

Model Evaluation & Validation

* Measures of Central Tendency

* Mode - Most frequently occurring

Median - Middle value when the data is ordered

- Half way between the two middle values
when there is an even number of points

Mean - Average = $\sum_{i=1}^n x_i / n$

* Mode

- distribution can have more than one mode (can be a range or set of values)
- not all samples of a distrⁿ contribute to it
- not affected by outliers

Mean

- one value
- all samples of a distrⁿ contribute to it
- affected by outliers

Median

- one value
- all samples of a distrⁿ contribute to it
- not affected by outliers

* Order of measures



mode < median < mean



mode = median = mean



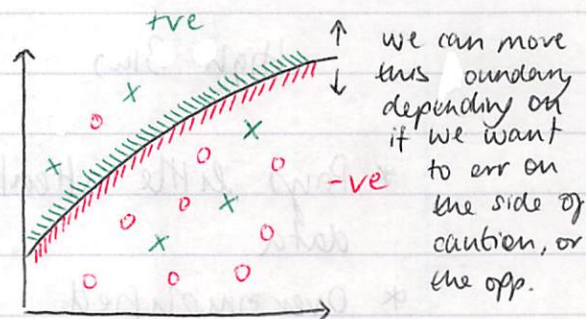
mode > median > mean

Evaluation Metrics

* Confusion Matrix Example:

	actual class	
	+ve	-ve
Predicted class	+ve	9 True +ve
	-ve	3 False -ve
	+ve	1 False +ve
	-ve	8 True -ve

False Alarm



* Precision & Recall

$$\text{Precision} = P(\text{true} | \text{+ve}) = \frac{\# \text{ True +ve}}{\# \text{ True +ve} + \# \text{ False +ve}}$$

$$\text{Recall} = P(\text{+ve} | \text{true}) = \frac{\# \text{ True +ve}}{\# \text{ True +ve} + \# \text{ False -ve}}$$

Terminology:

$\left\{ \begin{array}{l} \text{True} \\ \text{False} \end{array} \right\}$ refers to the correctness of the prediction
 $\left\{ \begin{array}{l} \text{+ve} \\ \text{-ve} \end{array} \right\}$ refers to the prediction result

* Example:

	PREDICTED					precision = 1	recall = 5/8		
ACTUAL	Ariel Sharon	[13	4	1	1	0	0	1]
	Colin Powell	[0	55	0	8	0	0	0]
	Donald Rumsfeld	[0	1	25	8	0	0	2]
	George W Bush	[0	3	0	123	0	0	1]
	Gerhard Schroeder	[0	1	0	7	14	0	4]
	Hugo Chavez	[0	3	0	2	1	10	0]
	Tony Blair	[0	0	1	7	0	0	26]

Causes of Error

* Bias Variance Dilemma Solution

High Bias

- * Pays little attention to data
- * Oversimplified
- * High error on training set (low r^2 , high SSE)

- * Few features used

High Variance

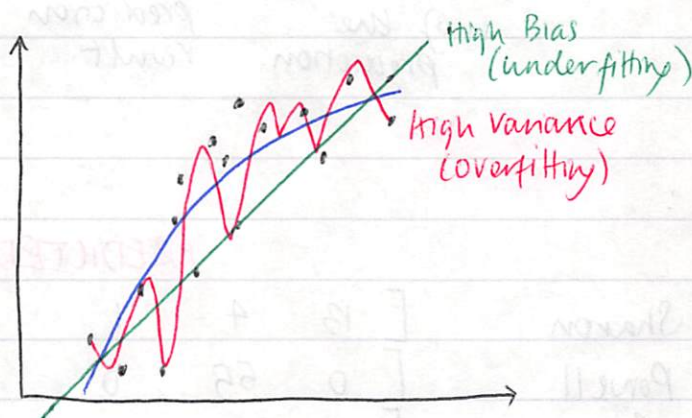
- * Pays too much attention to data (doesn't generalise well)

- * Overfit
- * Much higher error on test set than training set

- * Carefully minimised error over many features
→ optimised performance on training set

Want to use as few features as possible to

- maximise r^2
- minimise SSE



Data Types

* Numeric Data

- Discrete
- Continuous

* Categorical Data

- Can take on numerical values but these don't have mathematical meaning i.e. mean or variance don't make sense
- Ordinal data - Categories with ordering / rank

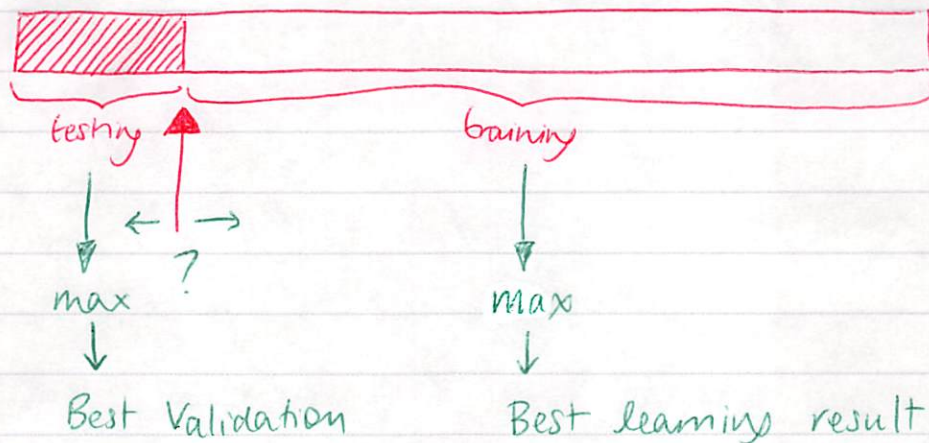
* Time Series Data

- Data point collected over time at

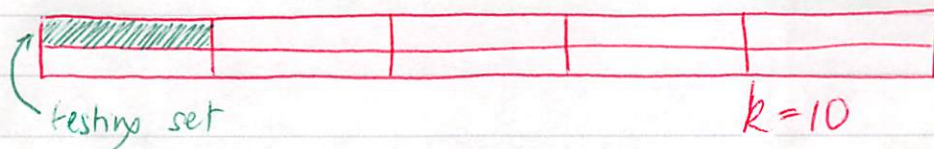
Cross Validation

* k -Fold Cross Validation

Problems with splitting into training and testing data



k fold cross validation



Partition your data into k equal parts and run k separate learning experiments. For each experiment

- pick a different testing set
- train on the remaining data

Average test results from those k experiments

This way we use all our data for both training and testing. Clearly training will take longer in this case than the simple train then test method of validation but k -fold cross validation results in better accuracy in our model.

Representative Power of a Model

* Curse of Dimensionality

As the number of number of features or equivalently dimensions grows, the amount of data we need to generalize accurately, grows exponentially

→ You're better off getting more data than you are adding more features!