ASSIGNMENT #3

18th Feb 2016

ADITYA GAUTAM
agautam1@andrew

Ques 1)

1) Kernel feature mapping.

$x = (x_1, x_2)^T$

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

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

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

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

corresponding kernel function = $\boxed{(x_1z_1 + x_2z_2)^2}$

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

a) mapping into feature space.

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

Total no. of multiplication mapping = 3

Similarly, for z mapping, we would need 3 multiplication

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

Total multiplication = 3+3+3 = $\boxed{9 \text{ multiplication}}$

Thus total multiplication

$\boxed{\text{Total Addition} = 2}$

b) computing through kernel function.

$(x_1z_1 + x_2z_2)^2 \rightarrow 2+1 \rightarrow \boxed{3 \text{ Multiplication}}$

$\boxed{\text{Addition} = 1}$

## 2) Perceptrons

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

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

**②** AND.

**Case 1** $x_1 = 0, x_2 = 0 \Rightarrow y = 0 \Rightarrow \phi\left(\sum x_i w_i + b\right) \leq 0$

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

$\Rightarrow$ as $x_1 = x_2 = 0$

$\boxed{b \leq 0}$ ——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 $\rightarrow 1$

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

$\Rightarrow w_1 + w_2 + b > 0$ ——④

Solving above 4 inequalities we get,

$b \leq 0, \ w_1 \leq b, \ w_2 \leq b, \ w_1 + w_2 > -b$

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

**①** AND table

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

3) OR contingency table

| X₁ | X₂ | OR |
|---|---|---|
| 0 | 0 | 0 |
| – | – | – |
| – | 0 | – |
| – | – | – |

4) Case 1) $x_1=0, x_2=0, y=0$

$y=0 \Rightarrow \phi(w_0+w_1x_1+w_2x_2) < 0$

$\Rightarrow \boxed{w_0 < 0}$ ——— ①

Case 2) $x_1=0, x_2=1, y=1$

$\Rightarrow w_0 + w_2\cdot 0 + w_1\cdot 0 \dots$

$\Rightarrow w_0 + w_2 > 0$ ——— ②

Case 3) $x_1=1, x_2=0, y=1$

$\Rightarrow w_0 + w_1 > 0$ ——— ③

Case 4) $x_1=1, x_2=1, y=1$

$\Rightarrow w_0 + w_1 + w_2 > 0$ ——— ④

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

one set of values

$$\boxed{w_0 = 2, \ w_2 = -2, \ b = -1}$$

## 5)

XOR Table

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

## 6)

Case 1) $x_1=0$, $x_2=0$  XOR $=0$

$y=0 \Rightarrow \phi(w_1\cdot 0 + w_2\cdot 0 + b) \le 0$
$\Rightarrow \phi(b) \le 0$  ... ①
$\boxed{b \le 0}$

Case 2) $x_1=0$, $x_2=1$,  XOR $=\phi$

$\Rightarrow \phi(w_2 + b) > 0$
$\Rightarrow w_2 + b > 0$  ... ②
$\Rightarrow f=1$

Case 3) $x_1=1$, $x_2=0$,  XOR $=1$

$\Rightarrow \phi(w_1 + b) > 0$
$\Rightarrow w_1 + b > 0$  ... ③

Case 4) $x_1=1$, $x_2=1$,  XOR $=0$

$\Rightarrow \phi(w_1 + w_2 + b) \le 0$
$\Rightarrow w_1 + w_2 + b \le 0$  ... ④

Above 4 equations are not possible to solve
thus single layer perceptron is not possible.
②&③&④ contradicts ①

However, two layer is possible, as XOR $=$ OR-AND
for logical XOR.

3) Regression Theory
===

3.1) Linear Regression

$\{x_1, y_1, \ldots$ $\{x_n, y_n\}$ $x \in \mathbb{R}^m$

$f(x_i) = \omega^T x$

a) Square loss error

error = misclassification $\Rightarrow f(x_i) \neq \hat{y}_i$ for $x_i$

$error = (y_i - f^i(x_i))$ for $i^{th}$ sample

$J(\omega) = (y_i - \hat{y}_i)^2$

loss func    for all samples.

$$J(\omega) = \frac{1}{n} \sum_{i=1}^{n} (\hat{y}_i - y_i)^2 = \frac{1}{n} \sum_{i=1}^{n} (\hat{y}_i - \omega^T x_i)^2$$

$$J(\omega) = \frac{1}{n} \sum_{i=1}^{n} (\hat{y}_i - \omega^T x_i)^2$$

b) Partial derivation of motivation to cost w.r.t $\omega^k$

$$\frac{\partial J(\omega)}{\partial \omega^k} = \frac{1}{n} \sum_{i=1}^{n} 2(y_i - \omega^T x_i) \times -x^k = \frac{2}{n} (y_i - \omega^T x_i) x - x_i^k$$

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

$$\frac{\partial J(\omega)}{\partial \omega^k} = \sum_{i=1}^{n} x_i^k (y_i - \omega^T x_i) \times$$

$$\sum_{i=1}^{n} (y_i - \omega^T x_i) x_i^k \qquad ①$$

c) Gradient descent

$\omega^k = \omega^k + \alpha \frac{\partial J(\omega)}{\partial \omega^k}$

$\omega^k_{new} = \omega^k + \alpha \frac{\partial J(\omega)}{\partial \omega^k}$

$\omega_{new}$    $\uparrow$ step size
old towards

from ①

$$\omega_{new}^k = \omega^k + \alpha \sum_{i=1}^{N} x_i^k (y_i - \omega^T x_i)$$

$\downarrow$ do this till convergence.

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

a) Conditional likelihood $L$

$$L = \frac{\prod P(y_i \mid x_i, w)}{\phantom{x}} = \prod_{i=1}^{n} P\left(\frac{y_i}{x_i, w}\right)$$

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

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

$$P(y/x, \omega) = \frac{1}{\sqrt{2\pi\sigma^2}} \, e^{-\frac{1}{2}\left(\frac{y - f(w^T x)}{\sigma}\right)^2}$$

$$\Rightarrow \boxed{\,L(w; x, y) = \text{argmax} \prod \frac{1}{\sqrt{2\pi\sigma^2}} \, e^{-\frac{1}{2}\left(\frac{y - f(w^T x)}{\sigma}\right)^2}\,}$$

b) log conditional likelihood.

$$\log \text{ conditional likelihood} = \text{argmax} \log \prod \frac{1}{\sqrt{2\pi\sigma^2}} \, e^{-\frac{1}{2}\left(\frac{y - f(w^T x)}{\sigma}\right)^2}$$

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

$$\log a \cdot b = \log a + \log b$$

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

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

$$= c - \frac{1}{2} \sum \left(\frac{y - f(w^T x)}{\sigma}\right)^2$$

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

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

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

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

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

$$\frac{\partial \log L(\omega;x,y)}{\partial\omega} = \frac{1}{2\sigma^2}\frac{\partial \sum(y-\omega^Tx_i)^2}{\partial\omega}$$

$$= \frac{1}{2\sigma^2}\sum\limits_{i=1}^{N} 2(y-\omega^Tx_i)\,x^i$$

$$\boxed{\frac{\partial \log L(\omega;x,y)}{\partial\omega} = \frac{1}{\sigma^2}\sum\limits_{i=1}^{N}(y-\omega^Tx_i)\,x_i}$$

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) $y = f(x) + \epsilon$;   mean $= 0$
   $\qquad \downarrow$   variance $\sigma^2$
   $\quad$ Noise
   defined by us

a) $E_D \left[ \iint (y - f(x))^2 \, p(y/x) \, p(x) \, dy \, dx \right]$

point is $x$.

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

$(\hat{y} = f_D(x))$   $\quad \ell = E_D(h(x))$

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

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

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

$= E[ (y - \ell)^2 ] + E[ (\ell - \hat{y})^2 ] + 2 E[\{\hat{y}\ell - \ell^2\} - E[\ell\hat{y}] - E[\ell\hat{y}]]$
$\qquad\qquad \downarrow \qquad\qquad\qquad\quad \downarrow \qquad\qquad\qquad\qquad -E[\ell\hat{y}]$
$\qquad\qquad \downarrow \qquad\qquad\qquad\quad \downarrow \qquad\qquad\qquad\qquad \downarrow$
$\qquad\quad$ Bais $\qquad\qquad$ variance $\;+\; \sigma^2$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad \downarrow$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ unavoidable error

## 3.2) Regularization

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

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

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

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

b) first Algorithm is more likely to give a sparse matrix as it penalised the high weighage of w(k+) weight with

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

In second algo, reduces parainbow ht fig by $\lambda$. Otherwise it will be high.

So, first Algorithm is more likely to give spane matrix with more w's set to zero

## 4.1) Logistic Regression

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

**1) logistic function.**

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

**2) conditional likelihood**

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

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

Conditional likelihood $= \prod P(Y^i | x^i, w)$

over all the samples training data

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

$$w_{mle}(w) = \underset{w}{\arg\max} \prod_i P(y^i | x^i, w)$$

**3) log conditional likelihood**

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

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

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

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

$$\qquad (\log ab = \log a + \log b)$$

$$= \sum_i y^i (w_0 + \sum_i w_i x_i^i) - \log\left(1 + e^{\left(w_0 + \sum_{i=1}^{n} w_i x_i^i\right)}\right)$$

4) Derivative w.r.t $w^k$

$$\boxed{\frac{\partial L}{\partial w_k} = \sum_i x_k^i \left(y^i - P(y^i=1|x^i, w)\right)}$$

$$\underset{f(w^T x_i)}{\uparrow}$$

5) gradient descent

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

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

6) Object = argmin $-l(w^i; x, y)$

Adding regularization.

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

Object = argmin

gradient w.r.t $w_i$

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

$$\boxed{\text{gradient w.r.t } w = -\sum_{j=1}^{N} x_i^i \left(y^i - f(w^T x_i)\right) + \lambda(w)}$$

$$w^i_{new} = w^k + \alpha \frac{\partial Object}{\partial w_i}$$

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

## 4.3 Analysis of results (Programming logistic regression)

1) The Accuracy of the training set increased with the application of Regularization (lambda = 26).

Accuracy = 98.5075 ⎤ without regularization
# Mis classified Image = 15 ⎦

Accuracy = 98.7065 ⎤ with L2 regularization
# Mis classified 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 anything set. However on test-data, accuracy improved a little bit.

2) Mis classified Image count = 13 (with L2 regularization) figure attached