

1?  $k(x, y) = \frac{1}{\sqrt{1 + \|x - y\|^2}}$

# Theoretical Assignment 3

Intro to Machine Learning, TAU

December 17, 2024

## Question 1

The Hyperbolic Tangent kernel is defined as  $k(x, y) = \tanh(x^T y + \gamma)$ . Where  $\gamma \in \mathbb{R}$  is some constant scalar.

Show that this is not a valid kernel according to Mercer's theorem.

Hint: Develop the conditions for a set of 2 samples:  $S = \{x_1, x_2\}$ .

## Question 2

The hinge loss  $\max(0, 1 - y(w^T x + b))$  is commonly used in SVMs. However, consider an alternative quadratic loss:

$$L_i(w, b) = [\max(0, 1 - y_i(w^T x_i + b))]^2.$$

- (a) In soft-SVM, what do the slack variables represent? What is the geometrical meaning of the different value ranges that the slack variables can have? How are the slack variables related to the hinge loss?
- (b) Write the primal optimization problem for a soft-margin SVM using this new quadratic loss. Hint: Consider about your answers to section a to understand how it would be different from the primal problem with the hinge loss.
- (c) Write the Lagrangian, KKT conditions, and the dual problem for a soft-margin SVM using this new quadratic loss.
- (d) Derive the gradient of the loss function  $L_i(w, b)$  with respect to  $w$  and  $b$ .
- (e) Discuss the implications of using this loss function instead of the hinge loss. How might this impact the sparsity of support vectors and the generalization performance?

## Question 3

- (a) What is the sample complexity of hard margin SVM in the linear case (without a kernel) and why? Explain your answer, don't just quote from class.
- (b) Briefly explain how will using an RBF kernel affect the sample complexity. For an SVM using the RBF kernel, discuss how the sample complexity is influenced by the kernel hyperparameter  $\sigma$ , which controls the spread of the RBF kernel.

## Question 4

Given a dataset with two classes,  $C_1$  and  $C_2$ , and the following class distributions:

Feature Value	Class $C_1$ Count	Class $C_2$ Count
$A$	8	2
$B$	3	7

- (a) Compute the entropy  $H(S)$  of the entire dataset and for the subsets  $S_A$   $S_B$  (the subsets corresponding to feature values  $A$  and  $B$ ).
- (b) Calculate the information gain  $G(S, A)$  when splitting the dataset on the feature value.
- (c) Interpret your results: What does the information gain tell you about the usefulness of splitting on this feature?
- (d) Compute the Gini impurity for the entire dataset  $S$ , and for the subsets  $S_A$  and  $S_B$  **before** splitting.
- (e) Calculate the weighted average Gini impurity **after** splitting the dataset by the feature values  $A$  and  $B$ .
- (f) Compare the results from entropy-based information gain and Gini impurity. Are the rankings of the splits the same? Explain why or why not.

## Question 5

Consider a binary decision tree used for binary classification, such as the ones we've seen in class, where each internal node splits the data based on a single feature threshold, and each leaf represents a class label. Suppose the dataset has  $n$  features, and each feature can take on real values.

- (a) Prove that a binary decision tree with depth  $d$  can divide the feature space into at most  $2^d$  regions.
- (b) Explain how this result relates to the expressiveness of decision trees and their capacity to overfit.
- (c) Using this result, determine the minimum depth  $d$  required for a decision tree to perfectly classify a dataset of  $m$  unique points (no two points share the same feature values).
- (d) Based on the previous sections, what is the VC dimension of a tree with depth  $d$ ?

# Question 1

The Hyperbolic Tangent kernel is defined as  $k(x, y) = \tanh(x^T y + \gamma)$ . Where  $\gamma \in \mathbb{R}$  is some constant scalar.

Show that this is not a valid kernel according to Mercer's theorem.

Hint: Develop the conditions for a set of 2 samples:  $S = \{x_1, x_2\}$ .

$$m=2$$

$$S = \{x_1, x_2\} \quad x_{1,2} \in \mathbb{R}^d$$

$$k_{i,j} = k(x_i, x_j)$$

$$K = \begin{pmatrix} k(x_1, x_1) & k(x_1, x_2) \\ k(x_2, x_1) & k(x_2, x_2) \end{pmatrix}$$

- **Mercer's theorem**

A function  $k : \mathbb{R}^d \times \mathbb{R}^d \mapsto \mathbb{R}$  is a kernel iff

- For any  $m \in \mathbb{N}$  and any  $x_1, \dots, x_m \in \mathbb{R}^d$ .
- Define the matrix  $K \in \mathbb{R}^{m \times m}$  where  $K_{i,j} = k(x_i, x_j)$  for all  $i, j = 1, \dots, m$ .
- Then  $K$  is symmetric and positive semi-definite.

$\forall c$  positive semi-definite  $K$

$$\nexists w \neq 0 \quad w^T K w \geq 0$$

$$w = \begin{pmatrix} a \\ b \end{pmatrix}, \quad w \neq 0, \quad w \in \mathbb{R}^2 \quad \text{if } \gamma > 0$$

$$w^T K w < 0$$

positive semi-definite  $\Leftrightarrow K \succeq 0$

$$(a \ b) \begin{pmatrix} k_{11} & k_{12} \\ k_{21} & k_{22} \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} = \begin{pmatrix} ak_{11} + bk_{21} & ak_{12} + bk_{22} \end{pmatrix} \begin{pmatrix} a \\ b \end{pmatrix} =$$

$$a^2 \tanh(\|x_1\|^2 + f) + 2ab \tanh(x_1^\top x_2 + f) + b^2 \tanh(\|x_2\|^2 + f)$$

$$b = -1, \quad a = 1 \quad \text{Ansatz}$$

$$a^2 \tanh(\|x_1\|^2 + f) - 2 \tanh(x_1^\top x_2 + f) + \tanh(\|x_2\|^2 + f)$$

$$x_1 = \begin{pmatrix} x_1 \\ n \end{pmatrix}, \quad n \rightarrow \infty \quad \text{Ansatz}$$

$$x_2 = \frac{1}{\sqrt{2}} \begin{pmatrix} 1 \\ 1 \end{pmatrix} \Rightarrow \|x_2\| = 1$$

$$\tanh(\infty) - 2 \tanh(\infty) + \tanh(1+f) = -1 + \tanh(1+f) \leftarrow 0$$

$$f \in \mathbb{R} \Rightarrow -\infty < 1+f < \infty$$

↓

$$-1 < \tanh(1+f) < 1$$



## Question 2

The hinge loss  $\max(0, 1 - y(w^T x + b))$  is commonly used in SVMs. However, consider an alternative quadratic loss:

$$L_i(w, b) = [\max(0, 1 - y_i(w^T x_i + b))]^2.$$

- (a) In soft-SVM, what do the slack variables represent? What is the geometrical meaning of the different value ranges that the slack variables can have? How are the slack variables related to the hinge loss?
  - (b) Write the primal optimization problem for a soft-margin SVM using this new quadratic loss. Hint: Consider about your answers to section a to understand how it would be different from the primal problem with the hinge loss.
  - (c) Write the Lagrangian, KKT conditions, and the dual problem for a soft-margin SVM using this new quadratic loss.
  - (d) Derive the gradient of the loss function  $L_i(w, b)$  with respect to  $w$  and  $b$ .
  - (e) Discuss the implications of using this loss function instead of the hinge loss. How might this impact the sparsity of support vectors and the generalization performance?

a)

1

- Local Stack Variables - ↗
  - ↳ Local variables are stacked fixed
  - ↳ Temporary variables - local → 300
  - ↳ Registers → variables local
  - ↳ Temporary variables local register
  - ↳ Local variables → stack variables local
  - ↳ Temporary variables → stack variables local

(1)  $\text{H}_2\text{O} \rightarrow \text{H}_2 + \text{O}_2$

$$-\log \left( \frac{1}{\rho} \right) = \text{const} \Rightarrow \rho \propto e^{-\lambda x}$$

(D)  $\sum_i \max(0, 1 - y_i(\omega^\top x_i + b)) \leq C$

- hinge loss  $\hat{y}_i$  slack  $\Rightarrow$  when  $\epsilon_i > 1$

if  $y_i(\omega^\top x_i + b) \geq 1 - \epsilon_i$  complementary slackness -N

$$y_i(\omega^\top x_i + b) = 1 - \epsilon_i$$

$\downarrow$

$$\epsilon_i = \max\{0, 1 - y_i(\omega^\top x_i + b)\}$$

$\Leftrightarrow$   
hinge loss  $\Rightarrow$   $\hat{y}_i$

b)

hinge loss BP Egret SVM if  $\epsilon_i = 0$

$$\min_{\omega, b, \epsilon} \frac{1}{2} \|\omega\|^2 + \frac{1}{n} \sum_{i=1}^n \epsilon_i$$

$\{\epsilon_i\}$

$$\ell_i = \text{hinge}(\omega, b, (x_i; y_i))$$

הנתקה מ- $\mathbb{R}^d$  ב- $\mathbb{R}^n$  נקראת גראן

$$L(w, b) = \sum_i \max(0, 1 - y_i(w^\top x_i + b))^2 = \sum_i \varepsilon_i^2$$

המינימום של פונקציית האפס-אחד מושג על ידי  $w = 0, b = 0$

$$\varepsilon_i^2 = \max(0, 1 - y_i(w^\top x_i + b))$$

המינימום של פונקציית האפס-אחד מושג על ידי  $w = 0, b = 0$

המינימום של פונקציית האפס-אחד מושג על ידי  $w = 0, b = 0$

המינימום של פונקציית האפס-אחד מושג על ידי  $w = 0, b = 0$

$$\min_{w,b,\varepsilon} \frac{1}{2} \|w\|^2 + \sum_i \varepsilon_i^2$$

$$y_i(w^\top x_i + b) \geq 1 - \varepsilon_i \quad \varepsilon_i \geq 0 \quad \forall i$$

המינימום של פונקציית האפס-אחד מושג על ידי  $w = 0, b = 0$

"מינימום" מוגן  $\Rightarrow$  מינימום מוגן

ענו על תבנית סינטטיית מילויים ב-  
השאלה פונקצייתית

$$L(w, b, \epsilon, \alpha, \mu) = \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \epsilon_i^2 - \sum_{i=1}^n \alpha_i (y_i (w^T u_i + b) - 1 - \epsilon_i) - \sum_{i=1}^n \mu_i \epsilon_i$$

$$L(w, b, \epsilon, \alpha, \mu) = \frac{\lambda}{2} \|w\|^2 + \frac{1}{n} \sum_{i=1}^n \epsilon_i^2 - \sum_{i=1}^n (\alpha_i y_i (w^T u_i + b) - \alpha_i - \epsilon_i) - \sum_{i=1}^n \mu_i \epsilon_i$$

$$\frac{\partial L}{\partial w} = \lambda w - \sum_{i=1}^n \alpha_i u_i y_i = 0$$

$$\frac{\partial L}{\partial \epsilon_i} = \frac{2}{n} \epsilon_i - \alpha_i - \mu_i = 0$$

$$\frac{\partial L}{\partial \mu} = - \sum_{i=1}^n \alpha_i y_i = 0$$

$$\begin{cases} w = \frac{1}{\lambda} \sum_{i=1}^n \alpha_i u_i y_i \\ \sum_{i=1}^n \alpha_i y_i = 0 \\ \epsilon_i = \frac{1}{2} (\alpha_i - \mu_i) \end{cases}$$

$$y_i(\omega^T u_i + b) \geq 1 - \varepsilon_i, \quad \varepsilon_i \geq 0 \quad \forall i \quad \text{--- optimal probability}$$

$$\alpha_i \geq 0, \quad \nu \geq 0 \quad \text{--- equal probability}$$

$$\alpha_i(1 - \varepsilon_i - y_i(\omega^T u_i + b)) = 0 \quad \nu_i \varepsilon_i = 0 \quad \text{--- complementary slackness}$$

$$\cdot \rightarrow \text{if } \beta_i > 0 \rightarrow \text{--- (c) bad}$$

$\nu_i \varepsilon_i = 0$  Complementary slackness (1) true

$$\text{for } \varepsilon_i > 0 \quad \nu_i = 0 \quad \nu_i \varepsilon_i = 0$$

$$\nu_i \varepsilon_i = 0 \leftarrow \varepsilon_i \geq 0 \text{ iff } \nu_i = 0$$

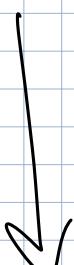
$$\cup_{i=1}^n \{x_i\} \subset \{x_i\}_{i=1}^n \subset N \subset \mathbb{R}^d$$

$$\varepsilon_i = \frac{1}{2} \alpha_i \quad \text{for } i = 1, 2, \dots, n$$

$$\rightarrow \text{the } \cup_{i=1}^n \{x_i\} \text{ is full } \Rightarrow 3)$$

$$\therefore \{x_i\} \cup \{x_j\} \cap \{x_k\} = \emptyset$$

$$\|\omega\|^2 = \left\| \frac{1}{\lambda} \sum_{i=1}^n \alpha_i u_i y_i \right\|^2 = \frac{1}{\lambda^2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (u_i^T u_j) \quad (\text{MKP})$$



$$\frac{1}{n} \sum_{i=1}^n \epsilon_i^2 = \frac{1}{n} \sum_{i=1}^n \left( \frac{y_i - \hat{y}_i}{\sigma} \right)^2 = \frac{1}{n} \sum_{i=1}^n \hat{\alpha}_i^2$$

- (100%) 73%) 100%

$$\max_{\alpha} \sum_{i=1}^n \alpha_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \alpha_i \alpha_j y_i y_j (u_i^T u_j) - \frac{n}{4} \sum_{i=1}^n \alpha_i^2$$

$$\text{s.t. } \sum_{j=1}^4 \alpha_j y_j = 0 \quad \alpha_i \geq 0$$

200%) 100% (6) 100% 73%) 100% (6)

- 100% 100%

$$\because f_i = 1 - y_i (\omega^T u_i + b) \geq 0 \quad (1)$$

$$\frac{\partial f_i}{\partial w} = 0 \quad \frac{\partial f_i}{\partial b} = 0$$

$$\therefore f_i = (1 - y_i (\omega^T u_i + b))^2 = 1 - y_i (\omega^T u_i + b) \geq 0 \quad (2)$$

$$\alpha_i < 1 - y_i (\omega^T u_i + b) + y_i^2 (\omega^T u_i + b)^2 =$$

$$1 - 2y_i \omega^T u_i - 2y_i b + y_i^2 (\omega^T u_i)^2 + 2y_i^2 (\omega^T u_i b) + y_i^2 b^2 =$$

$$1 - 2y_i \omega^T u_i - 2y_i b + y_i^2 (\omega^T u_i)^2 + y_i^2 (\omega^T u_i b) + y_i^2 b^2$$

$$\frac{\partial f_i}{\partial w} = -2y_i u_i - 2y_i^2 \omega^T u_i^2 - 2y_i^2 b - 2y_i u_i (-1 - y_i \omega^T) w + y_i b$$

$$= -2y_i u_i (1 - (y_i)(\omega^T u_i + b))$$

$$\frac{\partial L}{\partial \theta_i} = -y_i \log y_i + w^T u_i - \log 2 \quad \text{or} \quad -y_i(1 - (y_i)(w^T u_i - \log 2))$$

Now we have to find the force exerted by the string on the ball. This force is the resultant of the tension in the string and the weight of the ball. Let's call the angle between the string and the vertical  $\theta$ . Then, the horizontal component of the tension is  $T \sin \theta$  and the vertical component is  $T \cos \theta$ . The weight of the ball is  $mg$ . So, the horizontal component of the tension is equal to the horizontal component of the weight, i.e.,  $T \sin \theta = mg$ . The vertical component of the tension is equal to the vertical component of the weight, i.e.,  $T \cos \theta = mg$ . From the first equation, we get  $T = mg / \sin \theta$ . Substituting this value of  $T$  in the second equation, we get  $mg / \sin \theta \cos \theta = mg$ . Simplifying, we get  $\tan \theta = 1$ . Therefore,  $\theta = 45^\circ$ .

היום ה-10.02.2020 ברכבת מירושלים לוד

„Overfitting“ ist ein Begriff aus der Statistik und beschreibt die Tendenz eines Modells, an den gegebenen Daten zu gut zu passen. Dies führt dazu, dass das Modell bei neuen, unbekannten Daten nicht mehr präzise vorhersagen kann. Ein Modell ist überfittet, wenn es die Struktur der Störungen im Trainingssatz so genau wie möglich wiedergibt, was zu einem schlechten Generalisierungsleistung führt.

לפחות נול -tip-in כוונתנו מיל' ית'

רשותה של אוניברסיטת ירושלים לשלוח נייר

-היפר בעיה הינה מיל' כוונתנו מיל'

היפר בעיה הינה מיל' כוונתנו מיל'

### Question 3

- (a) What is the sample complexity of hard margin SVM in the linear case (without a kernel) and why? Explain your answer, don't just quote from class.
- (b) Briefly explain how will using an RBF kernel affect the sample complexity. For an SVM using the RBF kernel, discuss how the sample complexity is influenced by the kernel hyperparameter  $\sigma$ , which controls the spread of the RBF kernel.

(a)

hypothecis class -> a VC-dim ->  $\sigma$  (margin)  $\rightarrow$  10.0

2NCP k3n 1n1l & margin  $\rightarrow$  1n NnL

hyperplane  $\rightarrow$  1e 1.3N hard SVM

- mif  $\rightarrow$  2NQ, l' class -> 1.3N hard SVM  
- hyperplane  $\rightarrow$

$$\text{if } \omega^\top x_i + b = 0 \\ \downarrow$$

hyperplane  $\rightarrow$  1.3N bias

$\rightarrow$  1.3N  $\rightarrow$   $\|\omega\|$

$\rightarrow$  (if  $\|\omega\|$  is large)

separators  $\rightarrow$  1.3N if sample complexity  $\rightarrow$

$\rightarrow$  1.3N  $\rightarrow$   $\|\omega\|$

$$\text{err}(h_{\omega, b}) = O\left(\frac{\delta^{-1-\epsilon} \log \frac{1}{\sigma}}{n}\right)$$

3

$$n = O\left(\frac{\delta^{-1-\epsilon} \log \frac{1}{\sigma}}{\epsilon}\right)$$

הנחתה שוגר במאוד לאט וסביר

ולכן, הטענה מוגדרת כזאת שקיים מילוי למשתנה  $\underline{\lambda}$

$$\frac{\lambda}{\lambda + \delta \lambda} \geq \mu \quad \forall \lambda \in \mathbb{R}$$

וניהו  $\phi \in \mathcal{H}_{\omega, b}$  מושג RBF Kernel (b)

וניהו  $\delta \lambda$  sample complexity של  $\varphi$ , כלומר גודל מינימום

overfit מהילך צדקה מינימום סטטיסטי

$\sigma_1$  kernel hyperparameter  $\rightarrow$  מילוי

$$k(x, y) = e^{-\frac{\|x-y\|^2}{2\sigma^2}}$$

$x, y \in \mathbb{R}^d$  ו $k$ )

kernel  $\Gamma \vdash t : A$  (Fig)  $\vdash t : A$   $\vdash t : A$

לעכ. י.ר. ס. פ.ג.יכ נ.פ.ר. ז.ר. מ.ר. ק.ר. ל.ר. מ.ר. ג.ר. (ל.ר. מ.ר.)

overfitting  $\rightarrow$  fit first sample - then  $\rightarrow$  fit second sample

גג' RBF Kernel  $\rightarrow$  ס.  $\Gamma \rightarrow \Gamma$   $j(\gamma)$   $\rho(\gamma)$

71  $\int_{\text{min}}^{\text{max}} f(x) dx$   $\rightarrow$   $\int_{\text{min}}^{\text{max}} g(x) dx$

It is upon you to hope for me.

לְמִזְבֵּחַ וְלְמִזְבֵּחַ וְלְמִזְבֵּחַ

Training set  $\{x_i\}_{i=1}^n$   $\rightarrow$   $f(x)$   
and  $y_i \in \mathcal{Y}$ ,  $\hat{f}(x)$ , Overfitting

## Question 4

Given a dataset with two classes,  $C_1$  and  $C_2$ , and the following class distributions:

Feature Value	Class $C_1$ Count	Class $C_2$ Count
A	8	2
B	3	7

- (a) Compute the entropy  $H(S)$  of the entire dataset and for the subsets  $S_A, S_B$  (the subsets corresponding to feature values A and B).
- (b) Calculate the information gain  $G(S, A)$  when splitting the dataset on the feature value.
- (c) Interpret your results: What does the information gain tell you about the usefulness of splitting on this feature?
- (d) Compute the Gini impurity for the entire dataset S, and for the subsets  $S_A$  and  $S_B$  before splitting.
- (e) Calculate the weighted average Gini impurity after splitting the dataset by the feature values A and B.
- (f) Compare the results from entropy-based information gain and Gini impurity. Are the rankings of the splits the same? Explain why or why not.

$$a) H(S) = - \left( s \log_2(s) + (1-s) \log_2(1-s) \right)$$

dataset S

$$|S| = 20$$

number of 1's = 8  
number of 0's = 12

$$= - \left( \frac{8}{20} \log_2 \left( \frac{8}{20} \right) + \left( \frac{12}{20} \right) \log_2 \left( \frac{12}{20} \right) \right) = 0.992$$

$$|S| = 20$$

$$H(S_A) = - \left( \frac{8}{10} \log_2 \left( \frac{8}{10} \right) + \left( \frac{2}{10} \right) \log_2 \left( \frac{2}{10} \right) \right) = 0.722$$

$$|S_A| = 10$$

$$H(S_B) = - \left( \frac{3}{10} \log_2 \left( \frac{3}{10} \right) + \left( \frac{7}{10} \right) \log_2 \left( \frac{7}{10} \right) \right) = 0.881$$

$$|S_B| = 10$$

b) 'if' if condition information gain  $\Rightarrow$   
'not' if condition

A ∪ {3d})

$$G(S, A) = H(S) - H(A) \stackrel{?}{=} ?$$

$$H(A) = P_r(A) \cdot H(S_A) + P_r(B) \cdot H(S_B) =$$

$$\frac{10}{20} \cdot 0.722 + \frac{10}{20} \cdot 0.881 = 0.802$$



$$G(S, A) = H(S) - H(A) = 0.19$$

c) Széf a formulában gondolj fel, hogy miért írunk  
sziszemelési szintet a számításba (azaz 0-nak színezett)  
- ekkor minden szám két -k

d) Régi index = Val(q) =  $\sum q_i \log(1-q_i)$

$$H(S) = 2S(1-S) = 2 \cdot \frac{1}{20} \left( \frac{1}{20} \right) = 0.495$$

$$H(S_A) = 2 \cdot \frac{8}{10} \cdot \frac{2}{20} = 0.32$$

$$H(S_B) = 2 \cdot \frac{2}{10} \cdot \frac{7}{20} = 0.42$$

e)  $H(\text{after } \text{spf}) = P_r(A)H(S_A) + P_r(B)H(S_B) =$

(~~?~~)

$$\frac{1}{2} \cdot 0.32 + \frac{1}{2} \cdot 0.42 = 0.37$$

f)  $P(A \cap B) = P(A)P(B|A)$   
 $= 0.2 \cdot 0.16 = 0.032$

## Question 5

Consider a binary decision tree used for binary classification, such as the ones we've seen in class, where each internal node splits the data based on a single feature threshold, and each leaf represents a class label. Suppose the dataset has  $n$  features, and each feature can take on real values.

Tree 3

- (a) Prove that a binary decision tree with  $\underline{\text{depth}} d$  can divide the feature space into at most  $2^d$  regions.
- (b) Explain how this result relates to the expressiveness of decision trees and their capacity to overfit.
- (c) Using this result, determine the minimum depth  $d$  required for a decision tree to perfectly classify a dataset of  $m$  unique points (no two points share the same feature values).
- (d) Based on the previous sections, what is the VC dimension of a tree with depth  $d$ ?

a)



My tree  $d \sqrt{m}$

Binary tree ( $f(p)$ )  $\sim 2^d$  regions  $\sim m^{(d)}$  leaves  $\sim m^d$

$2^d$  regions  $\sim$   $\underline{2^d}$

GP L.  $\rightarrow$   $\sim 2^d$  regions  $\sim 2^d$  leaves  $\sim 2^d$  nodes

R-regions

- 0.07

$$\delta = 0$$

$$f_k(f_0) = 2^0 = 1$$

Split  $f_k$   $\circ \psi$

0

$\delta^{-1}$  กับ  $\mu^{\delta}$  ล้วน เปรียบเท่าๆ - 3f3

$$R(\delta^{-1}) = \delta^{\delta^{-1}}$$

( ผู้ที่  $\delta^{\delta^{-1}}$  นับว่า บีบ หรือ )

โดย จำกัด บน 0/1/1

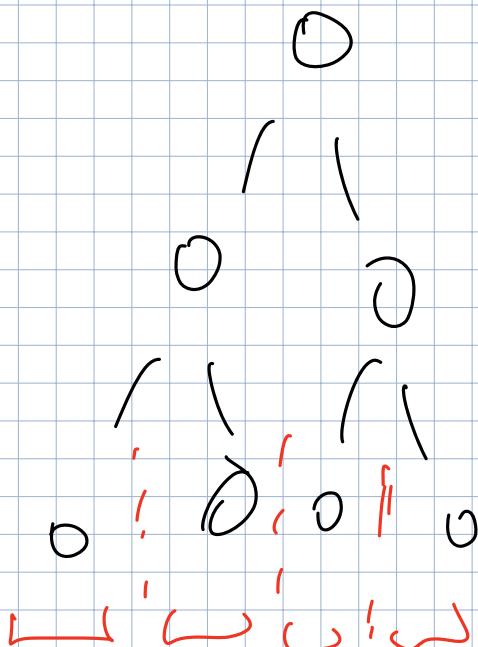
$$R(\delta) = \delta^\delta \quad \text{ผู้ที่}$$

( ผู้ที่  $\delta^\delta$  นับ )

$$R(\delta) = \delta^{\delta^{-1}} \cdot \delta = \delta^\delta$$

—  $\delta^\delta$  regions บน  $\mathbb{R}^d$  จด 0

c)



b)

## Question 5

Consider a binary decision tree used for binary classification, such as the ones we've seen in class, where each internal node splits the data based on a single feature threshold, and each leaf represents a class label. Suppose the dataset has  $n$  features, and each feature can take on real values.

Tree 3

- (a) Prove that a binary decision tree with  $\underline{\text{depth } d}$  can divide the feature space into at most  $2^d$  regions.
- (b) Explain how this result relates to the expressiveness of decision trees and their capacity to overfit.
- (c) Using this result, determine the minimum depth  $d$  required for a decision tree to perfectly classify a dataset of  $m$  unique points (no two points share the same feature values).
- (d) Based on the previous sections, what is the VC dimension of a tree with depth  $d$ ?

Decision tree  $\rightarrow$  regions

Each level splits the space into 2 regions

Number of regions  $\rightarrow 2^d$

Regions overlap  $\rightarrow$  overfitting

With more nodes, fewer regions

More nodes  $\rightarrow$  less regions

- (c) Using this result, determine the minimum depth  $d$  required for a decision tree to perfectly classify a dataset of  $m$  unique points (no two points share the same feature values).

- (d) Based on the previous sections, what is the VC dimension of a tree with depth  $d$ ?

c) To find the minimum depth  $d$  for a tree to classify  $m$  unique points.

Find  $\min(d)$  such that regions  $\rightarrow$   $\geq m$

$m$  unique points  $\rightarrow$   $\geq 2^d$

so  $2^d \geq m$  or  $d \geq \log_2(m)$

$$2^{(\text{min } d)} = m$$

↓

$$\min d = \log_2(m)$$

$d \geq \lceil \log_2(m) \rceil$  (round up to nearest integer)

d) VC dim =  $2^d$  - if a tree can classify  $2^d$  points

- if  $d$   $\geq n$  then it can't misclassify