Inf2 - Foundations of Data Science: k-Nearest neighbours and evaluation



THE UNIVERSITY of EDINBURGH informatics



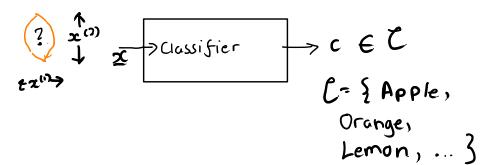
Recap of Monday's lecture

Supervised learning of classifications

	Fcatures	Fecture	vector	Lobel /cli
	Income 2	Housing	Employment	Repaid
_	20,000	Rent	Student	Yes
	50,000	Owns	Retiracl	No
	•	· ;	,	
-	30,000	Rent	Retired	Ż

Definition of a classifier

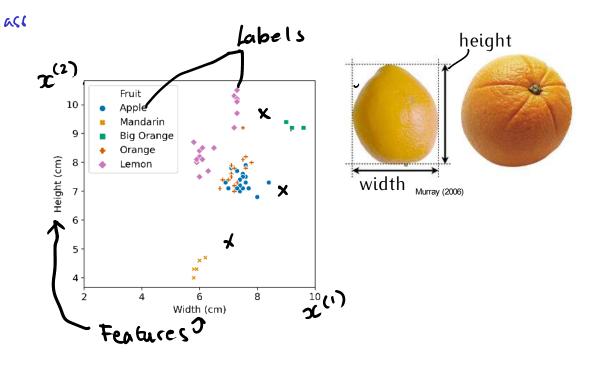
Feature vector
$$\underline{x} = (x^{(1)}, \dots, x^{(D)})^{+}$$



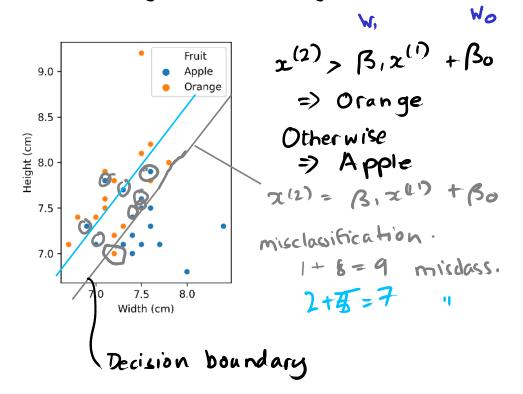
Properties

- 1. Decision Doundaries
- 2. Classification errors

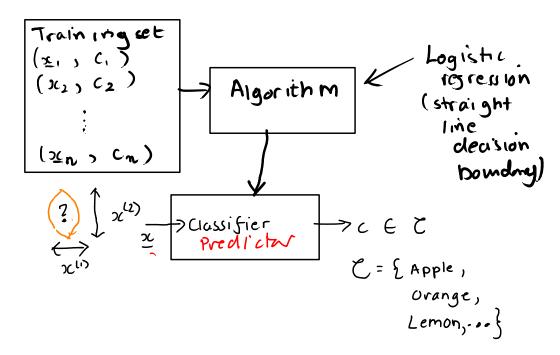
Visualising the classification problem



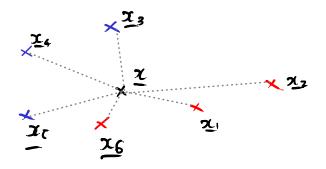
Constructing a classifier by hand

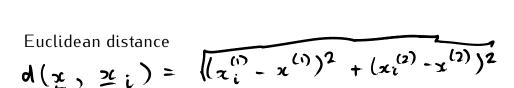


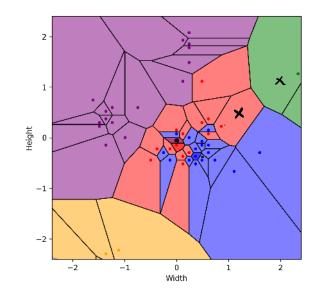
Constructing a classifer by supervised learning



Principle of nearest neighbour classification Nearest neighbour classification of fruit

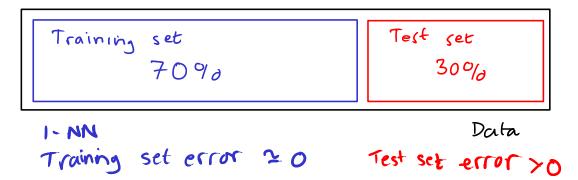






- Standardised variables
- Often better for distance-based algorithms

Evaluating generalisation to unseen data



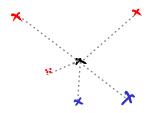
Never use the test set to train the classifier and then quote the error based on the test set

The estimate of the error will be too optimistic.

Principle of k-NN

To improve generalisation look at...

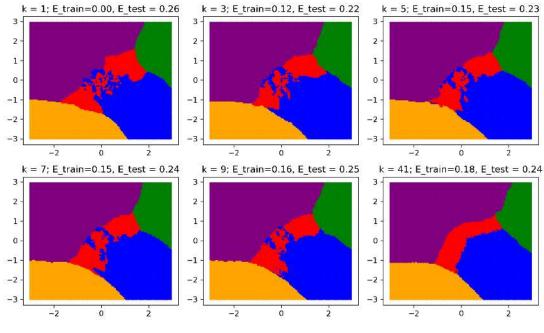
R- nearest neighbours



k	R	B	Winner
l	:	0	R
2	1	1	B/R
3	1	2	3
4	2	2	B/R
5	3	2	R

Ties: random resolution or weighting Odd numbers of k

Effect of k-NN



How to pick k?

k is a hyperparameter

over-generalising

0.25

0.20

0.15

0.00

Etrain
Etest

20

40

overfitting

under-generalising

Overview

- Metrics

k-NN regression

- Evaluation - K-fold cross validation

Inf2 - Foundations of Data Science: k-Nearest neighbours and evaluation Metrics







(Classification) accuracy

Accuracy and unbalanced classes

Definition =

Question: You want to test for a disease that you think around 5% of the population have. A company offers you a test that has 95% accuracy. Should you use it?

How to build a 95% accurate Covid classifier

Better measures for unbalanced classes

- 1. Sensitivity = Fraction of positives classified as paritive
- 2. Selectivity = Fraction of negatives classified as negative

Accuracy: 95%

```
Honly 5% of cases are truly positive then classifying all cases as negative gives:

Sensitivity: 0%

Selectivity: 100%
```

Confusion matrix

			Predic	ted class	1	Bad classifie	, (
			P	N		PN	
Actual	Positive	P					
CINSS	Negative	N					

Other metrics for classifiers

- Precision and recall
- F1
- In the context of a classifier with a adjustable threashold (e.g. Logistic regression, see later lectures)
 - Receiver operator characteristic
 - Area under the curve

Inf2 - Foundations of Data Science: k-Nearest neighbours and evaluation k-NN regression

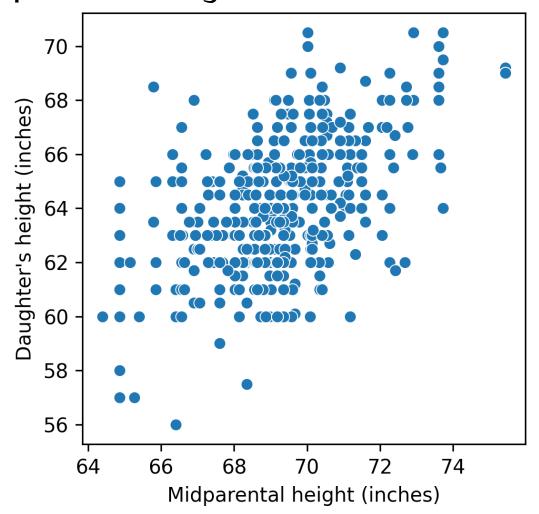


THE UNIVERSITY of EDINBURGH Informatics



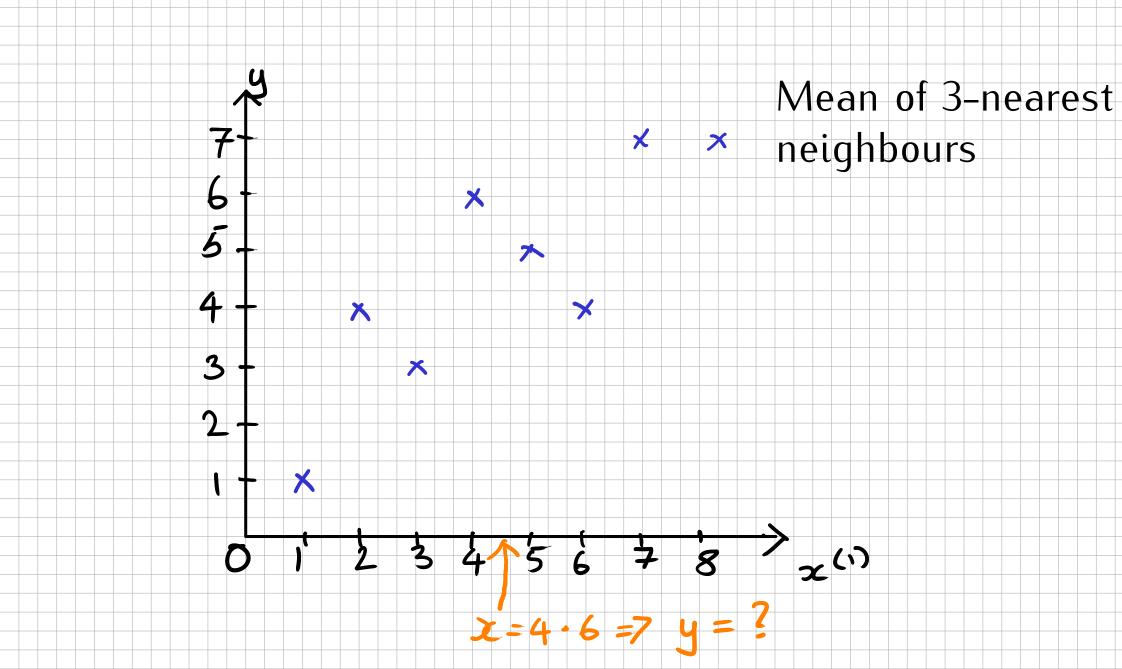
k-NNs for regression

Height of daughter versus parent height

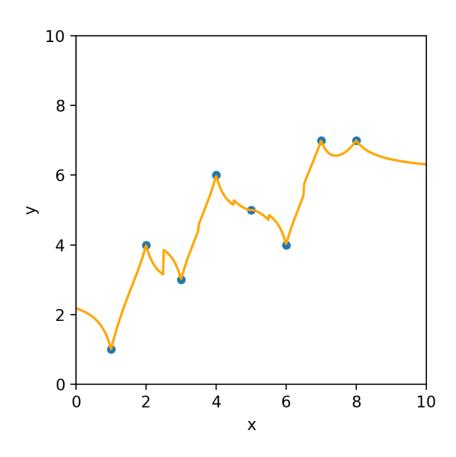


Could we use the principle of k-Nearest Neighbours to predict the daughter's height from the midparental height?

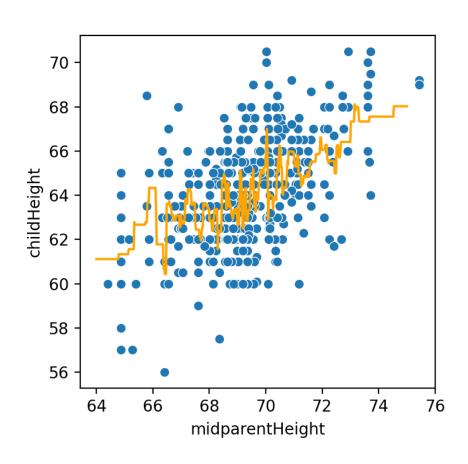
1-D example

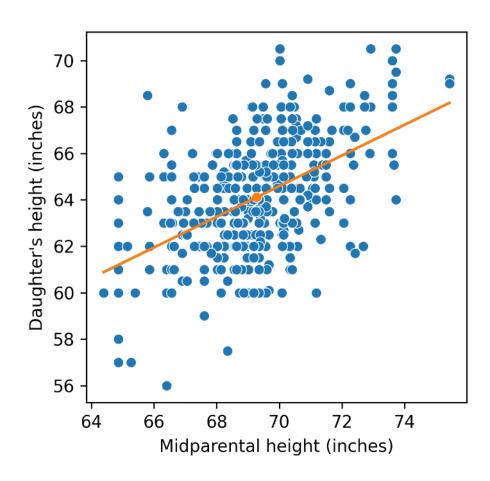


1-D example weighted by distance



k-NN regression versus linear regression





k-NN in more than one dimension

Na-set of Repoints newrest

Uniform weights:

$$y_{3}=9$$

$$y_{4}=6$$

$$y_{5}=9$$

$$y_{5}=2$$

$$x_{5}=7$$

$$y_{5}=7$$

$$y_{5}=7$$

Inverse-distance Weighted

Why use k-NN regression?

Metrics for regression

Same as for Linear regression!

- (Root) Mean squared error

- R-squared

Also:

- Mean Absolute Error (MAE)

Inf2 - Foundations of Data Science: k-Nearest neighbours and evaluation -K-fold cross-validation







The supervisedmachine learning paradigm

- 1. Split data into training, validation and test sets
- 2. Use training and validation data to select
 - (a) an alogorithm
 - (b) hyperparameters (if applicable)
- 3. Use test data to report performance of resulting predictor (classifier or regressor)

Limitations of the supervised machine learning paradigm

Data can change over time

- e.g. Covid symptoms changed over time

Selection of training, validation and test sets

- should be drawn from the same distribution
- should represent the real world to which the predictor will be applied c.f. Ethics!

Limited amount of data

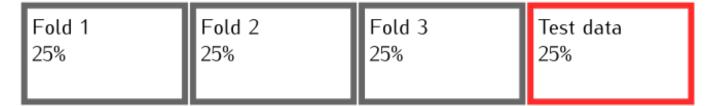
Training, validation and test data

Training 50% Validation duta Text data 25% Use to choose model & hyperparameters

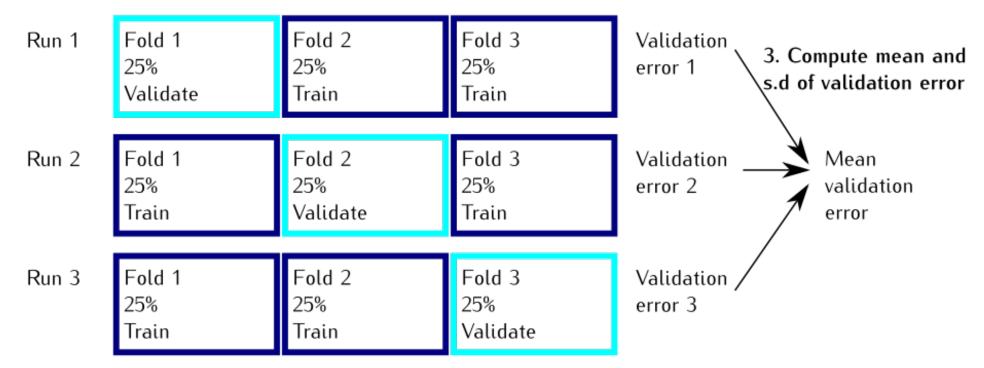
Problem: with small data, we're not able to get a very good estimate fo the validation loss/error

Solution: K-fold cross-validation

1. Split into testing data and non-testing data. Spilt non-testing data into 3 folds.



2. For each model and hyperparameter, train and compute validation error three times:



 Pick hyperparameter with lowest mean validation error. Train on all non-testing data and report performance on test data.

Training data	Test data
75%	25%

Model-fitting procedure

- la. Split data into training and testing sets
 - 16. Split training data into K folds
 - 2. for model in models
 - for hyperparameter in hyperparameters

for k in K

train model with hyperparameter on all folds apart from k

test trained model on ken fold

store mean and s.d of test errors with hyperpar.

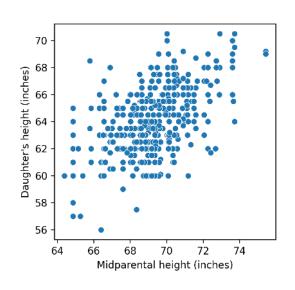
3. Choose model and hyperparameters with best Mean test error

Reporting test metric

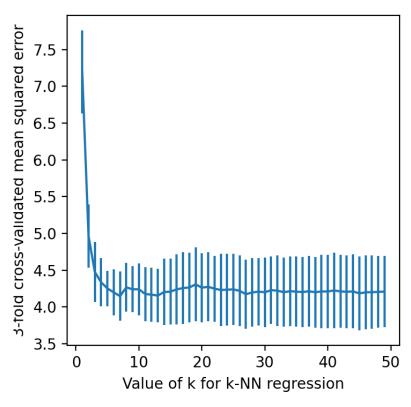
 Although validation metric gives estimate of expected test metric, still need to compute metric on separate test set

- Exception: in a competition, the organisers create training and testing sets, keeping test set secret

Example – applying cross-validation with k-NN vs linear regression



k-NN cross-validation

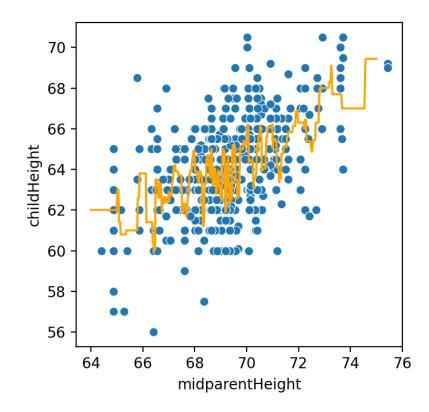


Linear regression: MSE mean (s.d) = 3.95 (0.21)

Test results

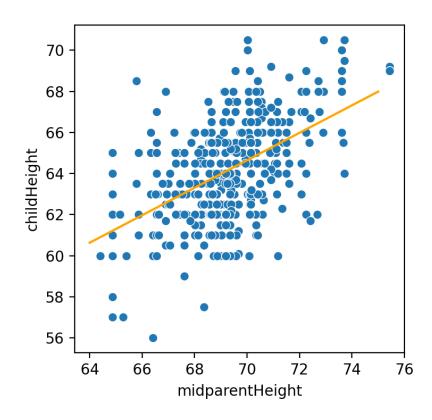
k-NN

MSE = 4.87 sq in RMSE = 2.21 in MAE = 1.68 in



Linear regression

MSE = 4.56 sq in RMSE = 2.14 in MAE = 1.66 in



Summary

- Classification metrics: careful choosing what to measure
 - accuracy can be misleading with unbalanced classes
- k-NN regression good for predicting nonlinear relationships between a numeric label and numeric predictors/features
 - Often worth trying linear regression!
- Cross-validation helps choose ML algorithms and hyperparameters
- Test data is to estimate how the predictor works in the real world not for training!!