Probability Density Estimation

Overview

Parametric Methods

- ✓ Maximum Likelihood Estimation Method
- ✓ Bayesian Inference Method
- □ Non-parametric Methods
 - ✓ Histogram
 - ✓ Kernel-based Approaches
 - ✓ K-Nearest Neighbor (KNN) Approach
- □ Kullback-Leibler (KL) Distance
- □ Semi-Parametric Methods : Mixture Models (Gaussian Mixture Models (GMMs))
 - ✓ Maximum Likelihood (Nonlinear Optimization)
 - ✓ The EM Algorithm

Overview (Cont..)

Parametric Methods

- ✓ Specific functional form is chosen
- \checkmark Evaluation with new data is faster
- □ Non-parametric Methods
 - \checkmark No specific functional form
 - ✓ Number of parameters grows with number of data points
 - ✓ Very slow for evaluation of new data
- □ Semi-Parametric Methods
 - \checkmark No specific functional form
 - \checkmark Number of parameters grows with complexity
 - ✓ Computationally intensive

Parametric Methods

EX: Nômal & Guussian distribution Estimation of parameters of a Model - Maximum licelyhood estimation Method - Bayesian Inference Method Normal Distribution (Single Variable) $P(x) = \frac{1}{(2\pi\sigma^{2})^{\frac{1}{2}}} \exp\left[-\frac{(x-\mu)^{2}}{2\sigma^{2}}\right]; \int_{-\infty}^{\infty} b(x) = 1$ M - Meun; J = Stundard Deviation; J = Variance $M = E[x] = \int x p(x) dx = F Expectation of x$ $\sigma = \varepsilon \left[(x - \mu)^{2} \right] = \int_{-\infty}^{\infty} (x - \mu)^{2} \dot{\rho}(x) dx$ E[.] = Eppectation



MAXIMUM LIKELIHOOD ESTIMATION

Optimum Values of the parameters by Hapimiting the Licelihood function desived from Training data parameter set a = [a1, a2, ... on] Data set $X = [x', x', \dots, x]$ bat x = b(x[o)) $pof(X) = \frac{N}{11} p(\tilde{x}|o) = L(o)$ L(Q) = Liculihool g O given X ML function: $p(X|\alpha) = \prod_{x=1}^{N} p(x|\alpha) = L(\alpha)$ Ervol function $E = -\ln L(\alpha) = \sum \ln \left| p(x'|\alpha) \right|$



$$\begin{aligned}
\left(\begin{array}{c} 0 \mid \chi \end{array}\right) &= \frac{\beta(\chi \mid \alpha)}{\beta(\chi)} \frac{\beta(\alpha)}{\beta(\chi)} &= \frac{\prod_{n \in I} \beta(\chi' \mid \alpha)}{\beta(\chi)} \frac{\beta(\alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \frac{\prod_{n \in I} \beta(\chi \mid \alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)}{\beta(\chi)} \\
&= \frac{\beta(\alpha)$$

Non-Parametric Methods

Histograms

- □ Kernel-based approaches
- □ K-Nearest Neighbor (KNN)
- □ Histograms (Simple Models)
 - Role of smoothing parameters
 - \checkmark Smoothing parameter \rightarrow Number of bins
 - ✓ Lower number of bins (smooth out crucial details)
 - ✓ Higher number of bins (spiky& noisy)
 - Advantages
 - \checkmark Can be constructed sequentially, entire data is not required at one shot. No storage issues.
 - Disadvantages
 - \checkmark Estimated density is not smooth
 - ✓ Curse of dimensionality with high-dimensionality data



Non-Parametric Methods : Density Estimation

Demosity Estimation in General P - prob of vector x drawn from whilerow density p(x) will fall instille some region R of X - Mule is $P = \int \varphi(x') Q x'$ Prob & K pointe fall within region R out of N prints druce $P_r(ic) = \frac{N!}{ic!(N-ic)!} P_i^{ic}(i-p) = Bimonial (aur$ Mean of IC points Juli in these regions is given by Mean - E [X] K/N] = P 2 K N Variance = $E\left[\left(\frac{k}{N} - P\right)^{V}\right] = \frac{P(I-P)}{N}$ $x \neq p(x)$ in continuum => $p = \int p(x') dx' = p(x) V$ =) $\frac{1}{N} = \left| \rho(x) V \right| = \frac{R}{|\rho(x)|} = \frac{1}{|NV|}$

Non-Parametric Methods : Density Estimation

Practicul Demnity Estimation

$$P \approx \frac{1}{N} \Rightarrow Region "R" (c be relatively large => P will be large.
P \approx P(x) v => p(x) = \frac{P}{V} = \frac{1}{NV} => & p(x)^{-1} = a(p) \frac{1}{NV} constant when "R"
is relatively should.
Choice of R for the best estimate of p(x)
K- Nearest Neigh (hor (ICNN))
Fip (the value of (C => Det the Corresponding value "U" from data
(Cernel - based durstly
Fip (the value of (C => Det "K" from data
Fip (the value of (C => Det "K" from data
Fip (the value of (C => Det "K" from data
For infinite "N" => Bath (CNN & Icernel band Mataras
will converge.$$

Non-Parametric Methods : Kernel-based Method

Kernel Barel Metwog
Region "R" to be hypercube with Aiden of lawyth h cartoned
on the point x in d-dimensional Africe
Volume of the hypercube
$$V = h^d$$

Kernel function $H(U) = \begin{cases} 1 & |u_3| < \frac{1}{2} & 3 = 1,1,...d. \\ 0 & 0 (thermine) \end{cases}$
 $H\left(\frac{x-x^2}{h}\right) = 1$; if x^n fuch within hypercube
 $K = \sum_{n=1}^{N} H\left(\frac{x-x^n}{h}\right)$
 $\int_{P(x)} = \frac{K}{NV} = \int_{N} \sum_{m=1}^{N} \frac{1}{h} H\left(\frac{x-x^n}{h}\right)$
Estimated Dentity = Superposition of N cubes of Analith, with each
cubes V3 bins in the print y hintery hum

Non-Parametric Methods : Kernel-based Method



Non-Parametric Methods : Kernel-based Method

Drawbacks

- \checkmark All data points to be stored
- \checkmark Evaluation of density is very slow if # data points is

large

 \checkmark It gives biased estimate of density

□Solution to fast density estimation

- ✓ Use of fewer kernel functions
- \checkmark Adapt their positions and width as per the data

Non-Parametric Methods : K-Nearest Neighbor

 \Box Problem with Kernel-based approach \rightarrow Fixed "h" for all data points

 $\Box \text{Remedy} \rightarrow \text{KNN} : \text{Fix "K" value & allow "V" to vary}$

Disadvantages of K-Nearest Neighbor (KNN)

- Estimated PD is not true density, because integral over all Xspace diverges.
- \checkmark All the training data points must be retained.
- ✓ Large computer storage is required to store the whole data
- Require large amount of processing to evaluate the density at new data points.

Non-Parametric Methods : KNN-Classifier



Non-Parametric Methods : Smoothing Parameters



Semi-Parametric Models : Mixture Models

Parametric Methods

- \checkmark Specific functional form is chosen
- \checkmark Evaluation with new data is faster
- □ Non-parametric Methods
 - \checkmark No specific functional form
 - ✓ Number of parameters grows with number of data points
 - ✓ Very slow for evaluation of new data
- □ Semi-Parametric Methods
 - \checkmark No specific functional form
 - ✓ Number of parameters grows with complexity
 - ✓ Computationally intensive

Semi-Parametric Models : Mixture Models (Cont..)



Semi-Parametric Models : Mixture Models (Cont..)

In complete Date : (on ponent labels one not (mound
clanification:
$$\beta(x|ce) - claim Aperite Holder
From Bayels Theolow: $P(T|x) = \frac{\beta(x|T)\beta(T)}{p(x)}$
 $\stackrel{H}{\underset{T=1}{\overset{H}{\underset{T=1}{\sum}}} = \frac{1}{(2\pi\sigma_T^2)^2} \exp\left\{-\frac{1|x-M_3|^2}{2\sigma_T^2}\right\}$$$



Mixture Models : Maximum Likelihood Technique

Mixture Models : ML Technique (Cont..)

$$\frac{M \ L \ Technique \left((ont...) \right)}{\frac{\partial E}{\partial H_{3}} = \sum_{m=1}^{N} P(\exists | x^{m}) \frac{(H_{1} - x^{m})}{G_{3}}}{\frac{\partial E}{\partial G_{7}} = \sum_{m=1}^{N} P(\exists | x^{m}) \left\{ \frac{d}{G_{7}} - \frac{||x^{m} - H_{7}||^{2}}{G_{7}^{2}} \right\}}{\frac{\partial E}{\partial G_{7}} = \sum_{m=1}^{N} P(\exists | x^{m}) \left\{ \frac{d}{G_{7}} - \frac{||x^{m} - H_{7}||^{2}}{G_{7}^{2}} \right\}}{\frac{\partial E}{\partial G_{7}} = -\sum_{m=1}^{N} \left\{ P(\exists | x^{m}) - P(\exists) \right\}}$$

$$\frac{\partial E}{\partial V_{7}} = -\sum_{m=1}^{N} \left\{ P(\exists | x^{m}) - P(\exists) \right\}}{\sum_{m=1}^{N} P(\exists | x^{m}) x^{m}} / \sum_{m=1}^{N} P(\exists | x^{m})}$$

$$\frac{\partial E}{\partial G_{7}} = -\sum_{m=1}^{N} \left\{ \frac{\sum_{m=1}^{N} P(\exists | x^{m}) - P(\exists)}{\sum_{m=1}^{N} P(\exists | x^{m})} \right\}}{\sum_{m=1}^{N} P(\exists | x^{m}) \prod_{m=1}^{N} \sum_{m=1}^{N} P(\exists | x^{m})}$$

Mixture Models : The EM Algorithm

The EM Algorithm No direct solm to compute parameter MJ, JE P(J) Highly non-lineur coupled egos => Difficult to solve. Iteratule Scheme to find minimum of E - Initiate with (quem) some parameters consumers old - Evaluate Ries of MT, GT & P(T) => New surameters - New paremeters => Errol fr. 1 - New paramean -> old parameter Refreat above steps until local anim in nucled. event - Eold = - E In [prove (x)] preur (x) 2 pole (x) => Densition evaluated using $e^{new} = -\sum_{n} \left[m \left[\frac{\sum_{j=1}^{new} (\tau_{j})}{p^{old} (\chi^{m})} + \frac{\sum_{j=1}^{new} (\chi^{m})}{p^{old} (\chi^{m}$



Mixture Models : The EM Algorithm (Example)



Mixture Models : The EM Algorithm (Incomplete Data)

Mixture Models : The EM Algorithm (Incomplete Data)

$$EM Algorithm - Grouplate Buta (ut-)$$

$$Minimi tution g E [E^{(u-b)}] =) Maximi tuto g E [L^{(u-b)}]$$

$$H-Step =) Maximi tuto g E [L^{(u-b)}]$$

$$E [E^{(u-b)}] = \sum_{2l=1}^{H} \sum_{2l=1}^{H}$$

Illustration of Use of GMMs : Isolated Word Rec

