

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

**Lecture 4**

**Confidence Intervals  
Predictions, and multidimensional data**

## Today's agenda

---

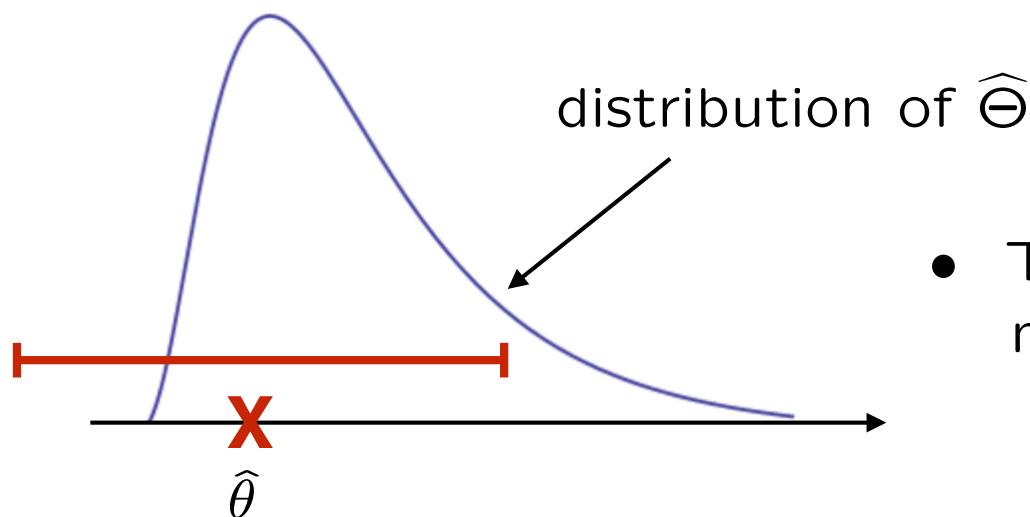
- Confidence intervals
  - definition and interpretation
  - based on normal (and other) approximations
  - via bootstrap

## Today's agenda

---

- Predicting an unknown random variable
  - expectation as a predictor
  - loss functions
  - empirical risk minimization (ERM)
- Multidimensional data
  - visualization
  - i.i.d. models
  - estimation (plugin, feature matching, ML)
- Predicting one variable from another
  - model-based,  $\mathbb{E}[Y | X]$
  - based on data  $(X_i, Y_i)$
- Structured (e.g., linear predictors)
  - Model-based; plugin; ERM
- Caution

## Confidence intervals (CIs)



- The value of an estimator  $\hat{\Theta}$  may not be informative enough

constructed using the data

- An  $1 - \alpha$  **confidence interval** is an interval  $[\hat{\Theta}^-, \hat{\Theta}^+]$ ,  
s.t.  $P^\theta(\hat{\Theta}^- \leq \theta \leq \hat{\Theta}^+) \geq 1 - \alpha$ , for all  $\theta$ 
  - often  $\alpha = 0.05$ , or  $0.025$ , or  $0.01$
- Interpretation (for  $\alpha = 0.05$ ). Suppose CI =  $[-3.2, 4.7]$ 
  - “There is probability 95% that  $\theta \in [-3.2, 4.7]$ ” **WRONG!!!**
  - if 100 statisticians produce 100 such CIs (for different data sets) we expect that 95 of them will capture  $\theta$
  - unless a certain extreme (5%) event has occurred, the CI captures  $\theta$

## Formulas for CIs – “nice” but common cases

- Assume sampling distribution:  $\widehat{\Theta} \approx N(\theta, se^2) \approx N(\theta, \widehat{se}^2)$

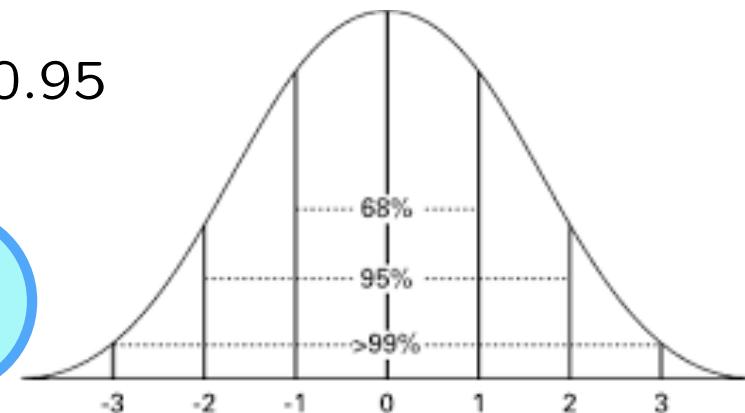
standard normal

exact,  
or CLT

bootstrap or other  
approximation (ML)

- Normal tables:  $\mathbb{P}\left(\frac{|\widehat{\Theta} - \theta|}{\widehat{se}} \leq 1.96\right) \approx 0.95$

$$\mathbb{P}(\widehat{\Theta} - 1.96 \widehat{se} \leq \theta \leq \widehat{\Theta} + 1.96 \widehat{se}) = 0.95$$

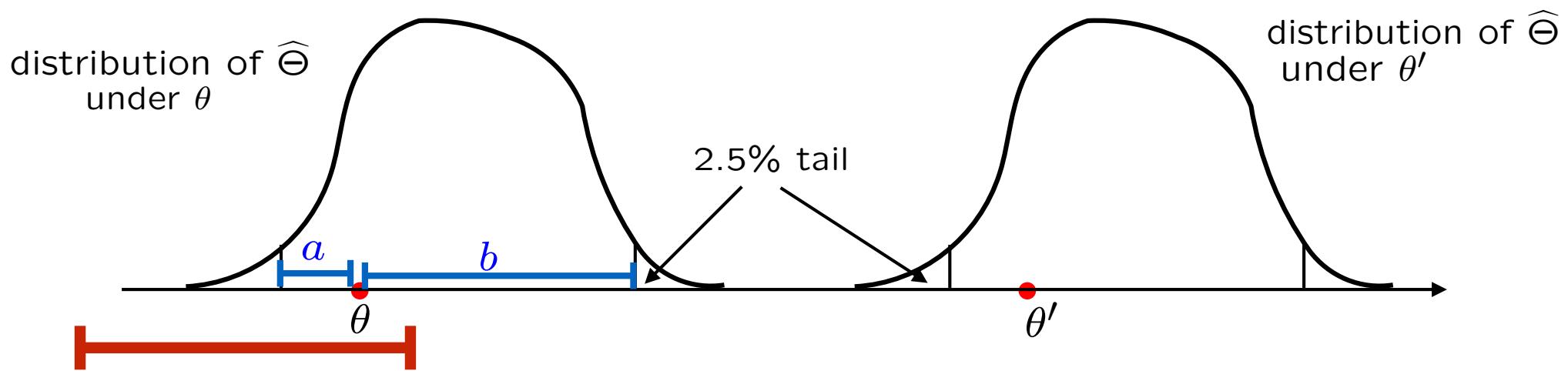


- Even if  $\widehat{\Theta}$  is normal (no approximation)  
the approximation  $se \approx \widehat{se}$   
may cause significant departure from normality ( $\rightarrow t$ -distribution)
- Caution:** We are dealing with the tails of the distributions;  
these can be very sensitive to approximations

## Constructing CIs more generally (assume $\alpha = 0.05$ )

- Want  $\mathbb{P}^\theta(\widehat{\Theta}^- \leq \theta \leq \widehat{\Theta}^+) \geq 0.95$ , for all  $\theta$  (Hard!)
- Data confine  $\theta$  and  $\mathbb{P}$  to a “small” range

**Assume:** sampling distribution does not change much within that small range



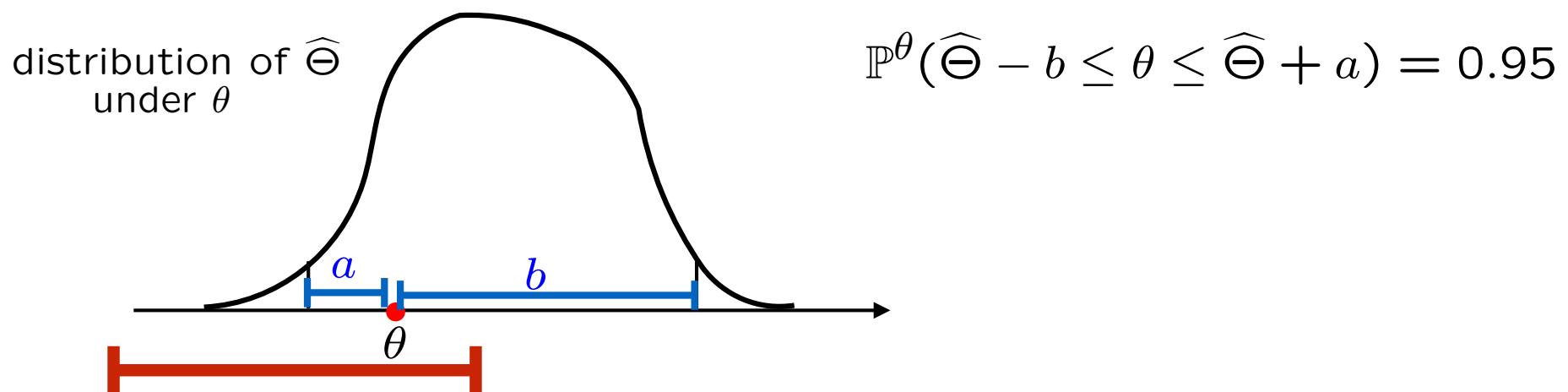
- unless a certain extreme (5%) event has occurred, the CI captures  $\theta$

$$\mathbb{P}^\theta(\theta - a \leq \widehat{\Theta} \leq \theta + b) = 0.95, \quad \text{for all } \theta$$

$$\mathbb{P}^\theta(\widehat{\Theta} - b \leq \theta \leq \widehat{\Theta} + a) = 0.95$$

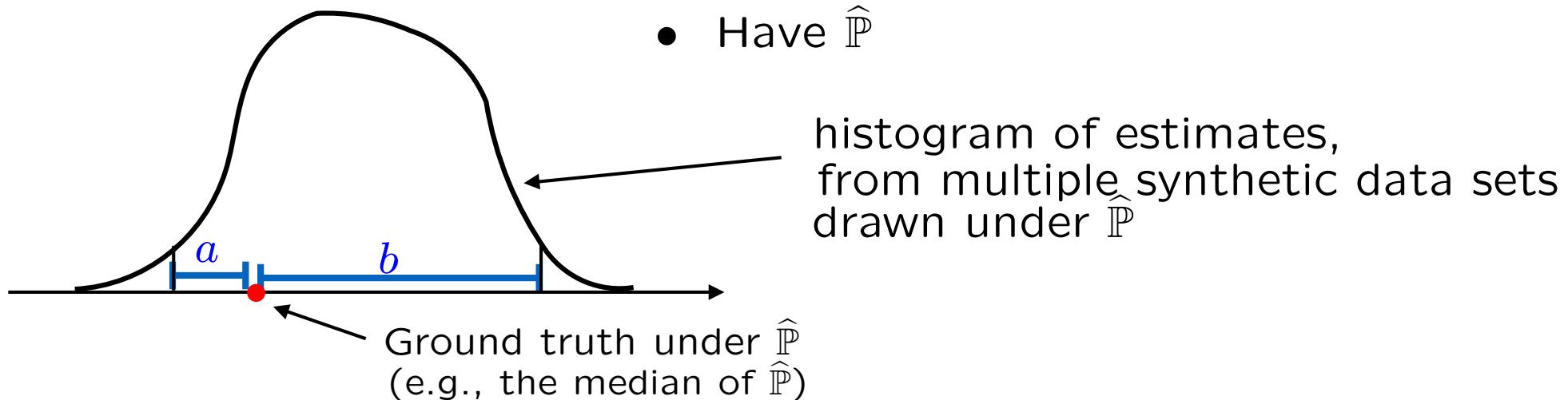
- earlier normal example:  $a = b = 1.96 \cdot \text{se}$

## Constructing CIs more generally (assume $\alpha = 0.05$ )



**Assume:** sampling distribution does not change much within that small range

- Simulate under some “reasonable”  $\theta$ : use  $\hat{\theta}$ , and simulate  $\mathbb{P}^{\hat{\theta}}$  (parametric bootstrap)
- General case; e.g., estimate the median of  $\mathbb{P}$ , know nothing about  $\mathbb{P}$ 
  - Have  $\hat{\mathbb{P}}$



## Recap

---

- Underlying assumptions (must always be questioned):
  - independent, identically distributed  $X_i$
  - finite mean and variance
- Key properties of an estimator  $\widehat{\Theta}$ :
  - bias, variance, standard error, sampling distribution  
(affected by the unknown true  $\theta$ )
  - need to know something about these properties in order to “interpret”  $\widehat{\Theta}$
  - estimating standard error and bias:  
exact analysis, CLT, parametric and general bootstrap
  - main focus is usually on the standard error (as bias is often smaller)

## Today's agenda

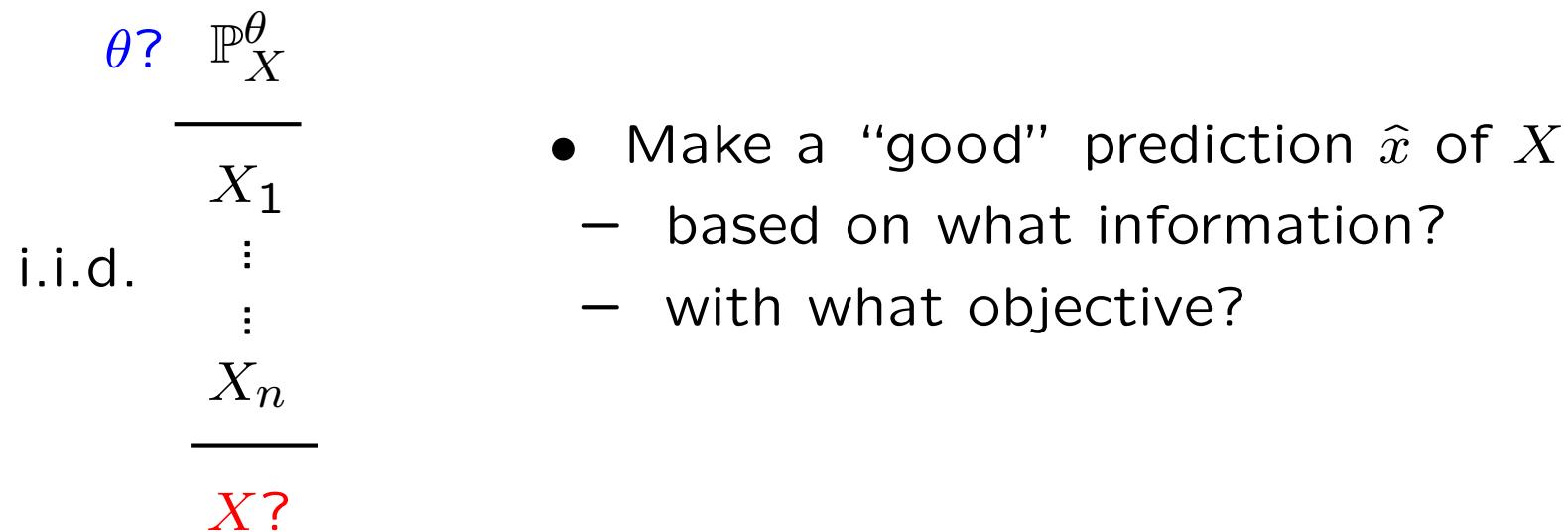
---

- Predicting an unknown random variable
  - expectation as a predictor
  - loss functions
  - empirical risk minimization (ERM)
- Multidimensional data
  - visualization
  - i.i.d. models
  - estimation (plugin, feature matching, ML)
- Predicting one variable from another
  - model-based,  $\mathbb{E}[Y | X]$
  - based on data  $(X_i, Y_i)$
- Structured (e.g., linear predictors)
  - Model-based; plugin; ERM
- Caution

## Why do we care about data or models?

---

- To make predictions [today, we focus on this]
- To understand how data were generated

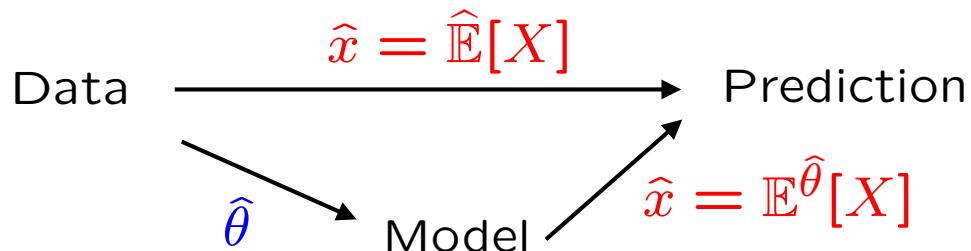


## Using data and/or models for predictions

$$\begin{array}{c} \theta? \quad \mathbb{P}_X^\theta \\ \hline X_1 \\ \vdots \\ X_n \\ \hline X? \end{array}$$

i.i.d.

- Make a “good” prediction  $\hat{x}$  of  $X$ 
    - based on what information?
    - with what objective?  $\min_{\hat{x}} \mathbb{E}^\theta[(X - \hat{x})^2]$
- A.** Have  $\mathbb{P}_X^\theta$ , know  $\theta$
- $$\hat{x} = \mathbb{E}[X] \quad (\text{data are irrelevant})$$
- B.** “Have”  $\mathbb{P}_X^\theta$ , don’t know  $\theta$
1. Model-based: estimate  $\hat{\theta}$  of  $\theta$   $\longrightarrow \hat{x} = \mathbb{E}^{\hat{\theta}}[X]$
  2. Plugin:  $\hat{x} = \hat{\mathbb{E}}[X] = (X_1 + \dots + X_n)/n$

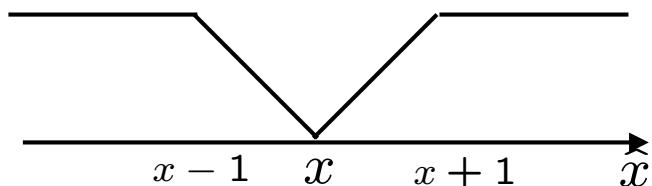


## Empirical risk minimization (ERM)

---

- Variation of the plugin idea
  - Model-based: minimize **risk**  $\mathbb{E}[(X - \hat{x})^2]$   $\hat{x} = \mathbb{E}[X]$
  - Data-based: minimize **empirical risk**  $\hat{\mathbb{E}}[(X - \hat{x})^2]$   $\hat{x} = \hat{\mathbb{E}}[X]$
- Approach extends to other risk functions (estimation objectives)
  - Model-based: minimize **risk**  $\mathbb{E}[|X - \hat{x}|]$   $\hat{x} = \text{median}(\mathbb{P})$
  - Data-based: minimize **empirical risk**  $\hat{\mathbb{E}}[|X - \hat{x}|]$   $\hat{x} = \text{median}(\hat{\mathbb{P}})$
- Objective formulated in terms of a loss function  $\ell(x, \hat{x})$

$$\text{minimize } \mathbb{E}[\ell(X - \hat{x})] \quad \hat{\mathbb{E}}[\ell(X - \hat{x})] = \frac{1}{n} \sum_{i=1}^n \ell(x_i - \hat{x})$$



- Approach has solid theory behind

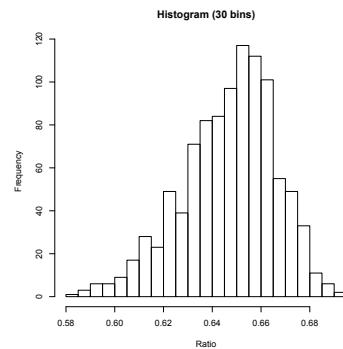
# Multidimensional data

---

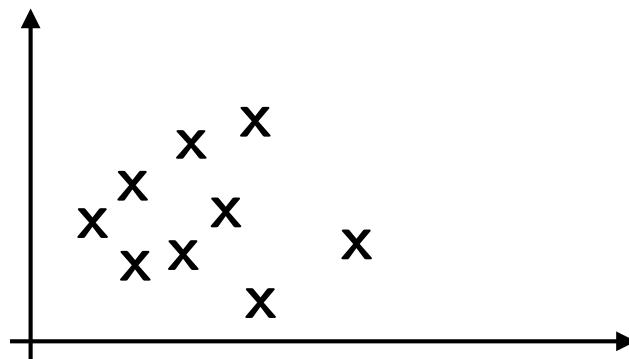
- $i$ th sample  $(x_i, y_i, z_i)$   
e.g., (gender, height, weight) of  $i$ th person

- Visualize:

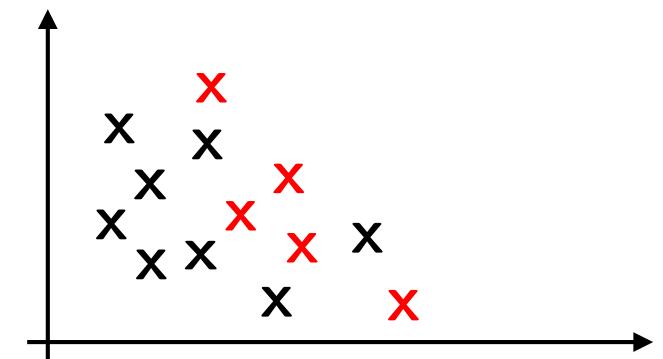
one at a time



two at a time



three, one categorical



## “Standard model:” Independent random sampling

---

- Assume two-dimensional data (only for exposition purposes)
- $\mathbb{P}$ : distribution on  $\mathbb{R}^2$  (discrete; continuous; mixed)

$$p_{X,Y}(x,y) \quad f_{X,Y}(x,y)$$

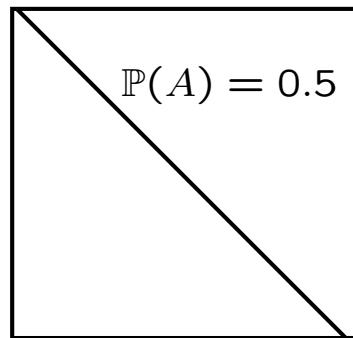
- $(X_i, Y_i) \sim \mathbb{P}$ , independent for different  $i$

$X_i$  and  $Y_i$  usually dependent

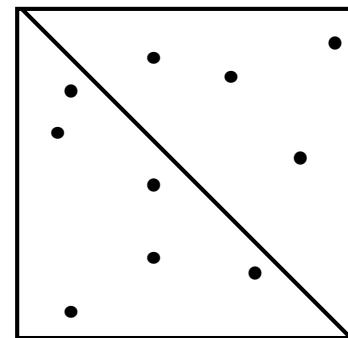
- $\hat{\mathbb{P}}$ : empirical distribution on  $\mathbb{R}^2$ 
  - weight  $1/n$  on each data point

$\hat{\mathbb{P}}(A) =$  fraction of data that belong to  $A$

uniform  $\mathbb{P}$



$\hat{\mathbb{P}}$



$\hat{\mathbb{P}}(A) = 4/10$

## Parameter estimation with multidimensional data

---

- Nothing new!

- **Plugin:**

$$\theta = \mathbb{E}[g(X, Y)] \quad \widehat{\Theta} = \widehat{\mathbb{E}}[g(X, Y)] = \frac{1}{n} \sum_{i=1}^n g(X_i, Y_i)$$

$$\theta = \text{Cov}(X, Y) = \mathbb{E}[XY] - \mathbb{E}[X]\mathbb{E}[Y]$$

$$\widehat{\Theta} = \widehat{\mathbb{E}}[XY] - \widehat{\mathbb{E}}[X]\widehat{\mathbb{E}}[Y]$$

## Parameter estimation with multidimensional data (ctd.)

---

- Feature matching and the method of moments
- Example:

$$\mathbb{P}^\theta \quad \theta = (\theta_1, \theta_2, \theta_3)$$

Choose  $\mathbb{E}^\theta[X]$ ,  $\mathbb{E}^\theta[Y]$ ,  $\mathbb{E}^\theta[XY]$  as features

$$\mathbb{E}^\theta[X] = \hat{\mathbb{E}}[X]$$

$$\mathbb{E}^\theta[Y] = \hat{\mathbb{E}}[Y] \qquad \text{solve for } \theta = (\theta_1, \theta_2, \theta_3)$$

$$\mathbb{E}^\theta[XY] = \hat{\mathbb{E}}[XY]$$

- easiest when we have closed-form expressions  
(as a function of  $\theta$ ) for left-hand side

## Parameter estimation with multidimensional data (ctd.)

---

- Maximum Likelihood

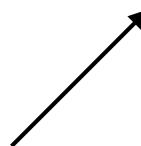
$$\max_{\theta} \prod_{i=1}^n \mathbb{P}^{\theta}(X_i = x_i) \quad \prod_{i=1}^n f_X^{\theta}(x_i)$$

$X_i, x_i$  are now vectors

$f_X$  is a joint PDF

- Nothing new
  - nice theoretical properties remain

- Sampling distribution, bootstrap, confidence intervals:  
again nothing new



one parameter at a time

## What is new? New types of prediction problems

---

$$\begin{array}{c} \theta? \quad \mathbb{P}_Y^\theta \\ \hline Y_1 \\ \vdots \\ Y_n \\ \hline Y? \end{array}$$

i.i.d.

$$\min_{\hat{y}} \quad \mathbb{E}^\theta \left[ (Y - \hat{y})^2 \right]$$

$$\begin{array}{c} \theta? \quad \mathbb{P}_{X,Y}^\theta \\ \hline X_1 \quad | \quad Y_1 \\ \vdots \quad | \quad \vdots \\ X_n \quad | \quad Y_n \\ \hline X \quad | \quad Y? \end{array}$$

$$\min_{g(\cdot)} \quad \mathbb{E}^\theta \left[ (Y - g(X))^2 \right]$$

$\hat{Y} = g(X)$  estimator

A. Have  $\mathbb{P}_Y^\theta$ , know  $\theta$

$$\hat{Y} = \mathbb{E}[Y]$$

Have  $\mathbb{P}_{X,Y}^\theta$ , know  $\theta$

$$\hat{Y} = \mathbb{E}[Y | X]$$

data  $X_1, Y_1, \dots, X_n, Y_n$  are irrelevant

## Some details

---

$$\mathbb{E}[(Y - \mathbb{E}[Y])^2] \leq \mathbb{E}[(Y - c)^2] \quad \text{for all } c$$

apply to conditional universe, where  $X = x$  has occurred

$$\mathbb{E}[(Y - \mathbb{E}[Y | X = x])^2 | X = x] \leq \mathbb{E}[(Y - c)^2 | X = x], \quad \text{for all } c$$

$$\mathbb{E}[(Y - \mathbb{E}[Y | X = x])^2 | X = x] \leq \mathbb{E}[(Y - g(x))^2 | X = x], \quad \text{for all } g(\cdot)$$

$$\mathbb{E}[(Y - \mathbb{E}[Y | X])^2 | X] \leq \mathbb{E}[(Y - g(X))^2 | X], \quad \text{for all } g(\cdot)$$

take expectation of both sides; use law of iterated expectations,  
and the property  $X \leq Y \Rightarrow \mathbb{E}[X] \leq \mathbb{E}[Y]$

$$\mathbb{E}[(Y - \mathbb{E}[Y | X])^2] \leq \mathbb{E}[(Y - g(X))^2], \quad \text{for all } g(\cdot)$$