

CIVIL-557

Decision-aid methodologies in transportation

Lecture 2: Deterministic methods and discrete metrics

Tim Hillel

**Transport and Mobility Laboratory TRANSP-OR
École Polytechnique Fédérale de Lausanne EPFL**

Motivation

How to **aid** informed **decisions** for **transportation** network management and investment?

Understand how people interact with a transport network

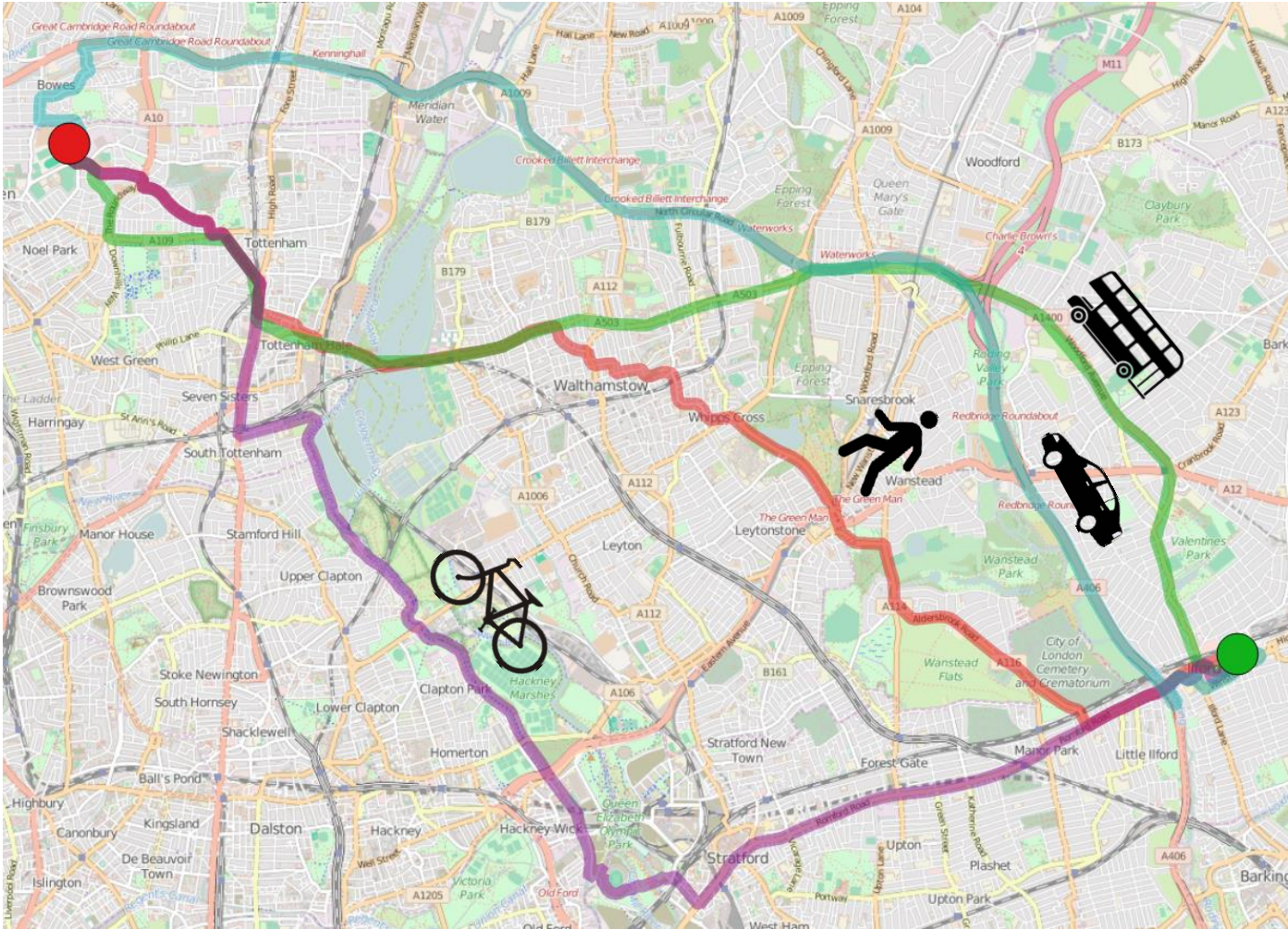
Case study



Transport for London



Mode choice



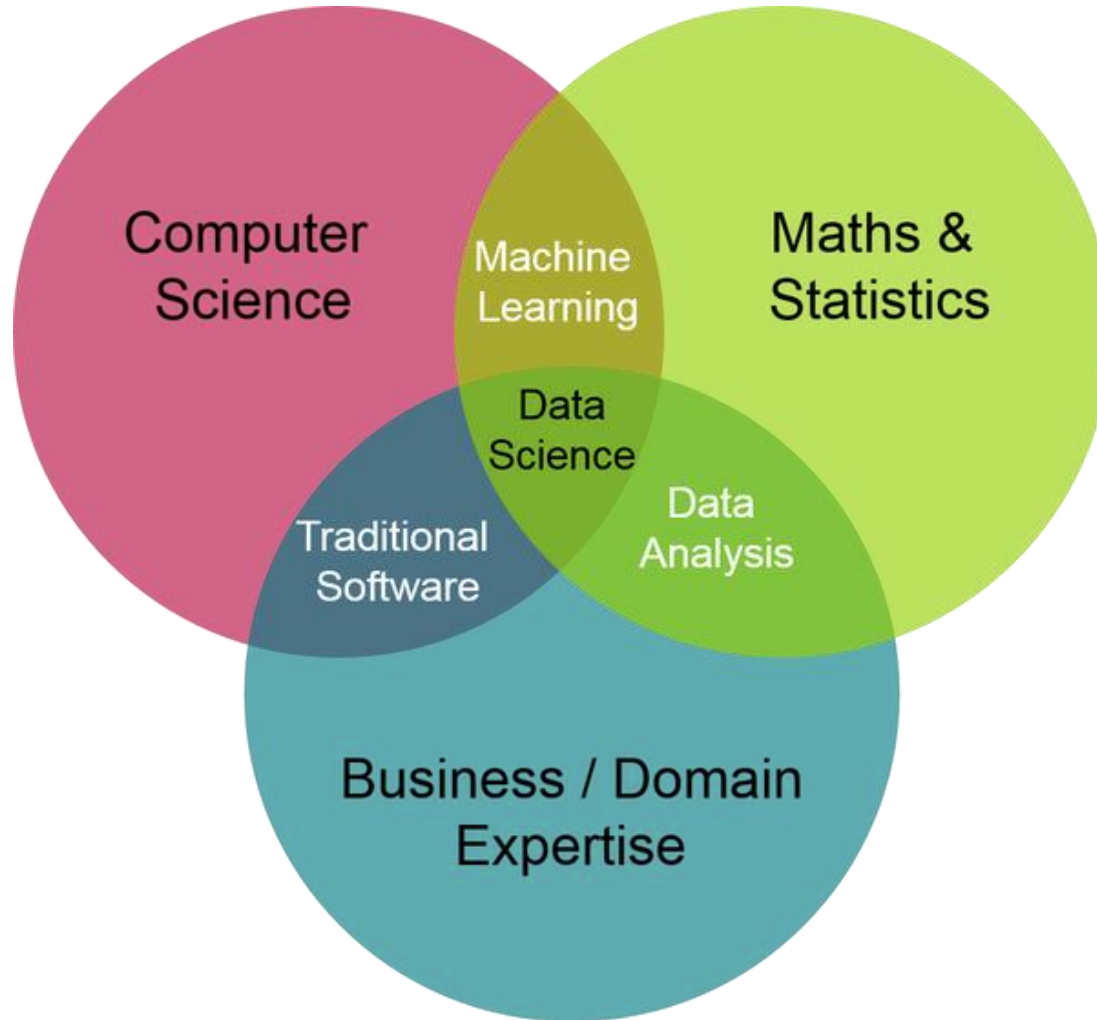
Last week

- Introduce and use tools
- Including:
 - Python
 - Jupyter
 - Numpy
 - Matplotlib
 - Pandas

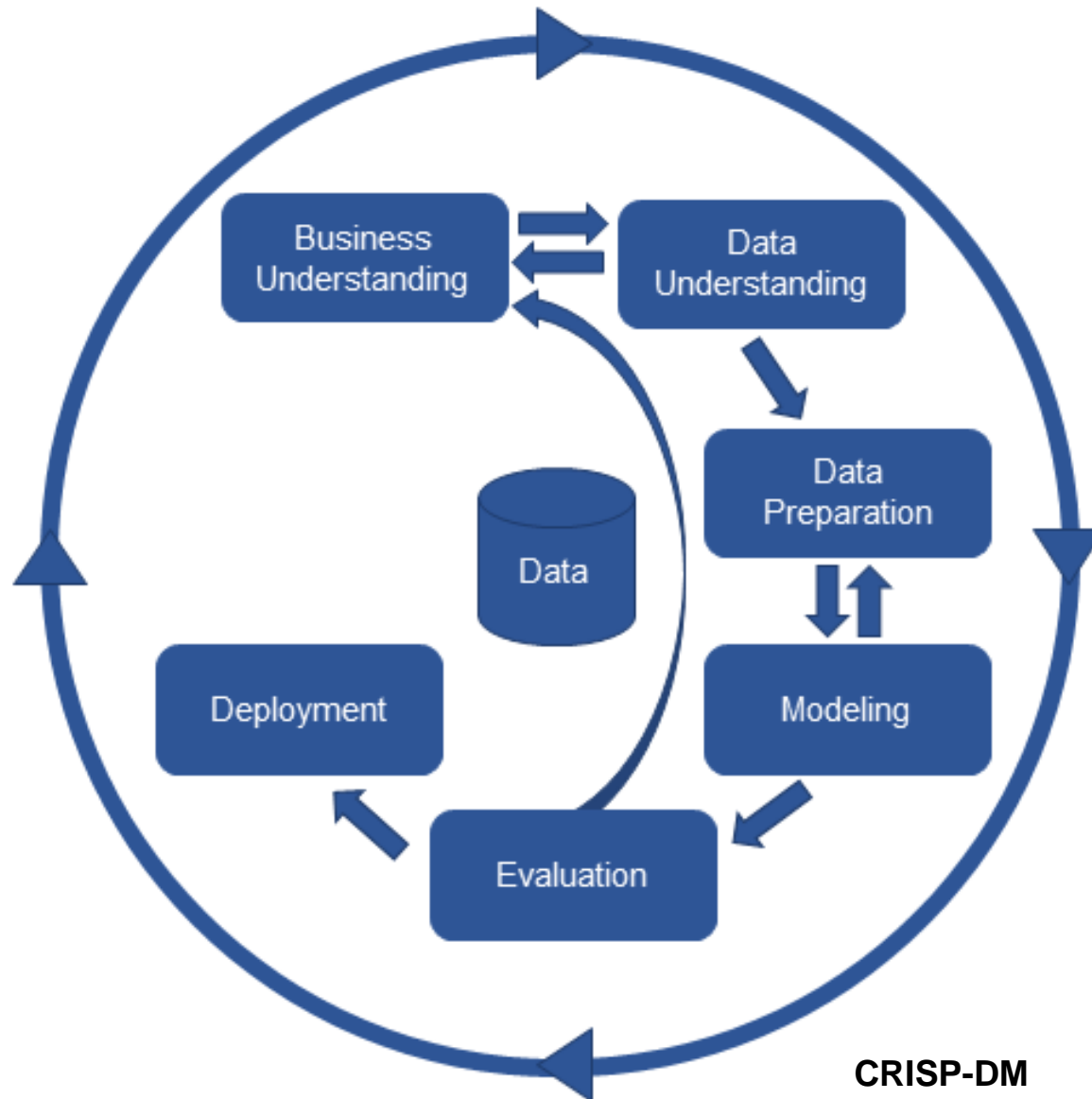
Today

- Data science process
- Dataset
- Deterministic methods
 - K-Nearest Neighbours (KNN)
 - Decision Tree (DT)
- Discrete metrics

Data science



Data science



CRISP-DM

Visualised last week in Pandas

Thoughts?

Tim Hillel, Mohammed Z E B Elshafie, and Ying Jin (2018). “Recreating Passenger Mode Choice-Sets for Transport Simulation: A Case Study of London, UK”. in: *Proceedings of the Institution of Civil Engineers - Smart Infrastructure and Construction* 171.1, pp. 29–42

Dataset building process



Historical trip data

London Travel Demand Survey (LTDS)

- Annual rolling **household travel survey**
- Each household member fills in **trip diary**

3 years of data (2012/13-2014/15)

- ~130,000 trips

Dataset building process

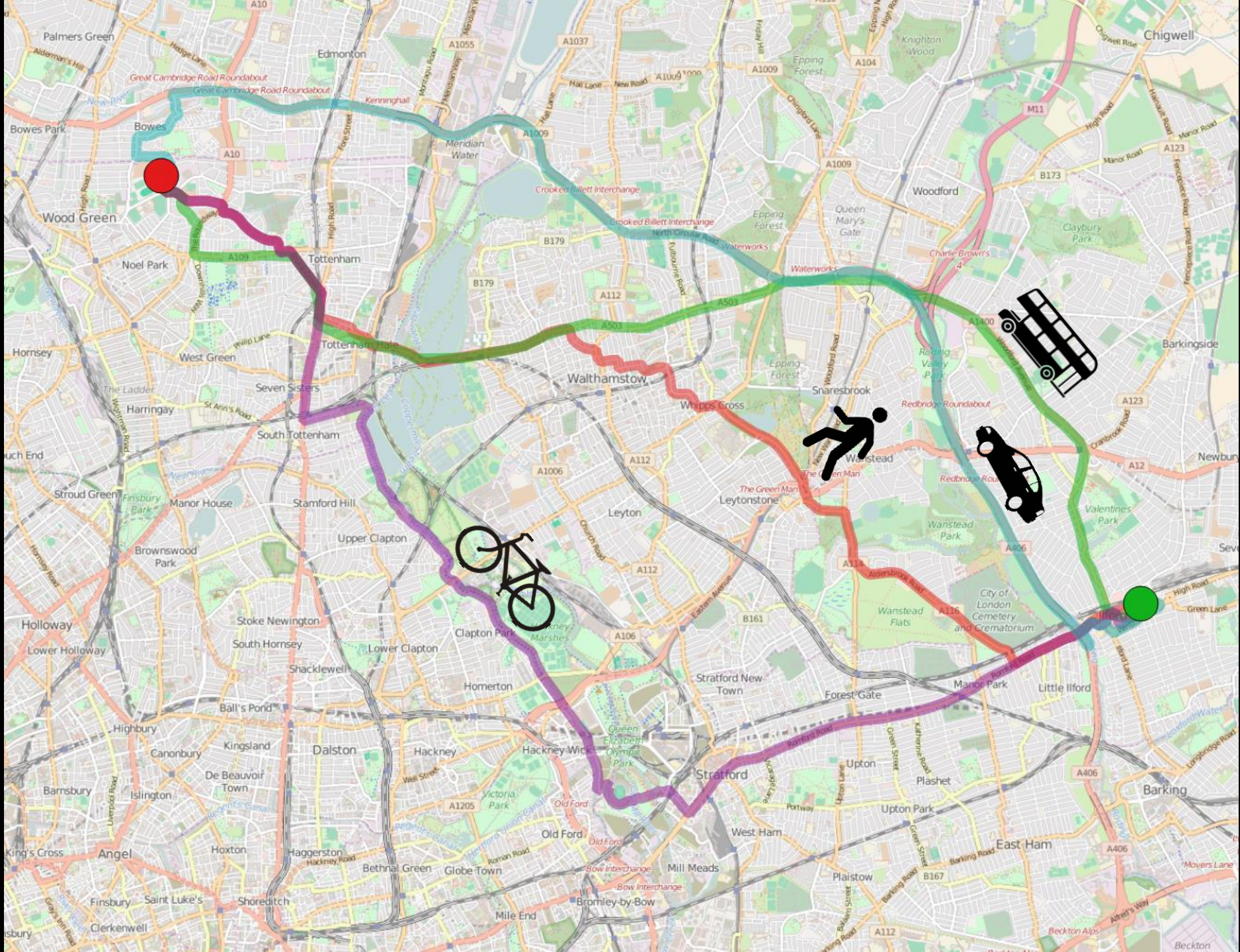


**Historical
trip data**

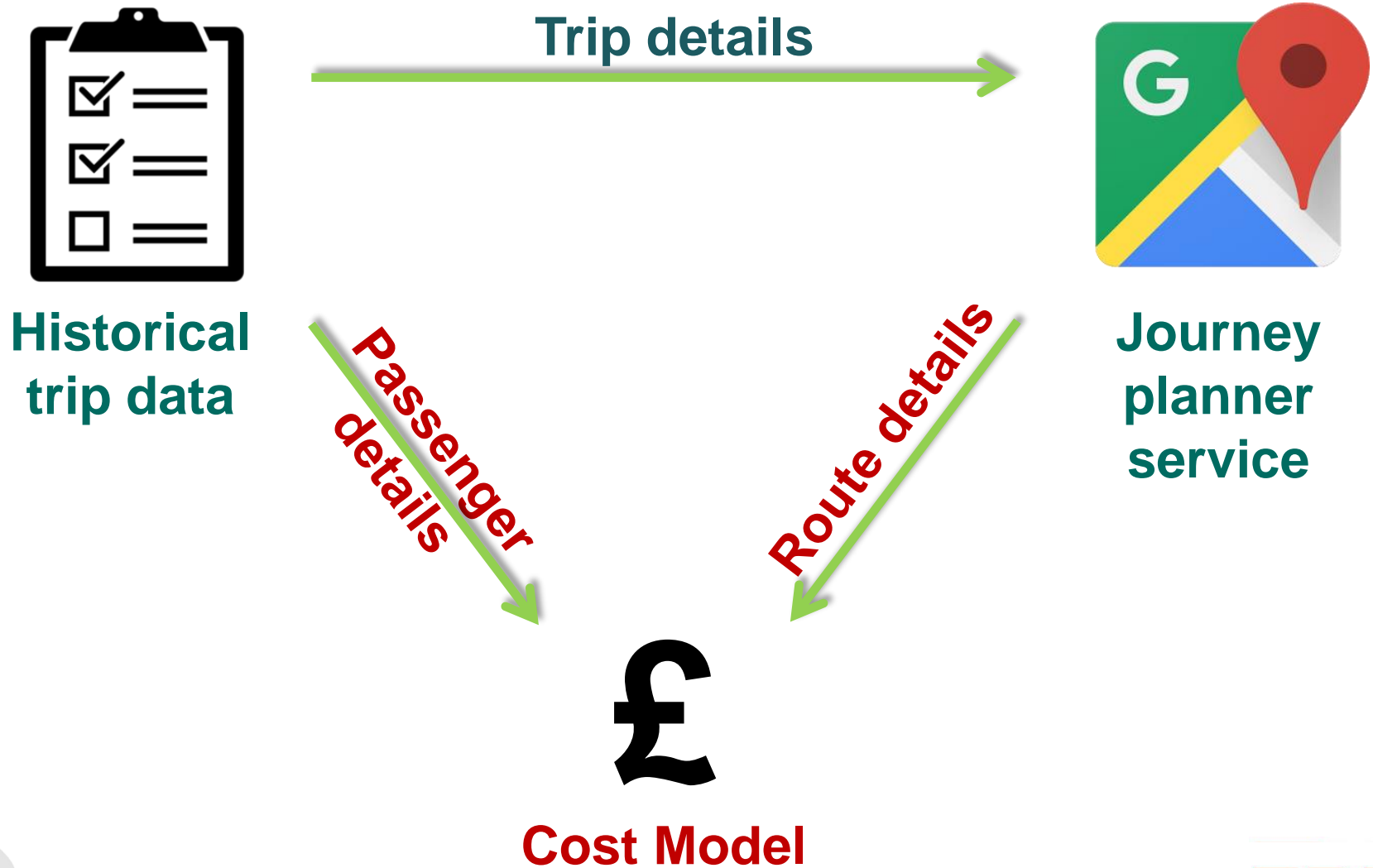
Trip details



**Journey
planner
service**



Dataset building process



Dataset

- 82,350 trips and alternatives for each major mode
 - walking, cycling, public transport (combined), rail, bus, driving:
 - Path
 - Duration
 - Traffic variability
 - Journey purpose
 - Socio-economic class
 - Departure time

Dataset

Feature vector  **Class label**

- Start time
 - Journey purpose
 - Vehicle ownership
 - Fare type
 - Alternative specific constants:
 - Duration
 - Cost
 - Typical traffic
 - Congestion charge
- Mode taken!

Machine learning

- Supervised learning
- Unsupervised learning
- Reinforcement learning
- ...and more:
 - Semi-supervised learning
 - Generative models
 - Active learning etc.

Supervised learning

- Regression – **Continuous value**
- Classification - **Discrete class**

Supervised classification

- Predict, for a feature vector x , the class label y
- Start with **binary** case
 - Was a trip made by car or not?
- **Discrete** classification
 - Predict 1 or 0

General approach

Fit model on some data (**train set**):

□ N instances (rows), comprising of:

- Features x (columns)
- Labels y

□ Predict for **unseen** data:

- y (unknown) from x (known)

General approach

$$y = f(x)$$

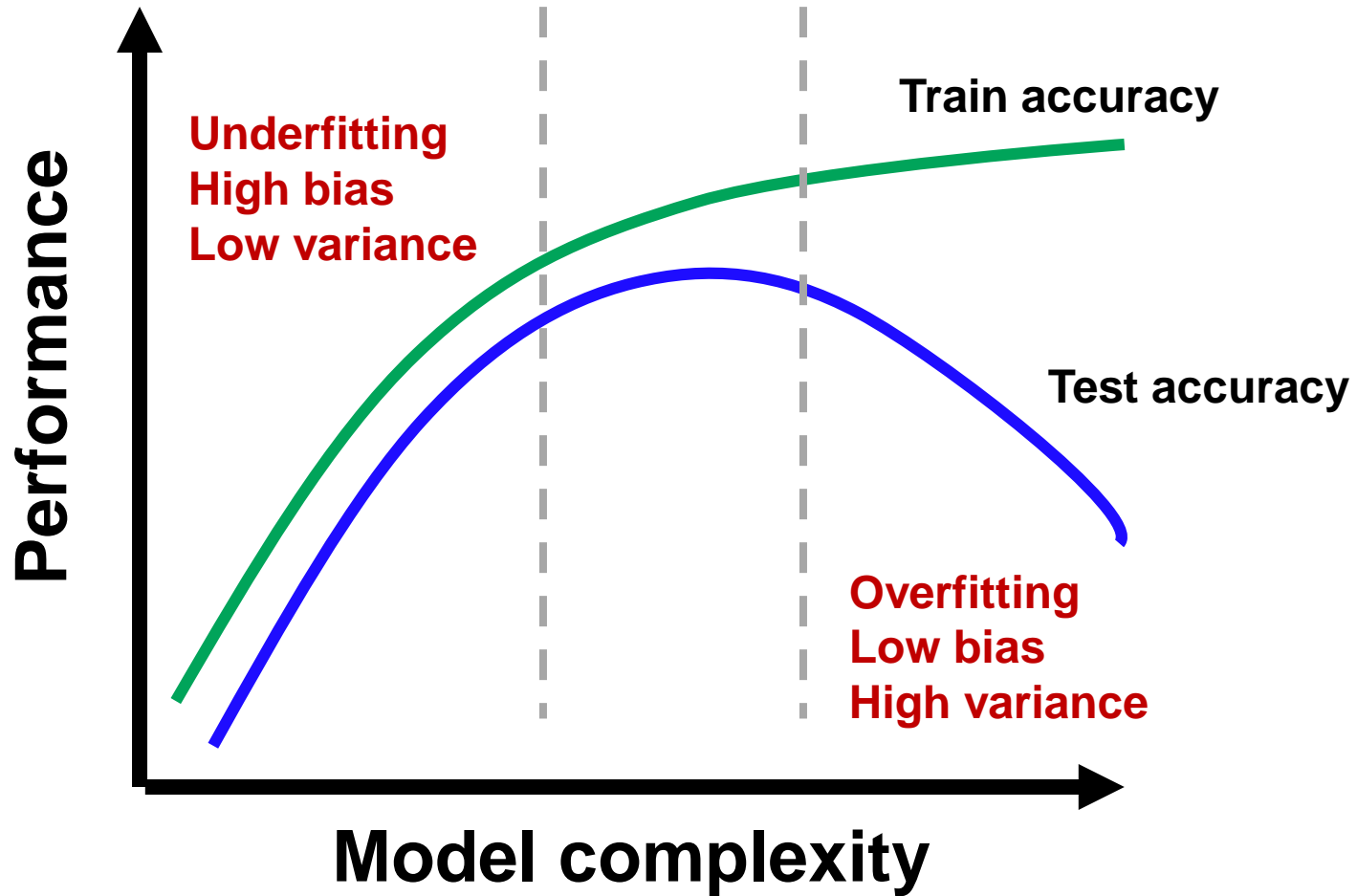
$$y = f'(x)$$

Model evaluation

How do we predict how model will **perform** on
unseen data?

Evaluate the model on an unseen **test set!**

Bias-variance trade-off



Train-test sampling

How do we define the **train** and **test** sets?

For now, we will use **random sampling** with 80:20
train:test split

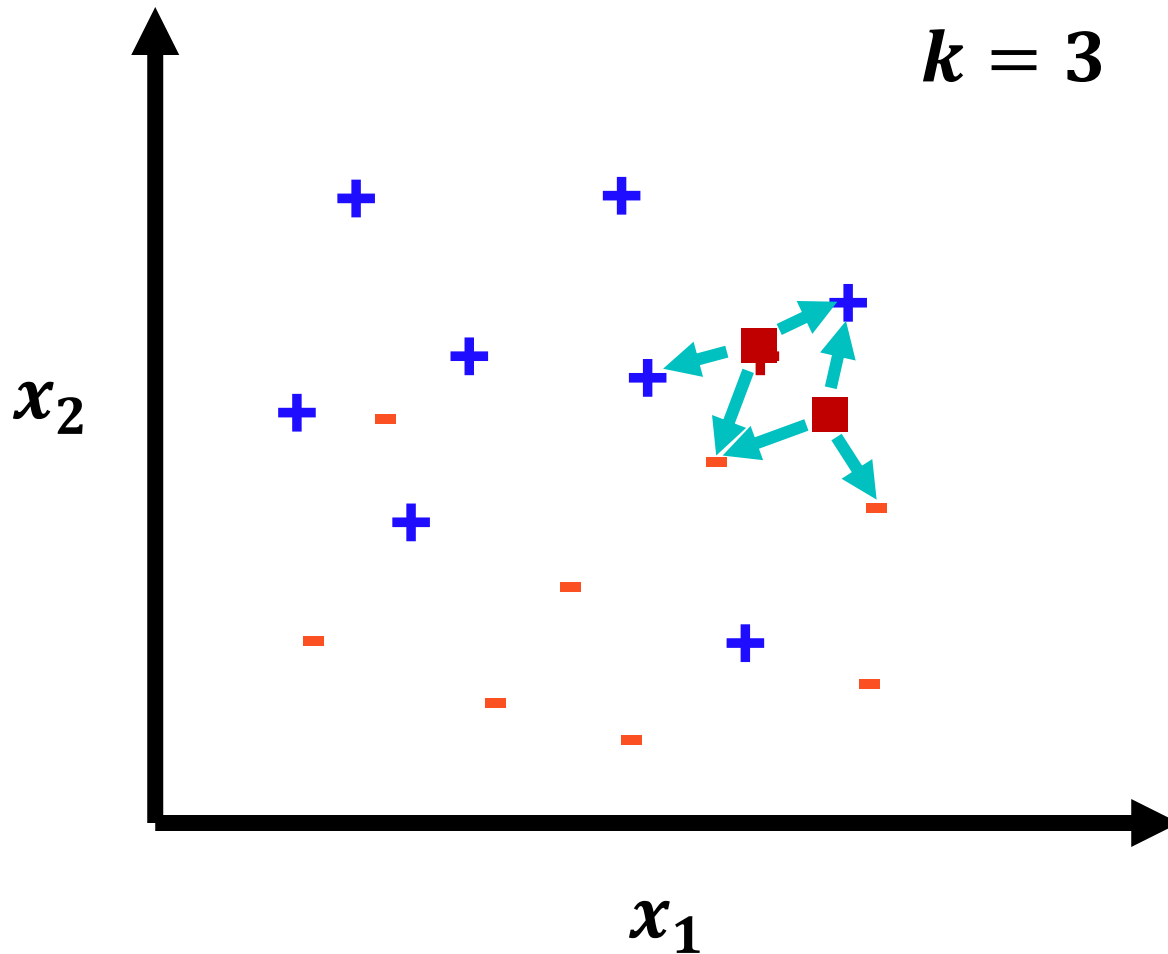
Metrics

How do we assess model performance?

For now, we consider **classification accuracy**:

What **proportion** did we get **right**?

K-Nearest Neighbours



*The
data is
our
model!*

In words...

1. Compute distance between the **candidate** point and all other neighbours in the **train** set:
 - For now, Euclidean distance:

$$d(p, q) = \sqrt{\sum_i (p_i - q_i)^2}$$

2. Select the k -nearest neighbours
3. Determine the candidate class from the k -nearest neighbours
 - For now, take the majority vote!

Notebook 1: Implementing k-NN

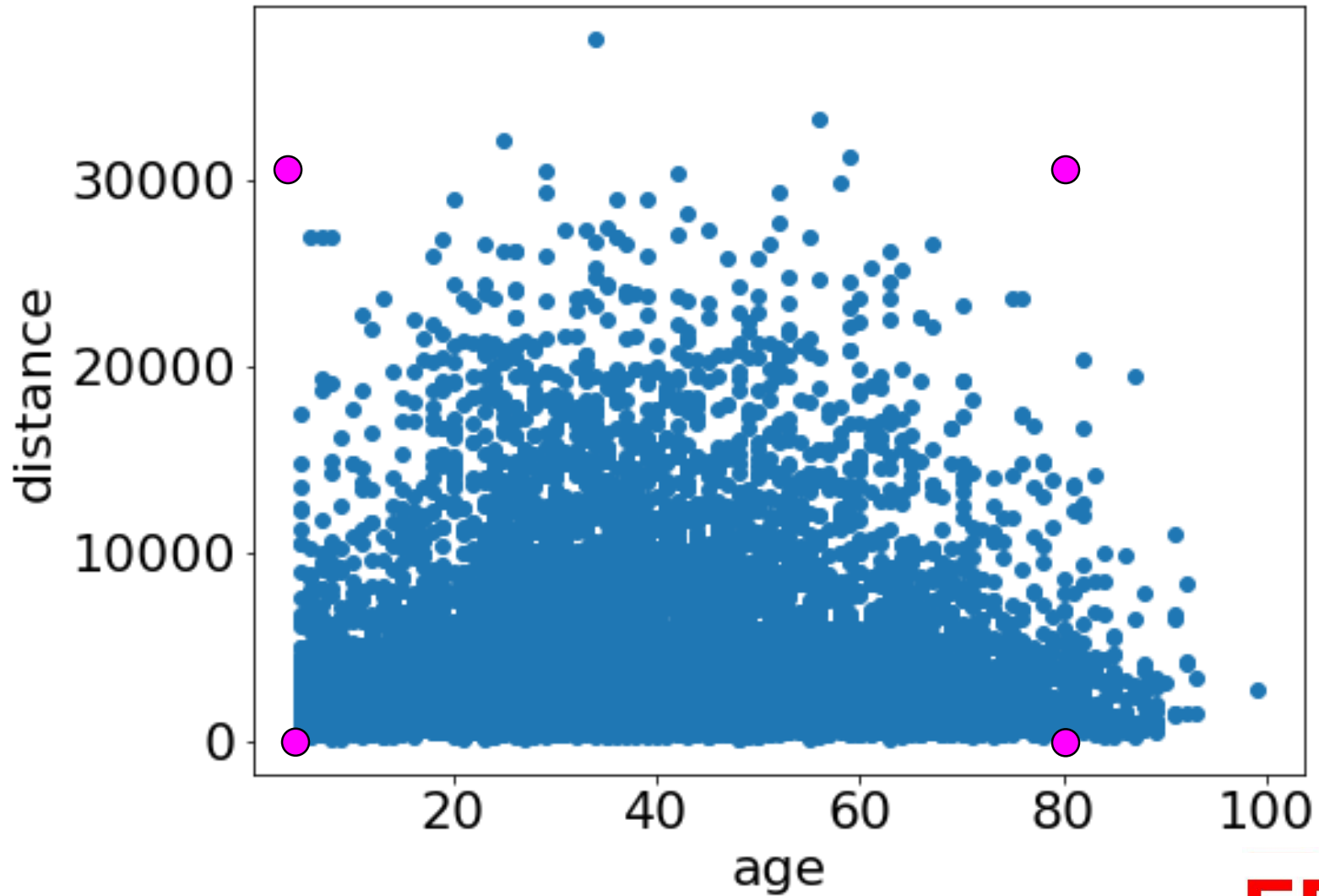
Areas for improvement

Ideas?

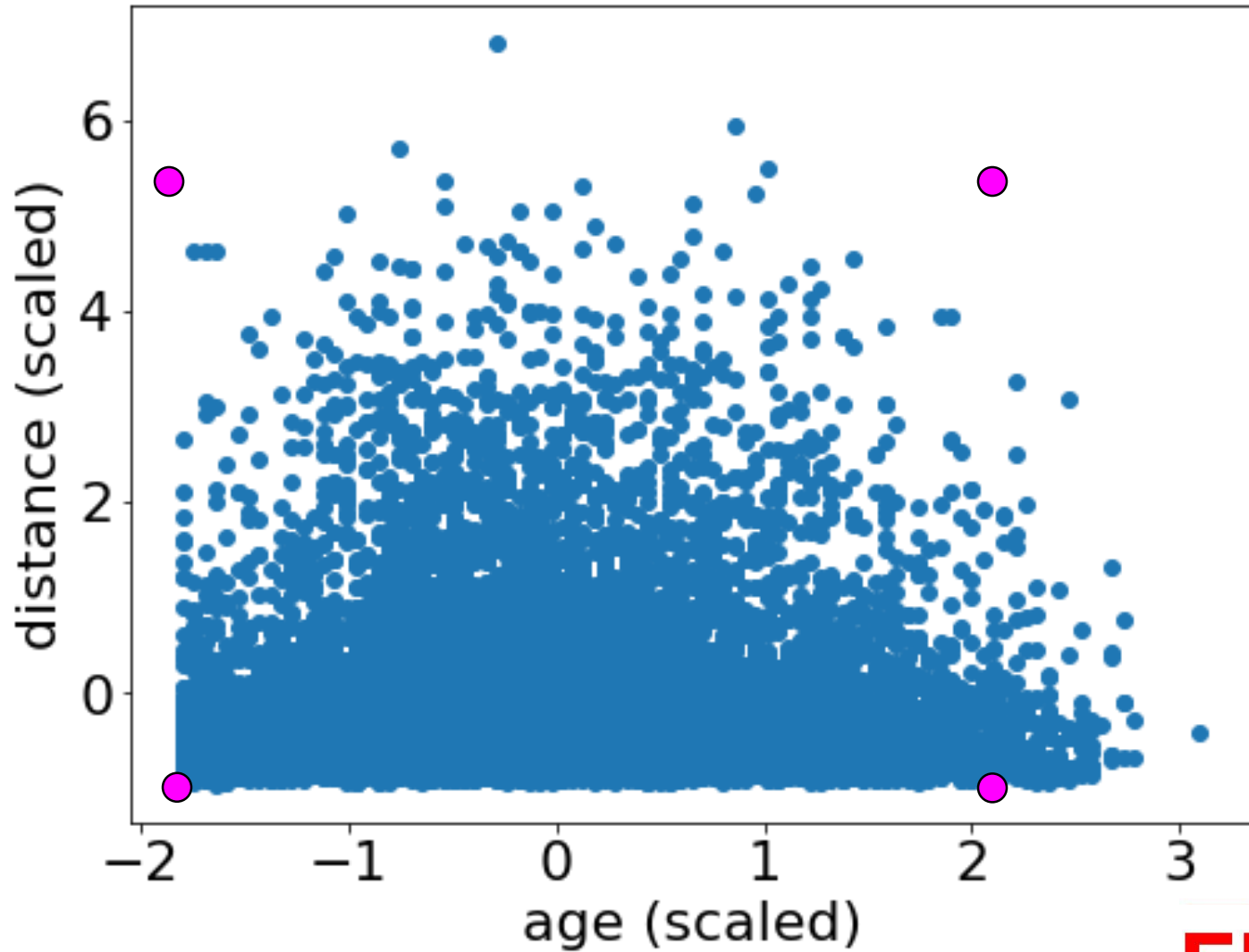
Areas for improvement

- Data preparation
 - Scaling the data
- Model optimisation
 - Choice of k
 - Distance metric
 - Improvements on majority vote
- Model evaluation
 - Alternatives to accuracy

Data scaling



Data scaling



Model optimisation

- Introduced three model *hyperparameters*
 - k
 - Distance metric
 - Decision rule

- All of these can be modified to improve *out-of-sample* performance
 - $k = 1, 2, 3, \dots, 100$? Bigger?
 - Manhattan distance? Minkowski distance?
 - Distance based weighting?

Model optimisation

How to select model **hyperparameters**?

Evaluate different values on **unseen data!**

Model evaluation

Is **accuracy** the best policy?

Consider the **confusion matrix**

Confusion Matrix

	PREDICTED CLASS		
ACTUAL CLASS		1	0
	1	TP	FN
	0	FP	TN

$$\text{Accuracy} = (TP+TN) / (TP+TN+FP+FN)$$

Confusion Matrix

	PREDICTED CLASS		
		1	0
ACTUAL CLASS	1	990	0
	0	10	0

Accuracy = 99% !

Known as *accuracy paradox*

Precision and recall

	PREDICTED CLASS		
		1	0
ACTUAL CLASS	1	TP	FN
	0	FP	TN

Precision (+ve) = $TP / TP+FP$

Recall (+ve) = $TP / TP+FN$

Precision and recall

	PREDICTED CLASS		
		1	0
ACTUAL CLASS	1	990	0
	0	10	0

Precision (+ve) = 0.99%

Recall (+ve) = 100%

Precision and recall

	PREDICTED CLASS		
		1	0
ACTUAL CLASS	1	TP	FN
	0	FP	TN

Precision (-ve) = $TN / (TN + FN)$

Recall (-ve) = $TN / (TN + FP)$

Precision and recall

	PREDICTED CLASS		
		1	0
ACTUAL CLASS	1	990	0
	0	10	0

Precision (-ve) = 0%

Recall (-ve) = 0%

Multinomial case - Precision

	PREDICTED CLASS					
	0	1	2	3	...	
ACTUAL CLASS	0	10				
	1	20	100	40	10	30
	2		5			
	3		15			
	...		20			

$$\begin{aligned}\text{Precision} &= TP / TP+FP \\ &= 100/150 \\ &= 67\%\end{aligned}$$

Multinomial case - Recall

	PREDICTED CLASS					
	0	1	2	3	...	
ACTUAL CLASS	0		10			
	1	20	100	40	10	30
	2		5			
	3		15			
	...		20			

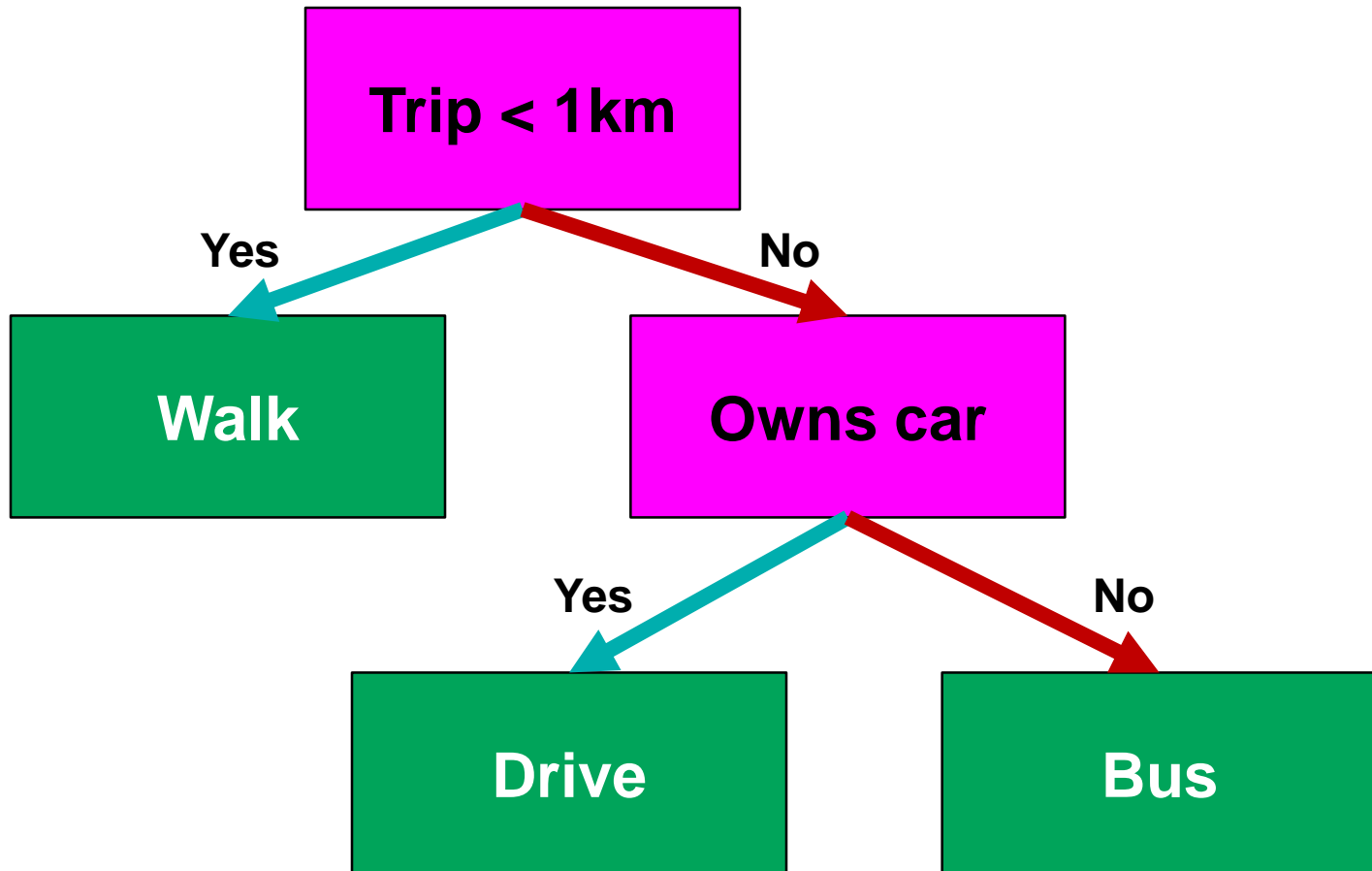
$$\begin{aligned}\text{Recall} &= \text{TP} / \text{TP} + \text{FP} \\ &= 100 / 200 \\ &= 50\%\end{aligned}$$

Machine learning in Python



Notebook 2: k-NN in scikit-learn

Decision trees



Decision trees

- Recursive hierarchical structure of *binary* splits

To fit (each split):

- For **every feature**:
 - Sort the data over the feature
 - For **every possible split point**:
 - ◆ Calculate the **gain** in separation
- Choose the split with the highest **gain**
- Repeat

Separation

- Separation is the reduction in the *shuffledness* of the data
- Two common metrics
 - GINI impurity
 - Entropy

GINI Impurity

$$G(p) = \sum_{i=1}^J p_i(1 - p_i) = 1 - \sum_{i=1}^J p_i^2$$

where p_i is the **proportion** of class i

Entropy

$$H(p) = - \sum_{i=1}^J p_i \log_2 p_i$$

where p_i is the **proportion** of class i

When to stop splitting

- Maximum depth
- Minimum leaf node size
- Minimum split size
- Minimum gain from split
- Maximum number of leaf nodes
- Etc...

Lots of hyperparameters!

Homework assignment

Notebook 3: Decision trees