

M378K: April 11<sup>th</sup>, 2025.

## Sufficient Statistics.

Consider two election candidates: (A) & (B)

Goal: To "predict" if A wins.

Let the random sample be:  $X_1, X_2, \dots, X_n$

Scenarios for book-keeping:

- $T_1 = (X_1, \dots, X_n)$
- $T_2 = X_1 + X_2 + \dots + X_n$
- $T_3 = \mathbb{1}_{\{\bar{Y} > 0.5\}}$

Formally,

- the conditional dist'n of  $T_1 = (X_1, \dots, X_n)$  given  $T_2 = X_1 + \dots + X_n$  does not depend on  $p$ .
- the conditional dist'n of  $T_1 = (X_1, \dots, X_n)$  given  $T_3 = \mathbb{1}_{\{\bar{Y} > 0.5\}}$  DOES DEPEND ON  $p$

In general, consider two r.v.s  $Y$  and  $T$ ,  
the conditional dist'n of  $Y$  given  $T$   
are these probabilities  $\mathbb{P}[Y=y | T=t]$

In the discrete case, take  $Y = (X_1, \dots, X_n)$ , we write

$$p_{X_1, \dots, X_n, T}(y_1, \dots, y_n, t) = \mathbb{P}[X_1 = y_1, \dots, X_n = y_n, T = t]$$

Def'n. The conditional joint pmf for  $t$  such that  $\mathbb{P}[T=t] > 0$  is

$$p_{X_1, \dots, X_n | T}(y_1, \dots, y_n | t) = \mathbb{P}[X_1 = y_1, \dots, X_n = y_n | T = t]$$

$$= \frac{\mathbb{P}[X_1 = y_1, \dots, X_n = y_n, T = t]}{\mathbb{P}[T = t]}$$

Analogously, in the continuous case:

$$f_{Y_1, \dots, Y_n | T}(y_1, \dots, y_n | t) = \frac{f_{Y_1, \dots, Y_n, T}(y_1, \dots, y_n, t)}{f_T(t)}$$

Def'n. A statistic  $T$  of a random sample  $(Y_1, \dots, Y_n)$  is said to be **SUFFICIENT** for an unknown parameter  $\theta$  if the **conditional dist'n of the sample  $(Y_1, \dots, Y_n)$  given  $T$  does not depend on  $\theta$ .**

Example. **Claim:**  $T_2$  is sufficient for  $p$

$$P_{Y_1, \dots, Y_n | T_2}(y_1, \dots, y_n | t) = \frac{P[Y_1 = y_1, \dots, Y_n = y_n, T_2 = t]}{P[T_2 = t]} \quad \checkmark$$

Cases:

1<sup>st</sup>  $y_1 + \dots + y_n \neq t \rightarrow$  we get 0

2<sup>nd</sup>  $y_1 + \dots + y_n = t \quad \checkmark \quad \{T_2 = t\} \supseteq \{Y_1 = y_1, \dots, Y_n = y_n\}$

$$P[Y_1 = y_1, \dots, Y_n = y_n, T_2 = t] = P[Y_1 = y_1, \dots, Y_n = y_n]$$

$$= P[Y_1 = y_1] \cdot P[Y_2 = y_2] \cdots P[Y_n = y_n]$$

$$= p^{y_1} \cdot (1-p)^{1-y_1} \cdot p^{y_2} (1-p)^{1-y_2} \cdots p^{y_n} (1-p)^{1-y_n}$$

$$= p^{\sum y_i} (1-p)^{n - \sum y_i}$$

$$= p^t (1-p)^{n-t}$$

$$P[T_2 = t] = \binom{n}{t} p^t (1-p)^{n-t}$$

$$P_{Y_1, \dots, Y_n | T_2}(y_1, \dots, y_n | t) = \frac{p^t (1-p)^{n-t}}{\binom{n}{t} p^t (1-p)^{n-t}} = \frac{1}{\binom{n}{t}}$$

Theorem.

## The Fisher-Neyman Factorization Criterion

Let  $Y_1, \dots, Y_n$  be a random sample w/ the likelihood function  $L(\theta; y_1, \dots, y_n)$ .

The statistic  $T$  is sufficient for  $\theta$  if and only if  $L$  can be expressed as

$$L(\theta; y_1, \dots, y_n) = g(\theta, T(y_1, \dots, y_n)) \cdot h(y_1, \dots, y_n)$$