

The background of the slide features a stylized globe on the left side, partially obscured by a dense pattern of binary code (0s and 1s) that recedes into the distance, creating a sense of depth and digital connectivity. The colors are primarily light blue and white.

# Foundations of Data Science:

## Hypothesis testing - eT      sting for goodness-of-fit

## Multiple categories

American Civil Liberties Union report in L6 jury selection in Alameda County, Ca (2010)

	Caucasian	Black/AA	Hispanic	Asian/PI	Other	Total
Population %	54	18	12	15	1	100
Observed panel numbers	780	117	114	384	58	1453
Expected panel numbers	784.62	261.54	174.36	217.95	14.53	1453.00
$\frac{(\text{Observed}-\text{Expected})^2}{\text{Expected}}$	0.03	79.88	20.90	126.51	130.05	357.36

$H_0$ : The panels were chosen by random selection from the population

$H_a$ : The panels were chosen by some other, unspecified method.

$k$  - groups

$p_i$  - population proportion in the  $i$ th group

$n$  - total size of population  $n = \sum_i n_i$

$n_i$  - number in  $i$ th group

$np_i$  - expected number in each group.

$$\chi^2 = \sum_{i=1}^k \frac{(n_i - np_i)^2}{np_i}$$

e.g.  $np_1 = 100$   $n_1 = 95$   $\overbrace{5}^{5\%}$   
 $np_2 = 10$   $n_2 = 5$   $\underbrace{5}_{50\%}$

"chi - squared"

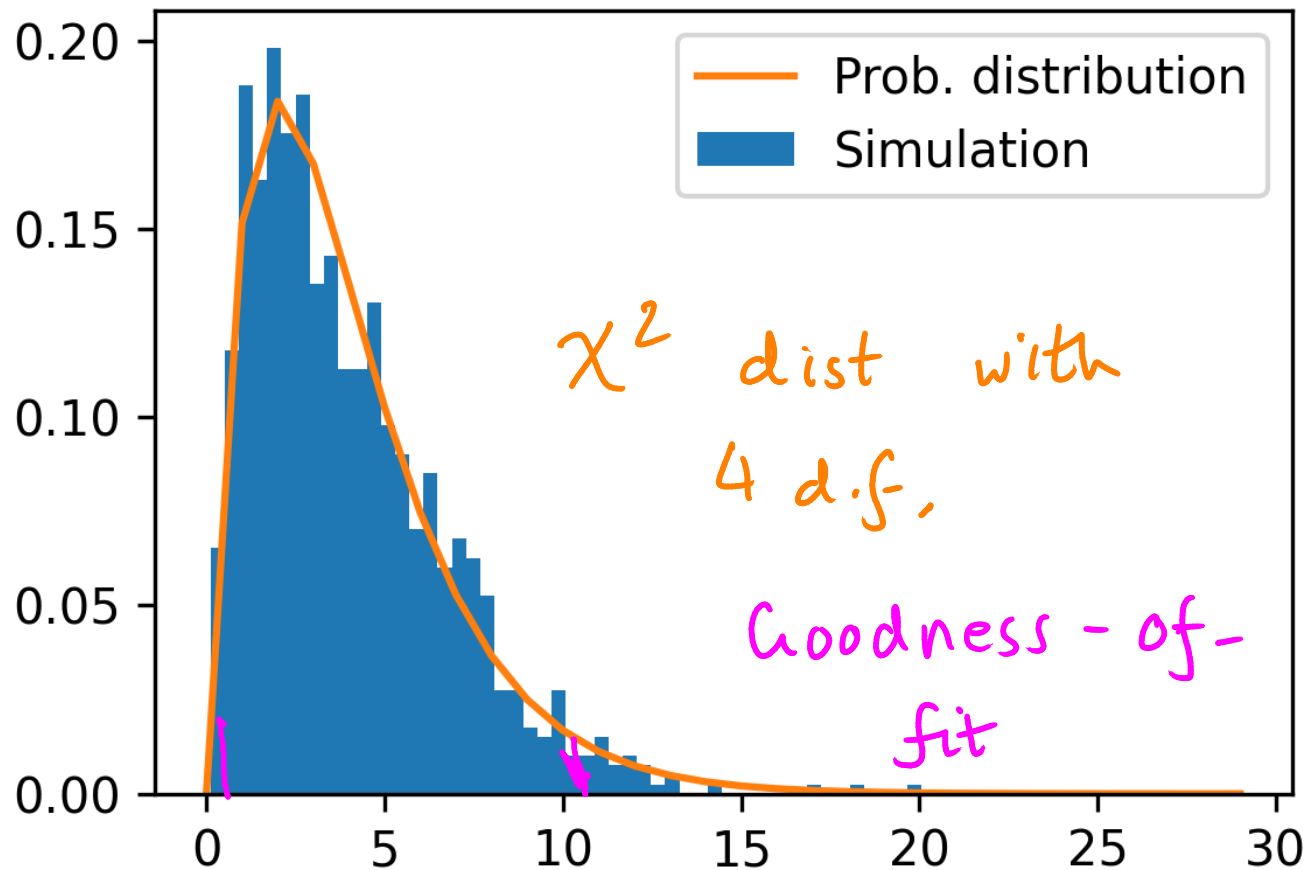
In e.g.  $\chi^2 = 357.36$

# Statistical simulations

$n = 1543$

$p_i$

$p \approx 0$



→  
357

$\chi^2$  stat is distributed according to  
a  $\chi^2$  distribution with  $k-1$  d.f.

## 2-way contingency tables

	Female	Male	Total
Depressed	30	12	42
Not depressed	2048	1663	3711
Total	2078	1675	3753

	Population 1	Population 2	Total
Category 1	$n_{11}$	$n_{12}$	$n_{1\bullet}$
Category 2	$n_{21}$	$n_{22}$	$n_{2\bullet}$
Total	$n_{\bullet 1}$	$n_{\bullet 2}$	$n_{\bullet\bullet}$

$H_0$  = Being severely depressed is independent of being female or male

$H_a$  = Some other, unspecified hypotheses.

$$H_0 = P(X=x | Y=y) = P(X=x) P(Y=y)$$

$$p_{i\bullet} = \frac{n_{i\bullet}}{n_{\bullet\bullet}} \quad p_{\bullet j} = \frac{n_{\bullet j}}{n_{\bullet\bullet}}$$

$$\Rightarrow \hat{E}_{ij} = n_{\bullet\bullet} p_{i\bullet} p_{\bullet j} = \frac{n_{i\bullet} n_{\bullet j}}{n_{\bullet\bullet}}$$

# Multiway contingency tables

$J=2$

$I=2$

	Female	Male	Total
Depressed	30	12	42
Not depressed	2048	1663	3711
Total	2078	1675	3753

	Population 1	Population 2	Total
Category 1	$n_{11}$	$n_{12}$	$n_{1\bullet}$
Category 2	$n_{21}$	$n_{22}$	$n_{2\bullet}$
Total	$n_{\bullet 1}$	$n_{\bullet 2}$	$n_{\bullet\bullet}$

	Female	Male
Depressed	23.25	18.75
Not depressed	2054.75	1656.25

	Population 1	Population 2
Category 1	$\hat{e}_{11}$	$\hat{e}_{12}$
Category 2	$\hat{e}_{21}$	$\hat{e}_{22}$

$$\chi^2 = \frac{(\text{Observed} - \text{Expected})^2}{\text{Expected}} = \sum_i \sum_j \frac{(n_{ij} - \hat{e}_{ij})^2}{\hat{e}_{ij}}$$

$$= 4.433$$

$$\# \text{ d.f.} = (I - 1)(J - 1) + 1 = 1$$

$$p = 0.035$$