

Unsupervised learning

Not
Assumed
any
model
for data

Representation learning

↳ PCA / kernel PCA

Clustering

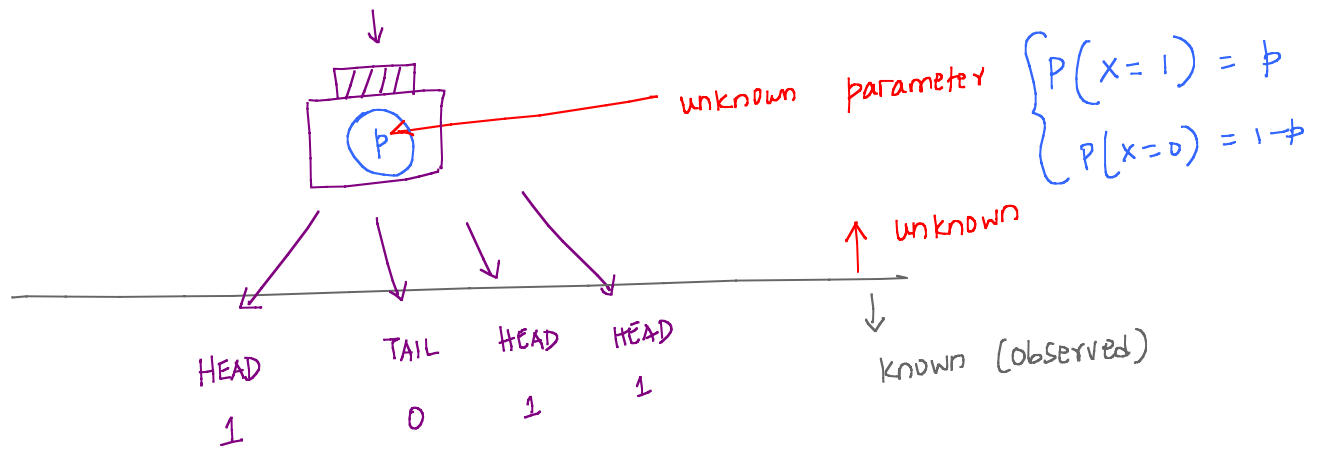
↳ Lloyd's / k-means.

Probabilistic
assumption

ESTIMATION

" There is some probabilistic mechanism that generates data about which we don't know "something". Given data, find/estimate what we don't know"

- OBSERVE data
- "ASSUME" a ^{probabilistic} model that generates data
- ESTIMATE unknown parameters using data.

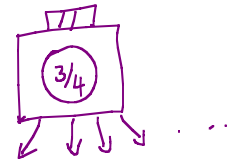


OBSERVE

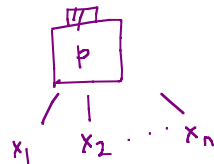
$\{1, 0, 0, 1, 1, 0, 1, 1, 1, 1, 1\}$

ESTIMATE:

$$9/12 = \boxed{3/4} \leftarrow$$



ASSUMPTIONS



OBSERVATIONS ARE

- (i) INDEPENDENT
- (ii) IDENTICALLY DISTRIBUTED

$$\text{Guess} = 2/3 ? \quad Y$$

$$= 0.0001 ? \quad Y$$

$$= 0 ? \quad N$$

$$= 1 ? \quad N$$

INDEPENDENCE

$$P(x_i / x_j) = P(x_i) \quad * i \neq j$$

IDENTICAL DISTRIBUTION

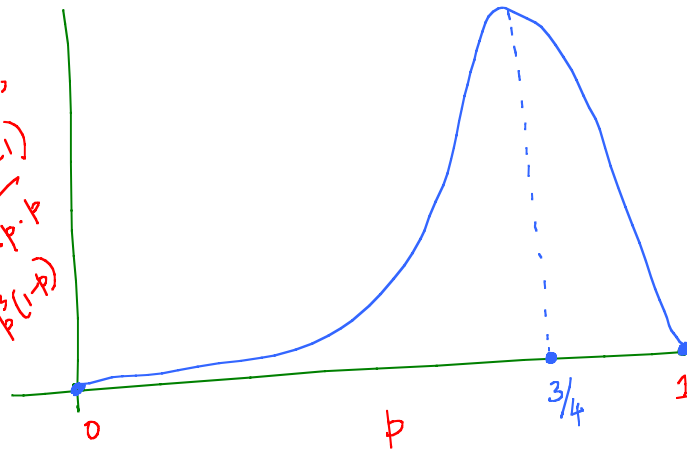
$$P(x_i = 1) = P(x_j = 1) = p \quad * i, j$$

LIKELIHOOD FUNCTION

$$(x) : \{1, 0, 1, 1\}$$

3/4

$$P(x_1=1, x_2=0, x_3=1, x_4=1) \\ = p \cdot (1-p) \cdot p \cdot p \\ = p^3(1-p)$$



FISHER'S PRINCIPLE OF MAXIMUM LIKELIHOOD

$$L(p; \{x_1, x_2, \dots, x_n\}) = P(x_1, x_2, \dots, x_n; p) \quad \text{underlying parameter.}$$

$$= P(x_1; p) \cdot P(x_2; p) \cdots P(x_n; p) \quad \text{[Independence]}$$

$$= \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} \rightarrow \begin{cases} \text{if } x_i = 1 \Rightarrow p^1 (1-p)^0 = p \\ \text{if } x_i = 0 \Rightarrow p^0 (1-p)^1 = 1-p \end{cases}$$

$$\hat{p}_{ML} = \arg \max_p \prod_{i=1}^n p^{x_i} (1-p)^{1-x_i}$$

$$= \arg \max_p \log \left(\prod_{i=1}^n p^{x_i} (1-p)^{1-x_i} \right)$$

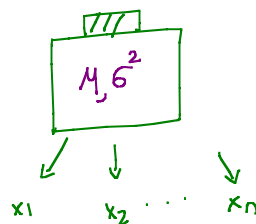
[log is monotonic increasing]

$$= \arg \max_p \sum_{i=1}^n [x_i \log p + (1-x_i) \log (1-p)]$$

Take derivative of $\log L(p)$, set it to 0 to get

$$\hat{p}_{ML} = \frac{1}{n} \sum_{i=1}^n x_i \leftarrow \text{Fraction of 1's}$$

$$\text{Data} = \{x_1, \dots, x_n\} \quad x_i \in \mathbb{R} \quad \forall i$$



$$x_i \sim \text{Gaussian}(\mu, \sigma^2) \quad \forall i$$

$\mu \rightarrow \text{unknown}; \quad \sigma^2 \rightarrow \text{known.}$

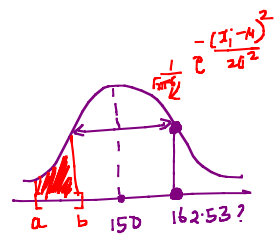
$$\begin{aligned} L(\mu, \sigma^2, \{x_1, \dots, x_n\}) &= P(x_1, \dots, x_n; \mu, \sigma^2) \\ &= \prod_{i=1}^n \underbrace{P(x_i; \mu, \sigma^2)}_0 \quad \times \end{aligned}$$

$$\underline{L(\mu, \sigma^2, \{x_1, \dots, x_n\})} = \int_{x_1, \dots, x_n} (x_1, \dots, x_n; \mu, \sigma^2)$$

$$= \prod_{i=1}^n f_{x_i}(\underline{x_i}; \underline{\mu}, \underline{\sigma^2})$$

$$= \prod_{i=1}^n \left[\frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x_i - \mu)^2}{2\sigma^2}} \right]$$

$$x_k = 162.53 \text{ cm.}$$



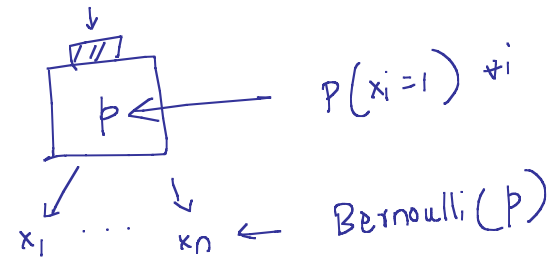
$$\log L(\mu, \sigma^2, \{x_1, \dots, x_n\}) = \sum_{i=1}^n \left[\underbrace{\log\left(\frac{1}{\sqrt{2\pi}\sigma}\right)}_x - \frac{(x_i - \mu)^2}{2\sigma^2} \right]$$

$$\hat{\mu}_{ML} = \operatorname{argmax}_{\mu} \sum_{i=1}^n - (x_i - \mu)^2$$

$$\boxed{\hat{\mu}_{ML} = \frac{1}{n} \sum_{i=1}^n x_i}$$

Consider the coin examples.

BAYESIAN
MODELING



Let's say
Someone says:

"I believe the bias p is somewhere close to 1"

We may have "HUNCH" about parameters

Goal: Incorporate "hunch/belief" about parameters of interest into the estimation procedure.

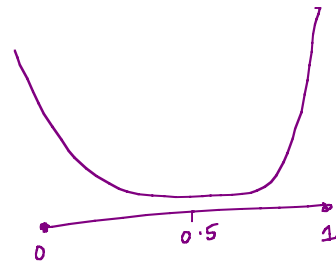
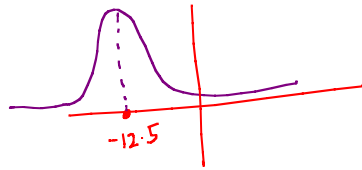
APPROACH: Think of the parameter to estimate as a "random" variable.

EARLIER

μ

p

Now



HUNCH → Codified using a probability
distribution over θ

$P(\theta)$
PRIOR

⇓ DATA

UPDATED
HUNCH

Codified using a prob.
distribution

$P(\theta/\text{DATA})$

POSTERIOR

Bayes law

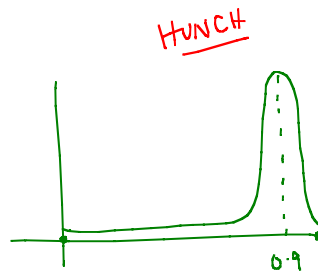
$$P(A|B) = \frac{P(B|A) \cdot P(A)}{P(B)}$$

A → Parameters θ

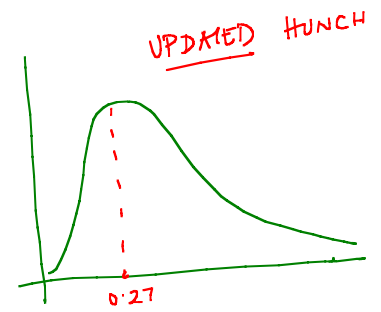
B → DATA $\{x_1, \dots, x_n\}$

$$\underbrace{P(\theta | \{x_1, \dots, x_n\})}_{\text{POSTERIOR}} = \left(\frac{\underbrace{P(\{x_1, \dots, x_n\} | \theta)}_{\text{LIKELIHOOD}}}{\underbrace{P(\{x_1, \dots, x_n\})}_{\text{EVIDENCE}} \uparrow \text{Does not depend on } \theta} \right) \cdot \underbrace{P(\theta)}_{\text{PRIOR}}$$

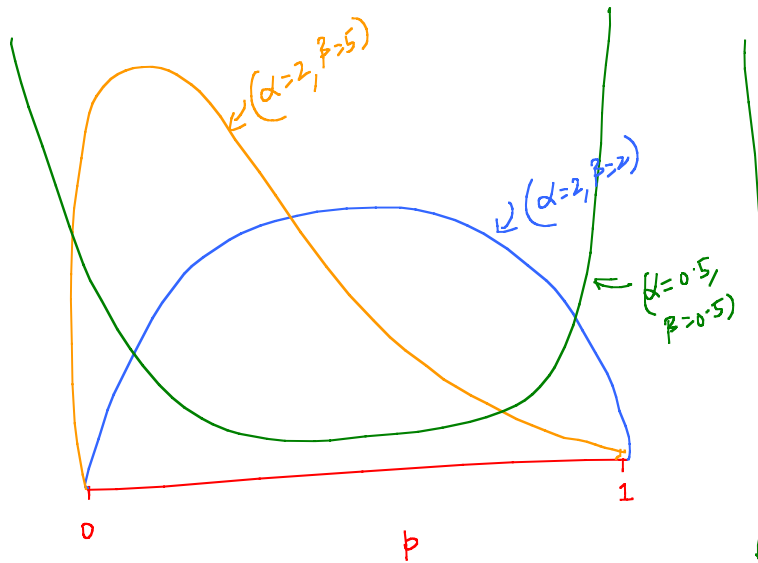
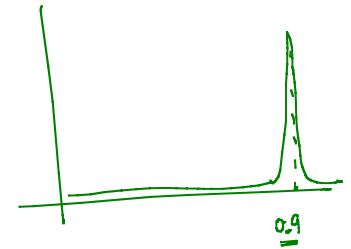
EXAMPLE



DATA
→
 $\{0, 0, 0, 0, 0, 1, 1, 0, 0, 0\}$



$\hookrightarrow \{ \overbrace{1,1,1,1,1}^{\Rightarrow}, \underbrace{0,1,1,1,1}_{} \}$



DATA - Bernoulli (p)

PRIOR ? $P(\theta)$

BETA PRIOR

$$f(p; \alpha, \beta) =$$

$$\left\{ \frac{p^{\alpha-1} (1-p)^{\beta-1}}{Z} \right.$$

$+ p \in [0, 1]$

$$P(\theta | \text{DATA}) \propto P(\text{DATA} | \theta) \cdot P(\theta)$$

$$f_{p/\text{DATA}}(p) \propto \underbrace{\left[\prod_{i=1}^n p^{x_i} (1-p)^{(1-x_i)} \right]}_{\text{LIKELIHOOD}} \cdot \underbrace{\left[\frac{p^{\alpha-1} (1-p)^{\beta-1}}{\text{PRIOR}} \right]}_{\text{PRIOR}}$$

$$f_{p/\text{DATA}}(p) \propto \underbrace{p^{\sum x_i + \alpha - 1} (1-p)^{\sum (1-x_i) + \beta - 1}}_{\text{}} \quad \underbrace{\quad}_{\text{}} \quad \underbrace{\quad}_{\text{}}$$

↳ same functional form as the PRIOR!!

$$\text{BETA PRIOR } (\alpha, \beta) \xrightarrow[\text{Bernoulli}]{\text{DATA}} \text{BETA POSTERIOR } (\alpha + n_h, \beta + n_t)$$

one possible
guess

$$\frac{\alpha + n_h}{\alpha + n_h + \beta + n_t} = \frac{\alpha + n_h}{(\alpha + \beta) + n}$$

$$\underline{\underline{E[\text{Posterior}]} = E[\text{Beta}(\alpha + n_h, \beta + n_t)] =$$

MAP Estimator - Maximum A posteriori Estimator
 \uparrow
 \hat{p}_{MAP}