

Lec 11

Thursday, October 3, 2019 10:58

Recap: Clustering

First clustering algo: k-means

Soft clustering / Gaussian Mixture Models

GMM: $Y \sim \text{Bernoulli}(\pi)$ (k=2 clusters
p=1 univariate x) $(X|Y=y) \sim \mathcal{N}(\mu_y, \sigma_y^2)$

This is a generative model for the data

w/ parameters $\theta = (\pi, \mu_0, \sigma_0^2, \mu_1, \sigma_1^2)$

$$\begin{aligned} \mathbb{P}(Y=1|X=x) &= \frac{\pi \varphi\left(\frac{x-\mu_1}{\sigma_1}\right)}{\pi \varphi\left(\frac{x-\mu_1}{\sigma_1}\right) + (1-\pi) \varphi\left(\frac{x-\mu_0}{\sigma_0}\right)} \\ &= \text{cluster responsibility} \end{aligned}$$

Estimate θ by MLE using EM algoGiven observations X_1, \dots, X_n , the log-lik is

$$l(\theta; \mathbf{X}) = \sum_i \log \left((1-\pi) \varphi\left(\frac{x_i-\mu_0}{\sigma_0}\right) + \pi \varphi\left(\frac{x_i-\mu_1}{\sigma_1}\right) \right)$$

↑
actually hard to optimize

Instead:

let's pretend that we also observe

$\gamma_1, \dots, \gamma_n$ (the true cluster membership)

$$l(\theta; \mathbf{X}, \mathbf{Y}) = \sum_{i=1}^n \left((1-\gamma_i) \log \mathcal{Q}\left(\frac{x_i - \mu_0}{\sigma_0}\right) + \gamma_i \log \mathcal{Q}\left(\frac{x_i - \mu_1}{\sigma_1}\right) \right) + \sum_{i=1}^n \left((1-\gamma_i) \log(1-\pi) + \gamma_i \log \pi \right)$$

max of this over π gives $\hat{\pi} = \frac{1}{n} \sum \gamma_i$

only involves $\mu_1, \sigma_1, \mu_0, \sigma_0$

only involves π

This is an easy convex opt problem

$$(*) = \underbrace{\sum_{i: \gamma_i=1} \log \mathcal{Q}\left(\frac{x_i - \mu_1}{\sigma_1}\right)}_{\text{max of this over } \mu_1, \sigma_1} + \underbrace{\sum_{i: \gamma_i=0} \log \mathcal{Q}\left(\frac{x_i - \mu_0}{\sigma_0}\right)}_{\text{Same but for the } \gamma_i=0 \text{ group}}$$

$$\hat{\mu}_1 = \frac{\sum_{i: \gamma_i=1} x_i}{\sum_{i: \gamma_i=1} 1}$$

$$\hat{\sigma}_1^2 = \frac{\sum_{i: \gamma_i=1} (x_i - \hat{\mu}_1)^2}{\sum_{i: \gamma_i=1} 1}$$

E-M trick: replace γ_i w/ currently estimated cluster responsibilities

Expectation - Maximization (EM) algo:

1. Take an initial guess $\hat{\theta} = (\hat{\pi}, \hat{\mu}_0, \hat{\mu}_1, \hat{\sigma}_1, \hat{\sigma}_0)$

2. E-step:

Compute the cluster responsibilities implied by $\hat{\theta}$:

$$\hat{\gamma}_i = \frac{\hat{\pi} \mathcal{Q}\left(\frac{x_i - \hat{\mu}_1}{\hat{\sigma}_1}\right)}{\hat{\pi} \mathcal{Q}\left(\frac{x_i - \hat{\mu}_1}{\hat{\sigma}_1}\right) + (1-\hat{\pi}) \mathcal{Q}\left(\frac{x_i - \hat{\mu}_0}{\hat{\sigma}_0}\right)}$$

Impute $y_i \leftarrow \hat{\delta}_i \in [0, 1]$

(impute unknown y_i by its conditional expectation)

3. M-step

Maximize the likelihood of $\hat{\theta}$ given a complete set of observations

\underline{X} (real obs)

\underline{Y} (imputed obs)

$$\hat{\mu}_1 = \frac{\sum_{i=1}^n \hat{\delta}_i X_i}{\sum_{i=1}^n \hat{\delta}_i}$$

$$\hat{\mu}_0 = \frac{\sum_i (1 - \hat{\delta}_i) X_i}{\sum_i (1 - \hat{\delta}_i)}$$

$$\hat{\sigma}_1^2 = \frac{\sum_i \hat{\delta}_i (X_i - \hat{\mu}_1)^2}{\sum_i \hat{\delta}_i}$$

$$\hat{\sigma}_0^2 = \frac{\sum_i (1 - \hat{\delta}_i) (X_i - \hat{\mu}_0)^2}{\sum_i (1 - \hat{\delta}_i)}$$

$$\hat{\pi} = \frac{1}{n} \sum_i \hat{\delta}_i$$

4. Repeat from step 2

Unsupervised learning w/ abstract distances

So far in unsupervised learning:

$X \in \mathbb{R}^p \rightsquigarrow \mathbb{Z}$ lower dim

either: $\mathbb{G} \mathbb{R}^2$ PCA

$\mathbb{G} \{1, \dots, K\}$ cluster

but we always started
w/ vector data $X \in \mathbb{R}^p$

Next: use more abstract descriptions of the data.

(Dis) Similarity Measures

1. Euclidean dist

$$x \in \mathbb{R}^p \quad x' \in \mathbb{R}^p$$

$$\|x - x'\|_2^2 = \sum_i (x_i - x'_i)^2$$

2. Chi-squared dist

$$\frac{1}{2} \sum_{k=1}^p \frac{(x_{1k} - x'_{1k})^2}{x_{1k} + x'_{1k}}$$

nonnegative vectors representing
a histogram distribution
or a ^{list of} frequencies

3. Cosine-similarity

measure the alignment b/w two vectors

let θ be the angle b/w two vectors,
 x & x' .

$$\text{Cosine Sim} = \cos(\theta) = \frac{x^T x'}{\|x\|_2 \|x'\|_2}$$

$\cos \theta = 1 \Rightarrow$ perfectly aligned

$= -1 \Rightarrow$ perfectly opposite

$$\text{Cos dissim} = 1 - \cos \theta$$

often used in conjunction w/ BoW.

4. Edit distance

Given two strings

$$X = A C G T C C A$$

$$X' = G G T C A C A$$

edit dist = how many insertions, deletions, & replacements minimally needed to go from X to X'

~~A~~ ~~C~~ G T C C A \Rightarrow edit dist = 3
 (Red annotations: "delete" with arrows pointing to A and C, "G" with an arrow pointing to the first G, and "A" with an arrow pointing to the second C)

Multi dimensional Scaling (MDS)

MDS: Given distance/dissimilarity matrix

$$D \in \mathbb{R}^{n \times n}$$

Q: What vectors $X \in \mathbb{R}^{n \times p}$

recover D in the closest way

$$\text{via } D_{ij} \approx \|x_i - x_j\|_2^2$$

I.e. X is embedding of the n abstract datapts in p real dimensions that maintains dissimilarities as vector distances.

If indeed $D_{ij} = \|x_i - x_j\|_2^2$ then


$$D_{ij} = \|x_i\|_2^2 + \|x_j\|_2^2 - 2x_i^T x_j$$

$$= T_{ii} + T_{jj} - 2T_{ij}$$

Where $T = \bar{X} \bar{X}^T \in \mathbb{R}^{n \times n}$
 (n x n matrix whose k_{ij} th entry
 contains inner prod of x_i & x_j)

Solve for T , then eigendecompose
 it to get \bar{X}

MDS: Eigendecompose T to get \bar{X}
 s.t. $T = \bar{X} \bar{X}^T$

Specifically: $T = U \Lambda U^T$  a soln
 up to
 rotation
 & translation
 set $\bar{X} = U_p \Lambda_p^{1/2}$

Fact: if D really came from Euclidean
 distances; then MDS is just PCA