

# CSCI 5521: Machine Learning Fundamentals

## (Spring 2022)<sup>1</sup>

### Final Exam

Due on Gradescope by 12pm, May 11

**Instructions:**

- The final exam has 5+1 questions, 100+2 points, on 8 pages, including one extra credit problem worth 2 points.
- Please write your name & ID on this cover page.
- For full credit, show how you arrive at your answers.

1. (30 points) In I-IV, select the correct option(s) (it is not necessary to explain).

I. Select all the option(s) that are **loss functions**:

- (a) Cross entropy (b) Maximum likelihood (c) Squared error (d) 0-1 loss (e) Hinge loss

II. Select all the option(s) that are **activation functions**:

- (a) Sigmoid (b) ReLu (c) Hyperbolic tangent (d) Softmax (e) Entropy

III. Select all the option(s) that are true when describing **kernel methods**:

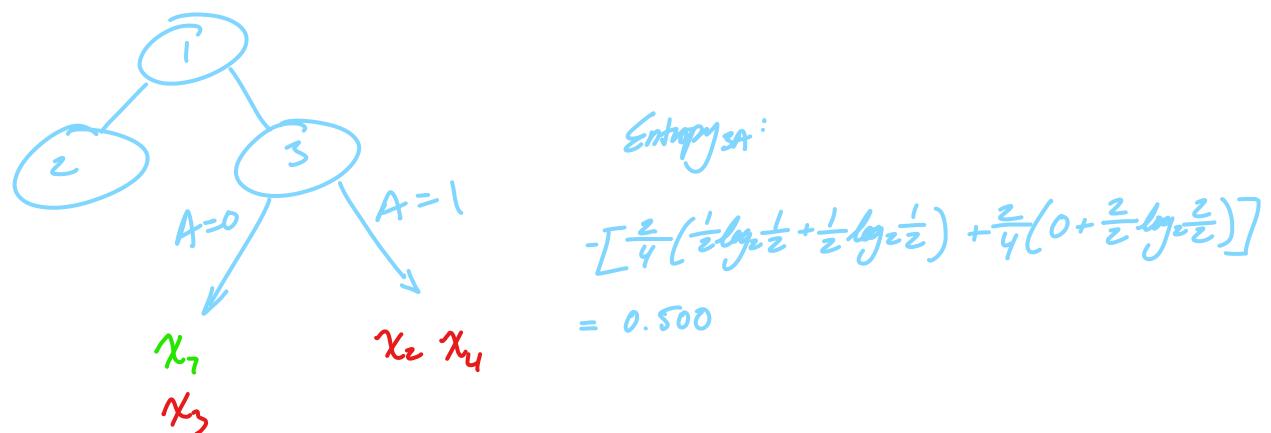
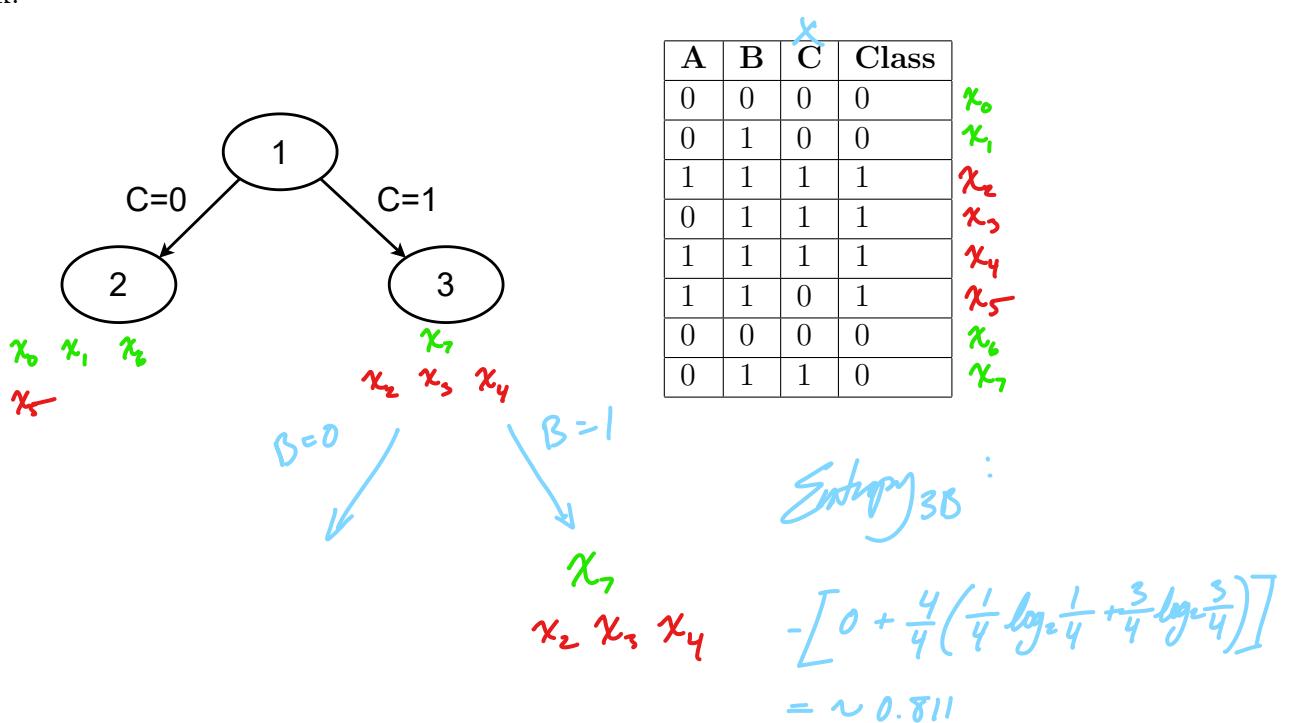
- (a) Kernel methods are designed to reduce overfitting.  
 (b) Kernel methods only work with Support Vector Machines (SVMs).  
 (c) Common kernels include polynomial, Gaussian, and Lagrange functions.  
 (d) Kernel methods are designed to map data into better representational space.  
 (e) Linear SVMs always work better than kernel SVMs.

IV. Select all the option(s) that are true about **overfitting**:

- (a) A bigger number of hidden nodes in a Multilayer Perceptron helps reduce overfitting with the same amount of data.  
 (b) A smaller number of layers in a Multilayer Perceptron helps reduce overfitting with the same amount of data.  
 (c) Pruning helps control overfitting in decision trees.  
 (d) Random forest reduces overfitting compared with decision tree.  
 (e) Data augmentation (e.g., creating copies of image data that are rotated or scaled versions of the original ones) helps reduce overfitting with the same amount of original data.

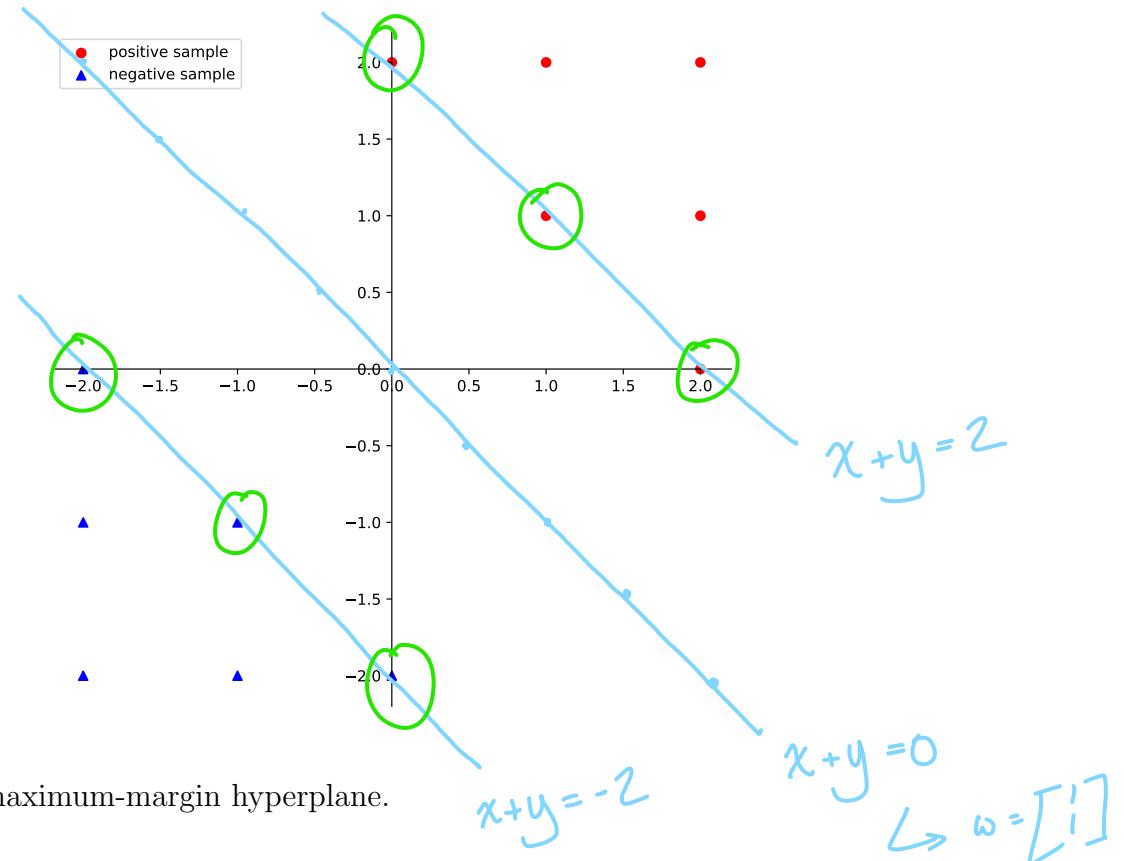
<sup>1</sup>Instructor: Catherine Qi Zhao. TA: Prithvi Raj Botcha, Shi Chen, Suzie Hoops, James Yang, Yifeng Zhang. Email: csci5521.s2022@gmail.com

2. (10 points) Given the decision tree in the figure below, the node 1 was split using feature C. Now suppose we wish to split node 3. What is the feature that you will be using to split? Show your work.



→ ∵  $\therefore$  Since  $\text{Entropy}_{3B} > \text{Entropy}_{3A}$  ( $\sim 0.811 > 0.500$ ), we would split node 3 by feature A

3. (15 points) Suppose we are training a linear SVM on a tiny dataset of 12 points shown in the figure below. Samples with positive labels are  $(1, 1), (2, 2), (1,2), (2,1), (2,0), (0,2)$  (denoted as red dots) and samples with negative labels are  $(-1, -1), (-2, -2), (-1, -2), (-2,-1), (-2,0), (0,-2)$  (denoted as blue triangles).



- (c) Pick one positive and one negative sample, and calculate their distances to the hyperplane.

$$\bullet : \frac{|1+1|}{\sqrt{1^2 + 1^2}} = \boxed{\frac{2}{\sqrt{2}}}$$

$$\blacktriangle : \frac{|-1+(-1)|}{\sqrt{1^2 + 1^2}} = \boxed{\frac{2}{\sqrt{2}}}$$

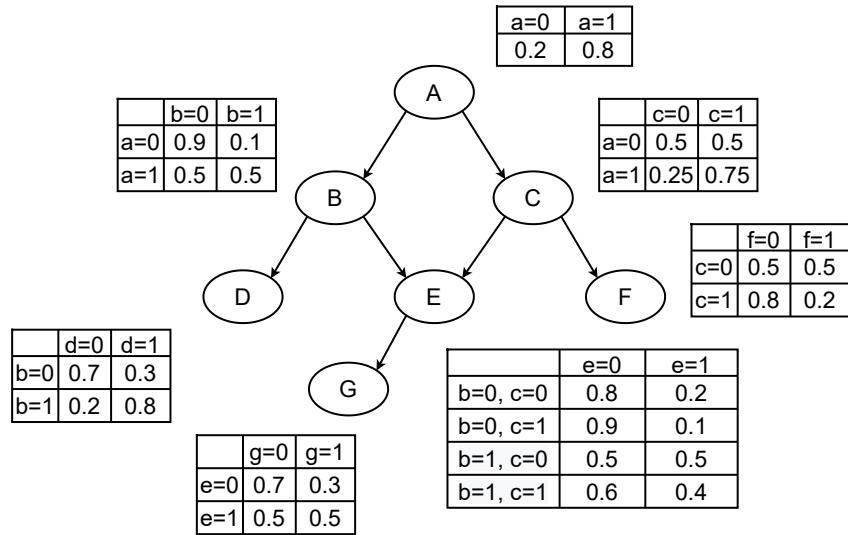
\*  $g(x) = \omega^T x$   
 $\& \text{dist} = \frac{|g(x)|}{\|\omega\|}$

- (d) If a new sample  $(-1, 1)$  comes as a negative sample on top of the original 12 points, select all the option(s) that are true:

- i. The decision boundary would change.
- ii. Kernel SVM would be a good option with the new data.

→ b/c it's linearly separable w/in the boundaries  
 → No, b/c the new data's still linearly separable.

4. (25 points) Consider the Bayesian Network below:



**Note:** The numerical values of the probabilities are for part (e). You do not need to use them for (a)-(d).

- (a) Find the joint probability  $P(A, B, C, D, E, F, G)$  as the product of conditional probabilities, according to the graphical model given above.

$$P(A, B, C, D, E, F, G) = \overbrace{P(A)P(B|A)P(C|A)P(D|B)P(E|B,C)P(F|C)P(G|E)}^{\text{Chain Rule}}$$

- (b) List all conditional independence given the graph.

$$\begin{aligned} & P(A)P(B|A)P(C|A,B)P(D|A,B,C)P(E|A,B,C,D)P(F|A,B,C,D,E)P(G|A,B,C,D,E,F) \\ & \quad \downarrow \qquad \downarrow \qquad \downarrow \qquad \downarrow \qquad \downarrow \qquad \downarrow \\ & C \perp\!\!\!\perp B | A \qquad \qquad D \perp\!\!\!\perp \{A, C\} | B \qquad E \perp\!\!\!\perp \{A, D\} | \{B, C\} \qquad F \perp\!\!\!\perp \{A, B, D, E\} | C \qquad G \perp\!\!\!\perp \{A, B, C, D, F\} | E \end{aligned}$$

*Chain Rule*

- (c) Show how to find the conditional probability  $P(A|C)$ .

$$\begin{aligned} P(A|C) &= \frac{P(C|A)P(A)}{P(C)} \quad \text{Baye's Theorem} \\ &= \frac{P(C|A)P(A)}{P(C|A)P(A) + P(C|\neg A)P(\neg A)} \end{aligned}$$

(d) Show how to find the marginal probability  $P(A, D)$ .

$$P(A, D) \Rightarrow \sum_B P(A, B, D) = \sum_B P(A) P(B|A) P(D|B)$$

$$= \boxed{P(A) \sum_B P(B|A) P(D|B)}$$

(e) Using the conditional probability distribution (CPD) tables in the figure, find:

i.  $P(a = 1|c = 0)$

$$\frac{P(c=0|a=1)P(a=1)}{P(c=0|a=1)P(a=1) + P(c=0|\neg a=1)P(\neg a=1)} = \frac{0.25 * 0.8}{(0.25 * 0.8) + (0.5 * 0.2)}$$

$$= \boxed{\frac{2}{3} \approx 0.667}$$

ii.  $P(a = 1, d = 0)$

$$P(A=1) \sum_B P(B|A=1) P(D=0|B) = P(A=1) (P(B=0|A=1)P(D=0|B=0) + P(B=1|A=1)P(D=0|B=1))$$

$$= 0.8 * (0.5 * 0.7 + 0.5 * 0.2)$$

$$= \boxed{\frac{9}{25} \approx 0.36}$$

iii.  $P(a = 1, b = 0, c = 0, d = 0, e = 0, f = 0, g = 1)$

$$P(A=1)P(B=0|A=1)P(C=0|A=1)P(D=0|B=0)P(E=0|B=0, C=0)P(F=0|C=0)P(G=1|E=0)$$

$$= \frac{0.8 * 0.5 * 0.25 * 0.7 * 0.8 * 0.5 * 0.3}{A \quad B \quad C \quad D \quad E \quad F \quad G}$$

$$= \boxed{0.0094 \approx \frac{21}{2500}}$$

iv.  $P(b = 0, c = 0, e = 0, f = 0, g = 1|a = 1, d = 0)$

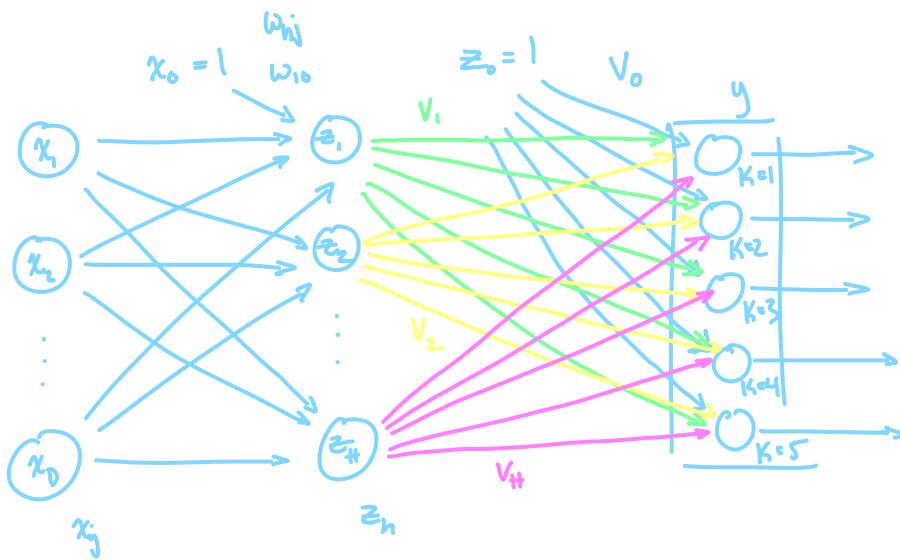
$$\frac{P(A=1, B=0, C=0, D=0, E=0, F=0, G=1)}{P(A=1, D=0)} = \frac{0.0084}{0.36} = \frac{\frac{21}{2500}}{\frac{9}{25}}$$

$$\approx \boxed{0.233 = \frac{7}{300}}$$

5. (20 points) Consider a Multi-layer Perceptron (MLP) for multi-class classification of  $K=5$  categories with 5 output units, where each hidden unit uses a hyperbolic tangent function such that  $z_h^t = \tanh(\sum_{j=1}^D w_{hj}x_j^t + w_{h0})$ . The output unit uses a softmax activation function such that  $y_i^t = \frac{\exp(\sum_h v_{ih}z_h^t + v_{i0})}{\sum_j \exp(\sum_h v_{jh}z_h^t + v_{j0})}$ . The error function is given below:

$$E(W, V | X) = - \underbrace{\sum_{t=1}^N \sum_{i=1}^K r_i^t \log y_i^t}_{E_1} + \underbrace{\frac{\lambda}{2} \sum_{h=1}^H \|w_h\|_2^2}_{E_2} + \underbrace{\frac{\sigma}{2} \sum_{k=1}^K \|v_k\|_2^2}_{E_3} \quad (1)$$

- (a) Draw the Multi-layer Perceptron showing: input values  $x_0 \dots x_D$ , output of the hidden units  $z_0 \dots z_H$ , weights  $W$  and  $V$ , and the outputs.



(b) Derive the Forward Step equation.

$$z_h^t = \tanh\left(\sum_{j=1}^D w_{hj} x_j^t + w_{h0}\right)$$

$$y_i^t = \text{softmax}\left(\frac{\exp(\sum_h v_{ih} z_h^t + v_{i0})}{\sum_j \exp(\sum_h v_{jh} z_h^t + v_{j0})}\right)$$

(c) Derive the Backward Step equation for  $w_{hj}$ .**Hint:**

- $\tanh'(x) = 1 - \tanh^2(x)$
- Given the softmax function  $f(\alpha_i) = \frac{\exp(\alpha_i)}{\sum_j \exp(\alpha_j)}$ , then  $\frac{\partial f(\alpha_i)}{\partial \alpha_j} = f(\alpha_i)(\delta_{ij} - f(\alpha_j))$ , in which  $\delta_{ij}$  is an indicator function, such that  $\delta_{ij} = 1$  if  $i = j$ , and 0 otherwise.

$$\frac{\partial E_t}{\partial w_{hj}} = \frac{\partial E_t}{\partial y_i^t} \frac{\partial y_i^t}{\partial z_h^t} \frac{\partial z_h^t}{\partial w_{hj}}$$

$= -\sum_t^N \sum_i^K \left( \frac{r_t^i}{y_i^t} * (y_i^t(1-y_i^t)) * (1 - \tanh^2(\sum_{j=1}^D w_{hj} x_j^t + w_{h0})) \right)$   
for  $i=j$   
 $= -\sum_t^N \sum_i^K \left( \frac{r_t^i}{y_i^t} * (-y_i^t y_i^t) * (1 - \tanh^2(\sum_{j=1}^D w_{hj} x_j^t + w_{h0})) \right)$   
for  $i \neq j$

$$\frac{\partial E_t}{\partial w_{hj}} = [\lambda w_{hj}]$$

$$\frac{\partial E_t}{\partial w_{hj}} = [0]$$

$$\therefore \Delta w_{hj} = -\gamma \frac{\partial E_t}{\partial w_{hj}} = -\left[ -\sum_t^N \sum_i^K \frac{r_t^i}{y_i^t} * (y_i^t(\delta_{ij} - y_i^t)) * (1 - \tanh^2(\sum_j^D w_{hj} x_j^t + w_{h0})) + \lambda w_{hj} \right]$$

\* converge  $i=j$  &  $i \neq j$  assumptions into  $\delta_{ij}$  b/c I have no more space

6. (**2 points, extra credit**) Charlotte is a machine learning scientist in a self-driving company. The company is developing an obstacle detection algorithm. Charlotte's team has access to over 1,000,000 images each with label, and the data are in the form of raw images without any features extracted.

- (a) Suggest a machine learning method she could use, and explain your suggestion.

A machine learning method that she could use would be the Graphical Model. The reason being so is because the relationships of the pixels between themselves could be modeled as a graph -- with similar and/or related pixels having stronger values between each other, holding as well for the opposite assumption. Since each image has a graph, the system could learn off of those like features and given labels, it would be able to produce appropriate outputs for each image.

- (b) What if only 1,000 out of the 1,000,000 images have labels? Suggest a method and explain your suggestion.

A method would be to develop some sort of layered, graphical network. Similar with answer (a), we add the extra layer of building graphs from the 1,000 images with labels to learn and predict the other 999,000 images to then be able to produce appropriate outputs during the validation stage. By just learning the 1,000 images and not incorporating the other 999,000, we risk underfitting. But by also using those 999k during training and validation(?) stages, we could increase the probability of a successful probability in the final stages.