

Statistics for Data Science: Week 9.2

Myrto Limnios and Rajita Chandak

Institute of Mathematics – EPFL

`rajita.chandak@epfl.ch`, `myrto.limnios@epfl.ch`



ÉCOLE POLYTECHNIQUE
FÉDÉRALE DE LAUSANNE

Distribution Theory of Least Squares

Gaussian Linear Model:

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

We have derived the estimators:

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

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

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

Theorem (Sampling Distribution of LSE under Gaussian Model)

Let $\mathbf{Y}_{n \times 1} = \mathbf{X}_{n \times p} \boldsymbol{\beta}_{p \times 1} + \boldsymbol{\varepsilon}_{n \times 1}$ with $\boldsymbol{\varepsilon} \sim \mathcal{N}_n(0, \sigma^2 \mathbf{I}_{n \times n})$ and assume that \mathbf{X} has full rank $p < n$. Then,

- 1 $\hat{\boldsymbol{\beta}} \sim \mathcal{N}_p\{\boldsymbol{\beta}, \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1}\};$
- 2 the random variables $\hat{\boldsymbol{\beta}}$ and S^2 are independent; and
- 3 $\frac{n-p}{\sigma^2} S^2 \sim \chi_{n-p}^2$, where χ_ν^2 denotes the chi-square distribution with ν degrees of freedom.

Proof.

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

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

2. If \mathbf{e} is independent of $\hat{\mathbf{Y}} = \mathbf{X}\hat{\beta}$, then $S^2 = \mathbf{e}^\top \mathbf{e} / (n - p)$ will be independent of $\hat{\beta}$ (why?). Now notice that:

- $\mathbf{e} = (\mathbf{I} - \mathbf{H})\mathbf{Y}$, with $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$
- $\hat{\mathbf{Y}} = \mathbf{H}\mathbf{Y}$
- $\mathbf{Y} \sim \mathcal{N}(\mathbf{X}\beta, \sigma^2 \mathbf{I})$

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

3. For the last part recall that

$$\mathbf{e} = (\mathbf{I} - \mathbf{H})\boldsymbol{\varepsilon} \implies (n - p)S^2 = (n - p) \frac{\mathbf{e}^\top \mathbf{e}}{n - p} = \boldsymbol{\varepsilon}^\top (\mathbf{I} - \mathbf{H})\boldsymbol{\varepsilon}$$

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

$$\frac{(n-p)}{\sigma^2} S^2 = (\sigma^{-1}\boldsymbol{\varepsilon})^\top (\mathbf{I} - \mathbf{H})(\sigma^{-1}\boldsymbol{\varepsilon}) \sim \chi_{n-p}^2.$$

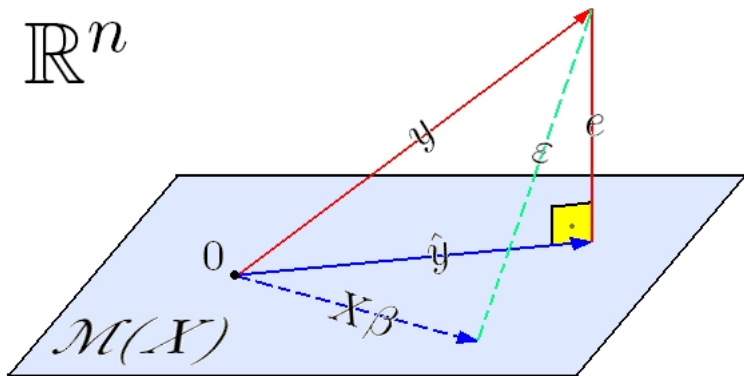
Recall that $\mathbf{H} = \mathbf{X}(\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top$.

Corollary

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

Corollary

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



Proof of the first Corollary.

Write $\mathbf{Y} = \mathbf{H}\mathbf{Y} + (\mathbf{I} - \mathbf{H})\mathbf{Y} = \hat{\mathbf{Y}} + \mathbf{e}$.

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

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

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



Proof of the second Corollary.

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



We have characterized various properties of LSE,
now what about characterizing whether the resulting linear model is *good*?

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

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

In conclusion:

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

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

$$\frac{\mathbf{c}^\top \hat{\beta} - \mathbf{c}^\top \beta}{\sqrt{S^2 \mathbf{c}^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{c}}} \sim t_{n-p}$$

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

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

- What about a $(1 - \alpha)$ CI for β_r ? (r th coordinate)
- Let $\mathbf{c}_r = (0, 0, \dots, 0, \underset{r^{th} \text{ position}}{1}, 0, \dots, 0)$
- Then $\beta_r = \mathbf{c}_r^\top \boldsymbol{\beta}$
- Therefore, base CI on

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

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

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

$$\hat{\beta}_r \pm t_{n-p}(\alpha/2) \sqrt{S^2 v_{rr}}.$$

A **prediction interval** aims to give confidence bounds on a **potential response**.

- Suppose we want to predict the value of Y_+ for an $\mathbf{x}_+ \in \mathbb{R}^p$
- Our model predicts Y_+ by $\mathbf{x}_+^\top \hat{\boldsymbol{\beta}}$.
- But $Y_+ = \mathbf{x}_+^\top \boldsymbol{\beta} + \varepsilon_+$ so a **prediction interval is different** from an interval for a linear combination $\mathbf{c}^\top \boldsymbol{\beta}$ (extra uncertainty due to ε_+):
 - $\mathbb{E}[\mathbf{x}_+^\top \hat{\boldsymbol{\beta}} + \varepsilon_+] = \mathbf{x}_+^\top \boldsymbol{\beta}$
 - $\text{var}[\mathbf{x}_+^\top \hat{\boldsymbol{\beta}} + \varepsilon_+] = \text{var}[\mathbf{x}_+^\top \hat{\boldsymbol{\beta}}] + \text{var}[\varepsilon_+] = \sigma^2[\mathbf{x}_+^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_+ + 1]$
- Base prediction interval on:

$$\frac{\mathbf{x}_+^\top \hat{\boldsymbol{\beta}} - Y_+}{\sqrt{S^2 \{1 + \mathbf{x}_+^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_+\}}} \sim t_{n-p}.$$

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

$$\mathbf{x}_+^\top \hat{\boldsymbol{\beta}} \pm t_{n-p}(\alpha/2) \sqrt{S^2 \{1 + \mathbf{x}_+^\top (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{x}_+\}}.$$

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

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

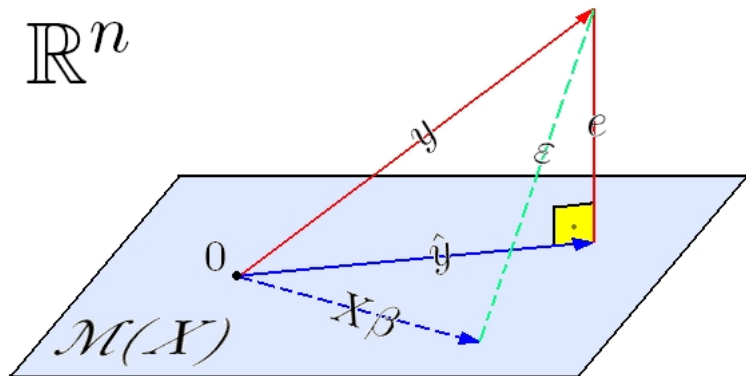
$$\frac{\|\mathbf{e}\|^2}{\|\mathbf{Y}\|^2} = \frac{\mathbf{e}^\top \mathbf{e}}{\mathbf{Y}^\top \mathbf{Y}} = \frac{\mathbf{Y}^\top (\mathbf{I} - \mathbf{H}) \mathbf{Y}}{\mathbf{Y}^\top \mathbf{Y}} = 1 - \frac{\hat{\mathbf{Y}}^\top \hat{\mathbf{Y}}}{\mathbf{Y}^\top \mathbf{Y}}$$

- Define

$$R_0^2 = \frac{\hat{\mathbf{Y}}^\top \hat{\mathbf{Y}}}{\mathbf{Y}^\top \mathbf{Y}} = \frac{\|\hat{\mathbf{Y}}\|^2}{\|\mathbf{Y}\|^2}$$

- Note that $0 \leq R_0^2 \leq 1$

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



“**Centred (in fact, usual) R^2** ”. Compares empirical variance of $\hat{\mathbf{Y}}$ to empirical variance of \mathbf{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 $\mathbf{e} \perp \mathbf{1}$ (here $\mathbf{1}$ is the vector of 1's = first column of design matrix \mathbf{X}) so $\sum_i e_i = 0$).

Note that

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

- R_0^2 mathematically more natural (does not treat first column of \mathbf{X} as special).
- R^2 statistically more relevant (expresses variance—the first column of \mathbf{X} usually *is* special, in statistical terms!).
- R_0^2 and R^2 may differ a lot when $\bar{Y} = \frac{1}{n} \sum_{i=1}^n Y_i$ large.

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

- $\mathbf{1}(\mathbf{1}^\top \mathbf{1})^{-1} \mathbf{1}^\top \mathbf{Y} = \mathbf{1} n^{-1} \sum_{i=1}^n Y_i = \mathbf{1} \bar{Y}$.
- $\mathbf{1}(\mathbf{1}^\top \mathbf{1})^{-1} \mathbf{1}^\top \hat{\mathbf{Y}} = \mathbf{1} n^{-1} \sum_{i=1}^n \hat{Y}_i = \mathbf{1} \bar{Y}$.

So

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

Intuition: Should not take into account the part of $\|\mathbf{Y}\|$ 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 $R^2 \leq R_0^2$.

The adjusted R^2 takes into account the number of variables employed. It is defined as:

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

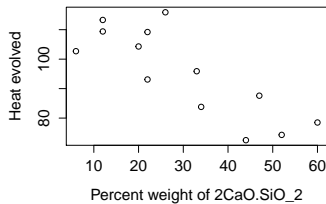
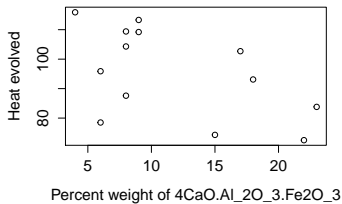
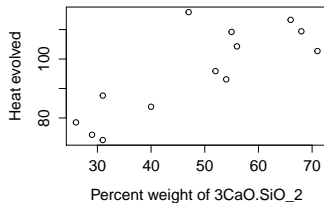
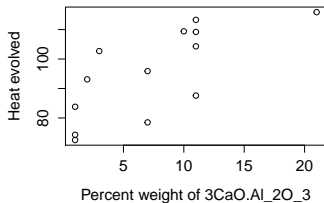
Corrects for the fact that we can always increase R^2 by adding variables. One can also correct the un-centred R_0^2 and take into account instead

$$R_{a0}^2 = 1 + (1 - R_0^2) \frac{n}{n - p}.$$

Example: Cement Heat Evolution

Case	$3\text{CaO} \cdot \text{Al}_2\text{O}_3$	$3\text{CaO} \cdot \text{SiO}_2$	$4\text{CaO} \cdot \text{Al}_2\text{O}_3 \cdot \text{Fe}_2\text{O}_3$	$2\text{CaO} \cdot \text{SiO}_2$	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~1+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
x1	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
x4	-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