

# K-nearest Neighbours

①

- handle both binary & multiclass
- No assumptions about parametric form of the decision boundary

KNN has no training phase

Normal Scenario:

- $N$  objects
- each rep by set of attris  $X_n$
- And a label  $t_n$

to classify  $X_{new}$  w/ KNN we first find the  $k$  number of training points that are closest to  $X_{new}$

$t_{new}$  is then decided based on the majority class among the  $n$  ~~neigh~~ neighbours



One issue with ~~the exact~~ knn is the example of tied class groups  
 i.e.  $k = 8$  neighbour class 1 = 4 class 2 = 4

Solutions -

- use odd  $k$
- weight neighbours, closer = more influential. e.g. euclidean

How to choose  $k$ ?

- too small = influenced by noise
  - This creates islands or clusters which represent overfitting
- Too big =
  - upto a point  $> k$  has a regularization affect and which reduces overfitting
  - Too far then we loose true fit of the data
  - data imbal w/ high  $k$  can completely shut out a class or just overrule



3

## K-cross validation

the most popular method for choosing  $k$  is  $k$ -cross validation

The metric we look @ to decide is the proportion of classification mistakes made

10-fold cross validation & repeat 100 times

Plot error result for each  $k$

Select the minima

With knn & 10-Fold, train on 9 segments & validate on the 1 to count the errors

repeat cycling through the holdout segments