Quiz 2 - Introduction to Unsupervised Learning Techniques

Name *

Dawei Chen

Email *

dawei.chen24@imperial.ac.uk

Decison Trees

What is the entropy of a node where the label distribution is uniform (i.e., each class 1 point occurs with equal probability)?



• Entropy is maximum

Entropy is minimum

Entropy is 0.5





Feature Importance Analysis

Algorithm Mean Decrease Impurity (MDI) for Random Forests

Require: Trained random forest model with T trees **Ensure:** Feature importance scores $\{MDI_j\}_{j=1}^d$ 1: Initialize importance scores: $IMP_i = 0$ for all features j 2: for each tree t = 1 to T in the forest do for each internal node *n* in tree *t* do 3: Identify the feature F_i used for splitting at node n4: Let w_n be the proportion of samples reaching node n5: Calculate information gain $IG_n(j, \tau_n)$ 6: Update importance: $IMP_{j} = IMP_{j} + w_{n} \cdot IG_{n}(j, \tau_{n})$ 7: end for 8: 9: end for 10: Compute $MDI_j = \frac{1}{T} \cdot \frac{IMP_j}{\sum_{k=1}^d IMP_k}$ for all j11: return $\{\mathsf{MDI}_j\}_{j=1}^d$

• It is an out-of-sample measure

) It is biased toward high-cardinality features

) It is specific to tree-based models

Algorithm Permutation Feature Importance (PFI)

Require: Fitted model *m*, validation data *D*, repetitions *K* **Ensure:** Feature importance scores $\{PFI_j\}_{j=1}^d$ with stds $\{\sigma_j\}_{j=1}^d$

- 1: Compute reference score s of model m on data D
- 2: for each feature F_j (column of D) do
- 3: Initialize array scores; of length K
- 4: for each repetition k in $1, \ldots, K$ do
- 5: Randomly shuffle column j of dataset D to generate corrupted version $\tilde{D}_{k,j}$
- 6: Compute score $s_{k,i}$ of model *m* on corrupted data $\tilde{D}_{k,i}$

7: Store in array: $scores_j[k] = s - s_{k,j}$

- 8: end for
- 9: Compute mean importance PFI_j and standard deviation σ_j from array scores_j

(日)

- 10: end for
- 11: return $\{PFI_j\}_{j=1}^d$ and $\{\sigma_j\}_{j=1}^d$

) It is agnostic to the evaluation metric

It doesn't suffer from the substitution effect when features are correlated

) It is an out-of-sample measure



Oluster 3 contains noise features that the model mistakenly used during training

Cluster 3 is highly informative, but PFI underestimates it due to randomness

PFI is incorrect, because feature importance cannot be negative

K-means and GMMs

What does the K-means algorithm aim to minimize?

1 point

Algorithm The K-means Algorithm

Require: A data set $X = \{x_1, \ldots, x_n\}$ $(x_i \in \mathbb{R}^p)$ **Ensure:** An assignment function Ψ^* and the associated centroids c_1^*, \ldots, c_K^* . 1: Initialization: Choose c_1, \ldots, c_K in X at random 2: repeat for $i = 1 \dots n$ do 3: $\Psi(x_i) \leftarrow \arg\min_{k \in \{1,\dots,K\}} \|x_i - c_k\|^2$ 4: end for 5: for $k = 1 \dots K$ do 6: $c_j \leftarrow \frac{1}{\sum\limits_{i=1}^{n} \mathbb{I}(\Psi(x_i)=k)} \sum_{i=1}^{n} \mathbb{I}(\Psi(x_i)=k) x_i$ 7: end for 8: 9: until convergence 10: return $\Psi^*, c_1^*, \ldots, c_K^*$

) The number of clusters

The within-cluster distortion

) The average silhouette score

What is the expected silhouette score of a point that lies exactly between two clusters 1 point (equidistant to both)?





The EM algorithm increases the evidence lower bound (denoted as $L(q, \theta)$) at each iteration.

It guarantees the optimal number of mixture components

In the programming session, after computing the distance matrix, applying PCA, and * 1 point then using GMMs, we apply Bayes' rule to compute posterior probabilities. What do these represent?



Hard cluster assignments, like K-means does

Soft cluster assignments – probabilities that each point belongs to each cluster

A transformation of the distance matrix into centroids

Questions ?

Is the pace of the course appropriate? Would a probability recap during office hours be useful?

This content is neither created nor endorsed by Google.

Google Forms