# Classification task

Pattern $\longrightarrow$ [ Alg ] $\longrightarrow$ label
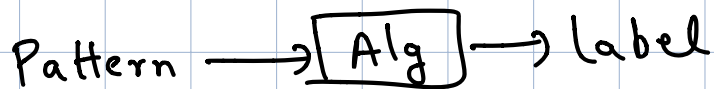
## Notation :

① $\mathcal{X}$ : Feature space = Set of all feature vectors

For e.g. $\mathcal{X} = \mathbb{R}^d$

② Classifier $h : \mathcal{X} \to \{1, ---, M\}$

Simple Case : $M = 2$, $\{0, 1\}$ or $\{-1, +1\}$

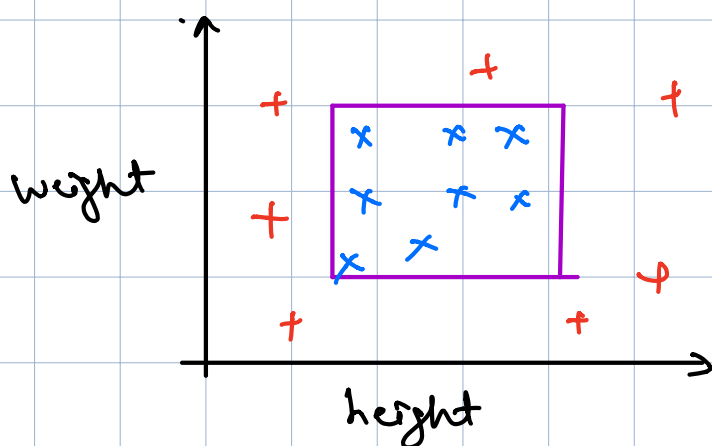How to design classifiers & how to judge their performance?

Input : $\{(x_i, y_i), i = 1 --- n\}$     Training dataset

$\nearrow$ feature vector    $\nwarrow$ class label

Using training data, learn an appropriate classifier $h$

Test : Test & validate the $h$ on "new" data
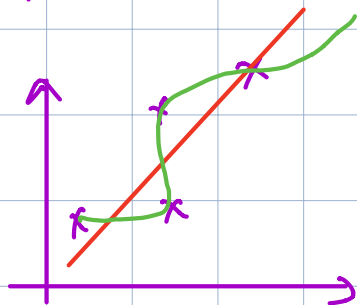
P TO

Example:    Medium-build



Regression task:

   Training data:  $\{(x_i, y_i), i = 1 \cdots n\}$
                    $x_i \in \mathbb{R}^d, \quad y_i \in \mathbb{R}$

   Need is to learn a function  $h: X \to \mathbb{R}$

      Examples: prediction of stock prices, etc

Curve-fitting:        $\{(x_1, y_1), ---- (x_n, y_n)\}$

         Find a function  $f$  s.t.  $y = f(x)$



Generalization:-   Obtain a classifier/regression function

using a training dataset s.t. the error on "unseen"
(test) data is   low.

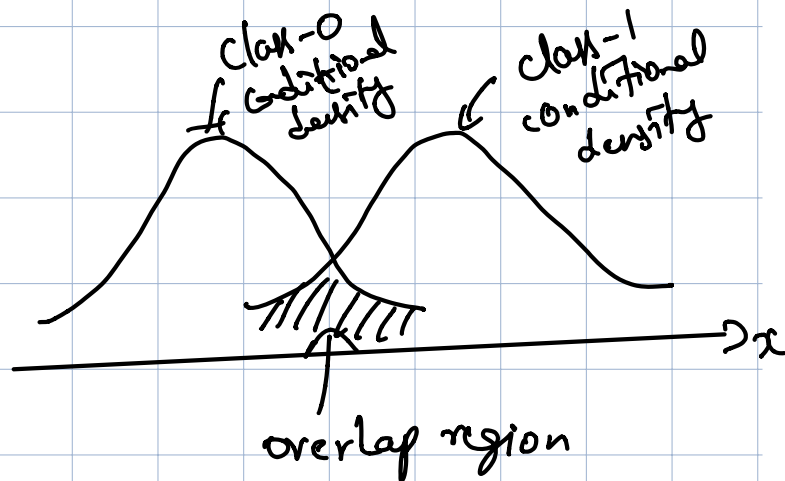Assumption: There is a distribution underlying the data.
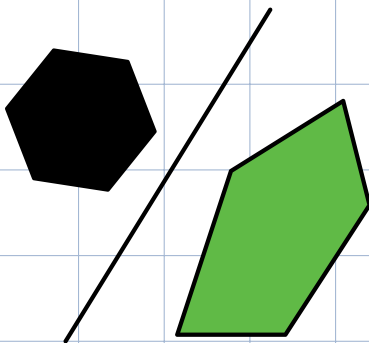
Formally,

$f_i$ : conditional density of features from class-i

Let $x = (x_1, \dots x_d) \in \mathbb{R}^d$ represent a feature vector

$f_i(x)$: joint density of $(x_1, \dots x_d)$ given that $x$ is in class-i

Note! A feature vector $x$ can belong to different classes with different probabilities



class-0 conditional density

class-1 conditional density

$\to x$

overlap region

P TO

$\mathcal{X}$: feature space $\qquad$ $\mathcal{Y} = \{0,1\}$ class-labels

$\mathcal{H} =$ set of classifiers $= \{ h \mid h: \mathcal{X} \to \mathcal{Y} \}$

Performance: $\qquad F(h) = P\left( h(x) \neq y(x) \right)$
metric

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad\nearrow\qquad\qquad\uparrow$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ label assigned $\quad$ true class
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ by classifier-h $\qquad$ label

$\qquad\qquad$ Probability of mis-classification

Goal: $\qquad h^* = \underset{h \in \mathcal{H}}{\arg\min} \; F(h)$

---

Bayes classifier

Let $\qquad P_i = P\left( y(x) = i \right) \leftarrow$ prior probabilities

$\qquad\qquad q_i(x) = P\left( y(x) = i \mid x=x \right)$

Bayes classifier $\qquad h_B(x) = \begin{cases} 0 & \text{if } \dfrac{q_0(x)}{q_1(x)} > 1 \\[2mm] 1 & \text{else} \end{cases}$

$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ P.T.O

From Bayes theorem,

$$\text{posterior} \propto \text{prior} \times \text{likelihood}$$

$$q_i(x) = \frac{p_i \, f_i(x)}{p_0 \, f_0(x) + p_1 \, f_1(x)}$$

$\hookrightarrow$ Normalizing constant

Bayes classifier:

Classify $x$ as "belonging to class "0"

if $\quad q_0(x) > q_1(x)$

$(\Rightarrow) \quad p_0 \, f_0(x) > p_1 \, f_1(x)$

Optimality of Bayes classifier:

Fix a classifier $h$

Define $R_i(h) = \{ x \in X \mid h(x) = i \}, i = 0,1$

$F(h) = P(h(x) \neq y(x))$

$\quad = P(x \in R_1(h), x \in \text{Class-0})$

$$+ P(X \in R_0(h), X \in C(\text{class-1}))$$

$$= P_0 \; P(X \in R_1(h) \mid X \in C(\text{class-0}))$$
$$+ P_1 \; P(X \in R_0(h) \mid X \in C(\text{class-1}))$$

$$= P_0 \int_{R_1(h)} f_0(x)\, dx + P_1 \int_{R_0(h)} f_1(x)\, dx$$
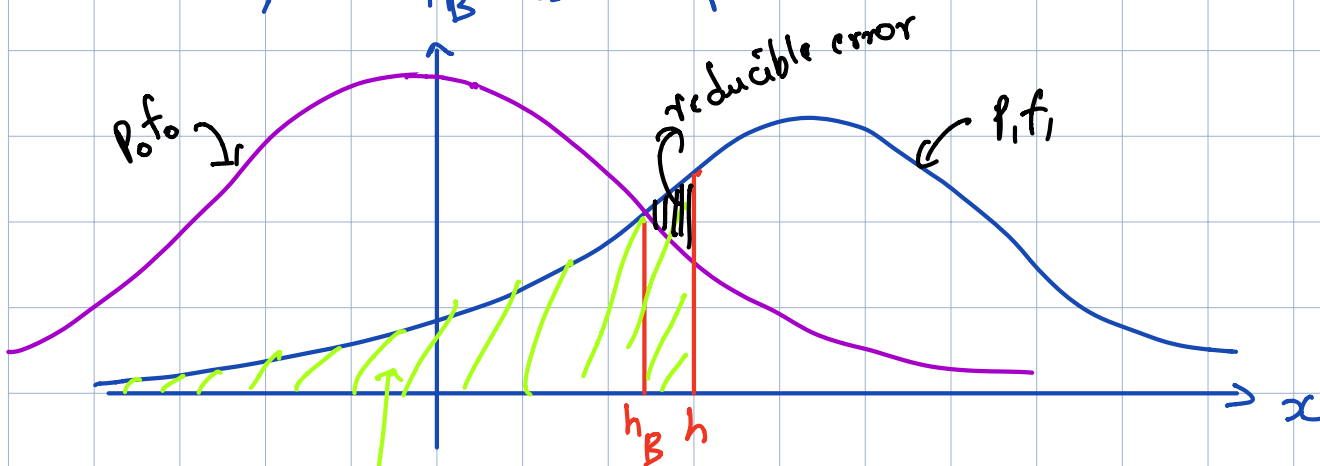
For the Bayes classifier $h_B$:

$$R_0(h_B) = \{ x \mid P_0 f_0(x) \geq P_1 f_1(x) \}$$
$$R_1(h_B) = \{ x \mid P_1 f_1(x) > P_0 f_0(x) \}$$

So, $$F(h_B) = \int_{R_1(h)} P_0 f_0(x)\, dx + \int_{R_0(h)} P_1 f_1(x)\, dx$$

$$= \int_{\mathcal{X}} \min\left( P_0 f_0(x),\; P_1 f_1(x) \right) dx$$

So, $h_B$ is optimal.



$P_0 f_0$    reducible error    $P_1 f_1$

$h_B$  $h$    $x$

$$\int_{R(h)} P_i f_i(x) dx$$

# Generalizing Bayes classifier to handle loss functions:-

Loss function $\quad L: Y \times Y \longrightarrow \mathbb{R}^+$

$L(h(x), y(x))$ is the loss suffered by $h$ on pattern $X$.

Performance metric $\quad F(h) = E\left[ L(h(x), y(x)) \right]$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad \longrightarrow (*)$

Special Case:- $\quad L(a,b) = \begin{cases} 0 & \text{if } a=b \\ 1 & \text{if } a \neq b \end{cases}$

(0-1) loss function

Then, $\quad F(h) = P(h(x) \neq y(x)) \quad$ for above $L$.

In general, $\quad L(0,1) \neq L(1,0) \quad [L(0,0) = L(1,1) = 0]$

$h_B$ for optimizing $(*)$ is

$$h_B(x) = \begin{cases} 0 & \text{if } \dfrac{\varphi_0(x)}{\varphi_1(x)} > \dfrac{L(0,1)}{L(1,0)} \\ \\ 1 & \text{else} \end{cases}$$

Note: If $L(1,0) = L(0,1)$, we recover the Bayes classifier for 0-1 loss function

## Extending Bayes clasifier to multi-class clasification problems!-

Let $\{0, 1, --- M-1\}$ be the M-class labels.

$L(i,j) \rightarrow$ loss function

Clasifier predicts $\nearrow$

true class label

$$R(h) = E\Big(L(h(x), y(x))\Big) \Leftarrow (**)$$

Goal: minimize $R(h)$

## Deriving the Bayes clasifier:

Let
$$R(i|x) = E\Big(L(h(x), y(x)) \,\big|\, h(x)=i, x\Big)$$

$$= E\Big(L(i, y(x)) \,\big|\, - a - \Big)$$

$$= \sum_{j=0}^{M-1} L(i,j)\, P\big(y(x)=j|x\big)$$

$$= \sum_{j=0}^{M-1} L(i,j)\, q_j(x)$$

$$R(h) = E\Big[E\big(L(h(x), y(x))\,|x\big]\Big]$$

$$= \int R(h(x)|x)\, f(x)\, dx$$
$\uparrow$

## Bayes Classifier:

$$h_B(x) = i \quad \text{if} \quad \sum_{j=0}^{M-1} L(i,j) q_j(x) \le \sum_{j=0}^{M-1} L(k,j) q_j(x), \quad \forall k$$

Claim: $h_B$ is optimal for $(**)$

---

Special Case for $M=2$:

$$h_B(x) = 0 \quad \text{if}$$

$$L(0,0) q_0(x) + L(0,1) q_1(x)$$
$$\le L(1,0) q_0(x) + L(1,1) q_1(x) \qquad \text{---} \quad (***)$$

If $L(0,0) = L(1,1) = 0$, then $(***)$ is equivalent to

$$\frac{q_0(x)}{q_1(x)} \ge \frac{L(0,1)}{L(1,0)}$$

H.W.: For 0-1 loss function, with a $M$-class classification problem, derive $h_B$.

Special Case: $X \in R$

$$f_i \sim N(\mu_i, \sigma_i^2), \quad i = 0, 1$$

$$f_i(x) = \frac{1}{\sqrt{2\pi} \, \sigma_i} \exp\left( - \frac{(x - \mu_i)^2}{2\sigma_i^2} \right), \quad i = 0, 1$$

$h_B(x) = 0$ if

$$P_0 \, f_0(x) \, L(1,0) > P_1 \, f_1(x) \, L(0,1)$$

$$\log(P_0 \, L(1,0)) + \log(f_0(x)) > \log(P_1 \, L(0,1)) + \log(f_1(x))$$
$$\overset{\llcorner}{(\cancel{\ast\ast\ast\ast})}$$

$$\log(P_0 \, L(1,0)) - \log(\sigma_0) - \frac{1}{2}\log(2\pi) - \frac{(x - \mu_0)^2}{2\sigma_0^2}$$

$$> \log(P_1 \, L(0,1)) - \log(\sigma_1) - \frac{1}{2}\log(2\pi) - \frac{(x - \mu_1)^2}{2\sigma_1^2}$$

$$\frac{1}{2} x^2 \left( \frac{1}{\sigma_1^2} - \frac{1}{\sigma_0^2} \right) + x \left( \frac{\mu_0}{\sigma_0^2} - \frac{\mu_1}{\sigma_1^2} \right)$$

$$+ \frac{1}{2} \left( \frac{\mu_1^2}{\sigma_1^2} - \frac{\mu_0^2}{\sigma_0^2} \right) + \log\left( \frac{\sigma_1}{\sigma_0} \right) + \log\left( \frac{P_0 \, L(1,0)}{P_1 \, L(0,1)} \right)$$
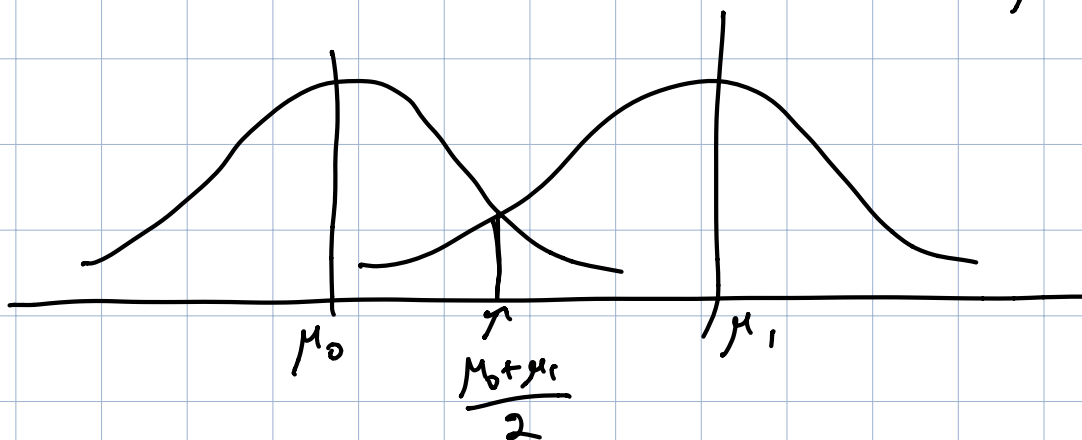
$$> 0$$

Special Cases:

① $\sigma_0 = \sigma_1 = \sigma$, $P_0 = P_1$, ⟵ equi-probable classes $\quad L(1,0) = L(0,1)$

Then, $\quad h_B(x) = 0$ if

$$\frac{x}{\sigma^2} \left( \mu_0 - \mu_1 \right) - \frac{1}{2\sigma^2} \left( \mu_0^2 - \mu_1^2 \right) > 0$$

i.e., $\quad x > \dfrac{\mu_0 + \mu_1}{2}$ $\qquad \left( \begin{array}{c} \text{assuming} \\ \mu_0 > \mu_1 \end{array} \right)$



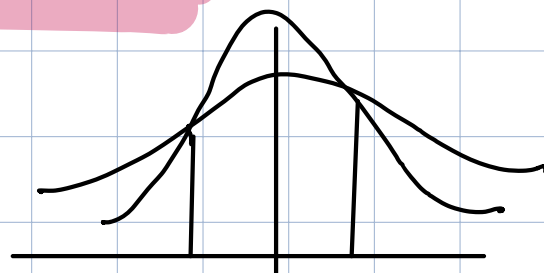$\mu_0 \qquad \dfrac{\mu_0 + \mu_1}{2} \qquad \mu_1$

② $\mu_0 = \mu_1 = 0$, $P_0 = P_1$, $\quad L(1,0) = L(0,1)$

$h_B(x) = 0$ if

$$\frac{x^2}{2} \left( \frac{1}{\sigma_1^2} - \frac{1}{\sigma_0^2} \right) - \log \left( \frac{\sigma_0}{\sigma_1} \right) > 0$$

$\sigma_0 > \sigma_1 \implies$ $\qquad \dfrac{x^2}{2} > \dfrac{\sigma_1^2 \sigma_0^2 \log \left( \sigma_0 / \sigma_1 \right)}{\left( \sigma_0^2 - \sigma_1^2 \right)}$

③ Generalization to multivariate class-conditional densities:

$$f_i(x) = \frac{1}{\sqrt{(2\pi)^d |\Sigma_i|}} \exp\left( -\frac{1}{2}(x-\mu_i)^T \Sigma_i^{-1}(x-\mu_i) \right), \quad i=0,1$$

$h_B(x) = 0$   if

CHECK!

$$\frac{1}{2} x^T \left( \Sigma_1^{-1} - \Sigma_0^{-1} \right) x + x^T \left( \Sigma_0^{-1}\mu_0 - \Sigma_1^{-1}\mu_1 \right)$$

$$+ \frac{1}{2}\left( \mu_1^T \Sigma_1^{-1}\mu_1 - \mu_0^T \Sigma_0^{-1}\mu_0 \right)$$

$$+ \log\left( \frac{p_0 \, L(1,0)}{p_1 \, L(0,1)} \right) + \frac{1}{2} \log\left( \frac{|\Sigma_1|}{|\Sigma_0|} \right) > 0$$
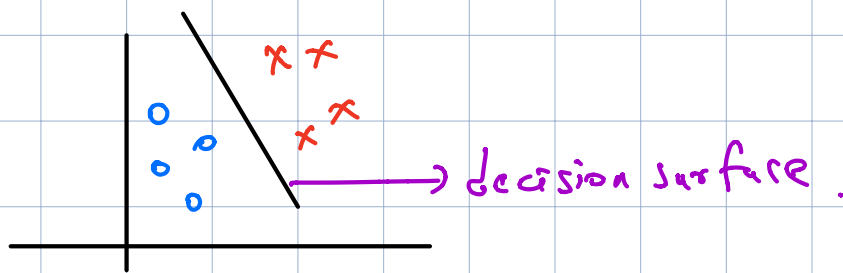
H.W: Work out the special cases as in the univariate setting.

---

## DISCRIMINANT FUNCTIONS:

$$h(x) = \begin{cases} 0 & \text{if } g(x) > 0 \\ 1 & \text{else} \end{cases}$$

Example: 0-1 loss function. Bayes classifier is based on the discriminant function $g(x) = q_0(x) - q_1(x)$

$g(x) = 0$ is the decision surface.



$\longrightarrow$ decision surface.

# Maximum likelihood estimation

① To implement Bayes classifier, we need class conditional densities + prior probabilities

② Suppose we are given i.i.d. samples from a class conditional distribution

$$\{ x_1, \dots x_n \}$$

③ We adopt a "parametric" approach to estimation

$\mathcal{D} = \{ x_1, \dots x_n \}$ iid from the distribution of r.v. $X$ parameterized by $\Theta$.

Aim: Use $\mathcal{D}$ to estimate $\Theta$

Example:-  $f(x|\Theta) \sim N(\Theta_1, \Theta_2)$

$(\Theta_1, \Theta_2) \longrightarrow$ unknown

$$f(x|\theta) = \frac{1}{\sqrt{2\pi} \, \theta_2} \exp\left(-\frac{(x-\theta_1)^2}{2\theta_2^2}\right)$$

$\hat{\theta}_n \leftarrow$ estimate of $\theta$.

**Example:-** Suppose $X$ is drawn uniformly at random from the set $\{1, ---, n\}$

Suppose "$n$" is unknown

Suppose we observe a sample "$k$".

Want an estimator of $n$ given $X = k$ is observed.

## Max-Likelihood (ML) idea:

The pmf of $X$ $\quad P_x(k) = \frac{1}{n} I\{1 \leq k \leq n\}$

Think of $P_x(k)$ as a function of $n$.

& find the maximizer of this function

$P_x(k)$ is $0$ for $n \leq k-1$, jumps to $\frac{1}{k}$ at $n=k$ & decreases beyond $k$. So, it is maximized at $n=k$.

$$\hat{n}_{MC}(k) = k.$$

---

Suppose we are given $\{x_1, \ldots, x_n\}$ iid $\sim f(\cdot | \theta)$

E.g. $x_i \sim N(\underset{\underset{\text{unknown}}{\uparrow}}{\theta}, 1) \in$ One-parameter distribution, $i = 1 \sim n$

$\hat{\theta}_n \leftarrow$ estimate of $\theta$

Sample mean $\rightarrow \hat{\theta}_n = \frac{1}{n} \sum_{i=1}^{n} x_i$, $\hat{\theta}_2 = \frac{x_1 + x_2}{2}$, $\hat{\theta}_1 = x_1$

Why choose $\hat{\theta}_n$ over $\hat{\theta}_2$ & $\hat{\theta}_1$?

Use a "mean-square error" objective

$$MSE(\hat{\theta}) = E\left( (\hat{\theta} - \theta)^2 \right)$$

$$= E\left[ \left( \hat{\theta} - E[\hat{\theta}] \right) + \left( E[\hat{\theta}] - \theta \right) \right]^2$$

$$= E\left[ (\hat{\theta} - E\hat{\theta})^2 \right] + \left( E(\hat{\theta}) - \theta \right)^2 + 2 \underbrace{E\left[ (\hat{\theta} - E\hat{\theta})(E\hat{\theta} - \theta) \right]}_{= (E\hat{\theta} - \theta) E(\hat{\theta} - E\hat{\theta}) = 0}$$

$$= E\left[(\hat{\theta} - E\hat{\theta})^2\right] + (E\hat{\theta} - \theta)^2$$

$$\underset{\text{Variance}}{\nearrow} \qquad \underset{\text{bias}}{\nearrow}$$
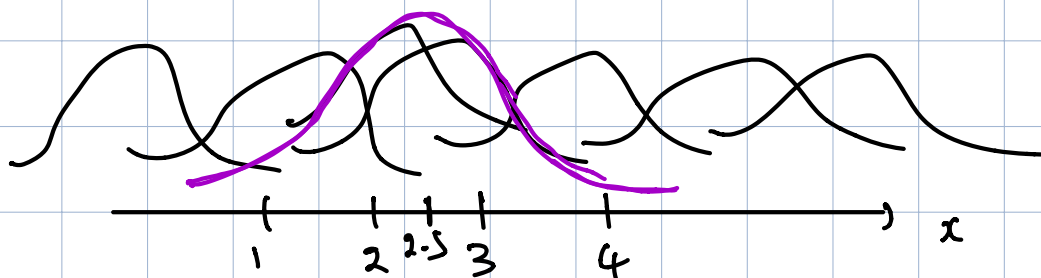
Estimator $\hat{\theta}$ unbiased if $E\hat{\theta} = \theta$

$$Var(\hat{\theta}_n) = \frac{1}{n}, \qquad Var(\hat{\theta}_2) = \frac{1}{2}, \qquad Var(\hat{\theta}_1) = 1$$

---

### Max- likelihood estimation:

$$\mathcal{D} = \{x_1 \text{---} x_n\}$$

E.g. $x_i \sim N(\theta, 1)$



Likelihood function $L(\theta) = \prod_{j=1}^{n} f(x_j | \theta)$

ML estimate: $\theta^* \in \underset{\theta}{arg\,max} \{L(\theta) = L(x_1 \text{--} x_n, \theta)\}$

$$L(\theta^*) \geq L(\theta) \quad \forall \theta$$

Log- likelihood $\ell(\theta) = Log\, L(\theta) = \sum_{j=1}^{n} log\, f(x_j | \theta)$

**Example:** Univariate normal
$$x_i \sim N(\theta_1, \theta_2^2), \quad i = 1 \cdots n$$

$$\ell(\theta) = \log L(\theta) = \sum_{j=1}^{n} \left( -\log \theta_2 - \frac{1}{2} \log 2\pi - \frac{(x_j - \theta_1)^2}{2\theta_2^2} \right)$$

$$\ell(\theta) = -n \log \theta_2 - \frac{n}{2} \log 2\pi - \sum_{j=1}^{n} \frac{(x_j - \theta_1)^2}{2\theta_2^2}$$

To find ML estimates $\hat{\theta}_1, \hat{\theta}_2$ of $\theta_1$ & $\theta_2$, do

(i) $\dfrac{\partial \ell}{\partial \theta_1} = 0$    &    (ii) $\dfrac{\partial \ell}{\partial \theta_2} = 0$

$$\hat{\theta}_1 = \frac{1}{n} \sum_{j=1}^{n} x_j$$

$$\hat{\theta}_2^2 = \frac{1}{n} \sum_{j=1}^{n} (x_j - \hat{\theta}_1)^2$$

↗ *not unbiased*

**H.W.:**  ① Assume $Exp(\lambda)$ likelihood & find the ML estimate

② Generalize to multivariate case
$$x_i \sim N(\mu, \Sigma), \quad i = 1 \cdots n, \quad x_i \in \mathbb{R}^d$$

Calculate ML estimates for $\mu, \Sigma$.
$$\hat{\mu}_n = \frac{1}{n} \sum_{i=1}^{n} x_i ; \quad \hat{\Sigma}_n = \frac{1}{n} \sum_{i=1}^{n} (x_i - \hat{\mu}_n)(x_i - \hat{\mu}_n)^T$$

$$\left(\text{Recall} \quad \Sigma = E(x-\mu)(x-\mu)^T\right.$$

Example:- Suppose a coin with bias $\theta$ is tossed $n$ times & $k$ heads occur.

$$P_{\hat{\theta}}(X=k) = \binom{n}{k}(\hat{\theta})^k((1-\hat{\theta})^{n-k} = L(\hat{\theta})$$

$$\hat{\theta}_{MC} = \underset{\hat{\theta}}{\arg\max} \; (\hat{\theta})^k ((1-\hat{\theta})^{n-k}, \quad n,k \text{ known}$$

Let $\quad 1 \le k \le n-1$

$$\frac{d}{d\hat{\theta}}\left(\hat{\theta}^k(1-\hat{\theta})^{n-k}\right) = \frac{k\hat{\theta}^k}{\hat{\theta}}\left(1-\hat{\theta}\right)^{n-k} - \frac{(n-k)}{(1-\hat{\theta})}\hat{\theta}^k(1-\hat{\theta})^{n-k}$$

$$= \left(\frac{k}{\hat{\theta}} - \frac{(n-k)}{(1-\hat{\theta})}\right)\hat{\theta}^k(1-\hat{\theta})^{n-k}$$

$$= \left(k-n\hat{\theta}\right)\hat{\theta}^{k-1}(1-\hat{\theta})^{n-k-1} \underline{\quad\quad}(\ast)$$

If $\quad \hat{\theta} \le \frac{k}{n}$, then $\quad (\ast) \geqslant 0$

if $\quad \hat{\theta} > \frac{k}{n}$, then $\quad (\ast) < 0$

So, $\quad \hat{\theta}_{MC} = \frac{k}{n}$

for $k=0$, $\quad L(\hat{\theta}) = (1-\hat{\theta})^n$ & $\quad \hat{\theta}_{MC} = 0$

for $k=n$, $\quad L(\hat{\theta}) = \hat{\theta}^n$ & $\quad \hat{\theta}_{MC} = 1$

# Bayesian estimation

| Frequentists | Bayesian |
|---|---|
| Probability = long run frequency | Probability = degree of belief (subjective) |
| Goal = algorithms with long-run freq. guarantees | Goal:- State & analyze beliefs |

### Bayesian approach:-

Parameter $\theta$ is given a prior density $\pi(\theta)$

Observe data $\{x_1, \dots x_n\}$ sampled from $f(x|\theta)$

Compute posterior density (Use Bayes rule)

$$f(\theta| x_1 \dots x_n) \propto L(\theta)\, \pi(\theta)$$

Example: $\mathcal{D}_n = \{x_1, \dots x_n\}$
$$x_i \sim Ber(\theta)$$

Uniform prior $\pi(\theta) = 1$

Posterior probability

$$f(\theta \mid \mathcal{D}_n) \propto \underset{\text{prior}}{\pi(\theta)} \; \underset{\text{likelihood}}{\mathcal{L}_n(\theta)}$$

$$\propto 1 \times \prod_{i=1}^{n} f(x_i \mid \theta)$$

$$\propto \prod_{i=1}^{n} \theta^{x_i} (1-\theta)^{1-x_i}$$

Denote $S_n = \sum_{i=1}^{n} x_i$

$$\propto \theta^{S_n} (1-\theta)^{n-S_n}$$

$$\propto \theta^{(S_n+1)-1} (1-\theta)^{(n-S_n+1)-1}$$

Beta density: $\pi_{\alpha,\beta}(\theta) = \dfrac{\overline{\alpha+\beta}}{\overline{\alpha}\,\overline{\beta}} \; \theta^{\alpha-1}(1-\theta)^{\beta-1}$

So, $f(\theta \mid \mathcal{D}_n) \sim \text{Beta}(S_n+1,\ n-S_n+1)$

Posterior - mean $= \dfrac{S_n+1}{n+2}$

H.W. Redo the Calculation above for $\pi(\theta) \sim \text{Beta}(\alpha,\beta)$

# Bayes estimation

Working procedure:

① Choose a prior density $\pi(\theta)$

② Observe data $\mathcal{D}_n = \{x_1, \dots x_n\}$ sampled iid

from $f(x|\theta)$ ← likelihood

③ Update beliefs:

Calculate posterior density

Bayes rule

$$f\left(\theta \mid \underbrace{x_1, \dots, x_n}_{\mathcal{D}_n}\right) = \frac{f(x_1 \dots x_n | \theta)\, \pi(\theta)}{f(x_1 \dots x_n)}$$

$$= \frac{L_n(\theta)\, \pi(\theta)}{\underset{\longrightarrow}{} \; C_n}$$

Normalization
Constant

$$\underset{\text{posterior}}{f(\theta | x_1 \dots x_n)} \quad \alpha \quad \underset{\text{likelihood}}{L_n(\theta)} \times \underset{\text{prior}}{\pi(\theta)}$$

Point estimates

① Posterior mean $\bar{\theta}_n = \int \theta\, f(\theta | \mathcal{D}_n)\, d\theta$

② Maximum aposteriori probability (MAP)

= mode of the posterior

Example:
$$\mathcal{D}_n = \{x_1, \dots x_n\}$$
$$x_i \sim Ber(\theta)$$

Beta prior ie., $\pi(\theta) \sim Beta(\alpha, \beta)$

Then,

$$f(\theta|\mathcal{D}_n) \propto \theta^{\alpha-1} (1-\theta)^{\beta-1} \theta^{S_n} (1-\theta)^{n-S_n}$$

$$\propto \theta^{\alpha+S_n-1} (1-\theta)^{\beta+n-S_n-1}$$

$$\Rightarrow \quad \theta|\mathcal{D}_n \sim Beta(\alpha+S_n, \beta+n-S_n)$$

Special case: $\alpha = \beta = 1$ leads to uniform (flat) prior

Posterior mean $\overline{\theta}_n = \dfrac{\alpha+S_n}{\alpha+\beta+n}$

$$= \left(\underbrace{\frac{n}{\alpha+\beta+n}}_{}\right)\left(\underbrace{\frac{S_n}{n}}_{\text{ML-ultimate}}\right) + \left(\frac{\alpha+\beta}{\alpha+\beta+n}\right)\left(\underbrace{\frac{\alpha}{\alpha+\beta}}_{\substack{\text{Prior} \\ \text{mean}}}\right)$$

— Sum to 1

Note: Conjugacy: Prior & posterior belong to the same family of parameterized distributions.
e.g. Bernoulli case: Beta prior $\Rightarrow$ conjugacy.

Example: $\mathcal{D}_n = \{x_1, --- x_n\}$

$$f(x|\theta) \sim N(\theta, \sigma^2)$$

$\nwarrow$ known

Prior $\quad \pi(\theta) \sim N(\theta_0, \sigma_0^2)$

$$f(\theta \mid \mathcal{D}_n) \propto \left[ \prod_{i=1}^{n} \exp\left(-\frac{1}{2}\left(\frac{x_i - \theta}{\sigma}\right)^2\right) \right] \exp\left(-\frac{1}{2}\left(\frac{\theta - \theta_0}{\sigma_0}\right)^2\right)$$

$$\propto \exp\left(-\frac{1}{2}\left[ \underbrace{\sum_{i=1}^{n} \left(\frac{x_i - \theta}{\sigma}\right)^2 + \left(\frac{\theta - \theta_0}{\sigma_0}\right)^2}_{\text{quadratic in } \theta} \right]\right) \qquad (*)$$

So, $\quad f(\theta \mid \mathcal{D}_n) \sim N(\theta_n, \sigma_n^2)$

$$f(\theta \mid \mathcal{D}_n) \propto \exp\left(-\frac{1}{2}\left(\frac{\theta - \theta_n}{\sigma_n}\right)^2\right) \qquad (**)$$

Equating the co-efficients of $\theta$ & $\theta^2$ in $(*)$ & $(**)$

$$\frac{1}{\sigma_n^2} = \left(\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}\right) \quad \Leftarrow \text{ coeff - of } \theta^2$$

$$\frac{\theta_n}{\sigma_n^2} = \frac{n}{\sigma^2}\frac{S_n}{n} + \frac{\theta_0}{\sigma_0^2} \quad \Leftarrow \text{ co-eff of } \theta$$

$$\Rightarrow \quad \theta_n = \left(\frac{n\sigma_0^2}{n\sigma_0^2 + \sigma^2}\right)\left(\frac{S_n}{n}\right) + \left(\frac{\sigma^2}{n\sigma_0^2 + \sigma^2}\right)\theta_0$$

$$\sigma_n^2 = \frac{\sigma_0^2 \sigma^2}{n\sigma_0^2 + \sigma^2}$$

Some observations:-

① $\Theta_n$ lies between $\frac{S_n}{n}$ & $\Theta_0$

② $\sigma_0 \neq 0 \implies \Theta_n \to \frac{S_n}{n}$ as $n \to \infty$

③ $\sigma_0 \to 0 \implies$ Prior dominates

④ $\sigma_0 >> \sigma \implies \Theta_n \approx \frac{S_n}{n}$

---

Ref:  Duda Hart Stork's book for Bayes Classifier &
       ML/Bayes estimation

---

Minimum mean square error estimation

Ref: Sec 4.9 of
Bruce Hajek's notes
on probability
(ECE 313 UIUC)

I Constant estimators:

Y is a r.v. & we wish to estimate Y
by a constant $\delta$.

$$MSE(\delta) = E(Y-\delta)^2$$

$$= \int_{-\infty}^{\infty} (y-\delta)^2 f_Y(y) \, dy$$

It is easy to see that $\delta^* = EY$ minimizes
the MSE.

So, $EY$ is the MMSE estimate.

II  Unconstrained estimators

$Y$ is a r.v. & we observe $X$.

$$MSE = E(Y - g(X))^2$$

↑ Want to find the $g$ that minimizes the MSE

Suppose you observe $X=10$.

$$E[Y \mid X=10] = \int_{-\infty}^{\infty} y \, f_{Y|X}(y|x) \, dy$$

Claim:  $E[Y \mid X=x]$ is the MMSE estimate

$$MSE = E[(Y - g(X))^2]$$

$$= \int_{-\infty}^{\infty} \left[ \int_{-\infty}^{\infty} (y - g(x))^2 f_{Y|X}(y|x) \, dy \right] f_X(x) \, dx$$

↓

minimized by $g^*(x) = E[Y \mid X=x]$

$$MMSE = E\left[\left(Y - E[Y|x]\right)^2\right]$$

show
that
$$\overset{\text{\textdownarrow}}{=} \quad E[Y^2] - E\left(\left[E[Y|x]\right]^2\right)$$

---

## Some problems

① Bayes estimation, Gaussian : unknown mean
known variance $\sigma^2$

**Aim:** Sequential estimation of posterior mean & variance

From class notes (see above),

① $\quad \dfrac{1}{\sigma_n^2} = \dfrac{n}{\sigma^2} + \dfrac{1}{\sigma_0^2}$

$\quad\quad = \dfrac{1}{\sigma^2} + \dfrac{n-1}{\sigma^2} + \dfrac{1}{\sigma_0^2}$

$\quad\quad = \dfrac{1}{\sigma^2} + \dfrac{1}{\sigma_{n-1}^2}$

② $\quad \dfrac{\theta_n}{\sigma_n^2} = \dfrac{S_n}{\sigma^2} + \dfrac{\theta_0}{\sigma_0^2}$

$\quad\quad = \dfrac{S_{n-1}}{\sigma^2} + \dfrac{x_n}{\sigma^2} + \dfrac{\theta_0}{\sigma_0^2}$

$$= \frac{\theta_{n-1}}{\sigma_{n-1}^2} + \frac{x_n}{\sigma^2}$$

**Pb 2)** $\quad x_1 = \begin{pmatrix} -1 \\ -1 \end{pmatrix} \quad x_2 = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad x_3 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$

$$\bar{x} = \begin{pmatrix} 0 \\ 0 \end{pmatrix}$$

$$C = \frac{1}{3} \sum_{i=1}^{3} x_i x_i^T = \frac{2}{3} \begin{bmatrix} 1 & 1 \\ 1 & 1 \end{bmatrix}$$

$u_1$

Top eigenvalue $= \frac{4}{3}$ & corresponding eigenvector $= \frac{1}{\sqrt{2}} \begin{bmatrix} 1 \\ 1 \end{bmatrix}$

$$\tilde{x}_i = (x_i^T u_1) u_1 + \underbrace{(\bar{x}^T u)}_{0} u_2 \quad 0$$

$$\tilde{x}_1 = \begin{pmatrix} -1 \\ -1 \end{pmatrix} \quad \tilde{x}_2 = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad \tilde{x}_3 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$$

**Pb 3)** Linear estimator

$$MSE = E\left[ (Y - (aX+b))^2 \right]$$

$$= E\left[ ((Y-aX) - b)^2 \right]$$

So, for a given $a$, $\quad b = E(Y-aX)$

$$= \mu_Y - a\mu_X$$

The estimator has the form $aX + \mu_Y - a\mu_X$

$$MSE = E\left[\left(Y - \mu_y - aX + a\mu_x\right)^2\right]$$

$$= E\left[\left[(Y - aX) - (\mu_y - a\mu_x)\right]^2\right]$$

$$= Var\left(Y - aX\right)$$

$$= Cov\left(Y - aX, Y - aX\right)$$

$$= Var(Y) - 2a\,Cov(Y,X) + a^2\,Var(X)$$
$$\underbrace{\qquad}_{(*)}$$

Differentiate $(*)$ wrt $a$ to obtain

$$a^* = \frac{Cov(Y,X)}{Var(X)}$$

Best linear estimator $= L^*(X) = a^* X + b^*$

$$= a^* X + \mu_y - a^* \mu_x$$

$$= \left(\frac{Cov(Y,X)}{Var\,X}\right)(X - \mu_x) + \mu_y$$

$$= \mu_y + \sigma_y\,\rho_{x,y}\left(\frac{X - \mu_x}{\sigma_x}\right)$$

$$MMSE = Var\left(Y - a^* X\right)$$

$$= \sigma_y^2 - \left(\frac{\sigma_y^2 \rho_{xy}^2}{\sigma_x^2}\right) \cdot \sigma_x^2$$

$$= \sigma_y^2\left(1 - \rho_{xy}^2\right)$$

**Pb 4)** ML estimation for Geometric $(\theta)$

$$L(\theta) = \prod_{i=1}^{n} (1-\theta)^{x_i-1} \theta \qquad \{x_1 - - x_n\}$$

$$S_n = \sum_{i=1}^{n} x_i$$

$$= (1-\theta)^{S_n-n} \theta^n$$

$$L'(\theta) = - (1-\theta)^{S_n-n-1}(S_n-n)\theta^n$$

$$+ (1-\theta)^{S_n-n} n\theta^{n-1}$$

$$= (1-\theta)^{S_n-n-1} \theta^{n-1} \left[ - (S_n-n)\theta + n(1-\theta) \right]$$

$$= (1-\theta)^{S_n-n-1} \theta^{n-1} \left[ n - S_n\theta \right]$$

$$L' > 0 \quad \text{if} \quad n > S_n\theta \quad \text{i.e.,} \quad \theta < \frac{1}{\left(\frac{S_n}{n}\right)}$$

$$= 0 \quad \text{if} \quad \theta = \frac{1}{\left(\frac{S_n}{n}\right)}$$

$$< 0 \quad \text{if} \quad \theta > \frac{1}{\left(\frac{S_n}{n}\right)}$$

So, $\qquad \hat{\theta}_{ML} = \frac{1}{\left(\frac{S_n}{n}\right)}$

---

**Pb 5)** Intuitive justification for ML estimate

of variance being biased

$$E\left(\frac{1}{n}\sum_{i=1}^{n}(x_i-\mu)^2\right) = \sigma^2$$

Recall $\quad \hat{\sigma}_{mc}^2 = \frac{1}{n}\sum_{i=1}^{n}(x_i-\hat{\mu}_n)^2$

Observe that

$$\sum_{i=1}^{n}(x_i-\hat{\mu}_n)^2 \leq \sum_{i=1}^{n}(x_i-\mu)^2 \quad - \quad (\divideontimes)$$

Why?

$$RHS = \sum_{i=1}^{n}\left((x_i-\hat{\mu}_n) + (\hat{\mu}_n-\mu)\right)^2$$

Cross term vanishes

$$= \sum_{i=1}^{n}(x_i-\hat{\mu}_n)^2 + \sum_{i=1}^{n}(\hat{\mu}_n-\mu)^2 + \quad 0$$

$$\geq \sum_{i=1}^{n}(x_i-\hat{\mu}_n)^2$$

from $(\divideontimes)$, $\quad E\left(\frac{1}{n}\sum_{i=1}^{n}(x_i-\hat{\mu}_n)^2\right) \leq \sigma^2$