**Q1 (a)** In the class we showed that the margin is $\frac{2c}{||w||}$ and we want to maximiz this margin.

The magnitude of $c$ is merely scales $w$ and $b$ in below

$$\max \frac{2c}{||w||} \quad s.t. \quad y^i(w^Tx^i+b) \geqslant c \quad \forall i$$

and because $c$ is constant we need to find the optimum value for $\frac{2c}{||w||}$, we can omit $c$ and "2".

hence we can find the statement can be replaced by

$$\frac{1}{||w||} \quad s.t. \quad y^i(w^Tx^i+b) \geqslant 1 \quad \forall i$$

**Q1 (b)** the Lagrangian dual Problem is

$$L(w,a,b) = \frac{1}{2}ww^T + \sum_{i=1}^{m} a_i(1-y^i(w^Tx^i+b))$$

$$\frac{\partial L(w,a,b)}{\partial w} = w + 0 - \sum_{i=1}^{m} a_i y^i x^i + 0 = 0$$

$$= w - \sum_{i=1}^{m} a_i y^i x^i = 0 \implies \boxed{w = \sum_{i=1}^{m} a_i y^i x^i}$$

in the HW Question $n = m$ & $y^i = y_i$ & $x^i = x_i$.

Q2

(a):

$x_i$ number of times word $V_i$ appears in $x$

$n = 15$ total Vocabulary in $V$

$x$ is of length 15. this is the input data

$y = 0$ when input data is NOT spam

$y = 1$ when input data is spam

$P(y=0) = P(NOT spam) = \frac{4}{7}$ $\longrightarrow$ Not spam $\longrightarrow$ total data

these are the prior

$P(y=1) = P(spam) = \frac{3}{7}$

(b):

million dollar offer        [0,1,0,0, 0,0,0,1,1, 0, 0, 0, 0,0,0]

secret offer today          [1,1,0,0, 0,0,1, 0, 0,0, 0, 0, 0,0,0]

secret is secret            [2,0, 0,0,0,0,0, 0,0, 0,1,0,0,0,0]

low price for valued customer   [0,0,1,1,1,1, 0, 0,0, 0, 0,1, 0, 0, 0]

play secret sports today    [1,0,0,0,0,0,1,0,0,1,0,0,1,0,0]

sports is healthy           [0,0,0,0, 0,0,0,0,0,1, 1,0,0,1,0]

low price pizza             [0,0,1,1,0,0, 0,0,0,0, 0, 0, 0,0,1]

(C) According to bayes rule $P(y|x) = \dfrac{P(x|y)\,P(y)}{P(x)}$

Prior

based on naive bayes rule $\Rightarrow P(y|x) = P(x|y)\,P(y)$

Posterior      likelihood

if $P(y=i|x) > P(y=j|x) \Rightarrow y=i$   else $\Rightarrow y=j$

Because we assume indeptndance for all data incidence Then we can write the probability for data X to be $y=1$ (spam) or $y=0$ (ham) can be written as:

$$P(y|X) = P(X,y) = P(x^1, y^1) \cdot P(x^2, y^2) \cdots P(x^m, y^m)$$

$$= \prod_{i=1}^{m} P(x^i, y^i) \quad \boxed{\text{Likelihood}}$$

$\boxed{\text{m is total number of records of data}}$

$$\underset{\text{naive}}{=} \prod_{i=1}^{m} P(x|y^i)\, P(y^i)$$

$$= \prod_{i=1}^{m} P(y^i) \prod_{k=1}^{N} \theta_{c,k}^{x_k^i}$$

(why in next page)$^{\pm}$

$\boxed{\text{we want to maximize this equation. this is constained optim}}$

for simplification we take the log (ln) from both sides. log is monotonic so the maximization result will hold

the constrained is: $\displaystyle\sum_{k=1}^{N} \theta_{c,k} = 1$   $\forall_0$

\* the answer to why from previous page:
  we are looking for maximizing the $P(X|y)$ likelihood.

to max $P(X|y) = \dfrac{P(y|X)}{P(y)}$ we can maximiz $P(y|X)$

because $\log P(y) < 1$. And it will vanish when taking Foc of lagrange.

Now the optimization problem, by taking log:

$$\log P(X,y) = \log \prod_{i=1}^{m} \left( P(y^i) \prod_{k=1}^{n} \theta_{c,k}^{x_k^i} \right) = \sum_{i=1}^{m} \log \left( P(y^i) \prod_{k=1}^{n} \theta_{c,k}^{x_k^i} \right) =$$

$$= \sum_{i=1}^{m} \left\{ \log (P(y^i)) + \sum_{k=1}^{n} \log \theta_{c,k}^{x_k^i} \right\}$$

with constraint $\displaystyle\sum_{k=1}^{n} \theta_{c,k} = 1 \quad \forall c$

$\theta_{c,k}^{x_k^i}$ is the likelihood of word $x_k$ from sentence $i$ that belongs to category $c$.

$y^i$ is the catigory indicator. Here $c$ is the category here. it can be $c$

$c$ is variable for category. here is $1$ or $0$

the lagrangian optimization problem will be:

$$L(P(X,y)) = \sum_{i=1}^{m} \left\{ \log (P(y^i)) + \sum_{k=1}^{n} \log \theta_{c,k}^{x_k^i} \right\} + \lambda \left( \sum_{k=1}^{n} \theta_{c,k} - 1 \right)$$

we'll take derivative from $\underline{a}$ $\theta_{c,k}$ which pertains to $\underline{a}$ $c$
this will make the derivative easier to see.

$$L(P(X,y)) = \sum_{i=1}^{m} \left\{ \log(P(y^i=c)) + \sum_{k=1}^{n} \log \theta_{c,k}^{x_k^i} \right\} + \lambda \left( \sum_{k=1}^{n} \theta_{c,k} - 1 \right)$$

$$\frac{\partial (P(X,y))}{\partial \theta_{c,k}} = \sum_{i=1}^{m} \left\{ 0 + \frac{1}{\theta_{c,k}} 1^{x_{ck}^i} \right\} + \lambda = 0$$

Because we are differentiating for a particular $\theta_{c,k}$ hence
only one term from $\lambda \sum_{k=1}^{n} \theta_{c,k}$ can be differenciated to $1$.
the rest are constants.

$1^{x_{ck}^i}$ : indicates the presence of word $k$ in group $c$ in sentence $i$
so if the word $k$ exist in sentence $i$ group $c$ the
value is $\underline{1}$ otherwise it is $\underline{0}$.

$$\frac{\partial (P(X,y))}{\partial \theta_{c,k}} = \frac{1}{\theta_{c,k}} \sum_{i=1}^{m} 1^{x_{ck}^i} + \lambda = 0 \Rightarrow \theta_{c,k} = \frac{-\sum_{i=1}^{m} 1^{x_{ck}^i}}{\lambda} \quad \circledast$$

since $\sum_{k=1}^{n} \theta_{c,k} = 1 \Rightarrow \sum_{k=1}^{n} \theta_{c,k} = \frac{-\sum_{k=1}^{n} \sum_{i=1}^{M} 1^{x_{ck}^i}}{\lambda} = 1$

$$\Rightarrow \lambda = - \sum_{k=1}^{n} \sum_{i=1}^{m} 1^{x_{ck}^i}$$

$$\Rightarrow \text{replace } \lambda \text{ in } \circledast$$

$$\theta_{c,k} = \frac{\sum_{i=1}^{m} 1^{x_{ck}^{i}}}{\sum_{k=1}^{n} \sum_{i=1}^{m} 1^{x_{ck}^{i}}}$$

→ this says count all word k in the category c

→ this says count all words from all sencences with category c

$$\theta_{0,1} = \frac{3.}{9 \, (8=0)} = \frac{1}{3} \qquad \theta_{1,1} = \frac{1}{15}$$

$$\theta_{0,7} = \frac{1}{9} \qquad\qquad \theta_{1,15} = \frac{1}{15}$$

(d) $P(\text{today is secret} | \text{spam}) = P((1,0,0,0,0,0,1,0,0,0,1,0,0,0,0) | \text{spam})$

$= \frac{3}{9} \times (1-\frac{2}{9}) \times 1 \times 1 \times 1 \times 1 \times \frac{1}{9} \times (1-\frac{1}{9}) \times (1-\frac{1}{9}) \times 1 \times \frac{1}{9} \times 1 \times 1 \times 1 \times 1$

$= 0.0025$

|          | Spam | Ham  |
|----------|------|------|
| secret   | 3/9  | 1/15 |
| after    | 2/9  | 0    |
| low      | 0    | 2/15 |
| price    | 0    | 2/15 |
| valued   | 0    | 1/15 |
| customer | 0    | 1/15 |
| today    | 1/9  | 1/15 |
| dollar   | 1/9  | 0    |
| million  | 1/9  | 0    |
| sports   | 0    | 2/15 |
| is       | 1/9  | 1/15 |
| for      | 0    | 1/15 |
| play     | 0    | 1/15 |

|        | Spam | ham  |
|--------|------|------|
| healty | 0    | 1/15 |
| Pizza  | 0    | 1/15 |

$P(\text{today is secret} | \text{ham}) = \frac{1}{13} \times 1 \times (1-\frac{2}{15}) \times (1-\frac{2}{15})$

$\times (1-\frac{1}{15}) \times (1-\frac{1}{15}) \times \frac{1}{15} \times \times 1 \times 1 \times 1 \times (1-\frac{2}{15}) \times \frac{1}{15}$

$\times (1-\frac{1}{15}) \times (1-\frac{1}{13}) \times (1-\frac{1}{15}) (1-\frac{1}{15}) = 0.00013$

$P(\text{ham} | \text{today is secret}) = \frac{0.0025 \times \frac{4}{7}}{0.0025 \times \frac{1}{7} + 0.000113 \times \frac{3}{7}}$

$= \frac{0.0014}{0.001497} = 0.967$

$P(\text{ham} | \text{today is secret}) = \frac{0.00005591}{0.001477} = 0.037$

It is Spam

# Q3

## Q3-a-i

When I ran the classification with the test data (20%), all three classifiers performed the same. As shown below.

```
LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
          intercept_scaling=1, max_iter=100, multi_class='warn',
          n_jobs=None, penalty='l2', random_state=None, solver='warn',
          tol=0.0001, verbose=0, warm_start=False)
              precision    recall  f1-score   support

        0.0       0.93      1.00      0.97        14
        1.0       1.00      0.95      0.97        20

   micro avg       0.97      0.97      0.97        34
   macro avg       0.97      0.97      0.97        34
weighted avg       0.97      0.97      0.97        34

[[14  0]
 [ 1 19]]
```

```
GaussianNB(priors=None, var_smoothing=1e-09)
              precision    recall  f1-score   support

        0.0       0.93      1.00      0.97        14
        1.0       1.00      0.95      0.97        20

   micro avg       0.97      0.97      0.97        34
   macro avg       0.97      0.97      0.97        34
weighted avg       0.97      0.97      0.97        34

[[14  0]
 [ 1 19]]
```

```
KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
          metric_params=None, n_jobs=None, n_neighbors=5, p=2,
          weights='uniform')
              precision    recall  f1-score   support

        0.0       0.93      1.00      0.97        14
        1.0       1.00      0.95      0.97        20

   micro avg       0.97      0.97      0.97        34
   macro avg       0.97      0.97      0.97        34
weighted avg       0.97      0.97      0.97        34

[[14  0]
 [ 1 19]]
```

But then I ran the classification on the entire data, the I noticed that logistic regression had higher precision. I believe that data in the divorce dataset is linearly separable and therefore logistic regression performed the best and Gaussian Naïve Bayes and KNN performed slightly worse than logistic regression.

## Logistic Regression

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0.0 | 0.99 | 1.00 | 0.99 | 86 |
| 1.0 | 1.00 | 0.99 | 0.99 | 84 |
| micro avg | 0.99 | 0.99 | 0.99 | 170 |
| macro avg | 0.99 | 0.99 | 0.99 | 170 |
| weighted avg | 0.99 | 0.99 | 0.99 | 170 |

```
[[86  0]
 [ 1 83]]
```

## Guassian Naïve Bayes

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0.0 | 0.96 | 1.00 | 0.98 | 86 |
| 1.0 | 1.00 | 0.95 | 0.98 | 84 |
| micro avg | 0.98 | 0.98 | 0.98 | 170 |
| macro avg | 0.98 | 0.98 | 0.98 | 170 |
| weighted avg | 0.98 | 0.98 | 0.98 | 170 |

```
[[86  0]
 [ 4 80]]
```

## KNearest Neighbor

|  | precision | recall | f1-score | support |
|---|---|---|---|---|
| 0.0 | 0.96 | 1.00 | 0.98 | 86 |
| 1.0 | 1.00 | 0.95 | 0.98 | 84 |
| micro avg | 0.98 | 0.98 | 0.98 | 170 |
| macro avg | 0.98 | 0.98 | 0.98 | 170 |
| weighted avg | 0.98 | 0.98 | 0.98 | 170 |

```
[[86  0]
 [ 4 80]]
```
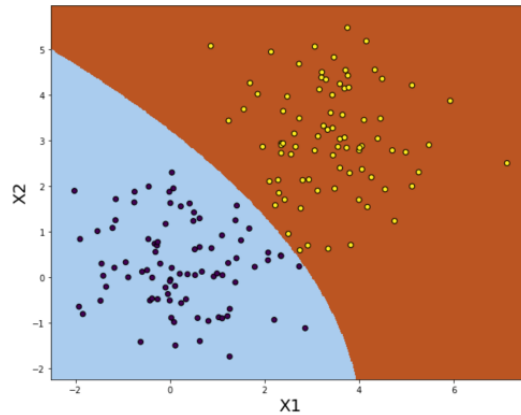
## Q3-a-ii

As can be seen from the figures below, logistic regression is classifiers works best for data that linearly separable. KNN shows more overfitting, and Naïve Bayes is more suitable for non- linearly separable data sets classification.
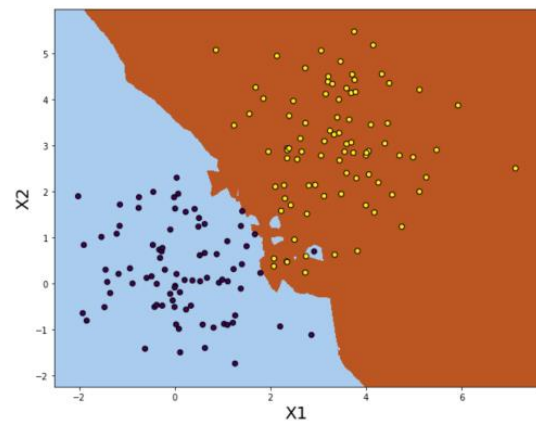
Logistic Regression Classification


Gaussian Naive Bayes Classification


KNN Classification

## Q3-b

For the part b, we had more data in our test dataset. The results are more impressive than section before. As it can be seen in figure below, KNN performed the best. And the reason for this performance is because the MNIST data set is less linearly separable and KNN performed the best here.

```
LogisticRegression(C=1.0, class_weight=None, dual=False, fit_intercept=True,
          intercept_scaling=1, max_iter=100, multi_class='warn',
          n_jobs=None, penalty='l2', random_state=None, solver='warn',
          tol=0.0001, verbose=0, warm_start=False)
              precision    recall  f1-score   support

           0       0.98      0.98      0.98       201
           1       0.98      0.98      0.98       197

   micro avg       0.98      0.98      0.98       398
   macro avg       0.98      0.98      0.98       398
weighted avg       0.98      0.98      0.98       398

[[197    4]
 [  3  194]]
```

```
GaussianNB(priors=None, var_smoothing=1e-09)
              precision    recall  f1-score   support

           0       0.98      0.63      0.77       201
           1       0.72      0.99      0.84       197

   micro avg       0.81      0.81      0.81       398
   macro avg       0.85      0.81      0.80       398
weighted avg       0.85      0.81      0.80       398

[[126  75]
 [  2 195]]




KNeighborsClassifier(algorithm='auto', leaf_size=30, metric='minkowski',
          metric_params=None, n_jobs=None, n_neighbors=5, p=2,
          weights='uniform')
              precision    recall  f1-score   support

           0       1.00      1.00      1.00       201
           1       0.99      1.00      1.00       197

   micro avg       1.00      1.00      1.00       398
   macro avg       1.00      1.00      1.00       398
weighted avg       1.00      1.00      1.00       398

[[200   1]
 [  0 197]]
```