# Statistical Learning II

## Lecture 4 - Least squares

_____

**Bruno Loureiro**

@ CSD, DI-ENS & CNRS

brloureiro@gmail.com

# Summary of ERM

Let $\mathscr{D} = \{(x_i, y_i) \in \mathscr{X} \times \mathscr{Y} : i = 1, \ldots, n\}$ denote training data sampled i.i.d. from $p$.

Given a choice of:

- Parametric hypothesis class $\mathscr{H} = \{f_\theta : \mathscr{X} \to \mathscr{Y} : \theta \in \Theta\}$

- Loss function $\ell : \mathscr{X} \times \mathscr{Y} \to \mathbb{R}_+$

Empirical Risk Minimisation consists of:

$$\min_{\theta \in \Theta} \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, f_\theta(x_i))$$

# Summary of ERM

Let $\mathcal{D} = \{(x_i, y_i) \in \mathcal{X} \times \mathcal{Y} : i = 1, \ldots, n\}$ denote training data sampled i.i.d. from $p$.

Given a choice of:

- Parametric hypothesis class $\mathcal{H} = \{f_\theta : \mathcal{X} \to \mathcal{Y} : \theta \in \Theta\}$

- Loss function $\ell : \mathcal{X} \times \mathcal{Y} \to \mathbb{R}_+$

Empirical Risk Minimisation consists of:

$$\min_{\theta \in \Theta} \ \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, f_\theta(x_i))$$

# Key questions

- What optimisation procedure to choose?

$$F(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, f_\theta(x_i))$$

Is typically a non-convex function of $\theta \in \Theta$.

# Key questions

- What optimisation procedure to choose?

$$F(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, f_\theta(x_i))$$

Is typically a non-convex function of $\theta \in \Theta$.

- How large $n$ needs to be (with respect to $p, d$) so that $\hat{\theta} \in \operatorname{argmin} F(\theta)$ has low training and/or test error?

# Key questions

- What optimisation procedure to choose?

$$F(\theta) = \frac{1}{n} \sum_{i=1}^{n} \ell(y_i, f_\theta(x_i))$$

Is typically a non-convex function of $\theta \in \Theta$.

- How large $n$ needs to be (with respect to $p, d$) so that $\hat{\theta} \in \mathrm{argmin}\ F(\theta)$ has low training and/or test error?

- What properties of the data distribution $p$ makes the problem easier / harder?

# Least-squares regression

# Least-squares regression

Let $\mathcal{D} = \{(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R} : i = 1, \ldots, n\}$ denote the training data.

Ordinary least-squares (OLS) regression is defined as:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \hat{\mathcal{R}}_n(\boldsymbol{\theta}) := \frac{1}{2n} \sum_{i=1}^{n} \left(y_i - \langle \boldsymbol{\theta}, \boldsymbol{x}_i \rangle\right)^2$$

# Least-squares regression

Let $\mathcal{D} = \{(x_i, y_i) \in \mathbb{R}^d \times \mathbb{R} : i = 1, \ldots, n\}$ denote the training data.

Ordinary least-squares (OLS) regression is defined as:

$$\min_{\boldsymbol{\theta} \in \mathbb{R}^d} \hat{\mathcal{R}}_n(\boldsymbol{\theta}) := \frac{1}{2n} ||\boldsymbol{y} - \boldsymbol{X\theta}||_2^2$$

Where we have defined the data matrix $\boldsymbol{X} \in \mathbb{R}^{n \times d}$ and label vector $\boldsymbol{y} \in \mathbb{R}^n$:

$$\boldsymbol{X} = \begin{bmatrix} - & \boldsymbol{x}_1 & - \\ - & \boldsymbol{x}_2 & - \\ & \vdots & \\ - & \boldsymbol{x}_n & - \end{bmatrix} \in \mathbb{R}^{n \times d} \qquad \boldsymbol{y} = \begin{bmatrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{bmatrix}$$

# Bayes risk for OLS

- This corresponds to an ERM problem on the class of linear functions:

$$\mathcal{H} = \{f_\theta(\boldsymbol{x}) = \langle \boldsymbol{\theta}, \boldsymbol{x} \rangle : \boldsymbol{\theta} \in \mathbb{R}^d\}$$

with the square loss functions:

$$\ell(y, f_\theta(\boldsymbol{x})) = \frac{1}{2}\left(y - f_\theta(\boldsymbol{x})\right)^2$$

# Bayes risk for OLS

Remarks:

- This corresponds to an ERM problem on the class of linear functions:

$$\mathcal{H} = \{f_\theta(\boldsymbol{x}) = \langle \boldsymbol{\theta}, \boldsymbol{x} \rangle : \boldsymbol{\theta} \in \mathbb{R}^d\}$$

  with the square loss functions:

$$\ell(y, f_\theta(\boldsymbol{x})) = \frac{1}{2}\left(y - f_\theta(\boldsymbol{x})\right)^2$$

- The Bayes predictor and risk are given by:

$$f_\star(\boldsymbol{x}) = \mathbb{E}[y \mid \boldsymbol{x}] \qquad \mathcal{R}_\star = \mathbb{E}\left[\frac{1}{2}(y - \mathbb{E}[y \mid \boldsymbol{x}])^2\right]$$

🤔 Exercise: show this.

# Intercept

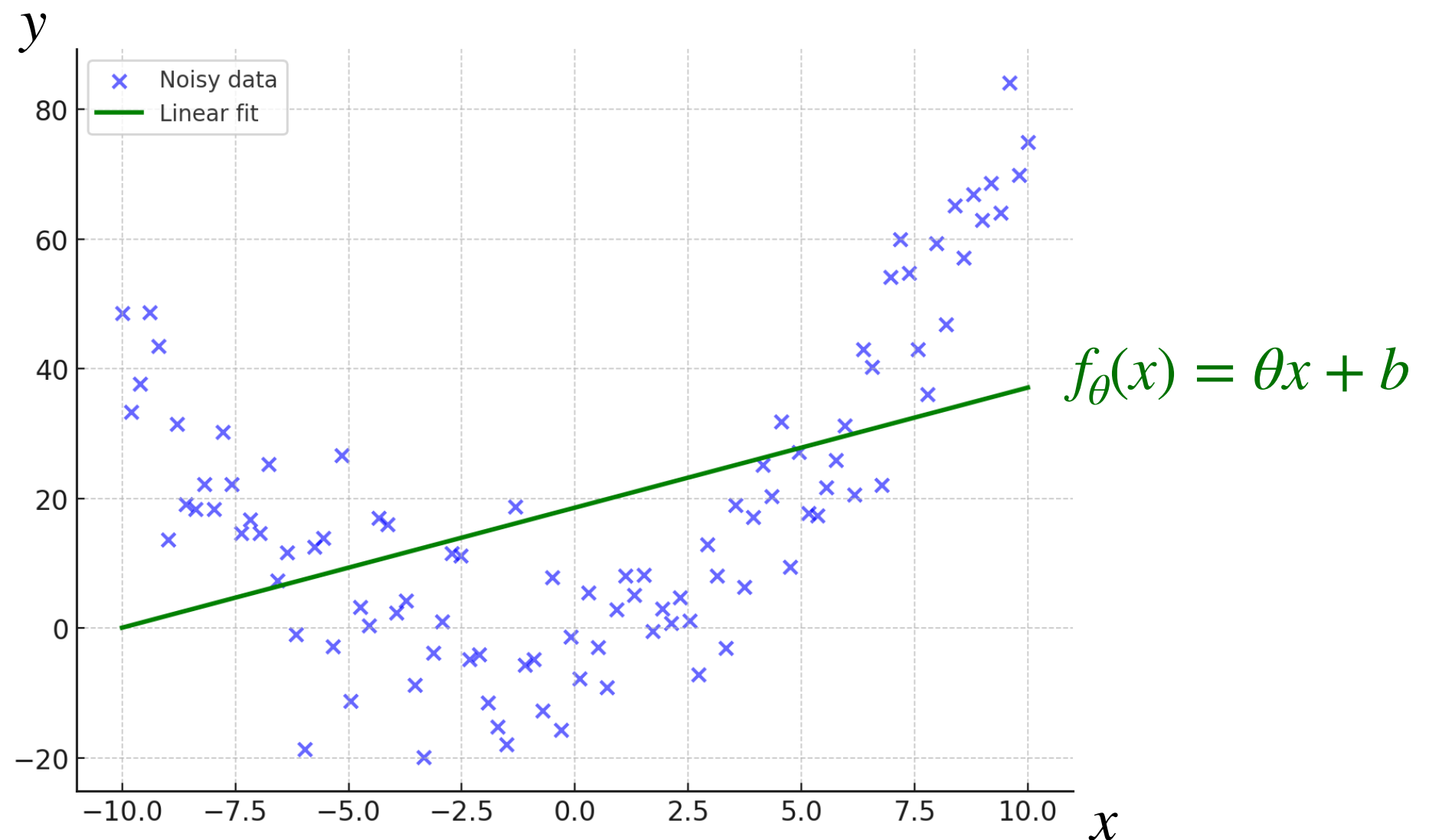- Without loss of generality, can add an intercept:

$$f_\theta(x) = \langle \theta, x \rangle + b$$

By redefining:

$$\tilde{X} = \begin{bmatrix} - & x_1 & - & 1 \\ - & x_2 & - & 1 \\ & \vdots & & \\ - & x_n & - & 1 \end{bmatrix} \in \mathbb{R}^{n \times (d+1)}$$

# Inductive bias of OLS

Remarks:

- <u>Inductive bias</u>: can only fit affine functions of $\boldsymbol{x} \in \mathbb{R}^d$



$$f_\theta(x) = \theta x + b$$

# Convexity of OLS

$$\hat{\mathcal{R}}_n(\boldsymbol{\theta}) := \frac{1}{2n} ||\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta}||_2^2$$

- Gradient: $\quad \nabla_{\boldsymbol{\theta}} \hat{\mathcal{R}}_n = -\frac{1}{n} \boldsymbol{X}^\top (\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta}) \in \mathbb{R}^d$

# Convexity of OLS

$$\hat{\mathscr{R}}_n(\boldsymbol{\theta}) := \frac{1}{2n}||\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta}||_2^2$$

- Gradient:  $\nabla_{\boldsymbol{\theta}}\hat{\mathscr{R}}_n = -\frac{1}{n}\boldsymbol{X}^\top(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta}) \in \mathbb{R}^d$

- Hessian:  $\nabla_{\boldsymbol{\theta}}^2\hat{\mathscr{R}}_n = \frac{1}{n}\boldsymbol{X}^\top\boldsymbol{X} \in \mathbb{R}^{d \times d}$  $(:= \hat{\boldsymbol{\Sigma}}_n)$

# Convexity of OLS

$$\hat{\mathscr{R}}_n(\boldsymbol{\theta}) := \frac{1}{2n} ||\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta}||_2^2$$

- Gradient: $\quad \nabla_{\boldsymbol{\theta}}\hat{\mathscr{R}}_n = -\frac{1}{n}\boldsymbol{X}^\top(\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta}) \in \mathbb{R}^d$

- Hessian: $\quad \nabla^2_{\boldsymbol{\theta}}\hat{\mathscr{R}}_n = \frac{1}{n}\boldsymbol{X}^\top\boldsymbol{X} \in \mathbb{R}^{d\times d} \quad (:= \hat{\boldsymbol{\Sigma}}_n)$

Since $\boldsymbol{X}^\top\boldsymbol{X} \succeq 0$, $\hat{\mathscr{R}}_n$ is convex over $\mathbb{R}^d$. This implies that any minimum of $\hat{\mathscr{R}}_n$ is a global minimum.

# Convexity of OLS

$$\hat{\mathcal{R}}_n(\boldsymbol{\theta}) := \frac{1}{2n} ||\boldsymbol{y} - \boldsymbol{X\theta}||_2^2$$

- **Gradient:** $\quad \nabla_{\boldsymbol{\theta}} \hat{\mathcal{R}}_n = -\frac{1}{n} \boldsymbol{X}^\top (\boldsymbol{y} - \boldsymbol{X\theta}) \in \mathbb{R}^d$

- **Hessian:** $\quad \nabla_{\boldsymbol{\theta}}^2 \hat{\mathcal{R}}_n = \frac{1}{n} \boldsymbol{X}^\top \boldsymbol{X} \in \mathbb{R}^{d \times d} \quad (:= \hat{\boldsymbol{\Sigma}}_n)$

Since $\boldsymbol{X}^\top \boldsymbol{X} \succeq 0$, $\hat{\mathcal{R}}_n$ is convex over $\mathbb{R}^d$. This implies that any minimum of $\hat{\mathcal{R}}_n$ is a global minimum.

For $n \geq d$, $\hat{\mathcal{R}}_n$ is strictly convex if and only if $\mathrm{rank}(\boldsymbol{X}^\top \boldsymbol{X}) = d$. This implies that $\hat{\mathcal{R}}_n$ can have at most one global minimum.

# Closed-form solution

- <u>Gradient:</u> $\nabla_{\boldsymbol{\theta}} \hat{\mathscr{R}}_n = -\dfrac{1}{n} X^{\top} (\boldsymbol{y} - \boldsymbol{X\theta}) \in \mathbb{R}^d$

If it exists, a minima must satisfy:

$$\nabla_{\boldsymbol{\theta}} \hat{\mathscr{R}}_n \overset{!}{=} 0$$

# Closed-form solution

- Gradient: $\quad \nabla_{\boldsymbol{\theta}} \hat{\mathscr{R}}_n = -\dfrac{1}{n} X^{\top}(\boldsymbol{y} - X\boldsymbol{\theta}) \in \mathbb{R}^d$

If it exists, a minima must satisfy:

$$\nabla_{\boldsymbol{\theta}} \hat{\mathscr{R}}_n \overset{!}{=} 0 \qquad \Leftrightarrow \qquad X^{\top} X \boldsymbol{\theta} = X^{\top} \boldsymbol{y}$$

# Closed-form solution

- <u>Gradient:</u>  $\nabla_{\boldsymbol{\theta}} \hat{\mathscr{R}}_n = -\dfrac{1}{n} X^{\top}(\boldsymbol{y} - X\boldsymbol{\theta}) \in \mathbb{R}^d$

If it exists, a minima must satisfy:

$$\nabla_{\boldsymbol{\theta}} \hat{\mathscr{R}}_n \overset{!}{=} 0 \qquad \Leftrightarrow \qquad X^{\top}X\boldsymbol{\theta} = X^{\top}\boldsymbol{y}$$

This is precisely the definition of the pseudo-inverse:

$$\hat{\boldsymbol{\theta}}_{OLS} = X^{+}\boldsymbol{y}$$

# Closed-form solution

- <u>Gradient:</u>   $\nabla_{\boldsymbol{\theta}} \hat{\mathscr{R}}_n = -\dfrac{1}{n} X^\top (\boldsymbol{y} - X\boldsymbol{\theta}) \in \mathbb{R}^d$

If it exists, a minima must satisfy:

$$\nabla_{\boldsymbol{\theta}} \hat{\mathscr{R}}_n \overset{!}{=} 0 \qquad\qquad \Leftrightarrow \qquad\qquad X^\top X\boldsymbol{\theta} = X^\top \boldsymbol{y}$$

This is precisely the definition of the pseudo-inverse:

$$\boxed{\hat{\boldsymbol{\theta}}_{OLS} = X^+ \boldsymbol{y}}$$

If $\mathrm{rank}(X) = \min(n, d)$:    $\hat{\boldsymbol{\theta}}_{OLS} = \begin{cases} (X^\top X)^{-1} X^\top \boldsymbol{y} & \text{if } n \geq d \\ X^\top (XX^\top)^{-1} \boldsymbol{y} & \text{if } n < d \end{cases}$

# Geometrical interpretation

This gives a natural interpretation of the OLS predictor as an orthogonal projection of the labels in the row space of $X$:

$$\hat{\boldsymbol{\theta}}_{OLS} = X^+ \boldsymbol{y} \qquad \Rightarrow \qquad \hat{\boldsymbol{y}}_{OLS} = X\hat{\boldsymbol{\theta}}_{OLS} = XX^+ \boldsymbol{y}$$



$\boldsymbol{y}$

$\hat{\boldsymbol{y}} = \boldsymbol{P}_{X^\top} \boldsymbol{y}$

$0$

$\text{im}(X) \subset \mathbb{R}^n$

$$\min_{z \in \text{im}(X)} ||\boldsymbol{y} - z||_2^2$$

# Two scenarios

From now on, let's assume $X$ is full-rank.

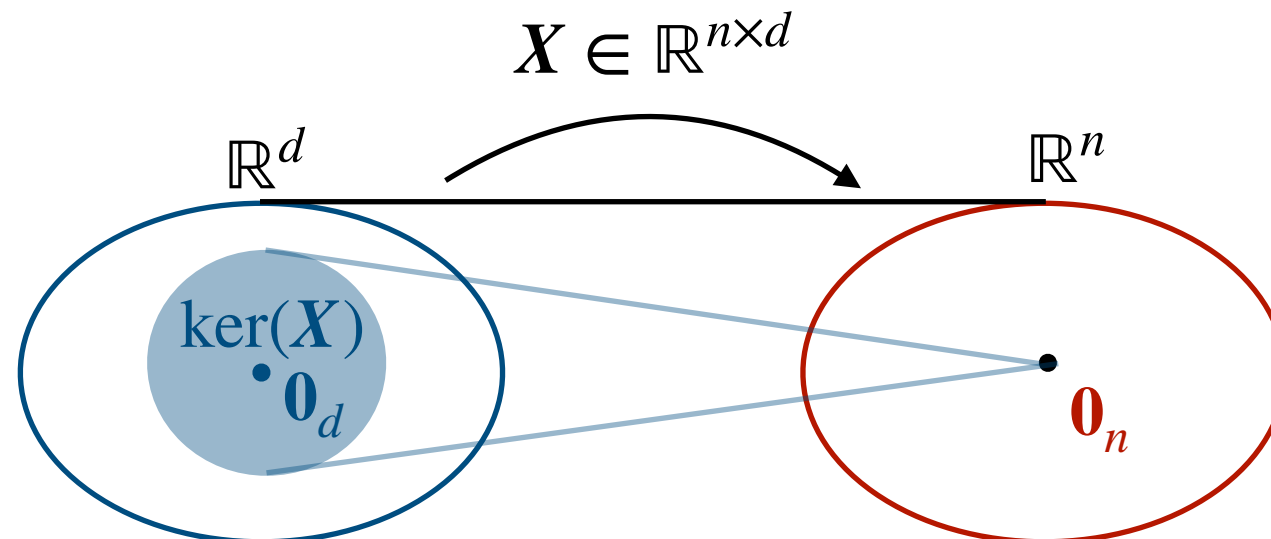- For $n \geq d$ : more equations than variables. $X\boldsymbol{\theta} = \boldsymbol{y}$ admits an unique solution.

# Two scenarios

From now on, let's assume $X$ is full-rank.

- For $n \geq d$ : more equations than variables. $X\theta = y$ admits an unique solution.

$$X \in \mathbb{R}^{n \times d}$$



$\mathbb{R}^d$        $\mathbb{R}^n$

$\mathbf{0}_d$

$\mathrm{im}(X)$

$\mathbf{0}_n$

- For $n < d$ : more variables than equations. $X\theta = y$ admits several solutions.

$$X \in \mathbb{R}^{n \times d}$$



$\mathbb{R}^d$        $\mathbb{R}^n$

$\ker(X)$

$\mathbf{0}_d$

$\mathbf{0}_n$

# OLS as least norm solution

Assume $\text{rank}(X) = n < d$. Then, OLS admits the following interpretation as the minimum $\ell_2$-norm solution:

$$\hat{\boldsymbol{\theta}}_{OLS} = \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\text{argmin}} \; ||\boldsymbol{\theta}||_2$$

$$\text{subject to} \quad X\boldsymbol{\theta} = y$$

# OLS as least norm solution

Assume $\text{rank}(X) = n < d$. Then, OLS admits the following interpretation as the minimum $\ell_2$-norm solution:

$$\hat{\boldsymbol{\theta}}_{OLS} = \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\text{argmin}} \ ||\boldsymbol{\theta}||_2$$

$$\text{subject to} \quad X\boldsymbol{\theta} = y$$

_Proof:_ Let $\hat{\boldsymbol{\theta}} \in \mathbb{R}^d$ denote a different solution from $\hat{\boldsymbol{\theta}}_{OLS}$.

# OLS as least norm solution

Assume $\text{rank}(X) = n < d$. Then, OLS admits the following interpretation as the minimum $\ell_2$-norm solution:

$$\hat{\boldsymbol{\theta}}_{OLS} = \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\text{argmin}} \; ||\boldsymbol{\theta}||_2$$

$$\text{subject to} \quad X\boldsymbol{\theta} = \boldsymbol{y}$$

_Proof:_ Let $\hat{\boldsymbol{\theta}} \in \mathbb{R}^d$ denote a different solution from $\hat{\boldsymbol{\theta}}_{OLS}$.

Then: $\langle \hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}, \hat{\boldsymbol{\theta}}_{OLS} \rangle = \langle \hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}, X^\top (XX^\top)^{-1}\boldsymbol{y} \rangle$

# OLS as least norm solution

Assume $\mathrm{rank}(X) = n < d$. Then, OLS admits the following interpretation as the minimum $\ell_2$-norm solution:

$$\hat{\boldsymbol{\theta}}_{OLS} = \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\mathrm{argmin}} \ ||\boldsymbol{\theta}||_2$$

$$\text{subject to} \quad X\boldsymbol{\theta} = y$$

*Proof:* Let $\hat{\boldsymbol{\theta}} \in \mathbb{R}^d$ denote a different solution from $\hat{\boldsymbol{\theta}}_{OLS}$.

Then: $\langle \hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}, \hat{\boldsymbol{\theta}}_{OLS} \rangle = \langle \hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}, X^\top(XX^\top)^{-1}y \rangle$

$$= \langle X(\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}), (XX^\top)^{-1}y \rangle$$

# OLS as least norm solution

Assume $\text{rank}(X) = n < d$. Then, OLS admits the following interpretation as the minimum $\ell_2$-norm solution:

$$\hat{\boldsymbol{\theta}}_{OLS} = \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\text{argmin}} \ ||\boldsymbol{\theta}||_2$$

$$\text{subject to} \quad X\boldsymbol{\theta} = y$$

_Proof:_ Let $\hat{\boldsymbol{\theta}} \in \mathbb{R}^d$ denote a different solution from $\hat{\boldsymbol{\theta}}_{OLS}$.

Then: $\langle \hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}, \hat{\boldsymbol{\theta}}_{OLS} \rangle = \langle \hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}, X^\top(XX^\top)^{-1}y \rangle$

$$= \langle X(\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}), (XX^\top)^{-1}y \rangle$$

$$= 0$$

# OLS as least norm solution

Assume $\text{rank}(X) = n < d$. Then, OLS admits the following interpretation as the minimum $\ell_2$-norm solution:

$$\hat{\boldsymbol{\theta}}_{OLS} = \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\text{argmin}} \; ||\boldsymbol{\theta}||_2$$

$$\text{subject to} \quad X\boldsymbol{\theta} = \boldsymbol{y}$$

_Proof:_  Let $\hat{\boldsymbol{\theta}} \in \mathbb{R}^d$ denote a different solution from $\hat{\boldsymbol{\theta}}_{OLS}$.

Then: $\quad \langle \hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}, \hat{\boldsymbol{\theta}}_{OLS} \rangle = \langle \hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}, X^\top (XX^\top)^{-1} \boldsymbol{y} \rangle$

$$= \langle X(\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}), (XX^\top)^{-1} \boldsymbol{y} \rangle$$

$$= 0$$

Therefore $\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS} \perp \hat{\boldsymbol{\theta}}_{OLS}$.

# OLS as least norm solution

Assume $\text{rank}(X) = n < d$. Then, OLS admits the following interpretation as the minimum $\ell_2$-norm solution:

$$\hat{\boldsymbol{\theta}}_{OLS} = \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\text{argmin}} \ ||\boldsymbol{\theta}||_2$$

$$\text{subject to} \quad X\boldsymbol{\theta} = y$$

*Proof:* Let $\hat{\boldsymbol{\theta}} \in \mathbb{R}^d$ denote a different solution from $\hat{\boldsymbol{\theta}}_{OLS}$.

Then: $\langle \hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}, \hat{\boldsymbol{\theta}}_{OLS} \rangle = \langle \hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}, X^\top(XX^\top)^{-1}y \rangle$

$$= \langle X(\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}), (XX^\top)^{-1}y \rangle$$

$$= 0$$

Therefore $\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS} \perp \hat{\boldsymbol{\theta}}_{OLS}$. Hence:

$$||\hat{\boldsymbol{\theta}}||_2^2 = ||\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS} + \hat{\boldsymbol{\theta}}_{OLS}||_2^2$$

# OLS as least norm solution

Assume $\mathrm{rank}(X) = n < d$. Then, OLS admits the following interpretation as the minimum $\ell_2$-norm solution:

$$\hat{\boldsymbol{\theta}}_{OLS} = \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\mathrm{argmin}} \ ||\boldsymbol{\theta}||_2$$

$$\text{subject to} \quad X\boldsymbol{\theta} = y$$

*Proof:*  Let $\hat{\boldsymbol{\theta}} \in \mathbb{R}^d$ denote a different solution from $\hat{\boldsymbol{\theta}}_{OLS}$.

Then:  $\langle \hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}, \hat{\boldsymbol{\theta}}_{OLS} \rangle = \langle \hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}, X^\top(XX^\top)^{-1}y \rangle$

$$= \langle X(\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}), (XX^\top)^{-1}y \rangle$$

$$= 0$$

Therefore $\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS} \perp \hat{\boldsymbol{\theta}}_{OLS}$. Hence:

$$||\hat{\boldsymbol{\theta}}||_2^2 = ||\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS} + \hat{\boldsymbol{\theta}}_{OLS}||_2^2 = ||\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}||_2^2 + ||\hat{\boldsymbol{\theta}}_{OLS}||_2^2$$

# OLS as least norm solution

Assume $\text{rank}(X) = n < d$. Then, OLS admits the following interpretation as the minimum $\ell_2$-norm solution:

$$\hat{\boldsymbol{\theta}}_{OLS} = \underset{\boldsymbol{\theta} \in \mathbb{R}^d}{\text{argmin}} \ ||\boldsymbol{\theta}||_2$$

$$\text{subject to} \quad X\boldsymbol{\theta} = y$$

_Proof:_   Let $\hat{\boldsymbol{\theta}} \in \mathbb{R}^d$ denote a different solution from $\hat{\boldsymbol{\theta}}_{OLS}$.

Then:   $\langle \hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}, \hat{\boldsymbol{\theta}}_{OLS} \rangle = \langle \hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}, X^\top (XX^\top)^{-1} y \rangle$

$$= \langle X(\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}), (XX^\top)^{-1} y \rangle$$

$$= 0$$

Therefore $\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS} \perp \hat{\boldsymbol{\theta}}_{OLS}$. Hence:

$$||\hat{\boldsymbol{\theta}}||_2^2 = ||\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS} + \hat{\boldsymbol{\theta}}_{OLS}||_2^2 = ||\hat{\boldsymbol{\theta}} - \hat{\boldsymbol{\theta}}_{OLS}||_2^2 + ||\hat{\boldsymbol{\theta}}_{OLS}||_2^2 \geq ||\hat{\boldsymbol{\theta}}_{OLS}||_2^2$$