# Statistical Learning II

Lecture 5 - Least squares (continued)
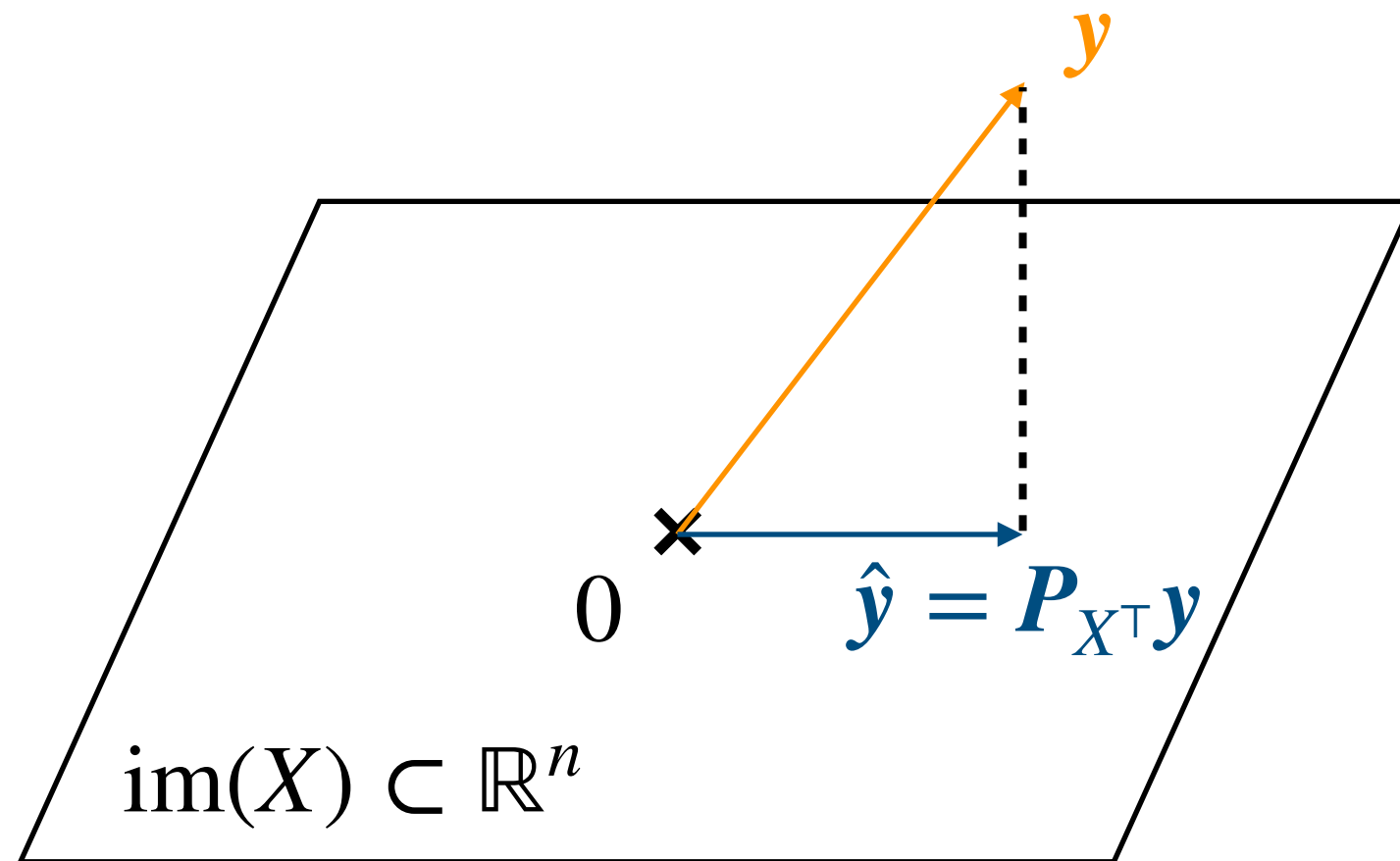
**Bruno Loureiro**
@ CSD, DI-ENS & CNRS

brloureiro@gmail.com

*DL3 IASO, Université Paris Dauphine-PSL*
*09.10.2024*

# Geometrical interpretation

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

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



$\boldsymbol{y}$

$0$

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

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

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

# Statistical analysis of OLS

Fixed-design analysis

# Assumptions

We now assume the following data generative model:

$$y_i = \langle \boldsymbol{\theta}_\star, \boldsymbol{x}_i \rangle + \varepsilon_i$$

With:
- Fixed $\boldsymbol{\theta}_\star \in \mathbb{R}^d$ and $\boldsymbol{x}_i \in \mathbb{R}^d$    "fixed design"
- $\mathbb{E}[\varepsilon_i] = 0$ and $\mathbb{E}[\varepsilon_i^2] = \sigma^2 < \infty$

# Assumptions

We now assume the following data generative model:

$$y_i = \langle \boldsymbol{\theta}_\star, \boldsymbol{x}_i \rangle + \varepsilon_i$$

With:
- Fixed $\boldsymbol{\theta}_\star \in \mathbb{R}^d$ and $\boldsymbol{x}_i \in \mathbb{R}^d$    "fixed design"
- $\mathbb{E}[\varepsilon_i] = 0$ and $\mathbb{E}[\varepsilon_i^2] = \sigma^2 < \infty$

Remarks:
- The Bayes predictor and error are given by

$$f_\star(\boldsymbol{x}) = \mathbb{E}[y \mid X = x] = \langle \boldsymbol{\theta}_\star, \boldsymbol{x} \rangle \qquad \mathscr{R}_\star = \mathscr{R}(\boldsymbol{\theta}_\star) = \sigma^2$$

# Assumptions

We now assume the following data generative model:

$$y_i = \langle \boldsymbol{\theta}_\star, \boldsymbol{x}_i \rangle + \varepsilon_i$$

With: 
- Fixed $\boldsymbol{\theta}_\star \in \mathbb{R}^d$ and $\boldsymbol{x}_i \in \mathbb{R}^d$    "fixed design"
- $\mathbb{E}[\varepsilon_i] = 0$ and $\mathbb{E}[\varepsilon_i^2] = \sigma^2 < \infty$

Remarks:
- The Bayes predictor and error are given by

$$f_\star(\boldsymbol{x}) = \mathbb{E}[y \,|\, \boldsymbol{X} = \boldsymbol{x}] = \langle \boldsymbol{\theta}_\star, \boldsymbol{x} \rangle \qquad \mathscr{R}_\star = \mathscr{R}(\boldsymbol{\theta}_\star) = \sigma^2$$

- In particular

$$f_\star \in \mathscr{H} = \{ f(\boldsymbol{x}) = \langle \boldsymbol{\theta}, \boldsymbol{x} \rangle : \boldsymbol{\theta} \in \mathbb{R}^d \}$$
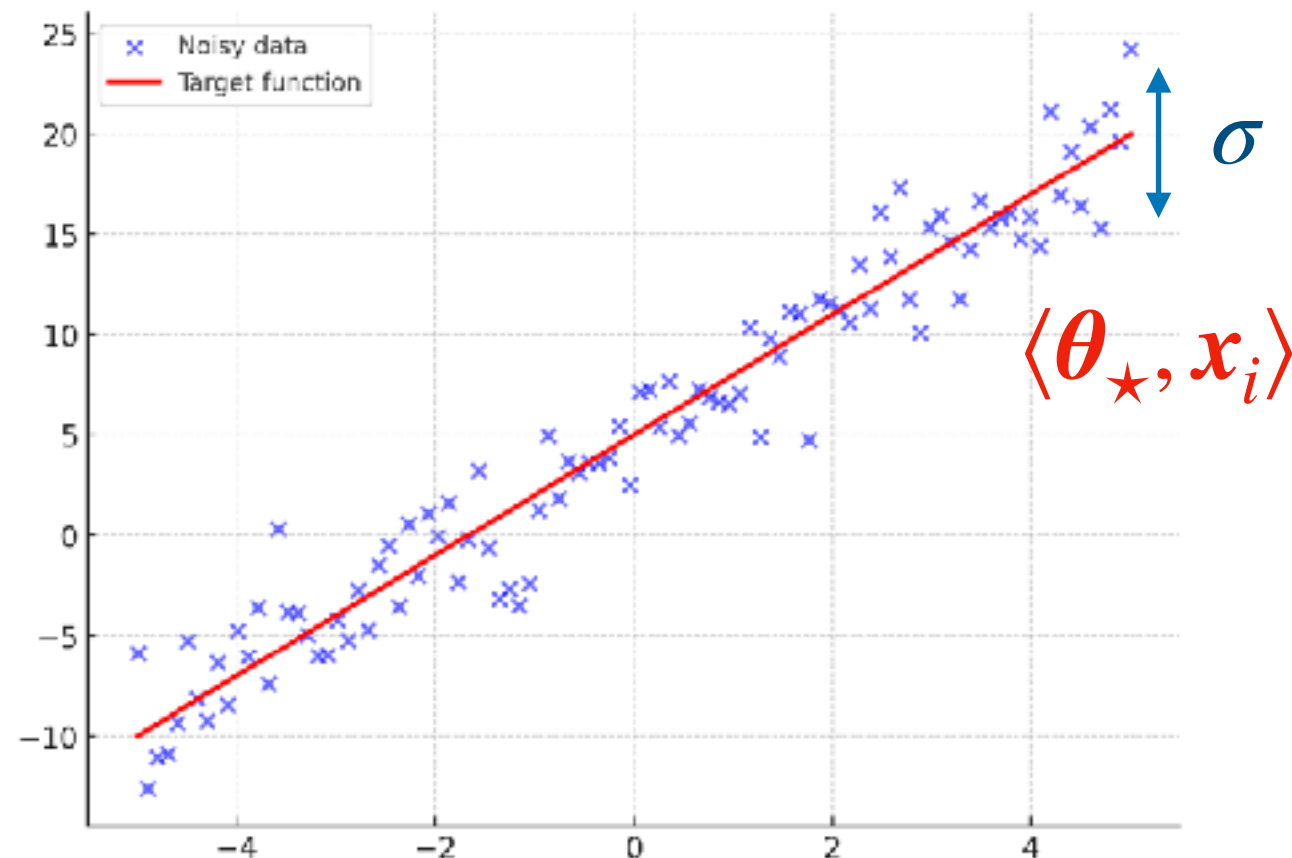
"Well-specified setting"

# Assumptions

We now assume the following data generative model:

$$y_i = \langle \boldsymbol{\theta}_\star, \boldsymbol{x}_i \rangle + \varepsilon_i$$

With:
- Fixed $\boldsymbol{\theta}_\star \in \mathbb{R}^d$ and $\boldsymbol{x}_i \in \mathbb{R}^d$    "fixed design"
- $\mathbb{E}[\varepsilon_i] = 0$ and $\mathbb{E}[\varepsilon_i^2] = \sigma^2 < \infty$
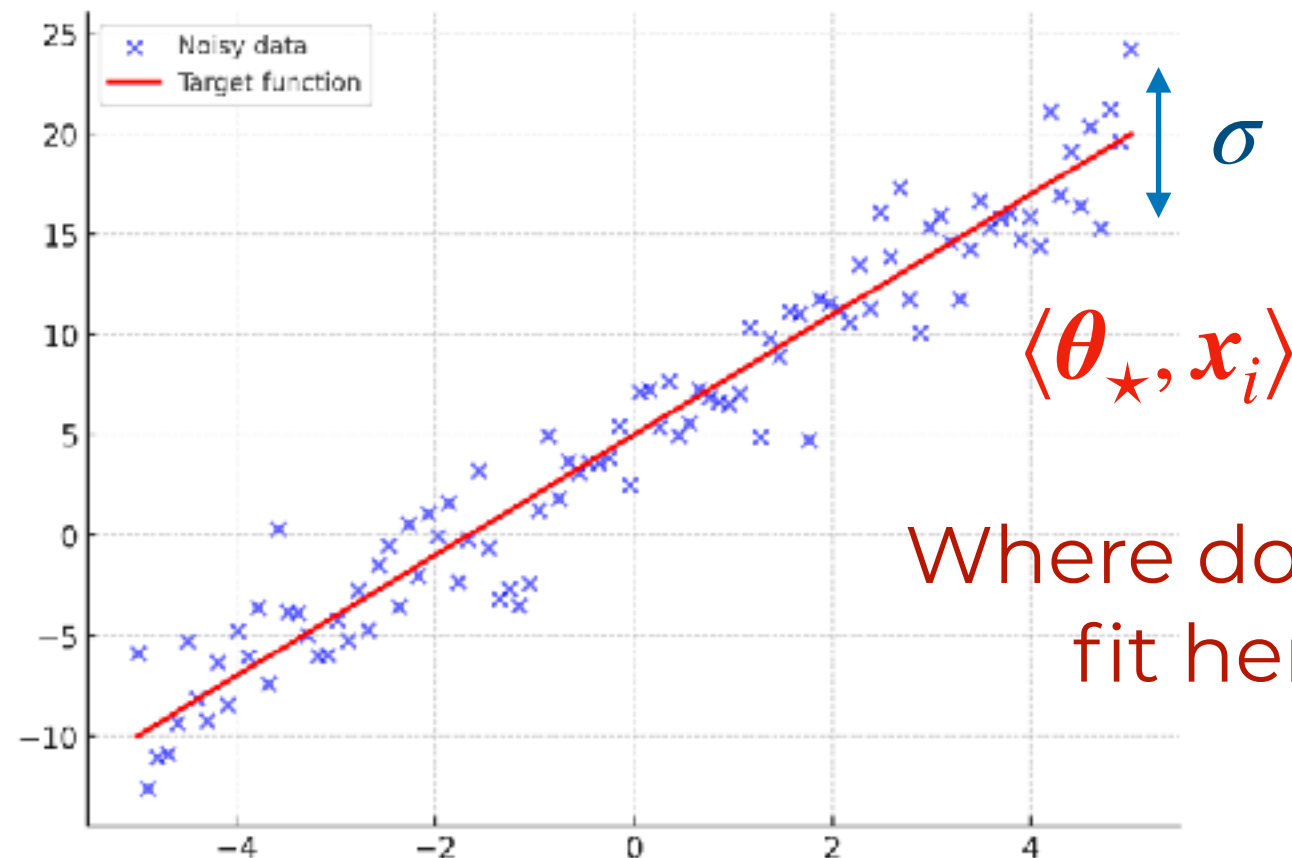
# Assumptions

We now assume the following data generative model:

$$y_i = \langle \boldsymbol{\theta}_\star, \boldsymbol{x}_i \rangle + \varepsilon_i$$

With:
- Fixed $\boldsymbol{\theta}_\star \in \mathbb{R}^d$ and $\boldsymbol{x}_i \in \mathbb{R}^d$    "fixed design"
- $\mathbb{E}[\varepsilon_i] = 0$ and $\mathbb{E}[\varepsilon_i^2] = \sigma^2 < \infty$



$\sigma$

$\langle \boldsymbol{\theta}_\star, \boldsymbol{x}_i \rangle$

Where does OLS fit here?

# Decomposition of OLS

Given a batch of data sampled from this model:

$$y = X\theta_\star + \varepsilon \in \mathbb{R}^n$$

Our goal is to understand the statistical properties of OLS.
For simplicity, assume that $\mathrm{rank}(X) = d$ $(n > d)$:

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

# Decomposition of OLS

Given a batch of data sampled from this model:

$$y = X\theta_\star + \varepsilon \in \mathbb{R}^n$$

Our goal is to understand the statistical properties of OLS. For simplicity, assume that $\text{rank}(X) = d$ $(n > d)$:

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

# Decomposition of OLS

Given a batch of data sampled from this model:

$$\textcolor{red}{y = X\theta_\star + \varepsilon \in \mathbb{R}^n}$$

Our goal is to understand the statistical properties of OLS.
For simplicity, $\textcolor{red}{\text{assume that } \mathrm{rank}(X) = d \ (n > d)}$:

$$\hat{\theta}_{OLS}(X, y) = (X^\top X)^{-1} X^\top \textcolor{red}{y} = (X^\top X)^{-1} X^\top \textcolor{red}{(X\theta_\star + \varepsilon)}$$

# Decomposition of OLS

Given a batch of data sampled from this model:

$$y = X\theta_\star + \varepsilon \in \mathbb{R}^n$$

Our goal is to understand the statistical properties of OLS.
For simplicity, assume that $\text{rank}(X) = d$ $(n > d)$:

$$\hat{\theta}_{OLS}(X, y) = (X^\top X)^{-1} X^\top y = (X^\top X)^{-1} X^\top (X\theta_\star + \varepsilon)$$
$$= (X^\top X)^{-1} X^\top X\theta_\star + (X^\top X)^{-1} X^\top \varepsilon$$

# Decomposition of OLS

Given a batch of data sampled from this model:

$$y = X\theta_\star + \varepsilon \in \mathbb{R}^n$$

Our goal is to understand the statistical properties of OLS. For simplicity, assume that $\mathrm{rank}(X) = d$ $(n > d)$:

$$\hat{\theta}_{OLS}(X, y) = (X^\top X)^{-1} X^\top y = (X^\top X)^{-1} X^\top (X\theta_\star + \varepsilon)$$

$$= (X^\top X)^{-1} X^\top X \theta_\star + (X^\top X)^{-1} X^\top \varepsilon$$

$$= \theta_\star + \frac{1}{n} \hat{\Sigma}_n^{-1} X^\top \varepsilon$$

Where we have defined $\hat{\Sigma}_n = \frac{1}{n} X^\top X \in \mathbb{R}^{d \times d}$ (Empirical covariance)

# Decomposition of OLS

$$\hat{\boldsymbol{\theta}}_{OLS}(\boldsymbol{X}, \boldsymbol{y}) = \boldsymbol{\theta}_{\star} + \frac{1}{n}\hat{\boldsymbol{\Sigma}}_n^{-1}\boldsymbol{X}^{\top}\boldsymbol{\varepsilon}$$

"signal"          "noise"

# Decomposition of OLS

$$\hat{\boldsymbol{\theta}}_{OLS}(\boldsymbol{X}, \boldsymbol{y}) = \boldsymbol{\theta}_\star + \frac{1}{n}\hat{\boldsymbol{\Sigma}}_n^{-1}\boldsymbol{X}^\top\boldsymbol{\varepsilon}$$

"signal"        "noise"

In particular:

- Bias:    $\mathbb{E}_{\boldsymbol{\varepsilon}}\left[\hat{\boldsymbol{\theta}}_{OLS}(\boldsymbol{X}, \boldsymbol{y})\right] = \boldsymbol{\theta}_\star$     "Unbiased"

# Decomposition of OLS

$$\hat{\boldsymbol{\theta}}_{OLS}(\boldsymbol{X}, \boldsymbol{y}) = \boldsymbol{\theta}_{\star} + \frac{1}{n}\hat{\boldsymbol{\Sigma}}_n^{-1}\boldsymbol{X}^{\top}\boldsymbol{\varepsilon}$$

"signal"   "noise"

In particular:

- Bias: $\quad \mathbb{E}_{\boldsymbol{\varepsilon}}\left[\hat{\boldsymbol{\theta}}_{OLS}(\boldsymbol{X}, \boldsymbol{y})\right] = \boldsymbol{\theta}_{\star}$   "Unbiased"

- Variance: $\quad \mathrm{Var}_{\boldsymbol{\varepsilon}}\left[\hat{\boldsymbol{\theta}}_{OLS}(\boldsymbol{X}, \boldsymbol{y})\right] = \frac{\sigma^2}{n}\hat{\boldsymbol{\Sigma}}_n^{-1}$

# Decomposition of OLS

$$\hat{\boldsymbol{\theta}}_{OLS}(\boldsymbol{X}, \boldsymbol{y}) = \boldsymbol{\theta}_{\star} + \frac{1}{n}\hat{\boldsymbol{\Sigma}}_n^{-1}\boldsymbol{X}^{\top}\boldsymbol{\varepsilon}$$

"signal"   "noise"

In particular:

- Bias:   $\mathbb{E}_{\boldsymbol{\varepsilon}}\left[\hat{\boldsymbol{\theta}}_{OLS}(\boldsymbol{X}, \boldsymbol{y})\right] = \boldsymbol{\theta}_{\star}$   "Unbiased"

- Variance:   $\mathrm{Var}_{\boldsymbol{\varepsilon}}\left[\hat{\boldsymbol{\theta}}_{OLS}(\boldsymbol{X}, \boldsymbol{y})\right] = \frac{\sigma^2}{n}\hat{\boldsymbol{\Sigma}}_n^{-1}$

Hence, informally:

$$\hat{\boldsymbol{\theta}}_{OLS} \to \boldsymbol{\theta}_{\star} \qquad \text{as } n \to \infty \qquad \text{"Consistency"}$$

# Risk of OLS

We now look at the population risk.

# Risk of OLS

We now look at the population risk. First, note that for fixed $\boldsymbol{\theta} \in \mathbb{R}^d$:

$$\mathcal{R}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{y}} \left[ \frac{1}{n} ||\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta}||_2^2 \right]$$

# Risk of OLS

We now look at the population risk. First, note that for fixed $\boldsymbol{\theta} \in \mathbb{R}^d$:

$$\mathscr{R}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{y}} \left[ \frac{1}{n} ||\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta}||_2^2 \right] = \mathbb{E}_{\boldsymbol{\varepsilon}} \left[ \frac{1}{n} ||\boldsymbol{X}\boldsymbol{\theta}_\star + \boldsymbol{\varepsilon} - \boldsymbol{X}\boldsymbol{\theta}||_2^2 \right]$$

# Risk of OLS

We now look at the population risk. First, note that for fixed $\boldsymbol{\theta} \in \mathbb{R}^d$:

$$\mathcal{R}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{y}} \left[ \frac{1}{n} ||\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta}||_2^2 \right] = \mathbb{E}_{\boldsymbol{\varepsilon}} \left[ \frac{1}{n} ||\boldsymbol{X}\boldsymbol{\theta}_\star + \boldsymbol{\varepsilon} - \boldsymbol{X}\boldsymbol{\theta}||_2^2 \right]$$

$$= (\boldsymbol{\theta}_\star - \boldsymbol{\theta})^\top \hat{\boldsymbol{\Sigma}}_n (\boldsymbol{\theta}_\star - \boldsymbol{\theta}) + \sigma^2$$

# Risk of OLS

We now look at the population risk. First, note that for fixed $\boldsymbol{\theta} \in \mathbb{R}^d$:

$$\mathcal{R}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{y}}\left[\frac{1}{n}||\boldsymbol{y} - \boldsymbol{X\theta}||_2^2\right] = \mathbb{E}_{\boldsymbol{\varepsilon}}\left[\frac{1}{n}||\boldsymbol{X\theta_\star} + \boldsymbol{\varepsilon} - \boldsymbol{X\theta}||_2^2\right]$$

$$= (\boldsymbol{\theta_\star} - \boldsymbol{\theta})^\top \hat{\boldsymbol{\Sigma}}_n (\boldsymbol{\theta_\star} - \boldsymbol{\theta}) + \sigma^2$$

In the fixed design setting, $\boldsymbol{X} \in \mathbb{R}^{n \times d}$ is the same for training and testing.

# Risk of OLS

We now look at the population risk. First, note that for fixed $\boldsymbol{\theta} \in \mathbb{R}^d$:

$$\mathcal{R}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{y}}\left[\frac{1}{n}||\boldsymbol{y} - \boldsymbol{X}\boldsymbol{\theta}||_2^2\right] = \mathbb{E}_{\boldsymbol{\varepsilon}}\left[\frac{1}{n}||\boldsymbol{X}\boldsymbol{\theta}_\star + \boldsymbol{\varepsilon} - \boldsymbol{X}\boldsymbol{\theta}||_2^2\right]$$

$$= (\boldsymbol{\theta}_\star - \boldsymbol{\theta})^\top \hat{\boldsymbol{\Sigma}}_n (\boldsymbol{\theta}_\star - \boldsymbol{\theta}) + \sigma^2$$

⚠️ In the fixed design setting, $\boldsymbol{X} \in \mathbb{R}^{n \times d}$ is the same for training and testing.

Therefore, for the OLS $\hat{\boldsymbol{\theta}}_{OLS}(\boldsymbol{X}, \boldsymbol{y}) = \boldsymbol{\theta}_\star + \frac{1}{n}\hat{\boldsymbol{\Sigma}}_n^{-1}\boldsymbol{X}^\top\boldsymbol{\varepsilon}$:

$$\mathcal{R}(\hat{\boldsymbol{\theta}}_{OLS}) - \sigma^2 = \frac{1}{n^2}\boldsymbol{\varepsilon}^\top \boldsymbol{X}\hat{\boldsymbol{\Sigma}}_n^{-1}\boldsymbol{X}^\top\boldsymbol{\varepsilon}$$

# Risk of OLS

We now look at the population risk. First, note that for fixed $\boldsymbol{\theta} \in \mathbb{R}^d$:

$$\mathcal{R}(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{y}}\left[\frac{1}{n}||\boldsymbol{y} - \boldsymbol{X\theta}||_2^2\right] = \mathbb{E}_{\boldsymbol{\varepsilon}}\left[\frac{1}{n}||\boldsymbol{X\theta_\star} + \boldsymbol{\varepsilon} - \boldsymbol{X\theta}||_2^2\right]$$

$$= (\boldsymbol{\theta_\star} - \boldsymbol{\theta})^\top \hat{\boldsymbol{\Sigma}}_n (\boldsymbol{\theta_\star} - \boldsymbol{\theta}) + \sigma^2$$

⚠️ In the fixed design setting, $X \in \mathbb{R}^{n \times d}$ is the same for training and testing.

Therefore, for the OLS $\hat{\boldsymbol{\theta}}_{OLS}(X, \boldsymbol{y}) = \boldsymbol{\theta_\star} + \frac{1}{n}\hat{\boldsymbol{\Sigma}}_n^{-1}X^\top\boldsymbol{\varepsilon}$:

$$\mathcal{R}(\hat{\boldsymbol{\theta}}_{OLS}) - \sigma^2 = \frac{1}{n^2}\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon}$$

⚠️ This is a random variable since $\hat{\boldsymbol{\theta}}_{OLS}$ is random!

# Risk of OLS

First, note we can rewrite:

$$\mathscr{R}(\hat{\boldsymbol{\theta}}_{OLS}) - \sigma^2 = \frac{1}{n^2} \boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon}$$

# Risk of OLS

First, note we can rewrite:

$$\mathcal{R}(\hat{\boldsymbol{\theta}}_{OLS}) - \sigma^2 = \frac{1}{n^2}\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon} = \frac{1}{n^2}\text{Tr}\left[\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon}\right]$$

# Risk of OLS

First, note we can rewrite:

$$\mathscr{R}(\hat{\boldsymbol{\theta}}_{OLS}) - \sigma^2 = \frac{1}{n^2}\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon} = \frac{1}{n^2}\text{Tr}\left[\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon}\right]$$

$$= \frac{1}{n^2}\text{Tr}\left[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top\right]$$

# Risk of OLS

First, note we can rewrite:

$$\mathcal{R}(\hat{\boldsymbol{\theta}}_{OLS}) - \sigma^2 = \frac{1}{n^2}\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon} = \frac{1}{n^2}\mathrm{Tr}\left[\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon}\right]$$

$$= \frac{1}{n^2}\mathrm{Tr}\left[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top\right]$$

Now taking the expectation:

# Risk of OLS

First, note we can rewrite:

$$\mathcal{R}(\hat{\boldsymbol{\theta}}_{OLS}) - \sigma^2 = \frac{1}{n^2}\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon} = \frac{1}{n^2}\mathrm{Tr}\left[\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon}\right]$$

$$= \frac{1}{n^2}\mathrm{Tr}\left[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top\right]$$

Now taking the expectation:

$$\mathbb{E}_{\boldsymbol{\varepsilon}}\left[\mathcal{R}(\hat{\boldsymbol{\theta}}_{OLS})\right] - \sigma^2 = \frac{1}{n^2}\mathrm{Tr}\left[\mathbb{E}_{\boldsymbol{\varepsilon}}\left[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top\right] X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top\right]$$

# Risk of OLS

First, note we can rewrite:

$$\mathscr{R}(\hat{\boldsymbol{\theta}}_{OLS}) - \sigma^2 = \frac{1}{n^2}\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon} = \frac{1}{n^2}\mathrm{Tr}\left[\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon}\right]$$

$$= \frac{1}{n^2}\mathrm{Tr}\left[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top\right]$$

Now taking the expectation:

$$\mathbb{E}_{\boldsymbol{\varepsilon}}\left[\mathscr{R}(\hat{\boldsymbol{\theta}}_{OLS})\right] - \sigma^2 = \frac{1}{n^2}\mathrm{Tr}\left[\mathbb{E}_{\boldsymbol{\varepsilon}}\left[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top\right] X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top\right]$$

$$= \frac{\sigma^2}{n^2}\mathrm{Tr}\left[X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top\right]$$

# Risk of OLS

First, note we can rewrite:

$$\mathscr{R}(\hat{\boldsymbol{\theta}}_{OLS}) - \sigma^2 = \frac{1}{n^2}\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon} = \frac{1}{n^2}\mathrm{Tr}\left[\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon}\right]$$

$$= \frac{1}{n^2}\mathrm{Tr}\left[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top\right]$$

Now taking the expectation:

$$\mathbb{E}_{\boldsymbol{\varepsilon}}\left[\mathscr{R}(\hat{\boldsymbol{\theta}}_{OLS})\right] - \sigma^2 = \frac{1}{n^2}\mathrm{Tr}\left[\mathbb{E}_{\boldsymbol{\varepsilon}}\left[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top\right] X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top\right]$$

$$= \frac{\sigma^2}{n^2}\mathrm{Tr}\left[X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top\right] = \frac{\sigma^2}{n^2}\mathrm{Tr}\left[X^\top X \hat{\boldsymbol{\Sigma}}_n^{-1}\right]$$

# Risk of OLS

First, note we can rewrite:

$$\mathscr{R}(\hat{\boldsymbol{\theta}}_{OLS}) - \sigma^2 = \frac{1}{n^2}\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon} = \frac{1}{n^2}\text{Tr}\left[\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon}\right]$$

$$= \frac{1}{n^2}\text{Tr}\left[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top\right]$$

Now taking the expectation:

$$\mathbb{E}_{\boldsymbol{\varepsilon}}\left[\mathscr{R}(\hat{\boldsymbol{\theta}}_{OLS})\right] - \sigma^2 = \frac{1}{n^2}\text{Tr}\left[\mathbb{E}_{\boldsymbol{\varepsilon}}\left[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top\right] X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top\right]$$

$$= \frac{\sigma^2}{n^2}\text{Tr}\left[X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top\right] = \frac{\sigma^2}{n^2}\text{Tr}\left[X^\top X \hat{\boldsymbol{\Sigma}}_n^{-1}\right]$$

$$= \frac{\sigma^2}{n}\text{Tr}[I_d]$$

# Risk of OLS

First, note we can rewrite:

$$\mathcal{R}(\hat{\boldsymbol{\theta}}_{OLS}) - \sigma^2 = \frac{1}{n^2}\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon} = \frac{1}{n^2}\text{Tr}\left[\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top \boldsymbol{\varepsilon}\right]$$

$$= \frac{1}{n^2}\text{Tr}\left[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top\right]$$

Now taking the expectation:

$$\mathbb{E}_{\boldsymbol{\varepsilon}}\left[\mathcal{R}(\hat{\boldsymbol{\theta}}_{OLS})\right] - \sigma^2 = \frac{1}{n^2}\text{Tr}\left[\mathbb{E}_{\boldsymbol{\varepsilon}}\left[\boldsymbol{\varepsilon}\boldsymbol{\varepsilon}^\top\right] X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top\right]$$

$$= \frac{\sigma^2}{n^2}\text{Tr}\left[X \hat{\boldsymbol{\Sigma}}_n^{-1} X^\top\right] = \frac{\sigma^2}{n^2}\text{Tr}\left[X^\top X \hat{\boldsymbol{\Sigma}}_n^{-1}\right]$$

$$= \frac{\sigma^2}{n}\text{Tr}[I_d] = \sigma^2 \frac{d}{n}$$

# Risk of OLS

Therefore, we have the following final result for the excess risk of OLS

$$\mathbb{E}_{\boldsymbol{\varepsilon}}\left[\mathscr{R}(\hat{\boldsymbol{\theta}}_{OLS})\right] - \sigma^2 = \sigma^2\frac{d}{n}$$

# Risk of OLS

Therefore, we have the following final result for the excess risk of OLS

$$\mathbb{E}_{\boldsymbol{\varepsilon}}\left[\mathscr{R}(\hat{\boldsymbol{\theta}}_{OLS})\right] - \sigma^2 = \sigma^2\frac{d}{n}$$

Remarks:

- Excess risk is proportional to the noise level $\mathbb{E}[\varepsilon^2] = \sigma^2$.

- Excess risk is proportional to the data dimension.

- To achieve excess risk $\Delta\mathscr{R} < \delta$, need:

$$n > \frac{\sigma^2 d}{\delta}$$

samples.