



University of
South Australia

COMP 2019

Week 7
ML Training

Learning Objectives

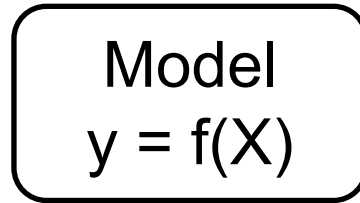
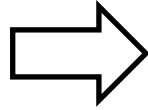
- Explain the machine learning process (CO3)
- Explain how data is prepared (CO3)
- Explain how ML models are trained (CO3)
- Explain how ML models are evaluated (CO3)



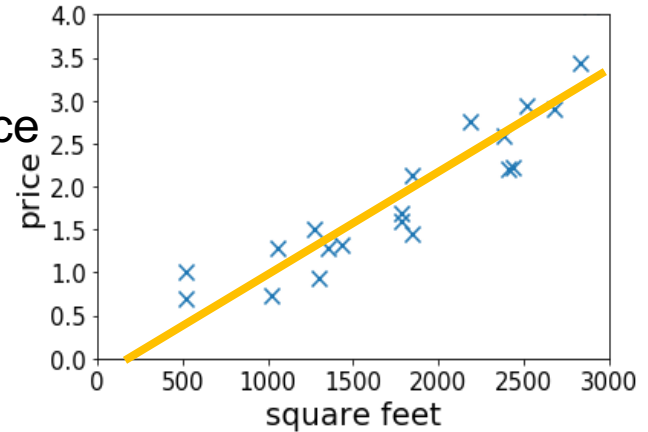
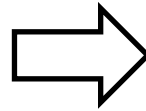
Supervised Learning from Data

x1	xn	y

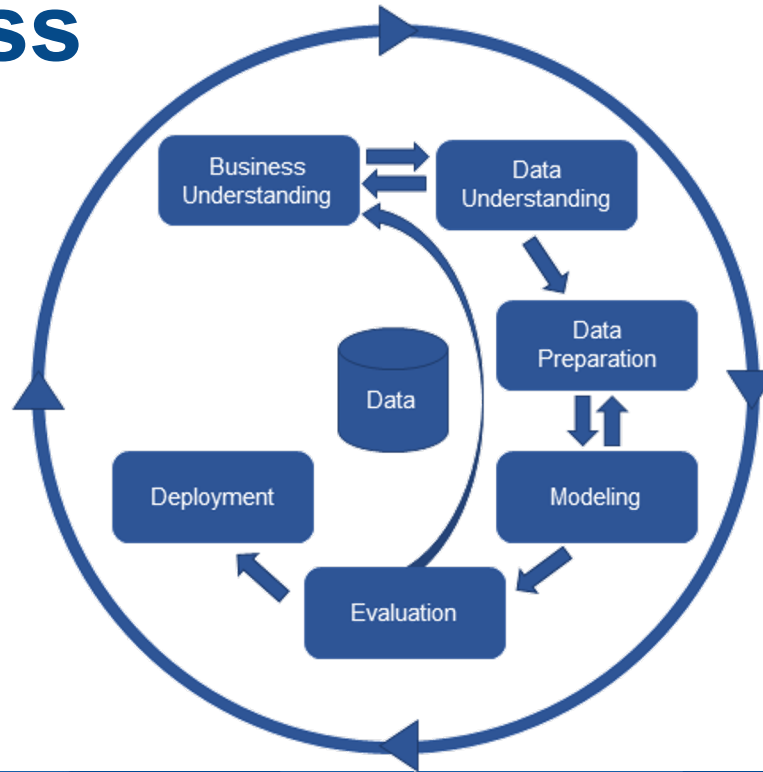
Training



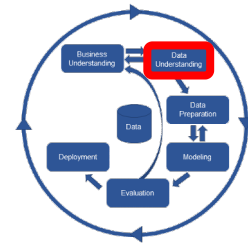
Inference



ML Process



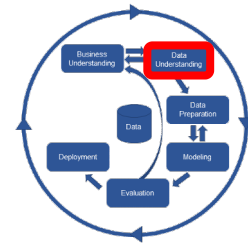
Understanding Data



- Key attribute distribution
- Label distribution
- Relationships between key attributes
- Attributes of important sub-populations
- Simple aggregation results
- Simple analysis of statistics



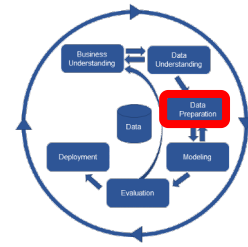
Examine the Quality of Data



- Is the data complete, covering all the required cases?
- Is the data correct?
 - How often do errors occur?
 - What is the nature of errors?
- Does the data contain missing values?
 - Where do they occur?
 - How they are they represented?
 - How frequent?
 - Systematic or random?



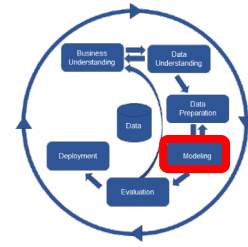
Preparing Data



- Cleaning the data
 - Rectify data quality issues
- Construct features
- Generate records
 - Negative cases may not be represented in the data
- Integration: combine data from multiple sources
- Aggregation



Data Sets



x1	...	xn	y

Never use the same data for training and testing

--	--	--	--

Training Dataset

x1	...	xn	y

Test Dataset

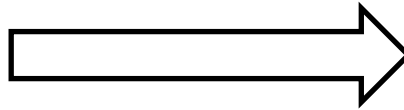
x1	...	xn	y



Model Fitting

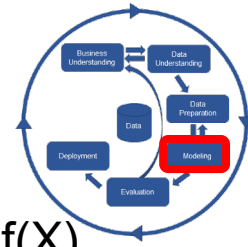
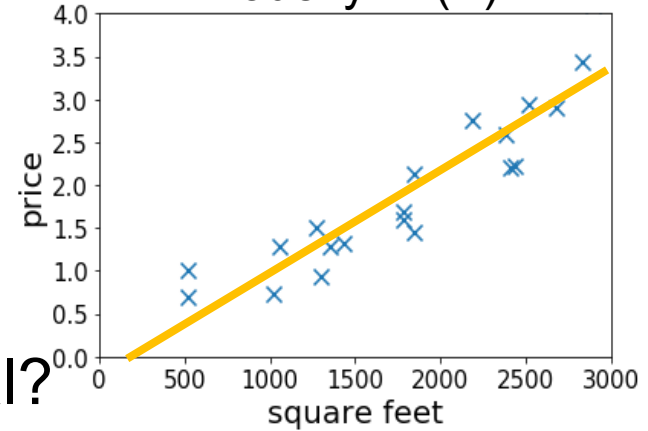
Training Dataset

x1	...	xn	y

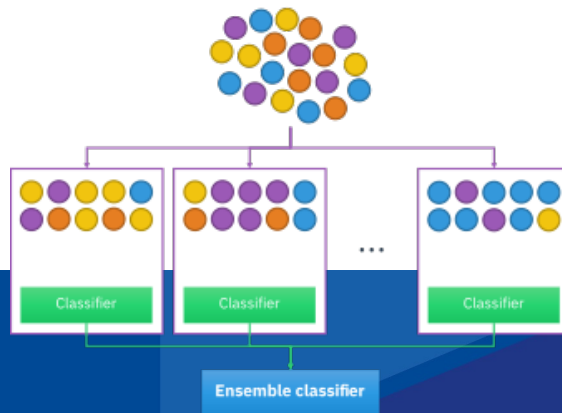
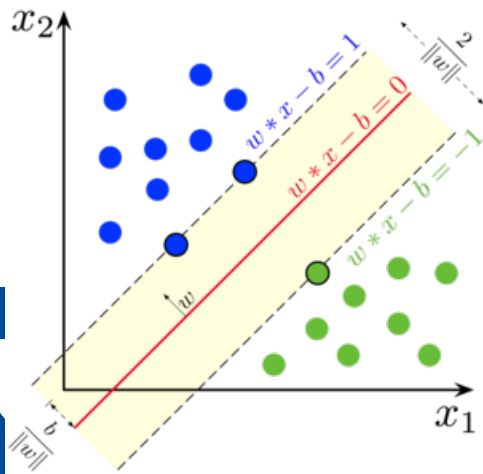
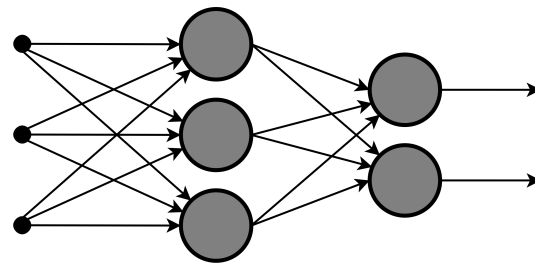
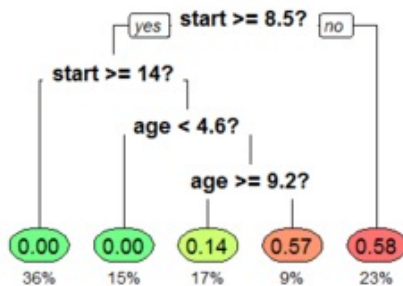
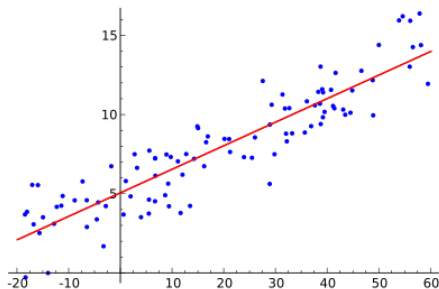
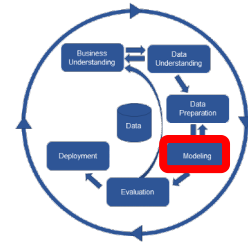


- Type of model?
- Optimisation goal?
- Learning algorithm?

Model $y = f(X)$



Types of Model



THIS IS YOUR MACHINE LEARNING SYSTEM?

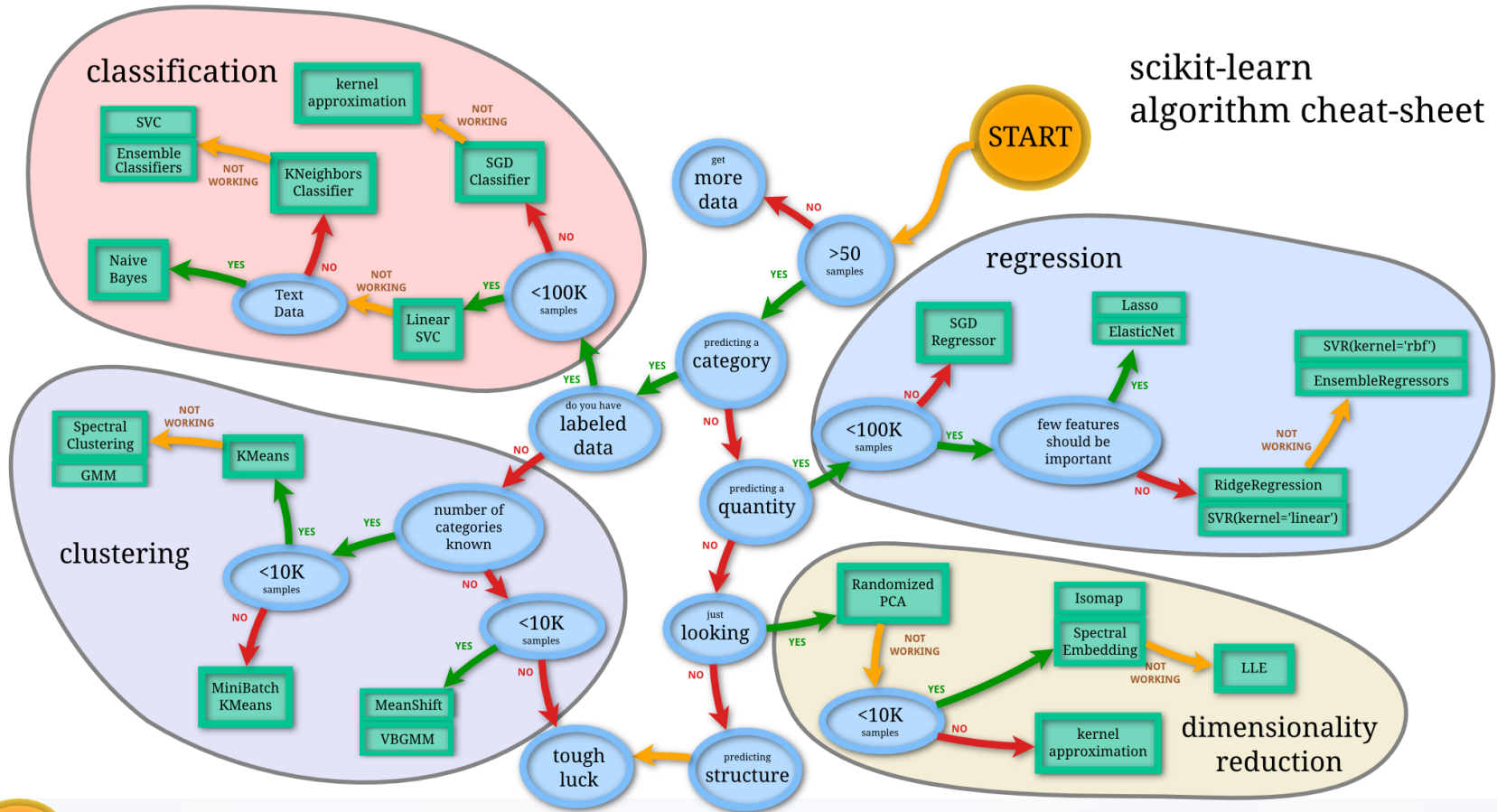
YUP! YOU POUR THE DATA INTO THIS BIG
PILE OF LINEAR ALGEBRA, THEN COLLECT
THE ANSWERS ON THE OTHER SIDE.

WHAT IF THE ANSWERS ARE WRONG?

JUST STIR THE PILE UNTIL
THEY START LOOKING RIGHT.



scikit-learn algorithm cheat-sheet



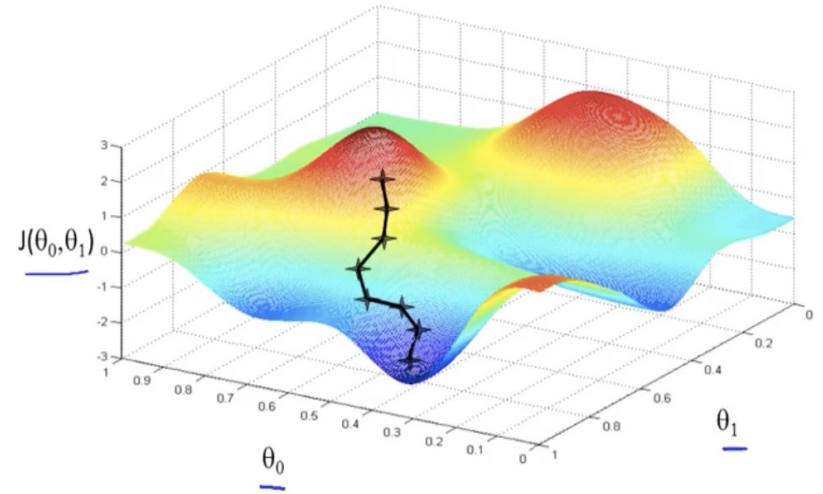
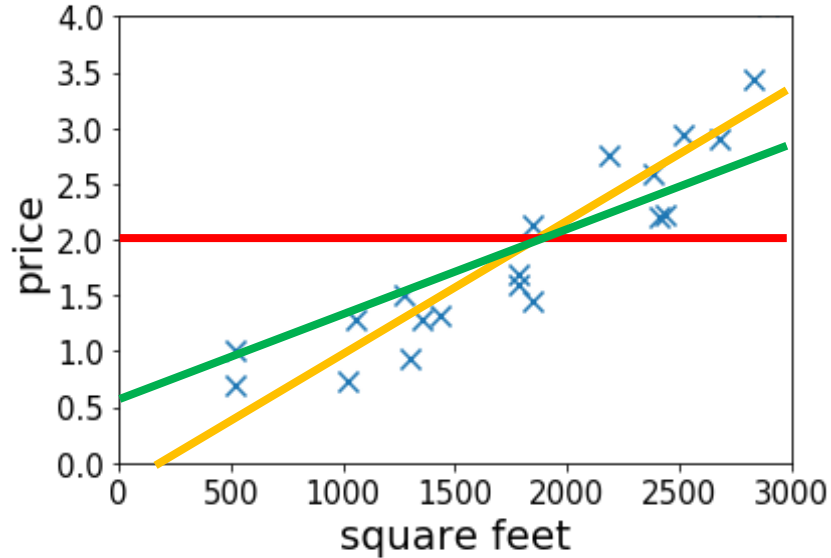
Back



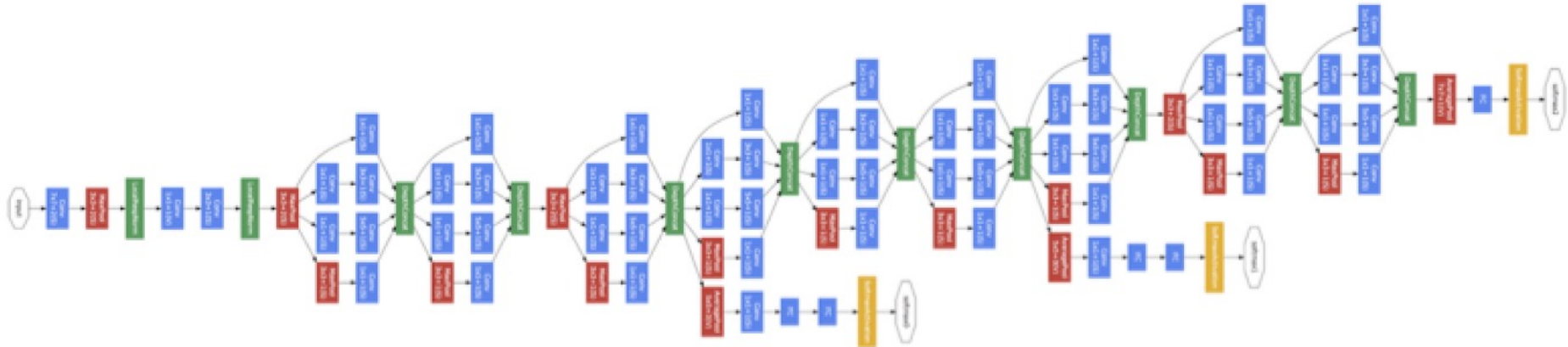
scikit
learn

Training as Optimisation

$$\text{price} = \theta_0 \times \text{sqFeet} + \theta_1$$



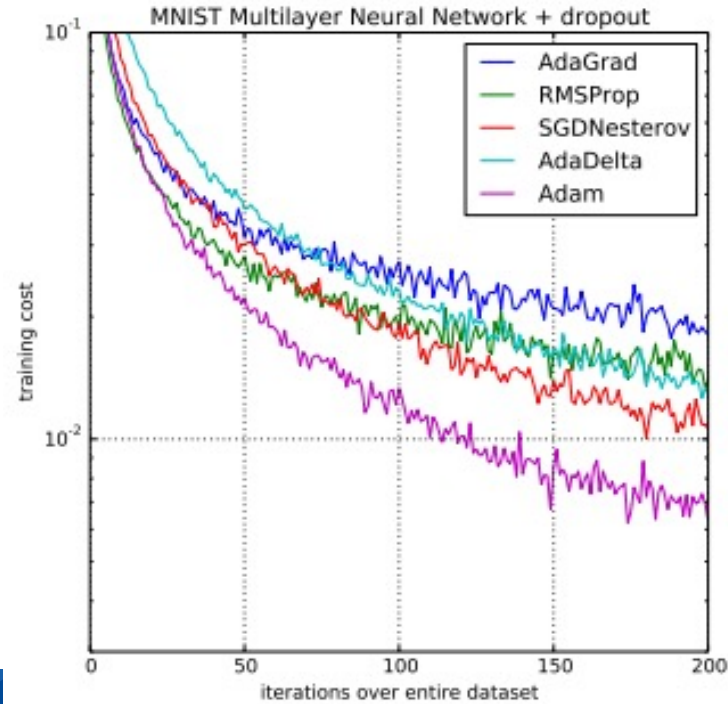
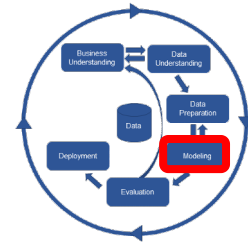
Deep Neural Net



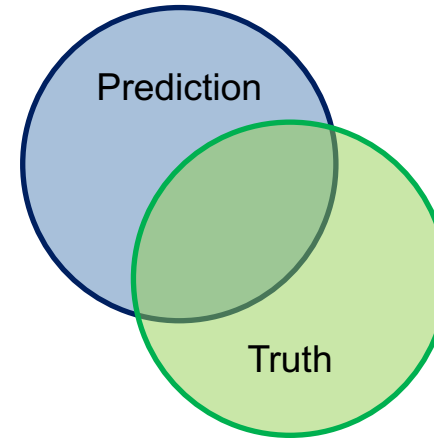
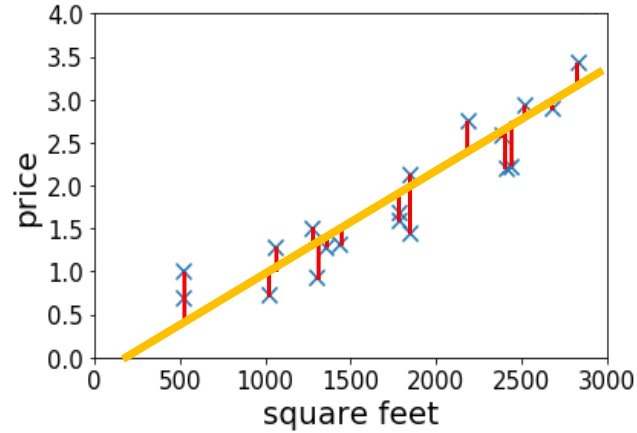
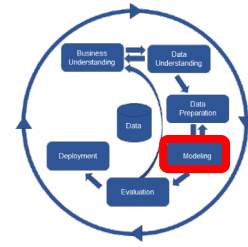
GoogLeNet/InceptionV1 (2014) has 22 layers and 4 million parameters



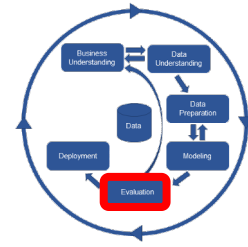
Optimisation Algorithm



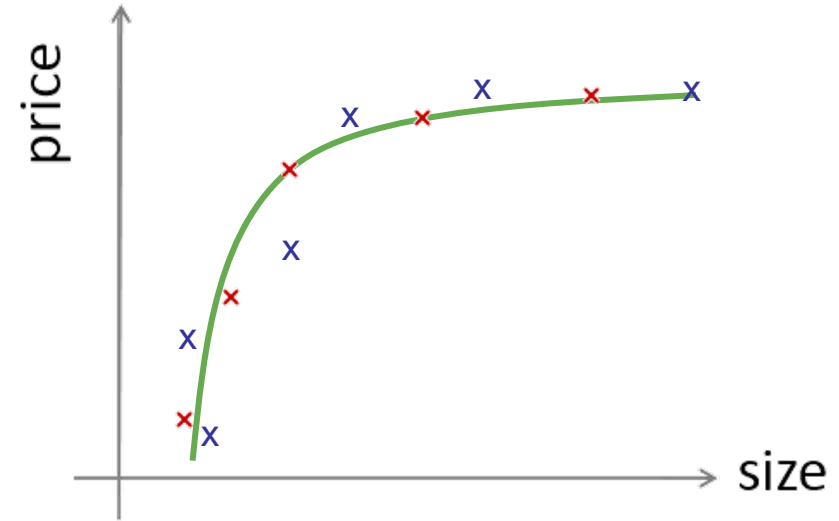
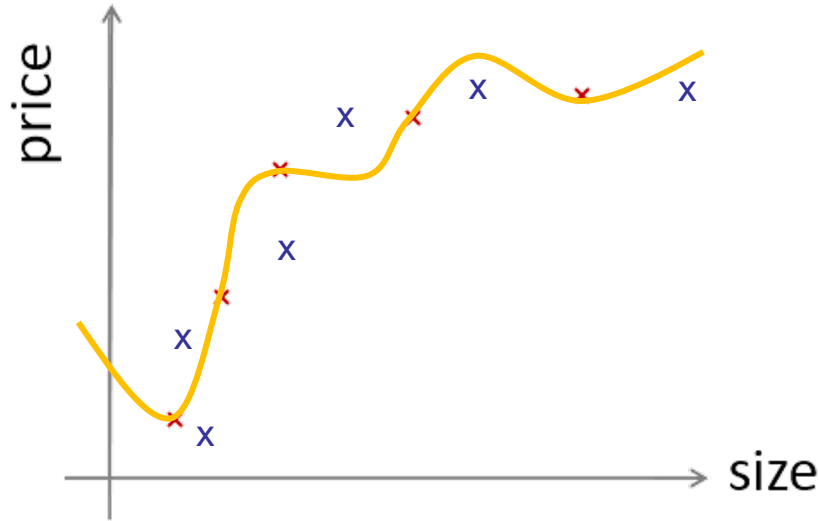
Metrics and Loss



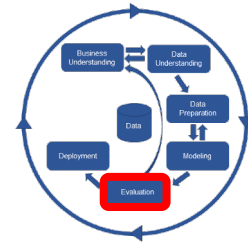
Learning = Generalisation



- How well does the model perform on **UNSEEN** data?

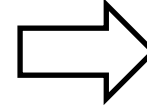
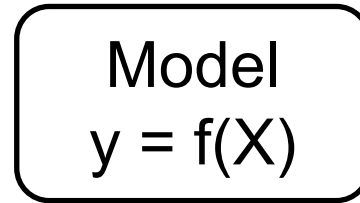
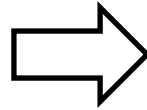


Model Evaluation

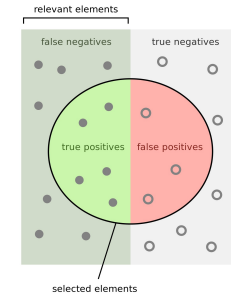


Test Dataset

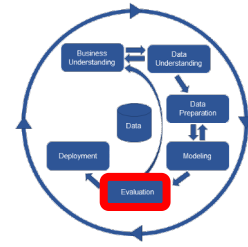
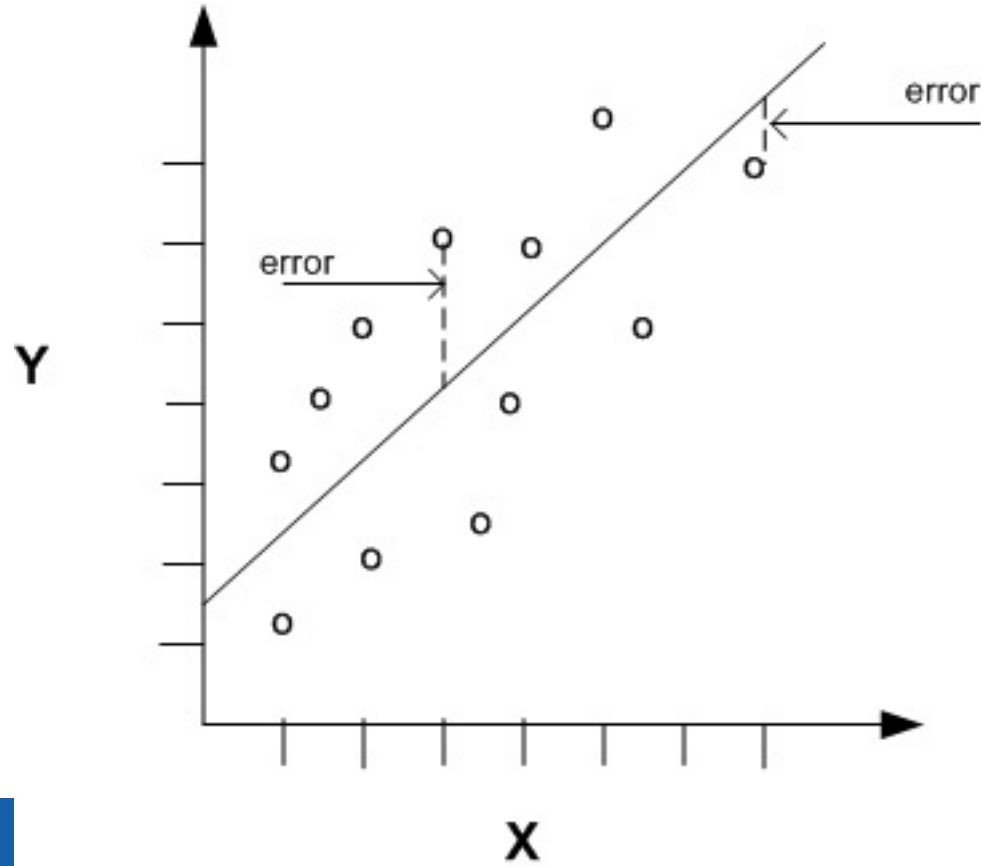
x1	...	xn	y



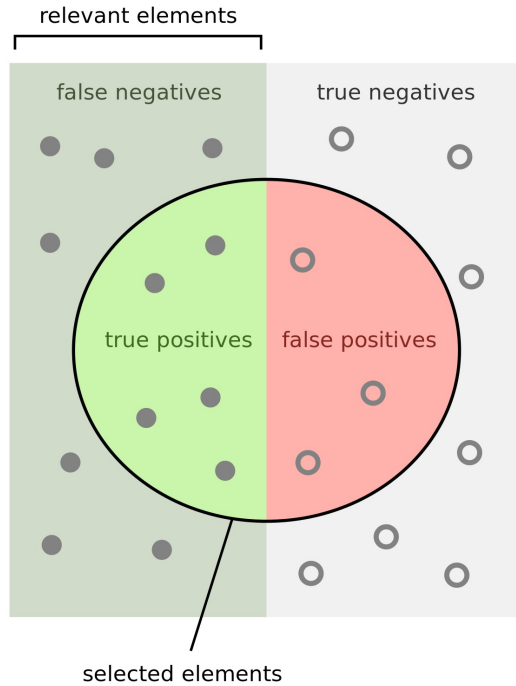
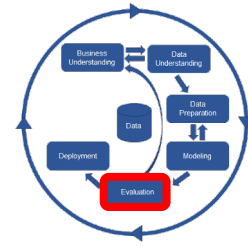
Metrics



(R)MSE



Accuracy and Precision



How many decisions
are correct?

$$\text{Accuracy} = \frac{\text{true positives} + \text{true negatives}}{\text{true positives} + \text{false positives} + \text{false negatives} + \text{true negatives}}$$

How many selected
items are relevant?

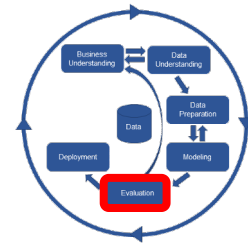
$$\text{Precision} = \frac{\text{true positives}}{\text{true positives} + \text{false positives}}$$

How many relevant
items are selected?

$$\text{Recall} = \frac{\text{true positives}}{\text{true positives} + \text{false negatives}}$$



Precision AND Recall?

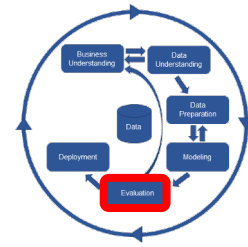


	Precision	Recall
Classifier 1	0.5	0.4
Classifier 2	0.7	0.1
Classifier 3	0.02	1.0

- Usually there is a trade-off between Precision and Recall
- Mean ($\frac{P+R}{2}$) is not meaningful
 - Classifier 3 has highest mean but predicts 'positive' all the time



F₁ Score

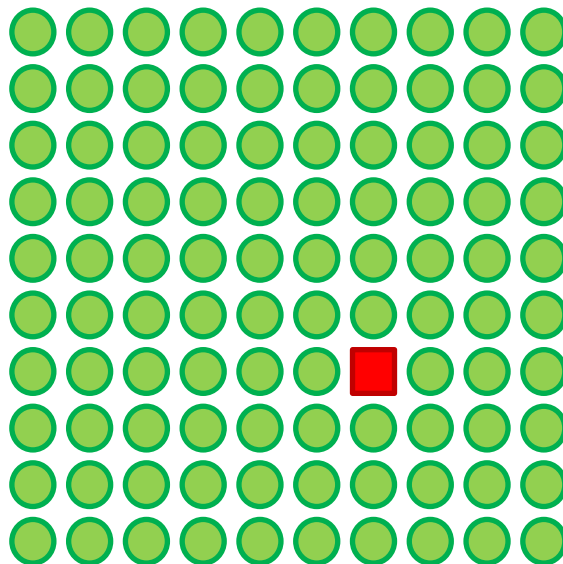
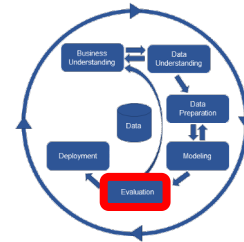


	Precision	Recall	F ₁
Classifier 1	0.5	0.4	0.444
Classifier 2	0.7	0.1	0.175
Classifier 3	0.02	1.0	0.039

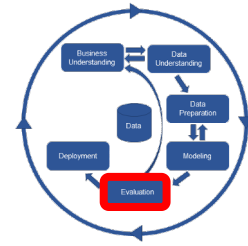
- $F_1 = 2 \frac{P \times R}{P + R}$
 - Higher is better
 - F_1 is 0 if $P=0$ or $R=0$
 - F_1 is 1 if $P=1$ and $R=1$



Accuracy Paradox



Contingency Tables: > 2 Classes



	Predicted C1	Predicted C2	Predicted C3
Actual C1	100	40	25
Actual C2	25	50	4
Actual C3	1	0	7



Summary

- ML systems are created in a 6 step process
 - Data collection & preparation is where most effort is spent
- ML Models are trained and evaluated on separate data sets
- Training means choosing parameters that optimise a metric
- Knowing the distribution of (unseen) data is important for selecting models and metrics





**University of
South Australia**

Questions?