# IDS 702

**Cross validation**

# Cross validation

- The train/test method of model validation is often referred to as cross validation

- Splitting the data randomly into two sets can have a big impact on the train and test MSE, particularly in small samples

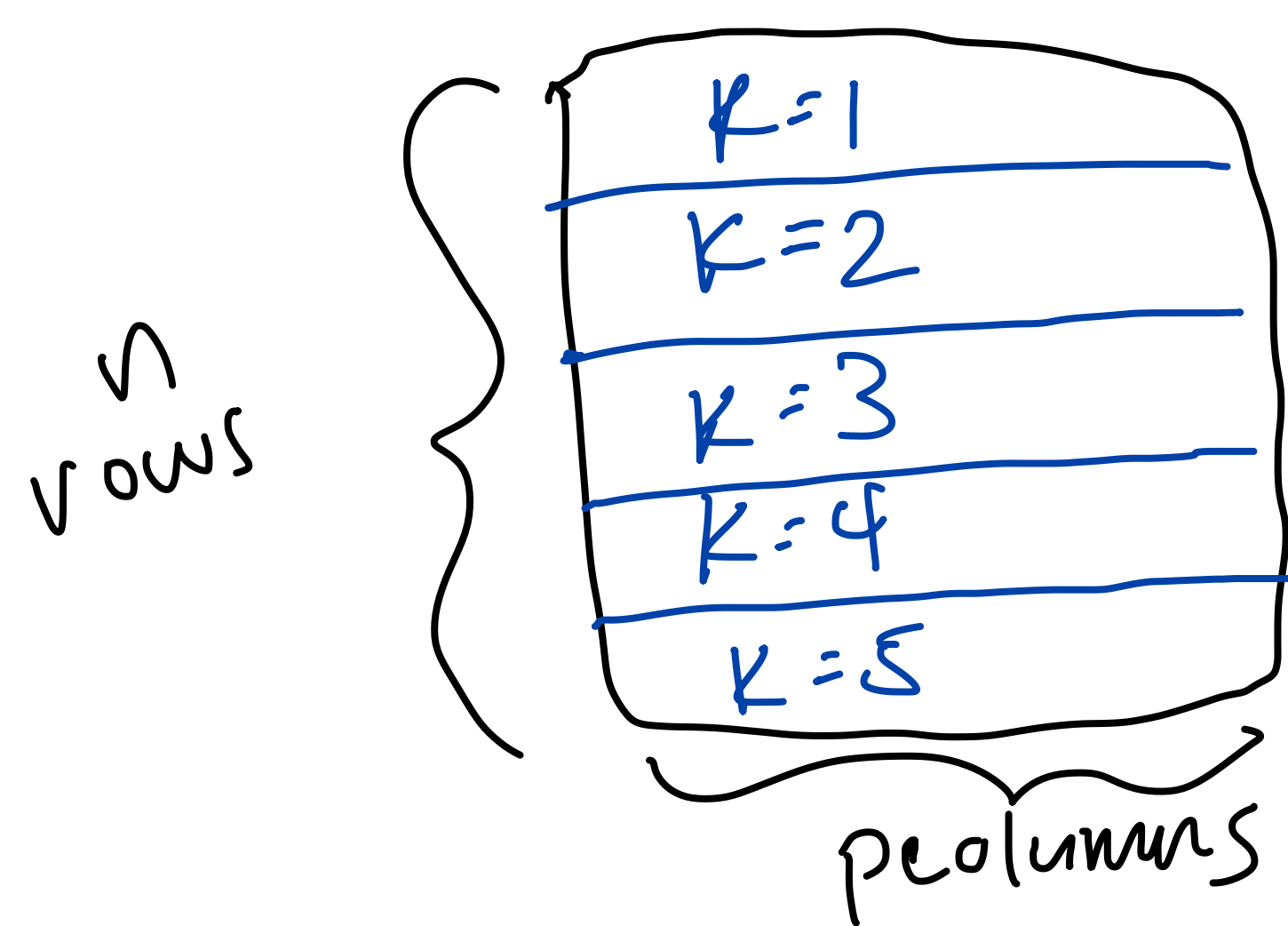- **K-fold cross validation** is a type of cross validation that aims to address this sensitivity

# K-fold cross validation

- Split the data into $K$ mutually exclusive groups (folds)

- For the $k$th fold, with $k = 1,...,K$, fit the model on all the remaining data excluding the $k$th fold (that is, all the other folds combined) and use the $k$th fold as the test set

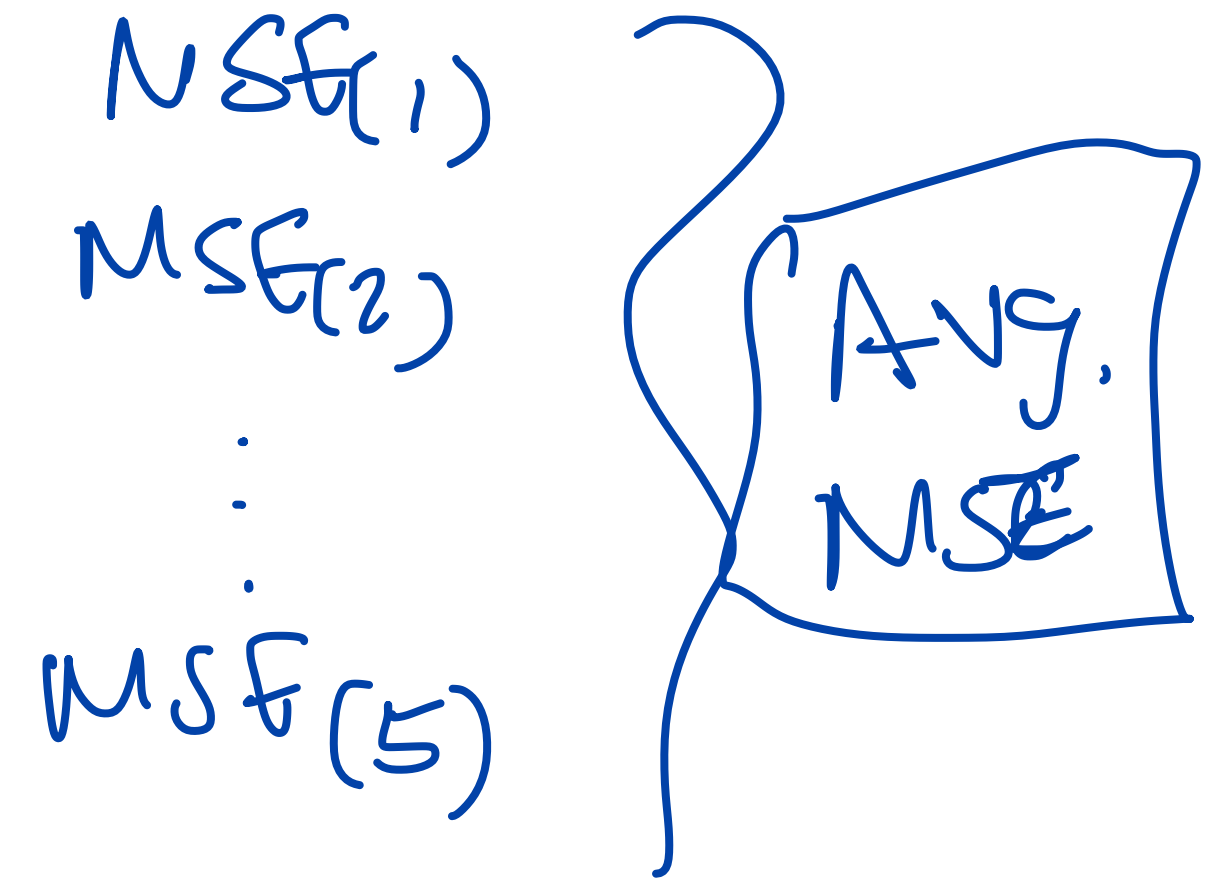- Repeat for each $k$, obtain MSE for each $k$, and summarize using the average over $K$

- $$AvgMSE = \frac{1}{K} \sum_{i=1}^{K} MSE_{test}^{(k)}$$

MSE is an appropriate metric for an average (p-values are valid)

# K-fold cross validation

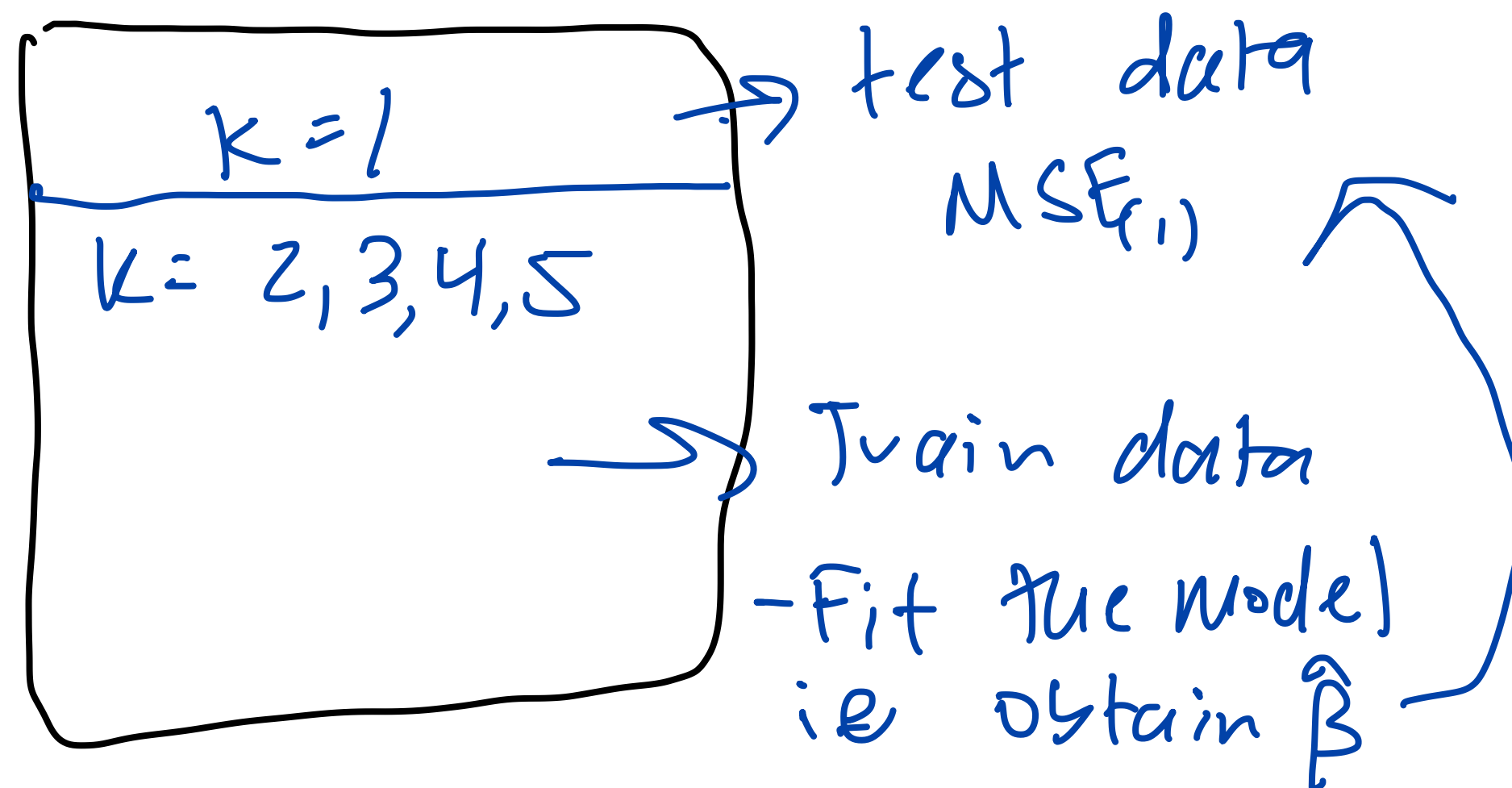$K = 5$

$n$ rows $\left\{ \begin{array}{c} \\ \\ \\ \\ \end{array} \right.$

| K=1 |
|---|
| K=2 |
| K=3 |
| K=4 |
| K=5 |

$\underbrace{\quad\quad\quad}_{p \text{ columns}}$

$MSE_{(1)}$
$MSE_{(2)}$
$\vdots$
$MSE_{(5)}$

Avg. MSE

**$k=1$:**

| K=1 | → test data $MSE_{(1)}$ |
|---|---|
| K= 2,3,4,5 | → Train data |

- Fit the model ie obtain $\hat{\beta}$

**$k=2$**

$MSE_{(2)}$

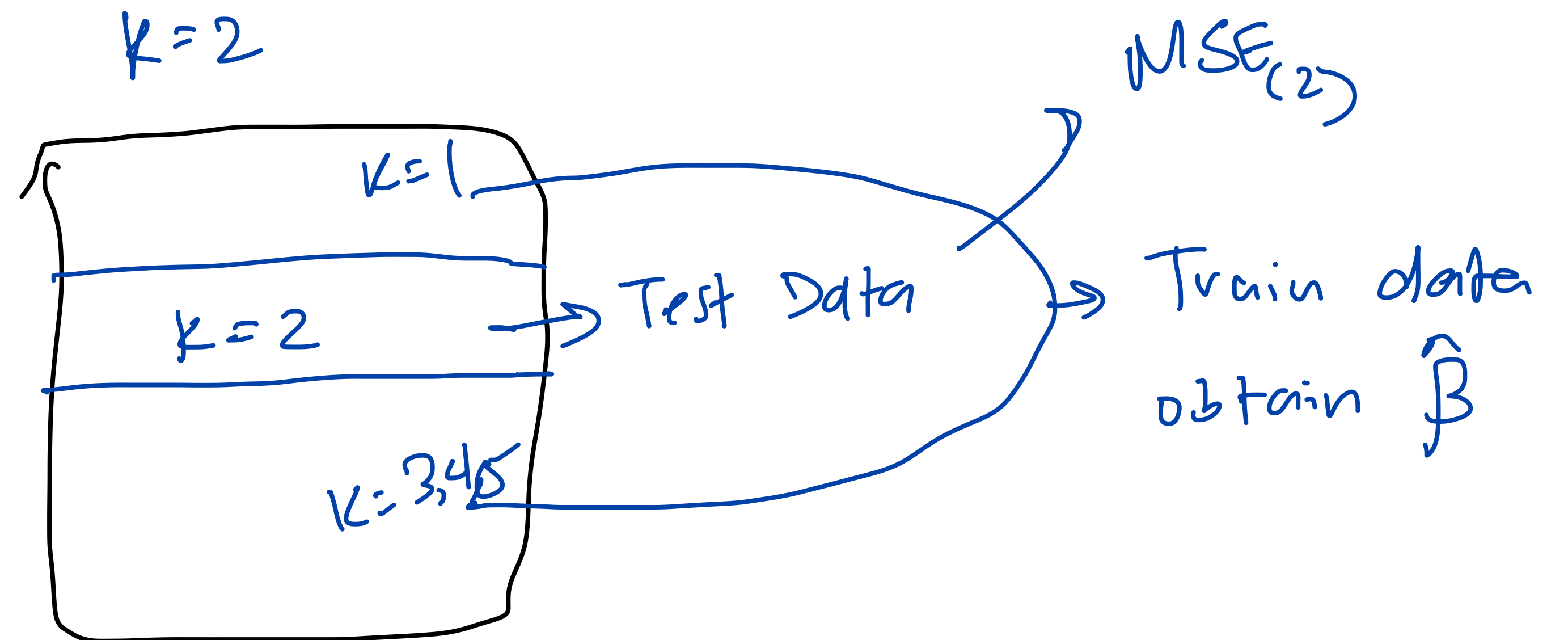| K=1 | → Train data |
|---|---|
| K=2 | → Test Data |
| K=3,4,5 | |

→ Train data obtain $\hat{\beta}$

# What should K be?

- Leave-one-out cross validation: $K = N$ (computationally intensive)

- $k = 5, k = 10$ are common choices

Consider:

- Sample size
- How many Models to compare (computation resources)