

(1)

1. (10 points) In class, we discussed how to represent XOR-like functions using quadratic features, since standard linear classifiers (such as perceptrons) are insufficient for this task. However, here we show that XOR-like functions can indeed be simulated using *multi-layer* networks of perceptrons. This example shows a glimpse of the expressive power of “deep neural networks”: merely increasing the depth from 1 to 2 layers can help reproduce nonlinear decision boundaries.

- a. Consider a standard two-variable XOR function, where we have 2-dimensional inputs

$$x_1, x_2 = \pm 1, \text{ and output } y = x_1(\text{XOR})x_2 = \begin{cases} -1 & \text{if } x_1 = x_2 \\ 1 & \text{otherwise} \end{cases}.$$

Geometrically argue why a single perceptron cannot be used to simulate the above function.

- b. Graphically depict, and write down the equation for, the optimal decision region for the following logical functions:

(i)  $x_1(\text{AND})(\text{NOT}(x_2))$

(ii)  $(\text{NOT}(x_1))(\text{AND})x_2$

- (iii)  $x_1(\text{OR})x_2$  Make note of the weights learned corresponding to the optimal decision boundary for each function.

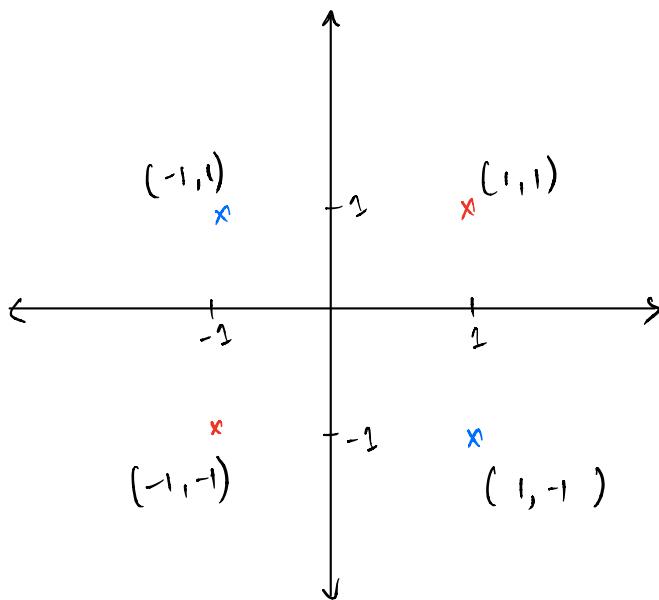
- c. Using the above information, simulate a multi-layer perceptron *network* for the XOR operation with the learned weights from Part (b).

(a)

$x_1$	$x_2$	XOR	Perceptron equation
-1	-1	-1	$w_0 + \sum_{i=1}^2 w_i x_i < 0 \rightarrow w_0 - w_1 - w_2 < 0$
-1	1	1	$w_0 + \sum_{i=1}^2 w_i x_i \geq 0 \rightarrow w_0 - w_1 + w_2 \geq 0$
1	-1	1	$w_0 + \sum_{i=1}^2 w_i x_i \geq 0 \rightarrow w_0 + w_1 - w_2 \geq 0$
1	1	-1	$w_0 + \sum_{i=1}^2 w_i x_i < 0 \rightarrow w_0 + w_1 + w_2 < 0$

Obtained by substituting  
the input values,  $x_1 \& x_2$

By looking at last column that contains equations that can be solved to obtain the weights, we can see some clear contradiction in the values of weights. This means that there are no definite weights that can be learnt for solving the XOR function.



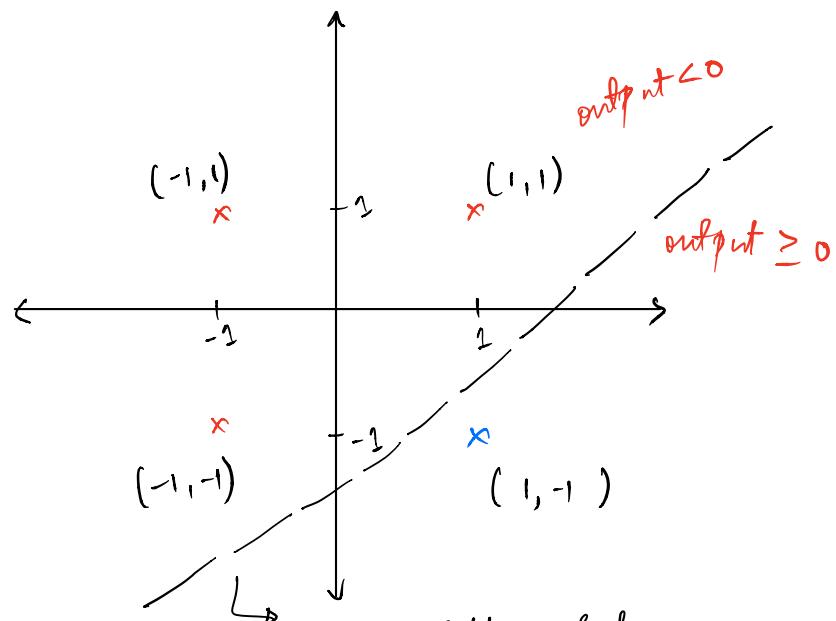
Geometrically, it is evident that no solution exists that can classify these 4 points according to their labels  $\rightarrow \{ -1, 1 \}$

$\text{x} \rightarrow \text{XOR output } 1$

$\text{x} \rightarrow \text{XOR output } -1$

(b) (i)  $y = x_1 \text{ (AND)} \quad \text{NOT}(x_2)$

$x_1$	$x_2$	$y$
-1	-1	-1 $\text{x}$
-1	1	-1 $\text{x}$
1	-1	1 $\text{x}$
1	1	-1 $\text{x}$



For label = 1: Equation for decision region

$$\hookrightarrow w_0 + \sum_{i=1}^2 w_i x_i \geq 0$$

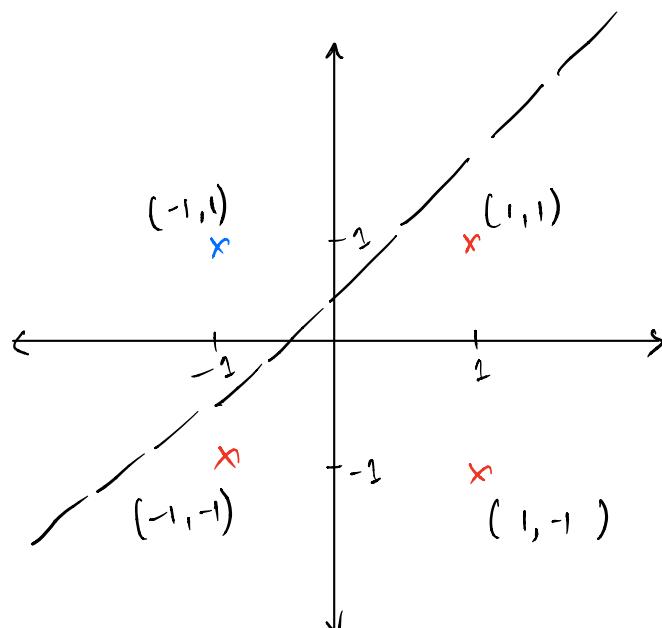
For label = -1:  $w_0 + \sum_{i=1}^2 w_i x_i \leq 0$

$$w_0 = -1.5, w_1 = 1, w_2 = -1$$

$$\text{Eq of perceptron} \rightarrow -1.5 + 1 - 1 = 0$$

(ii)  $y = \text{NOT}(x_1) \text{ AND } x_2$

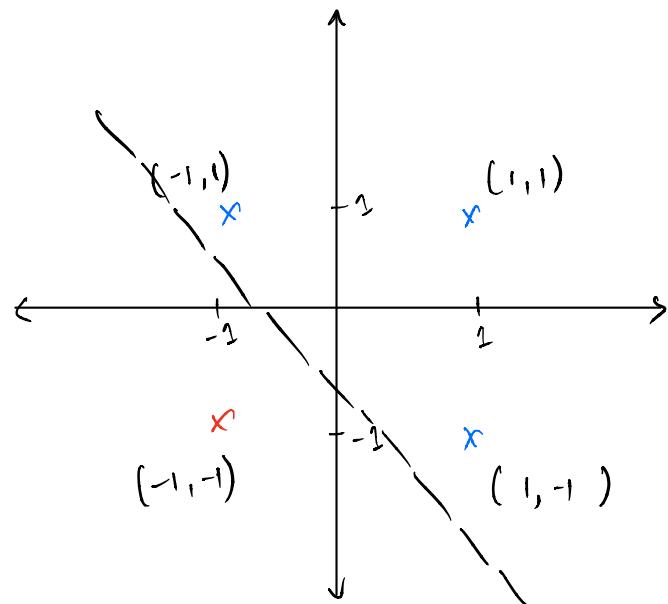
$x_1$	$x_2$	$y$
-1	-1	-1 $\times$
-1	1	1 $\times$
1	-1	-1 $\times$
1	1	-1 $\times$



<u>For label = 1</u> : $w_0 + \sum_{i=1}^2 w_i x_i \geq 0$	$w_0 = -0.5, w_1 = -1, w_2 = 1$ <u>Eq of perceptron is :</u> $-0.5 - x_1 + x_2 = 0$
--	---

(iii)  $y = x_1 \text{ OR } x_2$

$x_1$	$x_2$	$y$
-1	-1	-1 $\times$
-1	1	1 $\times$
1	-1	1 $\times$
1	1	1 $\times$

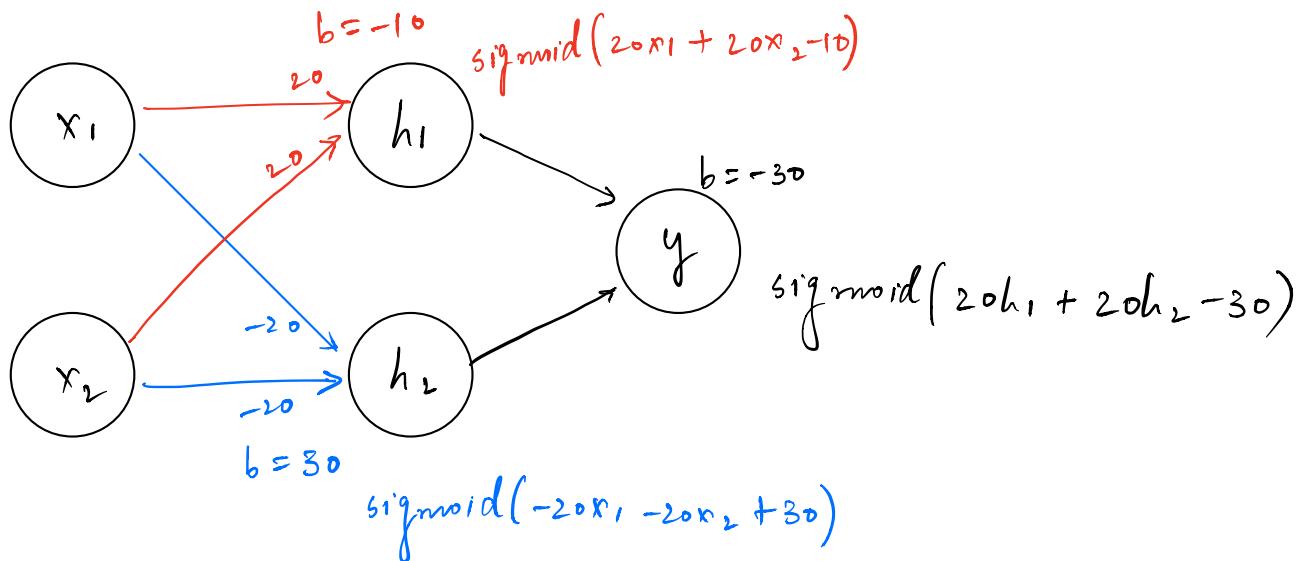


<u>For label = 1</u> : $w_0 + \sum_{i=1}^2 w_i x_i \geq 0$
--

<u>For label = -1</u> : $w_0 + \sum_{i=1}^2 w_i x_i < 0$
--

$$w_0 = 0.5, w_1 = 1, w_2 = 1, \text{ Eq of perceptron} = 0.5 + x_1 + x_2 = 0$$

c) To solve the XOR function consider the following multi-layer NN, initialized with some weights and biases.



Outputs of  $h_1$

$$\Gamma[(20)(0) + (20)(0) - 10] = 0$$

$$\Gamma[(20)(1) + (20)(1) - 10] = 1$$

$$\Gamma[(20)(0) + (20)(1) - 10] = 1$$

$$\Gamma[(20)(1) + (20)(0) - 10] = 1$$

Outputs of  $h_2$

$$\Gamma[(-20)(0) - (20)(0) + 30] = 1$$

$$\Gamma[(-20)(1) - (20)(1) + 30] = 0$$

$$\Gamma[(-20)(0) - (20)(1) + 30] = 1$$

$$\Gamma[(-20)(1) - (20)(0) + 30] = 0$$

Output at  $y$

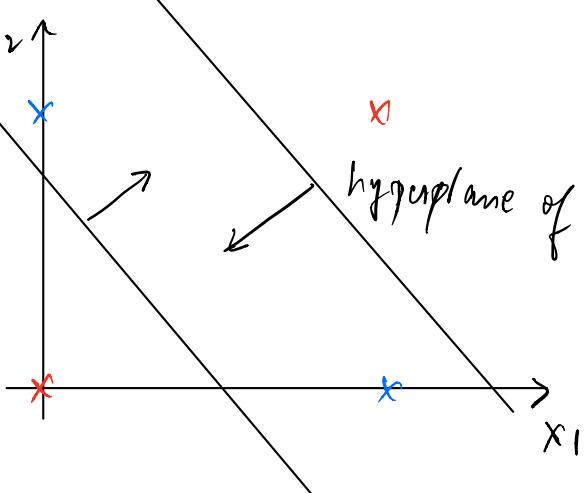
$$\Gamma[(20)(0) + (20)(1) - 30] = 0$$

$$\Gamma[(20)(1) + (20)(0) - 30] = 0$$

$$\Gamma[(20)(1) + (20)(1) - 30] = 1$$

$$\Gamma[(20)(1) + (20)(1) - 30] = 1$$

hyperplane  
of  $h_1$



(3)

3. (10 points) In Lectures 2 and 5, we derived a closed form expression for solving linear and ridge regression problems. This is great for finding linear behavior in data; however, if the data is nonlinear, just as in the classification case, we have to resort to the *kernel trick*, i.e., replace all dot products in the data space with kernel inner products. In this problem, we theoretically derive kernel regression. Suppose we are given training data  $\{(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)\}$  where each response  $y_i$  is a scalar, and each data point  $x_i$  is a vector in  $d$  dimensions.

- Assume that all data have been mapped into a higher dimensional space using the feature mapping  $x \mapsto \phi(x)$ , write down an expression for the squared error function using a linear predictor  $w$  in the high-dimensional space.
- Let  $\Phi$  be the matrix with  $n$  rows, where row  $i$  consists of the feature mapping  $\phi(x_i)$ . Write down a closed form expression for the optimal linear predictor  $w$  as a function of  $\Phi$  and  $y$ .
- For a new query data point  $z$ , the predicted value is given by  $f(z) = \langle w, \phi(z) \rangle$ . Plug in the closed form expression for  $w$  from the previous sub-problem to get an expression for  $f(z)$ .
- Suppose you are given black-box access to a kernel dot product function  $K$  where  $K(x, x') = \langle \phi(x), \phi(x') \rangle$ . Mathematically show that all the calculations in (b) and (c) can be performed by invoking the kernel dot product alone *without explicitly writing down  $\phi(x)$  ever*. You may want to use the Sherman-Morrison-Woodbury identity for matrices:

$$(A^{-1} + B^T C^{-1} B)^{-1} B^T C^{-1} = A B^T (B A B^T + C)^{-1}.$$

3a)  $L(w) = \min \sum_{i=1}^n (w^T x_i - y_i)^2$  - Loss function without Kernelization

$x$  is mapped to  $\phi(x)$  . . .  $\therefore L(w) = \min \sum_{i=1}^n (\underbrace{w^T \phi(x_i)} - y_i)^2$

↓  
The dimensions of this are  
very high. So, its hard to  
compute.

3b) closed form expression for linear regression

$$w = (X X^T)^{-1} X y$$

When expanded to higher dimensions using,  $x \rightarrow \phi(x)$

$$w = (\Phi \Phi^T)^{-1} \Phi y$$

By expressing ' $w$ ' as a linear combination of inputs of the training space, we get

$$w = \sum_{i=1}^n \alpha_i x_i = X \vec{\alpha}$$

$x \rightarrow \phi(x) \therefore x$  is replaced by  $\phi$

$$\phi \vec{\alpha} = w = (\phi \phi^T)^{-1} \phi y$$

Multiplying on both sides by  $\phi^T \phi \phi^T$

$$(\phi^T \phi)(\phi^T \phi) \vec{\alpha} = \phi^T (\phi \phi^T (\phi \phi^T)^{-1}) \phi y$$

Replace  $\phi^T \phi$  by  $K$

$$K^2 \vec{\alpha} = Ky$$

$$\boxed{\vec{\alpha} = K^{-1}y}$$

Kernel regression can be extended to Kernelized Ridge regression

where,  $\boxed{\vec{\alpha} = (K + \lambda I)^{-1}y}$

$\therefore$  The closed form expression for Kernelized Ridge regression is given by,

$$w = \phi \vec{\alpha}, \text{ where } \vec{\alpha} = (K + \lambda I)^{-1}y \quad \& \quad K = \phi \phi^T$$

$$\therefore \boxed{w = (\phi \phi^T + \lambda I)^{-1} \phi y}$$

3c) For a test point  $\tilde{z}$ ,  $f(\tilde{z}) = \langle w, \phi(\tilde{z}) \rangle$

$$f(\tilde{z}) = \langle (\phi\phi^T + \lambda I)^{-1}\phi\tilde{z}, \phi(\tilde{z}) \rangle$$

$$\boxed{f(\tilde{z}) = \phi(\tilde{z}) \phi ((\phi\phi^T + \lambda I)^{-1} \phi\tilde{z})}$$

3d) The answer obtained in '3c' is not an efficient way of computing the prediction for the label of the test point. Computing  $\phi$ ,  $\phi(\tilde{z})$  and the product of  $\phi$  &  $f(\tilde{z})$  is hard to compute. But we actually don't need to compute any of these things including  $\phi(x)$ .

The inner product  $\phi(x) \phi(\tilde{z})$  can be efficiently computed by,

$$\begin{aligned} \phi(x) \phi(\tilde{z}) &= 1 \cdot 1 + x_1 \tilde{z}_1 + x_2 \tilde{z}_2 + \dots + x_d \tilde{z}_d \\ &= \prod_{k=1}^d (1 + x_k \tilde{z}_k) \end{aligned}$$

The sum of  $2^d$  terms is converted to a product of  $d$  terms. Running time is reduced from  $O(2^d)$  to  $O(d)$ .

$$k(x_i, x_j) = \phi(x_i)^T \phi(x_j)$$

$\hookrightarrow$  Kernel function

For  $n$  training samples, all possible inner products can be precomputed and stored in a matrix  $K$ . With this kernel dot product function ' $K$ ', we can now perform a look-up to get any  $K_{ij}$  value.

$$\text{Here, } K_{ij} = \underline{\phi(x_i)^T \phi(x_j)}$$

While training in the  $\phi(x)$  space, to perform G.D we need to compute  $\langle w^T \cdot \phi(x) \rangle$ . Because 'w' can be expressed as a linear combination of the training samples,  $w = \sum_{j=1}^n \alpha_j \phi(x_j)$ .

$$\text{G.D} \rightarrow w_{t+1} = w_t - (\lambda r) \left( \frac{\partial l}{\partial w} \right)$$

$$\text{where, } \frac{\partial l}{\partial w} = 2 \left( w^T \phi(x_i) - y_i \right) x_i$$

$\frac{\partial l}{\partial w}$  becomes  $2 \left( \sum_{j=1}^n \alpha_j K_{ij} - y_i \right)$ , so the gradient update becomes

$$\alpha_i^{t+1} = \alpha_i^t - 2(\lambda r) \left( \sum_{j=1}^n \alpha_j^t K_{ij} - y_i \right)$$

So, in conclusion we never calculate  $\phi(x_i)$  we just calculate the dot product of  $\phi(x_i)$  with  $\phi(x_j)$  in  $O(d)$  time.