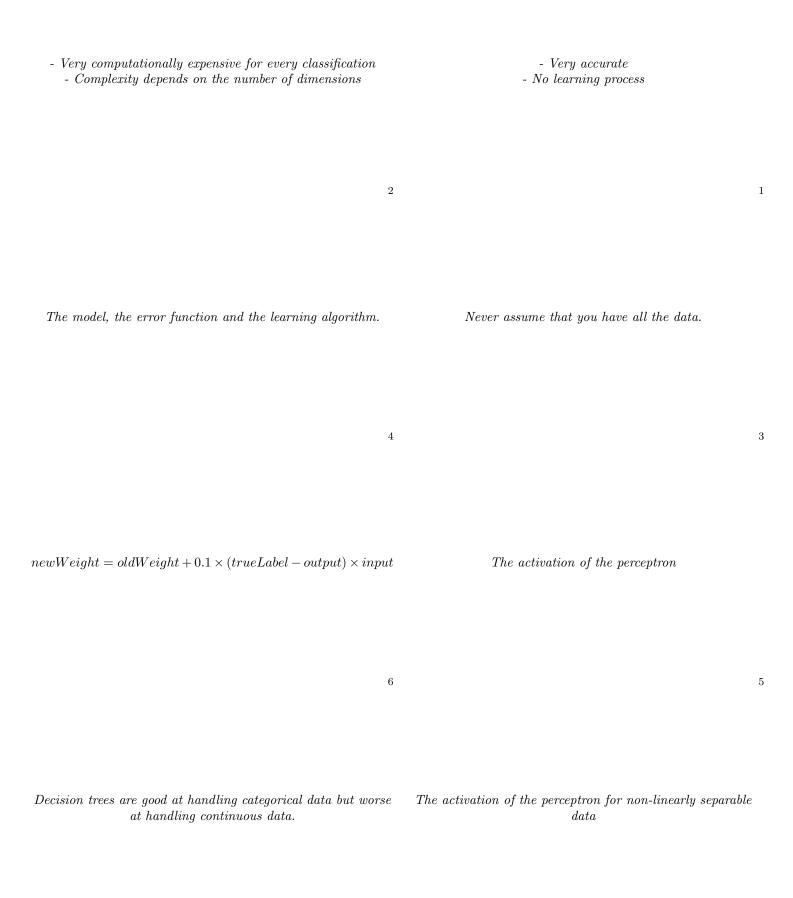
| What are the advantages of a nearest neighbour classifier?                                 | What are the <b>disadvantages</b> of a nearest neighbour classifier?  |
|--|---|
| What is the most important concept in machine learning?                                    | What are the three 'ingredients' of a machine learning algorithm?     |
| What does this equation calculate? $a = \sum_{i=1}^F x_i w_i$                              | What is the perceptron learning rule?                                 |
| What does this equation calculate? $a = \frac{1}{1 + exp(-\sum\limits_{i=1}^d w_i x_i)}$ 7 | Decision trees are good at handling data but  worse at handling data. |



| What does this equation calculate? $H(X) = -\sum_i p(x_i) \log_2 p(x_i)$    | The 'information' contained in a variable is called the            |
|---|--|
| Explain the process of cross validation.                                    | What factors should affect our decision on the best value of $k$ ? |
| What is the ensemble approach to machine learning?                          | Briefly describe bootstrapping                                     |
| On average, what is the percentage of data points that are left unselected? | Explain bagging  |

The 'information' contained in a variable is called the entropy.

The entropy of a variable X

10

9

- 1. Accuracy
- 2. Training time and space complexity
- 3. Testing time and space complexity
- 4. Interpretability

- 1. Break the data evenly into N chunks
- 2. Leave one chunk out
- 3. Train on the remaining N-1 chunks
- 4. Test on the chunk you leave out
- 5. Repeat until all chunks have been used to test
- 6. Plot the average and error bars for the N chunks

12

11

Bootstrapping is the process of generating multiple data sets from an original.

Select a class of models, fit multiple models to training data (called base learners), use the models as a committee to vote on testing data.

14

13

Generate m bootstraps and train a model on each one. When the testing data arrives a simple majority vote takes place.

36.8%

| $Explain\ boosting$   | What type of classifier models a classification rule directly and models the probability of class memberships based on input data? |
|---|--|
| What type of classifier makes a probabilistic model of data within each class?          | What type of classifier uses probabilities to classify data?   |
| What is the formula to work out $P(c X')$ Where $c$ is a class and $X'$ is an example?  | What is the formula to work out a Gaussian model?  |
| What are the two data representation methods that we talk about in clustering analysis? | What is the formula to work out Minkowski distance.  |

A discriminative classifier

Get a data set, take a bootstrap and train a model on it. See which examples the model got wrong then upweight those 'hard' examples and downweight the 'easy' ones. Now go back to training a model, but now you have a weighted bootstrap.

18

17

 $A\ probabilistic\ classifier$ 

A generative classifier

20

19

$$\frac{1}{\sigma\sqrt{2\pi}}exp(-\frac{(x-\mu)^2}{2\times\sigma^2})$$

$$P(c|X') = [P(x_1|c)P(x_2|c)...P(x_n|c)]P(c)$$

Where x is a feature in the example.

22

21

$$d(x,y) = \sqrt[p]{(x_1 - y_1)^p + (x_2 - y_2)^p \dots + (x_n - y_n)^p}$$

Data matrices and distance matrices.

| What is the formula for Manhattan distance?                                | What is the formula for Euclidean distance?                               |  |  |
|--|---|--|--|
| What is the cosine measure equation  | What is the formula for the distance between symmetric binary attributes? |  |  |
| What is the formula for the distance between asymmetric binary attributes? | Briefly explain the partitioning clustering approach                      |  |  |
| Briefly explain the hierarchical clustering approach                       | How does the K-means clustering algorithm work?                           |  |  |

$$d(x,y) = \sqrt{(x_1 - y_1)^2 + (x_2 - y_2)^2 \dots + (x_n - y_n)^2}$$

$$d(x,y) = |x_1 - y_1| + |x_2 - y_2| + |x_n - y_n|$$

26

25

$$d(x,y) = \frac{b+c}{a+b+c+d}$$

$$\frac{x_1y_1 + \dots + x_ny_n}{\sqrt{x_1^2 + \dots + x_n^2}\sqrt{y_1^2 + \dots + y_n^2}}$$

28

27

Various partitions are constructed and then evaluated by some criterion, r.g. minimizing the sum of square distance cost.

Typical methods are k-means, k-medoids, CLARANS.

$$d(x,y) = \frac{b+c}{a+b+c}$$

30

29

- 1. Pick K random seed points (from the data)
- 2. Assign each data object to the cluster with the nearest seed point
- 3. Compute the mean points of the clusters (aka centroids)
- 4. Go back to step 1, but the seed points are now the means calculated in step 3. Do this until the assigned sets don't change between iterations.

Create a hierarchical decomposition of the set of data (or objects) using some criterion. Typical methods are Agglomerative, Diana, Agnes, BIRCH, ROCK.

| What is the runtime of the K-means algorithm?                                | What are some problems with the K-means algorithm?                        |
|--|---|
| 33   | 34  |
| What are the two sequential strategies used to construct a tree of clusters? | Explain the single link method of measuring the distance between clusters |
| Explain the complete link method of measuring the distance between clusters  | Explain the average method of measuring the distance between clusters     |
| How does the agglomerative hierarchical algorithm work?                      | What are some weaknesses of the agglomerative approach?                   |

- The initial seed points can cause 'local optimum' clusters, which may not be representative of the whole data.
- If a cluster has a non-convex (i.e. concave) shape, then it won't be detected.
- Unable to classify categorical data (unless you modify the algorithm).
- It's hard to evaluate the performance of the algorithm,

## O(tKn)

Where t is the number of iterations
K is the number of clusters
n is the number of objects

34

Use the smallest distance between an element in one cluster and an element in another.

Agglomerative and divisive

36

Find the average distance between the points in the two clusters and use that.

Use the largest distance between an element in one cluster and an element in another.

38

- It is very sensitive to noise.
- It is less efficient than k-means clustering, with a runtime of  $O(n^2)$

- 1. Convert the object attributes into a distance matrix, and make each object a cluster (of size 1).
- 2. Merge the two closest clusters together.
- 3. Update the distance matrix to take into account the new cluster.
- 4. Repeat from step two until we've got the desired number of clusters.

40

| In order to | calculate the F-re | atio index u | ve divide t | the     |
|-------------|--------------------|--------------|-------------|---------|
|             | variance by the    |              | var         | riance. |

What is the formula to calculate the F-ratio index?

41 42

$$F(m) = \frac{mSSW(m)}{SSB(m)} = \frac{m \sum_{i=1}^{m} \sum_{j=1}^{n_i} d^2(x_{ij}, c_i)}{\sum_{i=1}^{m} n_i d^2(c_i, c)}$$

What is additive smoothing used for? Give an equation for it.

43

44

$$F(m) = \frac{mSSW(m)}{SSB(m)} = \frac{m \sum_{i=1}^{m} \sum_{j=1}^{n_i} d^2(x_{ij}, c_i)}{\sum_{i=1}^{m} n_i d^2(c_i, c)}$$

In order to calculate the F-ratio index we divide the intra-cluster variance by the inter-cluster variance.

41

42

When a  $P(x|\omega_i)$  has a probability of 0 (it's 0 for just class i, additive smoothing can ensure the probability of the data object isn't 0.

$$P(x|\omega_i) = \frac{N_{i,\omega} + \alpha}{N_i + \alpha d}$$

Where  $N_i$  is the count of that feature,  $N_{i,\omega}$  is the count for that class,  $\alpha$  is the smoothing parameter (1 for laplace smoothing) and d is the number of dimensions.

m The number of clusters generated by the algorithm.

 $n_i$  The number of data points in the ith cluster.

 $c_i$  The centroid for the ith cluster.

 $x_{ij}$  The jth datapoint in cluster  $c_i$ .

c The mean centroid for the whole data set.

d(x,y) The distance function we're using.