

**6.s077 — INTRODUCTION TO DATA SCIENCE  
EECS, MIT, Spring 2018**

**Lecture 3**

**Properties and assessment  
of parameter estimators**

## Today's agenda

---

- Estimators (incl. review)
- Sampling distributions
  - Bias, variance, standard error, MSE
  - for the sample mean
  - for ML
- Sampling distributions, in general
  - analytical
  - parametric bootstrap
  - general bootstrap
- Confidence intervals
  - definition and interpretation
  - based on normal (and other) approximations
  - via bootstrap

## Review

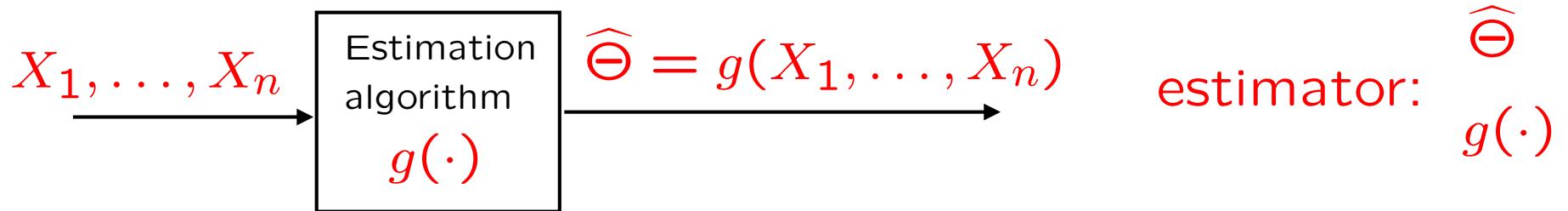
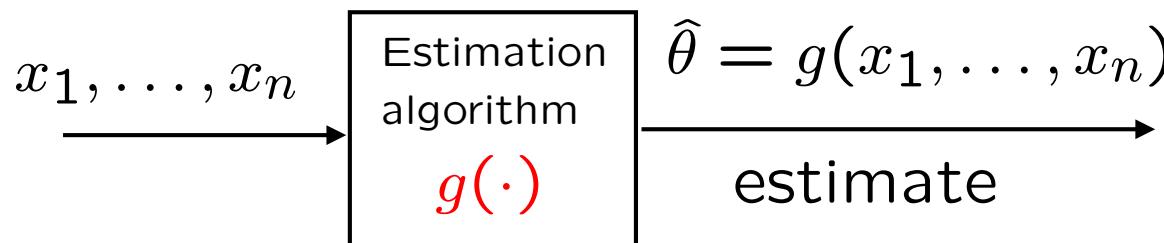
---

- data  $x_1, \dots, x_n$ 
  - observed values of i.i.d.  $X_1, \dots, X_n$ , distributed according to  $\mathbb{P}$   
unknown
- determine empirical distribution  $\hat{\mathbb{P}}$ 
  - probability  $1/n$  to each data point  $x_i$
- Three types (so far) of estimators
  - **Plugin:**  $a = h(\mathbb{P}) \rightarrow \hat{A} = h(\hat{\mathbb{P}})$
  - **Feature matching** (and method of moments): solve  $h(\mathbb{P}^\theta) = h(\hat{\mathbb{P}})$
  - **Maximum likelihood:** maximize  $\mathbb{P}^\theta(X = x)$

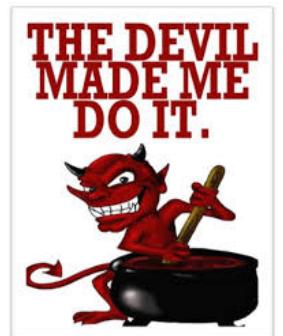
# Data and estimators

Data:  $x_1, \dots, x_n$

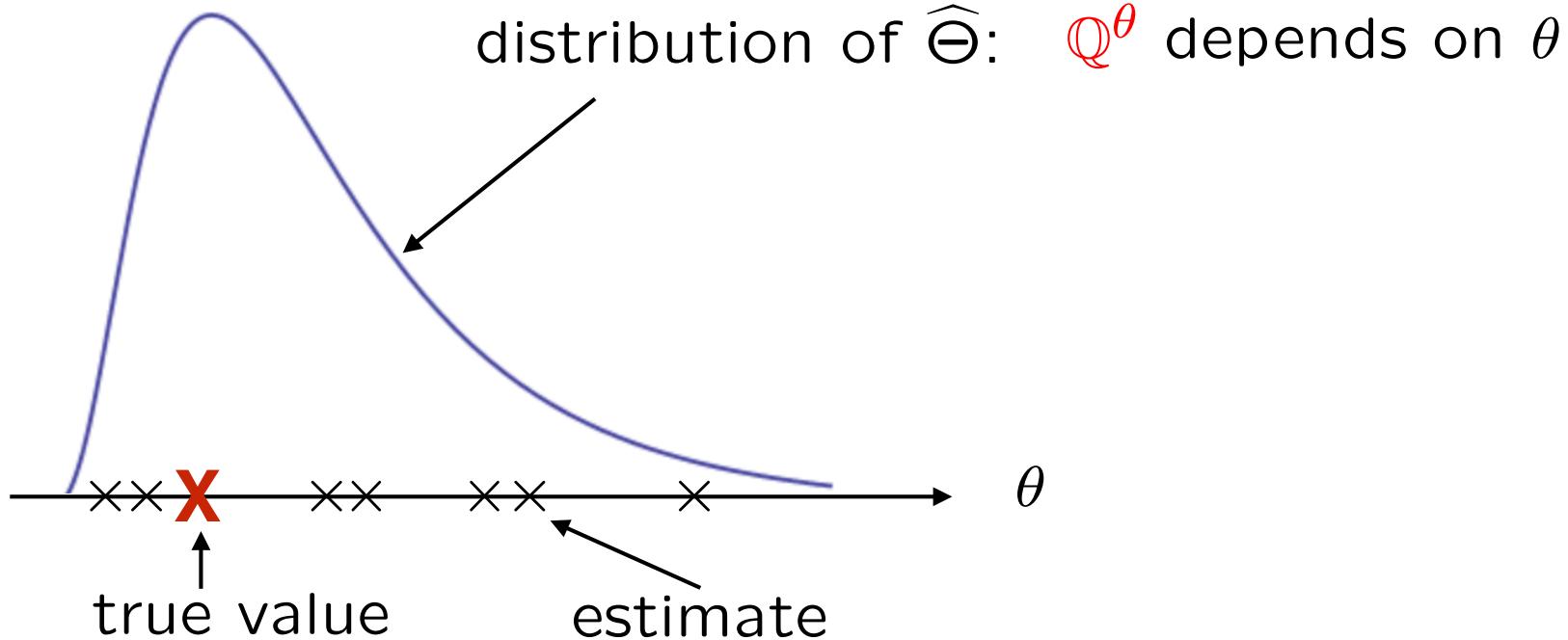
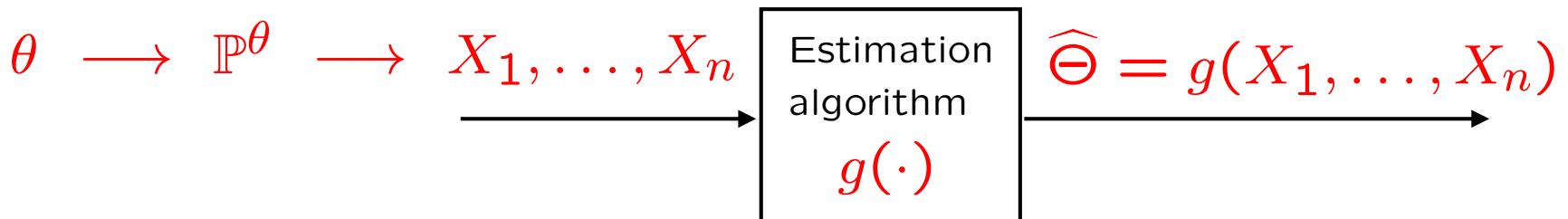
realized values of r.v.s:  
 $X_1, \dots, X_n$ , i.i.d.,  $\sim \mathbb{P}^\theta$



- Data scientist: design “good” estimator  $g(\cdot)$ 
  - before looking at the data
  - e.g., plugin, feature matching, ML, etc.
- What does “good” mean?

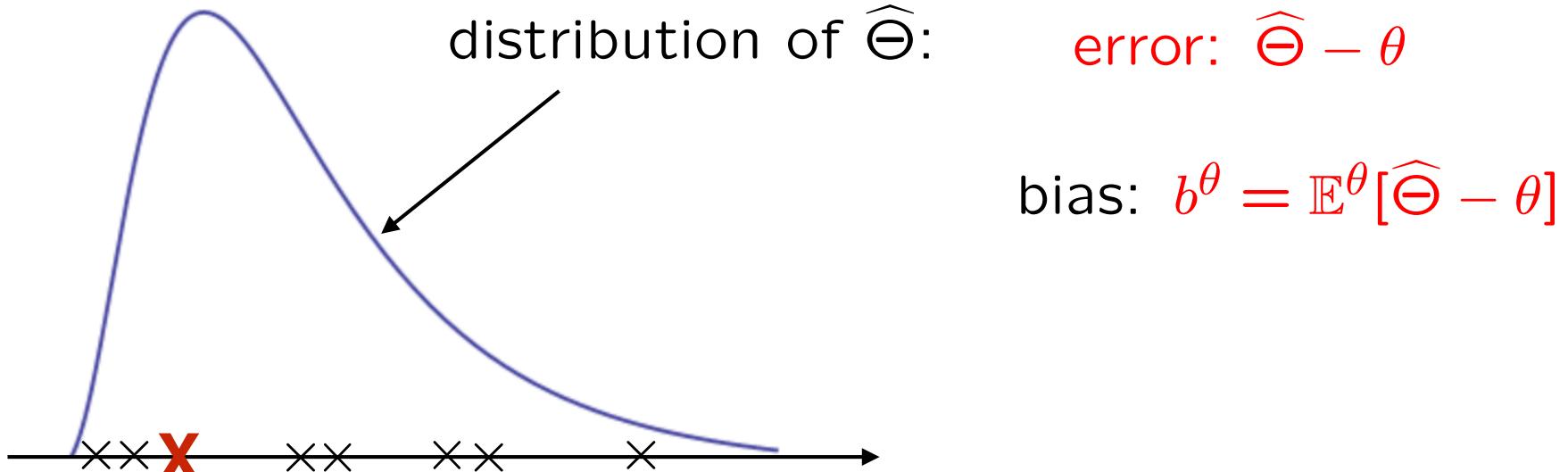


# Sampling distribution (of an estimator)



## Bias, variance, MSE, standard error

---



$$\text{variance: } v^\theta = E^\theta[(\hat{\Theta} - E^\theta[\hat{\Theta}])^2]$$

$$\text{standard error: } se^\theta = \sqrt{v^\theta}$$

$$\text{MSE}^\theta = E^\theta[(\hat{\Theta} - \theta)^2] = (b^\theta)^2 + v^\theta$$

- Usually:  $(b^\theta)^2$  is much smaller than  $v^\theta$   
there is a tradeoff

## Sampling distribution of the sample mean

$$X_i: \text{i.i.d, } \mu, \sigma^2 \quad \mu = \mathbb{E}[X_i] \quad \widehat{M} = \frac{1}{n} \sum_{i=1}^n X_i$$

$$\mathbb{E}[\widehat{M}] = \frac{1}{n} \sum_{i=1}^n \mu = \mu \quad \text{bias} = 0 \quad (\text{unbiased})$$

$$v = \text{var}(\widehat{M}) = \frac{1}{n^2} \text{var}\left(\sum_{i=1}^n X_i\right) = \frac{\sigma^2}{n} \xrightarrow{n \rightarrow \infty} 0$$

$$se = \frac{\sigma}{\sqrt{n}}$$

$$\bullet \text{ CLT: } \frac{\sqrt{n}(\widehat{M} - \mu)}{\sigma} \sim N(0, 1) \quad n \rightarrow \infty$$

Loosely speaking:  
 $\widehat{M} \sim N(\mu, \sigma^2/n)$

everything extends to estimation of  $\mathbb{E}[h(X)]$

# Sampling distribution of the sample mean

---

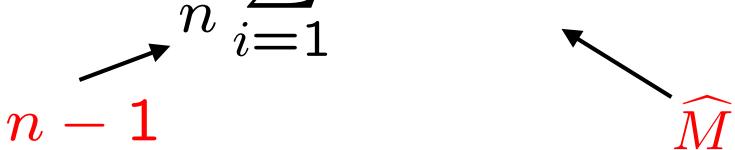
Loosely speaking:

$$\widehat{M} \sim N(\mu, \sigma^2/n)$$

Loosely speaking:

$$\widehat{M} \sim N(\mu, \widehat{\sigma}^2/n)$$

- Don't know  $\sigma$ ? Estimate it!

$$\text{plugin estimate: } \widehat{\sigma}^2 = \frac{1}{n-1} \sum_{i=1}^n (X_i - \bar{M})^2$$


- $\widehat{M}$  not quite normal for moderate  $n$ , especially if you are interested in the tails of the sampling distribution

# The wonderful properties of ML estimators

---

- Sampling distribution?

In general, can't avoid (difficult) problem-specific analysis or simulation

- **But**, under mild technical conditions, as  $n \rightarrow \infty$ , it is:

- **consistent**:  $\widehat{\Theta} \rightarrow \theta$

- **asymptotically normal**: 
$$\frac{\widehat{\Theta} - \theta}{\text{se}} \sim N(0, 1)$$

$\widehat{\text{se}}$ : estimate of se; there are good methods for obtaining it

- **efficient**: smallest possible se (asymptotically) among “reasonable” estimators

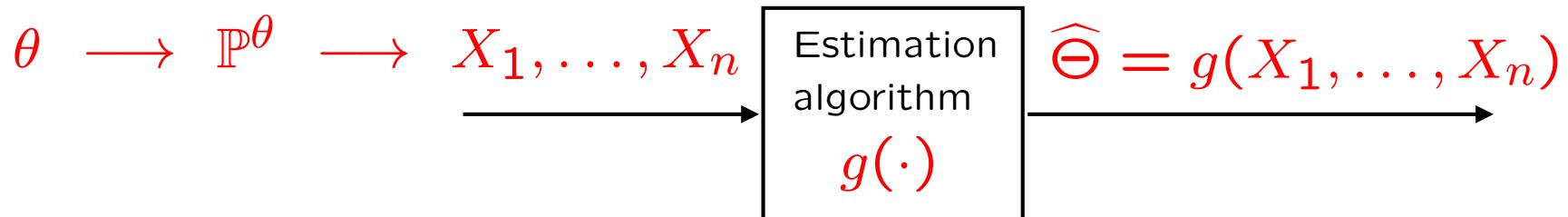
## Moving beyond the “friendly” cases

---

- For a general estimator  $\widehat{\Theta} = g(X_1, \dots, X_n)$ , how do we:
  - calculate/approximate the sampling distribution?
  - calculate/approximate the bias and the standard error?
  - use the above to obtain “confidence intervals” ?
- Approaches:
  - analytically
  - via the CLT
  - parametric bootstrap
  - general bootstrap

## Finding the sampling distribution: analytically or via the CLT

---



- We have the “model” structure  $\mathbb{P}^\theta$
  - Pretend  $\theta$  were known
  - A function of r.v's  $X_1, \dots, X_n$  is an r.v.
    - knowing  $\mathbb{P}^\theta$ , the distribution  $\mathbb{Q}^\theta$  of  $\widehat{\Theta}$  is completely determined
- 

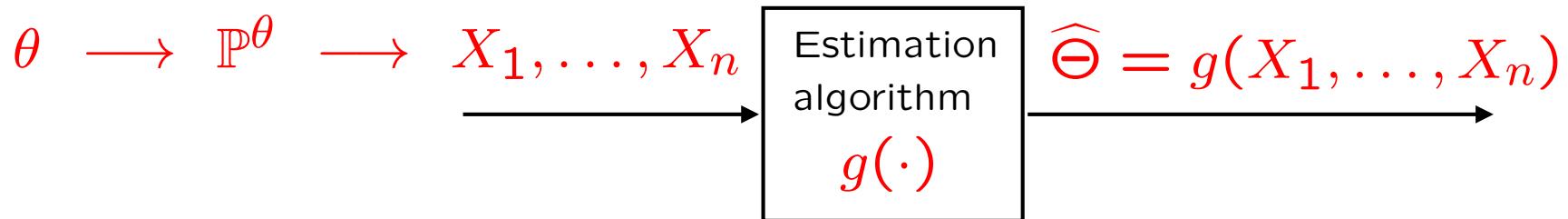
“derived distribution” exercise sometimes doable analytically

- Gives  $\mathbb{Q}^\theta$
- But!!! True  $\theta$  is not known

• Use  $\mathbb{Q}^{\widehat{\theta}}$

- 
- **CLT-based approach:** If theory tells us that  $\mathbb{Q}^\theta \sim N(\theta, se^\theta)$ , approximate:  $\mathbb{Q}^\theta \approx N(\widehat{\theta}, \widehat{se}^{\widehat{\theta}})$

## Finding the sampling distribution: Parametric bootstrap

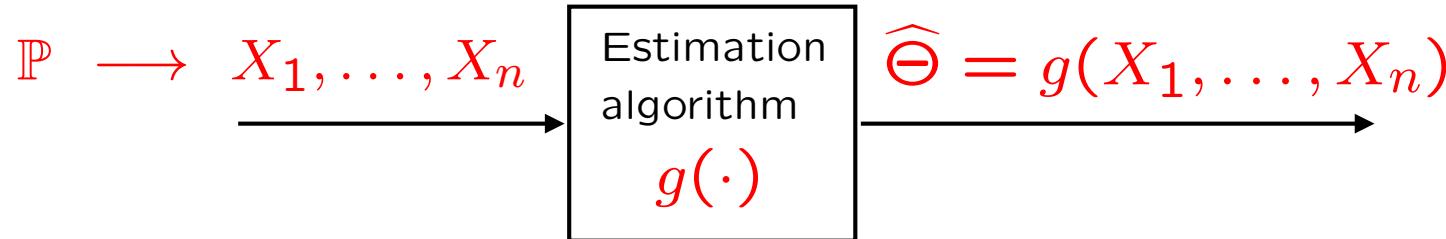


- We have the “model” structure  $P^\theta$
- Pretend  $\theta$  were known
- A function of r.v's  $X_1, \dots, X_n$  is an r.v.
  - knowing  $P^\theta$ , the distribution  $Q^\theta$  of  $\hat{\Theta}$  is completely determined and can be obtained by simulation
- But!!! True  $\theta$  is not known Rely on  $\hat{\theta} \approx \theta$

### Parametric bootstrap:

1. use actual data to compute  $\hat{\theta}$
2. generate synthetic data  $X_1, \dots, X_n$ , according to  $P^{\hat{\theta}}$
3. compute corresponding (new/synthetic)  $\hat{\theta}$
4. repeat steps 2-3 many times to generate a histogram/empirical distribution  $\hat{Q}$
5. Calculate bias, standard deviation of  $\hat{Q}$

## Finding the sampling distribution: General bootstrap



- Know little or nothing about the structure of  $\mathbb{P}$ 
  - $\theta$  need not be a parameter that determines the distribution, but just a feature/functional; e.g.,  $\theta = \text{median}$
- Want new samples from  $\mathbb{P}$ ; but cannot simulate  $\mathbb{P}$ 
  - but we know  $\hat{\mathbb{P}}$  (the actual data); and  $\hat{\mathbb{P}} \approx \mathbb{P}$   
use the samples that we have: draw  $n$  independent samples from  $\hat{\mathbb{P}}$

### General bootstrap:

1. use actual data  $x_1, \dots, x_n$  to compute actual  $\hat{\theta}$
2. draw  $n$  samples from the set  $\{x_1, \dots, x_n\}$ , uniformly, i.i.d. (with replacement)
3. use new  $n$  samples to compute corresponding (new)  $\hat{\theta}$
4. repeat steps 2-3 many times  $\rightarrow$  histogram/empirical distribution  $\hat{Q}$
5. Calculate bias, standard deviation of  $\hat{Q}$