

Homework 9

Submitted By: Puneet Singhal

1) Images at end. Sorry for formatting

2)

- (a) True
- (b) False
- (c) True
- (d) A
- (e) A, C

(f) Kernel SVM will be lower training error and lower test error. This is because we can choose kernels corresponding to transformations that best describe our model. This is given the Neural network in this question. If we start increasing layers in neural network and we can allow more non-linearities and hence can get better results with neural networks

3)

- (a)
 $x_1^2 + x_2^2$
- (b) images on new page. Sorry for formatting

4)

- (a)

The appropriate loss function is Hinge loss as it provides tight (convex) upper bound on 0-1 loss function. While this function is convex and continuous, it is not differentiable at 1 but this can be solved using subgradient descent methods.

- (b)

Mean squared error loss is appropriate as it square, absolute, continuous and differential function which will help us use techniques like gradient descent or stochastic gradient descent. This function will represent the actual cost of wrong predictions.

- (c)

- (a) SVM with no kernel (Soft margin): This will allow the error in classification but we will still not be able to classify 100% samples correctly.
- (b) SVM with quadratic kernel with penalty on slack variable $C=0$: The quadratic kernel will be able to classify each sample correctly as we can see that it is possible to draw an ellipse that divides the two samples. C does not matter as there are no green sample that is surrounded by red and vice versa.
- (c) SVM with quadratic kernel with penalty on slack variable $C=\infty$: The quadratic kernel will be able to classify each sample correctly as we can see that it is possible to draw an ellipse that divides the two samples. C does not matter as there are no green sample that is surrounded by red and vice versa.
- (d) Logistic regression (No kernel): Logistic regression with no kernel will have linear surface as separator which will not be able to classify all samples correctly.

Q.1)

$$y = w_0 + w_1 x + \epsilon \quad \text{where } \epsilon \sim N(0, \sigma^2)$$

$$1.) \quad P(w | y^t, x^t) = \frac{P(y^t | x^t, w) P(w_0) P(w_1)}{P(y^t | x^t)}$$

$$P(w_i) = \frac{1}{\sqrt{2\pi\lambda}} \exp\left(-\frac{w_i^2}{2\lambda}\right)$$

$$P(y^t | x^t, w) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y^t - w_0 - w_1 x^t)^2}{2\sigma^2}\right\}$$

$$P(y^t | x^t) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y^t)^2}{2\sigma^2}\right\}$$

$$\Rightarrow P(w | y^t, x^t) = \frac{1}{\sqrt{2\pi\lambda}} \exp\left\{-\frac{(y^t - w_0 - w_1 x^t)^2 + (y^t)^2 - w_0^2 - w_1^2}{2\sigma^2} - \frac{w_0^2 - w_1^2}{2\lambda}\right\}$$

$$2.) \quad P(w | y, x) = \frac{P(y | x, w) P(w_0) P(w_1)}{P(y | x)}$$

$$= \frac{1}{\sqrt{2\pi\lambda}} \exp\left\{-\frac{(y - x^T w)^T (y - x^T w)}{2\sigma^2} + \frac{y^T y}{2\sigma^2} - \frac{w_0^2}{2\lambda} - \frac{w_1^2}{2\lambda}\right\}$$

Here, $X \rightarrow \begin{bmatrix} 1 & x^0 \\ 1 & x^1 \\ \vdots & \vdots \\ 1 & x^N \end{bmatrix}, Y \rightarrow \begin{bmatrix} y^0 \\ y^1 \\ \vdots \\ y^N \end{bmatrix}, w = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}$

5) Collaboration Policy

- (a) No
- (b) No
- (c) No

$$3.) \quad -\frac{(Y-Xw)^T(Y-Xw)}{2\sigma^2} + \frac{Y^TY}{2\sigma^2} - \frac{w^Tw}{2\lambda}$$

$$4.) \quad d \frac{\log(P(w|y,x))}{dw} = 0$$

$$\Rightarrow -\frac{(Y-Xw)^TX}{\sigma^2} - \frac{w}{\lambda} = 0$$

$$5.) \quad \text{Gradient} = \frac{d \log(P(w|y,x))}{dw} = -\frac{(Y-Xw)^TX}{\sigma^2} - \frac{w}{\lambda}$$

$$w_{\text{new}} = w_{\text{old}} + \alpha \left(-\frac{(Y-X^Tw_{\text{old}})^TX}{\sigma^2} - \frac{w_{\text{old}}}{\lambda} \right)$$

$$6.) P(w|y, x) \propto \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left\{-\frac{(y-xw)^T(y-xw)}{2\sigma^2}\right\} \frac{1}{2\lambda} \exp\left\{-\frac{|w_0|}{\lambda}\right\}$$

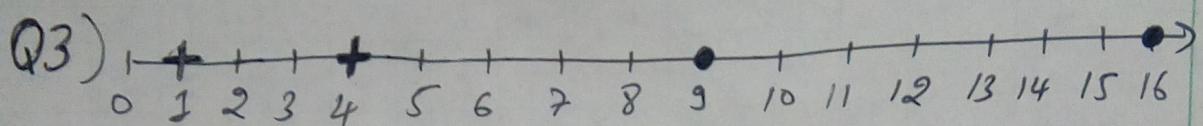
$$\times \frac{1}{2\lambda} \exp\left\{-\frac{|w_1|}{\lambda}\right\}$$

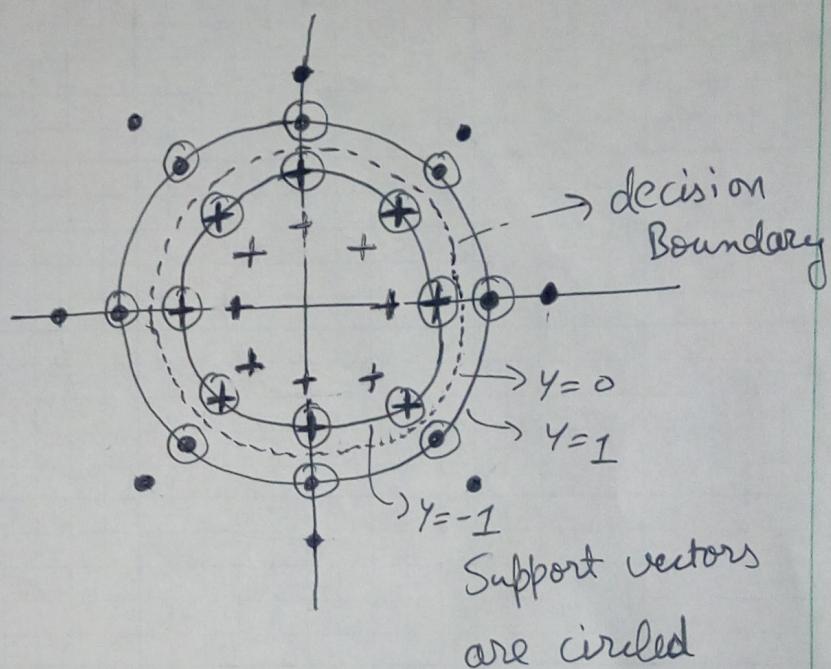
$$\log P(w|y, x) = \text{constant} - \frac{(y-xw)^T(y-xw)}{2\sigma^2} - \frac{|w_0|}{\lambda} - \frac{|w_1|}{\lambda}$$

$$\text{MAP: } \frac{P(w|y, x)}{dw} = 0$$

$$\frac{d \log P(w|y, x)}{dw_0} = 0$$

$$\frac{d \log(P(w|y, x))}{dw_1} = 0$$

Q3) 



3.) VC dimension $\rightarrow 2$