

$\frac{4}{8} \quad \frac{5}{2} \quad \boxed{\frac{6}{9}} \quad \frac{7}{6} \quad \frac{8}{4}$

①

August 8, 2019

# TOWER PROPERTY OF CONDITIONAL EXPECTATION:

A fair coin is tossed until two tails occur successively. Let  $N$  be the total number of tosses required to terminate the experiment. Compute  $\mathbb{E}[N]$ .

The essence of the tower property is the following:

$$\mathbb{E}[N] = \mathbb{E}[\mathbb{E}[N|X]]$$

$\uparrow$  w.r.t.  $X$        $\uparrow$  w.r.t.  $N$ , given  $X$

obviously,  $X$  must be wisely chosen to be helpful.

Let  $X = \begin{cases} 1 & \text{if the first toss results in tails} \\ 0 & \text{if the first toss results in heads} \end{cases}$

$\uparrow$  indicator r.v.  
dummy r.v.

$$\mathbb{E}[N] = \mathbb{E}[\mathbb{E}[N|X]] \quad (\text{work from outside, in})$$

$$= \mathbb{E}[N|X=1]P(X=1) + \mathbb{E}[N|X=0]P(X=0)$$

$$= \left( \left( \frac{1}{2} \right) (2 + \mathbb{E}[N]) + \left( \frac{1}{2} \right) (2) \right) \cdot \frac{1}{2} + (1 + \mathbb{E}[N]) \cdot \frac{1}{2}$$

$\uparrow$  heads  
~~heads~~  
on second flip

$\uparrow$  tails  
~~heads~~  
on second flip

$$\Rightarrow \mathbb{E}[N] = (1 + \frac{1}{2}\mathbb{E}[N] + 1) \cdot \frac{1}{2} + \frac{1}{2} + \frac{1}{2}\mathbb{E}[N]$$

$$\Rightarrow \mathbb{E}[N] = 1 + \frac{1}{2}\mathbb{E}[N] + \frac{1}{2} + \frac{1}{2}\mathbb{E}[N] \Rightarrow \frac{1}{2}\mathbb{E}[N] = \frac{3}{2} \Rightarrow \mathbb{E}[N] = 6$$

(2)

DEF: If  $\vec{X} \sim N(\vec{\mu}, \Sigma)$ , we say that  $\vec{X}$  is a multivariate Gaussian random vector with mean vector  $\vec{\mu}$  and variance-covariance matrix  $\Sigma$ .

NOTE: If  $\vec{x}$  is  $k$ -dimensional, the  $\vec{\mu} \in \mathbb{R}^k$  and  $\Sigma \in \mathbb{R}^{k \times k}$ . Moreover,  $\Sigma$  must be positive semi-definite.

RECALL:  $\Sigma$  is positive semi-definite iff  $\forall \vec{x} \in \mathbb{R}^k$ ,  
 $\vec{x}^T \Sigma \vec{x} \geq 0$ .

RECALL: If all eigenvalues are  $\geq 0$ , then  $\Sigma$  is positive & semi-definite.

Qn: What is the density function for  $\vec{X}$ ?

ANS:  $f(x_1, \dots, x_k) = f(\vec{x}) = (2\pi)^{-\frac{k}{2}} |\Sigma|^{-\frac{1}{2}} \exp\left\{-\frac{1}{2} (\vec{x} - \vec{\mu})^T \Sigma^{-1} (\vec{x} - \vec{\mu})\right\}$

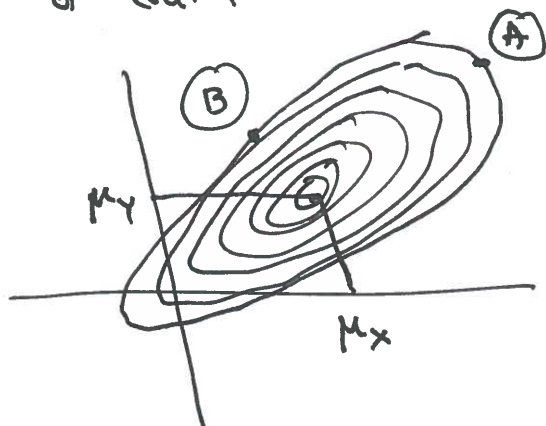
FUN FACT: Call  $d(\vec{x}, \vec{y}) = (\vec{x} - \vec{y})^T \Sigma^{-1} (\vec{x} - \vec{y})$ .

This distance is called Mahalanobis distance.

this is the variance-covariance matrix of some underlying set of data

$\Rightarrow$  key role in discriminant analysis  
 $\Rightarrow$  key role in  $k$ -means clustering

(A) is farther than (B) from  $(\mu_x, \mu_y)$ .  
 (A) and (B) are equidistant from  $(\mu_x, \mu_y)$  under Mahalanobis.



# FACTS ABOUT MULTIVARIATE GAUSSIAN RANDOM VECTORS

Suppose that  $\vec{X} = (X_1, \dots, X_k)$  with  $\vec{X} \sim N(\vec{\mu}, \Sigma)$ .

① Every linear combination of components of  $\vec{X}$  is again Gaussian. In other words,

$$a_1 X_1 + a_2 X_2 + \dots + a_k X_k = \vec{a}^T \vec{X} \quad \leftarrow \text{this is Gaussian}$$

$$\Rightarrow \mathbb{E}[\vec{a}^T \vec{X}] = \vec{a}^T \cdot \mathbb{E}[\vec{X}] = \vec{a}^T \cdot \vec{\mu}$$

$$\Rightarrow \text{Var}(\vec{X}) = \Sigma$$

by definition, i.e., this is how we will extend the notion of variance to random vectors

$$\Rightarrow \text{Var}(\vec{a}^T \vec{X}) = \vec{a}^T \Sigma \vec{a}$$

② For any such  $\vec{X}$ , there exists a matrix  $A$  so that

$$\vec{X} = A \cdot \vec{Z} + \vec{\mu}$$

where  $\vec{Z}$  is a  $k$ -dimensional vector of independent standard normal random variables.

clearly:  $A$  is related to a "square root" of  $\Sigma$ , which  $\Sigma$  always has because (1) it is positive semi-definite and (2) it is symmetric.

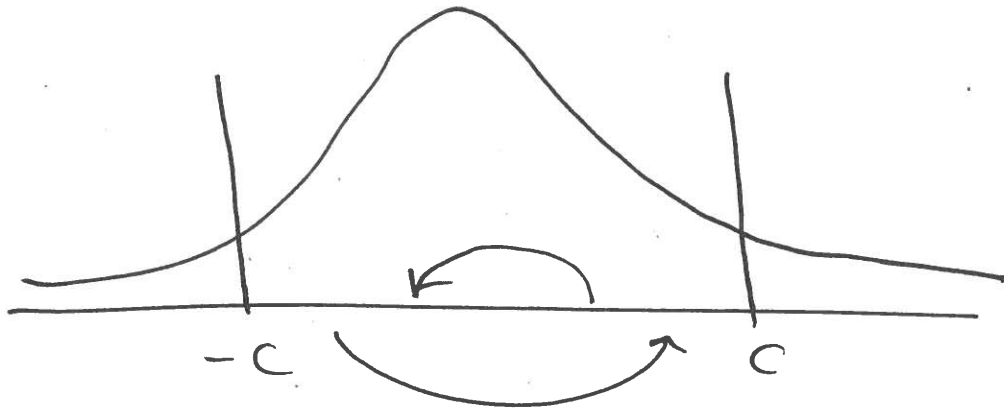
③ For any  $k$ , the level curves (or level sets) of the density  $f(\vec{x})$  are ellipsoids.

④ If  $X_i$  and  $X_j$  are such that  $\rho(X_i, X_j) = 0$ , then  $X_i$  and  $X_j$  are independent.

⑤ Two random variables that are normally distributed are not necessarily multivariate Gaussian.

EX: Let  $X \sim N(0, 1)$ .

Define  $Y = \begin{cases} X & \text{if } |X| > c \\ -X & \text{if } |X| \leq c \end{cases}$



Clearly  $Y \sim N(0, 1)$ .

Unfortunately,  $(X, Y)$  is NOT bivariate normal.

Qw: Why can't  $X$  and  $Y$  be jointly normal?

$$\Rightarrow \begin{aligned} \text{b/c } \text{CORR}(X, Y) &= +1 && \text{when } |X| > c \\ \text{CORR}(X, Y) &= -1 && \text{when } |X| \leq c \end{aligned}$$

(A multivariate Gaussian has the same covariances, or same correlations, throughout the support set of the density.)

## Asymptotic Distribution of an MLE

Suppose that  $X_1, \dots, X_n$  is a random sample from some density  $f(x; \vec{\theta})$ . Let  $\hat{\vec{\theta}}$  be an MLE vector for  $\vec{\theta}$ . Under a (fairly easy to satisfy) set of conditions mentioned on the homework,

$$\sqrt{n} (\hat{\vec{\theta}} - \vec{\theta}) \sim N(0, I^{-1})$$

↑  
inverse of the  
Fisher transformation;  
see more on the  
homework

In the case when  $\vec{\theta} = \theta$  is one-dimensional,

$$I = \mathbb{E} \left[ \frac{\partial}{\partial \theta} \log f(X; \theta) \cdot \frac{\partial}{\partial \theta} \log \left( \frac{\partial}{\partial x} f(x; \theta) \right) \Big| x = X \right]$$

Cramer-Rao Inequality  
- or - Cramer-Rao Lower Bound

BIG PICTURE POINT: whenever you are estimating a  $k$ -dimensional parameter vector  $\vec{\theta}$  (e.g., in time series class, linear regression class), that vector of MLEs  $\hat{\vec{\theta}}$  is asymptotically multivariate Gaussian.

↳ its variance-covariance matrix will be useful  
↳ door is opened to normality-based CIs and hypothesis tests