

SML - Formelsammlung

Probability distributions

Normal distribution pdf: $\frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2\right)$; if $p(x,y)$ is gaussian, then conditional $p(x|y)$ and marginal $p(x)$ are

Expectation: $E(X) = \int x f(x) dx$, Variance: $E((X-\mu)^2) = E(X^2) - E(X)^2$ also gaussian

Marginal distribution: $f_X(x) = \int_{-\infty}^{\infty} f(x,y) dy$; Conditional dist.: $f(x|y) = \frac{f(x,y)}{f_Y(y)}$

Moments: $M_n = E(X^n)$; Central Moments: $c_{mn} = E((x-\mu)^m)$; c_{m2} : variance, c_{m3} : skewness, c_{m4} : kurtosis

Bayes rule: $p(y|x) = \frac{p(x|y)p(y)}{p(x)}$; posterior or likelihood · prior; $p(x)$ = Normalization factor = $\sum_j p(x|c_j)p(c_j)$

KL-Divergence: $KL(p, q) = - \int p(x) \ln \frac{q(x)}{p(x)} dx$

Bayesian Decision Theory

Minimizing error: $p(\text{error}) = p(x \in R_1, C_2) + p(x \in R_2, C_1) = \int_{R_1} p(x|C_2)p(C_2) dx + \int_{R_2} p(x|C_1)p(C_1) dx$

Optimal decision: Decide for class C_1 if $p(C_1|x) > p(C_2|x)$ or $\frac{p(x|C_1)}{p(x|C_2)} > \frac{p(C_2)}{p(C_1)}$ (likelihood-Ratio-Test) Bayes optimal classifier

Risk minimization: loss function $\lambda(a_i|C_j) = \lambda_{ij}$ with decision a_i and true class C_j

Expected Loss: $R(\alpha_i|x) = \sum_j \lambda(a_i|C_j) \cdot p(C_j|x)$ class posterior

Decide for a_i if $R(\alpha_2|x) > R(\alpha_1|x)$ or $\frac{p(x|C_2)}{p(x|C_1)} > \frac{(r_{22} - r_{12})}{(r_{21} - r_{11})} \cdot \frac{p(C_2)}{p(C_1)}$

Probability Density Estimation

Is about finding the class conditional prob. $p(x|C_k)$

Parametric Models: Small number of parameters completely define pdf

Non-parametric: No explicit parameters (except hyperparameters) but every point is a parameter

Mixture Models: Combination of both

Maximum Likelihood:

PDF is given by parameters in θ : $p(x|\theta)$

We want to find θ that explains data the best

$$L(\theta) = p(D|\theta) = p(x_1, x_2, \dots | \theta)$$

$$= p(x_1|\theta) \cdot p(x_2|\theta) \cdots \Rightarrow \text{i.i.d. assumption}$$

$$= \prod_{i=1}^N p(x_i|\theta)$$

$$\text{or } \ln L(\theta) = \sum \ln p(x_i|\theta)$$

To find maximum, derive w.r.t. θ , set 0 and

solve

Does not work for single point,

provides point estimates

Clustering: Unsupervised learning, unlabeled data

Goal: Find similar datapoints and put them in discrete set of clusters

K-means clustering:

1. Initialization, pick k centroids

2. Assign each datapoint to closest centroid

3. Adjust centroids to be clusters means

4. Back to step 2 until no change

Problems: Needs specified k , sensitive to initialization, only local optimum, only spherical clusters

Mean shift clustering: Assign data to clusters by iteratively shifting datapoints to the mode (region with highest density) of a specified region. In the end, assign to points to the cluster they converge to.

Problems: Needs specified window size, hard to compute

Histograms: Split data to bins, count datapoints in each bin
Problem: Doesn't scale well to high dimensions.

Needs specified bin-width

Kernel-Density-Estimation: Fix volume V and determine number of datapoints N

$$p(x) \approx \frac{N}{NV} = \frac{1}{NV} \sum_{i=1}^N k\left(\frac{\|x-x_i\|}{V}\right) \text{ with Kernel } k,$$

Problem: Need to choose

$$\text{e.g. gaussian kernel } k(w) = \frac{1}{\sqrt{2\pi}} \exp\left(-\frac{w^2}{2}\right)$$

Kernell and bandwidth,

Choice of kernel not so important for many datapoints or chosen window width w = $\begin{cases} 1 & \text{if } N < 0.5 \\ 0 & \text{otherwise} \end{cases}$

K-Nearest Neighbor:

Calculate distance to k th neighbor. High distance means lower density

$$\text{Mixture models: } p(x) = \sum_{j=1}^M p(x|z_j) \cdot p(z_j)$$

(see EM algorithm) \rightarrow conditional prob. of x given its in cluster j ; \rightarrow Prior of clusters

EM-Clustering: Soft assignment to underlying prob. distribution

1. E-step: Assign each datapoint prob. of belonging to clusters j : $p(C_j|x) = \frac{p(x|C_j) \cdot p(C_j)}{\sum_{i=1}^M p(x|C_i) \cdot p(C_i)}$

2. M-step: Update priors, means and variances

$$p(C_j)^{\text{new}} = \frac{1}{N} \sum_{i=1}^N p(C_j|x_i)$$

$$\mu_j^{\text{new}} = \frac{\sum_{i=1}^N p(C_j|x_i) \cdot x_i}{\sum_{i=1}^N p(C_j|x_i)}$$

$$\sigma_j^{\text{new}} = \sqrt{\sum_{i=1}^N p(C_j|x_i) \cdot (x_i - \mu_j)^2 / \sum_{i=1}^N p(C_j|x_i)}$$

3. Repeat

Problems: Sensitive to initialization, only local optimum, needs specified number of mixture components

Evaluation

Ockham's Razor: choose smallest model that fits the data

Bias vs Variance: Too simple model has high bias, underfitting. Too complex model has high variance, overfitting

Maximum Likelihood Estimation is not always unbiased, e.g. Variance estimator of gaussian

