### Tests for independence
### + Permutation tests

Ex.

| X \ Y | Online | offline |
|-------|--------|---------|
| A | 25 ? | 15 |
| B | 10 | 10 |
| C | 15 | 10 |

2D cont. table

← obs-d freq.

Q. Are $X$ and $Y$ dependent?

$H_0$: $X$ and $Y$ are indep.
$H_A$: $X$ and $Y$ are dep.

**Any ideas?**

**Classic Pearson's Independence $\chi^2$ test:**

$$P.S. = \sum_{i,j} \frac{(N_{ij} - \hat{\mu}_{ij})^2}{\hat{\mu}_{ij}}$$

$$\sum_i \frac{(N_i - E(N_i))^2}{E(N_i)}$$

$\hat{\mu}_{ij}$ - is the expected value
     of $N_{ij}$ estimated under $H_0$ (independency)

idea: estimate $\hat{p}_A, \hat{p}_B, \hat{p}_C$

$\hat{p}_{online}, \hat{p}_{offline}$
under $H_0$

$$\hat{p}_A = \frac{N_{A,on} + N_{A,off}}{N}$$

$$\hat{p}_{on} = \frac{N_{A,on} + N_{B,on} + N_{C,on}}{N}$$

$$\hat{\mu}_{A,on} = N \cdot \hat{p}_A \cdot \hat{p}_{on}$$

$$\hat{p}_{i\cdot} = \frac{\sum_j N_{ij}}{N} \qquad \hat{p}_{\cdot j} = \frac{\sum_i N_{ij}}{N}$$

$$\hat{\mu}_{ij} = N \times \hat{p}_{i\cdot} \times \hat{p}_{\cdot j} = \frac{\left(\sum_i N_{ij}\right) \cdot \left(\sum_j N_{ij}\right)}{N}$$

$$P.S. = \sum_{ij} \frac{(N_{ij} - \hat{\mu}_{ij})^2}{\hat{\mu}_{ij}} \quad \xrightarrow[dist]{N \to \infty, H_0} \quad \chi^2_{df}$$
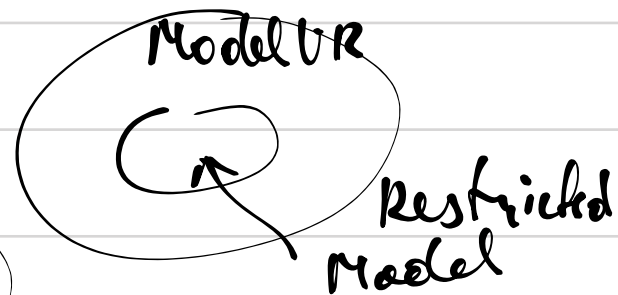
$$df = (C_x - 1) \cdot (C_y - 1) \qquad where \qquad \begin{array}{l} C_x - number \\ \quad of\ possible \\ \quad values\ of\ X \\ C_y \ -//- \ of\ Y \end{array}$$

---

## Likelihood Ratio test.

$$\boxed{LR = 2 \cdot \left( \ln L_{UR} - \ln L_R \right)}$$

nested models:
Model UR
Restricted
Model

$$\xrightarrow[\text{dist}]{H_0,\ N \to \infty} \quad \boxed{p_{UR} - p_R}$$

$p_R, p_{UR}$ - number of par-s.

$H_0$: Restricted model is TRUE
$H_A$: R-model is FALSE, but UR-model is TRUE

---

$C_y$ number of cat

| X\Y | $\alpha$ | $\beta$ | $\gamma$ .... |
|-----|---|---|---|
| A | · | · | · |
| B | · | · | · |
| C | · | · | · |
| ... | | | |

$C_x$ - number of cat.

R-model: X and Y are indep.
UR-model: maybe X and Y are depend., any probab-s are OK.

one constraint: $\sum_{ij} p_{ij} = 1$

$$p_{UR} \overset{?}{=} \underbrace{C_x \cdot C_y - 1}_{\substack{\text{number of cells in the table}}} \quad \text{free parameters.}$$

$$p_{UR} - p_R = C_x \cdot C_y - 1 - C_x - C_y + 2 = \boxed{= (C_x - 1) \cdot (C_y - 1)}$$

$$p_R \ ? = \underbrace{C_x - 1}_{\text{prob-s for X}} + \underbrace{C_y - 1}_{\text{prob-s for Y}} = C_x + C_y - 2$$

## UR - model :

| x \ Y | Online | offline |
|-------|--------|---------|
| A | (25) ? | 15 |
| B | 10 | 10 |
| C | 15 | 10 |

$$N = 85$$

| x \ Y | On | Off |
|-------|-----|-----|
| A | $\hat{p}_{11}$ | $\hat{p}_{12}$ |
| B | $\hat{p}_{21}$ | $\hat{p}_{22}$ |
| C | $\hat{p}_{31}$ | $\hat{p}_{32}$ |

( it is not a free parameter!

$$p_{UR} = 2 \cdot 3 - 1 = 5$$

$$\hat{p}_{11}^{UR} = \frac{25}{85} \qquad \hat{p}_{12}^{UR} = \frac{15}{85}$$

$$\hat{p}_{21}^{UR} = \frac{10}{85} \qquad \vdots$$

$$LR = 2 \cdot \left( \ln L_{UR} - \ln L_R \right) = [\text{use past lecture}]$$

$$= 2 \cdot \sum_{i,j} N_{ij} \cdot \left( \ln \hat{p}_{ij}^{UR} - \ln \hat{p}_{ij}^{R} \right) \xrightarrow[\text{dist}]{n \to \infty \ H_0} \chi^2_{df}$$

$$df = p_{UR} - p_R = (C_x - 1) \cdot (C_y - 1)$$

R - model : [$H_0$ of independency]

| x \ Y | Online | offline |
|-------|--------|---------|
| A | (25) ? | 15 |
| B | 10 | 10 |
| C | 15 | 10 |

$$\hat{p}_{1\cdot} = \frac{25 + 15}{85} \Big\}$$

$$\hat{p}_{2\cdot} = \frac{20}{85} \qquad \text{2 free params}$$

$$\hat{p}_{3\cdot} = \frac{25}{85} \qquad \text{It is not a free param!}$$

$$\hat{p}_{\cdot 1} = \frac{25 + 10 + 15}{85} \Big\} \text{1 free param}$$

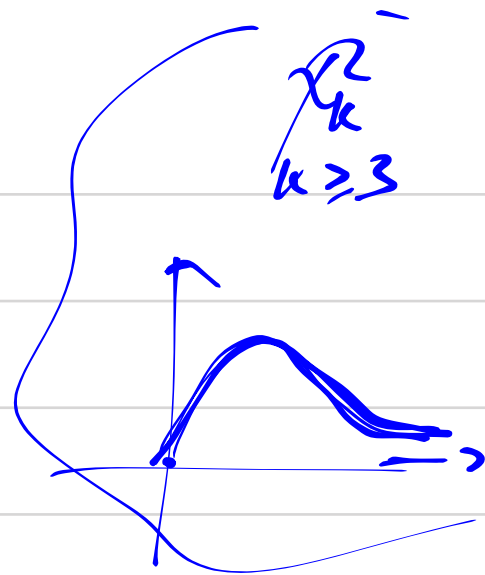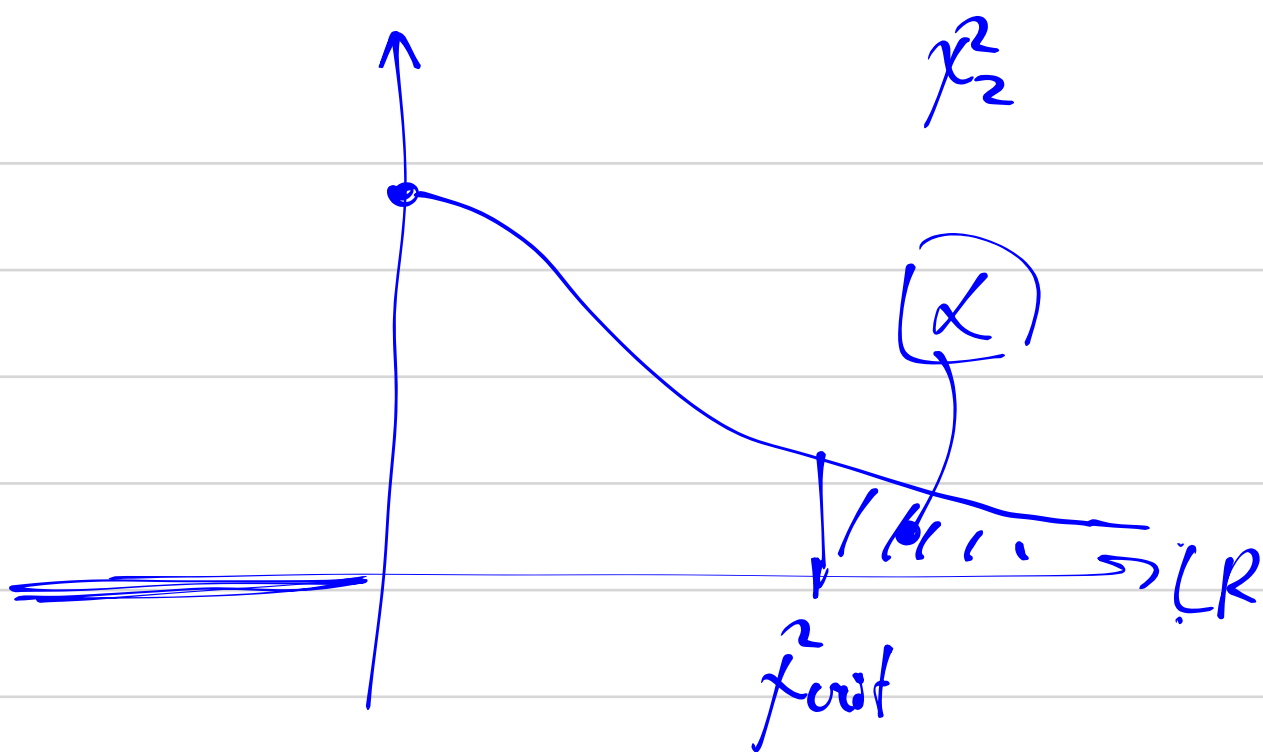$$\hat{p}_{\cdot 2} = \frac{15 + 10 + 10}{85}$$

$$p_R = 1 + 2 = 3 \text{ param.}$$

$H_0$ : X and Y are indep.

$$\hat{p}_{ij}^{R} = \hat{p}_{i\cdot} \times \hat{p}_{\cdot j}$$

$$\hat{p}_{11}^{R} = \frac{25 + 15}{85} \cdot \frac{25 + 10 + 15}{85}$$

in my ex : $\chi^2_{(3-1) \cdot (2-1)} = \chi^2_2$

$\chi^2_2$

$\chi^2_k$
$k \geq 3$



if $LR > \chi^2_{crit}$ then reject $H_0$.

## Permutation test. (Bootstrap.)

**idea:**
write a long
table, row =
= one observ.

| X | Y |
|---|---|
| A | On |
| B | On |
| C | Off |
| A | Off |
| ⋮ | ⋮ |
| C | On |

Calculate **any [!]**
reasonable measure
of similarity (correla-n)

**example:** [of a reas. measure for this case]

| You | $X_A$ | $X_B$ | $X_C$ |
|-----|-------|-------|-------|
| 1 | 1 | 0 | 0 |
| 1 | 0 | 1 | 0 |
| 0 | . | | |
| 0 | | | |
| ⋮ | | | |
| 1 | | | |

Step1 regress: $\hat{y}_i^{on} = \hat{\beta}_1 \cdot X_i^A + \hat{\beta}_2 \cdot X_i^B + \hat{\beta}_3 X_i^C$

Step 2:  calculate $R^2_{orig}$

① Randomly permute values of Y ⟹ obtain $R^2_{perm1}$
② — '' — ⟹ obtain $R^2_{perm2}$
   ⋮
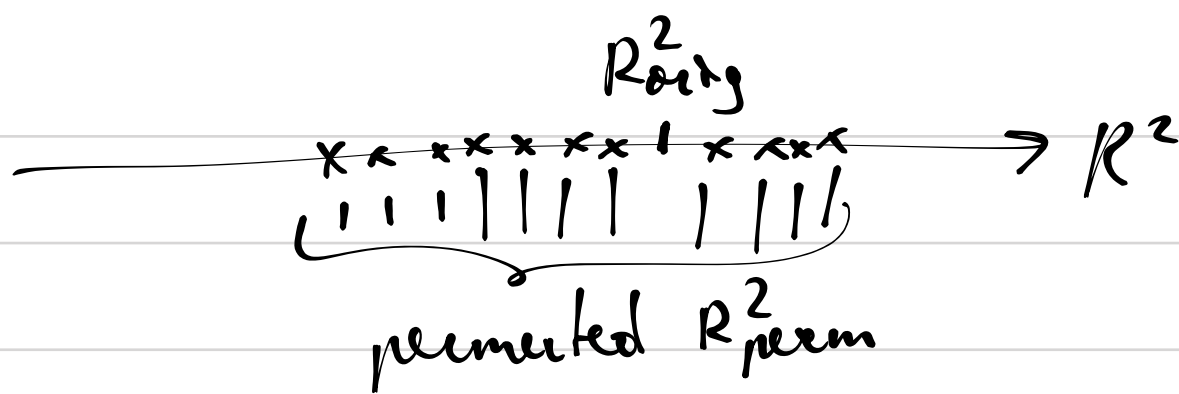10000 . — '' — ⟹ obtain $R^2_{perm\,10000}$

$R^2_{p1} .. \quad R^2_{p10000} \quad , R^2_{orig}$  → $R^2$

$$R^2_{orig}$$

x x xxx xx | x xx x $\longrightarrow R^2$
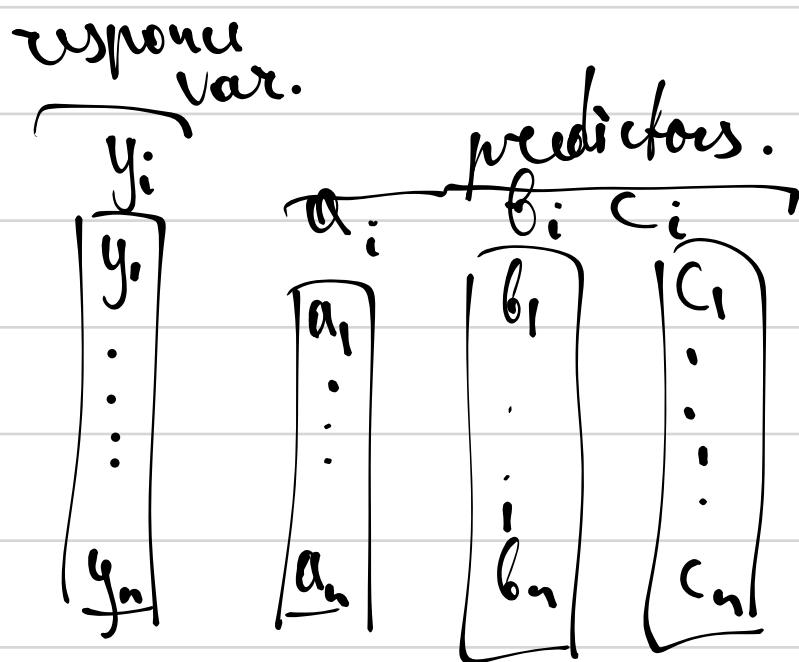
permuted $R^2_{perm}$

P-value = fraction of $R^2_{perm}$ higher than $R^2_{orig}$.

If p-value $< \alpha = 0.05$ then reject $H_0$

If p-value $\geq \alpha = 0.05$ then do not reject $H_0$.

<span style="color:red">Permutation tests in machine learning.</span>

response var.

$$\overbrace{\begin{matrix} y_i \\ y_1 \\ \vdots \\ \vdots \\ y_n \end{matrix}}$$ $$\overbrace{\begin{matrix} a_i & b_i & c_i \\ a_1 & b_1 & c_1 \\ \vdots & \vdots & \vdots \\ \vdots & \vdots & \vdots \\ a_n & b_n & c_n \end{matrix}}^{\text{predictors.}}$$

$\rightarrow$ Classification problem
  $y_i \in \{A, B, C\}$

$\rightarrow$ Gradient boosting

Q. Which predictors are really important?

Step 1. Split into train and test

Train $\longrightarrow$

Test

Step 2. <span style="color:red">(only one model est-n)</span> estimate the param.s of the algorithm on the training set.

$A_{orig}$ | · Step 3. Calculate accuracy on test sample (or any other reasonable quality measure)

(Step 4) Run permut-n test for variable _a_.
  $\rightarrow$ permute values of var-le a.
      obtain $A_{perm 1}$

[Step 5] calculate : $\overrightarrow{A_{orig} - A_{perm}}$  " obtain $A_{perm 2}$