

Linear Algebra

1. (15 pts) A symmetric matrix A over \mathbb{R} is called *positive semidefinite* (PSD) if for every vector v , $v^T A v \geq 0$.

(a) Show that a symmetric matrix A is PSD if and only if it can be written as $A = X X^T$, if and only if all of its eigenvalues are non-negative. ③

Hint: Recall that a real symmetric matrix A can be decomposed as $A = Q D Q^T$, where Q is an orthogonal matrix whose columns are eigenvectors of A and D is a diagonal matrix with eigenvalues of A as its diagonal elements.

פתרון

לנו: $1 \leq 2 \leq 3$

לנו: $1 \leq 3$

לנו: A היא PSD, כלומר לכל וקטור v מתקיים $v^T A v \geq 0$. בנוסף, A סימטרית.

$$0 \leq v^T A v = v^T Q D Q^T v = (Q^T v)^T D (Q^T v) = \sum_{i=1}^n \lambda_i (Q^T v)_i^2$$

כלומר, $\lambda_i \geq 0$

$$\leq \sum_{i=1}^n \lambda_i (Q^T v)_i^2$$

מכאן, פנימית על x הינו אי-שלילי.

3-2: יהי A סימטרית, כלומר $A = A^T$. נניח כי לכל $i \in [n]$, $\lambda_i \geq 0$. נכתוב $A = Q D Q^T$, כאשר Q היא מטריצה אורתוגונלית ו- D היא מטריצה דיאгонаלית.

$$A = Q D Q^T = Q \cdot \text{diag}(\lambda_1, \dots, \lambda_n) \cdot Q^T = (Q \sqrt{\lambda}) (\sqrt{\lambda} Q^T)$$

ניתן אף לסמלל:

$$(\sqrt{\lambda})_{i,j} = \begin{cases} \sqrt{\lambda_i} & i=j \\ 0 & \text{else} \end{cases}$$

על צורה אי-שלילית

ונתון סלילר, ויתר מוגדר היטב.

$$A = M M^T \quad \text{כאשר} \quad M = Q \sqrt{\lambda}$$

2-1: נניח A (ניתנת) סמי-נגזר באופן חיובי: $A = XX^T$. לנניח כי A היא PSD.

יהי v וקטור א-זר מתקין P -

$$v^T A v = v^T X X^T v = (v^T X)(X^T v) = \|v^T X\|^2 \geq 0$$

לוח: גריבול

□

שם ווקטוריות

הראו: $1 \leftarrow 2 \leftarrow 1$ כוונת C .

(b) Show that for all $\alpha, \beta \geq 0$ and PSD matrices $A, B \in \mathbb{R}^{n \times n}$, the matrix $\alpha A + \beta B$ is also PSD. Does this mean that the set of all $n \times n$ PSD matrices over \mathbb{R} is a vector space over \mathbb{R} ?

בתוכן: יהי $\alpha, \beta \geq 0$ וכן $A, B \in \mathbb{R}^{n \times n}$ PSD's. לנניח כי $\alpha A + \beta B$ היא PSD.

יהי וקטור v . נראה כי $v^T (\alpha A + \beta B) v \geq 0$.

$$v^T (\alpha A + \beta B) v = (v^T \alpha A + v^T \beta B) v = v^T \alpha A v + v^T \beta B v = \alpha (v^T A v) + \beta (v^T B v)$$

$\geq 0 \qquad \geq 0$

$$= \alpha \cdot \delta_1 + \beta \cdot \delta_2 \geq 0$$

↓

כזו נראה כי קבוצת המטריצות מאגל חזק את ספם לא מורח וקטור:

אגור $\alpha, \beta, \delta_1, \delta_2 \geq 0$ (בהמשך גבינה עזר) -

$$\alpha \cdot \delta_1 + \beta \cdot \delta_2 \geq 0, \quad \alpha, \beta \geq 0, \quad \delta_1, \delta_2 \geq 0$$

אנחנו לא אלה נקבל כי נניח בלתי-שלילי \leftarrow כבוד לא מקיף $v^T (\alpha A + \beta B) v$ כוונת C .

אני לא מורח וקטור: □

Calculus and Probability

1. (15 pts) Let X_1, \dots, X_n be i.i.d $U([0,1])$ (uniform) continuous random variables. Let $Y = \max(X_1, \dots, X_n)$.

(a) What is the PDF of Y ? Write the mathematical formula and plot the PDF as well. Compute $E[Y]$ and $\text{Var}[Y]$ - how do they behave as a function of n as n grows large?

$$P(Y \leq x) = P(X_1 \leq x, X_2 \leq x, \dots, X_n \leq x) = \prod_{i=1}^n P(X_i \leq x) = F_x^n(x)$$

צפייה

$$F_x(y) = \begin{cases} 0 & y < 0 \\ y^n & 0 \leq y \leq 1 \\ 1 & y > 1 \end{cases}$$

צפייה

(F_x של Y : צפייה) $f_x(x) = \begin{cases} 0 & y < 0 \\ n \cdot y^{n-1} & y \in [0,1] \\ 0 & y > 1 \end{cases}$ $y \in [0,1]$: PDF

$$E[Y] = \int_0^1 x \cdot f_Y(x) dx = \int_0^1 x \cdot n \cdot x^{n-1} dx = n \cdot \int_0^1 x^n dx = n \left[\frac{x^{n+1}}{n+1} \right]_0^1 = \frac{n}{n+1}$$

$$\text{Var}(Y) = E[Y^2] - (E[Y])^2$$

$$E[Y^2] = \int_0^1 x^2 f_Y(x) dx = \int_0^1 x^2 n x^{n-1} dx = n \int_0^1 x^{n+1} dx = n \left[\frac{x^{n+2}}{n+2} \right]_0^1 = \frac{n}{n+2}$$

$$\Rightarrow \text{Var}(Y) = \frac{n}{n+2} - \left(\frac{n}{n+1} \right)^2 = \frac{(n(n+1)^2 - n^2(n+2))}{(n+1)^2}$$

$$= \frac{n^3 + 2n^2 + n - n^3 - 2n^2}{(n+2)(n+1)^2}$$

$$= \frac{n}{(n+2)(n+1)^2}$$

כאן נחזיר את הבעיה.

סימטריה במרחב / סימטריה



Optimal Classifiers and Decision Rules

1. (15 pts)

- (a) Let X and Y be random variables where Y can take values in $\mathcal{Y} = \{1, \dots, L\}$. Let ℓ_{0-1} be the 0-1 loss function defined in class. Show that $h = \arg \min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}[\ell_{0-1}(Y, f(X))]$ is given by

$$h(x) = \arg \max_{i \in \mathcal{Y}} \mathbb{P}[Y = i | X = x]$$

$$P[X=x, Y=y] \ell_{0-1}(y, h(x))$$

הסתברות שיהיה $h(x)$ ויהיה x ויהיה y

$$P[X=x, Y=y] \cdot \ell_{0-1}(y, h(x)) =$$

הסתברות שיהיה $h(x)$ ויהיה x ויהיה y

$$\sum_{i \in \mathcal{Y}} P[X=x, Y=i] \cdot \ell_{0-1}(i, h(x)) = P[X=x] \left(\sum_{i \in \mathcal{Y}} P[Y=i | X=x] \cdot \ell_{0-1}(i, h(x)) \right)$$

$$= \begin{cases} \sum_{1 \leq j \leq L} P[Y=j | X=x] & h(x) = 1 \\ \vdots & \vdots \\ \sum_{L \neq j \leq L} P[Y=j | X=x] & h(x) = L \end{cases} \quad \leftarrow \text{(*)} \rightarrow \text{הסתברות שיהיה } h(x) \text{ ויהיה } x$$

הסתברות שיהיה $h(x)$ ויהיה x ויהיה y

$$h(x) = \arg \max_{i \in \mathcal{Y}} P[Y=i | X=x] \text{ כל } i \in \mathcal{Y} \Rightarrow P[Y=h(x) | X=x]$$

(b) Let X and Y be random variables where Y can take values in $\mathcal{Y} = \{0, 1\}$. Let Δ be the following asymmetric loss function:

$$\Delta(y, \hat{y}) = \begin{cases} 0 & y = \hat{y} \\ a & y = 0, \hat{y} = 1 \\ b & y = 1, \hat{y} = 0, \end{cases}$$

where $a, b \in (0, 1]$ (note that this loss function generalizes the 0-1 loss defined in class). Compute the optimal decision rule h for the loss function Δ , i.e. the decision rule which satisfies:

$$h = \arg \min_{f: \mathcal{X} \rightarrow \mathcal{Y}} \mathbb{E}[\Delta(Y, f(X))]$$

פתרון נתונה פונקציית אובדן, (b) פונקציית אובדן

נסה $x \in \mathcal{X}$ (נסה להיקלע את הביטוי)

$$P[X=x, Y=y] \Delta(y, h(x))$$

$$P[X=x, Y=y] \Delta(y, h(x)) = P[X=x, Y=0] \Delta(0, h(x)) + P[X=x, Y=1] \Delta(1, h(x)) =$$

$$= P[X=x] \cdot (P[Y=0 | X=x] \Delta(0, h(x)) + P[Y=1 | X=x] \Delta(1, h(x)))$$

$$= \begin{cases} a \cdot P[Y=0 | X=x] & h(x) = 1 \\ b \cdot P[Y=1 | X=x] & h(x) = 0 \end{cases}$$

$$h(x) = \arg \min_{y \in \{0, 1\}} (a \cdot y + b(1-y)) \cdot P[Y=1-y | X=x] \Leftarrow$$

$$\Rightarrow h(x) = \begin{cases} 1 & a \cdot P[Y=0 | X=x] \leq b \cdot P[Y=1 | X=x] \\ 0 & b \cdot P[Y=1 | X=x] \leq a \cdot P[Y=0 | X=x] \end{cases}$$

* נראה שזה נכון עבור כל $x \in \mathcal{X}$

2. (15 pts) Let X and Y be random variables where X can take values in some set \mathcal{X} and Y can take values in $\mathcal{Y} = \{0, 1\}$ (i.e. binary label space). Assume we wish to find a predictor $h : \mathcal{X} \rightarrow [0, 1]$ (note that the hypothesis can output any number between 0 and 1) which minimizes $\mathbb{E}[\Delta_{\log}(Y, h(X))]$, where Δ_{\log} is the following loss function known as the *log-loss*:

$$\Delta_{\log}(y, \hat{y}) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y}).$$

Find the predictor $h : \mathcal{X} \rightarrow [0, 1]$ which minimizes $\mathbb{E}[\Delta_{\log}(Y, h(X))]$.

Note: This loss function may seem odd at first, but it is very important and we'll discuss it further in the future.

פתרון

נחזיק $\hat{x} \in \mathcal{X}$ נקבע את הפונקציה

$$P[X = \hat{x} | Y = y] \cdot \Delta_{\log}(y, h(\hat{x}))$$

$$P[X = \hat{x}, Y = y] \cdot \Delta_{\log}(y, \hat{y}) =$$

לפי $h(\hat{x}) = \hat{y}$ (נקבע)

$$= P[X = x | Y = 0] \cdot \Delta_{\log}(0, \hat{y}) + P[X = x | Y = 1] \cdot \Delta_{\log}(1, \hat{y})$$

$$= P[X = x] \cdot (P[Y = 0 | X = x] \cdot \Delta_{\log}(0, \hat{y}) + P[Y = 1 | X = x] \cdot \Delta_{\log}(1, \hat{y}))$$

$$= P[Y = 0 | X = x] (-\log(\hat{y})) + P[Y = 1 | X = x] (-\log(1 - \hat{y}))$$

↓

עכשיו
נבדוק
כמה
(ענה)

$$\hat{y} = 1 \quad \text{אם} \quad P[Y = 0 | X = x] = 0$$

אם $P[Y = 0 | X = x] = 0$

$$\hat{y} = 0 \quad \text{אם} \quad P[Y = 1 | X = x] = 0$$

אם $P[Y = 1 | X = x] = 0$

$$\frac{P[Y = 0 | X = x]}{1 - \hat{y}} - \frac{P[Y = 1 | X = x]}{\hat{y}} = 0 \quad \text{נפתור את המשוואה}$$

$$\Rightarrow P[Y = 0 | X = x] \cdot \hat{y} = P[Y = 1 | X = x] (1 - \hat{y})$$

$$\Rightarrow \frac{(1 - \hat{y})}{\hat{y}} = \frac{P[Y = 0 | X = x]}{P[Y = 1 | X = x]} \Rightarrow \hat{y} = \left(\frac{P[Y = 0 | X = x]}{P[Y = 1 | X = x]} + 1 \right)^{-1}$$

$$= \left(\frac{P[Y=0|X=x] + P[Y=1|X=x]}{P[Y=1|X=x]} \right)^{-1} = P[Y=1|X=x]$$

$$= P[Y=0|X=x] \quad \square$$

3. (10 pts)

Let X and Y be random variables taking values in $\mathcal{X} = \mathbb{R}$ and $\mathcal{Y} = \{0,1\}$ respectively, and assume that given $Y = 0$, X is distributed normally with mean μ and variance σ_0^2 , i.e. $X \sim \mathcal{N}(\mu, \sigma_0^2)$, and similarly, given $Y = 1$, $X \sim \mathcal{N}(\mu, \sigma_1^2)$, where $\sigma_0 \neq \sigma_1$. Also assume $\Pr[Y=1] = p_1$.

Find the optimal decision rule for this distribution and the zero-one loss, i.e. find $h: \mathbb{R} \rightarrow \{0,1\}$ which minimizes $\mathbb{E}[\ell_{0-1}(Y, h(X))]$ where ℓ_{0-1} is the zero-one loss defined in class (write the decision rule only in terms of $x, \mu, \sigma_0, \sigma_1$ and p_1).

פתרון

ראינו כי ה-func decision הוא המקסימום של ה-0-1 loss במקרה זה:

$$h(x) = \begin{cases} 1 & P[Y=1|X=x] > P[Y=0|X=x] \\ 0 & \text{אחרת} \end{cases}$$

נחנא עמר אילו עובד המקסימום

$$P[Y=1|X=x] > P[Y=0|X=x]$$

$$\frac{f_X(x|Y=1)P[Y=1]}{f_X(x)} > \frac{f_X(x|Y=0)P[Y=0]}{f_X(x)}$$

\Leftrightarrow ארשר אירות:

$$\frac{f_X(x|Y=1)P[Y=1]}{f_X(x)} > \frac{f_X(x|Y=0)P[Y=0]}{f_X(x)}$$

$$\frac{f_X(x|Y=1)}{f_X(x|Y=0)} > \frac{P[Y=0]}{P[Y=1]} \quad \Leftrightarrow$$

\Leftrightarrow

$$\frac{\frac{1}{\sqrt{2\pi\sigma_1^2}} \cdot e^{-\frac{(x-\mu)^2}{2\sigma_1^2}}}{\sqrt{2\pi\sigma_1^2}} > \frac{1-P_1}{P_1} \quad (\Leftarrow)$$

$$\frac{1}{\sqrt{2\pi\sigma_1^2}} \cdot e^{-\frac{(x-\mu)^2}{2\sigma_1^2}}$$

$$\frac{\sqrt{2\pi\sigma_0^2}}{\sqrt{2\pi\sigma_1^2}} \cdot e^{-\frac{(x-\mu)^2}{2\sigma_1^2} + \frac{(x-\mu)^2}{2\sigma_0^2}} = \frac{\sigma_0}{\sigma_1} \cdot e^{-\frac{(x-\mu)^2}{2\sigma_1^2} + \frac{(x-\mu)^2}{2\sigma_0^2}} >$$

$$> \frac{1-P_1}{P_1} \quad (\Leftarrow) \quad \frac{(x-\mu)^2}{2\sigma_0^2} - \frac{(x-\mu)^2}{2\sigma_1^2} > \log\left(\frac{1-P_1}{P_1} \cdot \frac{\sigma_0}{\sigma_1}\right)$$

$$\Leftrightarrow (x-\mu)^2 \left(\frac{1}{2\sigma_0^2} - \frac{1}{2\sigma_1^2}\right) > \log\left(\left(\frac{1-P_1}{P_1}\right) \cdot \frac{\sigma_0}{\sigma_1}\right)$$

הנחה: $\sigma_0 < \sigma_1$

$$b) \sigma_0 < \sigma_1 \Rightarrow \frac{1}{2\sigma_0^2} - \frac{1}{2\sigma_1^2} > 0 \Rightarrow (x-\mu)^2 \left(\frac{1}{2\sigma_0^2} - \frac{1}{2\sigma_1^2}\right) >$$

$$> \log\left(\left(\frac{1-P_1}{P_1}\right) \cdot \frac{\sigma_0}{\sigma_1}\right) \Leftrightarrow (x-\mu)^2 > \log\left(\left(\frac{1-P_1}{P_1}\right) \cdot \frac{\sigma_0}{\sigma_1}\right) \left(\frac{1}{\frac{1}{2\sigma_0^2} - \frac{1}{2\sigma_1^2}}\right)$$

$$(x-\mu)^2 > \log\left(\left(\frac{1-P_1}{P_1}\right) \cdot \frac{\sigma_0}{\sigma_1}\right) \left(\frac{1}{\frac{1}{2\sigma_0^2} - \frac{1}{2\sigma_1^2}}\right) \quad \text{כאשר } h(x) = \dots$$

$h(x) = 0$ כאשר

$$b) \sigma_0 > \sigma_1 \Rightarrow \frac{1}{2\sigma_0^2} - \frac{1}{2\sigma_1^2} < 0 \Rightarrow (X-\mu)^2 \left(\frac{1}{2\sigma_0^2} - \frac{1}{2\sigma_1^2} \right) >$$

$$> \text{p.g.} \left(\left(\frac{1-p_1}{p_1} \right) \frac{\sigma_0}{\sigma_1} \right) \Leftrightarrow (X-\mu)^2 < \text{p.g.} \left(\left(\frac{1-p_1}{p_1} \right) \cdot \frac{\sigma_0}{\sigma_1} \right) \left(\frac{1}{2\sigma_0^2} - \frac{1}{2\sigma_1^2} \right)$$

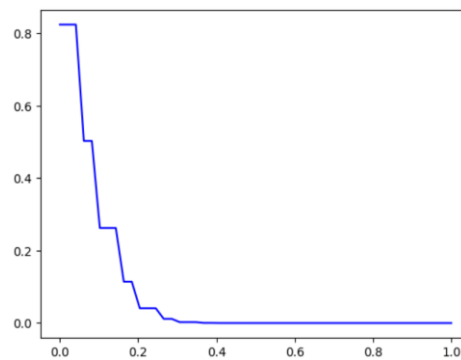
$$(X-\mu)^2 < \text{p.g.} \left(\left(\frac{1-p_1}{p_1} \right) \frac{\sigma_0}{\sigma_1} \right) \left(\frac{1}{2\sigma_0^2} - \frac{1}{2\sigma_1^2} \right) \text{ m.a. } h(x) = 1 \quad \text{if } |x - \mu| < \sqrt{\text{p.g.} \left(\left(\frac{1-p_1}{p_1} \right) \frac{\sigma_0}{\sigma_1} \right) \left(\frac{1}{2\sigma_0^2} - \frac{1}{2\sigma_1^2} \right)}$$

$$\square \quad , \quad h(x) = 0 \text{ m.a. } |x - \mu| > \sqrt{\text{p.g.} \left(\left(\frac{1-p_1}{p_1} \right) \frac{\sigma_0}{\sigma_1} \right) \left(\frac{1}{2\sigma_0^2} - \frac{1}{2\sigma_1^2} \right)}$$

חלק תכנותי:

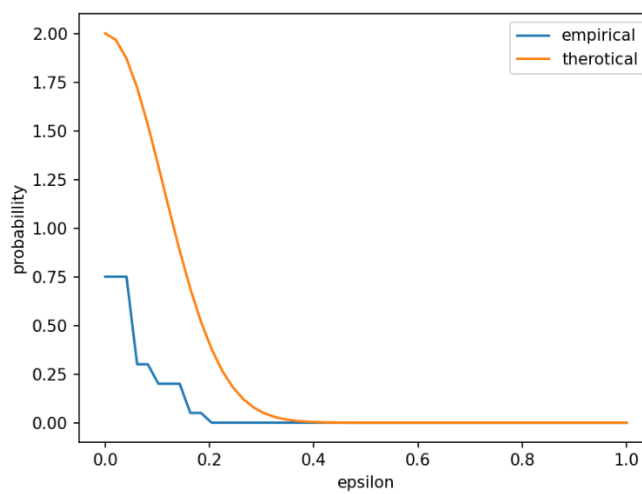
שאלה 1

סעיפים א+ב. מצורפת גרסה רציפה.



שאלה 1

סעיף ג'

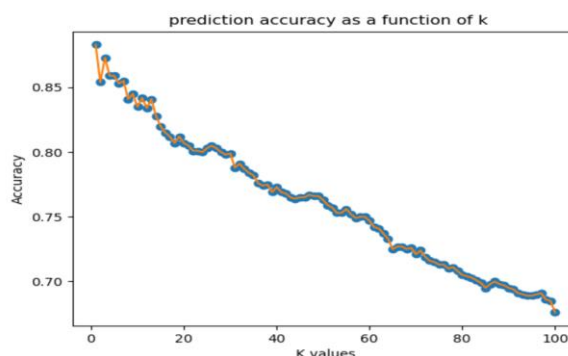


שאלה 2

עבור הפרמטרים איתם התבקשנו להריץ: $k=10$, $n=1000$. הדיוק שהתקבל הינו 0.846.

המשמעות הינה – 846 מתמונות המבחן קיבלו פרדיקציה נכונה. גודל המדגם של labels הינו 10 (הספרות 0-9) <- נצפה לדיוק של 0.1 ע"י predictor שיפעל באופן רדנומלי.

סעיף ג'



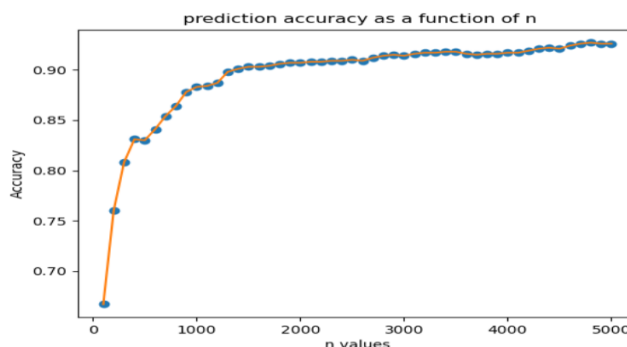
קל לראות כי ככל שערך K גדל <- רמת הדיוק של הפרדיקציה פוחתת.

נשים לב כי עבור $k=1$ אנו מקבלים את הדיוק הטוב ביותר.

הסבר אפשרי לכך: ככל שערך K גדל, תוצאת ההערכה (הפרדיקציה) תהיה מושפעת ע"י תמונות (ובפרט labels שלהן בהתאמה) אשר פחות קרובות לתמונה המקורית.

למשל עבור המקרה בו k הוא מספר כל התמונות, נקבל כי אין משמעות לקירוב הנל, שכן כל התמונות יהיו "קרובות" במידה שווה לתמונה שלנו.

סעיף ד'



לעומת הסעיף הקודם – כאן רמת הדיוק של הפרדיקציה גדלה ככל שערך n גדל.

הסבר מתבקש – ככל שמדגם המבחן שלנו גדול יותר (בדוגמא הנל יותר תמונות להתאמן עליהן) נקבל רמת דיוק גבוהה יותר.