

02/16/16

MACHINE LEARNING HOMEWORK-3

SAM DAVID CHRIST DOSS PUSHPAM
schristd@andrew.cmu.edu.

1. KERNEL FEATURE MAPPINGS :-

$$1) \quad x = (x_1, x_2)^T, \quad \phi(x) = (x_1^2, \sqrt{2} x_1 x_2, x_2^2)^T$$

$$\begin{aligned} K(x, z) &= \phi(x) \cdot \phi(z) = (x_1^2, \sqrt{2} x_1 x_2, x_2^2)^T \cdot (z_1^2, \sqrt{2} z_1 z_2, z_2^2)^T \\ &= x_1^2 z_1^2 + 2 x_1 z_1 x_2 z_2 + x_2^2 z_2^2 \\ &= (x_1 z_1 + x_2 z_2)^2 \\ &= (x \cdot z)^2 \end{aligned}$$

2) a) map to feature space and dot product.

$$K(x, z) = \phi(x) \cdot \phi(z) = x_1^2 z_1^2 + 2 x_1 z_1 x_2 z_2 + x_2^2 z_2^2$$

$$\begin{aligned} x_1^2 z_1^2 &= 4 \text{ operations } (x_1 \times x_1, x_2 \times x_2, z_1 \times z_1, z_2 \times z_2) \\ 2 x_1 z_1 x_2 z_2 &= 5 \text{ operations.} \\ x_2^2 z_2^2 &= 4 \text{ operations.} \end{aligned}$$

$$\begin{aligned} \phi(x) &= (x_1^2, \sqrt{2} x_1 x_2, x_2^2)^T \\ \left. \begin{aligned} x_1^2 &= 1 \\ \sqrt{2} x_1 x_2 &= 2 \\ x_2^2 &= 1 \end{aligned} \right\} \Rightarrow 4 \text{ operations for each } \phi \end{aligned}$$

dot product.

$$\begin{aligned}\phi(x) \cdot \phi(x) &= \left((x_1^2), (\sqrt{2}x_1x_2), (x_2^2) \right)^T \cdot \left((z_1^2), (\sqrt{2}z_1z_2), (z_2^2) \right)^T \\ &= \left(\underset{\substack{\downarrow \\ 1}}{(x_1^2)} \underset{\substack{\downarrow \\ 1}}{(z_1^2)} + \underset{\substack{\downarrow \\ 1}}{(\sqrt{2}x_1x_2)} \underset{\substack{\downarrow \\ 1}}{(\sqrt{2}z_1z_2)} + \underset{\substack{\downarrow \\ 1}}{(x_2^2)} \underset{\substack{\downarrow \\ 1}}{(z_2^2)} \right) \\ &\quad \underbrace{\hspace{10em}}_{5 \text{ operations.}}\end{aligned}$$

$\phi(x) \Rightarrow 2 \times (4 \text{ for each } \phi) + 5 \text{ operations while dot product}$

$\Rightarrow 13 \text{ operations}$

b) $K(x, z) = (x \cdot z)^2$

$$x \cdot z = x_1z_1 + x_2z_2 \Rightarrow 3 \text{ operations.}$$

$$(x \cdot z)^2 \Rightarrow 1 \text{ operations}$$

$\Rightarrow 4 \text{ operations totally.}$

2) Perceptrons :-

$$y = \psi \left(\sum_i x_i * w_i + b \right),$$

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

1) AND operation

x_1	x_2	AND
0	0	0
0	1	0
1	0	0
1	1	1

2) w_1 and w_2 for AND.

$$w_1 = 1$$

$$w_2 = 1$$

$$b = -1$$

3) OR:-

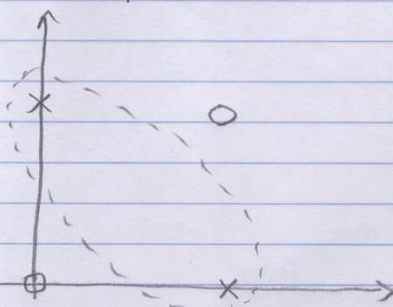
x_1	x_2	OR
0	0	0
0	1	1
1	0	1
1	1	1

4) $w_1 = 1$, $w_2 = 1$, $b = 0$

5) XOR :-

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

6) Single layer perceptron can't be used to create logical XOR.



The XOR function is not linearly separable, which is possible in AND, OR functions. The above graph shows the plot that linear line cannot be drawn to separate (1,1), (0,0) from (0,1), (1,0). So it is not possible to find a values for w_1 , w_2 , b in the line equation $w_1 x_1 + w_2 x_2 + b$

3. Regression theory :-

3.1 Linear regression :-

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

a) Squared loss, $J(w)$ associated with $f(x)$

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

b) Partial derivative w.r. to, w^k .

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

$$\frac{\partial}{\partial \omega_k} f(x_i) = \frac{\partial}{\partial \omega_k} \omega^T x_i \Rightarrow (x_i)^k$$

$$\Rightarrow \frac{\partial}{\partial \omega^k} J(\omega) = \frac{1}{2} \sum_i (y_i - f(x_i)) (x_i)^k = \sum_i (y_i - f(x_i)) x_i^k$$

c) update rule:-

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

$$= \omega^k + \alpha \left(\sum_i (y_i - f(x_i)) x_i^k \right)$$

$$\omega_{\text{new}}^k = \omega^k + \alpha \sum_i (y_i - f(x_i)) x_i^k$$

where, $\alpha < 0$

with, $\alpha < 0$

2) maximum likelihood:-

$$a) y_i \sim N(\omega^T x_i, \sigma^2)$$

Assuming that y_i is iid in normal distribution.

$$L(\omega; x, y) = \prod_i \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(y_i - f(x_i))^2}{2\sigma^2}}$$

\Rightarrow This is due to the condition that y_i is iid,

b) Log of conditional likelihood :-

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

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

To maximize $\ell(w; x, y)$, take derivative w.r.t. w^k .

$$\frac{\partial}{\partial w^k} \ell(w; x, y) = -\frac{1}{2\sigma^2} \frac{\partial \sum_i (y_i - f(x_i))^2}{\partial w^k}$$

$$= -\frac{1}{2\sigma^2} \sum_i (y_i - f(x_i)) \frac{\partial f(x_i)}{\partial w^k}$$

$$\frac{\partial f(x_i)}{\partial w^k} = x_i^k$$

$$= -\frac{1}{2\sigma^2} \sum_i ((y_i - f(x_i)) x_i^k)$$

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

maximizing ℓ is same as minimizing $-\ell$

So, minimizing $-l$ is written as.

$$(-l) = \sum_i \frac{(y_i - f(x_i))^2}{2\sigma^2} - n \log\left(\frac{1}{\sqrt{2\pi}\sigma}\right)$$

The second term is constant. This formula will give the same result as minimizing in (1b) which is,

$$\text{minimizing } \frac{1}{2} \sum_i (y_i - f(x_i))^2$$

So, maximizing $l(w; x, y)$ is same as minimizing least square error.

3) $y = f(x) + \epsilon$, ϵ has mean 0 and variance σ^2

a) Expected loss under squared error at some test point x .

$$E_{D, \epsilon}((y - \hat{f})^2) = E_{D, \epsilon}((y - f)^2),$$

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

The expected loss is taken over the dataset and the noise.

$$y = f(x) + \epsilon, \quad \hat{f} = h_D(x)$$

$\Rightarrow h_D$ is the linear regressor trained using D .

$$\begin{aligned}
b) \quad MSE &= E_{D, \epsilon} (y - \hat{\beta})^2 \\
&= E_{D, \epsilon} (y - E(\hat{\beta}) + E(\hat{\beta}) - \hat{\beta})^2 \\
&= E_{D, \epsilon} \left[(y - E(\hat{\beta}))^2 + (E(\hat{\beta}) - \hat{\beta})^2 + 2(y - E(\hat{\beta}))(E(\hat{\beta}) - \hat{\beta}) \right] \\
&= E_{D, \epsilon} (y - E(\hat{\beta}))^2 + E_{D, \epsilon} (E(\hat{\beta}) - \hat{\beta})^2 + 2 E_{D, \epsilon} \left\{ (y - E(\hat{\beta}))(E(\hat{\beta}) - \hat{\beta}) \right\} \\
&= E_{D, \epsilon} (y - E(\hat{\beta}))^2 + E_{D, \epsilon} (E(\hat{\beta}) - \hat{\beta})^2 + 2 (y - E(\hat{\beta}))(E(\hat{\beta}) - E(\hat{\beta})) \\
&= E_{D, \epsilon} ((\beta(x) + \epsilon - E(\hat{\beta}))^2) + E_{D, \epsilon} ((E(\hat{\beta}) - \hat{\beta})^2) + 0 \\
&= E_{D, \epsilon} ((\beta(x) - E(\hat{\beta}))^2) + E_{D, \epsilon} (\epsilon^2) - 2 E_{D, \epsilon} \left\{ (\beta(x) + \epsilon)(\epsilon) \right\} + E_{D, \epsilon} ((E(\hat{\beta}) - \hat{\beta})^2) \\
&= E_{D, \epsilon} ((\beta(x) - E(\hat{\beta}))^2) + \sigma^2 + E_{D, \epsilon} ((E(\hat{\beta}) - \hat{\beta})^2) - 2 (E(\beta(x)) + E(\epsilon))(E(\epsilon))
\end{aligned}$$

Since, $E(\epsilon) = 0 \Rightarrow$ mean is 0,
last term becomes 0.

$$\begin{aligned} \text{MSE} &= E_{\mathcal{D}} [f(x) - E(f(x))]^2 + E[(E(\hat{f}) - f)^2] + \sigma^2 \\ &= \text{Bias}^2 + \text{Variance} + \sigma^2 \end{aligned}$$

3.2 Regularization

$$1) \ a) \ L = \frac{1}{2} \sum_{i=1}^n (y_i - w^T x_i)^2 + \frac{\lambda}{2} \|w\|^2$$

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

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

b) The algorithm (ii) is most likely to give sparse w . The value of w is ~~normalized~~ normalized in that case, so it may suppress the value to 0 completely.

with ~~the~~ the normalization factor actually suppressing the value of w^k , (ii) tends to give sparse w .

4. Programming logistic regression.

4.1 Logistic regression.

1) Logistic function $f(x; w)$

$$f(x; w) = \frac{1}{1 + e^{(-w^T x)}}$$

2) Conditional likelihood,

$$L(w; x, y) = \begin{cases} f(x; w), & \text{if } y = 1 \\ 1 - f(x; w), & \text{if } y = 0 \end{cases}$$

~~$L(w; x, y) = f(x; w)^y (1 - f(x; w))^{1-y}$~~

3) Log conditional likelihood.

~~$\ell(w; x, y) = \log L(w; x, y)$~~

Log likelihood on one example.

$$\Rightarrow \begin{cases} \log(f(x; w)), & \text{if } y = 1 \\ \log(1 - f(x; w)), & \text{if } y = 0 \end{cases}$$

~~$\ell(w; x, y) = \log(f(x; w)^y (1 - f(x; w))^{1-y})$~~

4) Derivative of log likelihood.

$$\frac{\partial}{\partial w_k} \ell(w; x, y) = \begin{cases} \frac{\partial}{\partial w_k} \log(f(x; w)), & \text{if } y = 1 \\ \frac{\partial}{\partial w_k} \log(1 - f(x; w)), & \text{if } y = 0 \end{cases}$$

$$\frac{\partial}{\partial \omega_k} \ell(\omega; x, y) = \begin{cases} \frac{1}{f} \frac{\partial}{\partial \omega_k} f(x; \omega), & \text{if } y=1 \\ \frac{1}{1-f} \left(-\frac{\partial}{\partial \omega_k} f(x; \omega) \right), & \text{if } y=0 \end{cases}$$

$$f = \frac{1}{1 + e^{-\omega^T x}}, \quad 1-f = \frac{e^{-\omega^T x}}{1 + e^{-\omega^T x}}$$

$$\frac{\partial}{\partial \omega_k} f \Rightarrow (-1) (1 + e^{-\omega^T x})^{-2} \frac{\partial}{\partial \omega_k} (e^{-\omega^T x})$$

$$\Rightarrow (-1) (1 + e^{-\omega^T x})^{-2} (e^{-\omega^T x}) (x^k)$$

$$\Rightarrow \frac{1}{(1 + e^{-\omega^T x})} \cdot \frac{e^{-\omega^T x}}{(1 + e^{-\omega^T x})} x^k$$

$$\frac{\partial}{\partial \omega_k} (f) \Rightarrow (f(x; \omega))(1-f(x; \omega)) x^k$$

$$\frac{\partial}{\partial \omega_k} \ell(\omega; x, y) = \begin{cases} \frac{1}{f} \cdot f(1-f) x^k, & \text{if } y=1 \\ \left(\frac{1}{1-f} \right) (-f(1-f)) x^k, & \text{if } y=0 \end{cases}$$

$$= \begin{cases} (1-f) x^k, & \text{if } y=1 \\ -f x^k, & \text{if } y=0 \end{cases}$$

This can also be written as,

$$\frac{\partial}{\partial \omega_k} \ell(\omega; x, y) = \sum_{i=1}^n (y_i - f) x^k$$

5) The update rule is written as,

$$w_{\text{new}}^k = w^k - \alpha \frac{\partial}{\partial w_k} \ell(w; x, y)$$

$$w_{\text{new}}^k = w^k - \alpha \sum_{i=1}^n (y_i - f(w; x, y)) x_i^k$$

6) $\arg \min_w -\ell(w; x, y) + \frac{\lambda}{2} \|w\|^2$

Taking derivative w.r.to w_k gives

$$\Rightarrow -\frac{\partial}{\partial w_k} \ell(w; x, y) + \frac{\partial}{\partial w_k} \left(\frac{\lambda}{2} \|w\|^2 \right)$$

$$\Rightarrow -(y_i - f(w; x, y)) x_i^k + \lambda w_k$$

4.3)

1) For Unregularized LR,
Training accuracy :- 99.1575 %
Testing accuracy :- 98.607 %

For regularized LR,
Training accuracy :- 99.2236 %
Testing accuracy :- 98.607 %

Training accuracy increased slightly in, regularized LR, whereas the testing accuracy remains the same. The gap between training test error becomes wider as a result.

We would expect a ~~decrease~~ ~~in~~ increase in testing accuracy as well. This might be due to low value of the regularization parameter, that it does not affect the parameter values significantly.

5) Programming Kernel perception :-

5.2)

4. Perception :-

Training accuracy = 98.8932 %
Testing accuracy = 98.5075 %

Kernel perceptron:

Training accuracy : 99.7605%.

Testing accuracy : 99.602%.

Kernel perceptron has higher accuracy in both training and testing set. This is because the kernel perceptron actually works in higher dimensional space implicitly, which results in better classification.

4.3)

2)

