

Regression Methods

Myrto Linnios

Autumn 2025 - Week 1

General Information

- ▶ Lectures on Tuesdays, 10h15-12h, CM1221
- ▶ Exercises on Tuesdays, 8h15-10h with Marco Piccininni¹

- ▶ Slides and exercises sheets available each Monday evening
- ▶ Solutions available on the following Sunday
- ▶ Some exercises sessions will be practical and require coding, solutions will be given in **R**
- ▶ If you have any question regarding exercises, please come to the dedicated sessions

- ▶ No midterm, written final exam only (3h)
- ▶ No external material allowed
- ▶ We will not respond to emails related to the final exam after December 19th

- ▶ **November 18th: Recap session.** No exercise session, problem solving during the lecture session

The Linear Model - Introduction

What is a Regression Model?

Statistical model for:

y (random output) ^{whose law is influenced by} x (non-random input)

Aim: understand the effect of x on the distribution of random variable Y

General formulation²:

$$y_i \overset{\text{independent}}{\sim} \underbrace{\text{Distribution}\{g(x_i)\}}_{=\theta_i}, \quad i = 1, \dots, n.$$

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

- ▶ Inference
- ▶ Prediction
- ▶ Data compression (parsimonious representations)

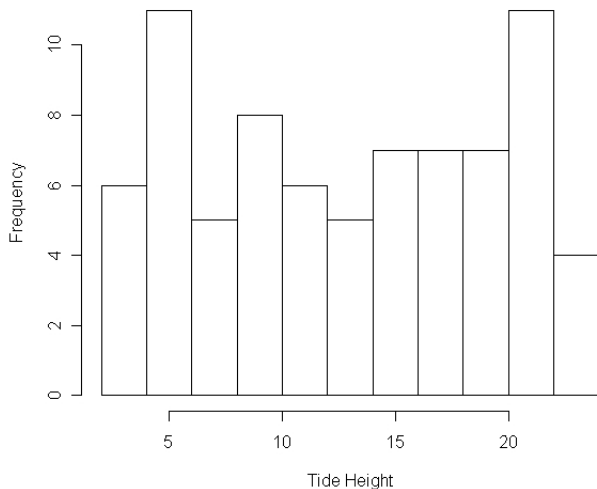
²Sometimes we write $y_i|x_i \overset{\text{independent}}{\sim} \text{Distribution}\{g(x_i) = \theta_i\}$ to highlight that the distribution of y depends on x , but without meaning that (x, y) are jointly random; such an assumption is unnecessary (e.g., in a designed experiment we choose values for x).

Absence or presence of covariates

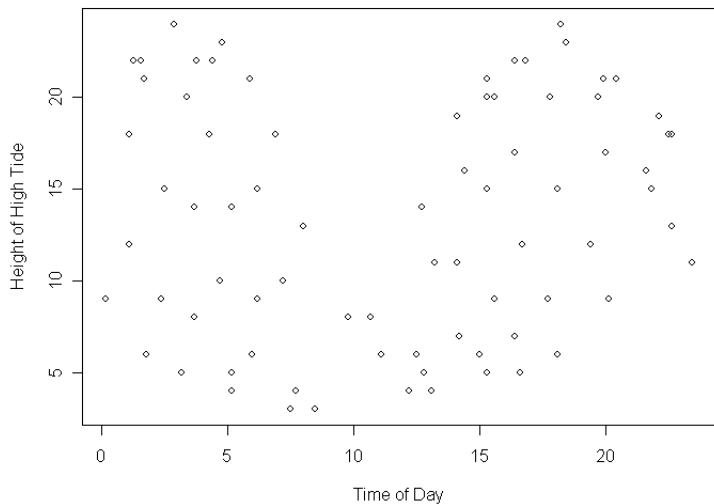
$(y_1, \dots, y_n)^\top$ has independent entries, each with **distribution $F(y; \theta_i)$ of the same family but with different parameters.**

- ▶ Each observation was generated under **slightly different experimental conditions.** They depend in a similar way on different θ_i .
- ▶ These θ_i correspond to different experimental conditions, say x_i .
- ▶ Each x_i is called a **covariate/feature**, and is an input that the experimenter can vary. They are **known**. The index i reminds us that it corresponds to the i th observation y_i .
- ▶ Usually θ_i is postulated to have a special relationship to x_i (through g), for example $\theta_i = \exp\{\alpha + \beta x_i\}$, for (α, β) unknown parameters.
- ▶ The point here is to understand the effect of varying the covariate/feature on the distribution of the observable.

How to model the height of Honolulu tides throughout the day - Histogram



Height of Honolulu tides as function of the time of day



A **bewildering variety** of models can be captured by the general specification

$$y_i \overset{\text{independent}}{\sim} \text{Distribution} \underbrace{\{g(x_i)\}}_{=\theta_i}, \quad i = 1, \dots, n.$$

x_i can be:

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

Distribution can be:

- ▶ Gaussian, Laplace, Bernoulli, Poisson, gamma, general exponential family, ...

Function $g(\cdot)$ can be:

- ▶ $g(x) = \beta_0 + \beta_1 x$, $g(x) = \sum_{k=-K}^K \beta_k e^{-ikx}$, cubic spline, neural net...

Table: A coarse classification of regression models we will consider

Distribution / Function g	$g(x_i^T) = x_i^T \beta$	g nonparametric
Gaussian	Linear Regression	Smoothing
Exponential Family	GLM	GAM

GLM: Generalized Linear Model and GAM: Generalized Additive Model

We start with a very standard model: **Linear Regression with $y|x$ being Gaussian.**

Fundamental case: Gaussian linear regression

- ▶ $y, x \in \mathbb{R}$, $g(x) = \beta_0 + \beta_1 x$

$$\begin{aligned} y | x &\sim \mathcal{N}(\beta_0 + \beta_1 x, \sigma^2) \\ &\Updownarrow \\ y &= \beta_0 + \beta_1 x + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2) \end{aligned}$$

The second version is useful for mathematical work, but is puzzling statistically, since we don't observe ϵ .

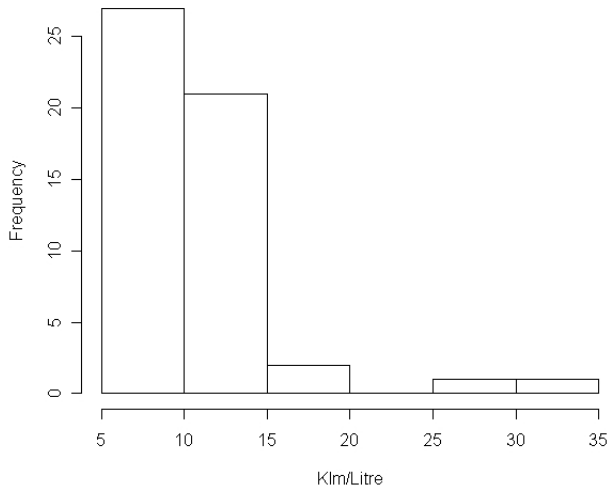
- ▶ Also, covariate could be vector ($y, \beta_0 \in \mathbb{R}$, $x \in \mathbb{R}^p$, $\beta \in \mathbb{R}^p$):

$$\begin{aligned} y | x &\sim \mathcal{N}(\beta_0 + \beta^\top x, \sigma^2) \\ &\Updownarrow \\ y &= \beta_0 + \beta^\top x + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2) \end{aligned}$$

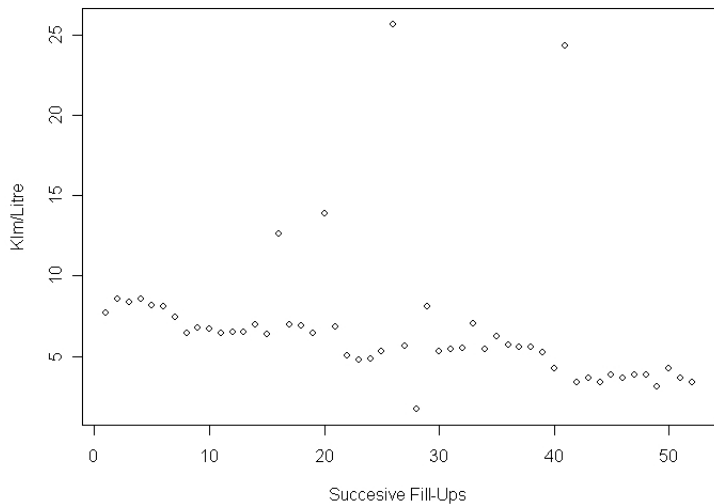
How to model my van's consumption of gas



Histogram of consumption of gas (km/L)



Gas consumption as function of successive fill-ups



The tools of the trade ...

Start from **Gaussian linear regression** then gradually generalise ...

Obviously: important features of Gaussian linear model are

- ▶ Gaussian distribution
- ▶ Linearity

These two **combine well** and give **geometric insights** to solve the estimation problem related to some **probabilistic linear algebra**...

- ▶ Subspaces and projection matrices
- ▶ Multivariate Gaussian Distribution
- ▶ Optimal dimension reduction
- ▶ Random quadratic forms

To sum up

- ▶ **Regression**: (statistics) a measure of the relation between the mean value of
 - ▶ one variable (e.g., output), denoted y (the **response variable**) and
 - ▶ corresponding values of other variables (e.g., time and cost), denoted x (**explanatory variables**).
- ▶ The explanatory variables are also called **covariates** or **features** (ML).
- ▶ We avoid the terms **dependent variable** (y) and **independent variable** (x) used in older books.
- ▶ Questions we try and answer:
 - ▶ (**description/explanation**) how does y depend on x ? How much of the variation of y is due to x ? Do I need all of x to explain the variation in y ?
 - ▶ (**prediction**) what will y be if $x = x_+$?
 - ▶ (**causation**) if I change x , what will happen to y ?
- ▶ The causation question presupposes that we can change (some of) x , which is not always true.

The Linear Model - First Results and Geometric Interpretations

Linear model

- ▶ Simplest explanation of y in terms of x is **linear model**:

$$y = g(x) = x_1\beta_1 + \dots + x_p\beta_p = x^T\beta,$$

where

$$y \in \mathbb{R}, \quad x^T = (x_1, \dots, x_p) \in \mathbb{R}^p, \quad \beta^T = (\beta_1, \dots, \beta_p) \in \mathbb{R}^p.$$

- ▶ The data consist of n **instances/examples/cases** (x_j, y_j) for $j = 1, \dots, n$, so

$$y_{n \times 1} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}, \quad X_{n \times p} = \begin{pmatrix} x_{11} & \dots & x_{1p} \\ \vdots & \vdots & \vdots \\ x_{n1} & \dots & x_{np} \end{pmatrix}, \quad \beta_{p \times 1} = \begin{pmatrix} \beta_1 \\ \vdots \\ \beta_p \end{pmatrix}$$

and we write

$$y = X\beta.$$

- ▶ **Key point**: linearity refers to linearity in β , not in terms of elements of X , which might be polynomials, or basis functions, or ...
- ▶ Sometimes we can transform to a linear model. For example, the multiplicative expression $y = \gamma x_1^{\beta_1} x_2^{\beta_2}$ becomes

$$\log y = \log \gamma + \beta_1 \log x_1 + \beta_2 \log x_2.$$

Notation

- ▶ Vectors are column vectors
- ▶ We write $X_{n \times p}$ to give the dimensions of a matrix or vector
- ▶ a^T (row vector) is the transpose of a (column vector)
- ▶ $j \in \{1, \dots, n\}$ (or sometimes i) indexes the rows of y (cases/examples)
- ▶ x_j^T is the j th row of X
- ▶ $r, s, t, \dots \in \{1, \dots, p\}$ indexes the columns of X (covariates/features)
- ▶ Roman letters (y, X, z, \dots) denote observed quantities, and may be the realisations of random variables
- ▶ Greek letters ($\beta, \gamma, \theta, \sigma, \dots$) denote unknown (often vector) parameters of models
- ▶ $\hat{\beta}$ denotes an estimate of β
- ▶ α denotes the level of significance tests and confidence intervals
- ▶ If Q is scalar (or a row vector) and β is a vector, then $\partial Q / \partial \beta$ denotes the vector (or matrix) the same shape as β with elements $\partial Q / \partial \beta_r$.
- ▶ If Q is scalar and β, γ are vectors, then $\partial^2 Q / \partial \beta \partial \gamma^T$ denotes the matrix with (r, s) element $\partial^2 Q / \partial \beta_r \partial \gamma_s$.
- ▶ $u \perp v$ means that the vectors u and v are orthogonal (i.e., $u^T v = 0$); ditto for matrices.
- ▶ $Y \perp\!\!\!\perp Z$ means that the random variables Y and Z are independent.

Useful matrix decompositions

- ▶ **Singular value decomposition (SVD)**: any real matrix X can be written in the form

$$X_{n \times p} = U_{n \times n} D_{n \times p} V_{p \times p}^T$$

where

- ▶ $U = (u_1, \dots, u_n)$ and $V = (v_1, \dots, v_p)$ are orthogonal (i.e., $U^T U = U U^T = I_n$, $V^T V = V V^T = I_p$) and D is $n \times p$ rectangular diagonal with real diagonal entries (**singular values**)
 $d_1 \geq \dots \geq d_m \geq 0$, where $m = \min(n, p)$,
 - ▶ if one or more $d_j = 0$, then X is singular, and
 - ▶ the u_j and v_r respectively span the column and row spaces of X .
- ▶ The SVD implies that the ranks of X , $X^T X$ and $X X^T$ are equal and at most m .
 - ▶ **Spectral theorem**: any real symmetric matrix H can be written as

$$H_{n \times n} = U_{n \times n} D_{n \times n} U_{n \times n}^T,$$

where

- ▶ $D = \text{diag}(d_1, \dots, d_n)$ contains the eigenvalues of H ;
- ▶ U is an orthogonal matrix whose columns are the corresponding eigenvectors; and
- ▶ if H is positive semi-definite then $d_1 \geq \dots \geq d_n \geq 0$.

Least squares fit

- ▶ Assume that

$$y = X\beta$$

and find the ‘best fit’ by choosing β to minimise the (squared) Euclidean distance between y and $X\beta$, i.e., the sum of squares

$$\|y - X\beta\|^2 = (y - X\beta)^T (y - X\beta) = \sum_{j=1}^n (y_j - x_j^T \beta)^2.$$

- ▶ In vector space terms, $y \in \mathbb{R}^n$ and $X\beta \in \text{span}(X) \subset \mathbb{R}^n$.
- ▶ The ‘best fit’ vector \hat{y} is the vector in $\text{span}(X)$ closest to y ; Pythagoras’ theorem (sketch) gives $\hat{y} \perp (y - \hat{y})$ (but see below).
- ▶ We call $\hat{y} \in \mathbb{R}^n$ the **fitted value(s)** and $e = y - \hat{y} \in \mathbb{R}^n$ the **residual (vector)**.
- ▶ Recall: $\underbrace{\mathcal{M}(X)}_{\text{Column Space}} = \text{span}(X) := \{X\gamma : \gamma \in \mathbb{R}^p\}$

Optimal regressor

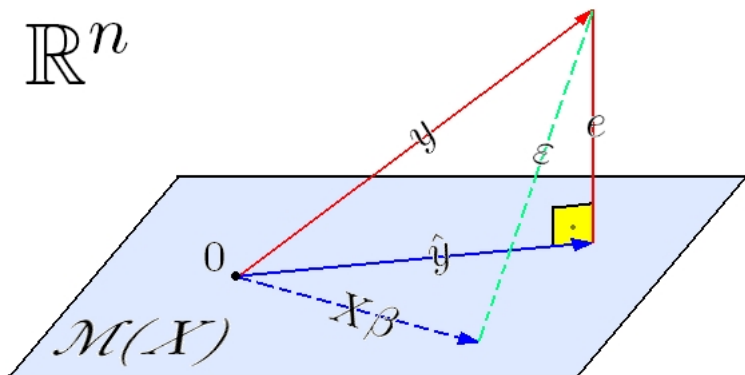
Lemma 1

When X has rank p and $n \geq p$ then $\hat{y} = X\hat{\beta} = Hy$, where

$$\hat{\beta} = (X^T X)^{-1} X^T y, \quad H = X(X^T X)^{-1} X^T.$$

The 'hat matrix' H has rank p , is symmetric and idempotent, and satisfies $HX = X$: it gives the orthogonal projection of \mathbb{R}^n onto $\text{span}(X)$.

The (column) geometry of least squares



Proof

- ▶ If X has rank p , so too does the $p \times p$ matrix $X^T X$, which is therefore invertible.
- ▶ The sum of squares

$$\begin{aligned} Q &= (y - X\beta)^T (y - X\beta) = y^T y - \beta^T X^T y - y^T X\beta + \beta^T X^T X\beta \\ &= y^T y - 2y^T X\beta + \beta^T X^T X\beta \end{aligned}$$

has first and second derivatives (respectively a $p \times 1$ vector and $p \times p$ matrix)

$$\frac{\partial Q}{\partial \beta} = -2X^T y + 2X^T X\beta, \quad \frac{\partial^2 Q}{\partial \beta \partial \beta^T} = 2X^T X$$

with respect to β . Setting $\partial Q / \partial \beta = 0$ implies that $(X^T X)\beta = X^T y$, and as $X^T X$ is invertible we can write

$$\hat{\beta} = (X^T X)^{-1} X^T y, \quad \hat{y} = X\hat{\beta} = X(X^T X)^{-1} X^T y = Hy,$$

say. The matrix $X^T X$ is positive definite, so $(y - X\beta)^T (y - X\beta)$ is minimised at $\hat{\beta}$.

Note to Lemma 1

- ▶ The $n \times n$ ‘hat matrix’ H (which ‘puts a hat’ on y) satisfies $\widehat{H}^T = H$, $H^2 = H$, so it is symmetric and idempotent, i.e., its eigenvalues equal 0 or 1, and their multiplicities must be $n - p$ and p , as its rank is p .

H is the matrix that projects \mathbb{R}^n orthogonally onto the span of the columns of X , $\text{span}(X)$.

- ▶ The inner product between \widehat{y} and $y - \widehat{y}$ equals zero, because $\widehat{y} = Hy$, $y - \widehat{y} = (I - H)y$, and $\widehat{y}^T(y - \widehat{y}) = y^T H^T(I - H)y = y^T(H - H)y = 0$.

Hence \widehat{y} and $y - \widehat{y}$ are orthogonal.

- ▶ Clearly $HX = X(X^T X)^{-1} X^T X = X$, so $H(X\beta) = X\beta$ for any $\beta \in \mathbb{R}^p$, i.e., a vector in $\text{span}(X)$ is left unchanged by multiplication by H .

Analysis of variance I

Lemma 2

Let $X_{n \times p} = (X_0, X_1, \dots, X_R)$ have rank p , where $p \leq n$, and let H_r denote the projection matrices formed using (X_0, \dots, X_r) , for $r = 0, \dots, R$; hence $H_R = H$.

Define $P_r = H_r - H_{r-1}$ for $r = 1, \dots, R$ and $P_{R+1} = I - H$.

Then:

- (i) $H_r H_s = H_r$ whenever $r \leq s$,
- (ii) $H_0 P_r = 0$ for any r ,
- (iii) the matrices P_r are symmetric and idempotent, with $P_r P_s = 0$ when $r \neq s$.

Rephrasing the question: Gaussian linear model

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

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

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

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

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:

$$y = X\beta + \varepsilon = (X_1 \ X_2) \begin{pmatrix} \beta_1 \\ \beta_2 \end{pmatrix} + \varepsilon = X_1\beta_1 + X_2\beta_2 + \varepsilon.$$

Our question can now be stated as:

- ▶ Is $\beta_2 = 0$?

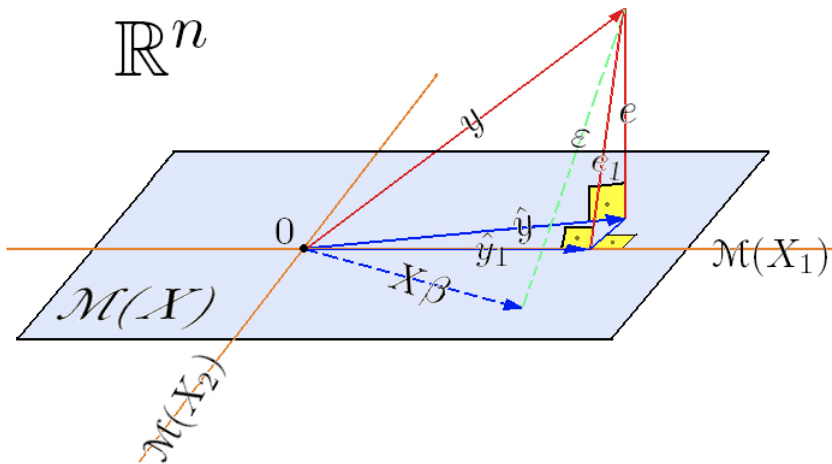


Figure: Geometry Revisited

Proof

- (i) Let $\mathcal{V}_0 \subset \dots \subset \mathcal{V}_R$ denote the *nested* linear spaces onto which \mathbb{R}^n is projected by $H_0, \dots, H_R = H$, and suppose that $r \leq s$. Now $H_r y \in \mathcal{V}_r$ for any $y \in \mathbb{R}^n$, so as $\mathcal{V}_r \subset \mathcal{V}_s$, $H_r y \in \mathcal{V}_s$. Hence $H_s H_r y = H_r y$ for any $y \in \mathbb{R}^n$, so $H_s H_r = H_r$. This implies that

$$H_s H_r = H_r = H_r^T = (H_s H_r)^T = H_r^T H_s^T = H_r H_s, \quad s \geq r.$$

- (ii) For $r = 1, \dots, R$, (i) yields $H_0 P_r = H_0 H_r - H_0 H_{r-1} = H_0 - H_0 = 0$, and $H_0 P_{R+1} = H_0(I - H_R) = 0$.
- (iii) The matrices P_1, \dots, P_R are symmetric because

$$P_r^T = (H_r - H_{r-1})^T = H_r^T - H_{r-1}^T = H_r - H_{r-1} = P_r,$$

and idempotent because (i) gives

$$\begin{aligned} P_r^2 &= (H_r - H_{r-1})(H_r - H_{r-1}) \\ &= H_r^2 - H_r H_{r-1} - H_{r-1} H_r + H_{r-1}^2 \\ &= H_r - H_{r-1} - H_{r-1} + H_{r-1} \\ &= H_r - H_{r-1} = P_r. \end{aligned}$$

Moreover if $r < s \leq R$, then

$$\begin{aligned} P_r P_s &= (H_r - H_{r-1})(H_s - H_{s-1}) \\ &= H_r H_s - H_r H_{s-1} - H_s H_{r-1} + H_{r-1} H_{s-1} \\ &= H_r - H_r - H_{r-1} + H_{r-1} \\ &= 0. \end{aligned}$$

The corresponding results for P_{R+1} are equally easy to check.

General decomposition

- ▶ In the setup of Lemma 2 suppose we fit the models with projection matrices $H_0, \dots, H_R = H$ and corresponding fitted values $\hat{y}_r = H_r y$. Then

$$\begin{aligned}y &= \hat{y}_0 + (\hat{y}_1 - \hat{y}_0) + \dots + (\hat{y}_R - \hat{y}_{R-1}) + (y - \hat{y}_R) \\ &= H_0 y + (H_1 - H_0) y + \dots + (H_R - H_{R-1}) y + (I - H) y \\ &= H_0 y + P_1 y + \dots + P_R y + P_{R+1} y,\end{aligned}$$

and Lemma 2 implies that the terms on the RHS are orthogonal, i.e.,

$$(H_0 y)^T (P_r y) = 0, \quad (P_s y)^T (P_r y) = 0, \quad r \neq s.$$

- ▶ Hence Pythagoras' theorem gives the **analysis of variance (ANOVA)** decomposition

$$\|y\|^2 = \|\hat{y}_0\|^2 + \sum_{r=1}^R \|\hat{y}_r - \hat{y}_{r-1}\|^2 + \|y - \hat{y}_R\|^2.$$

Analysis of variance II

- ▶ Usually $X_0 = 1_n$ (*intercept*); then $\hat{y}_0 = 1_n(1_n^T 1_n)^{-1} 1_n^T y = \bar{y} 1_n$ and

$$\|y\|^2 - \|\hat{y}_0\|^2 = \sum_{j=1}^n y_j^2 - \sum_{j=1}^n \bar{y}^2 = \sum_{j=1}^n (y_j - \bar{y})^2,$$

equals n times the empirical variance of y_1, \dots, y_n . Hence

$$\sum_{j=1}^n (y_j - \bar{y})^2 = \|y\|^2 - \|\hat{y}_0\|^2 = \sum_{r=1}^R \|\hat{y}_r - \hat{y}_{r-1}\|^2 + \|y - \hat{y}\|^2$$

decomposes ('analyses') the variability of y around its average \bar{y} into

- ▶ the contributions $\|\hat{y}_r - \hat{y}_{r-1}\|^2$ due to adding the columns of X_r to X_0, \dots, X_{r-1} ,
 - ▶ the **residual sum of squares** $\|y - \hat{y}\|^2$ left after fitting $X = (X_0, \dots, X_R)$.
- ▶ Large $\|\hat{y}_r - \hat{y}_{r-1}\|^2$ implies that X_r explains a lot of the variation of y even after allowing for that explained by X_0, \dots, X_{r-1} .
 - ▶ The
 - ▶ **degrees of freedom** of a fit is the rank ν_r of the corresponding H_r , and the
 - ▶ **residual degrees of freedom** is $n - \nu_R = n - p$.

Terms

- ▶ A constant column $X_0 = 1_n$ is almost always present in the design matrix, so

$$X\beta = (1_n \quad X_1 \quad \cdots \quad X_R) \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_R \end{pmatrix} = 1_n\beta_0 + X_1\beta_1 + \cdots + X_R\beta_R,$$

where the matrices X_1, \dots, X_R , the **terms**, are successively included.

- ▶ The baseline model with only 1_n has fitted value and residual vector

$$\hat{y}_0 = \bar{y}1_n, \quad y - \hat{y}_0 = y - \bar{y}1_n.$$

- ▶ Starting from the baseline we ask which terms lead to large reductions in the residual sum of squares, i.e., best explain the variation of y .
- ▶ The successive residual degrees of freedom, i.e., the ranks of the matrices $I - H_r$, are

$$n - 1 = n - \nu_0 \geq n - \nu_1 \geq \cdots \geq n - \nu_R.$$

- ▶ When the columns of X_{r+1} depend linearly on those of $1_n, X_1, \dots, X_r$, we have $\nu_{r+1} = \nu_r$, so inclusion of X_{r+1} does not change the fitted value or improve the fit.

Model formulae

- ▶ A mean vector such as $1_n\beta_0 + X_1\beta_1 + X_2\beta_2$ is often written as the right-hand side of

$$y \sim X1 + X2$$

where

- ▶ the columns of 1_n is (silently) included first by default,
 - ▶ $X1$ and $X2$ represent the vector subspaces of \mathbb{R}^n generated by the corresponding terms, and
 - ▶ $+$ represents addition of vector subspaces.
- ▶ Software generally drops any column of a design matrix that is linearly dependent on previous columns, and this affects which elements of β can be estimated and the meaning of estimates corresponding to later columns.
 - ▶ Carefully choosing the order of terms in a model can give easily interpreted estimates of the parameters of interest — for example, if X_2 is full-rank and a column of 1_n lies in $\text{span}(X_1) + \text{span}(X_2)$ then

$$y \sim X1 + X2, \quad y \sim X2 + X1 - 1,$$

span the same linear space but the second estimates the parameters of β_2 (unadjusted for the mean) and the parameters of β_1 , adjusted for the presence of X_2 .

ANOVA

Terms	Residual df	Residual SS	Term added	Reduction in residual df	Reduction in SS	Mean square
1_n	$n - \nu_0 = n - 1$	MS_0				
$1_n, X_1$	$n - \nu_1$	MS_1	X_1	$\nu_1 - \nu_0$	$MS_0 - MS_1$	$\frac{MS_0 - MS_1}{\nu_1 - \nu_0}$
$1_n, X_1, X_2$	$n - \nu_2$	MS_2	X_2	$\nu_2 - \nu_1$	$MS_1 - MS_2$	$\frac{MS_1 - MS_2}{\nu_2 - \nu_1}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
$1_n, X_1, \dots, X_R$	$n - \nu_R = n - p$	MS_R	X_R	$\nu_R - \nu_{R-1}$	$MS_{R-1} - MS_R$	$\frac{MS_{R-1} - MS_R}{\nu_R - \nu_{R-1}}$

- ▶ The sum of squares when including terms $1_n, X_1, \dots, X_r$ is

$$MS_r = \|y - \hat{y}_r\|^2.$$

- ▶ The ‘mean square’ for term X_r ,

$$MS_r = \frac{MS_{r-1} - MS_r}{\nu_r - \nu_{r-1}}$$

is the average reduction in MS_r per degree of freedom when X_r is added to the model.

- ▶ Usually show only the RHS of the table and the bottom line of its LHS (next slide).

ANOVA table

Term added	df	Reduction in SS	Mean square
X_1	$\nu_1 - \nu_0$	$\mathbf{MS}_0 - \mathbf{MS}_1$	$\mathbf{MS}_1 = (\mathbf{MS}_0 - \mathbf{MS}_1)/(\nu_1 - \nu_0)$
X_2	$\nu_2 - \nu_1$	$\mathbf{MS}_1 - \mathbf{MS}_2$	$\mathbf{MS}_2 = (\mathbf{MS}_1 - \mathbf{MS}_2)/(\nu_2 - \nu_1)$
\vdots	\vdots	\vdots	\vdots
X_R	$\nu_R - \nu_{R-1}$	$\mathbf{MS}_{R-1} - \mathbf{MS}_R$	$\mathbf{MS}_R = (\mathbf{MS}_{R-1} - \mathbf{MS}_R)/(\nu_R - \nu_{R-1})$
Residual	$n - \nu_R$	\mathbf{MS}_R	$\mathbf{MS}_{\text{Res}} = \mathbf{MS}_R/(n - \nu_R)$

- ▶ Used to screen which terms give the largest reductions, comparing \mathbf{MS}_r with the residual mean square \mathbf{MS}_{Res} .
- ▶ Judge ‘significance’ of reductions relative to residual using F -tests (later).
- ▶ Problem: the order of adding terms matters, so there is no unique reduction in general.

Coefficient of determination

- ▶ **Coefficient of determination R^2** measures reduction in variance of y as

$$R^2 = \frac{\|\hat{y} - \bar{y}1_n\|^2}{\|y - \bar{y}1_n\|^2} = \frac{\{(H - H_0)y\}^T (H - H_0)y}{\{(I - H_0)y\}^T (I - H_0)y} = \frac{y^T (H - H_0)y}{y^T (I - H_0)y},$$

where H_0 and H are the hat matrices for regression on 1_n and X , and $1_n \in \text{span}(X)$.

- ▶ $R^2 \in [0, 1]$ is the squared empirical correlation between y and \hat{y} , so $R^2 \approx 1$ implies that most of the variation in y is explained by \hat{y} .
- ▶ There is a geometric interpretation, as the terms on the right of

$$(I_n - H_0)y = (I_n - H)y + (H - H_0)y$$

are orthogonal (check this).

- ▶ Adding columns to X must increase R^2 , unlike the **adjusted R^2** ,

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

- ▶ If $1_n \notin \text{span}(X)$, use

$$R_0^2 = \frac{\hat{y}^T \hat{y}}{y^T y}, \quad R_{0,a}^2 = R_0^2 + (1 - R_0^2) \frac{n}{n - p}.$$

Comments

- ▶ We have supposed that $X_{n \times p}$ has rank p :
 - ▶ if X is rank-deficient, then a least squares algorithm usually drops columns that lie in the span of preceding ones, but care is needed to construct X so that the resulting $\hat{\beta}$ is easy to interpret;
 - ▶ if X is nearly rank-deficient, then regularisation may be needed. More later ...
- ▶ Everything so far as purely numerical:
 - ▶ least squares estimation is a numerical technique for using X to approximate y ;
 - ▶ $\hat{y} = X\hat{\beta}$ is the resulting approximation, which lies in $\text{span}(X)$;
 - ▶ $\hat{\beta}$ gives the coefficients of the columns of X for the best approximation;
 - ▶ the coefficient of determination R^2 measures how much of the overall variation of y was explained by X ; and
 - ▶ the ANOVA decomposition summarises how much of the variation in y is explained by different subsets of columns of X (terms).
- ▶ For statistics we need to add some distributional assumptions ... shortly ...
- ▶ First some reminders ...