

רק מליל  
322357542 - 5.7

# Theoretical Assignment 1

Intro to Machine Learning, TAU

November 18, 2024

## Question 1

Imagine you are classifying flowers as either *Iris-setosa* or *Iris-versicolor* based on two features: petal length and petal width. You have the following information:

Class	Mean Petal Length	Std Dev Petal Length	Mean Petal Width	Std Dev Petal Width
Iris-setosa	1.5	0.2	0.2	0.1
Iris-versicolor	4.7	0.3	1.4	0.2

Table 1: Summary statistics for petal length and width by flower class.

Assume that the features are independent and normally distributed.

1. Calculate the likelihoods  $P(\text{petal length} = 1.6 \mid \text{Iris-setosa})$  and  $P(\text{petal width} = 0.25 \mid \text{Iris-setosa})$  using the probability density function of a normal distribution, and plot them.
2. Based on these likelihoods, design a classifier that receives the petal width and petal length and predicts the the class.
3. Using the classifier from the previous section, determine which class (Iris-setosa or Iris-versicolor) is more likely for a flower with a petal length of 1.6 and a petal width of 0.25. Assume equal priors for both classes.

## Question 2

A medical test is used to classify patients into two categories: *Disease* or *No Disease*. However, due to the high costs associated with misclassification, there is also an option to reject a classification if confidence is too low. You are given the following information:

- $P(\text{Disease}) = 0.2$

- $P(\text{No Disease}) = 0.8$
- The test has a sensitivity (true positive rate) of 90%, meaning  $P(\text{Positive Test} \mid \text{Disease}) = 0.9$ .
- The test has a specificity (true negative rate) of 85%, meaning  $P(\text{Negative Test} \mid \text{No Disease}) = 0.85$ .

A patient receives a positive test result.

**Cost Structure:**

- The cost of a false positive (FP) error (i.e., diagnosing a healthy patient as having the disease) is \$1,000.
- The cost of a false negative (FN) error (i.e., failing to diagnose a patient with the disease) is \$5,000.
- The cost of rejecting the classification for further testing is \$2,000.

**Tasks:**

1. Calculate the probability that this patient actually has the disease given a positive test result,  $P(\text{Disease} \mid \text{Positive Test})$ .
2. Based on this probability and the cost structure, decide whether to classify the patient as *Disease*, *No Disease*, or *Reject*. Choose the option with the lowest expected cost. Use a confidence threshold of 70% for deciding whether to classify the patient as having the disease. If the probability of disease is below 70% but above 30%, classify as *Reject*; otherwise, classify as *No Disease*.

## Question 3

Given a real number  $R \geq 0$ , define the hypothesis  $h_R : \mathbb{R}^d \rightarrow \{0, 1\}$  as follows:

$$h_R(x) = \begin{cases} 1 & \text{if } \|x\|_2 \leq R, \\ 0 & \text{otherwise.} \end{cases}$$

Consider the hypothesis class  $H_{\text{ball}} = \{h_R \mid R \geq 0\}$ . Prove directly (without using the Fundamental Theorem of PAC Learning) that  $H_{\text{ball}}$  is PAC learnable in the realizable case. Assume for simplicity that the marginal distribution of  $X$  is continuous. How does the sample complexity depend on the dimension  $d$ ? Explain.

## Question 4

Given a polynomial  $P : \mathbb{R} \rightarrow \mathbb{R}$ , define the hypothesis  $h_P : \mathbb{R}^2 \rightarrow \{0, 1\}$  as follows:

$$h_P(x_1, x_2) = \begin{cases} 1 & \text{if } P(x_1) \geq x_2, \\ 0 & \text{otherwise.} \end{cases}$$

Determine the VC-dimension of  $H_{\text{poly}} = \{h_P \mid P \text{ is a polynomial}\}$ . You can use the fact that given  $n$  distinct values  $x_1, \dots, x_n \in \mathbb{R}$  and  $z_1, \dots, z_n \in \mathbb{R}$ , there exists a polynomial  $P$  of degree  $n - 1$  such that  $P(x_i) = z_i$  for every  $1 \leq i \leq n$ .

Imagine you are classifying flowers as either *Iris-setosa* or *Iris-versicolor* based on two features: petal length and petal width. You have the following information:

Assume that the features are independent and normally distributed.

 $\lambda$ 

$\text{Iris-setosa length} \sim N(1.5, 0.04)$   
 $\text{Iris-setosa width} \sim N(0.2, 0.01)$   
 $\text{Iris-versicolor length} \sim N(4.7, 0.09)$   
 $\text{Iris-versicolor width} \sim N(1.4, 0.04)$

$$\underline{x=1.6}$$

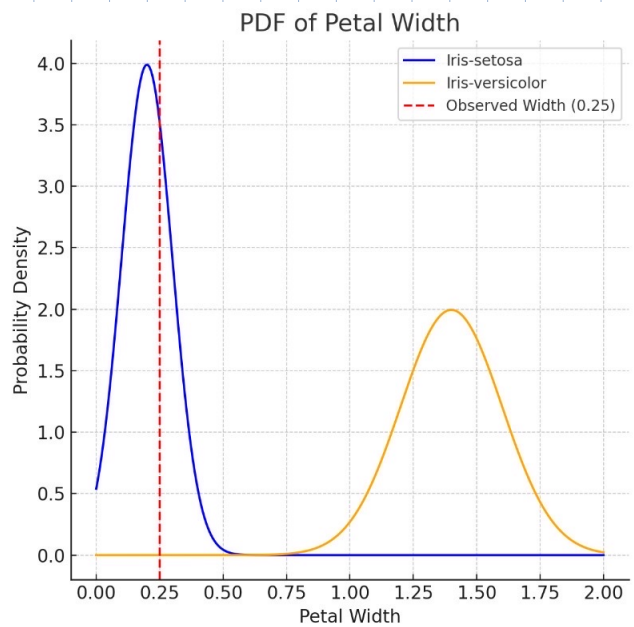
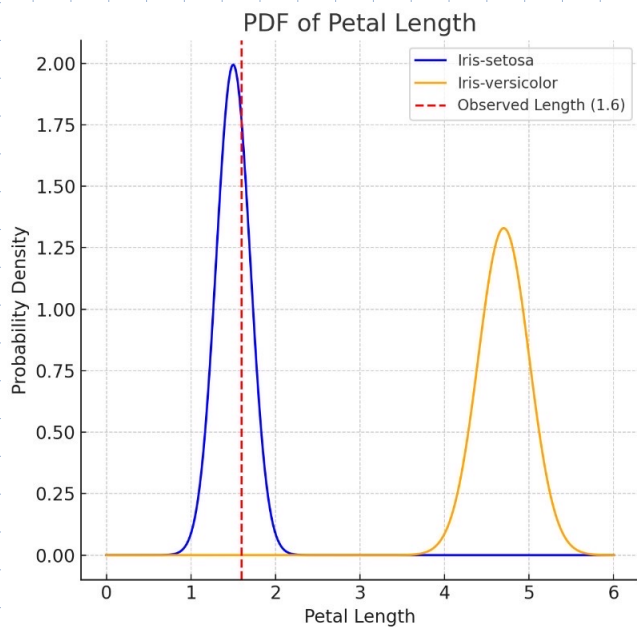
$$f(\text{petal length}=1.6 | \text{Iris-Setosa}, 0.2^2) = \frac{1}{\sqrt{2\pi \cdot 0.2^2}} e^{-\frac{(1.6-1.5)^2}{2 \cdot 0.2^2}}$$

$$= 1.76$$

$$f(\text{petal width}=0.5 | \text{Iris-Setosa}, 0.1^2) = \frac{1}{\sqrt{2\pi \cdot 0.1^2}} e^{-\frac{(0.5-0.2)^2}{2 \cdot 0.1^2}}$$

$$= 3.52$$

1



2

$$p \sim N(\mu, \Sigma)$$

$$\bar{\mu}_1 = \begin{pmatrix} 1.7 \\ 0.2 \end{pmatrix}$$

$$\Sigma_1 = \begin{pmatrix} 0.04 & 0 \\ 0 & 0.01 \end{pmatrix}$$

$$\bar{\mu}_2 = \begin{pmatrix} 4.7 \\ 1.4 \end{pmatrix}$$

$$\Sigma_2 = \begin{pmatrix} 0.09 & 0 \\ 0 & 0.04 \end{pmatrix}$$

$$\bar{\Sigma}_1 = \begin{pmatrix} 0.04 & 0 \\ 0 & 0.01 \end{pmatrix}$$

$$\bar{\Sigma}_2 = \begin{pmatrix} 0.09 & 0 \\ 0 & 0.04 \end{pmatrix}$$

$$p(\bar{x} | \mu_i) = \frac{1}{\sqrt{2\pi}^2 \sqrt{|\Sigma_i|}} e^{-\frac{-(\bar{x} - \mu_i)^T \Sigma_i^{-1} (\bar{x} - \mu_i)}{2}}$$

min  $\beta$  (classification — 1 2)

$$C = \arg \max_{i=1,2} (P(\bar{x} | a_i) \cdot P(a_i))$$

8)

length = 1.6 width = 0.25  
— 1 2 3 4 5 6 7

Iris - setosa length  $\sim N(1.5, 0.04)$

Iris - setosa width  $\sim N(0.2, 0.01)$

Iris - versicolor length  $\sim N(4.7, 0.09)$

Iris - versicolor width  $\sim N(1.4, 0.04)$

Iris - setosa -  $a_1$

Iris - versicolor -  $a_2$

$$\underline{x} = \begin{pmatrix} 1.6 \\ 0.25 \end{pmatrix}$$

—  $a_1$  1 2

$$P(\bar{x} | a_i) = \frac{1}{\sqrt{2\pi}^2 \sqrt{|\Sigma|}} e^{\frac{-(\bar{x} - \mu_i)^T \Sigma^{-1} (\bar{x} - \mu_i)}{2}}$$

$$\sqrt{|\Sigma_1|} = \sqrt{4 \cdot 10^{-9}} = 0.02$$

$$\Sigma_1^{-1} = \begin{pmatrix} 25 & 0 \\ 0 & 100 \end{pmatrix}$$

$$(\underline{x} - \underline{\mu}_1) = \begin{pmatrix} 1.6 \\ 0.25 \end{pmatrix} - \begin{pmatrix} 1.5 \\ 0.2 \end{pmatrix} = \begin{pmatrix} 0.1 \\ 0.05 \end{pmatrix}$$

$$(\underline{x} - \underline{\mu}_1)^T = (0.1 \quad 0.05)$$

$$(\underline{x} - \underline{\mu}_1)^T \cdot \Sigma_1^{-1} \cdot (\underline{x} - \underline{\mu}_1) = (0.1 \quad 0.05) \begin{pmatrix} 25 & 0 \\ 0 & 100 \end{pmatrix} \begin{pmatrix} 0.1 \\ 0.05 \end{pmatrix} = 0.7$$

$$P(x|a_1) = \frac{1}{\sqrt{2\pi} \cdot 0.02} \cdot e^{-\frac{0.7}{2}} \approx 8.74$$

= 0.12 n/r

$$\sqrt{|\Sigma_2|} = \sqrt{36 \cdot 10^{-9}} = 0.06$$



$$\Sigma_{\underline{g}}^{-1} = \begin{pmatrix} 11.1 & 0 \\ 0 & 2.5 \end{pmatrix}$$

$$(\underline{x} - \underline{\mu}_{\underline{g}}) = \begin{pmatrix} 1.6 \\ 0.25 \end{pmatrix} - \begin{pmatrix} 4.7 \\ 1.4 \end{pmatrix} = \begin{pmatrix} 3.1 \\ 1.15 \end{pmatrix}$$

$$(\underline{x} - \underline{\mu}_{\underline{g}})^T = (3.1 \quad 1.15)$$

$$(\underline{x} - \underline{\mu}_{\underline{g}})^T \cdot \Sigma_{\underline{g}}^{-1} \cdot (\underline{x} - \underline{\mu}_{\underline{g}}) = (3.1 \quad 1.15) \begin{pmatrix} 11.1 & 0 \\ 0 & 2.5 \end{pmatrix} \cdot \begin{pmatrix} 3.1 \\ 1.15 \end{pmatrix} = 139.733$$

$$P(\bar{x} | a_2) = \frac{1}{\sqrt{2} \pi \cdot 0.06} \cdot e^{-\frac{139.733}{2}} \rightarrow 0$$

$$P(\bar{x} | a_1) \gg P(\bar{x} | a_2)$$

$$P(a_2) = P(a_1) = \frac{1}{2}$$

Fris-Settera  $\rho$

## Question 2

A medical test is used to classify patients into two categories: *Disease* or *No Disease*. However, due to the high costs associated with misclassification, there is also an option to reject a classification if confidence is too low. You are given the following information:

- $P(\text{Disease}) = 0.2$

1

- $P(\text{No Disease}) = 0.8$
- The test has a sensitivity (true positive rate) of 90%, meaning  $P(\text{Positive Test} \mid \text{Disease}) = 0.9$ .
- The test has a specificity (true negative rate) of 85%, meaning  $P(\text{Negative Test} \mid \text{No Disease}) = 0.85$ .

A patient receives a positive test result.


### Cost Structure:

- The cost of a false positive (FP) error (i.e., diagnosing a healthy patient as having the disease) is \$1,000.
- The cost of a false negative (FN) error (i.e., failing to diagnose a patient with the disease) is \$5,000.
- The cost of rejecting the classification for further testing is \$2,000.

### Tasks:

1. Calculate the probability that this patient actually has the disease given a positive test result,  $P(\text{Disease} \mid \text{Positive Test})$ .
2. Based on this probability and the cost structure, decide whether to classify the patient as *Disease*, *No Disease*, or *Reject*. Choose the option with the lowest expected cost. Use a confidence threshold of 70% for deciding whether to classify the patient as having the disease. If the probability of disease is below 70% but above 30%, classify as *Reject*; otherwise, classify as *No Disease*.

$$1. P(\text{Disease} | \text{Positive Test}) = \frac{P(\text{Positive Test} | \text{Disease}) \cdot P(\text{Disease})}{P(\text{Positive Test})} = 0.2$$



∴ P(D) = 0.2

$$P(\text{Positive Test}) = P(\text{Positive Test} | \text{Disease}) P(\text{Disease}) + P(\text{Positive Test} | \text{No Disease}) P(\text{No Disease})$$

$$P(\text{Positive Test} | \text{No Disease}) + P(\text{Negative Test} | \text{No Disease}) = 1$$

↓

$$P(\text{Negative Test} | \text{No Disease}) = 1 - P(\text{Positive Test} | \text{No Disease})$$

$$1 - 0.85 = 0.15$$

0.15

$$P(\text{Positive Test}) = P(\text{Positive Test} | \text{Disease}) P(\text{Disease}) + P(\text{Positive Test} | \text{No Disease}) P(\text{No Disease})$$

$$= 0.9 \cdot 0.2 + 0.15 \cdot 0.8 = 0.3$$

0.3

$$P(\text{Disease} | \text{Positive Test}) = \frac{P(\text{Positive Test} | \text{Disease}) \cdot P(\text{Disease})}{P(\text{Positive Test})}$$

$$= \frac{0.9 \cdot 0.2}{0.3} = 0.6$$

2. threshold 70%.

$p < 0.3$  - No Disease

$0.3 < p < 0.7$  - Disease

$p > 0.7$  - reject

0.6 אכן חולה במחלה

0.3 < p < 0.7 חולה

reject  
החלטה שגויה  
החלטה נכונה

$L(a, c)$	cost - $L(c)$	$y_1$ disease	$y_2$ no disease
		$y_1$	$y_2$
$a_1$ reject		2000	2000
$a_1$ positive		0	1000
$a_2$ negative		1000	0

-positive result 2020

$$\text{cost}(a_1, y_1) = 0 + 1(a_1, y_2) P(y=y_2|a_1) =$$

$$1000(1 - P(y=y_1|a_1)) = 1000 \cdot 0.07 = 700$$

-reject 2020

$$\text{cost}(a_0, y_1) = 2000$$

$$\text{cost}(a_2, y_1) = 0 + 5000 \cdot P(y=y_1|a_2) =$$

$$0 + 5000 \cdot P(\text{Disease} | \text{negative}) = 5000 \cdot \frac{P(\text{Disease, negative})}{P(\text{negative})} =$$

$$5000 \cdot \frac{0.02}{0.7} = 142.857$$

142.857 < 2000 **reject** →  $C_0$  opt.

threshold → reject unless 2020 pr. only

## Question 4

Given a polynomial  $P : \mathbb{R} \rightarrow \mathbb{R}$ , define the hypothesis  $h_P : \mathbb{R}^2 \rightarrow \{0, 1\}$  as follows:

$$h_P(x_1, x_2) = \begin{cases} 1 & \text{if } P(x_1) \geq x_2, \\ 0 & \text{otherwise.} \end{cases}$$

Determine the VC-dimension of  $H_{\text{poly}} = \{h_P \mid P \text{ is a polynomial}\}$ . You can use the fact that given  $n$  distinct values  $x_1, \dots, x_n \in \mathbb{R}$  and  $z_1, \dots, z_n \in \mathbb{R}$ , there exists a polynomial  $P$  of degree  $n - 1$  such that  $P(x_i) = z_i$  for every  $1 \leq i \leq n$ .

$$\text{VC-dim}(H_{\text{poly}}) = \infty$$

$P(x_1) \geq x_2$  -  $(x_1, x_2)$  נמצא בתחתית של  $P$

$y_1 \geq x_2$  -  $(x_1, y_1)$  נמצא מעל  $P$  (  $y_1 = P(x_1)$  )

(mapping)

נבחר  $n$  נקודות  $(x_1, y_1), \dots, (x_n, y_n)$  ונמצא  $P$  כזה ש  $P(x_i) = y_i$

$x_1, \dots, x_n, x_{n+1}$  נבחרים כך ש  $x_{n+1} > x_i$

אם  $y_i = 1$  אז  $P(x_i) \geq x_{n+1}$  ונמצא מעל

$y_i = 0$  אז  $P(x_i) < x_{n+1}$  ונמצא בתחתית

לכן  $H_{\text{poly}}$  יכול להפריד כל קבוצת  $n$  נקודות. מכאן ש  $\text{VC-dim}(H_{\text{poly}}) \geq n$  לכל  $n$ .

$$1 \leq i \leq n \quad (x_1^i, x_2^i) \quad \text{נבחר}$$

$1 \leq i \leq n$ ,  $(x_1^i, x_2^i)$  נבחרים כך ש  $x_2^i > x_1^i$

נבחר  $y_i = 1$  או  $0$  לפי  $x_2^i > x_1^i$

$$W_i = \begin{cases} x_2^i - x_1^i & \text{if } y_i = 1 \\ x_1^i - x_2^i & \text{if } y_i = 0 \end{cases}$$

(ה) שני סדרים  $\{x_n\}$  ו- $\{y_n\}$  בלתי

שונים (אם  $n \in \mathbb{N}$  אז  $x_n \neq y_n$ )

יש  $N \in \mathbb{N}$  כזה שכל  $n > N$  מקיים  $|x_n - y_n| < \epsilon$

וההפך: אם  $\{x_n\}$  ו- $\{y_n\}$  הם שני סדרים

כאלה ש- $\lim_{n \rightarrow \infty} x_n = a$  ו- $\lim_{n \rightarrow \infty} y_n = b$  אז

אם  $a \neq b$  אז  $\lim_{n \rightarrow \infty} (x_n - y_n) = a - b \neq 0$

לכן  $\lim_{n \rightarrow \infty} (x_n - y_n) = 0$  אם ורק אם  $a = b$

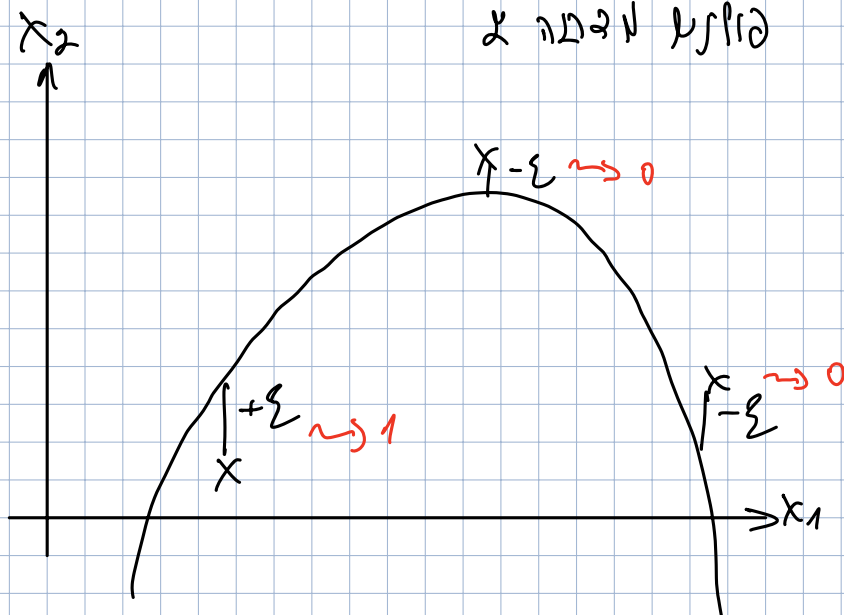
כלומר: שני סדרים מתכנסים לאותו הגבול

אם ורק אם ההפרש מתכנס ל-0

למשל:  $\lim_{n \rightarrow \infty} (x_n - y_n) = 0$  אם ורק אם  $\lim_{n \rightarrow \infty} x_n = \lim_{n \rightarrow \infty} y_n$

דוגמה:  $x_n = \frac{1}{n}$  ו- $y_n = \frac{1}{n^2}$ . אז  $\lim_{n \rightarrow \infty} x_n = 0$  ו- $\lim_{n \rightarrow \infty} y_n = 0$

באופן כללי:



### Question 3

Given a real number  $R \geq 0$ , define the hypothesis  $h_R : \mathbb{R}^d \rightarrow \{0, 1\}$  as follows:

$$h_R(x) = \begin{cases} 1 & \text{if } \|x\|_2 \leq R, \\ 0 & \text{otherwise.} \end{cases}$$

2

Consider the hypothesis class  $H_{\text{ball}} = \{h_R \mid R \geq 0\}$ . Prove directly (without using the Fundamental Theorem of PAC Learning) that  $H_{\text{ball}}$  is PAC learnable in the realizable case. Assume for simplicity that the marginal distribution of  $X$  is continuous. How does the sample complexity depend on the dimension  $d$ ? Explain.

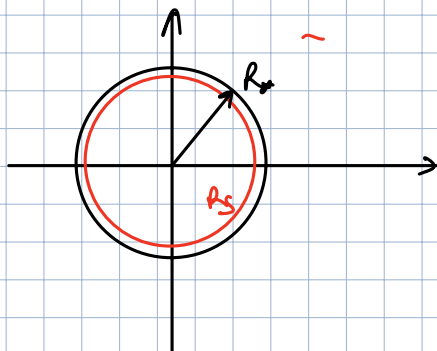
Realizable  $\Rightarrow h \in H$  r.r.

$$y = h(x) \quad \text{-- l } y$$

$$h(x) = \begin{cases} 1 & \|x\| \leq R \\ 0 & \text{o.w.} \end{cases} \quad \text{-- l } y \quad R \in \mathbb{R}^+ \quad \text{r.r.}$$

$$R_S = \max_{1 \leq i \leq n} \|x_i\|_2, \quad S = \{x_i\}_{i=1}^n \quad \text{r.r.} \quad \text{-- l } y \quad \text{r.r.}$$

$$(\text{realizable } \Rightarrow) R_S \leq R_0$$





$$\hat{h}(x_i) = h(x_i) \quad \text{if} \quad \|x_i\| \leq R_S \quad \text{and} \quad x_i \in \mathcal{X}$$

$$0 \leq \|x_i\| \leq R_S$$

$$x_j \in \mathcal{X} \quad \text{and} \quad \|x_j\| \leq R_S \quad \text{and} \quad x_j \in \mathcal{X}$$

$$R_S < \|x_j\| \leq R_{\infty} \quad \text{and} \quad x_j \in \mathcal{X}$$

$$P[\text{error}(\hat{h})] = P[R_S < \|x_j\| \leq R_{\infty}] = P[\|x_j\| \leq R_{\infty}] - P[\|x_j\| \leq R_S]$$

$$P[\text{error}(\hat{h})] < \epsilon \quad \text{if} \quad P[\|x_j\| \leq R_{\infty}] < \epsilon \quad \text{and} \quad P[\|x_j\| \leq R_S] = 0$$

$$P[\text{error}(\hat{h}) > \epsilon] \leq \prod_{i=1}^n (1 - \epsilon) = (1 - \epsilon)^n \leq e^{-\epsilon n} < \delta$$

$$-\epsilon n < \ln(\delta)$$

$$\epsilon n > -\ln(\delta)$$

$$\epsilon n > \ln\left(\frac{1}{\delta}\right)$$

$$n > \frac{1}{\epsilon} \ln\left(\frac{1}{\delta}\right)$$

$$n > \ln\left(\frac{1}{\delta}\right) \cdot \frac{1}{\epsilon}$$

↓

radius of ball