

ADITYA GAUTAM.

agautam1@andrew

MACHINE LEARNING

ASSIGNMENT #3

18th Feb 2016

Ques 1) Kernel Feature Mapping.

1) $x = (x_1, x_2)^T$

$$\phi(x) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)^T$$

$$K(x, z) = \phi(x) \cdot \phi(z)$$

$$= (x_1^2, \sqrt{2}x_1x_2, x_2^2) \cdot (z_1^2, \sqrt{2}z_1z_2, z_2^2)$$

$$= (x_1z_1 + x_2z_2)^2 = (x \cdot z)^2$$

Corresponding kernel function = $(x_1z_1 + x_2z_2)^2$

2) $K(x, z) \quad x, z \in R^2$

a) Mapping into feature space.

$$(x_1, x_2) \rightarrow \phi(x) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$$

Total no. of multiplication needed for above feature mapping = 3

Similarly, for z mapping, we would need 3 multiplications

To multiply in feature space we would need 3 multiplication (dot product)

$$\text{Thus, total Multiplication} = 3+3+3 = 9 \text{ Multiplication}$$

Total Addition = 2

b) Computing through kernel function.

$$(x_1z_1 + x_2z_2)^2 \rightarrow 2f1 \rightarrow 3 \text{ Multiplication}$$

Addition = 1

2) Perceptrons

$$y = \phi(\sum_i w_i x_i + b)$$

$$\phi(z) = \begin{cases} 0 & \text{if } z \leq 0 \\ 1 & \text{otherwise} \end{cases}$$

(2) AND

Case 1 $x_1=0, x_2=0 \Rightarrow y=0 \Rightarrow \phi(\sum_i w_i x_i + b) \leq 0$

$$\Rightarrow w_1 x_1 + w_2 x_2 + b \leq 0$$

$$\Rightarrow b \leq 0 \quad \text{as } x_1=x_2=0 \quad \textcircled{1}$$

Case 2 $x_1=0, x_2=1 \text{ AND } \rightarrow 0$

$$\Rightarrow 0 \cdot w_1 + 1 \cdot w_2 + b \leq 0$$

$$w_2 + b \leq 0 \quad \textcircled{2}$$

Case 3 $x_1=1, x_2=0 \text{ AND } \rightarrow 0$

$$\Rightarrow 1 \cdot w_1 + 0 \cdot w_2 + b \leq 0$$

$$w_1 + b \leq 0 \quad \textcircled{3}$$

Case 4 $x_1=1, x_2=1 \text{ AND } \rightarrow 1$

$$\Rightarrow 1 \cdot w_1 + 1 \cdot w_2 + b \geq 0$$

$$w_1 + w_2 + b \geq 0 \quad \textcircled{4}$$

Solving above 4 inequalities we get,
 $b \leq 0, w_1 \leq b, w_2 \leq -b, w_1 + w_2 > -b$

$$\boxed{b=-1, w_1=2/3, w_2=-2/3}$$

(3)

AND table

| x_1 | x_2 | AND |
|-------|-------|-----|
| 0 | 0 | 0 |
| 0 | 1 | 0 |
| 1 | 0 | 0 |
| 1 | 1 | 1 |

3) OR Contingency table

| x_1 | x_2 | OR |
|-------|-------|----|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

| x_1 | x_2 | y |
|-------|-------|-----|
| 0 | 0 | 0 |
| 0 | 1 | 1 |
| 1 | 0 | 1 |
| 1 | 1 | 1 |

4) Case 1 $x_1=0, x_2=0, y=0$
 $y=0 \Rightarrow w_0 + w_1 x_1 + w_2 x_2 + b \leq 0$
 $\Rightarrow b \leq 0 \quad \text{--- } \textcircled{1}$

Case 2 $x_1=0, x_2=1, y=1$
 $w_0 + w_2 x_2 + b \geq 0$
 $\Rightarrow b + w_2 \geq 0 \quad \text{--- } \textcircled{2}$

Case 3 $x_1=1, x_2=0, y=1$
 $w_0 + w_1 + b \geq 0$
 $w_0 + b \geq 0 \quad \text{--- } \textcircled{3}$

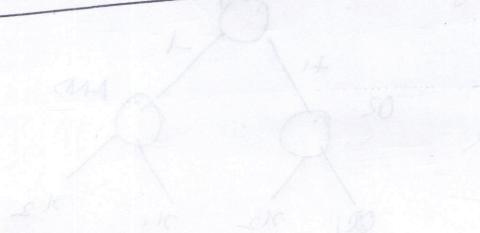
Case 4 $x_1=1, x_2=1, y=1$
 $w_0 + w_1 + w_2 + b \geq 0 \quad \text{--- } \textcircled{4}$

Solving above 4 inequalities, we get

a range of values.

one set of values

$$w_0 = 2, w_1 = 2, w_2 = 2, b = -1$$



5)

| x_1 | x_2 | XOR |
|-------|-------|-----|
| 0 | 0 | 0 |
| 1 | 0 | 1 |
| 0 | 1 | 1 |
| 1 | 1 | 0 |

XOR Table

6) Case 1) $x_1=0, x_2=0, \text{ XOR}=0$
 $y=0 \Rightarrow \nexists (\exists x_1 w_1 + b \leq 0)$
 $\Rightarrow x_1 w_1 + x_2 w_2 + b \leq 0$
 $b \geq 0 \quad \textcircled{1}$

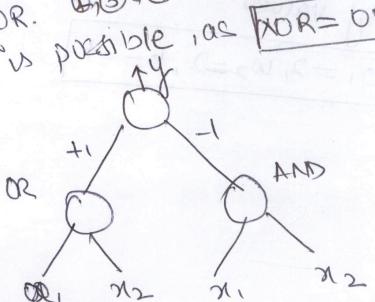
Case 2) $x_1=0, x_2=1, \text{ XOR}=\phi$
 $y=1 \Rightarrow \exists x_1 w_1 + b > 0$
 $\Rightarrow 0 \cdot w_2 + b > 0 \quad \textcircled{2}$

Case 3) $x_1=1, x_2=0, \text{ XOR}=1$
 $\Rightarrow 1 \cdot w_1 + 0 \cdot w_2 + b > 0$
 $w_1 + b > 0 \quad \textcircled{3}$

Case 4) $x_1=1, x_2=1, \text{ XOR}=0$
 $\Rightarrow 1 \cdot w_1 + 1 \cdot w_2 + b \leq 0$
 $w_1 + w_2 + b \leq 0 \quad \textcircled{4}$

Above 4 equations are not possible to solve
thus single layer perceptron is not possible.
for logical XOR. $\textcircled{1}, \textcircled{3}, \textcircled{4}$ contradicts $\textcircled{1}$.

However, two layer is possible, as $\text{XOR} = \text{OR-AND}$



3) Regression Theory

3.1 Linear Regression

a) $\{x_i, y_i\}, \dots, \{x_n, y_n\}$, $x \in \mathbb{R}^M$.

$$f(x) = w^T x$$

a) square loss error
 $\text{error} = \text{misclassification} \Rightarrow f(x_i) \neq y_i \text{ for } x_i$

$$J(w) = (y_i - f(x)) \text{ for } i^{\text{th}} \text{ sample}$$

loss func for all samples,

$$J(w) = \frac{1}{n} \sum_{i=1}^n (y_i - f(x))^2$$

$$J(w) = \frac{1}{n} \sum_{i=1}^n (y_i - w^T x_i)^2$$

b) Partial derivation of $J(w)$ w.r.t w^K

$$\frac{\partial J(w)}{\partial w^K} = \frac{1}{n} \frac{\partial (\sum (y_i - w^T x_i)^2)}{\partial w^K} = \frac{2}{n} (y_i - w^T x_i) x^K - \cancel{x^K}$$

$$= -\frac{2}{n} (y_i - w^T x_i) x^K$$

$$\Rightarrow \sum_{i=1}^n (y_i - w^T x_i) x^K$$

c) gradient descent

$$w^K_{\text{new}} = w^K + \alpha \frac{\partial J(w)}{\partial w^K}$$

old params step size

$$\frac{\partial J(w)}{\partial w^K} = \sum_{i=1}^n x^K (y_i - w^T x_i) \quad \text{--- ①}$$

$$w^K_{\text{new}} = w^K + \alpha \sum_{i=1}^n x^K (y_i - w^T x_i) \quad \text{from ①}$$

do this till convergence.

$$2) f(x) = w^T x$$

a) conditional likelihood L

$$\ell(w; x, y) := \prod_{i=1}^n P(y_i | x_i, w)$$

$$L(w; x, y) = \operatorname{argmax}_w \prod_{i=1}^n P\left(\frac{y_i}{x_i, w}\right)$$

$$P(y_i | x_i, w) = N(w^T x_i, \sigma^2)$$

$$P(y_i | x_i, w) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left(\frac{y_i - f(x_i; w)}{\sigma} \right)^2}$$

$$\Rightarrow L(w; x, y) = \operatorname{argmax}_w \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left(\frac{y_i - f(x_i; w)}{\sigma} \right)^2}$$

b) log conditional likelihood

$$\log(L(w; x, y)) = \operatorname{argmax}_w \left(\log \prod_{i=1}^n \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left(\frac{y_i - f(x_i; w)}{\sigma} \right)^2} \right)$$

$$\text{as } \log(x \cdot y) = \log x + \log y$$

$$\log(L(w; x, y)) = \sum_i \log \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2} \left(\frac{y_i - f(x_i; w)}{\sigma} \right)^2}$$

$$\begin{aligned} \log(L(w; x, y)) &= \sum_i \log \frac{1}{\sqrt{2\pi\sigma^2}} + \sum_i \log e^{-\frac{1}{2} \left(\frac{y_i - f(x_i; w)}{\sigma} \right)^2} \\ &= C - \frac{1}{2} \sum_i \left(\frac{y_i - f(x_i; w)}{\sigma} \right)^2 \end{aligned}$$

$$\log(L(w; x, y)) = C - \frac{1}{2} \sum_i \left(\frac{y_i - f(x_i; w)}{\sigma} \right)^2$$

$$C = \sum_i \log \frac{1}{\sqrt{2\pi\sigma^2}}$$

c) Maximizing the log likelihood
 \Rightarrow taking a derivative and putting it to zero

$$\frac{\partial \log L(\omega; x_i, y)}{\partial \omega} = \frac{\partial}{\partial \omega} \left(C - \frac{1}{2} \sum_{j=1}^n (y_j - f(x_j, \omega))^2 \right)$$

$$= 0 - \frac{1}{2C^2} \frac{\partial}{\partial \omega} \sum_{j=1}^n (y_j - f(x_j, \omega))^2$$

$$\frac{\partial \log L(\omega; x_i, y)}{\partial \omega} = -\frac{1}{2C^2} \sum_{j=1}^n \frac{\partial}{\partial \omega} (y_j - \omega^T x_j)^2$$

$$= -\frac{1}{2C^2} \sum_{j=1}^n (y_j - \omega^T x_j) x_j$$

$$\boxed{\frac{\partial \log L(\omega; x_i, y)}{\partial \omega} = \frac{1}{2C^2} \sum_{j=1}^n (y_j - \omega^T x_j) x_j}$$

Above function will give the derivative to zero at a point same as that of loss function proved earlier so, maximizing log likelihood is same as that of minimizing the least square error.

3) $y = f(x) + \epsilon$
 \downarrow
 Noise
 defined by us

ϵ : mean = 0
 variance σ^2

a) $E_D \left[\sum_{y_n} (h(x) - f(x))^2 p(y|x) p(x) dy dx \right]$
 point x .

b) $E_D [(y - h(x))^2]$

$l = E_D(h(x)) \quad l = E_D(h(x))$

$$E[(y - l)^2] = E[(y - l + l - \hat{y})^2]$$

$$= E[(y - l)^2 + (l - \hat{y})^2 + 2[y - l][l - \hat{y}]]$$

$$= E[(y - l)^2 + (l - \hat{y})^2 + 2[yl - l^2 + l\hat{y} - \hat{y}^2]]$$

$$= E[(y - l)^2] + E[(l - \hat{y})^2] + 2(E[yl] - E[l^2] - E[l\hat{y}])$$

$$\downarrow \quad \downarrow \quad \downarrow$$

Bias

Variance

+ σ^2

\downarrow
 Unavoidable error

3.2) Regularization.

$$L = \frac{1}{2} \sum_{i=1}^N (y_i - w^T x_i)^2 \quad (\text{original loss function})$$

$$\text{a) } L = \frac{1}{2} \sum_{i=1}^N (y_i - w^T x_i)^2 + \frac{\lambda}{2} \|w\|^2 \quad (\text{loss function with regularization})$$

$$\frac{\partial L}{\partial w_k} = \frac{2}{2} \sum_{i=1}^N (y_i - w^T x_i) x_{ik} + \frac{2\lambda}{2} w_k$$

$$\boxed{\frac{\partial L}{\partial w_k} = \sum_{i=1}^N (y_i - w^T x_i) x_{ik} + \lambda w_k}$$

b) First Algorithm is more likely to give a sparse matrix as it penalised the weight with high weightage of w_k .

$$w^{(k+1)} = w^{(k,t)} - \sum_{i=1}^n (y_i - w^T x_i) x_{ik} - \lambda w^{(k,t)}$$

In second algo, it will reduce by A. however in other it will be by the magnitude of $w^{(k+1)}$

so, First Algorithm is likely to give sparse matrix with more w_i s set to zero

4.1) logistic regression

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

1) logistic function.

$$f(x_i; w) = \frac{1}{1 + e^{-a}}, \text{ where } a = \sum_{j=1}^N w_j x_j, N \rightarrow \text{no. of feature.}$$

2) conditional likelihood

$$P(y=0|w) = \frac{1}{1 + e^{(w_0 + \sum w_i x_i)}}$$

$$P(y=1|w) = 1 - P(y=0|w) = \frac{e^{(w_0 + \sum w_i x_i)}}{1 + e^{(w_0 + \sum w_i x_i)}}$$

$$\text{Conditional likelihood} = \prod_{\text{Over all the samples}} P(y^i|x^i, w)$$

$$L(w; x, y) = \prod_i P(y^i|x^i, w)$$

$$L(w; x, y) = \underset{w}{\operatorname{argmax}} \prod_i P(y^i|w, x^i)$$

3) ~~Final~~ log conditional likelihood

$$\log L(w; x, y) = \log \prod_i P(y^i|x^i, w) \quad (\log ab = \log a + \log b)$$

$$= \sum_i \log P(y^i|x^i, w)$$

$$\log(L) = \sum_i y^i \log P(y=1|x^i, w) + (1-y^i) \log P(y=0|x^i, w)$$

$$= \sum_i y^i \log \frac{P(y=1|x^i, w)}{P(y=0|x^i, w)} + \log P(y=0|x^i, w)$$

$$= \sum_i y^i (w_0 + \sum_{j=1}^N w_j x_j^i) - \log(1 + e^{(w_0 + \sum_{j=1}^N w_j x_j^i)})$$

4) Derivative w.r.t w_k

$$\frac{\partial L}{\partial w_k} = \sum_i^N x_k (y_i - P(y=1|w^T x_i))$$

$\downarrow f(w^T x_i)$

5) Gradient descent

$$w_{\text{new}}^k = w^k + \alpha \frac{\partial L}{\partial w_k}$$

$$w_{\text{new}}^k = w^k + \alpha \sum_{i=1}^N x_k (y_i - f(w^T x_i))$$

6) Object = $\arg\min -l(w; x, y)$

Adding regularization.

$$\text{Object} = \arg\min -l(w; x, y) + \frac{\lambda}{2} \|w\|^2$$

$$\text{Object} = \arg\min -$$

Gradient w.r.t w_i

$$\frac{\partial (\text{Object})}{\partial w_i} = -\frac{\partial (-l(w; x, y))}{\partial w_i} + \frac{\lambda}{2} \frac{\partial \|w\|^2}{\partial w_i}$$

$$\text{Gradient w.r.t } w_i = \sum_{j=1}^N x_i^j (y_j - f(w^T x_j)) + \lambda w_i$$

$$w_{\text{new}}^i = w^i + \alpha \frac{\partial (\text{Object})}{\partial w_i}$$

$$= w^i + \alpha (\lambda w^i) + \sum_{j=1}^N x_i^j (y_j - f(w^T x_j))$$

4.3 Analysis of results (Programming logistic regression)

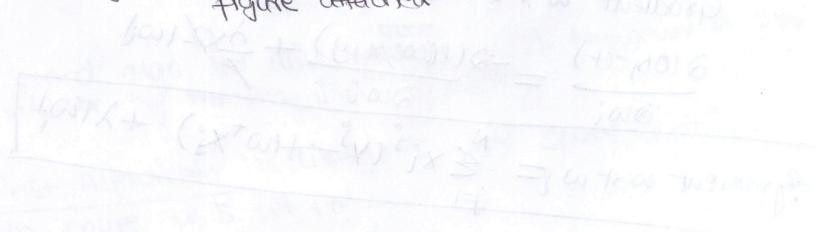
- 1) The Accuracy of the training set increased with the application of regularization ($\lambda = 25$).

accuracy = 98.5075] without regularization
Misclassified Image = 15]

accuracy = 98.7065] with L2 regularization
Misclassified Image = 13]

This is expected as penalizing the high weight features based on the training set would make the hypothesis much better i.e. reduce overfitting. Likewise, we see that L2 regularization improves the overall accuracy of any training set. However, on test data, accuracy improved a little bit.

- 2) Misclassified Image count = 13 (with L2 regularization)
figure attached



$$(b_1 + \sum_{j=1}^4 w_{1j}x_j + b_2 + \sum_{j=1}^4 w_{2j}x_j + b_3 + \sum_{j=1}^4 w_{3j}x_j) \geq 0$$

$$w_1^T x + b_1 \geq 0$$

Figure 1

