

Nose Oct 4th

Hard Clustering - 1 pc \rightarrow 1 cluster

Soft Clustering - helps w/ overlapping clustering clusters. to give each point a probability of being in a given cluster

$$\theta^* = \arg \max_{\theta} \sum_{i=1}^n \sum_{j=1}^k p(c_j) P(x_i | c_j)$$

Note: taking log transformation doesn't change

$$J(\theta^*) = \sum_{i=1}^n \log \sum_{j=1}^k p(c_j) P(x_i | c_j)$$

↓

$$\text{solve: } \frac{\partial}{\partial \theta} J(\theta) = 0$$

$$\frac{\partial}{\partial \theta} L(\theta) = 0$$

$$\Rightarrow \text{MLE estimate: } \hat{\mu}_j = \frac{\sum_{i=1}^n p(c_j | x_i) x_i}{\sum_{i=1}^n p(c_j | x_i)} \quad \left| \quad \hat{\Sigma}_j = \frac{\sum_{i=1}^n p(c_j | x_i) (x_i - \hat{\mu}_j)^T (x_i - \hat{\mu}_j)}{\sum_{i=1}^n p(c_j | x_i)} \right. \quad \left| \quad \hat{p}(c_j) = \frac{1}{n} \sum_{i=1}^n p(c_j | x_i) \right.$$

$$\text{Bayes rule: } p(c_j | x_i) = \frac{P(x_i | c_j) p(c_j)}{\sum_{j=1}^k P(c_j) P(x_i | c_j)}$$

Clustering Aggregation: Goal: ① Compare clusterings

② Combine info from multi clusterings \Rightarrow create a new clustering

Disagreement Distance:

$$D(p, c) = \sum_{x,y} \mathbb{I}_{p,c}(x, y)$$

Note Our 6th:

SVD : Singular value Decomposition:

- approx A w/ smaller B easier to store
- Dimensionality Reduction / Feature Extraction
- Anomaly Detection & Denoising

Linear Algebra review: Linear independent: linearly represent any cov as linear comb of other vectors in a set iff $A = [v_1, \dots, v_n]$

$$\text{Determinant: } \begin{bmatrix} a & b \\ c & d \end{bmatrix} \Rightarrow \det = ad - bc$$

$$3 \times 3 \Rightarrow \det = a \cdot \det \begin{pmatrix} a & b \\ f & g \end{pmatrix} - b \cdot \det \begin{pmatrix} a & c \\ f & h \end{pmatrix} + c \cdot \det \begin{pmatrix} a & b \\ g & h \end{pmatrix}$$

(Note: $m > n$ vectors in n -dim cannot be linearly independent.

Basis: of a vector space is a linearly independent subset of V that spans V .

Rank: max number of linearly independent cols of A . \rightarrow Full Rank iff $\text{rank}(A) = \min(m, n)$

Approximation: In practice, there are a lot of redundant data.

$$A_{(n \times m)} = U_{(n \times r)} \cdot V_{(r \times m)}$$

$\cdot \|A - A^{(k)}\|$ is small

- Frobenius Distance: $\|A - B\| = \sqrt{\sum_{ij} (a_{ij} - b_{ij})^2}$

when $k < \text{rank}(A)$ $A^{(k)} = \arg \min \{B \mid \text{rank}(B) = k\}$ def (A, B)

$$\text{SVD of rank } r \quad A_{(n \times m)} = U_{(n \times r)} \Sigma_{(r \times r)} V^T_{(r \times m)}$$

Occ 12th / 13th Note.

Classification: Given a training set where data is labeled w/ special attribute called class

Techniques:

- instance based
 - decision tree
 - Naive Bayes
 - Support Vector Machine (SVM)
 - Neural Networks.

kNN: used k closest records to perform classification.

→ limits:
 ↳ If k is too small, then it is sensitive to noise points + over fitting.

 ↳ If k is too big, then neighborhood may include points from other classes.

→ pros:
 ↳ simple to understand why a given unseen record was given a particular class.
 ↳ Add new attributes.

→ cons:
 ↳ expensive to classify points.

 ↳ kNN can be problematic in high dimensions.

Decision Tree

Algorithm: Need to specify terminating conditions: → Let D be set of training rec sets can reach node = .

→ Splitting based on normal attributes:

- multiway split → use as many partitions as values.
- Binary split → categorized as binary.

$$O \in 10^{56} / 2^{32}$$

Recommendation System:
 (user, user) similarity measure.
 (item, item) similarity measure

Pros \rightarrow no training required [Bias].

Cons \rightarrow users rate differently rating changing over time.

Content Based: feature to movie similarity.

Matrix: If we had user -> feature \Rightarrow feature to movie = user to movie

$$p^{(i)} = \underset{p}{\operatorname{argmin}} \frac{1}{\|u_i\|_F} \sum_{j \neq i} (p^T q^{(j)} - r_{ij})^2 + \lambda \|p\|_F^2$$

penalty on the size of the param p.

Collaborative Feature

can we learn P, Q. $R = PQ$ \Rightarrow can't learn using SVD as it is sparse.

$$\min_{p, q} \sum_{i,j} (r_{ij} - p^T q_j)^2 + \lambda (\|p\|_F^2 + \|q\|_F^2)$$

- Steps:
- 1. Score w/ random Q.
- 2. Get P
- 3. Improve Q.
- 4. Repeat 2 & 3

$$\text{Least Square: } \hat{\beta}_w = \underset{\beta}{\operatorname{argmin}} \sum_i d(h_{\beta}(x_i), y_i) = \underset{\beta}{\operatorname{argmin}} \|h_{\beta} - y\|_2^2 = \hat{\beta}_w.$$

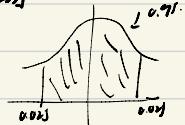
Maximum Likelihood: Since $\epsilon \sim N(0, \sigma^2)$ and $y = x\beta + \epsilon$, then $y \sim N(x\beta, \sigma^2)$, $\hat{\beta}_w = (x^T x)^{-1} x^T y$

Oct 27th. Note.

Evaluate Regression Model: $y_i \sim N(5, 25)$ for $1 \leq i \leq 100$ and $y_i = \mu + \epsilon$ where $\epsilon \sim N(0, 25)$

Then, the LSS for μ ($\mu_{\text{LS}} = \text{sample mean}$) $\hat{\mu} = \frac{1}{n} \sum y_i$.

$$\text{CL} = [\hat{\mu} - 1.96(\text{SE}(\mu_{\text{LS}})), \hat{\mu} + 1.96(\text{SE}(\mu_{\text{LS}}))]$$



$$z\text{-values: } \int_{-2}^2 \frac{1}{2\pi} e^{-\frac{1}{2}x^2} dx = 0.95$$

Q-Q plot: need to check our variance is distributed in the same way as variable following our target distribution

constant variance: we assume that our noise has constant variance

we can plot our fitted values against our residuals (noise estimates)