# Density Estimation:

$$x = \{x_i\}_{i=1}^{N} \quad N \text{ samples (data points)}$$

$$x_i \sim p(x) \quad \forall i \quad \Rightarrow \quad \text{probability distribution}$$

$\underbrace{p(x)}_{\text{unknown}}$

parameters (?)

DENSITY ESTIMATION

learn these parameters from training data



$$x_i \sim N(x; \mu, \sigma^2)$$

$\mu^{\#}$: the best $\mu$ parameter

$\sigma^{2\#}$: the best $\sigma^2$ parameter

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

class densities $\Rightarrow P(x \mid y = c) \rightsquigarrow$ density estimation

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

## BAYES RULE

$$P(B \mid A) = \frac{P(A, B)}{P(A)}$$

$$P(B \mid A) = \frac{P(A \mid B) \, P(B)}{P(A)}$$

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

a new data point $\overbrace{x_{N+1}}$

$$\Rightarrow P(y = c \mid x_{N+1}) \quad
\begin{array}{l}
\rightarrow P(y = 1 \mid x_{N+1}) \\
\rightarrow P(y = 2 \mid x_{N+1}) \\
\rightarrow P(y = 3 \mid x_{N+1})
\end{array}
\left.\right\} \text{Pick the maximum}$$

# MAXIMUM LIKELIHOOD ESTIMATION (MLE)

$$X = \{x_i\}_{i=1}^{N} \qquad x_i \sim p(x \mid \theta)$$

$\hookrightarrow$ unknown parameters

$\log(a^b) = b \cdot \log(a)$

$x_i$'s are i.i.d.

identically & independently distributed

$\log(a \cdot b \cdot c) = \log(a) + \log(b) + \log(c)$

$$\text{Likelihood} \equiv p(x_1, x_2, \dots, x_N \mid \theta)$$

$$L(\theta \mid X) \equiv p(x_1 \mid \theta)\, p(x_2 \mid \theta) \dots \dots p(x_N \mid \theta)$$

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

$$\theta^* = \arg\max_{\theta} L(\theta \mid X)$$

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

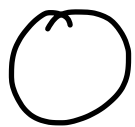$$= \sum_{i=1}^{N} \log\left[ p(x_i \mid \theta) \right]$$

Bernoulli density: $0 < p < 1$

$\hookrightarrow$ success probability

(H) success: $p \Rightarrow x = 1$

(T) failure: $1-p \Rightarrow x = 0$

$\dfrac{\partial \log(1-x)}{\partial x} = -\dfrac{1}{(1-x)}$

$\dfrac{\partial \log(x)}{\partial x} = \dfrac{1}{x}$

$\bigcirc \Rightarrow HTHHHT \cdots\cdots T \quad \Big/ \quad$ 70 heads, 30 tails

$\quad\quad\quad\quad\quad\quad x_1\ x_2\ x_3\ x_4\ x_5\ x_6\ x_7 \quad\quad\quad x_{100}$

$\quad\quad\quad\quad\quad\quad 1\ \ 0\ \ 1\ \ 1\ \ 1\ \ 1\ \ 0 \quad\quad\quad\quad 0$

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

$L(p \mid x) = \prod_{i=1}^{N} \left[ p^{x_i} \cdot (1-p)^{1-x_i} \right]$

$P(x_i = 1 \mid p) = p^1 \cdot (1-p)^{1-1} = p$

$P(x_i = 0 \mid p) = p^0 (1-p)^{1-0} = 1-p$

$\log L(p \mid x) = \sum_{i=1}^{N} \left[ x_i \log(p) + (1-x_i) \log(1-p) \right] \Rightarrow p = ?$

$\dfrac{\partial \log L(p \mid x)}{\partial p} = \sum_{i=1}^{N} \left[ x_i \cdot \dfrac{1}{p} - (1-x_i) \dfrac{1}{1-p} \right] = 0 \Rightarrow \hat{p} = \dfrac{\sum_{i=1}^{M} x_i}{N}$

\# of heads

## Gaussion Density:   $X = \{x_i\}_{i=1}^{N}$

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

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

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

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

$$\frac{\partial \log\text{-likelihood}}{\partial \mu} = 0 \qquad \& \qquad \frac{\partial \log\text{-likelihood}}{\partial \sigma^2} = 0$$

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

# Parametric Classification!

Input: A training dataset

Output: A classifier

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

test (unseen) ↑ data point → $y_i \in \{1, 2, \ldots, K\}$

$$\hat{y}_{N+1} = \arg\max_{c} g_c(X_{N+1})$$

↳ score function for class #c

$$P(y=c \mid x) = \frac{p(x \mid y=c)\, P(y=c)}{p(x)} \Rightarrow \text{independent of class labels}$$

$$P(y=c \mid x) \propto p(x \mid y=c)\, P(y=c)$$

↳ "proportional to"

constant ↗

$$\log P(y=c \mid x) = \log(p(x \mid y=c)) + \log(P(y=c)) - \log(p(x))$$

$$\overset{+}{=} \log(p(x \mid y=c)) + \log(P(y=c))$$

↳ "equal up to a constant"

$$g_c(x) = \log\left(p(x|y=c)\right) + \log\left(P(y=c)\right)$$

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

$$\overbrace{\text{frequency of class} \# c}^{\text{in our data set.}}$$

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

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

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

$$\mu_c^* = \frac{\sum\limits_{i=1}^{N}\left[x_i \cdot 1(y_i=c)\right]}{\sum\limits_{i=1}^{N} 1(y_i=c)} \qquad \sigma_c^{2*} = \frac{\sum\limits_{i=1}^{N}\left[(x_i - \mu_c^*)^2 \cdot 1(y_i=c)\right]}{\sum\limits_{i=1}^{N} 1(y_i=c)}$$

$\hookrightarrow$ # of data points that belong to class # c.

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

$$\mu_1^*, \mu_2^*, \ -------, \mu_K^* \quad \Big\} \ K$$

$$\sigma_1^{2*}, \sigma_2^{2*}, \ ------, \sigma_K^{2*} \quad \Big\} \ K$$

$$\hat{P}(y=1), \hat{P}(y=2), ----, \hat{P}(y=K) \Big\} \ K-1$$

$$\underbrace{\hspace{4cm}}$$

total # of parameters $= 3K - 1$.