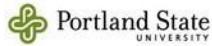




CS 445/545: Machine Learning



Midterm (49 pts. Possible),

Spring 2024 Rhodes

SRIHARI TANMAY KARTHIK TADALA

PSU ID - 918597835

Please show a sufficient amount of work for full credit on exercises. If the question prompts you to simply provide an answer, it is acceptable to simply submit the answer (without explanation, if none is necessary). The exam is open book and open notes.
However, you are not allowed to confer with fellow students or, nor are you permitted to use “Google” to seek out a solution.

*Email your exam solutions to our grader by the assigned due date. You can submit typed or hand-written solutions (or a combination of both if this is preferred); please make an effort to ensure that your solutions are clear and legible.

1. (1 pt.) Soft-margin SVMs, defined with slack variables ξ_i , always admit of a solution. **True False**

True

2. (1 pt.) A zero training set error necessarily indicates good generalization performance.
True False

False

The efficiency of a model is determined by its ability to accurately predict outcomes on new, unfamiliar data, rather than solely relying on its similarity to the data it was trained on.

3. (1 pt.) Recall the general expression for backprop weight updates:

$$\Delta w^t = \eta \delta_j x_i + \alpha \Delta w^{t-1} - \lambda w^{t-1}$$
. Explain the role of the following terms (just in a sentence):

$\alpha \Delta w^{t-1}$:

λw^{t-1} :

Solution-

The general expression for backpropagation weight updates with weight

decay (L2 regularization) can be represented as follows:

$$\Delta w^t = \eta \delta_{ji} x_i + a \Delta w^{t-1}_{ji} - \lambda w^{t-1}_{ji}$$

Where,

- $[\Delta w_j(t)]$: This represents weight updates for weight "i" at time "t". It's the change that will be applied to the weight to update its value.
- $[\eta]$: The learning rate, which controls the step size during weight updates.
- $[\partial j x_i]$: This term represents the gradient of the loss function "j" with respect to the input " x_i ".

It quantifies how much the loss changes concerning the input.

- $[a \Delta w_{ji}(t-1)]$: The term "a" represents the momentum coefficient. $\Delta w_{ji}(t-1)$ represents the change in weight at the previous time step. Momentum in the weight update takes into account the direction and velocity of previous weight changes, which aids in achieving convergence.
- $[\lambda w_{ji}(t-1)]$: The symbol " λ " denotes the weight decay or L2 regularization term. It imposes a penalty on significant weights by deducting a portion of the weight " w_{ji} " from the previous time step " $t-1$." Weight decay mitigates overfitting by promoting the regularization of weight values towards smaller magnitudes.

This formula combines the gradient of the loss with respect to the input, momentum from the previous weight update, and weight decay to calculate the weight update.

4. (1pt.) The PLA (perceptron learning algorithm) always converges in a finite number of steps. **True** **False**

False

5. (1 pt.) Given a finite data set, there exists a finite dimensional vector space for which the data is linearly separable. **True** **False**

True

6. (1 pt.) A model with infinite VC dimension can have a finite number of parameters.

True False

True

7. (1 pt.) Suppose we wish to calculate $P(H|X_1, X_2)$ and we have no conditional independence information. Which of the following sets of numbers are sufficient for the calculation? (circle)

(i) $P(X_1, X_2), P(H), P(X_1|H), P(X_2|H)$

(ii) $P(X_1, X_2), P(H), P(X_1, X_2|H)$

(iii) $P(H), P(X_1|H), P(X_2|H)$

8. (1 pt.) Consider the sigmoid function: $f(x) = \frac{1}{1+e^{-x}}$. Which expression is equal to _____

$f'(x)$? (circle)

(i) $f(x)\log(1-f(x))$

(ii) $f(x)(1-f(x))$

(iii) $\frac{1}{f(x)}$

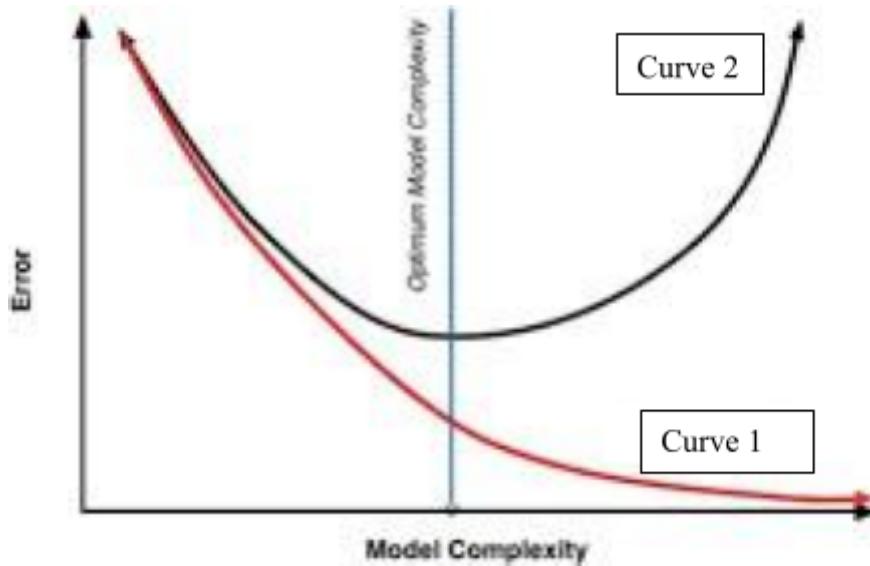
(iv) κ_0

Option 2 $f(x)(1-f(x))$

9. (1 pt.) Suppose you are given a dataset of cellular images from patients with and without cancer and that you are required to train a classifier that predicts the probability that the patient has cancer. The dataset has 900 cancer-free images and 100 images from cancer patients. If I train a classifier which achieves 95% accuracy on this dataset, should we consider this to be a good classifier? Explain your answer in a sentence or two.

Ans) Even if the classifier attains a 95% accuracy rate, it may still produce inaccurate outcomes due to the inadequate distribution of the dataset. With only 100 cancer images and 900 images of cancer-free patients, the dataset is not well-balanced.

10. (1 pt.) The plots of training and test error are shown as a function of model complexity below. Appropriately identify the plots.



Training Error: Curve 1 Curve 2 (circle)

Test Error: Curve 1 Curve 2

11. (2 pts.) Suppose that the prevalence of a disease is 2%. This disease can be screened by a medical test that is 85% accurate. This means that the test result is positive about 85% of the times when it is applied on patients who have the disease and that the test result is negative about 85% of the time when it is applied on patients who do not have the disease. Suppose that you take the test and the test shows a positive result. How likely is it that you have the disease?

- i) H_1 be the event of having disease
 H_0 be the event of not having disease
 D_+ be the event of testing positive
 D_- be the event of testing negative
Here we can find the possibility of having the disease given the positive test results which is $P(H_1/D_+)$

Probability of having disease $P(H_1) = 0.02$

Probability of not having disease $P(H_0) = 0.98$

Probability of getting positive of having the disease $P(D_+/H_1) = 0.85$

Probability of testing negative of not having a disease $P(D_-/H_0) = 0.85$.

$$\text{From Bayes Theorem: } P(H_1/D_+) = \frac{P(D_+/H_1) \times P(H_1)}{P(D_+)}$$

To find the value of $P(D_+)$

$$P(D_+) = P(D_+/H_1) \times P(H_1) + P(D_+/H_0) \times P(H_0)$$

$$\text{The value of } P(D_+/H_0) = 1 - P(D_-/H_0) = 1 - 0.85 = 0.15$$

$$\begin{aligned} \text{So } P(D_+) &= P(D_+/H_1) \times P(H_1) + P(D_+/H_0) \times P(H_0) \\ &= 0.85 \times 0.02 + 0.15 \times 0.98 \\ &= 0.184 \end{aligned}$$

$$\text{From Bayes theorem } P(H_1/D_+) = \frac{P(D_+/H_1) \times P(H_1)}{P(D_+)}$$

$$= \frac{0.85 \times 0.02}{0.184} = 0.103$$

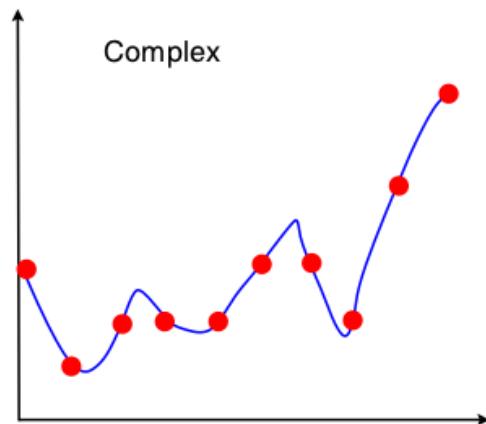
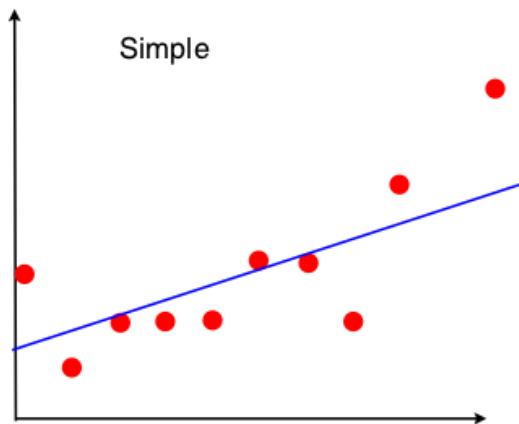
The chance of having disease is $0.103 \times 100 = 10.3\%$.
The objective is to minimize the

12. (2pts.) Recall that with LDA dimensionality reduction we wish to solve the following

optimization problem: $\arg \max_w \frac{w^T S_B w}{w^T S_w w}$. In a few simple sentences, explain the meaning of this expression with respect to LDA.

12) In LDA, the objective is to maximize the ratio $w^T S_B w / w^T S_w w$ where
 w is the weight vector
 S_B is the within class scatter matrix and
 S_w is the within class scatter matrix.
LDA's goal is to identify the projections that maintains close grouping of data within each class ($S_w(w)$) and also maximize the separation between high class S_B . The optimal w is the reduced dimensional space that encodes the difference b/w classes.

13. (2 pts.) Using the following mathematical models (i.e. the curves shown), explain briefly the idea of the “bias-variance tradeoff” in machine learning, how it relates to model *complexity* and the notions of *overfitting* and *underfitting*. In relation to these ideas, elaborate on what it means to have a “good” predictive model.



The diagram on the left depicts a model that is underfitting, meaning it does not accurately represent the trend in the data points. This model shows a high level of bias and a low level of variance. Conversely, the diagram on the right illustrates the phenomenon of overfitting, characterized by a model that exhibits low bias but high variance. This means that the model captures excessive noise from the data. An effective predictive model is one that achieves a good balance between bias and variance, resulting in minimal error. This model effectively captures the fundamental attributes of the data while also demonstrating strong performance on new, unseen data.

14. (1 pt.) “t-SNE” is an example of which type of general ML algorithm: (circle)

- (i) classification (ii) regression (iii) dimensionality reduction (iv) backpropagation

(iii) dimensionality reduction

15. (2 pts.) Let $\mathbf{x} = (x_1, x_2)$. Using the feature mapping

$$\begin{matrix} 1 & \sqrt{2} \\ 1 & 2 & 2 \end{matrix}$$

$$\Phi(\mathbf{x}) = \left(x^2, \quad \cdot x \ x, x^2 \right)$$

show that

$$\Phi((2, 3)) \cdot \Phi((4, 4)) = ((2, 3) \cdot (4, 4))$$

15. $x = (x_1, x_2)$ $y = (y_1, y_2)$
here we have, $\phi(x) = (x_1^2, \sqrt{2}x_1x_2, x_2^2)$
we need to substitute $(2, 3)$ and $(4, 4)$ in the
equation to get,

$$\begin{aligned}\phi(2, 3) \times \phi(4, 4) &= (2^2, \sqrt{2}(2)(3), 3^2) \\ &= (4, \sqrt{2}(4)(4), 4^2)\end{aligned}$$

$$\begin{aligned}\phi((2, 3) \cdot \phi(4, 4)) &= 64 + 192 + 144 \\ &= 400.\end{aligned}$$

Now, to prove it, we shall calculate the
value

$$\phi((2, 3)(4, 4))^2 = (8+12)^2 = 20^2 = 400$$

They are equal and hence proved!!!

16. (5 pts.) **Gradient Descent.** Consider the multivariate function: $f(x, y) = x^4 + y^2$

Devise an iterative rule using gradient descent that will iteratively move closer to the minimum of this function. Assume we start our search at an arbitrary point: (x_0, y_0) . Give your update rule in the conventional form for gradient descent, using η for the learning rate.

- (i) Write the explicit x-coordinate and y-coordinate updates for step $(i+1)$ in terms of the x- coordinate and y-coordinate values for the i th step.

$$x^{(i+1)} \leftarrow$$

$$y^{(i+1)} \leftarrow$$

$$16) \text{ i) Given } f(x, y) = x^4 + y^2$$

$$\frac{\partial f(x, y)}{\partial x} = 4x^3$$

$$\frac{\partial f(x, y)}{\partial y} = 2y$$

$$x^{i+1} = x^i - \eta \left(\frac{\partial f(x, y)}{\partial x} \right)^i$$

$$= x^i - 4x^3$$

$$y^{i+1} = y^i - \eta \left(\frac{\partial f(x, y)}{\partial y} \right)^i$$

$$= y^i - 2y$$

16) iv) Rewrite the rule using momentum (α)

$$x^{i+1} = x^i - \eta (4x^3)^i - \alpha [(4x^3)^i - (4x^3)]$$

$$y^{i+1} = y^i - \eta (2y)^i - \alpha [(2y)^i - (2y)^{i-1}]$$

$$x^{i+1} = x^i - \eta (4x^3)^i - \alpha [x^i - x^{i-1}]$$

$$y^{i+1} = y^i - \eta (2y)^i - \alpha [y^i - y^{i-1}]$$

(ii) Briefly explain how G.D. works, and the purpose of the learning rate.

Gradient Descent is an iterative algorithm employed to minimize a cost function by identifying the direction of steepest descent. Initially, the weights are assigned arbitrary values and subsequently adjusted through incremental changes referred to as the learning rate. The mentioned rate influences the magnitude of the increments made in each update to approach the local minimum. Augmenting the magnitude of the steps can accelerate the convergence process, although there is a potential risk of exceeding the intended objective. Conversely, a decreased learning rate results in a slower and more gradual progression towards the local minimum.

- (iii) Is your algorithm guaranteed to converge to the minimum of f (you are free to assume that the learning rate is sufficiently small)? Why or why not?

Answer: Yes, it will converge to the minimum of f if we use a small learning rate because it has only one local minima which is the global minimum.

- (iv) Re-write your rule from part (i) with a momentum term, including a momentum parameter α .

16) iv) Rewrite the rule using momentum (α)

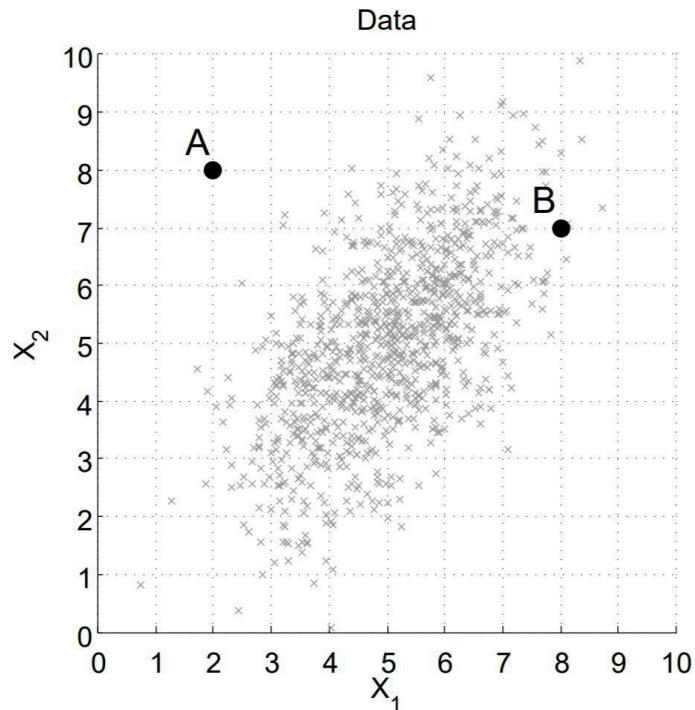
$$x^{i+1} = x^i - \eta (4x^3)^i - \alpha [(4x^3)^i - (4x^3)^{i-1}]$$

$$y^{i+1} = y^i - \eta (2y)^i - \alpha [(2y)^i - (2y)^{i-1}]$$

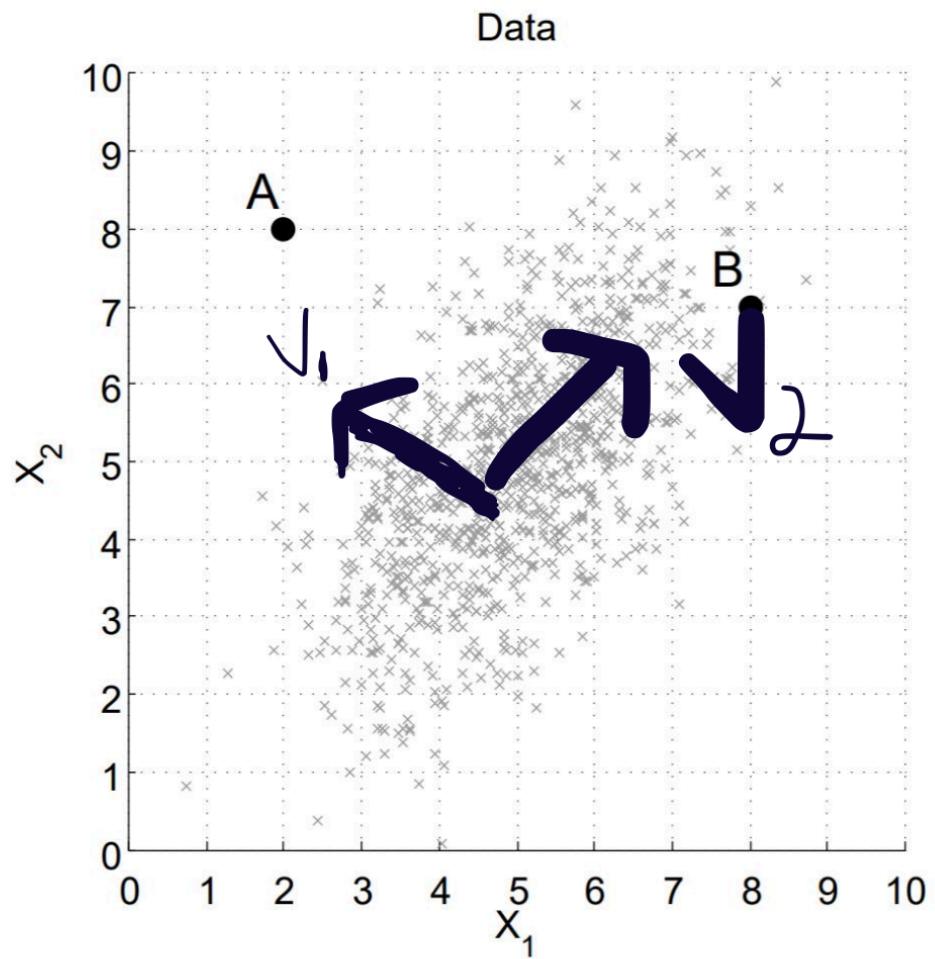
$$n^{i+1} = n^i - \eta (4n^3)^i - \alpha [n^i - n^{i-1}]$$

$$y^{i+1} = y^i - \eta (2y)^i - \alpha [y^i - y^{i-1}]$$

17. (3 pts.) PCA. The plot below shows a sample drawn from a two dimensional multivariate Normal (Gaussian) distribution. Define vectors \mathbf{v}_1 and \mathbf{v}_2 as the directions of the first and second principal components, after applying PCA to the dataset, where
- $$\|\mathbf{v}_1\| = \|\mathbf{v}_2\| = 1.$$



- (i) Sketch and label \mathbf{v}_1 and \mathbf{v}_2 in the figure above. The arrows should originate from the mean of the distribution. You do not need to compute the actual PCA procedure, instead simply visually estimate the directions of the arrows.

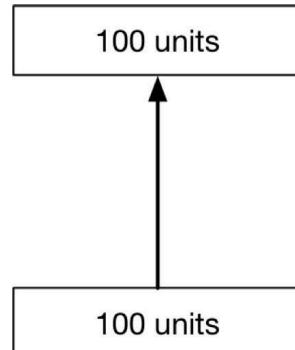


- (ii) Which point (A or B) would have the higher reconstruction error after projecting onto the first principal component direction v₁? Circle one:

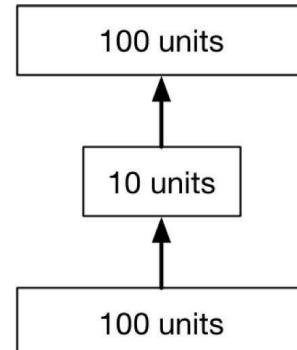
Point A Point B

Answer: Point A

18. (2 pts.) Consider the following two MLPs, where all of the layers use linear activation functions.



Network A



Network B

- (i) Give one advantage of Network A over Network B.

Implementing Network A is easier than Network B due to its smaller number of units and lower computational complexity.

- (ii) Give one advantage of Network B over Network A.

Network B has the capacity to acquire a greater number of features compared to Network A. Thus, the likelihood of overfitting is lower compared to Network A.

19. (2 pts.) Suppose that we have trained an MLP on the MNIST dataset, and pass a datum \mathbf{x} (the encoding of a handwritten digit) through our network and the corresponding output neuron activations consist of the following values:

$$\{0.10, 0.40, 0.45, 0.10, 0.20, 0.90, 0.05, 0.20, 0.10, 0.35\}$$

Apply the softmax transformation to this set of activations and list the transformed values (preserving the order of the original activations).

In practice, what is the purpose of using the softmax transformation?

(3)

19) We know that

$$y_n = g(h_n) = \frac{\exp(h_n)}{\sum_{k=1}^N \exp(h_k)}$$

$$\{0.10, 0.40, 0.45, 0.10, 0.20, 0.90, 0.05, 0.20, \\ 0.10, 0.25\}$$

$$\exp(0.10) = 1.1052$$

$$\exp(0.40) = 1.4918$$

$$\exp(0.45) = 1.5683$$

$$\exp(0.10) = 1.1052$$

$$\exp(0.20) = 1.2214$$

$$\exp(0.30) = 2.4596$$

$$\exp(0.05) = 1.0513$$

$$\exp(0.20) = 1.2214$$

$$\exp(0.10) = 1.1052$$

$$\exp(0.35) = 1.4191$$

$$\sum_{k=1}^N \exp(h_k) \rightarrow \frac{\exp[0.10 + 0.40 + 0.45 + 0.10 + 0.20 + \\ 0.98 + 0.05 + 0.20]}{\exp[0.10] + \exp[0.35]}$$

$$\text{softmax}(h(x)) = \left[\frac{1.1052}{13.748}, \frac{1.4918}{13.748}, \frac{1.5683}{13.748}, \right. \\ \left. \frac{1.1052}{13.748}, \frac{1.2214}{13.748}, \frac{1.0513}{13.748}, \frac{1.2214}{13.748}, \frac{1.1052}{13.748}, \frac{1.4191}{13.748} \right]$$

$$\text{Softmax}(h(n)) = \begin{bmatrix} 0.0804, 0.1085, 0.1140, 0.0804, 0.0888, \\ 0.1789, 0.0765, 0.0888, 0.0804, 0.103 \end{bmatrix}$$

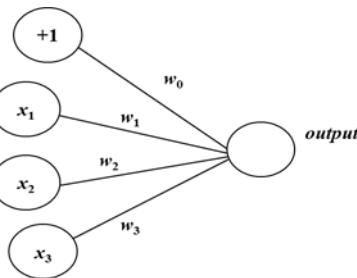
The softmax function rescales the outputs by calculating the exponential of the inputs to the neuron.

It is basically used for normalizing values which helps in predicting the probability

of each class

20. (4 pts) Consider the training set, perceptron, and initial weights given below.

x_1	x_2	x_3	Class
.1	-.2	.2	+1
.4	.4	.1	+1
-.3	.1	.3	-1
-.2	.2	-.1	-1



Initial Weights:

$$w_0 = .1$$

$$w_1 = .1$$

$$w_2 = .1$$

$$w_3 = .1$$

- (i) Using stochastic gradient descent (i.e., update the weights after processing each training example), and learning rate $\eta = 0.1$, simulate, by hand, the perceptron learning algorithm for one epoch, showing the new values of the weights after completing the epoch.

2b) Given $w_0 = 0.1, w_1 = 0.1, w_2 = 0.1, w_3 = 0.1$
 learning rate (η) = 0.1 for 1 epoch
 for $x_1 = 0.1, x_2 = -0.2, x_3 = 0.2, b = 1 (t^k = +1)$

$$y^k = \sum w_i x_i + b$$

$$= (0.1 \times 0.1) + (0.1 \times -0.2) + 0.1(0.2) + 1(0.1)$$

$$= 0.01 - 0.02 + 0.02 + 0.1 = 0.11$$

$y^k > 0$, hence output +1

$y^k = t^k$ here no weight change

For $x_1 = 0.4, x_2 = 0.4, x_3 = 0.1, b = 1 (t^k = +1)$

$$y^k = \sum w_i x_i + b$$

$$= 0.1(0.4) + 0.1(0.4) + 0.1(0.1) + 0.1(1)$$

$$= 0.19 > 0 \text{ (hence output is +1)}$$

④

$$y^k = t^k \text{ (weight change not required here b/c)} \quad 0$$

$$\text{For } x_1 = 0.3, x_2 = 0.1, x_3 = 0.3, b = -1 \quad (t^k = -1) \quad 2$$

$$y^k = \sum w_i x_i + b$$

$$= 0.1(0.3) + 0.1(0.1) + 0.1(0.3) + 0.1$$

$$= 0.16 > 0$$

perception fires +1

$$y^{k+1}, \text{ update weights}$$

$$\Delta w = 0.1(-1 - 1)(1) = 0.1(-2)$$

$$\boxed{\Delta w_0 = -2}$$

$$\Delta w_1 = (0.1)(-2)(-0.3) = -0.06$$

$$\Delta w_2 = (0.1)(-1)(0.1) = -0.02$$

$$\Delta w_3 = (0.1)(-1)(0.3) = -0.06$$

New weights are

$$w_0 = w_0 + \Delta w_0 = 0.1 - 0.2 = -0.1$$

$$w_1 = w_1 + \Delta w_1 = 0.1 + 0.06 = 0.16$$

$$w_2 = w_2 + \Delta w_2 = 0.1 - 0.02 = 0.08$$

$$w_3 = w_3 + \Delta w_3 = 0.1 - 0.06 = 0.04$$

Now consider the y^k output

$$x_1 = 0.2, x_2 = 0.2, x_3 = 0.1, b = 1 \quad (t^k = -1)$$

$$y^k = \sum w_i x_i + b$$

$$(0.16)(-0.02) + (0.08)(0.2) + (0.04)(-0.1) = -0.12$$

$$- (-0.1)(1)$$

$$= 0.12 > 0$$

perceptron fires -1 $y^k = f_k$ no weight change

- (ii) At the end of the epoch, how would your resulting perceptron classify a new example with $x_1 = 0.8, x_2 = -0.2, x_3 = -0.1$?

$$\text{2011)} \quad x_1 = 0.8, x_2 = -0.2, x_3 = -0.1 \quad y^k = \sum w_i x_i + b$$

$$= 0.16(0.8) + (0.08)(-0.2) + 0.04(-0.1)$$

$$= 0.128 - 0.016 - 0.004 - 0.1$$

$$= 0.0008 > 0$$

perceptron fires]

As the value is greater than zero the perceptron will return 1 as the class

21. (2 pts.) Suppose you have the following short DNA sequences in your training set:

$$s_1 = \text{ACGGT} \quad s_2 = \text{ATTGT} \quad s_3 = \text{CGCCT}$$

Suppose you are using these to train an SVM with a "match-count" kernel, where $k(s_i, s_j)$ returns the number of locations strings s_1 and s_2 at which the symbols match.

Give the Kernel matrix K for this training set and kernel function.

21) Given $s_1 = \text{ACGAT}$
 $s_2 = \text{ATTGT}$
 $s_3 = \text{CGCTT}$

Now the kernel function $K(s_i, s_j) = \sum (s_i, s_j)$

$$\begin{bmatrix} K(s_1, s_1) & K(s_1, s_2) & K(s_1, s_3) \\ K(s_2, s_1) & K(s_2, s_2) & K(s_2, s_3) \\ K(s_3, s_1) & K(s_3, s_2) & K(s_3, s_3) \end{bmatrix}$$

$$= \begin{bmatrix} 5 & 3 & 1 \\ 3 & 5 & 1 \\ 1 & 1 & 5 \end{bmatrix}$$

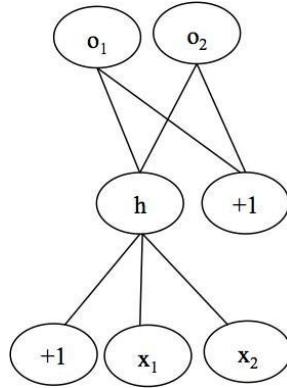
$s_1, s_1 = 5$
 $s_1, s_2 = 3$
 $s_1, s_3 = 1$
 $s_2, s_1 = 3$
 $s_2, s_2 = 5$
 $s_2, s_3 = 1$
 $s_3, s_1 = 1$
 $s_3, s_2 = 1$
 $s_3, s_3 = 5$

21. (2 pts.) Explain in a few sentences why kernel functions allow SVMs to classify non-linearly separable problems in an efficient way.

Kernel functions enable SVMs to efficiently classify non-linearly separable problems by transforming the dataset into a higher-dimensional space. In this transformed space, the previously inseparable classes become linearly separable, allowing for the use of a linear

classifier. The effectiveness of this approach heavily depends on the selection of the kernel function, as it determines the characteristics of the high-dimensional space in which the data is organized, thereby influencing the performance of the SVM.

22. (4 pts.) Consider the multilayer neural network given below.



Suppose all the weights are initialized to 0.1. Assume the sigmoid activation function for hidden and output nodes:

$$o = \sigma(\mathbf{w} \cdot \mathbf{x}), \text{ where } \sigma(z) = \frac{1}{1 + e^{-z}}.$$

- (a) Given input $\mathbf{x} = (3, 2)$, what is the activation at output node o_1 ?

⑤

22)B) 22) a)

Given sigmoid function $O = \sigma(w \cdot x)$ and $\sigma(z) = \frac{1}{1+e^{-z}}$

and all the weights are initialized to 0.1

Given input is (3, 2) and now we calculate the
weights hidden node $a = (0.1)(3) + (0.1)(2) + 0.1(1)$
 $= 0.3 + 0.2 + 0.1 = 0.6$

Now sigmoid function $= \frac{1}{1+e^{-z}} \times (0.6) = \frac{1}{1+e^{-0.6}}$

WKT the value $e = 2.718$

So the $\frac{1}{1+2.718^{0.6}} = 0.6457$

As the value we got greater than 0, so the output
would be 1

Now we calculate the value from input to
the output (0, 1)

$\sum w_i x_i + b = 0.1(1) + 0.1 = 0.2$

Now $O = \frac{1}{1+e^{-0.2}} = \frac{1}{1+(2.718)^{-0.2}} = 0.5498$

which is > 0, so output returns 1

(b) Recall that the weight update rule for backpropagation is

$$\Delta w_{ji} = \eta \delta_j x_{ji} + \text{momentum-term}$$

Let $\eta = 1$ and momentum = 0.

Suppose that after the weights are initialized to 0.1 and the input in part (a) is given, the error term for o_1 is calculated to be $\delta_{o1} = 0.2$.

What is the new value of w_{o1h} , the weight from the hidden unit to o_1 ?

22)b) Given weights = 0.1, n=1, momentum = 0
 $Dw_{ji} = n \delta_j \gamma_{ij} + \text{momentum}$
 $w_{oh} = w_{o1h} + Dw_{o1h} + \alpha(\text{momentum})$
 $= 0.1 + [0](0.2)(0.65) + 0$
 $= 0.1 + 0.13 + 0 = 0.23$

23. (3 pts.) **Definition.** A real-valued symmetric matrix A is positive semi-definite if:

$$\mathbf{x}^T A \mathbf{x} \geq 0 \text{ for all real-valued vectors } \mathbf{x} \neq \mathbf{0}.$$

(i) Let's consider a simple matrix as an example; let

$$A = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}. \text{ Prove that } A \text{ in this case}$$

is positive semi-definite. In other words, you need to show that the necessary condition holds in all cases.

23) Given $x^T A x \geq 0$ $A = \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix}$

Now assume $x = \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$ $x^T = (x_1 \ x_2)$

Now $x^T A x = \begin{bmatrix} x_1 \ x_2 \end{bmatrix} \begin{bmatrix} 1 & -1 \\ -1 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$= [x_1 \ x_2] \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}$

$= [x_1^2 - x_1 x_2 - x_1 x_2 + x_2^2]$

$= [x_1^2 - 2x_1 x_2 + x_2^2] \geq 0$

As A is positive semi-definite

(ii) Briefly, explain the significance of positive semi-definite matrices in relation to the *Mercer's Theorem* in ML.

Positive semi-definite matrices play a critical role in machine learning, as emphasized by Mercer's Theorem. This theorem is fundamental for assessing the appropriateness of kernel functions in support vector machines (SVMs) and other related applications. According to Mercer's Theorem, a kernel function is deemed suitable for use if the resulting kernel matrix is positive semi-definite. The condition is crucial as it guarantees the mathematical integrity

of the kernel method, enabling the creation of dependable and efficient machine learning models. The effectiveness of the kernel method, which is essential for producing accurate machine learning outcomes, depends on the utilization of specific matrices to ensure the validity of the underlying mathematics.

24. (3 pts.) Recall that the classification formula used by support vector machines is

$$h(\mathbf{x}) = \text{sgn} \left| \sum_{k=1}^m a_k (\mathbf{v}_k \cdot \mathbf{x}) + b \right|$$

for support vectors $\{\mathbf{v}_1, \mathbf{v}_2, \dots, \mathbf{v}_m\}$.

Let $\mathbf{v}_1 = (0, -1)$, $\mathbf{v}_2 = (1, 1)$, and $\mathbf{v}_3 = (-1, 0)$ be three support vectors defining a separating line, and let the corresponding coefficients and bias be $a_1 = -0.2$, $a_2 = 0.4$, $a_3 = -0.2$, and $bias = -0.4$.

(i) Using the formula for h above, give the class the h would assign to $\mathbf{x} = (-3, 1)$. Show your work.

$$\begin{aligned}
 h(n) &= \text{sigmoid} \left[\sum_{k=1}^m a_k (\mathbf{v}_k \cdot \mathbf{n}) + b \right] \\
 \text{given us } \mathbf{v}_1 &= (0, -1), \mathbf{v}_2 = (1, 1), \mathbf{v}_3 = (-1, 0) \\
 a_1 &= -0.2, a_2 = 0.4, a_3 = -0.2, bias = -0.4 \\
 \text{For } \mathbf{x} &= (-3, 1) \\
 h(n) &= \text{sgmd} \left(\sum_{k=1}^m a_k (\mathbf{v}_k \cdot \mathbf{x}) + b \right) \\
 &= \text{sgmd} \left(-0.2[(0)(-1)](-3, 1) + 0.4[(1, 1)(-3, 1)] \right. \\
 &\quad \left. -0.2[(-1, 0)(-3, 1)] + (-0.4) \right) \\
 &= \text{sgmd} \left[-0.2(-1) + 0.4(-3) - 0.2[3 + 0 - 0.4] \right] \\
 &= \text{sgmd} \left[-0.2(-1) + 0.4(-3) - 0.2(2.6) \right] \\
 \mathbf{x} = (-3, 1) &\text{ is assigned to class -1}
 \end{aligned}$$

(ii) Letting $\mathbf{x} = (x_1, x_2)$, give the equation of the separating line in the form $x_2 = m x_1 + b$, where m is the slope of the line and b is the vertical-axis intercept. (Note, here, b denotes the vertical-axis intercept, **not** the bias). Sketch the line described by this equation.

⑦

24) ii) Let $x = (x_1, x_2)$ $x_2 = mx_1 + b$

$$\begin{aligned} \text{WKT } w &= \sum x_i n_i \\ &= (-0.2)[(0, -1)] + 0.4[(1, 1)] \\ &= -0.2[-1/10] \\ &= [0 + 0.4 + 0.2], [0.2 + 0.4 + 0] = [0.6, 0.6] \end{aligned}$$

from the above points

$$\begin{aligned} w_1 x_1 + b + w_2 x_2 &= 0 \\ 0.6x + 0.6x_2 - 0.4 &= 0 \\ x_1 + x_2 &= \frac{0.4}{0.6} = \frac{2}{3} \end{aligned}$$

The graph for above equation $(x_1 = 0, x_2 = 0.66)$
 $x_2 = 0, x_1 = 0.66$

