

04/15

MLA

K-Nearest Neighbours

- Classification that is non parametric / distribution free method of assigning group membership
- Make no assumption about spread of data within each class
- Consequence of this is that class assignment is fixed, with no measurement of uncertainty concerning any particular assignment
- Classification techniques that do not make distributional assumptions at least allow quantification of the uncertainty in group membership
- K nearest neighbors simply look at the K-closest points of known origins to the point of unknown origin
- The point is then classified as belonging to the group which contains the most of these K points
- Results are not invariant to the scaling of the original variables, nor to the method in which distance is calculated, as these are likely to change the nearest neighbor for any particular point

Advantage: No assumptions, so assumptions can't be wrong

Disadvantage: Cannot express probabilistic uncertainty about our classification i.e. class assignment is fixed e.g. point A is in class A as opposed to 65% point is in class A.

We look at the K nearest point of known label to our new point and classify it as belonging to the class which is most prevalent.

Choosing K

- The classification varies with K.

- One approach for choosing K is to split labelled data into 3 parts:

→ Training - points whose labels are used to classify unlabelled points

→ Test - points we know the labels for but consider unlabelled in order to find best K to classify them

→ Validation - Remaining labelled are considered unlabelled in order to estimate classification error for best K in test set

- Plot as function of k the proportion incorrectly classified

Why Validate?

- The correct classification rate for test data typically overestimates the percentage correct classifications in validation data. Because value of k is chosen specifically for the test set, and may not be representative for another unlabeled sample.

- Validation data is not used at any stage of the model fitting, and hence offers a more reasonable estimate of the correct classification rate.

Cross-Validation

- In this problem (leave-one-out) cross validation would be used as follows:

→ For each value of k remove each data point and determine if other data point would be correctly classified knowing the label of all other data points.
e.g. 100 points, making 100 classifications of 1 point based on labels of other 99.

→ Example. 3 data points (x_1, x_2, x_3). See if we correctly classify x_1 given labels for x_2, x_3 , x_2 for x_1, x_3 , x_3 for x_1, x_2 .

Each value of k will either get non correct, one correct, two or all correct.

→ Select value of k that has best classification rate.

k-MEANS

- Aim is to divide the data into k distinct groups so that observations within a group are similar, whilst observations between groups are different.

- k -mean clustering is an iterative, rather than hierarchical, clustering algorithm.

- Means that at each stage of algorithm data point will be assigned to a fixed number of clusters. (Contrasts with hierarchical clustering where the number of clusters ranges from the number of data points to a single cluster).

- Can select k from expert knowledge.

- Or use previous results from preliminary data exploration.

- Simple and computationally efficient, can be sensitive to selection of starting points.

- Running k -mean several times with different start points can help check robustness of results.

MLA

3

24/15

K-MEANS

Pseudo Code

1. Choose the number of clusters k and designate cluster centers
2. Assign each data point to the cluster whose center is closest.
3. For cluster i , calculate its centroid $C_i^T = (C(i)_1, C(i)_2, \dots, C(i)_m)$ where m is number of variables in observation \rightarrow these are found by averaging variable scores for data points within the cluster.
4. Calculate the sum of squared distance to each object to its cluster centroid:
$$SS = \sum_{i=1}^k \sum_{j=1}^n (x_{ij} - C(i)_j)^2$$

Assume a total of n observations. Minimizing SS value

5. Re-assign each observation to the cluster whose centroid is closest.
6. Repeat (3) \rightarrow (5) until convergence

- Randomly assign centers of each group

- Initial partition can be constructed in several ways:

- A random selection of k observations

- Specify selection based on prior knowledge

- By using results from an existing hierarchical clustering algorithm

- label points according to which center is closest

- update values for the three centers (prototypes)

- repeat

- k -means has converged when no points are moved between groups on an iteration

- Once this happens, centers will no longer change, points either

- the convergence criteria might not be suitable in some cases e.g. if n is very large, and alternatives are possible, e.g. ~~within~~ within cluster sum of squares does not change over 3 iterations etc

Choosing k for k

- Euclidean
- Generally shall run k mean for different number of clusters
- Why not choose k that minimises SS? The more clusters there are fitted to smaller the SS
- General rule is to plot k against SS and look for a 'kink' in the curve. If there is no kink, then there is a trade off between additional complexity by increasing k and better fit by reducing the SS

Local Minima

- k mean algorithm can give different answers when initialised at different starting values
- Mean algo will not always find min value for the Total Within Sum of Squares
- Choose lowest this

k -Mean Clustering - Model Based

- Model based clustering makes use of statistical method and is the parametric model
- Several possible extensions
- we used euclidean distance, and new centres are mean value of points assigned to that group
- Choice of dissimilarity and choice of center are related

Extensions

- Can use different measure of dissimilarity
- Could assign x_i to group k so that $d(x_i, \mu_k)$ is minimised but we use different measure of dissimilarity than euclidean
- Can help algorithm from forming circular clusters
- The new centres for the group could be computed by minimising $\sum_{i=1}^n d(x_i, \mu_k) I_{ik}$ where $I_{ik} = 1$ if observation i is assigned to cluster k and 0 otherwise
- The partitioning around medoid algorithm does this

MLA

5

K-MEANS CLUSTERING

Cluster Medoids:

- A medoid is a representative object of a cluster so that its average dissimilarity to all the data points in the cluster is minimal
- Unlike the mean or centroid used in k-means, a medoid has to be an actual data point within the data.
- Like mean v median
- Medoid useful in applications where a mean or centroid are conceptually difficult to conceptually understand, or it may not even be defined e.g. 3D-coordinates

Partitioning around Medoid (PAM) Pseudo Code

1. Select a dissimilarity metric to be used
2. Initialize by selecting k of the n data points to be medoids
3. Cluster each data point to belong to the same group as the medoid it is closest to under the dissimilarity metric selected
4. For each medoid x^* :
 - for each non-medoid point x , swap x^* with x and compute the total dissimilarity cost of the configuration
5. Select the configuration with lowest total dissimilarity cost.
6. Repeat 3 to 5 until convergence i.e. no change in medoids.

Mixture Model

- Suppose we have data $x_i^T (x_{i1}, x_{i2}, \dots, x_{in})$ which is known to come from one of k populations
- Within each population (cluster) the data follow a density $f(x_i | \theta_j)$ for $j=1 \dots k$
- Where θ_j are the parameters governing population j
- Suppose the probability that a data point is from population j is π_j
- Then the data can be modelled using a mixture model:

$$P(x_i) = \sum_{j=1}^k P(x_i | \theta_j) \pi_j$$

$$p(x_i) = \sum_{j=1}^k \pi_j f(x_i | \theta_j)$$

- Mixture model can be used to form a model based clustering technique

Mixture Model: $m=1$

- The π_j could be called the mixing proportions, and $F(\cdot | \theta_j)$ is known as the j th component density
- Model offer good model flexibility by allowing both π_j and the model parameters within each population to vary
- Common form of $m=1$ is the mixture of normals, where each component density is a normal distribution density (or multivariate normal density)

$$\text{In the univariate case: } p(x_i) = \sum_{j=1}^K \pi_j \frac{1}{\sqrt{2\pi\sigma_j^2}} \exp\left[-\frac{(x_i - \mu_j)^2}{2\sigma_j^2}\right]$$

- Assumes that the data in each population is normal with mean μ_j and variance σ_j^2
- Similar idea to QDA, LDA

(Linear Shape)

- K mean lines for circular shape
- LDA fits ellipse of the same shape and direction (equal covariance matrix assumption)
- QDA allows for different cov matrices between groups (i.e. different shapes and directions)
- Model based clustering allows different shapes and directions as well as varying π_j

Multivariate Mixture Model: $m \geq 1$

- Normal MM can be easily extended to a multivariate mixture model
- In this case it is assumed that the data within group j follow a multivariate normal distribution with mean μ_j and cov matrix Σ_j
- We could constrain cov matrix in different ways to allow model flexibility

Decomposing Cov Matrix

- Can decompose any eigenvalue decomposition: $\Sigma = \lambda D A D'$

λ = a constant

D = orthogonal matrix of eigenvectors

A = diagonal matrix with entries proportional to eigenvalues

14/15

MLA

2

K-MEANS

Decomposing Cov Mat

- Seen that the control of the shape for a MVN from ellipse, the shape being controlled by cov mat Σ .

A - control Shape of ellipse

D - control orientation of ellipse

k - control size of ellipse

- If $A=D=1$ ellipse becomes circle
- If $D=1$ and A is unimodal, ellipse becomes diagonal with axes
- When we move to mixture of normal, the flexibility increased. This is the hall for model based clustering

- Main aim is to find clusters in data where K is unknown (unsupervised)
- Idea is to fit a mixture of ~~probable~~ normal to the data set for a range of possible values for k and for a range of different possibilities for the Σ 's
- Use a MLE approach via Expectation Maximization (EM), this is generally what happens in K-means and PAM algorithms
- For a given cluster number k and a likelihood model $L(\theta|X, z)$ for the probability of parameter θ on given data set X and cluster assignment z the pseudo EM algorithm iterates between the following:
 - E-step: find expected value of θ as a function of z using given likelihood model
 - M-step: Change cluster assignment z to max the expected likelihood from E-step

How to choose optimum model?

- Use Bayesian Information Criterion
- $BIC = -2 \times \text{maximized likelihood of data} + (\log N \times \text{\# of parameters})$
- take model with smallest BIC
- First term reward good fit, second term penalizes complexity (penalizes)
- Model based clustering will return optimal # of groups in data and initial optimal centroid decomposition of the Σ 's (can incorporate k into likelihood model). Can return estimates of group membership of each point and an estimate of the parameters in the group assignment