

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

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)

Thus, total multiplication = $3 + 3 + 3 =$ 9 Multiplication

Total Addition = 2

b) Computing through kernel function.

$$(x_1z_1 + x_2z_2)^2 \rightarrow 2 + 1 \rightarrow$$

Addition = 1

3 Multiplication

2) Perceptrons

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

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

② AND

Case 1 $x_1=0, x_2=0 \Rightarrow y=0 \Rightarrow \phi(\sum x_i w_i + b) \leq 0$
 $\Rightarrow x_1 w_1 + x_2 w_2 + b \leq 0$
 $\Rightarrow \boxed{b \leq 0}$ as $x_1=x_2=0$ ①

Case 2) $x_1=0, x_2=1$ AND $\rightarrow 0$
 $\Rightarrow 0 \cdot w_1 + 1 \cdot w_2 + b \leq 0$
 $w_2 + b \leq 0$ ②

Case 3) $x_1=1, x_2=0$ AND $\rightarrow 0$
 $\Rightarrow 1 \cdot w_1 + 0 \cdot w_2 + b \leq 0$
 $w_1 + b \leq 0$ ③

Case 4) $x_1=1, x_2=1$ AND $\rightarrow 1$
 $\Rightarrow 1 \cdot w_1 + 1 \cdot w_2 + b \geq 0$
 $\Rightarrow w_1 + w_2 + b \geq 0$ ④

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

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

⑦

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

4) Case 1) $x_1=0, x_2=0, y=0$
 $y=0 \Rightarrow \cancel{\phi} (\leq w_1 x_1 + b) \leq 0$
 $\Rightarrow \boxed{b \leq 0} \quad \text{--- ①}$

Case 2) $x_1=0, x_2=1, y=1$
 $\Rightarrow w_1 \cdot 0 + w_2 x_2 + b \geq 0$
 $\Rightarrow b + w_2 \geq 0 \quad \text{--- ②}$

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

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

Solving above 4 inequalities, we get a range of values.

one set of values

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

5)

x_1	x_2	XOR
0	0	0
1	0	1
0	1	1
1	1	0

XOR Table

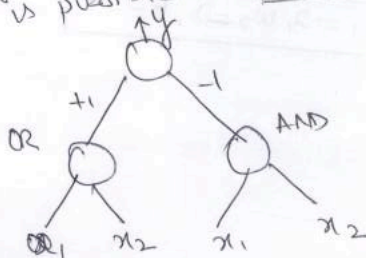
6) Case 1) $x_1=0, x_2=0, \text{ XOR}=0$
 $y=0 \Rightarrow \sum x_i w_i + b \leq 0$
 $\Rightarrow 0 \cdot w_1 + 0 \cdot w_2 + b \leq 0$
 $[b \leq 0] \text{ --- } \textcircled{1}$

Case 2) $x_1=0, x_2=1, \text{ XOR}=1$
 $y=1 \Rightarrow \sum x_i w_i + b > 0$
 $\Rightarrow 0 \cdot w_1 + 1 \cdot w_2 + b > 0$ --- $\textcircled{2}$

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

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

Above 4 equations are not possible to solve
 thus single layer perceptron is not possible.
 for logical XOR. $\textcircled{2}, \textcircled{3}, \textcircled{4}$ contradicts $\textcircled{1}$.
 However, two layer is possible, as XOR = OR-AND



3) Regression Theory

3.1) Linear Regression

g) $\{x_1, y_1\}, \dots, \{x_n, y_n\}$. $x \in \mathbb{R}^M$.
 $f(x) = w^T x$.

a) Square loss error.

error = misclassification $\Rightarrow f(x) \neq y_i$ for x_i

loss func $J(w) = (y_i - f(x))$ for i^{th} sample
 for all samples,

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

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

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

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

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

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

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

c) Gradient descent

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

new param \downarrow old params step size

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

$$w_{\text{new}}^k = w^k + \alpha \sum_{i=1}^n x_{ik} (y_i - w^T x_i) \quad \text{--- from ①}$$

do this till convergence.

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

a) Conditional likelihood L

$$L(w; x, y) = \underset{w}{\operatorname{argmax}} \prod_{i=1}^n P(y_i | x_i, w)$$

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

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

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

b) log conditional likelihood.

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

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

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

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

$$= c + \sum \log e^{-\frac{1}{2} \left(\frac{y_i - f(x_i, w)}{\sigma} \right)^2}$$

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

$$c = \sum \log \frac{1}{\sqrt{2\pi}\sigma^2}$$

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

$$\frac{\partial \log L(w; x, y)}{\partial w} = 2 \left(c - \frac{1}{2} \sum \left(\frac{y - f(x, w)}{\sigma} \right)^2 \right)$$

$$= 0 \quad - \frac{1}{2\sigma^2} \frac{\partial \sum (y - f(x, w))^2}{\partial w}$$

$$\frac{\partial \log L(w; x, y)}{\partial w} = - \frac{1}{2\sigma^2} \frac{\partial \sum_i (y - w^T x_i)^2}{\partial w}$$

$$= - \frac{1}{2\sigma^2} \sum_{i=1}^N (y - w^T x_i) x_i$$

$$\boxed{\frac{\partial \log L(w; x, y)}{\partial w} = - \frac{1}{2\sigma^2} \sum_{i=1}^N (y - w^T x_i) x_i}$$

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

3)

$$y = f(x) + \epsilon$$

↓
defined by us

↓
Noise

ϵ : mean = 0
variance σ^2

a)

$$E_D \left[\int_{y_n} \int_{x_n} (y - f(x))^2 p(y|x) p(x) dy dx \right]$$

point is x .

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

~~$E(h(x))$~~

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

$$E[(y - y')^2] = E[(y - l + l - y')^2]$$

$$= E[(y - l)^2 + (l - y')^2 + 2[y - l][l - y']]$$

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

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

↓

Bais

↓

variance

+ σ^2

↓

unavoidable error

3.2) Regularization.

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

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

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

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

b) First Algorithm is more likely to give a sparse matrix as it penalises the weight with high weightage of $w_{k,t}$.

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

In second algo, it will ~~reduce~~ ^{increase} by λ . However in other it will be by the magnitude of $w_{k,t}$.

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.

$$f(x, w) = \frac{1}{1 + e^{-a}}, \text{ where } a = \sum_{i=1}^N x_i w_i + w_0, N \rightarrow \text{no. of feature.}$$

2) conditional likelihood

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

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

$$\text{Conditional Likelihood} = \prod_{\text{over all the samples/training data}} P(y^i | x^i, w)$$

~~$L(w; x, y)$~~

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

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

3) ~~Cond~~ 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

$$\frac{\partial L}{\partial w^k} = \sum_i x_k^i (y^i - p(y^i | x^i; w))$$

\downarrow
 $f(w^T x_i)$

5) Gradient descent

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

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

6) Object = $\text{argmin} -l(w; x, y)$

Adding regularization.

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

~~$$\text{Object} = \text{argmin}$$~~

Gradient w.r.t w_i

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

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

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

$$= w^k + \alpha (\lambda w^k) + \sum_{j=1}^N x_k^j (y^j - f(w^T x_j))$$

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
Misclassified Image = 15 } without regularization

accuracy = 98.7065
Misclassified Image = 13 } with L2 regularization

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

