

Parzen - Rosenblatt Estimator

- Nonparametric Density Estimation

Kernel Window Function:

$$k(\vec{x}_i, \vec{x}) = \begin{cases} 1, & \frac{|\vec{x}_i - \vec{x}|}{h} \leq 1/2 \\ 0, & \text{otherwise} \end{cases}$$

(Box Kernel)

$$p_R = p(\vec{x}) \cdot V_R \\ = \frac{k_R}{N} \cdot V_R$$

$$\Rightarrow p(\vec{x}) = \frac{k_R}{N \cdot V_R}$$

→ use box kernel →

$$p(\vec{x}) = \frac{1}{N \cdot h^d} \cdot \sum_{i=1}^N k_h(\vec{x}_i, \vec{x})$$

Cross Validation

$$\hat{h} = \underset{h}{\operatorname{argmax}} \alpha(h) \stackrel{ML}{=} \underset{h}{\operatorname{argmax}} \frac{1}{N} \sum_{j=1}^N p_{h, N-1}^j(\vec{x}_j)$$

$$S = \{x_1, \dots, x_N\}$$

$$S_j = S \setminus \{x_j\} \leftarrow \text{test sample}$$

$$= \underset{h}{\operatorname{argmax}} \sum_{j=1}^N \log \left[p_{h, N-1}^j(\vec{x}_j) \right]$$

Mean Shift Algorithm

1. Compute Mean Shift Vector:

$$\sum_{i=1}^N k'_h(\|\vec{x}_i - \vec{x}\|^2) \cdot \vec{x}_i$$

$$\sum_{i=1}^N k'_h(\|\vec{x}_i - \vec{x}\|^2)$$

$$- \vec{x} = \vec{0}$$

2. Update \vec{x} :

$$\vec{x}^{(t+1)} = \vec{x}^{(t)} + m(\vec{x}^{(t)})$$

$$= \frac{\sum_{i=1}^N k'_h(\|\vec{x}_i^{(t)} - \vec{x}^{(t)}\|^2) \cdot \vec{x}_i^{(t)}}{\sum_{i=1}^N k'_h(\|\vec{x}_i^{(t)} - \vec{x}^{(t)}\|^2)}$$

$$+ \vec{x}^{(t)} - \vec{x}^{(t)}$$

Epanechnikov Kernel:

$$k_E(x) = \begin{cases} c \cdot (1 - x^T x), & x^T x \leq 1 \\ 0, & \text{otherwise} \end{cases}$$

K-Means

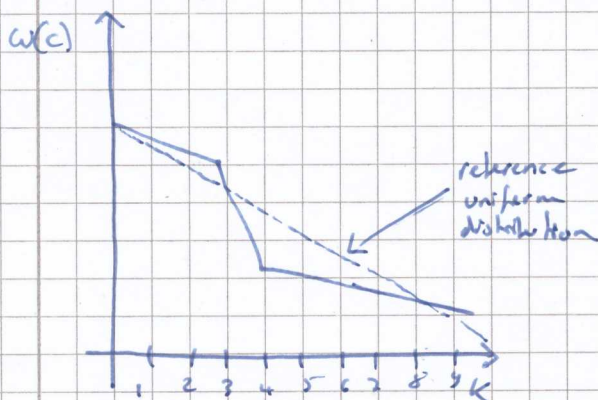
1. For a given cluster assignment C : compute μ
2. Assign each observation to the closest μ
3. Repeat 1. and 2. until convergence

$$W(C) = \frac{1}{2} \sum_{k=1}^K \sum_{d(i)=k} \sum_{C(i')=k} d(x_i, x_{i'})$$

Using medoids increases

$$O(N) \rightarrow O(N^2) \quad (\text{additional pair checking})$$

K model selection problem



Tibshirani: Gap Statistics

Tibshirani Significance:

- 1.) Repeat sampling B times
- 2.) calculate standard deviation of $W(C)$ for each K of the synthetic curve

$$K^* = \underset{K}{\operatorname{argmin}} \left\{ K \mid G(K) \geq G(K+1) - S'_{K+1} \right\}$$

\uparrow
 $E(\log(W_K) - \log(W_{K+1}))$

Hierarchical Clustering

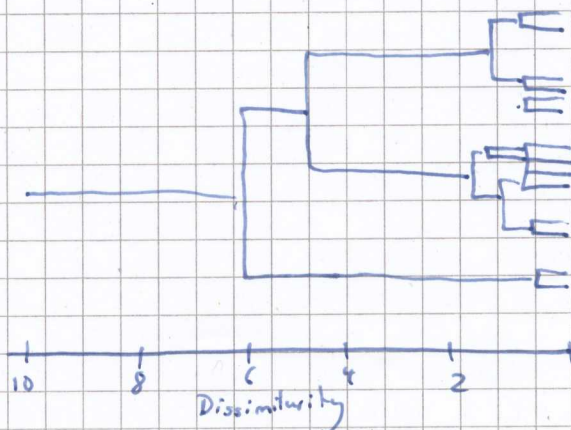
- top-down: "divisive clustering"
- bottom-up: "agglomerative clustering"

Single Linkage: Merge least dissimilar points
 \rightarrow chains may form

Complete Linkage: merge where largest sample distance is small

Group Average: merge where avg distance is small

Dendrogram:



Model selection for GMMs

- Dirichlet Process Mixture
- Prior distribution of GMM Parameters

$$p(\vec{x}) = \sum_{k=1}^K \pi_k \mathcal{N}(\vec{x} | \vec{\mu}_k, \Sigma_k)$$

$$p_{ik} \equiv \gamma(z_{nk})$$

$$p_k = \pi_k$$

GMM Expectation Maximization: (classic)

0.) Initialize $\pi_k, \vec{\mu}_k, \Sigma_k$

1.) "Expectation":

$$\gamma(z_{nk}) = \frac{\pi_k \cdot \mathcal{N}(\vec{x}_n | \vec{\mu}_k, \Sigma_k)}{p(\vec{x}_n)}, \quad p(\vec{x}_n) = \sum_{j=1}^K \pi_j \mathcal{N}(\vec{x}_n | \vec{\mu}_j, \Sigma_j)$$

2.) "Maximization":

$$\vec{\mu}_k^{\text{new}} = \frac{1}{N_k} \cdot \sum_{n=1}^N \gamma(z_{nk}) \cdot \vec{x}_n$$

weighted average

\sum_k \otimes 10% \otimes 100% \otimes 90%

$$\Sigma_k^{\text{new}} = \frac{1}{N_k} \cdot \sum_{n=1}^N \gamma(z_{nk}) \cdot (\vec{x}_n - \vec{\mu}_k^{\text{new}}) \cdot (\vec{x}_n - \vec{\mu}_k^{\text{new}})^T$$

$$\pi_k^{\text{new}} = \frac{N_k}{N}, \quad N_k = \sum_{n=1}^N \gamma(z_{nk})$$

Gibbs sampler:

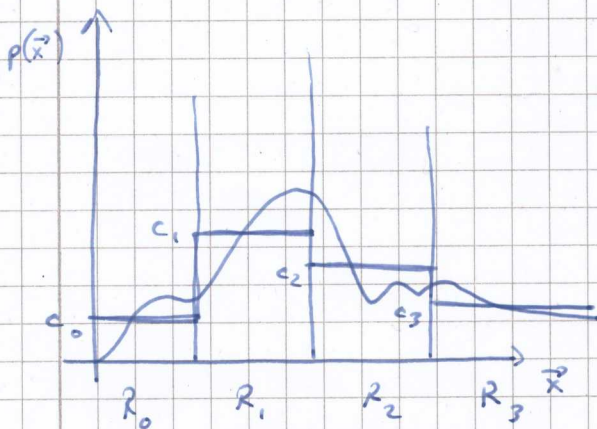
0.) Assign initial clusters

1.) Randomly choose sample \vec{x}_i , remove from clustering

2.) Assignment to a new sent:

$$\underbrace{p(z_i = k)}_{\text{"rich get richer"}} \cdot \underbrace{p(\vec{x}_i | \mu_k, \Sigma_k)}_{\text{"do I like the food?"}} = \frac{N_k}{N + \alpha_0} \cdot \mathcal{N}(\vec{x}_i | \vec{\mu}_k, \Sigma_k)$$

$$\underbrace{\frac{\alpha_0}{N + \alpha_0}}_{\text{for a new cluster}} \cdot \mathcal{N}(\vec{x}_i | \mu_0, \Sigma_0)$$

Regression Tree $R_j \hat{=}$ leaf node j

$$\hat{c}_j = \operatorname{argmin}_{c_j} \int_{R_j} (p(\vec{x}) - \hat{p}(\vec{x}))^2 d\vec{x}$$

$$= \operatorname{argmin}_{c_j} \int_{R_j} \left(p(\vec{x}) - \sum_{j=1}^K c_j I(\vec{x} \in R_j) \right)^2 d\vec{x}$$

$$c_j = \frac{\int_{R_j} p(\vec{x}) d\vec{x}}{\int_{R_j} 1 d\vec{x}}$$

 $\Rightarrow c_j$ is the mean of the values in R_j Mitigate overfittingBagging: Train many trees on random subsets of dataBoosting: Train a sequence of classifiers with increasing difficulty. Make a weighted vote (AdaBoost, Viola & Jones)Random Forests: Randomize Training Process :

- random data subset
- random dimension for split
- finite subset of parameters for split function

Entropy:

Classification

$$H(S) = - \sum_{c \in C} p(\vec{x} \in c) \cdot \log(p(\vec{x} \in c))$$

 C : all class labels

Regression

$$H(S) = - \frac{1}{|S|} \sum_{x \in S} \int_Y p(y|x) \cdot \log(p(y|x)) dy$$

 c : class label

$$p(y|x) = \mathcal{N}(y; \bar{y}(x), \bar{\sigma}_y^2(x))$$

