

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

**Lecture 2**

**Estimating distributions and parameters  
(one-dimensional data)**

## Today's agenda

---

- Visualization
- Probability model/distribution  $\mathbb{P}$
- Empirical distribution  $\hat{\mathbb{P}}$
- Estimating a parameter
  - Plugin estimator
  - Feature matching (and method of moments)
  - Maximum likelihood (ML)
  - Empirical risk minimization
  - Bayesian methods
- Estimation of multiple parameters
- Representation matters

} Later lectures

## Data

---

- A collection of  $n$  records
- $i$ th record:  $x_i$ 
  - or  $(x_i, y_i, z_i)$
  - e.g.,  $x_i$  =height,  $y_i$  =weight,  $z_i$  =age
- How do we...
  - visualize?
  - think about them? (Probability!)
  - get something out of them? (This class)

$$\begin{matrix} x_1 & y_1 & z_1 \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ x_n & y_n & z_n \end{matrix}$$

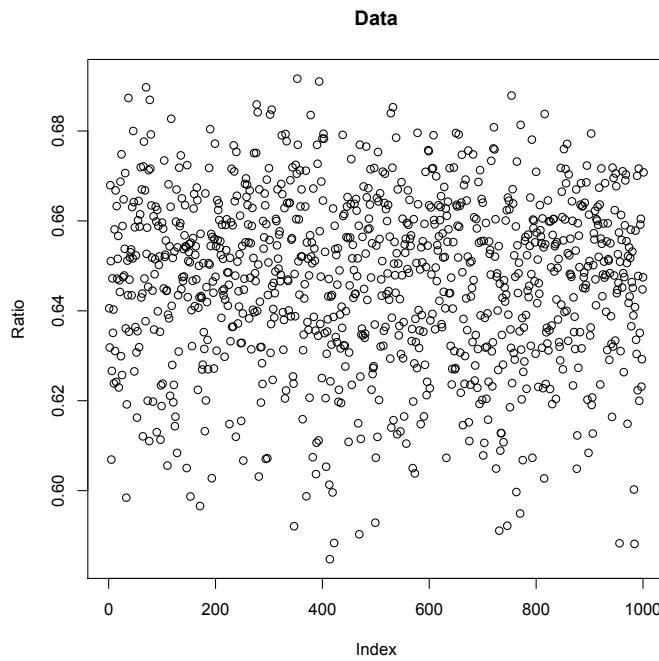
# Visualizing 1-dimensional data



$$x = \frac{\text{forehead breadth}}{\text{body length}}$$

$n = 1000$

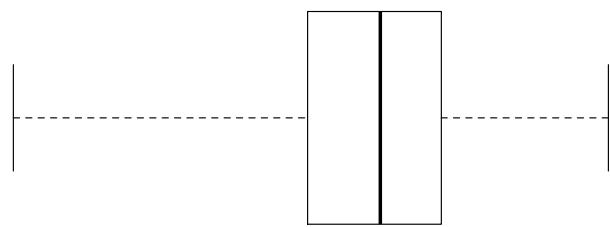
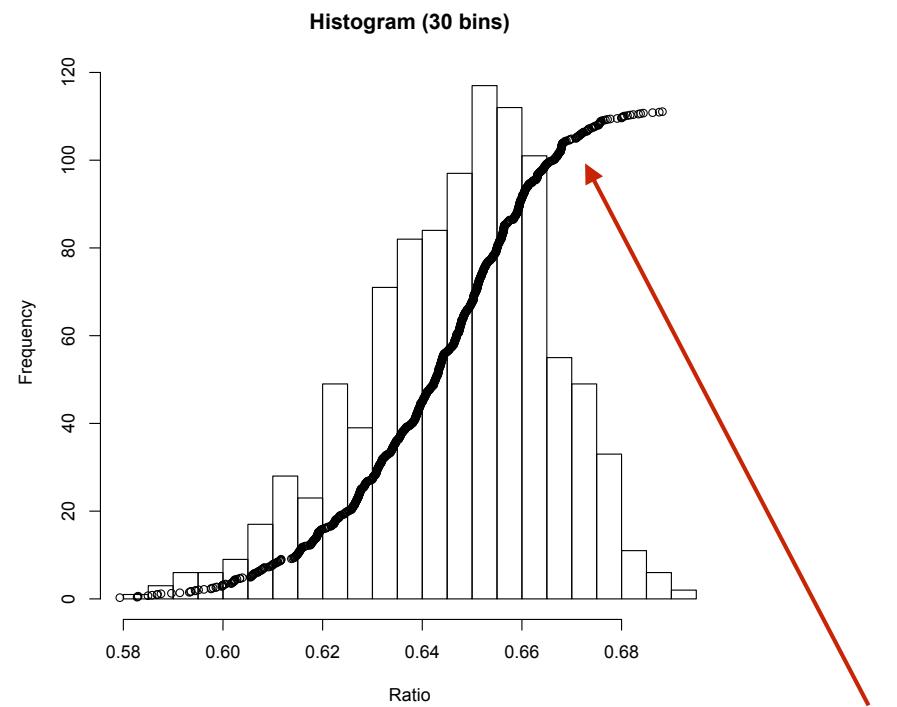
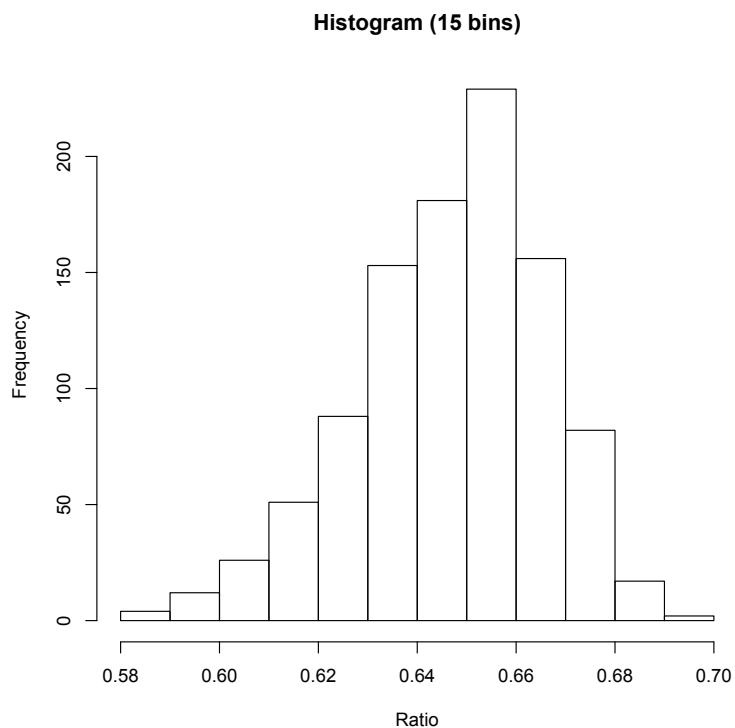
Naples crabs [W.F.R. Weldon; K. Pearson, 1894]



raw data ("scatter plot")

ratio	freq
0.5835	1
0.5875	3
0.5915	5
0.5955	2
0.5995	7
0.6035	10
0.6075	13
0.6115	19
0.6155	20
0.6195	25
0.6235	40
0.6275	31
0.6315	60
0.6355	62
0.6395	54
0.6435	74
0.6475	84
0.6515	86
0.6555	96
0.6595	85
0.6635	75
0.6675	47
0.6715	43
0.6755	24
0.6795	19
0.6835	9
0.6875	5
0.6915	1

## More insightful visualizations



box plot  
(shows the quartiles)

Empirical CDF,  $\hat{F}$

$$\hat{F}(\alpha) = \frac{\#\text{i s.t. } x_i \leq \alpha}{n}$$

## Numerical summaries of the data

---

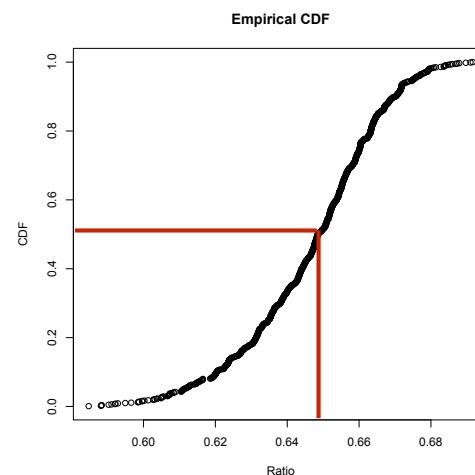
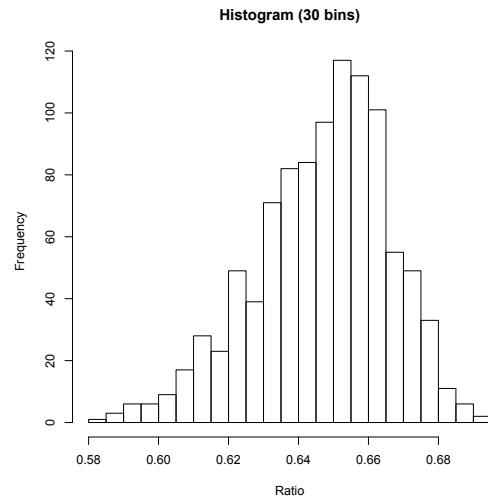
- **Statistic:** A real number that summarizes some aspect of the data

sample mean:  $\bar{x} = \frac{x_1 + \cdots + x_n}{n}$

sample standard deviation:  $\hat{\sigma} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \bar{x})^2}$

sample median:

$$\hat{F}(\alpha) = \frac{1}{2}$$



## How to think about the data

---

- How were the data generated?
  - Arbitrary/know nothing about the mechanism: can't do much
  - According to a partially known mechanism
    - for example:  $x_{k+1} = g(x_k, w_k)$ 

noise
- estimate  $g$ ,  $\mathbb{P}$
- model using probabilities  $\mathbb{P}$

- Start simple: no dependence

## Simplest data generation: independent random sampling

---

$\mathbb{P}$ : discrete or continuous distribution on  $\mathbb{R}$



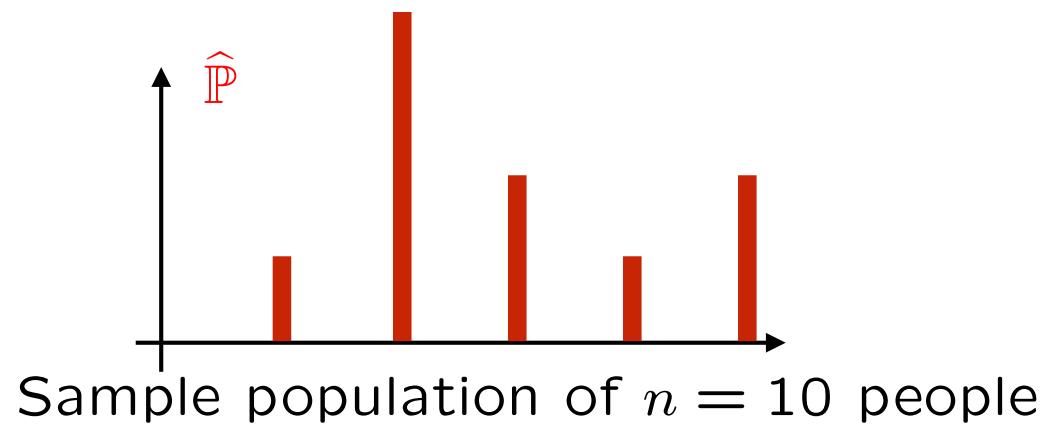
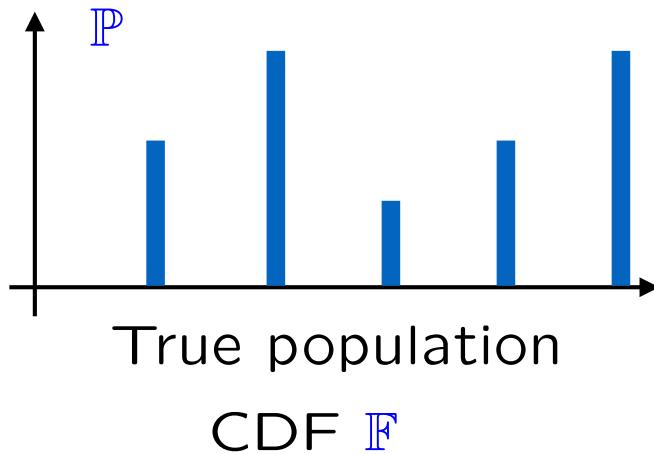
$X, X_1, \dots, X_n$ : random variables

$$\left. \begin{array}{l} X_i \sim \mathbb{P} \\ \text{independent} \end{array} \right\} \text{i.i.d.}$$

$$\mathbb{P}(X_1 \in A_1, \dots, X_n \in A_n) = \prod_{i=1}^n \mathbb{P}(X_i \in A_i) = \prod_{i=1}^n \mathbb{P}(X \in A_i)$$

## Independent random sampling

- **Example:** Finite population with 5 types of individuals
  - each person equally likely to be picked
  - independence: sampling “with replacement”



$\hat{\mathbb{P}}$  is a random “object”  
**“empirical distribution”**  
associated empirical CDF  $\hat{F}$

- Continuous  $X$ : as if infinite underlying population

**Central fact:** as  $n \rightarrow \infty$ , we have  $\hat{\mathbb{P}} \rightarrow \mathbb{P}$   
(see Lecture notes/Problem Set 1)

$\hat{F}(x) \rightarrow F(x)$ , for all  $x$

## Estimating parameters or numerical properties of $\mathbb{P}$

---

- Plugin estimator
- Feature matching (and method of moments)
- Maximum likelihood (ML)
- Empirical risk minimization
- Bayesian methods

} Later lectures

- Same/similar answers? Good
- If not, need ways to choose (performance analysis/assessment)

## Plugin estimators

---

$\mathbb{E}$ ,  $\widehat{\mathbb{E}}$ : expectations w.r.t.  $\mathbb{P}$ ,  $\widehat{\mathbb{P}}$

$$\mu = \mathbb{E}[X] \quad \widehat{M} = \frac{1}{n} \sum_{i=1}^n X_i = \widehat{\mathbb{E}}[X] \quad \xrightarrow[n \rightarrow \infty]{\longrightarrow} \mu \quad (\text{LLN})$$

$$\gamma = \mathbb{E}[X^3] \quad \widehat{\Gamma} = \frac{1}{n} \sum_{i=1}^n X_i^3 = \widehat{\mathbb{E}}[X^3]$$

$$\phi = \mathbb{E}[g(X)] \quad \widehat{\Phi} = \frac{1}{n} \sum_{i=1}^n g(X_i) = \widehat{\mathbb{E}}[g(X)]$$

$$v = \mathbb{E}\left[\left(X - \mathbb{E}[X]\right)^2\right] \quad \widehat{V} = \frac{1}{n} \sum_{i=1}^n (X_i - \widehat{M})^2$$

$$\widehat{V} = \frac{1}{n} \sum_{i=1}^n (X_i - \mu)^2 \quad (\text{if we know } \mu)$$

$$a = \text{median}(\mathbb{P}) \quad \widehat{A} = \text{median}(\widehat{\mathbb{P}})$$

$$a = h(\mathbb{P}) \quad \widehat{A} = h(\widehat{\mathbb{P}})$$

## Feature matching

---

- “**Feature**”: a property of a distribution  
e.g.: mean, variance, median, 70th percentile, etc.

$$\mathbb{P} \longrightarrow h \quad h = g(\mathbb{P})$$

- Distribution determined by parameter  $\theta$ :  $\mathbb{P}^\theta$   
feature depends on the parameter:  $h^\theta = g(\mathbb{P}^\theta)$

- Use data to calculate empirical value of the feature  $\hat{h} = g(\hat{\mathbb{P}})$
- Find (analytically)  $h_\theta$  as a function of  $\theta$
- Match the two: **solve, for  $\theta$ ,  $h_\theta = \hat{h}$**

- Example:
  - Find median  $\hat{h}$  of the data
  - Find  $\theta$  such that the median of  $\mathbb{P}^\theta$  agrees with  $\hat{h}$
- Food for thought:  
What would you do if you have one parameter  
but consider two (or more) features?

## Maximum likelihood

---

- Assume  $X_i$  discrete, i.i.d.;  $X = (X_1, \dots, X_n)$

**Likelihood function:**  $L^\theta(x_1, \dots, x_n) = \mathbb{P}^\theta(X = x) = \prod_{i=1}^n \mathbb{P}^\theta(X_i = x_i)$

- Find  $\theta$  under which the observed  $x$  is most likely to have been observed

having seen  $x$ :  $\max_\theta L^\theta(x) \quad \max_\theta \sum_{i=1}^n \log \left( \mathbb{P}^\theta(X_i = x_i) \right)$

Sometimes analytically

Usually numerically

- Example:  $X \sim \text{Binomial}(n, \theta)$

$$\max_\theta \binom{n}{x} \theta^x (1 - \theta)^{1-x} \quad \text{algebra} \longrightarrow \hat{\theta} = \frac{x}{n}$$

## Maximum likelihood — continuous data

---

- $X_i$ , independent PDF  $\sim f_{X_i}^\theta(\cdot)$        $X = (X_1, \dots, X_n)$

Likelihood function:  $L^\theta(x_1, \dots, x_n) = f_X^\theta(x) = \prod_{i=1}^n f_{X_i}^\theta(x_i)$

having seen  $x$ :  $\max_\theta L^\theta(x)$        $\max_\theta \sum_{i=1}^n \log(f_{X_i}^\theta(x_i))$

- Example:  $X_i \sim \text{Exp}(\theta)$ ,  $f^\theta(x) = \theta e^{-\theta x}$ ,  $x \geq 0$

## What comes next?

---

- Understand the different methods
  - study their properties
  - do they have desirable properties?
  - understand their errors
- Generalize
  - Each  $X_i$  may be a vector [later]
  - Parameter  $\theta$  may be a vector [now: straightforward]

## The case of a vector of parameters — Plugin and ML

---

- Examples:  $X \sim N(\mu, v)$   $\theta = (\mu, v)$   
unfair die,  $X \in \{1, \dots, 6\}$   $\theta = (p_1, \dots, p_6)$
- **Plugin:** Apply to each component of  $\theta$  separately
- **ML:**  $\max_{\theta} L^{\theta}(x)$   $\nabla_{\theta} \log L^{\theta}(x) = 0$   
computation can be hard

## Vector of parameters — Feature matching

---

- For concreteness: consider two parameters  $\theta = (\theta_1, \theta_2)$ 
  - will use two features  $h_i = g_i(\mathbb{P})$ ,  $i = 1, 2$

e.g.,  $h_1 = \text{mean}$   $h_2 = \text{median}$

$$\left. \begin{array}{l} h_1^\theta = g_1(\mathbb{P}^\theta) \\ h_2^\theta = g_2(\mathbb{P}^\theta) \end{array} \right\} \text{find expressions, as functions of } \theta$$

$$\text{data} \longrightarrow \hat{\mathbb{P}} \longrightarrow \left\{ \begin{array}{l} \hat{h}_1 = g_1(\hat{\mathbb{P}}) \\ \hat{h}_2 = g_2(\hat{\mathbb{P}}) \end{array} \right.$$

Solve  $2 \times 2$  system of equations  $\left\{ \begin{array}{l} h_1^\theta = \hat{h}_1 \\ h_2^\theta = \hat{h}_2 \end{array} \right.$

## Vector of parameters — The method of moments

---

- Suppose  $\theta$  involves  $k$  parameters
- Use  $\mathbb{E}[X], \dots, \mathbb{E}[X^k]$  as your  $k$  features
- Find  $\theta$  such that the first  $k$  moments of  $\mathbb{P}^\theta$  match the first  $k$  empirical moments  $\hat{\mathbb{E}}[X], \dots, \hat{\mathbb{E}}[X^k]$
- Crab example: Pearson (1894)
  - fit weighted average (“mixture”) of two normal PDFs
  - $k = 5$ , by hand!

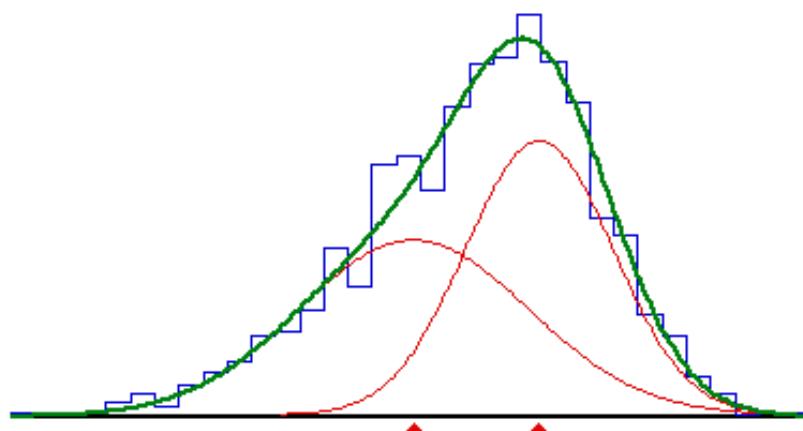


figure credit: Peter Macdonald

## Representation matters

---

- **Data:**  $X$  has the same information as  $Y = X^3$ 
  - Use  $X$  or  $Y$ ?
- **Parameters:**
  - variance  $v$  has same information as standard deviation:  $v = \sigma^2$
  - Estimate  $v$  or  $\sigma$ ?
- **Does it make a difference?**
  - in general, it does!
  - for ML it does not!