

Handwritten
hardcopy approved by
Dr. Yan Liu

Anjali Krishna Prasada
MORNING SESSION
USC id: 3400313602

Assignment 3

Q1.a)

$$\hat{\beta}_\lambda = \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n (y_i - x_i^T \beta)^2 + \lambda \|\beta\|_2^2 \right\}$$
$$= \arg \min_{\beta} \left\{ \frac{1}{n} \sum_{i=1}^n (x_i^T \beta - y_i)^2 + \lambda \|\beta\|_2^2 \right\}$$

Converting the above equation to vectorized form.

$$\hat{\beta}_\lambda = (\beta X - y)^T (\beta X - y) + \lambda \beta^T \beta$$

Finding $\nabla \hat{\beta}_\lambda = 0$

$$\nabla_\beta = X^T (\beta X - y) + \lambda \beta = 0$$

$$X^T \beta X - X^T y + \lambda \beta = 0$$

$$\beta (X^T X + \lambda I) = X^T y$$

$$\beta = (X^T X + \lambda I)^{-1} X^T y$$

closed form solution.

$$\beta = [X^T X + \lambda I]^{-1} X^T y$$

we know $y = X^T \beta^* + \epsilon$

$$\hat{\beta}_\lambda = [X^T X + \lambda I]^{-1} X^T (X^T \beta^* + \epsilon)$$

$$= \underset{P}{\underbrace{[X^T X + \lambda I]}^{-1}} \underset{B}{\underbrace{[X X^T B^k + X^T E]}} \underset{A}{\uparrow}$$

$$\text{where } P = [X^T X + \lambda I]^{-1}$$

$$A = X^T$$

$$B = X X^T B^*$$

$$= P(B + A(t))$$

Using ~~background~~ Affine transformation :-
where $X \sim N(0, 1)$

$$Ax + B \in N(Au + B, A^T \circ A)$$

where ϵ is taken from $N(0_g)$

we get

$$N(x x^T B^*, x x^T)$$

AS

~~case~~ $\mu=0$ & $\sigma=1$

$$\Rightarrow P \left[N(XX^T B^*, X\hat{X}^T) \right]$$

substituting P for ~~pos~~ its value.

$$[X^T X + \lambda I]^{-1} [N(X X^T \beta^*, X X^T)]$$

we again use affine transformation to get:

$$\beta \sim N \left[\left(X^T X + \lambda I \right)^{-1} X X^T \beta^*, \sigma^2 I \right]$$

$$\left(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} \right)^{-1} \mathbf{X} \mathbf{X}^T \left(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I} \right)^{-1}$$

b). In the previous solution we got

$$\hat{B} \sim N \left[\left(X^T X + \lambda I \right)^{-1} X X^T \beta^*, \left(X^T X + \lambda I \right)^{-1} X X^T \left(X^T X + \lambda I \right)^{-1} \right]$$

here mean is equivalent to ~~not~~ bias

∴ bias term is

$$\left(X^T X + \lambda I \right)^{-1} X X^T \beta^*$$

c) In a part of this answer we got

$$\hat{B} \sim N \left[\left(X^T X + \lambda I \right)^{-1} X X^T \beta^*, \left(X^T X + \lambda I \right)^{-1} X X^T \left(X^T X + \lambda I \right)^{-1} \right]$$

Here 2nd term is variance.

$$\therefore \text{variance term} = \left(X^T X + \lambda I \right)^{-1} X X^T \left(X^T X + \lambda I \right)^{-1}$$

d) without regularization, we can get variance very high. When we use regularization & keep λ small, there is a slight increase in bias but the variance comes down significantly. If we keep increasing λ , bias increases significantly.

We can see in variance equation

$$\left(X^T X + \lambda I \right)^{-1} X X^T \left(X^T X + \lambda I \right)^{-1}$$

As λ increase variance decreases.

$\therefore \lambda \uparrow \quad \cancel{\text{variance}} \uparrow$
 $\lambda \downarrow \quad \text{bias} \uparrow$

Q2a) To prove $k_3(x, x') = a_1 k_1(x, x') + a_2 k_2(x, x')$

where $a_1, a_2 \geq 0$.

Let $\langle x, x' \rangle$ be euclidean inner product.

$k_1(x, x')$ is a kernel therefore the semi-definite i.e.

$$\langle v, k_1(x, x') v \rangle > 0$$

a_1 is a scalar

$$\langle v, a_1 k_1(x, x') v \rangle = a_1 \underbrace{\langle v, k_1(x, x') v \rangle}_{\text{This is the semidefinite.}}$$

$$\langle v, a_2 k_2(x, x') v \rangle = a_2 \underbrace{\langle v, k_2(x, x') v \rangle}_{\text{This is also the semidefinite.}}$$

Combining the two.

$$\langle v, \{a_1 k_1(x, x') + a_2 k_2(x, x')\} v \rangle$$

By property of inner product.

$$a_1 \langle v, k_1(x, x') v \rangle + a_2 \langle v, k_2(x, x') v \rangle > 0$$

For all non zero vectors x

since $a_1 \langle v, k_1(x, x') v \rangle > 0$ and

$$a_2 \langle v, k_2(x, x') v \rangle > 0$$

Thus by definition of matrix.

$a_1 k_1(x, x') + a_2 k_2(x, x')$ is positive semidefinite.

Proving symmetric nature.

~~Since~~ $a_1 k_1(x, x')$ and $a_2 k_2(x, x')$ are symmetric as

$$a A = a[A] = a[A^T]$$

We know $A = A^T$ when symmetric & k_1 & k_2 are kernels so they are symmetric.

To prove: $a_1 k_1(x, x') + a_2 k_2(x, x')$ is symmetric.

$$\text{Let } a_1 k_1(x, x') = A \\ a_2 k_2(x, x') = B.$$

$$(A+B)^T = A^T + B^T$$

$$\text{But } A^T = A \quad \& \quad B^T = B$$

$$\therefore (A+B)^T = A+B$$

Hence $a_1 k_1(x, x') + a_2 k_2(x, x')$ is symmetric

$\therefore k_3(x, x')$ is also a kernel.

2 b) $k_4(x, x') = f(x) f(x')$ To prove

$$k_4(x, x') = \langle \phi(x), \phi(x') \rangle$$

where $\phi: x \rightarrow f(x)$

$$\phi: x' \rightarrow f(x').$$

$$k_4(x, x') = \langle f(x), f(x') \rangle$$

$\therefore f(x) f(x')$ is symmetric and positive definite.

2c) To prove $k_s(x, x') = k_1(x, x') k_2(x, x')$

$$\phi(x)_{ij} = \phi_1(x)_i \cdot \phi_2(x)_j.$$

i.e. the i, j component of the new feature vector is the product of i^{th} of $\phi_1(x)$ & j^{th} of $\phi_2(x)$.

$$k_s(x, x') = \langle \phi(x), \phi(x') \rangle$$

$$= \sum_{i,j} \phi(x)_{ij} \phi(x')_{ij}$$

$$= \sum_{i,j} \phi_1(x)_i \phi_2(x)_j \cdot \phi_1(x')_i \phi_2(x')_j$$

$$= \left(\sum_i \phi_1(x)_i \phi_2(x'_i) \right) \left(\sum_j \phi_1(x_j) \phi_2(x'_j) \right)$$

$$= k_1(x, x') k_2(x, x')$$

Hence proved.

$$3a) \min_{\omega} \sum_i (y_i - \omega^T x_i)^2 + \lambda \|\omega\|_2^2$$

Converting the above equation in vector form, we get.

$$(x\omega - y)^T (x\omega - y) + \lambda \omega^T \omega.$$

Taking gradient wrt ω . $\nabla_{\omega} = 0$

$$\nabla_{\omega} = x^T (x\omega - y) + \lambda \omega = 0.$$

$$x^T x \omega - x^T y + \lambda \omega = 0$$

$$\omega (x^T x + \lambda I_D) = x^T y.$$

$$\omega = [x^T x + \lambda I_D]^{-1} x^T y.$$

Therefore we can say

Optimal ω can be written as.

$$[x^T x + \lambda I_D]^{-1} x^T y.$$

3 b) Applying non linear feature mapping to each sample. $x_i \rightarrow \phi_i = \phi(x) \in R^T$
 s.t. $\phi \in R^{N \times T}$

then

$$\sum_n [w^T \phi(x_n) - y_n]^2 + \lambda \|w\|_2^2$$

changing this to vector form

$$\phi^T [\phi w - y] [\phi w - y]^T + \lambda w^T w$$

taking gradient wrt w. $\nabla_w = 0$

$$\nabla_w = \phi^T [\phi w - y] + \lambda w = 0$$

$$\phi^T [\phi w] - \phi^T y + \lambda w = 0$$

$$w^T [\phi^T \phi + \lambda I_t] = \phi^T y$$

$$w^* = [\phi^T \phi + \lambda I_t]^{-1} \phi^T y \quad -①$$

we have to prove,

$$(\phi^T \phi + \lambda I_t)^{-1} \phi^T y = \phi^T (\phi \phi^T + \lambda I_N)^{-1} y$$

lets use the following property.

$$(P^{-1} + B^T R^{-1} B)^{-1} B^T R^{-1} = P B^T (B P B^T + R)^{-1}$$

Let $P = I$

$$B = \Phi$$

$$R = \lambda I$$

lets keep aside y for now

$$(P^{-1} + B^T R^{-1} B) B^T R^{-1}$$

$$= (I + \frac{\phi^T I \phi}{\lambda})^{-1} \frac{\phi^T I}{\lambda}$$

$$\text{as } A^{-1} B^{-1} = (BA)^{-1}$$

$$\therefore \left(\frac{\phi^T I}{\lambda} \right)^{-1} = \lambda I \phi^{-T}$$

$$\Rightarrow \left[\lambda I \phi^{-T} I + \lambda I \phi^{-T} \frac{\phi^T I}{\lambda} \phi \right]^{-1}$$

= ~~$\lambda I \phi^{-T} I$~~ *

$$= [\phi^{-T} \lambda I + \phi^{-T} \phi^T \phi]^{-1}$$

$$= [\lambda I + \phi^T \phi]^{-1} \phi^T$$

Multiplying the above eqⁿ by y we get ①

$$[\lambda I + \phi^T \phi]^{-1} \phi^T y \Rightarrow [\phi^T \phi + \lambda I] \phi^T y$$

This means we can use the matrix identity to get the result.

[using P , B & R values]

$$\begin{aligned}
 & PB^T (B P B^T + R)^{-1} \\
 &= I \phi^T (\phi I \phi^T + \lambda I)^{-1} \\
 &= \phi^T (\phi \phi^T + \lambda I)^{-1}
 \end{aligned}$$

Multiplying the above equation by y
we get:

$$= \phi^T (\phi \phi^T + \lambda I)^{-1} y.$$

Hence proved.

3c) we know

$$\omega^* = \phi^T (\phi \phi^T + \lambda I_N)^{-1} y.$$

given To prove:

$$\hat{y} = y^T (K + \lambda I_N)^{-1} k(x)$$

Given:

$$\begin{aligned}
 \hat{y} &= \omega^{*T} \phi(x) \\
 &= y^T (\phi \phi^T + \lambda I_N) \cancel{\phi(x)} \phi^T \\
 &= y^T (K + \lambda I_N)^{-1} k(x)
 \end{aligned}$$

3d) Linear Ridge regression has a complexity
of ~~order~~ $O(CD^3)$

Kernel Ridge regression has the complexity
of ~~order~~ $O(N^3)$
~~same same~~

linear ridge regression: $(X^T X + \lambda I_D)^{-1} X^T y$.

$X^T X$ is a ~~DxN~~ $D \times N$ and $N \times D$ matrix

$(X^T X + \lambda I_D)^{-1}$ is a $D \times D$ inversion.

$\therefore O(D^3)$ is complexity.

* Kernel Ridge regression:

$$\phi^T (\phi \phi^T + \lambda I_N)^{-1} y$$

$\phi \phi^T$ is a ~~TxN~~ $N \times T$ and $T \times N$ matrix

$(\phi \phi^T + \lambda I_N)^{-1}$ is a $N \times N$ matrix inversion

- $O(N^3)$ is complexity of kernel ridge regression.

Linear ridge regression complexity depends on dimension of the features & kernel ridge regression depends on no of examples x.

Q4a) Positive examples are not linearly separable from the negative examples in the original space as there is no line present that can divide the positive and negative examples from each other.

b). We have 2 positive points $(1, 1)$, $(-1, -1)$ and we have 2 negative points $(-1, 1)$, $(1, -1)$

Our transformation function is

$\phi(x) = [1, x_1, x_2]$. Applying this on all 4 points we get.

Positive points : $(1, 1)$ transforms to

$$[1, 1, 1, 1],$$

$$(-1, 1) \rightarrow [1, -1, -1, 1]$$

negative points :

$$(1, -1) \rightarrow [1, 1, -1, -1]$$

$$(-1, 1) \rightarrow [1, -1, 1, -1]$$

we can clearly see that only terms that decides that if the example is positive or negative is the term $x_1 x_2$

Therefore we can set w coefficients $[0, 0, 0, 1]$

4c) If we have a positive example $(2, -2)$, the total five examples can no longer be linearly separated in the feature space $\phi(x)$. As multiplying Φ the transformation we get.

$$[1, 2, -2, -4]$$

with the coefficients we got in the previous answer for w .

We get -4 . which classifies the example to be negative. But it is actually positive.

4d) $K(x, x) = \phi(x)^\top \phi(x)^T$

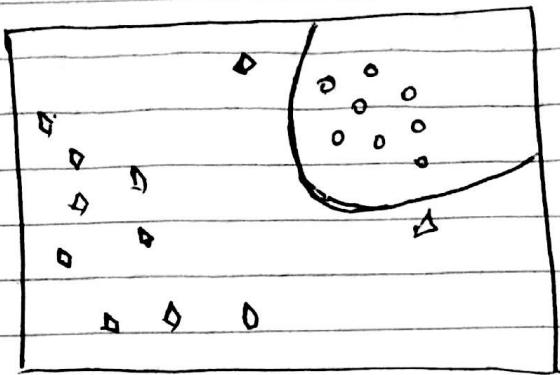
$$\phi(x) = [1, x_1, x_2, x_1 x_2]$$

$$\phi(x)^T = [1, x_1, x_2, x_1 x_2]$$

$$K(x, x) = 1 + x_1^2 + x_2^2 + x_1^2 x_2^2$$

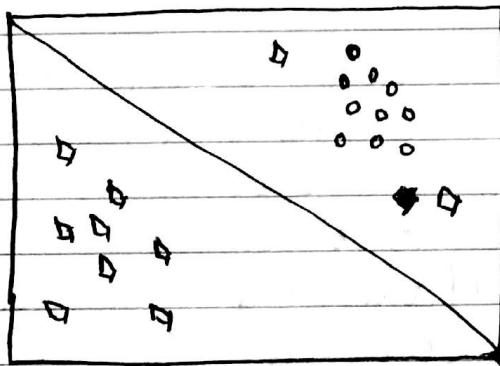
$K(x, x)$ is a polynomial kernel.

Q5a)



When $C \rightarrow \infty$, we get hard margin in SVM which results in overfitting. Therefore I drew a parabola to separate the 2 classes without misclassifying any points.

b)

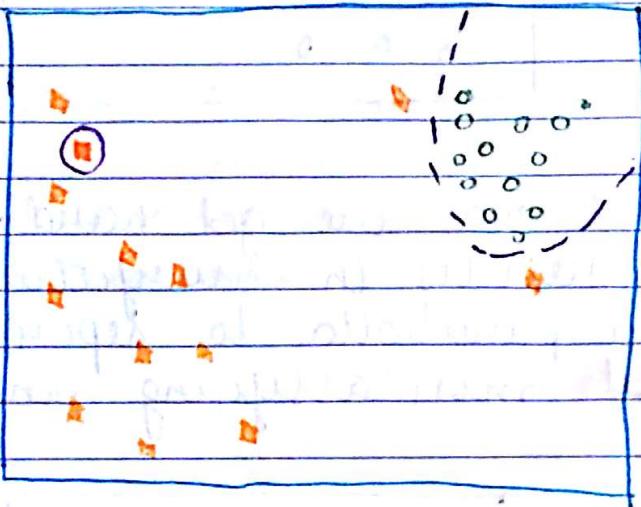


For ~~$C \rightarrow \infty$~~ $C \approx 0$, we get soft margin that means misclassification is allowed in this case. For a very tiny value of C , we should get misclassified examples.

- c) I expect the second one ($C \approx 0$) to work better than the 1st one ($C \rightarrow \infty$) as it is more generic. We can clearly see overfitting in the first ~~case~~ case.

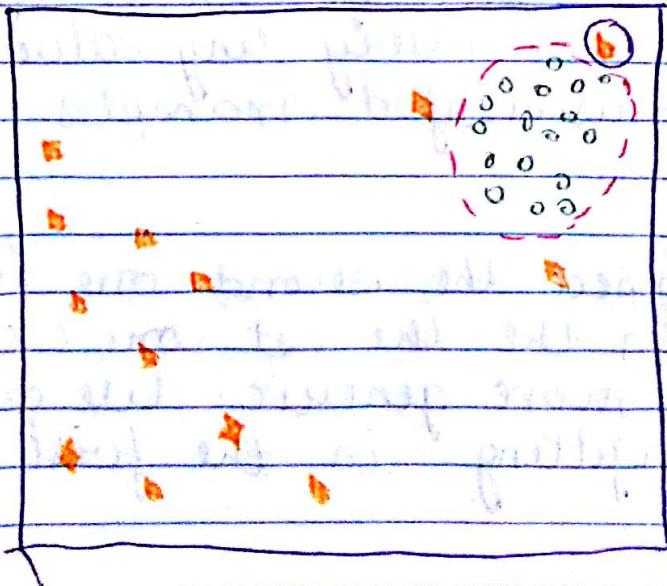
In the second case, only 2 points are misclassified but the model looks more generic.

d)



The point which I have circled, is the new point I have added which will not change the decision boundary when $C \rightarrow \infty$. As when $C \rightarrow \infty$, we get Hard margin. Hard Margin restricts misclassification. and according to the hard margin in this case, the new point will not be misclassified.

e)



The circled point is the new point I have added which changes the decision boundary for very large $C \rightarrow \infty$. Hence, it is a hard margin that does not allow any misclassification. Therefore according to the old decision boundary, the new point I added will be misclassified. Therefore the decision boundary will change as shown in the figure.

Programming question:

Extension by 1 week for programming question approved by Dr. Yashas on medical grounds.