

# SUMMARY

- RECAP ON PROBABILITY: DISTRIBUTIONS,  $E[\cdot], V[\cdot]$ , see  $\rightsquigarrow$

STATISTICAL METHODS

ex.

RANDOM EXPERIMENT: "  $n = 10$  RANDOM PERSON WITH REPLACEMENT AND COUNT HOW MANY BORN IN ITALY" =  $X$

$$\cdot X \sim \text{Bin}(n=10, p = \frac{24}{36} = \frac{2}{3})$$

$$\rightarrow P(X=3) = \binom{10}{3} \left(\frac{2}{3}\right)^3 \left(\frac{1}{3}\right)^{10-3}$$

- IN PYTHON: # 10 random number in [1, 37)  
`np.random.choice(range(1,37), size=10)`

BORN	COUNT
ITALY	24
ABROAD	12
TOT	36

ex. MULTINOMIAL DISTRIBUTION

BORN	COUNT
P PIEDMONT	10
I ITALY NOT PIEDMONT	14
E EU NOT ITALY	5
A OUT OF EU	7
TOT	36

$X_1 = \#$  BORN IN PIEDMONT  
 $X_2 = \#$  BORN IN ITALY NOT PIEDMONT  
 $X_3 = \#$  BORN IN EU NOT ITALY  
 $X_4 = \#$  BORN OUT OF EU

$$\rightarrow P(X_1=2, X_2=3, X_3=0, X_4=5) = ?$$

... BUT FIRST, 1 POSSIBLE DISPOSITION

$$\begin{aligned}
 &P(" \text{SEEDING } P, I, I, A, A, A, I, P, A, A") = \\
 &= \frac{10}{36} \cdot \frac{14}{36} \cdot \frac{14}{36} \cdot \frac{7}{36} \cdot \frac{7}{36} \cdot \frac{7}{36} \cdot \frac{14}{36} \cdot \frac{70}{36} \cdot \frac{7}{36} \cdot \frac{7}{36}
 \end{aligned}$$

BUT  $\forall$  SYMBOL  $\rightarrow$  MULTIPLE COMBINATIONS: 2P, 3I, 0E, 5A

$$\rightarrow P(X_1=2, X_2=3, X_3=0, X_4=5) = \binom{10}{2} \binom{8}{3} \binom{5}{0} \binom{5}{5} \cdot \left(\frac{10}{36}\right)^2 \left(\frac{14}{36}\right)^3 \left(\frac{7}{36}\right)^0 \left(\frac{7}{36}\right)^5 =$$

$$= \frac{10!}{2! 3! 0! 5!} \cdot \left(\frac{10}{36}\right)^2 \left(\frac{14}{36}\right)^3 \left(\frac{7}{36}\right)^0 \left(\frac{7}{36}\right)^5$$

- MULTINOMIAL DISTRIBUTION:  $\rightarrow$  GENERALIZATION OF BINOMIAL EXPERIMENT: 1 CLASS  $\longmapsto$   $k+1$  CLASSES

UNION  $K+1$  CLASSES AS POSSIBLE OUTCOMES, IF WE MAKE  $n$  CHOICES

INDEPENDENTLY WITH THE SAME CLASS PROBABILITIES  $\pi_1, \dots, \pi_{K+1}$  /  $\prod_{k=1}^{K+1} = 1 - \sum_{j=1}^K \pi_j$

$$\rightarrow P(X_1 = x_1, \dots, X_{K+1} = x_{K+1}) = \frac{n!}{x_1! \dots x_{K+1}!} \cdot \prod_{i=1}^K \pi_i^{x_i} \quad / \quad n = \sum_{j=1}^K x_j + x_{K+1}$$

R.V. OUTCOME  
MULTINOMIAL RANDOM VECTOR

MULTINOMIAL COEFFICIENT

$\rightarrow$  ex. BINOMIAL:  $K=1 \rightarrow P, (1-p)$

- $\left\{ \begin{array}{l} K+1: n^{\circ} \text{ OF MUTUALLY EXCLUSIVE EVENTS } (P(A \cap B) = 0) \\ n: n^{\circ} \text{ OF TRIALS} \end{array} \right\} / K: n^{\circ} \text{ R.V.s}$

ex.

UNION 3 POSSIBLE DEFECT A, B, C /  $X_i = \# \text{ OF } i \text{ DEFECTIVE OUTCOME}$

$\rightarrow$  WHEN RANDOM SAMPLING  $\rightarrow$  MULTINOMIAL  $(n, \pi_A, \pi_B)$

- RANDOM VECTORS:

$$X = \begin{pmatrix} X_1 \\ \vdots \\ X_d \end{pmatrix} / E(X) = \begin{pmatrix} E(X_1) \\ \vdots \\ E(X_d) \end{pmatrix} \rightsquigarrow d \times 1 \text{ VECTOR}$$

- $V(X_i)$  AND  $Cov(X_i, X_j)$  ARE COLLECTED IN A  $d \times d$  MATRIX  
CALLED "VARIANCE-COVARIANCE MATRIX"

$$Var(Cov(\bar{X})) = \begin{pmatrix} V(X_1) & Cov(X_1, X_2) & \dots & Cov(X_1, X_d) \\ Cov(X_1, X_2) & V(X_2) & & \\ \vdots & & \ddots & \\ Cov(X_1, X_d) & & & V(X_d) \end{pmatrix}$$

$$\cdot P(\bar{x} \in A) = \iint \dots \int_A f(x_1, \dots, x_d) dx$$

es.

IN THE PLACE OF BIRTH EXAMPLE, P OF SEEING NO ITACONS AT ALL?

$$\cdot P(X_1=0, X_2=0, X_3=\text{ANY}) = f(0,0,0) + f(0,0,1) + \dots + f(0,0,10)$$

- PROPERTY: THE MARGINAL DISTRIBUTION OF A SUBSET OF COMPONENTS IS AGAIN MULTINOMIAL

$$\hookrightarrow \begin{pmatrix} X_1 \\ \vdots \\ X_d \end{pmatrix} \sim \text{MULTINOMIAL}(n, \bar{\pi}_1, \dots, \bar{\pi}_d) \quad / \quad X_1 \sim \text{Binom}(n, \bar{\pi}_1)$$

$$\rightarrow \begin{pmatrix} X_1 \\ X_2 \end{pmatrix} \sim \text{MULTINOMIAL}(n, \bar{\pi}_1, \bar{\pi}_2)$$

- THE DISTRIBUTION OF PARTIAL SUMS OF COMPONENTS ARE MULTINOMIAL;

$$\rightarrow (X_1, X_2 + \dots + X_d) \sim \text{MULTINOMIAL}(n, \bar{\pi}_1, \bar{\pi}_2 + \dots + \bar{\pi}_d)$$

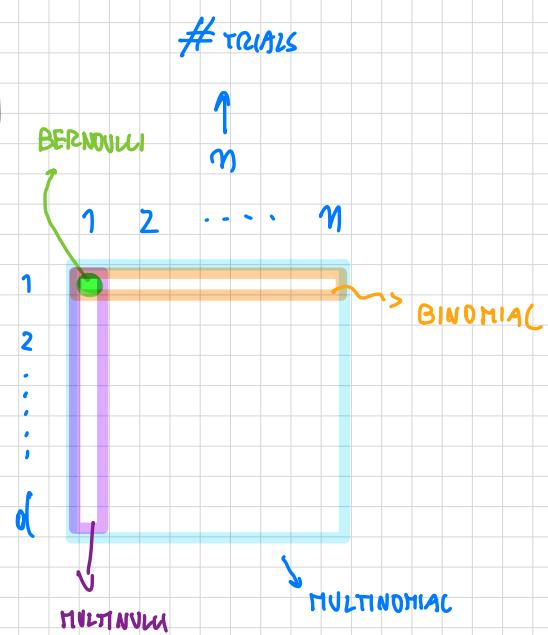
- IF  $n=1 \rightarrow \underline{\text{MULTINOMIAL}}$ :

$$\cdot X \sim \text{MULINUML}(\bar{\pi}_1, \dots, \bar{\pi}_d)$$

$$\cdot E(\bar{X}) = \begin{pmatrix} E(X_1) \\ \vdots \\ E(X_d) \end{pmatrix} = \begin{pmatrix} \bar{\pi}_1 \\ \vdots \\ \bar{\pi}_d \end{pmatrix}$$

$$\cdot V(X_i) = \bar{\pi}_i (1 - \bar{\pi}_i)$$

# RVs  $\leftarrow k$



SUMMARY:

- $\bar{X} \sim \text{MULTINOMIAL}(n, \tilde{\pi}_1, \dots, \tilde{\pi}_{k+1})$

$$\hookrightarrow f(x_1, \dots, x_k) = \binom{n}{x_1, \dots, x_{k+1}} \tilde{\pi}_1^{x_1} \cdot \tilde{\pi}_2^{x_2} \cdots \tilde{\pi}_k^{x_k} \cdot \left(1 - \sum_i \tilde{\pi}_i\right)^{x_{k+1}}$$

$\rightarrow$  1 R.V.

$$\cdot K=1 : f(x_1) = \binom{n}{x_1} \cdot \tilde{\pi}_1^{x_1} \cdot \left(1 - \tilde{\pi}_1\right)^{n-x_1} \rightarrow \text{BINOMIAL DISTRIBUTION}$$

# EXTRactions IN EXPERIMENT = 1

$$\cdot n=1 : f(x_1, \dots, x_n) = \tilde{\pi}_1^{x_1} \cdots \tilde{\pi}_k^{x_k} \cdot \left(1 - \sum_j \tilde{\pi}_j\right)^{1-\sum x_j} \rightarrow \text{MULTINOMIAL DISTRIBUTION}$$

$$\cdot K=1, n=1 : f(x_1) = \begin{cases} 1, p = \tilde{\pi}_1 \\ 0, p = 1 - \tilde{\pi}_1 \end{cases} \rightarrow \text{BERNOULLI DISTRIBUTION}$$

- CONSIDER  $\bar{X} \sim \text{MULTINOMIAL}(\tilde{\pi}_1, \dots, \tilde{\pi}_k)$  :  $\rightarrow$  OUTCOMES:  $\{0, \dots, 1, \dots, 0\}$  OF SHAPE

$$\rightarrow \text{SINCE } X_i \sim \text{Bernoulli}(\tilde{\pi}_i) \rightarrow E[X_i] = \tilde{\pi}_i, V[X_i] = \tilde{\pi}_i(1 - \tilde{\pi}_i)$$

$$\hookrightarrow C_oV(X_i, X_j) = \underset{i \neq j}{E[X_i X_j]} - E[X_i] E[X_j] = 0 - \tilde{\pi}_i \cdot \tilde{\pi}_j = -\tilde{\pi}_i \tilde{\pi}_j$$

PROPERTIES:

$$\cdot \text{IF } \bar{X} \sim \text{MULTINOMIAL}(n, \tilde{\pi}_1, \dots, \tilde{\pi}_{k+1}) \rightarrow \bar{X} = \sum_{i=1}^k n_i / n, n_i \sim \text{MURNOULI}(\tilde{\pi}_1, \dots, \tilde{\pi}_k)$$

$$\cdot E(\bar{X}) = \sum_{i=1}^k E(n_i) = n \begin{pmatrix} \tilde{\pi}_1 \\ \vdots \\ \tilde{\pi}_k \end{pmatrix}, V(X_i) = n \tilde{\pi}_i (1 - \tilde{\pi}_i)$$

$$\cdot C_oV(X_i, X_j) = -n \tilde{\pi}_i \tilde{\pi}_j$$

$$\cdot E[\bar{a} \bar{X} + \bar{b}] = \bar{a} E[\bar{X}] + \bar{b}$$

$$\cdot V_{\text{ar}} C_oV(A \cdot \bar{X}) = A \cdot V_{\text{ar}} C_oV(\bar{X}) \cdot A^T$$

$\rightarrow$

so IF

$$\bar{X} \sim \text{MULTINOMIAL}(n, \tilde{\pi}_1, \dots, \tilde{\pi}_k)$$

$\rightarrow$

$$V_{\text{ar}} C_oV(\bar{X}) = n \cdot \begin{pmatrix} \tilde{\pi}_1(1 - \tilde{\pi}_1) & -\tilde{\pi}_1 \tilde{\pi}_2 & \cdots \\ -\tilde{\pi}_1 \tilde{\pi}_2 & \tilde{\pi}_2(1 - \tilde{\pi}_2) & \cdots \\ \vdots & \vdots & \ddots \end{pmatrix}$$

## MULTIVARIATE NORMAL:

$$\mu = \begin{pmatrix} 0 \\ 0 \end{pmatrix}, \Sigma = \begin{pmatrix} 1 & 0 \\ 0 & 1 \end{pmatrix}$$

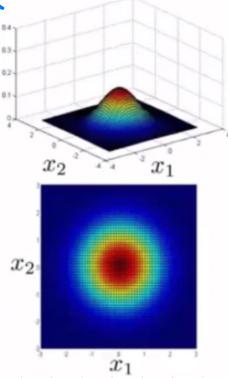
1. DEFINE THE MULTIVARIATE STD NORMAL:

$$\bar{z} = \begin{pmatrix} z_1 \\ \vdots \\ z_k \end{pmatrix} / z_1, \dots, z_k \stackrel{iid}{\sim} N(0, 1)$$

$$f(z_1, \dots, z_k) = \prod_{i=1}^k \frac{1}{\sqrt{2\pi}} e^{-\left(\frac{z_i^2}{2}\right)} = (2\pi)^{-\frac{k}{2}} e^{-\frac{1}{2}\|z\|^2}$$

$$\cdot E(\bar{z}) = \bar{0}^T, \text{Var}(\text{CoV}(\bar{z})) = \begin{pmatrix} 1 & 0 & \dots & 0 \\ 0 & 1 & & : \\ \vdots & & \ddots & \\ 0 & \dots & \dots & 1 \end{pmatrix} = I_k$$

IDENTITY MATRIX  
 $k \times k$



2. DEFINE ANY MULTIVARIATE NORMAL VIA LINEAR TRANSFORMATION

• STARTING WITH  $\bar{z} \sim N_k(\bar{0}, I_k)$ ,  $\bar{\mu}$ ,  $A$ :

$$\rightarrow \text{MULTIVARIATE NORMAL: } \bar{X} = A \cdot \bar{z} + \bar{\mu}$$

$$\cdot E(\bar{X}) = A \cdot E(\bar{z}) + \bar{\mu} = \bar{\mu}$$

$$\cdot \text{Var}(\text{CoV}(\bar{X})) = \text{Var}(\text{CoV}(A \cdot \bar{z})) = A \cdot \underbrace{\text{Var}(\text{CoV}(\bar{z}))}_{I} \cdot A^T = A \cdot A^T$$

so,

$$\mu = (0 \ 1)^T, A = \begin{pmatrix} 1 & 0 \\ 2 & 1 \end{pmatrix} / k = n = 2$$

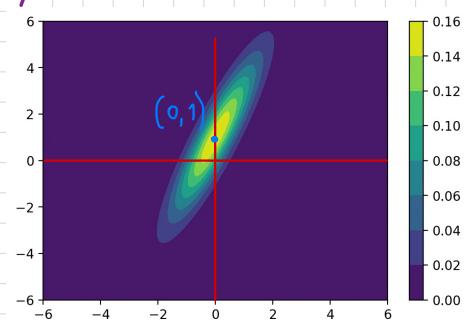
$$\rightarrow \begin{pmatrix} X \\ Y \end{pmatrix} = A \cdot \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} + \mu = \begin{pmatrix} 1 & 0 \\ 2 & 1 \end{pmatrix} \begin{pmatrix} z_1 \\ z_2 \end{pmatrix} + \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \begin{pmatrix} z_1 \\ 2z_1 + z_2 + 1 \end{pmatrix}$$

• WHERE  $z_1, z_2 \stackrel{iid}{\sim} N(0, 1)$

$$\rightarrow E\left(\begin{pmatrix} X \\ Y \end{pmatrix}\right) = \bar{\mu} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

$$\rightarrow \text{Var}(\text{CoV}\left(\begin{pmatrix} X \\ Y \end{pmatrix}\right)) = \begin{pmatrix} 1 & 0 \\ 2 & 1 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 2 \\ 2 & 5 \end{pmatrix}$$

$$g = \frac{\text{CoV}(X, Y)}{\sigma_X \cdot \sigma_Y} = \frac{2}{\sqrt{5}}$$



• PROOF THAT  $\text{Cov}(X_l, X_m) = -n \tilde{\Sigma}_l \tilde{\Sigma}_m$ , IF  $m \neq l$  :

$\tilde{X} = \sum M_i$ ,  $M_1, \dots, M_n \sim \text{MVN}(0, \tilde{\Sigma})$

$\hookrightarrow$   $n$  MULTIVARIATE  $(M_1, \tilde{\Sigma}_1, \dots, \tilde{\Sigma}_n)$

$\rightarrow \text{Cov}(X_l, X_m) = \text{Cov}\left(\sum_{i=1}^n M_{il}, \sum_{i=1}^n M_{im}\right) =$

$$= \sum_{i=1}^n \sum_{j=1}^n \text{Cov}(M_{il}, M_{jm}) \stackrel{\text{IF } i=j \Rightarrow \text{Cov} = 0}{=} \sum_{i=1}^n \text{Cov}(M_{il}, M_{im})$$

$$= n (-\tilde{\Sigma}_l \tilde{\Sigma}_m) = -n \tilde{\Sigma}_l \tilde{\Sigma}_m$$

### DENSITY OF A MULTIVARIATE NORMAL:

$$\cdot \bar{X} \sim \underset{n \times 1}{A} \cdot \underset{n \times K}{\bar{\Sigma}} + \underset{n \times 1}{\bar{\mu}} \sim \mathcal{N}(\bar{\mu}, \underset{n \times n}{\tilde{\Sigma}} = A \cdot A^T)$$

IF  $\tilde{\Sigma}^{-1}$  EXISTS, WE CAN APPLY THE JACOBIAN METHOD TO PROVE  
THAT THE DENSITY OF  $\bar{X}$  ALSO EXISTS, THAT IS:

$$\mathbb{R}^n \rightsquigarrow f(\bar{X}) = (2\pi)^{-\frac{n}{2}} \cdot |\det(\tilde{\Sigma})|^{-\frac{1}{2}} \cdot \exp\left\{-\frac{1}{2} (\bar{X} - \bar{\mu})^T \tilde{\Sigma}^{-1} (\bar{X} - \bar{\mu})\right\}$$

ELLIPSOIDE IN  $\mathbb{R}^n$

$\mathbb{R}^2 \rightsquigarrow$  IF  $n=2$ :

$$\cdot \bar{X} = \begin{pmatrix} X \\ Y \end{pmatrix} \rightarrow \tilde{\Sigma} = \begin{pmatrix} \tilde{\sigma}_x^2 & \rho \tilde{\sigma}_x \tilde{\sigma}_y \\ \rho \tilde{\sigma}_x \tilde{\sigma}_y & \tilde{\sigma}_y^2 \end{pmatrix}$$

$$\rightarrow \tilde{\Sigma}^{-1} = \frac{1}{\det(\tilde{\Sigma})} \begin{pmatrix} \tilde{\sigma}_x^2 & -\rho \tilde{\sigma}_x \tilde{\sigma}_y \\ -\rho \tilde{\sigma}_x \tilde{\sigma}_y & \tilde{\sigma}_y^2 \end{pmatrix} = \frac{1}{\underbrace{\tilde{\sigma}_x^2 \tilde{\sigma}_y^2 - \rho^2 \tilde{\sigma}_x^2 \tilde{\sigma}_y^2}_{G_x^2 G_y^2 (1-\rho^2)}} \begin{pmatrix} \tilde{\sigma}_x^2 & -\rho \tilde{\sigma}_x \tilde{\sigma}_y \\ -\rho \tilde{\sigma}_x \tilde{\sigma}_y & \tilde{\sigma}_y^2 \end{pmatrix} =$$

$$= \frac{1}{1-\rho^2} \begin{pmatrix} 1/\tilde{\sigma}_x^2 & -\rho/\tilde{\sigma}_x \tilde{\sigma}_y \\ -\rho/\tilde{\sigma}_x \tilde{\sigma}_y & 1/\tilde{\sigma}_y^2 \end{pmatrix}$$

(... AFTER SOME ALGEBRA), FOR  $\rho \neq \pm 1$

$\rightarrow$

$$f(x, y) = (2\pi)^{-1} \left( (1-\rho^2) \tilde{\sigma}_x^2 \tilde{\sigma}_y^2 \right)^{-\frac{1}{2}} \cdot \exp\left\{-\frac{1}{2(1-\rho^2)} \left[ \left(\frac{X-\mu_x}{\tilde{\sigma}_x}\right)^2 - 2\rho \left(\frac{X-\mu_x}{\tilde{\sigma}_x}\right) \left(\frac{Y-\mu_y}{\tilde{\sigma}_y}\right) + \left(\frac{Y-\mu_y}{\tilde{\sigma}_y}\right)^2 \right]\right\}$$

ELLIPSE IN  $\mathbb{R}^2$

ex.

$$\mu = (0 \ 1)^T, A = \begin{pmatrix} 1 & 0 \\ 2 & 1 \end{pmatrix} \quad / \quad k = n = 2$$

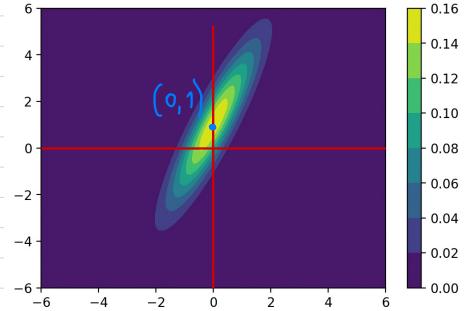
$$\rightarrow \begin{pmatrix} X \\ Y \end{pmatrix} = A \cdot \begin{pmatrix} Z_1 \\ Z_2 \end{pmatrix} + \mu = \begin{pmatrix} 1 & 0 \\ 2 & 1 \end{pmatrix} \begin{pmatrix} Z_1 \\ Z_2 \end{pmatrix} + \begin{pmatrix} 0 \\ 1 \end{pmatrix} = \begin{pmatrix} Z_1 \\ 2Z_1 + Z_2 + 1 \end{pmatrix}$$

• WHERE  $Z_1, Z_2 \stackrel{\text{iid}}{\sim} N(0, 1)$

$$\rightarrow E\begin{pmatrix} X \\ Y \end{pmatrix} = \bar{\mu} = \begin{pmatrix} 0 \\ 1 \end{pmatrix}$$

$$\rightarrow \text{Var}(\text{Cov}\begin{pmatrix} X \\ Y \end{pmatrix}) = \begin{pmatrix} 1 & 0 \\ 2 & 1 \end{pmatrix} \begin{pmatrix} 1 & 2 \\ 0 & 1 \end{pmatrix} = \begin{pmatrix} 1 & 2 \\ 2 & 5 \end{pmatrix}$$

$$\cdot \beta = \frac{2}{\sqrt{5}} \neq \pm 1$$



$$\cdot \text{DENSITY: } f(x, y) = \frac{1}{2\pi\sigma_X\sigma_Y\sqrt{1-\rho^2}} \exp\left(-\frac{1}{2(1-\rho^2)} \left[ \left(\frac{x-\mu_X}{\sigma_X}\right)^2 - 2\rho \left(\frac{x-\mu_X}{\sigma_X}\right) \left(\frac{y-\mu_Y}{\sigma_Y}\right) + \left(\frac{y-\mu_Y}{\sigma_Y}\right)^2 \right]\right)$$

$$\begin{aligned} f(x, y) &= (2\pi)^{-1} \cdot \left(1 - \left(\frac{2}{\sqrt{5}}\right)^2 \cdot 1 \cdot 5\right)^{-\frac{1}{2}} \cdot \exp\left\{-\frac{1}{2 \cdot \left(1 - \frac{4}{5}\right)} \left[x^2 - 2 \cdot \frac{2}{\sqrt{5}} \cdot x \cdot \left(\frac{y-1}{\sqrt{5}}\right) + \left(\frac{y-1}{\sqrt{5}}\right)^2\right]\right\} \\ &= (2\pi)^{-1} \cdot \exp\left\{-\frac{5}{2} \left[x^2 - \frac{4}{5}x + (y-1)^2 + \left(\frac{y-1}{\sqrt{5}}\right)^2\right]\right\} \end{aligned}$$

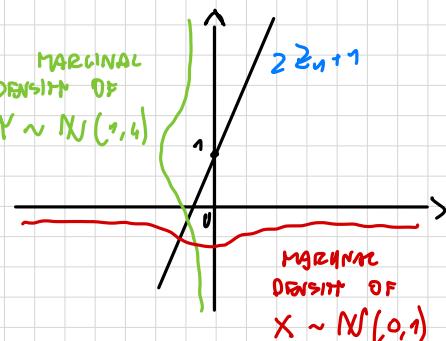
ex. DEGENERATE CASE,  $\beta = \pm 1$

$$\begin{pmatrix} X \\ Y \end{pmatrix} = \begin{pmatrix} Z_1 \\ 2Z_1 + 1 \end{pmatrix} \rightarrow \text{LINEAR CORRECTION}$$

$$\cdot E\begin{pmatrix} X \\ Y \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \quad \rightarrow \beta = \frac{2}{\sqrt{1+4}} = 1$$

$$\cdot \text{Var}(\text{Cov}\begin{pmatrix} X \\ Y \end{pmatrix}) = \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix}$$

$\rightarrow$  SINCE  $Y = 2X + 1 \rightarrow$  NO 2D DENSITY



• PROPERTIES :

MARGINAL DISTRIBUTIONS OF SUBVECTORS OF MULTIVARIATE NORMAL VECTORS  
ARE MULTIVARIATE NORMAL THEMSELVES

$$\rightarrow \bar{X} = \begin{pmatrix} \bar{X}_1 \\ \bar{X}_2 \end{pmatrix}$$

$$\cdot E(\bar{X}) = \begin{pmatrix} E(\bar{X}_1) \\ E(\bar{X}_2) \end{pmatrix} = \begin{pmatrix} \bar{\mu}_1 \\ \bar{\mu}_2 \end{pmatrix} = \bar{\mu}$$

$$\cdot \text{Var}(\bar{X}) = \Sigma = \begin{pmatrix} \sum_{11} & \sum_{12} \\ \sum_{21} & \sum_{22} \end{pmatrix}$$

$$\rightarrow \bar{X}_1 \sim N_{m_1} (\bar{\mu}_1, \Sigma_{11})$$

$$\rightarrow \bar{X}_2 \sim N_{m_2} (\bar{\mu}_2, \Sigma_{22})$$

• N.B.:  $\Sigma$  IN GENERAL CAN HAVE  $m_1 \neq m_2$

$$\text{so, } m_1 = 1, m_2 = n-1$$

$$m_1 = 1 \left\{ \begin{array}{c|c} \sigma_{x_1}^2 & \text{Cov}(x_1, x_2) \dots \text{Cov}(x_1, x_n) \\ \hline \text{Cov}(x_1, x_2) & \vdots \\ & \text{Cov}(x_1, x_n) \end{array} \right\}$$

$$m_2 = n-1 \left\{ \begin{array}{c|c} \Sigma_{22} & \end{array} \right\}$$

## CONDITIONAL NORMAL DISTRIBUTIONS:

IT CAN BE PROVED THAT CONDITIONAL DISTRIBUTIONS OF NORMAL ARE NORMAL

$\rightarrow$  IF  $\Sigma_{11}$  IS NOT SINGULAR  $\left( \Sigma^{-1}_{11} \text{ EXISTS} \right)$

$$\rightarrow X_2 | X_1 = x_1 \sim N_{m_2} \left( \bar{\mu}_2 + \Sigma_{21} \cdot \Sigma_{11}^{-1} (\bar{x}_1 - \bar{\mu}_1), \Sigma_{22} - \Sigma_{21} \cdot \Sigma_{11}^{-1} \Sigma_{12} \right)$$

IF  $m = 2$ ,  $m_1 = 1$ ,  $m_2 = 1$ :  $\Leftrightarrow \rho \neq \pm 1$

$$\cdot \begin{pmatrix} X \\ Y \end{pmatrix} : Y | X = x \sim N \left( \mu_y + \rho \frac{\sigma_y}{\sigma_x} (x - \mu_x), (1 - \rho^2) \sigma_y^2 \right)$$

$$\cdot P(X \leq x_1 \wedge Y \leq y_1) = \int \begin{pmatrix} X \\ Y \end{pmatrix} \sim \begin{pmatrix} z_1 \\ z_2 + z_1 + 1 \end{pmatrix}$$

$$= \iint_A f(x, y) dx dy = \int_{-\infty}^{x_1} \int_{-\infty}^{y_1} f(x, y) dx dy :$$

$$= \int_{-\infty}^{x_1} f_X(x) \underbrace{\int_{-\infty}^{y_1} f_{Y|X}(y|x) dy}_{P(x,y) = P(y)P(x|y)} dx$$

```
library(mvtnorm)

# Mean vector: (mu_1, ..., mu_n)
mean <- c(0, 1)
# Covariance matrix
cov <- matrix(c(1, 2,
                2, 4), nrow = 2, ncol = 2)

lower_bounds <- c(-Inf, -Inf)
# Upper bounds
upper_bounds <- c(2, 1)

# Compute the probability
probability <- pmvnorm(
  lower = lower_bounds,
  upper = upper_bounds,
  mean = mean,
  sigma = cov
)
# paste(): concatenate float and str
print(paste("Probability:", probability))
```

$\hookrightarrow$  [1] "Probability: 0.5"

CDF OF A NORMAL

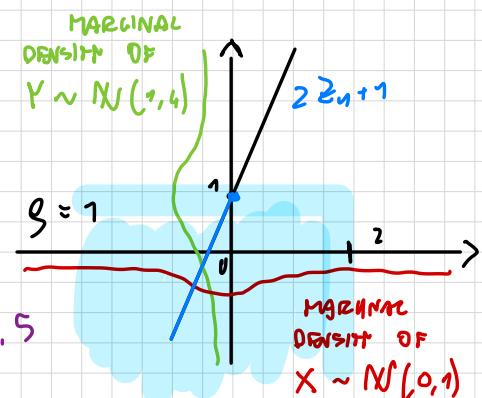
$$\cdot \text{IN GENERAL: } P(X_1 \leq x_1, \dots, X_n \leq x_n) = \int_{-\infty}^{x_1} \dots \int_{-\infty}^{x_n} f(t_1, \dots, t_n) dt_1 \dots dt_n$$

$\rightsquigarrow$  SOLVABLE  
IN PRIMAR, R

$$\text{ex. } \begin{pmatrix} X \\ Y \end{pmatrix} = \begin{pmatrix} z_1 \\ z_2 + 1 \end{pmatrix} \rightarrow \begin{matrix} E(X) = 0 \\ E(Y) = 1 \end{matrix} \rightarrow \text{Var}(\text{Cov}) \begin{pmatrix} X \\ Y \end{pmatrix} = \begin{pmatrix} 1 & 2 \\ 2 & 4 \end{pmatrix}$$

$$\cdot P(X \leq 2 \wedge Y \leq 1) =$$

$$= P \left[ \begin{pmatrix} X \\ Y \end{pmatrix} \in A \right] = P[X \leq 0] = 0.5$$



- SIMULATION :

THE BASE IS TO EXTRACT  $x_0$  FROM  $X \sim U[0, 1]$

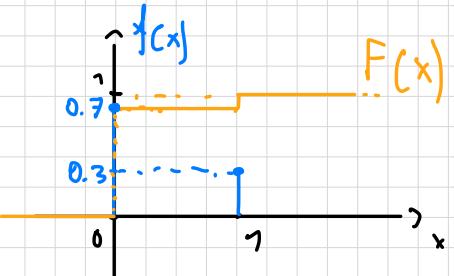
$\hookrightarrow \text{np.random.random(1)}$

$\rightarrow$  EVERY EXTRACTION OF  $X \sim f(x)$  CAN BE DONE BASED ON THE EXTRACTION FROM THE UNIFORM DISTRIBUTION :

ex.

$X \sim \text{Bernoulli}(p_0)$  :

$\rightarrow$  RETURN  $\begin{cases} 0 & \text{if } X_0 \leq p_0 \\ 1 & \text{if } X_0 > p_0 \end{cases}$



ex.

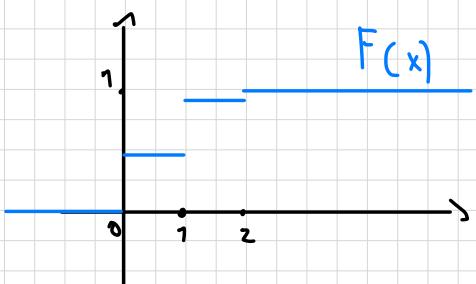
$X \sim \text{Binomial}(n = n_0, p = p_0)$

$\rightarrow$  COMPUTE CDF AND APPLY SAME TRICK

$$f_X(x_0) = \begin{cases} k=0 : \binom{n_0}{0} p_0^0 (1-p_0)^{n_0} = 0.4^0 \\ k=1 : \binom{n_0}{1} p_0^1 (1-p_0)^{n_0-1} = 0.42 \\ \dots \end{cases}$$

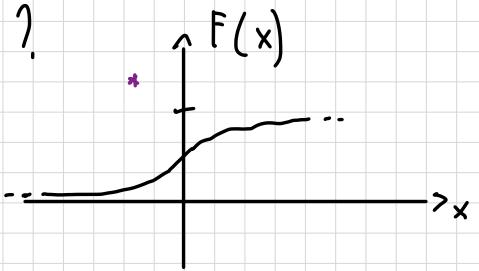
$n_0=2, p_0=0.3$

$$\rightarrow F_X(x_0) = \begin{cases} 0 & , x \leq 0 \\ 0.4^0 & , 0 < x \leq 1 \\ 0.97 & , 1 < x \leq 2 \\ 1 & , x \geq 2 \end{cases}$$



$$\rightarrow \text{RETURN} \begin{cases} 0 & , p = 0.4^0 \\ 1 & , p = 0.42 \\ 2 & , p = 0.04 \end{cases}$$

- WHAT IF  $X$  has density  $f_X(x)$  CONTINUOUS?
- $F_X(x_0) = P(X \leq x_0)$
- $\text{es. } Z \sim N(0, 1)$ \*

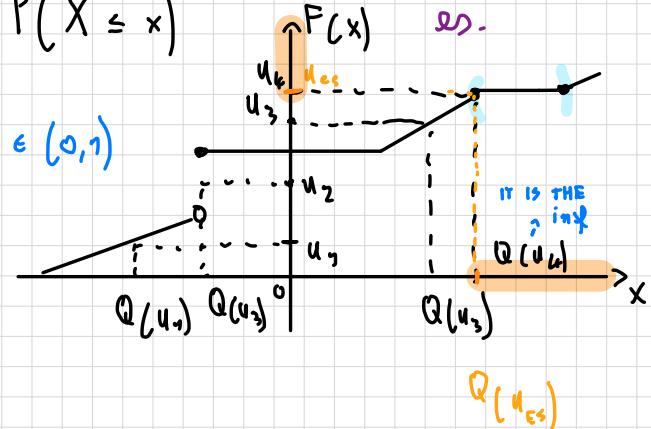


-> WE NEED TO DEFINE A (GENERALIZED) INVERSE OF  $F$ :

- IF  $X$  has CDF  $F(x) = P(X \leq x)$

-> QUANTILE FUNCTION:

$$Q(u) = \inf \{x : F(x) \geq u\}$$



• THEOREM :

$$Q(u) \leq x \Leftrightarrow u \leq F(x)$$

$$\forall u \in (0,1), \forall x \in \mathbb{R}$$

• PROOF:

- IF  $u \leq F(x) \rightarrow Q(u)$  is AN INFINITE:  $Q(u) \leq x$
- IF  $Q(u) \leq x \rightarrow F(Q(u)) \leq F(x)$

$$\rightarrow u \leq F(Q(u)) \leq F(x)$$

• THEOREM :

IF  $X$  is a R.V. with QUANTILE FUNCTION  $Q(u)$ ,  $u \in (0,1)$   
AND IF  $U \sim \text{UNIFORM}(0,1)$

->  $Q(u)$  AND  $X$  HAVE SAME DISTRIBUTION

• PROOF:

DEFINITION  $F_X(x) = P(X \leq x)$ , CDF of  $X$

-> CDF of  $Q(U)$ :  $F_{Q(U)}(x) = P(Q(U) \leq x) = P(U \leq F(x)) = F(x)$

THEOREM ABOVE

SINCE  $U$  IS UNIFORM

→ THIS RESULTS SIMPLIFIES THE WENBOG ALGORITHM TO  
SIMULATE AN UNIVARIATE R.V.:

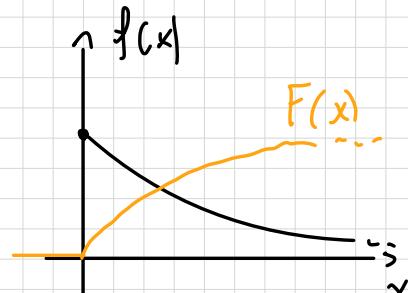
1. OBTAIN  $U \sim \text{UNIFORM}(0, 1)$

2. RETURN  $X = Q(U)$

ex.

$$X \sim \text{Exp}(\lambda) \rightarrow f_X(x) = \lambda e^{-\lambda x}, x > 0$$

$$\cdot F_X(x) = \int_0^x \lambda e^{-\lambda t} dt = -e^{-\lambda t} \Big|_0^x = 1 - e^{-\lambda x}$$



$$(Y = 1 - e^{-\lambda x} = 1 - V = e^{-\lambda x} : \ln(1-V) = -\lambda x \rightarrow x = -\frac{\ln(1-V)}{\lambda})$$

$$\cdot Q(u) = F_x^{-1}(u) \quad [NOT SURE AS NOTATION]$$

→ SIMULATION:

1. GET  $U \sim \text{UNIFORM}(0, 1)$

2. RETURN  $X = -\frac{\ln(1-U)}{\lambda}$  TO SIMULATE  $X \sim \text{Exp}(\lambda)$

• HOW DO WE SIMULATE U EXTRACTION (STEP 1.)

↳ USE OF PSEUDO-RANDOM NUMBER GENERATION

• ONE CLASS OF ALGORITHM TO DO THIS: LINEAR CONGRUENTIAL GENERATORS

→ THEY EXPLOIT THE UNPREDICTABILITY OF THE REMAINDER OF A DIVISION BY 2 INTEGERS

• ALGORITHM:

• CHOOSE 4 POSSIBLY LARGE NUMBERS:  $M, a, c$

↗ USED VALUES ARE:  $a = 7^5$   
 $M = 2^{31}-1$   
 $c = 0$

• CHOOSE SEED  $L_0$

• DEFINE FOR  $i = 1, 2, \dots$ :  $L_i = (aL_{i-1} + c) \bmod M$

→ EXTRACT  $u_i = \frac{L_i}{M}$ : RATIONAL NUMBER IN  $[0, 1]$

## MULTIVARIATE SIMULATION

HOW DO WE SIMULATE FROM

$$\bar{X} = \begin{pmatrix} X_1 \\ \vdots \\ X_d \end{pmatrix}, \text{ WITH A COVARIANCE } \Sigma \text{ AND A DENSITY } f(x_1, \dots, x_d)$$

->

- TAKE  $\bar{X} \sim N_d(\bar{\mu}, \Sigma)$  AND

$$\text{FIND } A / \Sigma = A \cdot A^T, \text{ FOR } \text{SO. : } A \stackrel{\text{CHOLESKY}}{\underbrace{\cdot}} \text{chol}(\Sigma)$$

- THEN:

$$1) \text{ SIMULATE } \bar{Z} = \begin{pmatrix} Z_1 \\ \vdots \\ Z_d \end{pmatrix}, \text{ iid } N(0, 1)$$

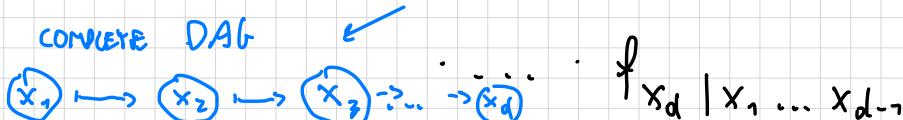
$$2) \text{ SET } \bar{X} = \bar{\mu} + A \cdot \bar{Z}, \text{ IN THIS WAY } \bar{X} \sim N_d(\bar{\mu}, AA^T)$$

THIS WORKS ONLY WITH MULTIVARIATE NORMAL

→ PROBLEM CAN BE REDUCED TO A SERIES OF UNIVARIATE SIMULATION

- A通用ALGORITHM TO APPLY FOR ANY DISTRIBUTION, TO GENERATE  $\bar{X}$ :

$$f(x_1, \dots, x_d) = f_{x_1}(x_1) \cdot f_{x_2|x_1}(x_2|x_1) \cdot f_{x_3|x_1, x_2}(x_3|x_1, x_2)$$



- ALGORITHM:

1. GENERATE  $x_1$  FROM THE MARGINAL  $f_{x_1}$

2. "  $x_2$  " "

:

d. "  $x_d$  " "

$f_{x_2|x_1}$

$f_{x_3|x_1, x_2}$

$f_{x_d|x_1, \dots, x_{d-1}}$

GIBBS  
SAMPLER

```

A <- matrix(c(
  4, 12, -16,
  12, 37, -43,
  -16, -43, 98
), nrow = 3, ncol = 3)
L <- chol(A) # A = L * L^T
A = L %*% t(L)

```

ex.

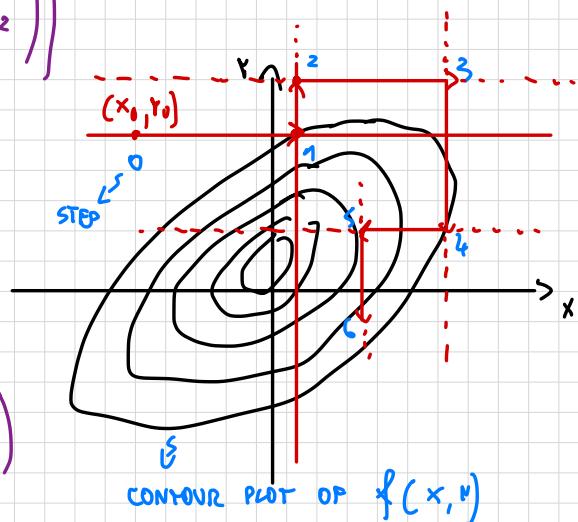
$$\begin{pmatrix} X \\ Y \end{pmatrix} \sim N_2 \left( \begin{pmatrix} \mu_X \\ \mu_Y \end{pmatrix}, \begin{pmatrix} \sigma_X^2 & \sigma_{XY} \\ \sigma_{YX} & \sigma_Y^2 \end{pmatrix} \right)$$

0. START AT  $\begin{pmatrix} x_0 \\ y_0 \end{pmatrix}$

FOR  $j = 1, \dots, B$ :

1. GENERATE  $X_{j+1}$  FROM  $f_{X|Y}(x | y_j)$

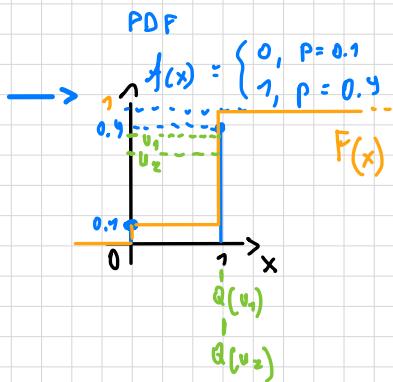
2. GENERATE  $y_{j+1}$  FROM  $f_{Y|X}(y | x_{j+1})$



→ AT iteration B :  $(x_B, y_B)$  will have density  $f(x, y)$

→  $(x_B, y_B)$  will be a genuine simulated value from  $f(x, y)$

### • UTILITIY:



SIMULATION :

```
x <- round(runif(1), 2)
```

$$u_1 \sim U[0,1] = 0.86 \rightarrow 1$$

$$u_2 = 0.71 \rightarrow 1$$

$$u_3 = \dots$$

### • CONTOUR PLOT OF A MULTIVARIATE NORMAL:

## • MONTE CARLO SIMULATION:

GIVEN  $X_1, Y$  iid  $\sim U[-1, 1]$

$$f_X(x) = f_Y(y) = \frac{1}{2} \quad (-1 \leq x \leq 1)$$

$$\rightarrow f_{(X,Y)}(x,y) = \frac{1}{4}$$

• what is  $P\left[\begin{pmatrix} X \\ Y \end{pmatrix} \in C\right]?$

$$\rightarrow P\left[\begin{pmatrix} X \\ Y \end{pmatrix} \in C\right] = \iint_C f_{(X,Y)}(x,y) dx dy = \frac{1}{4} \iint_C dx dy =$$

$$= \frac{1}{4} (\text{area of } C) = \frac{1}{4} \cdot (\pi \cdot 1^2) = \frac{\pi}{4}$$

## • MONTE CARLO APPROXIMATION:

1. GENERATE  $(X_1, Y_1), (X_2, Y_2), \dots, (X_B, Y_B)$  iid  $\begin{pmatrix} X \\ Y \end{pmatrix}$  FOR large  $B$

2. SET  $y(x_i, y_i) = \begin{cases} 1 & \text{IF } \begin{pmatrix} X \\ Y \end{pmatrix} \in C \\ 0 & \text{IF } \begin{pmatrix} X \\ Y \end{pmatrix} \notin C \end{cases}$

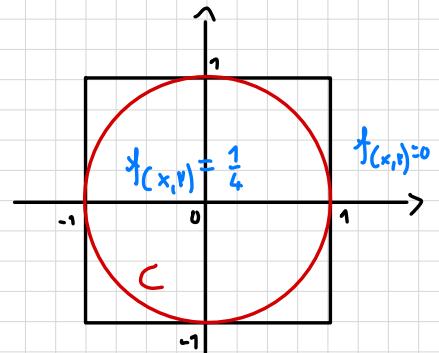
• BY LN :  $\frac{1}{B} \sum_{i=1}^B y(x_i, y_i) = \frac{\# \text{TIMES } \begin{pmatrix} X \\ Y \end{pmatrix} \in C}{B} \approx P\left(\begin{pmatrix} X \\ Y \end{pmatrix} \in C\right)$

$\hookrightarrow \frac{1}{B} \sum_{i=1}^B y(x_i, y_i) \approx E[y(x_i, y_i)]$

• TO APPROXIMATE  $\pi$ :

$$\rightarrow \frac{1}{B} \sum y(x_i, y_i) \approx \frac{\pi}{4}$$

$$\rightarrow \frac{4}{B} \sum y(x_i, y_i) \approx \pi$$



```
# Number of random points
n <- 10000000

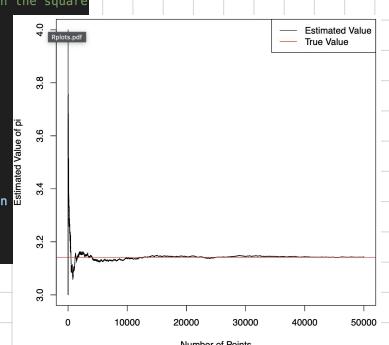
# Generate random x and y coordinates within the square
x <- runif(n, -1, 1)
y <- runif(n, -1, 1)

# Calculate distance from the origin
dist <- sqrt(x^2 + y^2)

# Check if points fall within the circle
points_in_circle <- dist <= 1

# Estimate pi
pi_estimate <- 4 * sum(points_in_circle) / n

# Print the estimated value of pi
print(pi_estimate)
```



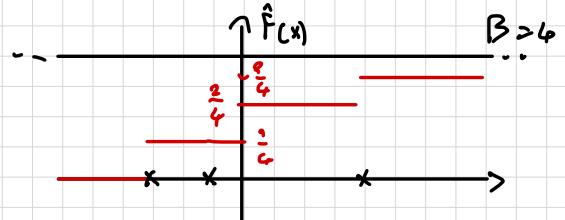
## 2. EMPIRICAL CDF

SUPPOSE  $X$  UNIVARIATE R.V. /  $F(X) = \int_{-\infty}^x f(t) dt = P(X \leq x)$

• HOW TO APPROX.  $F(x_0)$ :

1. TENDENCE  $x_1, \dots, x_B$  iid  $X$

$$2. \hat{F}(x_0) = \frac{1}{B} \sum_{i=1}^B (x_i \leq x_0)$$



• EMPIRICAL CDF DO NOT PROVIDE AN ESTIMATE OF  $f(x)$

$\hookrightarrow \frac{d}{dx} \hat{F}(x)$  DO NOT HAVE A SENSITIVE

→ AN HISTOGRAM IS A DENSITY ESTIMATE, BUT VERY PRIMITIVE

• HOW TO MEASURE VARIABILITY IN A MEASURE?

$$\cdot \bar{y} = \frac{1}{B} \sum_{i=1}^B y(\tilde{x}_i^*) = \text{SAMPLE MEAN ESTIMATE } E[y(\bar{x})]$$

$\hookrightarrow \tilde{x}_1^*, \dots, \tilde{x}_B^*$  ARE SYNTHETIC OBSERVATIONS

• IF  $V(y(\bar{x}))$  IS ASSUMED TO EXIST  $\rightarrow$  CLT APPLIES:

$$\frac{\bar{y} - E[y(\bar{x})]}{\sqrt{\frac{V(y(\bar{x}))}{B}}} \xrightarrow[B \rightarrow \infty]{d} N(0, 1)$$

$$\rightarrow \text{SAMPLE VARIANCE: } s_y^2 = \frac{\sum_{i=1}^B (y(\tilde{x}_i^*) - \bar{y})^2}{B-1}$$

$$\rightarrow \text{APPROXIMATE ASYMPTOTIC C.I. OF UNEC } 1-\alpha: \bar{y} \pm z_\alpha \sqrt{\frac{s_y^2}{B}}$$

• SOMETIMES WE CANNOT DO iid SIMULATION  $\rightarrow$  iid MONTE CARLO NOT POSSIBLE:

MARSHAL CHAIN  
MONTE CARLO

$\rightarrow$  2 WAYS:

- OR
1. START MANY MARKOV CHAINS AND USE LAST OBSERVATION  $\tilde{x}_i^*$
  2. USE A SINGLE VERY LONG MARKOV CHAIN AND SAMPLE IT

THERE WILL BE SIMILAR LN FOR MARKOV CHAINS  
ENSURENCE  $\bar{y}$  IS A GOOD ESTIMATE OF  $E(y(\bar{x}))$

ESSENCE OF MCMC:

## • LINEAR MODELS :

A QUANTITATIVE RESPONSE VARIABLE  $Y$  HAS TO BE EXPLAINED IN TERMS OF 1 OR MORE PREDICTORS (FEATURES)

$$\bar{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} 1 & X_{1,1} & \dots & X_{1,p-1} \\ \vdots & \vdots & & \vdots \\ 1 & X_{n,1} & \dots & X_{n,p-1} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_{p-1} \end{pmatrix} + \begin{pmatrix} \epsilon_0 \\ \vdots \\ \epsilon_n \end{pmatrix}$$

INTERCEPT COLUMN       $p-1$  FEATURES      VECTOR OF COEFFICIENTS      ERROR TERM

RANDOM VECTOR OF OBSERVATIONS

$$\rightarrow Y_i = \beta_0 + \beta_1 \cdot X_{i,1} + \dots + \beta_{p-1} \cdot X_{i,p-1} + \epsilon_i$$

• ONCE OBSERVED  $Y_1, \dots, Y_n$  BECOMES NUMBERS  $y_1, \dots, y_n$

•  $\epsilon_1, \dots, \epsilon_n$  ARE ERRORS / iid NORMAL  $(0, \sigma^2)$

•  $\beta_0, \dots, \beta_{p-1}$  ARE UNKNOWN COEFFICIENTS

•  $X_{i,1}, \dots, X_{i,p-1}$  WILL BE A VECTOR OF CONSTANT FEATURES

ASSOCIATED TO THE  $i$ -TH RESPONSE

ERROR VARIANCE

$$\rightarrow Y_i \sim N(\beta_0 + \beta_1 \cdot X_{i,1} + \dots + \beta_{p-1} \cdot X_{i,p-1}, \sigma^2)$$

ex.

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_m \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ \vdots & \vdots \\ 1 & 0 \\ \vdots & \vdots \\ 1 & 1 \end{pmatrix}_{m \times m} \cdot \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} + \begin{pmatrix} \epsilon_1 \\ \vdots \\ \epsilon_m \end{pmatrix}$$

1 BINARY  
FEATURE ONLY

$$\cdot Y_i = \beta_0 + \epsilon_i, \quad i = 1, 2, \dots, m_1$$

$$\cdot Y_i = \beta_0 + \beta_1 + \epsilon_i, \quad i = m+1, \dots, m$$

$\rightarrow$

$$Y_1, \dots, Y_{m_1} \text{ iid } N\left(\frac{\beta_0}{2}, \sigma^2\right)$$

$$Y_{m+1}, \dots, Y_m \text{ iid } N\left(\frac{\beta_0 + \beta_1}{2}, \sigma^2\right)$$

• THE PREVIOUS EXAMPLE IS THE TWO-NORMAL-SAMPLE PROBLEM  
IN TWO GROUP EXPERIMENT:

WE ARE COMPARING TWO MEANS  $\mu$  AND  $\nu$  / IN LINEAR MODEL :  $\begin{cases} \mu = \beta_0 \\ \nu = \beta_0 + \beta_1 \end{cases}$

$\rightarrow \nu - \mu = \beta_1$  : DIFFERENCE OF 2 NORMAL MEANS, EACH GROUP SAME  $V(\cdot)$

$\hookrightarrow t\text{-TEST AND RELATED } t\text{-C.I.}$

2). 1 SAMPLE MODEL, NORMAL MODEL

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \beta_0 + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix} \rightarrow Y_1, \dots, Y_n \text{ iid } N\left(\frac{\mu}{\beta_0}, \sigma^2\right)$$

3). SIMPLE LINEAR REGRESSION  $\rightarrow p=2$ , 1 QUANTITATIVE PREDICTOR

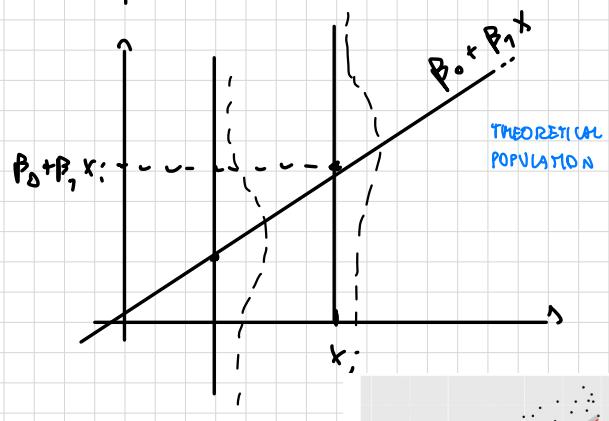
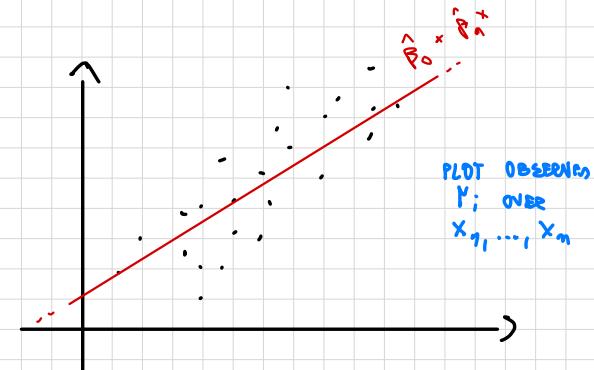
$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$$\rightarrow Y_i \sim N(\beta_0 + \beta_1 x_i, \sigma^2)$$

• WE DO NOT KNOW  $\beta_0, \beta_1, \sigma^2$

$\rightarrow$  WE WILL ESTIMATE THEM

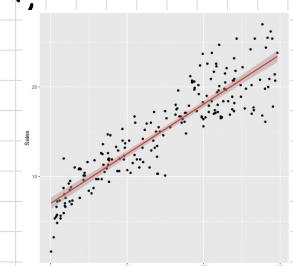
BASED ON DATA



• IN R:

```
advertising <- read_csv("./datasets/advertising.csv")
```

```
simple_reg <- lm(Sales ~ TV, data = advertising)
```



## • ESTIMATION IN LINEAR MODELS:

LINEAR MODEL:

$$\bar{Y} = X \cdot \beta + \varepsilon$$

$m \times 1$      $m \times p$      $p \times 1$      $m \times 1$

SINCE  $\varepsilon \sim N_m(\bar{0}, \sigma^2 I_m)$ :

•  $\bar{Y} = X \bar{\beta} + \bar{\varepsilon}$  IS A LINEAR COMBINATION OF MULTIVARIATE NORMAL

$$\rightarrow \bar{Y} \sim N(E(\bar{Y}), \text{Var}(oV(\bar{Y}))$$

$$\cdot E(\bar{Y}) = X \bar{\beta} + E(\varepsilon) = X \bar{\beta}$$

$$\cdot \text{Var}(oV(\bar{Y})) = \text{Var}(oV(X \bar{\beta} + \bar{\varepsilon})) = \text{Var}(oV(\bar{\varepsilon})) = \sigma^2 I$$

$$\rightarrow \bar{Y} \sim N(X \bar{\beta}, \sigma^2 I)$$

$\hookrightarrow \bar{Y}$  ARE INDEPENDENT, BUT NOT iid

• MLE TO ESTIMATE PARAMETERS  $\beta_0, \dots, \beta_{p-1}, \sigma^2$ :

$$\mathcal{L}(\bar{\beta}, \sigma^2; Y_1, \dots, Y_m) = \prod_{i=1}^m \left[ \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp \left\{ -\frac{1}{2\sigma^2} \left( Y_i - \sum_{j=1}^{p-1} \beta_j \cdot x_{i,j} \right)^2 \right\} \right] =$$

$$= (2\pi\sigma^2)^{-\frac{m}{2}} \cdot \exp \left\{ -\frac{1}{2\sigma^2} \sum_{i=1}^m \underbrace{\left( Y_i - \sum_{j=1}^{p-1} \beta_j \cdot x_{i,j} \right)^2}_{(\bar{Y} - X \bar{\beta})^T \cdot (\bar{Y} - X \bar{\beta})} \right\}$$

IN MATRIX FORM

• MAXIMIZING  $\log \mathcal{L}$ :

$$\log \mathcal{L}(\beta, \sigma^2, Y_1, \dots, Y_m) = -\frac{m}{2} \log(2\pi\sigma^2) - \frac{1}{2\sigma^2} (\bar{Y} - X \bar{\beta})^T (\bar{Y} - X \bar{\beta})$$

$$\cdot \underline{N.B.}: \max \log L \Leftrightarrow \min (\bar{Y} - X\bar{\beta})^T (\bar{Y} - X\bar{\beta})$$

$$\rightarrow \min_{\beta} (\bar{Y} - X\bar{\beta})^T (\bar{Y} - X\bar{\beta}) = \min_{\beta} \left( Y^T Y - 2 \cdot \underbrace{\beta^T X^T Y}_{1 \times p \quad p \times n \quad n \times 1} + \beta^T X^T X \beta \right)$$

N.B.: IMPORTANT RULES FOR MATRIX DIFFERENTIATION

$$\cdot \frac{\delta}{\delta x} A\bar{x} = A^T \quad , \quad \cdot \frac{\delta}{\delta x} x^T A x = 2Ax$$

$$\rightarrow \frac{\delta}{\delta \beta} \left( Y^T Y - 2 \cdot \underbrace{\beta^T X^T Y}_{1 \times 1} + \beta^T X^T X \beta \right) = -2X^T Y + 2X^T X \bar{\beta} = 0$$

$$\rightarrow \bar{x}^T \cdot \bar{x} \cdot \bar{\beta} = \bar{x}^T \cdot \bar{Y}$$
NORMAL EQUATION

• PROVIDED THAT  $(\bar{x}^T \bar{x})^{-1}$  EXISTS :

$$\rightarrow \hat{\beta} = (\bar{x}^T \bar{x})^{-1} \cdot \bar{x}^T \cdot \bar{Y}$$

$\rightarrow$

$$\left. \begin{array}{l} \text{ESTIMATED COEFFICIENTS} \\ \cdot \hat{\beta} = (X^T X)^{-1} \cdot X^T \cdot \bar{Y} \\ \cdot \hat{\beta} = (X^T X)^{-1} \cdot X^T \cdot \bar{Y} \end{array} \right. \begin{array}{l} \text{LS ESTIMATE OF } \beta, \text{ CORRESPONDING TO OUR DATA} \\ \text{LS ESTIMATE OF } \bar{\beta}, \text{ A RANDOM VECTOR} \end{array}$$

$$\cdot \hat{Y} = X \cdot \hat{\beta}$$
PROJECTION OF OUR DATA ONTO SPAN(X)

$$\cdot \bar{Y} = X \cdot \bar{\beta}$$
CORRESPONDING RANDOM VECTOR

$$\left. \begin{array}{l} \text{RESIDUALS} \\ \cdot \bar{e} = \bar{Y} - \hat{Y} \\ \cdot \bar{E} = \bar{Y} - \hat{Y} \end{array} \right. \begin{array}{l} \text{DIFFERENCE BETWEEN OBSERVED DATA AND THEIR PROJECTION} \\ \text{RESIDUALS} \end{array}$$

## en. NULL MODEL

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_m \end{pmatrix} = \underbrace{\begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix}}_{\mathbf{x}} \cdot \beta_0 + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_m \end{pmatrix}$$

$\leadsto Y_i = \beta_0 + \varepsilon_i$

$$\rightarrow \beta_0 = (\bar{x}^T \bar{x})^{-1} \cdot \bar{x}^T \bar{Y} = \begin{pmatrix} (1 \dots 1) & (1) \\ \vdots & \vdots \\ 1 & 1 \end{pmatrix}, \quad (1 \dots 1) \begin{pmatrix} Y_1 \\ \vdots \\ Y_m \end{pmatrix} =$$

$$= \dots = \frac{1}{n} \sum Y_i$$

$\rightarrow$  LSSE OF  $\beta_0$  IS THE SAMPLE MEAN

## en. LINEAR REGRESSION

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_m \end{pmatrix} = \begin{pmatrix} 1 & X_1 \\ \vdots & \vdots \\ 1 & X_m \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_m \end{pmatrix}$$

$\leadsto Y_i = \beta_0 + \beta_1 X_i + \varepsilon_i$

$$\cdot \hat{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} = (\bar{x}^T \cdot \bar{x})^{-1} \cdot \bar{x}^T \cdot \bar{Y} = \left( \begin{pmatrix} 1 & \dots & 1 \\ X_1 & \dots & X_m \end{pmatrix} \begin{pmatrix} 1 & X_1 \\ \vdots & \vdots \\ 1 & X_m \end{pmatrix} \right)^{-1} \cdot \begin{pmatrix} 1 & \dots & 1 \\ X_1 & \dots & X_m \end{pmatrix} \begin{pmatrix} Y_1 \\ \vdots \\ Y_m \end{pmatrix} =$$

$$= \begin{pmatrix} n & \sum X_i \\ \sum X_i & \sum X_i^2 \end{pmatrix}^{-1} \cdot \begin{pmatrix} 1 & \dots & 1 \\ X_1 & \dots & X_m \end{pmatrix} \begin{pmatrix} Y_1 \\ \vdots \\ Y_m \end{pmatrix} =$$

$$= \dots = \begin{cases} \hat{\beta}_0 = \bar{Y} - \hat{\beta}_1 \bar{X} \\ \hat{\beta}_1 = \frac{\sum (x_i - \bar{x}) Y_i}{\sum (x_i - \bar{x})^2} \end{cases}$$

MEAN

• WHAT ABOUT  $\hat{\sigma}^2$  ?

$$\rightarrow \frac{S}{\hat{\sigma}^2} \log \mathcal{L}(\hat{\beta}, \hat{\sigma}^2) = 0$$

$$\rightarrow \frac{S}{\hat{\sigma}^2} \left( -\frac{n}{2} \log(2\pi\hat{\sigma}^2) - \frac{1}{2\hat{\sigma}^2} (\mathbf{r} - \mathbf{x}\hat{\beta})^\top (\mathbf{r} - \mathbf{x}\hat{\beta}) \right) = 0$$

$$\rightarrow -\frac{n}{2} \frac{1}{2\hat{\sigma}^2} 2\hat{\sigma}^2 - \frac{1}{2} (\mathbf{r} - \mathbf{x}\hat{\beta})^\top (\mathbf{r} - \mathbf{x}\hat{\beta}) \frac{1}{(\hat{\sigma}^2)^2} = 0$$

• IF WE CAN EXCLUDE THE SOLUTION  $\hat{\sigma}^2 = 0$  :

$$\rightarrow \hat{\sigma}^2 = \frac{(\mathbf{r} - \mathbf{x}\hat{\beta})^\top (\mathbf{r} - \mathbf{x}\hat{\beta})}{n} = \frac{\mathbf{e}^\top \mathbf{e}}{n} \quad / \begin{array}{l} \text{RESIDUAL } \bar{\mathbf{e}} \\ \mathbf{e}^\top \mathbf{e} : \text{SUM OF SQUARES} \end{array}$$

$\rightarrow$  IT IS ALSO THE MLE OF  $\hat{\sigma}^2$

so. NULL MODEL

$$\begin{pmatrix} \mathbf{r}_1 \\ \vdots \\ \mathbf{r}_m \end{pmatrix} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \cdot \beta_0 + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_m \end{pmatrix} \quad \sim \quad \mathbf{r}_i = \beta_0 + \varepsilon_i$$

$$\rightarrow \hat{\sigma}^2 = \frac{\mathbf{e}^\top \mathbf{e}}{n} = \frac{\sum_i (\mathbf{r}_i - \hat{\beta}_0)^2}{n} = \frac{\sum_i (\mathbf{r}_i - \bar{\mathbf{r}})^2}{n} = \frac{\sum_i (\mathbf{r}_i - \bar{\mathbf{r}})^2}{n-1}$$

• ACTUALLY A BETTER ESTIMATE OF  $\hat{\sigma}^2$  :  $S^2 = \frac{\sum_i (\mathbf{r}_i - \hat{\mathbf{r}})^2}{n-p}$

... OR :  $MSR = \frac{\mathbf{e}^\top \mathbf{e}}{n-p} \quad / p: \text{DEGREE OF FREEDOMS}$

$\rightarrow \sim N_m(E(\hat{\beta}), \text{Var}(Cov(\hat{\beta}))$

•  $\hat{\beta}$  IS A R.V. ITSELF : IT'S ALSO NORMAL SINCE IS A LINEAR TRANS. OF  $\mathbf{r}$

$$\cdot E(\hat{\beta}) = E((\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \mathbf{r}) = (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \cdot E(\mathbf{r}) = (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \mathbf{x} \beta = \beta$$

$$\cdot \text{Var}(Cov(\hat{\beta})) = \text{Var}(Cov((\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \mathbf{r})) = (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \cdot \text{Var}(Cov(\mathbf{r})) \cdot ((\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top)^\top = \text{Var}(Cov(\mathbf{A}^\top \mathbf{X})) = \mathbf{A} \cdot \text{Var}(Cov(\mathbf{X})) \cdot \mathbf{A}^\top = (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x}^\top \cdot \hat{\sigma}^2 \mathbf{I} \cdot (\mathbf{x}^\top \mathbf{x})^{-1} \mathbf{x} = (\mathbf{x}^\top \mathbf{x})^{-1} \cdot \hat{\sigma}^2 \mathbf{I}$$

SQUARED STANDARD ERROR OF  $\hat{\beta}_i$

$\rightarrow$  A SINGLE COMPONENT :  $\hat{\beta}_i \sim N(\beta_i, \frac{\hat{\sigma}^2 (\mathbf{x}^\top \mathbf{x})^{-1}_{ii}}{\text{IF WILL BE ESTIMATED WITH MSR}})$ ,  $i = 0, \dots, p-1$

## CONFIDENCE INTERVALS FOR $\beta$ .

RECALL :  $\bar{Y} = X \cdot \beta + \epsilon \rightarrow \text{LINEAR MODEL}$

GIVEN  
 $\bar{Y} = \sum_{i=1}^n Y_i \sim N(\bar{\beta}, \sigma^2 I)$

$\Rightarrow R.V. , \text{ OBSERVED INTO VECTORS } \bar{Y}$

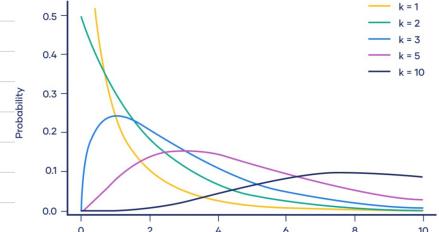
$n: \# \text{SAMPLES}$   
 $p: \# \text{FEATURES}$

AN ESTIMATE (MLE) IS :  $\hat{\beta} \sim \text{ESTIMATOR: } \hat{\beta} = \frac{E^T E}{n}$

$\hat{\sigma}^2 = \frac{e^T e}{n} = \text{SUM OF SQUARED RESIDUALS} = \frac{\sum e_i^2}{n} / e_i : r_i - \hat{Y}_i$

IT CAN BE PROVEN THAT :  $\frac{E^T E}{\hat{\sigma}^2} \sim \chi^2_{n-p}$

IF R.V  $X \sim \chi^2_k \rightarrow E[X] = k$



$\rightarrow E\left[\frac{E^T E}{\hat{\sigma}^2}\right] = n-p \rightarrow E\left[\frac{E^T E}{n-p}\right] = E[MSE] = \hat{\sigma}^2$

$\Rightarrow MSE$  IS A BETTER ESTIMATOR THAN  $\frac{E^T E}{n}$ , WHICH IS BIASED

### C. I.:

$\hat{\beta}_i \sim N(\beta_i, \sigma^2 (X^T X)_{ii}^{-1})$

THEOREM:  $\frac{N(0,1)}{\sqrt{\frac{\chi^2_k}{k}}} \sim t_k$

$\rightarrow$  WE CAN STANDARDIZE IT:

$$\frac{\hat{\beta}_i - \beta_i}{\sqrt{\sigma^2 (X^T X)_{ii}^{-1}}} = \frac{\hat{\beta}_i - \beta_i}{\sqrt{\frac{E^T E}{n-p} \cdot \frac{\sigma^2}{\hat{\sigma}^2} \cdot (X^T X)_{ii}^{-1}}} = \frac{\hat{\beta}_i - \beta_i}{\sqrt{\frac{\chi^2_{n-p}}{n-p}}} \sim t_{n-p}$$

THEOREM

• C.I. FOR  $\beta$ :

• t-student C.I. from ELEMENTARY STATISTICS:  $SE(\hat{\beta}_i) = \sqrt{G^2(X^T X)^{-1}_{ii+1, ii+1}} = \sqrt{\frac{s^2}{n}}$

• SINCE:  $P\left[\beta_i - t_{\frac{\alpha}{2}, n-p} \cdot SE(\hat{\beta}_i) \leq \hat{\beta}_i \leq \beta_i + t_{\frac{\alpha}{2}, n-p} \cdot SE(\hat{\beta}_i)\right] = 1-\alpha$

$$\rightarrow C.I._{1-\alpha} = \left[ \beta_i - t_{\frac{\alpha}{2}, n-p} \cdot SE(\hat{\beta}_i); \beta_i + t_{\frac{\alpha}{2}, n-p} \cdot SE(\hat{\beta}_i) \right]$$

• C.I. FOR  $G^2$ :

$$C.I._{1-\alpha} : P\left[\chi^2_{1, \frac{n}{2}, n-p} < \frac{E^T E}{G^2} < \chi^2_{\frac{n}{2}, n-p}\right] = 1-\alpha$$

2. TWO SAMPLE PROBLEM

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_m \end{pmatrix} = \begin{pmatrix} 1 & 0 \\ \vdots & \vdots \\ 0 & 1 \\ \vdots & \vdots \\ 1 & 1 \end{pmatrix}_{m \times m} \cdot \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

1 BINARY  
FEATURE ONLY

$$\cdot Y_i = \beta_0 + \varepsilon_i, i = 1, 2, \dots, m, \text{ iid } N\left(\tilde{\beta}_0, G^2\right)$$

$$\cdot Y_i = \beta_0 + \beta_1 + \varepsilon_i, i = m+1, \dots, n \rightarrow Y_{m+1}, \dots, Y_n \text{ iid } N\left(\tilde{\beta}_0 + \tilde{\beta}_1, G^2\right)$$

$\rightarrow$

$$\text{DIFFERENCE OF MEANS: } \beta_1 = \mu - \nu$$

$\rightarrow p=2 \rightarrow 2 \text{ FEATURES}$

$$\rightarrow C.I._{1-\alpha} = \left[ \hat{\beta}_1 - t_{\frac{\alpha}{2}, n-2} \cdot SE(\hat{\beta}_1); \hat{\beta}_1 + t_{\frac{\alpha}{2}, n-2} \cdot SE(\hat{\beta}_1) \right]$$

## INTUITIVE GRAPHS :

- $\text{Y}$  : RESPONSE
- $X_1, X_2, \dots; Z_1, Z_2, \dots$  : QUANTITATIVE FEATURES
- $b_1, b_2, \dots$  : BINARY FEATURES
- $a_1, a_2, \dots$  : CATEGORICAL FEATURES

->

- $\text{Y}$       NULL MODEL
- $\text{Y} \leftarrow b$  . 2 NORMAL SAMPLE PROBLEM
- $\text{Y} \leftarrow X$  : SIMPLE LINEAR REGRESSION
- $\text{Y} \leftarrow \begin{matrix} X_1 \\ \vdots \\ X_n \end{matrix}$  . MULTIPLE LINEAR REGRESSION
- $\text{Y} \leftarrow \begin{matrix} X \\ \vdots \\ X^k \end{matrix}$  : POLYNOMIAL REGRESSION  $\rightsquigarrow$  BUT LINEARITY IS IN  $\beta$
- $\text{Y} \leftarrow a$  : SINGLE CATEGORICAL FEATURE (ONEWAY ANOVA)
- $\text{Y} \leftarrow \begin{matrix} a \\ \vdots \\ c \end{matrix}$  : TWO-WAY ANOVA

## TESTS ON $\beta_i$ :

- NULL HYPOTHESIS  $\rightarrow H_0: \beta_i = \beta_i^0$

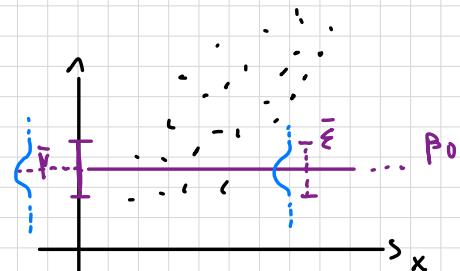
- IF WE TEST  $H_0: \beta_i = \beta_i^0 \leq 0 \rightarrow$  WE ARE TESTING

IF FEATURE  $X_i$  HAS INFERENCE ON  $\gamma$ :

$$\bar{Y} = \begin{pmatrix} 1 & X_{11} & \dots & X_{1,p-1} \\ \vdots & \vdots & & \vdots \\ 1 & X_{n1} & \dots & X_{n,p-1} \\ \hline X_{(..,i)} & = 0 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_{p-1} \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

ex. SIMPLE REGRESSION

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} 1 & X_1 \\ \vdots & \vdots \\ 1 & X_m \\ \hline \vdots & = 0 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 = 0 \end{pmatrix} + \bar{\varepsilon}$$



$\rightarrow$  IF  $\beta_1 = 0 \rightarrow \bar{Y} = \beta_0 + \bar{\varepsilon} \rightarrow \bar{X}$  NO INFERENCE ON  $\gamma$

ex. 2 NORMAL SAMPLES

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_m \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} 1 & 0 & \dots & 0 \\ \vdots & \vdots & & \vdots \\ 1 & 0 & \dots & 0 \\ \vdots & \vdots & & \vdots \\ 1 & 0 & \dots & 0 \end{pmatrix}_{m \times m} \cdot \begin{pmatrix} \beta_0 \\ \beta_1 \\ \vdots \\ \beta_p \end{pmatrix} + \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_m \\ \vdots \\ \varepsilon_n \end{pmatrix}$$

$\rightarrow$  EQUIVALENT TO COMPARING 2 NORMAL MEANS:  $\beta_1 = \mu_x - \mu_y$

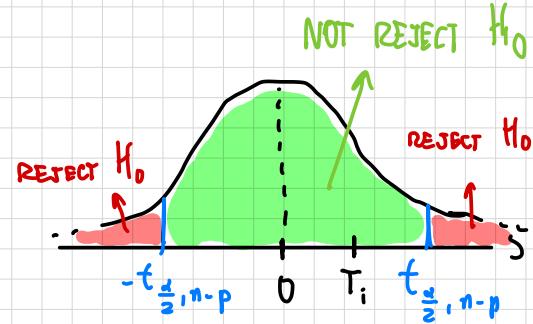
$\rightarrow H_0: \beta_1 = 0 \equiv H_0: \mu_x - \mu_y = 0 \equiv H_0: \mu_x = \mu_y$

- IN GENERAL: IF  $H_0: \beta_i = 0$  IS TRUE  $\rightarrow X_i$  IS USELESS IN EXPLAINING  $\gamma$

## • HOW TO REJECT / NOT REJECT $H_0$ ?

- REJECT  $H_0$  IF :

$$T_i = \left| \frac{\hat{\beta}_i - \beta_i^0}{SE(\hat{\beta}_i)} \right| > t_{\frac{\alpha}{2}, n-p}$$



## Ex. ADVERTISING DATASET

- FIT:  $Y$ : SALES,  $X$ : RADIO

$$\begin{aligned} Y &= \hat{\beta}_0 + \hat{\beta}_1 \cdot \text{RADIO} = \\ &= 12.2357 + 0.1246 \cdot \text{RADIO} \end{aligned}$$

$$\cdot T_1 = \frac{\hat{\beta}_1}{SE(\hat{\beta}_1)} = 5.257$$

$$\cdot SE(\hat{\beta}_1) = \sqrt{\frac{e^T e}{n-p}} (x^T x)^{-1}_{i+1, i+1} = 0.0237$$

→ suppose we choose  $\alpha = 0.05$ :

LARGE DATASET

$$t_{\frac{\alpha}{2} = 0.025, n-p = 198} = 1.972 \quad (\approx Z_{0.025} = 1.96)$$

$$\rightarrow T_1 = 5.257 > t_{\frac{\alpha}{2}, n-p} = 1.972 \rightarrow H_0 \text{ REJECTED}$$

• p-VALUE:

• ↓ p-VALUE → WE BECOME LESS IN  $H_0$ : OUR OBSERVED VALUE IS EXTREME

→ { p-VALUE  $\leq \alpha \rightarrow H_0$  REJECTED)

{ p-VALUE  $> \alpha \rightarrow H_0$  NOT REJECTED)

Residuals:				
Min	1Q	Median	3Q	Max
-15.5632	-3.5293	0.6714	4.2504	8.6796
Coefficients:				
Estimate Std. Error t value Pr(> t )				
(Intercept)	12.2357	$\beta_0$	0.6535	18.724 < 2e-16 ***
Radio	0.1244	$\beta_1$	0.0237	5.251 3.88e-07 ***
---				
Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1				
Residual standard error: 4.963 on 198 degrees of freedom				
Multiple R-squared: 0.1222, Adjusted R-squared: 0.1178				
F-statistic: 27.57 on 1 and 198 DF, p-value: 3.883e-07				

AN OLR NOTHING TEST:

- $H_0: \beta_1 = \beta_2 = \dots = \beta_{p-1} = 0$

$\hookrightarrow$  opposite  $H_p: H_p: \exists \beta_i \neq 0$

$\rightarrow H_0$  refers to NULL MODEL  $\rightarrow \beta_D = k \neq 0, \beta_i = 0, i \neq 0, k \neq 0$

- IN THIS CASE:  $E(\bar{Y}) = \beta_0, \hat{\beta}_0 = \bar{Y} \rightarrow \begin{pmatrix} Y_1 \\ \vdots \\ Y_m \end{pmatrix} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \beta_0 + \bar{E} \sim N(0, \sigma^2)$

VECTOR OF FITTED VALUES:

$$\hat{Y} = \begin{pmatrix} \hat{Y}_1 \\ \vdots \\ \hat{Y}_m \end{pmatrix} = X \cdot \hat{\beta} = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \cdot \hat{\beta}_0 = \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \cdot \bar{Y} = \begin{pmatrix} Y_1 \\ \vdots \\ Y_m \end{pmatrix}$$

$\rightarrow$  NULL MODEL  $\subset$  FULL MODEL:

- FULL MODEL:  $\beta_i$  may not be 0

- NULL MODEL:  $\beta_1 = \dots = \beta_{p-1} = 0, \beta_0 = k$

- TESTING  $H_0$   $\equiv$  NULL MODEL IS ENOUGH TO EXPLAIN  $Y$

- IF NULL MODEL IS ENOUGH ( $H_0$ )  $\rightarrow E_0^T E_0 \cdot E^T E \mapsto$  SMALL NUMBER ?

$$/ E_0^T E_0 - E^T E = \| \hat{Y} - \hat{Y}_0 \|$$

EXTENSION OF PYTHAGORAS'S THEOREM

$\rightarrow$  WE REJECT  $H_0$  IF  $E_0^T E_0 \cdot E^T E$  IS TOO LARGE ?

## F TEST :

- IT CAN BE PROVEN THAT :

$$\frac{E_0^T E_0 - E^T E}{S^2} \sim \chi^2_{p-1}$$

AND IF  $H_0$  IS TRUE :

$$F_{\text{STATISTICAL}} = \frac{\frac{E_0^T E_0 - E^T E}{S^2(p-1)}}{\frac{E^T E}{S^2(n-p)}} \sim F_{p-1, n-p}$$

$S^2$  HERE CANCELS

$F_{\alpha, p-1, n-p}$

F DISTRIBUTION

$$\cdot \text{BY PARMALORG'S THEOREM : } F = \frac{\frac{\|\hat{Y} - \hat{Y}_0\|^2}{p-1}}{\frac{\|\hat{Y} - \hat{Y}\|^2}{n-p}} = \frac{\frac{\sum (\hat{Y}_i - \bar{Y})^2}{p-1}}{\frac{\sum (Y_i - \bar{Y})^2}{n-p}}$$

EXPLAINED VARIANCE

UNEXPLAINED VARIANCE

$\cdot H_0$  REJECTED ( $\Rightarrow F > F_{\alpha, p-1, n-p}$ )

L> FROM COMPUTERS : DETECT IF P-VALUE  $< \alpha$

## 2. ONE WAY ANOVA (ANALYSIS OF VARIANCE)

ONE WAY ANOVA : LINEAR MODEL WITH CATEGORICAL FEATURE WITH MODE

THEN 2 CLASSES, ONE-HOT ENCODING USE

• USING  $p=4$  CLASSES :

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_n \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ \vdots & 0 & 1 & 0 \\ 1 & 0 & 1 & 0 \end{pmatrix} \quad \rightarrow \text{A COLUMN NEEDS TO BE ERASED TO HAVE FULL RANK() }$$

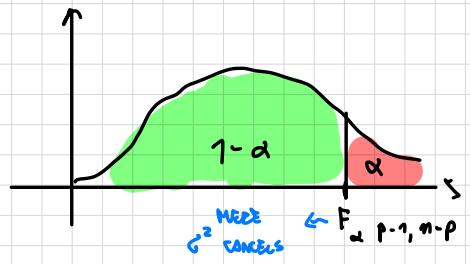
• INTERPRETATION :

$P_0$  : AVG PRICE FOR A

$P_1$  : AVG DIFF. IN PRICE B - A

,  $P_2$  : AVG DIFF. IN PRICE C - A

$P_3$  : AVG DIFF. IN PRICE D - A



→ TESTING  $H_0 : \beta_1 = \beta_2 = \beta_3 = 0$  MEANS ALL GROUPS HAVE SAME EXPECTATION

→ WE REJECT IF  $F = \frac{\text{VARIANCE BETWEEN GROUPS}}{\text{VARIANCE WITHIN GROUPS}}$  IS LARGE.

• AN F TEST ON A SMALL MODEL VS LARGE MODEL

CAN BE DONE IN R USING : anova (NESTED LINEAR MODELS)  
small / large)

Ex.

$Y$ : FINAL LAUREA SCORE (110)

PREDICTORS:

$X_1$ : MURIBITA,  $X_2$ : ANALYSIS I SCORE

$X_3$ : CHEMISTRY SCORE,  $X_4$ : PHYSICS SCORE

→

SMALL MODEL:  $m_1 \leftarrow \text{lm}(Y \sim X_1 + X_2)$

LARGE MODEL:  $m_2 \leftarrow \text{lm}(Y \sim X_1 + X_2 + X_3 + X_4)$

→ anova ( $m_1, m_2$ )

• IN R:

$$\cdot \text{DOF}_1 = n - p_1 - 1 = 200 - 2 - 1 = 197$$

$$\cdot \text{DOF}_2 = n - p_2 - 1 = 200 - 3 - 1 = 196$$

$$\cdot \text{RSS} = \sum (Y_i - \hat{Y}_i)^2 = e^T e$$

$$\cdot Df = \text{DOF}_2 - \text{DOF}_1$$

$$\cdot \text{SUM OF Sq} = \text{RSS}_2 - \text{RSS}_1$$

$$\cdot F = \frac{\frac{\sum (\hat{Y}_i - \bar{Y})^2}{p-1}}{\frac{\sum (Y_i - \hat{Y})^2}{n-p}} = 0,0034$$

$$\rightarrow P[X > 0,034] \approx 0,45 = p\text{-value} > \alpha \rightarrow H_0 \text{ NOT REJ}$$

```
# Fit the full model
full_model <- lm(Sales ~ TV + Radio + Newspaper, data = advertising)
# Fit the reduced model (without the Newspaper predictor)
reduced_model <- lm(Sales ~ TV + Radio, data = advertising)
# Perform the F-test
f_test <- anova(reduced_model, full_model)

# Print the F-test results
print(f_test)
```

#### Analysis of Variance Table

		Model 1: Sales ~ TV + Radio		Model 2: Sales ~ TV + Radio + Newspaper			
		Res.Df	RSS	Df	Sum of Sq	F	Pr(>F)
1		197	541.21				
2		196	541.20	1	0.0092861	0.0034	0.9538

## CONFIDENCE AND PREDICTION INTERVALS:

SUPPOSE  $[X_f]$ : NEW COMBINATION OF PREDICTORS

e.g. LINEAR REGRESSION  $X_f = [1 \ X_f^o]$

FITTING VALUE OF  $X_f$  IS:  $\hat{Y}_f = X_f \cdot \hat{\beta}$

WHERE  $E[\hat{Y}_f \cdot \hat{\beta}] = X_f \cdot \beta$

C.I FOR.  $\hat{Y}_f = X_f \cdot \hat{\beta} / X_f \cdot \hat{\beta} = X_f (X^T X)^{-1} X^T Y$ :

$$E[\hat{Y}_f] = E[X_f (X^T X)^{-1} X^T Y] = X_f \cdot (X^T X)^{-1} X^T \cdot X \beta = X_f \beta$$

$$\begin{aligned} V[\hat{Y}_f] &= X_f (X^T X)^{-1} X^T \cdot V(\text{Var}(Y)) \cdot X (X^T X)^{-1} X^T \\ &\approx \sigma^2 \cdot X_f (X^T X)^{-1} (X^T X) (X^T X)^{-1} X^T \end{aligned}$$

$$\rightarrow \hat{Y}_f \sim N(X_f \beta, \sigma^2 \cdot X_f (X^T X)^{-1} X^T)$$

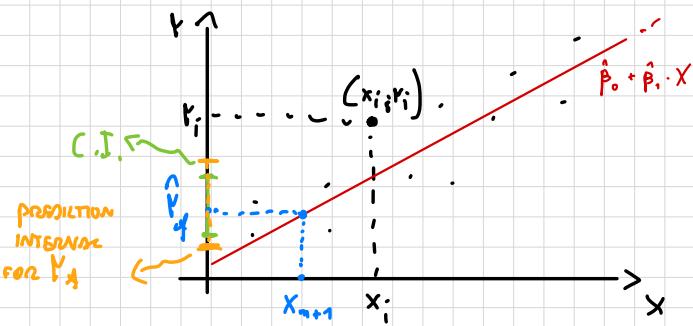
$\rightarrow$

$$C.I_{1-\alpha}: X_f \cdot \hat{\beta} \pm t_{\frac{\alpha}{2}, n-p} \sqrt{\frac{\sigma^2}{n-p} \cdot X_f (X^T X)^{-1} X^T}$$

## 20. SIMPLE LINEAR REGRESSION

$$X = \begin{bmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{bmatrix}$$

$$\bar{x}_y = [1 \ x_{n+1}]$$



- WHAT WILL BE THE NEW VALUE OF  $\bar{Y}_f$  CORRESPONDING TO  $x_f$ ?

→ POINT PREDICTION:  $\hat{Y}_f = \bar{x}_f \cdot \hat{\beta}$

- PREDICTION INTERVAL:

$$\text{IN GENERAL } \bar{Y}_f \sim N(x_f \hat{\beta}, \sigma^2)$$

$$\rightarrow \text{PREDICTION ERROR: } \bar{Y}_f - \hat{Y}_f = \bar{Y}_f - \bar{x}_f \hat{\beta} \sim N$$

$$\cdot E(\bar{Y}_f - \hat{Y}_f) = E(\bar{Y}_f) - E(\hat{Y}_f) = \bar{x}_f \bar{\beta} - \bar{x}_f \bar{\beta} = 0$$

$$\cdot V(\bar{Y}_f - \hat{Y}_f) = V(\bar{Y}_f) + V(\hat{Y}_f) + 2 \text{cov}(\bar{Y}_f, \hat{Y}_f) =$$

$$= \sigma^2 + \sigma^2 x_f (x^T x)^{-1} x_f^T = \sigma^2 \left( 1 + x_f (x^T x)^{-1} x_f^T \right)$$

$$\overset{\text{PREDICT}}{\rightarrow} \frac{\bar{Y}_f - x_f \hat{\beta}}{\sqrt{\frac{E^T E}{n-p} \left( 1 + x_f (x^T x)^{-1} x_f^T \right)}} \sim t_{n-p}$$

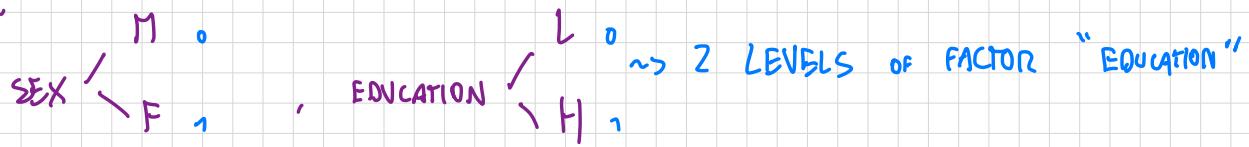
→ PREDICTION INTERVAL:

$$x_f \hat{\beta} \pm t_{\frac{\alpha}{2}, n-p} \cdot \sqrt{\frac{E^T E}{n-p} \left( 1 + x_f (x^T x)^{-1} x_f^T \right)}$$

A BIT LARGER THAN C.I.

- 2 BINARY PREDICTORS:

ex.



- IN R CATEGORICAL FEATURES ARE CALLED "FACTORS", AND THEIR VALUES ARE CALLED "LEVELS"

•  $\text{lm}(\text{INCOME} \sim \text{SEX}) \rightarrow \text{NO MEANING}$

$\hookrightarrow \text{lm}(\text{INCOME} \sim \text{ASPECTORS (SEX)}) \rightarrow$

$$\text{M/F} \rightarrow 0/1$$

$$\begin{pmatrix} 1 & 0 \\ 1 & 0 \\ \vdots & \vdots \\ 1 & 1 \end{pmatrix}$$

- GENERIC ADDITIVE MODEL WITH 2 BINARY FEATURES:

2 BINARY FEATURES

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_m \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 \\ \vdots & \vdots & \vdots \\ 1 & 0 & 1 \\ \vdots & \vdots & \vdots \\ 1 & 1 & 0 \\ 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \end{pmatrix} + \bar{\epsilon}$$

0 0  $\mu_{M,L} = \beta_0$  = MEAN RESPONSE IF SEX = M, EDUC = L

1 0  $\mu_{F,L} = \beta_0 + \beta_1$  = MEAN RESPONSE IF SEX = F, EDUC = L

0 1  $\mu_{M,H} = \beta_0 + \beta_2$  = MEAN RESPONSE IF SEX = M, EDUC = H

1 1  $\mu_{F,H} = \beta_0 + \beta_1 + \beta_2$  = MEAN RESPONSE IF SEX = F, EDUC = H

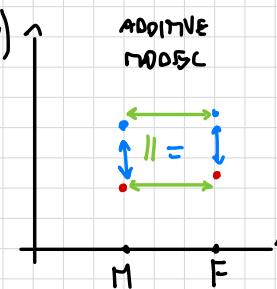
• IN R,  $\text{lm}(Y \sim X_1 + X_2)$ : ADDITIVE MODEL (NO INTERACTION)

$\hookrightarrow$  NO INTERACTION: DIFFERENCE BETWEEN 2 LEVELS IT'S THE SAME, NO MATTER WHAT OTHER FACTOR IS

→

$$\mu_{F,L} - \mu_{M,L} = \mu_{F,H} - \mu_{M,H} = \beta_1$$

$$\mu_{F,H} - \mu_{F,L} = \mu_{M,H} - \mu_{M,L} = \beta_2$$



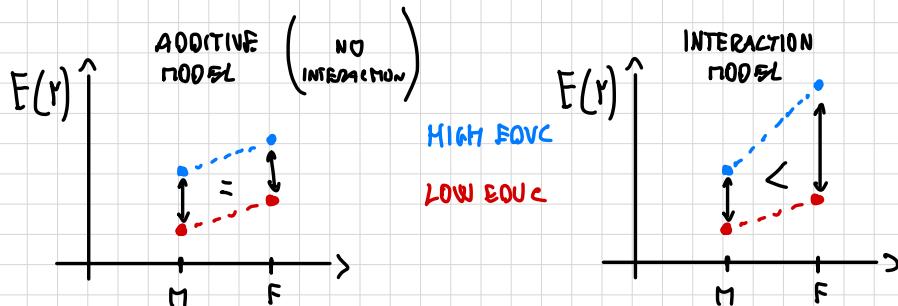
HIGH  
LOW

• GENERAL INTERACTION MODEL:

$$\begin{pmatrix} Y_1 \\ \vdots \\ Y_m \end{pmatrix} = \begin{pmatrix} 1 & 0 & 0 & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 0 & 0 & 1 & 0 \\ \vdots & \ddots & \ddots & \vdots \\ 1 & 0 & 0 & 1 \\ \vdots & 1 & 1 & 0 \\ 1 & 1 & 1 & 1 \end{pmatrix} \begin{pmatrix} \beta_0 \\ \beta_1 \\ \beta_2 \\ \beta_3 \end{pmatrix}$$

INTERACTION  
MODEL

• GRAPHICALLY:



∴



→

$$\begin{array}{ll}
 \mu_{M,L} = \beta_0 & = E(Y) \quad \text{IF } M \text{ AND } L \\
 \mu_{F,L} = \beta_0 + \beta_1 & = E(Y) \quad \text{IF } F \text{ AND } L \\
 \mu_{M,H} = \beta_0 + \beta_2 & = E(Y) \quad \text{IF } M \text{ AND } H \\
 \mu_{F,H} = \beta_0 + \beta_1 + \beta_2 + \beta_3 & = E(Y) \quad \text{IF } F \text{ AND } H
 \end{array}$$

• IN THIS CASE:

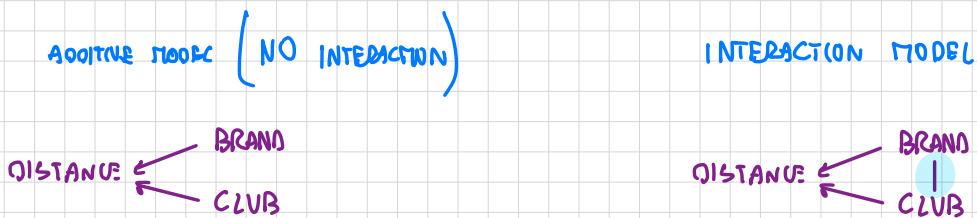
$$\mu_{F,L} - \mu_{M,L} = \beta_1 \neq \mu_{F,H} - \mu_{M,H} = \beta_1 + \beta_3$$

*Different from ADDITIVE MODEL*

↳ WHEN INTERACTION IS PRESENT, EACH COMBINATION OF FACTOR 1 AND A LEVEL OF FACTOR 2 HAS NO INTERACTION

## es. GOLF, R IN VS CODE

IN R:



- ANOVA (small, large) gives a test for null  $H_0$ ,  $H_0: \text{all interactions} = 0$ 
  - IN this example, p-value = 0.00179 → small →  $H_0$  REJECTED
  - ↳  $\exists$  EVIDENCE THAT SOME INTERACTIONS  $\neq 0$
- DIAGNOSTIC PLOTS:

• THEY ARE USED TO BE SURE WE ARE USING THE "RIGHT" MODEL

→ ONE WAY IS TO LOOK AT RESIDUAL PLOTS:  $E_i = Y_i - \hat{Y}_i$

↳  $E_i$  SHOULD BEHAVE AS NOISE, similar to  $\xi$

→ ALSO A PLOT  $E_i$  VS NEW FEATURE  $W_i$  CAN TELL IF  $W_i$  SHOULD BE INCLUDED IN MODEL: IF PATTERNS ARE EVIDENT (e.g. LINEARITY)

→  $W_i$  SHOULD BE INCLUDED IN THE MODEL

## MEASURES OF FIT:

GIVEN  $\vec{Y}$ : VECTOR OF TRUE VALUES,  $\vec{\hat{Y}}_0$ : VECTOR FROM NULL MODEL  
 $\vec{\hat{Y}}$ : PREDICTED VALUES BY THE MODEL

$\rightarrow$  THE FURTHER AWAY IS  $\vec{\hat{Y}}$  FROM  $\vec{\hat{Y}}_0 \rightarrow$  BETTER EXPLANATION OF DATA BY THE MODEL

### $R^2$ :

$$R^2 = \frac{\text{NUM}}{\text{DEN}} = \frac{\|\vec{Y} - \vec{\hat{Y}}_0\|^2}{\|\vec{Y} - \vec{\bar{Y}}\|^2} = \frac{\sum (\hat{Y}_i - \bar{Y})^2}{\sum (Y_i - \bar{Y})^2} = \frac{\text{TSS} - \text{RSS}}{\text{TSS}} = 1 - \frac{\text{RSS}}{\text{TSS}}$$

PTITAGORAS EXTRADIDA A  $12^n$  BT MEAN OF DATA

$$\cdot \text{NUM: } \|\vec{Y} - \vec{\hat{Y}}_0\|^2 = \|\vec{Y} - \vec{\hat{Y}}_0\|^2 - \|\vec{Y} - \vec{\bar{Y}}\|^2 = \text{TSS} - E^T E$$

$$\cdot \text{DEN: } \|\vec{Y} - \vec{\bar{Y}}\|^2 = E_0^T E_0 = \text{TSS}$$

TOTAL SUM OF SQUARES

RSS

RESIDUAL SUM OF SQUARES

$$\cdot 0 \leq R^2 \leq 1$$

$\cdot$  PROBLEM: IF A PREDICTOR IS ADDED  $\rightarrow R^2$ , shouldn't happen

$$\hookrightarrow \text{SOL: } R_{\text{adj}}^2$$

### $R_{\text{adj}}^2$ :

$$R_{\text{adj}}^2 = 1 - \frac{n-1}{n-p} \cdot \frac{\text{RSS}}{\text{TSS}}$$

n: # samples  
 p: # predictors

### AKAIKE INFORMATION CRITERION:

IT'S A GENERALIZATION OF  $R^2$  AND  $R_{\text{adj}}^2$  TO NON-LINEAR MODELS

BT ALSO FOR LINEAR MODELS

$$\cdot \text{AIC} = 2k - 2 \cdot \ln(\hat{L})$$

p: # estimated parameters in the model

$\hat{L}$ : maximum value of likelihood function

### BAYESIAN INFORMATION CRITERION:

$$\text{BIC} = k \cdot \ln(n) - 2 \cdot \ln(\hat{L})$$

n: # data points

$\rightarrow \downarrow \text{AIC}, \text{BIC} \rightarrow$  BETTER MODEL  $\rightsquigarrow$  UNLIKE  $R^2, R_{\text{adj}}^2$

## R ANALYSIS :

```
> advertising <- read_csv("./datasets/advertising.csv")
> simple_reg <- lm(Sales ~ TV, data = advertising)
> summary(simple_reg)
```

Residuals:

	Min	1Q	Median	3Q	Max
	-6.4438	-1.4857	0.0218	1.5042	5.6932

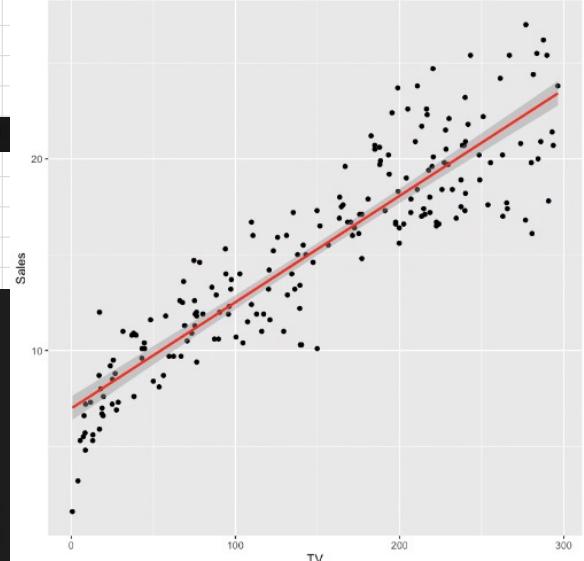
3

Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	6.974821 <sup>1</sup>	0.322553	21.62	<2e-16 ***
TV	0.055465 <sup>2</sup>	0.001896	29.26	<2e-16 ***

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.296 on 198 degrees of freedom  
 Multiple R-squared: 0.8122, Adjusted R-squared: 0.8112  
 F-statistic: 856.2 on 1 and 198 DF, p-value: < 2.2e-16



- SIMPLE REG:  $\bar{Y} = X \bar{\beta} + \tilde{\epsilon}$  /  $\bar{\beta} = \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix}$ ,  $X = \begin{pmatrix} 1 & x_{1,1} \\ \vdots & \vdots \\ 1 & x_{1,m} \end{pmatrix}$

-  $\hat{\beta}_0^1 = 6.974821$  → INTERCEPT OF UNFAR REGRESSION

-  $\hat{\beta}_1^2 = 0.055465$  → COEFFICIENT RELATED TO "TV" PREDICTOR,  $x_1$

$\rightarrow Y_i = (1 \ x_{1,i}) \begin{pmatrix} \beta_0 \\ \beta_1 \end{pmatrix} + \tilde{\epsilon} \sim N(0, G^2)$   $\uparrow$  TV predictor

## RESIDUALS:

RESIDUAL ESTIMATOR:  $\hat{Y} - X \hat{\beta}$

→ TO COMPUTE THEM:  $e_i = Y - X \hat{\beta}$

• 3:

$\min\{e_i\} = -6.4438$ ,  $\max\{e_i\} = 5.6932$

.  $1Q = -1.4857$  :  $P[e_i < -1.4857] = 25\%$   
 .  $3Q = 1.5042$  :  $P[e_i < 1.5042] = 75\%$ .

MEDIAN = 0.0218 :  $P[e_i < 0.0218] = 50\%$ .

```
library(tibble)
y <- advertising$Sales
x <- cbind(1, advertising$TV)
e <- y - x %*% simple_reg$coefficients
head(tibble(
  lm_res = simple_reg$residuals,
  manual_res = as.vector(e)))
```

A tibble: 6 x 2	
lm_res	manual_res
<dbl>	<dbl>
2.3627348	2.3627348
0.9569962	0.9569962
4.0711845	4.0711845

300 RESIDUALS

## COEFFICIENTS :

$$\hat{\beta}_0 = 6.974821, \hat{\beta}_1 = 0.055465$$

```
beta_hat <- solve(t(x) %*% x) %*% t(x) %*% y
head(tibble(
```

```
  lm_coeff = simple_reg$coefficients,
  manual_coeff = as.vector(beta_hat)
))
```

lm_coeff	manual_coeff
<dbl>	<dbl>
6.97482149	6.97482149
0.05546477	0.05546477

### Residuals:

Min	1Q	Median	3Q	Max
-6.4438	-1.4857	0.0218	1.5042	5.6932

3

### Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	6.974821	0.322553	21.62	<2e-16 ***
TV	0.055465	0.001896	29.26	<2e-16 ***

8

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 2.296 on 198 degrees of freedom

Multiple R-squared: 0.8122, Adjusted R-squared: 0.8112

F-statistic: 856.2 on 1 and 198 DF, p-value: < 2.2e-16

$$\cdot \text{COEFFICIENTS ESTIMATOR: } \hat{\beta} = (X^T X)^{-1} X^T Y$$

$$\cdot \hat{\beta}_i : SE_i = \sqrt{\hat{G}_n^{-2} (X^T X)^{-1}_{ii}}$$

$$\rightarrow \downarrow SE_i \rightarrow \uparrow \hat{\beta}_i \text{ ACCURATE}$$

$$\cdot \hat{G}_n^{-2} = \frac{e^T e}{n-p} = \frac{\sum_{i=1}^n (Y_i - \hat{Y}_i)^2}{n-p} = RMS = \text{RESIDUAL MEAN SQUARES} = S^2$$

$$\cdot SE_0 = 0.322553$$

$$\cdot SE_1 = 0.01896 \xrightarrow{\text{an.}} SE_1 = \sqrt{\frac{e^T e}{n-p}}$$

```
n <- nrow(x)
p <- ncol(x)
rms <- t(e) %*% e / (n - p)

# SE for TV
tv_se <- sqrt(rms * solve(t(x) %*% x)[2, 2])
intercept_se <- sqrt(rms * solve(t(x) %*% x)[1, 1])
tv_se
intercept_se
```

A matrix: 1 x 1  
of type dbl  
0.001895551

A matrix: 1 x 1  
of type dbl  
0.3225535

## T-VALUES:

$$\text{GIVEN } H_0: \hat{\beta}_1^0 = 0 :$$

$$\cdot T_i = \left| \frac{\hat{\beta}_i - \hat{\beta}_i^0}{SE(\hat{\beta}_i)} \right| \rightarrow \begin{cases} T_0 = \frac{\hat{\beta}_0 - 0}{SE(\hat{\beta}_0)} = \frac{\textcircled{1}}{\textcircled{4}} = 21.62 \\ T_1 = \frac{\hat{\beta}_1 - 0}{SE(\hat{\beta}_1)} = \frac{\textcircled{2}}{\textcircled{5}} = 29.26 \end{cases}$$

$$\rightarrow \text{IF } T_i > t_{\frac{\alpha}{2}, n-p} \xrightarrow{\text{200-2}} H_0 \text{ REJECTED}$$

```
# for TV
t_val <- simple_reg$coefficients[2] / tv_se
t_val
```

A matrix: 1 x 1 of type dbl  
6 29.2605

## P-VALUES:

• FOR 2-SIDED TEST :

$$P\text{-VALUE}_i = 2 \cdot P[\hat{\beta}_i \geq |T_i|] \rightarrow \begin{cases} \text{IF } p\text{-value}_i < \alpha \rightarrow H_0 \text{ REJECT} \\ \text{IF } p\text{-value}_i > \alpha \rightarrow H_0 \text{ NOT REJECT} \end{cases}$$

$$\rightarrow P\text{-VALUE}_1 = 7.9279 \cdot 10^{-74} \approx 2.2 \cdot 10^{-16} \Rightarrow H_0 \text{ REJECT}$$

```
# multiply by two because the alternative hyp is two-sided
p_val <- 2 * pt(t_val, n - 2, lower.tail = FALSE)
```

A matrix: 1 x 1 of type dbl  
7.927912e-74

8

## SIGNIF. CODES:

THEY INDICATES THE

STRONGNESS TO REJECT  $H_0$ ,  
BASED ON P-VALUE;

- \*\*\* : STRONGLY REJECT  $H_0$
- :
- 'NULL' : NOT REJECT  $H_0$

↓ P-VALUE

### Residuals:

	Min	1Q	Median	3Q	Max
	-6.4438	-1.4857	0.0218	1.5042	5.6932

3

### Coefficients:

	Estimate	Std. Error	t value	Pr(> t )
(Intercept)	6.974821	0.322553	21.62	<2e-16 ***
TV	0.055465	0.001896	29.26	<2e-16 ***

---

Signif. codes: 0 '\*\*\*' 0.001 '\*\*' 0.01 '\*' 0.05 '.' 0.1 ' '

← STRONGNESS

Residual standard error: 2.296 on 198 degrees of freedom

Multiple R-squared: 0.8122, Adjusted R-squared: 0.8112

F-statistic: 856.2 on 1 and 198 DF, p-value: < 2.2e-16

## LOGISTIC REGRESSION AS GLM: $\Rightarrow$ GENERALIZED LINEAR MODELS

SO FAR WE ASSUMED:  $\bar{Y} \sim N(X\bar{\beta}, \sigma^2 I)$

, BUT SOME ASSUMPTIONS MAY NOT BE REALISTIC:

1. IF  $\sigma^2 \neq \text{const.}$   $\rightarrow$  TRY  $\log(\bar{Y})$  OR OTHER TRANSFORMATIONS

2. IF RESIDUALS DO NOT  $\sim N(\dots, \dots)$   $\rightarrow$  TRY SIMPLE TRANSFORMATIONS (  $\log(\bar{Y})$ ,  $\log(k)$ ,  $x^2$ , ... )  
OR ENRICH MODEL WITH NEW PREDICTORS

3. IF OBSERVATIONS CLUSTER INTO 2 GROUPS  $\rightarrow$  INTRODUCE GROUP LABEL AS A PREDICTOR

4. RESIDUALS OR  $\hat{Y}_i$ 'S  $\rightarrow$  USE SERIES MODEL (e.g. MARKOVIAN) NEXT COURSE

⑤  $Y$  CANNOT BE ASSUMED TO BE NORMAL  $\rightarrow$  GLMs

• ASSUME:  $Y_i \sim \text{Bernoulli}(p_i)$ ,  $Y_1, \dots, Y_m$  INDEPENDENT

$$\rightarrow X \cdot \beta = \begin{pmatrix} 1 & x_{1,1} & \dots & x_{1,p-1} \\ \vdots & \ddots & \ddots & \vdots \\ 1 & x_{m,1} & \dots & x_{m,p-1} \end{pmatrix} \begin{pmatrix} \beta_0 \\ \vdots \\ \beta_{p-1} \end{pmatrix}, \quad E(\bar{Y}) = \begin{pmatrix} p_1 \\ \vdots \\ p_m \end{pmatrix}$$

• PROBLEM:  $p_i \in [0, 1] \rightarrow$  WE SOMETIMES CAN USE  $\hat{p} = X\bar{\beta}$   
 ↳ SOL: IDEA OF GLM's  $\rightarrow$  TRANSFORM  $p: [0, 1] \xrightarrow{\text{LINK FUNCTION}} [-\infty, +\infty]$

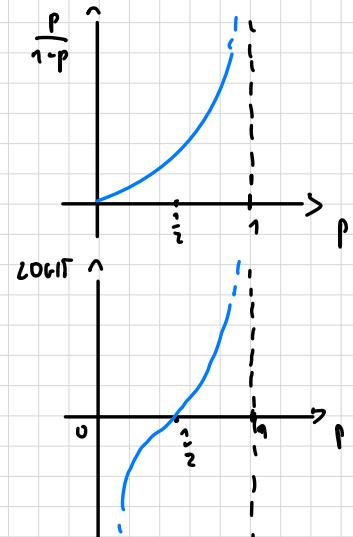
• Odds:

$$\cdot \frac{p}{1-p} : [0, 1] \mapsto [0, +\infty]$$

ex.  $p = 0.05 \rightarrow \text{Odds} = \frac{1}{20} : \text{"ONE TO TWENTY"}$

• LOGIT (LOG ODDS):

$$\cdot \log\left(\frac{p}{1-p}\right) : [0, 1] \mapsto [-\infty, +\infty]$$



- GLM for p:

$$\begin{pmatrix} \text{logit}(p_1) \\ \vdots \\ \text{logit}(p_n) \end{pmatrix} = X \cdot \tilde{\beta} \rightarrow \text{NO "SIGNAL + ERROR" STRUCTURE}$$

$\because Y_i \sim \text{Bernoulli}(p_i)$

- $\text{logit}(p) \mapsto p : l = \log\left(\frac{p}{1-p}\right) \rightarrow e^l = \frac{p}{1-p} \rightarrow (1-p)e^l = p$   
 $[-\infty; +\infty] \mapsto [0, 1]$

$$\rightarrow p = \frac{e^l}{1+e^l} \quad \text{INVERSE LOGIT TRANSFORM.}$$

$$P_i = \frac{e^{X \beta_{(i)}}}{1 + e^{X \beta_{(i)}}}$$

- IN R: <https://stats.oarc.ucla.edu/r/dae/logit-regression/>

```
my_fit <- glm(y ~ predictors, family=binomial)
```

Family	Default link function
binomial	(link = "logit") <span style="color: red;">→ logistic regression</span>
gaussian	(link = "identity") <span style="color: red;">→ LINEAR MODELS</span>
gamma	(link = "inverse")
inverse.gaussian	(link = "1/mu^2")
poisson	(link = "log") <span style="color: red;">→ for count variables, like number of accidents</span>
quasi	(link = "identity", variance = "constant")
quasibinomial	(link = "logit")
quasipoisson	(link = "log")

- METHODS FROM `lm()` CLASS:  $\rightsquigarrow$  ex. `summary(my_fit)`

Function	Description
<code>summary()</code>	Displays detailed results for the fitted model
<code>coefficients()</code> , <code>coef()</code>	Lists the model parameters (intercept and slopes) for the fitted model
<code>confint()</code>	Provides confidence intervals for the model parameters (95 percent by default)
<code>residuals()</code>	Lists the residual values for a fitted model
<code>anova()</code>	Generates an ANOVA table comparing two fitted models
<code>plot()</code>	Generates diagnostic plots for evaluating the fit of a model
<code>predict()</code>	Uses a fitted model to predict response values for a new dataset

$$\hat{P}_f = \frac{e^{\bar{x}_f \cdot \hat{\beta}}}{1 + e^{\bar{x}_f \cdot \hat{\beta}}} : \text{ESTIMATES PROBABILITY OF } Y_f = 1 \text{ CORRESPONDING TO A NEW CONFIGURATION OF PREDICTORS } \bar{x}_f$$

$\rightarrow$  IN ORDER TO ASSIGN CLASS  $0/1$  TO  $\bar{x}_f$   
WE CAN SET A THRESHOLD:  $\hat{P}_f > \text{threshold}$

ex.  $t = \text{threshold} = 0.5$



- CONFUSION MATRIX:

		PREDICTED	
		1	0
TRUE	1	TP	FN
	0	FP	TN

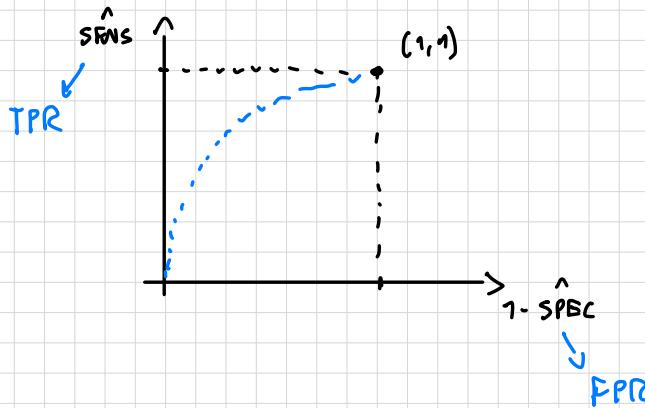
- SENSITIVITY = RECALL( $c=1$ ) =  $P(\hat{Y}=1 | Y=1)$

$$\hookrightarrow \hat{\text{SENSITIVITY}} = \frac{\text{TP}}{\text{TP} + \text{FN}} = \text{TPR}$$

- SPECIFICITY = RECALL( $c=0$ ) =  $P(\hat{Y}=0 | Y=0)$

$$\hookrightarrow \hat{\text{SPECIFICITY}} = \frac{\text{TN}}{\text{TN} + \text{FP}} = 1 - \text{FPR}$$

- IF WE VARY  $t \rightarrow$  CHANGE OF SENSITIVITY AND SPECIFICITY
- ROC CURVE:



ROC CURVE:

$$\hat{\text{SENS}} = \sqrt{1 - \hat{\text{SPEC}}}$$

• INTRODUCTION TO BAYESIAN STATISTICS :

• BAYES THEOREM :  $P(B_n | A) = \frac{P(B_n) P(A | B_n)}{\sum_i P(B_i) P(A | B_i)}$

$\rightarrow P[H_p: H_k, \text{ given evidence } E] : P(H_k | E) = \frac{P(H_k) P(E | H_k)}{\sum_i P(H_i) P(E | H_i)}$

ex. DIAGNOSTIC TEST

SUPPOSE 2  $H_p$  :  $\begin{cases} D: \text{DISEASED} \\ \bar{D}: \text{NOT DISEASED} \end{cases}$

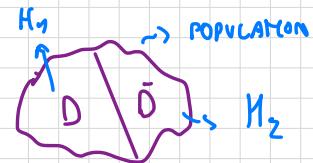
EVIDENCE WILL BE A TEST :  $\begin{cases} T+: \text{says you have } D \\ T-: \text{says you have NOT } D \end{cases}$

$\rightarrow \text{Prob} \approx 1 \quad GT \approx 1$

$$\cdot P(T+ | D) = \text{SENSITIVITY} = TPR$$

$$\cdot P(T+ | \bar{D}) = 1 - \text{SPECIFICITY} = 1 - P(T- | \bar{D}) = FPR$$

$$\rightarrow P(D | T+) = \frac{P(D) \cdot P(T+ | D)}{P(T+)} = \frac{P(D) \cdot P(T+ | D)}{P(D) P(T+ | D) + P(\bar{D}) P(T+ | \bar{D})}$$



BAYES THEOREM FOR R.V.:

$\rightarrow$  R.VEC.

LET  $X, Y$  BE RANDOM VECTORS WITH JOINT DENSITY  $f(x, y)$

- CONDITIONAL DENSITY OF  $X$  GIVEN  $Y = y$ ;

$$f_{X|Y}(x = x | Y = y) = f_{X|Y}(x | y) = \frac{f(x, y)}{f_Y(y)} \xrightarrow{\substack{\uparrow \\ \text{MARGINAL DENSITY OF } Y}} \int f(x, y) dx$$

- CONDITIONAL DENSITY OF  $Y$  GIVEN  $X = x$ ;

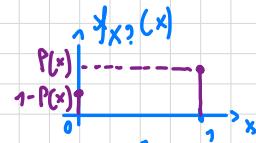
$$f_{Y|X}(y = y | X = x) = f_{Y|X}(y | x) = \frac{f(x, y)}{f_X(x)} \xrightarrow{\substack{\uparrow \\ \text{MARGINAL DENSITY OF } X}} \int f(x, y) dy$$

- APPLYING BAYES THEOREM:

$$\begin{aligned} f_{X|Y=y}(x | y) &= \frac{f_X(x) \cdot f_{Y|X=x}(y | x)}{f_Y(y)} \stackrel{\substack{\uparrow \\ \text{LAW OF} \\ \text{TOTAL PROB.}}}{=} \frac{f_X(x) \cdot f_{Y|X=x}(y | x)}{\int f_X(t) \cdot f_{Y|X=t}(y | t) dt} \\ &\quad \left( \int f(x, y) dx \right) \end{aligned}$$

- INDICATOR OF EVENTS:

$$X: \text{EVENT} \rightarrow X? = \begin{cases} 1, & \text{IF } X \text{ IS TRUE} \\ 0, & \text{IF } X \text{ IS FALSE} \end{cases}, \quad f_{X?}(x) = \begin{cases} P(X), & \text{IF } x = 1 \\ 1 - P(X), & \text{IF } x = 0 \end{cases}$$



Ex. DIAGNOSTIC EXAMPLE, suppose: SENS = 0.9, SPEC = 0.95, PREVALENCE =  $P(D) = 0.01$

$$f_{T?}(x) = \begin{cases} 1, & \text{IF } T+ \\ 0, & \text{IF } T- \end{cases}$$

0/1

$$\text{JOINT DENSITY OF } T? \text{ AND } D? : f_{(T?, D?)}(x, y) = f_{T?|D?}(y | x) \cdot f_{D?}(x)$$

$$\text{EXAMPLE: } f_{(D?, T?)}(0, 1) = f_{D?}(1) \cdot f_{T?|D?}(0 | 1) = 0.01 \cdot (1 - \underset{0.1}{\text{SENS}}) = 0.009$$

2. Consider a binary screening situation formulated in terms of the following random variables:

- $D?$  = indicator whether a randomly sampled person is diseased
- $T?$  = indicator whether the person has a single positive test
- $T_1?, T_2?$  = similar indicators for two conditionally independent tests

and compute the following probabilities.

- (a) Consider a single test. Assuming  $P(D? = 1) = 0.01$  (prevalence),  $P(T? = 1|D? = 1) = 0.90$  (sensitivity) and  $P(T? = 0|D? = 0) = 0.95$  (specificity), compute

$$P(D? = 1|T? = 1) \quad (\text{positive predictive value})$$

$$\text{PREV} = 0.01$$

$$\text{SENS} = 0.9$$

$$\text{SPEC} = 0.95$$

en. 2 a

$$P(D? = 1 | T? = 1) :$$

P OF HAVING DISEASE

GIVEN A POSITIVE RESULT ON A TEST

~

$$P(D? = 1 | T? = 1) = \frac{P(D? = 1, T? = 1)}{P(T? = 1)}$$

Sensitivity

Prevalence

$$\cdot P(D? = 1, T? = 1) = \underbrace{P(T? = 1 | D? = 1) P(D? = 1)}_{\text{Sensitivity}} = 0.90 \cdot 0.01 = 9 \cdot 10^{-3}$$

$$\begin{aligned} \cdot P(T? = 1) &= \underbrace{P(T? = 1 | D? = 0) P(D? = 0)}_{1 - \text{Specificity}} + \underbrace{P(T? = 1 | D? = 1) P(D? = 1)}_{\text{Sensitivity}} \\ &= (1 - 0.95) (1 - 0.01) + 0.9 \cdot 0.01 = 0.0585 \end{aligned}$$

$$\rightarrow P(D? = 1 | T? = 1) = \frac{9 \cdot 10^{-3}}{0.0585} = 0.1538 \sim \text{RELATIVELY LOW}$$

en. 2 b

- (b) Consider repeating two conditionally independent test on the same person. Assume sensitivity and specificity of  $T_1?$  and  $T_2?$  are the same as  $T?$ . Compute

$$P(D? = 1 | T_1? = 1, T_2? = 1)$$

$$\rightarrow P(D? = 1 | T_1? = 1, T_2? = 1) = \frac{P(D? = 1, T_1? = 1, T_2? = 1)}{P(T_1? = 1, T_2? = 1)}$$

INDEP.

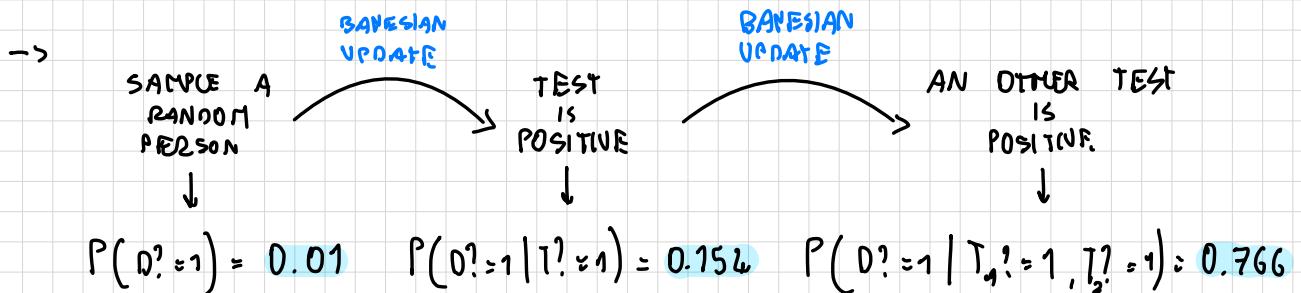
$$\begin{aligned} \cdot P(D? = 1, T_1? = 1, T_2? = 1) &= P(T_1? = 1, T_2? = 1 | D? = 1) P(D? = 1) = \\ &= P(T_1? = 1 | D? = 1) P(T_2? = 1 | D? = 1) P(D? = 1) \\ &= 0.9 \cdot 0.9 \cdot 0.01 \end{aligned}$$

$$\begin{aligned} \cdot P(T_1? = 1, T_2? = 1) &= P(T_1? = 1, T_2? = 1 | D? = 1) P(D? = 1) + \\ &\quad + P(T_1? = 1, T_2? = 1 | D? = 0) P(D? = 0) = \dots = \\ &= 0.9 \cdot 0.9 \cdot 0.1 + (1 - 0.95)(1 - 0.95) \cdot (1 - 0.01) \end{aligned}$$

$$\rightarrow P(D? = 1 | T_1? = 1, T_2? = 1) = \frac{0.9 \cdot 0.9 \cdot 0.01}{0.9^2 \cdot 0.1 + 0.05^2 \cdot 0.99} = 0.766$$

## BAYESIAN UPDATE :

IT IS POSSIBLE TO UPDATE  $P(D)$  GIVEN NEW EVIDENCE OF T?



ex. 3

- $(1, 1)$  : ACCURS SARS  
"PERSON HAS DISEASE"
- $(0, 0)$  : NEVER SARS  
"PERSON HAS NO DISEASE"

→

TPR : SENSITIVITY

FPR : 1 - SPECIFICITY

•  $t_1$  :  $\left\{ \begin{array}{l} TPR = \frac{TP}{TP + FN} = \frac{45}{45 + 5} = 0.9 \\ FPR = \frac{FP}{FP + TN} = \frac{5}{5 + 95} = 0.05 \end{array} \right.$

$$P1 = (0.05, 0.9)$$

•  $t_2$  :  $\left\{ \begin{array}{l} TPR = \frac{35}{35 + 15} = 0.7 \\ FPR = \frac{15}{15 + 85} = 0.17 \end{array} \right.$

$$P2 = (0.17, 0.7)$$

∴  $\rightarrow P3 = (0.1, 0.5)$

•  $t_4$  :  $\left\{ \begin{array}{l} TPR = \frac{15}{50} = 0.3 \\ FPR = \frac{5}{50} = 0.1 \end{array} \right.$

$$P4 = (0.1, 0.3)$$

3. The following example is adapted from the *Encyclopedia of Biostatistics*, Wiley (2005). Suppose 100 non-diseased ( $D^-$ ) patients and 50 diseased ( $D^+$ ) patients have also been classified on the basis of radiological exams over five ordered levels

--- = very mild

- = mild

+- = neutral

+ = serious

++ = very serious

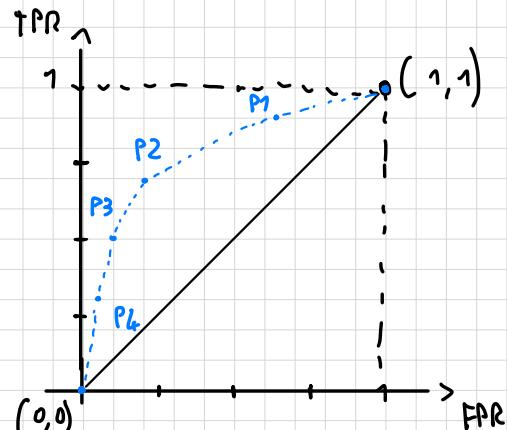
Here are the results:

	$t_0$	$t_1$	$t_2$	$t_3$	$t_4$	$t_5$	
$D^-$	---	-	+-	+	++		total
$D^+$	35	45	10	5	5	100	

	$t_0$	$t_1$	$t_2$	$t_3$	$t_4$	$t_5$	
$D^-$	5	10	10	10	15		50
$D^+$							1

Using the only four possible cut-off choices and the default points  $(0,0)$  and  $(1,1)$ , draw an empirical ROC curve.



## MULTIVARIATE DISTRIBUTIONS AND DAGs:

CONSIDER A R.VEC IN  $\mathbb{R}^d$   $(x_1, \dots, x_d)$ , WITH JOINT DENSITY  $f(x_1, \dots, x_d)$

USING THE DEFINITION OF CONDITIONAL DENSITY:

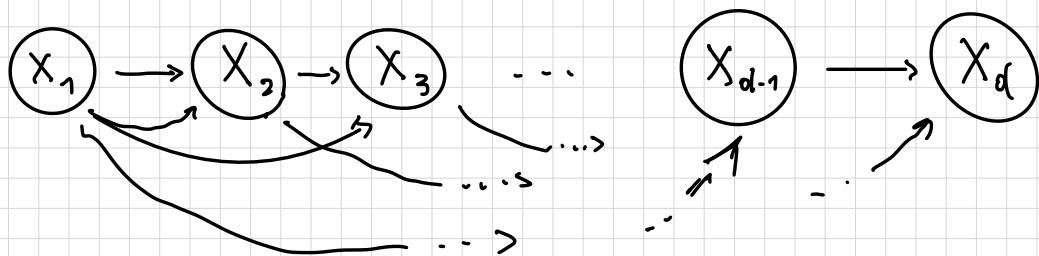
$$\cdot f(x_1, \dots, x_d) =$$

$$= f_{x_1, \dots, x_{d-1}}(x_1, \dots, x_{d-1}) \cdot f_{x_d | x_1, \dots, x_{d-1}}(x_d | x_1, \dots, x_{d-1}) =$$

$$= f_{x_1, \dots, x_{d-2}}(x_1, \dots, x_{d-2}) \cdot f_{x_{d-1} | x_1, \dots, x_{d-2}}(x_{d-1} | x_1, \dots, x_{d-2}) \cdot f_{x_d | x_1, \dots, x_{d-1}}(x_d | x_1, \dots, x_{d-1})$$

$$= \dots = f_{x_1}(x_1) \cdot f_{x_2 | x_1}(x_2 | x_1) \cdot f_{x_3 | x_1, x_2}(x_3 | x_1, x_2) \cdot \dots \cdot f_{x_d | x_1, \dots, x_{d-1}}(x_d | x_1, \dots, x_{d-1})$$

$\rightarrow$  DAG:



• FOR SOME DISTRIBUTIONS  $\rightarrow$  ~~for~~ SOME ARROWS  $\rightarrow$  JOINT DENSITY SIMPLIFIES

ex.

$$X \rightarrow Y \rightarrow Z \rightarrow \dots \rightarrow f(x, y, z) = f_X(x) \cdot f_{Y|X}(y|x) \cdot f_{Z|Y,X}(z|y,x)$$

• NOW SUPPOSE  $f_{Z|X,Y}(z|x,y) \stackrel{*}{=} f_{Z|Y}(z|y) \rightsquigarrow$  ~~DOES NOT DEPEND ON X~~

$$\cdot f_{X,Y,Z}(x,y,z) = f_{X|Y}(x|y) \cdot f_{Z|X,Y}(z|x,y) \stackrel{*}{=} f_{X|Y}(x|y) \cdot f_{Z|Y}(z|y)$$

$$\rightarrow f(x, y, z) = f_X(x) \cdot f_{Y|X}(y|x) \cdot f_{Z|Y}(z|y)$$

- MARKOVIAN CHAINS :

$(X_1) \rightarrow (X_2) \rightarrow \dots \rightarrow (X_d)$  :

$$f(x_1, \dots, x_d) = f_{X_1}(x_1) \cdot f_{X_2|X_1}(x_2|x_1) \cdot f_{X_3|X_2}(x_3|x_2) \cdot \dots \cdot f_{X_d|X_{d-1}}(x_d|x_{d-1})$$

- INDEPENDENT RVs:

$(X_1) \quad (X_2) \quad (X_d)$  :  $f(x_1, \dots, x_d) = \prod_i f_{X_i}(x_i)$

- FACTORIZATION THEOREM FOR DAGs :

$$f(x_1, \dots, x_d) = \prod_i f_{X_i | \text{PARENTS}(X_i)}(x_i | \text{PARENTS}(X_i))$$

ex.  $X \rightarrow Y \leftarrow Z$        $Y = \text{parents}(w, z)$        $X = \text{parents}(y)$   
 $f(x) = 0 \Rightarrow f(x, r, w, z) = f_{Z|Y}(z|y) \cdot f_{W|Y}(w|y) \cdot f_{Y|X}(r|x) \cdot f_X(x)$

ex. MARKOVIAN CHAIN :  $X_1 \rightarrow X_2 \rightarrow \dots \rightarrow X_d$

$$\rightarrow f(x_1, \dots, x_d) = f_{X_1}(x_1) \cdot f_{X_2|X_1}(x_2|x_1) \cdot \dots \cdot f_{X_d|X_{d-1}}(x_d|x_{d-1})$$

DAG :

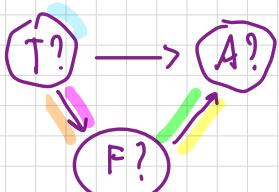
$$\cdot T? = \begin{cases} 1, & \text{CARD HOLDER IS TRAVELING} \\ 0, & \text{-- IS NOT --} \end{cases}$$

$$\cdot F? = \begin{cases} 1, & \text{TRANSACTION IS FREQUENT} \\ 0, & \text{-- IS NOT --} \end{cases}$$

$$\cdot A? = \begin{cases} 1, & \text{PURCHASE IS ABROAD} \\ 0, & \text{-- IS NOT --} \end{cases}$$

→

DAG:



BINARN BAYESIAN NETWORK

↗

$$\cdot f(x, v, z) = f_{T?}(x) \cdot f_{F?|T?}(v|x) \cdot f_{A?|T?, F?}(z|v, x)$$

• $T?$ : $T? = 0$	$ $	$T? = 1$
$P$ : 0.95		0.05

$$\cdot F? | T? : \quad T? = 0 \quad T? = 1$$

$F? = 0$	0.99	0.98
$F? = 1$	0.01	0.02

$$\cdot A? | T?, F?$$

$T? = 0$		$T? = 1$	
$F? = 0$		$F? = 1$	$F? = 1$
$A? = 0$	0.99	0.9	0.1
$A? = 1$	0.01	0.7	0.9

$$\rightarrow f(0, 0, 0) = 0.95 \cdot 0.99 \cdot 0.99$$

Suppose you are working for a financial institution and you are asked to build a fraud detection system. You plan to use the following information. When the card holder is traveling abroad, fraudulent transaction are more likely since tourists are prime targets for thieves. More precisely, 2% of transactions are fraudulent when the card holder is traveling, whereas only 1% of the transactions are fraudulent when he is not traveling. On average, 5% of all transactions happen while card holder is traveling. If a transaction is fraudulent, then the likelihood of a purchase abroad increases, unless the card holder happens to be traveling. More precisely, when the card holder is not traveling, 10% of the fraudulent transactions are abroad purchases, whereas only 1% of the legitimate transactions are abroad purchases.

## • CONDITIONAL EXPECTATIONS:

GIVEN A PAIR OF R.V.  $X, Y$ :

$$E[Y|X=x] = \begin{cases} \sum y f_{Y|X}(y|x) & \text{DISCRETE} \\ \int y f_{Y|X}(y|x) dy & \text{CONTINUOUS} \end{cases} / \int f_{Y|X}(y|x) dy = \frac{\int f_{Y|X}(y|x)}{\int f_X(x)}$$

ex. BIVARIATE NORMAL

$$\begin{pmatrix} X \\ Y \end{pmatrix} \sim N_2 \left( \begin{pmatrix} \mu_x \\ \mu_y \end{pmatrix}, \begin{pmatrix} \sigma_x^2 & \rho \sigma_x \sigma_y \\ \rho \sigma_x \sigma_y & \sigma_y^2 \end{pmatrix} \right) / \rho = \frac{\text{Cov}(X, Y)}{\sigma_x \sigma_y}$$

$$\rightarrow E[Y|X=x] = \mu_y + \rho \frac{\sigma_y}{\sigma_x} (x - \mu_x) \times$$

$\nwarrow$  IT'S A NUMBER  
NOT A R.V

$\hookrightarrow$  LINEAR FUNCTION OF  $X$

$X$  = HEIGHT OF A FATHER,  $Y$  = HEIGHT OF HIS SON

SUPPOSE  $X = 1.95$ ,  $\mu_x = 1.75 = \mu_y$ ,  $\rho = 0.9$ ,  $\sigma_x = \sigma_y = 0.1$

$\rightarrow$

$$E[Y|X=1.95] = 1.75 + 0.9 \frac{0.1}{0.1} (1.95 - 1.75) = 1.93$$

$\rightarrow$  NOT AS TALL AS HIS FATHER

$$\frac{E[Y|X=x] - \mu_y}{\sigma_y} = \rho \frac{x - \mu_x}{\sigma_x}$$

$\overbrace{\quad \quad \quad}$  HOW MANY STD THE SON IS DEVIATING FROM MEAN

$\overbrace{\quad \quad \quad}$  HOW MANY STD THE FATHER IS DEVIATING FROM MEAN

✓  
REGRESSION TO THE MEAN EFFECT (WALTON, 1895)

• PROPERTIES:

1.  $E[E[Y|X]] = E[Y]$

MAIN PROPERTY

PROOF:

$$\begin{aligned} E \left[ \int_Y f_{Y|X}(y|x) dy \right] &= \int_X \int_Y f_{Y|X}(y|x) f_X(x) dx = \\ &= \int_X \int_Y y \cdot \frac{f_{X,Y}(x,y)}{f_X(x)} f_X(x) dx = \int_X y \cdot \underbrace{\int_X f_{X,Y}(x,y) dx}_{d_Y(y)} dy = E[Y] \end{aligned}$$

⇒ BIUARIAVE NORMAL

$$\begin{aligned} E[E[Y|X]] &= E\left[\mu_Y + \rho \frac{\sigma_Y}{\sigma_X} (X - \mu_X)\right] = \\ &= \mu_Y + \rho \frac{\sigma_Y}{\sigma_X} E[X - \mu_X] = \mu_Y = E[Y] \end{aligned}$$

2.  $E(a|X) = a \quad / a = \text{const.}$

3.  $E(aX + bY|Z) = aE(X|Z) + bE(Y|Z)$

4.  $E(g(X) \cdot Y|X) = g(X) \cdot E(Y|X)$

• CONDITIONAL VARIANCE:

$$V(Y|X) = E[(Y - E(Y|X))^2 | X]$$

5.  $E[V(Y|X)] = V(Y) - V(E(Y|X))$

$\hookrightarrow V(Y) = E[V(Y|X)] + V(E(Y|X))$  OTHER WAY TO COMPUTE Var

PROOF:

$$\begin{aligned} E(V(Y|X)) &= E\left(E((Y - E(Y|X))^2 | X)\right) = E\left(E(Y^2 + E^2(Y|X) - 2Y \cdot E(Y|X) | X)\right) \stackrel{\text{LINEARITY}}{=} E\left(E(Y^2 | X) + E(E(E(Y|X) | X)) - 2E(E(Y|X) | X)\right) \\ &\stackrel{\text{MAIN PROPERTY}}{=} E(Y^2) + E(E^2(Y|X) \cdot \underbrace{E(1|X)}_{=1}) - 2E(E(Y|X) \cdot E(M|X)) = E(Y^2) - E(E^2(Y|X)) \left( + E^2(Y) - E^2(Y) \right) = \\ &= E(Y^2) \cdot E^2(Y) + E^2(Y) - E(E^2(Y|X)) = V(Y) - (E(E^2(Y|X)) - E^2(Y)) = V(Y) - V(E(Y|X)) \end{aligned}$$

## BAYESIAN STATISTICS:

- $\theta$  is a parameter itself, a R.V. (or R.VEC.)

as.  $\theta = \begin{pmatrix} M \\ \sigma^2 \end{pmatrix}$

- BEFORE SEEING ANY DATA,  $\theta$  HAS A PRIOR DISTRIBUTION

- AFTER SEEING DATA, WE UPDATE OUR PRIOR DISTRIBUTION ON  $\theta$  INTO A POSTERIOR DISTRIBUTION, USING BAYES THEOREM

- GIVEN  $\theta$  A VECTOR WITH PRIOR DENSITY  $\tilde{\Pi}(\theta)$ :

$$\underbrace{\tilde{\Pi}(\theta | \text{DATA})}_{\substack{\text{POSTERIOR} \\ \hookrightarrow \text{OPEN JUST } \tilde{\Pi}(\theta | \text{DATA})}} = \frac{\cancel{\times}(\text{DATA} | \theta) \cdot \underbrace{\tilde{\Pi}(\theta)}_{\substack{\text{PRIOR}}}}{\int \tilde{\Pi}(\theta) \cdot \cancel{\times}(\text{DATA} | \theta) d\theta}$$

The diagram shows a circle containing  $\theta$  with an arrow pointing down to a bracket labeled  $\{\text{DATA}\}$ . To the right of this, the text "BAYES THEOREM" is written.

- SUPPOSE :

$$X_1, \dots, X_n | \theta \stackrel{\text{iid}}{\sim} \mathcal{N}(\theta, \sigma^2)$$

WE ASSUME  $\theta \sim N(\theta_0, \sigma_0^2) = \underbrace{\tilde{\Pi}(\theta)}_{\substack{\text{PRIOR}}} / \theta_0, \sigma_0^2 \text{ ARE KNOWN}$

$$\rightarrow \underbrace{\times(x_1, \dots, x_n | \theta)}_{\substack{\text{DATA} \\ \rightarrow \text{NORMAL DENSITY}}} = \prod_{i=1}^n \times(x_i | \theta) = \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x_i - \theta)^2} =$$

$$\sum (x_i - \theta)^2 = \sum (x_i^2 - 2\theta x_i + \theta^2) = (2\pi\sigma^2)^{-\frac{m}{2}} \cdot \exp \left\{ -\frac{1}{2\sigma^2} \sum (x_i - \theta)^2 \right\} =$$

$$= \sum_i x_i^2 - 2\theta \sum_i x_i + m\theta^2 = \dots = (2\pi\sigma^2)^{-\frac{m}{2}} \cdot \exp \left\{ -\frac{1}{2\sigma^2} \left( \sum x_i^2 - 2\theta \sum x_i + m\theta^2 \right) \right\}$$

$$\therefore \sum_i x_i^2 - 2\theta \sum x_i - m\theta^2$$

$$n\bar{x} = n \cdot \frac{\sum x_i}{m}$$

• APPLING BAYES THEOREM:

$$\tilde{\Pi}(\theta | \text{DATA}) = \frac{\mathbb{P}(\text{DATA} | \theta) \cdot \tilde{\Pi}(\theta)}{\int \tilde{\Pi}(t) \cdot \mathbb{P}(\text{DATA} | t) dt} \propto \mathbb{P}(\text{DATA} | \theta) \cdot \tilde{\Pi}(\theta)$$

$\Rightarrow \text{const.}$

•  $\mathbb{P}(\text{DATA} | \theta) \cdot \tilde{\Pi}(\theta) =$

$$= \frac{1}{\sqrt{2\pi\sigma_0^2}} \cdot \exp\left\{-\frac{1}{2\sigma_0^2} (\theta - \theta_0)^2\right\} \cdot \left(\frac{1}{2\pi\sigma^2}\right)^{\frac{n}{2}} \cdot \exp\left\{-\frac{1}{2\sigma^2} \left(\sum_i x_i^2 - 2\theta n\bar{x} + n\theta^2\right)\right\}$$

$\Rightarrow \theta^2 - 2\theta\theta_0 + \theta_0^2$

$\text{const.}$

$$\propto \exp\left\{-\frac{1}{2\sigma_0^2} (\theta - \theta_0)^2\right\} \cdot \exp\left\{-\frac{1}{2\sigma^2} \left(\sum_i x_i^2 - 2\theta n\bar{x} + n\theta^2\right)\right\}$$

• GIVEN  $\begin{cases} \tilde{\tau} = \frac{1}{\sigma^2} & : \text{PRECISION} \\ \tilde{\tau}_0 = \frac{1}{\sigma_0^2} & : \text{PRIOR PRECISION} \end{cases}$

CONST.  $\rightarrow$  REMOVED FOR PROPORTIONALITY

$$\rightarrow \exp\left\{ \dots \cdot \exp\left\{ \dots \right\} = \exp\left\{ -\frac{1}{2} \tilde{\tau}_0 (\theta^2 - 2\theta\theta_0 + \theta_0^2) - \frac{1}{2} \tilde{\tau} \left( \sum_i x_i^2 - 2\theta n\bar{x} + n\theta^2 \right) \right\} \right\}$$

$$= \exp\left\{ -\frac{1}{2} \left( \tilde{\tau}_0 \theta^2 - 2\tilde{\tau}_0 \theta \theta_0 - 2\tilde{\tau} \theta n\bar{x} + \tilde{\tau} n\theta^2 \right) \right\} =$$

$$= \exp\left\{ -\frac{1}{2} \left( (\tilde{\tau}_0 + n\tilde{\tau}) \theta^2 - 2(\tilde{\tau}_0 \theta_0 + n\tilde{\tau} \bar{x}) \theta \right) \right\} =$$

$$= \exp\left\{ -\frac{1}{2} (\tilde{\tau}_0 + n\tilde{\tau}) \left( \theta^2 - 2 \frac{\tilde{\tau}_0 \theta_0 + n\tilde{\tau} \bar{x}}{\tilde{\tau}_0 + n\tilde{\tau}} \theta \right) \right\} \quad \left( \begin{array}{l} \text{TRICK TO COMPUTE THE SQUARE} \\ \text{+ } \left( \frac{\tilde{\tau}_0 \theta_0 + n\tilde{\tau} \bar{x}}{\tilde{\tau}_0 + n\tilde{\tau}} \right)^2 \end{array} \right) =$$

$$\propto \exp\left\{ -\frac{1}{2} (\tilde{\tau}_0 + n\tilde{\tau}) \left( \theta - \frac{\tilde{\tau}_0 \theta_0 + n\tilde{\tau} \bar{x}}{\tilde{\tau}_0 + n\tilde{\tau}} \right)^2 - \left( \frac{\tilde{\tau}_0 \theta_0 + n\tilde{\tau} \bar{x}}{\tilde{\tau}_0 + n\tilde{\tau}} \right)^2 \cdot \left( -\frac{1}{2} (\tilde{\tau}_0 + n\tilde{\tau}) \right) \right\}$$

$$\propto \exp\left\{ -\frac{1}{2} (\tilde{\tau}_0 + n\tilde{\tau}) \left( \theta - \frac{\tilde{\tau}_0 \theta_0 + n\tilde{\tau} \bar{x}}{\tilde{\tau}_0 + n\tilde{\tau}} \right)^2 \right\} \propto \tilde{\Pi}(\theta | x_1, \dots, x_m)$$

MEAN OF A NORMAL DISTRIBUTION

$$\rightarrow \theta | x_1, \dots, x_m \sim N\left(\frac{\tilde{\tau}_0 \theta_0 + n\tilde{\tau} \bar{x}}{\tilde{\tau}_0 + n\tilde{\tau}}, \frac{1}{\tilde{\tau}_0 + n\tilde{\tau}}\right)$$

## BAYESIAN ANALYSIS OF BINARY RANDOM SAMPLE:

- RANDOM SAMPLES MAY BE BINARY, e.g. BERNoulli:

$$X_1, \dots, X_n \mid p \stackrel{\text{ciel}}{\sim} \text{Bernoulli}(p)$$

PRIOR

- $\tilde{H}(p)$ :

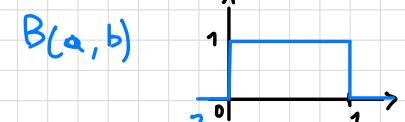
- WE COULD THINK IT AS A DISCRETE R.V.  $\rightarrow$  DISCRETE BNs
- ... OR AS A CONTINUOUS R.V., e.g. WITH A BETA DISTRIBUTION:

$$p \sim \text{Beta}(a, b) \rightarrow \tilde{H}(p) = \frac{P^{a-1} (1-p)^{b-1}}{\int_0^1 t^{a-1} (1-t)^{b-1} dt}, \quad 0 < p < 1$$

P<sup>a-1</sup> (1-p)<sup>b-1</sup>  
∫<sub>0</sub><sup>1</sup> t<sup>a-1</sup> (1-t)<sup>b-1</sup> dt  
NORMALIZING CONSTANT

- SPECIAL CASE,  $a = b = 1$ :

$$\tilde{H}(p) = \frac{P^{1-1} (1-p)^{1-1}}{\int_0^1 dt} = \frac{1}{\int_0^1 dt} = \tilde{H}_{[0,1]}(p)$$



- $E[p] = \int_0^1 p \cdot \frac{P^{a-1} (1-p)^{b-1}}{B(a, b)} dp = \frac{B(a+1, b)}{B(a, b)} = \frac{a}{a+b}$

- $V[p] = \frac{ab}{(a+b)^2 (a+b+1)}$

- N.B.: EUCLER'S BETA FUNCTION:  $B(a, b) = \int_0^1 t^{a-1} (1-t)^{b-1} dt$

- $B(a, b) = \frac{T(a) T(b)}{T(a+b)} \quad / \quad \zeta \Gamma(x) = \int_0^\infty t^{x-1} e^{-t} dt$

- PROPERTY:  $\Gamma(x) = (x-1) \Gamma(x-1) \quad / \quad \text{IF } X = n \in \mathbb{Z}: \quad \Gamma(n) = (n-1)!$

- $X_1, \dots, X_n | p \stackrel{\text{conditional}}{\sim} \text{Bernoulli}(p) / p \sim \text{Beta}(a, b)$

- CONDITIONAL DENSITY:

$$f(x_1, \dots, x_n | p) = \prod_{i=1}^n f_{X_i|p}(x_i | p) = \prod_{i=1}^n p^{x_i} (1-p)^{n-x_i}$$

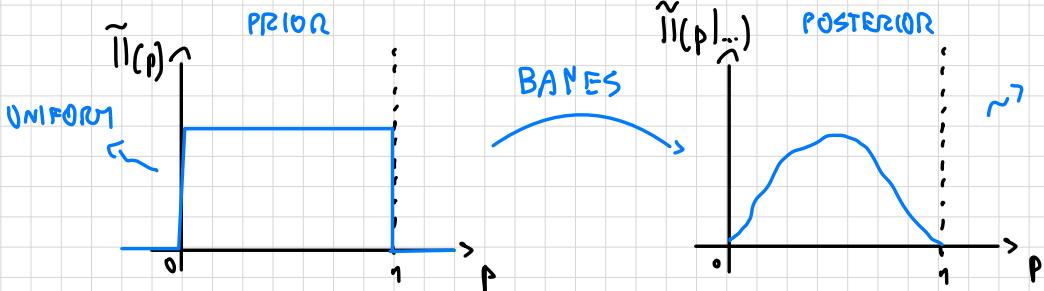
↳  $L(p, x_1, \dots, x_n)$   
 ↳ LIKELIHOOD FUNCTION =  $p^{\sum x_i} \cdot (1-p)^{n-\sum x_i}$

- BAYES POSTERIOR:

$$\begin{aligned} \overline{f}(p | x_1, \dots, x_n) &= \frac{\tilde{f}(p) \cdot f(x_1, \dots, x_n | p)}{\int_0^1 \tilde{f}(q) \cdot f(x_1, \dots, x_n | q) dq} = \\ &\underset{\text{CONST. PROP.}}{=} \frac{p^{a-1} (1-p)^{b-1} \cdot p^{\sum x_i} \cdot (1-p)^{n-\sum x_i}}{\left[ \int_0^1 t^{a-1} (1-t)^{b-1} dt \right]} \propto p^{a-1} (1-p)^{b-1} \cdot p^{\sum x_i} (1-p)^{n-\sum x_i} \\ &\underset{\text{CONST. PROP.}}{=} \frac{\left[ \int_0^1 \tilde{f}(q) \cdot f(x_1, \dots, x_n | q) dq \right]}{\text{IT IS } \rightarrow \text{BETA AGAIN}} \\ &\cdot p^{a-1} (1-p)^{b-1} \cdot p^{\sum x_i} (1-p)^{n-\sum x_i} = p^{a + \sum x_i - 1} (1-p)^{b + n - \sum x_i - 1} \end{aligned}$$

$$\rightarrow \overline{f}(p | x_1, \dots, x_n) = \frac{p^{a + \sum x_i - 1} (1-p)^{b + n - \sum x_i - 1}}{\int_0^1 t^{a + \sum x_i - 1} (1-t)^{b + n - \sum x_i - 1} dt}$$

$$\rightarrow p | x_1, \dots, x_n \sim \text{Beta}(a + \sum x_i, b + n - \sum x_i)$$

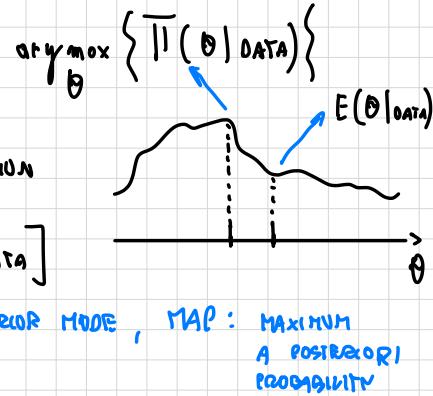


## BAYESIAN STATISTICS CORE:

### POINT ESTIMATION:

IT IS AN APPROXIMATE SUMMARY OF THE POSTERIOR DISTRIBUTION

$\hookrightarrow$  ex. POSTERIOR EXPECTED VALUE:  $E[\theta | \text{DATA}]$



### NORMAL RANDOM SAMPLE:

$$E(\mu | x_1, \dots, x_n) = \frac{\tilde{\tau}_0 \mu_0 + n \tilde{x}}{\tilde{\tau}_0 + n \tilde{\tau}} = \frac{\tilde{\tau}_0}{\tilde{\tau}_0 + n \tilde{\tau}} \mu_0 + \frac{n \tilde{x}}{\tilde{\tau}_0 + n \tilde{\tau}} \bar{x} =$$

WEIGHTED  
AVG OF:

•  $\mu_0$ : prior mean  
 •  $\bar{x}$ : sample mean  $\hookrightarrow = \left(1 - \frac{n \tilde{\tau}}{\tilde{\tau}_0 + n \tilde{\tau}}\right) \mu_0 + \frac{n \tilde{\tau}}{\tilde{\tau}_0 + n \tilde{\tau}} \bar{x} / w = \frac{n \tilde{\tau}}{\tilde{\tau}_0 + n \tilde{\tau}}$

•  $\downarrow \tilde{\tau}_0 \quad (\rightarrow \text{HIGHER } \tilde{\sigma}_0^2 : \text{ HIGHER UNCERTAINTY}) \rightarrow \uparrow \text{ DATA IMPORTANCE FOR ESTIMATION}$

• IF  $n \rightarrow \infty$ :  $\lim_{n \rightarrow \infty} \frac{n \tilde{\tau}}{\tilde{\tau}_0 + n \tilde{\tau}} = 1 \rightarrow \text{ONLY DATA COUNTS}$

$$\hookrightarrow n \rightarrow \infty : \underbrace{\left(1 - \frac{n \tilde{\tau}}{\tilde{\tau}_0 + n \tilde{\tau}}\right)}_{\rightarrow 0} \mu_0 + \underbrace{\frac{n \tilde{\tau}}{\tilde{\tau}_0 + n \tilde{\tau}}}_{\rightarrow 1} \bar{x}$$

### BINARY DATA CASE, Beta(a, b):

$$E(p | x_1, \dots, x_m) = \frac{a + \sum x_i}{a + \sum x_i + b + n - \sum x_i} = \frac{a + \sum x_i}{a + b + m} =$$

$$= \frac{a + b}{a + b + m} \cdot \frac{a}{a + b} + \frac{m}{a + b + m} \cdot \frac{\sum x_i}{m} =$$

SAME PROPERTIES AS ABOVE  $\hookrightarrow = \left(1 - \frac{m}{a + b + m}\right) \frac{a}{a + b} + \frac{m}{a + b + m} \bar{x} \hookrightarrow \hat{p} : \text{MLE OF } p$

• IF  $a = b = 1$ :  $E(p | x_1, \dots, x_m) = \frac{1 + \sum x_i}{2 + m} = \hat{p}$  BAMES ESTIMATION  
 $\hookrightarrow$  IT UNIFORM ON P: GOOD PRIOR IGNORANCE

BTW:  $x_1 = \dots = x_m = 1 \rightarrow \hat{p} = \frac{1 + m}{2 + m}$  LAPLACE SUCCESSION RULE