

10601 Machine Learning HW3
Hsueh Lin Huang
hsuehlinh

1. Kernel feature mappings

1.

$$K(x, z) = \phi(x) \cdot \phi(z) = (x_1^2, \sqrt{2}x_1x_2, x_2^2) \cdot (z_1^2, \sqrt{2}z_1z_2, z_2^2) \\ = (x_1z_1 + x_2z_2)^2$$

2.

(a) $\mathbb{R}^2 \rightarrow \mathbb{R}^3$ $\begin{bmatrix} \end{bmatrix}_{3 \times 2} \begin{bmatrix} \end{bmatrix}_{2 \times 1} \rightarrow (2 \text{ multiplication} + 1 \text{ addition}) \times 3$

$$\text{mul} = 2 \times 3 \times 2 + 3 = 15$$

$$\text{add} = \underbrace{3 \times 2}_{\text{map to feature space}} + (3 - 1) = 8$$

map to
feature space

(b) $\mathbb{R}^2 \rightarrow \mathbb{R}^3$

$$K(x, z) = (x_1^2, x_2^2, \sqrt{2}x_1x_2)$$

$$\text{mul} = 3$$

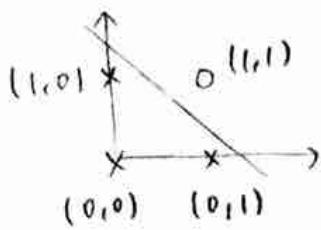
$$\text{add} = 1$$

2. Perceptrons

1.

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

2.



$$\sum w_i x_i = -b$$

$$w_1 x_1 + w_2 x_2 = -b$$

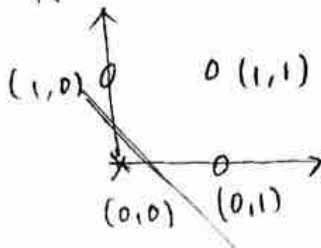
$$x_2 = \frac{-b - w_1 x_1}{w_2}$$

$$w_1 = 1, w_2 = 1, b = -1.5$$

3.

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

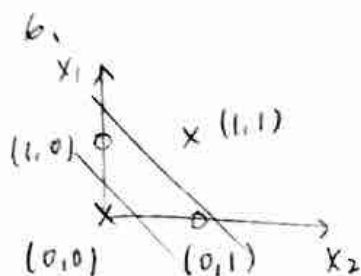
4.



$$w_1 = 1, w_2 = 1, b = -0.5$$

5. XOR

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



We need two straight lines to separate the different outputs.

→ we can only solve this either change the transfer function so that it has more than one decision boundary, or use a more complex network that is able to generate more complex decision boundary.

3. Regression Theory

3.1 Linear Regression

1. (a)

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

$$\text{Squared error loss } J(w) = \|y - f(x)\|^2 = \sum_{i=1}^n (y_i - f(x_i))^2$$

(b)

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

(c)

$$w_{\text{new}}^k = w^k + \underset{\substack{\downarrow \\ \text{Step}}}{\alpha} \cdot \frac{\partial J(w)}{\partial w^k} = w^k + \alpha \sum_{i=1}^n (y_i - f(x_i)) x_i^k$$

2.

(a)

$$\begin{aligned} L(w | x, y) &= p(y | x, w) = \prod_{i=1}^n p(y_i | x_i, w) \\ &= \prod_{i=1}^n (2\pi\sigma^2)^{-\frac{1}{2}} \exp\left(-\frac{1}{2} \frac{(y_i - f(x_i))^2}{\sigma^2}\right) \\ &= (2\pi\sigma^2)^{-\frac{n}{2}} \exp\left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - f(x_i))^2\right) \end{aligned}$$

$$\begin{aligned}
 (b) \quad \ell(w|x, y) &= \ln(L(w|x, y)) = \ln \left((2\pi\sigma^2)^{-\frac{n}{2}} \exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - f(x_i))^2 \right) \right) \\
 &= \ln \left((2\pi\sigma^2)^{-\frac{n}{2}} \right) + \ln \left(\exp \left(-\frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - f(x_i))^2 \right) \right) \\
 &= -\frac{n}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - f(x_i))^2 \\
 &= -\frac{n}{2} \ln(2\pi) - \frac{n}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^n (y_i - f(x_i))^2
 \end{aligned}$$

$$\begin{aligned}
 (c) \quad \frac{\partial \ell(w|x, y)}{\partial w_k} &= -\frac{1}{\sigma^2} \sum_{i=1}^n (y_i - f(x_i)) x_{i,k} \\
 \frac{\partial J(w)}{\partial w_k} &= -\sum_{i=1}^n (y_i - f(x_i)) x_{i,k}
 \end{aligned}$$

minimize negative log likelihood function
is the same as minimize the loss function.

3.

$$\epsilon \sim \mathcal{N}(0, \sigma^2)$$

$$(a) \quad E \left[\sum_{i=1}^n (y_i - \hat{f}(x_i))^2 \right]_{\epsilon_i}$$

→ The expectation is taken over by ϵ

(b)

$$\begin{aligned}
 E[(y - \hat{f}(x))^2] &= E[\hat{f}(x)^2 - 2\hat{f}(x)y + y^2] = E[\hat{f}(x)^2] - 2E[\hat{f}(x)]E[y] + E[y^2] \\
 &= E[(\hat{f}(x) - E[\hat{f}(x)])^2] + E[\hat{f}(x)]^2 - 2E[\hat{f}(x)]f(x) + E[(y - f(x))^2] + f(x)^2 \\
 &= E[(\hat{f}(x) - E[\hat{f}(x)])^2] + (\text{variance}) \\
 &\quad + (E[\hat{f}(x)] - f(x))^2 + (\text{bias}^2) \\
 &= E[(y - f(x))^2] + (\sigma^2)
 \end{aligned}$$

3.2 Regularization

1. (a)

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

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

(b)

(i) is more likely to give sparse w , because it has additional constraint that $w^{k,(t+1)} < 0$, set $w^{k,(t+1)} = 0$

4. Programming logistic regression

1. $f(x|w) = \frac{1}{1 + e^{-x \cdot w}}$

2. $L(w|x, y) = \prod_{i=1}^n P(y_i | x_i, w) = \prod_{i=1}^n f(x_i|w)^{y_i} (1 - f(x_i|w))^{1-y_i}$

3.
$$\begin{aligned} \ell(w|x, y) &= \sum_{i=1}^n y_i \log f(x_i|w) + (1 - y_i) \log (1 - f(x_i|w)) \\ &= \sum_{i=1}^n \log (1 - f(x_i|w)) + \sum_{i=1}^n y_i \log \frac{f(x_i|w)}{1 - f(x_i|w)} \\ &= \sum_{i=1}^n -\log (1 + e^{w_0 + x_i \cdot w}) + \sum_{i=1}^n y_i (w_0 + x_i \cdot w) \end{aligned}$$

4.
$$\begin{aligned} \frac{\partial \ell(w|x, y)}{\partial w^k} &= - \sum_{i=1}^n \frac{1}{1 + e^{w_0 + x_i \cdot w}} e^{w_0 + x_i \cdot w} x_i^k + \sum_{i=1}^n y_i x_i^k \\ &= \sum_{i=1}^n (y_i - f(x_i|w)) x_i^k \end{aligned}$$

$$5. \quad w_{\text{new}}^k = w^k + \alpha \frac{\partial \ell(w|x, y)}{\partial w^k} = w^k + \alpha \sum_{i=1}^n (y_i - f(w^T x_i)) x_i^k$$

$$6. \quad \frac{\partial \left(-\ell(w|x, y) + \frac{\lambda}{2} \|w\|^2 \right)}{\partial w^k} = - \sum_{i=1}^n (y_i - f(x|w)) x_i^k + \lambda w_i^k$$

4 Programming Logistic Regression

4.3 Analysis of Results

w/o Regularization:

Your training accuracy is:99.2318%

Your testing accuracy is:98.7065%

Regularization:

Your training accuracy is:97.5799%

Your testing accuracy is:96.8159%

1. The accuracy drops after adding the regularization term. The purpose of adding regularization term is to avoid overfitting problem. However, if lambda is too large, then it might underfit the data because of high bias.
- 2.



5 Programming Kernel Perceptron

5.1 Perceptron

Your training accuracy is:98.4224%

Your test accuracy is:97.7114%

5.2 Kernel Perceptron

Your training accuracy is:97.7699%

Your test accuracy is:97.7114%

It seems like the accuracy between kernel perceptron and perceptron are the same.