### Linear Models

#### Victor Panaretos

Institut de Mathématiques - EPFL

 ${\tt victor.panaretos@epfl.ch}$ 



What is a Regression Model?

Statistical model for:

• Y (random variable)  $\stackrel{\text{depending on}}{\longleftarrow} x$  (non-random variable)

Aim: understand the effect of x on the random quantity Y

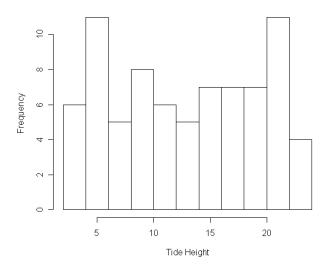
General formulation<sup>1</sup>:

$$Y \sim \mathsf{Distribution}\{g(x)\}$$

Statistical Problem: Estimate (learn)  $g(\cdot)$  from data  $\{(x_i, y_i)\}_{i=1}^n$ . Use for:

- Description
- Inference
- Prediction
- Data compression (parsimonious representations)
- . . .

 $<sup>^1</sup>$ Often books/people write  $Y\mid x\sim {\sf Distribution}\{g(x)\}$  but this implies that (X,Y) have a joint distribution; this assumption is unnecessary (e.g., in a designed experiment we choose values for x). Despite this, we write  $Y\mid x$  to remind ourselves that the distribution of Y depends on x.



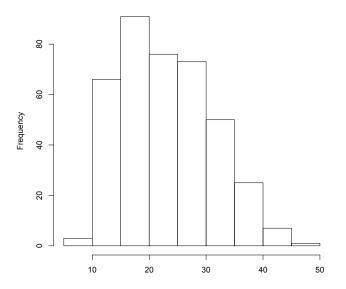
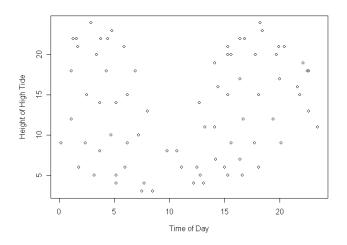
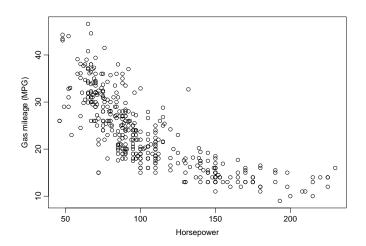


Figure: Miles per gallon for 392 car models





#### Great Variety of Models

Remember general model:

$$Y \sim \mathsf{Distribution}\{g(x)\}$$

#### x can be:

- continuous, discrete, categorical, vector . . .
- arrive randomly, or be chosen by experimenter, or both
- ullet however x arises, we treat it as constant in the analysis

#### Distribution can be:

• Gaussian (Normal), Laplace, binomial, Poisson, gamma, General exponential family, . . .

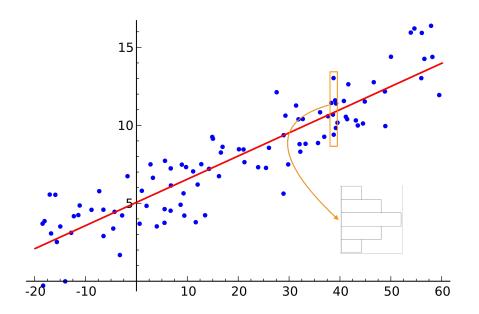
#### Function $g(\cdot)$ can be:

$$ullet g(x)=eta_0+eta_1 x$$
,  $g(x)=\sum_{k=-K}^Keta_k e^{-ikx}$ , Cubic spline,  $\dots$ 

ullet  $Y,x\in\mathbb{R},\ g(x)=eta_0+eta_1x$ , Distribution = Gaussian

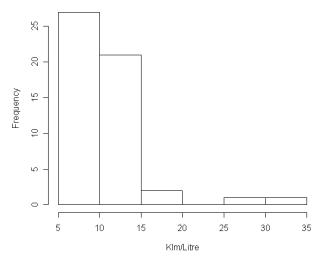
The second verson is useful for mathematical work, but is puzzling statistically, since we don't observe  $\epsilon$ .

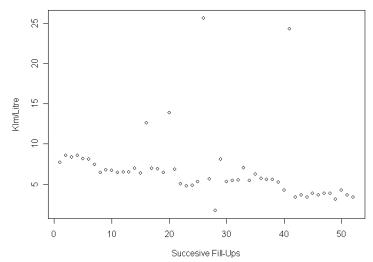
• Also, x could be vector  $(Y, \beta_0 \in \mathbb{R}, x \in \mathbb{R}^p, \beta \in \mathbb{R}^p)$ :





10 / 309





Tools of the trade ...

Start from Normal linear model  $\longrightarrow$  gradually generalise ... Important features of Normal linear model:

- Gaussian distribution
- Linearity

These two combine well and give geometric insights to solve the estimation problem. Thus we need to revise some linear algebra and probability ...

Will base course on the Gaussian assumption, but relax linearity later:

- linear Gaussian regression
- nonlinear Gaussian regression
- nonparametric Gaussian regression

Many further generalisations are possible . . .

Projections, Spectra, Gaussian Law

If Q is an  $n \times p$  real matrix, we define the *column space* (or *range*) of Q to be the set spanned by its columns:

$$\mathcal{M}(Q) = \{ y \in \mathbb{R}^n : \exists \beta \in \mathbb{R}^p, \ y = Q\beta \}.$$

- Recall that  $\mathcal{M}(Q)$  is a subspace of  $\mathbb{R}^n$ .
- ullet The columns of Q provide a coordinate system for the subspace  $\mathfrak{M}(Q)$
- If Q is of full column rank (p), then the coordinates  $\beta$  corresponding to a  $y \in \mathcal{M}(Q)$  are unique.
- Allows interpretation of system of linear equations

$$Q\beta = y$$
.

[existence of solution  $\leftrightarrow$  is y an element of  $\mathcal{M}(Q)$ ?] [uniqueness of solution  $\leftrightarrow$  is there a unique coordinate vector  $\beta$ ?]

Two further important subspaces associated with a real  $n \times p$  matrix Q:

ullet the null space (or kernel),  $\ker(Q)$ , of Q is the subspace defined as

$$\ker(Q) = \{x \in \mathbb{R}^p : Qx = 0\};$$

ullet the orthogonal complement of  $\mathcal{M}(\mathcal{Q}),\,\mathcal{M}^{\perp}(\mathcal{Q}),$  is the subspace defined as

$$\mathcal{M}^{\perp}(Q) = \{ y \in \mathbb{R}^n : y^{\top}Qx = 0, \ \forall x \in \mathbb{R}^p \}$$
  
=  $\{ y \in \mathbb{R}^n : y^{\top}v = 0, \ \forall v \in \mathcal{M}(Q) \}.$ 

The orthogonal complement may be defined for arbitrary subspaces by using the second equality.

## Theorem (Singular Value Decomposition)

Any  $n \times p$  real matrix can be factorised as

$$Q_{n \times p} = U_{n \times n} \sum_{n \times p} V_{p \times p}^{\top},$$

where U and  $V^{\top}$  are orthogonal with columns called left singular vectors and right singular vectors, respectively, and  $\Sigma$  is diagonal with non-negative real entries called singular values.

- **①** The left singular vectors corresponding to non-zero singular values form an orthonormal basis for  $\mathcal{M}(Q)$ .
- **2** The left singular vectors corresponding to zero singular values form an orthonormal basis for  $\mathcal{M}^{\perp}(Q)$ .

Since the statement is invariant to transposition, assume wlog that  $n \geq p$ . We will prove the statement by induction on p. Assume that p=1 so that Q is a column vector. Then the statement holds true trivially, by taking

$$U_{n\times 1} = Q/\|Q\|$$
  $\Sigma_{1\times 1} = \|Q\|$   $V^{\top} = V = 1.$ 

Thus the statement is true for all  $n \ge p$  when p = 1. This is the base case for our induction. For the inductive step, assume that the statement is true for some p > 1 and all  $n \ge p$ . Let us prove that it is also true for p + 1 and all  $n \ge p + 1$ .

Let  $\mathbb{S}^{p+1}=\{x\in\mathbb{R}^{p+1}:\|x\|=1\}$  and  $q(x)=\|Qx\|$ . Since  $q(\cdot)$  is continuous and  $\mathbb{S}^{p+1}$  is compact, we have that q(x) is bounded over  $\mathbb{S}^{p+1}$  and attains its bounds. So there exists  $v_1\in S^{p+1}$  such that

$$q(v_1) = \max_{x \in \mathbb{S}^{p+1}} q(x) = \sigma_1 < \infty.$$

and let  $v_1 \in \mathbb{S}^{p+1}$  be maximiser of q(x), i.e. such that  $q(v_1) = \max_{x \in \mathbb{S}^{p+1}} q(x)$ . Define  $u_1 = \sigma_1^{-1} Q v_1$  so  $\|u_1\| = 1$ . Given any orthonormal bases  $\{u_j\}_{j=2}^n$  for  $\operatorname{span}^{\perp}(u_1)$  and  $\{v_j\}_{j=2}^p$  for  $\operatorname{span}^{\perp}(v_1)$  define U and V to be orthogonal matrices

$$U = (u_1 \ u_2 \dots u_n) = (u_1 \ U_1)$$
 &  $V = (v_1 \ v_2 \dots v_n) = (v_1 \ V_1)$ .

Using block matrix multiplication, we see that

$$\begin{array}{c} U_{n\times n}^{\top} Q V_{n\times n} = \begin{pmatrix} u_{1}^{\top} \\ U_{1}^{\top} \end{pmatrix} Q \begin{pmatrix} v_{1} & V_{1} \end{pmatrix} = \begin{pmatrix} u_{1}^{\top} Q v_{1} & u_{1}^{\top} Q V_{1} \\ U_{1}^{\top} Q v_{1} & U_{1}^{\top} Q V_{1} \end{pmatrix} \\ = \begin{pmatrix} \sigma_{1} & \theta^{\top} \\ 1\times 1 & 1\times p \\ 0 & Z \\ (p-1)\times 1 & (p-1)\times p \end{pmatrix}. \end{array}$$

Now we claim that  $\theta = 0$ . To see this, first observe that

$$\sigma_1 = \max_{x \in \mathbb{S}^{p+1}} \|Qx\| = \max_{x \in \mathbb{S}^{p+1}} \|U^ op Qx\| = \max_{x \in \mathbb{S}^{p+1}} \|U^ op QVx\|.$$

Next, let's consider the norm of  $U^{\top}QV\left( \begin{array}{c} \sigma_1 \\ \theta \end{array} \right)$ ,

$$\begin{aligned} \left\| \begin{pmatrix} \sigma_1 & \theta^\top \\ \mathbf{0} & Z \end{pmatrix} \begin{pmatrix} \sigma_1 \\ \theta \end{pmatrix} \right\| &= \left\| \begin{pmatrix} \sigma_1^2 + \theta^\top \theta \\ Z \theta \end{pmatrix} \right\| &= \sqrt{(\sigma_1^2 + \theta^\top \theta)^2 + \|Z \theta\|^2} \\ &\geq \sigma_1^2 + \theta^\top \theta = (\sigma_1^2 + \theta^\top \theta)^{1/2} \left\| \begin{pmatrix} \sigma_1 \\ \theta \end{pmatrix} \right\|. \end{aligned}$$

Dividing across by  $\|(\sigma_1 \ \theta)^\top\|$ , we see that we must necessarily have  $(\sigma_1^2 + \theta^\top \theta)^{1/2} \leq \max_{x \in \mathbb{R}^{n+1}} \|U^\top Q V x\| = \sigma_1 = (\sigma_1^2 + 0)^{1/2}.$ 

and so it must be that 
$$\theta^{\top}\theta=0$$
. We conclude that

$$U^{ op} \, Q \, V = \left( egin{array}{cc} \sigma_1 & \mathbf{0}_{1 imes p} \ \mathbf{0}_{(n-1) imes 1} & Z \end{array} 
ight) \, \stackrel{\mathit{thus}}{\Longrightarrow} \, \, Q = U \left( egin{array}{cc} \sigma_1 & \mathbf{0}_{1 imes p} \ \mathbf{0}_{(n-1) imes 1} & Z \end{array} 
ight) \, V^{ op}.$$

But Z is an  $(n-1)\times p$  matrix, and since  $n\geq p+1$  it holds that  $n-1\geq p$ . So by our inductive hypothesis

 $Z_{(n-1)\times p} = W_{(n-1)\times (n-1)}\Omega_{(n-1)\times p}R_{n\times p}^{\top}$ 

where W, R are orthogonal and  $\Omega$  is diagonal. Thus

$$Q_{n\times p} = U_{n\times n} \begin{pmatrix} \sigma_1 & \mathbf{0}_{1\times p} \\ \mathbf{0}_{(n-1)\times 1} & W\Omega R^{\top} \end{pmatrix} V_{p\times p}^{\top} =$$

$$= \underbrace{U \begin{pmatrix} 1 & \mathbf{0}_{1\times (n-1)} \\ \mathbf{0}_{(n-1)\times 1} & W_{(n-1)\times (n-1)} \end{pmatrix}}_{W_{(n-1)\times (n-1)}} \underbrace{\begin{pmatrix} \sigma_1 & \mathbf{0}_{1\times p} \\ \mathbf{0}_{(n-1)\times 1} & \Omega_{(n-1)\times p} \end{pmatrix}}_{C_{(n-1)\times 1}} \underbrace{\begin{pmatrix} 1 & \mathbf{0}_{1\times p} \\ \mathbf{0}_{p\times 1} & R_{p\times p}^{\top} \end{pmatrix}}_{V} V^{\top}$$

orthogonal diagonal orthogonal

Victor Panaretos (EPFL) Linear Models 20 / 309

## Theorem (Spectral Theorem)

A  $p \times p$  matrix A is symmetric if and only if there exists a  $p \times p$  orthogonal matrix U and a real diagonal matrix  $\Lambda$  such that

$$A = U \Lambda U^{\top}.$$

In particular:

• the columns of  $U = (u_1 \cdots u_p)$  are **eigenvectors** of A, i.e.

$$Au_j=\lambda_ju_j, \qquad j=1,\ldots,p$$

where  $diag(\lambda_1, \dots, \lambda_p) = \Lambda$  are the corresponding (real) **eigenvalues** of A.

- $oldsymbol{0}$  the rank of A is the number of non-zero eigenvalues.
- if the eigenvalues are distinct, the eigenvectors are unique (up to re-ordering and sign flips).

### Proof.

If A = 0, the statement holds trivially, so let  $A = A^{\top} \neq 0$ .

First note that the SVD of A guarantees the existence of a singular vector pair (u,v) with non-zero singular value  $\sigma$ , so that

$$A(v+u) = Av + Au = Av + A^{\top}u = \sigma u + \sigma v = \sigma(u+v).$$

hence w = u + v is an eigenvector of A with real eigenvalue  $\sigma$ .

Now the theorem is obviously true for  $1 \times 1$  matrices (scalars). So use induction. Assume any non-zero  $p \times p$  symmetric matrix satisfies the theorem statement.

Let  $A = A^{\top} \neq 0$  be  $(p+1) \times (p+1)$ . By (1), A has at least one eigenvector  $w \in \mathbb{R}^p$  with real eigenvalue  $\sigma \neq 0$ .

Let  $W=(w\ R)$  where R has p orthonormal columns spanning  $\operatorname{span}^{\perp}(w)$ . Then

$$W^{\top}AW = \left(\begin{array}{c} w^{\top} \\ R^{\top} \end{array}\right) A \left(\begin{array}{cc} w & R \end{array}\right) = \left(\begin{array}{cc} w^{\top}Aw & w^{\top}AR \\ R^{\top}Aw & R^{\top}AR \end{array}\right)$$

$$= \left( \begin{array}{cc} \sigma^2 & (Aw)^\top R \\ R^\top Aw & R^\top AR \end{array} \right) = \left( \begin{array}{cc} \sigma^2 & \mathbf{0}_{1\times p} \\ \mathbf{0}_{p\times 1} & R^\top AR \end{array} \right) = \left( \begin{array}{cc} \sigma^2 & \mathbf{0}_{1\times p} \\ \mathbf{0}_{p\times 1} & B \end{array} \right)$$

where  $B = R^{\top}AR$  is a symmetric  $p \times p$  matrix.

Victor Panaretos (EPFL) Linear Models 22 / 309

Since B is symmetric, we have  $B=V\Omega V^{\top}$  for  $V_{p\times p}$  orthogonal and  $\Omega_{p\times p}$  diagonal by our induction hypothesis. In summary

$$A = W \begin{pmatrix} \sigma^2 & \mathbf{0}_{1 \times p} \\ \mathbf{0}_{p \times 1} & B \end{pmatrix} W^\top$$

$$= \underbrace{W \begin{pmatrix} 1 & \mathbf{0}_{1 \times p} \\ \mathbf{0}_{p \times 1} & V_{p \times p} \end{pmatrix}}_{\text{orthogonal}} \underbrace{\begin{pmatrix} \sigma_1 & \mathbf{0}_{1 \times p} \\ \mathbf{0}_{p \times 1} & \Omega_{p \times p} \end{pmatrix}}_{\text{diagonal}} \underbrace{\begin{pmatrix} 1 & \mathbf{0}_{1 \times p} \\ \mathbf{0}_{p \times 1} & V_{p \times p}^\top \end{pmatrix}}_{\text{orthogonal}} W_{p \times p}^\top$$

$$= U \Lambda U^\top$$

Combining the SVD and the spectral theorem, we notice that:

- The left singular vectors of Q are eigenvectors of  $A = QQ^{\top}$ .
- **2** The right singular vectors of Q are eigenvectors of  $A = Q^{T}Q$ .
- **1** The squared singular values of Q are eigenvalues of both  $QQ^{\top}$  and  $Q^{\top}Q$ .

(ロ) (部) (室) (室) (室) (室) (の)

Victor Panaretos (EPFL) Linear Models 23 / 309

A matrix Q is called *idempotent* if  $Q^2 = Q$ .

An orthogonal projection (henceforth projection) onto a subspace  $\mathcal V$  is a symmetric idempotent matrix H such that  $\mathcal M(H)=\mathcal V$ .

## Proposition

The only possible eigenvalues of a projection matrix are 0 and 1.

### Proposition

Let  $\mathcal V$  be a subspace and H be a projection onto  $\mathcal V$ . Then I-H is the projection matrix onto  $\mathcal V^\perp$ .

### Proof.

$$(I-H)^{\top} = I - H^{\top} = I - H$$
 since  $H$  is symmetric and,  
 $(I-H)^2 = I^2 - 2H + H^2 = I - H$ . Thus  $I-H$  is a projection matrix.

It remains to identify the column space of I-H. Let  $H=U\Lambda U^{\top}$  be the spectral decomposition of H. Then  $I-H=UU^{\top}-U\Lambda U^{\top}=U(I-\Lambda)U^{\top}$ .

Hence the column space of I-H is spanned by the eigenvectors of H corresponding to zero eigenvalues of H, which coincides with  $\mathcal{M}^{\perp}(H) = \mathcal{V}^{\perp}$ .

### Proposition

Let  $\mathcal V$  be a subspace and H be a projection onto  $\mathcal V$ . Then Hy=y for all  $y\in \mathcal V$ .

### Proposition

If P and Q are projection matrices onto a subspace V, then P=Q.

### Proposition

If  $x_1, \ldots, x_p$  are linearly independent and are such that  $span(x_1, \ldots, x_p) = \mathcal{V}$ , then the projection onto  $\mathcal{V}$  can be represented as

$$H = X(X^{\top}X)^{-1}X^{\top}$$

where X is a matrix with columns  $x_1, \ldots, x_p$ .

## Proposition

Let V be a subspace of  $\mathbb{R}^n$  and H be a projection onto V. Then

$$||x - Hx|| \le ||x - v||, \quad \forall v \in \mathcal{V}.$$

#### Proof

Let  $H = U\Lambda U^{\top}$  be the spectral decomposition of H,  $U = (u_1 \cdots u_n)$  and  $\Lambda = \operatorname{diag}(\lambda_1, \dots, \lambda_n)$ . Letting  $p = \dim(\mathcal{V})$ ,

- $\mathbf{0} \ \lambda_1 = \cdots = \lambda_p = 1 \ \text{and} \ \lambda_{p+1} = \cdots = \lambda_n = 0,$
- $u_1, \ldots, u_n$  is an orthonormal basis of  $\mathbb{R}^n$ ,
- $u_1, \ldots, u_p$  is an an orthonormal basis of  $\mathcal{V}$ .

### (proof continued)

$$\begin{split} \|x-Hx\|^2 &= \sum_{i=1} (x^\top u_i - (Hx)^\top u_i)^2 \qquad \text{[orthonormal basis]} \\ &= \sum_{i=1}^n (x^\top u_i - x^\top H u_i)^2 \qquad [H \text{ is symmetric]} \\ &= \sum_{i=1}^n (x^\top u_i - \lambda_i x^\top u_i)^2 \qquad [u' \text{s are eigenvectors of } H] \\ &= 0 + \sum_{i=p+1}^n (x^\top u_i)^2 \qquad \text{[eigenvalues 0 or 1]} \\ &\leq \sum_{i=1}^p (x^\top u_i - v^\top u_i)^2 + \sum_{i=p+1}^n (x^\top u_i)^2 \qquad \forall v \in \mathcal{V} \\ &= \|x-v\|^2. \end{split}$$

### Proposition

Let  $\mathcal{V}_1 \subseteq \mathcal{V} \subseteq \mathbb{R}^n$  be two nested linear subspaces. If  $H_1$  is the projection onto  $\mathcal{V}_1$  and H is the projection onto  $\mathcal{V}$ , then

$$HH_1 = H_1 = H_1 H.$$

#### Proof.

First we show that  $HH_1=H_1$ , and then that  $H_1H=HH_1$ . For all  $y\in\mathbb{R}^n$  we have  $H_1y\in\mathcal{V}_1$ . But then  $H_1y\in\mathcal{V}$ , since  $\mathcal{V}_1\subset\mathcal{V}$ .

Therefore  $HH_1y=H_1y$ . We have shown that  $(HH_1-H_1)y=0$  for all  $y\in\mathbb{R}^n$ , so that  $HH_1-H_1=0$ , as its kernel is all  $\mathbb{R}^n$ . Hence  $HH_1=H_1$ .

(Or, take n linearly independent vectors  $y_1, \ldots, y_n \in \mathbb{R}^n$ , and use them as columns of the  $n \times n$  matrix Y. Now Y is invertible, and  $(HH_1 - H_1)Y = 0$ , so  $HH_1 - H_1 = 0$ , giving  $HH_1 = H_1$ .)

To prove that  $H_1H=HH_1$ , note that symmetry of projection matrices and the first part of the proof give

$$H_1 H = H_1^{\top} H^{\top} = (HH_1)^{\top} = (H_1)^{\top} = H_1 = HH_1.$$

Victor Panaretos (EPFL) Linear Models 28 / 309

# Definition (Non-Negative Matrix – Quadratic Form Definition)

A  $p \times p$  real symmetric matrix  $\Omega$  is called non-negative definite (written  $\Omega \succeq 0$ ) if and only if  $x^\top \Omega x \geq 0$  for all  $x \in \mathbb{R}^p$ . If  $x^\top \Omega x > 0$  for all  $x \in \mathbb{R}^p \setminus \{0\}$ , then we call  $\Omega$  positive definite (written  $\Omega \succ 0$ ).

An equivalent definition is:

# Definition (Non-Negative Matrix - Spectral Definition)

A  $p \times p$  real symmetric matrix  $\Omega$  is called non-negative definite (written  $\Omega \succeq 0$ ) if and only the eigenvalues of  $\Omega$  are non-negative. If the eigenvalues of  $\Omega$  are strictly positive, then  $\Omega$  is called positive definite (written  $\Omega \succ 0$ ).

### Lemma (Exercise)

Prove that the two definitions are equivalent.

### Definition (Covariance Matrix)

Let  $Y=(Y_1,\ldots,Y_n)^{\top}$  be a random  $n\times 1$  vector such that  $\mathbb{E}\|Y\|^2<\infty$ . The covariance matrix of Y, say  $\Omega$ , is the  $n\times n$  symmetric matrix with entries

$$\Omega_{ij} = \operatorname{cov}(Y_i, Y_j) = \mathbb{E}[(Y_i - \mathbb{E}[Y_i])(Y_j - \mathbb{E}[Y_j])], \quad 1 \leq i \leq j \leq n.$$

That is, the covariance matrix encodes the variances of the coordinates of Y (on the diagonal) and the covariances between the coordinates of Y (off the diagonal). If we write

$$\mu = \mathbb{E}[\,Y] = (\mathbb{E}[\,Y_1], \ldots, \mathbb{E}[\,Y_n])^ op$$

for the mean vector of Y, then the covariance matrix of Y can be written as

$$\mathbb{E}[(Y-\mu)(Y-\mu)^{\top}] = \mathbb{E}[YY^{\top}] - \mu\mu^{\top}.$$

Whenever Y is a random vector, we will write cov(Y) or var(Y) for the covariance matrix of Y.

#### Covariance Matrices

#### Lemma

Let Y be a random  $d \times 1$  vector such that  $\mathbb{E}||Y||^2 < \infty$ . Let  $\mu$  be the mean vector and  $\Omega$  be the covariance matrix of Y. If A is a  $p \times d$  real matrix, the mean vector and covariance matrix of AY are  $A\mu$  and  $A\Omega A^{\top}$ , respectively.

### Proof.

Exercise.

# Corollary (Covariance of Projections)

Let Y be a random  $d \times 1$  vector such that  $\mathbb{E}||Y||^2 < \infty$ . Let  $\beta, \gamma \in \mathbb{R}^d$  be fixed vectors. If  $\Omega$  denotes the covariance matrix of Y,

- the variance of  $\beta^{\top} Y$  is  $\beta^{\top} \Omega \beta$ ;
- the covariance of  $\beta^{\top} Y$  with  $\gamma^{\top} Y$  is  $\gamma^{\top} \Omega \beta$ .

Non-negative Matrices  $\equiv$  Covariance Matrices

# Proposition (Non-Negative and Covariance Matrices)

Let  $\Omega$  be a real symmetric matrix. Then  $\Omega$  is non-negative definite if and only if  $\Omega$  is the covariance matrix of some random variable Y.

### Proof.

Exercise.



### Principal Component Analysis

- Let Y be a random vector in  $\mathbb{R}^d$  with covariance matrix  $\Omega$ .
- Find direction  $v_1 \in \mathbb{S}^{d-1}$  such that the projection of Y onto  $v_1$  has maximal variance.
- For  $j=2,3,\ldots,d$ , find direction  $v_j\perp v_{j-1}$  such that projection of Y onto  $v_j$  has maximal variance.

Solution: maximise  $\mathsf{Var}(v_1^ op Y) = v_1^ op \Omega v_1$  over  $\|v_1\| = 1$ 

$$v_1^ op \Omega v_1 = v_1^ op U \Lambda U^ op v_1 = \|\Lambda^{1/2} U^ op v_1\|^2 = \sum_{i=1}^d \lambda_i (u_i^ op v_1)^2$$
 [change of basis]

Now  $\sum_{i=1}^d (u_i^ op v_1)^2 = \|v_1\|^2 = 1$  so we have a convex combination of the  $\{\lambda_j\}_{j=1}^d$ ,

$$\sum_{i=1}^d p_i \lambda_i, \qquad \sum_i p_i = 1, \quad p_i \geq 0, \quad i = 1, \ldots, d.$$

But  $\lambda_1 \geq \lambda_i \geq 0$  so clearly this sum is maximised when  $p_1=1$  and  $p_j=0$   $\forall j \neq 1$ , i.e.  $v_1=\pm u_1$ .

Iteratively,  $v_j=\pm u_j$ , i.e. principal components are eigenvectors of  $\Omega_{i,j}$ 

33 / 309

## Theorem (Optimal Linear Dimension Reduction Theorem)

Let Y be a mean-zero random variable in  $\mathbb{R}^n$  with  $n \times n$  covariance  $\Omega$ . Let H be the projection matrix onto the span of the first k eigenvectors of  $\Omega$ . Then

$$\mathbb{E}||Y - HY||^2 \le \mathbb{E}||Y - QY||^2$$

for any  $n \times n$  projection operator Q or rank at most k.

Intuitively: if you want to approximate a mean-zero random variable taking values  $\mathbb{R}^n$  by a random variable that ranges over a subspace of dimension at most  $k \leq n$ , the optimal choice is the projection of the random variable onto the space spanned by its first k principal components (eigenvectors of the covariance). "Optimal" is with respect to the mean squared error.

For the proof, use lemma below (follows immediately from spectral decomposition)

#### Lemma

Q is a rank k projection matrix if and only if there exist orthonormal vectors  $\{v_j\}_{j=1}^k$  such that  $Q = \sum_{j=1}^k v_j v_j^{\top}$ .

### Optimal Linear Dimension Reduction.

Write  $Q = \sum_{i=1}^k v_i v_i 1^{\top}$  for some orthonormal  $\{v_i\}_{i=1}^k$ . Then,

$$\begin{split} \mathbb{E} \| \, Y - QY \|^2 &= \mathbb{E} \left[ \, Y^\top (I - Q)^\top (I - Q) \, Y \, \right] = \mathbb{E} \left[ \operatorname{tr} \{ (I - Q) \, YY^\top (I - Q)^\top \} \right] \\ &= \operatorname{tr} \{ (I - Q) \mathbb{E} \left[ \, YY^\top \right] (I - Q)^\top \} = \operatorname{tr} \{ (I - Q)^\top (I - Q) \Omega \} \\ &= \operatorname{tr} \{ (I - Q) \Omega \} = \operatorname{tr} \{ \Omega \} - \operatorname{tr} \{ \, Q\Omega \} = \sum_{i=1}^n \lambda_i - \operatorname{tr} \left\{ \sum_{j=1}^k v_j \, v_j^\top \Omega \right\} \\ &= \sum_{i=1}^n \lambda_i - \sum_{j=1}^k \operatorname{tr} \left\{ v_j \, v_j^\top \Omega \right\} = \sum_{i=1}^n \lambda_i - \sum_{j=1}^k v_j^\top \Omega v_j \\ &= \sum_{i=1}^n \lambda_i - \sum_{j=1}^k \operatorname{Var} [v_j^\top Y] \end{split}$$

If we can minimise this expression over all  $\{v_j\}_{j=1}^k$  with  $v_j^\top v_{j'} = \mathbf{1}\{j=j'\}$ , then we're done. By PCA, this is done by choosing the top k eigenvectors of  $\Omega$ .

## Corollary

Let  $\{x_1,...,x_p\}\subset\mathbb{R}^n$  be such that  $x_1+...+x_p=0$ , and let X be the  $n\times p$  matrix with columns  $\{x_j\}_{j=1}^p$ . The best approximating k-hyperplane to the points  $\{x_1,...,x_p\}$  is given by the span of the k leading eigenvectors of the matrix  $XX^\top$ , i.e. if H is the projection onto this span, it holds that

$$\sum_{j=1}^p \|x_j - Hx_j\|^2 \le \sum_{j=1}^p \|x_j - Qx_j\|^2$$

for any  $n \times n$  projection operator Q or rank at most k.

### Proof.

Define a discrete random vector Y by  $\mathbb{P}[Y=x_j]=1/p, j\in\{1,...,p\}$  and observe that  $\mathbb{E}[h(Y)]=p^{-1}\sum_{j=1}^p h(x_j)$ , for any vector-valued (or matrix-valued) deterministic map h. Now use the optimal linear dimension reduction theorem.

# Definition (Multivariate Gaussian Distribution)

A random vector Y in  $\mathbb{R}^d$  has the multivariate normal distribution if and only if  $\beta^\top Y$  has the univariate normal distribution,  $\forall \beta \in \mathbb{R}^d$ .

Observation: From the definition if follows that Y must have some well-defined mean vector  $\mu$  and some well defined covariance matrix  $\Omega$ .

To see this note that since  $\mathbb{E}\{(\beta^\top Y)^2\} < \infty$  for all  $\beta$ , then we can successively pick  $\beta$  to be equal to each canonical basis vector and conclude that each coordinate has finite variance and thus  $\mathbb{E}\|Y\|^2 < \infty$ .

So all the means, variances and covariances of its coordinates are well defined.

Then, the mean vector (say)  $\mu$  and covariance matrix (say)  $\Omega$  can be (uniquely) determined entrywise by equating

$$\mu_i = \mathbb{E}[e_i^{ op} Y] \qquad \& \qquad \Omega_{ij} = \mathsf{cov}\{e_i^{ op} Y, \ e_j^{ op} Y\}.$$

where  $e_j$  is the jth canonical basis vector

$$e_j = (0 \ , \ 0 \ , \dots, \ 1 \ , \dots, \ 0 \ , \ 0)^ op$$

$$j_{th} \ \ \text{position}$$

#### How can we use this definition to determine basic properties?

The moment generating function (MGF) of a random vector W in  $\mathbb{R}^d$  is defined as

$${M}_{W}( heta) = \mathbb{E}[e^{ heta^ op W}], \qquad heta \in \mathbb{R}^d,$$

provided the expectation exists. When the MGF exists *it characterises the distribution of the random vector*. Furthermore, two random vectors are independent if and only if their joint MGF is the product of their marginal MGF's, i.e.

$$X_{n \times 1}$$
 independent of  $Y_{m \times 1}$ 

$$\iff$$

$$ig| \mathbb{E}[e^{eta^ op X + m{\gamma}^ op Y}] = \mathbb{E}[e^{m{eta}^ op X}] imes \mathbb{E}[e^{m{\gamma}^ op Y}], \qquad orall \, eta \in \mathbb{R}^n \ \& \, m{\gamma} \in \mathbb{R}^m$$

#### Gaussian Vectors and Affine Transformations

#### Useful facts:

**①** Moment generating function of  $Y \sim \mathcal{N}(\mu, \Omega)$ :

$$M_Y(u) = \exp\left(u^ op \mu + rac{1}{2} u^ op \Omega u
ight).$$

- $\mathbf{2} \ \ Y \sim \mathcal{N}(\mu_{p \times 1}, \Omega_{p \times p}) \ \text{and given } B_{n \times p} \ \text{and } \theta_{n \times 1}, \ \text{then}$   $\theta + BY \sim \mathcal{N}(\theta + B\mu, B\Omega B^{\top}).$
- **3**  $\mathcal{N}(\mu, \Omega)$  density, assuming  $\Omega$  nonsingular:

$$f_Y(y) = rac{1}{\left(2\pi
ight)^{p/2} |\Omega|^{1/2}} \exp\left\{-rac{1}{2}(y-\mu)^ op \Omega^{-1}(y-\mu)
ight\}.$$

- Constant density isosurfaces are ellipsoidal
- Marginals of Gaussian are Gaussian (converse NOT true).
- **1**  $\Omega$  diagonal  $\Leftrightarrow$  independent coordinates  $Y_j$ .
- $m{O}$  If  $Y \sim \mathcal{N}(\mu_{p \times 1}, \Omega_{p \times p})$ ,

AY independent of  $BY \iff A\Omega B^{\top} = 0$ .

# Proposition (Property 1: Moment Generating Function)

The moment generating function of  $Y \sim \mathcal{N}(\mu, \Omega)$  is

$$M_Y(u) = \exp\left(u^ op \mu + rac{1}{2}u^ op \Omega u\right)$$

## Proof.

Let  $v \in \mathbb{R}^d$  be arbitrary. Then  $v^\top Y$  is scalar Gaussian with mean  $v^\top \mu$  and variance  $v^\top \Omega v$ . Hence it has moment generating function:

$$M_{v^{ op}\,Y}(t) = \mathbb{E}\left(\left.e^{tv^{ op}\,Y}
ight) = \exp\left\{t(v^{ op}\mu) + rac{t^2}{2}(v^{ op}\Omega v)
ight\}.$$

Now take t = 1 and observe that

$$M_{v^{\top}Y}(1) = \mathbb{E}\left(e^{v^{\top}Y}\right) = M_Y(v).$$

Combining the two, we conclude that

$$M_Y(v) = \exp\left(v^ op \mu + rac{1}{2}v^ op \Omega v
ight), \quad v \in \mathbb{R}^d.$$

## Proposition (Property 2: Affine Transformation)

For  $Y \sim \mathcal{N}(\mu_{p \times 1}, \Omega_{p \times p})$  and given  $B_{n \times p}$  and  $\theta_{n \times 1}$ , we have

$$\theta + BY \sim \mathcal{N}(\theta + B\mu, B\Omega B^{\top})$$

#### Proof.

$$\begin{split} M_{\theta+BY}(u) &= & \mathbb{E}\left[\exp\{u^{\top}(\theta+BY)\}\right] = \exp\left\{u^{\top}\theta\right\} \mathbb{E}\left[\exp\{(B^{\top}u)^{\top}Y\}\right] \\ &= & \exp\left\{u^{\top}\theta\right\} M_Y(B^{\top}u) \\ &= & \exp\left\{u^{\top}\theta\right\} \exp\left\{(B^{\top}u)^{\top}\mu + \frac{1}{2}u^{\top}B\Omega B^{\top}u\right\} \\ &= & \exp\left\{u^{\top}\theta + u^{\top}(B\mu) + \frac{1}{2}u^{\top}B\Omega B^{\top}u\right\} \\ &= & \exp\left\{u^{\top}(\theta+B\mu) + \frac{1}{2}u^{\top}B\Omega B^{\top}u\right\} \end{split}$$

And this last expression is the MGF of a  $\mathcal{N}(\theta + B\mu, B\Omega B^{\top})$  distribution.

# Proposition (Property 3: Density Function)

Let  $\Omega_{p \times p}$  be nonsingular. The density of  $\mathcal{N}(\mu_{p \times 1}, \Omega_{p \times p})$  is

$$f_Y(y) = rac{1}{(2\pi)^{p/2} |\Omega|^{1/2}} \exp\left\{-rac{1}{2} (y-\mu)^ op \Omega^{-1} (y-\mu)
ight\}$$

## Proof.

Let  $Z = (Z_1, \ldots, Z_p)^{\top}$  be a vector of iid  $\mathcal{N}(0, 1)$  random variables. Then, because of independence,

(a) the density of Z is

$$f_Z(z) = \prod_{i=1}^p f_{Z_i}(z_i) = \prod_{i=1}^p rac{1}{\sqrt{2\pi}} \exp\left(-rac{1}{2}z_i^2
ight) = rac{1}{(2\pi)^{p/2}} \exp\left(-rac{1}{2}z^ op z
ight).$$

(b) The MGF of Z is

$$M_Z(u) = \mathbb{E}\left\{\exp\left(\sum_{i=1}^p u_i Z_i
ight)
ight\} = \prod_{i=1}^p \mathbb{E}\{\exp(u_i Z_i)\} = \exp(u^ op u/2),$$

which is the MGF of a p-variate  $\mathcal{N}(0, I)$  distribution.

Victor Panaretos (EPFL) 42 / 309 Linear Models

## proof continued

$$\overset{(a)+(b)}{\Longrightarrow}$$
 the  $\mathcal{N}(0,I)$  density is  $f_Z(z) = rac{1}{(2\pi)^{p/2}} \exp\left(-rac{1}{2}z^ op z
ight)$ .

By the spectral theorem,  $\Omega$  admits a square root,  $\Omega^{1/2}$ . Furthermore, since  $\Omega$  is non-singular, so is  $\Omega^{1/2}$ .

Now observe that from our Property 2, we have  $Y \stackrel{d}{=} \Omega^{1/2}Z + \mu \sim \mathcal{N}(\mu, \Omega)$ . By the change of variables formula,

$$\begin{split} f_Y(y) &= f_{\Omega^{1/2}Z + \mu}(y) \\ &= |\Omega^{-1/2}| f_Z \{ \Omega^{-1/2}(y - \mu) \} \\ &= \frac{1}{(2\pi)^{p/2} |\Omega|^{1/2}} \exp \left\{ -\frac{1}{2} (y - \mu)^\top \Omega^{-1}(y - \mu) \right\}. \end{split}$$

[Recall that to obtain the density of W = g(X) at w, we need to evaluate  $f_X$  at  $g^{-1}(w)$  but also multiply by the Jacobian determinant of  $g^{-1}$  at w.]

# Proposition (Property 4: Isosurfaces)

The isosurfaces of a  $\mathcal{N}(\mu_{p\times 1},\Omega_{p\times p})$  are (p-1)-dimensional ellipsoids centred at  $\mu$ , with principal axes given by the eigenvectors of  $\Omega$  and with anisotropies given by the ratios of the square roots of the corresponding eigenvalues of  $\Omega$ .

## Proof.

Exercise: Use Property 3, and the spectral theorem.

# Proposition (Property 5: Coordinate Distributions)

Let  $Y=(Y_1,\ldots,Y_p)^{ op}\sim\mathcal{N}(\mu_{p imes 1},\Omega_{p imes p}).$  Then  $Y_j\sim\mathcal{N}(\mu_j,\Omega_{jj})$  .

## Proof.

Observe that  $Y_j = (0, 0, \dots, 1, \dots, 0, 0) Y$  and use Property 2.

# Proposition (Property 6: Diagonal $\Omega \iff$ Independence)

Let  $Y = (Y_1, \dots, Y_p)^\top \sim \mathcal{N}(\mu_{p \times 1}, \Omega_{p \times p})$ . Then the  $Y_i$  are mutually independent if and only if  $\Omega$  is diagonal.

## Proof.

Suppose that the  $Y_j$  are independent. Property 5 yields  $Y_j \sim \mathcal{N}(\mu_j, \sigma_j^2)$  for some  $\sigma_j > 0$ . Thus the density of Y is

$$f_{Y}(y) = \prod_{j=1}^{p} f_{Y_{j}}(y_{j}) = \prod_{i=1}^{p} \frac{1}{\sigma_{j}\sqrt{2\pi}} \exp\left\{-\frac{1}{2} \frac{(y_{j} - \mu_{j})^{2}}{\sigma_{j}^{2}}\right\}$$

$$= \frac{1}{\left(2\pi\right)^{p/2} |\mathsf{diag}(\sigma_1^2, \dots, \sigma_p^2)|^{1/2}} \exp\left\{-\frac{1}{2}(y-\mu)^\top \mathsf{diag}(\sigma_1^{-2}, \dots, \sigma_p^{-2})(y-\mu)\right\}.$$

Hence  $Y \sim \mathcal{N}\{\mu, \operatorname{diag}(\sigma_1^2, \dots, \sigma_p^2)\}$ , i.e. the covariance  $\Omega$  is diagonal.

Conversely, assume  $\Omega$  is diagonal, say  $\Omega=\operatorname{diag}(\sigma_1^2,\ldots,\sigma_p^2)$ . Then we can reverse the steps of the first part to see that the joint density  $f_Y(y)$  can be written as a product of the marginal densities  $f_{Y_j}(y_j)$ , thus proving independence.

# Proposition (Property 7: $AY, BY \text{ indep } \iff A\Omega B^{\top} = 0$ )

If  $Y \sim \mathcal{N}(\mu_{p \times 1}, \Omega_{p \times p})$ , and  $A_{m \times p}$ ,  $B_{d \times p}$  be real matrices. Then,

AY independent of  $BY \iff A\Omega B^{\top} = 0$ .

#### Proof

It suffices to prove the result assuming  $\mu=0$  (and it simplifies the algebra). First assume  $A\Omega B^{\top}=0$ . Let  $W_{(m+d)\times 1}=\binom{AY}{BY}$  and  $\theta_{(m+d)\times 1}=\binom{u_{m\times 1}}{y_{d\times 1}}$ .

$$egin{aligned} M_W( heta) &=& \mathbb{E}[\exp\{W^ op heta\}] = \mathbb{E}\left[\exp\left\{Y^ op A^ op u + Y^ op B^ op v
ight\}
ight] \ &=& \mathbb{E}\left[\exp\left\{Y^ op (A^ op u + B^ op v)
ight\}
ight] = M_Y(A^ op u + B^ op v) \ &=& \exp\left\{rac{1}{2}(A^ op u + B^ op v)^ op \Omega(A^ op u + B^ op v)
ight\} \ &=& \exp\left\{rac{1}{2}\left(u^ op A\Omega A^ op u + v^ op B\Omega B^ op v + u^ op \underline{A\Omega B^ op} v + v^ op \underline{B\Omega A^ op} u
ight) \ &=& M_{AY}(u)M_{BY}(v), \end{aligned}$$

i.e., the joint MGF is the product of the marginal MGFs, proving independence.

For the converse, assume that AY and BY are independent. Then,  $\forall u, v$ ,

$$M_W(\theta) = M_{AY}(u) M_{BY}(v), \quad \forall u, v,$$

$$\implies \exp\left\{\frac{1}{2}\left(u^{\top}A\Omega A^{\top}u + v^{\top}B\Omega B^{\top}v + u^{\top}A\Omega B^{\top}v + v^{\top}B\Omega A^{\top}u\right)\right\}$$
$$= \exp\left\{\frac{1}{2}u^{\top}A\Omega A^{\top}u\right\} \exp\left\{\frac{1}{2}v^{\top}B\Omega B^{\top}v\right\}$$

$$\implies \exp\left\{\frac{1}{2} \times 2u^{\top} A \Omega B^{\top} v\right\} = 1$$

$$\implies u^{\top} A \Omega B^{\top} v = 0, \qquad \forall \ u \in \mathbb{R}^d, v \in \mathbb{R}^m,$$

 $\implies$  the orthocomplement<sup>a</sup> of the column space of  $A\Omega B^{+}$  is the whole of  $\mathbb{R}^{m}$ .

$$\implies A\Omega B^{\top} = 0.$$

Victor Panaretos (EPFL)

Linear Models

arecall that for  $Q_{m imes d}$  we have  $\mathfrak{M}^\perp(Q) = \{y \in \mathbb{R}^m : y^ op Qx = 0, \ orall x \in \mathbb{R}^d\}$ 

Gaussian Quadratic Forms and the  $\chi^2$  Distribution

# Definition ( $\chi^2$ distribution)

Let  $Z \sim \mathcal{N}(0, I_{p \times p})$ . Then  $||Z||^2 = \sum_{j=1}^p Z_j^2$  is said to have the chi-square  $(\chi^2)$  distribution with p degrees of freedom; we write  $||Z||^2 \sim \chi_p^2$ .

[Thus,  $\chi_p^2$  is the distribution of the sum of squares of p real independent standard Gaussian random variates.]

## Definition (F distribution)

Let  $V \sim \chi_p^2$  and  $W \sim \chi_q^2$  be independent random variables. Then (V/p)/(W/q) is said to have the F distribution with p and q degrees of freedom; we write  $(V/p)/(W/q) \sim F_{p,q}$ .

## Proposition (Gaussian Quadratic Forms)

• If  $Z \sim \mathcal{N}(0_{p \times 1}, I_{p \times p})$  and H is a projection of rank  $r \leq p$ ,

$$Z^{\top}HZ \sim \chi_r^2$$
.

②  $Y \sim \mathcal{N}(\mu_{p \times 1}, \Omega_{p \times p})$  with  $\Omega$  nonsingular  $\Longrightarrow$ 

$$(Y - \mu)^{\top} \Omega^{-1} (Y - \mu) \sim \chi_p^2$$

Exercise: Prove these results.

What if the random vector is not Gaussian? Here's a CLT<sup>2</sup> that helps:

# Theorem (Hajék-Sidak Weighted Sum CLT)

Let  $\{X_n\}$  be an i.i.d sequence of real random variables, with common mean 0 and variance 1. Let  $\{\gamma_n\}$  be a sequence of real constants. Then,

$$\sup_{1 \leq j \leq n} \frac{\gamma_j^2}{\sum_{i=1}^n \gamma_i^2} \overset{n \to \infty}{\longrightarrow} 0 \Longrightarrow \frac{1}{\sqrt{\sum_{i=1}^n \gamma_i^2}} \sum_{i=1}^n \gamma_i X_i \overset{d}{\to} N(0,1).$$

- Supremum condition amounts to saying that, in the limit, any single component contributes a negligible proportion of the total variance.
- Coefficient sequence  $\{\gamma_n\}$  might very well diverge, without contradicting the negligibility condition (e.g.  $\gamma_k = \sqrt{k}$ )

Victor Panaretos (EPFL) Linear Models 50 / 309

<sup>&</sup>lt;sup>2</sup>Consequence of Lyapunov's CLT, see e.g. Sen & Singer, "Large Sample Methods in Statistics", Chapman & Hall, pp. 108-119.

Linear Models: Likelihood and Geometry

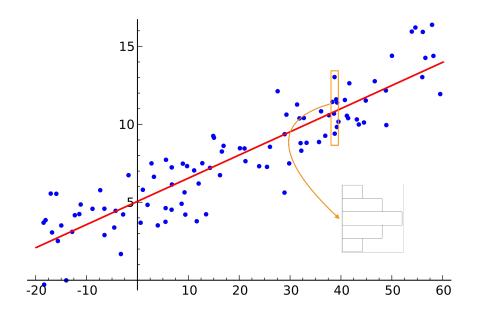
#### General formulation:

$$Y_i|x_i \stackrel{ind}{\sim} \mathsf{Distribution}\{g(x_i)\}, \quad i=1,\ldots,n.$$

### Simple Normal Linear Regression:

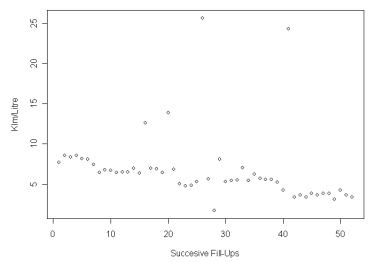
$$\left\{egin{array}{l} \mathsf{Distribution} = \mathcal{N}\{g(x), \sigma^2\} \ g(x) = eta_0 + eta_1 x \end{array}
ight.$$

## Resulting Model:



Example: Professor's Van

Fillup	Km/L
1	7.72
2	8.54
3	8.35
4	8.55
5	8.16
6	8.12
7	7.46
8	6.43
9	6.74
10	6.72



#### Simple Normal Linear Regression

Jargon: *Y* is response variable and *x* is explanatory variable (or covariate) Linearity: Linearity is in the parameters, not the explanatory variable.

Example: Flexibility in what we define as explanatory:

$$Y_j = eta_0 + eta_1 \underbrace{\sin(x_j)}_{x_j^*} + arepsilon_j, \quad arepsilon_j \stackrel{iid}{\sim} \mathsf{Normal}(0, \sigma^2).$$

Example: Sometimes a transformation may be required:

$$Y_j = eta_0 \, e^{eta_1 x_j} \eta_j, \quad \eta_j \stackrel{iid}{\sim} \mathsf{Lognormal}$$
  $\log(\cdot) \downarrow \qquad \uparrow \exp(\cdot)$ 

 $\log Y_j = \log eta_0 + eta_1 x_j + \log \eta_j, \quad \log \eta_j \stackrel{iid}{\sim} \mathsf{Normal}$ 

#### Data Structure:

For  $i=1,\ldots,n$ , pairs

$$(x_i, y_i) \longrightarrow \left\{egin{array}{l} x_i ext{ fixed values of } x \ y_i ext{ treated as a realisation of } Y_i ext{ at } x_i \end{array}
ight.$$

◆ロ → ◆団 → ◆ 豆 → ● ・ ● ・ りへの

Instead of  $x_i \in \mathbb{R}$  could have  $x_i^{\top} \in \mathbb{R}^q$ ):

$$Y_i = eta_0 + eta_1 x_{i1} + eta_2 x_{i2} + \ldots + eta_q x_{iq} + arepsilon_i, \quad arepsilon_i \stackrel{ind}{\sim} \mathcal{N}(0, \sigma^2).$$

Letting p = q + 1, this can be summarised via matrix notation:

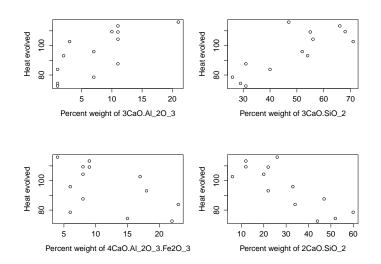
$$\underbrace{\begin{pmatrix} Y_1 \\ Y_2 \\ \vdots \\ Y_n \end{pmatrix}}_{Y} = \underbrace{\begin{pmatrix} 1 & x_{11} & \dots & x_{1q} \\ 1 & x_{21} & & x_{2q} \\ \vdots & \vdots & & \vdots \\ 1 & x_{n1} & \dots & x_{nq} \end{pmatrix}}_{X} \underbrace{\begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_q \end{pmatrix}}_{\beta} + \underbrace{\begin{pmatrix} \varepsilon_1 \\ \varepsilon_2 \\ \vdots \\ \varepsilon_n \end{pmatrix}}_{\varepsilon}$$

$$\Longrightarrow \underbrace{Y}_{n \times 1} = \underbrace{X}_{n \times p} \underbrace{\beta}_{n \times 1} + \underbrace{\varepsilon}_{n \times 1}, \quad \varepsilon \sim \mathcal{N}_n(0, \sigma^2 I)$$

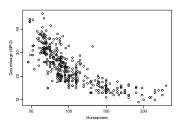
X is called the *design matrix*.

4日ト 4回ト 4 差ト 4 差ト 差 9000

Case	3 CaO. Al <sub>2</sub> O <sub>3</sub>	3 Ca O . Si O 2	4Cao. Al <sub>2</sub> O <sub>3</sub> . Fe <sub>2</sub> O <sub>3</sub>	2 CaO. SiO2	Heat
1	7.00	26.00	6.00	60.00	78.50
2	1.00	29.00	15.00	52.00	74.30
3	11.00	56.00	8.00	20.00	104.30
4	11.00	31.00	8.00	47.00	87.60
5	7.00	52.00	6.00	33.00	95.90
6	11.00	55.00	9.00	22.00	109.20
7	3.00	71.00	17.00	6.00	102.70
8	1.00	31.00	22.00	44.00	72.50
9	2.00	54.00	18.00	22.00	93.10
10	21.00	47.00	4.00	26.00	115.90
11	1.00	40.00	23.00	34.00	83.80
12	11.00	66.00	9.00	12.00	113.30
13	10.00	68.00	8.00	12.00	109.40



Example: polynomial terms for MPG vs Horsepower



Perhaps more fitting than

$$Y_j = \beta_0 + \beta_1 x_j + \varepsilon_j$$

would be

$$Y_j = \beta_0 + \beta_1 x_j + \beta_2 x_j^2 + \varepsilon_j$$

Still a linear model but now with 2 covariates:  $x_j$  and  $x_j^* = x_j^2$ 

- Normally would require a (hyper)plane to visualise dependence of mean on 2 or more covariates
- When additional covariates are variable transformation, can visualise mean dependence via a non-linear curve, even though model is linear

#### Model is:

$$Y_i = eta_0 + eta_1 x_{i1} + eta_2 x_{i2} + \dots + eta_q x_{iq} + arepsilon_i, \quad arepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

$$\updownarrow$$

$$Y = X \beta + arepsilon, \quad arepsilon \sim \mathcal{N}_n(0, \sigma^2 I)$$

Observe:  $y = (y_1, \dots, y_n)^{\top}$  for given fixed design matrix X, i.e.:

$$(y_1, x_{11}, \ldots, x_{1q}), \ldots, (y_i, x_{i1}, \ldots, x_{iq}), \ldots, (y_n, x_{n1}, \ldots, x_{nq})$$

### Likelihood and Loglikelihood

$$egin{aligned} L(eta,\sigma^2) &= rac{1}{(2\pi\sigma^2)^{n/2}} \exp\left\{-rac{1}{2\sigma^2}(y-Xeta)^ op (y-Xeta)
ight\} \ \ell(eta,\sigma^2) &= -rac{1}{2} \left\{n\log 2\pi + n\log \sigma^2 + rac{1}{\sigma^2}(y-Xeta)^ op (y-Xeta)
ight\} \end{aligned}$$

Whatever the value of  $\sigma$ , the log-likelihood is maximised when  $(y - X\beta)^{\top}(y - X\beta)$  is minimised. Hence, the MLE of  $\beta$  is:

$$\hat{\beta} = \argmax_{\beta} \left\{ -(y - X\beta)^{\top} (y - X\beta) \right\} = \arg\min_{\beta} (y - X\beta)^{\top} (y - X\beta)$$

Obtain minimum by solving:

$$0 = \frac{\partial}{\partial \beta} (y - X\beta)^{\top} (y - X\beta)$$

$$0 = \frac{\partial (y - X\beta)}{\partial \beta} \frac{\partial (y - X\beta)^{\top} (y - X\beta)}{\partial (y - X\beta)} \quad \text{(chain rule)}$$

$$0 = X^{\top} (y - X\beta) \quad \text{(normal equations)}$$

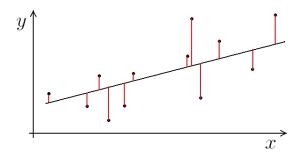
$$X^{\top} X\beta = X^{\top} y$$

$$\hat{\beta} = (X^{\top} X)^{-1} X^{\top} y \quad \text{(if } X \text{ has rank } p)$$

 $\hat{\beta}$  is called the *least squares estimator* because it is a result of minimising

$$(y - X\beta)^{\top}(y - X\beta) = \underbrace{\sum_{i=1}^{n} (y_i - \beta_0 - \beta_1 x_{i1} - \beta_2 x_{i2} - \dots - \beta_q x_{iq})^2}_{\text{sum of squares}}.$$

Thus we are trying to find the  $\beta$  that gives the hyperplane with minimum sum of squared vertical distances from our observations.



<u>Residuals</u>:  $e = y - X\hat{\beta}$ , so that  $e = (e_1, \dots, e_n)^{\top}$ , with

$$e_i = y_i - \hat{eta}_0 - \hat{eta}_1 x_{i1} - \hat{eta}_2 x_{i2} - \dots - \hat{eta}_q x_{iq}$$

"Regression Line" is such that  $\sum e_i^2$  is minimised over all  $\beta$ .

<u>Fitted Values</u>:  $\hat{y} = X \hat{\beta}^{\top}$ , so that  $\hat{y} = (\hat{y}_1, \dots, \hat{y}_n)^{\top}$ , with

$$\hat{y}_i = \hat{eta_0} + \hat{eta_1} x_{i1} + \cdots + \hat{eta_q} x_{iq}$$

Since the MLE of  $\beta$  is  $\hat{\beta} = (X^{\top}X)^{-1}X^{\top}y$  for all values of  $\sigma^2$ , we have

$$\begin{split} \hat{\sigma}^2 &= \underset{\sigma^2}{\arg\max} \left\{ \underset{\beta}{\max} \, \ell(\beta, \sigma^2) \right\} \\ &= \underset{\sigma^2}{\arg\max} \, \ell(\hat{\beta}, \sigma^2) \\ &= \underset{\sigma^2}{\arg\max} \, -\frac{1}{2} \left\{ n \log \sigma^2 + \frac{1}{\sigma^2} (y - X \hat{\beta})^\top (y - X \hat{\beta}) \right\}. \end{split}$$

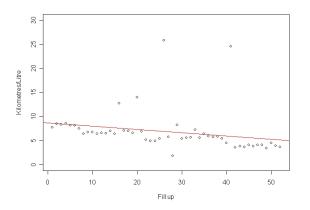
Differentiating and setting equal to zero yields

$$\hat{\sigma}^2 = rac{1}{n} (y - X \hat{eta})^ op (y - X \hat{eta}).$$

Next week we will see that a better (unbiased) estimator is

$$S^2 = rac{1}{n-p}(y-X\hat{oldsymbol{eta}})^{ op}(y-X\hat{oldsymbol{eta}}).$$

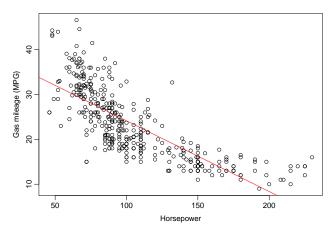
4 D > 4 B > 4 E > 4 E > 9 Q C



$$\hat{\beta}_0 = 8.6 \quad \hat{\beta_1} = -0.068 \quad S^2 = 17.4$$

◆ロト ◆団 ト ◆ 豆 ト ◆ 豆 ・ り Q (~)

### Model with linear term only

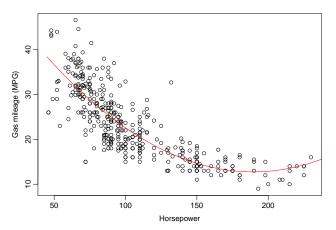


Parameter estimates:  $\hat{\beta}_0 = 39.94$  and  $\hat{\beta}_1 = -0.16$  and  $S^2 = 24.06$ .

◆ロト ◆団ト ◆豆ト ◆豆ト 豆 からぐ

Victor Panaretos (EPFL) Linear Models 67 / 309

#### Model with linear quadratic terms



Parameter estimates:  $\hat{\beta}_0 = 56.90$ ,  $\hat{\beta}_1 = -0.47$  and  $\hat{\beta}_2 = 0.0012$  and  $S^2 = 19.13$ .

Victor Panaretos (EPFL) Linear Models 68 / 309

There are two <u>dual</u> geometrical viewpoints that one may adopt:

$$\left( egin{array}{c} Y_1 \ Y_2 \ dots \ Y_n \end{array} 
ight) = \left( egin{array}{ccccc} 1 & x_{11} & x_{12} & \dots & x_{1q} \ 1 & x_{21} & x_{22} & & x_{2q} \ dots & dots & dots \ 1 & x_{(n-1)1} & x_{(n-1)2} & \dots & x_{(n-1)q} \ 1 & x_{n1} & x_{n2} & \dots & x_{nq} \end{array} 
ight) \left( egin{array}{c} eta_0 \ eta_1 \ dots \ eta_q \end{array} 
ight) + \left( egin{array}{c} arepsilon_1 \ dots \ dots \ eta_n \end{array} 
ight)$$

- Row geometry: focus on the *n* OBSERVATIONS
- Column geometry: focus on the p EXPLANATORIES

Both are useful, usually for different things:

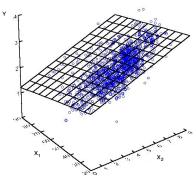
- Row geometry useful for exploratory analysis.
- Column geometry useful for theoretical analysis.

Both geometries give useful, but different, intuitive interpretations of the least squares estimators.

Victor Panaretos (EPFL) Linear Models 69 / 309

Corresponds to the "scatterplot geometry" – (data space)

- n points in  $\mathbb{R}^p$
- each corresponds to an observation
- least squares parameters give parametric equation for a hyperplane
- hyperplane has property that it minimizes the sum of squared vertical distances of observations from the plane itself over all possible hyperplanes



 Fitted values are vertical projections (NOT orthogonal projections!) of observations onto plane, residuals are signed vertical distances of observations from plane. Column Geometry (Variables)

#### Adopt the dual perspective:

- ullet Consider the entire vector y as a single point living in  $\mathbb{R}^n$
- ullet Then consider each variable (column) as a point also in  $\mathbb{R}^n$

What is the interpretation of the p-dimensional vector  $\hat{\beta}$ , and the n-dimensional vectors  $\hat{y}$  and e in this dual space?

Turns out there is another important plane here: the plane spanned by the variable vectors (the column vectors of X).

Recall that this is the *column space* of X, denoted by  $\mathfrak{M}(X)$ .

#### Column Geometry (Variables)

Recall: 
$$\underbrace{\mathcal{M}(X)}_{\mathsf{Column}} := \{X\gamma : \gamma \in \mathbb{R}^p\}$$

Q: What does  $Y = X\beta + \varepsilon$  mean?

A: Y is [some element of  $\mathfrak{M}(X)$ ] + [Gaussian disturbance].

Any realisation y of Y will lie outside  $\mathcal{M}(X)$  (almost surely). MLE estimates  $\beta$  by minimising

$$(y - X\beta)^{\top}(y - X\beta) = ||y - X\beta||^2$$

Thus we search for a  $\beta$  giving the element of  $\mathfrak{M}(X)$  with the minimum distance from y .

Hence  $\hat{y} = X\hat{\beta}$  is the projection of y onto  $\mathcal{M}(X)$ :

$$\hat{y} = X\hat{\beta} := \underbrace{X(X^{\top}X)^{-1}X^{\top}}_{H} y = Hy.$$

H is the hat matrix (puts hat on y!)

## Column Geometry (Variables)

## Another derivation of the MLE of $\beta$ :

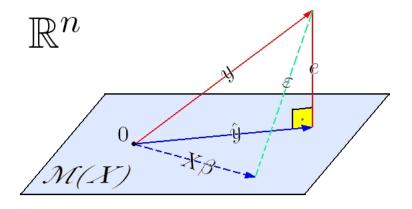
ullet Choose  $\hat{eta}$  to minimise  $(y-Xeta)^{ op}(y-Xeta)=\|y-Xeta\|^2$ , so

$$\hat{\beta} = \arg\min \|y - X\beta\|^2.$$

- $\bullet \ \min_{\beta \in \mathbb{R}^p} \|y X\beta\|^2 = \min_{\gamma \in \mathcal{M}(X)} \|y \gamma\|^2$
- ullet But the unique  $\gamma$  that yields  $\min_{\gamma \in \mathcal{M}(X)} \|y \gamma\|^2$  is  $\gamma = Py$ .
- Here P is the projection onto the column space of X,  $\mathfrak{M}(X)$ .
- Since X is of full rank,  $P = X(X^{\top}X)^{-1}X^{\top}$ .
- So  $\gamma = X(X^\top X)^{-1}X^\top y$
- $\hat{\beta}$  will now be the unique (since X non-singular) vector of coordinates of  $\gamma$  with respect to the basis of columns of X.
- So

$$X\hat{\beta} = \gamma = X(X^{\top}X)^{-1}X^{\top}y,$$

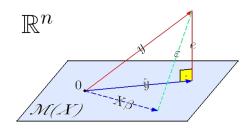
which implies that  $\hat{\beta} = (X^\top X)^{-1} X^\top y$ 



## The (Column) Geometry of Least Squares

# So what is $\hat{\beta}$ ?

- If X columns linearly independent, they are a (non-orthogonal) basis for  ${\mathfrak M}$
- Hence for any  $z \in \mathcal{M}(X)$ , there exists a unique  $\gamma \in \mathbb{R}^p$  such that  $z = X\gamma$



- ullet So  $\gamma$  contains coordinates of z with respect to the X-column basis
- $\bullet$  Consequently,  $\hat{\beta}$  contains coordinates of  $\hat{y}$  with respect to the X-column basis
- But  $\hat{y} = Hy = X\underbrace{(X^{\top}X)^{-1}X^{\top}y}_{u} = Xu$ , so u is the unique vector that gives coordinates of y with respect to the X-column basis
- ullet Hence we must have  $\hat{eta} = u = (X^{ op}X)^{-1}X^{ op}y$

◆ロト ◆昼ト ◆昼ト ◆ ■ めへで

## The (Column) Geometry of Least Squares

#### Facts:

- $\bullet e = (I H)y = (I H)\varepsilon.$
- ②  $\hat{y}$  and e are orthogonal, i.e.  $\hat{y}^{\top}e = 0$
- $\textbf{ 9 Pythagoras: } y^\top y = \hat{y}^\top \hat{y} + e^\top e = y^\top H y + \varepsilon^\top (I H) \varepsilon$

#### Derivation:

- $\bullet e = y X\hat{\beta} = y Hy = (I H)y = (I H)(X\beta + \varepsilon) = (I H)X\beta + (I H)\varepsilon = (I H)\varepsilon$
- $e = y \hat{y} = (I H)y \implies \hat{y}^{\top} e = y^{\top} H^{\top} (I H)y = 0$

Assume slightly different model:

$$egin{aligned} Y_i &= eta_0 + eta_1 x_{i1} + eta_2 x_{i2} + \cdots + eta_q x_{iq} + rac{arepsilon_i}{\sqrt{w_i}}, \quad arepsilon_i \stackrel{ind}{\sim} \mathcal{N}(0,\sigma^2), \quad w_i > 0 \end{aligned}$$
  $\diamondsuit$ 
 $Y_i \stackrel{ind}{\sim} N\left(eta_0 + eta_1 x_{i1} + eta_2 x_{i2} + \cdots + eta_q x_{iq}, rac{\sigma^2}{w_i}
ight).$ 

With the  $w_j$  known weights (example: each  $Y_j$  is an average of  $w_j$  measurements).

Arises often in practice (e.g., in sample surveys), but also arises in theory.

◆ロト ◆昼 ◆ 単 ◆ 単 ・ り へ ○

**Transformation:** 

$$y' = W^{1/2}y, \quad X' = W^{1/2}X$$

with

$$W_{n\times n} = \mathsf{diag}(w_1,\ldots,w_n)$$

Leads to usual scenario. In this notation we obtain:

$$\hat{\beta} = [(X')^{\top} X']^{-1} (X')^{\top} y'$$

$$= (X^{\top} WX)^{-1} X^{\top} Wy$$

Similarly:

$$S^2 = rac{1}{n-p} y^{ op} \left[ W - WX(X^{ op}WX)^{-1}X^{ op}W 
ight] y$$

# Distribution Theory of Least Squares

Gaussian Linear Model:

$$Y_{n\times 1} = X_{n\times p}\beta_{p\times 1} + \varepsilon_{n\times 1}, \quad \varepsilon \sim \mathcal{N}_n(0, \sigma^2 I)$$

We have derived the estimators:

- $\bullet \ \hat{\beta} = (X^\top X)^{-1} X^\top y$
- $\bullet \ \hat{\sigma}^2 = \frac{1}{n} (y X \hat{\beta})^\top (y X \hat{\beta}) = \frac{1}{n} ||\hat{y} y||^2$
- $S^2 = \frac{1}{n-p} ||\hat{y} y||^2$

We need to study the distribution of these estimators for the purpose of:

- Understanding their precision
- Building confidence intervals
- Testing hypotheses
- Comparing them to other candidate estimators
- ...



## **Theorem**

Let  $Y_{n\times 1} = X_{n\times p}\beta_{p\times 1} + \varepsilon_{n\times 1}$  with  $\varepsilon \sim \mathcal{N}_n(0,\sigma^2I)$  and assume that X has full rank p < n. Then,

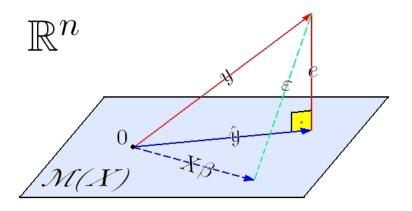
- $\bullet \ \hat{\beta} \sim \mathcal{N}_p\{\beta, \sigma^2(X^\top X)^{-1}\};$
- **2** the random variables  $\hat{\beta}$  and  $S^2$  are independent; and
- $\frac{n-p}{\sigma^2}S^2 \sim \chi^2_{n-p}$ , where  $\chi^2_{\nu}$  denotes the chi-square distribution with  $\nu$  degrees of freedom.

# Corollary

Let  $Y_{n\times 1} = X_{n\times p}\beta_{p\times 1} + \varepsilon_{n\times 1}$  with  $\varepsilon \sim \mathcal{N}_n(0, \sigma^2 I)$ . The statistic Hy is sufficient for the parameter  $\beta$ . If X has full rank p < n, then  $\hat{\beta}$  is also sufficient for  $\beta$ .

## Corollary

 $S^2$  is unbiased whereas  $\hat{\sigma}^2$  is biased (so we prefer  $S^2$ ).



#### Proof of the Theorem.

1. Recall our results for linear transformations of Gaussian variables:

$$\left. \begin{array}{l} \hat{\beta} = (X^\top X)^{-1} X^\top Y \\ Y \sim \mathcal{N}_n(X\beta, \sigma^2 I) \end{array} \right\} \implies \hat{\beta} \sim \mathcal{N}_p \{ \beta, \sigma^2 (X^\top X)^{-1} \}$$

- 2. If e is independent of  $\hat{y} = X\hat{\beta}$ , then  $S^2 = e^{\top}e/(n-p)$  will be independent of  $\hat{\beta}$  (why?). Now notice that:
  - $\bullet$  e = (I H)y
  - $\hat{y} = Hy$
  - $y \sim \mathcal{N}(X\beta, \sigma^2 I)$

Therefore, from the properties of the Gaussian distribution e is independent of  $\hat{y}$  since  $(I-H)(\sigma^2 I)H = \sigma^2 (I-H)H = 0$ , by idempotency of H.

## proof cont'd.

3. For the last part recall that

$$e = (I-H)arepsilon \implies (n-p)S^2 = (n-p)rac{e^+e}{n-p} = arepsilon^ op (I-H)arepsilon$$

by idempotency of H. But recall that  $\varepsilon \sim \mathcal{N}_n(0, \sigma^2 I_n)$  so  $\sigma^{-1}\varepsilon \sim \mathcal{N}_n(0, I_n)$ . Therefore, by the properties of normal quadratic forms (slide 40),

$$rac{(n-p)}{\sigma^2}S^2=(\sigma^{-1}arepsilon)^ op(I-H)(\sigma^{-1}arepsilon)\sim \chi^2_{n-p}.$$



# Proof of the first Corollary.

Write  $y = Hy + (I - H)y = \hat{y} + e$ .

If we can show that the conditional distribution of the 2n-dimensional vector  $W=(\hat{y},e)^{\top}$  given  $\hat{y}$  does not depend on  $\beta$ , then we will also know that the conditional distribution of  $y=\hat{y}+e$  given  $\hat{y}$  does not depend on  $\beta$  either, proving the proposition.

But we have proven that  $\hat{y}$  is independent of e. Therefore, conditional on  $\hat{y}$ , e always has the same distribution  $\mathcal{N}(0,(I-H)\sigma^2)$ . It follows that, conditional on  $\hat{y}$ , the vector W has a distribution whose first n coordinates equal  $\hat{y}$  almost surely, and whose last n coordinates are  $\mathcal{N}(0,(I-H)\sigma^2)$ . Neither of those two depend on  $\beta$ , and the proof is complete.

When X has full rank,  $\hat{\beta}$  is a 1-1 function of Hy, and is also sufficient for  $\beta$ .

# Proof of the second Corollary.

Recall that if  $Q \sim \chi_d^2$ , then  $\mathbb{E}[Q] = d$ .

How to construct  $1 - \alpha$  CI for a linear combination of the parameters,  $c^{\top}\beta$ ?

- Have  $c^{\top}\hat{\beta} \sim \mathcal{N}_1(c^{\top}\beta, \sigma^2 c^{\top}(X^{\top}X)^{-1}c) = \mathcal{N}_1(c^{\top}\beta, \sigma^2\delta)$
- Therefore  $Q = (c^{\top}\hat{\beta} c^{\top}\beta)/(\sigma\sqrt{\delta}) \sim \mathcal{N}_1(0,1)$
- Hence  $Q^2 \sim \chi_1^2$
- ullet and  $Q^2$  is independent of  $S^2$  (since  $\hat{eta}$  is independent of  $S^2$ )
- ullet while  $rac{n-p}{\sigma^2}S^2\sim\chi^2_{n-p}$ .

In conclusion:

$$\frac{\frac{Q^2}{1}}{\frac{(n-p)}{n-p}S^2} \sim F_{1,n-p} \Rightarrow \frac{\frac{(c^\top \hat{\beta} - c^\top \beta)^2}{\sigma^2 \delta}}{\frac{S^2}{\sigma^2}} = \left(\frac{c^\top \hat{\beta} - c^\top \beta}{\sqrt{S^2 c^\top (X^\top X)^{-1} c}}\right)^2 \sim F_{1,n-p}$$

• But for real W,  $W^2 \sim F_{1,n-p} \iff W \sim t_{n-p}$ , so base CI on:

$$rac{c^{ op}\hat{eta}-c^{ op}eta}{\sqrt{S^2c^{ op}(X^{ op}X)^{-1}c}}\sim t_{n-p}$$

• We obtain  $(1 - \alpha) \times 100\%$  CI:

$$c^{\top}\hat{\beta} \pm t_{n-p}(1-\alpha/2)\sqrt{S^2c^{\top}(X^{\top}X)^{-1}c}.$$

- What about a  $(1 \alpha)$  CI for  $\beta_r$ ? (rth coordinate)
- Let  $c_r = (0, 0, \dots, 0, 1, 0, \dots, 0)$
- ullet Then  $eta_r = c^ op eta$
- Therefore, base CI on

$$\frac{c_r^\top \hat{\beta} - c_r^\top \beta}{\sqrt{S^2 c_r^\top (X^\top X)^{-1} c_r}} = \frac{\hat{\beta}_r - \beta_r}{\sqrt{S^2 v_{r,r}}} \sim t_{n-p},$$

where  $v_{r,s}$  is the r,s element of  $(X^{\top}X)^{-1}$ .

• Obtain  $(1 - \alpha) \times 100\%$  CI:

$$\hat{eta} \pm t_{n-p} (1-lpha/2) \sqrt{S^2 v_{rr}}$$

#### Confidence and Prediction Intervals

- ullet Suppose we want to predict the value of  $y_+$  for an  $x_+ \in \mathbb{R}^p$
- Our model predicts  $y_+$  by  $x_+^{\top} \hat{\beta}$ .
- But  $y_+ = x_+^\top \beta + \varepsilon_+$  so a prediction interval is DIFFERENT from an interval for a linear combination  $c^\top \beta$  (extra uncertainty due to  $\varepsilon_+$ ):
  - ullet  $\mathbb{E}[x_+^ op \hat{eta}_+ arepsilon_+] = x_+^ op eta$
  - $\bullet \ \operatorname{var}[x_+^\top \hat{\beta} + \varepsilon_+] = \operatorname{var}[x_+^\top \hat{\beta}] + \operatorname{var}[\varepsilon_+] = \sigma^2 [x_+^\top (X^\top X)^{-1} x_+ + 1]$
- Base prediction interval on:

$$rac{x_+^{ op}\hat{eta} - y_+}{\sqrt{S^2\{1 + x_+^{ op}(X^{ op}X)^{-1}x_+\}}} \sim t_{n-p}.$$

• Obtain  $(1 - \alpha)$  prediction interval:

$$x_+^{ op} \hat{eta} \pm t_{n-p} (1-lpha/2) \sqrt{S^2 \{1 + x_+^{ op} (X^{ op} X)^{-1} x_+ \}}.$$

 $R^2$  is a *measure of fit* of the model to the data.

- We are trying to best approximate y through an element of the column-space of X.
- How successful are we? Squared error is  $e^{\top}e$ .
- How large is this, relative to data variation? Look at

$$\frac{\|e\|^2}{\|y\|^2} = \frac{e^\top e}{y^\top y} = \frac{y^\top (I - H)y}{y^\top y} = 1 - \frac{\hat{y}^\top \hat{y}}{y^\top y}$$

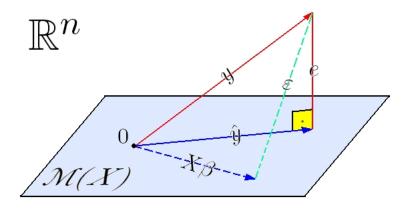
Define

$$R_0^2 = rac{\hat{y}^ op \hat{y}}{y^ op y} = rac{\|\hat{y}\|^2}{\|y\|^2}$$

• Note that  $0 \le R_0^2 \le 1$ 

Interpretation: what proportion of the squared norm of y does our fitted value  $\hat{y}$  explain?

◆ロト ◆団ト ◆重ト ◆重ト 重 める()



"Centred (in fact, usual)  $R^2$ ". Compares empirical variance of  $\hat{y}$  to empirical variance of y, instead of the empirical norms. In other words:

$$R^2 = \frac{\frac{1}{n} \sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\frac{1}{n} \sum_{i=1}^n (y_i - \bar{y})^2} = \frac{\sum_{i=1}^n (\hat{y}_i - \bar{y})^2}{\sum_{i=1}^n (y_i - \bar{y})^2} = \frac{\sum_{i=1}^n \hat{y}_i^2 - n\bar{y}^2}{\sum_{i=1}^n y_i^2 - n\bar{y}^2}.$$

(note that  $\frac{1}{n}\sum_{i=1}^n \hat{y}_i = \frac{1}{n}\sum_{i=1}^n (y_i - e_i) = \bar{y}$  because  $e \perp 1$  (recall that 1 is the vector of 1's = first column of design matrix X) so  $\sum_i e_i = 0$ .

Note that

$$R^2 = \frac{\|\hat{y}\|^2 - \|\bar{y}\mathbf{1}\|^2}{\|y\|^2 - \|\bar{y}\mathbf{1}\|^2}.$$

- $R_0^2$  mathematically more natural (does not treat first column of X as special).
- $R^2$  statistically more relevant (expresses variance—the first column of X usually *is* special, in statistical terms!).
- $R_0^2$  and  $R^2$  may differ a lot when  $\bar{y}$  large.

Geometrical interpretation of  $\mathbb{R}^2$ : project y and  $\hat{y}$  on orthogonal complement of 1, then compare the norms (of the projections):

- $\mathbf{1}(\mathbf{1}^{\top}\mathbf{1})^{-1}\mathbf{1}^{\top}y = \mathbf{1}n^{-1}\sum_{i=1}^{n}y_{i} = \mathbf{1}\bar{y}.$
- $ullet \ \mathbf{1}(\mathbf{1}^{\top}\mathbf{1})^{-1}\mathbf{1}^{\top}\hat{y} = \mathbf{1}n^{-1}\sum_{i=1}^{n}\hat{y}_{i} = \mathbf{1}\bar{y}.$

So

$$R^{2} = \frac{\|\hat{y}\|^{2} - \|\bar{y}\mathbf{1}\|^{2}}{\|y\|^{2} - \|\bar{y}\mathbf{1}\|^{2}} = \frac{\|(I - \mathbf{1}(\mathbf{1}^{\top}\mathbf{1})^{-1}\mathbf{1})\hat{y}\|^{2}}{\|(I - \mathbf{1}(\mathbf{1}^{\top}\mathbf{1})^{-1}\mathbf{1})y\|^{2}}$$

Intuition: Should not take into account the part of ||y|| that is explained by a constant, we only want to see the effect of the explanatory variables.

NOTE: Statistical packages (e.g., R) provide  $R^2$  (and/or  $R_a^2$ , see below), not  $R_0^2$ .

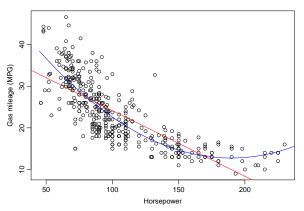
Exercise: Show that  $R^2 = [\text{corr}(\{\hat{y}_i\}_{i=1}^n, \{y_i\}_{i=1}^n)]^2$ .

Exercise: Show that  $R^2 \leq R_0^2$ .

4 D > 4 B > 4 B > B 9 9 9

 $\mathbb{R}^2$  coefficients for the linear and quadratic models:

	$R_0^2$	$R^2$
linear	0.96	0.61
quadratic	0.97	0.69



#### Different Versions of $R^2$

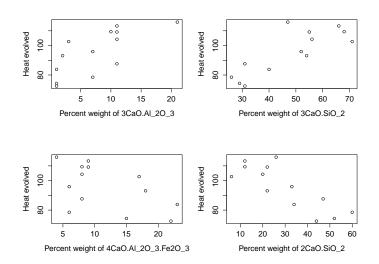
The adjusted  $\mathbb{R}^2$  takes into account the number of variables employed. It is defined as:

$$R_a^2 = R^2 - (1 - R^2) \frac{n-1}{n-p}.$$

Corrects for the fact that we can always increase  $\mathbb{R}^2$  by adding variables. One can also correct the un-centred  $\mathbb{R}^2_0$  by evaluating

$$R_0^2 - (1 - R_0^2) \frac{n}{n - p}$$
.

Case	3 CaO. Al <sub>2</sub> O <sub>3</sub>	3 Ca O . Si O 2	4 Cao. Al <sub>2</sub> O <sub>3</sub> . Fe <sub>2</sub> O <sub>3</sub>	2 CaO . SiO2	Heat
1	7.00	26.00	6.00	60.00	78.50
2	1.00	29.00	15.00	52.00	74.30
3	11.00	56.00	8.00	20.00	104.30
4	11.00	31.00	8.00	47.00	87.60
5	7.00	52.00	6.00	33.00	95.90
6	11.00	55.00	9.00	22.00	109.20
7	3.00	71.00	17.00	6.00	102.70
8	1.00	31.00	22.00	44.00	72.50
9	2.00	54.00	18.00	22.00	93.10
10	21.00	47.00	4.00	26.00	115.90
11	1.00	40.00	23.00	34.00	83.80
12	11.00	66.00	9.00	12.00	113.30
13	10.00	68.00	8.00	12.00	109.40



#### Example: Cement Heat Evolution

- > cement.lm<-lm( $y\sim1+x1+x2+x3+x4$ ,data=cement)
- > summary(cement.lm)

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	62.4054	70.0710	0.89	0.3991
×1	1.5511	0.7448	2.08	0.0708
x2	0.5102	0.7238	0.70	0.5009
x3	0.1019	0.7547	0.14	0.8959
×4	-0.1441	0.7091	-0.20	0.8441

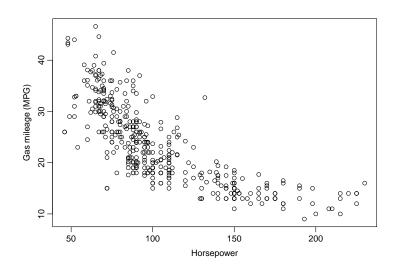
Residual standard error: 2.446 on 8 degrees of freedom R-Squared: 0.9824

> x.plus
[1] 25 25 25 25
predict(cement.lm,x.plus,interval="confidence",
se.fit=T,level=0.95)

Fit	Lower	Upper
112.8	97.5	128.2

predict(cement.lm,x.plus,interval="prediction",
se.fit=T,level=0.95)

Fit	Lower	Upper
112.8	96.5	129.2



### Horsepower and MPG of cars

- > auto.lm <- lm(mpg~1+horsepower+I(horsepower<sup>2</sup>),data=Auto)
- > summary(auto.lm)

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	56.9000	1.8004	31.60	$< 2 \times 10^{-16}$
horsepower	-0.4662	0.0311	-14.98	$< 2 \times 10^{-16}$
I(horsepower <sup>2</sup> )	0.0012	0.0001	10.08	$< 2 \times 10^{-16}$

Residual standard error: 4.374 on 389 degrees of freedom

R-Squared: 0.6876

> x.plus horsepower 120

> predict(auto.lm, x.plus, interval="confidence", se.fit=T, level=0.95)

Fit	Lower	Upper
18.68	18.03	19.33

> predict(auto.lm, x.plus, interval="prediction",
se.fit=T, level=0.95)

Fit	Lower	Upper
18.68	10.05	27.30

Gauss-Markov & Optimal Estimation

Gaussian Linear Model: Efficiency of LSE (Optimality)

Q: Geometry suggests that the LSE  $\hat{\beta}$  is a sensible estimator. But is it the best we can come up with?

A: Yes,  $\hat{\beta}$  is the unique minimum variance unbiased estimator of  $\beta$ .

(To be seen in Statistical Theory course, since  $\hat{\beta}$  is sufficient and complete)

Thus, in the Gaussian Linear model, the LSE are the best we can do as far as unbiased estimators go.

(actually can show  $S^2$  is optimal unbiased estimator of  $\sigma^2$ , by similar arguments)

Second Order Assumptions: Optimality in a weaker setting?

The crucial assumption so far was:

• Normality:  $\varepsilon \sim \mathcal{N}_n(0, \sigma^2 I)$ 

What if we drop this strong assumption and assume something weaker?

• Uncorrelatedness:  $\mathbb{E}[\varepsilon] = 0 \& \text{var}[\varepsilon] = \sigma^2 I$ 

(notice we do not assume any particular distribution.)

How well do our LSE estimators perform in this case?

(note that in this setup the observations may not be independent — uncorrelatedness implies independence only in the Gaussian case.)

For a start, we retain unbiasedness:

#### Lemma

If we only assume both

$$\mathbb{E}[\epsilon] = 0 \quad var[\epsilon] = \sigma^2 I$$

instead of

$$\varepsilon \sim \mathcal{N}(0, \sigma^2 I),$$

then the following remain true:

- $\bullet \ \mathbb{E}[\hat{\beta}] = \beta;$
- **2**  $Var[\hat{\beta}] = \sigma^2(X^{\top}X)^{-1};$
- **3**  $\mathbb{E}[S^2] = \sigma^2$ .

But what about optimality properties?

#### **Theorem**

Let  $Y_{n\times 1} = X_{n\times p}\beta_{p\times 1} + \varepsilon_{n\times 1}$ , with p < n, X having rank p, and

- $\mathbb{E}[\varepsilon] = 0$ ,
- $var[\varepsilon] = \sigma^2 I$ .

Then,  $\hat{\beta} = (X^{\top}X)^{-1}X^{\top}Y$  is the best linear unbiased estimator of  $\beta$ , that is, for any linear unbiased estimator  $\tilde{\beta}$  of  $\beta$ , it holds that

$$var(\hat{\beta}) - var(\hat{\beta}) \succeq 0.$$

## Proof.

Let  $\tilde{\beta}$  be linear and unbiased, in other words:  $\begin{cases} \tilde{\beta} = AY, & \text{for some } A_{p \times n}, \\ \mathbb{E}[\tilde{\beta}] = \beta, & \text{for all } \beta \in \mathbb{R}^p. \end{cases}$ 

These two properties combine to yield,

$$eta = \mathbb{E}[\tilde{eta}] = \mathbb{E}[AY] = \mathbb{E}[AX\beta + Aarepsilon] = AXeta, \quad eta \in \mathbb{R}^p$$

$$\implies (AX - I)\beta = 0, \, \forall \beta \in \mathbb{R}^p.$$

We conclude that the null space of (AX - I) is the entire  $\mathbb{R}^p$ , and so AX = I.

$$\begin{aligned} \operatorname{var}[\tilde{\beta}] - \operatorname{var}[\hat{\beta}] &= A\sigma^2 I A^\top - \sigma^2 (X^\top X)^{-1} \\ &= \sigma^2 \{AA^\top - AX(X^\top X)^{-1} X^\top A^\top \} \\ &= \sigma^2 A (I - H) A^\top \\ &= \sigma^2 A (I - H) (I - H)^\top A^\top \\ \succeq 0. \end{aligned}$$

## Large Sample Distribution of $\hat{\beta}$

If  $\mathbb{E}[\varepsilon] = 0$  and  $\operatorname{cov}[\varepsilon] = \sigma^2 I$ 

 $\hookrightarrow$  Gauss-Markov says  $\hat{\beta}$  optimal linear unbiased estimator, regardless of whether or not  $\varepsilon$  is Gaussian.

**Question**: What can we say about the distribution of  $\hat{\beta}$  when  $cov(\varepsilon) = \sigma^2 I$ , but  $\varepsilon$  is not necessarily Gaussian?

Note that we can always write

$$\hat{\beta} - \beta = (X^{\top}X)^{-1}X^{\top}\varepsilon.$$

- Since there is a huge variety of candidate distributions for  $\varepsilon$  that would be compatible with the property  $\text{cov}(\varepsilon) = \sigma^2 I$ , we cannot say very much about the exact distribution of  $\hat{\beta} \beta = (X^\top X)^{-1} X^\top \varepsilon$ .
- Can we at least hope to say something about this distribution asymptotically, as the sample becomes large?

◆ロト ◆個 ト ◆ 恵 ト ◆ 恵 ・ り へ ○

# Large Sample Distribution of $\hat{\beta}$

Large sample  $\iff$  increasing number of observations.

- We let  $n \to \infty$  (# rows of X tend to infinity)
- # columns of X, i.e., p, (held fixed).

# Theorem (Large Sample Distribution of $\hat{\beta}$ )

Let  $\{X_n\}_{n\geq 1}$  be a sequence of  $n\times p$  design matrices and  $Y_n=X_n\beta+\varepsilon_n$ . If

- **1**  $X_n$  is of full rank p for all  $n \geq 1$
- $\varepsilon_n$  is a zero mean n-vector with i.i.d. coordinates of variance  $\sigma^2$ , then the least squares estimator  $\hat{\beta}_n = (X_n^\top X_n)^{-1} X_n^\top Y_n$  satisfies

$$(X_n^{\top} X_n)^{1/2} (\hat{\beta}_n - \beta) \stackrel{d}{\longrightarrow} \mathcal{N}_p(0, \sigma^2 I).$$

4 나 ▶ 4 년 ▶ 4 년 ▶ 년 \*) 역(

<sup>&</sup>lt;sup>a</sup>Gottfried Noether (not Emmy Noether), Ann. Math. Stat., 1949.

Theorem's conclusion can be interpreted as:

for 
$$n$$
 "large enough",  $\hat{oldsymbol{eta}} \overset{d}{\simeq} \mathcal{N}\{oldsymbol{eta}, \sigma^2(X_n^{ op}X_n)^{-1}\}$ 

- ullet i.e. distribution of  $\hat{eta}$  gradually becomes what it would be if arepsilon were Gaussian
- ullet ... provided design matrix X satisfies Noether's condition (2).
- This equivalent to: diagonal elements of  $H_n=X_n(X_n^\top X_n)^{-1}X_n^\top$  , say  $h_{jj}(n)$  converge to zero uniformly in j as  $n\to\infty$

Because  $x_i^\top (X^\top X)^{-1} x_i = (e_i^\top X) (X^\top X)^{-1} (e_i^\top X)^\top = e_i^\top H e_i = h_{ii}$  where  $e_i$  is the ith canonical basis vector for  $\mathbb{R}^n$ .

• Note that trace(H) = p, so that the average  $\sum h_{jj}(n)/n \to 0$  — the question is do all the  $h_{jj}(n) \to 0$  uniformly?

Has a very clear interpretation in terms of the form of the design that we will see when we discuss the notions of *leverage* and *influence*.

4□ > 4□ > 4□ > 4□ > 4□ > 4□

Victor Panaretos (EPFL)

To understand Condition (2), consider simple linear model

$$Y_i = \beta_0 + \beta_1 t_i + \varepsilon_i, \qquad i = 1, \ldots, n.$$

Here, p = 2. Can show that

$$h_{jj}(n) = rac{1}{n} + rac{(t_j - \overline{t}\,)^2}{\sum_{k=1}^n (t_k - \overline{t}\,)^2}$$

ullet Suppose  $t_i=i$ , for  $i=1,\ldots,n$  (regular grid). Then

$$h_{jj}(n) = rac{1}{n} + rac{\{j-(n+1)/2\}^2}{(n^3-n)/12}$$

so 
$$\max_{1\leq j\leq n}h_{jj}(n)=h_{nn}(n)=rac{1}{n}+rac{3(n-1)}{n^2(n-1)}\stackrel{n o\infty}{\longrightarrow} 0.$$

Victor Panaretos (EPFL)

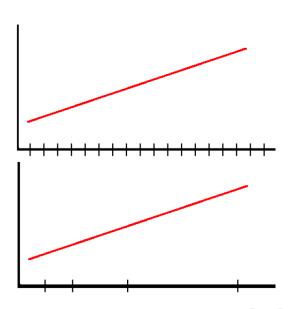
## Large Sample Distribution of $\hat{\beta}$

• Now consider  $t_i = 2^i$  (grid points spread apart as n grows). The centre of mass and sum of squares of the grid points is now

$$\overline{t} = \frac{2(2^n - 1)}{n}, \quad \sum_{i=1}^n (t_i - \overline{t})^2 = \frac{4^{n+1} - 4}{3} - \frac{4^{n+1} + 4 - 2^{n+3}}{n}$$

and so

$$\max_{1 \leq j \leq n} h_{jj}(n) = h_{nn}(n) \stackrel{n o \infty}{\longrightarrow} rac{3}{4}.$$



#### Proof.

Recall that  $\hat{\beta}_n - \beta = (X_n^\top X_n)^{-1} X_n^\top \varepsilon_n$ . We will show that for any unit vector u,

$$u^{\top}(X_n^{\top}X_n)^{-1/2}X_n^{\top}\varepsilon_n\stackrel{d}{
ightarrow} N(0,\sigma^2),$$

and then the theorem will be proven by the Cramér-Wold device<sup>a</sup>. Now notice that

$$u^{ op}(X_n^{ op}X_n)^{-1/2}X_n^{ op}arepsilon_n=\gamma_n^{ op}arepsilon_n$$

where:

$$\bullet \hspace{0.1cm} \boldsymbol{\gamma}_n = (\gamma_{n,1}, \ldots, \gamma_{n,n})^\top = \left( u^\top (X_n^\top X_n)^{-1/2} x_1, \ldots, u^\top (X_n^\top X_n)^{-1/2} x_n \right)^\top$$

Consequently, the result follows from the weighted sum CLT upon noticing:

$$\max_{1 \leq i \leq n} \gamma_{n,i}^2 \left/ \sum_{k=1}^n \gamma_{n,k}^2 \leq \max_{1 \leq i \leq n} x_i^ op (X_n^ op X_n)^{-1} x_i 
ightarrow 0$$

<sup>a</sup>Cramér-Wold:  $\xi_n \overset{d}{\to} \xi$  in  $\mathbb{R}^d$  if and only if  $u^\top \xi \overset{d}{\to} u^\top \xi$  in  $\mathbb{R}$  for all unit vectors u.

# Diagnostics

Four basic assumptions inherent in the Gaussian linear regression model:

- Linearity:  $\mathbb{E}[Y]$  is linear in X.
- Homoskedasticity:  $var[\varepsilon_j] = \sigma^2$  for all j = 1, ..., n.
- Gaussian Distribution: errors are Normally distributed.
- Independent Errors:  $\varepsilon_i$  independent of  $\varepsilon_j$  for  $i \neq j$ .

When one of these assumptions fails clearly, then Gaussian linear regression is inappropriate as a model for the data.

Isolated problems, such as <u>outliers</u> and <u>influential observations</u> also deserve investigation. They *may or may not* decisively affect model validity.

Victor Panaretos (EPFL)

How do we check these assumptions?

Scientific reasoning: impossible to *validate* model assumptions.

Cannot *prove* that the assumptions hold. Can only provide evidence in favour (or against!) them.

### Strategy:

- Find implications of each assumption that we can check graphically (mostly concerning residuals).
- Construct appropriate plots and assess them (requires experience).

"Magical Thinking": Beware of overinterpreting plots!

Residuals e: Basic tool for checking assumptions.

Recall: 
$$e = y - \hat{y} = y - X\hat{\beta} = (I - H)y = (I - H)\varepsilon$$

Intuition: the residuals represent the aspects of y that cannot be explained by the columns of X.

Since  $\varepsilon \sim \mathcal{N}_n(0, \sigma^2 I)$ , if the model is correct we should have  $e \sim \mathcal{N}_n\{0, \sigma^2(I-H)\}.$ 

So if assumptions hold 
$$ightarrow \left\{egin{array}{l} e_i \sim \mathcal{N}\{0,\sigma^2(1-h_{ii})\} \ \operatorname{\mathsf{cov}}(e_i,e_j) = -\sigma^2 h_{ij} \end{array}
ight.$$

Note the residuals are correlated, and that they have unequal variances. Define the standardised residuals:

$$r_i:=rac{e_i}{s\sqrt{1-h_{ii}}}, \quad i=1,\ldots,n.$$

These are still correlated but have variance  $\approx 1$ .

(can decorrelate by  $U^{\top}e$ , where  $H=U\Lambda U^{\top}$ ) – why?

Is  $\mathbb{E}[Y]$  entirely specified as linear functional of X? Did we leave variables out? Is it also a linear functional of non-linear transformations of X-columns?

A first impression can be drawn by looking at plots of the response against each of the explanatory variables. Other plots to look at? Notice that, by construction of e = (I - H)y we have

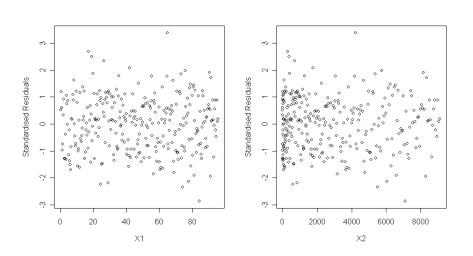
$$X^{\top}e=0.$$

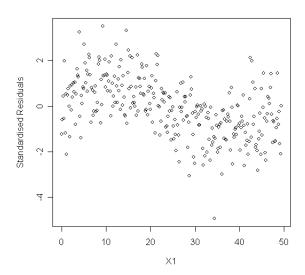
Hence, no correlation will appear between explanatory variables and residuals.

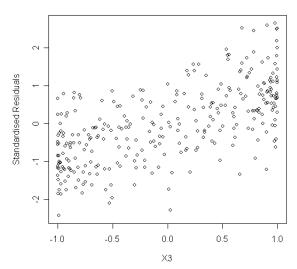
- ullet Plot stand. residuals r against each explanatory variable (columns of X).
  - → No systematic (non-linear) patterns should appear in these plots. A systematic pattern would suggest incorrect dependence of the response on the particular explanatory (e.g. need to add a transformation of that explanatory as an additional variable).

Also, no correlation **should** appear between unused explanatory variables and residuals.

- $\bullet$  Plot standardised residuals r against explanatories left out of the model.







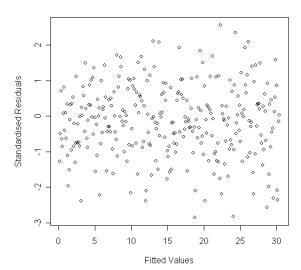
#### Checking for Homoskedasticity

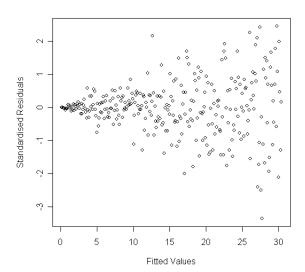
$$\mathsf{Homoskedastic} = \underbrace{\delta\mu\,o}_{\mathit{same}} + \underbrace{\sigma\,\kappa\varepsilon\delta\,\alpha\sigma\,\mu\grave{\mathsf{o}}\varsigma}_{\mathit{spread}}$$

According to our model assumptions, the variance of the errors  $\varepsilon_j$  should be the same across indices:

$$\mathsf{var}(arepsilon_j) = \sigma^2$$

- Plot r against the fitted values  $\hat{y}$ . (why not against y?)
  - $\hookrightarrow$  A random scatter should appear, with approximately constant spread of the values of r for the different values of  $\hat{y}$ . "Trumpet" or "bulging" effects indicate failure of the homoskedasticity assumption.
  - $\hookrightarrow$  Since  $\hat{y}^{\top}e=0$ , this plot can also be used to check linearity, as before.





#### Checking for Normality

<u>Idea</u>: compare the distribution of standardised residuals against a Normal distribution.

How?

Compare the empirical with the theoretical quantiles . . .

The p-quantile  $(p \in [0,1])$  of a distribution F is the value  $\delta$  defined as

$$\delta := \inf \{ \alpha \in \mathbb{R} : F(\alpha) \ge p \}.$$

Notation:  $\delta = F^{-1}(p)$  (although the inverse may not be well defined) Given a sample  $X_1, \ldots, X_n$ , the *empirical p quantile* is the value  $\gamma$  defined as

$$\gamma = \inf \left\{ lpha \in \mathbb{R} : rac{\#\{X_i \leq lpha\}}{n} \geq p 
ight\}.$$

Notation:  $\gamma = \hat{F}_n^{-1}(p)$ 

◆ロ → ◆団 → ◆ 豆 → ○ ● ・ ○ へ ○ ○

#### Checking for Normality

A quantile plot for a given sample plots certain empirical quantiles against the corresponding theoretical quantiles (i.e. those under the assumed distribution). If the sample at hand originates from F, then we expect that the points of the plot fall close to the  $45^{\circ}$  line.

• Plot the empirical  $\{k/n\}_{k=1}^n$  quantiles of standardised residuals

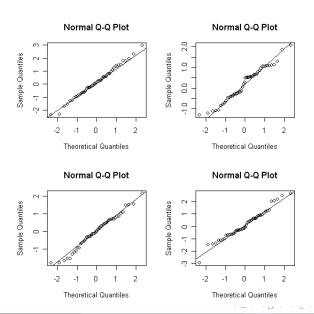
$$r_{(1)} \leq r_{(2)} \leq \cdots \leq r_{(n)}$$

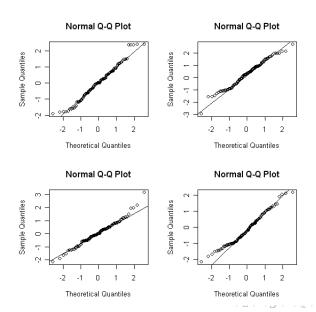
against theoretical quantiles  $\Phi^{-1}\{1/(n+1)\},\ldots,\Phi^{-1}\{n/(n+1)\}$  of a  $\mathcal{N}(0,1)$  distribution.

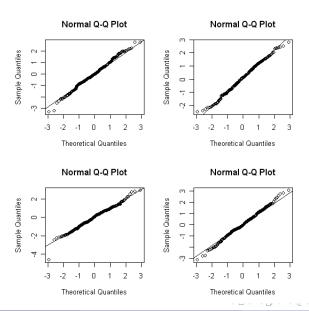
- $\hookrightarrow$  Think why we pick  $\Phi^{-1}\left(\frac{k}{n+1}\right)$  instead of  $\Phi^{-1}\left(\frac{k}{n}\right)$ .

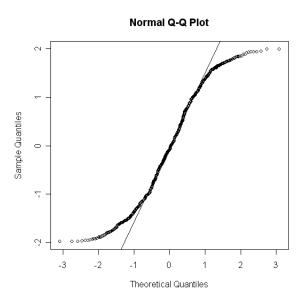
Beware of overinterpretation when n is small!

◆ロト ◆部ト ◆草ト ◆草ト 草 り900









#### Checking for Independence

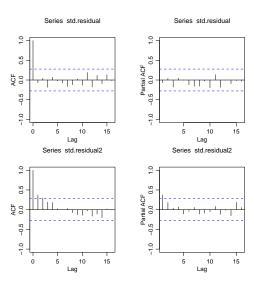
- It is assumed that  $var[\varepsilon] = \sigma^2 I$ .
- Under assumption of normality this is equivalent to independence

Difficult to check this assumption in practice.

- One thing to check for is clustering, which may suggest dependence.
  - $\hookrightarrow$  e.g. identifying groups of related individuals with correlated responses
- When observations are time-ordered can look at correlation  $\operatorname{corr}[r_t, r_{t+k}]$  or partial correlation  $\operatorname{corr}[r_t, r_{t+k} | r_{t+1}, \dots, r_{t+k-1}]$ . When such correlations exist, we enter the domain of *time series*.

#### Existence of dependence:

- seriously affects estimator reliability
- inflates standard errors



An influential observation can usually be categorised as an:

- outlier (relatively easier to spot by eye)
   OR
- leverage point (not as easy to spot by eye)

#### Influential observations

- May or may not decisively affect model validity.
- Require scrutiny on an individual basis and consultation with the data expert.

David Brillinger (Berkeley): You will not find your Nobel prize in the fit, you will find it in the outliers!

Influential observations may reveal unanticipated aspects of the scientific problem that are worth studying, and so must not simply be scorned as "non-conformists"!

#### Outliers

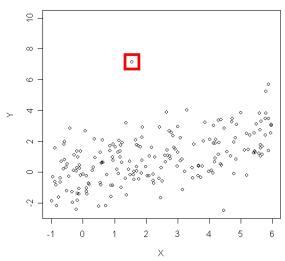
An *outlier* is an observation that stands out in some way from the rest of the observations, causing *Surprise!* Exact mathematical definition exists (Tukey) but we will not pursue it.

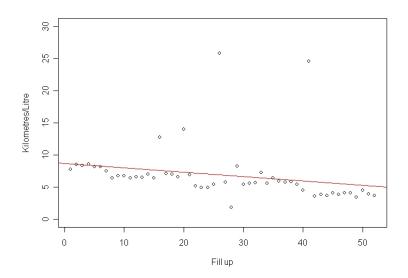
- In regression, outliers are points falling far from the cloud surrounding the regression line (or surface).
- They have the effect of "pulling" the regression line (surface) toward them.

Outliers can be checked for visually through:

- The regression scatterplot.
- Residual Plots.
  - $\hookrightarrow$  Points that fall beyond (-2,2) in the  $(\hat{y},r)$  plot.

Outliers may result from a data registration error, or a single extreme event. They can, however, result because of a deeper inadequacy of our model (especially if there are many!).





- Outliers may be influential: they "stand out" in the "y-dimension".
- However an observation may also be influential because of unusual values in the "x-dimension".
- Such influential observations cannot be so easily detected through plots.

Call  $(x_j, y_j)$  the *j-th* case and notice that

$$\operatorname{\mathsf{var}}(y_j - \hat{y}_j) = \operatorname{\mathsf{var}}(e_j) = \sigma^2(1 - h_{jj}).$$

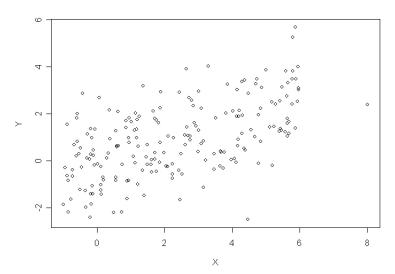
If  $h_{jj} \approx 1$ , then the model is constrained so  $\hat{y}_j = x_j^\top \hat{\beta} \simeq y_j!$  (i.e., need a separate parameter entirely devoted to fitting this observation!)

- $h_{ij}$  is called the *leverage* of the *j*-th case.
- since trace $(H) = \sum_{i=1}^{n} h_{jj} = p$ , cannot have low leverage for all cases
- ullet a good design corresponds to  $\mathit{h}_{jj} \simeq \mathit{p}/\mathit{n}$  for all j

(i.e. assumption  $\max_{j \le n} h_{jj} \stackrel{n \to 0}{\to} 0$  satisfied in asymptotic thm).

Leverage point: (rule of thumb) if  $h_{jj}>2\,p/n$  observation needs further scrutiny—e.g., fitting again without j-th case and studying effect.

 $Outlier+Leverage\ Point=TROUBLE$ 



- How to find cases having strong effect on fitted model?
- ullet Idea: see effect when case j , i.e.,  $(x_j,\,y_j)$  , is dropped.
- Let  $\hat{\beta}_{-j}$  be the LSE when model is fitted to data without case j, and let  $\hat{y}_{-j} = X \hat{\beta}_{-j}$  be the corresponding fitted value.
- Define Cook's distance

$$C_j = rac{1}{ps^2} (\hat{y} - \hat{y}_{-j})^ op (\hat{y} - \hat{y}_{-j}),$$

which measures scaled distance between  $\hat{y}$  and  $\hat{y}_{-j}$ .

Can show that

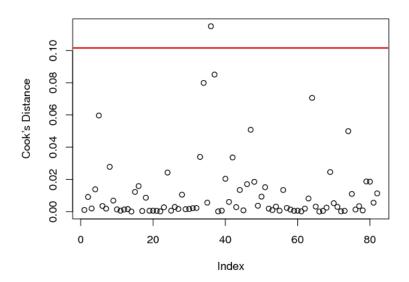
$$C_j = rac{r_j^2 h_{jj}}{p(1-h_{jj})},$$

so large  $C_j$  implies large  $r_j$  and/or large  $h_{jj}$ .

- Cases with  $C_i > 8/(n-2p)$  worth a closer look (rule of thumb)
- Plot  $C_j$  against index  $j=1,\ldots,n$  and compare with 8/(n-2p) level.

4□ > 4□ > 4∃ > 4∃ > ∃
9

Victor Panaretos (EPFL)



#### Summary

#### Diagnostic plots usually constructed:

- ullet y against columns of X
  - $\hookrightarrow$  check for linearity and outliers
- ullet standardized residual r against columns of X
- r against explanatories not included
  - $\hookrightarrow$  check for variables left out
- ullet r against fitted value  $\hat{y}$
- Normal quantile plot
  - $\hookrightarrow \ \mathsf{check} \ \mathsf{for} \ \mathsf{normality}$
- Cook's distance plot
  - $\hookrightarrow$  check for influential observations

Detour: Reminder on Hypothesis Tests

Detour: Very brief Reminder on Testing Hypotheses

- Scientific theories lead to assertions that are testable using empirical data.
- Data may discredit the theory (call it the *hypothesis*) or not (i.e., empirical findings reasonable under hypothesis).
- Example: Theory of "luminoferous aether" in late 19th century to explain light travelling in vacuum. Discredited by Michelson-Morley experiment.
- Similarities with the logical/mathematical concept of a *necessary condition*.

# Hypothesis Testing Setup

- $H_0$ : The <u>null hypothesis</u>
  - $\hookrightarrow$  scientific theory under scrutiny
- $\begin{array}{c} \bullet & \left\{ \begin{array}{l} Y, & \mathsf{data} \\ T(\cdot), & \mathsf{test} \ \mathsf{statistic}, \ \mathsf{assumed} \ \mathsf{positive} \end{array} \right. \\ \hookrightarrow & \mathsf{the} \ \mathsf{experimental} \ \mathsf{setup} \ \mathsf{to} \ \mathsf{test} \ \mathsf{theory} \end{array}$

#### INTUITION:

- The null hypothesis would predict a certain plausible range of values for T(Y) (plausible results of the experiment).
- ullet We would say that the assertion made by the null hypothesis (theory) is not supported by the data if T(Y) is an extreme (unlikely) observation given the range of plausible values predicted by the hypothesis (if the experimental evidence appears to be inconsistent with the theory).

Plausibility of different values of  $T(\cdot)$  under the theory  $H_0$ 

 $\rightarrow$  described by the distribution of T(Y) under the null hypothesis:

$$\mathbb{P}_{H_0}[T(Y) \in \cdot]$$

Suppose that we perform the experiment T(Y) and the result is T(Y) = t. The result t is judged to be incompatible with the hypothesis when

$$p = \mathbb{P}_{H_0}[T(Y) \geq t]$$

is small. The value p is called the p-value.

- Small values of p suggest that we have observed something which is unlikely to happen if  $H_0$  holds true.
- Large values of p suggest that what we have observed is plausible if  $H_0$  holds true.
- (Choice of T often guided by an alternative hypothesis  $H_1$ , under which T should be large)

Thus we reject the null hypothesis when p is small.

146 / 309

Victor Panaretos (EPFL) Linear Models

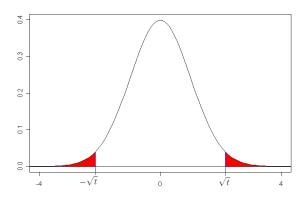
Example: Mean of a Normal Distribution

- Let  $X \sim \mathcal{N}(\mu, \sigma^2)$ , unknown mean, known variance
- $H_0: \mu = 0$
- Data:  $Y = (X_1, \ldots, X_n)$ ,  $X_i \stackrel{d}{=} X$ ,  $X_i$  indep  $X_j$  for  $i \neq j$ .
- Test statistic:  $T(Y) = \left(\frac{\sum_i X_i}{\sigma \sqrt{n}}\right)^2$ . (tends to be large when  $\mu \neq 0$ ).
- ullet Perform experiment (i.e., obtain values  $y=(x_1,\ldots,x_n)$ ) and observe T(y)=t.

Under the null hypothesis:  $T(Y) \stackrel{H_0}{\sim} \chi_1^2$ . Hence:

$$\begin{array}{rcl} p & = & \mathbb{P}_{H_0}[T(Y) \geq t] \\ & = & \mathbb{P}[\chi_1^2 \geq t] \\ & = & \mathbb{P}[\{\mathcal{N}(0,1) \leq -\sqrt{t}\} \cup \{\mathcal{N}(0,1) \geq \sqrt{t}\}]. \end{array}$$

Usually reject when p < 0.05.



- For continuous test statistics with everywhere positive densities, if we reject  $H_0$  whenever  $p<\alpha$ , then our (type I) error probability is  $\alpha$ .
  - $\hookrightarrow$  The probability of rejecting  $H_0$  when in fact  $H_0$  is true is  $\alpha$
- There is a close link with confidence intervals.
  - $\,\hookrightarrow\,$  We will only illustrate this link in a specific example

- 4日 4 4 日 4 日 4 日 4 日 4 日 4 日 9 4 (で

Example: Testing for  $c^{\top}\beta = 0$  in a Gaussian Regression

- Let  $Y \sim \mathcal{N}(X\beta, \sigma^2 I)$ , unknown  $\beta$ , unknown variance
- $H_0: c^{\top}\beta = 0$
- Data: (y, X).
- Test statistic:  $T(Y) = \left(\frac{c^{\top}\hat{\beta}}{S\sqrt{c^{\top}(X^{\top}X)^{-1}c}}\right)^2$

Suppose we observe  $\,T(y)= au\,$  and let  $\,W\,\sim\,t_{n-p}\,.\,$  Then,

$$p = \mathbb{P}_{H_0}[T(Y) \ge \tau] = \mathbb{P}[\{W \le -\sqrt{\tau}\} \cup \{W \ge \sqrt{\tau}\}].$$

Reject the null hypothesis if  $p < \alpha$ , some small  $\alpha$ .

• Identical to building a  $1-\alpha$  confidence interval for  $c^\top\beta$  based on  $\frac{c^\top\beta-c^\top\beta}{s\sqrt{c^\top(X^\top X)^{-1}c}}$  and rejecting the hypothesis  $H_0$  if and only if the interval does not contain zero.

Many many issues remain (this was just a reminder!)

- The role of an alternative hypothesis.
- How do we choose a test statistic?
- Are there optimal tests in a given situation?
- Simple and composite hypotheses.
- One and two-sided tests.
- Limitations of hypothesis testing . . .
- ...
- Review your 2nd year Probability/Statistics course!

# Nested Model Selection & ANOVA

Consider the model:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \varepsilon.$$

This will always have higher  $R^2$  than the sub-model:

$$y=\beta_0+\beta_1x_1+\varepsilon.$$

- Why? (think of geometry...)
- The question is: is the first model *significantly* better than the second one?
  - $\hookrightarrow$  i.e. does the first model explain the variation adequately enough, or should we incorporate extra explanatory variables? Need a quantitative answer.

Model is  $y = X\beta + \varepsilon$  with  $\varepsilon \sim \mathcal{N}_n(0, \sigma^2 I)$ . Estimate:

$$\hat{\beta} = (X^\top X)^{-1} X^\top y.$$

Interpretation:  $\hat{y}=X\hat{\beta}=Hy$  is the projection of y into the column space of X,  $\mathcal{M}(X)$ . This subspace has dimension p, when X is of full column rank p. Now for q< p write X in block notation as

$$X = (\begin{array}{cc} X_1 & X_2 \\ n \times q & n \times (p-q) \end{array}).$$

Interpretation:  $X_1$  is built by the first q columns of X and  $X_2$  by the rest. Similarly write  $\beta = (\beta_1 \ \beta_2)^{\top}$  so that:

$$y = X\beta + \varepsilon = (X_1 \ X_2) {eta_1 \choose eta_2} + \varepsilon = X_1\beta_1 + X_2\beta_2 + \varepsilon.$$

Our question can now be stated as:

• Is  $\beta_2 = 0$ ?

## Residual Sums of Squares

Let  $H_1 = X_1(X_1^{ op}X_1)^{-1}X_1^{ op}$ , and  $\hat{y}_1 = H_1y$ ,  $e_1 = y - \hat{y}_1$ .

Pythagoras tells us that:

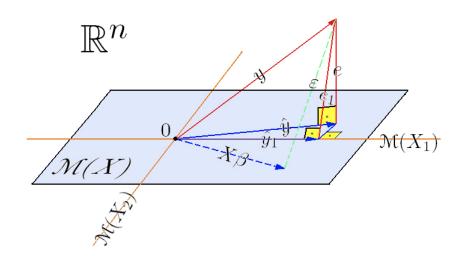
$$\underbrace{ \| \underline{y} - \hat{y}_1 \|^2 }_{RSS(\hat{\beta}_1) = \|e_1\|^2} = \underbrace{ \| \underline{y} - \hat{y} \|^2 }_{RSS(\hat{\beta}) = \|e\|^2} + \underbrace{ \| \hat{y} - \hat{y}_1 \|^2 }_{RSS(\hat{\beta}_1) - RSS(\hat{\beta}) = \|e - e_1\|^2}$$

Notice that  $RSS(\hat{\beta}_1) \geq RSS(\hat{\beta})$  always (think why!)

So the idea is simple: to see if it is worthwhile to include  $\beta_2$  we will compare how much larger  $RSS(\hat{\beta}_1)$  is compared to  $RSS(\hat{\beta})$ .

- ullet Equivalently, we can look at a ratio like  $\{RSS(\hat{eta}_1) RSS(\hat{eta})\}/RSS(\hat{eta})$
- To construct a test based on this quantity, we need to figure out distributions

. .



# **Theorem**

We have the following properties:

- (A)  $e e_1 \perp e$ ;
- (B)  $||e||^2 = RSS(\hat{\beta})$  and  $||e_1 e||^2 = RSS(\hat{\beta}_1) RSS(\hat{\beta})$  are independent;
- (C)  $||e||^2 \sim \sigma^2 \chi_{n-p}^2$ ;
- (D) under the hypothesis  $H_0: \beta_2 = 0$ ,  $||e_1 e||^2 \sim \sigma^2 \chi_{p-q}^2$ .

# Proof.

(A) holds since  $e-e_1=y-\hat{y}-y+\hat{y}_1=-\hat{y}+\hat{y}_1\in \mathcal{M}(X_1,X_2)$  but  $e\in [\mathcal{M}(X_1,X_2)]^{\perp}$ .

To show (B), we notice that

$$e_1 = (I - H_1)y = (I - H_1H)y$$

because  $\mathcal{M}(X_1) \subset \mathcal{M}(X_1, X_2)$ .

## proof continued

Therefore,

$$e - e_1 = (I - H)y - (I - H_1H)y = y - Hy - y + H_1Hy = (H_1 - I)Hy.$$

But recall that  $y \sim \mathcal{N}(X\beta, \sigma^2 I)$ . Therefore, to prove independence of  $e - e_1 = (H_1 - I)Hy$  and e = (I - H)y, we need to show that

$$(H_1 - I)H[\sigma^2 I](I - H)^{\top} = 0.$$

This is clearly the case since H(I - H) = 0, proving (B).

(C) follows immediately, since we have already proven last time that  $\forall \beta$  (even when  $\beta_2 = 0$ )

$$RSS(\hat{\beta}) \sim \sigma^2 \chi_{n-p}^2$$

# proof continued.

To prove (D), we note that

$$e-e_1=(H_1-I)Hy\sim \mathcal{N}\{(H_1-I)HXeta,\sigma^2\underbrace{(H_1-I)HH^{\top}(H_1-I)^{\top}}_{=H-H_1}\}.$$

But  $HX = X(X^{\top}X)^{-1}X^{\top}X = X$ . So, in block notation,

$$e - e_1 \sim \mathcal{N}((H_1 - I)X_1\beta_1 + (H_1 - I)X_2\beta_2, \sigma^2(H - H_1)).$$

Now  $(I-H_1)X_1eta_1=0$  always, since  $I-H_1$  projects onto  $\mathfrak{M}^\perp(X_1)$ . Therefore,

$$e-e_1\sim \mathcal{N}(0,\sigma^2(H-H_1)),$$
 when  $eta_2=0.$ 

Now observe that  $(H-H_1)^{\top}=(H-H_1)$  and  $(H-H_1)^2=(H-H_1)$  (because  $\mathcal{M}(X_1)\subset\mathcal{M}(X_1,X_2)$ ). Thus,

$$egin{aligned} e-e_1 &\sim \mathcal{N}(0,\sigma^2(H-H_1)^2) \implies e-e_1 &\stackrel{d}{=} & (H-H_1)arepsilon \ &\Longrightarrow RSS(\hat{eta}_1)-RSS(\hat{eta}) = \|e-e_1\|^2 &\stackrel{d}{=} & arepsilon^ op (H-H_1)arepsilon \sim \sigma^2\chi^2_{p-q}. \end{aligned}$$

since  $(H - H_1)$  is symmetric idempotent with trace p - q and  $\varepsilon \sim \mathcal{N}(0, \sigma^2 I_n)$ .

158 / 309

# Corollary

We conclude that, under the hypothesis  $\beta_2 = 0$ ,

$$rac{\left(rac{RSS(\hat{eta}_1) - RSS(\hat{eta})}{p-q}
ight)}{\left(rac{RSS(\hat{eta})}{n-p}
ight)} \sim F_{p-q,n-p}$$

#### The F-Test

Distributional results suggest the following test:

- Have  $Y \sim \mathcal{N}(X_1\beta_1 + X_2\beta_2, \sigma^2 I)$
- $H_0: \beta_2 = 0$
- Data:  $(y, X_1, X_2)$ .

$$ullet$$
 Test statistic:  $T=rac{\left(rac{RSS(\hat{eta}_1)-RSS(\hat{eta})}{p-q}
ight)}{\left(rac{RSS(\hat{eta})}{n-p}
ight)}$ 

Then, under  $H_0$ , it holds that  $T \sim F_{p-q,n-p}$ . Suppose we observe  $T = \tau$ . Then,

$$p = \mathbb{P}_{H_0}[T(Y) \ge au] = \mathbb{P}[F_{p-q,n-p} \ge au]$$

Reject the null hypothesis if  $p < \alpha$ , some small  $\alpha$ , usually 0.05.

Example: Nested Models in Cement Data

►We fitted the model:

$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \beta_3 x_3 + \beta_4 x_4 + \varepsilon$$

▶ But would the following simpler model be in fact adequate?

$$y = \beta_0 + \beta_1 x_1 + \varepsilon$$

- ▶ Intuitively: is the extra explanatory power of the "larger" model significant enough in order to justify its use instead of a simpler model? (i.e., is the residual vector for the "larger" model significantly smaller than that of the simpler model?)
- ▶ In this case, n = 13, p = 5, q = 2 and

$$RSS(\hat{\beta}) = 47.86, \qquad RSS(\hat{\beta}_1) = 1265.7$$

yielding

$$\tau = \frac{(1265.7 - 47.86)/(5 - 2)}{(47.86)/(13 - 5)} = 67.86$$

▶ $p = \mathbb{P}[F_{3,8} \ge 67.86] = 4.95 \times 10^{-6}$ , so we reject the hypothesis  $H_0: \theta_2 = \theta_3 = \theta_4 = 0$ .

►We can fit the quadratic model:

$$\texttt{MPG} = \beta_0 + \beta_1 \texttt{horsepower} + \beta_2 \texttt{horsepower}^2 + \varepsilon$$

▶ But would the model only with linear term suffice?

$$exttt{MPG} = eta_0 + eta_1 ext{horsepower} + arepsilon$$

- ▶ Intuitively: is the reduction of RSS afforded by the "complex" model substantial enough in order to justify its use instead of a simpler model?
- ▶ In this case, n = 392, p = 3, q = 2 and

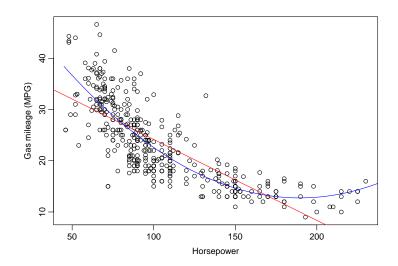
$$RSS(\hat{\beta}) = 7442, \qquad RSS(\hat{\beta}_1) = 9385.9$$

yielding

$$\tau = \frac{(9385.9 - 7442)/(3 - 2)}{7442/(392 - 3)} = 101.6$$

 $ightharpoonup p = \mathbb{P}[F_{1,389} \ge 101.6] = 2.2 \times 10^{-21}$ , so we reject the hypothesis  $H_0: \beta_2 = 0$ .

4 D > 4 P > 4 E > 4 E > 9 Q P



▶ Let  $\mathbf{1}, X_1, \ldots, X_r$  be groups of columns of X (the "terms"), such that

We have

$$y = X\beta + \varepsilon = \mathbf{1}\beta_0 + X_1\beta_1 + \cdots + X_r\beta_r + \varepsilon$$

- ▶ Would like to do the same "F-test investigation", but this time do it term-by-term. That is, we want to look at the following sequence of nested models:
  - $y = 1\beta_0 + \varepsilon$
  - $y = \mathbf{1}\beta_0 + X_1\beta_1 + \varepsilon$
  - $y = \mathbf{1}\beta_0 + X_1\beta_1 + X_2\beta_2 + \varepsilon$ 
    - :
  - $\bullet \ y = \mathbf{1}\beta_0 + X_1\beta_1 + X_2\beta_2 + \cdots + X_r\beta_r + \varepsilon$

Proceed similarly as before. Define:

- $ullet \ X_0 := m{1} \ ext{and} \ \mathcal{X}_k = (X_0 \ X_1 \ X_2 \ \dots \ X_k), \quad k \in \{0,\dots,r\}$
- $ullet \ \mathcal{H}_k := \mathcal{X}_k (\mathcal{X}_k^ op \mathcal{X}_k)^{-1} \mathcal{X}_k^ op, \quad k \in \{0,\dots,r\}$
- $\bullet \ \hat{y}_k := \mathcal{H}_k y, \quad k \in \{0, \ldots, r\}$
- $\bullet \ e_k = y \hat{y}_k, \quad k \in \{0, \ldots, r\}$
- Note that  $\hat{y}_0 = \bar{y} \mathbf{1}$ .
- ► As before, Pythagoras implies

$$\underbrace{ \frac{\|y - \hat{y}_0\|^2}{\|e_0\|^2}} = \underbrace{ \frac{\|y - \hat{y}_r\|^2}{\|e_r\|^2}} + \underbrace{ \frac{\|\hat{y} - \hat{y}_{r-1}\|^2}{\|e_r - e_{r-1}\|^2}} + \dots + \underbrace{ \frac{\|\hat{y}_1 - \hat{y}_0\|^2}{\|e_1 - e_0\|^2}}$$

$$= \underbrace{ \frac{\|e_r\|^2}{\|e_r\|^2}} + \sum_{k=0}^{r-1} \underbrace{ \frac{\|e_{k+1} - e_k\|^2}{\|e_{k+1} - e_k\|^2}}$$

with  $RSS_k$  the residual sum of squares for  $\hat{y}_k$ , with  $\nu_k$  degrees of freedom.

### Some observations:

- $RSS_k RSS_{k+1}$  is the reduction in residual sum of squares caused by adding  $X_{k+1}$ , when the model already contains  $X_0, \ldots, X_k$ .
- $RSS_r$  and  $\{RSS_k RSS_{k+1}\}_{k=0}^{r-1}$  are all mutually independent.
- Obviously,  $\nu_0 \geq \nu_1 \geq \nu_2 \geq \cdots \geq \nu_r$
- $\nu_{k+1} = \nu_k$  if  $X_{k+1} \in \mathcal{M}(\mathcal{X}_k)$ .
- ► Given this information, we want to see how adding each term in the model sequentially, affects the explanatory capacity of the model.

Terms	df	Residual	Terms	df	Reduction	F-test
		RSS	added		in RSS	
1	n-1	$RSS_0$				
$1, X_1$	$ u_1$	$RSS_1$	$X_1$	$n-1-\nu_1$	$RSS_0-RSS_1$	
$1, X_1, X_2$	$\nu_2$	$RSS_2$	$X_2$	$ u_1 -  u_2$	$RSS_1 - RSS_2$	
	-	•		•	•	
•	•	:		•	•	
$1, X_1, \ldots, X_r$	$ u_r$	$RSS_r$	$X_r$	$\nu_{r-1} - \nu_r$	$RSS_{r-1} - RSS_r$	

The F-statistic for testing the significance of the reduction in RSS when  $X_k$  is added to the model containing terms  $1, X_1, \ldots, X_k$  is

$$F_k = rac{(RSS_{k-1} - RSS_k)/(
u_{k-1} - 
u_k)}{RSS_r/
u_r},$$

and  $F_k \sim F_{
u_{k-1}u_k,
u_r}$  under the null hypothesis  $H_0: eta_k = 0$ .

Large values of  $F_k$  relative to the null distribution are evidence against  $H_0$ .

◆ロト ◆団ト ◆豆ト ◆豆 ・ りゅぐ

Example: Nested Sequence in Cement Data

- Reductions in overall sum of squares when sequentially entering terms  $x_1$ ,  $x_2$ ,  $x_3$  and  $x_4$ .
- Does adding extra variables improve model significantly?

	Df	Red Sum Sq	F value $( au)$	<i>p</i> -value
$\overline{x_1}$	1	1450.08	242.37	$2.88 \times 10^{-7}$
$x_2$	1	1207.78	201.87	$5.86 \times 10^{-7}$
$x_3$	1	9.79	1.64	0.2366
$x_4$	1	0.25	0.04	0.8441
Residual SSq	8	47.86		

▶ In this case, each term is a single column (variable).

- Significance of entering a term depends on how the sequence is defined: when entering terms in different order get different results! (why?)
- When a term is entered "early" and is significant, this does not tell us much (why?)
- When a term is entered "late" is significant, then this is quite informative (why?)
- ▶ Why is this true? Are there special cases when the order of entering terms doesn't matter?

▶ Consider terms  $X_0 = 1, X_1, X_2$  from X, so

$$X = (\underset{n \times 1}{X_0} \underset{n \times q_1}{X_1} \underset{n \times q_2}{X_2}), \quad \beta = (\underset{1 \times 1}{\beta_0} \underset{1 \times q_1}{\beta_1} \underset{1 \times q_2}{\beta_2})^\top$$

Assume orthogonality of terms, i.e.  $X_i^{\top} X_j = 0$ ,  $i \neq j$ Notice that in this case

$$\hat{\beta} = \begin{pmatrix} X_0^\top X_0 & 0 & 0 \\ 0 & X_1^\top X_1 & 0 \\ 0 & 0 & X_2^\top X_2 \end{pmatrix}^{-1} \begin{pmatrix} X_0 & X_1 & X_2 \end{pmatrix}^\top y$$

$$\implies \hat{\beta}_0 = \bar{y}, \ \hat{\beta}_1 = (X_1^\top X_1)^{-1} X_1^\top y, \ \hat{\beta}_2 = (X_2^\top X_2)^{-1} X_2^\top y$$

It follows that the reductions of sums of squares are unique, in the sense that they do not depend upon the order of entry of the terms in the model. (show this!) Intuition:  $X_i$  contains completely independent linear information from  $X_j$  for y,  $i \neq j$ 

Model Selection / Collinearity / Shrinkage

# Theory VS Practice

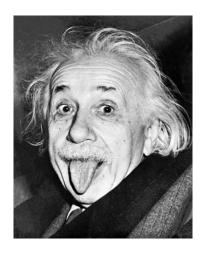
▶ **Theory:** We are given a relationship

$$y = X\beta + \varepsilon$$

and asked to provide estimators, tests, confidence intervals, optimality properties . . .

...and we can do it with complete success!

- **Practice**: We are given data (y, X) and suspect a linear relationship between y and some of the columns of X. We don't know a priori which exactly!
  - $\hookrightarrow$  Need to select a "most appropriate" subset of the columns of X
  - General principle: parsimony (Latin parsimonia: sparingness; simplicity and least number of requisites and assumptions; economy or frugality of components and associations).



'Everything should be made as simple as possible, but no simpler.'

 ♦ □ ▷ ▷ ♦ □ ▷ ▷ ♦ □ ▷ ♦



Occam's razor: *It is vain to do with more what can be done with fewer.*Given several explanations of the same phenomenon, we should prefer the simplest.

Graphical exploration → provides initial picture:

- plots of y against candidate variables;
- plots of transformations of y against candidate variables;
- plots of transformations of certain variables against *y*;
- plots of pairs of candidate variables.

This will often provide a starting point, but:

- Automatic Model Selection: Need objective model comparison criteria, as a screening device.
  - $\hookrightarrow$  We saw how to do an F-test, but what if models to be compared are not nested?
- Automatic Model Building: Situations when p large, so there are lots of possible models.

Consider design matrix X with p variables.

- 2<sup>p</sup> possible models!
- Denote set of all models generated by X by  $2^X$  (model powerset)
- ullet If wish to consider k different transformations of each variable, then p becomes (1+k)p
- Fast algorithms (branch and bound, leaps in R) exist to fit them, but they don't work for *large* p, and anyway . . .
- ... need criterion for comparison.

So given a collection of models, we need an automatic (objective) way to pick out a "best" one (unfortunately cannot look carefully at all of them, BUT NOTHING replaces careful scrutiny of the final model by an experienced researcher).

Many possible choices, none universally accepted. Some (classical) possibilities:

- Prediction error based criteria (CV)
- Information criteria (AIC, BIC, ...)
- Mallow's  $C_p$  statistic

Before looking at these, let's introduce terminology: Suppose that the truth is

$$y = X\beta + \varepsilon$$
 but with  $\beta_r = 0$  for some subset  $\beta_r$  of  $\beta$ .

- The true model contains only the columns for which  $\beta_r \neq 0$ 
  - $\hookrightarrow$  Equivalently, the true model uses  $X \oslash$  as the design matrix, the latter being the matrix of columns of X corresponding to non-zero coefficients.
- A correct model is the true model plus extra columns.
  - $\hookrightarrow$  Equivalently, a correct model has a design matrix  $X_{\diamondsuit}$ , such that  $\mathcal{M}(X_{\heartsuit}) \subset \mathcal{M}(X_{\diamondsuit})$ .
- A wrong model is a model that does not contain all the columns of the true model.
  - $\hookrightarrow$  Equivalently, a wrong model has a design matrix  $X_{\diamondsuit}$ , such that  $\mathcal{M}(X_{\heartsuit}) \cap \mathcal{M}(X_{\diamondsuit}) \neq \mathcal{M}(X_{\heartsuit})$ .

### **Expected Prediction Error**

- ► We may wish to choose a model by minimising the error we make on average, when predicting a future observation given our model.

  Our "experiment is":
  - Design matrix X
  - ullet response y at X

Every model  $f \in 2^X$ , will yield fitted values  $\hat{y}(f) = H_f y$ . And suppose we now obtain new independent responses  $y_+$  for the same "experimental setup" X. Then, one approach is to select the model

$$f^* = rg \min_{f \in 2^X} \underbrace{rac{1}{n} \mathbb{E}\left\{ \|y_+ - \hat{y}(f)\|^2 
ight\}}_{\Delta(f)},$$

where expectation is taken over both y and  $y_+$ .

◆ロト ◆節ト ◆意ト ◆意ト 意 めなべ

Let X be a design matrix, and let  $X_{\diamondsuit}$   $(n \times p)$  and  $X_{\heartsuit}$   $(n \times q)$  be matrices built using columns of X. Suppose that the true relationship between y and X is

$$y = \underbrace{X_{\heartsuit}\beta}_{\mu} + \epsilon$$

but we use the matrix  $X_{\diamondsuit}$  instead of  $X_{\heartsuit}$  (i.e., we fit a different model). Therefore our fitted values are

$$\hat{y} = (X_{\diamondsuit}^{\top} X_{\diamondsuit})^{-1} X_{\diamondsuit}^{\top} y = H_{\diamondsuit} y.$$

Now suppose that we obtain new observations  $y_+$  corresponding to the same design  $\boldsymbol{X}$ 

$$y_+ = X_{\odot}\beta + \varepsilon_+ = \mu + \varepsilon_+.$$

Then, observe that

$$y_{+} - \hat{y} = \mu + \varepsilon_{+} - H_{\Diamond}(\mu + \varepsilon)$$
  
=  $(I - H_{\Diamond})\mu + \varepsilon_{+} - H_{\Diamond}\varepsilon$ .

It follows that

$$\begin{split} \|y_{+} - \hat{y}\|^2 &= (y_{+} - \hat{y})^{\top} (y_{+} - \hat{y}) \\ &= \mu^{\top} (I - H_{\diamondsuit}) \mu + \varepsilon^{\top} H_{\diamondsuit} \varepsilon + \varepsilon_{+}^{\top} \varepsilon_{+} + [\text{cross terms}]. \end{split}$$

Since  $\mathbb{E}[\text{cross terms}] = 0$  (why?), we observe that

$$\Delta = \left\{ \begin{array}{ll} n^{-1}\mu^\top (I-H_\diamondsuit)\mu + (1+p/n)\sigma^2, & \text{if model wrong,} \\ (1+p/n)\sigma^2, & \text{if model correct,} \\ (1+q/n)\sigma^2, & \text{if model true.} \end{array} \right.$$

- Selecting a *correct model* instead of the *true model* brings in additional variance, because q < p.
- Selecting a wrong model instead of the true model results in bias, since  $(I H_{\diamondsuit})\mu \neq 0$  when  $\mu$  is not in the column space of  $X_{\diamondsuit}$ .
- Must find a balance between small variance (few columns in the model) and small bias (all columns in the model).

←ロト→□ト→重ト→重・のQで

#### Cross Validation

▶ Impossible to calculate  $\Delta$  (depends on unknown  $\mu$  and  $\sigma^2$ ), so we must find a proxy (estimator)  $\widehat{\Delta}$ .

Suppose that n is large so that we can split the data in two pieces:

- $X^*$ ,  $y^*$  used to estimate the model
- $\bullet$  X', y' used to estimate the prediction error for the model

The estimator of the prediction error will be

$$\widehat{\Delta} = (n')^{-1} ||y' - X' \widehat{\beta}^*||^2.$$

In practice n can be small and we often cannot afford to split the data (variance of  $\hat{\Delta}$  is too large).

Instead we use the *leave-one-out cross validation* sum of squares:

$$n\widehat{\triangle}_{CV} = CV = \sum_{j=1}^n (y_j - x_j^{ op} \hat{\pmb{\beta}}_{-j})^2,$$

where  $\hat{\beta}_{-j}$  is the estimate produced when dropping the jth case.

No need to perform n regressions since

$$CV = \sum_{j=1}^n rac{(y_j - x_j^{ op} \hat{oldsymbol{eta}})^2}{(1 - h_{jj})^2},$$

so the full regression may be used (show this!). Alternatively one may use a more stable version:

$$GCV = \sum_{j=1}^n rac{(y_j - x_j^ op \hat{oldsymbol{eta}})^2}{(1 - \mathsf{trace}(H)/n)^2},$$

where "G" stands for "generalised", and we guard against any  $h_{jj} \approx 1$ .

It holds that:

$$\mathbb{E}[GCV] = rac{\mu^+(I-H)\mu}{(1-p/n)^2} + rac{n\sigma^2}{1-p/n} pprox n\Delta.$$

▷ Suggests strategy: pick variables to minimise (G)CV.

Criteria can be obtained based on the notion of information (relative entropy).

• Same basic idea as for prediction error: aim to choose candidate model f(y) to minimise *information distance*:

$$\int \log \left\{ rac{g(y)}{f(y)} 
ight\} g(y) dy \geq 0,$$

where g(y) represents true model—equivalent to maximising expected log likelihood

$$\int \log f(y)g(y)dy.$$

Can show that (apart from constants) information distance is estimated by

$$\mathsf{AIC} = -2\hat{\ell} + 2p \quad (\equiv n \log \hat{\sigma}^2 + 2p \text{ in linear model})$$

where  $\hat{\ell}$  is maximised log likelihood for given model, and p is number of parameters.

Improved (corrected) version of AIC for regression problems:

$$\mathsf{AIC}_c \equiv \mathsf{AIC} + rac{2p(p+1)}{n-p-1}.$$

Also can use Bayes' information criterion

$$\mathsf{BIC} = -2\hat{\ell} + p\log n.$$

Mallows suggested

$$C_p = \frac{SS_p}{s^2} + 2p - n,$$

where  $SS_p$  is RSS for fitted model and  $s^2$  estimates  $\sigma^2$ .

- Comments:
  - AIC tends to choose models that are too complicated, buts AIC<sub>c</sub> cures this somewhat;
  - BIC is model selection consistent—if the true model is among those fitted, BIC chooses it with probability  $\to 1$  as  $n \to \infty$  (for fixed p).

←ロト ←団ト ← 豆ト ← 豆 ・ りへで

Victor Panaretos (EPFL)

Linear Models

#### Simulation Experiment

For each  $n \in \{10, 20, 40\}$  we construct  $20 \ n \times 7$  design matrices. We multiply each of these design matrices from the right with  $\beta = (1, 2, 3, 0, 0, 0, 0, 0)^{\top}$  and we add a  $n \times 1$  Gaussian error. We do this independently 50 times, obtaining 1000 regressions with p = 3. Selected models with 1 or 2 covariates have a bias term, and those with 4 or more covariates have excess variance.

n				Numbe	er of co	ovariate	S	
		1	2	3	4	5	6	7
10	$C_p$		131	504	91	63	83	128
	BIC		72	373	97	83	109	266
	AIC		52	329	97	91	125	306
	$AIC_c$	15	398	565	18	4		
20	$C_p$		4	673	121	88	61	53
	вĺС		6	781	104	52	30	27
	AIC		2	577	144	104	76	97
	$AIC_c$		8	859	94	30	8	1
40	$C_p$			712	107	73	66	42
	BIC			904	56	20	15	5
	AIC			673	114	90	69	54
	$AIC_c$			786	105	52	41	16

### Automatic Model Building

▶ We saw so far:

Automatic Model Selection: build a set of models and select the "best" one.

▶ Now look at different philosophy:

**Automatic Model Building**: construct a single model in a way that would hopefully provide a good one.

There are three standard methods for doing this:

- Forward Selection
- Backward Elimination
- Stepwise Selection

CAUTION: Although widely used, these have no theoretical basis. Element of arbitrariness . . .

- Forward selection: starting from the model with constant only,
  - add each remaining term separately to the current model;
  - if none of these terms is significant, stop; otherwise
  - **1** update the current model to include the most significant new term; go to step 1.
- Backward elimination: starting from the model with all terms,
  - if all terms are significant, stop; otherwise
  - ② update current model by dropping the term with the smallest *F* statistic; go to step 1.
- Stepwise: starting from an arbitary model,
  - consider three options—add a term, delete a term, swap a term in the model for one not in the model, and choose the most significant option;
  - 2 if model unchanged, stop; otherwise go to step 1.

### Some thoughts:

- Each procedure may produce a different model.
- Systematic search minimising Prediction Error, AIC or similar over all possible models is preferable— BUT not always feasible (e.g., when p large).
- Stepwise methods can fit 'highly significant' models to purely random data! Main problem is lack of objective function.
- Can be improved by comparing Prediction Error/AIC for different models at each step — uses objective function, but no systematic search.

Example: Nuclear Power Station Data

Data on light water reactors (LWR) constructed in the USA. The covariates are date (date construction permit issued), T1 (time between application for and issue of permit), T2 (time between issue of operating license and construction permit), capacity (power plant capacity in MWe), PR (=1 if LWR already present on site), NE (=1 if constructed in north-east region of USA), CT (=1 if cooling tower used), BW (=1 if nuclear steam supply system manufactured by Babcock–Wilcox), N (cumulative number of power plants constructed by each architect-engineer), PT (=1 if partial turnkey plant).

	cost	date	$T_1$	$T_2$	capacity	PR	NE	CT	BW	N	PT
1	460.05	68.58	14	46	687	0	1	0	0	14	0
2	452.99	67.33	10	73	1065	0	0	1	0	1	0
3	443.22	67.33	10	85	1065	1	0	1	0	1	0
4	652.32	68.00	11	67	1065	0	1	1	0	12	0
5	642.23	68.00	11	78	1065	1	1	1	0	12	0
6	345.39	67.92	13	51	514	0	1	1	0	3	0
7	272.37	68.17	12	50	822	0	0	0	0	5	0
8	317.21	68.42	14	59	457	0	0	0	0	1	0
:											
32	270.71	67.83	7	80	886	1	0	0	1	11	1

Example: Nuclear Power Station Data

	Full model			Back	ward	Forward		
	Est	t		Est	t		Est	$\overline{t}$
Int.	-14.24	-3.37	_	-13.26	-4.22	-	-7.62	-2.66
date	0.2	3.21		0.21	4.91		0.13	3.38
logT1	0.092	0.38						
logT2	0.29	1.05						
logcap	0.694	5.10		0.72	6.09		0.67	4.75
PR	-0.092	-1.20						
NE	0.25	3.35		0.24	3.36			
CT	0.12	1.82		0.14				
BW	0.033	0.33						
log(N)	-0.08	-1.74		-0.08	-2.11			
PT	-0.22	-1.83		-0.22	-1.99	-	-0.49	-4.77
s (df)	0.164	(21)		0.159	(25)		0.195	5 (28)

Recall:  $\hat{y}$  is projection of y onto  $\mathfrak{M}(X)$ 

 $\hookrightarrow$  Adding more variables (columns) into X "enlarges"  $\mathcal{M}(X)$  ... IF the rank increases by the # of new variables

#### Consider two extremes

- Adding a new variable  $X_{p+1} \in \mathcal{M}^{\perp}(X)$ 
  - $\hookrightarrow$  Gives us completely "new" information.
- Adding a new variable  $X_{p+1} \in \mathcal{M}(X)$ 
  - Gives no "new" information cannot even do least squares (why not?)

What if we are between the two extremes? What if

$$X_{p+1} \notin \mathcal{M}(X)$$
 but  $X(X^{\top}X)^{-1}X^{\top}X_{p+1} = HX_{p+1} \simeq X_{p+1}$ ?

We can certainly fit the regression, but what will happen?

◆ロト ◆節ト ◆意ト ◆意ト 意 めなべ

Using block matrix properties, have

$$\mathsf{var}(\hat{oldsymbol{eta}}) = \sigma^2 \left[ (X \ X_{p+1})^ op (X \ X_{p+1}) 
ight]^{-1}$$

with

$$\left[ (X \ X_{p+1})^{\top} (X \ X_{p+1}) \right]^{-1} = \left[ \begin{array}{cc} A & B \\ C & D \end{array} \right]$$

where

$$A = (X^{\top}X)^{-1} + (X^{\top}X)^{-1}X^{\top}X_{p+1} \\ \times (X_{p+1}^{\top}X_{p+1} - X_{p+1}^{\top}HX_{p+1})^{-1}X_{p+1}^{\top}X(X^{\top}X)^{-1},$$

$$B = -(X^{\top}X)^{-1}X^{\top}X_{p+1}(X_{p+1}^{\top}X_{p+1} - X_{p+1}^{\top}HX_{p+1})^{-1},$$

$$C = -(X_{p+1}^{\top}X_{p+1} - X_{p+1}^{\top}HX_{p+1})^{-1}X_{p+1}^{\top}X(X^{\top}X)^{-1},$$

$$D = (X_{p+1}^{\top}X_{p+1} - X_{p+1}^{\top}HX_{p+1})^{-1}.$$

## Problem of Multicollinearity

 $\frac{\mbox{Multicollinearity:}}{\mbox{dimension } q < p} \mbox{ when } p \mbox{ explanatories concentrate around a subspace of }$ 

[simplest case: pairs of variables that are correlated]

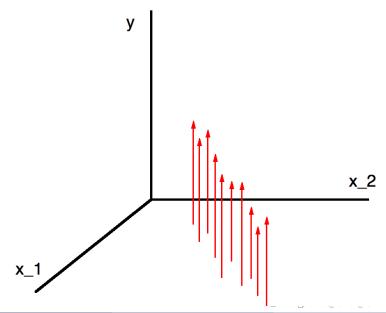
BUT: might exist even if pairs of variables appear uncorrelated!

## Can be caused by:

- Poor design [can try designing again],
- Inherent relationships [other remedies needed].

#### So what are the results?

- Huge variances of the estimators!
  - $\hookrightarrow$  Can even flip signs for different data, to give the impression of inverse effects.
- Individual coefficients insignificant:
  - $\hookrightarrow$  *t*-test *p*-values inflated.
- But global *F*-test might give significant result!



Simple first steps:

- Look at scatterplots,
- Look at correlation matrix of explanatories,

Might not reveal more complex linear constraints, though.

• Look at the variance inflation factors:

$$VIF_j = rac{\mathsf{var}(\hat{eta}_j) \|X_j\|^2}{\sigma^2} = \|X_j\|^2 \left[ (X^{\top} X)^{-1} \right]_{jj}.$$

Can show that

$$VIF_j = \frac{1}{1 - R_j^2}$$

where  $R_j^2$  is the coefficient of determination for the regression

$$X_j = \beta_{0,j} + \beta_{1,j} X_1 + \cdots + \beta_{j-1,j} X_{j-1} + \beta_{j+1,j} X_{j+1} + \cdots + \beta_{p,j} X_p + \varepsilon,$$

measuring linear dependence of  $X_j$  on the other columns of X.

(ロ) (固) (量) (量) (量) の(で)

Let  $X_{-j}$  be the design matrix without the j-th variable. Then

$$R_{j}^{2} = \frac{\|X_{-j}(X_{-j}^{\top}X_{-j})^{-1}X_{-j}^{\top}X_{j}\|^{2}}{\|X_{j}\|^{2}} \in [0, 1]$$

is close to 1 if  $\underbrace{X_{-j}(X_{-j}^{\top}X_{-j})^{-1}X_{-j}}_{H_{-j}}X_j \simeq X_j.$ 

Large values of  $VIF_j$  indicate that  $X_j$  is linearly dependent on the other columns of the design matrix.

Interpretation: how much the variance is inflated when including variable j as compared to the variance we would obtain if  $X_j$  were orthogonal to the other variables—how much worse are we doing as compared to the ideal case.

Rule of thumb:  $VIF_j > 5$  or  $VIF_j > 10$  considered to be "large".

Consider the spectral decomposition of  $X^\top X$ ,  $X^\top X = U \Lambda U^\top$  with  $\Lambda = \text{diag}\{\lambda_1,\ldots,\lambda_p\}$  and  $U^\top U = I$ . Then

$$\operatorname{\mathsf{rank}}(X^{ op}X) = \#\{j: \lambda_j 
eq 0\}, \qquad \det(X^{ op}X) = \prod_{j=1}^p \lambda_j.$$

Hence "small"  $\lambda_j$ 's mean "almost" reduced rank, revealing the effect of collinearity. Measure using *condition index*:

$$\mathit{CI}_j(X^{\top}X) := \sqrt{\lambda_{\mathsf{max}}/\lambda_j}$$

Global "instability" measured by the condition number,

$$CN(X^{\top}X) = \sqrt{\lambda_{\mathsf{max}}/\lambda_{\mathsf{min}}}$$

Rule of thumb: CN > 30 indicates moderate to significant collinearity, CN > 100 indicates severe collinearity (choices vary).

#### Remedies?

If design faulty, may redesign.

- Otherwise? Inherent relationships between explanatories.
  - Variable deletion attempt to remove problematic variables
    - $\rightarrow$  E.g., by backward elimination.
  - ullet Choose an orthogonal basis for  $\mathcal{M}(X)$  and use its elements as explanatories
    - $\rightarrow$  Use columns of U from spectrum,  $X^{\top}X = U\Lambda U^{\top}$
    - $\rightarrow\,$  OK for prediction
    - ightarrow Problem: lose interpretability

Other approaches?

Example: Body Fat Data

Body fat is measure of health  $\rightarrow$  not easy to measure! Collect 252 measurements on body fat and some explanatory variables.

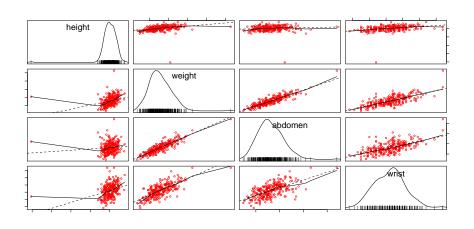
Can we use measuring tape and scales only to find body fat? Explanatory variables:

- age
- weight
- height
- biceps

- neck
- chest
- abdomen
- forearm

- hip
- thigh
- knee a
- ankle
- wrist

Some Scatterplots [library(car);scatterplot.matrix( $\dots$ )]



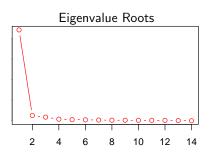
Looks like we're in trouble. Let's go ahead and fit anyway . . .

Estimate	Std. Error	t value	$\Pr(> t )$
-18.1885	17.3486	-1.05	0.2955
0.0621	0.0323	1.92	0.0562
-0.0884	0.0535	-1.65	0.0998
-0.0696	0.0960	-0.72	0.4693
-0.4706	0.2325	-2.02	0.0440
-0.0239	0.0991	-0.24	0.8100
0.9548	0.0864	11.04	0.0000
-0.2075	0.1459	-1.42	0.1562
0.2361	0.1444	1.64	0.1033
0.0153	0.2420	0.06	0.9497
0.1740	0.2215	0.79	0.4329
0.1816	0.1711	1.06	0.2897
0.4520	0.1991	2.27	0.0241
-1.6206	0.5349	-3.03	0.0027
	-18.1885 0.0621 -0.0884 -0.0696 -0.4706 -0.0239 0.9548 -0.2075 0.2361 0.0153 0.1740 0.1816 0.4520	-18.1885         17.3486           0.0621         0.0323           -0.0884         0.0535           -0.0696         0.0960           -0.4706         0.2325           -0.0239         0.0991           0.9548         0.0864           -0.2075         0.1459           0.2361         0.1444           0.0153         0.2420           0.1740         0.2215           0.1816         0.1711           0.4520         0.1991	-18.1885         17.3486         -1.05           0.0621         0.0323         1.92           -0.0884         0.0535         -1.65           -0.0696         0.0960         -0.72           -0.4706         0.2325         -2.02           -0.0239         0.0991         -0.24           0.9548         0.0864         11.04           -0.2075         0.1459         -1.42           0.2361         0.1444         1.64           0.0153         0.2420         0.06           0.1740         0.2215         0.79           0.1816         0.1711         1.06           0.4520         0.1991         2.27

 $R^2 = 0.749$ , F-test:  $p < 2.2 \times 10^{-16}$ .

	Estimate	D <sub>*</sub> ( >  ±  )	Estimate	D <sub>*</sub> ( >  ± )
	Estimate	Pr(> t )		Pr(> t )
(Intercept)	-32.6564	0.1393	-1.2221	0.9730
age	0.1048	0.0153	0.0256	0.6252
weight	-0.1285	0.0502	-0.0237	0.8223
height	-0.0666	0.5207	-0.1005	0.7284
neck	-0.5086	0.0721	-0.4619	0.2635
chest	0.0168	0.9002	-0.0910	0.5877
abdomen	0.9750	0.0000	0.8924	0.0000
hip	-0.2891	0.1265	-0.0265	0.9130
thigh	0.3850	0.0565	0.0334	0.8793
knee	0.2218	0.5111	-0.1310	0.7366
ankle	0.4377	0.0694	-0.5037	0.3516
biceps	-0.1297	0.5485	0.4458	0.1179
forearm	0.8871	0.0174	0.2247	0.3750
wrist	-1.7378	0.0309	-1.5902	0.0560

	VIF			CI
age	2.25	_	1	1.00
weight	33.51		2	17.47
height	1.67		3	25.30
neck	4.32		4	58.61
chest	9.46		5	83.59
abdomen	11.77		6	100.63
hip	14.80		7	137.90
thigh	7.78		8	175.29
knee	4.61		9	192.62
ankle	1.91	-	10	213.01
biceps	3.62	-	11	228.16
forearm	2.19	-	12	268.21
wrist	3.38	-	13	555.67



Condition Number  $\simeq 556$ !

### Variable Deletion: Backward Elimination

Multiple R-Squared: 0.7466, F-statistic p-value: < 2.2e-16

	Estimate	Std. Error	t value	Pr(> t )	VIF
(Intercept)	-22.6564	11.7139	-1.93	0.0543	
age	0.0658	0.0308	2.14	0.0336	2.05
weight	-0.0899	0.0399	-2.25	0.0252	18.82
neck	-0.4666	0.2246	-2.08	0.0388	4.08
abdomen	0.9448	0.0719	13.13	0.0000	8.23
hip	-0.1954	0.1385	-1.41	0.1594	13.47
thigh	0.3024	0.1290	2.34	0.0199	6.28
forearm	0.5157	0.1863	2.77	0.0061	1.94
wrist	-1.5367	0.5094	-3.02	0.0028	3.09

Define Z=XU as design matrix.  $R^2$ =0.749, F-test p-value<2.2  $\times$   $10^{-16}$ 

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	-18.1885	17.3486	-1.05	0.2955
Z[, 1]	-0.1353	0.0619	-2.19	0.0297
Z[, 2]	-0.0168	0.0916	-0.18	0.8546
Z[, 3]	0.2372	0.1070	2.22	0.0276
Z[, 4]	-0.7188	0.0571	-12.58	0.0000
Z[, 5]	0.0248	0.0827	0.30	0.7649
Z[, 6]	0.4546	0.1001	4.54	0.0000
Z[, 7]	0.5903	0.1366	4.32	0.0000
Z[, 8]	-0.1207	0.1742	-0.69	0.4890
Z[, 9]	-0.0836	0.1914	-0.44	0.6627
Z[, 10]	0.5043	0.2082	2.42	0.0162
Z[, 11]	-0.5735	0.2254	-2.54	0.0116
Z[, 12]	0.3007	0.2628	1.14	0.2536
Z[, 13]	1.5168	0.5447	2.78	0.0058

- Eigenvector approach rotates space so as to "free" the dependence of one coefficient  $\beta_j$  on others  $\{\beta_i\}_{i\neq j}$
- $\hookrightarrow$  Imposes constraint on X (orthogonal columns)

```
Problem: lose interpretability! (prediction OK)
```

- Example: most significant "rotated" term in fat data: Z[,4]=-0.01\*age -0.058\*weight -0.011\*height +0.46\*neck -0.144\*chest -0.441\*abdomen +0.586\*hip +0.22\*thigh -0.197\*knee -0.044\*ankle -0.07\*biceps -0.33\*forearm -0.249\*wrist
- Other approach to reduce this strong dependence?
  - $\hookrightarrow$  Impose constraint on  $\beta$ ! How? (introduces bias)

Multicollinearity problem is that  $\det \left[ (X^\top X)^{-1} \right] \approx 0$  [i.e.  $X^\top X$  almost not invertible]

A Solution: add a "small amount" of a full rank matrix to  $X^{\top}X$ .

For reasons to become clear soon, we standardise the design matrix:

- Write  $X = (\mathbf{1} \ W), \ \beta = (\beta_0 \ \gamma)^{\top}$
- ullet Recentre/rescale the covariates defining:  $Z_j = \frac{\sqrt{n}}{\operatorname{sd}(W_j)} (W_j \mathbf{1}\overline{W}_j)$ 

  - $\hookrightarrow$  Interpretation of  $\beta_j$  slightly different: not "mean impact on response per unit change of explanatory variable", but now "mean impact on response per unit deviation of explanatory variable from its mean, measured in units of standard deviation"
- ullet The  $Z_j$  are all orthogonal to  ${f 1}$  and are of unit norm.

- Since  $Z_j \perp \mathbf{1}$  for all, j, we can estimate  $\beta_0$  and  $\gamma$  by two separate regressions (orthogonality).
- Least squares estimators become

$$\hat{\beta}_0 = \overline{Y}, \quad \hat{\gamma} = (Z^\top Z)^{-1} Z^\top Y.$$

ullet Ridge regression replaces  $Z^{\top}Z$  by  $Z^{\top}Z + \lambda I$  (i.e. adds a "ridge")

$$\hat{eta}_0 = \overline{Y}, \quad \hat{\gamma} = (Z^{\top}Z + \lambda I)^{-1}Z^{\top}Y$$

Adding  $\lambda I$  to  $Z^{\top}Z$  makes inversion more stable  $\hookrightarrow \lambda$  called *ridge parameter*.

Ridge Regression: Shrinkage Viewpoint

 $\rightarrow$  Ridge term  $\lambda I$  seems slightly ad-hoc. Motivation?

 $\hookrightarrow$  Can see that  $(\hat{\beta}_0 \quad \hat{\gamma}) = (\overline{Y} \quad (Z^\top Z + \lambda I)^{-1} Z^\top Y)$  minimizes

$$||Y - \beta_0 \mathbf{1} - Z\gamma||_2^2 + \lambda ||\gamma||_2^2$$

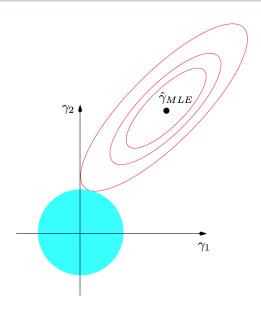
or equivalently

$$\|Y-eta_0\mathbf{1}-Z\gamma\|_2^2$$
 subject to  $\sum_{j=1}^{p-1}\gamma_j^2=\|\gamma\|_2^2\leq r(\lambda)$ 

instead of least squares estimator which minimizes

$$||Y-\beta_0\mathbf{1}-Z\gamma||_2^2.$$

Idea: in the presence of collinearity, coefficients are ill-defined: a wildly positive coefficient can be cancelled out by a largely negative coefficient (many coefficient combinations can produce the same effect). By imposing a size constraint, we limit the possible coefficient combinations!



## Proposition

Let  $Z_{n\times q}$  be a matrix of rank  $r\leq q$  with centred column vectors of unit norm. Given  $\lambda>0$ , the unique minimiser of

$$Q(\hat{eta}_0, \hat{\gamma}) = \|y - \hat{eta}_0 \mathbf{1} - Z\hat{\gamma}\|_2^2 + \lambda \|\hat{\gamma}\|_2^2$$

is

$$(\hat{\beta}_0, \hat{\gamma}) = (\overline{y}, (Z^\top Z + \lambda I)^{-1} Z^\top y).$$

## Proof.

Write

$$y = \underbrace{\left(y - \bar{y}\mathbf{1}\right)}_{=y^* \in \mathcal{M}^{\perp}(\mathbf{1})} + \underbrace{\bar{y}\mathbf{1}}_{\in \mathcal{M}(\mathbf{1})}$$

Note also that by assumption  $1 \in \mathcal{M}^{\perp}(Z)$ . Therefore by Pythagoras' theorem

∈M<sup>⊥</sup>(1)

$$||y - \hat{\beta}_0 \mathbf{1} - Z\hat{\gamma}||_2^2 = ||\underbrace{(\bar{y} - \hat{\beta}_0)\mathbf{1}}_2 + \underbrace{(y^* - Z\hat{\gamma})}_2||_2^2 = ||(\bar{y} - \hat{\beta}_0)\mathbf{1}||_2^2 + ||(y^* - Z\hat{\gamma})||_2^2.$$

$$\text{Therefore, } \min_{\hat{\beta}_0, \hat{\gamma}} Q(\hat{\beta}_0, \hat{\gamma}) = \min_{\hat{\beta}_0} \|(\bar{y} - \hat{\beta}_0)\mathbf{1}\|_2^2 + \min_{\hat{\gamma}} \left\{ \|(y^* - Z\hat{\gamma})\|_2^2 + \lambda \|\hat{\gamma}\|_2^2 \right\}$$

Clearly,  $\arg\min_{\hat{\theta}_0}\|(\bar{y}-\hat{\beta}_0)\mathbf{1}\|_2^2=\bar{y}$  while the second component can be written

$$\min_{\hat{m{\gamma}} \in \mathbb{R}^q} \left\| inom{Z}{\sqrt{\lambda} \, I_{q imes q}} \hat{m{\gamma}} - inom{y^*}{n imes 1} igg 0_{q imes 1} 
ight\}_2$$

using block notation. This is the usual least squares problem with solution

$$\left[ (Z^\top, \sqrt{\lambda} I_{q \times q}) \binom{Z}{\sqrt{\lambda} I_{q \times q}} \right]^{-1} (Z^\top, \sqrt{\lambda} I_{q \times q}) \binom{y^*}{\mathbf{0}_{q \times 1}} = (Z^\top Z + \lambda I)^{-1} Z^\top y^*$$

Note that  $Z^{\top}Z + \lambda I$  is indeed invertible. Writing  $Z^{\top}Z = U\Lambda U^{\top}$ , we have

$$Z^{\top}Z + \lambda I = U \Lambda U^{\top} + U(\lambda I_{g \times g}) U^{\top} = U(\Lambda + \lambda I_{g \times g}) U^{\top}$$

$$\text{ and } \Lambda = \operatorname{diag}\{\underbrace{\lambda_1, \dots, \lambda_r}, \underbrace{\lambda_{r+1}, \dots, \lambda_q}\} \; (Z^\top Z \succeq 0 \; \& \; \operatorname{rank}(Z^\top Z) = \operatorname{rank}(Z)).$$

To complete the proof, observe that 
$$Z^{\top}y^* = Z^{\top}y - \bar{y}Z^{\top}\mathbf{1} = Z^{\top}y$$
.

To complete the proof, observe that  $Z \cdot y^* = Z \cdot y - \overline{y}Z \cdot 1 = Z \cdot y$ .

### The Effect of Shrinkage

Note that if the SVD of Z is  $Z=V\Omega U^{\top}$ , last steps of previous proof may be used to show that

$$\hat{\gamma} = \sum_{j=1}^q rac{\omega_j}{\omega_j^2 + \lambda} (v_j^ op y) u_j,$$

where the  $v_i$ s and  $u_i$ s are the columns of V and U, respectively.

Compare this to the ordinary least squares solution, when  $\lambda = 0$ :

$$\hat{\gamma} = \sum_{j=1}^q rac{1}{\omega_j} (\mathit{v}_j^ op y) \mathit{u}_j,$$

which is not even defined if Z is of reduced rank.

Role of  $\lambda$  is to reduce the size of  $1/\omega_j$  when  $\omega_j$  becomes very small.

## **Proposition**

Let  $\hat{\gamma}$  be the ridge regression estimator of  $\gamma$ . Then

$$\mathit{bias}(\hat{\gamma}, \gamma) = -\lambda \left(Z^{\top}Z + \lambda I_q\right)^{-1} \gamma$$

and

$$cov(\hat{\gamma}) = \sigma^2 (Z^\top Z + \lambda I)^{-1} Z^\top Z (Z^\top Z + \lambda I)^{-1}.$$

## Proof.

Since  $\mathbb{E}(\hat{\gamma}) = (Z^{\top}Z + \lambda I)^{-1}Z^{\top}\mathbb{E}(y) = (Z^{\top}Z + \lambda I)^{-1}Z^{\top}Z\gamma$ , the bias is

$$\begin{aligned} \mathsf{bias}(\hat{\gamma}, \gamma) &= & \mathbb{E}(\hat{\gamma}) - \gamma = \{(Z^\top Z + \lambda I)^{-1} Z^\top Z - I\} \gamma \\ &= & \mathbb{E}(\hat{\gamma}) - \gamma = \{(Z^\top Z + \lambda I)^{-1} Z^\top Z - I\} \gamma \\ &= & \left\{ \left(\frac{1}{\lambda} Z^\top Z + I\right)^{-1} \left(\frac{1}{\lambda} Z^\top Z + I - I\right) - I\right\} \gamma \\ &= & \left\{ I - \left(\frac{1}{\lambda} Z^\top Z + I\right)^{-1} - I\right\} \gamma = -\left(\frac{1}{\lambda} Z^\top Z + I\right)^{-1} \gamma. \end{aligned}$$

The covariance term is obvious.

# Corollary

Assume that  $\operatorname{rank}(Z_{n \times q}) = q$  and let

$$\hat{\gamma} = (Z^{\top}Z)^{-1}Z^{\top}y$$
 &  $\hat{\gamma}_{\lambda} = (Z^{\top}Z + \lambda I)^{-1}Z^{\top}y$ 

be the least squares estimator and ridge estimator, respectively. Then,

$$\mathbb{E}\left\{ (\hat{\gamma} - \gamma)(\hat{\gamma} - \gamma)^{\top} \right\} - \mathbb{E}\left\{ (\hat{\gamma}_{\lambda} - \gamma)(\hat{\gamma}_{\lambda} - \gamma)^{\top} \right\} \succeq 0$$

for all  $\lambda \in (0, 2\sigma^2/||\gamma||^2)$ .

Ridge estimator uniformly better than least squares! How can this be? (What happened to Gauss-Markov?)

- Gauss-Markov only covers unbiased estimators but ridge estimator biased.
- A bit of bias can improve the MSE by reducing variance.
- ullet Also, there is a catch: the "right" range for  $\lambda$  depends on unknowns.
- $\bullet$  Choosing a good  $\lambda$  is all about balancing bias and variance.

Victor Panaretos (EPFL) Linear Models 215 / 309

### Proof.

From our bias/variance calculations on the ridge estimator, we have

$$\mathbb{E}\left\{(\boldsymbol{\hat{\gamma}}-\boldsymbol{\gamma})(\boldsymbol{\hat{\gamma}}-\boldsymbol{\gamma})^\top\right\}-\mathbb{E}\left\{(\boldsymbol{\hat{\gamma}}_{\lambda}-\boldsymbol{\gamma})(\boldsymbol{\hat{\gamma}}_{\lambda}-\boldsymbol{\gamma})^\top\right\}=$$

$$\sigma^{2}(\boldsymbol{Z}^{\top}\boldsymbol{Z})^{-1} - (\boldsymbol{Z}^{\top}\boldsymbol{Z} + \lambda\boldsymbol{I})^{-1}\sigma^{2}\boldsymbol{Z}^{\top}\boldsymbol{Z}(\boldsymbol{Z}^{\top}\boldsymbol{Z} + \lambda\boldsymbol{I})^{-1} - \lambda^{2}\left(\boldsymbol{Z}^{\top}\boldsymbol{Z} + \lambda\boldsymbol{I}\right)^{-1}\boldsymbol{\gamma}\boldsymbol{\gamma}^{\top}\left(\boldsymbol{Z}^{\top}\boldsymbol{Z} + \lambda\boldsymbol{I}\right)^{-1}$$

$$= \lambda(\boldsymbol{Z}^{\top}\boldsymbol{Z} + \lambda\boldsymbol{I})^{-1}\left(\sigma^{2}(2\boldsymbol{I} + \lambda(\boldsymbol{Z}^{\top}\boldsymbol{Z})^{-1}) - \lambda\boldsymbol{\gamma}\boldsymbol{\gamma}^{\top}\right)(\boldsymbol{Z}^{\top}\boldsymbol{Z} + \lambda\boldsymbol{I})^{-1}.$$

To go from 2nd to 3rd line, we wrote

$$\sigma^{2}(\boldsymbol{Z}^{\top}\boldsymbol{Z})^{-1} = \sigma^{2}(\boldsymbol{Z}^{\top}\boldsymbol{Z} + \lambda \boldsymbol{I})^{-1}(\boldsymbol{Z}^{\top}\boldsymbol{Z} + \lambda \boldsymbol{I})(\boldsymbol{Z}^{\top}\boldsymbol{Z})^{-1}(\boldsymbol{Z}^{\top}\boldsymbol{Z} + \lambda \boldsymbol{I})(\boldsymbol{Z}^{\top}\boldsymbol{Z} + \lambda \boldsymbol{I})^{-1}$$
$$= (\boldsymbol{Z}^{\top}\boldsymbol{Z} + \lambda \boldsymbol{I})^{-1}(\sigma^{2}\boldsymbol{Z}^{\top}\boldsymbol{Z} + 2\sigma^{2}\lambda \boldsymbol{I} + \sigma^{2}\lambda^{2}(\boldsymbol{Z}^{\top}\boldsymbol{Z})^{-1})(\boldsymbol{Z}^{\top}\boldsymbol{Z} + \lambda \boldsymbol{I})^{-1}$$

and did the tedious (but straighforward) algebra. The final term can be made positive definite if

$$2\sigma^2 I + \sigma^2 \lambda (Z^\top Z)^{-1} - \lambda \gamma \gamma^\top \succeq 0.$$

Noting that we can always write

$$I = rac{oldsymbol{\gamma}oldsymbol{\gamma}^ op}{\|oldsymbol{\gamma}\|^2} + \sum_{j=1}^{q-1} oldsymbol{ heta}_j oldsymbol{ heta}_j^ op$$

for  $\{\gamma/||\gamma||, \theta_1, ..., \theta_{q-1}\}$  an orthonormal basis of  $\mathbb{R}^q$  we see that  $\lambda \in (0, 2\sigma^2/||\gamma||^2)$  suffices for positive definiteness to hold true.

Victor Panaretos (EPFL) Linear Models 216 / 309

Role of  $\lambda$ : Regulates Bias-Variance tradeoff

- $\lambda \uparrow$  decreases variance (e.g. due to collinearity) but increases bias
- $\lambda \downarrow$  decreases bias but variance inflated if collinearity exists

#### Recall:

$$\mathbb{E}||\hat{\gamma} - \gamma||^2 = \underbrace{\mathbb{E}||\hat{\gamma} - \mathbb{E}\hat{\gamma}||^2}_{Variance = trace[cov(\hat{\gamma})]} + \underbrace{\|\mathbb{E}\hat{\gamma} - \gamma\|^2}_{Bias^2} + \underbrace{2(\mathbb{E}\hat{\gamma} - \gamma)^{\top}\mathbb{E}[\hat{\gamma} - \mathbb{E}\hat{\gamma}]}_{=0}$$

Note that if  $Z^{\top}Z = U\Omega U^{\top}$ , then  $trace(cov(\hat{\gamma})) = \sum_{j=1}^q \frac{\omega_i}{\omega_i^2 + \lambda} \sigma^2$ 

So choose  $\lambda$  so as to optimally  $\underline{\text{increase bias}}/\underline{\text{decrease variance}}$ 

Use cross validation!



Motivated from Ridge Regression formulation can consider:

min! 
$$\|Y - \beta_0 \mathbf{1} - Z\gamma\|_2^2$$
 subject to  $\sum_{j=1}^{p-1} |\gamma_j| = \|\gamma\|_1 \le r(\lambda)$   $\iff$  min!  $\|Y - \beta_0 \mathbf{1} - Z\gamma\|_2^2 + \lambda \|\gamma\|_1$ .

Shrinks coefficient size by different version of magnitude.

- ullet Resulting estimator non–linear in Y
- $\bullet$  No explicit form available (unless  $Z^\top Z = I)$  , needs quadratic programming algorithm

Why choose a different type of norm? Because  $L^1$  penalty (almost) produces a "continuous" model selection!

When the explanatory variables are orthogonal (i.e.  $Z^{\top}Z = I$ ), then the LASSO exactly performs model selection via thresholding:

## **Theorem**

Consider the linear model

$$Y_{n\times 1} = \beta_0 \frac{1}{1\times 1} + Z_{n\times (p-1)(p-1)\times 1} + \varepsilon_{n\times 1}$$

where  $Z^{\top}\mathbf{1}=\mathbf{0}$  and  $Z^{\top}Z=I$  . Let  $\hat{\gamma}$  be the least squares estimator of  $\gamma$  ,

$$\hat{\gamma} = (Z^\top Z)^{-1} Z^\top Y = Z^\top Y.$$

Then, the unique solution to the LASSO problem

$$\min_{oldsymbol{eta}_0 \in \mathbb{R}, oldsymbol{\gamma} \in \mathbb{R}^{p-1}} \left\{ \|Y - oldsymbol{eta}_0 \mathbf{1} - Z oldsymbol{\gamma} \|_2^2 + \lambda \|oldsymbol{\gamma}\|_1 
ight\}$$

is given by  $(\hat{eta}_0,\check{\gamma})=(eta_0,\check{\gamma}_1,\ldots,\check{\gamma}_{p-1})$ , defined as

$$\hat{eta}_0 = ar{Y} \qquad \& \qquad \check{\gamma}_i = sgn(\hat{\gamma}_i) \left( |\hat{\gamma}_i| - rac{\lambda}{2} 
ight)_{\perp}, \quad i = 1,...,p-1.$$

## Proof.

Note that since  $Z^{\top}\mathbf{1}=0$  and since  $\beta_0$  does not appear in the  $L^1$  penalty, we have

$$\hat{\beta}_0 = (\mathbf{1}^\top \mathbf{1})^{-1} \mathbf{1} Y = \bar{Y}.$$

Thus, the LASSO problem reduces to

$$\min_{\beta_0\in\mathbb{R},\gamma\in\mathbb{R}^{p-1}}\left\{\|Y-\beta_0\mathbf{1}-Z\gamma\|_2^2+\lambda\|\gamma\|_1\right\}=\min_{\gamma\in\mathbb{R}^{p-1}}\left\{\|u-Z\gamma\|_2^2+\lambda\|\gamma\|_1\right\}.$$

where  $u=Y-ar{Y}\mathbf{1}$  for tidiness. Expanding the squared norm gives

$$\|u - Z\hat{\gamma}\|_2^2 = u^\top u - 2u^\top Z\gamma + \gamma^\top \underbrace{\left(\underline{Z}^\top Z\right)}_{=I} \gamma = u^\top u - 2\underbrace{Y}^\top \underline{Z}\gamma + 2\,\bar{Y}\underbrace{1}^\top \underline{Z}\gamma + \gamma^\top \gamma$$

Since  $u^{ op}u$  does not depend on  $\gamma$ , we see that the LASSO objective function is

$$-2\hat{\gamma}^{\top}\gamma + ||\gamma||_{2}^{2} + \lambda ||\gamma||_{1}.$$

Clearly, this has the same minimizer if multiplied across by 1/2, i.e.

$$-\hat{\gamma}^{\top}\gamma + \frac{1}{2}\|\gamma\|_{2}^{2} + \frac{1}{2}\lambda\|\gamma\|_{1} = \sum_{i=1}^{p-1} \left(-\hat{\gamma}_{i}\gamma_{i} + \frac{1}{2}\gamma_{i}^{2} + \frac{\lambda}{2}|\gamma_{i}|\right).$$

Notice that we now have a sum of p-1 objective functions, each depending only on one  $\gamma_i$ . We can thus optimise each separately. That is, for any given  $i \leq p-1$ , we must minimise

$$-\hat{\gamma}_i\gamma_i+\frac{1}{2}\gamma_i^2+\frac{\lambda}{2}|\gamma_i|.$$

We distinguish 3 cases:

- Case  $\hat{\gamma}_i = 0$ . In this case, the objective function becomes  $\frac{1}{2}\gamma_i^2 + \frac{\lambda}{2}|\gamma_i|$  and it is clear that it is minimised when  $\gamma_i = 0$ . So in this case  $\check{\gamma}_i = 0$ .
- ② Case  $\hat{\gamma}_i > 0$ . In this case, the objective function  $-\hat{\gamma}_i \gamma_i + \frac{1}{2} \gamma_i^2 + \frac{\lambda}{2} |\gamma_i|$  is minimised somewhere in the range  $\gamma_i \in [0, \infty)$  because the term  $-\hat{\gamma}_i \gamma_i$  is negative there (and all other terms are positive). But when  $\gamma_i \geq 0$ , the objective function becomes

$$-\hat{\gamma}_i\gamma_i+\frac{1}{2}\gamma_i^2+\frac{\lambda}{2}\gamma_i=\left(\frac{\lambda}{2}-\hat{\gamma}_i\right)\gamma_i+\frac{1}{2}\gamma_i^2.$$

If  $\frac{\lambda}{2} - \hat{\gamma}_i \geq 0$ , then the minimum over  $\gamma_i \in [0, \infty)$  is clearly at  $\gamma_i = 0$ . Otherwise, when  $\frac{\lambda}{2} - \hat{\gamma}_i < 0$ , we differentiate and find the minimum at  $\gamma_i = \hat{\gamma}_i - \lambda/2 > 0$ . In summary,  $\tilde{\gamma}_i = (\hat{\gamma}_i - \lambda/2)_+ = sgn(\hat{\gamma}_i)(|\hat{\gamma}_i| - \lambda/2)_+$ .

• Case γ̂<sub>i</sub> < 0. In this case, the objective function  $-\hat{\gamma}_i \gamma_i + \frac{1}{2} \gamma_i^2 + \frac{\lambda}{2} |\gamma_i|$  is minimised somewhere in the range  $\gamma_i \in (-\infty, 0]$  because the term  $-\hat{\gamma}_i \gamma_i$  is negative there (and all other terms are positive). But when  $\gamma_i \leq 0$ , the objective function becomes

$$-\hat{\gamma}_i\gamma_i+\frac{1}{2}\gamma_i^2+\frac{\lambda}{2}(-\gamma_i)=\left(\frac{\lambda}{2}+\hat{\gamma}_i\right)(-\gamma_i)+\frac{1}{2}\gamma_i^2=\left(\frac{\lambda}{2}-|\hat{\gamma}_i|\right)(-\gamma_i)+\frac{1}{2}\gamma_i^2.$$

If  $\frac{\lambda}{2}-|\hat{\gamma}_i|\geq 0$ , then the minimum over  $\gamma_i\in (-\infty,0]$  is clearly at  $\gamma_i=0$ , since  $-\gamma_i$  ranges over  $[0,\infty)$ . Otherwise, when  $\frac{\lambda}{2}-|\hat{\gamma}_i|<0$ , we differentiate and find the minimum at  $\gamma_i=-|\hat{\gamma}_i|+\lambda/2<0$ , which we may re-write as:

$$-|\hat{\gamma}_i| + \lambda/2 = -(|\hat{\gamma}_i| - \lambda/2) = sgn(\hat{\gamma}_i)(|\hat{\gamma}_i| - \lambda/2).$$

In summary,  $\check{\gamma}_i = sgn(\hat{\gamma}_i)(|\hat{\gamma}_i| - \lambda/2)_+$ .

The proof is now complete, as we can see that all three cases yield

$$oldsymbol{\check{\gamma}}_i = sgn(\hat{\gamma}_i) \left( |\hat{\gamma}_i| - rac{\lambda}{2} 
ight)$$
 ,  $i=1,...,p-1$ .

How can we interpret the LASSO in terms of ANOVA in the orthogonal case?

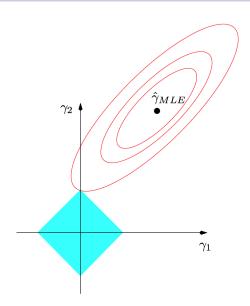
## Corollary

In the context of the previous theorem, and assuming that  $\varepsilon \sim \mathcal{N}(0, \sigma^2 I)$ , model selection using the LASSO tuned by  $\lambda > 0$  is equivalent to including only coefficients significant at level  $\alpha = 2(1 - G_{t_{n-p}}(\lambda/(2S)))$ , where  $G_{t_{n-p}}$  is the CDF of Student's t-distribution.

## Proof.

Remember that a coefficient  $\gamma_j$  is pronounced statistically significant at level  $\alpha$  if  $\{H_0: \gamma_j=0\}$  is rejected at level  $\alpha$ . Under the setting when  $\varepsilon \sim \mathcal{N}(0,\sigma^2I)$ , this happens when  $|\hat{\gamma}_j| > t_{n-p}(1-\alpha/2)S$ . So equating

$$\frac{\lambda}{2} = t_{n-p}(1-\alpha/2)S \implies 1-\frac{\alpha}{2} = G_{t_{n-p}}(\lambda/(2S)) \implies \alpha = 2(1-G_{t_{n-p}}(\lambda/(2S)))$$



Intuition:  $L_1$  norm induces "sharp" balls!

- Balls more concentrated around the axes
- Induces model selection by regulating the lasso (through  $\lambda$ )

Extreme case:  $L^0$  "Norm", gives best subsets selection!

$$\|\gamma\|_0 = \sum_{j=1}^{p-1} |\gamma_j|^0 = \sum_{j=1}^{p-1} \mathbf{1}_{\{\gamma_j \neq 0\}} = \#\{j : \gamma_j \neq 0\}$$

Generally:  $\|\gamma\|_p^p = \sum_{j=1}^{p-1} |\gamma_j|^p$ , sharp balls for 0

$$q = 4$$

q=2

q=1



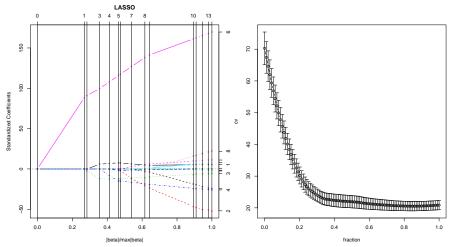
q = 0.5



q = 0.1



## LASSO and CV for different values of $r(\lambda)/\|\hat{\gamma}\|_1$

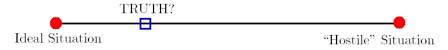


# Robust Linear Modeling

The "success" of the LSE in a regression model depends on "assumptions":

- Normality (LSE optimal in this case)
- Not many "extreme" observations (LSE affected from "extremities")

#### Picture:



- Resistant procedure: not strongly affected by changes to data.
- Robust procedure: not strongly affected by departures from distribution.
  - Often: Robust ⇔ Resistant

Motivating Example: Estimating a Mean

Let  $X_1,\ldots,X_n\stackrel{iid}{\sim} F$ , estimate  $\mu=\int_{-\infty}^{\infty}xF(\,dx)$  by

$$ar{x} = rac{1}{n} \sum_{i=1}^n x_i = rg \min_{\gamma \in \mathbb{R}} \sum_{i=1}^n (x_i - \gamma)^2$$

#### Some observations:

- Average  $\bar{x}$  is optimal (MLE) when F is Normal.
- Extremely sensitive to outliers (low breakdown point).
- Blows up from a single value:  $x \mapsto x + \epsilon \implies \bar{x} \mapsto \bar{x} + \epsilon/n$ .
- If  $\epsilon$  large relative to  $n \to \text{disaster} \dots$
- ullet May not be optimal for other possible F's ...

Can we "cure" sensitivity by using different distance function?

$$m=rg\min_{\gamma\in\mathbb{R}}\sum_{i=1}^n|x_i-\gamma|=\left\{egin{array}{cc}x_{(k+1),}&n=2k+1,\ rac{x_{(k)}+x_{(k+1)}}{2},&n=2k.\end{array}
ight.$$

- Median much less sensitive to bad values.
- Higher breakdown point: must blow up at least 50% of obs to blow m up.
- Median is optimal (MLE) when F is Laplace.
- But how well does m perform when  $F \simeq \text{Normal } (\textit{relative efficiency})?$

## Remember picture:



Other alternatives?

 $\triangleright \alpha$ -Trimmed mean: throw away most extreme observations:

$$trm = rac{1}{|E^c|} \sum_{i 
otin E} x_i,$$

E being subset of  $\alpha \times n$  most extreme observations from each end. Both m and trm may 'throw away' information. View as special cases of the  $\blacktriangleright$  Weighted estimate:

$$wm = \frac{\sum_{i=1}^n w_i x_i}{\sum_{i=1}^n w_i}.$$

- Weights downplaying certain observations (i.e., give less weight to extremes . . . )
- How to objectively/automatically choose weight?

Regression situation is similar. Have:

$$Y = X\beta + \varepsilon$$
,  $\varepsilon \sim F$ ,  $\mathbb{E}[\varepsilon] = 0$ ,  $\operatorname{cov}[\varepsilon] = \sigma^2 I$ 

LSE for  $\beta$  given by

$$\hat{eta} = (X^ op X)^{-1} X^ op y = rg \min_{oldsymbol{\gamma} \in \mathbb{R}^p} \sum_{k=1}^n (y_i - x_i^ op oldsymbol{\gamma})^2$$

- ullet Optimal at F = Normal
- Disastrous if  $y_i \mapsto y_i + c$  with c large:

$$\hat{eta} \mapsto \hat{eta} + (X^T X)^{-1} x_i c$$

- Gauss-Markov: optimal linear for any F
  - $\hookrightarrow$  May not be overall optimal for other F's

### Robust/Resistant Alternatives

- ullet  $L^1$  regression:  $ilde{eta} = rg \min_{oldsymbol{\gamma} \in \mathbb{R}^p} \sum_{k=1}^n |y_i x_i^ op oldsymbol{\gamma}|$
- ullet Trimmed least squares:  $\check{eta} = rg \min_{\gamma \in \mathbb{R}^p} \sum_{i=1}^K (y_i x_i^ op \gamma)_{(i)}^2$ , where we set K = |n/2| + |(p+1)/2|
- Weighted least squares:  $\check{\beta} = (X^\top V^{-1} X)^{-1} X^\top V^{-1} Y$  for a diagonal weight matrix V (recall earlier lecture):

$$V = \left( egin{array}{cccc} w_1 & & & 0 \ & & w_2 & & \ & & \ddots & \ & & & \ddots & \ 0 & & & w_n \end{array} 
ight).$$

Would like to **formalise** the concept of robust/resistant estimation

ightarrow Find a general formulation of which above are special cases.

◆ロト ◆部ト ◆差ト ◆差ト を めらぐ

Seek a unifying approach:

• Instead of  $(\cdot)^2$  or  $|\cdot|$ , consider a more general distance function  $\rho(\cdot)$ .

MLE when errors are Gaussian is obtained as maximising loglikelihood kernel

$$\hat{eta} = rg \max_{oldsymbol{\gamma} \in \mathbb{R}^p} - rac{1}{2} \sum_{i=1}^n \left( rac{y_i - x_i^ op oldsymbol{\gamma}}{\sigma} 
ight)^2$$

Replacing  $\rho(u)=u^2$  by general  $\rho(\cdot)$  yields:

$$\widehat{eta} := rg\min_{oldsymbol{\gamma} \in \mathbb{R}^p} \sum_{i=1}^n 
ho\left(rac{y_i - x_i^ op oldsymbol{\gamma}}{\sigma}
ight)$$

Call this an M(aximum likelihood like)-Estimator.

Obtaining  $\argmin_{\gamma \in \mathbb{R}^p} \sum_{i=1}^n \rho\left( \frac{y_i - x_i^{\top} \gamma}{\sigma} \right)$  reduces to solving

$$\sum_{i=1}^n x_i^ op \psi\left(rac{y_i-x_i^ op \gamma}{\sigma}
ight)=0$$

with  $\psi(t)=d
ho(t)/dt$ . Letting  $w(u)=\psi(u)/u$  this reduces to

$$\sum_{i=1}^n w_i x_i^ op (y_i - x_i^ op \gamma) = 0, \quad ext{where } w_i = w\left(rac{y_i - x_i^ op \gamma}{\sigma}
ight).$$

But this is simply the weighting scenario!

▶ Robust Regression can be written as a Weighted Regression, but the weights depend on the data.

Distance functions are in 1-1 correspondence with loss functions.

Idea: choose  $\rho$  to have desirable properties (reduce/eliminate impact of outliers) — same as choosing weight function.

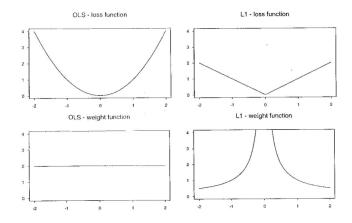
## Some typical examples are:

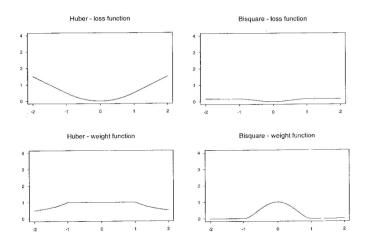
• 
$$\rho(z) = z^2$$
  $\Leftrightarrow w(u) = 2$ 

$$\bullet \ \rho(z) = |z| \qquad \Leftrightarrow w(u) = 1/|u|$$

• Huber: 
$$ho(z) = \left\{ egin{array}{l} z^2, & \mbox{if } |z| \leq H \\ 2H|z| - H^2, & \mbox{otherwise} \end{array} \right.$$

$$\bullet \ \, \mathsf{Bisquare:} \ \, \rho(z) = \left\{ \begin{array}{l} \frac{1}{6}B^2 \left[1 - \left\{1 - \left(z/B\right)^2\right\}^3\right], \quad |z| \leq B, \\ \frac{1}{6}B^2, \ \, \mathsf{otherwise}. \end{array} \right.$$





## Computing a Regression M-Estimator

- Explicit expression for LSE
- ▶ M-Estimation: non-linear optimisation problem use iterative approach
- ▶ Iteratively re-weighted least squares:
  - $oldsymbol{0}$  Obtain initial estimate  $\hat{oldsymbol{eta}}^{(0)}$
  - $\textbf{ 9} \text{ Form normalised residuals } u_i^{(0)} = (y_i x_i^\top \hat{\beta}^{(0)}) / \mathsf{MAD}(y_i x_i^\top \hat{\beta}^{(0)})$
  - $lackbox{0}$  Obtain  $w_i^{(0)} = w(u_i^{(0)})$  for the chosen weight function  $w(\cdot)$
  - lacktriangledown Perform weighted least squares with  $V^{(0)}= ext{diag}\{w_1^{(0)},\ldots,w_n^{(0)}\}$
  - **3** Obtain updated estimate  $\hat{\beta}^{(1)}$
  - Iterate until convergence (?)

## (Asymptotic) Distribution of M-Estimators

▶ Obtained M-Estimator as the solution to the system

$$X^{\top} \psi(\gamma) = 0$$

instead of  $X^{\top}(y - X\gamma) = 0$ . Here we defined

$$\psi(\gamma) = \left(\psi\left(rac{y_1 - x_1^ op \gamma}{\sigma}
ight), \ldots, \psi\left(rac{y_n - x_n^ op \gamma}{\sigma}
ight)
ight)^ op$$

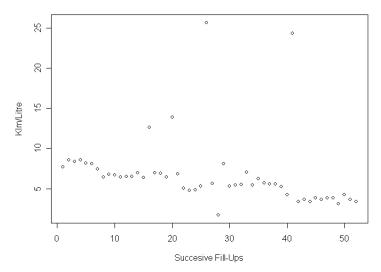
▶If these estimating equations are unbiased, i.e.,

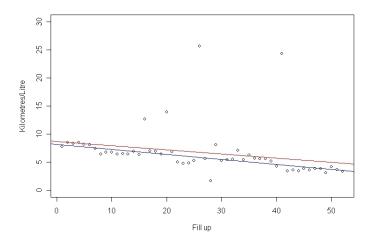
$$\mathbb{E}_{\beta}\left[X^{\top}\Psi(\beta)\right] = 0, \quad \forall \beta \in \mathbb{R}^{p},$$

then under mild regularity conditions, as  $n \to \infty$ , we can show that

$$\hat{\beta}_n \overset{d}{\approx} \mathcal{N}_p \left( \beta, \left\{ \mathbb{E}[X^\top \nabla \psi] \right\}^{-1} X^\top \mathbb{E}[\psi \psi^\top] X \left\{ \mathbb{E}[X^\top \nabla \psi] \right\}^{-1} \right).$$

4□ > 4□ > 4 = > 4 = > = 90



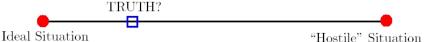


$$\hat{eta}=-0.07$$
 (with  $p=0.06$ ) while  $ilde{eta}=-0.09$  (with  $p\simeq 0$ )

Victor Panaretos (EPFL) Linear Models 242 / 309

## Asymptotic Relative Efficiency (ARE)

## Remember our picture:



- ARE measures quality of one estimator of  $\theta_{p\times 1}$  relative to another, often the MLE  $\hat{\theta}$ , for which  $\text{var}(\hat{\theta}) = I(\theta)^{-1}$ , for large sample size.
- Generally ARE of  $\tilde{\theta}$  relative to  $\hat{\theta}$  is less than 1 (100%): low ARE is bad, high ARE is good.
- ARE of  $\tilde{\theta}$  relative to  $\hat{\theta}$  is

$$\left\{\frac{|\mathsf{var}(\hat{\theta})|}{|\mathsf{var}(\tilde{\theta})|}\right\}^{1/p} \quad (\times 100\%).$$

ullet ARE of  $ilde{ heta}_r$  relative to  $\hat{ heta}_r$  is

$$rac{\mathsf{var}( ilde{ heta}_r)}{\mathsf{var}( ilde{ heta}_r)} \quad ( imes 100\%).$$

#### ARE in the Linear model

- Linear model  $y = X\beta + \varepsilon$ , with  $\varepsilon_j \stackrel{iid}{\sim} g(\cdot)$ ; assume  $\text{var}(\varepsilon_j) = \sigma^2 < \infty$  is known.
- Assume MLE is regular, with

$$i_g = \int -rac{\partial^2 \log g(u)}{\partial u^2} g(u) du = \int \left\{rac{\partial \log g(u)}{\partial u}
ight\}^2 g(u) du.$$

• ARE of LSE of  $\beta$  relative to MLE of  $\beta$  is

$$\frac{1}{\sigma^2 i_g}$$

#### Examples:

- ARE at  $g(\cdot)$  Gaussian: 1
- ARE at  $g(\cdot)$  Laplace: 1/2
- ARE of Huber at  $g(\cdot)$  Gaussian is 95% with H=1.345

Mallow's Rule

A simple and useful strategy is to perform one's analysis both robustly and by standard methods and to compare the results. If the differences are minor, either set may be presented. If the differences are not minor, one must perforce consider why not, and the robust analysis is already at hand to guide the next steps.

- Perform analysis both ways and compare results.
- Plot weights to see which observations were downweighted.
- Try to understand why.

# Nonlinear and Nonparametric Models

Recall most general version of regression given in Week 1:

$$Y_i \mid x_i^{ op} \stackrel{ind}{\sim} \mathsf{Dist}\{g(x_i^{ op})\}, \quad i = 1, \dots, n.$$

So far we have investigated what happens when

$$\left\{egin{array}{ll} g(x^ op) = x^ opeta, \ \operatorname{Dist} = \mathcal{N}(x^ opeta,\sigma^2). \end{array}
ight. egin{array}{ll} eta \in \mathbb{R}^p, \end{array}$$

We now consider a more general situation:

$$Y_i \mid x_i^{ op ind} \overset{ind}{\sim} \mathcal{N}\{\eta(x_i^{ op};eta), \sigma^2\}, \quad i=1,\ldots,n,$$

where  $\eta(x_i^{ op};eta)$ 

- is a KNOWN function,
- ullet that depends on a parameter  $eta \in \mathbb{R}^p$ ,
- but is **not** linear in  $\beta$ .

◆ロ > ◆母 > ◆豆 > ◆豆 > 豆 り Q @

## Example: Logistic Growth

- Decennial population data from US, for 1790–1990.
- y is population in millions, x is time.

## Regression model:

$$Y_i = rac{eta_1}{1 + \exp(eta_2 + eta_3 x_i)} + arepsilon_i, \quad arepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2) \,\, i = 1, \ldots, n.$$

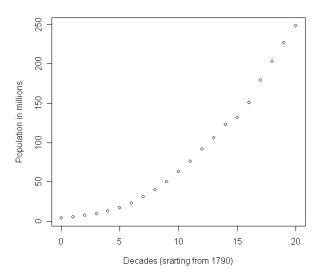
Here

$$\eta(x;eta) = rac{eta_1}{1+\exp(eta_2+eta_3x)}.$$

- Distribution remains Gaussian.
- Cannot transform into a linear regression problem.
- Coefficient interpretation different than in a linear model.
- Related to the differential equation

$$rac{d}{dx}\eta(x) = C imes \eta(x)\{1-\eta(x)\}.$$

◆□ → ◆□ → ◆ = → ◆ = ・ へ Q へ ○



- Still assume independent random variables  $Y_1, \ldots, Y_n$ , with observed values  $y_1, \ldots, y_n$ , and explanatories  $x_1, \ldots, x_n$ .
- Distribution still Gaussian.

#### Introduce notation:

- $ullet y=(y_1,\ldots,y_n)^{ op}\in\mathbb{R}^n$ ,
- $\eta(\beta)=(\eta_1(\beta),\ldots,\eta_n(\beta))^{ op}=(\eta(x_1^{ op},\beta),\ldots,\eta(x_n^{ op},\beta))^{ op}$ , i.e.,

$$\eta(eta): \mathbb{R}^p 
ightarrow \mathbb{R}^n \qquad eta \in \mathbb{R}^p \mapsto \eta(eta) \in \mathbb{R}^n$$

- Therefore  $\eta(\beta)$  is a vector-valued function.
- Analogy with linear case:  $\eta(\beta)$  plays the role of  $X\beta$  but is no longer linear in  $\beta$ .

#### Model now is:

$$y = \underbrace{\eta(eta)}_{n imes 1} + \mathop{arepsilon}_{n imes 1}, \qquad eta \in \mathbb{R}^p, \quad arepsilon \sim \mathcal{N}_n(0, \sigma^2 I).$$

◆ロ → ◆団 → ◆ 豆 → ○ ● ・ ○ へ ○ ○

Since  $\varepsilon \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2)$ , have

$$y \sim \mathcal{N}\{\eta(\beta), \sigma^2\},\$$

so likelihood and loglikelihood are

$$L(eta,\sigma^2) = rac{1}{(2\pi\sigma^2)^{n/2}} \exp\left\{-rac{1}{2\sigma^2}(y-\eta(eta))^{ op}(y-\eta(eta))
ight\},$$

$$\ell(eta,\sigma^2) = -rac{1}{2}\left\{n\log 2\pi + n\log \sigma^2 + rac{1}{\sigma^2}(y-\eta(eta))^{\top}(y-\eta(eta))
ight\}.$$

... exactly as in linear case, but with  $\eta(\beta)$  replacing  $X\beta$ . Hence, suggests *least* squares estimators,

$$\begin{cases} \hat{\beta} = \arg\min_{\beta \in \mathbb{R}^p} \lVert y - \eta(\beta) \rVert^2 & \text{(assuming identifiability),} \\ \hat{\sigma}^2 = \frac{1}{n} \lVert y - \eta(\hat{\beta}) \rVert^2. \end{cases}$$

4□ > 4□ > 4□ > 4□ > 4□ > 4□

Main problem is *non-linearity* — cannot obtain closed form solution in general.

 $\hookrightarrow$  Idea: linearise locally, assuming that  $\eta$  is sufficiently smooth.

First-order Taylor expansion: approximate as

$$\eta(\beta) \simeq \eta(\beta^{(0)}) + \underbrace{\left[\nabla_{\beta}\eta\right]_{\beta=\beta^{(0)}}}_{n\times p} \underbrace{\left(\beta-\beta^{(0)}\right)}_{p\times 1}$$

where  $\beta$  is sufficiently close to  $\beta^{(0)}$ .

• We dropped higher order terms by appealing to smoothness of  $\eta$  (smoothness  $\iff$  "close to zero" higher derivatives).

Linearised representation suggests Newton–Raphson iteration:

- ullet Suppose an initial estimate  $eta^{(0)}$  is available  $(\|eta^{(0)} \hat{eta}\| < \epsilon)$ .
- Let  $D^{(0)} = [\nabla_{\beta} \eta]_{\beta = \beta^{(0)}}$  and  $\beta = u^{(0)} + \beta^{(0)}$ .
- Taylor expansion yields

$$y - \eta(\beta^{(0)}) \approx D^{(0)}\underbrace{(\beta - \beta^{(0)})}_{u^{(0)}} + \varepsilon.$$

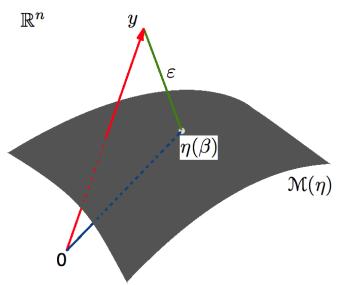
To get  $\beta$  we need  $u^{(0)}$ . Consider the following iteration:

- Initialise with  $\beta^{(0)}$ .
- ② Let  $u^{(1)} = \underset{u \in \mathbb{R}^p}{\arg \min} \|y \eta(eta^{(0)}) D^{(0)}u\|^2$
- :-) (but this is just a linear least squares problem, with  $y^{(0)}=y-\eta(\beta^{(0)})$  and  $X^{(0)}=D^{(0)}!)$
- **3** Thus set  $u^{(1)} = ([D^{(0)}]^{\top} D^{(0)})^{-1} [D^{(0)}]^{\top} \{ y \eta(\beta^{(0)}) \}.$
- Let  $\beta^{(1)} = \beta^{(0)} + u^{(1)}$  and iterate until convergence criterion satisfied. Return last  $\beta^{(k)}$  as  $\hat{\beta}$ .

As  $\beta$  ranges over  $\mathbb{R}^p$ ,  $\eta(\beta)$  traces a p-dimensional differentiable manifold (smooth surface) in  $\mathbb{R}^n$ ,

$$\mathcal{M}(\eta) = \{ \eta(\beta) : \beta \in \mathbb{R}^p \}.$$

- $\bullet$   $\beta$  provides the intrinsic coordinates on that manifold.
- y is obtained by selecting a point  $\eta(\beta)$  on the manifold, and adding a mean zero Gaussian vector  $\varepsilon$ .
- Regression asks to find the coordinates of the point on the manifold that generated y.
- Would like to project y on the manifold, but do not have a closed form expression!



200

255 / 309

Newton-Raphson algorithm is interpretable via differential geometry:

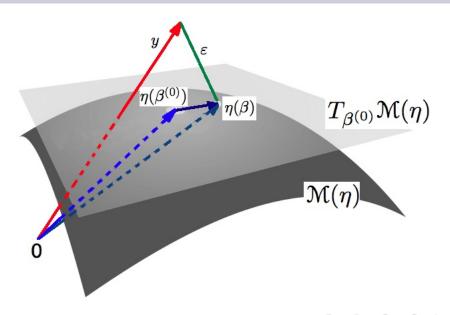
- The *p*-dimensional tangent plane at a point  $\eta(\beta^{(0)}) \in \mathcal{M}(\eta)$  is spanned by  $\eta(\beta^{(0)}) + [\nabla_{\beta}\eta(\beta)]_{\beta=\beta^{(0)}} u$ ,  $u \in \mathbb{R}^p$ .
- Hence we may write that

$$T_{oldsymbol{eta}^{(0)}} \mathcal{M}(\eta) = \{ \eta(oldsymbol{eta}^{(0)}) + D^{(0)}u : u \in \mathbb{R}^p \}$$

- In other words, the p columns of  $D^{(0)}$ , when translated by  $\eta(\beta^{(0)})$ , form a basis for the tangent plane at  $\eta(\beta^{(0)})$ .
- Taylor expansion merely says that if  $\beta$  is close to  $\beta^{(0)}$ , we approximately have  $\eta(\beta) \eta(\beta^{(0)}) \in T_{\beta^{(0)}} \mathcal{M}(\eta)$ . This is equivalent to the expression

$$\eta(eta) - \eta(eta^{(0)}) pprox \underbrace{\left[
abla_{eta}\eta\right]_{eta=eta^{(0)}}}_{D^{(0)}}\underbrace{\left(eta-eta^{(0)}
ight)}_{u^{(0)}}.$$

- Therefore,  $y \eta(\beta^{(0)}) \approx D^{(0)} u^{(0)} + \varepsilon$  means that  $\mathbb{E}[y]$  approximately lies in  $T_{\beta^{(0)}} \mathcal{M}(\eta)$ .
- Newton-Raphson algorithm ≡ iterated projection on approximating linear subspaces.



- Summarising, suppose we consider  $\eta(\beta^{(0)})$  as the origin of space (i.e., now the tangent space is a subspace).
- Then  $y \eta(\beta^{(0)})$  is approximately the response obtained when adding  $\varepsilon$  to an element  $D^{(0)}(\beta \beta^{(0)}) \in T_{\beta_{(0)}} \mathcal{M}(\eta)$ .
- So, approximately, we have our usual linear problem, and we can use orthogonal projection to solve it.
- Amounts to approximating the manifold  $\mathcal{M}(\eta)$  by a plane  $T_{\beta_{(0)}}\mathcal{M}(\eta)$  locally around  $\eta(\beta^{(0)})$ .

Once initial value  $\beta^{(0)}$  is updated to  $\beta^{(1)}$ , use a new tangent plane approximation and repeat the whole procedure.

But how do we obtain our initial  $\beta^{(0)}$ ?

# Choosing $\beta^{(0)}$

Successful linearisation depends on good initial value.

- Occasionally, can find initial values by inspection in simple problems.
- More generally, it takes some experimentation.
  - E.g., one can try fitting polynomial models to data.
  - Use these to find fitted values at fixed design points.
  - Solve a system of equations to get initial values.

Example: consider the model  $y_j = \beta_0 + \beta_1 \exp\{(-x_j/\theta)\} + \varepsilon_j$ 

- Fit a polynomial regression to data
- ② Find fitted values  $\tilde{y}_0$ ,  $\tilde{y}_1$ ,  $\tilde{y}_2$  at  $x_0$ ,  $x_0 + \delta$ ,  $x_0 + 2\delta$ .
- Equate fitted values with model expectation:

$$\tilde{y}_k = \beta_0 + \beta_1 \exp\{-(x_0 + k\delta)/\theta\}, \quad k = 0, 1, 2.$$

- System yields initial estimate  $heta^{(0)} = \delta/\log\left[( ilde{y}_0 ilde{y}_1)/( ilde{y}_1 ilde{y}_2)
  ight]$
- **3** Get initial values for  $\beta_0, \beta_1$  by linear regression, once  $\theta^{(0)}$  is at hand.

Under smoothness conditions on  $\eta$ , one can in general prove that

$$oxed{S^{-1} \left\{ 
abla_eta \eta(\hat{eta})^ op 
abla_eta \eta(\hat{eta}) 
ight\}^{1/2} (\hat{eta} - eta) \stackrel{d}{pprox} N_p(0, I_p)}$$

for large n, where  $S = (n - p)^{-1} ||e||^2$ . May thus mimic linear case:

$$c^{ op}\hat{eta} \overset{d}{pprox} \mathcal{N}_1 \left[ c^{ op}eta, S^2c^{ op} \left\{ 
abla_eta\eta(\hat{eta})^{ op}
abla_eta\eta(\hat{eta}) 
ight\}^{-1} c 
ight].$$

So base confidence intervals (and tests) on

$$\frac{c^{\top}\hat{\beta} - c^{\top}\beta}{\sqrt{S^2c^{\top}\left\{\nabla_{\beta}\eta(\hat{\beta})^{\top}\nabla_{\beta}\eta(\hat{\beta})\right\}^{-1}c}} \overset{d}{\approx} N(0,1),$$

which gives a  $(1 - \alpha) \times 100\%$  CI:

$$oxed{c^{ op}\hat{eta}\pm z_{lpha/2}\sqrt{S^2c^{ op}\left\{
abla_{eta}\eta(\hat{eta})^{ op}
abla_{eta}\eta(\hat{eta})
ight\}^{-1}c}}$$

Until today we have discussed the following setup:

$$Y_i \mid x_i \overset{ind}{\sim} \mathsf{Dist}[y \mid heta_i] 
ightarrow \left\{egin{array}{l} heta_i = g(x_i; eta), \ eta \in \mathbb{B} \subset \mathbb{R}^p, \end{array}
ight.$$

with  $g(\cdot; \beta)$  known up to  $\beta$  to be estimated from data, e.g.

- Dist $(\cdot \mid \mu) = \mathcal{N}(\cdot \mid \mu)$  and  $\mu = g(x \mid \beta) = x^{\top}\beta$ ,
- Dist $(\cdot \mid \mu) = \mathcal{N}(\cdot \mid \mu)$  and  $\mu = g(x \mid \beta) = \eta(x; \beta)$ .

Would now like to extend model to a more flexible dependence:

$$Y_i \mid x_i \overset{ind}{\sim} \mathsf{Dist}[y \mid heta_i] 
ightarrow \left\{egin{array}{l} heta_i = g(x_i), \ g \in \mathfrak{F} \subset L^2(\mathbb{R}^p) \ ext{(say)}, \end{array}
ight.$$

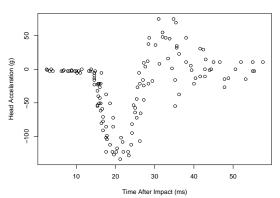
with g unknown, to be estimated given data  $\{(y_i, x_i)\}_{i=1}^n$ .

- A *nonparametric* problem (parameter ∞-dimensional)!
- How to estimate g in this context?
- $\bullet$   $\mathfrak F$  is usually assumed to be a class of smooth functions (e.g.,  $C^k$ ).

### Start from simplest problem:

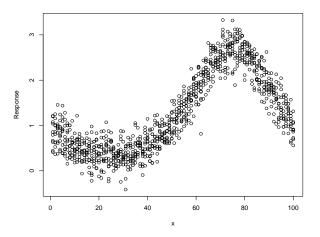
$$egin{aligned} \mathsf{Dist} &\equiv \mathcal{N}(\mu,\sigma^2) \ x_i \in \mathbb{R} \end{aligned} 
ight. egin{aligned} & \longrightarrow & Y_i = g(x_i) + arepsilon_i, & arepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0,\sigma^2) \end{aligned}$$

## Figure: Motorcycle Accident Data



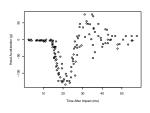
### **Exploiting Smoothness**

ullet Ideally: multiple y's at each  $x_i$  ( $n o \infty$  and large *covariate classes*):



- ullet Then average y's at each  $x_i$  and interpolate . . .
- But this is never the case ...

• Usually unique  $x_i$  distinct:



- Here is where the smoothness assumption comes in
- ullet Since have unique y at each  $x_i$ , need to borrow information from nearby  $\dots$
- ... use continuity!!! (or even better, smoothness)
- ▶ Recall: A function  $g : \mathbb{R} \to \mathbb{R}$  is *continuous* if:

$$orall \; \epsilon > 0 \; \exists \; \delta > 0: \; |x-x_0| < \delta \implies |g(x)-g(x_0)| < \epsilon.$$

- ▶ So maybe average  $y_i$ 's corresponding to  $x_i$ 's in a  $\delta$ -neighbourhood of x as  $\hat{g}(x)$ ?
- ▶ Motivates the use of a kernel smoother . . .

◆ロト ◆団ト ◆豆ト ◆豆 ・ りへで

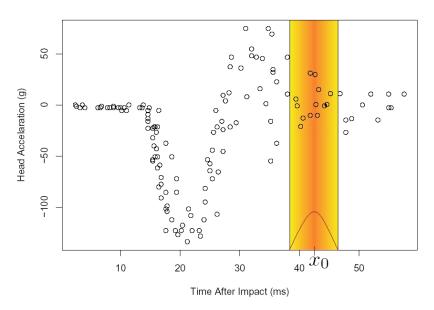
Naive idea:  $\hat{g}(x_0)$  should be the average of  $y_i$ -values with  $x_i$ 's "close" to  $x_0$ .

$$\hat{g}(x_0) = rac{1}{\sum_{i=1}^n \mathbf{1}\{|x_i - x_0| \leq \lambda\}} \sum_{i=1}^n y_i \mathbf{1}\{|x_i - x_0| \leq \lambda\}.$$

A weighted average! Choose other weights? Kernel estimator:

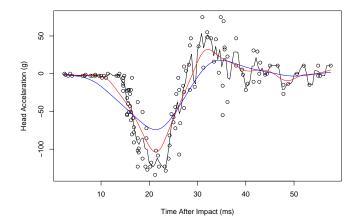
$$\hat{g}(x_0) = rac{1}{\sum_{i=1}^n K\left(rac{x_i-x_0}{\lambda}
ight)} \sum_{i=1}^n y_i K\left(rac{x_i-x_0}{\lambda}
ight).$$

- K is a weight function (kernel), e.g. a pdf
  - $\,\hookrightarrow\,$  Usually symmetric, non-negative, decreasing away from zero
- ullet  $\lambda$  is the bandwidth parameter
  - $\hookrightarrow$  small  $\lambda$  gives local behaviour, large  $\lambda$  gives global behaviour
- ullet Choice of K not so important, choice of  $\lambda$  very important!
- The resulting fitted values are linear in the responses, i.e.,  $\hat{y} = S_{\lambda} y$ , where the smoothing matrix  $S_{\lambda}$  depends on  $x_1, \ldots, x_n$ , K and  $\lambda$ . Analogous to a projection matrix in linear regression, but  $S_{\lambda}$  is NOT a projection.



### Motorcycle Data Kernel Smooth

- > plot(time,accel,xlab="Time After Impact (ms)",ylab="Head Accelaration (g)")
- > lines(ksmooth(time,accel,kernel="normal",bandwidth=0.7))
- > lines(ksmooth(time,accel,kernel="normal",bandwidth=5),col="red")
- > lines(ksmooth(time,accel,kernel="normal",bandwidth=10),col="blue")



# Find $g \in C^2$ that minimises

$$\underbrace{\sum_{i=1}^{n} \{y_i - g(x_i)\}^2}_{\text{Fit Penalty}} \ + \ \underbrace{\lambda \int_{I} \{g''(t)\}^2 dt}_{\text{Roughness Penalty}}$$

- $\lambda$  to balance fidelity to the data and smoothness of the estimated h.

# Remarkably, problem has unique explicit solution!

- $\hookrightarrow$  Natural Cubic Spline with knots at  $\{x_i\}_{i=1}^n$ :
  - piecewise polynomials of degree 3,
  - with pieces defined at the knots,
  - with two continuous derivatives at the knots,
  - and linear outside the data boundary.

Can represent splines via a basis  $B_j$ , as

$$s(t) = \sum_{j=1}^{n} \gamma_j B_j(t).$$

For example, one basis (the natural basis) is

$$egin{array}{lcl} B_1(t) &=& 1 \ B_2(t) &=& t \ B_{m+2}(t) &=& \delta_m(t) - \delta_{n-1}(t), & m=1,\ldots,n-2 \ \delta_k(t) &=& rac{(t-x_k)_+^3 - (t-x_n)_+^3}{t_n-t_k}, & k=1,\ldots,n-1 \end{array}$$

where  $x_m$  are the knot locations and

$$(\cdot)_+ = \max\{\cdot, 0\}$$

is the positive part of any function.

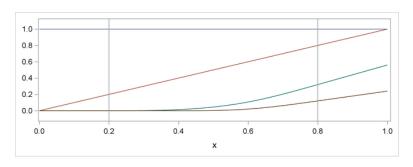


Figure: The n=4 natural spline basis functions for knots at  $x_1=0.2,\ x_2=0.4,\ x_3=0.6$  and  $x_4=0.8$ 

We wish to find a basis for natural cubic splines with knot locations  $\{x_i\}_{i=1}^n$ 

• Observe that any piecewise polynomial  $PP_3(t)$  of order 3 with 2 cts derivatives at the knots can be expanded in the truncated power series basis

$$PP_3(t) = \sum_{j=0}^3 \phi_j \, t^j \, + \sum_{i=1}^n heta_k (t-x_i)_+^3$$

- The n+4 coefficients  $\{\phi_j\}_{j=0}^3 \cup \{\theta_i\}_{i=1}^n$  must satisfy constraints to ensure linearity beyond boundary knots:
  - $\phi_2 = 0 \& \phi_3 = 0$
  - $\bullet \ \sum_{i=1}^n \theta_i = 0$
  - $\bullet \ \sum_{i=1}^n \theta_i x_i = 0$
- Can then use relations re-express basis in form on previous slide, with only n (rather than n+4) basis functions, and unconstrained coefficients.

Letting  $\gamma = (\gamma_1, \ldots, \gamma_n)^\top$ ,

$$g(t) = \sum_{i=1}^n \gamma_i B_i(t), \quad B = \{B_{ij}\} = \{B_j(x_i)\}, \quad \Omega_{ij} = \int B_i''(t) B_j''(t) \, dt,$$

our penalised likelihood

$$\sum_{i=1}^{n} \{y_i - g(x_i)\}^2 + \lambda \int_{I} \{h''(t)\}^2 dt$$

becomes

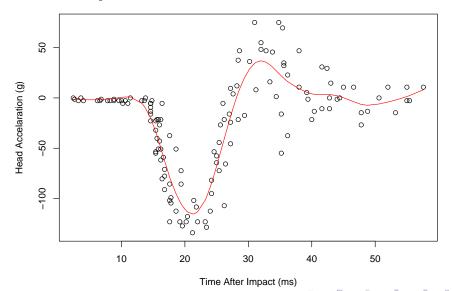
$$\left\{ (y-B\gamma)^{ op}(y-B\gamma) + \lambda \gamma^{ op}\Omega\gamma 
ight\}$$
 .

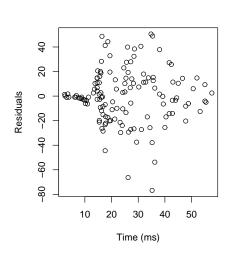
Differentiating and equating with zero yields

$$(B^{\top}B + \lambda\Omega)\hat{\gamma} = B^{\top}y \implies \hat{\gamma} = (B^{\top}B + \lambda\Omega)^{-1}B^{\top}y.$$

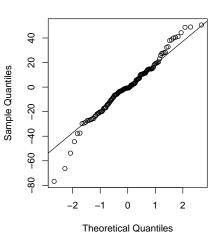
- The smoothing matrix is  $S_{\lambda} = B(B^{\top}B + \lambda\Omega)^{-1}B^{\top}$ .
- The natural cubic spline fit is approximately a kernel smoother.

lines(smooth.spline(time,accel),col="red")





### Normal Q-Q Plot



4ロト 4回 ト 4 重 ト 4 重 ト 9 Q (~)

• Least squares estimation:  $y = X_{n \times p} \beta + \varepsilon$ , we have  $\hat{y} = Hy$ , with trace(H) = p, in terms of the projection matrix  $H = X(X^{\top}X)^{-1}X^{\top}$ . Here

$$\hat{y} = \underbrace{B(B^{\top}B + \lambda\Omega)^{-1}B^{\top}}_{S_{\lambda}} y.$$

• Idea: define equivalent degrees of freedom of smoother

$$\operatorname{trace}(S_{\lambda}) = \sum_{j=1}^{n} \frac{1}{1 + \lambda \eta_{j}}$$

where  $\eta_i$  are eigenvalues of  $K = (B^\top B)^{-1/2} \Omega(B^\top B)^{-1/2}$ .

- Hence  $\operatorname{trace}(S_{\lambda})$  is monotone decreasing in  $\lambda$ , with  $\operatorname{trace}(S_{\lambda}) \to 2$  as  $\lambda \to \infty$  (K will have twos zero eigenvalues) and  $\operatorname{trace}(S_{\lambda}) \to n$  as  $\lambda \to 0$ . Note 1–1 map  $\lambda \leftrightarrow \operatorname{trace}(S_{\lambda}) = \operatorname{df}$ , so usually determine roughness using df (interpretation easier).
- Each eigenvalue of  $S_{\lambda}$  lies in (0,1), so this is a smoothing, NOT a projection, matrix.

Focus on the fit for the given grid  $x_1, \ldots, x_n$ :

$$\hat{\mathbf{g}}=(\hat{g}(x_1),\ldots,\hat{g}(x_n)),\quad \mathbf{g}=(g(x_1),\ldots,g(x_n))$$

Consider the mean squared error:

$$\mathbb{E}(\|\mathbf{g} - \hat{\mathbf{g}}\|^2) = \underbrace{\mathbb{E}\{\|\mathbb{E}(\hat{\mathbf{g}}) - \hat{\mathbf{g}}\|^2\}}_{\text{variance}} + \underbrace{\|\mathbf{g} - \mathbb{E}(\hat{\mathbf{g}})\|^2}_{\text{bias}^2}.$$

When estimator potentially biased, need to worry about both! In the case of a linear smoother, for which  $\hat{\mathbf{g}}=S_\lambda y$ , we find that

$$\mathbb{E}(\|\mathbf{g} - \hat{\mathbf{g}}\|^2) = \frac{\mathsf{trace}(S_{\lambda}S_{\lambda}^{\top})}{n}\sigma^2 + \frac{(\mathbf{g} - S_{\lambda}\mathbf{g})^{\top}(\mathbf{g} - S_{\lambda}\mathbf{g})}{n},$$

SO

- $\lambda \uparrow \implies$  variance  $\downarrow$  but bias  $\uparrow$ ,
- $\lambda \downarrow \implies$  bias  $\downarrow$  but variance  $\uparrow$ .
- Would like to choose  $\lambda$  to find optimal bias-variance tradeoff:

 $\hookrightarrow$  Unfortunately, optimal  $\lambda$  will generally depend on unknown g!

### Choosing $\lambda$

- Fitted values are  $\hat{y} = S_{\lambda} y$ .
- ullet Fitted value  $\hat{y}_{i}^{-}$  obtained when  $y_{j}$  is dropped from fit is

$$S_{jj}(\lambda)(y_j - \hat{y}_j^-) = \hat{y}_j - \hat{y}_j^-.$$

Cross-validation sum of squares is

$$\mathsf{CV}(\lambda) = \sum_{j=1}^n (y_j - \hat{y}_j^-)^2 = \sum_{j=1}^n \left\{ \frac{y_j - \hat{y}_j}{1 - S_{jj}(\lambda)} \right\}^2,$$

and generalised cross-validation sum of squares is

$$\mathsf{GCV}(\lambda) = \sum_{j=1}^n \left\{ rac{y_j - \hat{y}_j}{1 - \mathsf{trace}(S_\lambda)/n} 
ight\}^2,$$

where  $S_{jj}(\lambda)$  is (j,j) element of  $S_{\lambda}$ .

Orthogonal Series: "Parametrising" The Problem

Depending on what  $\mathcal{F} \ni g(\cdot)$  is (Hilbert space) can write:

$$g(x) = \sum_{k=1}^{\infty} eta_k \psi_k(x)$$
 (in an appropriate sense),

with  $\{\psi\}_{k=1}^\infty$  known (orthogonal) basis functions for  $\mathcal{F}$ , e.g.,

- $\bullet \ \mathcal{F} = L^2(-\pi,\pi),$
- $ullet \left\{ \psi_k 
  ight\} = \{e^{-ikx}\}_{k \in \mathbb{Z}}, \; \psi_i \perp \psi_j, \; i 
  eq j.$
- Gives Fourier series expansion,  $\beta_k = \frac{1}{2\pi} \int_{-\pi}^{\pi} g(x) e^{-ikx} dx$ .

Idea: if truncate series, then have simple linear regression!

$$Y_i = \sum_{k=1}^{ au} eta_k \psi_k(x_i) + arepsilon_i, \quad au < \infty$$

Notice: truncation has implications, e.g., in Fourier case:

- Truncating implies assume  $g \in \mathcal{G} \subset L^2$ .
- ullet Interpret this as a smoothness assumption on g.
- How to choose  $\tau$  optimally?

Easy exercise in Fourier analysis:

$$\sum_{k=- au}^{ au} eta_k \, e^{-ikx} = rac{1}{2\pi} \int_{-\pi}^{\pi} g(y) D_{ au}(x-y) \, dy$$

with the *Dirichlet kernel* of order  $\tau$ ,  $D_{\tau}(u) = \sin \{(\tau + 1/2) u\}/\sin(u/2)$ . Recall kernel smoother:

$$\hat{g}(x_0) = \sum_{i=1}^n rac{y_i K_\lambda(x_i - x_0)}{\sum_{i=1}^n K_\lambda(x_i - x_0)} = rac{1}{c} \int_I y(x) K_\lambda(x - x_0) dx,$$

with

$$y(x) = \sum_{i=1}^n y_i \delta(x - x_i).$$

- ullet So if K is the Dirichlet kernel, we can do series approximation via kernel smoothing.
- Works for other series expansions with other kernels (e.g., Fourier with convergence factors)

4ロ → 4回 → 4 呈 → 4 呈 → 9 へ ○

Victor Panaretos (EPFL) Linear Models 279 / 309

The Dirichlet kernel

From  $x \in \mathbb{R}$  to  $(x_1, \ldots, x_d) \in \mathbb{R}^p$ 

So far: how to estimate  $g:\mathbb{R} \to \mathbb{R}$  (assumed smooth) in

$$Y_i = g(x_i) + \varepsilon_i, \quad \varepsilon_i \stackrel{iid}{\sim} \mathcal{N}(0, \sigma^2), \text{ given data } \{(y_i, x_i)\}_{i=1}^n.$$

- ► Generalise to include multivariate explanatories?
- ▶ "Immediate" Generalisation:  $g: \mathbb{R}^p \to \mathbb{R}$  (smooth)

$$Y_j = g(x_{j1}, \ldots, x_{jp}) + arepsilon_j, \quad arepsilon_j \overset{iid}{\sim} \mathcal{N}(0, \sigma^2)$$

- ▶ Estimation by (e.g.) multivariate kernel method.
- ▶ Two basic drawbacks of this approach . . .
- Shape of kernel? (definition of *local*)

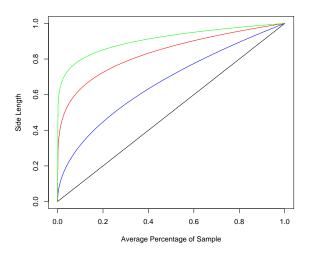
#### What is "local" in $\mathbb{R}^p$ ?

- → Need some definition of "local" in the space of explanatories
- $\hookrightarrow$  Use some metric on  $\mathbb{R}^p \ni (x_1, \dots, x_p)$ !

#### But which one?

- Choice of metric ←⇒ choice of geometry
  - $\hookrightarrow$  e.g., curvature reflects intertwining of dimensions
- Geometry  $\implies$  reflects structure in the explanatories
  - potentially different units of measurement (variable stretching of space)
  - g may be of higher variation in some dimensions (need finer neighbourhoods there)
  - statistical dependencies present in the explanatories
     ("local" should reflect these)

# Curse of Dimensionality $(\mathcal{U}[0,1]^p)$



$$p = 1$$
,  $p = 2$ ,  $p = 5$ ,  $p = 10$ 

### Curse of Dimensionality

"neighbourhoods with a fixed number of points become less local as the dimensions increase"

Bellman (1961)

- Hence to allow for reasonably small bandwidths
   → Density of sampling must increase.
- → Density of Sampling must increase.
- Need to have ever larger samples as dimension grows.

Attempt to find a link/compromise between:

- our mastery of 1D case (at least we can do that well ...),
- and higher dimensional explanatories (and associated difficulties).

One approach: Projection-Pursuit Regression

$$Y = \sum_{k=1}^K h_k(artheta_k^ op \mathbf{x}) + arepsilon, \quad \|artheta_k\| = 1, \,\, arepsilon \sim \mathcal{N}(0,\sigma^2).$$

- ullet Additively decomposes g into smooth functions  $h_k:\mathbb{R} \to \mathbb{R}.$
- Each function depends on a global feature

   → a linear combination of the explanatories,
- projections directions chosen for best fit

   ⇒ similarities to tomography.
- ullet Each  $h_k$  is a ridge function of  ${f x}$ : varies only in the direction defined by  $artheta_k$

4 D > 4 A > 4 B > 4 B > B = 900

285 / 309

Victor Panaretos (EPFL) Linear Models

How is the model fitted to data?

Assume only one term, K=1 and consider penalized likelihood:

$$\min_{h \in \mathcal{C}^2, \|\boldsymbol{\vartheta}\| = 1} \qquad \left\{ \sum_{i=1}^n \{y_i - h_1([\boldsymbol{\vartheta}^\top \mathbf{x}]_i)\}^2 + \int_I \{h_1''(t)\}^2 dt \right\}.$$

# Two steps:

- *Smooth*: Given a direction  $\vartheta$ , fitting  $g_1(\vartheta^\top \mathbf{x})$  is done via 1D smoothing splines.
- Pursue: Given  $h_1$ , have a non-linear regression problem w.r.t.  $\vartheta$ .

Hence, iterate between the two steps

- $\hookrightarrow$  Complication is that  $h_1$  not explicitly known, so need numerical derivatives.
- → Further terms added in forward stepwise manner.

←□ → ←□ → ← = → ○ へ ○ ○

### Projection pursuit:

- (+) Can uniformly approximate  $C^1(\text{compact}[\mathbb{R}^p])$  function arbitrarily well as  $K \to \infty$  (very useful for prediction)
- (-) Interpretability? What do terms mean within problem?

# Need something that can be interpreted variable-by-variable

► Compromise: Additive Model

$$Y_j = lpha_j + \sum_{k=1}^p f_k(x_{jk}) + arepsilon_j, \quad arepsilon_j \overset{iid}{\sim} \mathcal{N}(0, \sigma^2),$$

•  $f_j$ 's univariate smooth functions,  $\sum_j f_k(x_{jk}) = 0$ . In our standard setting, have:

$$Y_j \mid \widetilde{x}_j \overset{ind}{\sim} \mathsf{Dist}(\cdot \mid heta_j) 
ightarrow \left\{egin{array}{l} \mathsf{Dist} = \mathcal{N}(\mu_j, \sigma^2), \ heta_j = \mu_j = lpha_j + \sum_{k=1}^p f_k(x_{jk}). \end{array}
ight.$$

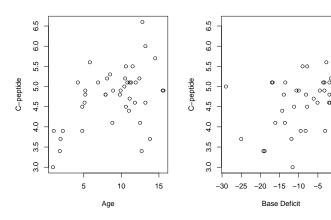
## The Backfitting Algorithm

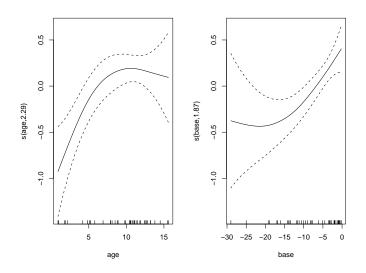
- ▶ How to fit additive model?
- $\hookrightarrow$  Know how to fit each  $f_k$  separately quite well
- $\hookrightarrow$  Take advantage of this . . .
- ▶ Motivation: Fix *j* and drop it for ease:

$$\mathbb{E}\left[Y-lpha-\sum_{m
eq k}f_m(x_m)
ight]=f_k(x_k)$$

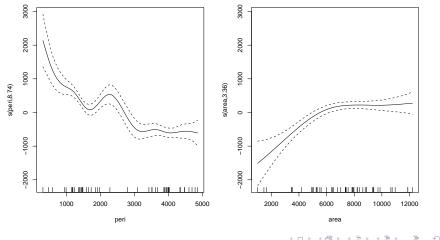
- ► Suggests the *Backfitting Algorithm*:
- (1) Initialise:  $lpha={\sf ave}\{y_j\}$ ,  $f_k=f_k^{\,0}$ ,  $k=1,\ldots,p$ .
- (2) Cycle:  $f_k = \mathcal{S}_k(y-\alpha-\sum_{m \neq k} \mathbf{f}_m)$   $k=1,\ldots,p,1,\ldots,p,\ldots$
- (3) Stop: when individual functions don't change
- $\triangleright$  S is arbitrary scatterplot smoother

◆ロ > ◆母 > ◆き > ◆き > き の Q @





Measurements on 48 rock samples from a petroleum reservoir: rock.gam<-gam(perm 1+s(peri)+s(area),family=gaussian)</pre>



```
Family: gaussian
Link function: identity
Formula:
perm ~1 + s(peri) + s(area)
Parametric coefficients:
           Estimate Std. Error t value Pr(>|t|)
(Intercept) 415.45 27.18 15.29 <2e-16 ***
___
Signif. codes: 0 '*** 0.001 '** 0.01 '* 0.05 '.' 0.1 ' 1
Approximate significance of smooth terms:
         edf Est.rank F p-value
s(peri) 8.739 9 18.286 9.49e-11 ***
s(area) 3.357 7 6.364 7.41e-05 ***
Signif. codes: 0 '*** 0.001 '** 0.01 '* 0.05 '.' 0.1 ' ' 1
R-sq.(adj) = 0.815 Deviance explained = 86.3\%
```

More on Splines

#### We want to rigorously show:

- The penalized least squares problem admints a natural cubic spline as a unique solution
- ② That any natural cubic spline on n distinct knots can be expanded in a basis of n elements  $\{B_1, ..., B_n\}$
- **1** That the matrix inversion involved in the expression  $(B^\top B + \lambda \Omega)^{-1} B^\top y$  is well-defined

### En route, we would also like to

• Construct at least one example of an explicit basis  $\{B_1,...,B_n\}$ .

To analyse spline smoothing we will need to first analyse spline interpolation.

Victor Panaretos (EPFL)

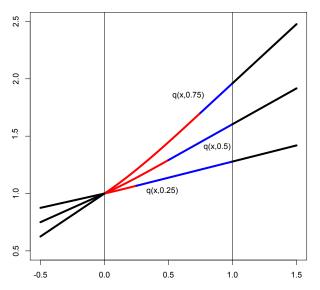
Our analysis will hinge on a very carefully chosen kernel:

$$q(x, y) = 1 + xy + k(x, y), \quad (x, y) \in [0, 1]^2,$$

where

$$k(x,y) = egin{cases} x^2y/2 - x^3/6 & ext{ for } x \leq y \ xy^2/2 - y^3/6 & ext{ for } x > y \end{cases}, \quad (x,y) \in [0,1]^2.$$

- We will write  $q_y(x)$  or  $k_y(x)$  whenever we want to emphasise that the second argument is taken fixed and we view the kernel as a function of the first argument.
- ullet In this light,  $q_y(x)$  is piecewise polynomial with two pieces:
  - **1** a cubic piece (for  $0 \le x \le y$ ), and
  - 2 a linear piece (for  $y \leq x \leq 1$ ).



Recall,  $q_y(x)$  is piecewise polynomial with two pieces:

- **1** a cubic piece (for  $0 \le x \le y$ )
- 2 a linear piece (for  $y \leq x \leq 1$ ).

# Theorem (Positive Definiteness)

Given any  $1 \le t_1 \le t_2 \le \ldots \le t_n \le 1$  we have

$$\sum_{i=1}^n \sum_{j=1}^n lpha_i lpha_j q(t_i,t_j) \geq 0 \qquad orall \, lpha = (lpha_1,...,lpha_n)^ op \in \mathbb{R}^n,$$

in other words  $Q = \{q(t_i, t_j)\}_{i,j=1}^n$  is nonnegative definite. When all the  $t_j$ 's are distinct,

$$0 \le t_1 < t_2 < \ldots < t_n \le 1,$$

the displayed inequality is strict unless  $lpha\in\mathbb{R}^n\setminus\{0\}$ , and so Q is positive definite.

## Proof.

Let  $K = \{k(t_i, t_j)\}_{i=1}^n$ ,  $t = (t_1, \dots, t_n)^{\top}$ ,  $1 = (1, \dots, 1) \in \mathbb{R}^n$  and note that

$$Q = \{q(t_i, t_j)\}_{i,j=1}^n = \{1 + t_i t_j + k(t_i, t_j)\}_{i,j=1}^n = \mathbf{1}\mathbf{1}^\top + t t^\top + K.$$

Thus, if we can verify that  $K \succeq 0$  we will obtain that  $Q \succeq 0$ , being the sum of three non-negative definite matrices.

297 / 309

Given any pair  $(t_i, t_j)$  with  $t_i \leq t_j$  (say), observe that

$$\int_0^1 k_{t_i}''(u)k_{t_j}''(u)\,du = \int_0^{t_i} (t_j-u)(t_i-u)\,du = t_i^2t_j/2 - t_i^3/6 = k(t_i,t_j). \quad (*)$$

Therefore, we may substitute the integral expression for  $k(t_i, t_j)$  into  $\alpha^{\top} Q \alpha$  to manifest a square:

$$\sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j \, k(t_i,t_j) = \sum_{i=1}^n \sum_{j=1}^n \int_0^1 \alpha_i k_{t_i}''(u) \alpha_j \, k_{t_j}''(u) \, du = \int_0^1 \left( \sum_{i=1}^n \alpha_i k_{t_i}''(u) \right)^2 \, du.$$

This shows that  $lpha^ op Qlpha \geq$  0, and so K (and hence Q) is always nonnegative.

Now suppose that the  $\{t_i\}$  are all distinct. Remark that each function  $k''_{t_i}(u)$  is supported on  $[0,t_i)$  and is linear thereon. We distinguish two cases:

•  $t_1 > 0$ . Then all n supports are disjoint non-empty intervals and the  $\{k_{t_j}^{\prime\prime}\}_{j=1}^n$  are linearly independent. Consequently the sum can be zero only if  $\alpha_1 = \alpha_2 = \ldots = \alpha_n = 0$ , and K (and hence Q) is strictly positive.

Victor Panaretos (EPFL)

Linear Models

•  $t_1=0$ . Then  $k_{t_1}''=0$ , so only the n-1 functions  $\{k_{t_j}''\}_{j=2}^n$  are linearly independent. In this case, first row/column of K will be uniformly zero, and only the bottom right  $(n-1)\times (n-1)$  submatrix

$$\boldsymbol{K}_{n-1} = \{k(t_i, t_j)\}_{j=2}^n$$

of K will be positive definite. Thus K is of reduced rank n-1. However, the first column of  $\mathbf{11}^{\top}$  is now linearly independent of all columns of K, and so  $Q = \mathbf{11}^{\top} + tt^{\top} + K$  is of full rank n.

In summary, when  $0 \leq t_1 < t_2 < \ldots < t_n \leq 1$ , the matrix Q is positive definite.

Notice that the calculation (\*) was the crucial ingredient. We will use this again when proving that  $\Omega$  is nonnegative.

Why go into all this trouble? It turns out that this property will give us both:

- A solution to the spline interpolation problem.
- A basis for natural cubic splines.

## Theorem (Spline Interpolation: Uniqueness and Optimality))

- . Let  $0 = t_1 < t_2 < \ldots < t_n = 1$  be distinct nodes, with  $n \ge 2$ , and  $y = (y_1, \ldots, y_n) \in \mathbb{R}^n$  be associated responses.
  - **1** There exists a unique natural cubic spline  $s:[0,1]\to\mathbb{R}$  with knots at  $\{t_j\}$  that interpolates  $\{(t_i,y_j)\}_{j=1}^n$ , and can be explicitly constructed as

$$s(x) = \sum_{j=1}^n heta_i \, q(x,t_j), \qquad ext{with } oldsymbol{ heta} = oldsymbol{Q}^{-1} oldsymbol{y}$$

where  $\boldsymbol{\theta} = (\theta_1, \dots, \theta_n)^{\top}$ , and  $\boldsymbol{Q} = \{q(t_i, t_j)\}_{i,j=1}^n$  is bone fide invertible.

**②** for any  $C^2$  function  $f:[0,1] \to \mathbb{R}$  that also interpolates  $\{(t_j,y_j)\}_{j=1}^n$ ,

$$\mathcal{C}(f) \equiv \int_0^1 [f''(u)]^2 du \geq \int_0^1 [s''(u)]^2 du \equiv \mathcal{C}(s).$$
 (I)

• The inequality in (I) is strict unless f(u) = s(u) everywhere on [0, 1].

4ロト 4団 ト 4 豆 ト 4 豆 ト 豆 り Q ○

#### Proof.

Notice that Q is indeed invertible by our previous theorem, so s(x) is well-defined and indeed a natural cubic spline by definition.

To verify that it interpolates  $\{(t_j,y_j)\}_{j=1}^n$ , write  $s=(s(t_1),\ldots,s(t_n))^{ op}$  and note

$$s(t_i) = \sum_{i=1}^n \theta_i q(t_i, t_j),$$
 and so  $s = Q\theta = QQ^{-1}y = y.$ 

This establishes existence of at least one interpolating cubic spline, constructible explicitly via the stated form. To establish that this is the unique interpolating spline, we will:

- prove that (2) and (3) hold for any interpolating spline (not s specifically).
- using this, we will show that there can only be one interpolating spline thus closing our proof loop.

Let f be an arbitrary  $C^2$  interpolant and w(x) be an interpolating cubic spline, not necessarily equal to s. Define  $\delta(x)=f(x)-w(x)$  and remark that  $\delta(t_j)=0$  for all j since w interpolates f at the nodes. Now, expand the square to write

$$\mathcal{C}(f) = \mathcal{C}(w+\delta) = \mathcal{C}(w) + \mathcal{C}(\delta) + \int_a^b w''(u) \delta''(u) du.$$

Victor Panaretos (EPFL) Linear Models 301 / 309

We claim that the last term vanishes. Using integration by parts

$$\int_a^b w''(u)\delta''(u) \, du = w''\delta' \Big|_0^1 - \int_0^1 w'''(u)\delta'(u) \, du = - \int_0^1 w'''(u)\delta'(u) \, du$$

because w''(0) = w''(1) = 0 by the natural boundary constraint. Breaking the integration over the knot partition and using integration by parts a second time,

$$\int_0^1 w'''(u) \delta'(u) du = \sum_{j=1}^{n-1} \int_{t_j}^{t_{j+1}} w'''(u) \delta'(u) du =$$
 $= \sum_{j=1}^{n-1} \left( w''' \delta \Big|_{t_j}^{t_{j+1}} - \int_{t_j}^{t_{j+1}} w''''(u) \delta'(u) du 
ight) = 0$ 

because on each partition subinterval w''' is a constant and w'''' vanishes, whereas  $\delta(t_j)=0$  by the interpolation constraint.

This establishes that for any  $\,C^{\,2}$  interpolant f and any interpolating natural cubic spline  $\,w,$  we must have

$$C(f) = C(w) + C(\delta) \ge C(w)$$
.

The inequality

$$C(f) = C(w) + C(\delta) \ge C(w).$$

becomes an equality if and only if  $\mathcal{C}(\delta)=0$ . But if if  $\mathcal{C}(\delta)=0$ , it must be that  $\delta''=0$  because  $\delta''$  is continuous (by w'' and f'' being so). Hence,  $\delta$  is linear everywhere on [0,1], and so must be uniformly zero on [0,1] since  $\delta(t_j)=0$ .

In summary, for any interpolating spline w and any  $C^2$  interpolant,

$$C(f) \ge C(w)$$
, unless  $f = w$ . (C)

Let us use this conclusion to establish uniqueness in (1). Let  $s_1(x)$  and  $s_2(x)$  be two natural cubic splines that interpolate  $\{(t_j,y_j)\}_{j=1}^n$ . Apply conclusion (C) to  $s_1$  and  $s_2$  twice, each time reversing their roles:

- First, take  $s_2$  as an interpolating spline and  $s_1$  as some  $C^2$  interpolant. We must have  $C(s_1) > C(s_2)$  unless  $s_1 = s_2$ .
- Second, take  $s_1$  as an interpolating spline and  $s_2$  as some  $C^2$  interpolant. We must have  $C(s_2) > C(s_1)$  unless  $s_2 = s_1$ .

The only way for the two conclusions to hold simultaneously is for  $s_1 = s_2$ , which proves uniqueness in (1) and completes the proof.

## Corollary

Given distinct nodes  $0 = t_1 < t_2 < \ldots < t_n = 1$ , the set  $S(t_1, \ldots, t_n)$  of natural cubic splines with knots  $\{t_j\}_{j=1}^n$  is a vector space of dimension n, and

$$arphi_i(x) = q_{t_i}(x) = q(x,t_i), \qquad i=1,...,n$$

forms a basis for  $S(t_1, ..., t_n)$ .

### Proof.

It is immediate that  $S(t_1,...,t_n)$  is a vector space by the definition of a natural cubic spline. And, for any  $y=(y_1,...,y_n)^{\top}\in\mathbb{R}^n$  there is a unique  $s_y\in S(t_1,...,t_n)$  that interpolates  $\{(t_j,y_j)\}_{j=1}^n$ . This establishes a bijection between  $\mathbb{R}^n$  and  $S(t_1,...,t_n)$ , and proves that the dimension of  $S(t_1,...,t_n)$  is n.

To show that the collection of n functions  $\{\varphi_i\}_{i=1}^n$  is linearly independent, we need to show that if  $\theta_1\varphi_1(x)+\ldots\theta_n\varphi_n(x)=0$ , then  $\boldsymbol{\theta}=(\theta_1,\ldots,\theta_n)^\top=\mathbf{0}$ 

Note that  $0 = \sum_{j=1}^n \theta_j \varphi_j(x) \equiv \sum_{j=1}^n \theta_j q_j(x, t_j)$ , is the unique natural cubic spline that interpolates  $\{(0, t_j)\}_{j=1}^n$ . Hence, we must have that  $Q\theta = 0$ , for  $Q = \{q(t_i, t_j)\}_{i,j=1}^n$  strictly positive definite, and so  $\theta = 0$ .

## Corollary

If  $\{B_i\}_{i=1}^n$  is a basis for natural cubic splines on n distinct nodes  $0=t_1<\ldots< t_n=1$ , then the  $n\times n$  matrix  $\mathbf{B}=\{B_i(t_j)\}_{i,j=1}^n$  is invertible and the  $n\times n$  matrix  $\mathbf{\Omega}=\{\int_0^1 B_m''(x)B_k''(x)dx\}_{m,k=1}^n$  is nonnegative definite.

## Proof.

The matrix B is invertible if and only if the equation

$$B\gamma = y$$

has a unique solution with respect to  $\gamma \in \mathbb{R}^n$  for any  $\boldsymbol{y} = (y_1, \dots, y_n)^\top \in \mathbb{R}^n$ . Notice, however, that as

$$\left\{\sum_{i=1}^n \gamma_i B_i : \boldsymbol{\gamma} \in \mathbb{R}^n\right\} = \mathcal{S}(t_1, ..., t_n)$$

since  $\{B_i\}$  is a basis of  $\mathcal{S}$ . Hence the matrix statement is equivalent to asking whether for any y, there exists a unique  $s \in \mathcal{S}(t_1, ..., t_n)$  such that

$$s(t_j) = y_j, \qquad j = 1, ..., n.$$

This is guaranteed by the unique interpolation theorem.

Victor Panaretos (EPFL)

Linear Models

To show  $\Omega \succeq 0$ , note that each  $B_m$  can be expanded in the basis  $\{q_{t_i}(x)\}_{i=1}^n$  as

$$B_m(x) = \sum_{i=1}^n heta_{i,m} \, q_{t_i}(x).$$

Therefore,

$$B_m''(x) = \sum_{i=1}^n heta_{i,m} \, q_{t_i}''(x) = \sum_{i=1}^n heta_{i,m} \, k_{t_i}''(x), \quad m = 1,...,n.$$

Consequently, we can make use of our earlier calculation (\*) to get

$$\int_0^1 B_m''(x) B_k''(x) \, dx = \sum_{i=1}^n \sum_{j=1}^n heta_{i,m} heta_{j,k} \int_0^1 k_{t_i}''(x) k_{t_j}''(x) \, dx = \ \stackrel{(*)}{=} \sum_{i=1}^n heta_{i,m} \sum_{i=1}^n k(t_i,t_j) heta_{j,k}.$$

Equivalently, 
$$\Omega = \Theta^{\top} K \Theta$$
, for  $\Theta = \{\theta_{i,m}\}_{i,m=1}^n$  so  $\Omega \succeq 0$ .

# Theorem (Splines Minimise Penalised Least Squares)

Given covariates  $0 = x_1 < \ldots < x_n = 1$  and responses  $\{y_i\}_{i=1}^n$ , the functional

$$\mathcal{L}(f) = \sum_{i=1}^n ig(y_i - f(x_i)ig)^2 + \lambda \int_0^1 (f''(u))^2 du$$

is uniquely minimised at a natural cubic spline  $\hat{f}(x)$  with knots  $\{x_j\}_{j=1}^n$  expressed as

$$\hat{f}(x) = \sum_{j=1}^n \hat{\gamma}_j B_j(x),$$

with

$$(\hat{\gamma}_1,\dots,\hat{\gamma}_n)^{ op}=\hat{oldsymbol{\gamma}}=(oldsymbol{B}^{ op}oldsymbol{B}+\lambdaoldsymbol{\Omega})^{-1}oldsymbol{B}^{ op}oldsymbol{y},$$

where

- ullet  $\{B_j(x)\}_{j=1}^n$  is any basis for natural cubic spline basis with knots  $\{x_j\}_{j=1}^n$ 
  - $\mathbf{v} = (y_1, \dots, y_n)^{\top}$
  - $B = \{B_i(x_i)\}_{i,i=1}^n$  is invertible.
  - $\Omega = \left\{ \int_0^1 B_i''(t) B_j''(t) dt \right\}_{i,j=1}^n$  is non-negative definite.

#### Proof.

Let  $f \in C^2$  be a candidate minimiser, and let s(x) be the unique element of  $S(t_1,...,t_n)$  that interpolates  $\{(t_j,f(t_j)\}_{j=1}^n$ . Then,

$$egin{array}{lcl} \mathcal{L}(f) & = & \sum_{i=1}^n ig(y_i - f(x_i)ig)^2 + \lambda \int_0^1 (f''(u))^2 du \ & = & \sum_{i=1}^n ig(y_i - s(x_i)ig)^2 + \lambda \int_0^1 (f''(u))^2 du \ & \geq & \sum_{i=1}^n ig(y_i - s(x_i)ig)^2 + \lambda \int_0^1 (s''(u))^2 du \ & = & \mathcal{L}(s). \end{array}$$

with equality only if f is itself a spline. Therefore, minimisation of  $\mathcal{L}$  over all of  $C^2$ , reduces to minimisation of  $\mathcal{L}$  over the vector space  $\mathcal{S}(t_1,...,t_n)$ . Since  $\{B_1,...,B_n\}$  is a basis for  $\mathcal{S}(t_1,...,t_n)$ , our problem is equivalent to minimising

$$\Im(\gamma) = \sum_{j=1}^n \left(y_j - \sum_{i=1}^n \gamma_i B_i(x_j) 
ight)^2 + \lambda \int_0^1 \left(\sum_{i=1}^n \gamma_i B_i''(u) 
ight)^2 du$$

over  $\boldsymbol{\gamma} = (\gamma_1, ..., \gamma_n)^{\top} \in \mathbb{R}^n$ .

In matrix notation, we want to minimize w.r.t.  $\gamma$  the expression

$$(\boldsymbol{y} - \boldsymbol{B} \boldsymbol{\gamma})^{ op} (\boldsymbol{y} - \boldsymbol{B} \boldsymbol{\gamma}) + \lambda \boldsymbol{\gamma}^{ op} \boldsymbol{\Omega} \boldsymbol{\gamma}.$$

This is a ridge regression problem, and will admit the unique solution

$$\hat{oldsymbol{\gamma}} = (oldsymbol{B}^ op oldsymbol{B} + \lambda oldsymbol{\Omega})^{-1} oldsymbol{B}^ op oldsymbol{y}$$

provided the matrix  $B^{\top}B + \lambda\Omega$  is indeed invertible. This follows from the fact that B is invertible and  $\Omega$  is nonnegative definite, as per our last corollary.