

Overview of parameter estimation

- We have a density function $f(X; \theta)$ whose parameters θ are unknown.
- We have a dataset D of n independent observations from f

- D is random variable denoting X_1, X_2, \dots, X_n where each $X_i \sim f(X; \theta)$

- We use any method to get an estimate $\hat{\theta} = A(D)$ as some function of D . Thus $\hat{\theta}$ is also a R.V.

- Alternative notations $\hat{\theta}_n, \hat{\theta}_D$ to stress that estimate depends on D .

Example: $f(x; \theta) \equiv N(x; \theta = (\mu, \sigma^2))$ $x \equiv$ height of people.

$$D \equiv \{X_1, X_2, \dots, X_{10}\}$$

$$\hat{\theta}_{\text{pop}} \equiv A(D) = \frac{v_1 + v_2 + \dots + v_{10}}{10}$$

$$\underbrace{f(x; \theta)}_{\text{hidden}} \xrightarrow{\quad D \quad} \begin{matrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{matrix}$$

$$x_1 = v_1, x_2 = v_2, \dots, x_n = v_n$$

$$\hat{\theta} \equiv \underline{A}(D)$$
$$\hat{\theta}_D \quad \hat{\theta}_n$$

Risk of an estimate — Theoretical exercise -

- Let $\underline{\theta}, \hat{\theta}_n$ be actual and estimated quantities.
- Risk = Expected square error

$$\underline{E_D [(\hat{\theta}_n - \theta)^2]}$$

$$\int (\hat{\theta}_n - \theta)^2 P(D) dD \quad D = \{x_1, \dots, x_n\}$$

$$\int_{x_1} \int_{x_2} \dots \int_{x_n} \underbrace{\int_{\hat{\theta}_n}^{\infty} (A(x_1, x_2, \dots, x_n) - \theta)^2 f(x_1) f(x_2) \dots f(x_n) dx_1 \dots dx_n}_{\text{Risk}} \quad \hat{\theta}^\infty = 6.0 \text{ ft}$$

Bias and Variance

$\hat{\theta}_n$ is a R.V

- Expected value of $E_D[\hat{\theta}_n] =$

$$\int_D A(D)f(D) = \int_{X_1} \dots \int_{X_n} A(X_1, \dots X_n) \prod_i f(X_i) dX_1 dX_2 \dots dX_n$$

Bias $E_D[\hat{\theta}_n] - \theta$

Variance

$$E_D\{(\hat{\theta}_n - E_D[\hat{\theta}_n])^2\} = E_D[\hat{\theta}_n^2] - (E_D[\hat{\theta}_n])^2$$

$$\text{Risk} = \text{Bias}^2 + \text{Variance}$$

$$\overline{\mathbb{E}_D[(\hat{\theta}_n - \theta)^2]} = (\mathbb{E}_D[\hat{\theta}_n] - \theta)^2 + \text{Var}(\hat{\theta}_n)$$

Proof:

$$\begin{aligned} \mathbb{E}_D[(\hat{\theta}_n - \theta)^2] &= \mathbb{E}_P[\hat{\theta}_n - \underbrace{\mathbb{E}_P[\hat{\theta}_n]}_{-b} + \underbrace{\mathbb{E}_P[\hat{\theta}_n] - \theta}_{\text{Bias}}]^2 \\ &= \mathbb{E}_D[(\hat{\theta}_n - \mathbb{E}_D[\hat{\theta}_n])^2] + \mathbb{E}_D[(\mathbb{E}_D[\hat{\theta}_n] - \theta)^2] \\ &\quad - 2 \mathbb{E}_D[(\hat{\theta}_n - \mathbb{E}_P[\hat{\theta}_n])(\mathbb{E}_P[\hat{\theta}_n] - \theta)] \\ &= \text{Var}(\hat{\theta}_n) + (\mathbb{E}_D[\hat{\theta}_n] - \theta)^2 - [2(\mathbb{E}_D[\hat{\theta}_n] - \mathbb{E}_D[\hat{\theta}_n]))(-)] \\ &= 0 \end{aligned}$$

Estimating CDF of any scalar random variable

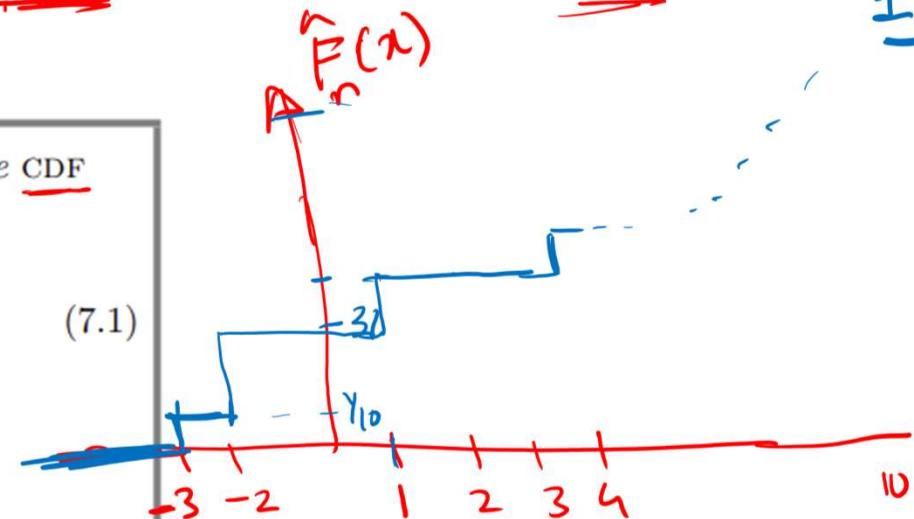
- Given sample: $D = \{X_1, X_2, \dots, X_n\}$ sampled i.i.d from an unknown $f(X)$
- Goal: estimate CDF function using D

7.1 Definition. The empirical distribution function \hat{F}_n is the CDF that puts mass $1/n$ at each data point X_i . Formally,

$$\hat{P}(X \leq x) \quad \hat{F}_n(x) = \frac{\sum_{i=1}^n I(X_i \leq x)}{n}$$

where

$$I(X_i \leq x) = \begin{cases} 1 & \text{if } X_i \leq x \\ 0 & \text{if } X_i > x. \end{cases}$$



Example: $n = 10, D = \{-3, -2, -2, 1, 2, 2, 2, 4, 6, 8, 10\}$

Analyzing Bias, Variance, and Risk of empirical CDF

Bias: at a x

$$E_D[\hat{F}_n(x)]$$

$\hat{F}_n(x)$ is unbiased.

$$\begin{aligned}
 &= E_D \left[\frac{1}{n} \sum_{i=1}^n I(X_i \leq x) \right] \\
 &= \frac{1}{n} \sum_{i=1}^n E_D [I(X_i \leq x)] = \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{x_i} f(x_i) dx_i \\
 &= \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^{x_i} I(X_i \leq x) \underbrace{\int_{-\infty}^{x_1} f(x_1) dx_1}_{x_1=1} \underbrace{\int_{-\infty}^{x_2} f(x_2) dx_2}_{x_2=1} \cdots \underbrace{\int_{-\infty}^{x_n} f(x_n) dx_n}_{x_n=1} = 1 \\
 &= \frac{1}{n} \sum_{i=1}^n \underbrace{\int_{-\infty}^{x_i} I(X_i \leq x) f(x_i) dx_i}_{\text{except } i} = \frac{1}{n} \sum_{i=1}^n \int_{-\infty}^x f(x_i) dx_i = \frac{n}{n} F(x) = \underline{F(x)}
 \end{aligned}$$

Variance $\hat{F}_n(x)$

$$\text{Var}\left(\left[\sum_{i=1}^n I(X_i \leq x)\right]\right) = \frac{1}{n^2} \sum_{i=1}^n \text{Var}_{X_i}(I(X_i \leq x))$$

$$= \frac{1}{n^2} \cdot n \text{Var}_{X_1}(I(X_1 \leq x))$$

$$\int_{x_1=-\infty}^x f(x_1) dx_1$$

$$= \frac{1}{n} \left(E_{X_1}([I(X_1 \leq x)^2]) - F(x)^2 \right)$$

$$= \frac{1}{n} \left(\int_0^x F(x) - F(x)^2 \right)$$

$$= \frac{1}{n} \left(F(x) - F(x)^2 \right)$$

$$= -\frac{1}{n} \underbrace{\int_0^x F(x)}_{F(x)(1-F(x))}$$

Non-parametric density estimation

Reading material: https://faculty.washington.edu/yenchic/18W_425/Lec6_hist_KDE.pdf

Motivation

- $D = \{X_1, X_2, \dots, X_n\}$ sampled i.i.d from an unknown $f(X)$
- Estimate $\hat{f}(x)$ without committing on a specific parametric form of $f(X)$
- Why not use empirical CDF?
 - Too inefficient to maintain. Need to store entire data
 - Too jerky. Non-zero density at observed points, zero elsewhere.
- Density estimation: assume some form of smoothness of $f(X)$