# 1. Comparisons of two proportions

**Recitation problem statement**

You are interested in comparing the proportions of people in their 20's that smoke in France and in the US. After you sample randomly and independently $n$ people in their 20's in both countries, you observe that $N_{US}$ sampled US Americans and $N_F$ sampled French are smokers. Based on such an experiment, how would you test whether there is a significant difference between the proportions of smokers in both countries ?
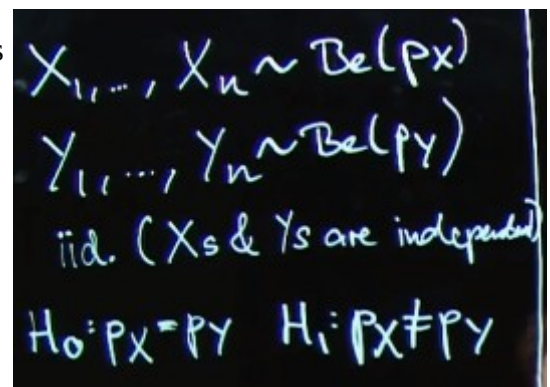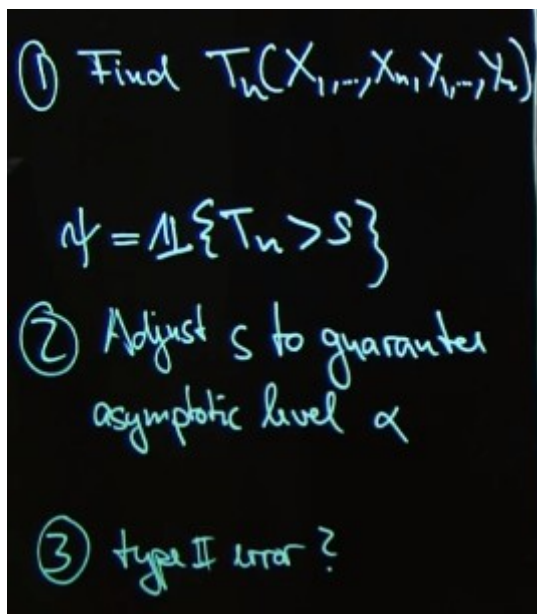
**Note:** In the following videos, we introduce a new term called "pivot". The formal definition of a pivotal quantity (or a pivot) is as follows. Let $X_1, \ldots, X_n$ be random samples and let $T_n$ be a function of $X$ and a parameter vector $\theta$. That is, $T_n$ is a function of $X_1, \ldots, X_n, \theta$. Let $g(T_n)$ be a random variable whose distribution is the same for all $\theta$. Then, $g$ is called a pivotal quantity or a pivot.

For example, let $X$ be a random variable with mean $\mu$ and variance $\sigma^2$. Let $X_1, \ldots, X_n$ be iid samples of $X$. Then,

$$g_n \triangleq \frac{\overline{X_n} - \mu}{\sigma}$$

is a pivot with $\theta = \begin{bmatrix} \mu & \sigma^2 \end{bmatrix}^T$ being the parameter vector. The notion of a parameter vector here is not to be confused with the set of paramaters that we use to define a statistical model.

Have a hypothesis test to check if the number of smokers is the same or different in the US and France

## Comparison of two proportions

(1) $\hat{p}_X = \frac{1}{n}\sum_{i=1}^{n} X_i \xrightarrow[n\to\infty]{\mathbb{P}} p_X$, $\quad \hat{p}_Y = \frac{1}{n}\sum_{i=1}^{n} Y_i \xrightarrow[n\to\infty]{\mathbb{P}} p_Y$

Consider $\hat{p}_X - \hat{p}_Y = g(\hat{p}_X, \hat{p}_Y)$, $\quad g(x,y) = x - y$

CLT: $\sqrt{n}\left(\begin{pmatrix} \hat{p}_X \\ \hat{p}_Y \end{pmatrix} - \begin{pmatrix} p_X \\ p_Y \end{pmatrix}\right) \xrightarrow[n\to\infty]{D} \mathcal{N}(0, \Sigma)$

$= \frac{1}{n}\sum_{i=1}^{n}\begin{pmatrix} X_i \\ Y_i \end{pmatrix}$

$\Sigma = \begin{pmatrix} p_X(1-p_X) & 0 \\ 0 & p_Y(1-p_Y) \end{pmatrix}$

Delta Method:

$\sqrt{n}\left(g(\hat{p}_X, \hat{p}_Y) - g(p_X, p_Y)\right) \to \mathcal{N}\left(0, \nabla g(p_X,p_Y)^T \Sigma \nabla g(p_X,p_Y)\right)$

$\nabla g(x,y) = \begin{pmatrix} 1 \\ -1 \end{pmatrix} \Rightarrow \boxed{\sigma_g^2} = \begin{pmatrix} 1 & -1 \end{pmatrix}\begin{pmatrix} p_X(1-p_X) & 0 \\ 0 & p_Y(1-p_Y) \end{pmatrix}\begin{pmatrix} 1 \\ -1 \end{pmatrix}$

$\underbrace{= p_X(1-p_X) + p_Y(1-p_Y)}$

Want an expression that does not depend on any of the underlying parameters\\

## Comparison of two proportions $\hat{p}_X = \frac{1}{n}\sum_{i=1}^{n} X_i$, $\hat{p}_Y = \frac{1}{n}\sum_{i=1}^{n} Y_i$

$\sqrt{n}\left(\hat{p}_X - \hat{p}_Y - (p_X - p_Y)\right) \xrightarrow[n\to\infty]{D} \mathcal{N}\left(0, p_X(1-p_X) + p_Y(1-p_Y)\right)$

$\Rightarrow \quad \dfrac{\sqrt{n}\cdot\left(\hat{p}_X - \hat{p}_Y - p_X + p_Y\right)}{\sqrt{p_X(1-p_X) + p_Y(1-p_Y)}} \xrightarrow[n\to\infty]{D} \mathcal{N}(0,1)$

Simplify by setting p_x = p_y for H_0 and using a plugin estimator since we do not currently know what p_x hat and p_y hat are.

$$\text{For } H_0, \text{ set } p_X = p_Y = p \in (0,1): \quad p_X(1-p_X) + p_Y(1-p_Y)$$
$$= 2p \cdot (1-p)$$
$$\hat{p} = \frac{1}{2}(\hat{p}_X + \hat{p}_Y) \xrightarrow[n\to\infty]{\mathbb{P}} p$$
$$\text{Slutsky's method: } \sqrt{n}\,\frac{\hat{p}_X - \hat{p}_Y}{\sqrt{2\hat{p}(1-\hat{p})}} \xrightarrow[n\to\infty]{D} \mathcal{N}(0,1)$$

This is the test statistic, Tn, that we want. The absolute value of it as it's a 2 sided test.

$$\textcircled{2} \quad T_n = \left| \sqrt{n}\,\frac{\hat{p}_X - \hat{p}_Y}{\sqrt{2\hat{p}(1-\hat{p})}} \right| \xrightarrow[n\to\infty]{D} \mathcal{N}(0,1) \text{ under } H_0\,(p_X = p_Y)$$

Cutting off at Tn > S to guarantee an asymptotic level of alpha

$$\mathbb{P}(T_n > s) = \mathbb{P}\left( \left| \sqrt{n}\,\frac{\hat{p}_X - \hat{p}_Y}{\sqrt{2\hat{p}(1-\hat{p})}} \right| > s \right) \xrightarrow[n\to\infty]{} \underbrace{\mathbb{P}(|Z| > s)}_{\sim \mathcal{N}(0,1)}$$
$$= 2(1 - \Phi(s)) \qquad\qquad \overset{!}{=} \alpha$$
$$\implies s = q_{\alpha/2}, \ 1-\alpha/2 \text{ quantile of } \mathcal{N}(0,1)$$

$$\textcircled{3} \quad \hat{p}_X \xrightarrow[n\to\infty]{\mathbb{P}} p_X, \ \hat{p}_Y \xrightarrow[n\to\infty]{\mathbb{P}} p_Y, \ \hat{p} = \frac{1}{2}(\hat{p}_X + \hat{p}_Y) \xrightarrow[n\to\infty]{\mathbb{P}} \frac{1}{2}(p_X + p_Y) =: \tilde{p}$$
$$T_n = \left| \sqrt{n}\,\frac{\overbrace{\hat{p}_X - \hat{p}_Y}^{\to p_X - p_Y \neq 0 \text{ under } H_1}}{\underbrace{\sqrt{2\hat{p}(1-\hat{p})}}_{\to \sqrt{2\tilde{p}(1-\tilde{p})}}} \right| \xrightarrow[n\to\infty]{\mathbb{P}} +\infty \implies \text{type II error} \xrightarrow[n\to\infty]{\mathbb{P}} 0$$

If the null hypothesis does not hold then the statistic test exceeds the limits with overwhelming probability (approaches infinity)

Remarks:

Different sample sizes

$$④ \text{ Different sample sizes}: X_1, \ldots, X_{n_1}, Y_1, \ldots, Y_{n_2}$$

$$\text{CLT}: Z_1, \ldots, Z_n \overset{\in \mathbb{R}}{\text{ iid}}, \quad \mathbb{E}[Z_i] = \mu, \quad Var(Z_i) = \sigma^2$$

$$\sqrt{n} \; \frac{\bar{Z}_n - \mu}{\sqrt{\sigma^2}} \xrightarrow[n \to \infty]{D} N(0,1); \quad \bar{Z}_n = \frac{1}{n}\sum_{i=1}^{n} Z_i$$

$$\mathbb{E}[\bar{Z}_n] = \frac{1}{n}\sum_{i=1}^{n}\mathbb{E}[Z_i]$$

$$= \frac{n}{n}\mu = \mu$$

$$Var(\bar{Z}_n) = \frac{1}{n^2}\cdot\sum_{i=1}^{n} Var(Z_i)$$

$$= \frac{n}{n^2}\cdot Var(Z_i) = \frac{1}{n}\sigma^2$$

$$\sqrt{n} \; \underbrace{\frac{\bar{Z}_n - \mu}{\sqrt{\sigma^2}}}_{\mathbb{E}[-n-]=0, \; Var(-n-)=1} \xrightarrow[n \to \infty]{D} N(0,1)$$