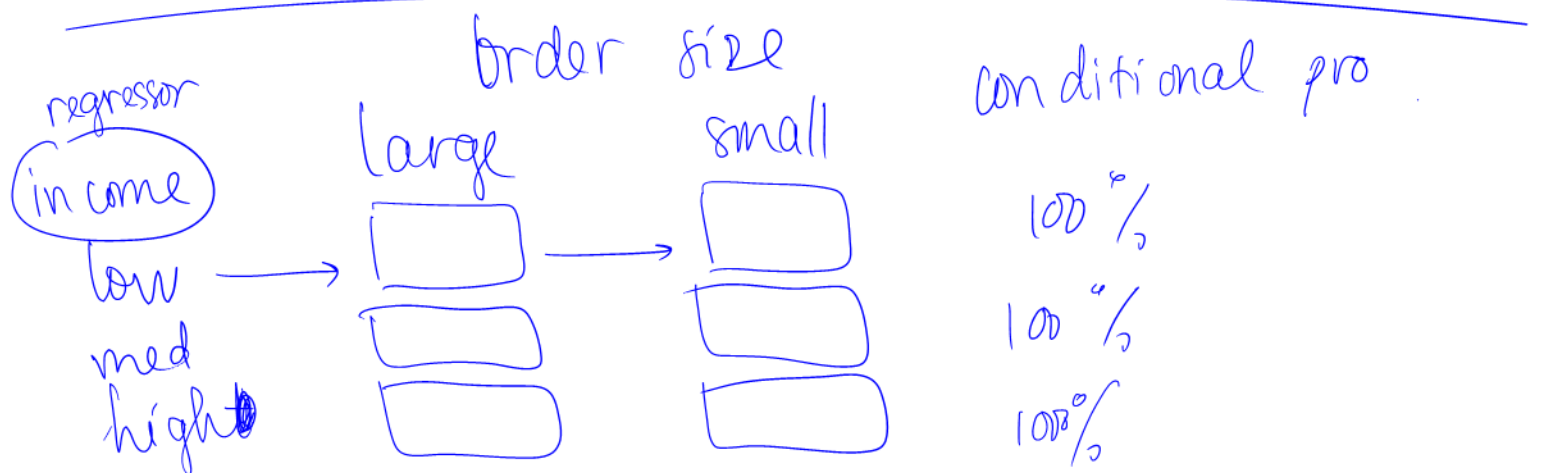


3 cate \Rightarrow $\left\{ \begin{array}{l} \text{plot} \\ \text{contingency table} \\ \text{not OR} \end{array} \right.$ (OR is only for 2x2 table).

$$\begin{aligned} \Pr(\underline{\text{large}} \mid M) &= \approx 3.7\% \\ \Pr(\text{large} \mid F) &= \sim 2.8\% \end{aligned}$$

conditional prop / prob could help.

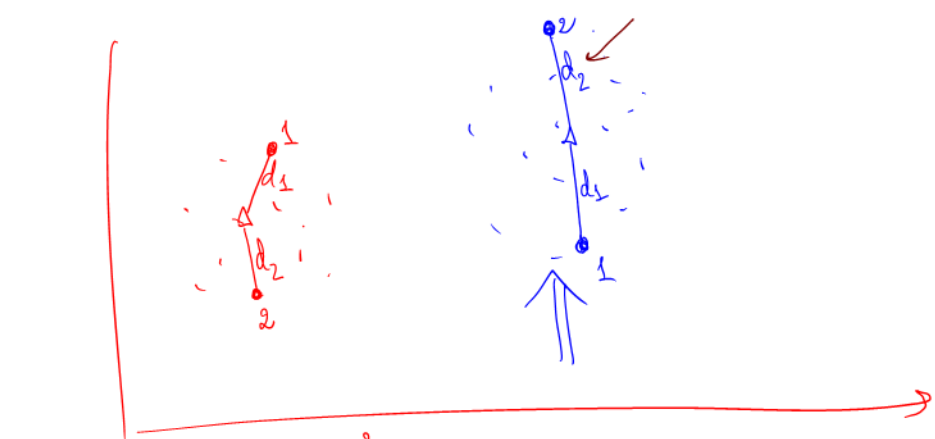


$$\begin{aligned} \text{Prop}(\text{large} \mid \text{low}) &= 5\% \\ \text{Prop}(\text{large} \mid \text{med}) &= 10\% \\ \text{Prop}(\text{large} \mid \text{high}) &= 20\% \end{aligned}$$

$$\frac{X}{3} \times Y$$

$\rightarrow \left\{ \begin{array}{l} KNN \\ DT \\ NB \\ LO.R \end{array} \right\}$
 $\rightarrow \left\{ \begin{array}{l} L \\ M \end{array} \right\}$

a x b



SS for blues:

d_1^2 : point 1 $\rightarrow \Delta$

d_2^2 : point 2 $\rightarrow \Delta$

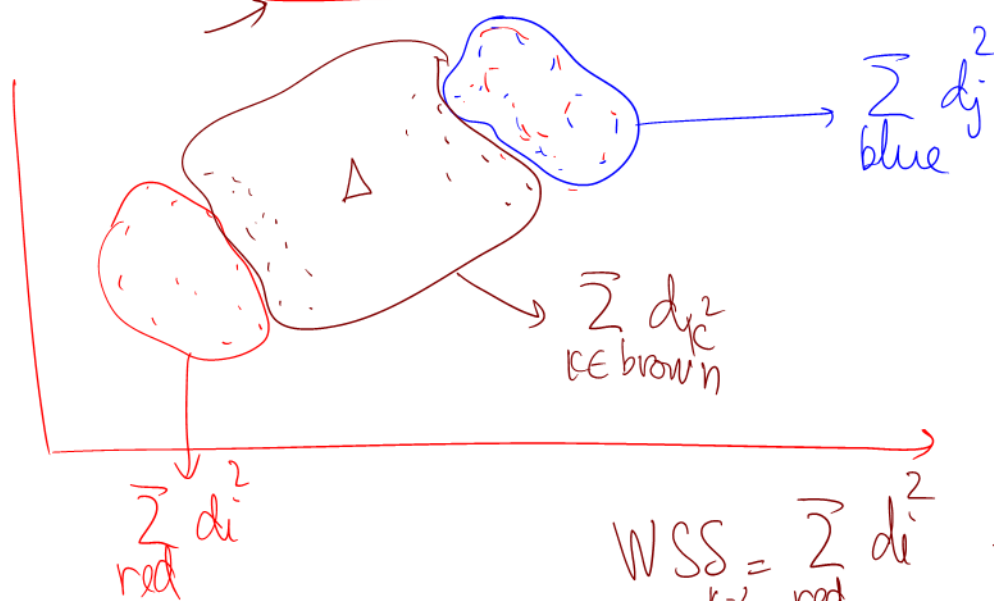
SS for reds: d_1^2
 d_2^2

all blue points: $\sum_{i \in \text{blue}} d_i^2$

all red points $\Rightarrow \sum_{j \in \text{reds}} d_j^2$

$WSS_{k=2}$

$$= \sum_{\text{blue}} d_i^2 + \sum_{\text{red}} d_j^2$$



$$WSS_{k=3} = \sum_{\text{red}} d_i^2 + \sum_{\text{blue}} d_j^2 + \sum_{\text{brown}} d_k^2$$

compare: $WSS_{k=2}$ vs $WSS_{k=3}$

value of k that gives smaller WSS is better. However we would prefer smaller k .

hdb : try with $k = 1 \rightarrow WSS_1$
 $k = 2 \rightarrow WSS_2$
 \vdots
 quant $\dots k = 10 \rightarrow WSS_{10}$

$\boxed{Y} \sim x_1 + x_2$ linear
 \downarrow
 response

$\log \frac{p}{1-p} = \dots$ logistic $\rightarrow Y = \text{response} \begin{cases} 0 \\ 1 \end{cases}$
 $p = \text{Prob}(Y = 1)$

to find best k for KNN:

fit KNN with $k = 1 \rightarrow$ check g.o.f. \rightarrow test set
 for train data (80%)

fit 2-NN \rightarrow check g.o.f. \rightarrow for test set
 for train data (80%)

N fold for 80% (train data)

\downarrow
 N folds \Rightarrow find best k
 \rightarrow check for test set.