#### Nearest Neighbor Classification

#### Morteza H. Chehreghani morteza.chehreghani@chalmers.se

Chalmers University of Technology

April 11, 2019

#### Reference

The content and the slides are adapted from

S. Rogers and M. Girolami, A First Course in Machine Learning (FCML), 2nd edition, Chapman & Hall/CRC 2016, ISBN: 9781498738484

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

## Introduction

#### Supervised learning

- Regression
  - Minimised loss (least squares)

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

- Maximised likelihood
- Bayesian approach
- Classification
- Unsupervised learning
  - Clustering
  - Projection

# Classification



A set of N objects with attributes (usually vector)  $\mathbf{x}_n$ .

Each object has an associated response (or label) t<sub>n</sub>.

(日) (日) (日) (日) (日) (日) (日) (日)

- Binary classification:  $t_n = \{0, 1\}$  or  $t_n = \{-1, 1\}$ ,
  - (depends on algorithm).
- Multi-class classification:  $t_n = \{1, 2, \dots, K\}$ .

# Classification syllabus

- ► 4 classification algorithms.
- Of which:
  - 2 are probabilistic.
    - Bayes classifier
    - Logistic regression.
  - 2 are non-probabilistic.
    - K-nearest neighbours
    - Support Vector Machines.

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

There are many others!

Classifier is trained on  $\mathbf{x}_1, \ldots, \mathbf{x}_N$  and  $t_1, \ldots, t_N$  and then used to classify  $\mathbf{x}_{new}$ .

Probabilistic classifiers produce a probability of class membership P(t<sub>new</sub> = k|x<sub>new</sub>, X, t)

• e.g. binary classification:  $P(t_{\text{new}} = 1 | \mathbf{x}_{\text{new}}, \mathbf{X}, \mathbf{t})$  and  $P(t_{\text{new}} = 0 | \mathbf{x}_{\text{new}}, \mathbf{X}, \mathbf{t})$ .

Which to choose depends on application....

Classifier is trained on  $\mathbf{x}_1, \ldots, \mathbf{x}_N$  and  $t_1, \ldots, t_N$  and then used to classify  $\mathbf{x}_{new}$ .

Probabilistic classifiers produce a probability of class membership P(t<sub>new</sub> = k|x<sub>new</sub>, X, t)

• e.g. binary classification:  $P(t_{\text{new}} = 1 | \mathbf{x}_{\text{new}}, \mathbf{X}, \mathbf{t})$  and  $P(t_{\text{new}} = 0 | \mathbf{x}_{\text{new}}, \mathbf{X}, \mathbf{t})$ .

Non-probabilistic classifiers produce a hard assignment
e.g. t<sub>new</sub> = 1 or t<sub>new</sub> = 0.

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

Which to choose depends on application....

Probabilities provide us with more information –  $P(t_{\text{new}} = 1) = 0.6$  is more useful than  $t_{\text{new}} = 1$ .

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

Tells us how sure the algorithm is.

Probabilities provide us with more information –  $P(t_{\text{new}} = 1) = 0.6$  is more useful than  $t_{\text{new}} = 1$ .

- ► Tells us how **sure** the algorithm is.
- Particularly important where cost of misclassification is high and imbalanced.
  - e.g. Diagnosis: telling a diseased person they are healthy is much worse than telling a healthy person they are diseased.

Probabilities provide us with more information –  $P(t_{\text{new}} = 1) = 0.6$  is more useful than  $t_{\text{new}} = 1$ .

- ► Tells us how **sure** the algorithm is.
- Particularly important where cost of misclassification is high and imbalanced.
  - e.g. Diagnosis: telling a diseased person they are healthy is much worse than telling a healthy person they are diseased.

・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・
・

- Extra information (probability) often comes at a cost.
- For large datasets, might have to go with non-probabilistic.

# Algorithm 1: K-Nearest Neighbours

Non-probabilistic.

Can do binary or multi-class.

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ のへぐ

► No 'training' phase.

# Algorithm 1: K-Nearest Neighbours

Non-probabilistic.

- Can do binary or multi-class.
- No 'training' phase.
- How it works:
  - Choose K
  - For a test object x<sub>new</sub>:
  - Find the K closest points from the training set.

- Find majority class of these *K* neighbours.
- (Assign randomly in case of a tie)



◆□▶ ◆□▶ ◆三▶ ◆三▶ ○三 のへぐ



◆□▶ ◆□▶ ◆ 臣▶ ◆ 臣▶ ○ 臣 ○ の Q @



Find K = 6 nearest neighbours.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで



Class one has most votes – classify  $\mathbf{x}_{new}$  as belonging to class 1.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ



Second example - class 2 has most votes.



Binary data.

▲□▶ ▲□▶ ▲臣▶ ★臣▶ = 臣 = のへで



- 1-Nearest Neighbour.
- Line shows decision boundary.
- Too complex should the islands exist?



<ロト < 部ト < 注入</p>

- 2-Nearest Neighbour.
- What's going on?



- 2-Nearest Neighbour.
- What's going on?
- Lots of ties random guessing.



<ロト <回ト < 注ト < 注ト

æ

- ► 5-Nearest Neighbour.
- Much smoother.



<ロト <回ト < 回ト < 回ト

æ

▶ 19-Nearest Neighbour.

Very smooth.



Binary data.

▲ロト▲舂▶▲恵▶▲恵▶ 恵 のへで



Non-smooth – too complex again?

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 - のへで



Random effects again...



► Getting smoother.

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ ○臣 - の々ぐ



Smoother still.

## Problems with KNN

#### Class imbalance

- ► As K increases, small classes will disappear!
- Imagine we had only 5 training objects for class 1 and 100 for class 2.

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ

For  $K \ge 11$ , class 2 will **always** win!

## Problems with KNN

#### Class imbalance

- As K increases, small classes will disappear!
- Imagine we had only 5 training objects for class 1 and 100 for class 2.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ■ ●の00

- For  $K \ge 11$ , class 2 will **always** win!
- How do we choose K?
  - Right value of K will depend on data.
  - Cross-validation!

## Cross-validation for classification

- E.g. to find *K* in KNN:
- Exactly the same as we have seen before.
- Split the data up use some to train, some to validation.

▲□▶ ▲□▶ ▲□▶ ▲□▶ ▲□ ● ● ●

- Need a measure of 'goodness'.
- Use number of mis-classifications.....
- ....and use K that minimises it!

## Remember...



Average number of misclassifications over the C folds.

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

#### Example – 5 classes



<ロト <回ト < 回ト < 回ト

э

- ▶ 5 classes.
- Smallest has 20 instances, biggest 120.

#### Example – 5 classes



- Curve shows average misclassification error for 10-fold CV.
- Minimum at approximately K = 30.

#### Example – 5 classes



(日)

æ

- As K increases, classes 'disappear'
- Causes the 'steps' in error.

# KNN – summary

- Non-probabilistic.
- Fast.
- Only one parameter to tune (K).

▲□▶ ▲□▶ ▲ □▶ ▲ □▶ □ のへぐ

- Important to tune it well....
- …can use CV.

# KNN – summary

- Non-probabilistic.
- Fast.
- Only one parameter to tune (K).
- Important to tune it well....
- …can use CV.
- There is a probabilistic version.
  - Not covered in this course.

▲□▶ ▲□▶ ▲ 三▶ ▲ 三▶ 三 のへぐ