

Problem Set #2 Solutions: Supervised Learning II

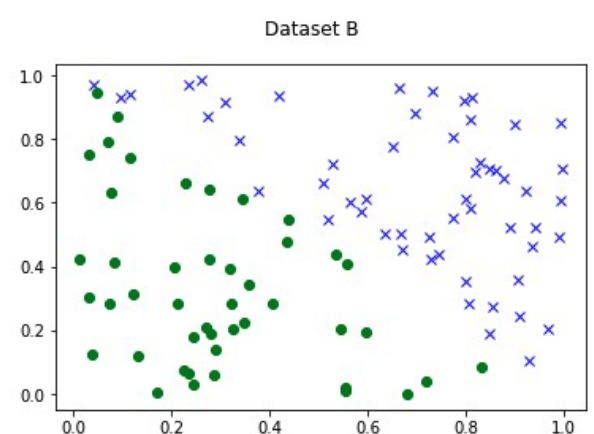
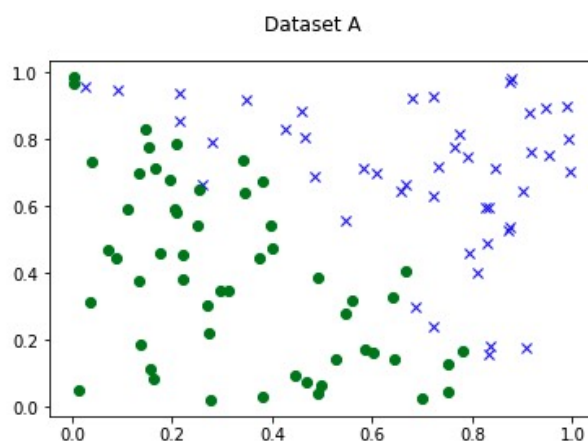
1.

(a)

Logistic regression converged on dataset A , but didn't converge on dataset B . Run `python P01_lr.py`

```
$ python p01_lr.py
==== Training model on data set A ====
Finished 10000 iterations
Finished 20000 iterations
Finished 30000 iterations
Converged in 30357 iterations
=====
==== Training model on data set B ====
Finished 10000 iterations
Finished 20000 iterations
Finished 30000 iterations
Finished 40000 iterations
Finished 50000 iterations
Finished 60000 iterations
Finished 70000 iterations
Finished 80000 iterations
Finished 90000 iterations
Finished 100000 iterations
```

(b)



At first glance, the only difference between the two datasets seems to be that dataset B is linearly separable, whereas dataset A is not.

```

def calc_grad(X, Y, theta):
    """Compute the gradient of the loss with respect to theta."""
    m, n = X.shape

    margins = Y * X.dot(theta)
    probs = 1. / (1 + np.exp(margins))
    grad = -(1./m) * (X.T.dot(probs * Y))

    return grad

def logistic_regression(X, Y):
    """Train a logistic regression model."""
    m, n = X.shape
    theta = np.zeros(n)
    learning_rate = 10

    i = 0
    while True:
        i += 1
        prev_theta = theta
        grad = calc_grad(X, Y, theta)
        theta = theta - learning_rate * grad
        if i % 10000 == 0:
            print('Finished %d iterations' % i)
        if np.linalg.norm(prev_theta - theta) < 1e-15:
            print('Converged in %d iterations' % i)
            break
    return

```

Recall that in SVM the functional margin $\hat{\gamma}$ is

$$\hat{\gamma}^{(i)} = y^{(i)}(w^T x^{(i)} + b)$$

Because there is no constraint on w (such as $\|w\|_2 = 1$), we can scale the w and b to increase the functional margin without changing the decision boundary.

In this problem, the labels y are $\{-1, +1\}$ instead of $\{0, 1\}$. So the loss function $J(\theta)$ is

$$J(\theta) = \frac{1}{m} \sum_i^m \log(1 + \exp\{-y^{(i)} \theta^T x^{(i)}\}) \quad \checkmark$$

Notice there is $y^{(i)} \theta^T x^{(i)}$ in the expression above, this has a similar property like the functional margin.

When the dataset is linearly separable, $y^{(i)} \theta^T x^{(i)} > 0$ for all training examples. So we can scale θ to make $J(\theta)$ smaller (close to 0). However, when the dataset is not linearly separable, $y^{(i)} \theta^T x^{(i)}$ could be greater or smaller than 0. So we can't arbitrarily scale θ to reduce $J(\theta)$.

(c)

i.

No, using a different learning rate will not help to reduce the value of θ .

ii.

~~No~~ Yes, using learning rate decay (e.g. by factor $1/t^2$) will make $\alpha \nabla_{\theta} J(\theta) \leq 10^{-15}$ in a few iterations.

iii.

No, as you can see in the pictures the input feature are already scaled.

iv.

Yes, adding L_2 regularization will help reduce the value of θ .

v.

Yes, adding noise could make the dataset becomes not linearly separable.

But how to control the scale of noise to avoid losing accuracy?

(d) η

SVM use hinge loss is not vulnerable to linearly separable dataset like B .

Here is the hinge loss

$$J(\hat{y}) = \max(0, 1 - y \cdot \hat{y}), \text{ where } \hat{y} = w^T x + b$$

Assume that the dataset is linearly separable, so $y \cdot \hat{y} > 0$.

When we increase w and b to make $|\hat{y}| \geq 1$, then $J(\hat{y}) = 0$.

2.

(a)

The log likelihood is

$$\ell(\theta) = \log L(\theta) = \sum_{i=1}^m y^{(i)} \log h(x^{(i)}) + (1 - y^{(i)}) \log(1 - h(x^{(i)}))$$

After training, the gradients are equal to 0

$$\frac{\partial \ell(\theta)}{\partial \theta_j} = \sum_{i=1}^m (y^{(i)} - h(x^{(i)})) x_j^{(i)} = 0$$

Set $j = 0$. Because $x_0^{(i)} = 1$, so

$$\sum_{i=1}^m (y^{(i)} - h(x^{(i)})) = 0$$

$$\sum_{i=1}^m h(x^{(i)}) = \sum_{i=1}^m y^{(i)}$$

$$h(x^{(i)}) = P(y^{(i)} = 1 | x^{(i)}; \theta), \quad y^{(i)} = \mathbb{I}\{y^{(i)} = 1\}$$

$$\sum_{i=1}^m P(y^{(i)} = 1 | x^{(i)}; \theta) = \sum_{i=1}^m \mathbb{I}\{y^{(i)} = 1\}$$

When $(a, b) = (0, 1)$, $I_{a,b} = \{x^{(i)}, y^{(i)}\}_{i=1}^m$ and $|\{i \in I_{a,b}\}| = m$

$$\frac{\sum_{i \in I_{a,b}} P(y^{(i)} = 1 | x^{(i)}; \theta)}{|\{i \in I_{a,b}\}|} = \frac{\sum_{i \in I_{a,b}} \mathbb{I}\{y^{(i)} = 1\}}{|\{i \in I_{a,b}\}|}$$

(b)

The model is perfectly calibrated doesn't necessarily imply that the model achieves perfect accuracy.

The converse is also not necessarily true.

①

Assume that $(a, b) = (0.5, 1)$.

When the model achieves perfect accuracy, the predictions are all correct, i.e.

$$\sum_{i \in I_{a,b}} \mathbb{I}\{y^{(i)} = 1\} = |\{i \in I_{a,b}\}|$$

For all $i \in I_{a,b}$

$$0.5 < P(y^{(i)} = 1 | x^{(i)}; \theta) < 1$$

So

$$\frac{\sum_{i \in I_{a,b}} P(y^{(i)} = 1 | x^{(i)}; \theta)}{|\{i \in I_{a,b}\}|} < \frac{\sum_{i \in I_{a,b}} \mathbb{I}\{y^{(i)} = 1\}}{|\{i \in I_{a,b}\}|}$$

However, when the model is perfectly calibrated, the following property always hold

$$\frac{\sum_{i \in I_{a,b}} P(y^{(i)} = 1 | x^{(i)}; \theta)}{|\{i \in I_{a,b}\}|} = \frac{\sum_{i \in I_{a,b}} \mathbb{I}\{y^{(i)} = 1\}}{|\{i \in I_{a,b}\}|}$$

So model is perfectly calibrated doesn't mean model achieves perfect accuracy. The converse neither.

(c)

When adding L_2 regularization, θ is not the maximum likelihood parameter learned after training.

Furthermore, the loss function is

$$J(\theta) = - \sum_{i=1}^m y^{(i)} \log h(x^{(i)}) + (1 - y^{(i)}) \log(1 - h(x^{(i)})) + \frac{1}{2} \lambda \|\theta\|_2^2$$

After training, the gradients are equal to 0

$$\frac{\partial J(\theta)}{\partial \theta_j} = \sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) x_j^{(i)} + \lambda \theta_j = 0$$

Set $j = 0$. Because $x_0^{(i)} = 1$, so

$$\sum_{i=1}^m (h(x^{(i)}) - y^{(i)}) + \lambda \theta_0 = 0$$

$$\sum_{i=1}^m h(x^{(i)}) + \lambda \theta_0 = \sum_{i=1}^m y^{(i)}$$

$$\sum_{i=1}^m P(y^{(i)} = 1 | x^{(i)}; \theta) + \lambda \theta_0 = \sum_{i=1}^m \mathbb{I}\{y^{(i)} = 1\}$$

So the model will not be well-calibrated.

② A model having perfect accuracy does not mean it is perfectly calibrated. Because:

$$\frac{\sum_{i \in I_{a,b}} P(y^{(i)} = 1 | x^{(i)}; \theta)}{|\{i \in I_{a,b}\}|} = 1 > \frac{\sum_{i \in I_{a,b}} \mathbb{I}\{y^{(i)} = 1\}}{|\{i \in I_{a,b}\}|}$$

3.

(a)

$$p(\theta|x, y) = \frac{p(x, y, \theta)}{p(x, y)} = \frac{p(y|x, \theta)p(x, \theta)}{p(x, y)} = \frac{p(y|x, \theta)p(\theta)p(x)}{p(x, y)}$$

Assume that $p(\theta) = p(\theta|x)$, then

$$p(\theta|x, y) = \frac{p(y|x, \theta)p(\theta)p(x)}{p(x, y)} = p(y|x, \theta)p(\theta) \cdot \frac{p(x)}{p(x, y)}$$

$$\theta_{\text{MAP}} = \arg \max_{\theta} p(\theta|x, y) = \arg \max_{\theta} p(y|x, \theta)p(\theta) \left(\frac{p(x)}{p(x, y)} \right) \xrightarrow{\text{固定常数}} \arg \max_{\theta} p(y|x, \theta)p(\theta)$$

(b)

$$\begin{aligned} \theta_{\text{MAP}} &= \arg \max_{\theta} p(y|x, \theta)p(\theta) \\ &= \arg \max_{\theta} \log(p(y|x, \theta)p(\theta)) \\ &= \arg \max_{\theta} \log p(y|x, \theta) + \log p(\theta) \\ &= \arg \min_{\theta} -\log p(y|x, \theta) - \log p(\theta) \end{aligned}$$

$$\theta \sim \mathcal{N}(0, \eta^2 I)$$

$$p(\theta) = \frac{1}{(2\pi)^{d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(\theta - \mu)^T \Sigma^{-1}(\theta - \mu)\right)$$

$$\log p(\theta) = -\frac{1}{2} \theta^T \Sigma^{-1} \theta$$

$$\begin{aligned} \theta_{\text{MAP}} &= \arg \min_{\theta} -\log p(y|x, \theta) - \log p(\theta) \\ &= \arg \min_{\theta} -\log p(y|x, \theta) + \frac{1}{2\eta^2} \|\theta\|_2^2 \end{aligned}$$

$$\lambda = \frac{1}{2\eta^2}$$

(c)

$$\epsilon^{(i)} \sim \mathcal{N}(0, \sigma^2)$$

$$y^{(i)} = \theta^T x^{(i)} + \epsilon^{(i)}$$

$$y^{(i)}|x^{(i)}, \theta \sim \mathcal{N}(\theta^T x^{(i)}, \sigma^2)$$

$$p(y^{(i)}|x^{(i)}, \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2} (y^{(i)} - \theta^T x^{(i)})^2\right\}$$

$$\begin{aligned} p(\vec{y}|X, \theta) &= \prod_{i=1}^m p(y^{(i)}|x^{(i)}, \theta) \\ &= \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left\{-\frac{1}{2\sigma^2} (y^{(i)} - \theta^T x^{(i)})^2\right\} \\ &= \frac{1}{(2\pi)^{m/2} \sigma^m} \exp\left\{-\frac{1}{2\sigma^2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2\right\} \\ &= \frac{1}{(2\pi)^{m/2} \sigma^m} \exp\left\{-\frac{1}{2\sigma^2} \|X\theta - \vec{y}\|_2^2\right\} \end{aligned}$$

$$\log p(\vec{y}|X, \theta) = -\frac{m}{2} \log(2\pi) - m \log \sigma - \frac{1}{2\sigma^2} \|X\theta - \vec{y}\|_2^2$$

$$\theta_{\text{MAP}} = \arg \min_{\theta} -\log p(y|x, \theta) + \frac{1}{2\eta^2} \|\theta\|_2^2$$

$$= \arg \min_{\theta} \frac{1}{2\sigma^2} \|X\theta - \vec{y}\|_2^2 + \frac{1}{2\eta^2} \|\theta\|_2^2$$

$$\nabla_{\theta} J(\theta) = \nabla_{\theta} \left(\frac{1}{2\sigma^2} (\vec{y} - X\theta)^T (\vec{y} - X\theta) + \frac{1}{2\eta^2} \|\theta\|_2^2 \right)$$

$$J(\theta) = \frac{1}{2\sigma^2} \|X\theta - \vec{y}\|_2^2 + \frac{1}{2\eta^2} \|\theta\|_2^2 = \frac{1}{2\sigma^2} (\vec{y} - X\theta)^T (\vec{y} - X\theta) + \frac{1}{2\eta^2} \|\theta\|_2^2$$

$$= \frac{1}{2\sigma^2} (\theta^T X^T X \theta - 2\vec{y}^T X \theta + \vec{y}^T \vec{y}) + \frac{1}{2\eta^2} \theta^T \theta$$

$$\nabla_{\theta} J(\theta) = \frac{1}{\sigma^2} (X^T X \theta - X^T \vec{y}) + \frac{1}{\eta^2} \theta = 0$$

$$= \frac{1}{\sigma^2} (X^T X \theta - X^T \vec{y}) + \frac{1}{\eta^2} \theta = 0$$

$$\therefore (X^T X + \frac{\sigma^2}{\eta^2}) \theta = X^T \vec{y}$$

$$\theta_{\text{MAP}} = \arg \min_{\theta} J(\theta) = (X^T X + \frac{\sigma^2}{\eta^2} I)^{-1} X^T \vec{y}$$

$$(d) \quad \theta = (X^T X + \frac{\sigma^2}{\eta^2})^{-1} X^T \vec{y}$$

$$\theta_{\text{MAP}} = \arg \min_{\theta} -\log p(y|x, \theta) - \log p(\theta)$$

$$\therefore \theta_{\text{MAP}} = \frac{1}{2\sigma^2} \sum_{i=1}^m (y^{(i)} - \theta^T x^{(i)})^2 + \sum_{i=1}^n \frac{1}{b} |\theta_i|$$

$$(\vec{y} - X\theta)^T (\vec{y} - X\theta)$$

$$\theta \sim \mathcal{L}(0, bI) \quad y = \theta^T x + \epsilon, \quad \epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$\therefore y \sim \mathcal{N}(\theta^T x, \sigma^2)$$

$$\therefore p(y|x, \theta) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - \theta^T x^{(i)})^2}{2\sigma^2}\right)$$

$$p(\theta) = \frac{1}{2^b b} \exp\left(-\frac{|\theta_i|}{b}\right)$$

$$\theta_{\text{MAP}} = \arg \min_{\theta} \frac{1}{2\sigma^2} \|X\theta - \vec{y}\|_2^2 - \log p(\theta)$$

$$= \arg \min_{\theta} \frac{1}{2\sigma^2} \|X\theta - \vec{y}\|_2^2 + \frac{1}{b} \|\theta\|_1$$

$$\therefore J(\theta) = \|X\theta - \vec{y}\|_2^2 + \gamma \|\theta\|_1$$

$$\theta_{\text{MAP}} = \arg \min_{\theta} J(\theta)$$

$$\gamma = \frac{2\sigma^2}{b}$$

4.

(a)

$\star \rightarrow \vec{z}$

Yes, K_1 and K_2 are both PSD, so $K_1 + K_2$ is PSD.

$$z^T K z = z^T (K_1 + K_2) z = z^T K_1 z + z^T K_2 z \geq 0$$

(b)

No, although K_1 and K_2 are both PSD, $K_1 - K_2$ may not be PSD.

For example, $K_2 = 2K_1$

$$z^T K z = z^T (K_1 - K_2) z = z^T (K_1 - 2K_1) z = -z^T K_1 z \leq 0$$

(c)

Yes, K_1 is PSD, so aK_1 ($a \in \mathbb{R}^+$) is PSD.

$$z^T K z = z^T a K_1 z = a \cdot z^T K_1 z \geq 0$$

(d)

No, K_1 is PSD, so $-aK_1$ ($a \in \mathbb{R}^+$) is not PSD.

$$z^T K z = z^T (-aK_1) z = -a \cdot z^T K_1 z \leq 0$$

(e)

Yes, $K_1 K_2$ is PSD.

$$\begin{aligned} z^T K z &= \sum_i \sum_j z_i K_{ij} z_j \\ &= \sum_i \sum_j z_i K_1(x^{(i)}, x^{(j)}) K_2(x^{(i)}, x^{(j)}) z_j \\ &= \sum_i \sum_j z_i \phi_1(x^{(i)})^T \phi_1(x^{(j)}) \phi_2(x^{(i)})^T \phi_2(x^{(j)}) z_j \\ &= \sum_i \sum_j z_i \sum_a \phi_{1a}(x^{(i)}) \phi_{1a}(x^{(j)}) \sum_b \phi_{2b}(x^{(i)}) \phi_{2b}(x^{(j)}) z_j \\ &= \sum_a \sum_b \sum_i \sum_j z_i \phi_{1a}(x^{(i)}) \phi_{1a}(x^{(j)}) \phi_{2b}(x^{(i)}) \phi_{2b}(x^{(j)}) z_j \\ &= \sum_a \sum_b \sum_i \left(z_i \phi_{1a}(x^{(i)}) \phi_{2b}(x^{(i)}) \right)^2 \geq 0 \end{aligned}$$

(f)

No, K is not a Mercer kernel.

$f: \mathbb{R}^n \mapsto \mathbb{R}$ is a real-valued function,

A counter example would be $f: y \mapsto \text{sign}(y)$, and choosing $(x, z) = (-1, 1)$

(g) It is straightforward to prove K is a Mercer kernel, since K_3 is a Mercer kernel. This is independent of the chosen map ϕ

Yes, $K_3(\phi(x), \phi(z))$ is a valid kernel, no matter what the inputs are.

(h)

Yes, $p(K_1)$ is a valid kernel.

$p(x)$ is a polynomial function with coefficients $c_k > 0, k = 0, 1, \dots, n$

$$p(x) = \sum_{k=0}^n c_k x^k$$

$$K(x, z) = p(K_1(x, z)) = \sum_{k=0}^n c_k \left(K_1(x, z) \right)^k$$

From (e) we know $K(x, z) = K_1(x, z)K_2(x, z)$ is a valid kernel, so $K(x, z) = \left(K_1(x, z)\right)^k$ is valid.

From (a) and (c), we know $K(x, z) = K_1(x, z) + K_2(x, z)$ and $K(x, z) = aK_1(x, z), a \in \mathbb{R}^+$ are both valid.

So $K(x, z) = \sum_{k=0}^n c_k \left(K_1(x, z)\right)^k$ is a valid kernel.

5.

(a) Recall the update rule

i.

$$\theta^{(i+1)} := \theta^{(i)} + \alpha \left(y^{(i+1)} - h_{\theta^{(i)}}(\phi(x^{(i+1)})) \right) \phi(x^{(i+1)})$$

We can easily figure out that $\theta^{(i)}$ is a linear combination of $\phi(x^{(1)}), \dots, \phi(x^{(i)})$, that is

$$\theta^{(i)} = \sum_{j=1}^i \beta_j \phi(x^{(j)})$$

ii.

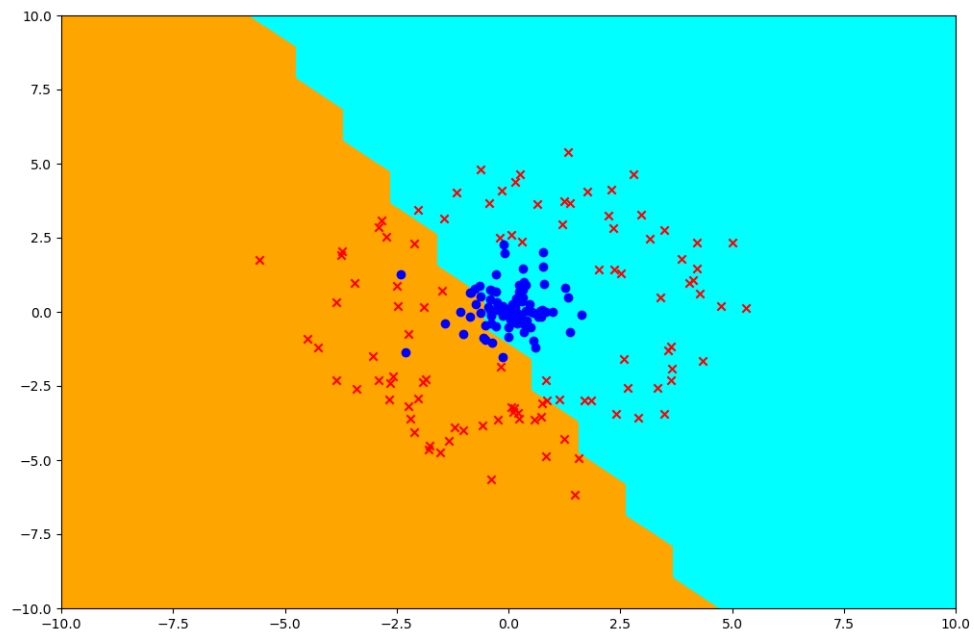
$$\theta^{(0)} = \vec{0}$$

$$\begin{aligned} h_{\theta^{(i)}}(\phi(x^{(i+1)})) &= g(\theta^{(i)T} \phi(x^{(i+1)})) \\ &= \text{sign}(\theta^{(i)T} \phi(x^{(i+1)})) \\ &= \text{sign}\left(\sum_{j=1}^i \beta_j \phi(x^{(j)})^T \phi(x^{(i+1)})\right) \\ &= \text{sign}\left(\sum_{j=1}^i \beta_j \langle \phi(x^{(j)}), \phi(x^{(i+1)}) \rangle\right) \\ &= \text{sign}\left(\sum_{j=1}^i \beta_j K(x^{(j)}, x^{(i+1)})\right) \end{aligned}$$

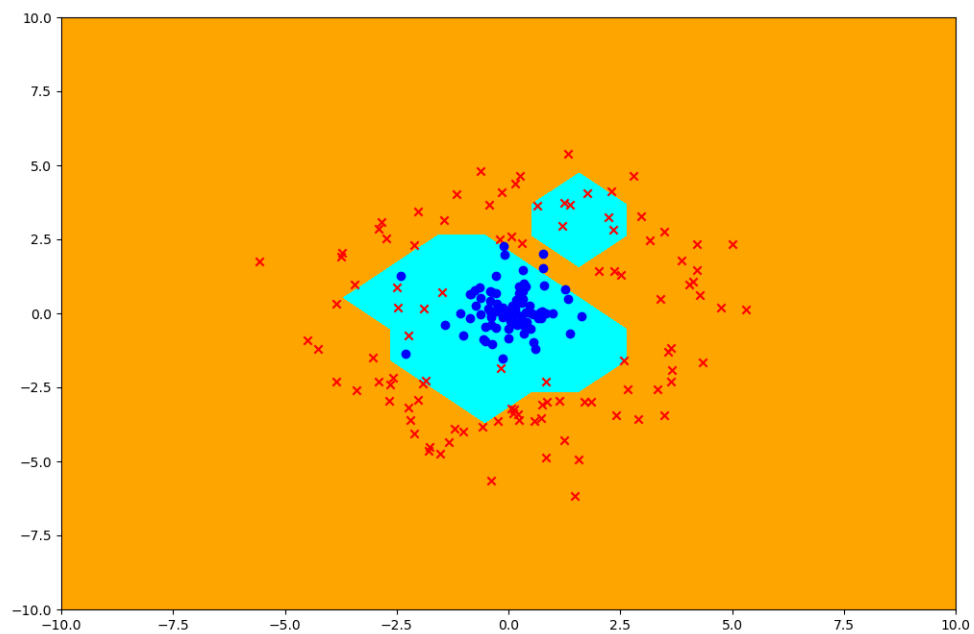
iii.

$$\begin{aligned} \theta^{(i+1)} &:= \theta^{(i)} + \alpha \left(y^{(i+1)} - h_{\theta^{(i)}}(\phi(x^{(i+1)})) \right) \phi(x^{(i+1)}) \\ &= \sum_{j=1}^i \beta_j \phi(x^{(j)}) + \alpha \left(y^{(i+1)} - \text{sign}\left(\sum_{j=1}^i \beta_j K(x^{(j)}, x^{(i+1)})\right) \right) \phi(x^{(i+1)}) \\ \beta_{i+1} &= \alpha \left(y^{(i+1)} - \text{sign}\left(\sum_{j=1}^i \beta_j K(x^{(j)}, x^{(i+1)})\right) \right) \end{aligned}$$

(c)



dot kernel



rbf kernel

Dot kernel performs poorly than rbf kernel. Because dot kernel doesn't do feature mapping $\phi(x) = x$, and the dataset is not linearly separable.

6.
(b)

$$\begin{aligned}\mathcal{L}(\phi_y, \phi_{k|y=1}, \phi_{k|y=0}) &= \sum_{i=1}^m \log p(x^{(i)}, y^{(i)}; \phi_y, \phi_{k|y=1}, \phi_{k|y=0}) \\ &= \sum_{i=1}^m \log p(x^{(i)} | y^{(i)}; \phi_{k|y=1}, \phi_{k|y=0}) p(y^{(i)}; \phi_y)\end{aligned}$$

where $x^{(i)} | y^{(i)=1}; \phi_{k|y=1} \sim \text{Mult}(n, \phi_{k|y=1})$, $\phi_{k|y=1}$ is the probability that word k appears given $y = 1$. Similar for $\phi_{k|y=0}$.

Therefore, to maximize \mathcal{L} w.r.t $\phi_{k|y=1}$ is equivalent to maximizing

$$\sum_{i=1}^m 1\{y^{(i)} = 1\} \sum_{k=1}^n x_k^{(i)} \log \phi_{k|y=1}$$

subject to the constraints $\phi_{k|y=1} \geq 0$ and $\sum_{k=1}^n \phi_{k|y=1} = 1$, where $x_k^{(i)}$ is the number of times that word k appears in the i -th message.

By introducing a Lagrange multiplier $\lambda \in \mathbb{R}$ corresponding to the second constraint, we have

$$\mathcal{L}(\lambda, \phi_{k|y=1}) = \sum_{i=1}^m 1\{y^{(i)} = 1\} \sum_{k=1}^n x_k^{(i)} \log \phi_{k|y=1} - \lambda \left(\sum_{k=1}^n \phi_{k|y=1} - 1 \right)$$

By solving the Lagrange multiplier problem and applying Laplace smoothing,

$$\phi_{k|y=1} := \frac{1 + \sum_{i=1}^m 1\{y^{(i)} = 1\} x_k^{(i)}}{n + \sum_{i=1}^m 1\{y^{(i)} = 1\} \sum_{j=1}^n x_j^{(i)}}$$

Similarly,

$$\phi_{k|y=0} := \frac{1 + \sum_{i=1}^m 1\{y^{(i)} = 0\} x_k^{(i)}}{n + \sum_{i=1}^m 1\{y^{(i)} = 0\} \sum_{j=1}^n x_j^{(i)}}$$

and ϕ_y is still $(1/m) \sum_{i=1}^m 1\{y^{(i)} = 1\}$.

To make a prediction, instead of computing the probabilities explicitly, we compare the logarithm of $p(y=1|x)/p(y=0|x)$ with 0.

$$\begin{aligned}\log \frac{p(y=1|x)}{p(y=0|x)} &= \log \frac{p(x|y=1)p(y=1)}{p(x|y=0)p(y=0)} \\ &= \log \frac{\left(\prod_{k=1}^n p(x_k|y=1) \right) p(y=1)}{\left(\prod_{k=1}^n p(x_k|y=0) \right) p(y=0)} \\ &= \sum_{k=1}^n \frac{\log \phi_{k|y=1}}{\log \phi_{k|y=0}} + \log \frac{\phi_y}{1-\phi_y} \\ &= \sum_{k=1}^n x_k (\log \phi_{k|y=1} - \log \phi_{k|y=0}) + \log \frac{\phi_y}{1-\phi_y}\end{aligned}$$

$p(y=1) = \phi_y$ $p(x_k | y=1) = \phi_{k|y=1}$
 $p(y=0) = 1 - \phi_y$ $p(x_k | y=0) = \phi_{k|y=0}$

(c). The top 5 indicative words for Naive Bayes are: ['claim', 'won', 'prize', 'tone', 'urgent!']

(d). The optimal SVM radius was 0.1.

The SVM model had an accuracy of 0.9695... on the testing set.