

CSCI 567 Machine Learning HW1

- Ankit Kothari
USC ID:4138854460

1 DENSITY ESTIMATION

(a)(i)

Beta Distribution $f(x) = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{B(\alpha, \beta)}$ where $B(\alpha, \beta) = \frac{\Gamma(\alpha)\Gamma(\beta)}{\Gamma(\alpha + \beta)}$

and $\Gamma(t) = (t-1)!$

Therefore $Beta(\alpha, \beta) = \frac{1}{\alpha}$

$f(x) = x^{\alpha-1} \alpha$

$L(\alpha|x) = \prod_1^N x^{\alpha-1} \cdot \alpha$

Taking log likelihood

$l(\alpha|x) = \sum_1^N \log(x^{\alpha-1} \cdot \alpha)$

$l(\alpha|x) = \sum_1^N \log(x^{\alpha-1}) + N \log(\alpha)$

Taking derivative wrt α

$= \frac{dl(\alpha)}{d\alpha} \sum_1^N \frac{1}{x^{\alpha-1}} \cdot x^{\alpha-1} \log(x) + \frac{N}{\alpha} = 0$

$= \sum_i^N \log(x) + \frac{N}{\alpha} = 0$

therefore, $\alpha = \frac{-N}{\sum_i^N \log(x)}$

(a)(ii)

$$f(x) = \frac{1}{\sqrt{2\pi\sigma^2}} \cdot \exp \frac{-(x-\mu)^2}{2\sigma^2}$$

Given $\mu = \Theta$, $\sigma = \Theta$

$$f(x) = \frac{1}{\sqrt{2\pi\theta}} \cdot \exp \frac{-(x-\theta)^2}{2\theta}$$

Likelihood function is given by,

$$L(\theta|x) = \prod_1^N f(X_i|\theta)$$

Log Likelihood function is given by,

$$l(\theta) = \sum_{i=1}^N \log \left(\frac{1}{\sqrt{2\pi\theta}} \cdot \exp \frac{-(x_i - \theta)^2}{2\theta} \right)$$

Differentiating wrt θ ,

$$\frac{dl(\theta)}{d\theta} = \frac{-1}{2} \left[-\frac{N}{\theta} - \sum \frac{x^2}{\theta^2} + N \right]$$

$$\frac{dl(\theta)}{d\theta} = \frac{-1}{2} \left[-\frac{N}{\theta} - \sum \frac{x^2}{\theta^2} + N \right] = 0$$

$$\theta = \frac{-N \pm \sqrt{N^2 + 4N \sum x^2}}{2N}$$

$$\therefore \theta = \frac{1}{2} \left[-1 \pm \sqrt{1 + \frac{4 \sum x^2}{N}} \right]$$

(b)

(i)

$$\text{To Show } E_{x_i, \dots, n}[\hat{f}(x)] = \frac{1}{h} \int K\left(\frac{x-t}{h}\right) f(t) dt$$

$$\hat{f}(c) = \frac{1}{n} \sum_{i=1}^n \frac{1}{h} K\left(\frac{x-X_i}{h}\right)$$

$$E\hat{f} = \frac{1}{nh} \sum_{i=1}^n E\left[K\left(\frac{x-X_i}{h}\right)\right] \quad \text{As } E \sum_{x_i} = \text{sum} E(x_i)$$

$$= \frac{1}{nh} [nE(K(\frac{x-X_i}{h}))]$$

$$= \frac{1}{h} \int K\left(\frac{x-X_i}{h}\right) f(X_i) dX_i \quad E(g(x)) = \int g(x) f(x) dx$$

$$= \frac{1}{h} \int K\left(\frac{x-t}{h}\right) f(t) dt$$

(ii)

$$E[\hat{f}(x)] = \frac{1}{n} \sum_{i=1}^n nE\left[\frac{1}{h} K\left(\frac{x-X_i}{h}\right)\right]$$

$$= E\left[\frac{1}{h} K\left(\frac{x-X}{h}\right)\right]$$

$$= \int dt \frac{1}{h} K\left(\frac{x-t}{h}\right) f(t)$$

$$\text{since } z = \frac{x-t}{h} \quad \therefore -h dz = dt$$

$$= \int dz K(z) f(x-hz) \dots (\text{changing limits and then eliminating } h \text{ and changing limits again})$$

$$= \int dz K(z) [f(x) - hzf'(x) + \frac{h^2 z^2}{2} f''(x) - \frac{h^3 z^3}{3!} f'''(x) + \dots]$$

$$= f(x) + \frac{h^2 f''(x)}{2} \int dz K(z) z^2 + \frac{h^3 f'''(x)}{3!} \int dz K(z) z^3 + \dots$$

by definition $\int dz K(z) = 1$ and $\int dz K(z) z = 0$. If we call $\int dz K(z) z^2 = \sigma_k^2$, then bias of KDE is

$$E[\hat{f}(x)] - f(x) = \frac{h^2 \sigma_k^2 f''(x)}{2} + o(h^2)$$

2 DENSITY ESTIMATION

(a)

$$P(X_1 \dots X_n | Y) = \prod_{i=1}^n P(X_i | Y)$$

Using Naive Bayes Thm

$$P(Y = y_k | X_1 \dots X_n) = \frac{P(Y = y_k) \prod_{i=1}^n P(X_i | Y = y_k)}{\sum P(Y = y_j) \prod_{i=1}^n P(X_i | Y = y_j)}$$

$$P(Y = 1 | X) = \frac{\pi \prod_{i=1}^n P(X_i | Y = 1)}{\pi \prod_{i=1}^n P(X_i | Y = 1) + (1 - \pi) \prod_{i=1}^n P(X_i | Y = 0)}$$

$$= \frac{1}{1 + \frac{(1 - \pi) \prod_{i=1}^n P(X_i | Y = 0)}{\pi \prod_{i=1}^n P(X_i | Y = 1)}}$$

Given $P(X_i | Y = y_k)$ follows a Gaussian distribution of the form $N(\mu_{jk}, \sigma_j^2)$

for $X_{i=0}$

$$P(X_0 | Y = 0) = \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp^{-\left(\frac{(x - \mu_{00})^2}{2\sigma_0^2}\right)} dx$$

$$P(X_0 | Y = 1) = \frac{1}{\sqrt{2\pi\sigma_0^2}} \exp^{-\left(\frac{(x - \mu_{01})^2}{2\sigma_0^2}\right)} dx$$

For X_i

$$\frac{\prod P(X_i | Y = 0)}{\prod P(X_i | Y = 1)} = \exp^{-\left[\sum \frac{(x_i - \mu_{i0})^2 - (x_i - \mu_{i1})^2}{2\sigma_i^2}\right]}$$

$$= \exp^{-\left[\sum \frac{2x_i(\mu_{i1}^2 - \mu_{i0}^2)}{2\sigma_i^2} - \log\left(\frac{1 - \pi}{\pi}\right)\right]}$$

$$= \exp^{-\left[\log\left(\frac{1 - \pi}{\pi}\right) - \sum \frac{x_i(\mu_{i1} - \mu_{i0})}{\sigma_i^2} + \sum \frac{(\mu_{i1}^2 - \mu_{i0}^2)}{2\sigma_i^2}\right]}$$

$$\therefore w_0 = -\left[\frac{\log(1 - \pi)}{\pi}\right] + \sum_{i=1}^n \frac{(\mu_{i1}^2 - \mu_{i0}^2)}{2\sigma_i^2}$$

$$w_i = \frac{(\mu_{i0} - \mu_{i1})}{\sigma_i^2} \quad w = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_k \end{bmatrix} \quad 4$$

(b)

$$\begin{aligned} P(X = x, Y = y) &= P(Y = y)P(X = x|Y = y) \\ &= P(Y = y) \prod_{d=1}^D DP(X_d = x_d|Y = y) \end{aligned}$$

$P(Y = k) = \pi_k$ for $k \in 0, 1$ given y follows Bernoulli Distribution

$$P(x_j|Y = y_k) \sim N(\mu_{jk}, \sigma_{jk})$$

$$\log(P(x_i|Y = k)) = \frac{\log(2\pi\sigma_{jk})}{2} - \frac{(x_j - \mu_{jk})^2}{2\sigma_{jk}}$$

Using Naive Bayes

$$\begin{aligned} P(X_i = x, Y = y) &= P(Y_i = y)P(X_{i1} = x_1, X_{i2} = x_2 \dots X_{iD} = x_D|Y = y_k) \\ &= P(Y_i = y) \prod_{j=1}^D P(X_{ij} = x_{ij}|Y = y_i) \end{aligned}$$

$$\begin{aligned} \log \text{likelihood function } \log(L) &= \sum_{i=1}^N \log(P(Y_i = y_i)) + \sum_{i=1}^N \sum_{j=1}^D \log(P(X_i = x_{ij}|Y = y_i)) \\ &= \sum_{k=1}^K \pi_k P(Y = k) N_k + \sum_{k=1}^K \sum_{j=1}^D \log(P(X_{ij} = x_{ij}|Y = k)) N_k \end{aligned}$$

Differentiating w.r.t the params π , μ and σ

$$\text{For } \pi, \sum_k \pi_k = 1, \text{ we get } \pi_k \frac{N_k}{N}$$

$$\frac{\partial \sum_{k=1}^K \sum_{j=1}^D \log(P(X_i = x_{ij}|Y = k)) N_k}{\partial \mu_{jk}} = 0$$

$$\mu_{jk} = \frac{\sum_{i; Y_i=k} x_{ij}}{N_k}$$

$$\frac{\partial \sum_{k=1}^K \sum_{j=1}^D \log(P(X_i = x_{ij}|Y = k)) N_k}{\partial \sigma_{jk}} = 0$$

$$\therefore \sigma = \frac{\sum_{i; Y_i=k} (x_{ij} - \mu_{jk})^2}{N_k}$$

3 NEAREST NEIGHBOR

(a)

$$\bar{x} = \frac{1}{N} \sum_{i=1}^{13} x_i = \frac{166}{13} = 12.769$$

$$\bar{y} = \frac{1}{N} \sum_{i=1}^{13} y_i = \frac{160}{13} = 12.307$$

$$\sigma(x) = \frac{1}{N-1} \sum_{i=1}^{13} (x_i - \bar{x})^2 = 20.717$$

$$\sigma(y) = \frac{1}{N-1} \sum_{i=1}^{13} (y_i - \bar{y})^2 = 25.931$$

Subject	Normalized x	Normalized y	L2 distance	L1 distance
Mathematics	-0.616	1.415	1.899	2.585
	-0.954	0.959	1.622	2.269
	-1.051	1.338	2.088	2.943
Electrical Engineering	0.783	-0.012	0.474	0.627
	1.749	0.721	1.680	2.326
	1.170	0.991	1.453	2.017
Computer Science	-0.230	-0.127	0.585	0.657
	0.011	-0.513	0.453	0.646
	-0.906	-0.590	1.316	1.64
	-1.630	-0.012	1.990	2.172
Economics	0.687	-1.709	1.538	1.842
	0.300	-1.014	0.801	0.858
	0.687	-1.246	1.094	1.379

Distance column is the distance between new point (20,7)

normalized to (0.349,-0.205) other points (Normalized)

From the above table, For L2 distance, the new point is classified as Computer Science for K=1 and 5

For L1 distance, it is classified as Electric Engineering for K=1 and Computer Science fro K=5

(b)

(i)

$$\begin{aligned}P(x) &= \sum_{i=1}^N P(x|Y=c)P(Y=c) \\&= \frac{K_1}{N_1 V} \cdot \frac{N_1}{N} + \frac{K_2}{N_2 V} \cdot \frac{N_2}{N} + \frac{K_3}{N_3 V} \cdot \frac{N_3}{N} + \dots + \frac{K_c}{N_c V} \cdot \frac{N_c}{N} \\&= \frac{1}{VN} [\sum_c K_c] \\&\therefore p(x) = \frac{K}{VN}\end{aligned}$$

(ii)

$$P(Y=c|x) = \frac{P(x|Y=c)P(Y=c)}{P(x)}$$

$$\frac{\frac{K_c}{N_c V} \cdot \frac{N_c}{N}}{\frac{K}{VN}}$$

$$\therefore P(Y=c|x) = \frac{K_c}{K}$$

4 DECISION TREE

(i)

$$H(x) = - \sum_{i=1}^n p_i \log(p_i)$$

$$H(y|x) = - \sum_{i=1}^n p_{X=X_i} H(Y|x = x_i)$$

$$= \sum_{m,n} p_{X_i, y_i} \log(p(y_i|x_i))$$

$$I(y; x) = H(y) - H(y|x)$$

$$H(AccRate) = -[\frac{73}{100} \log(\frac{73}{100}) + \frac{27}{100} \log(\frac{27}{100})] = 0.841465$$

$$I(Acc; Weather) = H(Acc) - H(Acc|Weather)$$

$$I(Acc; Traffic) = H(Acc) - H(Acc|Traffic)$$

$$H(Acc|Weather) = -[\frac{23}{28} \log(\frac{23}{28}) + \frac{5}{28} \log(\frac{5}{28})] \frac{28}{100} + \frac{72}{100} [\frac{50}{72} \log(\frac{50}{72}) + \frac{22}{72} \log(\frac{22}{72})]$$

$$= 0.82889$$

$$I(Acc; weather) = 0.0126$$

$$H(Acc; Traffic) = -[\frac{27}{100} (\frac{27}{27} \log(\frac{27}{27}))] + [\frac{73}{100} (\frac{73}{73} \log(\frac{73}{73}))]$$

$$= 0$$

$$\therefore I(Acc; Traffic) = 0.841$$

HencesplitonTrafficastheinformationgainis max

(ii)

The decision tree with the normalized data will give the same out come as the one with raw data. This is because normalization does not affects the number of objects class/category it just scales the data to follow a normal curve. This does not affect the data and hence the splitting of normalized tree is same as that of the normal tree.

(iii)

$$GiniIndex = \sum_{k=1}^k p_k(1 - p_k)$$

$$CrossEntropy = - \sum_{k=1}^k p_k \log(p_k)$$

$$\text{To prove, } \sum_{k=1}^k p_k(1 - p_k) \leq - \sum_{k=1}^k p_k \log(p_k)$$

i.e $p_k(1 - p_k) \leq -p_k \log(p_k)$ for each term in the summation

value of p_k in the range $[0,1]$

consider a function $p(x)=1-x\log(x)$ defined in the range $[0,1]$

to get maximum, $p'(x)=0$

$$\therefore, x - \frac{1}{x} = 0$$

$$\therefore, x = 1 \text{ with max value}=0$$

thus, $p(x) \leq 0$

$$1-x+\log(x) \leq 0$$

$$\therefore, 1-x \leq -\log(x)$$

i.e $x(1-x) \leq -x\log(x)$...which is what we want to prove

5 PROGRAMMING

- i. There are total 10 attributes including row id which is the first column.(Class attribute is not being considered as attribute)
- ii. No not all attributes are meaning full as the first record contains the row id which does not have any impact on classification as its just a number
- iii. There are total of 7 classes but class 4 has no records
- iv. The majority is for class 2 with 73 records. Its not a uniform distribution.

Table 1: Accuracy

Knn For K	Training L1	Training L2	Testing L1	Testing L2
1	75.0000	71.4286	66.6667	61.1111
3	74.4898	72.4490	72.2222	66.6667
5	70.9184	72.4490	61.1111	55.5556
7	69.8980	68.8776	55.5556	55.5556

* For KNN in case of a tie between classes, the point with minimum distance is considered for classification

Bayes	Training	Testing L2
1	54.5918	33.3333

As clearly seen by the outputs, KNN gives better accuracy than Naive Bayes as KNN is being non parametric makes no assumption about the data where as Naive Bayes assumes attributes are conditionally independent to each other given class values and are normally distributed. As a result Naive Bayes can have only linear boundary compared to flexible boundary KNN can has.

This assignment has been collaborated between me, Bhavya Gandhi, Dhaval Shah, Junaid Hundekar and Sagar Makwana