

PROBLEM 1

1. Maximum likelihood estimate for lambda

Referred from CASE -2 of <http://www2.imm.dtu.dk/courses/02711/lecture3.pdf>

$$X \sim \frac{\lambda^x e^{-\lambda}}{x!} \quad \text{for some real valued } \underline{\lambda > 0}$$

$$\Rightarrow \lambda_{MLE} = \underset{\lambda}{\operatorname{argmax}} P(D|\lambda) \rightarrow \text{Maximizing likelihood function}$$

$$P(D|\lambda) = \prod_{i=1}^m \frac{\lambda^{x^{(m)}} e^{-\lambda}}{x^{(m)}!} \rightarrow \text{MAXIMIZED s.t. } \underline{\lambda > 0}$$

$$\max_{\lambda} \prod_{i=1}^m \frac{\lambda^{x^{(m)}} e^{-\lambda}}{x^{(m)}!} \Rightarrow L(\lambda, \mu) = \frac{\lambda^{\sum x^{(m)} - M\lambda}}{\prod x^{(m)}!} - \mu(-\lambda)$$

$$\text{s.t. } -\lambda < 0$$

Applying KKT conditions gives us:

$$\text{STATIONARITY: } \frac{\partial}{\partial \lambda} \left(\frac{\lambda^{\sum x^{(m)} - M\lambda}}{\prod x^{(m)}!} + \mu \lambda \right) = 0 \quad \text{i.e. } \frac{\partial L}{\partial \lambda} = 0.$$

$$\text{COMPLEMENTARY SLACKNESS: } \mu(\lambda) = 0 \quad \swarrow$$

Two cases arise in order to satisfy
this constraint.

$$\text{PRIMAL FEASIBILITY: } -\lambda < 0$$

$$\textcircled{1} \mu = 0 \quad \textcircled{2} \mu \neq 0$$

$$\text{DUAL FEASIBILITY: } \mu > 0$$

CASE-1 : $\mu = 0$.

↙ [Two complementary conditions]

$$L(\lambda) \rightarrow \max_{\lambda} \frac{\sum z^{(m)} - M\lambda}{\prod z^{(m)}!}$$

$$\begin{aligned} \frac{\partial L}{\partial \lambda} = 0 &\Rightarrow \frac{\partial}{\partial \lambda} \left(\frac{\sum z^{(m)} - M\lambda}{\prod z^{(m)}!} \right) = 0 \\ \Rightarrow \frac{\sum z^{(m)}}{\prod z^{(m)}!} &+ \frac{(-M)}{\prod z^{(m)}!} = 0 \end{aligned}$$

$$\Rightarrow \sum z^{(m)} - M\lambda = 0 \Rightarrow \lambda = \frac{\sum z^{(m)}}{M} \quad L > 0$$

[for this λ]

CASE-2 : $\mu > 0, \lambda = 0$.

$\lambda = 0 \rightarrow$ does not satisfy the primal feasibility condition ($\lambda \geq 0$)

of the considered KKT conditions.

∴ Optimal Solution occurs when

$$\lambda = \frac{\sum z^{(m)}}{M}$$

2) Prior probability $P(\lambda) = \frac{1}{5} \max \left\{ 1 - \frac{\lambda}{10}, 0 \right\}$ is introduced

$$P(D|\lambda) = \frac{\lambda^{\sum x^{(m)}} \cdot e^{-M\lambda}}{\prod (x^{(m)})!} \Rightarrow P(\lambda|D) = P(D|\lambda) P(\lambda)$$

[MAP ESTIMATE]

s.t. $\lambda > 0, \lambda < 10$

$$\lambda_{MAP} = \underset{\lambda}{\operatorname{argmax}} P(\lambda|D) \quad \Rightarrow \quad \underset{\lambda}{\operatorname{max}} \frac{\lambda^{\sum x^{(m)}} \cdot e^{-M\lambda}}{\prod (x^{(m)})!} \cdot \left(\frac{1}{5} \left(1 - \frac{\lambda}{10} \right) \right)$$

s.t. $-\lambda < 0, \lambda - 10 < 0$

Let $\sum_m x^{(m)} = s, \prod_m (x^{(m)})! = P$

$$L(\lambda, \mu_1, \mu_2) = \frac{\lambda^s \cdot e^{-M\lambda}}{P} \cdot \left(\frac{10-\lambda}{50} \right) - \mu_1(-\lambda) - \mu_2(\lambda - 10)$$

Applying KKT conditions.

STATIONARITY : $\frac{\partial L}{\partial \lambda} = 0$

PRIMAL FEASIBILITY : $\lambda > 0, \lambda < 10$

COMPLEMENTARY : $\mu_1(\lambda) = 0$

DUAL FEASIBILITY : $\mu_1, \mu_2 > 0$

SLACKNESS : $\mu_2(10-\lambda) = 0$

① $\mu_1 = 0, \mu_2 = 0$

Two conditions leads to 4 cases

② $\mu_1 \neq 0, \mu_2 = 0$

③ $\mu_1 = 0, \mu_2 \neq 0$

④ $\mu_1 \neq 0, \mu_2 \neq 0$

CASE-1 : $\mu_1 = 0, \mu_2 = 0$.

By stationarity condition

$$\frac{\partial L}{\partial \lambda} = \frac{\partial}{\partial \lambda} \left(\frac{1}{50P} \left(\lambda^s \cdot e^{-M\lambda} (10-\lambda) \right) \right) = 0.$$

$$\lambda^s \left(\frac{\partial}{\partial \lambda} \left(e^{-M\lambda} (10-\lambda) \right) \right) + e^{-M\lambda} (10-\lambda) s \cdot \lambda^{s-1} = 0.$$

$$= \cancel{\lambda^s} \left(\cancel{-e^{-M\lambda}} (-1) + (10-\lambda) \cancel{(-M)} \cancel{e^{-M\lambda}} \right) + \cancel{-e^{-M\lambda}} (10-\lambda) s \cdot \cancel{\lambda^{s-1}} = 0.$$

$$= \lambda (M\lambda - 10M - 1) + s (10 - \lambda) = 0.$$

$$\Rightarrow M\lambda^2 - 10M\lambda - \lambda + 10s - s\lambda = 0 \Rightarrow M\lambda^2 - \lambda(10M + 1 + s) + 10s = 0.$$

$$\lambda = \frac{(10M + s + 1) \pm \sqrt{(10M + s + 1)^2 - 40MS}}{2M}.$$

$$2M \cdot \underline{\hspace{1cm}} \quad \textcircled{1}$$

$$\underline{\hspace{1cm}} \quad \textcircled{2}$$

$$\text{Consider the term } (10M + s + 1)^2 - 40MS = (10M + s)^2 + 1 + 2(10M + s) - 40MS$$

$$= (10M - s)^2 + 1 + 2(10M - s)$$

$$= (10M - s)^2 + 1 + 2(10M - s) + 4s$$

$$= (10M - s + 1)^2 + 4s \quad \text{--- } \textcircled{1}$$

For $\textcircled{+}$ term

$$\lambda > \frac{(10M + s + 1) + (10M - s + 1)}{2M}$$

$$\lambda > \frac{20M + 2}{2M} \Rightarrow \boxed{\lambda > 10} \rightarrow \text{Primal feasibility not satisfied.}$$

So, only $\textcircled{-}$ term exists satisfying all KKT conditions

$$\text{CASE-2: } \mu_1=0, \mu_2 \neq 0 \Rightarrow \lambda = 10.$$

Primal feasibility not satisfied. ($\lambda < 10$)

Value of likelihood function shall also be $\Rightarrow \frac{10^S \cdot e^{-10m}}{P} \left(\frac{1}{5} \left(1 - \frac{10}{10} \right) \right) \rightarrow \boxed{0}$

$$\text{CASE-3: } \mu_2=0, \mu_1 \neq 0 \Rightarrow \lambda = 0.$$

Primal Feasibility not satisfied ($\lambda > 0$)

Value of function is also $\boxed{0}$

$$\text{CASE-4: } \mu_1 \neq 0, \mu_2 \neq 0 \Rightarrow \lambda = 0, \lambda = 10.$$

Primal Feasibility is not satisfied for both ($\lambda > 0, \lambda < 10$)

For both the values Likelihood function shall also be $\boxed{0}$

Thus only candidate for local maximum is

)

$$\lambda = \frac{(10m+s+1) - \sqrt{(10m+s+1)^2 - 40ms}}{2M}.$$

Satisfying all the KKT conditions

3. We generally choose conjugate priors for analysis as that leads to
- Posterior distributions having the same functional form as prior distributions
 - So that posterior distribution can act as prior for the subsequent data observed
 - Greatly simplifies Bayesian Analysis

For the choice of prior distribution $P(\lambda) = \frac{1}{5} \max \left\{ 1 - \frac{\lambda}{10}, 0 \right\}$ — (I)

Considering our data likelihood function

$$P(D|\lambda) = \frac{\lambda^{\sum_m x^{(m)}} e^{-M\lambda}}{\prod_m (x^{(m)}!)} \quad \text{— (II)}$$

(I) is not a conjugate distribution for (II)

Thus choosing (I) doesn't result in POSTERIOR having same functional as PRIOR

As data likelihood has. product of λ & exponential in λ

A good choice of prior here should also be of the same form

Considering the [GAMMA DISTRIBUTION]

$$G(\lambda | a, b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda)$$

Thus resulting POSTERIOR shall be of the same of PRIOR

PROBLEM 2

1. Log likelihood of the data observations

For positive real valued r.v X , distributed according to LOG-NORMAL

$$X \sim \frac{1}{x \sigma \sqrt{2\pi}} \exp\left(-\frac{(\ln x - \mu)^2}{2\sigma^2}\right) \text{ for real valued parameters}$$

$$P(D | \mu, \sigma^2) = \prod_{m=1}^M p(x^{(m)} | \mu, \sigma^2). \leftarrow [\text{Likelihood function}]$$

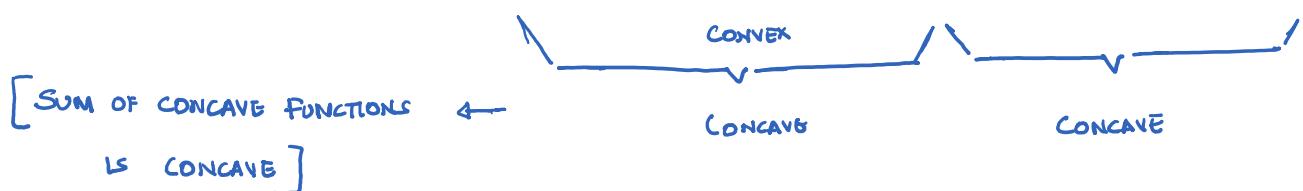
Log likelihood function $\rightarrow \ln P(D | \mu, \sigma^2)$

$$\sum_{m=1}^M \ln p(x^{(m)} | \mu, \sigma^2) = \sum_{m=1}^M \left(-\frac{(\ln x^{(m)} - \mu)^2}{2\sigma^2} - \ln(x^{(m)} \sigma \sqrt{2\pi}) \right)$$

This Log likelihood (μ, σ^2) is not concave but

$L(\mu, \lambda)$ where $[\lambda = \frac{1}{\sigma}]$ is concave $[\sigma \neq 0]$.

$$\Rightarrow L(\mu, \lambda) = \sum_{m=1}^M -\left(\frac{\lambda^2}{2} (\ln x^{(m)} - \mu)^2 \right) - \ln\left(\frac{x^{(m)} \sqrt{2\pi}}{\lambda}\right)$$



So, $L(\mu, \lambda)$ is concave

$$L(\mu, \lambda) = \sum_{m=1}^M - \left(\frac{\lambda}{2} (\ln x^{(m)} - \mu)^2 \right) - \ln \left(\frac{x^{(m)} \sqrt{2\pi}}{\lambda} \right)$$

$$\Rightarrow \frac{\partial L}{\partial \mu} = \sum_{m=1}^M - \left(\cancel{\lambda} (\ln x^{(m)} - \mu) \cancel{(}) \right) = 0.$$

$$\Rightarrow \sum_{m=1}^M \ln x^{(m)} - M\mu = 0 \Rightarrow \mu = \frac{\sum_{m=1}^M \ln x^{(m)}}{M}.$$

$$\frac{\partial L}{\partial \lambda} = \sum_{m=1}^M - \left(\lambda (\ln x^{(m)} - \mu)^2 \right) + \frac{1}{\left(\frac{x^{(m)} \sqrt{2\pi}}{\lambda} \right)} \left(+ \frac{x^{(m)} \sqrt{2\pi}}{\lambda} \right)$$

$$= \sum_{m=1}^M \left(-\lambda (\ln x^{(m)} - \mu)^2 + \frac{1}{\lambda} \right) = 0$$

$$\Rightarrow \frac{M}{\lambda} = \lambda \sum_{m=1}^M (\ln x^{(m)} - \mu)^2$$

$$\Rightarrow \frac{1}{\lambda^2} = \frac{\sum_{m=1}^M (\ln x^{(m)} - \mu)^2}{M} \Rightarrow \boxed{5}$$

3) Are the maximum likelihood estimators for μ , σ^2 biased?

Obtained sample mean and uncorrected sample variance of μ , σ^2 from maximum likelihood estimation are as follows.

$$\hat{\mu}_{MLE} = \frac{\sum_{m=1}^M \ln z^{(m)}}{M}, \quad \hat{\sigma}_{MLE}^2 = \frac{\sum_{m=1}^M (\ln z^{(m)} - \hat{\mu}_{MLE})^2}{M}$$

Let μ , σ^2 be true mean and variances of distribution.

$$\begin{aligned} \Rightarrow E[\hat{\mu}_{MLE}] &= E\left[\frac{1}{M} \sum_{m=1}^M \ln z^{(m)}\right] \quad (\text{By linearity of expectation}) \\ &= \frac{1}{M} \sum_{m=1}^M E[\ln z^{(m)}] = E[\ln \bar{x}] \\ &\qquad\qquad\qquad \text{Since each } z^{(m)} \text{ is IID } \bar{x} \end{aligned}$$

$$\begin{aligned} \text{By LOTUS} \rightarrow E[\ln \bar{x}] &= \int_0^\infty \ln x p(x) dx \leq \int_0^\infty x p(x) dx. \\ E[\hat{\mu}_{MLE}] &\leq \mu \quad \leftarrow \text{E}[x] \end{aligned}$$

[BIASED ESTIMATOR OF MEAN]

A simple intuition by looking at the estimates they are like

estimates for gaussian with r.v's $\underline{\ln z_1}, \underline{\ln z_2}, \underline{\ln z_3} \dots \underline{\ln z_n}$

In general $\underline{\sigma^2}$ for gaussian is biased. Here variance can be BIASED

Now checking for $\underline{\sigma_{MLE}^2}$

$$\sigma^2 = E \left[\frac{1}{n} \sum_i (x_i - \mu)^2 \right]$$

↑
[TRUE VARIANCE] ↑
[TRUE MEAN]

$$E \left[\sigma_{MLE}^2 \right] = E \left[\frac{1}{M} \sum_{i=1}^M (\ln x^{(m)} - \mu_{MLE})^2 \right]$$

$$= E \left[\frac{1}{M} \sum_{i=1}^M (\ln x^{(m)} - x^{(m)} + x^{(m)} - \mu + \mu - \mu_{MLE})^2 \right] \quad \text{By linearity}$$

$$= E \left[\frac{1}{M} \sum \left[(\ln x^{(m)} - x^{(m)})^2 + (x^{(m)} - \mu)^2 + (\mu - \mu_{MLE})^2 \right. \right. \\ \left. \left. + 2(x^{(m)} - \mu)(\ln x^{(m)} - x^{(m)}) + 2(x^{(m)} - \mu)(\mu - \mu_{MLE}) + 2(\ln x^{(m)} - x^{(m)})(\mu - \mu_{MLE}) \right] \right]$$

$$= \sigma^2 + E \left[\frac{1}{M} \sum (\ln x^{(m)} - x^{(m)})^2 \right] + E \left[(\mu - \mu_{MLE})^2 \right] \\ + 2 E \left[\frac{1}{M} \sum \underbrace{(x^{(m)} - \mu)(\ln x^{(m)} - x^{(m)})}_{\textcircled{1}} \right] + 2 E \left[\frac{1}{M} \sum \underbrace{(x^{(m)} - \mu)(\mu - \mu_{MLE})}_{\textcircled{2}} \right] \\ + 2 E \left[\frac{1}{M} \sum \underbrace{(\ln x^{(m)} - x^{(m)})(\mu - \mu_{MLE})}_{\textcircled{3}} \right]$$

Observing terms $\textcircled{1}, \textcircled{2}, \textcircled{3}$

$$(x^{(m)} - \mu)(\ln x^{(m)} - x^{(m)}) = x^{(m)} (\ln x^{(m)} - x^{(m)}) - \mu (\ln x^{(m)} - x^{(m)})$$

$$(\ln x^{(m)} - x^{(m)})(\mu - \mu_{MLE}) = \mu (\ln x^{(m)} - x^{(m)}) - \mu_{MLE} (\ln x^{(m)} - x^{(m)})$$

$$(x^{(m)} - \mu)(\mu - \mu_{MLE}) = x^{(m)} (\mu - \mu_{MLE}) - \mu (\mu - \mu_{MLE})$$

$$\Rightarrow 2E \left[\frac{1}{M} \sum x^{(m)} \ln x^{(m)} - (\bar{x}^{(m)})^2 - \mu_{MLE} \ln \bar{x}^{(m)} + \mu \bar{x}^{(m)} - \mu^2 + \mu \mu_{MLE} \right] + \textcircled{I}$$

$$\Rightarrow 2E \left[\frac{1}{M} \sum x^{(m)} \ln x^{(m)} - (\bar{x}^{(m)})^2 - \mu_{MLE} \ln \bar{x}^{(m)} + \mu \mu_{MLE} \right] + \textcircled{I}$$

$$\begin{aligned} &\Rightarrow 2E \left[\frac{1}{M} \sum x^{(m)} \ln x^{(m)} \right] - 2E \left[\frac{1}{M} \sum (\bar{x}^{(m)})^2 \right] - 2E \left[\frac{1}{M} \sum \mu_{MLE} \ln \bar{x}^{(m)} \right] \\ &+ E \left[(\mu - \mu_{MLE})^2 \right] + 2E \left[\frac{1}{M} \sum \mu \mu_{MLE} \right] + \sigma^2 + E \left[\frac{1}{M} \sum (\ln x^{(m)} - \bar{x}^{(m)})^2 \right] \end{aligned}$$

$$\rightarrow E[\hat{\sigma}_{MLE}^2] < \sigma^2 \quad [\text{TRUE-VARIANCE}]$$

Thus, $\hat{\sigma}^2$ is BIASED

4>

Log normal distribution insists rv \textcircled{X}
to be positive, real valued

$E[x]$ shall also be POSITIVE (MEAN)

Now, choosing a GUASSIAN PRIOR over $\boxed{\mu}$ results in
accounting for $\boxed{\sigma^2}$ values of $\boxed{\mu}$ which the data distribution
is unlikely to produce

Thus Gaussian Prior for μ may not be a good option

PROBLEM 3

- Fitting logistic regression model to training data no regularization involved

ACCURACIES :

TRAINING ACCURACY : 100%

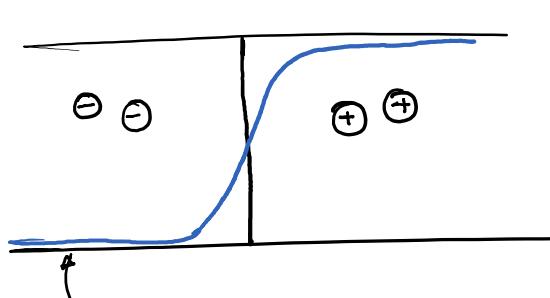
VALIDATION ACCURACY : 82.69230769230769%

TEST ACCURACY : 75.0%

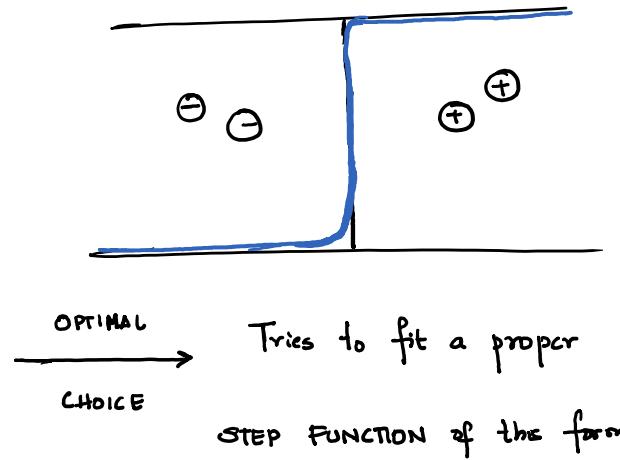
If there is a feature that would perfectly separate the two classes, the logistic regression model can no longer be trained. This is because the weight for that feature would not converge, because the optimal weight would be infinite

This problem of complete separation can be solved by introducing penalization of the weights or defining a prior probability distribution of weights.

Consider this linearly separable dataset



As logistic regression tries to fit
a curve like this



OPTIMAL CHOICE → Tries to fit a proper STEP FUNCTION of the form

To resolve this

→ Add REGULARIZATION

→ (or) Add PRIORS on WEIGHTS.



Thus throwing weights to

$$\boxed{\infty}$$

Although accuracy is converged

2. Fitting logistic regression with L2 penalty on the weights following are the accuracies for different values of C

C	TRAINING	VALIDATION	TEST
1.00E-06	0.509615	0.615385	0.5
1.00E-05	0.509615	0.615385	0.5
0.0001	0.509615	0.615385	0.5
0.001	0.509615	0.615385	0.5
0.01	0.634615	0.711538	0.596154
0.1	0.740385	0.75	0.634615
0.5	0.826923	0.826923	0.75
1	0.826923	0.826923	0.769231
10	0.884615	0.807692	0.807692
100	0.942308	0.769231	0.826923
1000	0.990385	0.807692	0.788462
10000	1	0.807692	0.75
100000	1	0.846154	0.75

Best validation accuracy is obtained for C = 100000

Best weights and bias calculated for this C are

Accuracy on the test set : 75%

Bias = array([30.29822813])

```
W = array([-109.3506758, -72.03251953, 76.77956541, -60.14462796,
-28.63187635, 44.91359683, -12.00867791, 131.98871562,
-123.65074068, 10.24641001, 5.80717194, -11.5652987,
-28.15630323, 25.15885572, -22.72300216, 18.97688102,
61.08025052, -42.25998256, -7.59122594, -30.79315092,
42.17821559, -38.71056678, 17.50353478, -36.67590325,
26.031912, -37.90732617, 46.35116147, -23.66595003,
14.15241658, -54.45700234, 77.7426917, -30.22867364,
-15.10334468, 60.51284367, -62.07836708, 74.12762222,
-26.27115732, 1.99257962, -1.49141903, 27.48790619,
-56.44536433, 43.33197411, -54.61678608, 4.3863678,
-21.31217387, 36.24032036, -12.75947559, -88.54028366,
-69.0936089, 73.24709733, -4.03034888, -31.76883203,
-28.13468045, -20.46467365, 7.90493342, -1.50289267,
13.38501368, -23.24556838, -55.26823679, -33.43392126]))
```

3. Fitting logistic regression with L1 penalty on weights, following are the accuracies

C	TRAINING	VALIDATION	TEST
1.00E-06	0.509615	0.615385	0.5
1.00E-05	0.509615	0.615385	0.5
0.0001	0.509615	0.615385	0.5
0.001	0.490385	0.384615	0.5
0.01	0.509615	0.615385	0.5
0.1	0.509615	0.615385	0.5
0.5	0.75	0.692308	0.711538
1	0.817308	0.788462	0.807692
10	0.923077	0.769231	0.788462
100	1	0.788462	0.788462
1000	1	0.769231	0.75
10000	1	0.807692	0.75
100000	1	0.826923	0.75

Best validation accuracy is obtained for C = 100000

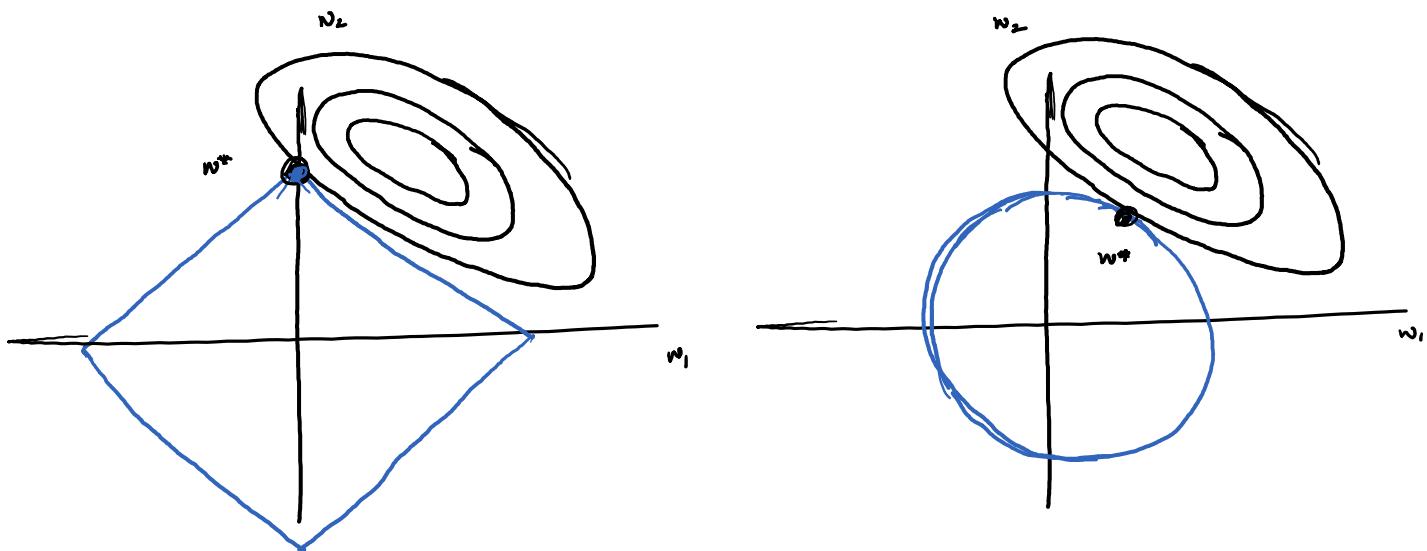
Corresponding accuracy on the test set = 75%

Bias = array([43.30640427])

Weights = array([-160.1945443, -105.80725621, 114.78170216, -86.88794287, -42.02491905, 66.32214868, -20.88374954, 192.12512919, -179.13769098, 16.61122584, 4.07186831, -14.68757473, -39.91093397, 38.52451869, -32.37626728, 24.87398742, 87.88732102, -59.82707773, -11.21993887, -45.73345418, 63.28176764, -56.05864627, 25.42596941, -52.31893448, 37.18580769, -57.09291681, 68.92534319, -34.28130793, 19.0614109, -78.46307405, 112.76784387, -42.26988118, -22.77675412, 89.07934674, -91.25545354, 106.80698489, -36.60797823, 0.95900351, -1.11895256, 37.68713966, -82.83541057, 60.53515299, -77.88782372, 6.99943983, -29.43999232, 53.28326561, -16.30656688, -123.65912787, -96.67104192, 107.267997, -3.76916167, -45.78622287, -40.56999317, -29.12157315, 8.95330377, -0.32430307, 18.73394906, -31.1382513, -80.83139619, -48.64800047]))

4) L_1 tends to produce sparser weights.

INTUITION



FORMALIC INTUITION

$$L_1: L(w) = \lambda \|w\|_1$$

$$\frac{\partial L_1}{\partial \vec{w}} : \frac{\partial L}{\partial \vec{w}} = \lambda$$

$$L_2: L(w) = \frac{\lambda}{2} \|w\|_2^2$$

$$\frac{\partial L_2}{\partial \vec{w}} : \frac{\partial L}{\partial \vec{w}} = \lambda \vec{w}$$

As weights get smaller

L_1 update weights by same $\boxed{\lambda}$ but L_2 update by $\boxed{\lambda w}$

Net change in $w \rightarrow \Delta w$ ↑ in L_1 than L_2 as weights get smaller

So, L_1 tends to produce sparser weights than L_2

PROBLEM 4: GAUSSIAN NAÏVE BAYES

1. Log likelihood of the Gaussian Naïve Bayes model, compute MLE for each of the parameters

Naïve Bayes makes the following conditional independence assumption

$$P(\mathbf{x} | Y=c_i) = \prod_{j=1}^n P(x_j | Y=c_i)$$

Gaussian Naïve Bayes model description :

- 1> Class priors : $P(Y)$
- 2> Conditional distributions that are gaussian.

$$P(x_j=x_j | Y=c_i) = N(\mu_{i,j}, \sigma_{i,j}^2)$$

Given a dataset with m continuous features.

$$\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}) \dots (x^{(m)}, y^{(m)})\}$$

Log likelihood of the data:

$$\log \left(\prod_{m=1}^M P(x=x^{(m)} | Y=y^{(m)}) \right)$$

(Curved arrow pointing from the first term to the sum)

$$\sum_{m=1}^M \log P(x=x^{(m)} | Y=y^{(m)})$$

$$\sum_{m=1}^M \sum_{j=1}^n \log p(x_j = x_j^{(m)} | Y = y^{(m)})$$

II

[ESTIMATING GAUSSIAN PARAMETERS]

II

$$\sum_{m=1}^M \sum_{j=1}^n \log p(x_j = x_j^{(m)} | Y = y^{(m)})$$

[INDICATOR RANDOM VARIABLE]

$$= \sum_{m=1}^M \sum_{j=1}^n \sum_{i=1}^k \underbrace{\log p(x_j = x_j^{(m)} | Y = c_i)}_{\text{Gaussian (By our assumption). } N(\mu_{ij}, \sigma_{ij}^2)} \cdot \mathbf{1}(y^{(m)} = c_i)$$

L

$$\underset{\mu_{ij}, \sigma_{ij}^2}{\operatorname{argmax}} \sum_{m=1}^M \sum_{j=1}^n \sum_{i=1}^k \left(-\frac{1}{2} \left(\frac{x_j^{(m)} - \mu_{ij}}{\sigma_{ij}} \right)^2 - \log(\sigma_{ij} \sqrt{2\pi}) \right) \cdot \mathbf{1}(y^{(m)} = c_i)$$

$$\frac{\partial L}{\partial \mu_{ij}} = \sum_{m=1}^M \frac{\partial}{\partial \mu_{ij}} \left[\left(-\frac{1}{2} \left(\frac{x_j^{(m)} - \mu_{ij}}{\sigma_{ij}} \right)^2 - \log(\sigma_{ij} \sqrt{2\pi}) \right) \cdot \mathbf{1}(y^{(m)} = c_i) \right]$$

$$\sum_{m=1}^M \left(\frac{\mu_{ij} - x_j^{(m)}}{\sigma_{ij}} \right) \cdot \mathbf{1}(y^{(m)} = c_i) = 0 \Rightarrow \mu_{ij} = \frac{\sum_{m=1}^M x_j^{(m)} \cdot \mathbf{1}(y^{(m)} = c_i)}{\sum_{m=1}^M \mathbf{1}(y^{(m)} = c_i)}$$

$$\sum_{m=1}^M \sum_{j=1}^n \sum_{i=1}^k \left(-\frac{1}{2} \left(\frac{x_j^{(m)} - \mu_{i,j}}{\sigma_{i,j}} \right)^2 - \log(\sigma_{i,j} \sqrt{2\pi}) \right) \cdot \mathbb{1}(y^{(m)} = c_i)$$

$$\frac{\partial L}{\partial \sigma_{i,j}} = \sum_{m=1}^M \frac{\partial}{\partial \sigma_{i,j}} \left[\left(-\frac{1}{2} \left(\frac{x_j^{(m)} - \mu_{i,j}}{\sigma_{i,j}} \right)^2 - \log(\sigma_{i,j} \sqrt{2\pi}) \right) \cdot \mathbb{1}(y^{(m)} = c_i) \right]$$

$\cancel{\sigma} = 0$

$$= \sum_{m=1}^M \left(\cancel{-\frac{1}{2}} \cancel{(-2)} \frac{1}{\sigma_{i,j}^3} (x_j^{(m)} - \mu_{i,j})^2 - \frac{1}{\sigma_{i,j}} \right) \cdot \mathbb{1}(y^{(m)} = c_i)$$

$$= \sum_{m=1}^M \left(\frac{(x_j^{(m)} - \mu_{i,j})^2}{\sigma_{i,j}^3} - \frac{1}{\sigma_{i,j}} \right) \cdot \mathbb{1}(y^{(m)} = c_i) = 0.$$

$$\Rightarrow \bar{\sigma}_{i,j}^2 = \frac{\sum_{m=1}^M (x_j^{(m)} - \mu_{i,j})^2 \cdot \mathbb{1}(y^{(m)} = c_i)}{\sum_{m=1}^M \mathbb{1}(y^{(m)} = c_i)}$$

$\bar{\sigma}_{i,j}$ is a BIASED ESTIMATOR.

UNBIASED VERSION

$$\hookrightarrow \bar{\sigma}_{i,j}^2 = \frac{\sum_{m=1}^M (x_j^{(m)} - \mu_{i,j})^2 \cdot \mathbb{1}(y^{(m)} = c_i)}{\sum_{m=1}^M \mathbb{1}(y^{(m)} = c_i) - 1}$$

2. Test accuracy of the GNB model : 69.23076923076923

3. PRIOR ON NAÏVE BAYES MODEL :

a. We generally choose a conjugate prior so that posterior also has same functional form

i. If we choose a Gaussian Prior, it will be the conjugate distribution of the likelihood function as the corresponding posterior will be product of two exponentials of quadratic functions and hence shall be Gaussian

$$P(X|\mu) = \prod_{m=1}^M P(x_m|\mu) = \frac{1}{(2\pi\sigma^2)^{N/2}} \exp \left\{ -\frac{1}{2\sigma^2} \sum_{m=1}^M (x_m - \mu)^2 \right\}$$

[LIKELIHOOD FUNCTION]

(Probability of observed data given μ) (EXponential OF QUADRATIC FORM.)

If we choose $p(\mu)$ by GAUSSIAN \rightarrow [CONJUGATE DISTRIBUTION]

Posterior has same functional form as PRIOR

$$P(X|\lambda) = \prod_{m=1}^M N(x_m|\mu, \lambda^{-1}) \propto \lambda^{M/2} \exp \left\{ -\frac{\lambda}{2} \sum_{m=1}^M (x_m - \mu)^2 \right\}$$

Likelihood function for λ

So corresponding conjugate prior \rightarrow product of a power of λ exponential of a linear function of λ

Corresponds to a GAMMA DISTRIBUTION

$$\text{Gam}(\lambda|a,b) = \frac{1}{\Gamma(a)} b^a \lambda^{a-1} \exp(-b\lambda)$$

4. I don't think Naïve Bayes is reasonable here for the following reasons

Naïve Bayes assumes that features are conditionally independent given classes

- a. Logistic regression being a discriminative classifier without any assumptions produced higher test accuracies than the Naïve Bayes model for this dataset
 - i. Logistic Regression (No penalty) - 75%
 - ii. Logistic Regression (L1 penalty) - 75%
 - iii. Logistic Regression (L2 penalty) - 75%

All are higher than Naïve Bayes