

Name:- Prateek Mahajan

EE511 : Homework 3

Problem 1:-

- a) I have read and understood the general instructions at the top of HW3 and I formally declare that all work I turn in for everything in this course will not contain or involve any cheating at all.

Problem 2:-

- a) Consider a binary naive Bayes classifier with multivariate class conditional distribution $p(x|y)$

Let class variances be equal.

Bayes assumption : - $p(x|y) = \prod_{j=1}^m p(x_j|y)$

Since they are gaussian :-

$$p(x|y) = \prod_{i=1}^m \frac{1}{(2\pi\sigma_{y,i}^2)^{1/2}} \exp\left(-\frac{1}{2} \frac{(x_i - \mu_{y,i})^2}{\sigma_{y,i}^2}\right)$$

for "m" univariate gaussians

As mentioned in the slides:-

~~Log odds ratio for the binary case:-~~

Posterior probability, $p(y=y_1|x) =$

$$\frac{1}{1 + \exp(-\log(\frac{p(x|y=y_1)}{p(x|y=y_0)} \cdot \frac{p(y=y_1)}{p(y=y_0)}))}$$

$$\Rightarrow p(y=y_1|x) = g\left(\sum_{i=1}^m \log\left(\frac{p(x_i|y=y_1)}{p(x_i|y=y_0)}\right) + \log\left(\frac{p(y=y_1)}{p(y=y_0)}\right)\right),$$

where the g inside is the log ratio

Consider the term $\frac{p(x_i | y = y_1)}{p(x_i | y = y_0)}$

$$\textcircled{1} \quad p(x_i | y = y_1) = \frac{1}{(2\pi\sigma_{y_1,i}^2)^{1/2}} \cdot \exp\left(-\frac{1}{2} \frac{(x_i - \mu_{y_1,i})^2}{\sigma_{y_1,i}^2}\right)$$

$$\textcircled{2} \quad p(x_i | y = y_0) = \frac{1}{(2\pi\sigma_{y_0,i}^2)^{1/2}} \cdot \exp\left(-\frac{1}{2} \frac{(x_i - \mu_{y_0,i})^2}{\sigma_{y_0,i}^2}\right)$$

But class variances are equal, $\sigma_{y_0,i}^2 = \sigma_{y_1,i}^2 = \sigma_i^2$

$$\begin{aligned} \Rightarrow \frac{p(x_i | y = y_1)}{p(x_i | y = y_0)} &= \exp\left(-\frac{1}{2\sigma_i^2} ((x_i - \mu_{y_1,i})^2 - (x_i - \mu_{y_0,i})^2)\right) \\ &= \exp\left(-\frac{1}{2\sigma_i^2} (2x_i - \mu_{y_1,i} - \mu_{y_0,i})(\mu_{y_0,i} - \mu_{y_1,i})\right) \\ &= \exp\left(-\frac{1}{\sigma_i^2} \frac{(x_i - (\mu_{y_1,i} + \mu_{y_0,i}))(\mu_{y_0,i} - \mu_{y_1,i})}{2}\right) \\ \log\left(\frac{p(x_i | y = y_1)}{p(x_i | y = y_0)}\right) &= \cancel{\mu_{y_1,i} - \mu_{y_0,i}} \frac{x_i}{\sigma_i^2} + \frac{\mu_{y_1,i}^2 - \mu_{y_0,i}^2}{2\sigma_i^2} \end{aligned}$$

Now consider $\log\left(\frac{p(y=1)}{p(y=0)}\right)$

Clearly, the above term is a constant & independent of X

\therefore the log ratio

$$\begin{aligned} &\sum_{i=1}^m \log\left(\frac{p(x_i | y=1)}{p(x_i | y=0)}\right) + \log\left(\frac{p(y=1)}{p(y=0)}\right) \\ &= \sum_{i=1}^m \cancel{\frac{(\mu_{y_1,i} - \mu_{y_0,i})}{\sigma_i^2} x_i} + \frac{\mu_{y_0,i}^2 - \mu_{y_1,i}^2}{2\sigma_i^2} + \log\left(\frac{p(y=1)}{p(y=0)}\right) \end{aligned}$$

$= b_i x_i + c_i$ where :-

$$b_i = \cancel{\frac{(\mu_{y_1,i} - \mu_{y_0,i})}{\sigma_i^2}}$$

$$c_i = \frac{\mu_{y_0,i}^2 - \mu_{y_1,i}^2}{2\sigma_i^2} + \log\left(\frac{p(y=1)}{p(y=0)}\right)$$

b) If the class variances are unequal, consider eq ① & ② from 2a).

From ① & ②:-

$$\frac{p(x_i | y=y_1)}{p(x_i | y=y_0)} = \frac{\sigma_{y_1, i}}{\sigma_{y_0, i}} \cdot \exp\left(\frac{-1}{2} \left[\frac{1}{\sigma_{y_1, i}^2} - \frac{1}{\sigma_{y_0, i}^2} \right]\right)$$

$$\exp\left[\sigma_{y_0, i}^2 \left(n_i^2 + \mu_{y_1, i}^2 - 2\mu_{y_1, i} n_i \right) - \sigma_{y_1, i}^2 \left(n_i^2 + \mu_{y_0, i}^2 - 2\mu_{y_0, i} n_i \right)\right]$$

$$\Rightarrow \frac{p(x_i | y=y_1)}{p(x_i | y=y_0)} = \frac{\sigma_{y_1, i}}{\sigma_{y_0, i}} \cdot \exp\left[\frac{-1}{2} \cdot \frac{1}{\sigma_{y_1, i}^2} \cdot \left[n_i^2 (\sigma_{y_0, i}^2 - \sigma_{y_1, i}^2) + n_i^2 [2\mu_{y_0, i} \sigma_{y_1, i}^2 - 2\mu_{y_1, i} \sigma_{y_0, i}^2] + \sigma_{y_0, i}^2 \mu_{y_1, i}^2 - \sigma_{y_1, i}^2 \mu_{y_0, i}^2 \right]\right]$$

Taking log on both sides:-

$$\log\left(\frac{p(x_i | y=y_1)}{p(x_i | y=y_0)}\right) = \log\left(\frac{\sigma_{y_1, i}}{\sigma_{y_0, i}}\right) + \frac{1}{2} \cdot \frac{n_i^2 (\sigma_{y_0, i}^2 - \sigma_{y_1, i}^2)}{\sigma_{y_1, i}^2 \sigma_{y_0, i}^2} + 2n_i (\mu_{y_1, i} \sigma_{y_0, i}^2 - \mu_{y_0, i} \sigma_{y_1, i}^2) + \sigma_{y_1, i}^2 \mu_{y_0, i}^2 - \sigma_{y_0, i}^2 \mu_{y_1, i}^2$$

Now consider $\log\left(\frac{p(y=1)}{p(y=0)}\right)$

Clearly, as in 2a, this is constant & independent of x .

\therefore the log ratio =

$$\sum_{i=1}^n \log\left(\frac{p(x_i | y=1)}{p(x_i | y=0)}\right) + \log\left(\frac{p(y=1)}{p(y=0)}\right)$$

$$\begin{aligned}
 &= \pi_i^2 \left(\frac{\sigma_{y_{1,i}}^2 - \sigma_{y_{0,i}}^2}{2 \sigma_{y_{1,i}}^2 \sigma_{y_{0,i}}^2} \right) + \lambda \pi_i (\mu_{y_{1,i}} \sigma_{y_{1,i}}^2 - \mu_{y_{0,i}} \sigma_{y_{0,i}}^2) \\
 &\quad + \frac{\sigma_{y_{1,i}}^2 \cdot \mu_{y_{0,i}}^2}{2 \sigma_{y_{1,i}}^2 \cdot \sigma_{y_{0,i}}^2} - \frac{\sigma_{y_{0,i}}^2 \mu_{y_{1,i}}^2}{2 \sigma_{y_{1,i}}^2 \cdot \sigma_{y_{0,i}}^2} + \log \left(\frac{\sigma_{y_{0,i}}}{\sigma_{y_{1,i}}} \right) \\
 &\quad + \log \left(\frac{p(y=1)}{p(y=0)} \right) \\
 &= a_i \pi_i^2 + b_i \pi_i + c_i
 \end{aligned}$$

where :-

$$a_i = \frac{\sigma_{y_{1,i}}^2 - \sigma_{y_{0,i}}^2}{2 \sigma_{y_{1,i}}^2 \cdot \sigma_{y_{0,i}}^2}$$

$$b_i = \frac{\mu_{y_{1,i}} \sigma_{y_{0,i}}^2 - \mu_{y_{0,i}} \sigma_{y_{1,i}}^2}{\sigma_{y_{1,i}}^2 \cdot \sigma_{y_{0,i}}^2}$$

$$c_i = \frac{\sigma_{y_{1,i}}^2 \mu_{y_{0,i}}^2 - \sigma_{y_{0,i}}^2 \mu_{y_{1,i}}^2}{2 \sigma_{y_{1,i}}^2 \cdot \sigma_{y_{0,i}}^2} + \log \left(\frac{\sigma_{y_{0,i}}}{\sigma_{y_{1,i}}} \right) + \log \left(\frac{p(y=1)}{p(y=0)} \right)$$

c) As mentioned in 2a :-

$$p(y=y_1 | x) = \frac{1}{1 + \exp(-\log(p(x|y=y_1) \cdot p(y=y_1)) - \log(p(x|y=y_0) \cdot p(y=y_0)))}$$

$$= g \left(\sum_{i=1}^m \log \left(\frac{p(x_i | y=y_1)}{p(x_i | y=y_0)} \right) + \log \left(\frac{p(y=y_1)}{p(y=y_0)} \right) \right),$$

where $g = \frac{1}{1 + e^{-z}}$ \rightarrow logistic function

& $z \rightarrow \log \text{ ratio of gaussian class conditional distributions}$

this classifier is essentially a logistic regression in the log ratio of the ~~more~~ name Bayes classifier.

For equal class variance (and assuming g is a logistic f^n): -

$$P(y=1|x) = g(\theta^T x + \text{const}) \rightarrow \text{linear logistic regression}$$

For unequal class variances: -

$$P(y=1|x) = g\left(\sum_{i=1}^m a_i x_i^2 + b_i x_i + c_i\right)$$

→ quadratic logistic regression,
which can be made linear by consider the x_i^2 terms as new features

Problem 3:-

a) For $m=2$: -

$$V_2(r) = \pi r^2$$

$$S_1(r) = 2\pi r$$

For $m=3$: -

$$V_3(r) = \frac{4}{3} \pi r^3$$

$$S_2(r) = 4\pi r^2$$

b) In my opinion, the core intuition to remember here is that surface area represents the "boundary" or "skin" of a given volume. "dr" represents a small change in the radius of a volume.

∴ if you change the ~~outer surface~~ of a volume radius of a volume by "dr" the proportional change in volume, " dV " will be the surface ~~area~~ of the volume (especially since $dr \rightarrow 0$, which means ~~the~~ surface area can be considered to be constant).

∴ change in volume (infinitely small) = surface area change in radius (infinitely small)

$$\Rightarrow \frac{dV_m(r)}{dr} = m S_{m-1}(r) \quad (\text{if } dr \rightarrow 0)$$

Now take $m=2$: -

$$V_2(r) = \pi r^2$$

$$\Rightarrow \frac{dV_2(r)}{dr} = 2\pi r = S_1(r) \quad \text{and hence ref}$$

$$\text{For } m=3: - \quad (\text{area} + \text{ref}) p = (r(1-p))$$

$$V_3 = \frac{4}{3} \pi r^3$$

$$\Rightarrow \frac{dV_3(r)}{dr} = 4\pi r^2 = S_2(r) \quad p = (r(1-p))$$

Clearly, this eqn holds for $m=2 \& 3$!

c) $V_m(r) = K r^m$ (as mentioned in question),
where K is a constant

$$\therefore S_{m-1}(r) = m \cdot K \cdot r^{m-1} = dV(r)/dr$$

$$\text{If } r=1: -$$

$$S_{m-1}(r) = m \cdot K = \bar{S}_{m-1}$$

$$\Rightarrow K = \frac{\bar{S}_{m-1}}{m}$$

$$S_{m-1}(r) = \frac{\bar{S}_{m-1} \cdot r^m}{m}$$

$$\text{wherever } V_m(r) = \frac{\bar{S}_{m-1} \cdot r^m}{m}, \text{ remains true if } (d)$$

$$\text{the surface area term is now "defined"}$$

$$\text{when } S_{m-1}(r) = \frac{\bar{S}_{m-1} \cdot r^{m-1}}{m}$$

d) $f_m(r) = \int_{r \in S_{m-1}(r)} p(x) \cdot d\mu \text{ where } \mu \text{ is surface measure}$

$$p(x) = \frac{1}{(2\pi\sigma^2)^{m/2}} \exp\left(-\frac{\|x\|_2^2}{2\sigma^2}\right)$$

$$\|x\|_2 = r + \text{min } S_{m-1}(r) \text{ or in words}$$

(distance) when on sphere

$\therefore f_m(r) = \text{probability density of sampled points lying on surface of } S_{m-1}(r)$

$$= \int_{x \in S_{m-1}(r)} p(x) \cdot d\mu$$

$$\text{But } p(x) = \frac{1}{(2\pi\sigma^2)^{m/2}} \exp\left(-\frac{\|x\|_2^2}{2\sigma^2}\right)$$

$$\& \text{since } \|x\|_2 = r \therefore$$

$$p(x) = \frac{1}{(2\pi\sigma^2)^{m/2}} \exp\left(-\frac{r^2}{2\sigma^2}\right)$$

$$\therefore f_m(r) = \frac{1}{(2\pi\sigma^2)^{m/2}} \exp\left(-\frac{r^2}{2\sigma^2}\right) \cdot \underbrace{\int_{x \in S_{m-1}(r)} dx}_{\text{int } S_{m-1}(r)}$$

note that
this is because
the gaussian
is isotropic

$$= \frac{1}{(2\pi\sigma^2)^{m/2}} \exp\left(-\frac{r^2}{2\sigma^2}\right) \cdot \overbrace{S_{m-1}(r)}^{\text{this is essentially}}$$

$$= \frac{1}{(2\pi\sigma^2)^{m/2}} \exp\left(-\frac{r^2}{2\sigma^2}\right) \cdot \overbrace{S_{m-1} \cdot r^{m-1}}^{\text{int } S_{m-1}(r)}$$

e) Using the above $f_m(r)$,

$$\text{set } \frac{d}{dr} f_m(r) = 0$$

$$\frac{d}{dr} f_m(r) = \frac{d}{dr} \frac{S_{m-1} \cdot r^{m-1}}{(2\pi\sigma^2)^{m/2}} \exp\left(-\frac{r^2}{2\sigma^2}\right)$$

$$= \frac{S_{m-1}}{(2\pi\sigma^2)^{m/2}} \cdot \left[(m-1)r^{m-2} \cdot \exp\left(-\frac{r^2}{2\sigma^2}\right) - \frac{2r^m}{2\sigma^2} \exp\left(-\frac{r^2}{2\sigma^2}\right) \right]$$

$$= \frac{S_{m-1}}{(2\pi\sigma^2)^{m/2}} \cdot r^{m-2} \exp\left(-\frac{r^2}{2\sigma^2}\right) \cdot \left[(m-1) - \frac{r^2}{\sigma^2} \right]$$

Setting $\frac{d}{dr} f_m(r) = 0$

$$\Rightarrow (m-1) - \frac{r^2}{\sigma^2} = 0 \Rightarrow r^2 = (m-1)\sigma^2 \Rightarrow r_+ = \pm \sqrt{(m-1)\sigma^2}$$

But $m > 2$ $\Rightarrow m-1 \approx m$ (i.e. m is large)

$$\Rightarrow r \approx \pm \sqrt{m}\sigma$$

Consider $\frac{d^2}{dr^2} f_m(r)$.

$$= \frac{\overline{S_{m-1}}}{(2\pi\sigma^2)^{m/2}} \cdot \exp\left(-\frac{r^2}{2\sigma^2}\right) \left[(m-2)r^{m-3} \cdot \left(\frac{(m-1)-r^2}{\sigma^2}\right) - \frac{3r^{m-1}}{\sigma^2} \right]$$

$$= \frac{\overline{S_{m-1}}}{(2\pi\sigma^2)^{m/2}} \cdot \exp\left(-\frac{r^2}{2\sigma^2}\right) r^{m-3} \left[(m-2)(m-1) - \frac{(m+1)}{\sigma^2} r^2 \right]$$

For $r = \sqrt{m}\sigma$, clearly $\frac{d^2}{dr^2} f_m(r) < 0$

~~(as $m(m+1)\sigma^2 = m(m+1)$)~~ (as $(m+1)m\sigma^2$)

which is greater than $(m-2)(m-1)$

$\therefore r = \sqrt{m}\sigma$ is a maximum value.

Clearly $\frac{d}{dr} f_m(r) = 0$ only at $r = \sqrt{m}\sigma$

$\Rightarrow f_m(r)$ has a single maximum value at $\hat{r} \approx \sqrt{m}\sigma$

f) As proven above, $\hat{r} = \sqrt{m}\sigma$

$$f(r) = \frac{\overline{S_{m-1}}}{(2\pi\sigma^2)^{m/2}} r^{m-1} \cdot \exp\left(-\frac{r^2}{2\sigma^2}\right)$$

$$\therefore f(\hat{r} + \epsilon) = \frac{\overline{S_{m-1}}}{(2\pi\sigma^2)^{m/2}} (\hat{r} + \epsilon)^{m-1} \cdot \exp\left(-\frac{(\hat{r} + \epsilon)^2}{2\sigma^2}\right)$$

Now consider $(\hat{r} + \epsilon)^{m-1}$

Using multinomial expansion:

$$(\hat{r} + \epsilon)^{m-1} = {}^{m-1}C_0 \cdot (\hat{r})^{m-1} \cdot \epsilon^0 + {}^{m-1}C_1 \cdot (\hat{r})^{m-2} \cdot \epsilon + \dots$$

$$+ {}^{m-1}C_{m-2} \cdot \hat{r}^{m-2} \cdot \epsilon^{m-2} + {}^{m-1}C_{m-1} \cdot (\hat{r})^0 \cdot \epsilon^{m-1}$$

Clearly, since $\epsilon \ll \sigma \ll \hat{r}$, all terms are much smaller than $m^{-1} C_0 \cdot (\hat{r})^{m-1} = (\hat{r}^1)^{m-1}$

$$\therefore (\hat{r} + \epsilon)^{m-1} \approx (\hat{r}^1)^{m-1}$$

$$\begin{aligned}\therefore g(\hat{r} + \epsilon) &\approx \frac{S_{m-1}}{(2\pi\sigma^2)^{m/2}} \cdot (\hat{r}^1)^{m-1} \cdot \exp\left(\frac{-(\hat{r}^1)^2 - \epsilon^2 - 2\hat{r}\epsilon}{2\sigma^2}\right) \\ &\approx \frac{S_{m-1}}{(2\pi\sigma^2)^{m/2}} \cdot (\hat{r}^1)^{m-1} \cdot \exp\left(\frac{-(\hat{r}^1)^2}{2\sigma^2}\right) \cdot \exp\left(\frac{-\epsilon^2 - 2\hat{r}\epsilon}{2\sigma^2}\right) \\ &\quad \text{this is } g(\hat{r}^1) \\ &\approx g(\hat{r}^1) \cdot \exp\left(\frac{-\epsilon^2}{2\sigma^2}\right) \cdot \exp\left(-\frac{\hat{r}\epsilon}{\sigma^2}\right) \\ &\approx g(\hat{r}^1) \cdot \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right) \cdot \exp\left(-\frac{\sqrt{m}\epsilon}{\sigma}\right)\end{aligned}$$

Using Taylor series expansion: —

$$\exp\left(-\frac{\sqrt{m}\epsilon}{\sigma}\right) = 1 - \frac{\sqrt{m}\epsilon}{\sigma} + \frac{m\epsilon^2}{2\sigma^2} - \frac{m^{3/2}\epsilon^3}{3\sigma^3} + \dots$$

But $\epsilon \ll \hat{r} \Rightarrow \epsilon \ll \sqrt{m}$ (as σ is constant)
 $\Rightarrow \sqrt{m}\epsilon \approx 0$.

$$\therefore \exp\left(-\frac{\sqrt{m}\epsilon}{\sigma}\right) \approx 1.$$

$$\Rightarrow g(\hat{r} + \epsilon) \approx g(\hat{r}^1) \cdot \exp\left(-\frac{\epsilon^2}{2\sigma^2}\right) \cdot 1$$

- q) If we take a finite ~~sample~~ number of samples from a high dimensional gaussian distribution, most of the points would reside at a distance/radius, $r = \sqrt{m}\sigma$

As shown in 3f), moving away from this point exponentially drops the mass density \propto

probability of sampling. This holds true whether we move from $r = \sqrt{m}\sigma$ towards both the origin & towards infinity.

For a low dimensional gaussian distribution, ~~the mass~~ most points reside near $r = (\sqrt{m-1})\sigma$.

\Rightarrow For $m=1$, they reside at the origin

For other low values of m as well, most points would be at small multiples of σ from the origin, and the drop would not be as dramatic as we move from the maxima

$$n) \text{ Probability density } f_m(r) = \frac{S_{m-1}}{(2\pi\sigma^2)^{m/2}} r^{m-1} \exp\left(-\frac{r^2}{2\sigma^2}\right)$$

$$\therefore f_m(0) = 0$$

$$f_m(r + S_{m-1}(r)) = f_m(r) = \frac{S_{m-1}}{(2\pi\sigma^2)^{m/2}} \cdot (S_m\sigma)^{m-1} \cdot \exp\left(-\frac{(r + S_m\sigma)^2}{2\sigma^2}\right)$$

$$= \frac{S_{m-1}}{(2\pi)^{m/2}} \cdot \frac{m^{(m-1)/2}}{\sigma} \cdot e^{-m/2}$$

clearly, probability density = 0 at the origin & is ~~very~~ very high at $r = \sqrt{m}\sigma$ (which is another consequence of dimensionality)

I wanted to show below since it's clear

(i) ~~Please refer to the pdf at the end~~

i) Please refer to the end of the pdf for the graph & code.

8. Please see the attached file for the

Clearly, the mean of the radius increases as m increases, which is in line with our earlier observations, which show that the mass of the gaussian goes to a higher radius as the dimensionality increases (as depicted by this increase in mean).

Further, it was observed that even when the mass increased, the probability density at the maximum "n" dropped exponentially as you moved away from it in higher dimensions \Rightarrow the standard deviations shouldn't increase & be \approx constant, as is deserved in this graph.

Problems:-

- a) Based on the scatter plot of the question, the decision boundary seems to be a non-linear & non-exponential one. Therefore, I believe a low-dimensional linear SVM (not a higher dimensional may work) would not be a great choice, which is also the case for logistic regression, as the decision boundary is not linear in nature.

\therefore I would choose the gaussian kernel SVM.

(since it is probably better at representing a non-linear decision boundary; like the provided one)

- b) Please refer to the end of this pdf.

c) Please refer to the end of this pdf for the accuracy answers.

The gaussian kernel SVM is better at representing non-linear decision boundaries than the others which are for linear boundaries. More so than anything else, this is a property of their mathematical formulations (which is "linear" in nature)

d) logistic regression: linear Decision Boundary, as it is a linear mathematical formulation (i.e. it relies on a linear combinations of input features)

linear SVM: linear Decision Boundary as it is a linear mathematical formulation (as it relies on a linear combinations of input features)

Gaussian Kernel SVM: Non linear decision boundary as is a non-linear mathematical formulation

e) Please refer to the end of this pdf for visualisations. The experiment confirmed my hypothesis.

a) Please refer to the end of the pdf for the code & plots. I took the solⁿ for $\lambda = 25.0$

- b) Please refer to the end of this pdf for this qn.
- c) Please refer to the end of this pdf for this qn..

In the higher σ case, the initial weights & the change in weights is less consistent & more random.

~~Also, The~~ As a consequence, the zero features take longer to converge to sparsity. This occurs due to an increase in randomness in data, ^{and} drop in consistency. ~~and~~

- d)
- For $n=50, m=75, \lambda = 25.0$
 - For $n=50, m=150, \lambda \approx 50.0$
 - For $n=50, m=1000, \lambda = 70.0$
 - For $n=100, m=75, \lambda = 25.0$
 - For $n=100, m=150, \lambda = 50.0$.

I think $n = O(m)$ is most ^{the} probable here.

- e) Please refer to the end of this pdf file for this qn.