

M358K : November 16<sup>th</sup>, 2020.

## $\chi^2$ test of independence

$\left. \begin{matrix} X \\ Y \end{matrix} \right\}$  categorical random variables

$X \backslash Y$	$y_1$	$y_2$	...	$y_j$	...	$y_c$
$x_1$						
$x_2$						
$\vdots$						
$x_i$				$p_{ij}$		$p_{X}(x_i)$
$\vdots$						
$x_r$						
				$p_Y(y_j)$		$1$

dimension:  $r \times c$

Joint probability  
mass function

$$\mathbb{P}[X=x_i, Y=y_j] =: p_{ij} \text{ for all } i, j$$

If  $X$  and  $Y$  are INDEPENDENT:

$$p_{ij} = \underbrace{\mathbb{P}[X=x_i]}_{p_X(x_i)} \cdot \underbrace{\mathbb{P}[Y=y_j]}_{p_Y(y_j)} \text{ for all } i, j$$

marginal dist's

Two-way tables (an empirical version of the joint pmf table)

	" $y_j$ "	
" $x_i$ "	$n_{i,j}$	$(r_i)$ ← the total count for $i$
	$c_j$ ↑ total count in column $j$	$n$ ... total sample size

$n_{i,j}$ ... # of observed cases w/ the combination " $x_i$ ", " $y_j$ "

Empirically:

Q: What is the probability that you land in row  $i$ ?

$$\frac{r_i}{n}$$

Q: What is the probability that you land in column  $j$ ?

$$\frac{c_j}{n}$$

If the row and column effects (random variables) are independent, then the probability of landing in cell  $(i,j)$  is  $\frac{r_i}{n} \cdot \frac{c_j}{n} = \frac{r_i \cdot c_j}{n^2}$ .

⇒ The expected count in cell  $(i,j)$  if the two effects are independent is

$$n \cdot \frac{r_i \cdot c_j}{n^2} = \frac{r_i \cdot c_j}{n}$$