

Density Estimation

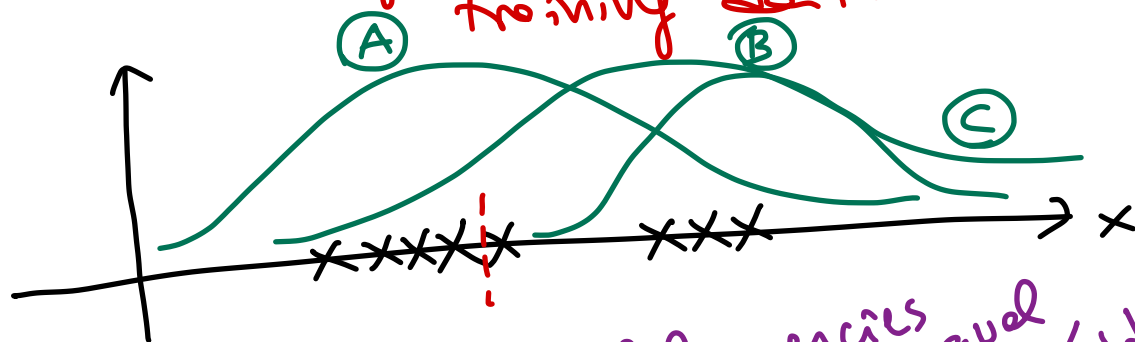
$$\mathcal{X} = \{x_i\}_{i=1}^N$$

N samples
 N data points

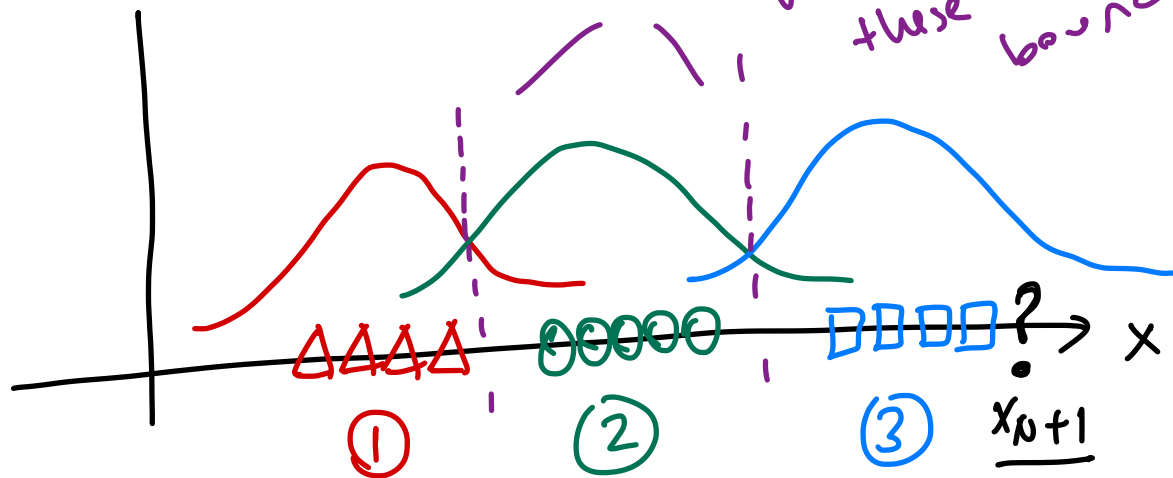
$x_i \sim p(x_i)$ $\forall i \Rightarrow$ probability distribution
 \Downarrow
 parameters (?)

DENSITY ESTIMATION

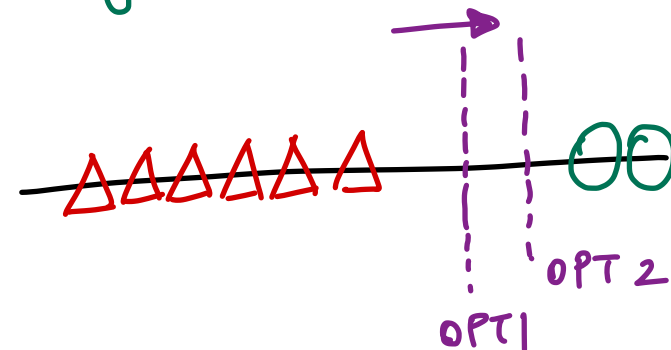
learning these parameters from training data



if frequencies were equal these would be boundaries



$x_i \sim N(x_i; \mu, \sigma^2)$
 μ^* : the best μ parameter
 σ^{*2} : the best σ^2 parameter



$$\mathcal{X} = \{ (x_i, y_i) \}_{i=1}^N \quad x_i \in \mathbb{R}^1 \quad y_i \in \{1, 2, 3\}$$

class densities $\Rightarrow p(x | y=c) \Rightarrow$ density estimation

prior distribution $\Rightarrow P(y=c) \Rightarrow$ class frequencies

BAYES RULE: $P(B|A) = \frac{P(A, B)}{P(A)} \rightarrow$ joint distribution

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

$$\overbrace{P(y=c | x)}^{\text{posterior}} = \frac{p(x | y=c) P(y=c)}{p(x)}$$

$$\begin{array}{|c|} \hline 1003 \\ \hline 1004 \\ \hline 1002 \end{array} \quad \begin{array}{|c|} \hline 3 \\ \hline 4 \\ \hline 2 \end{array}$$

test data point
 \Downarrow

x_{N+1}

$$\Rightarrow P(y=c | x_{N+1})$$

$$\rightarrow P(y=1 | x_{N+1})$$

$$\rightarrow P(y=2 | x_{N+1})$$

$$\rightarrow P(y=3 | x_{N+1})$$

} pick the maximum one

MAXIMUM LIKELIHOOD ESTIMATION (MLE)

$$\mathcal{X} = \{x_i\}_{i=1}^N$$

$$x_i \sim p(x_i | \theta)$$

unknown parameters

x_i 's are i.i.d.

↳ identically & independently distributed

Likelihood $\equiv p(x_1, x_2, x_3, \dots, x_N | \theta) \Rightarrow$ full joint

$$L(\theta | \mathcal{X}) \equiv p(x_1 | \theta) \cdot p(x_2 | \theta) \cdot \dots \cdot p(x_N | \theta)$$

$$\equiv \prod_{i=1}^N p(x_i | \theta)$$

$$\Rightarrow \underset{\substack{\uparrow \\ \text{best } \theta}}{\theta^*} = \arg \max_{\theta} L(\theta | \mathcal{X})$$

$$\log \text{likelihood} = \log \left[\prod_{i=1}^N p(x_i | \theta) \right]$$

$$= \sum_{i=1}^N \log [p(x_i | \theta)]$$

$$\begin{aligned} \log(a^b) &= b \cdot \log(a) \\ \log(a \cdot b) &= \log(a) + \log(b) \end{aligned}$$

Bernoulli Density:

$$0 < \pi < 1$$

→ success probability

$$\left[\begin{array}{l} \frac{\partial \log(x)}{\partial x} = \frac{1}{x} \\ \frac{\partial \log(1-x)}{\partial x} = -\frac{1}{1-x} \end{array} \right]$$

(H) success: $\pi \Rightarrow x=1$
(T) failure: $1-\pi \Rightarrow x=0$

(A) Head/Tail = $(1/2, 1/2)$
(B) Head/Tail = $(7/10, 3/10)$



⇒ H T H H H H T ... T
 $x_1 \ x_2 \ x_3 \ x_4 \ x_5 \ x_6 \ x_7 \ \dots \ x_{100}$
1 0 1 1 1 1 0 ... 0 } 70 heads
30 tails

$$p(x_i | \pi) = \pi^{x_i} \cdot (1-\pi)^{1-x_i}$$

$$P(x_i=1 | \pi) = \pi^1 \cdot (1-\pi)^{1-1} = \pi$$

$$P(x_i=0 | \pi) = \pi^0 \cdot (1-\pi)^{1-0} = 1-\pi$$

$$\mathcal{L}(\pi | \mathcal{X}) = \prod_{i=1}^N \pi^{x_i} (1-\pi)^{1-x_i}$$

$$\log \mathcal{L}(\pi | \mathcal{X}) = \sum_{i=1}^N [x_i \cdot \log(\pi) + (1-x_i) \log(1-\pi)] \Rightarrow \pi^* = ?$$

$$\frac{\partial \log \mathcal{L}(\pi | \mathcal{X})}{\partial \pi} = \sum_{i=1}^N \left[x_i \cdot \frac{1}{\pi} + (1-x_i) \cdot \frac{(-1)}{(1-\pi)} \right] = 0 \Rightarrow \pi^* = \frac{\sum_{i=1}^N x_i}{N}$$

of heads

Gaussian Density: $\mathcal{X} = \{x_i\}_{i=1}^N$

$$x_i \sim N(x_i; \mu, \sigma^2) \Rightarrow \mu^* = ? \quad \sigma^{2*} = ?$$

$$\sim \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left[-\frac{(x_i - \mu)^2}{2\sigma^2}\right] \quad -\infty < x_i < +\infty$$

$$\log \text{Likelihood} \equiv \log \prod_{i=1}^N \left[\frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp\left[-\frac{(x_i - \mu)^2}{2\sigma^2}\right] \right]$$

$$\log \text{likelihood} = \sum_{i=1}^N \left[-\frac{1}{2} \log(2\pi\sigma^2) + \left[-\frac{(x_i - \mu)^2}{2\sigma^2} \right] \right]$$

$$\frac{\partial \log \text{-likelihood}}{\partial \mu} = 0$$

$$\frac{\partial \log \text{likelihood}}{\partial \sigma^2} = 0$$

$$\mu^* = \frac{\sum_{i=1}^N x_i}{N}$$

sample mean

$$\sigma^{2*} = \frac{\sum_{i=1}^N (x_i - \mu^*)^2}{N}$$

sample variance

Parametric Classification:

Input: A training data set

Output: A classifier

$$\mathcal{X} = \{(x_i, y_i)\}_{i=1}^N$$

$$\hat{y}_{N+1} = \arg \max_c g_c(x_{N+1})$$

test data point

score function for class # c.

$$P(y=c|x) = \frac{p(x|y=c) \cdot P(y=c)}{p(x)}$$

$p(x) \Rightarrow$ independent of class labels.

$$P(y=c|x) \propto p(x|y=c) P(y=c)$$

\hookrightarrow "proportional to"

$$\log P(y=c|x) = \log[p(x|y=c)] + \log[P(y=c)] - \log[p(x)]$$

constant

$$= \log[p(x|y=c)] + \log[P(y=c)]$$

\hookrightarrow "equal up to a constant"

$$g_c(x) = \log[p(x|y=c)] + \log[P(y=c)]$$

$$N(x; \mu_c, \sigma_c^2)$$

frequency of class #c
in our training data set

$$= \log\left[\frac{1}{\sqrt{2\pi\sigma_c^2}} \cdot \exp\left[-\frac{(x-\mu_c)^2}{2\sigma_c^2}\right]\right] + \log[P(y=c)]$$

$$\mu_c^* = ? \quad \sigma_c^{2*} = ?$$

$$\frac{N_c}{N} = \frac{\sum_{i=1}^N 1(y_i=c)}{N}$$

of data points
that belong to class
#c.

$$\mu_c^* = \frac{\sum_{i=1}^N [x_i \cdot 1(y_i=c)]}{\sum_{i=1}^N [1(y_i=c)]}$$

$$\sigma_c^{2*} = \frac{\sum_{i=1}^N [(x_i - \mu_c^*)^2 \cdot 1(y_i=c)]}{\sum_{i=1}^N [1(y_i=c)]}$$

"one" function $\Rightarrow 1(\cdot) = \begin{cases} 1 & \text{if } \cdot \text{ is TRUE} \\ 0 & \text{otherwise} \end{cases}$

$$\begin{aligned}
 & \left[\begin{array}{c} \mu_1^*, \mu_2^*, \dots, \mu_K^* \\ \sigma_1^{2*}, \sigma_2^{2*}, \dots, \sigma_K^{2*} \\ \hat{P}(y=1), \hat{P}(y=2), \dots, \hat{P}(y=K) \end{array} \right] \Rightarrow K \\
 & \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \Rightarrow K \\
 & \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \qquad \Rightarrow K-1
 \end{aligned}$$

$$\text{Total \# of parameters} = \frac{+}{3K - 1}$$

$$H = 1/2 \quad T = 1/2$$

$$H = 3/4 \quad T = 1/4$$

$$H H H T \rightarrow \text{likelihood} \Rightarrow \frac{1}{2} \frac{1}{2} \frac{1}{2} \frac{1}{2} = \frac{8}{64}$$

$$H H H T \rightarrow \text{likelihood} \Rightarrow \frac{3}{4} \frac{3}{4} \frac{3}{4} \frac{1}{4} = \frac{27}{64}$$