
Oleksandr Romanko, Ph.D.

Associate Director, Financial Risk Quantitative Research, SS&C Algorithmics
Adjunct Professor, University of Toronto

MIE1624H – Introduction to Data Science and Analytics Lecture 7 – Machine Learning

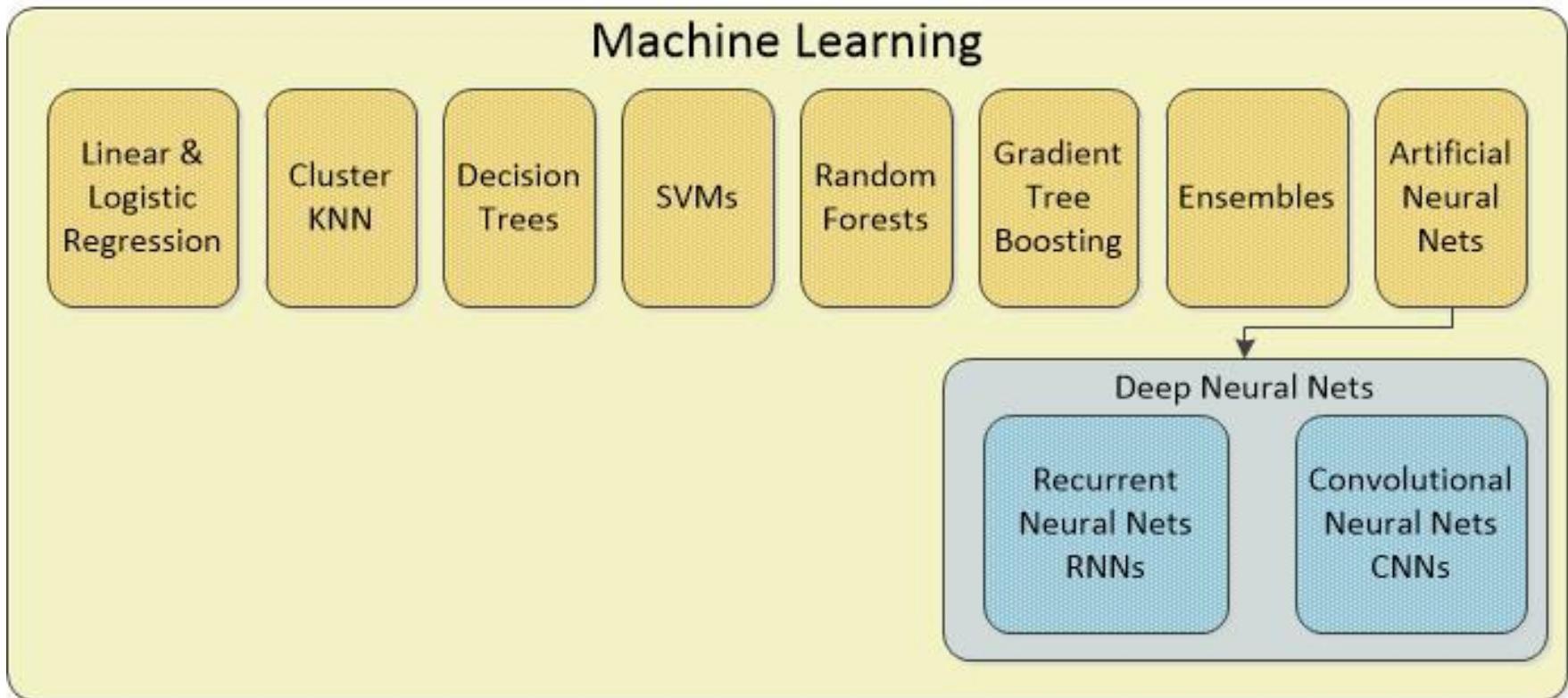
Machine learning

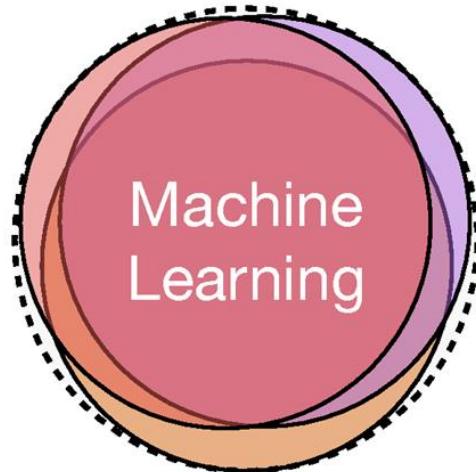
Machine learning gives computers the ability to learn without being explicitly programmed

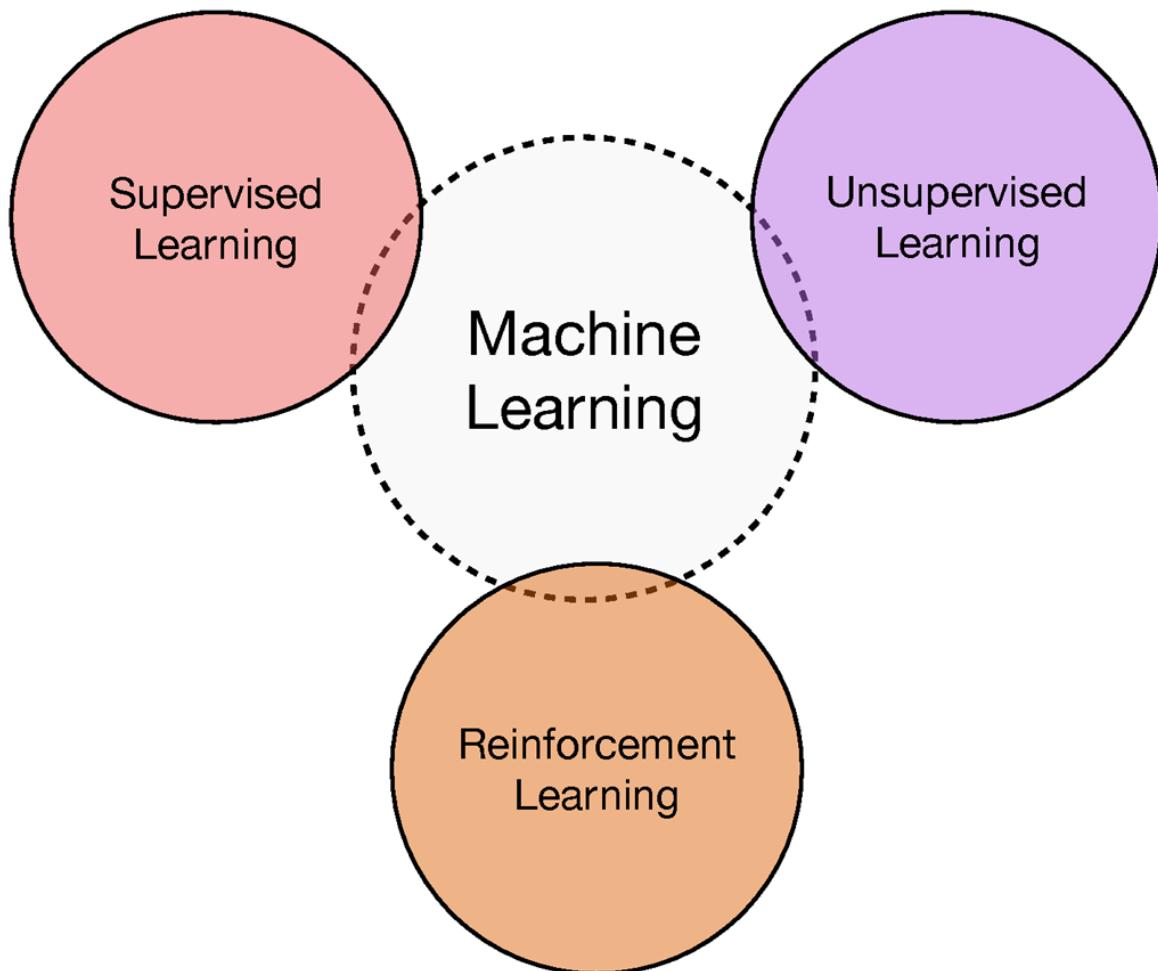
- **Supervised learning**: decision trees, ensembles (bagging, boosting, random forests), k-NN, linear regression, Naive Bayes, neural networks, logistic regression, SVM
 - Classification
 - Regression (prediction)
- **Unsupervised learning**: k-means, c-means, hierarchical clustering, DBSCAN
 - Clustering
 - Dimensionality reduction (PCA, LDA, factor analysis, t-SNE)
 - Association rules (market basket analysis)
- **Reinforcement learning**
 - Dynamic programming
- **Neural nets**: deep learning, multilayer perceptron, recurrent neural network (RNN), convolutional neural network (CNN), generative adversarial network (GAN)

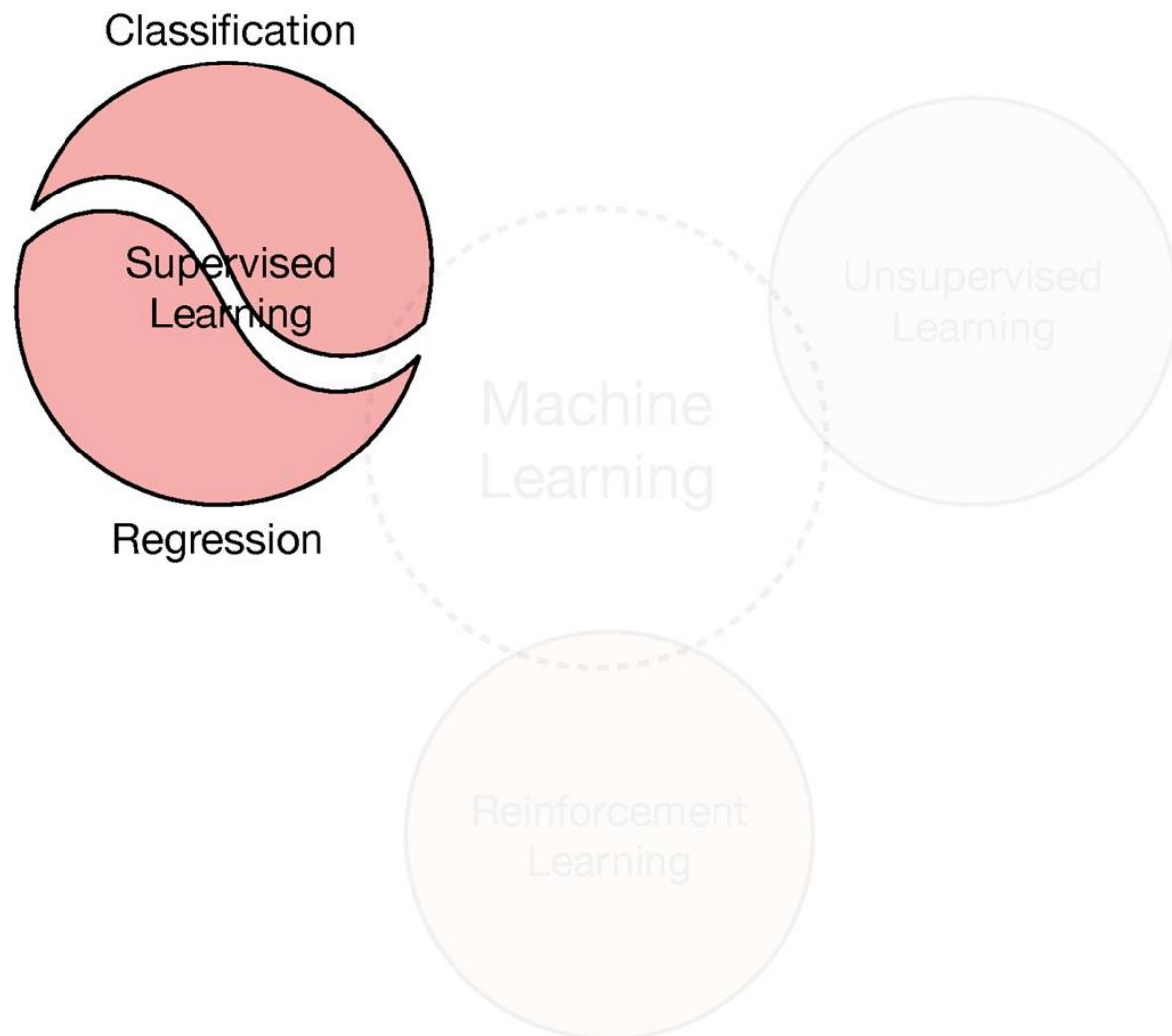
Machine learning

Machine learning gives computers the ability to learn without being explicitly programmed

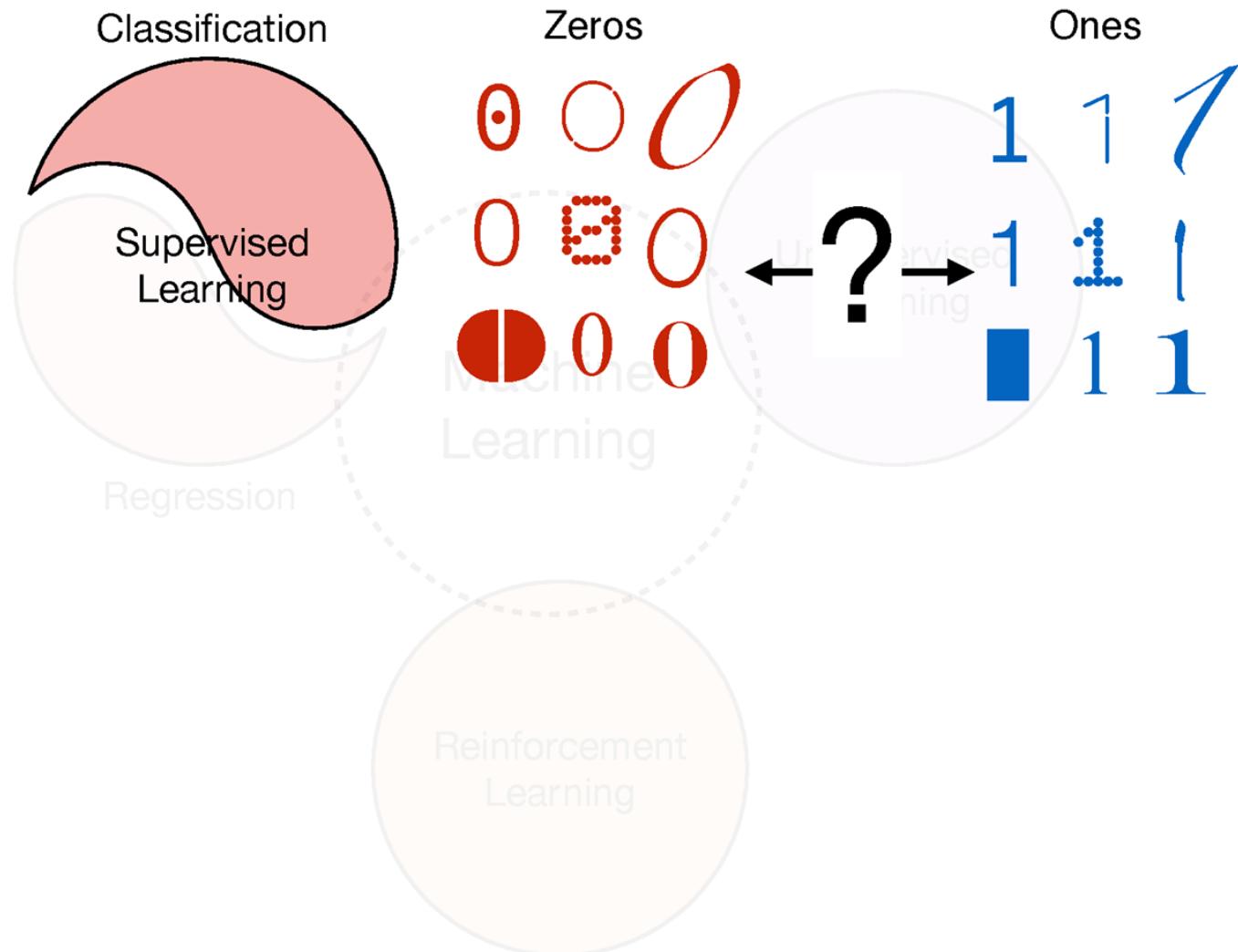








Binary Classification



Binary Classification

Classification

Supervised Learning

Regression

Zeros

0 0 0
0 0 0
0 0 0

Ones

1 1 1
1 1 1
1 1 1

→?←

Multiclass Classification

Reinforcement Learning

? →

0 0 0 0
1 1 1 1
2 2 2 2
3 3 3 3



Regression

Supervised Learning

7\$

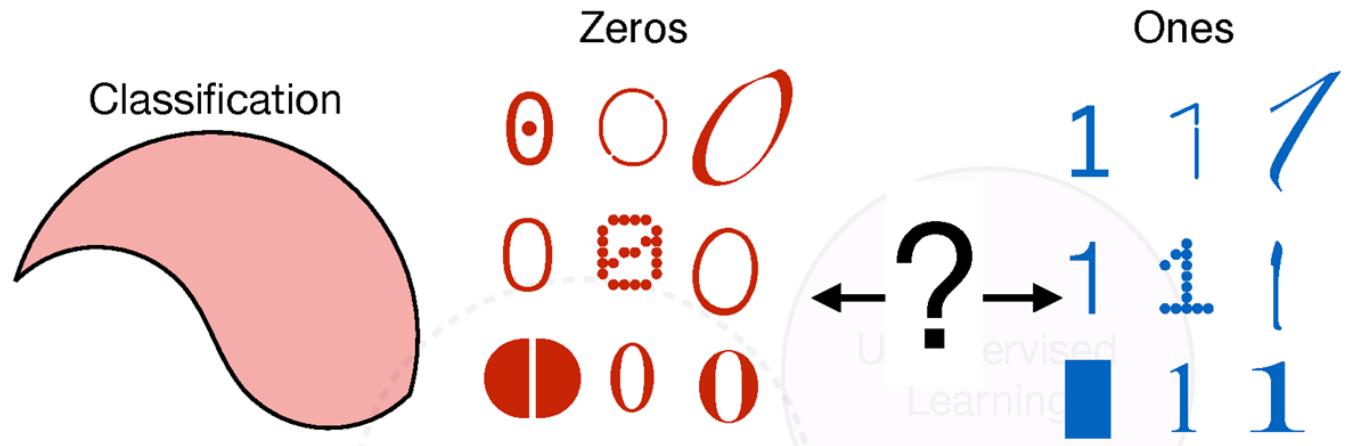
5\$

3.5\$

?\$

Predicting **continuous** values

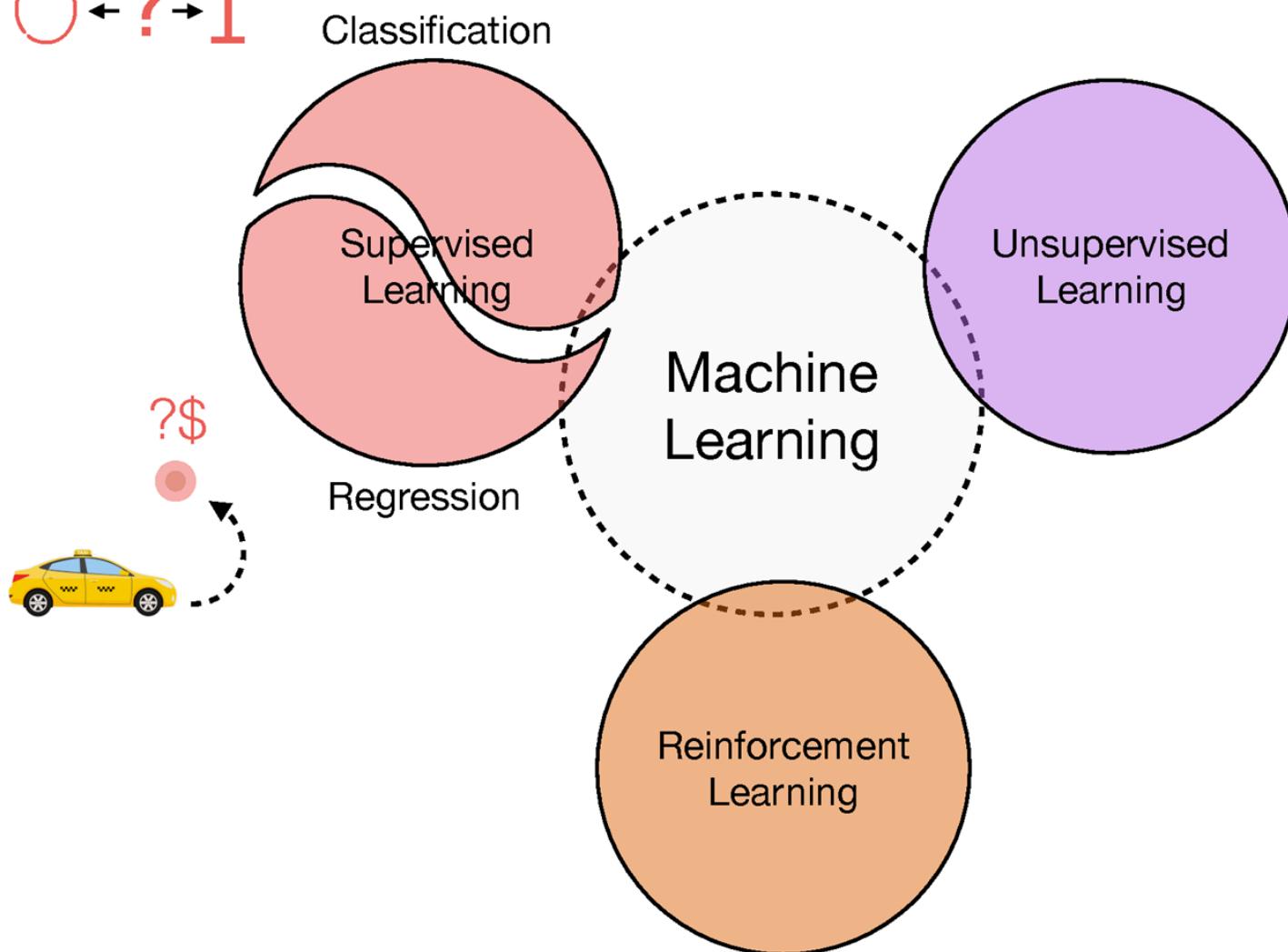
Classification



What is the main similarity/difference between these two classes of **supervised learning**?

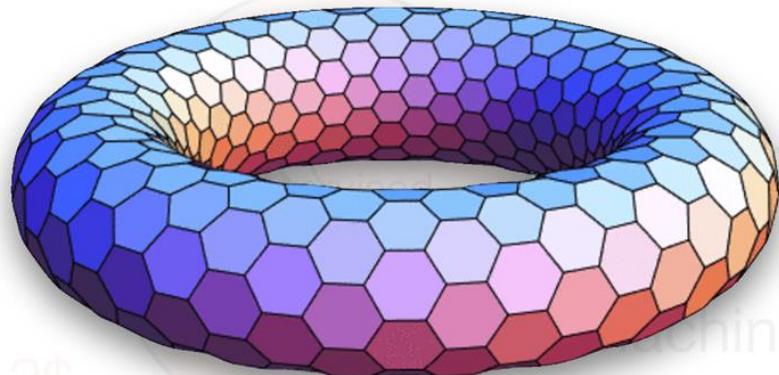


0 ← ? → 1

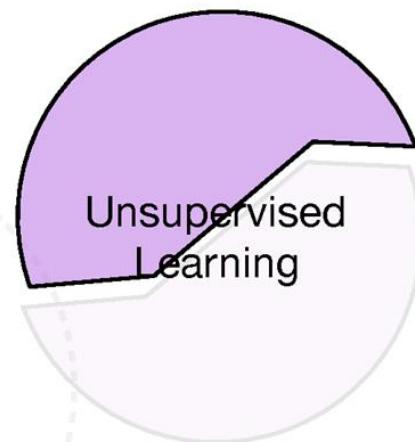




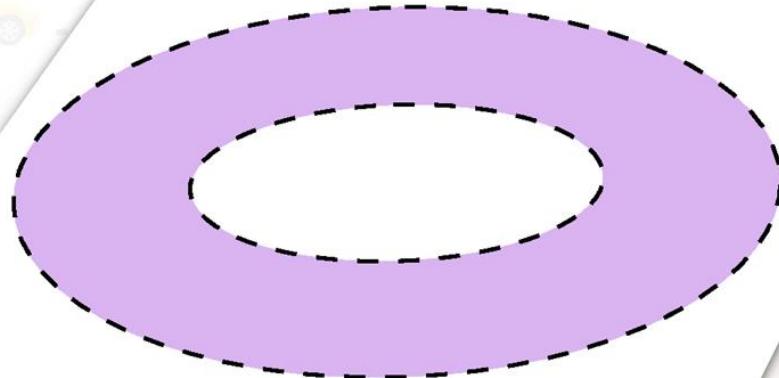
High Dimensional Space



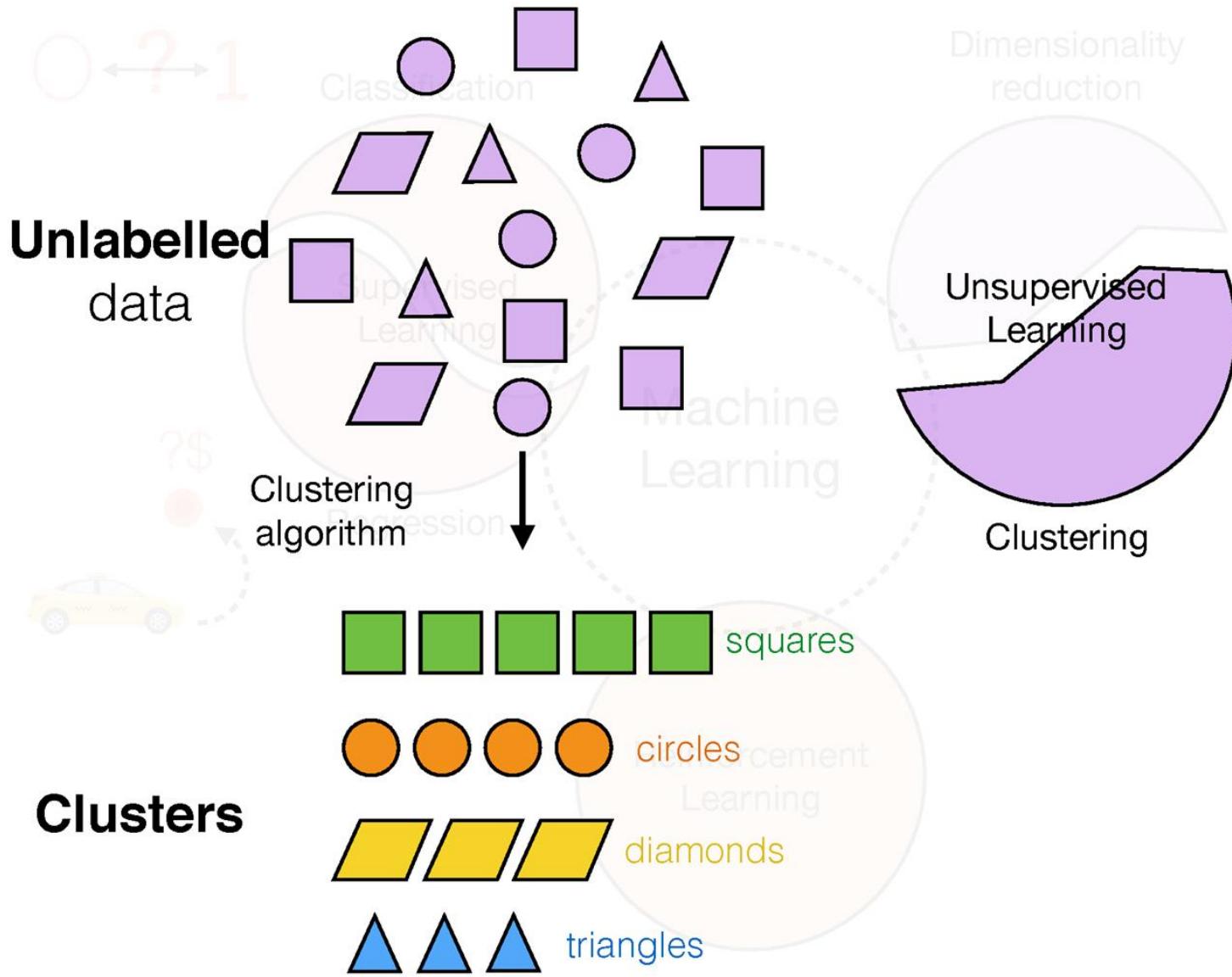
Dimensionality reduction



Dimensionality reduction



Low Dimensional Space



○ ← ? → 1

Classification

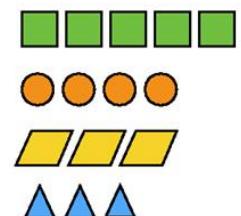
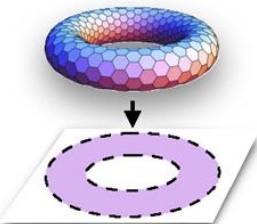
Supervised
Learning

Regression

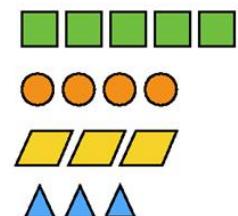
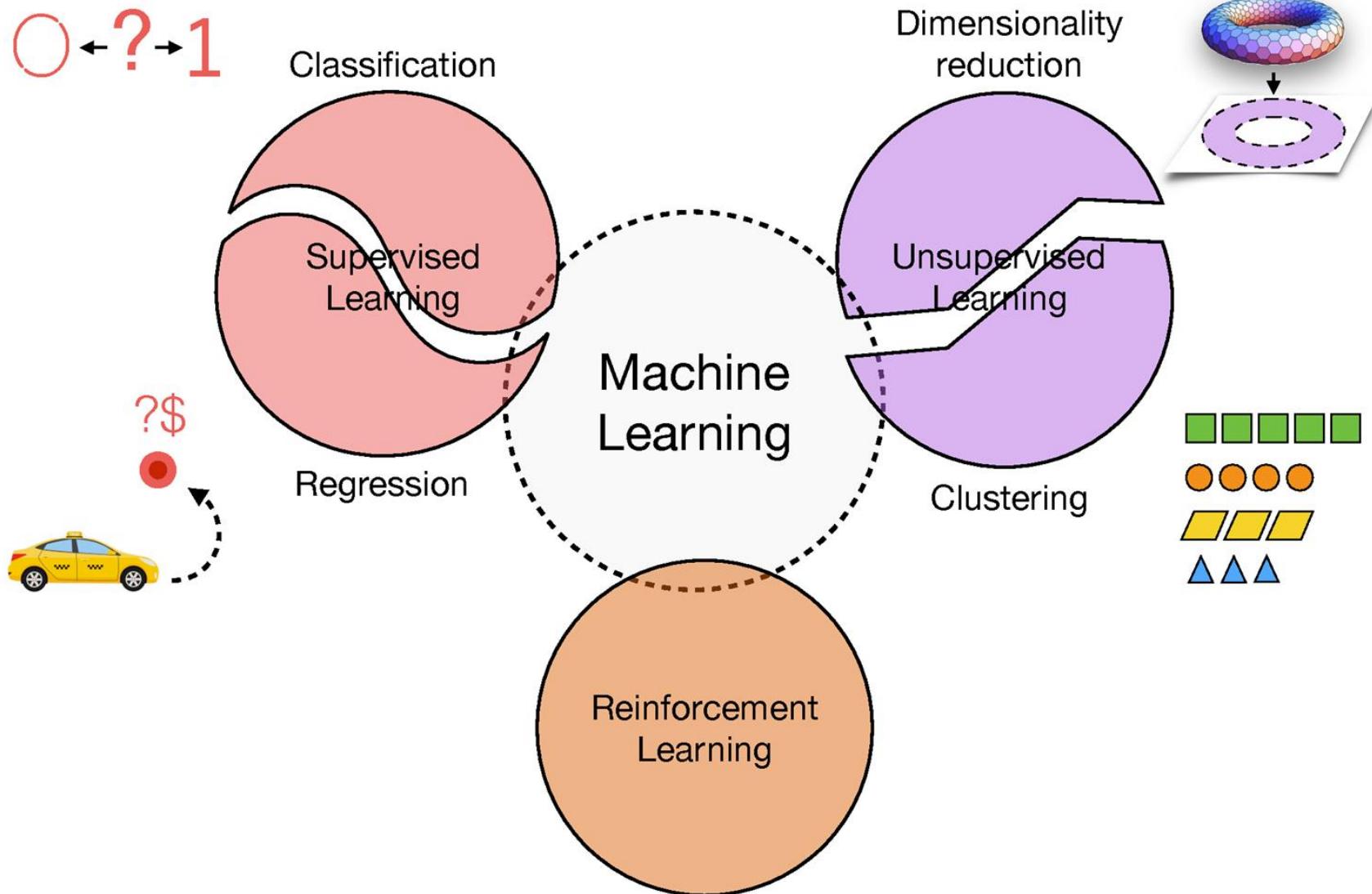
Dimensionality
reduction

Unsupervised
Learning

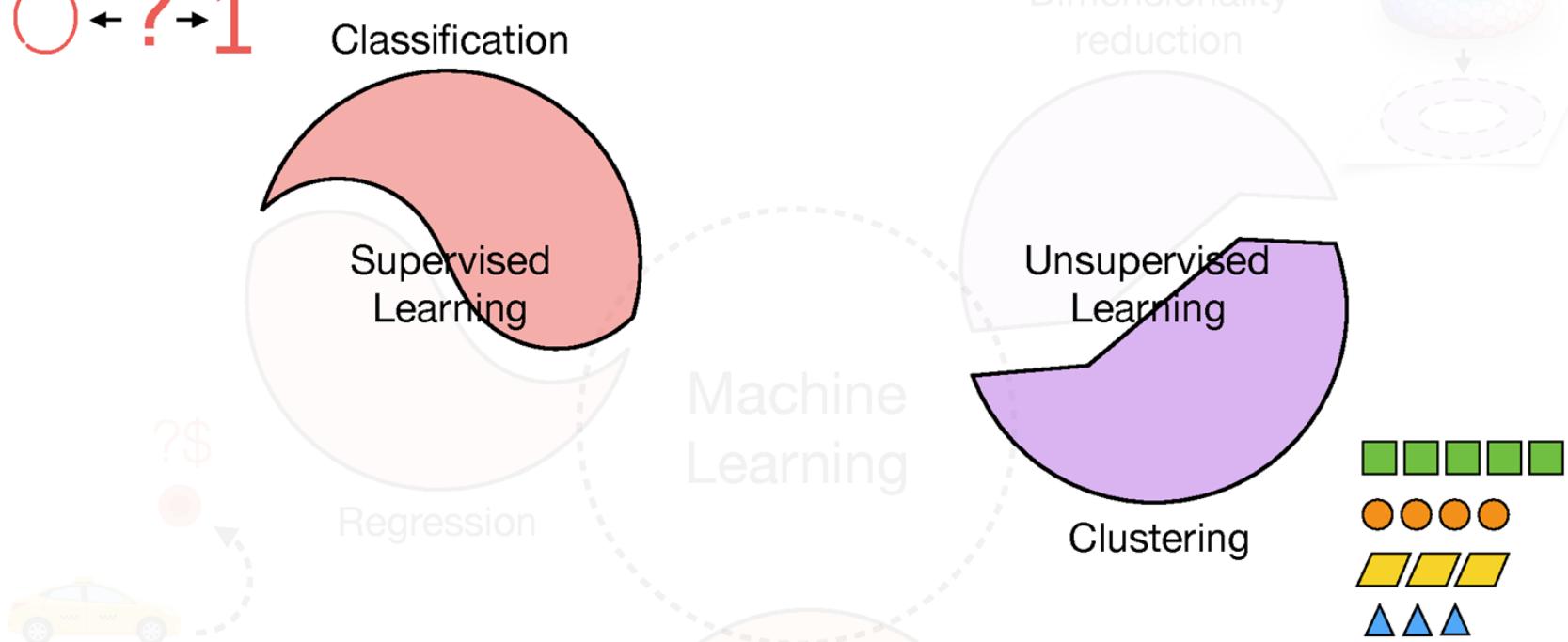
Clustering



0 ← ? → 1



0 ← ? → 1



What is the main similarity/difference between
classification and **clustering**?

○ ← ? → 1

Classification

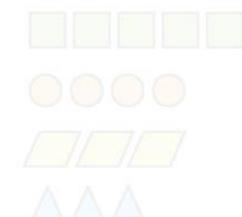


Regression

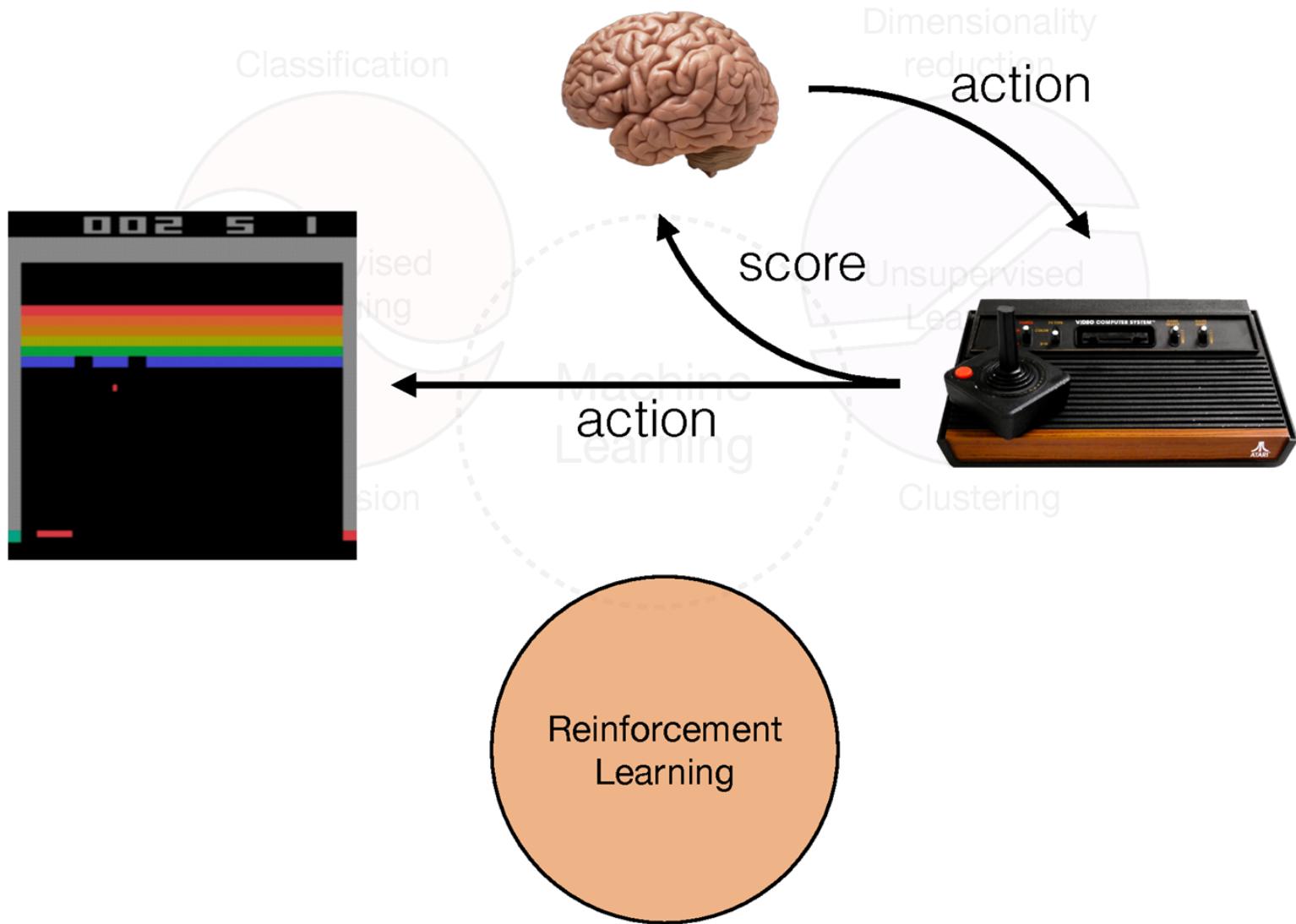
Dimensionality reduction

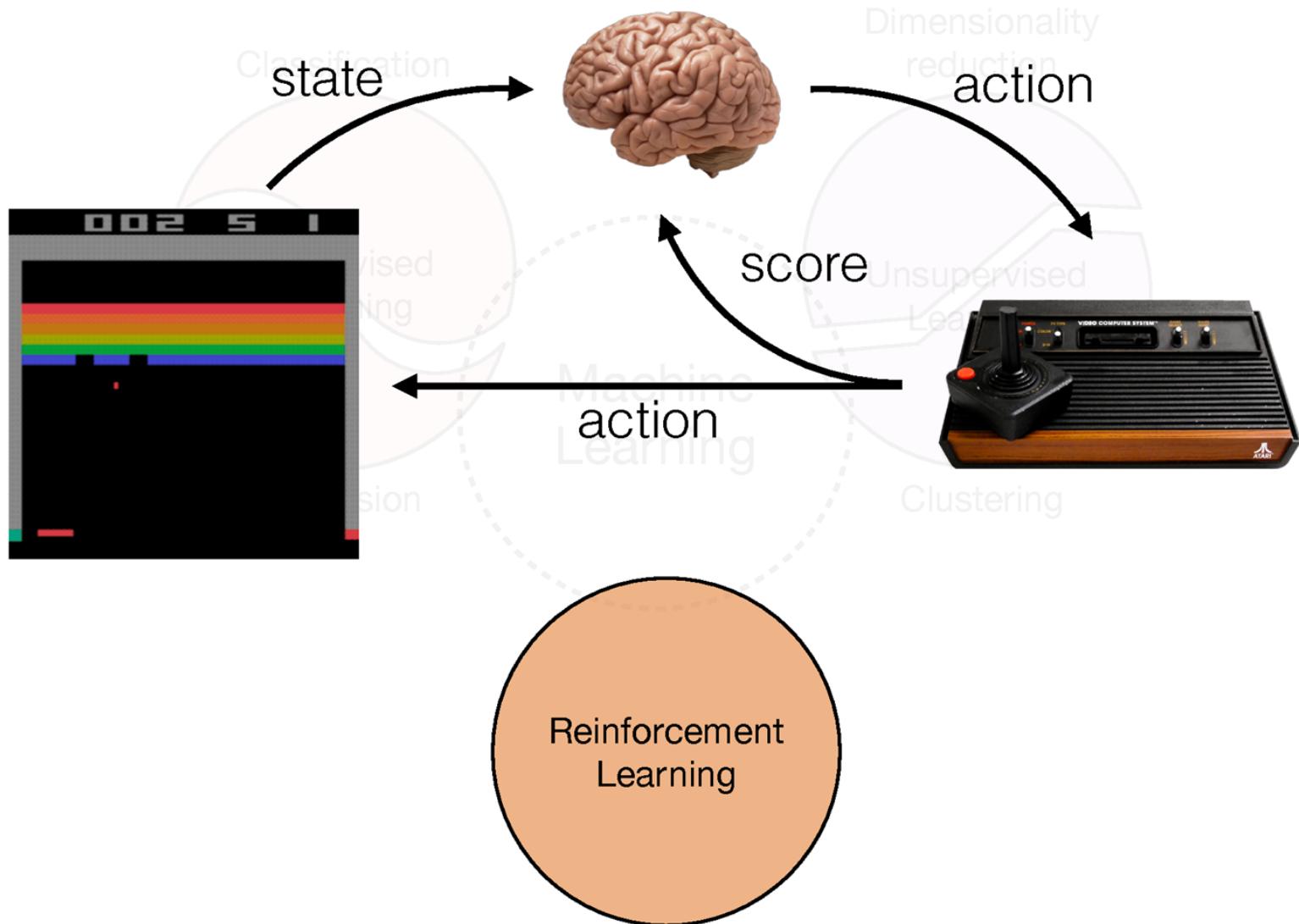


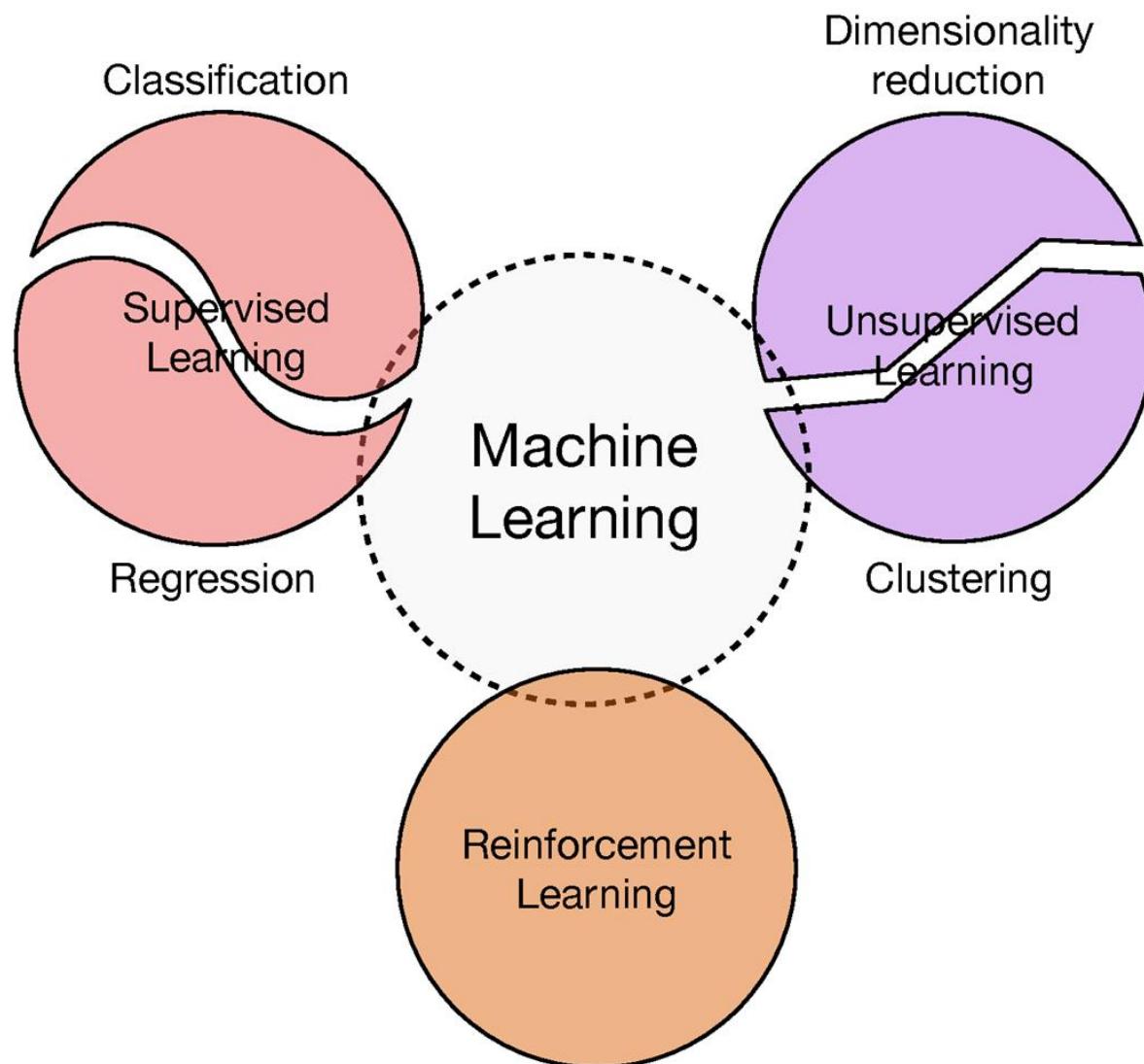
Clustering



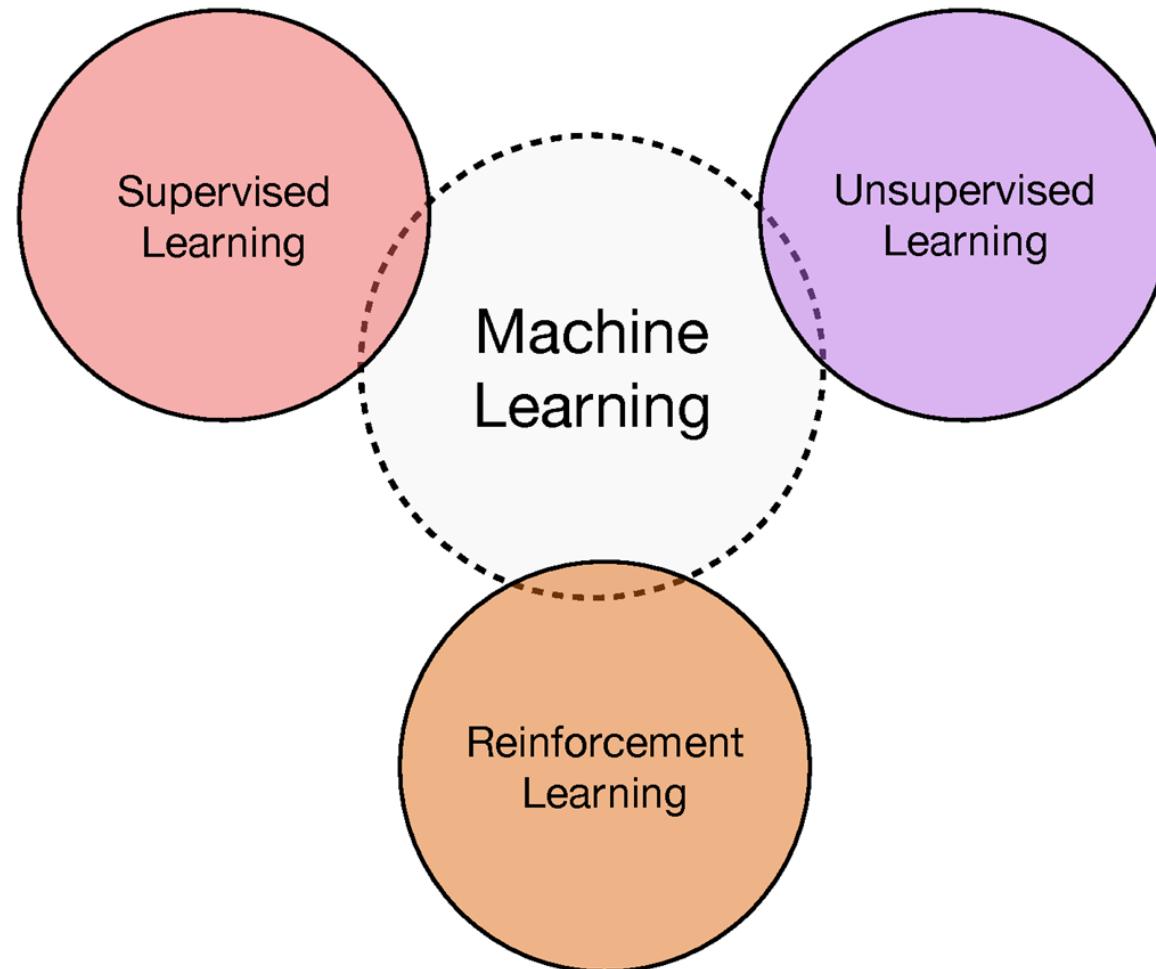
Reinforcement
Learning







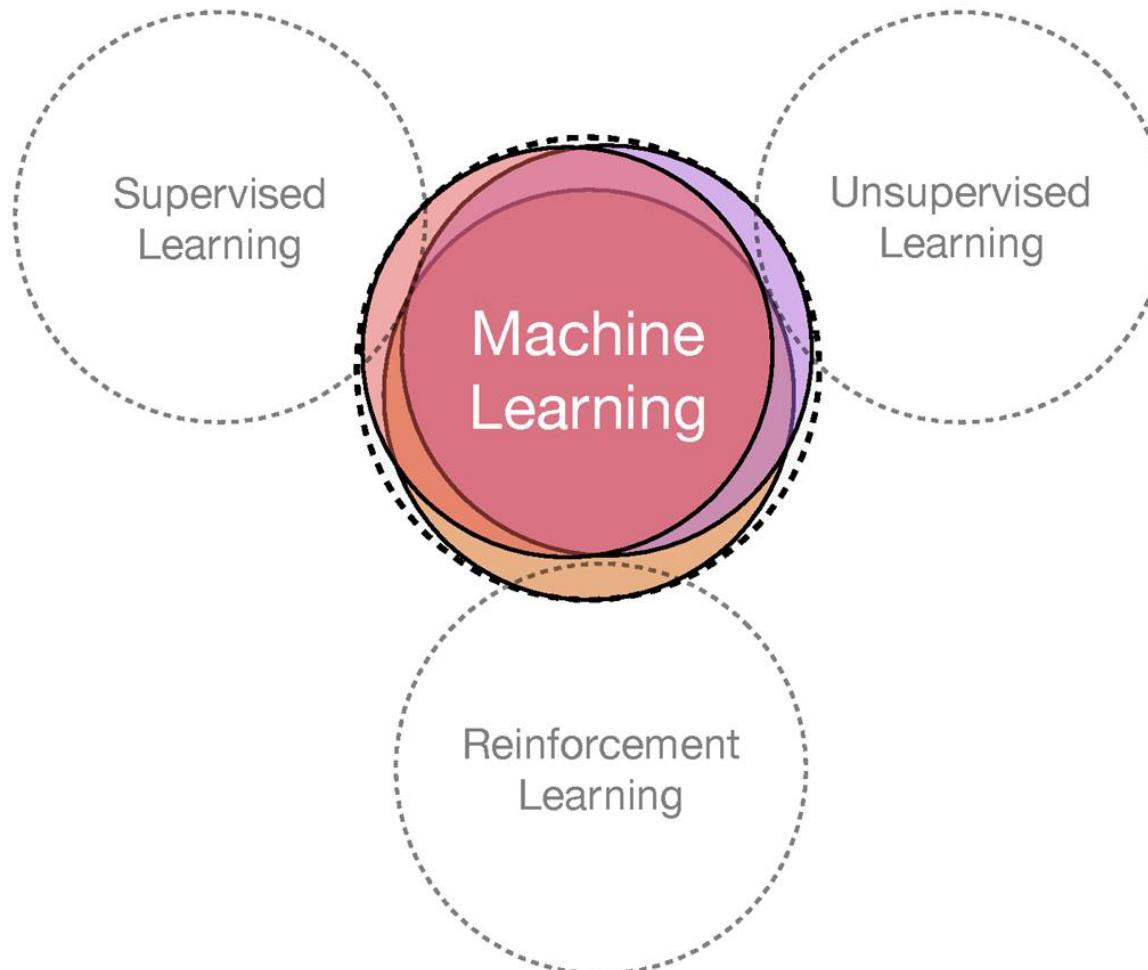
Lecture #5, #7, #8



Lecture #8

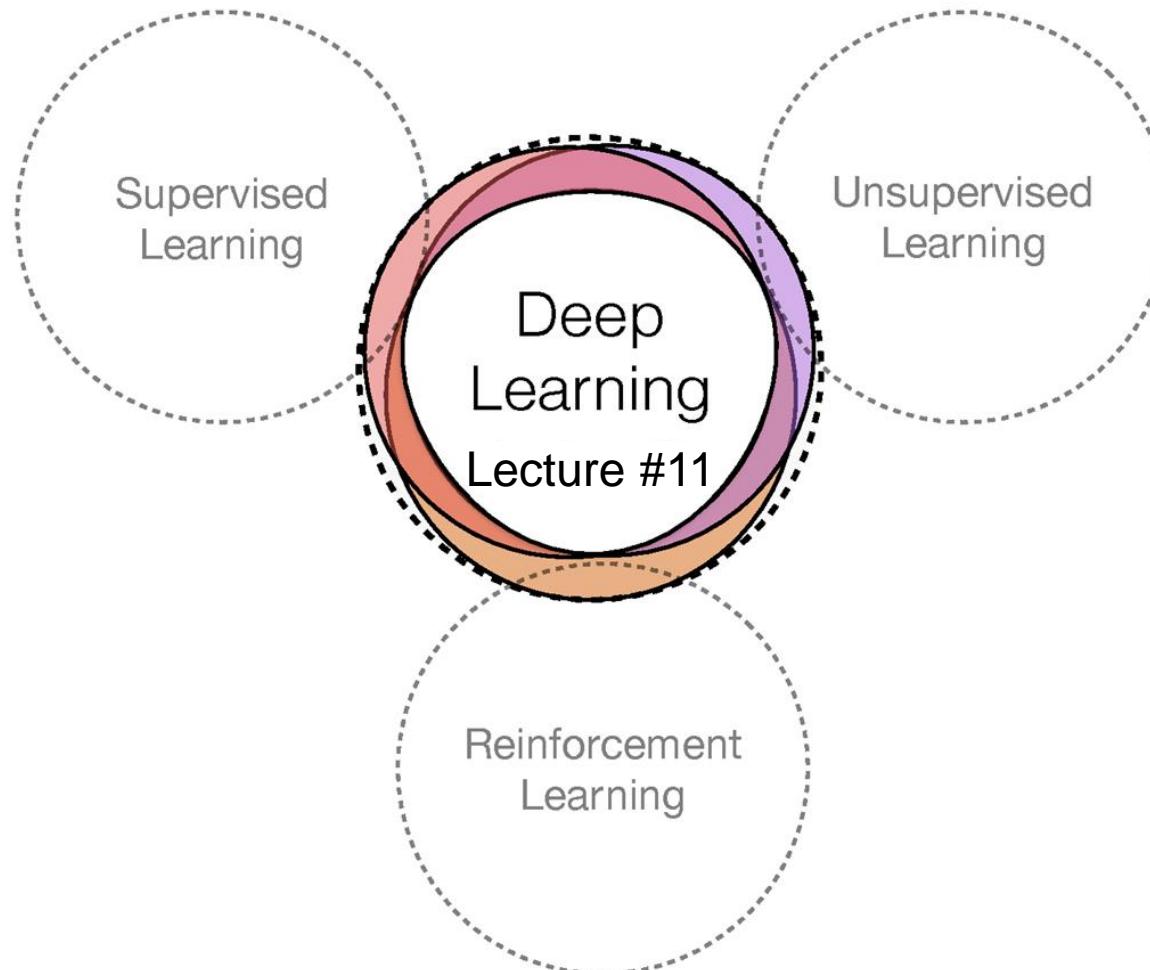
Lecture #5, #7, #8

Lecture #8



Lecture #5, #7, #8

Lecture #8



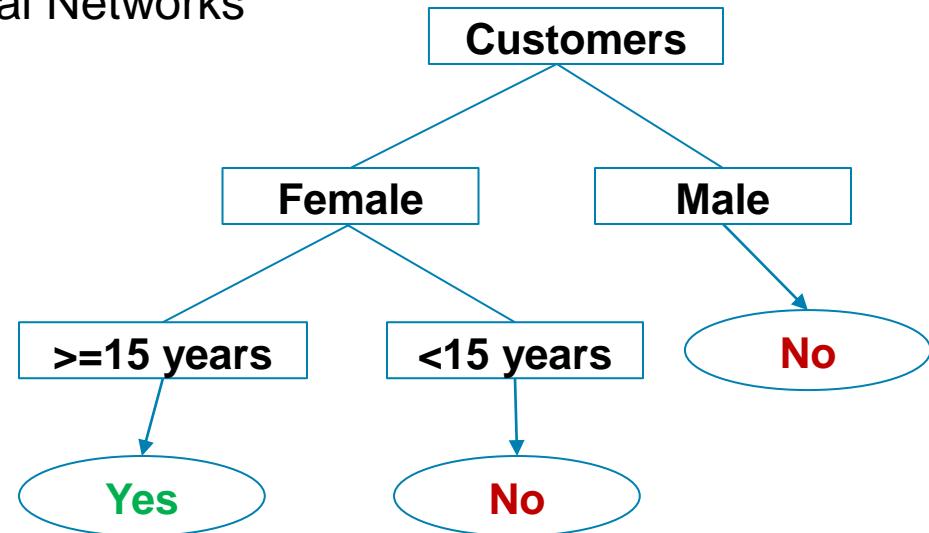


Supervised Machine Learning – Decision Trees

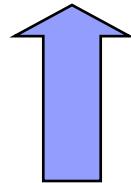
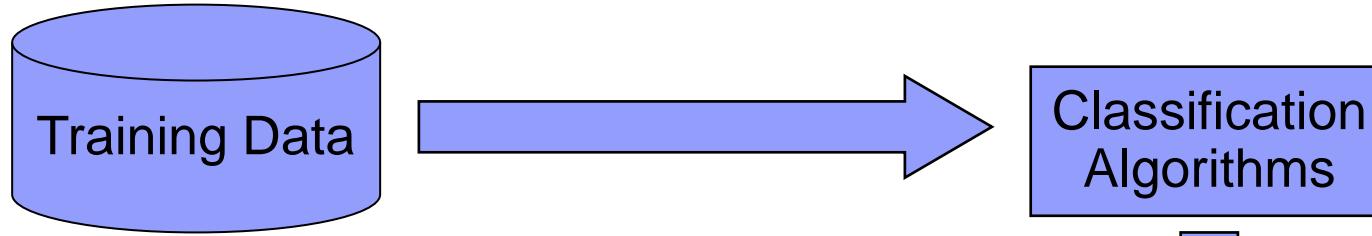
Classification

- Classification is a supervised learning technique, which maps data into predefined classes or groups
- Training set contains a set of records, where one of the records indicates class
- Modeling objective is to assign a class variable to all of the records, using attributes of other variables to predict a class
- Data is divided into test / train, where “train” is used to build the model and “test” is used to validate the accuracy of classification
- Typical techniques: Decision Trees, Neural Networks

Gender	Age	Lipstick
Female	21	Yes
Male	30	No
Female	14	No
Female	35	Yes
Male	17	No
Female	16	Yes

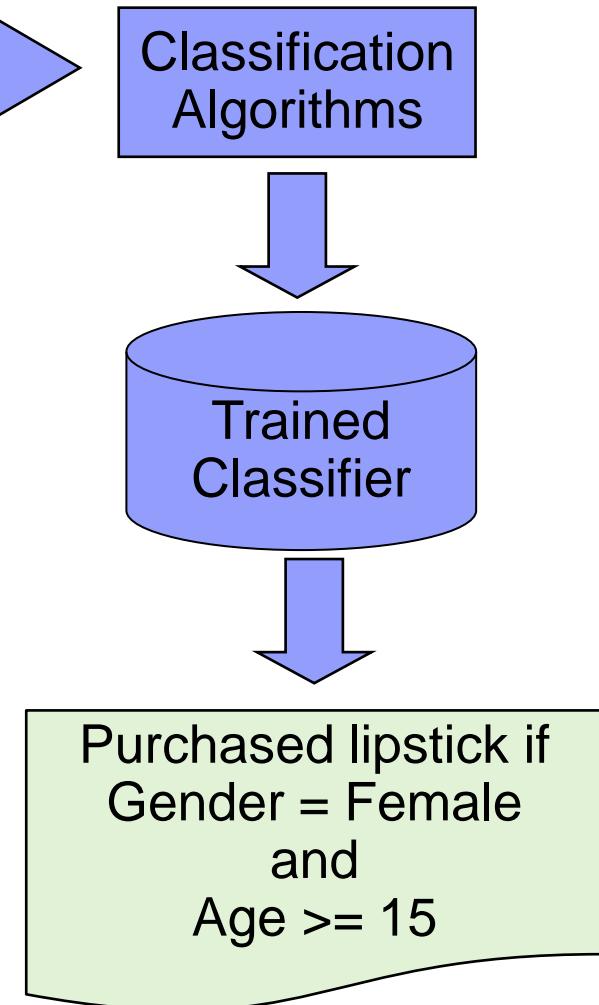


Classification: creating model



Works with both interval and categorical variables

Gender	Age	Lipstick
Female	21	Yes
Male	30	No
Female	14	No
Female	35	Yes
Male	17	No
Female	16	Yes



Classification: applying rules

Gender	Age	Lipstick
Female	27	?
Male	55	?
Female	47	?
Male	39	?
Female	27	?
Male	19	?

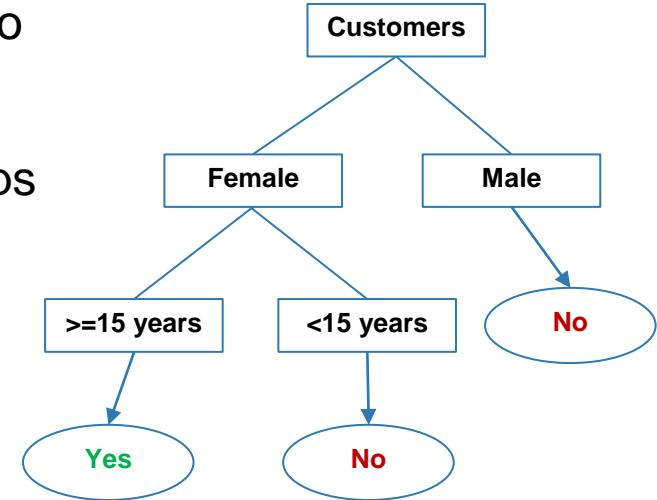
Apply
Scoring

Gender	Age	Lipstick
Female	27	P Yes
Male	55	P No
Female	47	P Yes
Male	39	P No
Female	27	P Yes
Male	19	P No

If
Gender = Female
and
Age ≥ 15 then
Purchase lipstick = YES

Decision (classification) trees

- A tree can be "learned" by splitting the source set into subsets based on an attribute value test
- Tree partitions samples into mutually exclusive groups by selecting the best splitting attribute, one group for each terminal node
- The process is repeated recursively for each derived subset, until the stopping criteria is reached



- Works with both interval and categorical variables
- No need to normalize the data
- Intuitive if-then rules are easy to extract and apply
- Best applied to binary outcomes

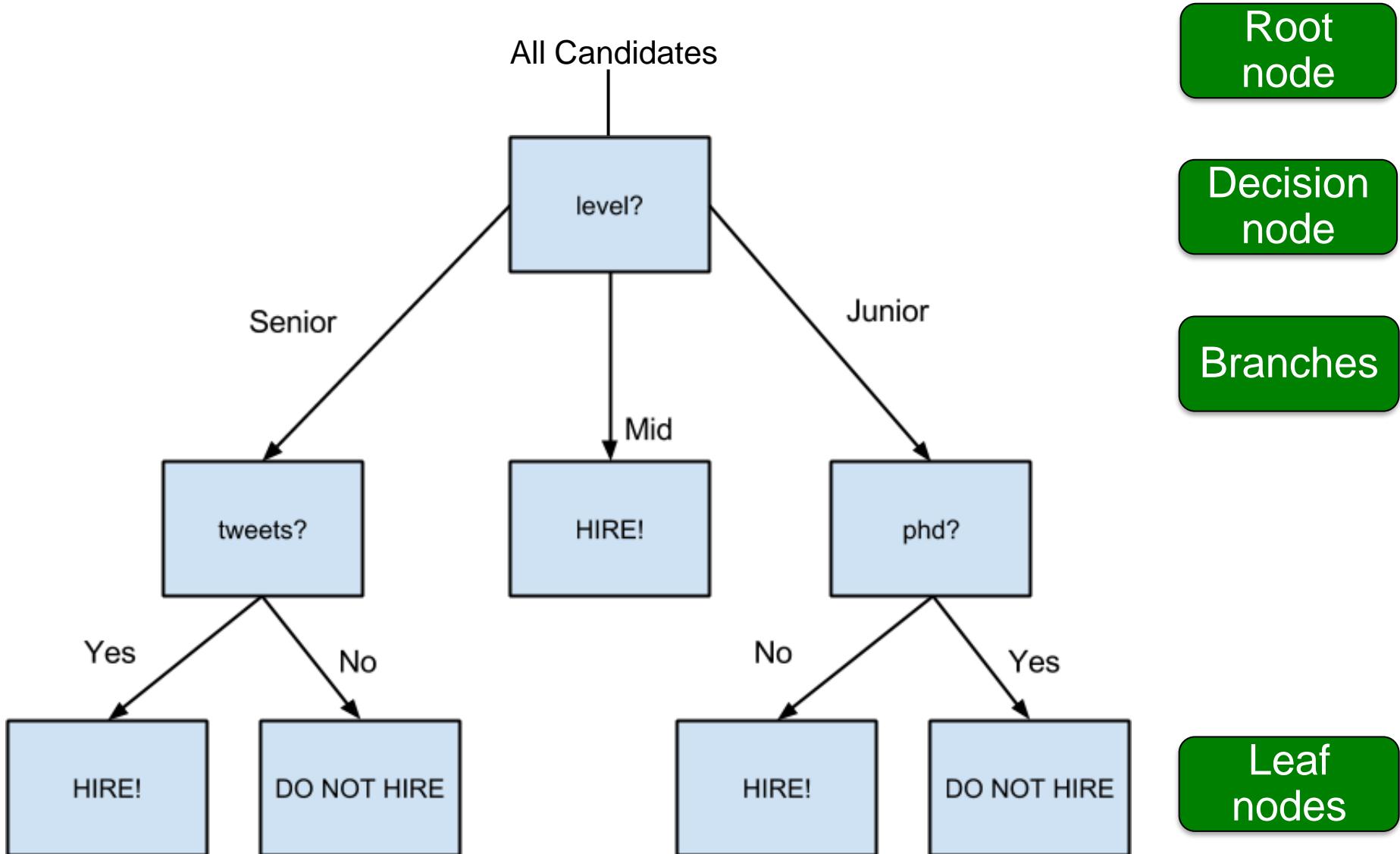
- Decision trees can be used to support multiple modeling objectives
 - Customer segmentation
 - Investment / portfolio decisions
 - Issuing a credit card or loan
 - Medical patient / disease classification

Decision (classification) trees

```
inputs = [
    ({'level':'Senior', 'lang':'Java', 'tweets':'no', 'phd':'no'},      False),
    ({'level':'Senior', 'lang':'Java', 'tweets':'no', 'phd':'yes'},       False),
    ({'level':'Mid', 'lang':'Python', 'tweets':'no', 'phd':'no'},        True),
    ({'level':'Junior', 'lang':'Python', 'tweets':'no', 'phd':'no'},      True),
    ({'level':'Junior', 'lang':'R', 'tweets':'yes', 'phd':'no'},         True),
    ({'level':'Junior', 'lang':'R', 'tweets':'yes', 'phd':'yes'},        False),
    ({'level':'Mid', 'lang':'R', 'tweets':'yes', 'phd':'yes'},          True),
    ({'level':'Senior', 'lang':'Python', 'tweets':'no', 'phd':'no'},      False),
    ({'level':'Senior', 'lang':'R', 'tweets':'yes', 'phd':'no'},         True),
    ({'level':'Junior', 'lang':'Python', 'tweets':'yes', 'phd':'no'},      True),
    ({'level':'Senior', 'lang':'Python', 'tweets':'yes', 'phd':'yes'},   True),
    ({'level':'Mid', 'lang':'Python', 'tweets':'no', 'phd':'yes'},       True),
    ({'level':'Mid', 'lang':'Java', 'tweets':'yes', 'phd':'no'},        True),
    ({'level':'Junior', 'lang':'Python', 'tweets':'no', 'phd':'yes'},  False)
]
```

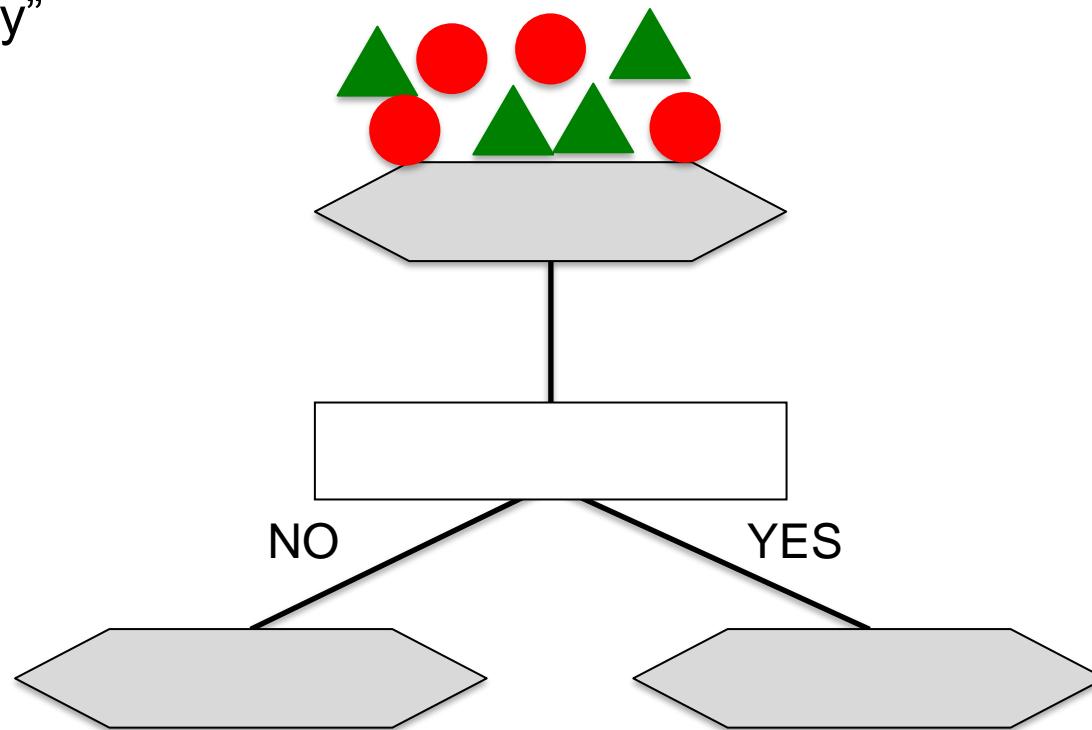
```
('level',
  {'Junior': ('phd', {'no': True, 'yes': False}),
   'Mid': True,
   'Senior': ('tweets', {'no': False, 'yes': True}))
```

Decision (classification) trees



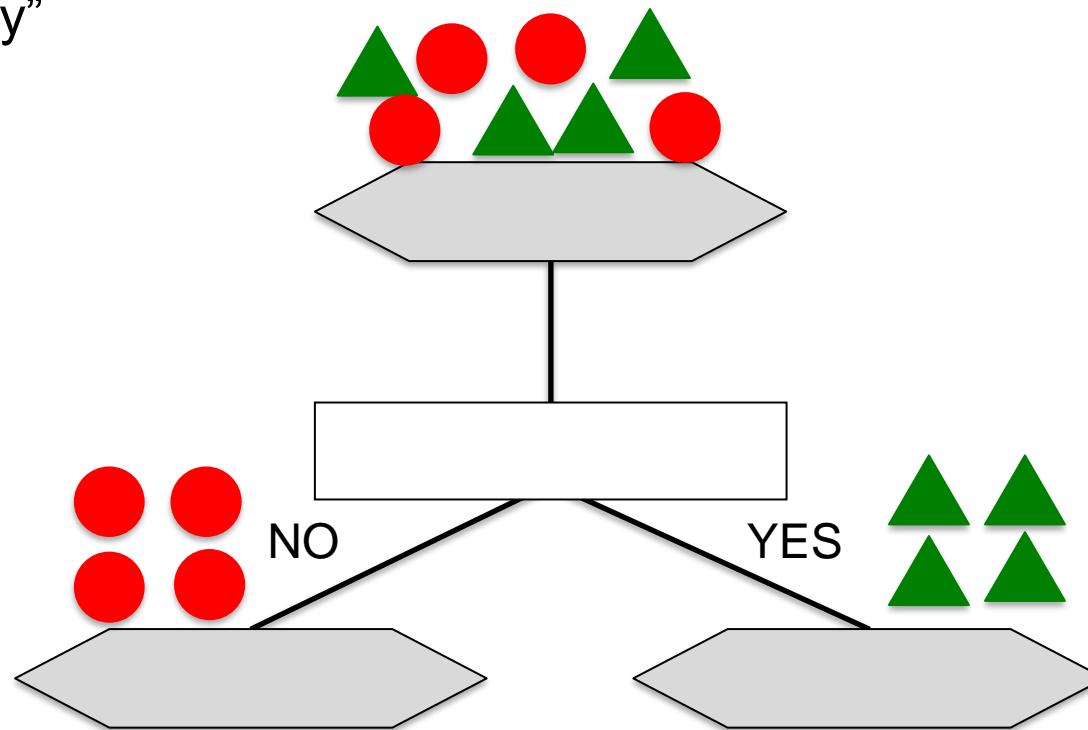
Understanding decision trees

- Decision trees are built using recursive partitioning to classify the data
- The algorithm chooses the most predictive feature to split the data on
- “Predictiveness” is based on decrease in entropy (gain in information) or “impurity”



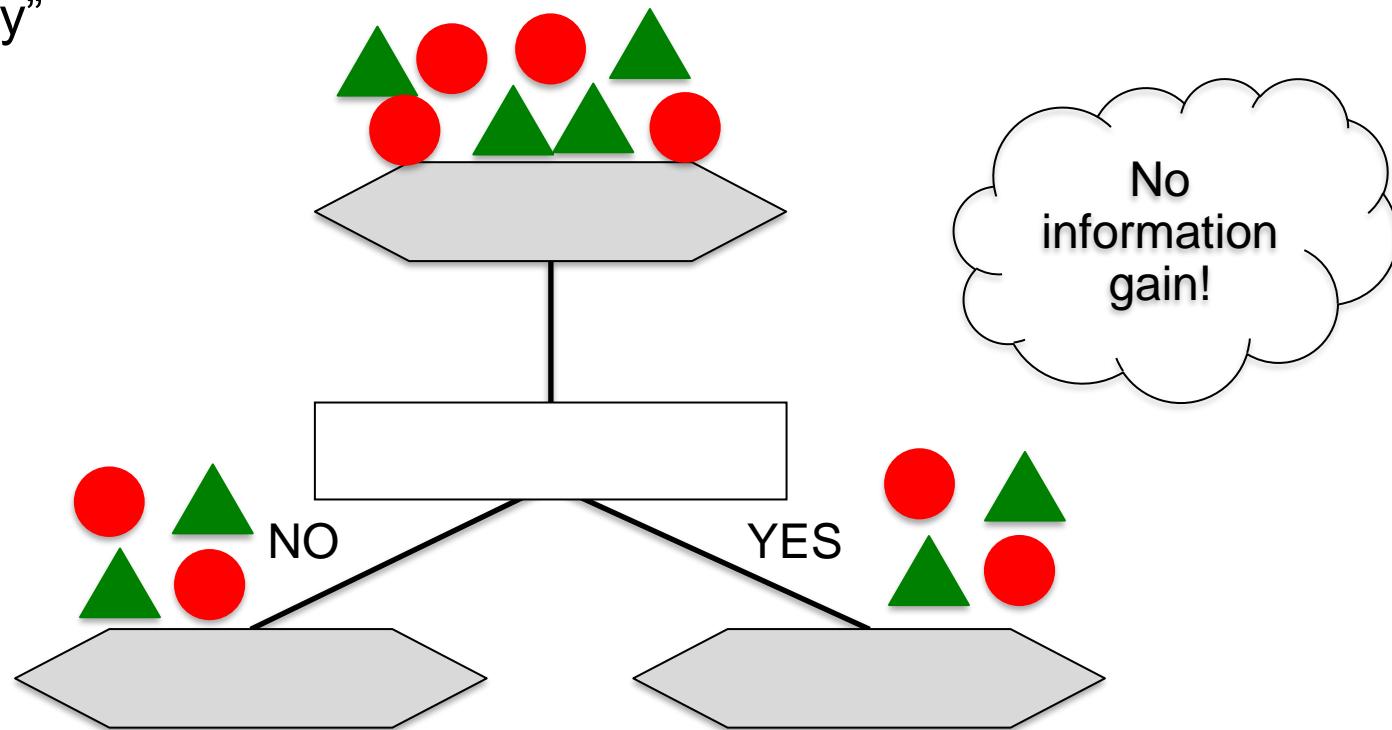
Understanding decision trees

- Decision trees are built using recursive partitioning to classify the data
- The algorithm chooses the most predictive feature to split the data on
- “Predictiveness” is based on decrease in entropy (gain in information) or “impurity”



Understanding decision trees

- Decision trees are built using recursive partitioning to classify the data
- The algorithm chooses the most predictive feature to split the data on
- “Predictiveness” is based on decrease in entropy (gain in information) or “impurity”



Understanding decision trees

- Root node partitions the data using the feature that provides the most information gain

- Information gain tells us how important a given attribute of the feature vectors is:

$$\text{Information Gain} = \text{entropy}(\text{parent}) - \text{average entropy}(\text{children})$$

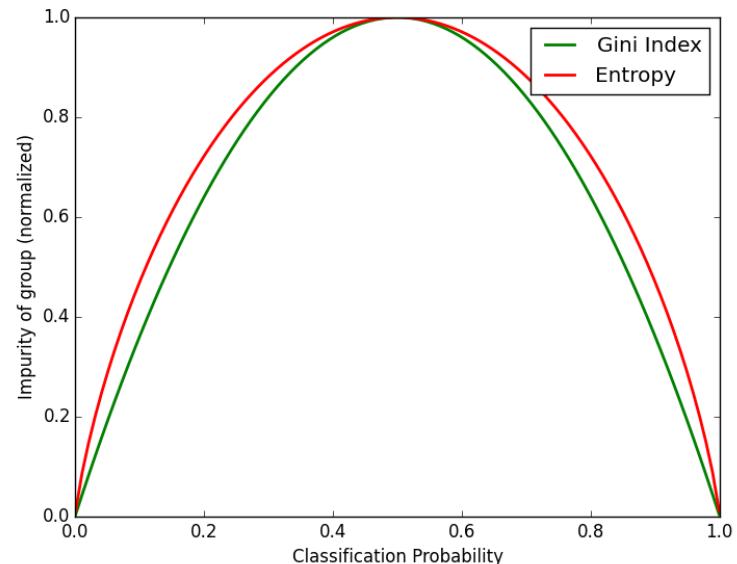
- Entropy is a common measure of target class impurity (i is each of the target classes, p_i is proportion of the number of elements in class 0 or 1):

$$\text{Entropy} = \sum_i -p_i \log_2 p_i$$

- Gini Index is another measure of impurity:

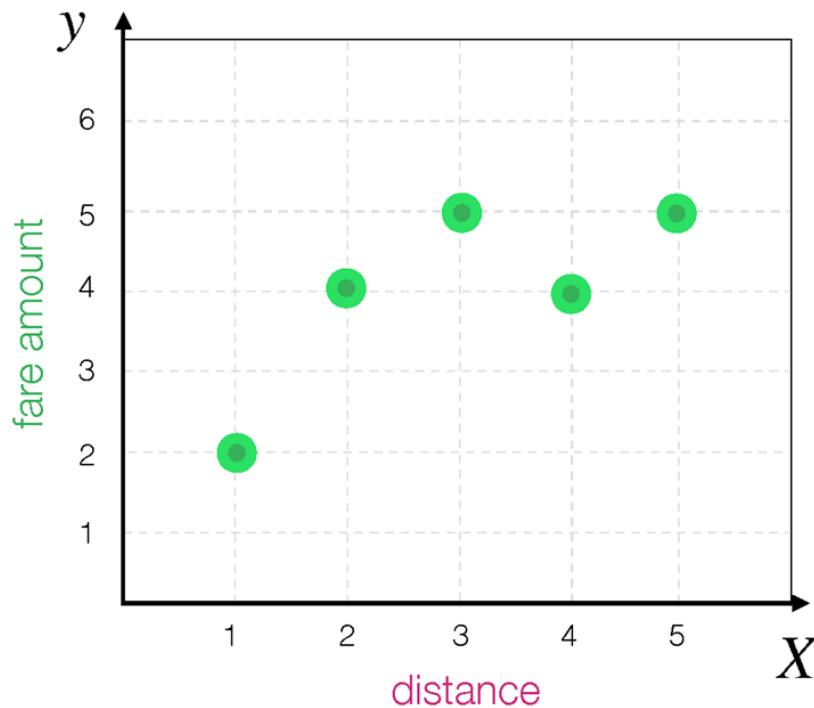
$$\text{Gini} = 1 - \sum_i p_i^2$$

Gini impurity is computationally faster as it doesn't require calculating logarithmic functions, though in reality which of the two methods is used rarely makes too much of a difference.



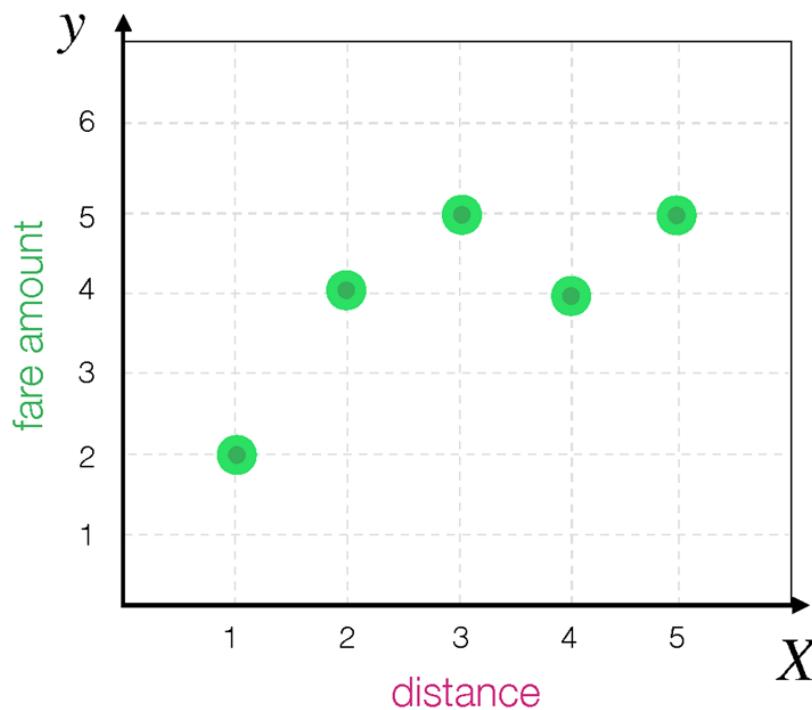
Decision Tree Algorithm

By asking a simple **question** about value of **independent variable** it tries to predict a value of **dependent variable**



Decision Tree Algorithm

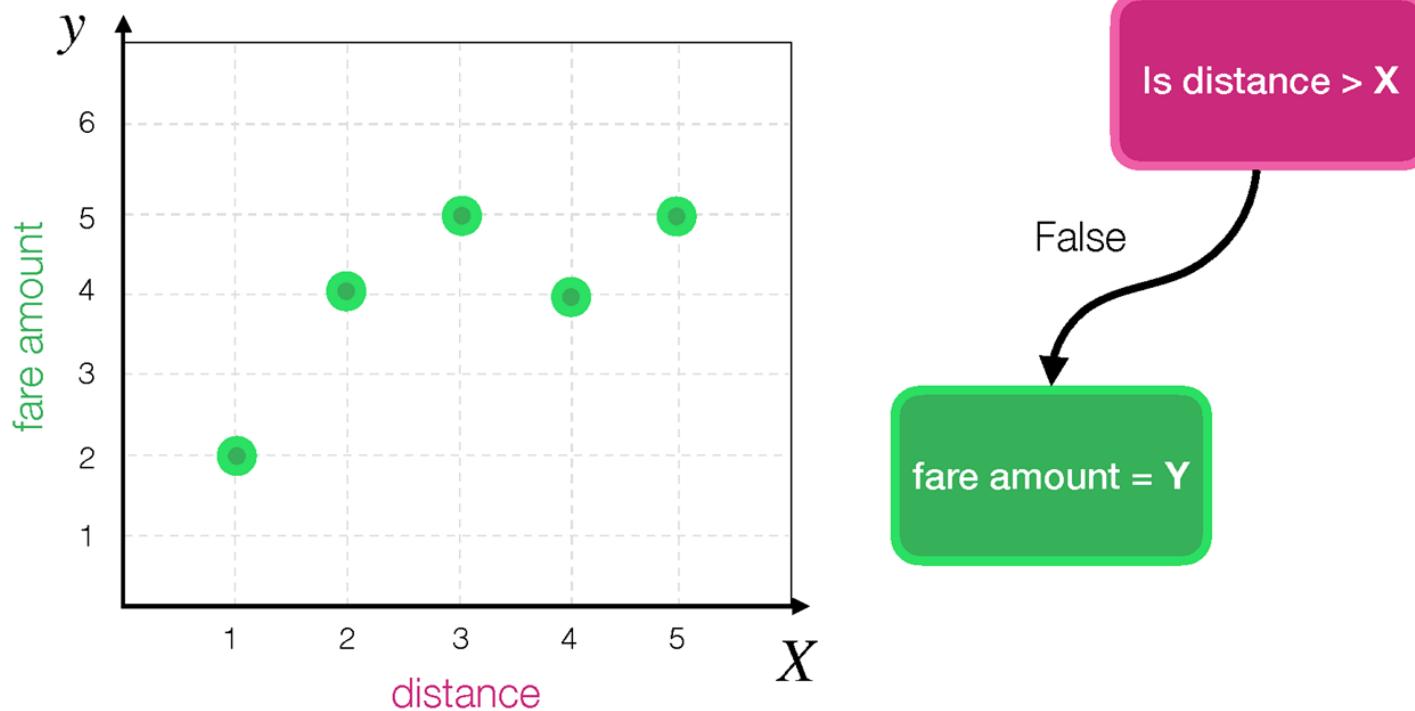
By asking a simple **question** about value of **independent variable** it tries to predict a value of **dependent variable**



Is distance > X

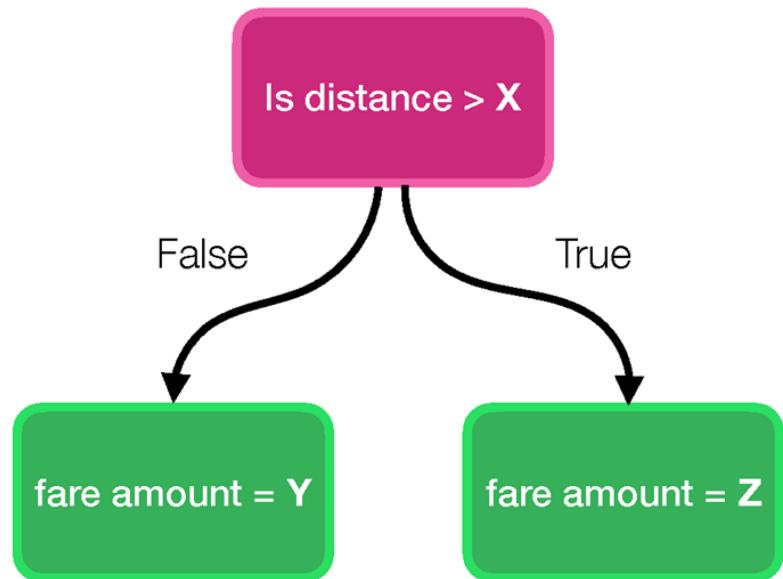
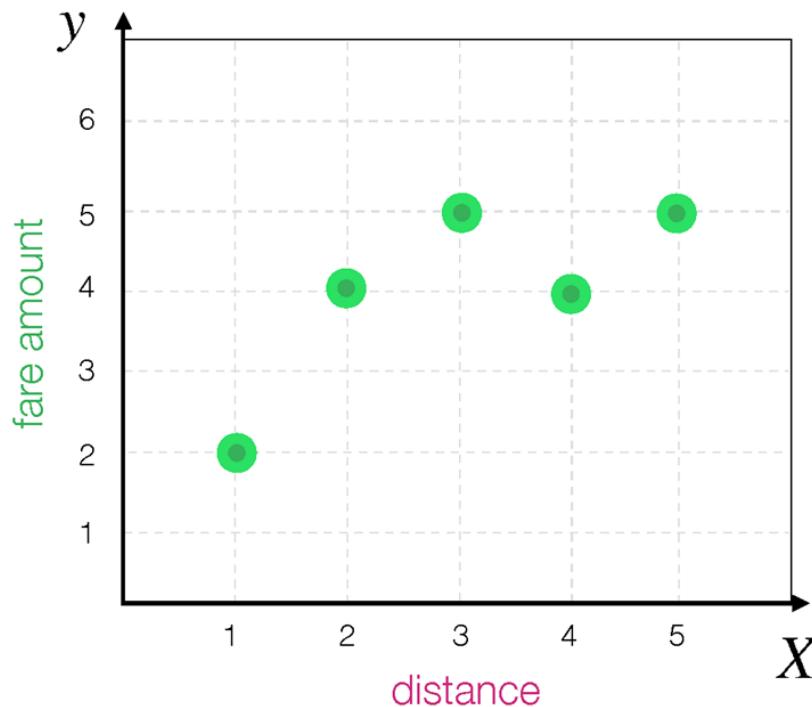
Decision Tree Algorithm

By asking a simple **question** about value of **independent variable** it tries to predict a value of **dependent variable**



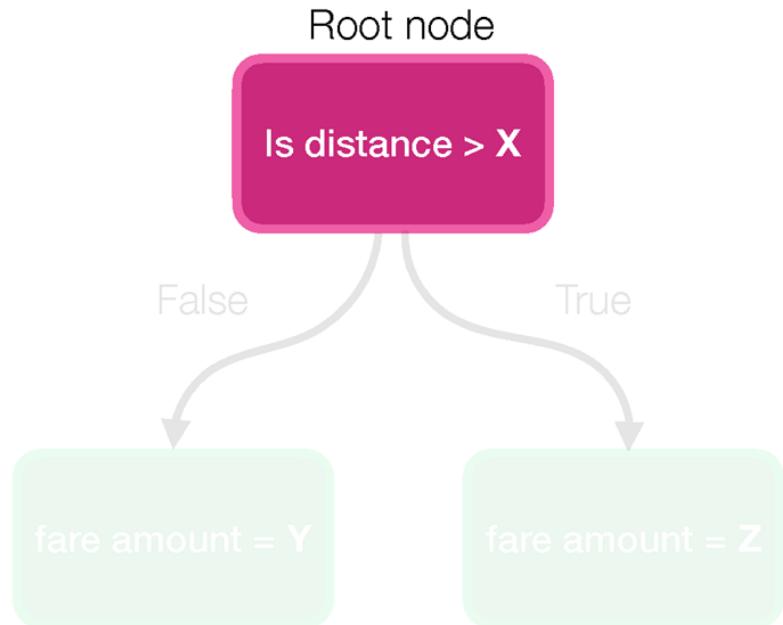
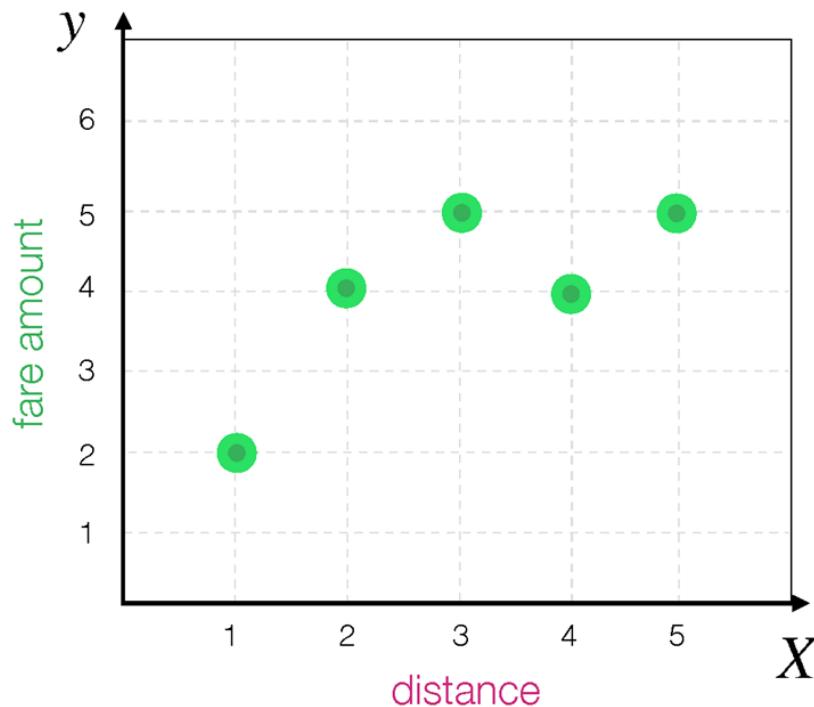
Decision Tree Algorithm

By asking a simple **question** about value of **independent variable** it tries to predict a value of **dependent variable**



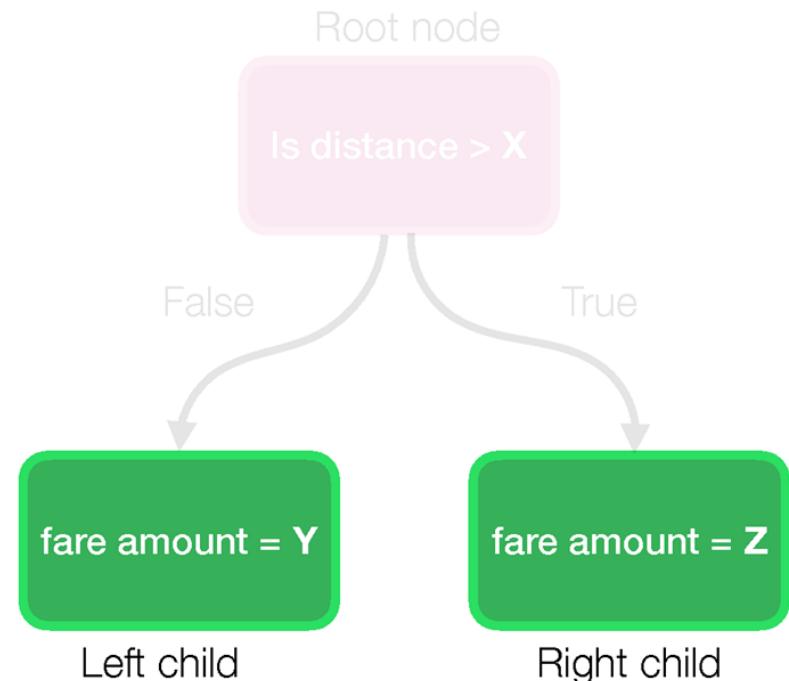
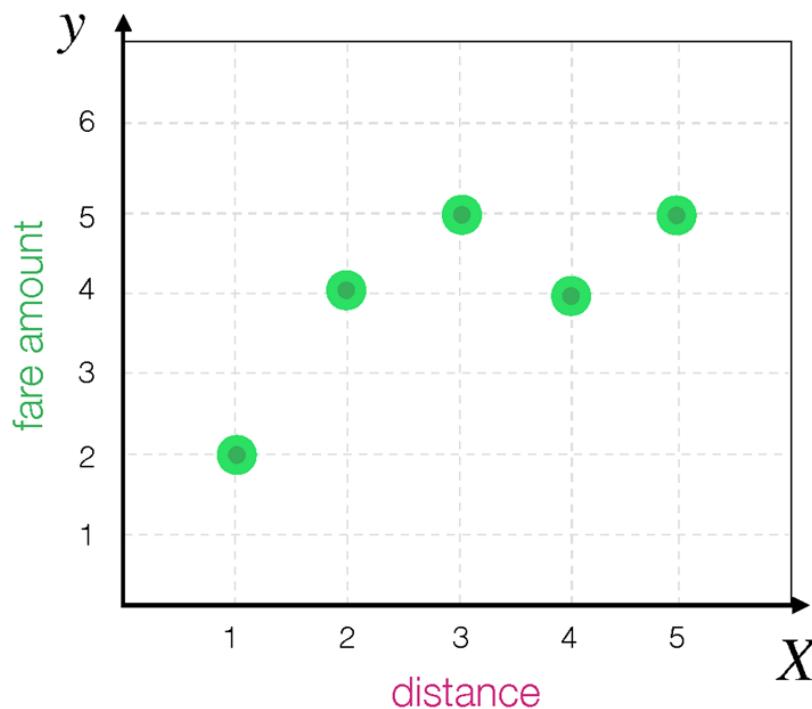
Decision Tree Algorithm

By asking a simple **question** about value of **independent variable** it tries to predict a value of **dependent variable**



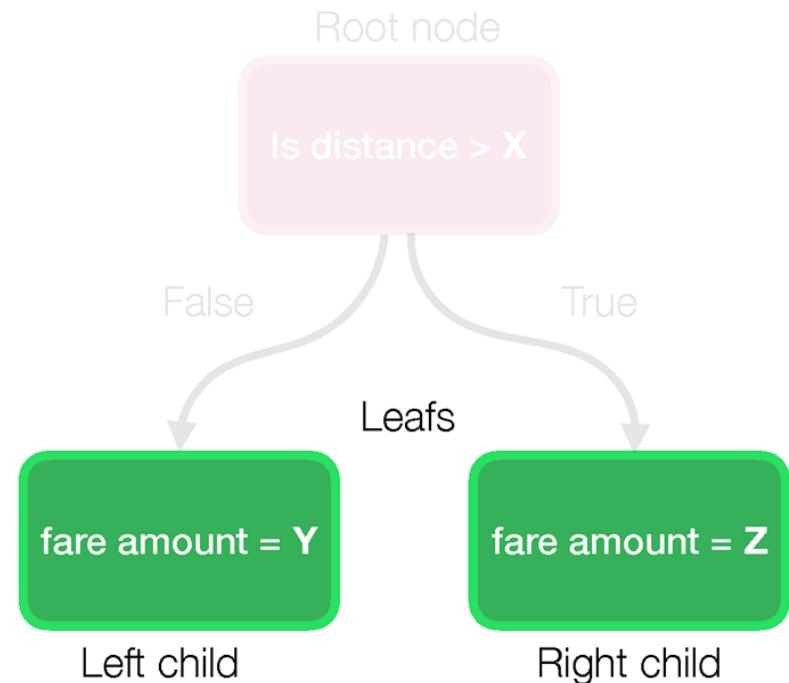
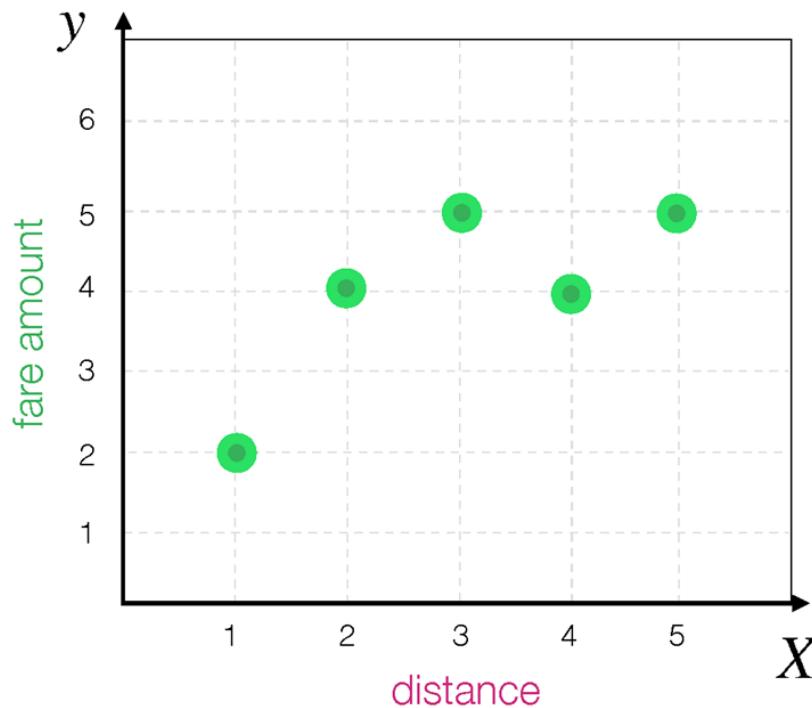
Decision Tree Algorithm

By asking a simple **question** about value of **independent variable** it tries to predict a value of **dependent variable**



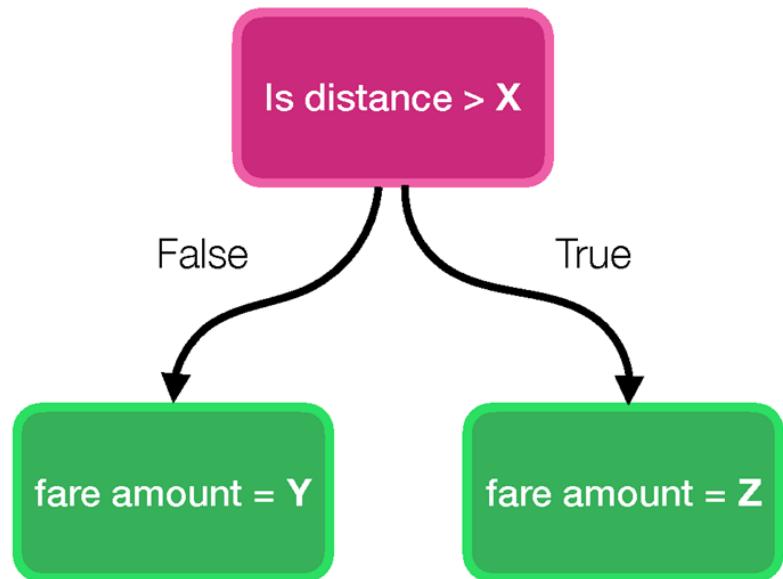
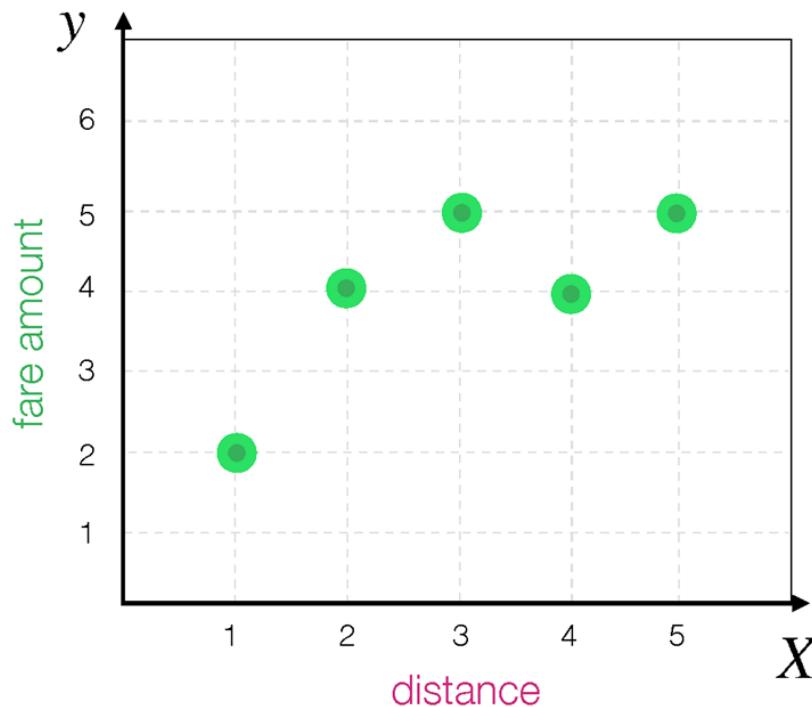
Decision Tree Algorithm

By asking a simple **question** about value of **independent variable** it tries to predict a value of **dependent variable**



Decision Tree Algorithm

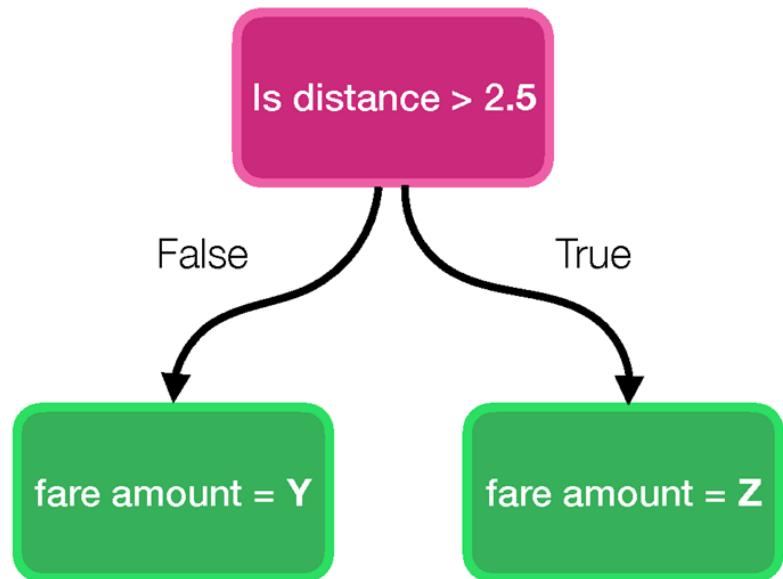
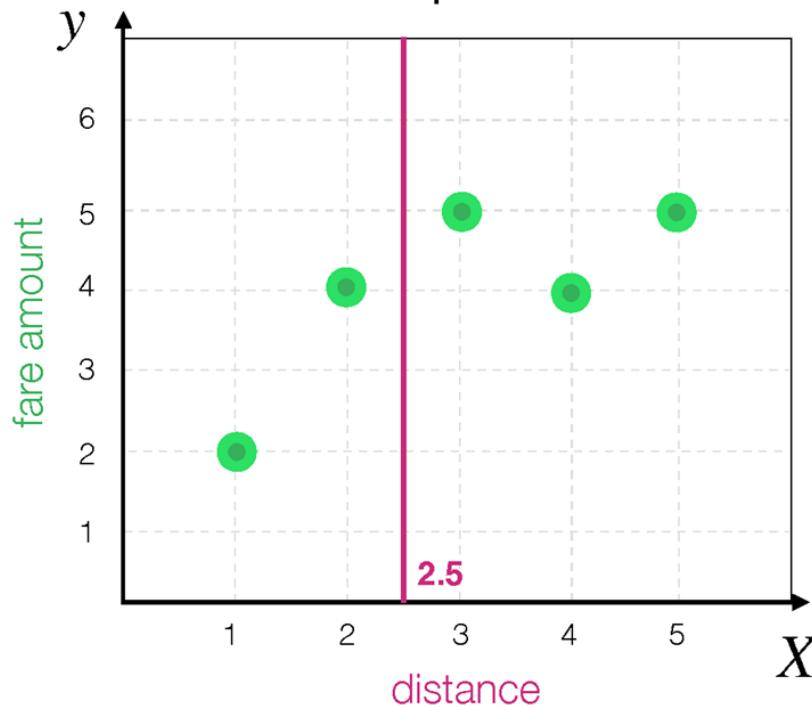
By asking a simple **question** about value of **independent variable** it tries to predict a value of **dependent variable**



Decision Tree Algorithm

Here, **X** may correspond to any **vertical** line.

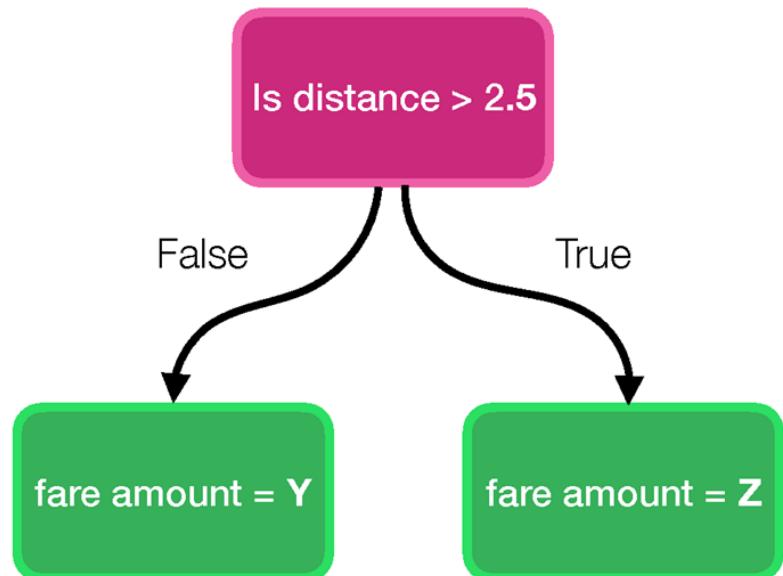
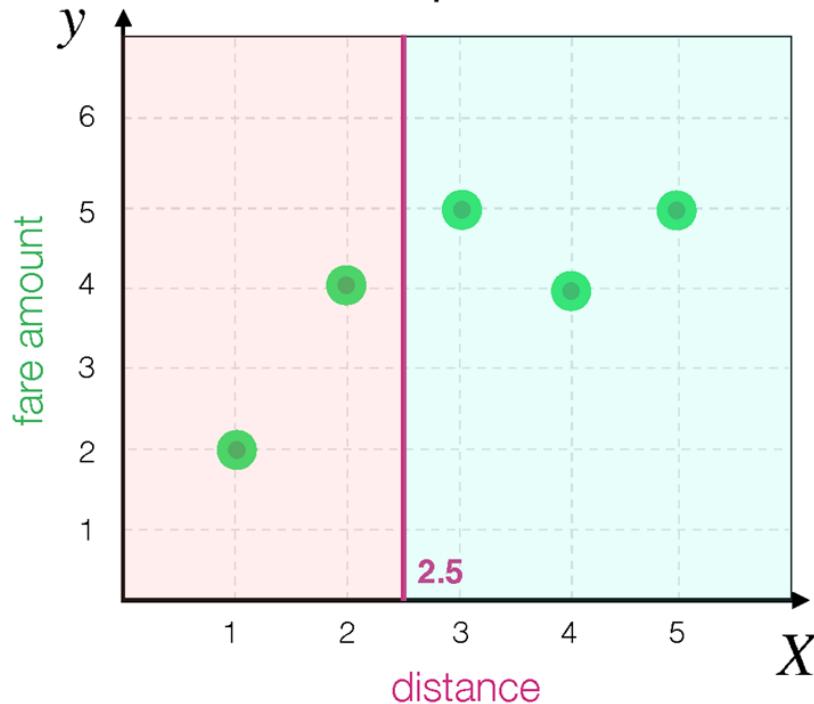
For example if $X = 2.5$:



Decision Tree Algorithm

Here, **X** may correspond to any **vertical** line.

For example if $X = 2.5$:

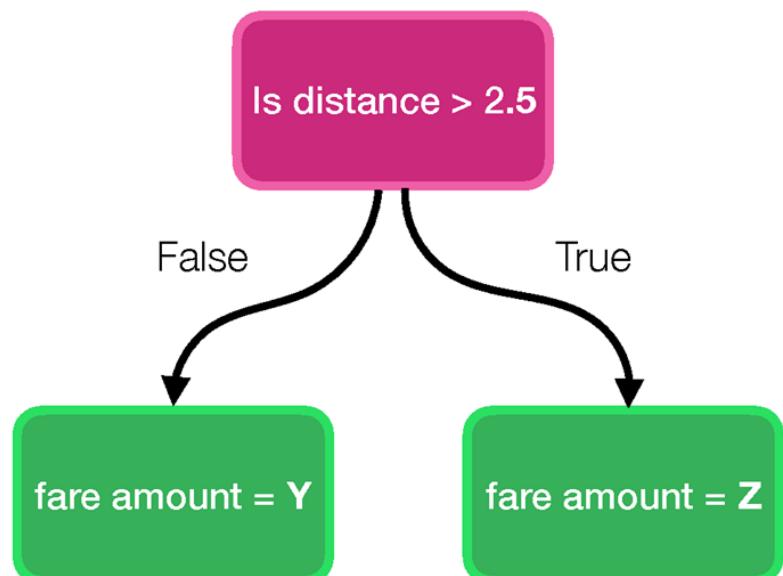
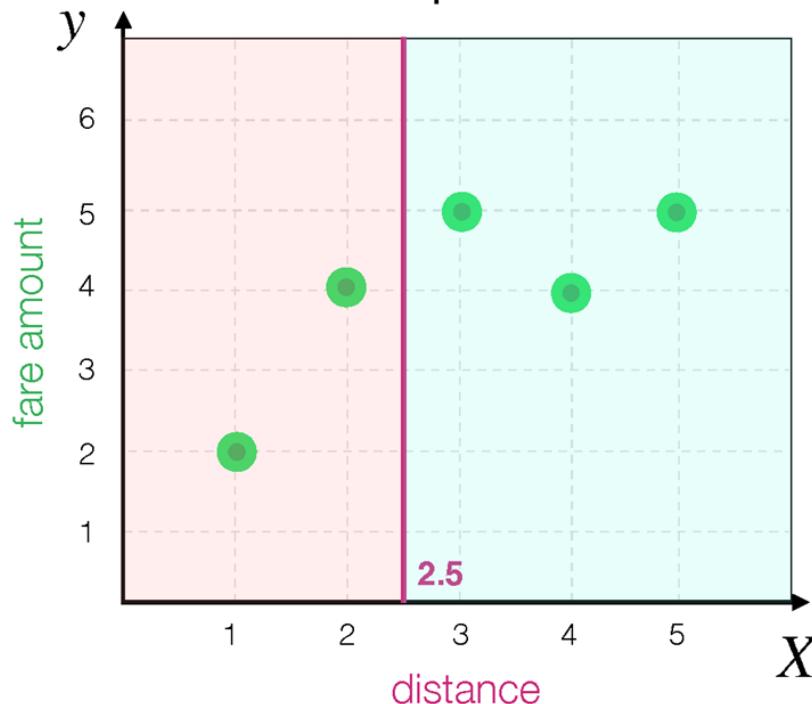


What are most reasonable values
for **Y** and **Z**?

Decision Tree Algorithm

Here, \mathbf{X} may correspond to any **vertical** line.

For example if $X = 2.5$:

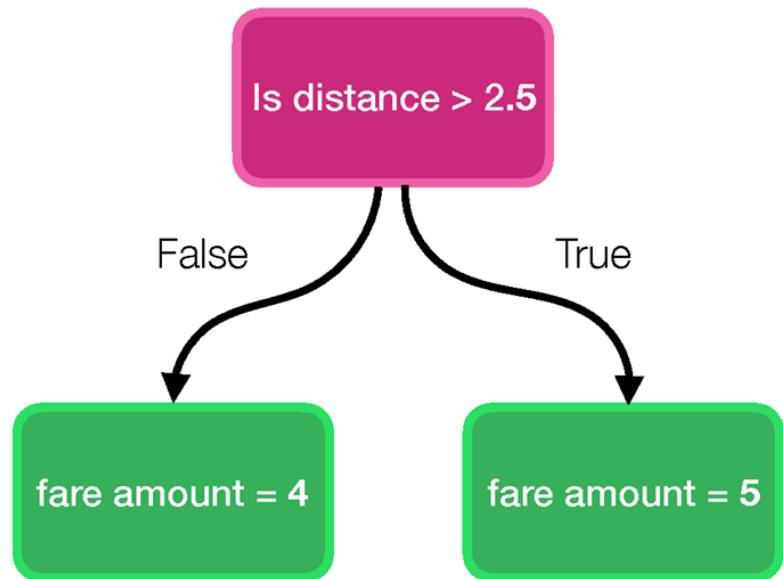
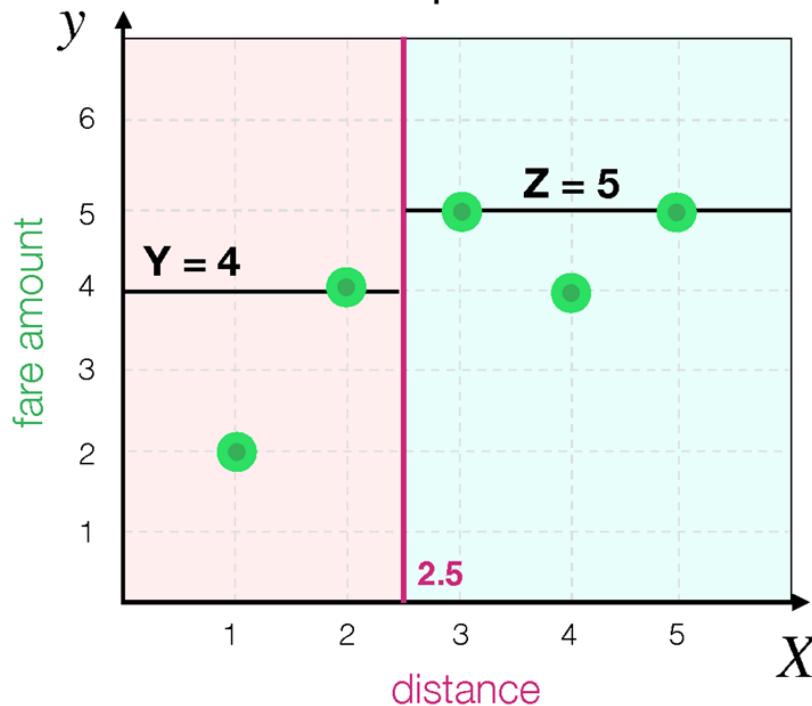


What are most reasonable values for **Y** and **Z** (**that minimise total MSE**)?

Decision Tree Algorithm

What would be MSE if $\mathbf{Y = 4}$ and $\mathbf{Z = 5}$?

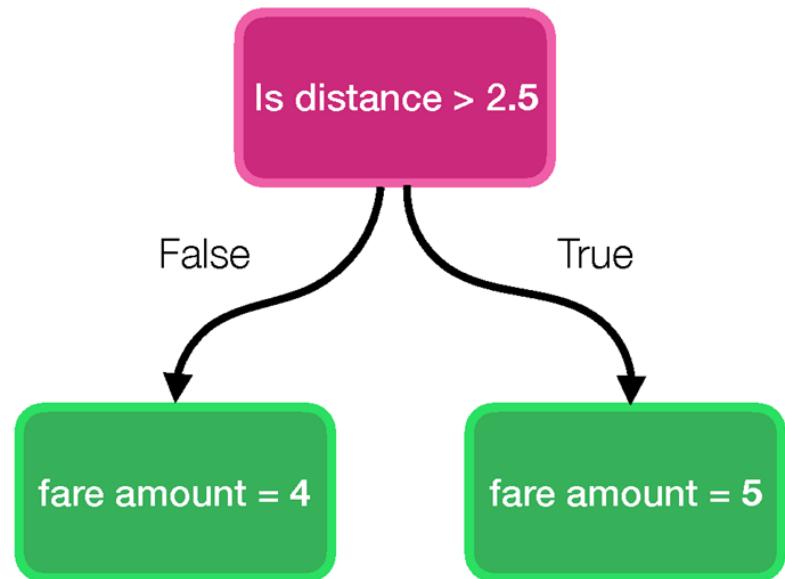
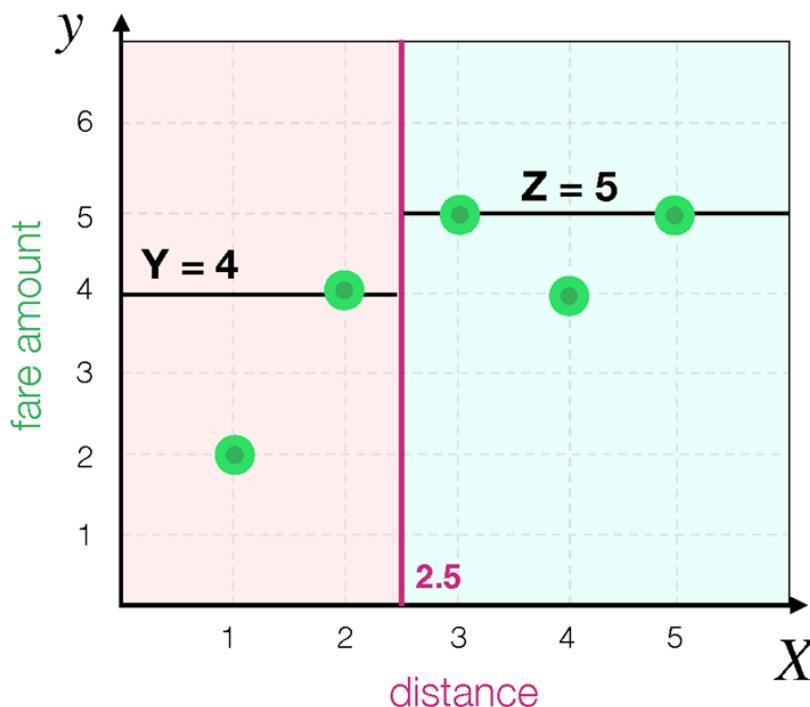
For example if $X = 2.5$:



What are most reasonable values for \mathbf{Y} and \mathbf{Z} (that minimise total MSE)?

Decision Tree Algorithm

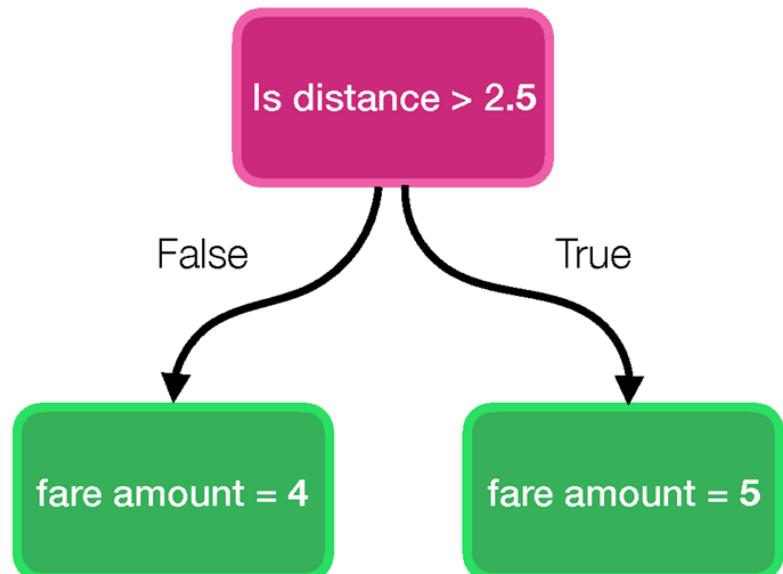
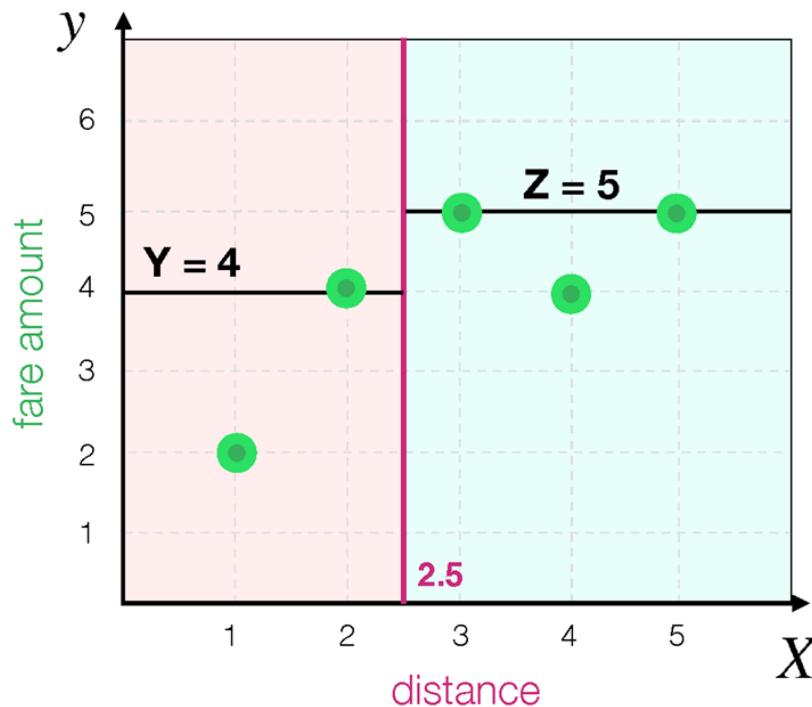
$$MSE = \frac{1}{n} \sum_{i=1}^n (\text{real value}_i - \text{predicted value}_i)^2$$



What are most reasonable values for \mathbf{Y} and \mathbf{Z} (that minimise total MSE)?

Decision Tree Algorithm

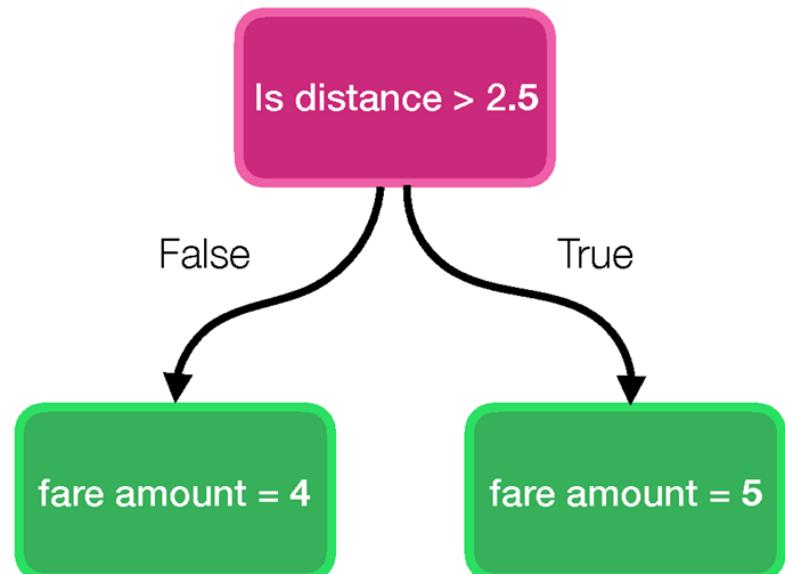
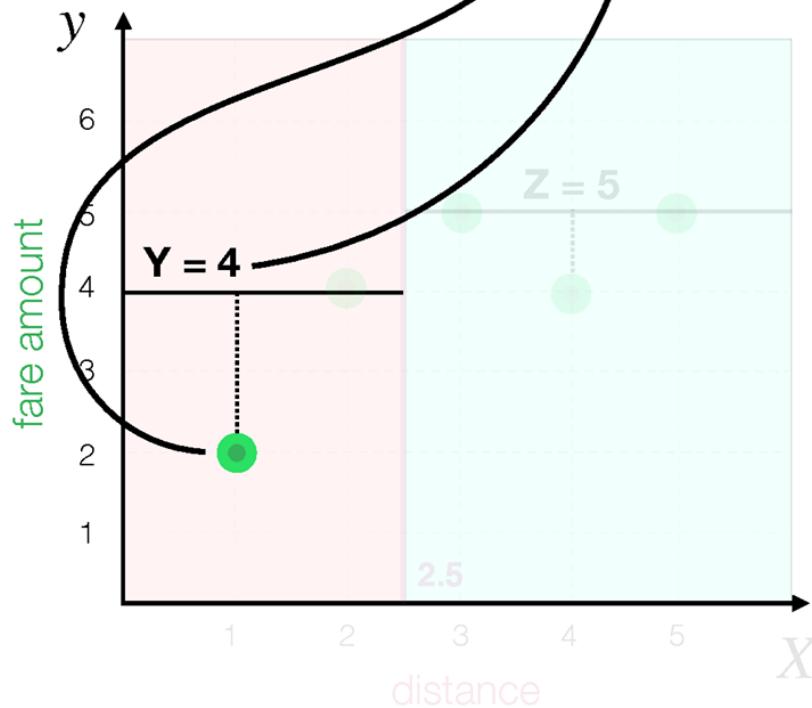
$$MSE = \frac{1}{n} \sum_{i=1}^n (\textcolor{teal}{y}_i - \hat{y}_i)^2 = \frac{(y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + (y_3 - \hat{y}_3)^2 + (y_4 - \hat{y}_4)^2 + (y_5 - \hat{y}_5)^2}{5}$$



What are most reasonable values for $\textcolor{teal}{Y}$ and Z (that minimise total MSE)?

Decision Tree Algorithm

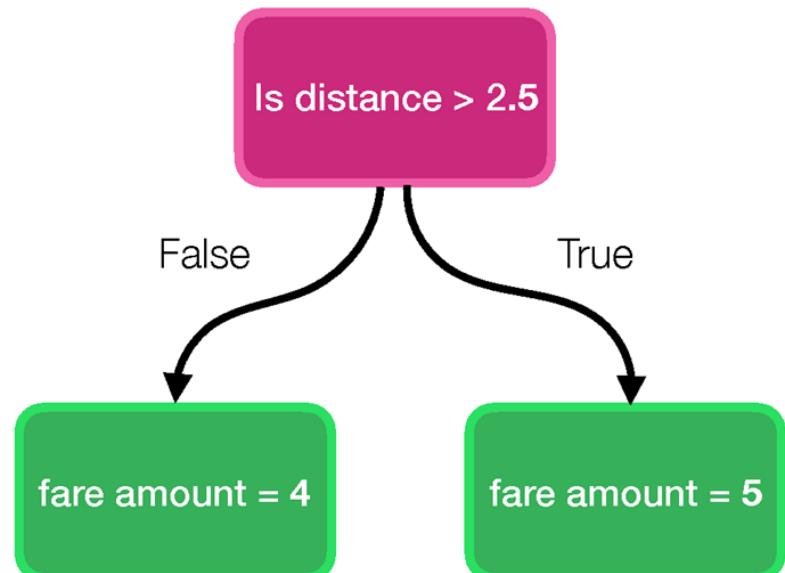
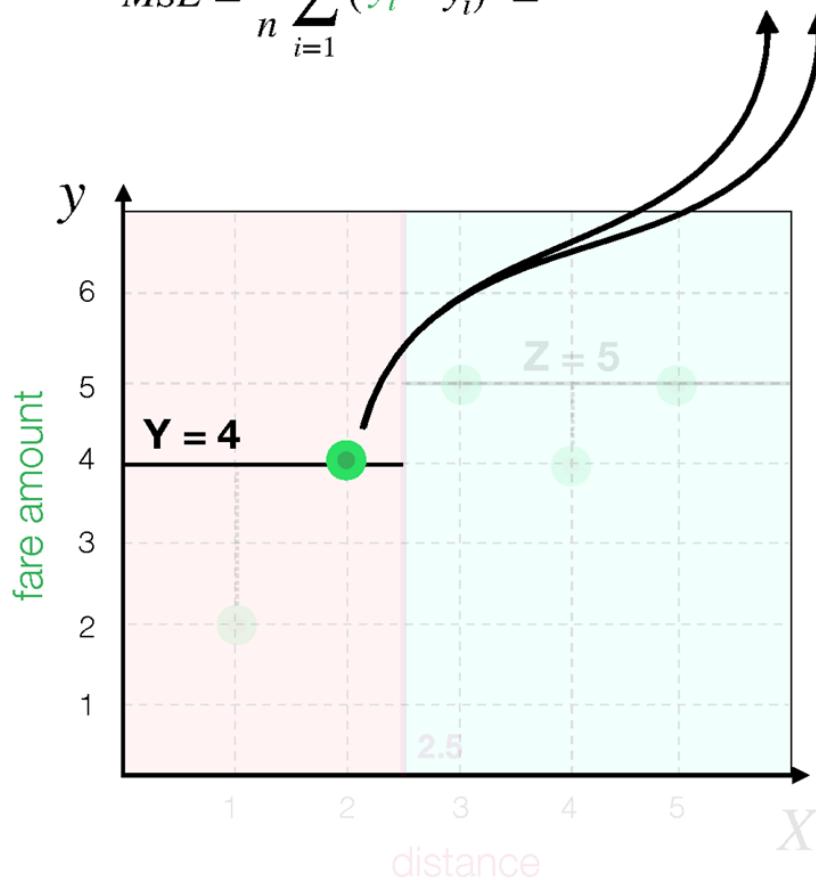
$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{(2 - 4)^2 + (y_2 - \hat{y}_2)^2 + (y_3 - \hat{y}_3)^2 + (y_4 - \hat{y}_4)^2 + (y_5 - \hat{y}_5)^2}{5}$$



What are most reasonable values for **Y** and **Z** (that minimise total MSE)?

Decision Tree Algorithm

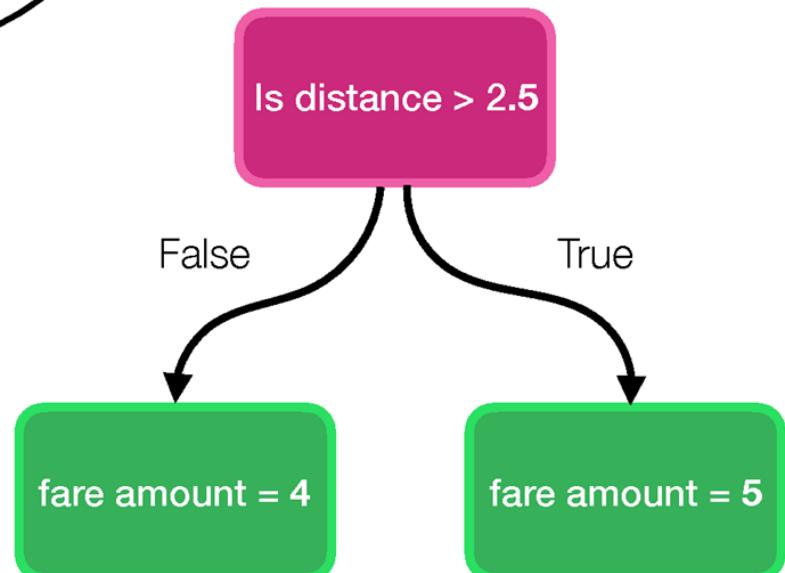
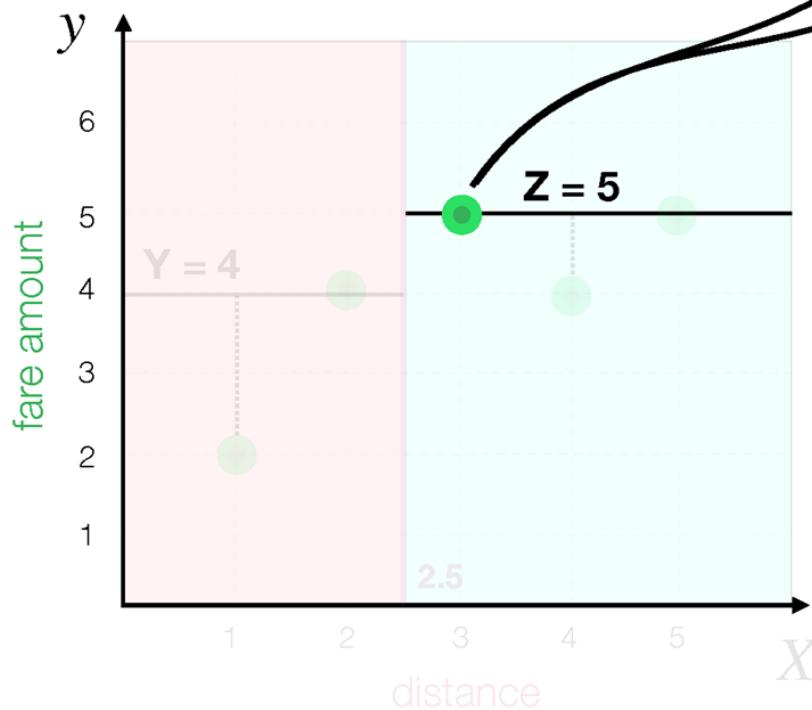
$$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = \frac{(2 - 4)^2 + (4 - 4)^2 + (y_3 - \hat{y}_3)^2 + (y_4 - \hat{y}_4)^2 + (y_5 - \hat{y}_5)^2}{5}$$



What are most reasonable values for **Y** and **Z** (that minimise total MSE)?

Decision Tree Algorithm

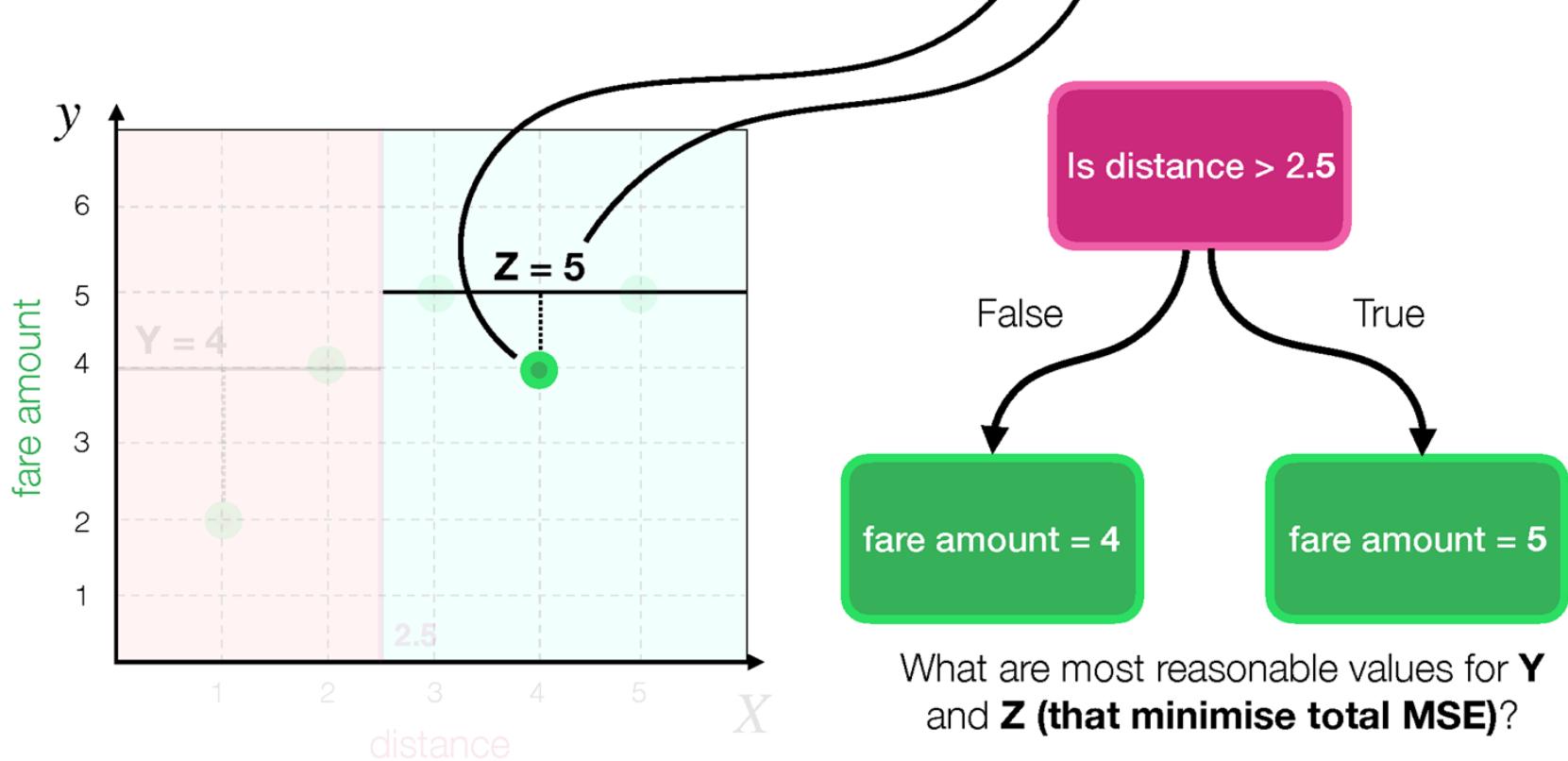
$$MSE = \frac{1}{n} \sum_{i=1}^n (\textcolor{teal}{y}_i - \hat{y}_i)^2 = \frac{(2 - 4)^2 + (4 - 4)^2 + (5 - 5)^2 + (y_4 - \hat{y}_4)^2 + (y_5 - \hat{y}_5)^2}{n}$$



What are most reasonable values for $\textcolor{teal}{Y}$ and $\textcolor{violet}{Z}$ (that minimise total MSE)?

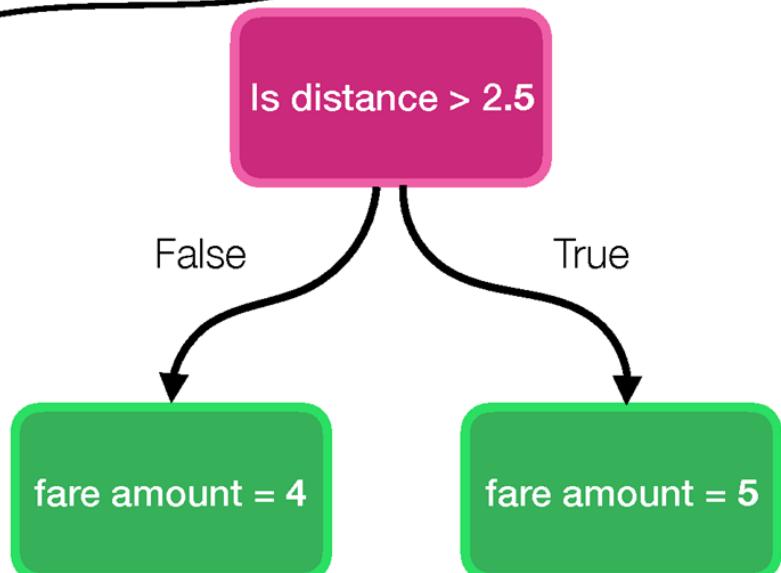
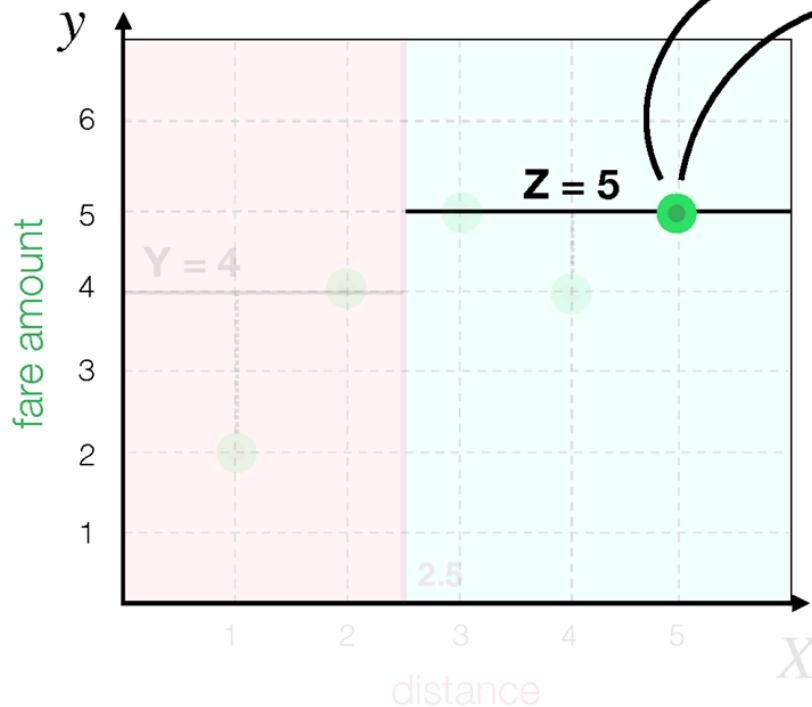
Decision Tree Algorithm

$$MSE = \frac{1}{n} \sum_{i=1}^n (\textcolor{teal}{y}_i - \hat{y}_i)^2 = \frac{(2-4)^2 + (4-4)^2 + (5-5)^2 + (4-5)^2 + (y_5 - \hat{y}_5)^2}{5}$$



Decision Tree Algorithm

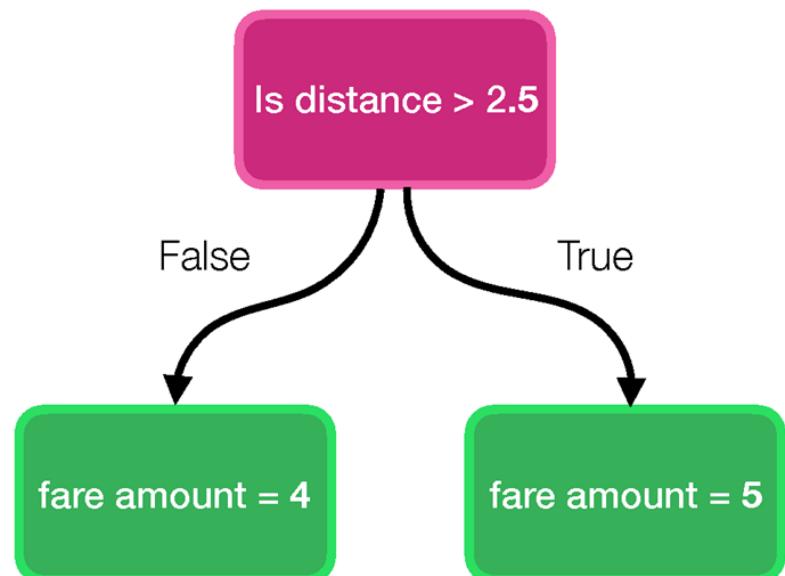
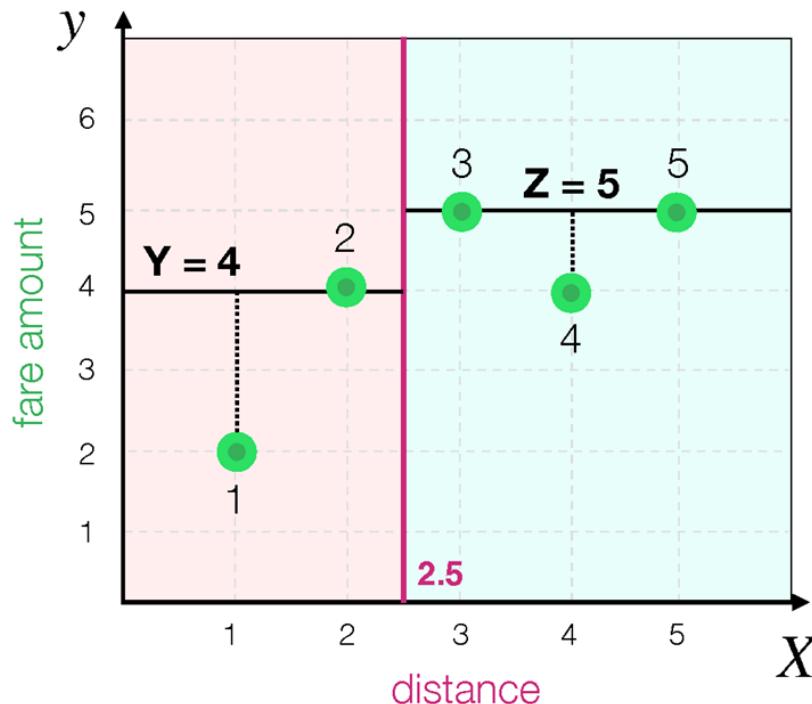
$$MSE = \frac{1}{n} \sum_{i=1}^n (\textcolor{teal}{y}_i - \hat{y}_i)^2 = \frac{(2-4)^2 + (4-4)^2 + (5-5)^2 + (4-5)^2 + (5-5)^2}{5}$$



What are most reasonable values for **Y** and **Z** (that minimise total MSE)?

Decision Tree Algorithm

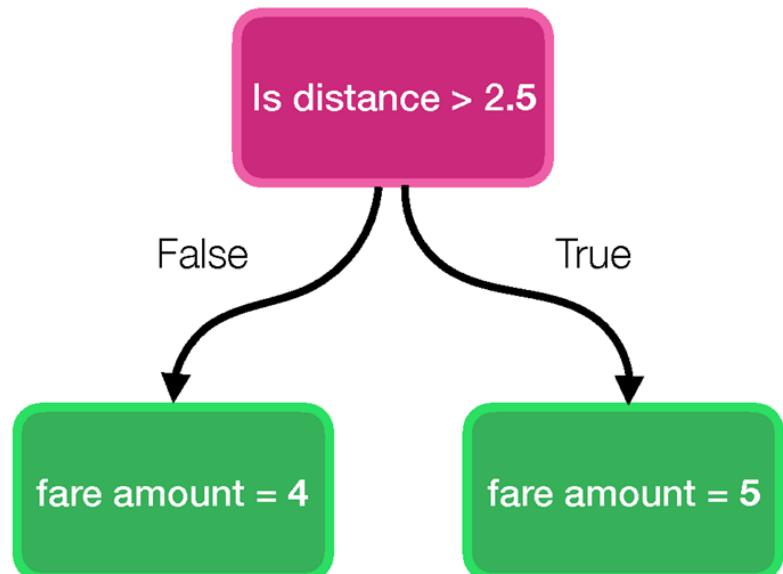
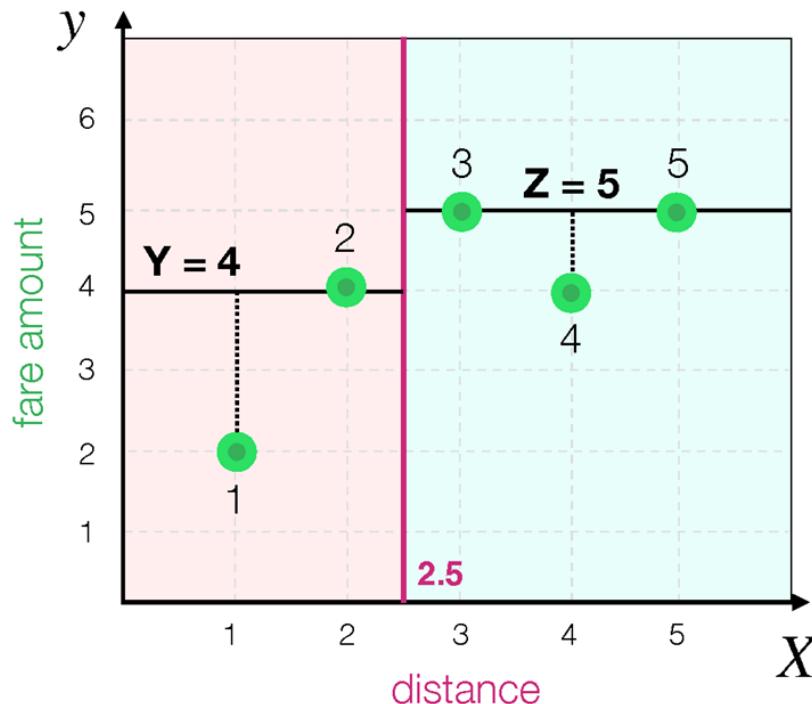
$$MSE = \frac{1}{n} \sum_{i=1}^n (\textcolor{teal}{y}_i - \hat{y}_i)^2 = \frac{(-2)^2 + (0)^2 + (0)^2 + (-1)^2 + (0)^2}{5}$$



What are most reasonable values for **Y** and **Z** (that minimise total MSE)?

Decision Tree Algorithm

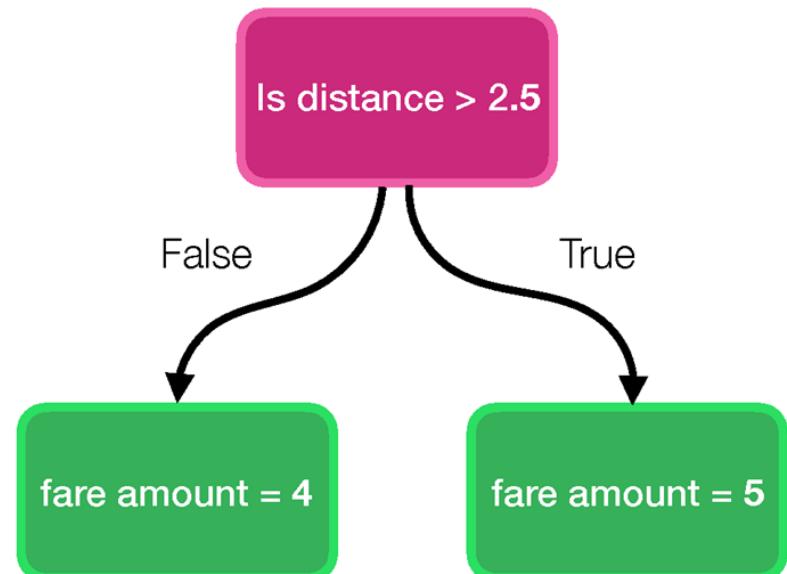
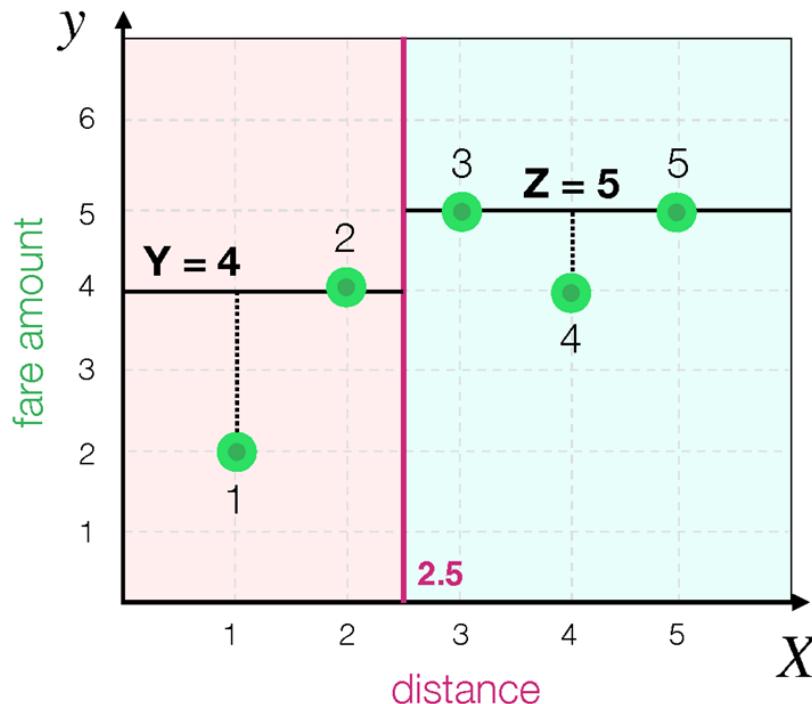
$$MSE = \frac{1}{n} \sum_{i=1}^n (\textcolor{teal}{y}_i - \hat{y}_i)^2 = \frac{4 + 0 + 0 + 1 + 0}{5}$$



What are most reasonable values for **Y** and **Z** (that minimise total MSE)?

Decision Tree Algorithm

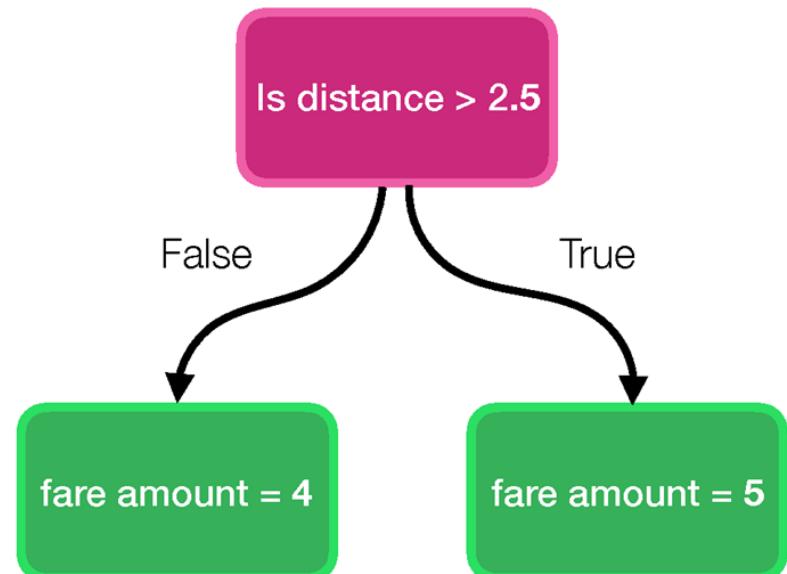
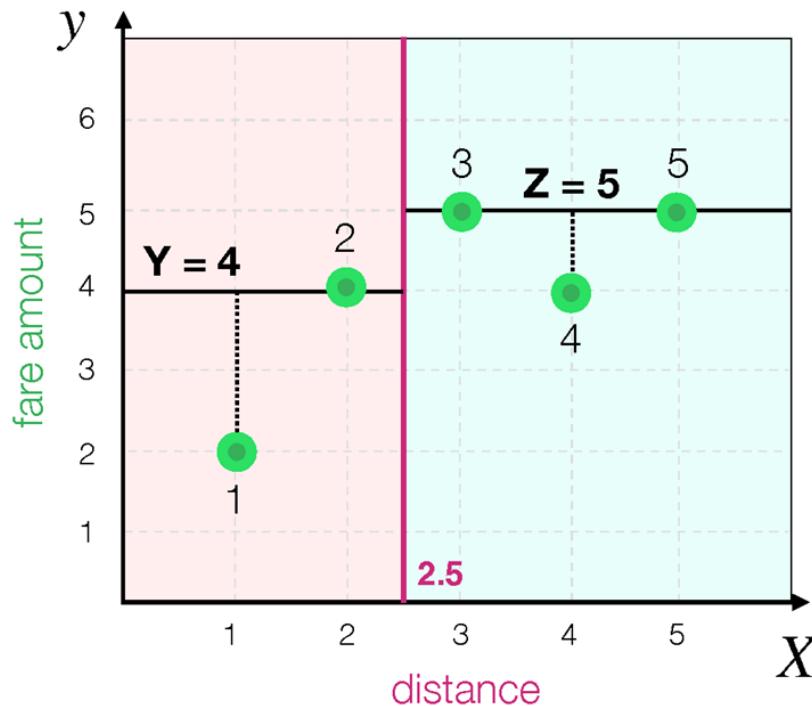
$$MSE = \frac{1}{n} \sum_{i=1}^n (\textcolor{teal}{y}_i - \hat{y}_i)^2 = \frac{5}{5}$$



What are most reasonable values for **Y** and **Z** (that minimise total MSE)?

Decision Tree Algorithm

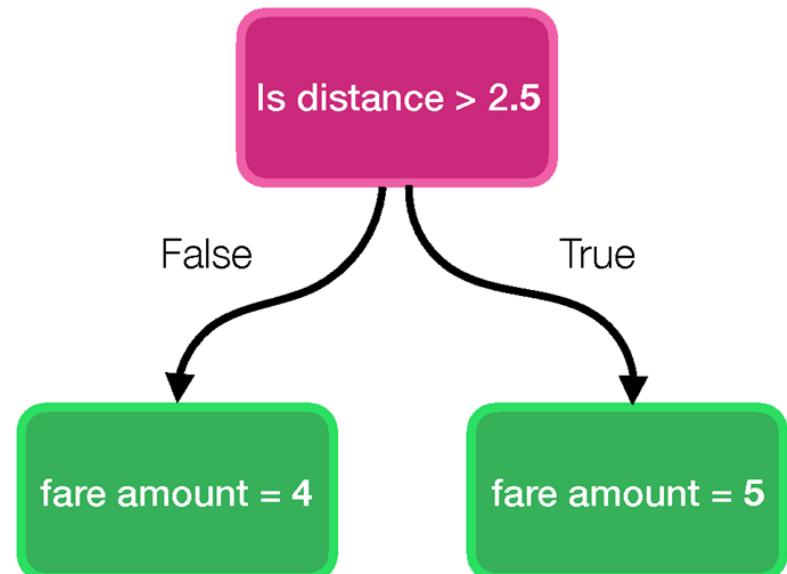
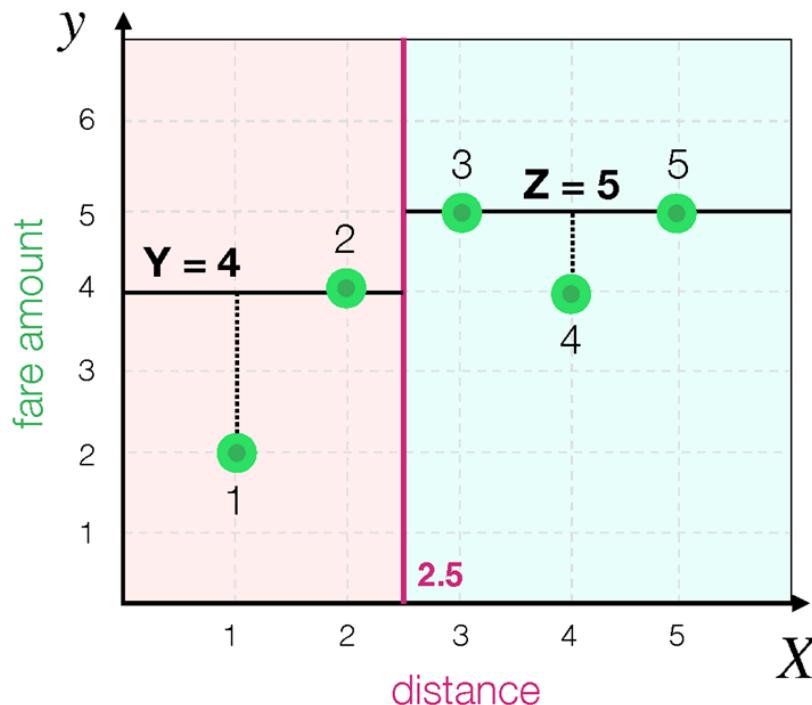
$$MSE = \frac{1}{n} \sum_{i=1}^n (\textcolor{teal}{y}_i - \hat{y}_i)^2 = 1$$



What are most reasonable values for **Y** and **Z** (that minimise total MSE)?

Decision Tree Algorithm

$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = 1$ so, if $\mathbf{X} = 2.5$, $\mathbf{Y} = 4$ and $\mathbf{Z} = 5$, MSE is 1

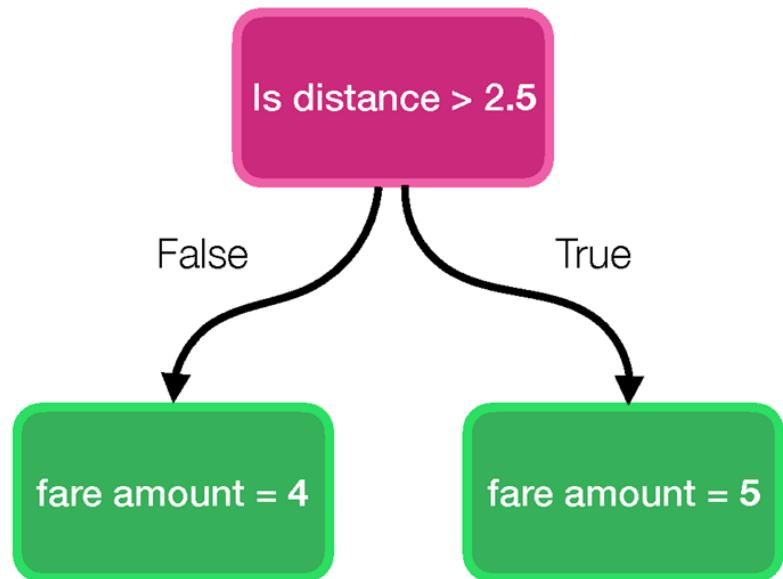
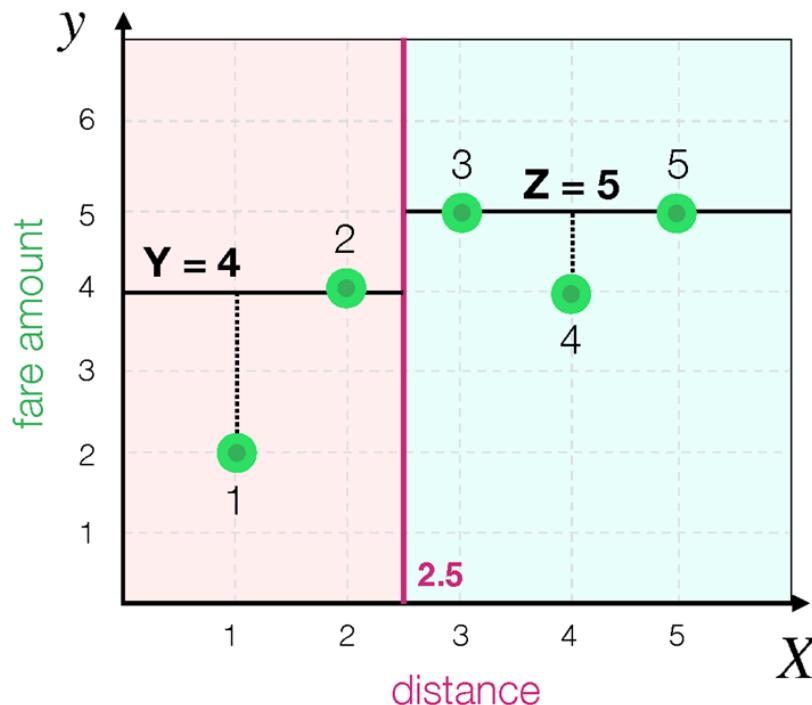


What are most reasonable values for \mathbf{Y} and \mathbf{Z} (that minimise total MSE)?

Decision Tree Algorithm

$MSE = \frac{1}{n} \sum_{i=1}^n (y_i - \hat{y}_i)^2 = 1$ so, if $\mathbf{X} = 2.5$, $\mathbf{Y} = 4$ and $\mathbf{Z} = 5$, MSE is 1

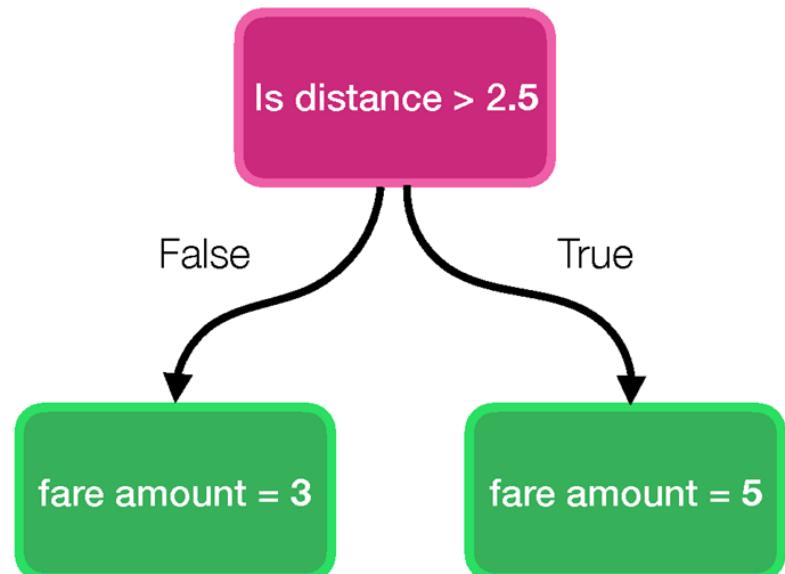
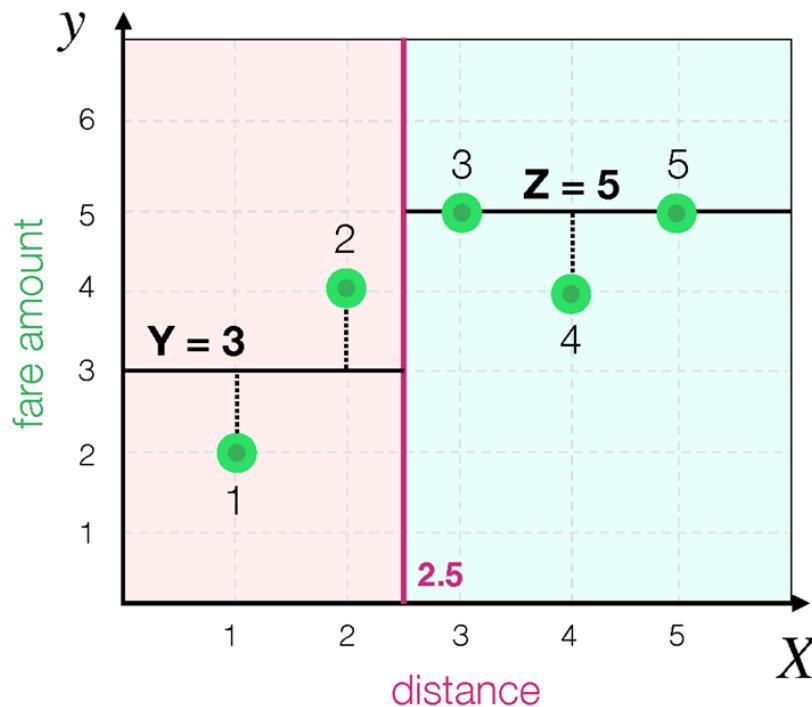
Can we find better \mathbf{Y} and \mathbf{Z} ?



What are most reasonable values for \mathbf{Y} and \mathbf{Z} (that minimise total MSE)?

Decision Tree Algorithm

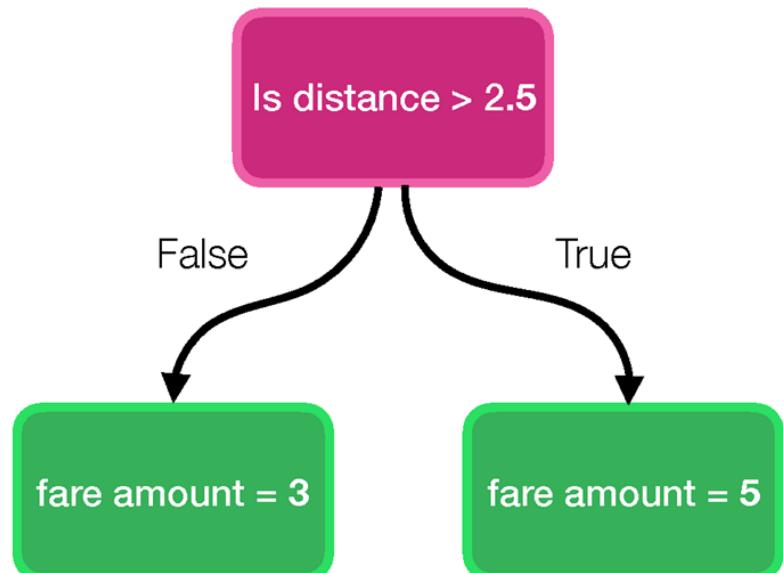
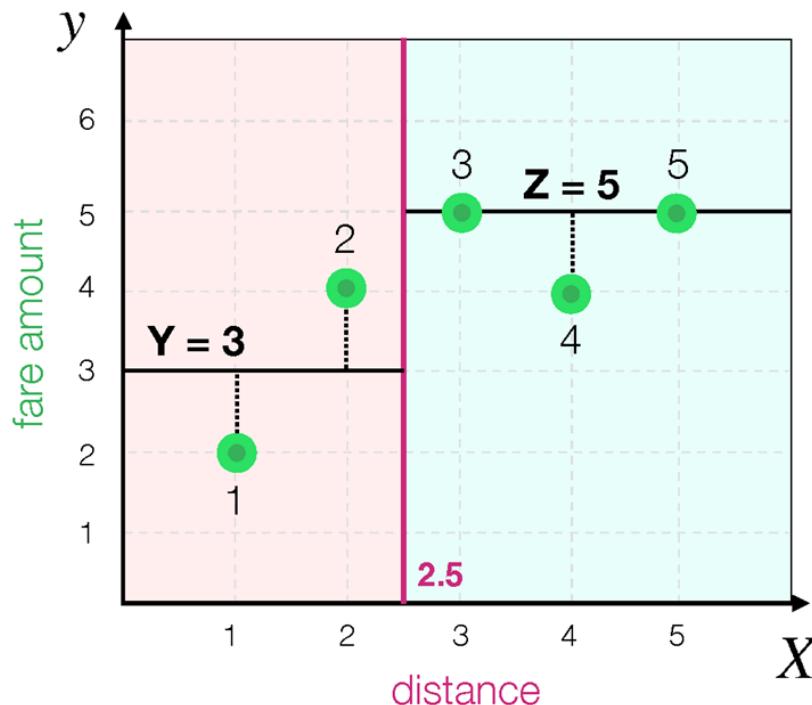
$$MSE = \frac{1}{n} \sum_{i=1}^n (\textcolor{teal}{y}_i - \hat{y}_i)^2 = \frac{(y_1 - \hat{y}_1)^2 + (y_2 - \hat{y}_2)^2 + (y_3 - \hat{y}_3)^2 + (y_4 - \hat{y}_4)^2 + (y_5 - \hat{y}_5)^2}{5}$$



What are most reasonable values for **Y** and **Z** (that minimise total MSE)?

Decision Tree Algorithm

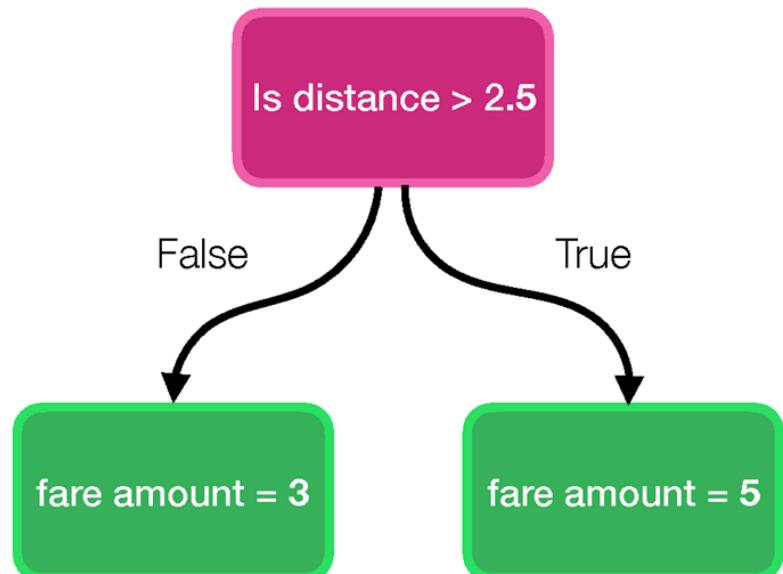
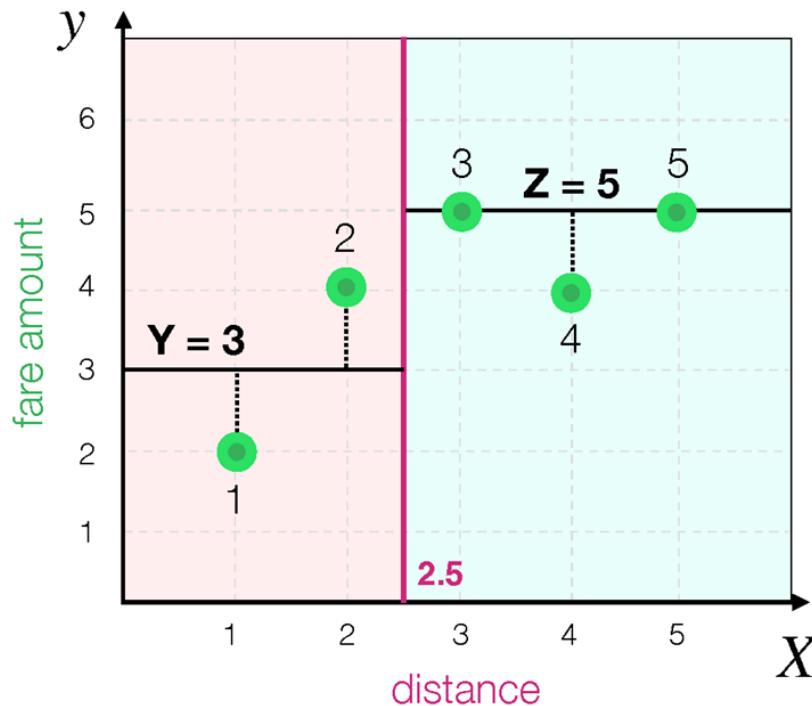
$$MSE = \frac{1}{n} \sum_{i=1}^n (\textcolor{teal}{y}_i - \hat{y}_i)^2 = \frac{(2-3)^2 + (4-3)^2 + (5-5)^2 + (4-5)^2 + (5-5)^2}{5}$$



What are most reasonable values for **Y** and **Z** (that minimise total MSE)?

Decision Tree Algorithm

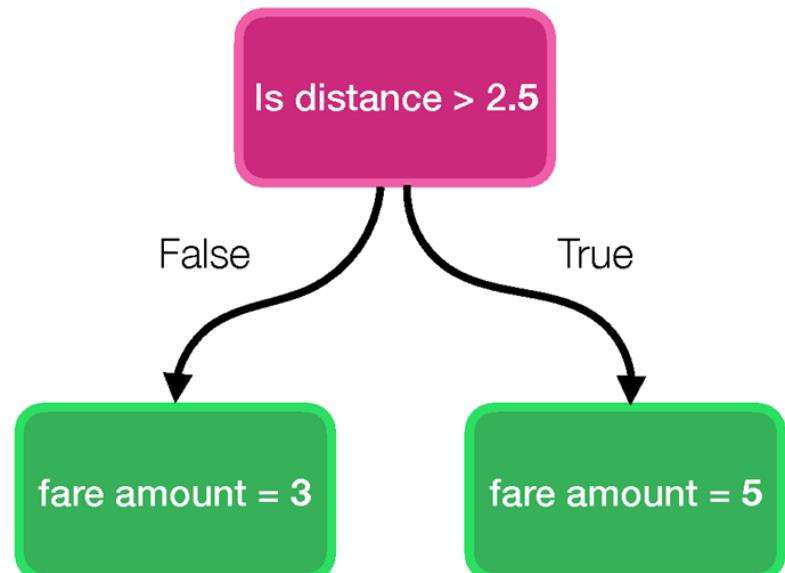
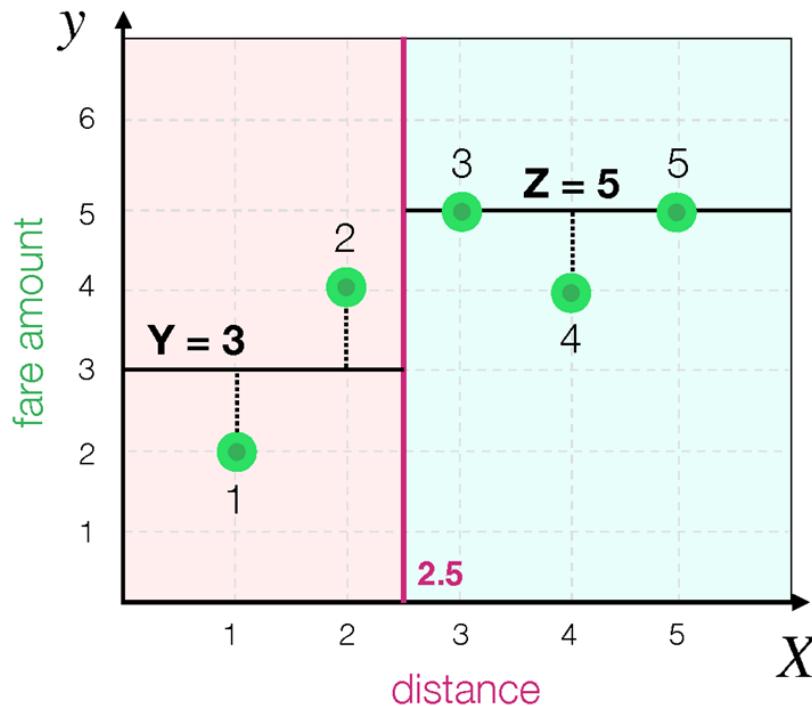
$$MSE = \frac{1}{n} \sum_{i=1}^n (\textcolor{teal}{y}_i - \hat{y}_i)^2 = \frac{1 + 1 + 0 + 1 + 0}{5}$$



What are most reasonable values for **Y** and **Z** (that minimise total MSE)?

Decision Tree Algorithm

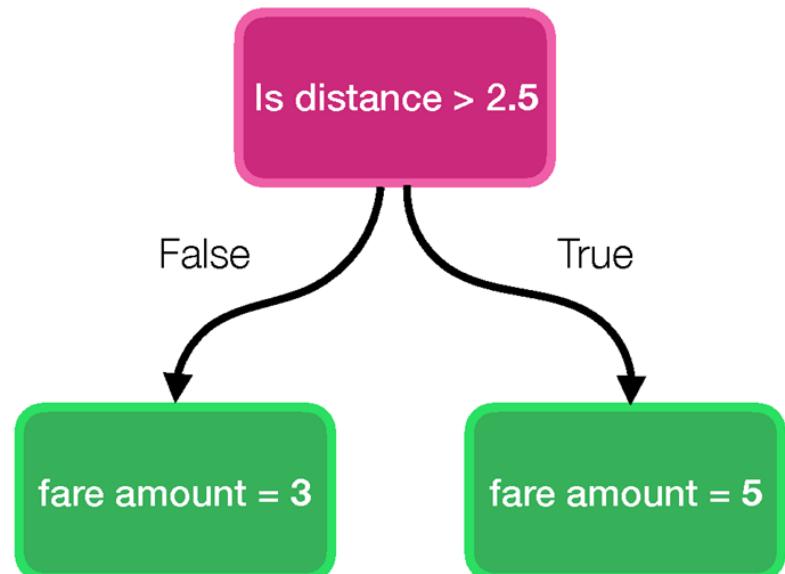
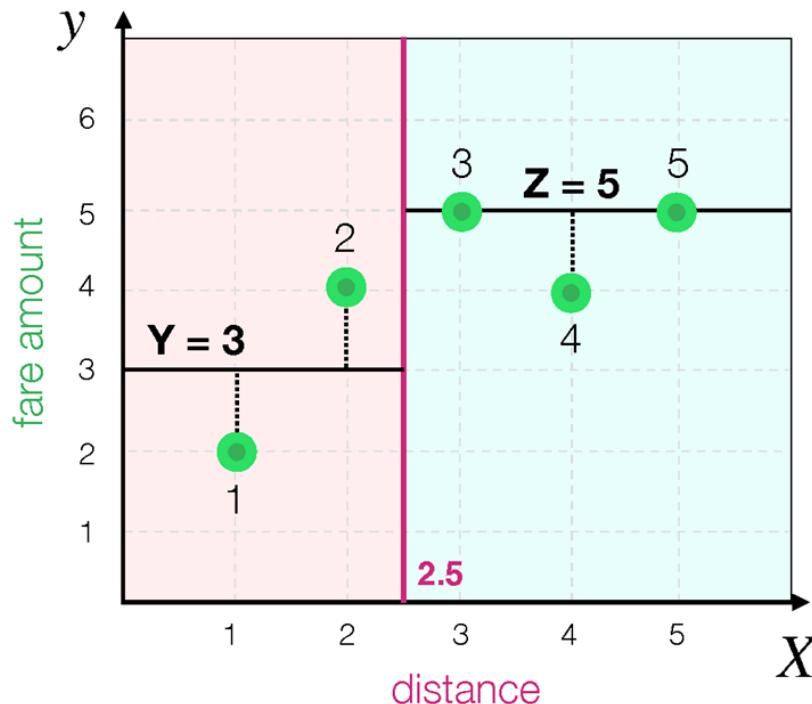
$$MSE = \frac{1}{n} \sum_{i=1}^n (\textcolor{teal}{y}_i - \hat{y}_i)^2 = \frac{3}{5} = 0.6$$



What are most reasonable values for \mathbf{Y} and \mathbf{Z} (that minimise total MSE)?

Decision Tree Algorithm

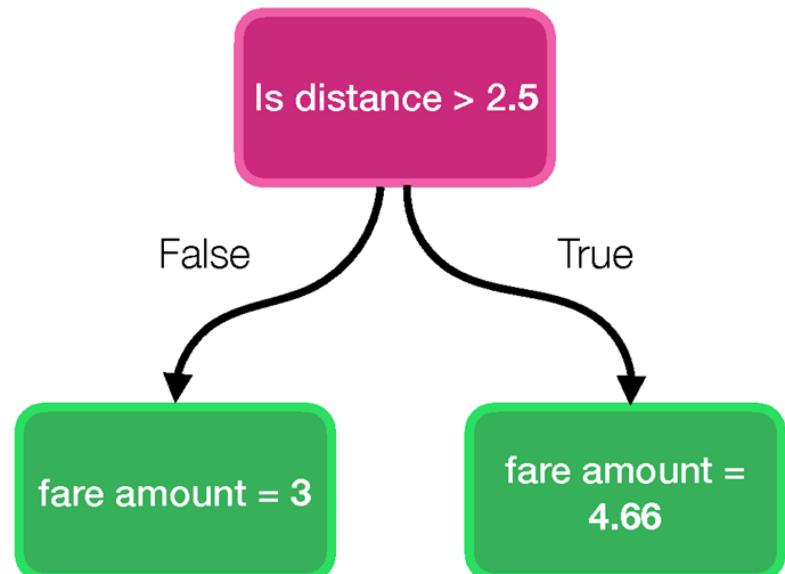
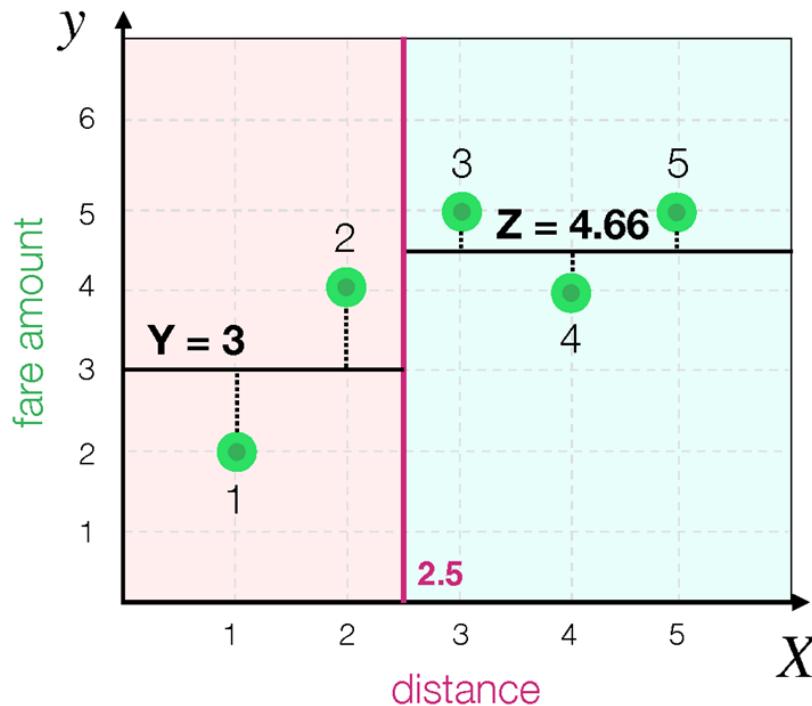
$$MSE = \frac{1}{n} \sum_{i=1}^n (\textcolor{teal}{y}_i - \hat{y}_i)^2 = \frac{3}{5} = 0.6 \quad \text{so, if } \mathbf{X} = \mathbf{2.5}, \mathbf{Y} = \mathbf{3} \text{ and } \mathbf{Z} = \mathbf{5}, \\ \text{MSE is } \mathbf{0.6}$$



What are most reasonable values for **Y** and **Z** (that minimise total MSE)?

Decision Tree Algorithm

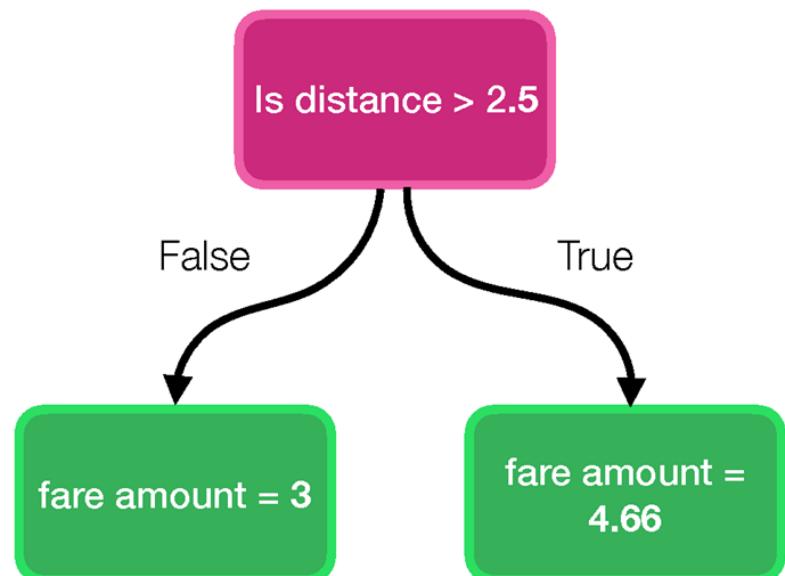
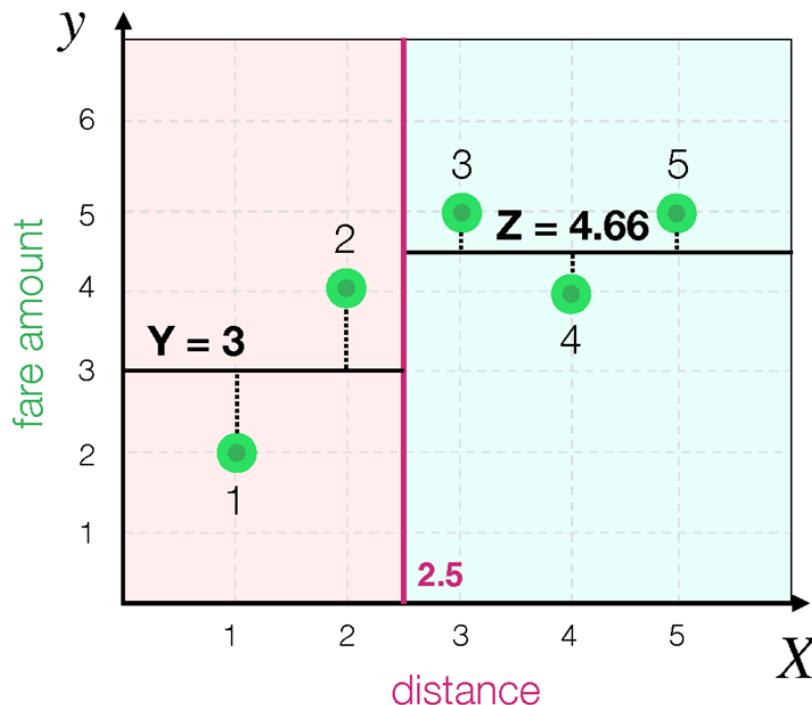
$$MSE = \frac{1}{n} \sum_{i=1}^n (\textcolor{teal}{y}_i - \hat{y}_i)^2 = \frac{(2 - 3)^2 + (4 - 3)^2 + (5 - 4.66)^2 + (4 - 4.66)^2 + (5 - 4.66)^2}{5}$$



What are most reasonable values for $\textcolor{teal}{Y}$ and $\textcolor{teal}{Z}$ (that minimise total MSE)?

Decision Tree Algorithm

$$MSE = \frac{1}{n} \sum_{i=1}^n (\textcolor{teal}{y}_i - \hat{y}_i)^2 = \frac{1 + 1 + 0.12 + 0.43 + 0.12}{5}$$

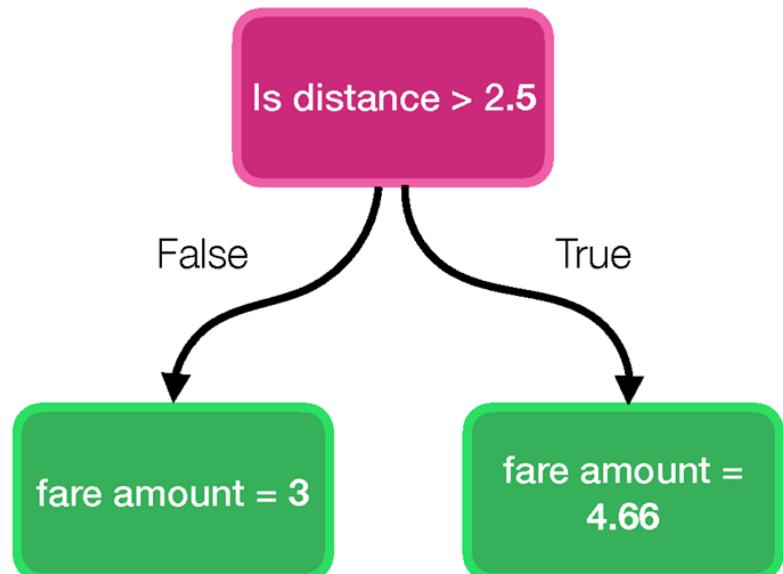
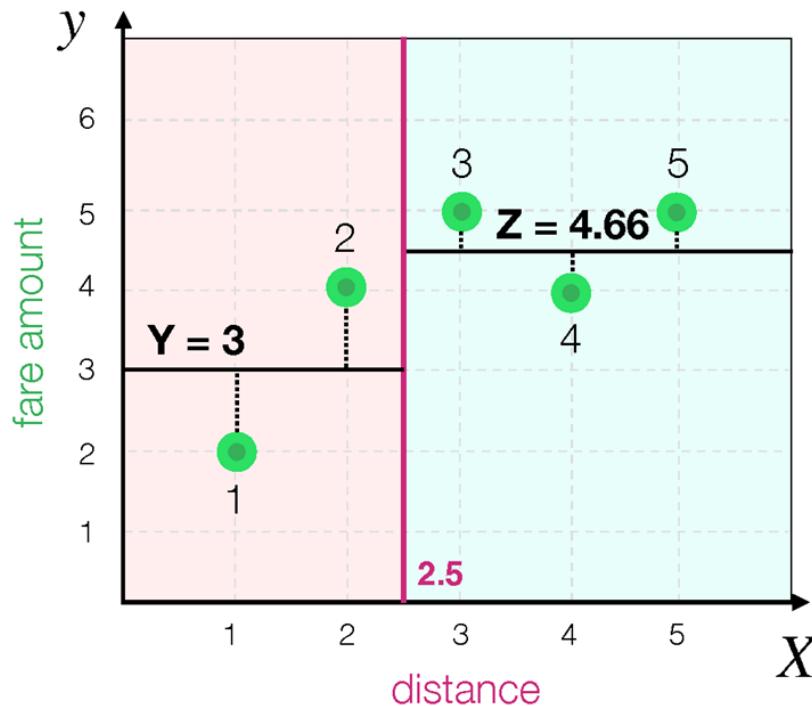


What are most reasonable values for **Y** and **Z** (that minimise total MSE)?

Decision Tree Algorithm

$$MSE = \frac{1}{n} \sum_{i=1}^n (\textcolor{teal}{y}_i - \hat{y}_i)^2 = \frac{2.67}{5} = 0.53$$

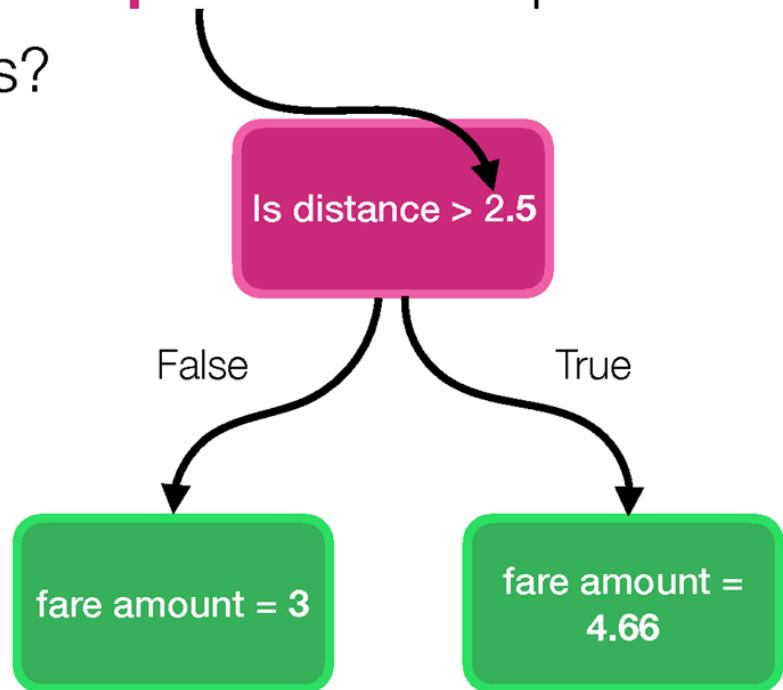
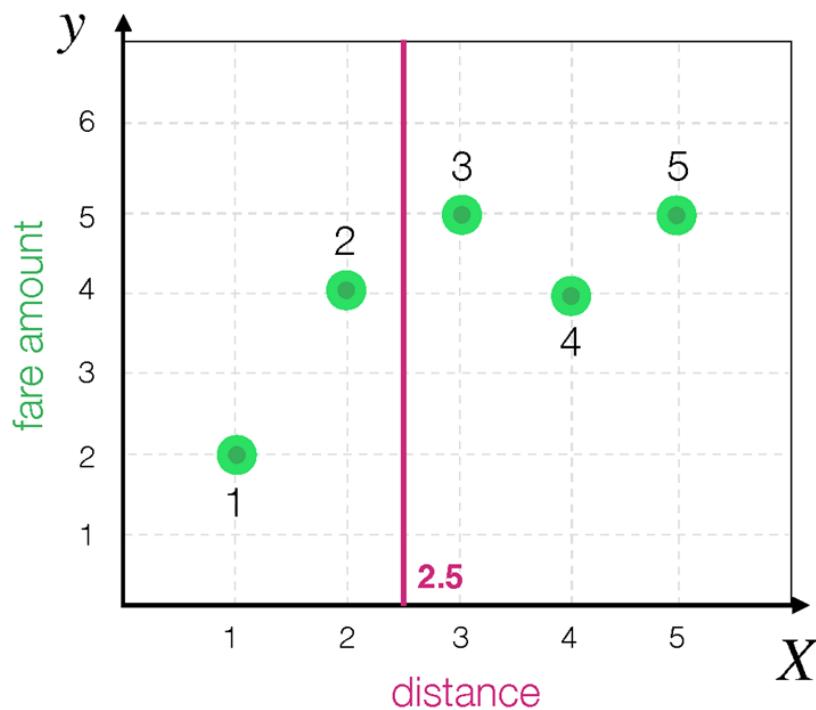
so, if **Y = 3** and **Z = 4.66**,
MSE is **smallest**



Are we happy?

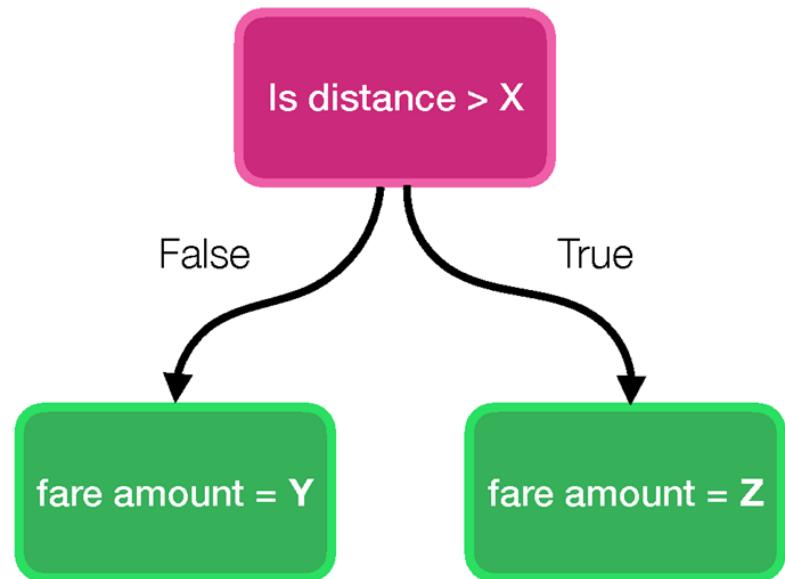
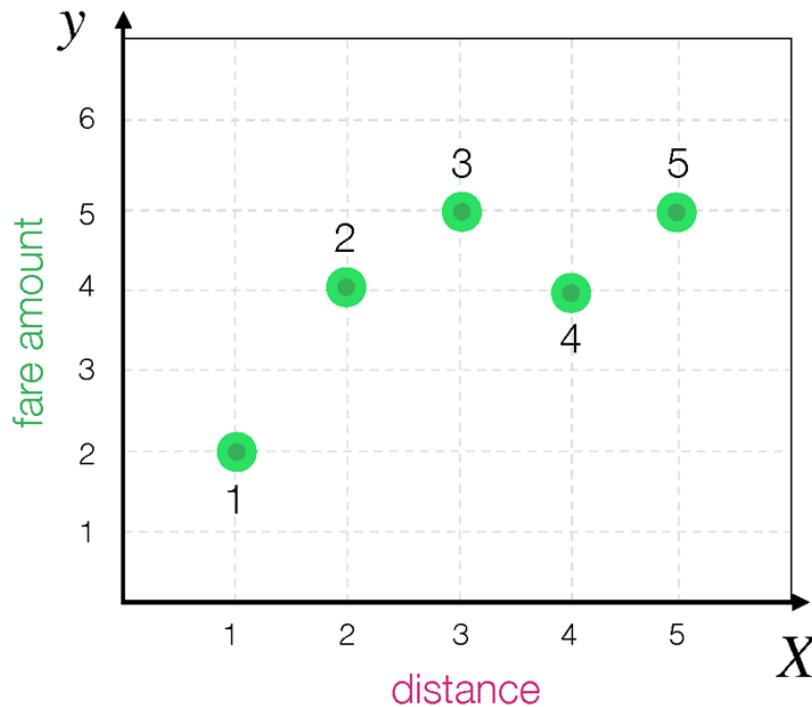
Decision Tree Algorithm

Hold on, how did we choose this **split** on the first place?
Maybe there are better options?



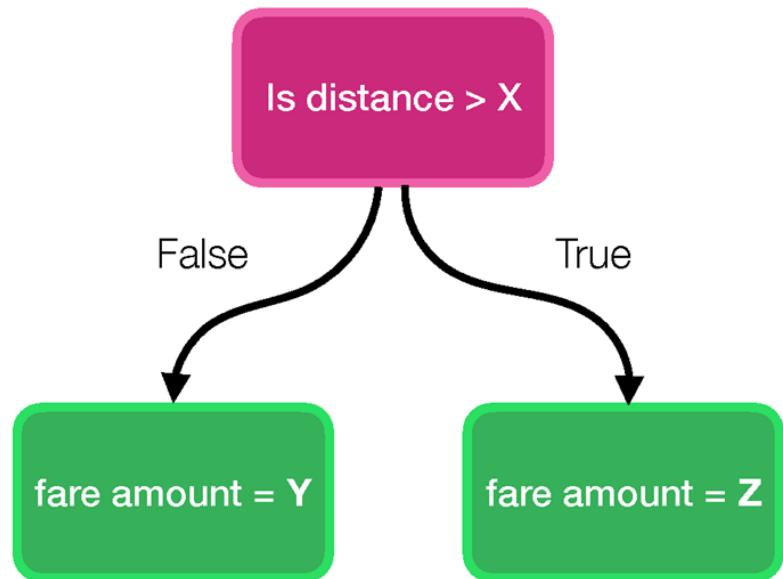
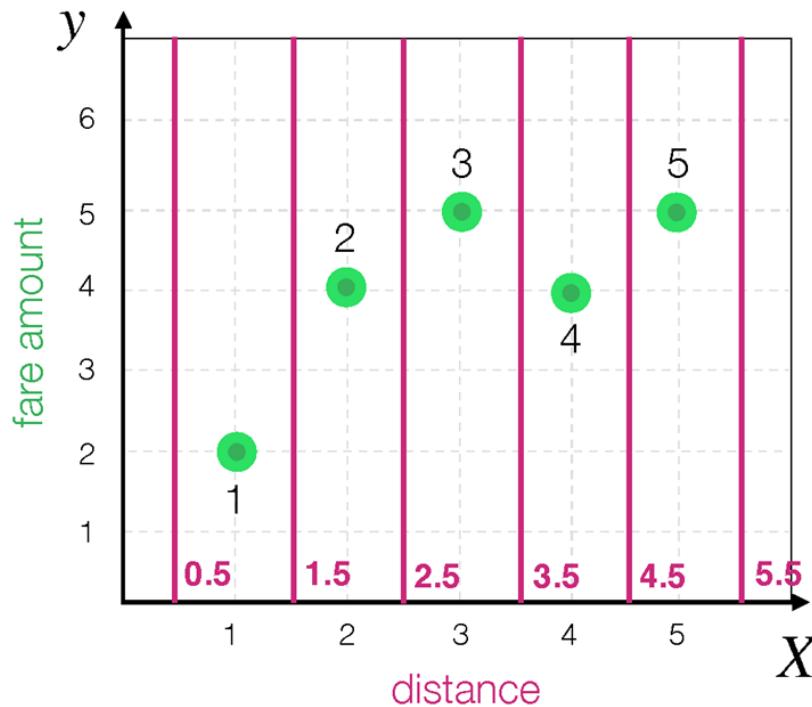
Decision Tree Algorithm

What are the possible **split** options in this case?



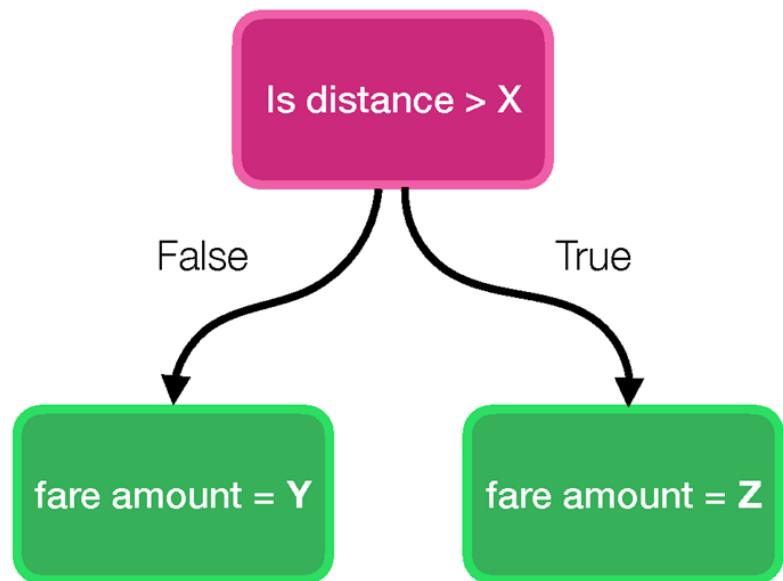
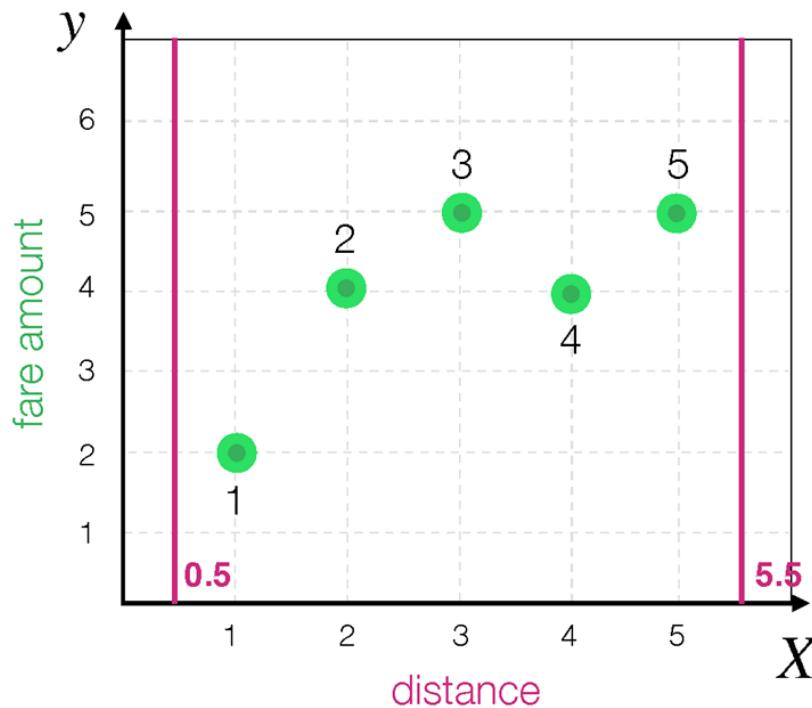
Decision Tree Algorithm

What are the possible **split** options in this case?



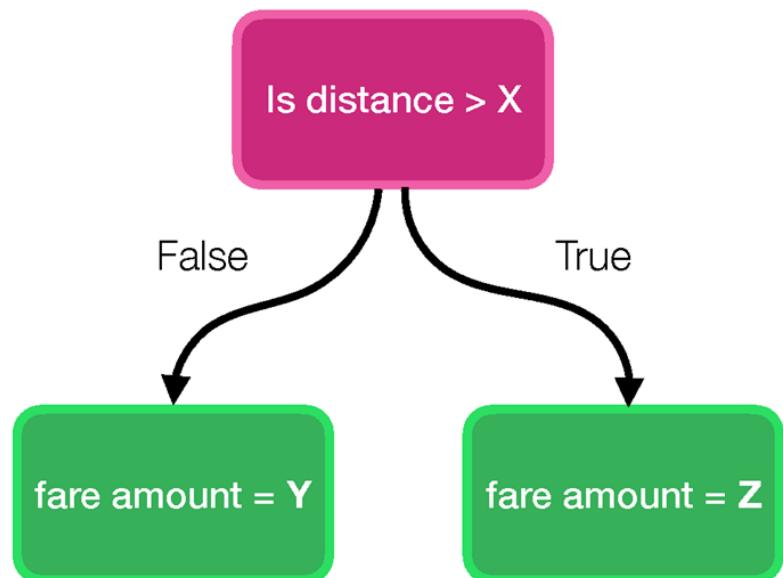
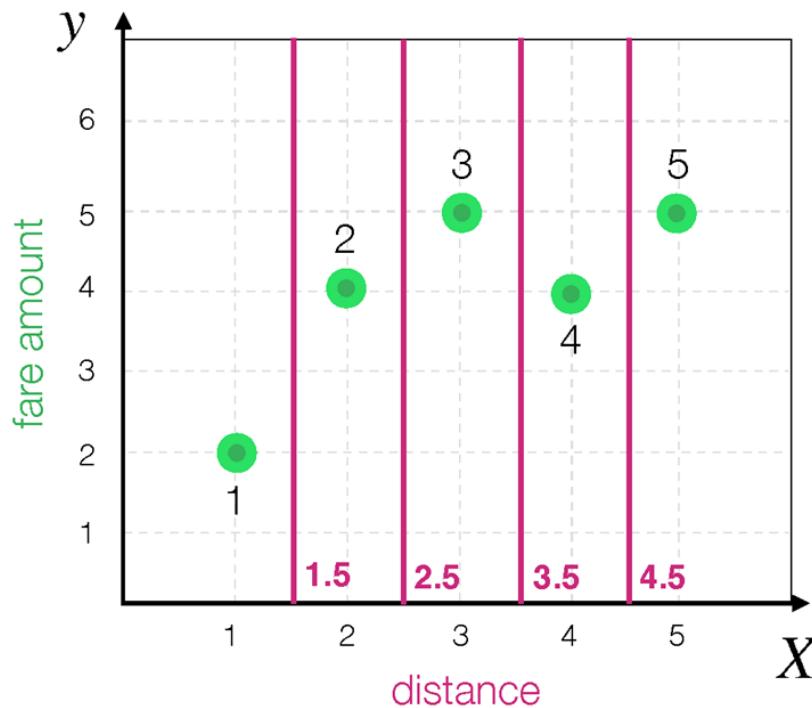
Decision Tree Algorithm

Are these meaningful?



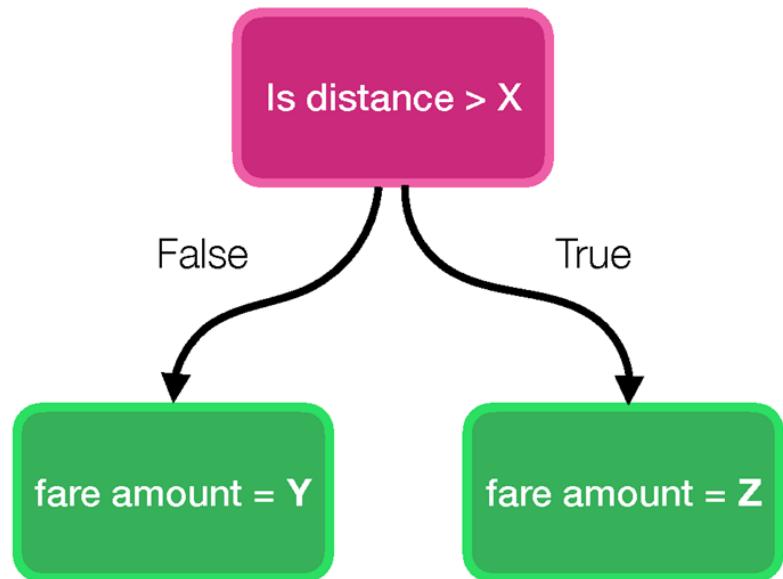
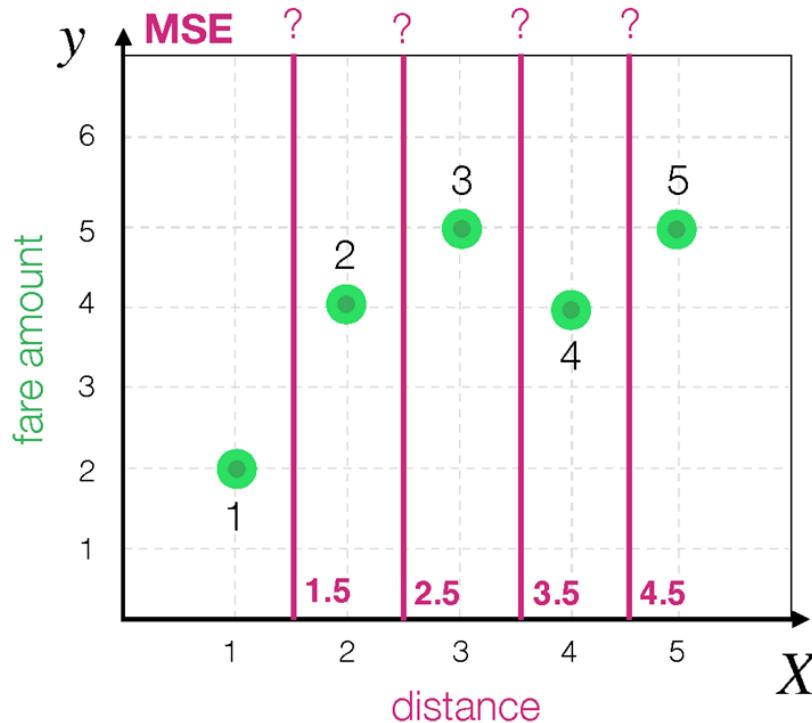
Decision Tree Algorithm

How to compare remaining?



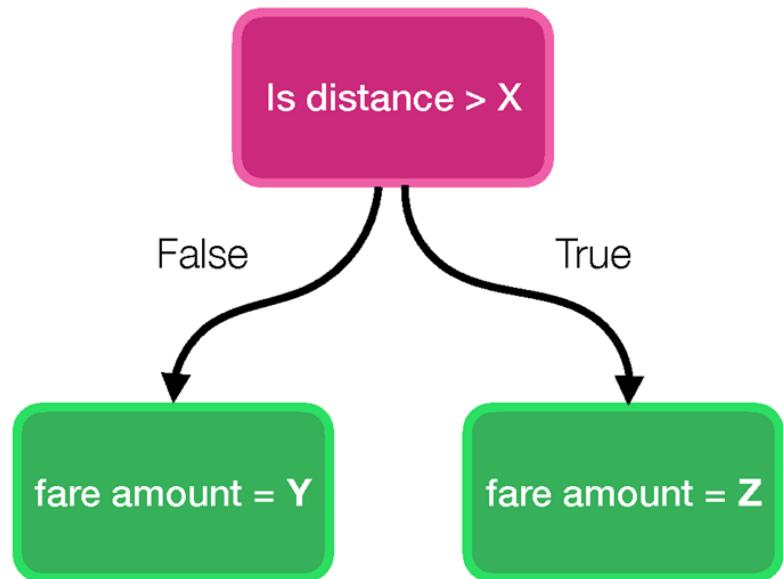
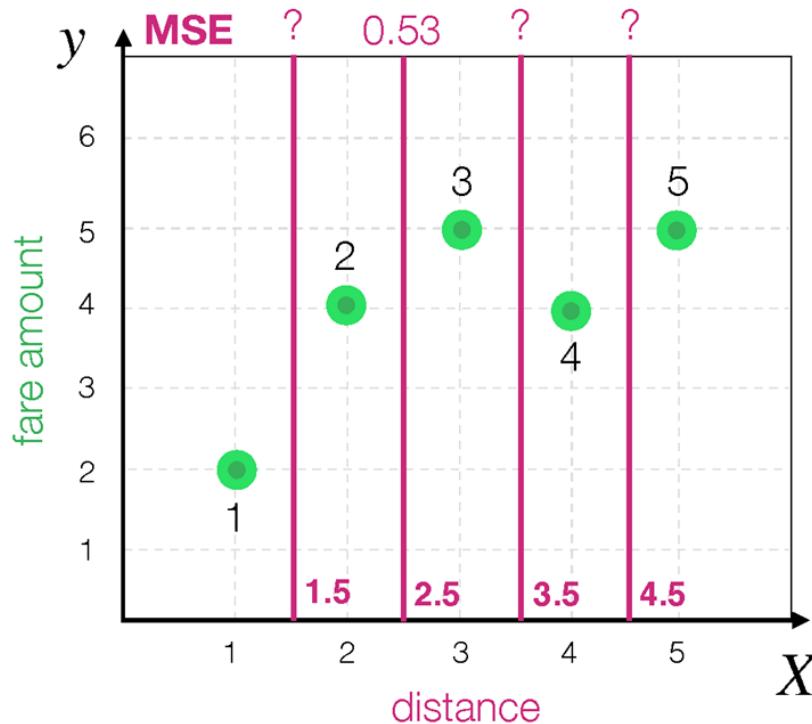
Decision Tree Algorithm

How to compare remaining?
For each one we can compute MSE

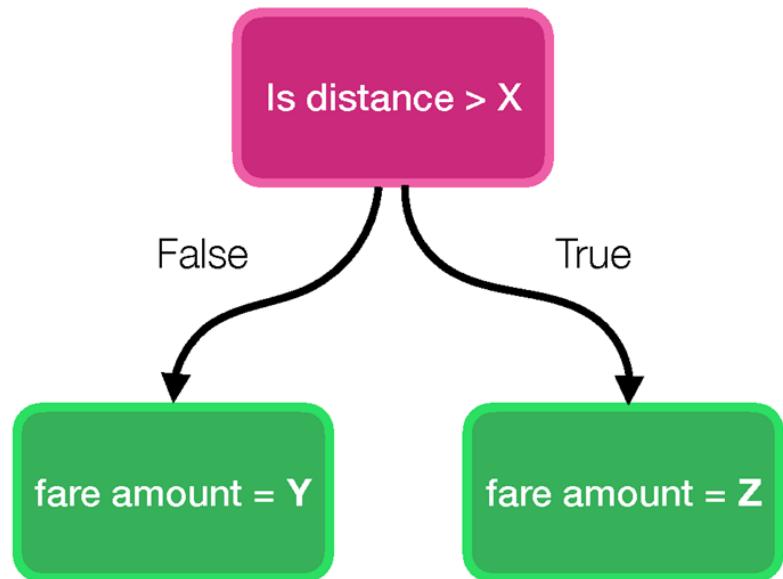
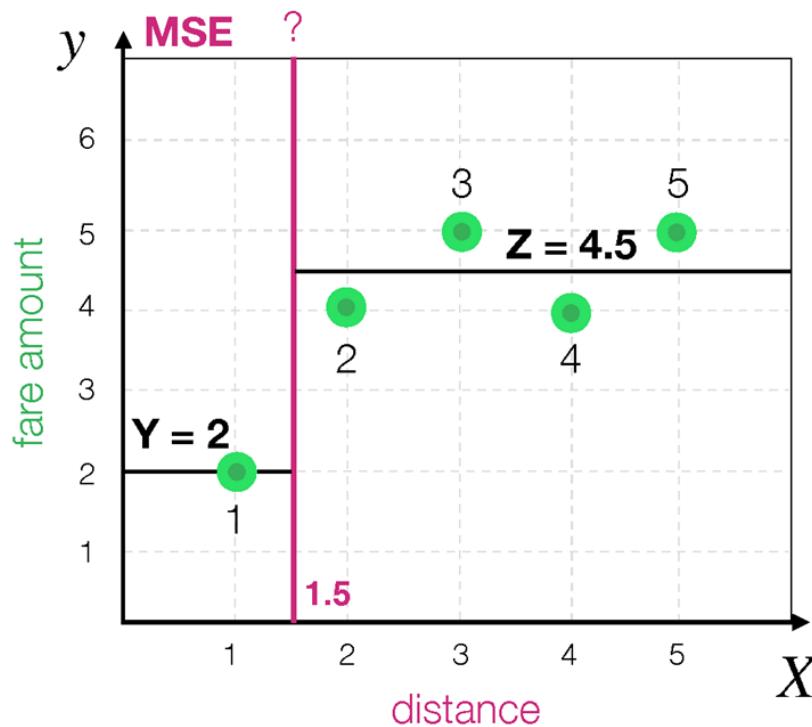


Decision Tree Algorithm

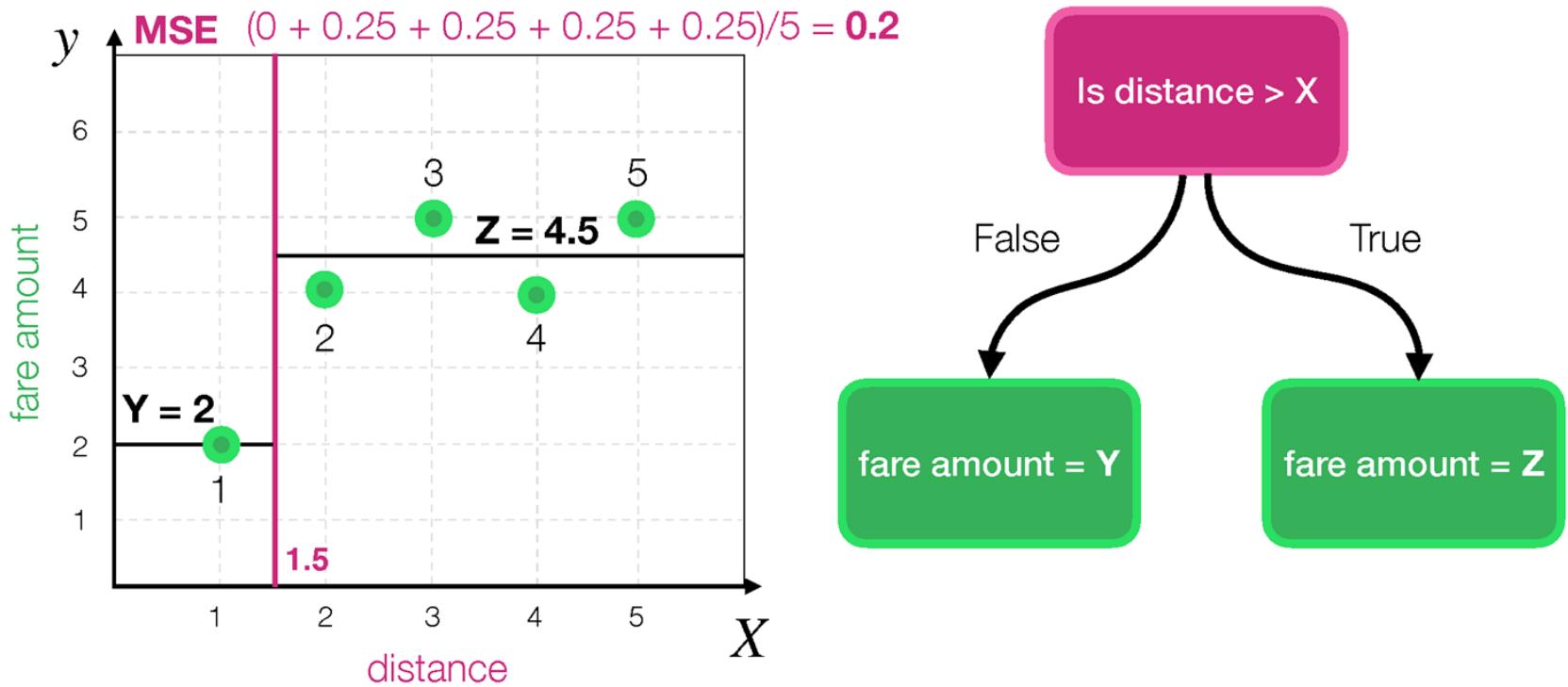
How to compare remaining?
For each one we can compute MSE



Decision Tree Algorithm

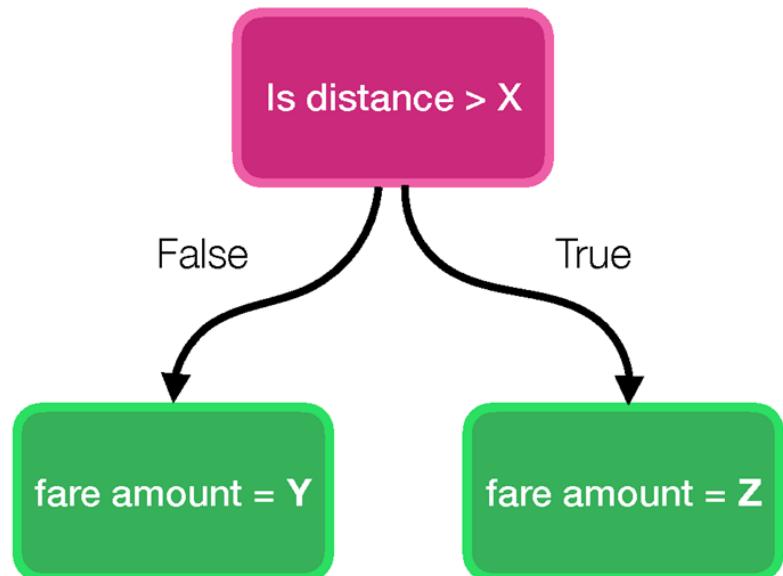
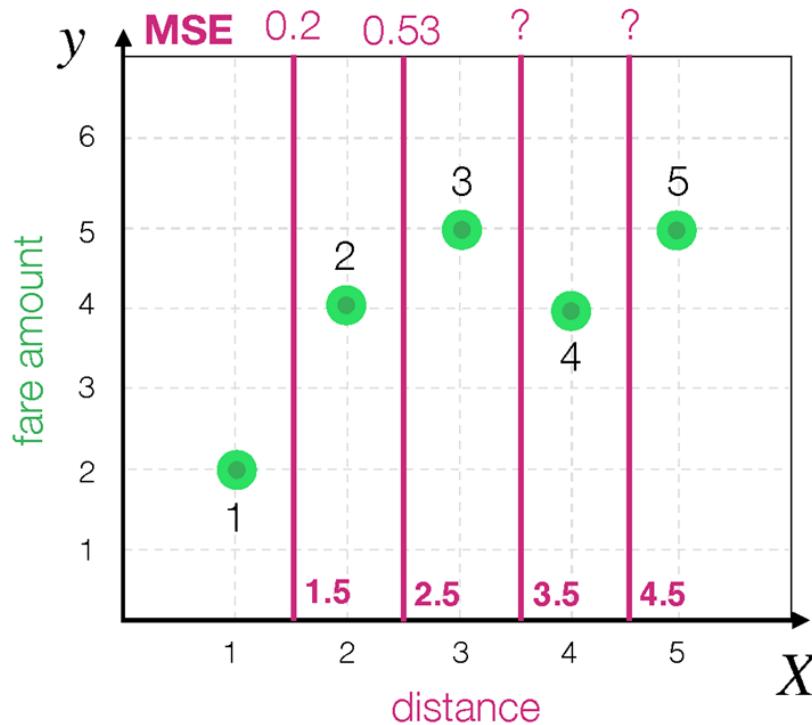


Decision Tree Algorithm

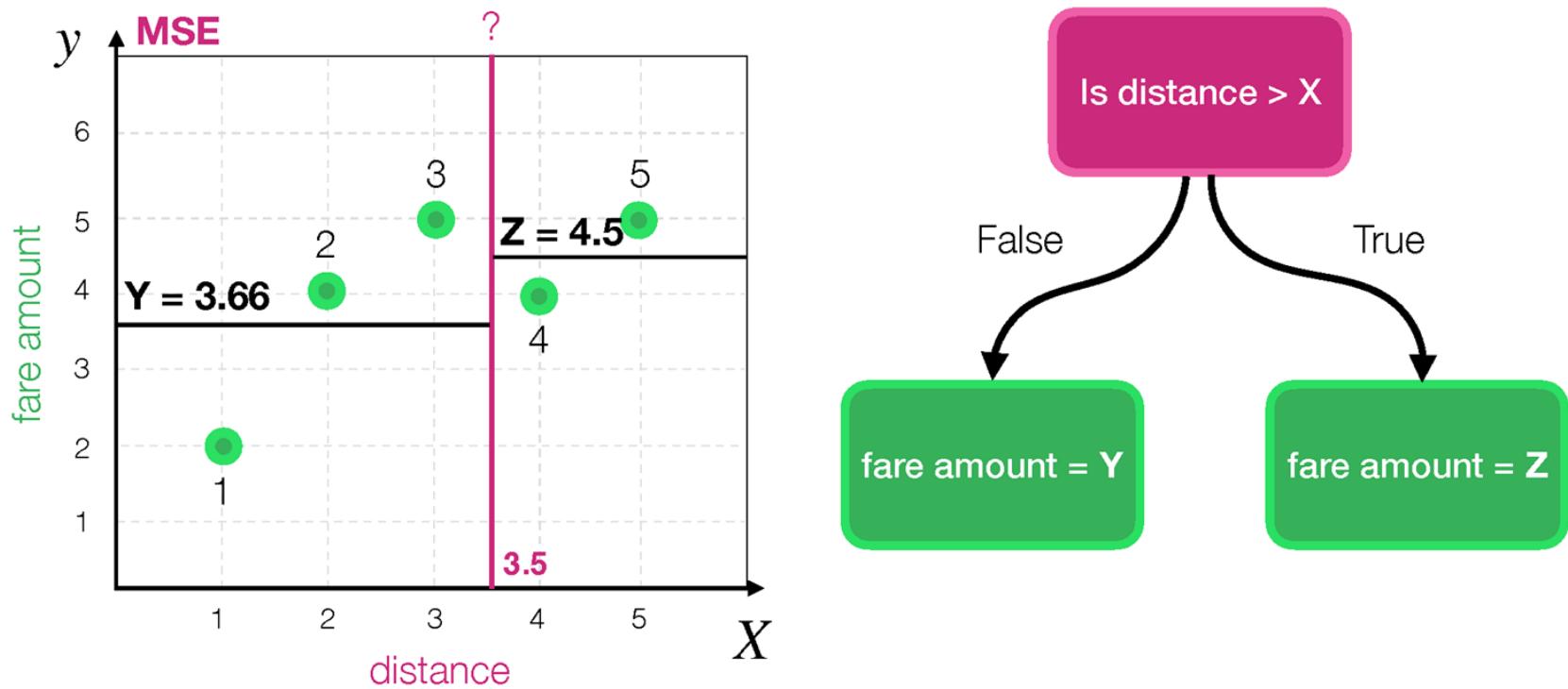


Decision Tree Algorithm

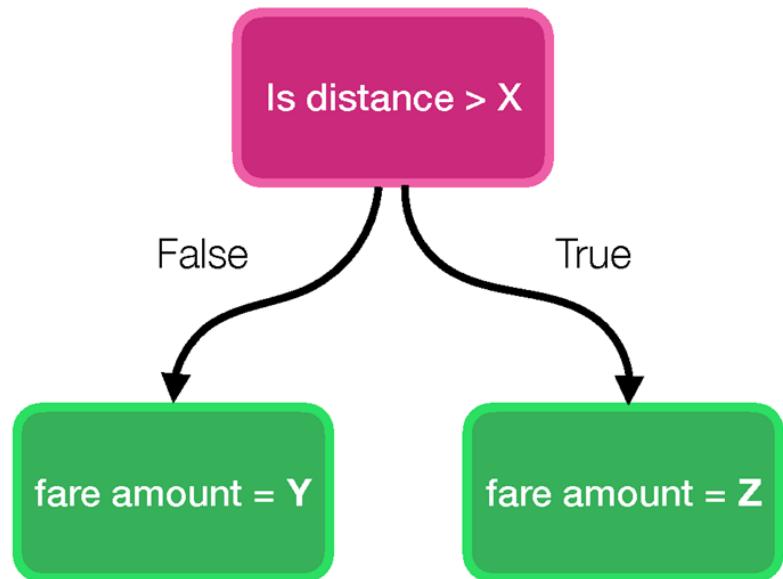
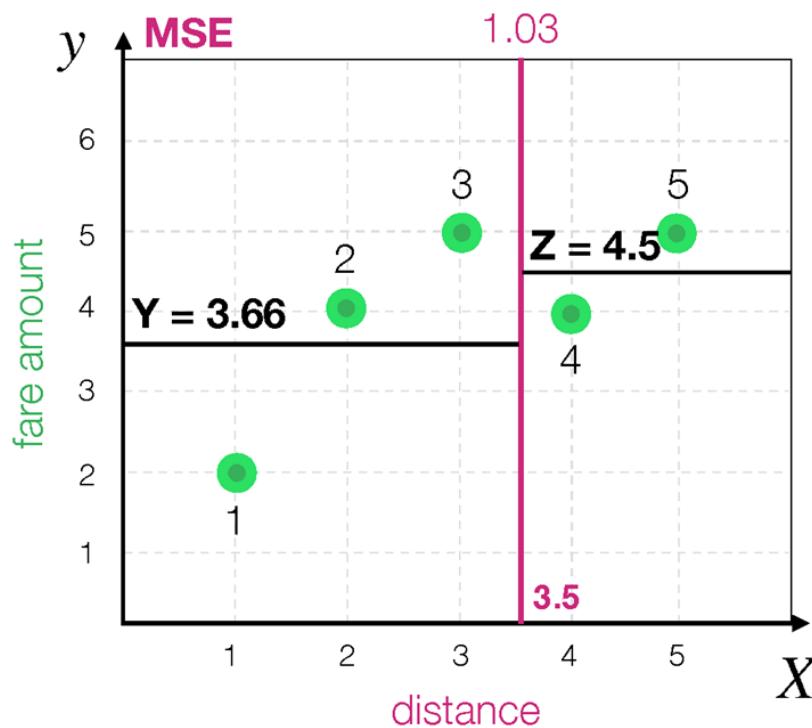
How to compare remaining?
For each one we can compute MSE



Decision Tree Algorithm

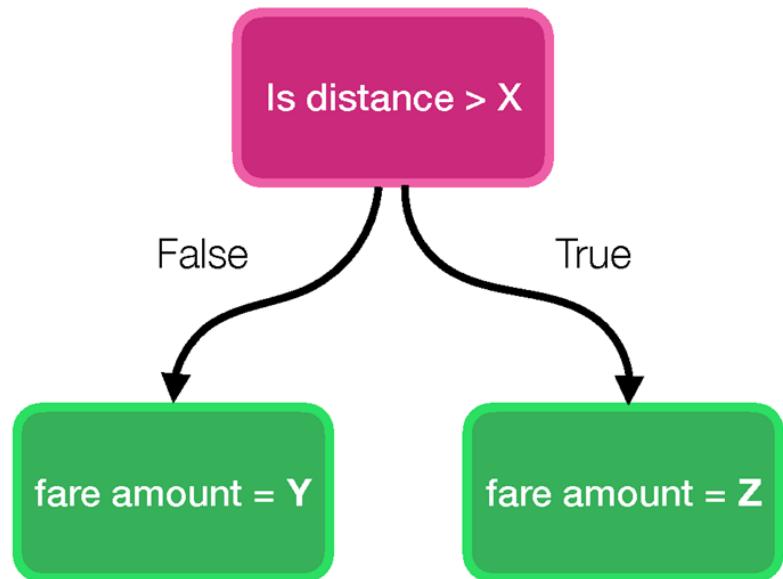
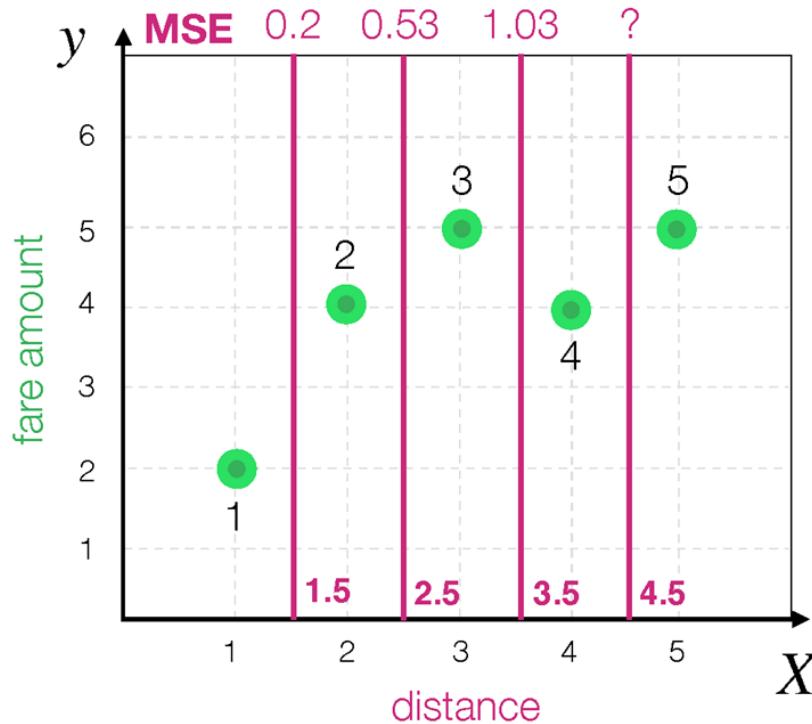


Decision Tree Algorithm

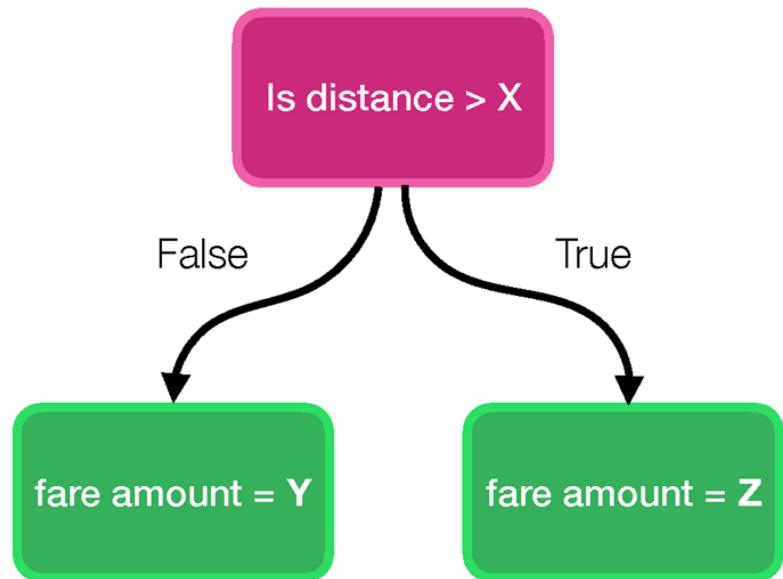
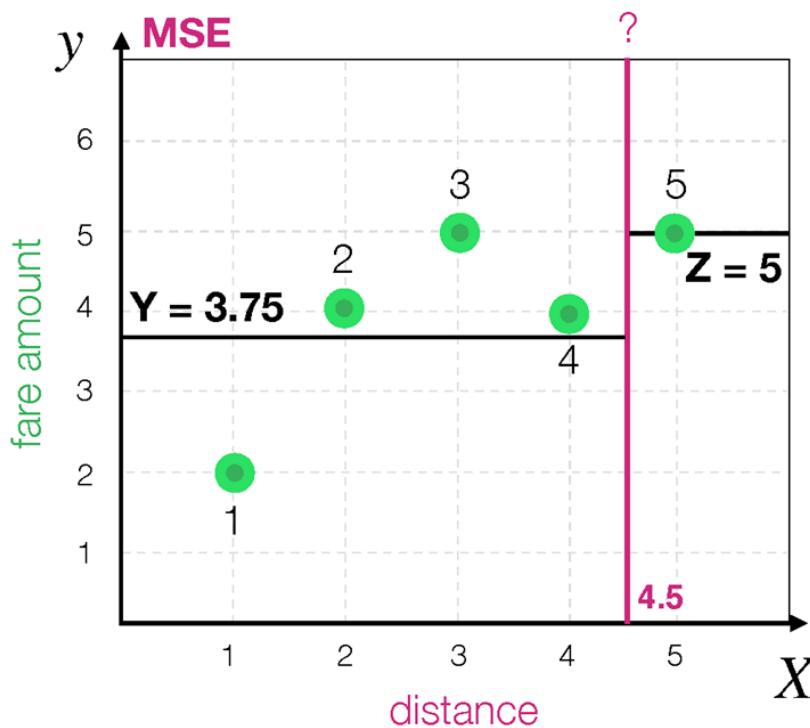


Decision Tree Algorithm

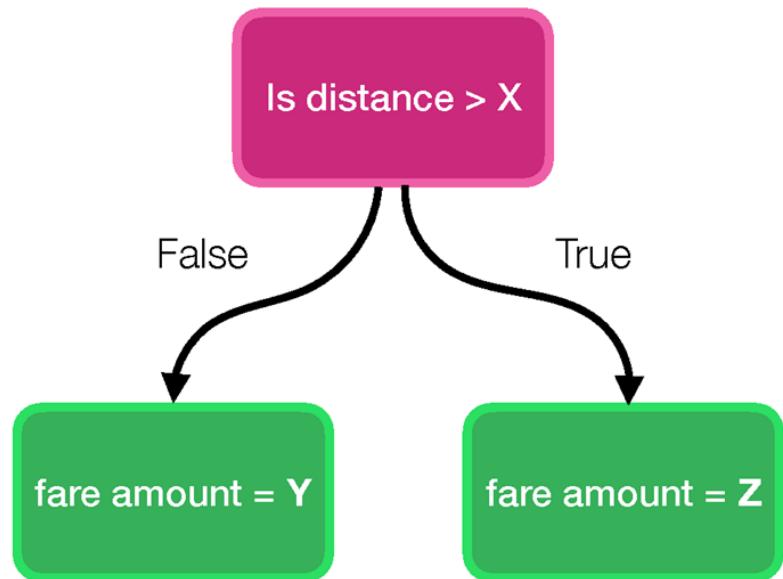
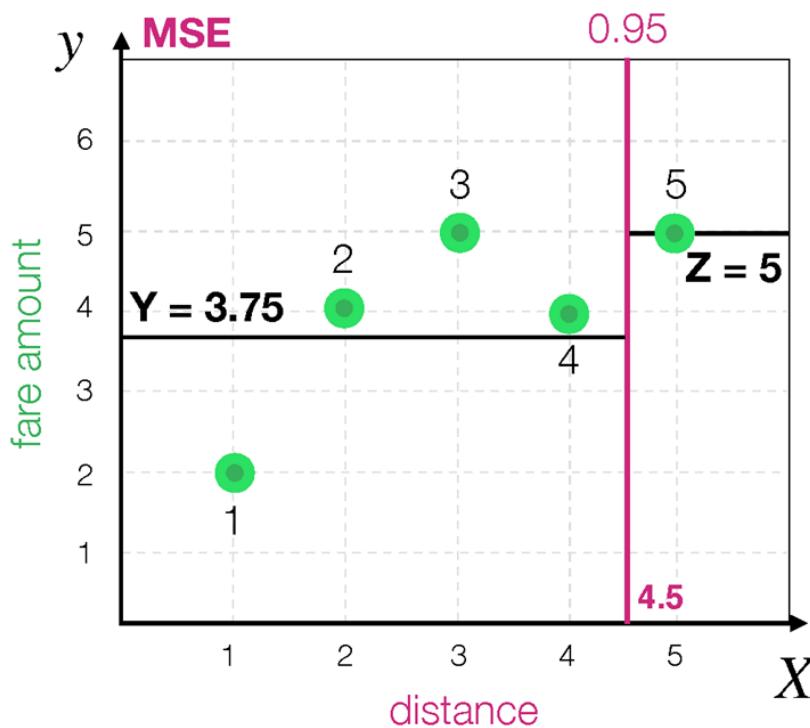
How to compare remaining?
For each one we can compute MSE



Decision Tree Algorithm

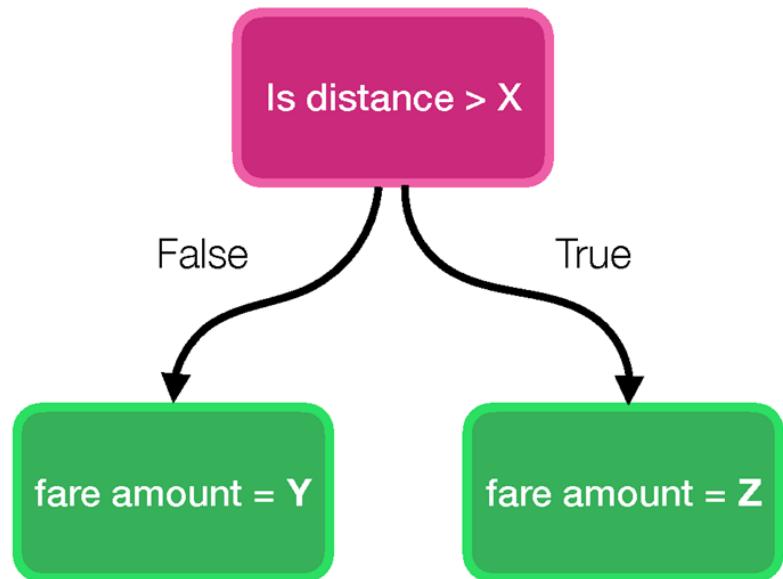
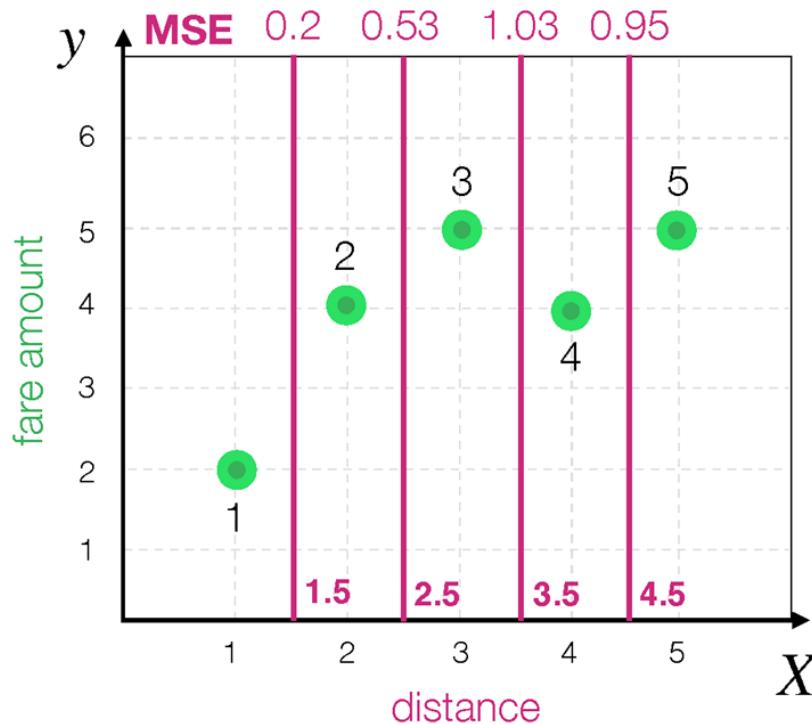


Decision Tree Algorithm



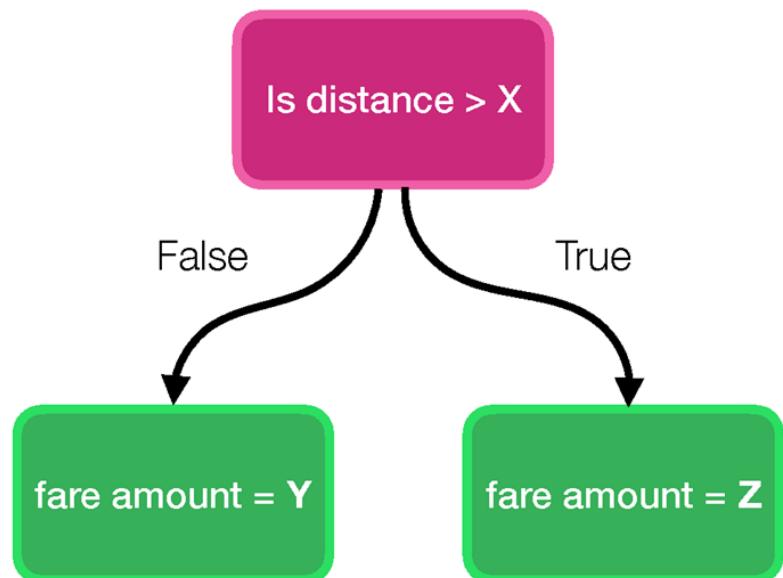
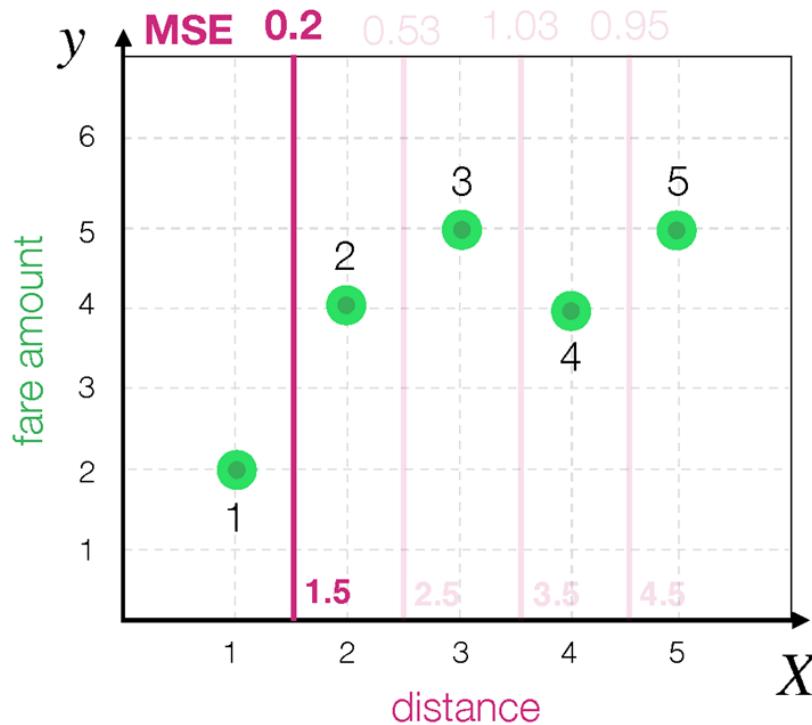
Decision Tree Algorithm

How to compare remaining?
For each one we can compute MSE



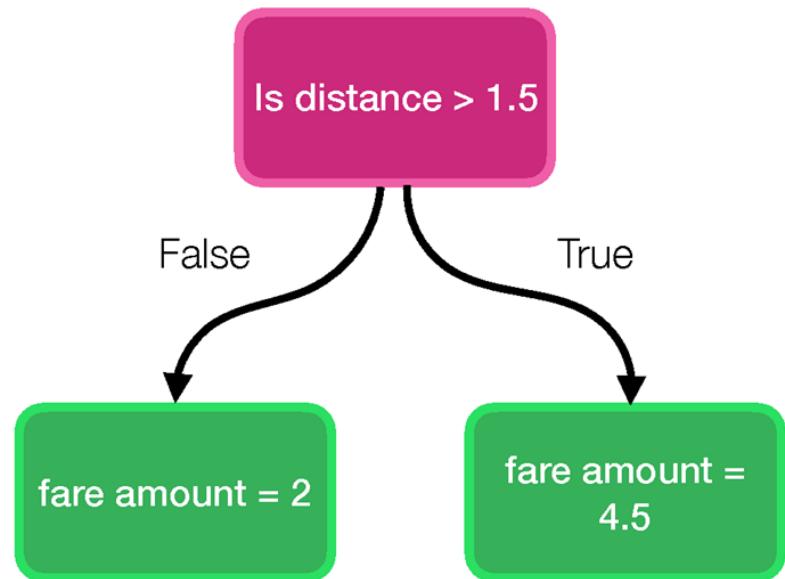
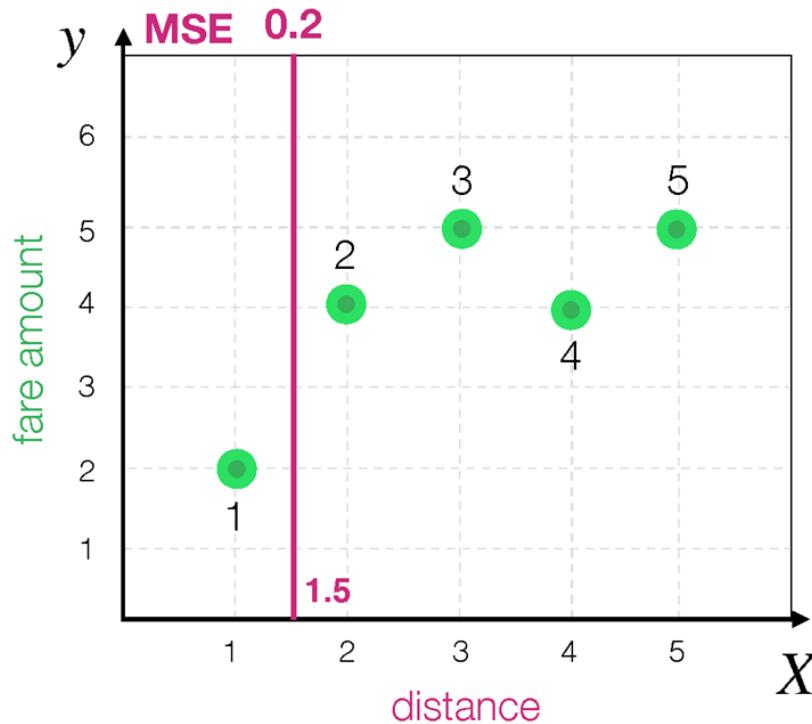
Decision Tree Algorithm

We choose the **split** that minimises total MSE



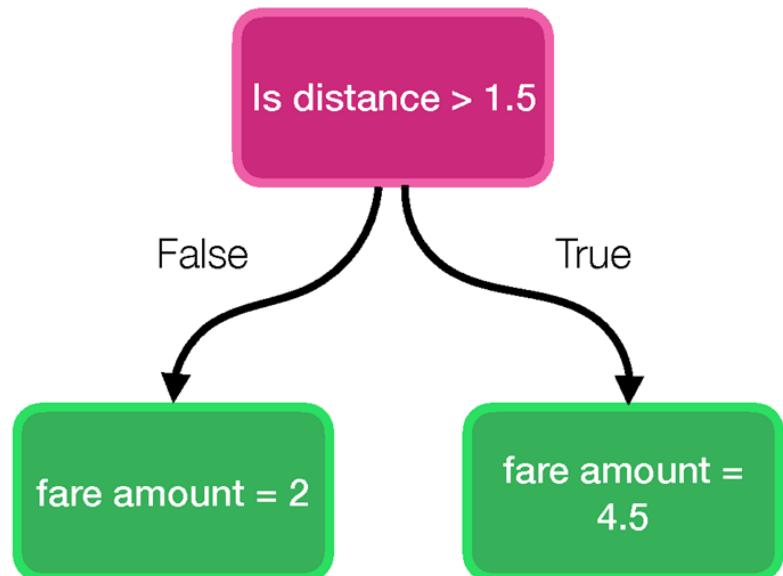
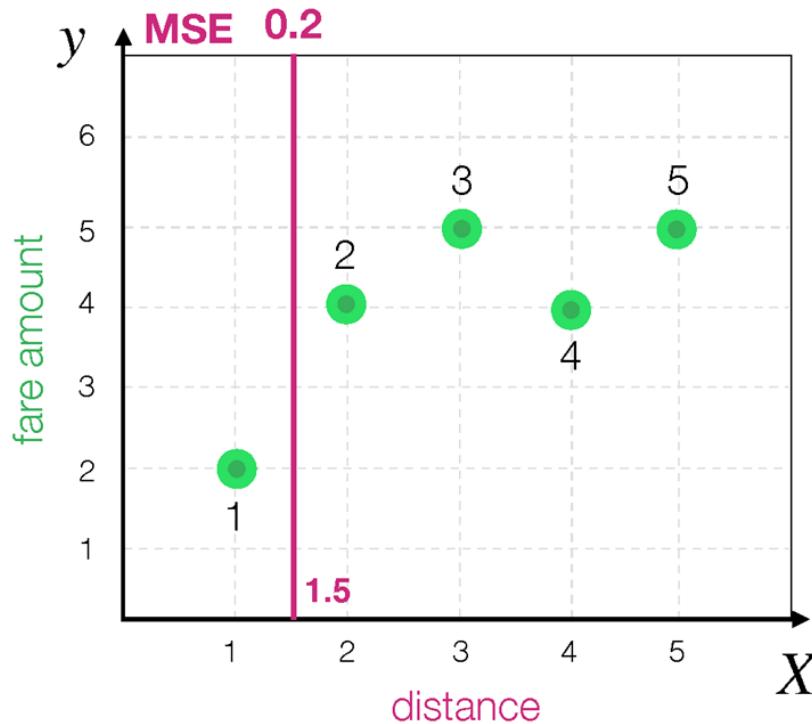
Decision Tree Algorithm

Thus, the resulting tree:



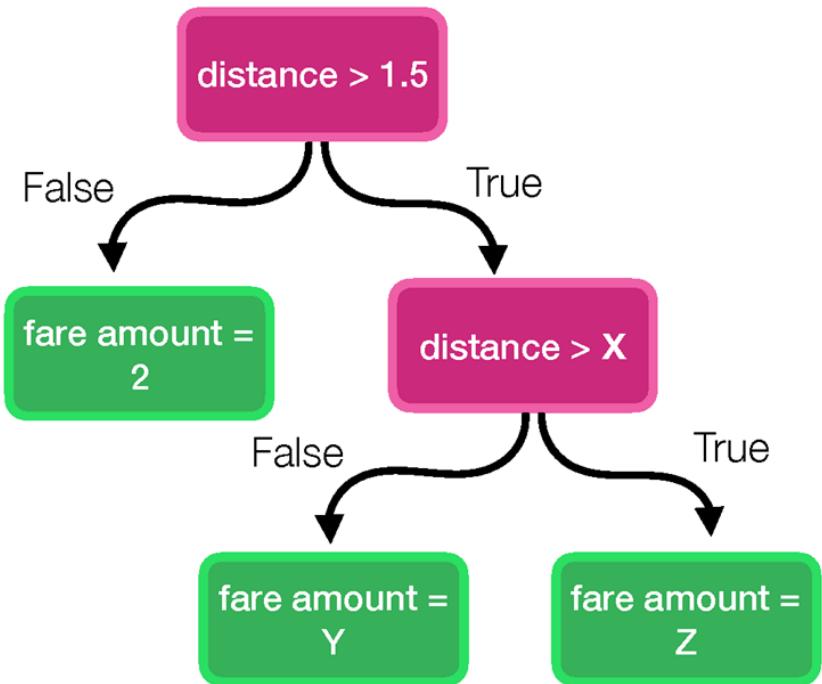
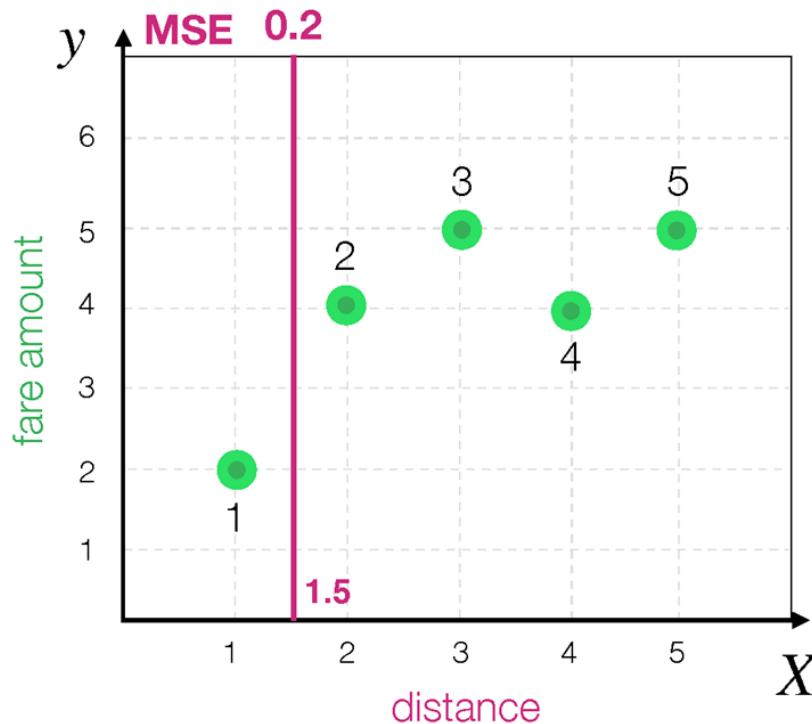
Decision Tree Algorithm

Can we make our decision tree more accurate?



Decision Tree Algorithm

Can we make our decision tree more accurate?
Yes, by going **deeper!**



Characteristics of decision trees

Pros	Cons
Easy to interpret	Easy to overfit or underfit the model
Can handle numeric or categorical features	Cannot model interactions between features
Can handle missing data	Large trees can be difficult to interpret
Uses only the most important features	
Can be used on very large or small data	

A tree stops growing at a node when...

- pure or nearly pure
- no remaining variables on which to further subset the data
- the tree has grown to a preselected size limit

Bias-variance tradeoff

$$\text{error} = \text{bias} + \text{variance}$$

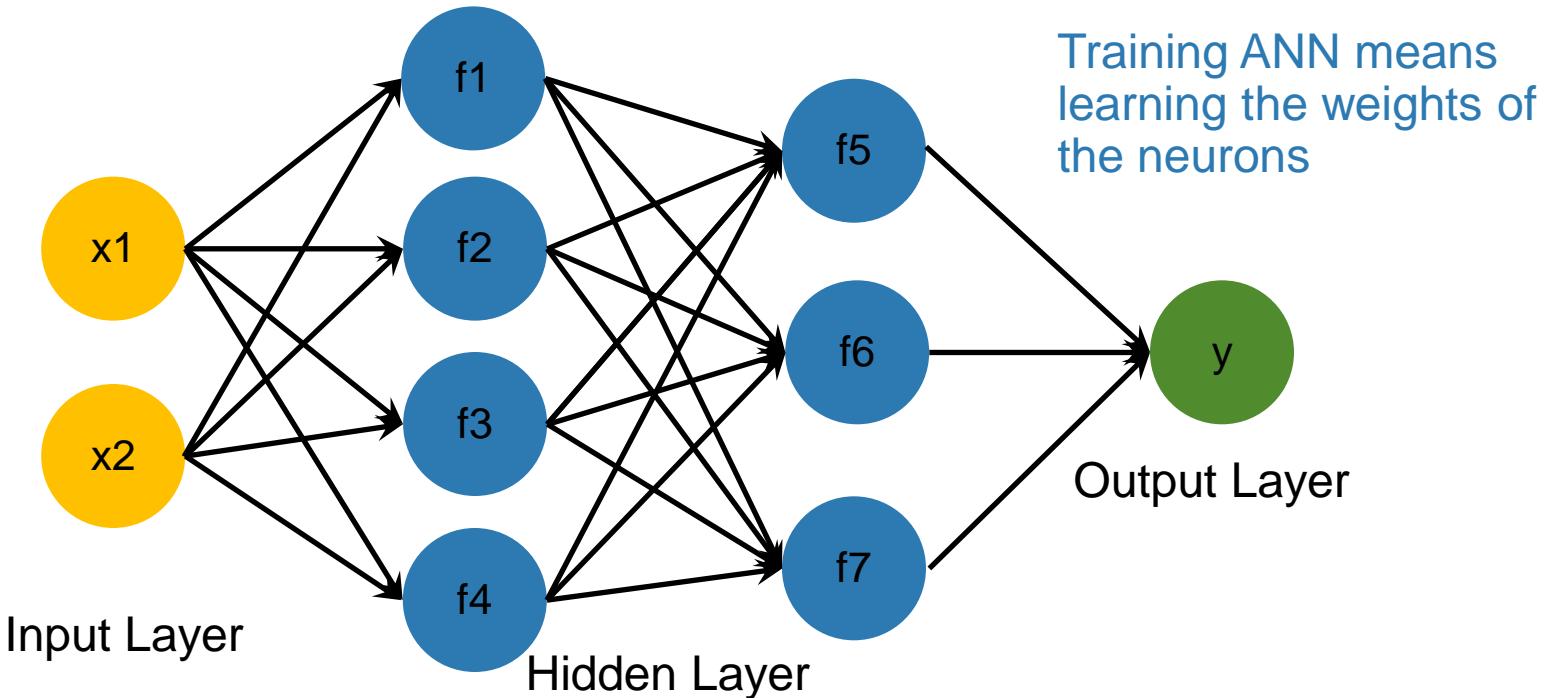
- **Bias-variance tradeoff** is the problem of simultaneously minimizing two sources of error that prevent supervised learning algorithms from generalizing beyond their training set:
 - **Bias** is error from erroneous assumptions in the learning algorithm, high bias can cause an algorithm to miss the relevant relations between features and target outputs (**underfitting**)
 - **Variance** is error from sensitivity to small fluctuations in the training set, high variance can cause **overfitting**, i.e., modeling the random noise in the training data, rather than the intended outputs
- **Ensemble tree methods:**
 - **Gradient Boosting** (GBoost) is based on weak learners (high bias, low variance). In terms of decision trees, weak learners are shallow trees, sometimes even as small as decision stumps (trees with two leaves). Boosting reduces error mainly by reducing bias.
 - **Random Forest** uses fully grown decision trees (low bias, high variance). It tackles the error reduction task by reducing variance. The trees are made uncorrelated to maximize the decrease in variance, but the algorithm cannot reduce bias (which is slightly higher than the bias of an individual tree in the forest).



Supervised Machine Learning – Neural Networks

Artificial neural networks (ANNs)

- Based loosely on computer models of how brains work
- Model is an assembly of inter-connected neurons (nodes) and weighted links
- Each neuron applies a nonlinear function to its inputs to produce an output
- Output node sums up each of its input value according to the weights of its links
- Used for classification, pattern recognition, speech recognition
- “Black Box” model – no explanatory power, very hard to interpret the results
- ANNs with more than one hidden layer are called Deep Learning networks



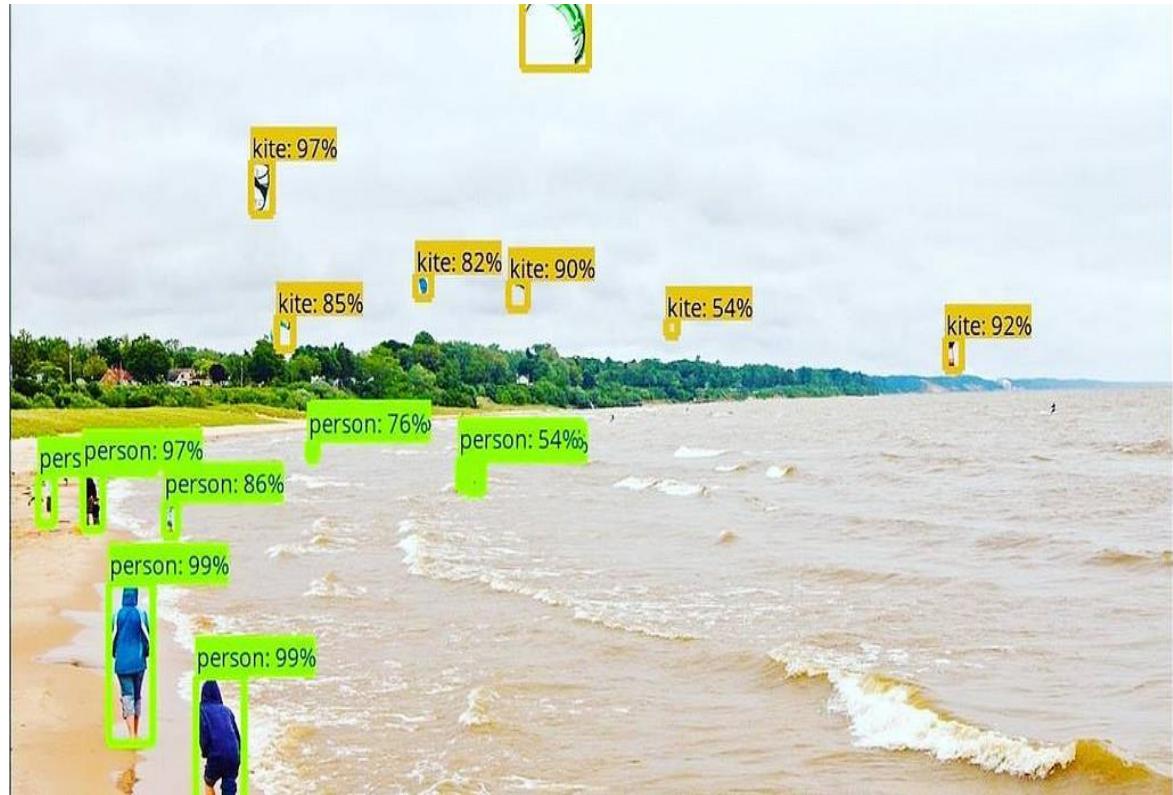
Applications of Deep Learning – Computer Vision

- **Image classification**
- Object detection
- Type: Feedforward Neural Network
- Type: Convolutional Neural Network



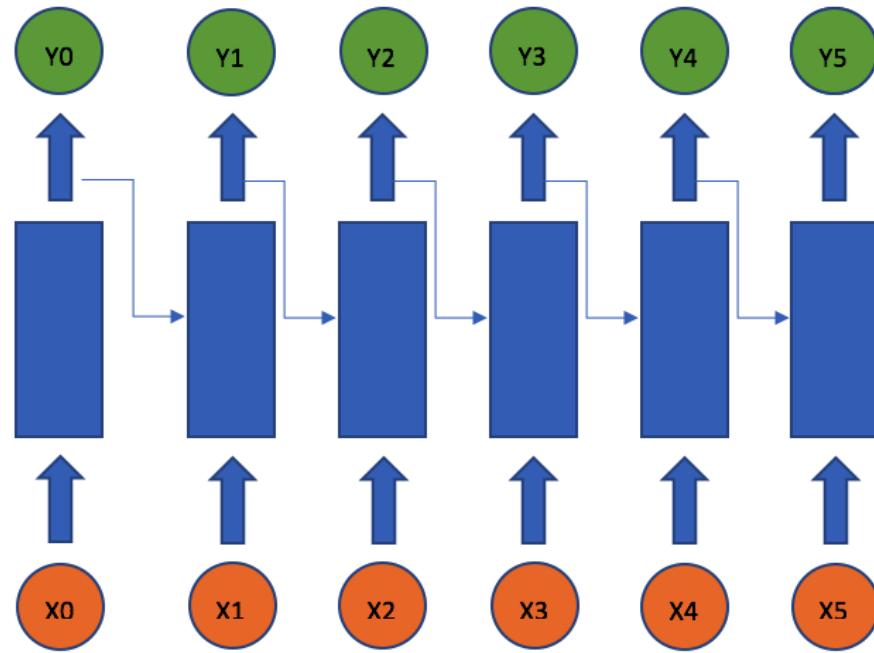
Applications of Deep Learning – Computer Vision

- Image classification
- **Object detection**
- Type: Feedforward Neural Network
- Type: Convolutional Neural Network



Applications of Deep Learning – Natural Language Processing and Time Series Analysis

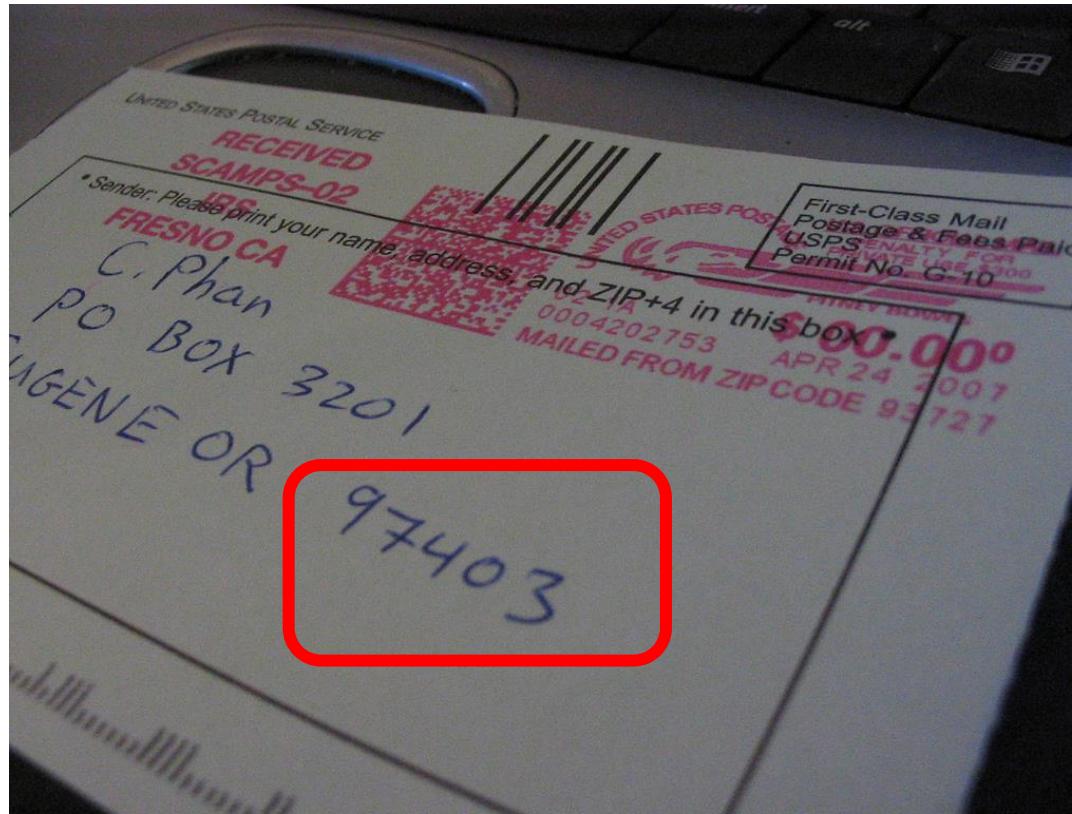
- Speech recognition
- Automatic machine translation
- Time series prediction
- Type: Recurrent Neural Network



<https://blog.statsbot.co/time-series-prediction-using-recurrent-neural-networks-lstms-807fa6ca7f>

Applications in Computer Vision

- Application: automatically classifying **postal codes**

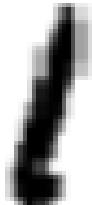


Target variables in MNIST dataset – multi-class classification

$y = 0$



$y = 1$



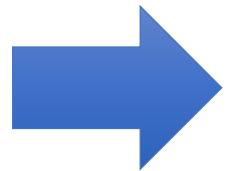
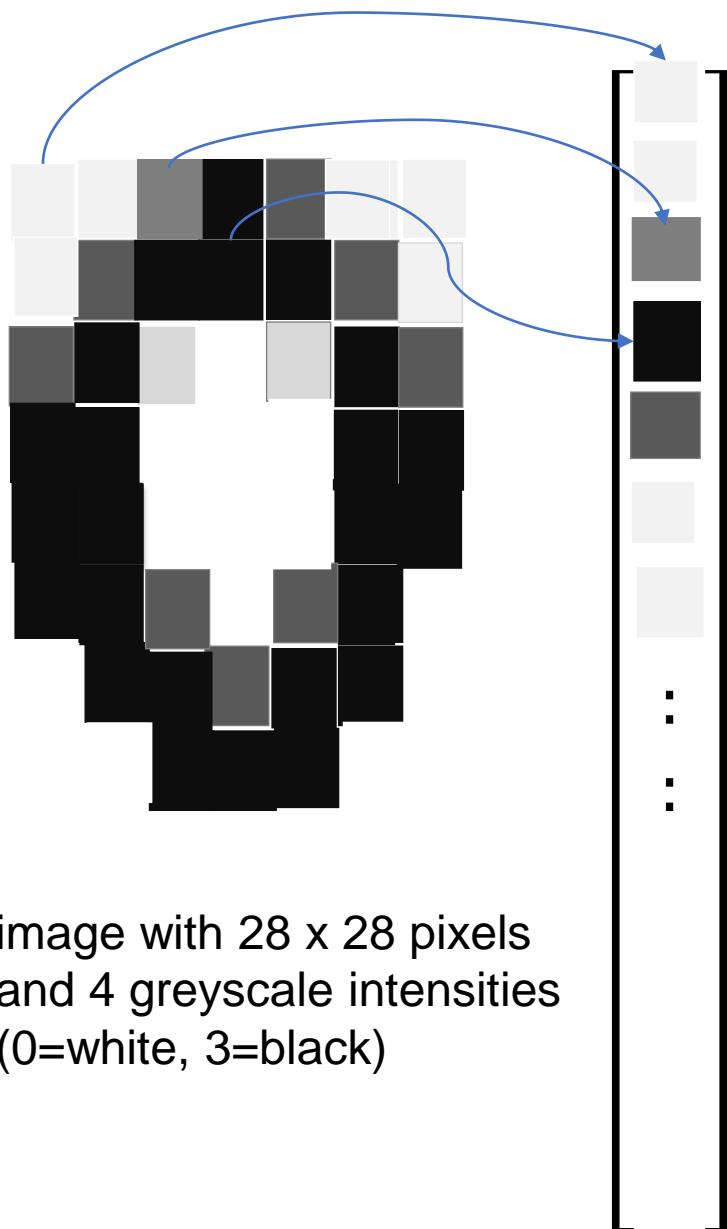
$y = 2 \dots$



$y = 9$



...



$$\mathbf{X} = \begin{bmatrix} 1 \\ 1 \\ 2 \\ 3 \\ 2 \\ \vdots \\ \vdots \end{bmatrix}$$

A blue bracket on the right side of the vector indicates its total length, labeled "784" in red text.

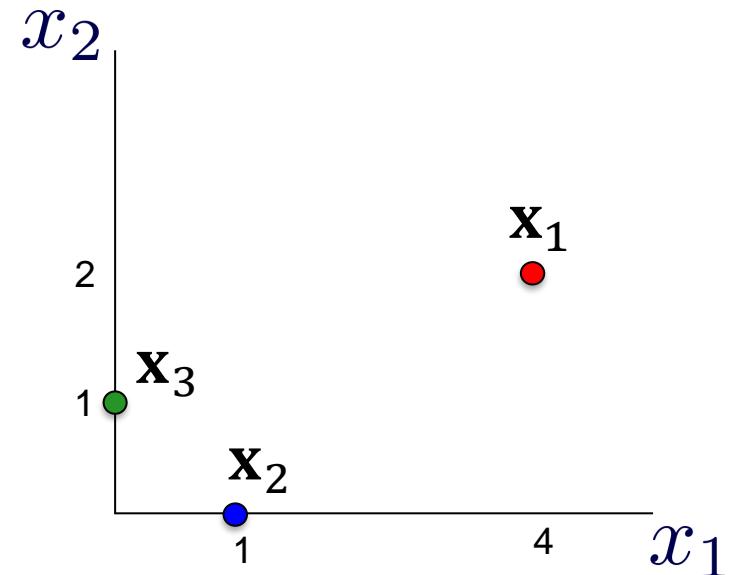
Features and targets: 3-class example with 2 features

Colors are used to indicates the class

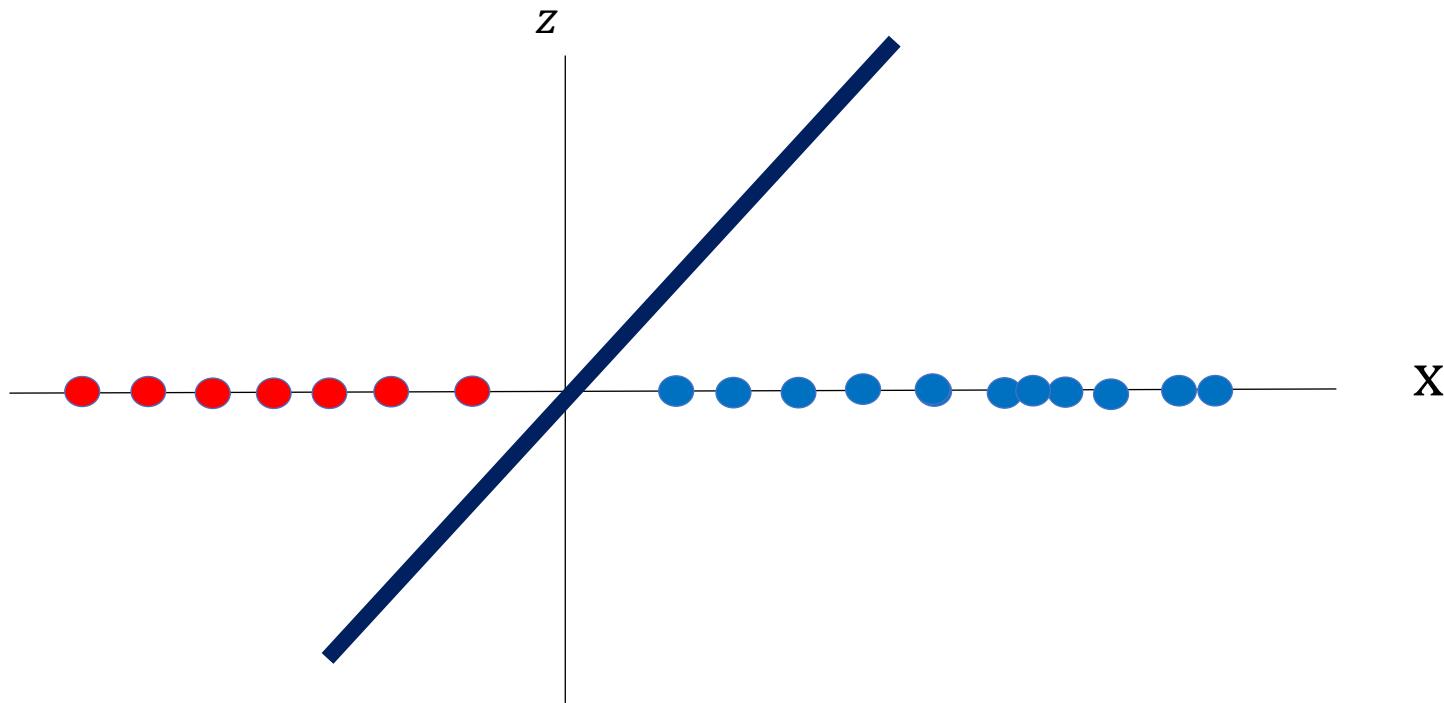
1. ($y_1 = \textcolor{red}{0}$, $\mathbf{x}_1 = [4, 2]$)

2. ($y_2 = \textcolor{blue}{1}$, $\mathbf{x}_2 = [1, 0]$)

3. ($y_3 = \textcolor{green}{2}$, $\mathbf{x}_3 = [0, 1]$)



Linear classifier



Linear classifier

- The equation of a line in 1 dimension is given by:

$$wx + b$$

- This generalizes in D dimensions to :

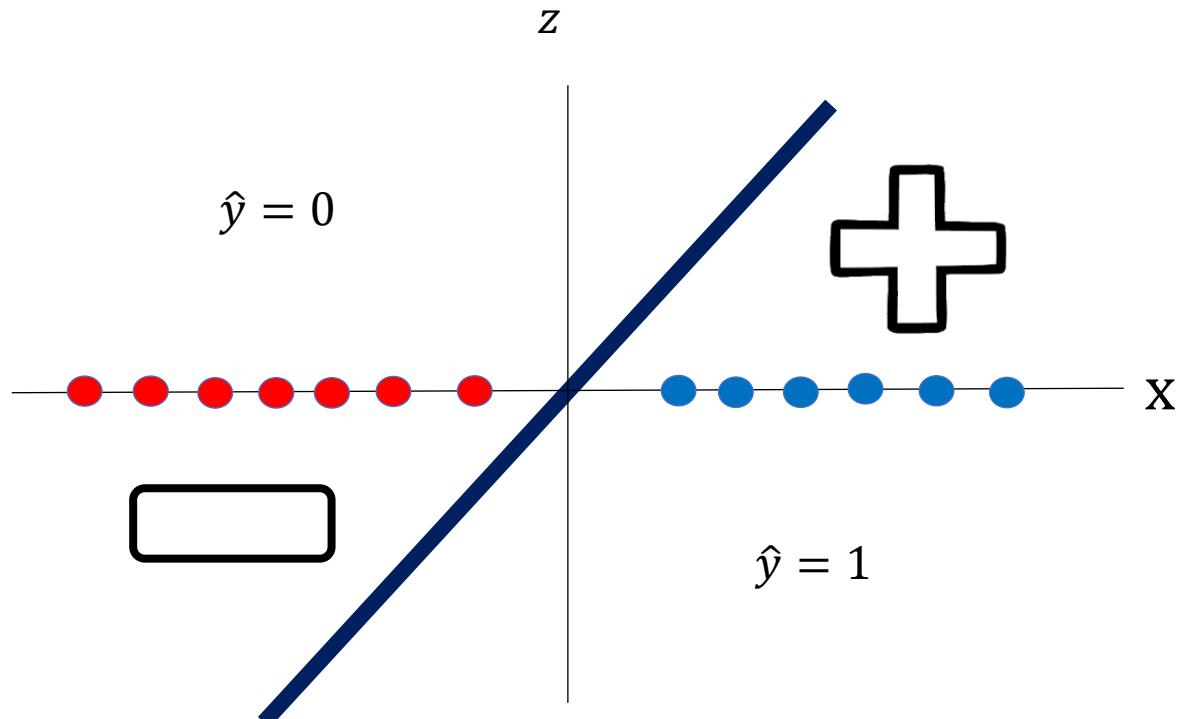
$$\mathbf{w}^T \mathbf{x} + b$$

- Let's see what happens for different values of \mathbf{x}

Linear classifier

- Consider the following data set
- If we can separate the data with a line, we can use that line to classify the sample

$$\begin{aligned} z &= w\mathbf{x} + b \\ &= \mathbf{1}\mathbf{x} - 1 \end{aligned}$$



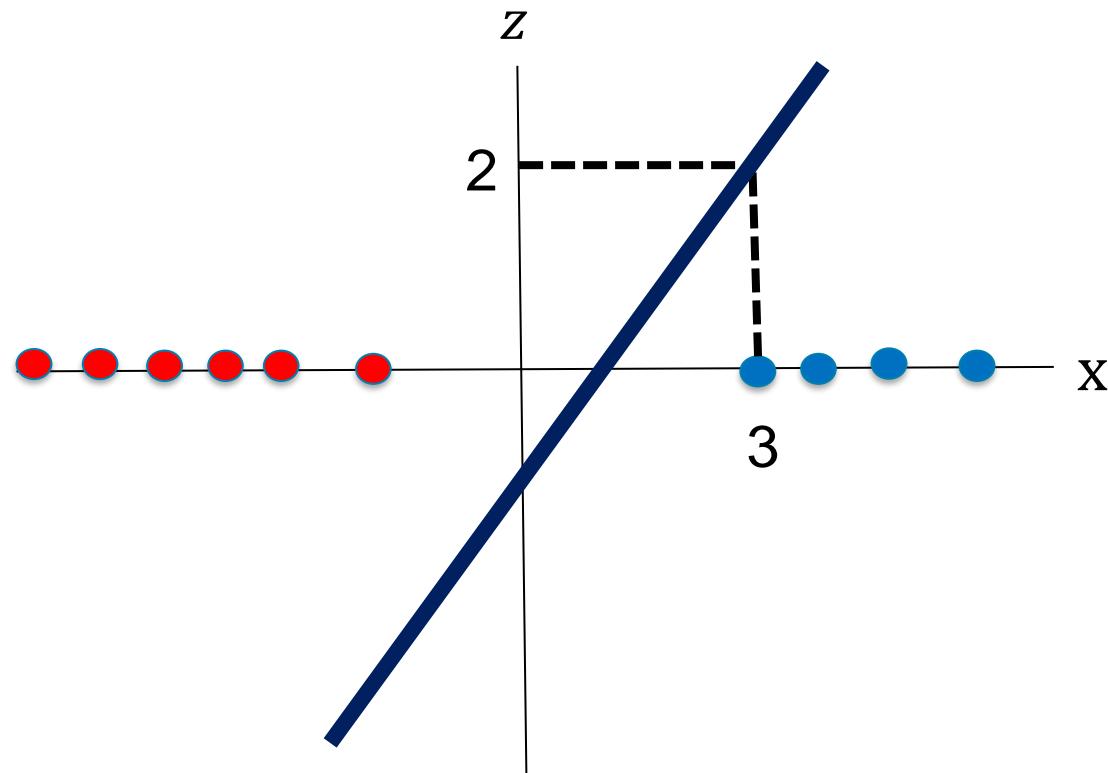
Linear classifier

- If x is on the left side of the line we get a positive number

$$z = 1x - 1$$

$$z = 1(3) - 1$$

$$= 2$$



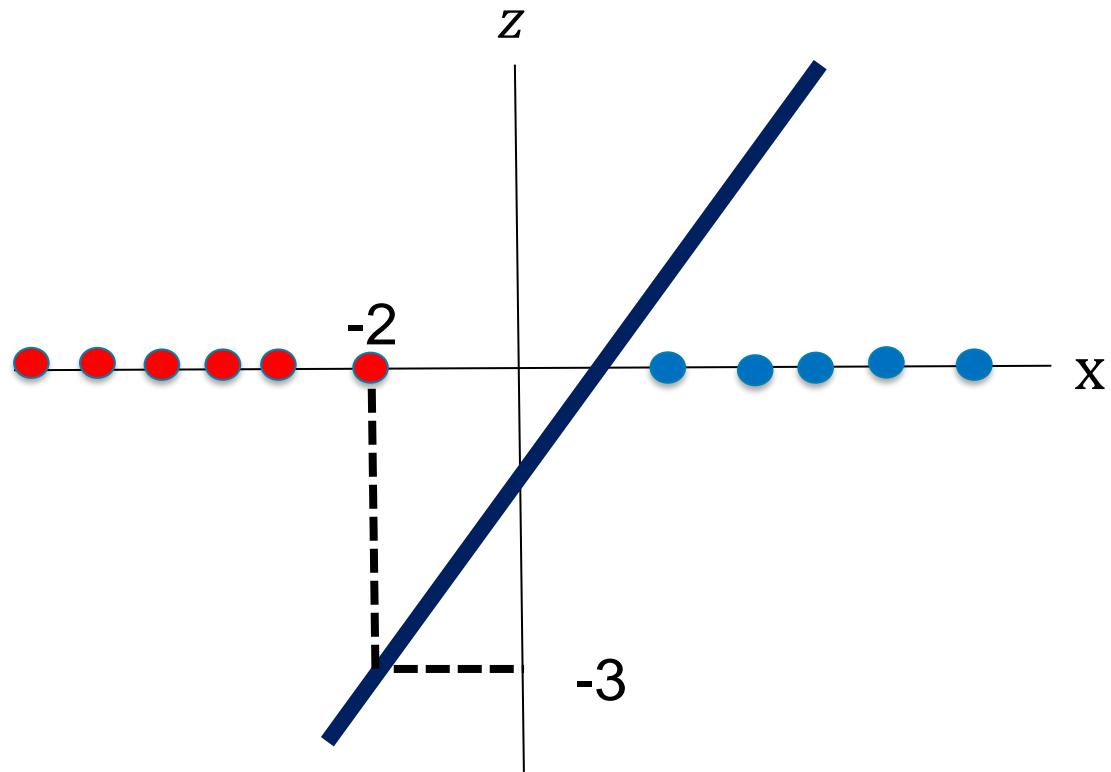
Linear classifier

- If x is on the right side of the line we get a positive number

$$z = 1x - 1$$

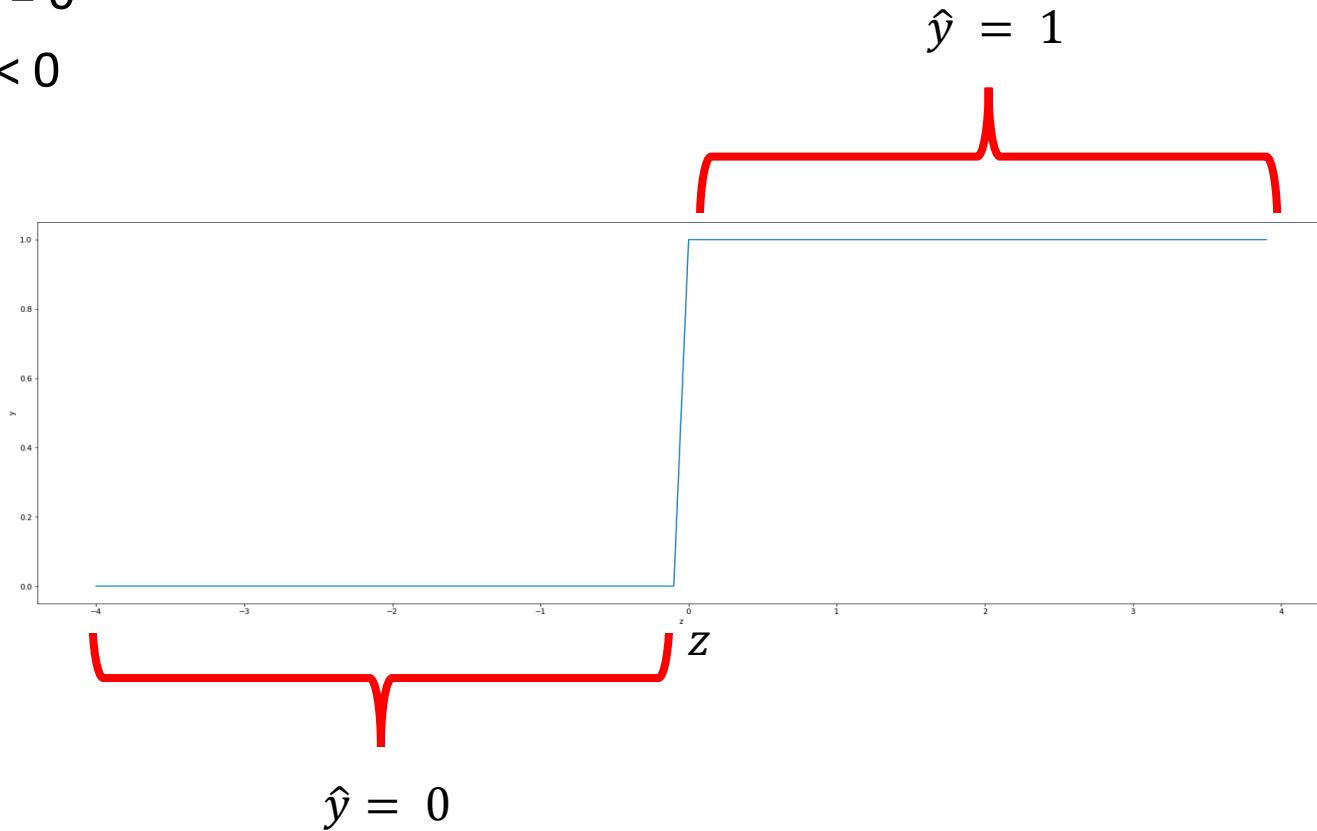
$$z = 1(-2) - 1$$

$$= -3$$



Linear classifiers

- If we use the line to calculate the class of a point, it always returns a positive or negative number, such as 3, -2, and so on.
- But, we need class 0 and class 1. How we can convert the numbers into 0 and 1?
- $\hat{y} = 1$, if $z \geq 0$
- $\hat{y} = 0$, if $z < 0$



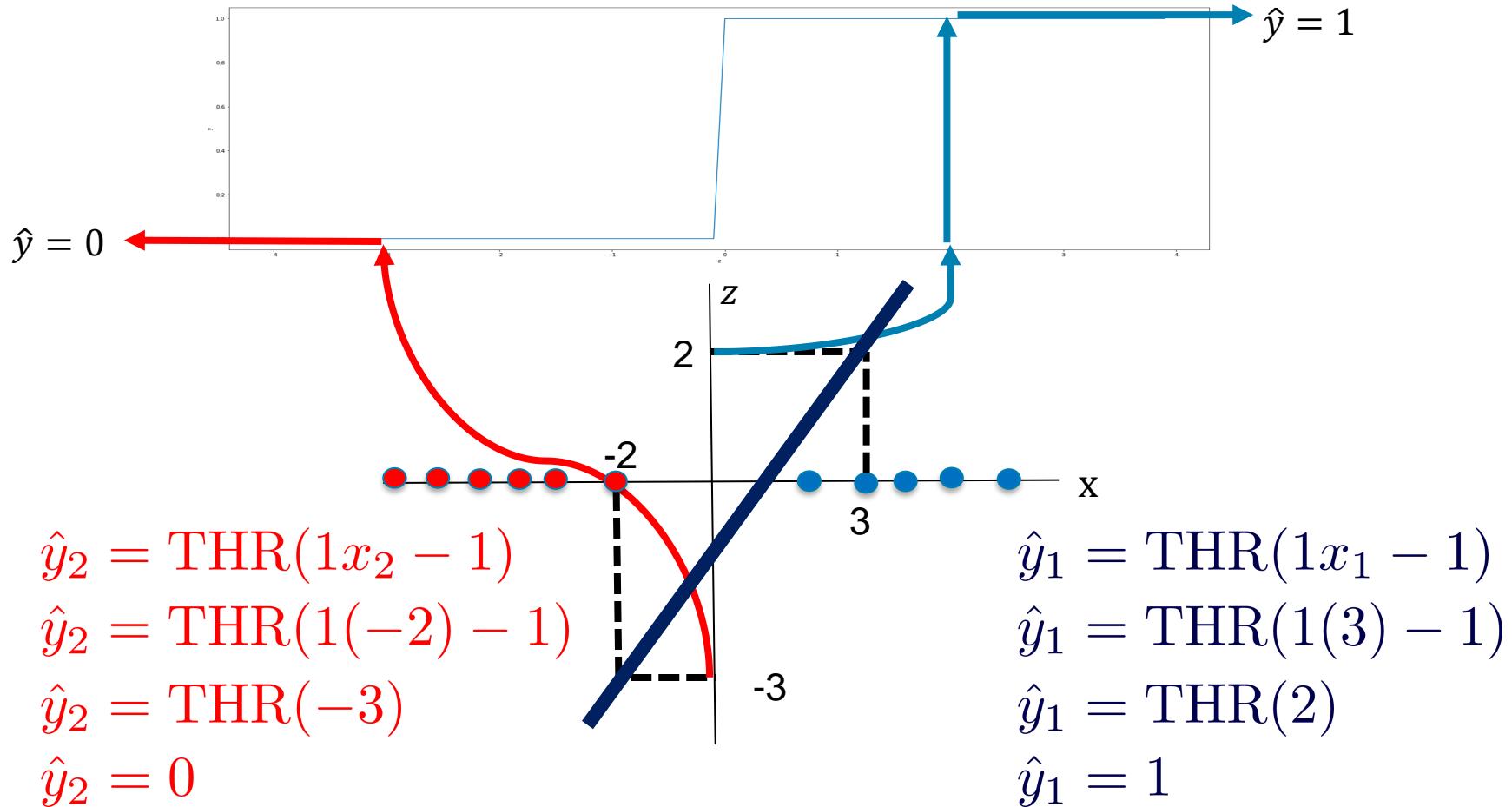
Linear classifier

$$\hat{y}_n = f(x_n)$$

parameters: w, b

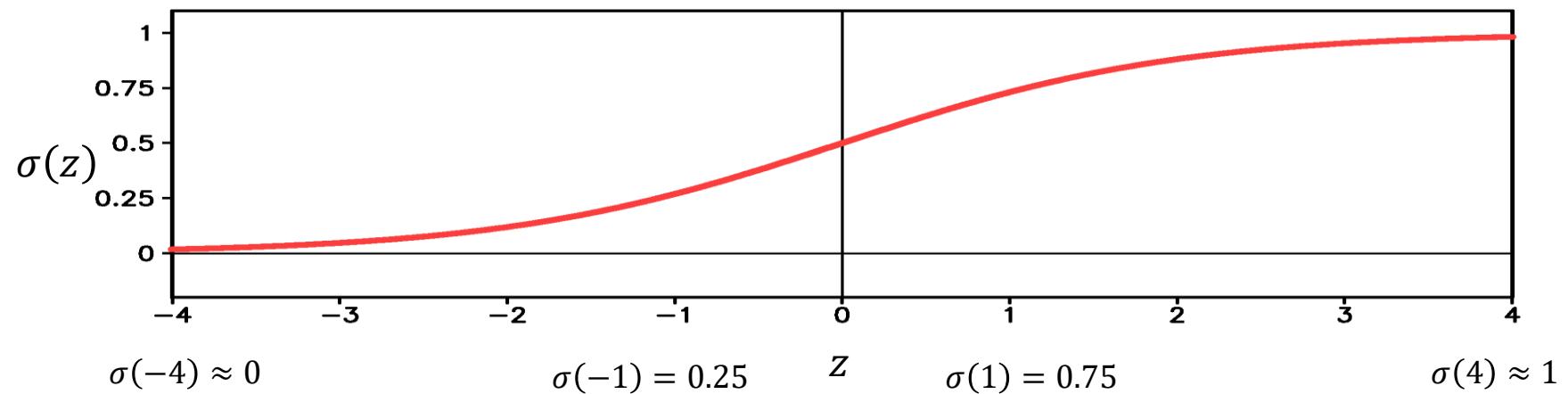
$$\hat{y}_n = \text{THR}(wx_n + b)$$

Linear classifier: Threshold Function



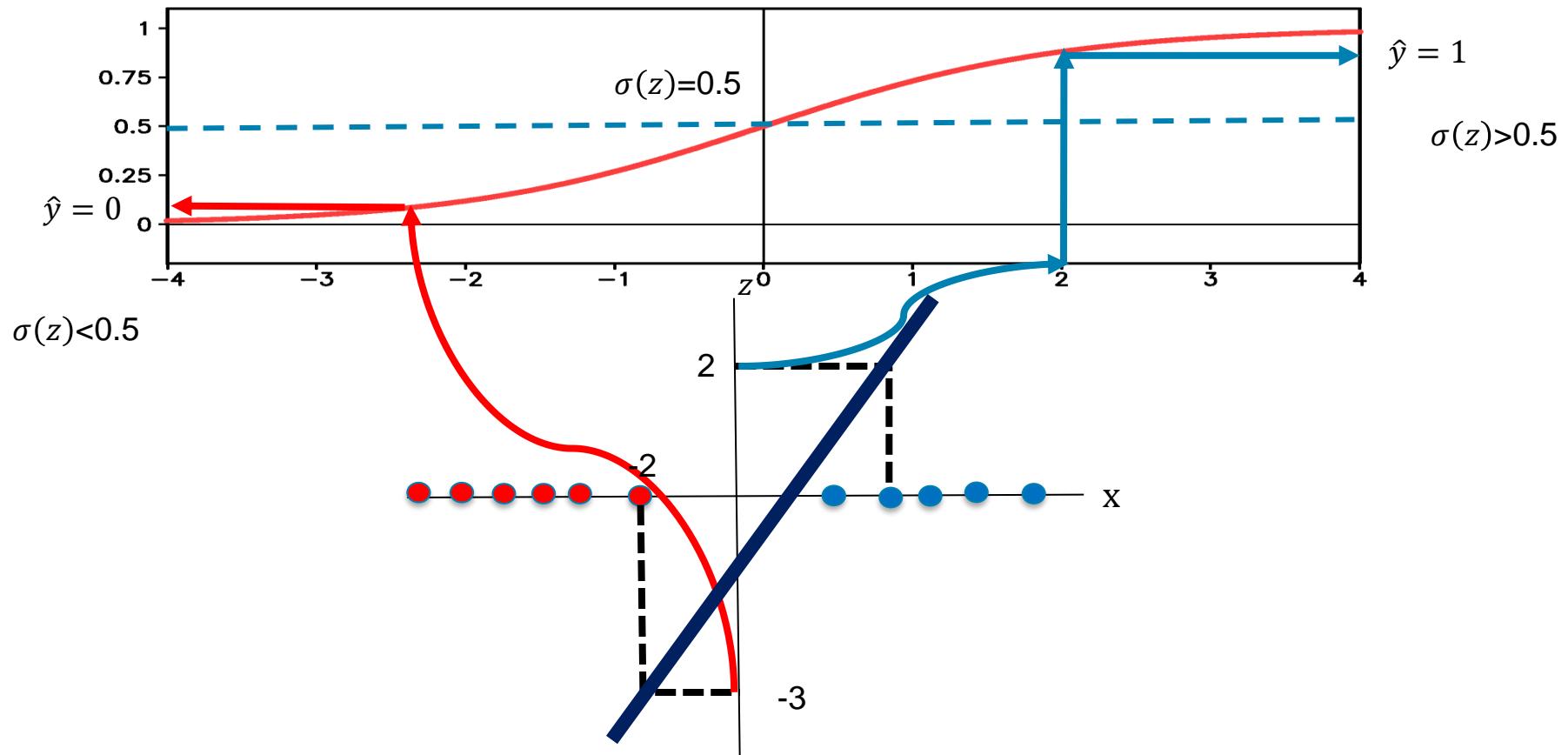
Linear classifier: Logistic Regression

Logistic function



$$\sigma(z) = \frac{1}{1 + e^{-z}}$$

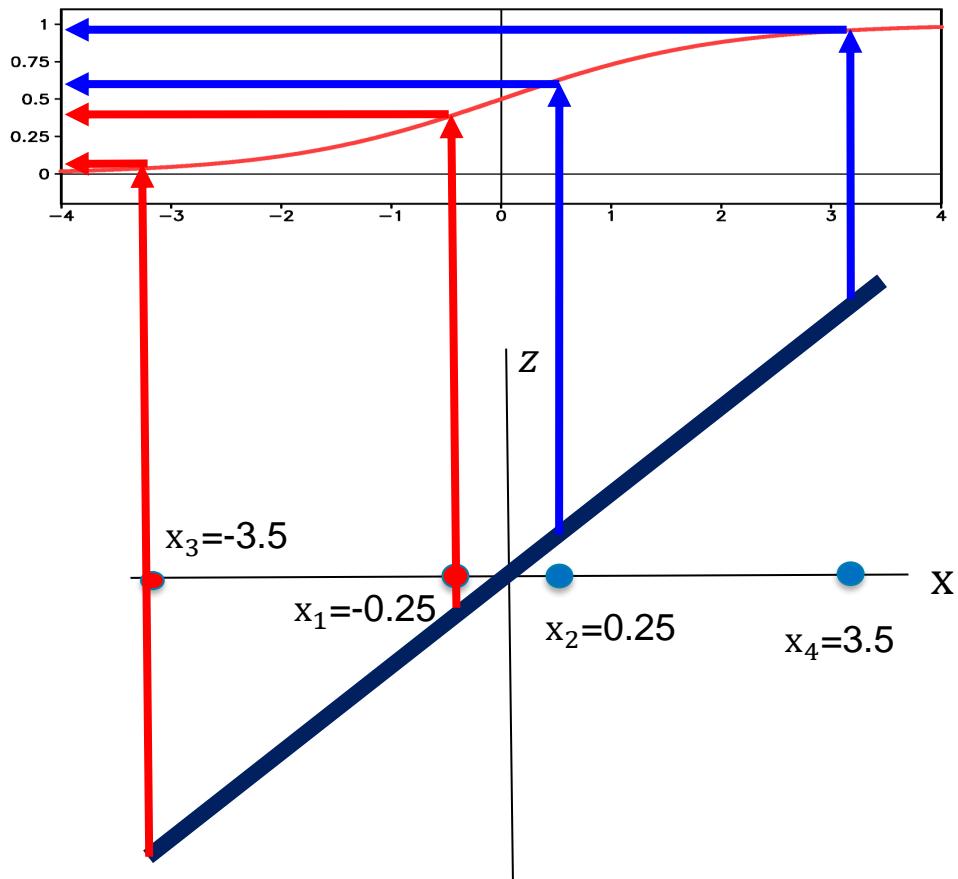
Linear classifier: Logistic Regression



Linear classifier: Logistic Regression

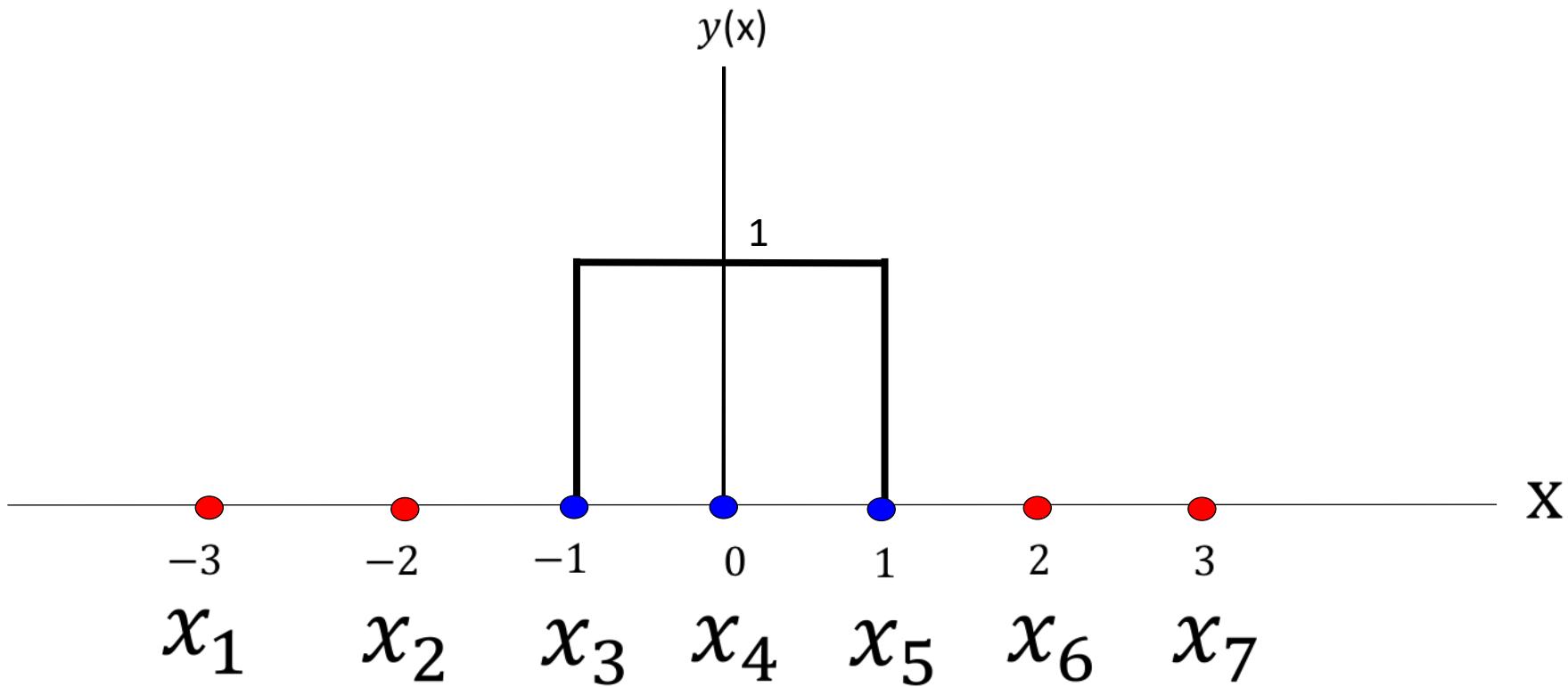
- Samples are near the line, we get values close to 0.5
- If they are far, we get values close to 0 or 1

x	y	$\sigma(z)$	\hat{y}
x_1	0	0.44	0
x_2	1	0.56	1
x_3	0	0.02	0
x_4	1	0.98	1



Linear classifiers vs. neural networks – example

- It is helpful to view the sample y as a decision function of x

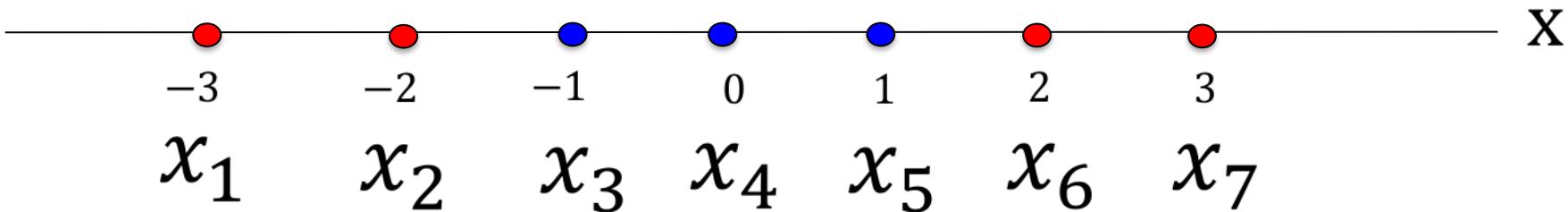


Linear classifiers vs. neural networks – example

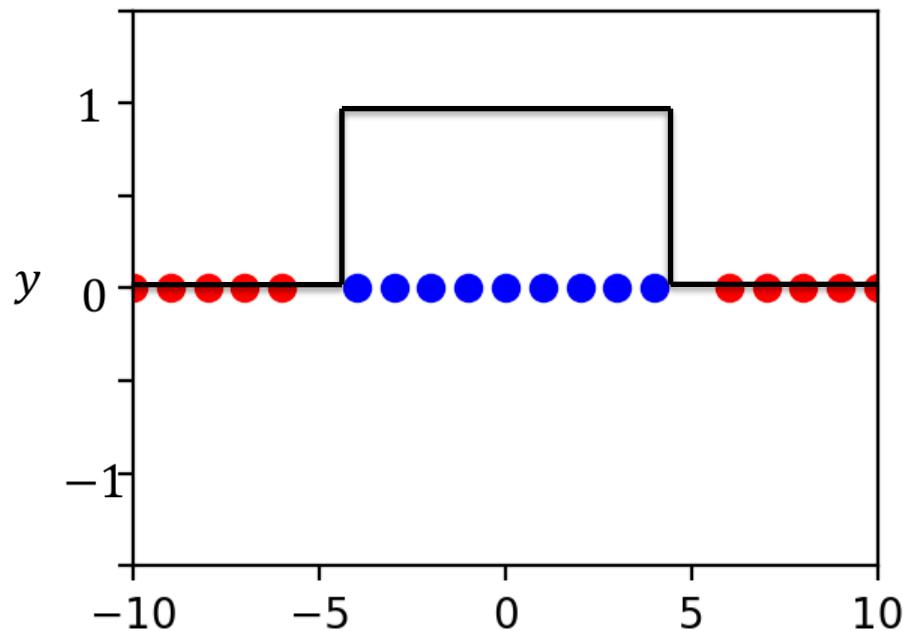
Colors are used to indicate the class

1. $y_1 = 0, x_1 = -3$
2. $y_2 = 0, x_2 = -2$
3. $y_3 = 1, x_3 = -1$
4. $y_4 = 1, x_4 = 0$
5. $y_5 = 1, x_5 = 1$
6. $y_6 = 0, x_6 = 2$
7. $y_7 = 0, x_7 = 3$

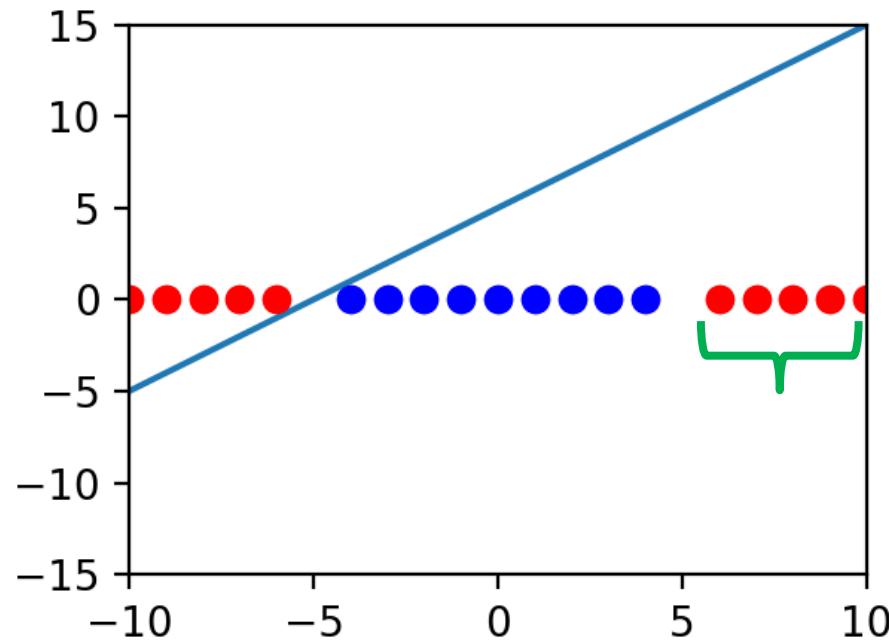
This dataset is not
linearly separable



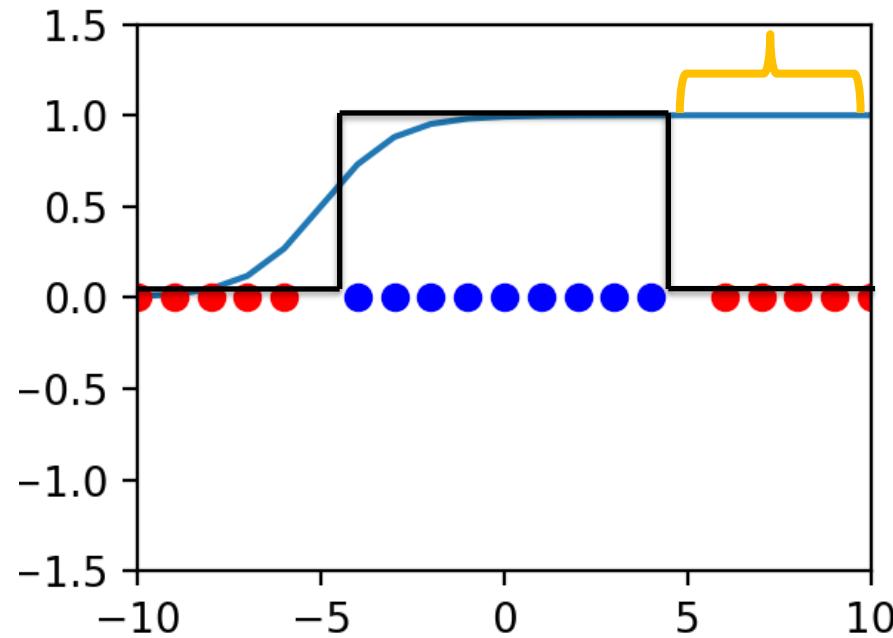
Linear classifiers vs. neural networks



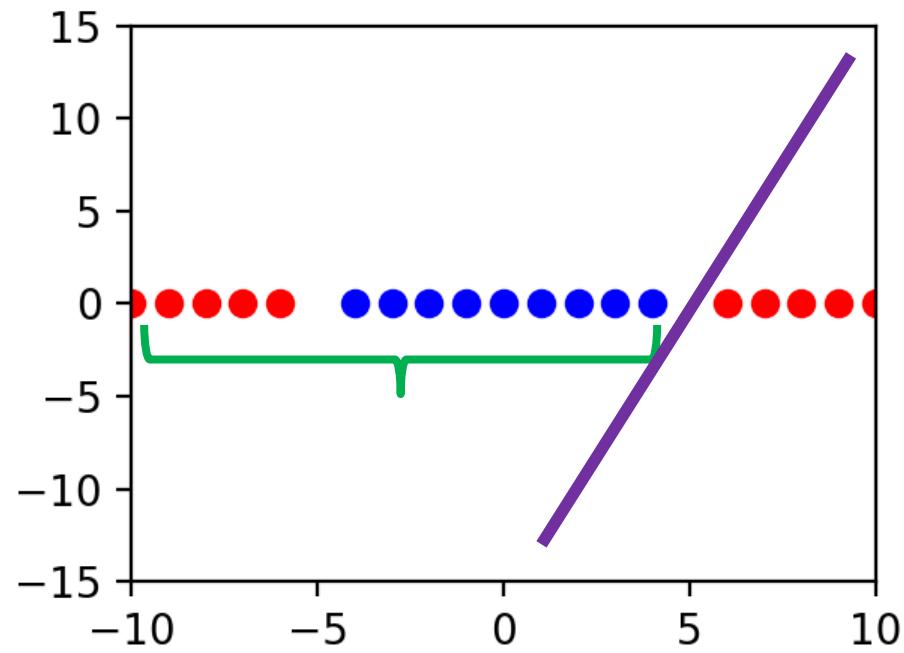
Neural networks



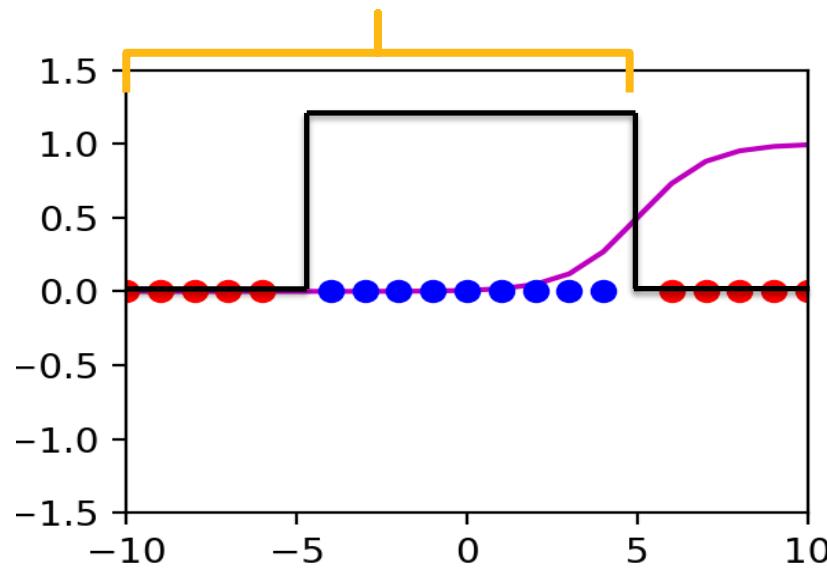
Neural networks



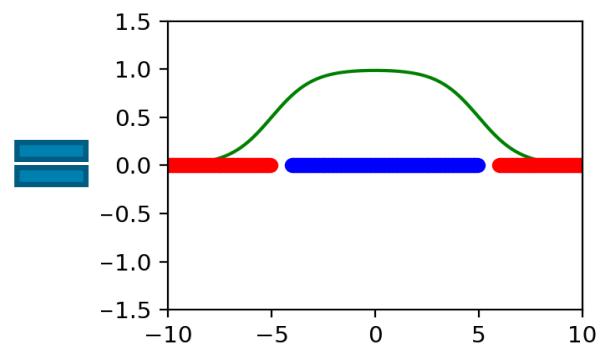
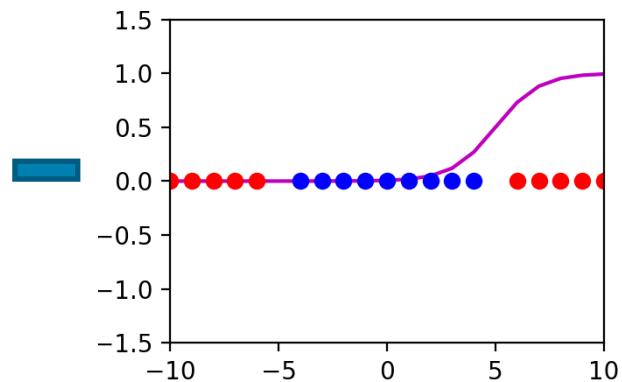
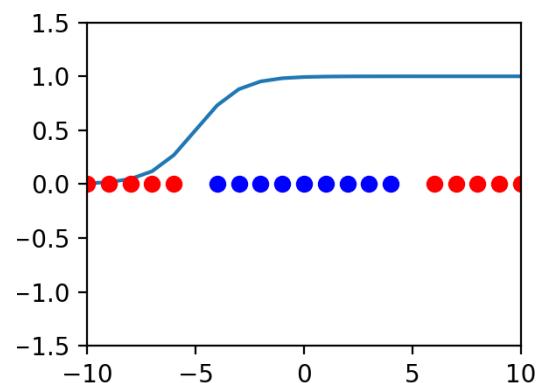
Neural networks



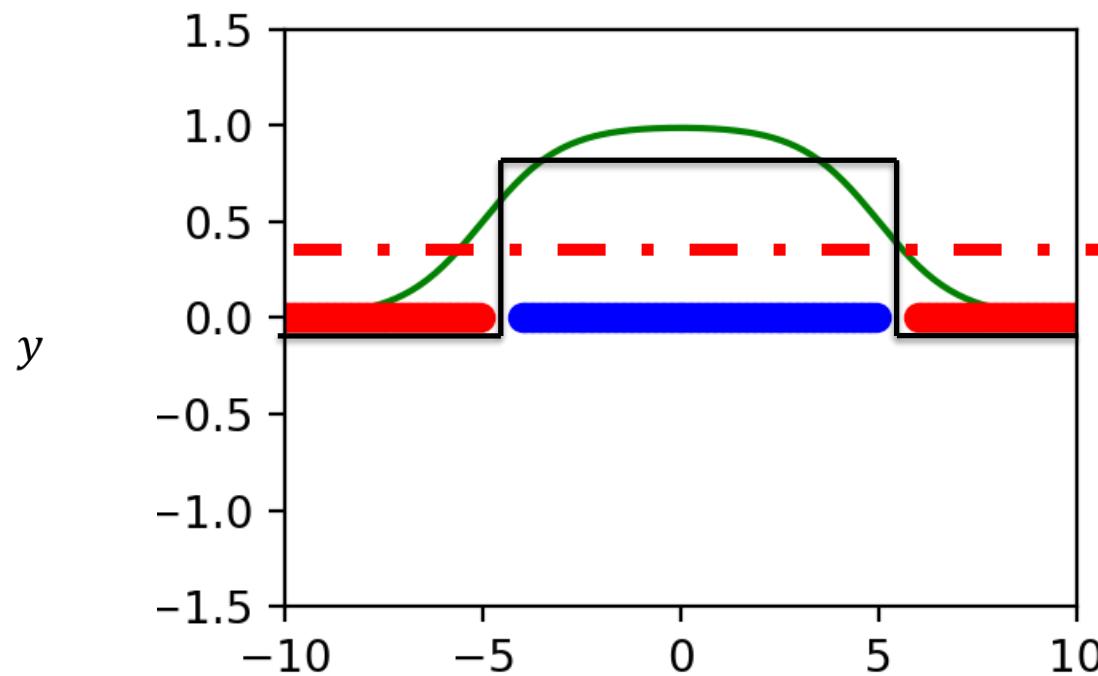
Neural networks

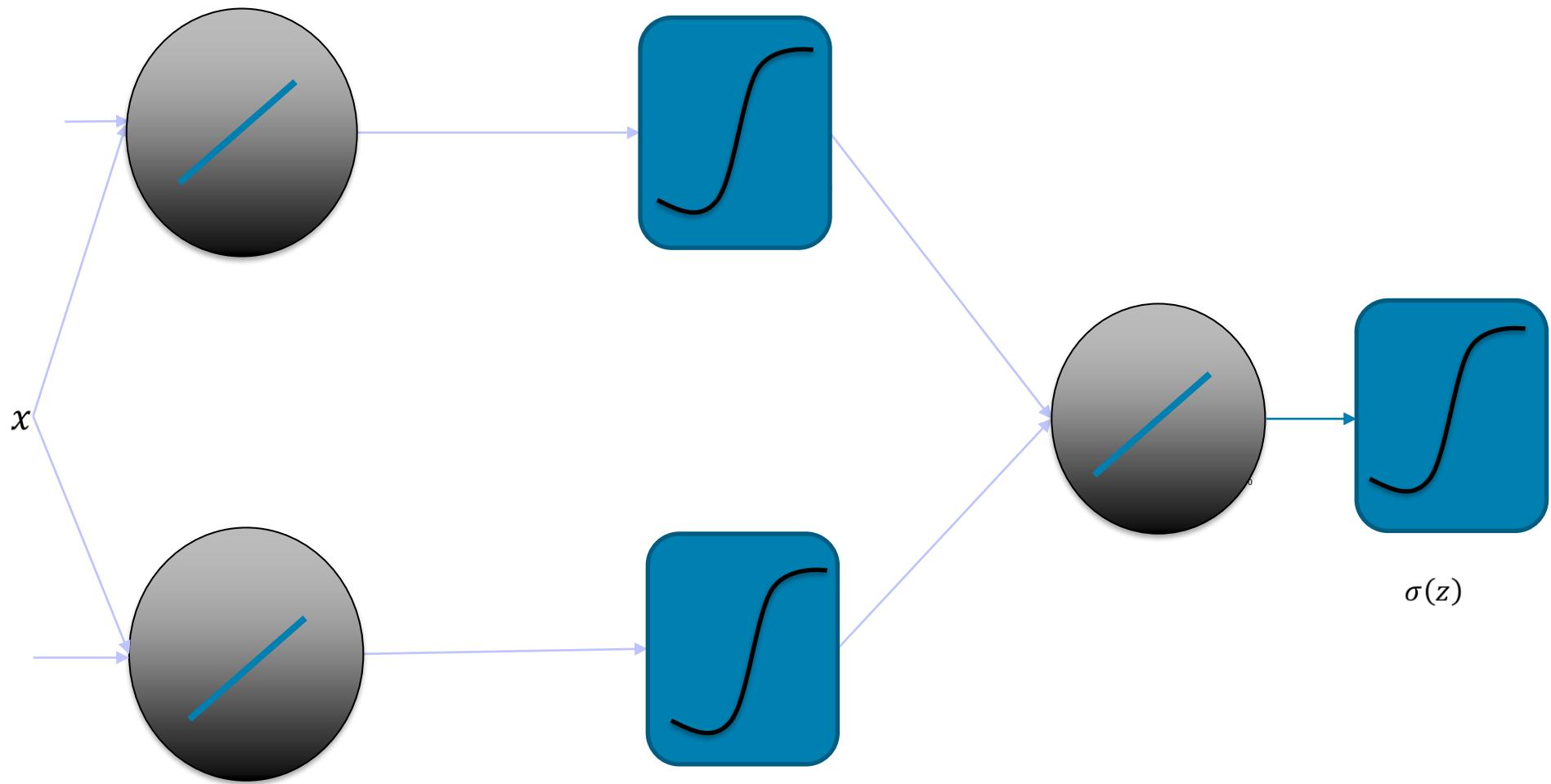


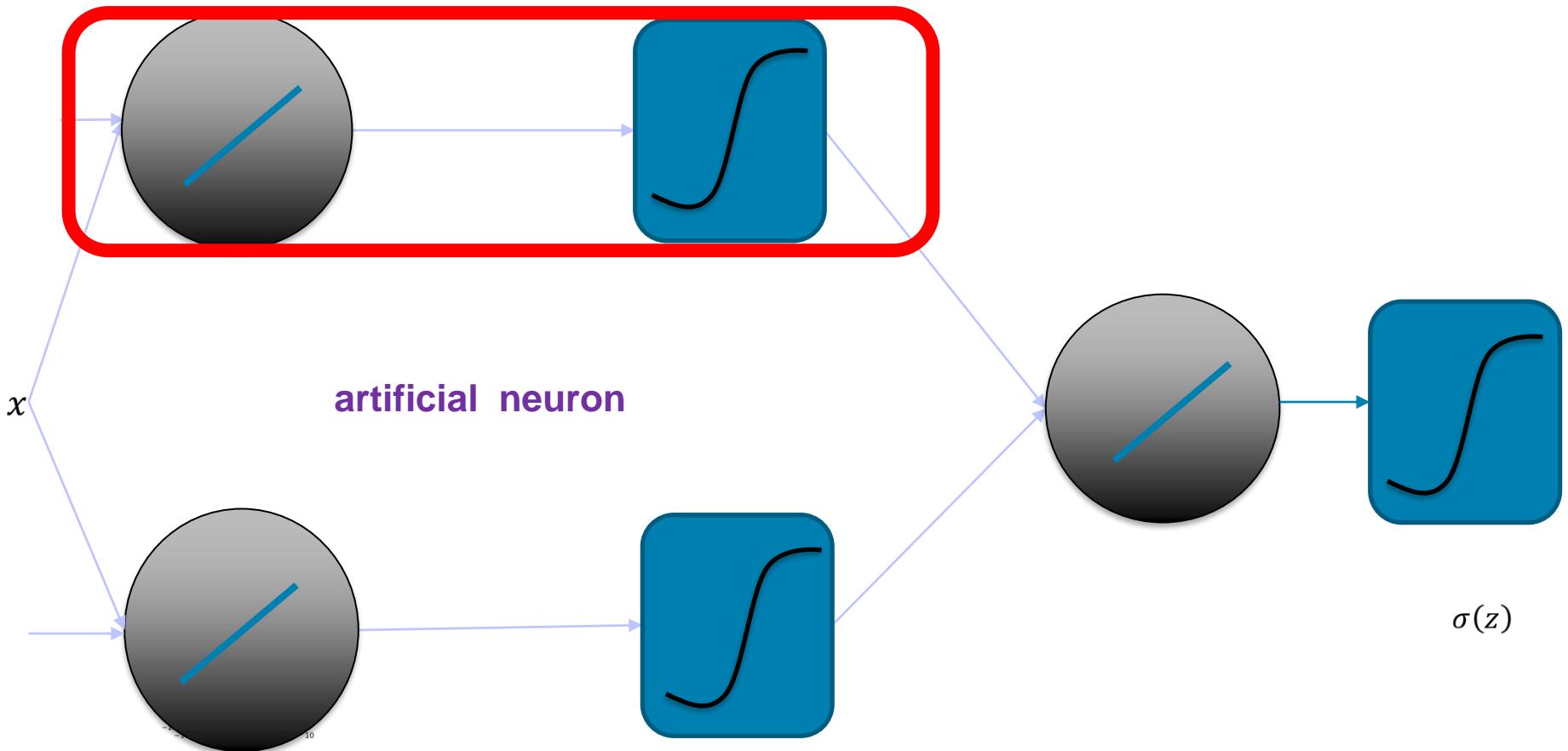
Neural networks

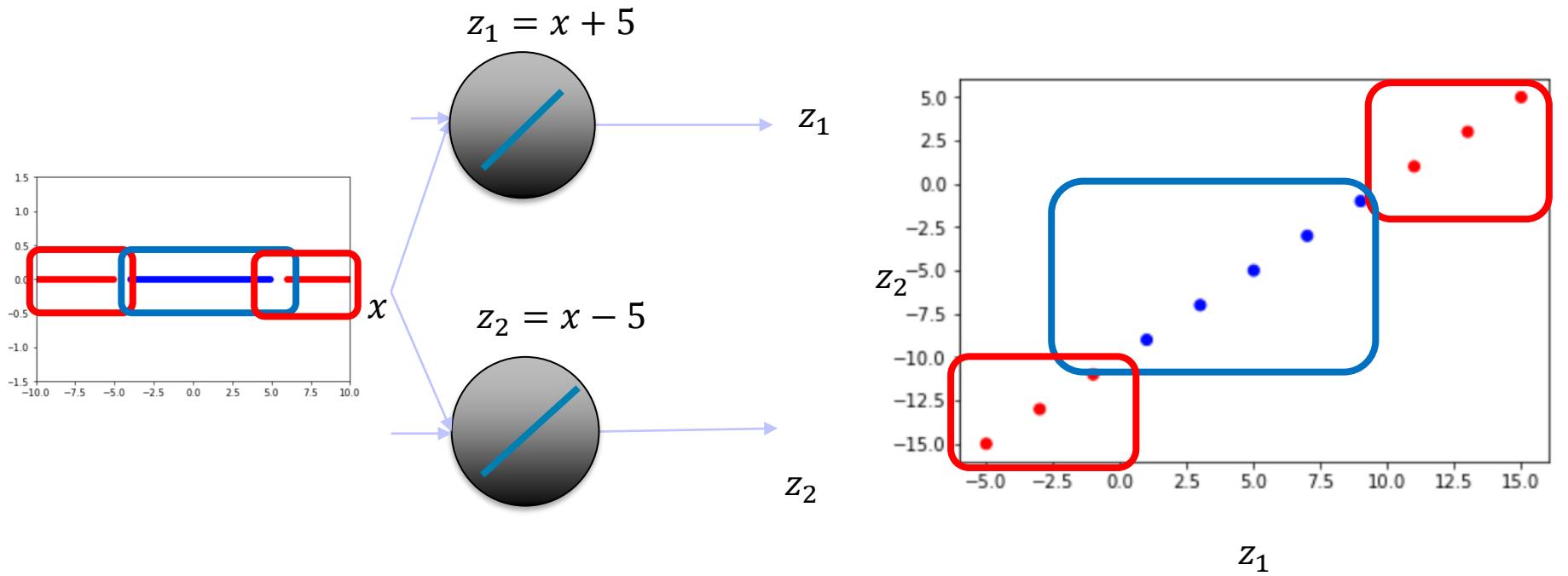


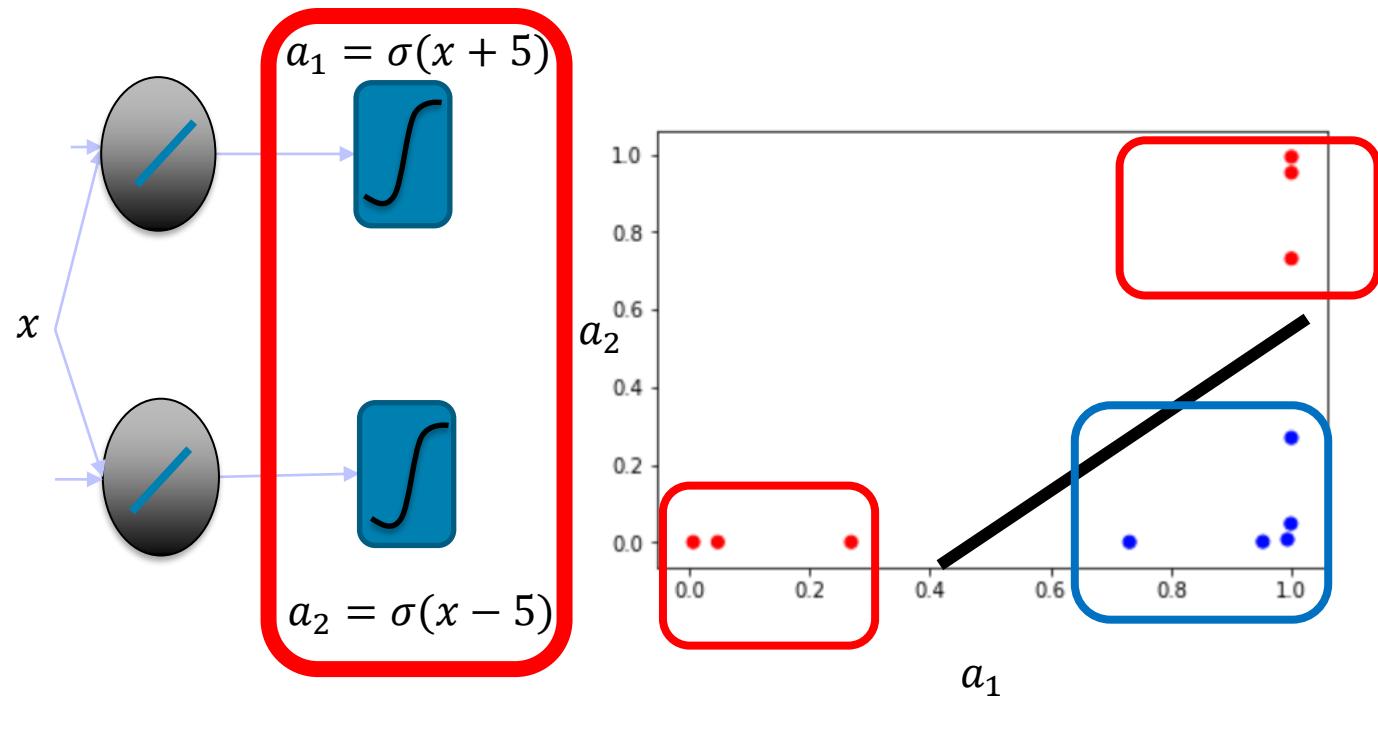
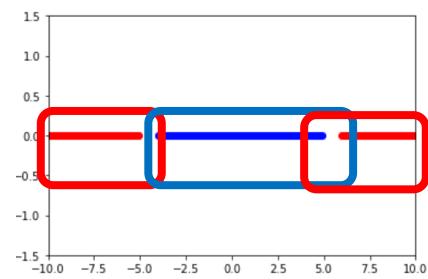
Neural networks



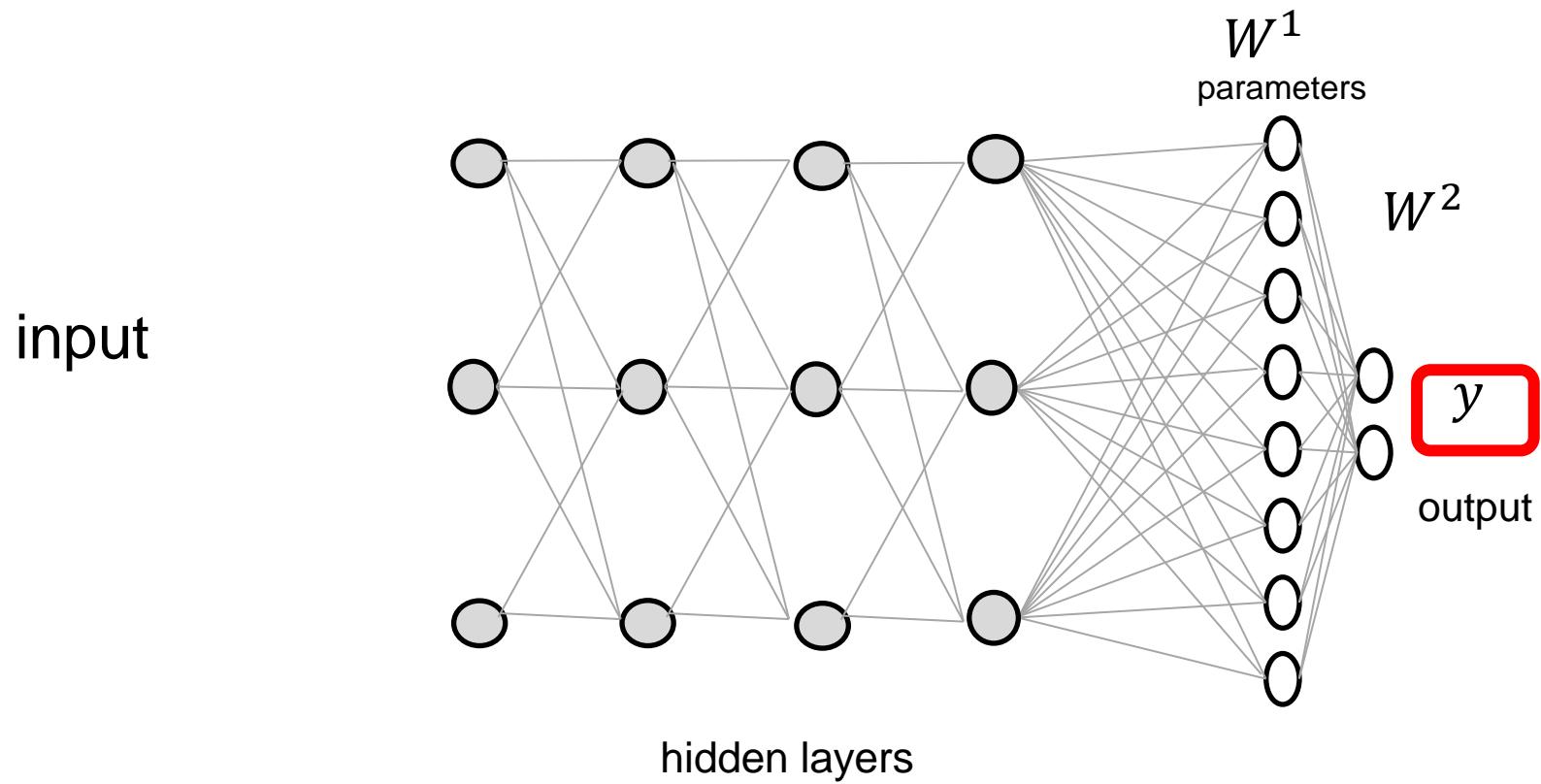








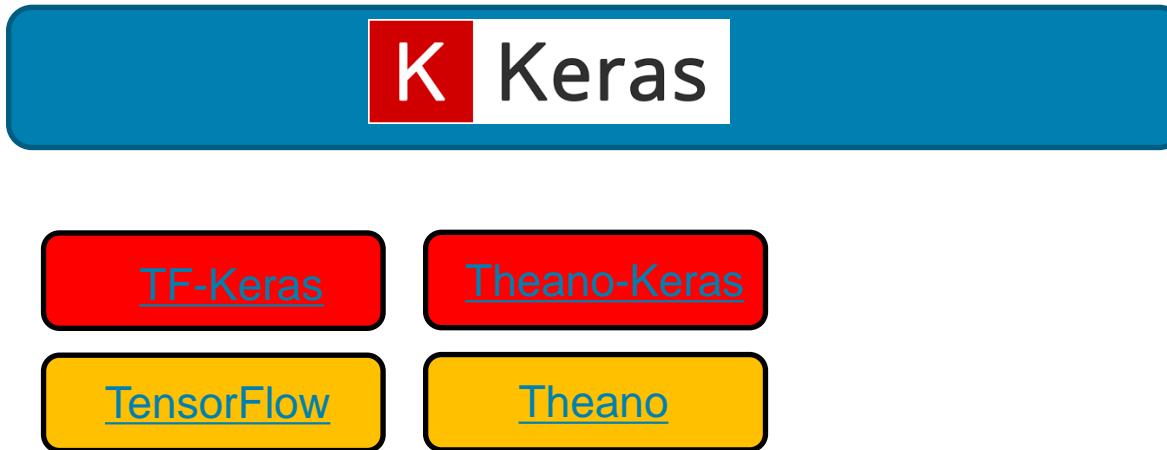
Neural networks and tensors



 TensorFlow	 PyTorch
TensorFlow is used in large scale production and deployment	Versatility of this Pythonic framework allows researchers to test out ideas with almost zero friction

*The Battle: TensorFlow vs. PyTorch
by Eitan Rosenzvraig*

Interface



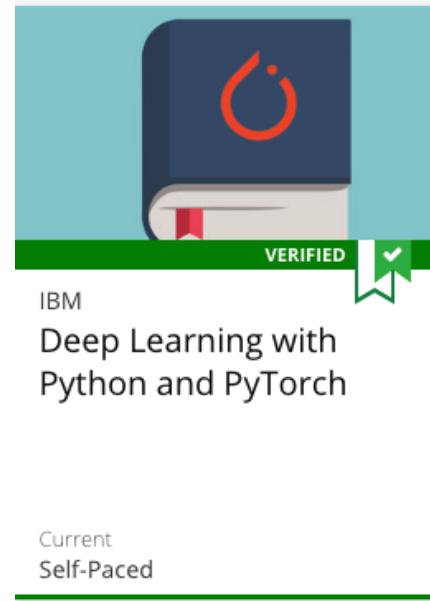
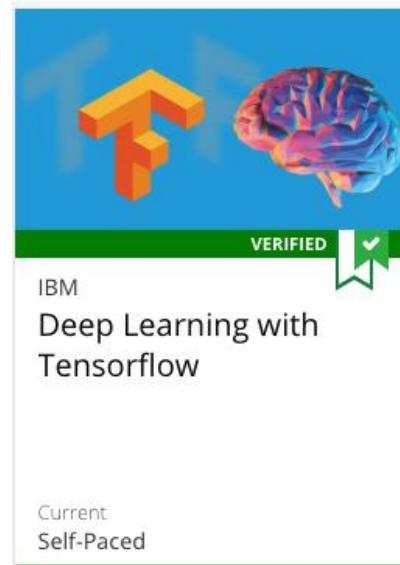
 Keras	 PyTorch
Keras may be easier to get into and experiment with standard layers, in a plug & play spirit	Low-level environment for experimentation, giving the user more freedom to write custom layers and look under the hood of numerical optimization tasks.

Keras or PyTorch as your first deep learning framework: by Piotr Migdal and Rafał Jakubanis

Deep Learning on EdX



Keras



<https://www.edx.org/professional-certificate/ibm-deep-learning>