

מערכות לומדות תרגיל 3

28 באפריל 2022

חלק תיאורטי

1. הוכיחו כי בעיית האופטימיזציה Hard-SVM היא בעיית תכנות ריבועית:

$$\operatorname{argmin}_{(w,b)} \|w\|^2 \quad \text{s.t.} \quad \forall i \ y_i (\langle w, x_i \rangle + b) \geq 1$$

כלומר, מצאו מטריצות Q ו- A ווקטורים a ו- d כך שהבעיה הנ"ל יכולה להיכתב כך:

$$\operatorname{argmin}_{v \in \mathbb{R}^n} \frac{1}{2} v^T Q v + a^T v \quad \text{s.t.} \quad A v \leq d$$

רמז: שימו לב כי $\|w\|^2 = w^T I w$

$$\begin{aligned} & \operatorname{argmin}_{(w,b)} \|w\|^2 \quad \text{s.t.} \quad \forall i \ y_i (\langle w, x_i \rangle + b) \geq 1 \\ & \operatorname{argmin}_{(w,b)} \begin{pmatrix} w \\ b \end{pmatrix}^T I \begin{pmatrix} w \\ b \end{pmatrix} \quad \text{s.t.} \quad \begin{pmatrix} (y_1 x_1) & y_1 \\ \vdots & \vdots \\ (y_m x_m) & y_m \end{pmatrix} \begin{pmatrix} w \\ b \end{pmatrix} \geq \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \\ & \operatorname{argmin}_{(w,b)} \begin{pmatrix} w \\ b \end{pmatrix}^T I \begin{pmatrix} w \\ b \end{pmatrix} \quad \text{s.t.} \quad \begin{pmatrix} (y_1 x_1) & y_1 \\ \vdots & \vdots \\ (y_m x_m) & y_m \end{pmatrix} \begin{pmatrix} w \\ b \end{pmatrix} \geq \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} \end{aligned}$$

אם נכפול ונחלק את הביטוי הראשון ב-2 (אין שינוי), נגדיר $a = \vec{0}$ ונכפיל את הביטוי השני בשתי צדדיו ב-1 נקבל:

$$\begin{aligned} & \operatorname{argmin}_{(w,b)} \frac{1}{2} \begin{pmatrix} w \\ b \end{pmatrix}^T 2 \cdot I \begin{pmatrix} w \\ b \end{pmatrix} + 0^T \begin{pmatrix} w \\ b \end{pmatrix} \quad \text{s.t.} \quad - \begin{pmatrix} (y_1 x_1) & y_1 \\ \vdots & \vdots \\ (y_m x_m) & y_m \end{pmatrix} \begin{pmatrix} w \\ b \end{pmatrix} \leq \begin{pmatrix} -1 \\ \vdots \\ -1 \end{pmatrix} \\ & .d = \begin{pmatrix} -1 \\ \vdots \\ -1 \end{pmatrix} \quad \text{ו-} \quad A = - \begin{pmatrix} (y_1 x_1) & y_1 \\ \vdots & \vdots \\ (y_m x_m) & y_m \end{pmatrix}, Q = 2 \cdot I, v = \begin{pmatrix} w \\ b \end{pmatrix} \end{aligned}$$

לכן קיבלנו:

2. בעיית Soft-SVM אופטימיזציה:

$$\operatorname{argmin}_{w, \{\varepsilon_i\}} \frac{\lambda}{2} \|w\|^2 + \frac{1}{m} \sum \varepsilon_i \quad \text{s.t.} \quad \forall i \quad y_i \langle w, x_i \rangle \geq 1 - \varepsilon_i \wedge \varepsilon_i \geq 0$$

סמנו את פונקציית $hinge-loss$ כ- $\ell^{hinge}(a) := \max\{0, 1 - a\}$. הראו כי בעיית האופטימיזציה Soft-SVM היא שקולה לבעיית האופטימיזציה ללא אילוצים הבאה:

$$\operatorname{argmin}_{w, \{\varepsilon_i\}} \frac{\lambda}{2} \|w\|^2 + \frac{1}{m} \sum_i \ell^{hinge}(y_i \langle w, x_i \rangle)$$

נשים לב כי התנאי $1 - \varepsilon_i \wedge \varepsilon_i \geq 0$ הוא בדיוק התנאי בפונקציה ℓ^{hinge} . אם $\varepsilon_i \geq 0$ ומתקיים $y_i \langle w, x_i \rangle \geq 1 - \varepsilon_i$ אזי $\ell^{hinge}(y_i \langle w, x_i \rangle) = 1 - y_i \langle w, x_i \rangle = \varepsilon_i$. ואחרת יהיה 0.

3. א. הניחו $x \in \mathbb{R}^d$. בהינתן $\text{trainset } S = \{(x_i, y_i)\}_{i=1}^m$, התאימו מסווג Gaussian Naive Bayes הפותר את (5) תחת הנחות (6).

נרצה למצוא θ הממקסמת את הנראות.

$$\begin{aligned} L(\theta | X, y) &= f_{X, y | \theta}(\{(x_i, y_i)\}_{i=1}^m) \stackrel{iid}{=} \prod_{i=1}^m f_{X, y | \theta}(X_i, y_i) \\ &= \prod_{i=1}^m f(y_i | \theta) \cdot f(x_i | y = y_i) = \prod_{i=1}^m \pi_{y_i} \cdot N(X_i | \mu_{y_i}, \sigma_{y_i}^2) \\ &= \prod_{i=1}^m \pi_{y_i} \cdot \frac{1}{\sqrt{2\pi(\sigma_{y_i})^2}} \cdot \exp\left(-\frac{(X_i - \mu_{y_i})^2}{2(\sigma_{y_i})^2}\right) \end{aligned}$$

נפעיל \log על הביטוי כדי לפשט אותו.

$$\begin{aligned} &= \sum_{i=1}^m \log(\pi_{y_i}) + \log\left(\frac{1}{\sqrt{2\pi\sigma_{y_i}^2}}\right) - \frac{(X_i - \mu_{y_i})^2}{2\sigma_{y_i}^2} \\ &= \sum_{i=1}^m \log(\pi_{y_i}) - \frac{1}{2} \log(2\pi\sigma_{y_i}^2) - \frac{(X_i - \mu_{y_i})^2}{2\sigma_{y_i}^2} \\ &= \sum_{i=1}^m \log(\pi_{y_i}) - \frac{1}{2} \log(2\pi) - \frac{1}{2} \log(\sigma_{y_i}^2) - \frac{(X_i - \mu_{y_i})^2}{2\sigma_{y_i}^2} \end{aligned}$$

נרצה לחשב את θ הממקסמת את הנראות, לכן

$$\begin{aligned}\operatorname{argmax}_{\theta} \ell(\theta, S) &= \operatorname{argmax}_{\theta} \sum_{i=1}^m \log(\pi_{y_i}) - \frac{1}{2} \log(\sigma_{y_i}^2) - \frac{(X_i - \mu_{y_i})^2}{2\sigma_{y_i}^2} \\ &= \operatorname{argmax}_{\theta} \sum_k \left(n_k \cdot \log(\pi_k) - \frac{1}{2} n_k \cdot \log(\sigma_k^2) - \sum_{i:y_i=k} \frac{(X_i - \mu_k)^2}{2\sigma_k^2} \right)\end{aligned}$$

נגזור את הביטוי וכל פעם עבור משתנה אחר: $\{\pi_k\}, \{\mu_k\}, \{\sigma_k\}$
 $\hat{\pi}_k = \frac{n_k}{m}$: נקבל מהגזירה אותו דבר כמו החישוב בתרגול (כי יש את אותו ביטוי עם π_k)
 $\hat{\mu}_k = \frac{1}{n_k} \sum_{i:y_i=k} X_i$: נקבל מהגזירה אותו דבר כמו החישוב בתרגול (כי יש את אותו ביטוי עם μ_k)
 $\hat{\sigma}_k^2 = \frac{1}{n_k} \sum_{i:y_i=k} \frac{(X_i - \mu_k)^2}{2}$: נקבל מהגזירה אותו דבר כמו החישוב בתרגול (כי יש את אותו ביטוי עם σ_k^2)

$$\begin{aligned}\frac{\partial \ell}{\partial \mu_k} &= \sum_{i:y_i=k} \frac{-2(X_i - \mu_k)}{2\sigma_k^2} = 0 \\ \sum_{i:y_i=k} -2(X_i - \mu_k) &= 0 \\ \sum_{i:y_i=k} X_i - n_k \cdot \mu_k &= 0 \\ \mu_k &= \frac{1}{n_k} \cdot \sum_{i:y_i=k} X_i\end{aligned}$$

: σ_k^2

$$\begin{aligned}\frac{\partial \ell}{\partial \sigma_k^2} &= -\frac{1}{2} \frac{n_k}{\sigma_k^2} - \sum_{i:y_i=k} \frac{(X_i - \mu_k)^2}{2} \cdot \left(-\frac{1}{\sigma_k^4} \right) = 0 \\ \frac{n_k}{\sigma_k^2} &= \sum_{i:y_i=k} (X_i - \mu_k)^2 \cdot \frac{1}{\sigma_k^4} \\ n_k &= \sum_{i:y_i=k} (X_i - \mu_k)^2 \cdot \frac{1}{\sigma_k^2} \\ \hat{\sigma}_k^2 &= \sum_{i:y_i=k} \frac{(X_i - \mu_k)^2}{n_k}\end{aligned}$$

ב. הניחו $x \in \mathbb{R}^d$ (לכל דגימה יש d פיצ'רים). בהניתן $\text{trainset } S = \{(X_i, y_i)\}_{i=1}^m$, התאימו מסווג Gaussian Naive Bayes הפותר את (5) תחת הנחות (6).

נשנה את ההנחה ל- $X \in \mathbb{R}^d$. כלומר, כעת כל דגימה היא בעלת d פיצ'רים, כאשר כל פיצ'ר X_{ij} מתפלג $X_{ij} \mid y_i = k \sim N(\mu_{kj}, \sigma_{kj}^2)$. בנוסף $y \sim \text{Mult}(\pi)$, $\sum \pi_i = 1$, $\pi = (\pi_1, \dots, \pi_n)$.
 נכתוב שוב את פונקציית הנראות עבור $\theta = (\{\pi_k\}, \{\mu_k\}, \{\sigma_{kj}^2\})$

$$\begin{aligned}\ell(\theta | \{X_i, y_i\}) &= f_{X,y|\theta}(\{X_i, y_i\}_{i=1}^m) \stackrel{iid}{=} \prod_{i=1}^m f_{X,y|\theta}(X_i, y_i) \\ &= \prod_{i=1}^m \prod_{j=1}^d f_{X_i, y|\theta}(X_{ij}, y_i) = \prod_{i=1}^m \prod_{j=1}^d f_y(y_i, \theta) \cdot f(X_{ij} | y = y_i) \\ &= \prod_{i \in [m], j \in [d]} f_y(y_i | \theta) \cdot f(X_{ij} | y = y_i)\end{aligned}$$

נגדיר $S' = \{(X_{ij}, y_i) : i \in [m], j \in [d]\}$: π_k

אם נסמן ב- n_k את כמות הדגימות בסט S המקורי ששייכות למחלקה k אז בקבוצה S' קיימות $n_k \cdot d$ דגימות כאלה, שכן כל פיצ'ר של כל אחת מהדגימות הנ"ל קיבל כעת את התווית y_i .

לכן ראינו $\hat{\pi}_k = \frac{|\{y_i=k\}|}{\text{num of samples}}$ לכן עבור הסט S' נקבל $\hat{\pi}_k = \frac{n_k \cdot d}{m \cdot d} = \frac{n_k}{m}$: μ_{k_j}

$$\hat{\mu}_{k,j} = \frac{1}{n_{k,j}} \sum_{i:y_i=k} X_{ij} = \frac{1}{n_k} \sum_{i:y_i=k} X_{ij}$$

: $\sigma_{k_j}^2$

$$\sigma_{k_j}^2 = \frac{1}{n_{k,j}} \cdot \sum_{i:y_i=k} (X_{ij} - \mu_{k,j})^2 = \frac{1}{n_k} \cdot \sum_{i:y_i=k} (X_{ij} - \mu_{k_j})^2$$

4. א. הניחו $x \in \mathbb{R}^d$. בהינתן $\text{trainset } S = \{(x_i, y_i)\}_{i=1}^m$, התאימו מסווג Gaussian Naive Bayes הפותר את (5) תחת הנחות (7).

נרצה למצוא θ הממקסמת את הנראות.

$$\begin{aligned}L(\theta | X, y) &= f_{X,y|\theta}(\{(x_i, y_i)\}_{i=1}^m) \stackrel{iid}{=} \prod_{i=1}^m f_{X,y|\theta}(X_i, y_i) \\ &= \prod_{i=1}^m f(y_i | \theta) \cdot f(x_i | y = y_i) = \prod_{i=1}^m \pi_{y_i} \cdot \text{Poi}(X_i | \lambda_{k,i}) \\ &= \prod_{i=1}^m \pi_{y_i} \cdot \frac{e^{-\lambda_{k,i}} \lambda_{k,i}^{y_i}}{y_i!}\end{aligned}$$

נפעיל \log כדי לפשט את הביטוי

$$\begin{aligned} \prod_{i=1}^m \pi_{y_i} \cdot \frac{e^{-\lambda} \lambda^k}{k!} &= \log(\pi_{y_i}) + \log\left(\frac{e^{-\lambda_{k,i}} \lambda_{k,i}^{y_i}}{y_i!}\right) \\ &= \sum_{i=1}^m \log(\pi_{y_i}) - \lambda_{k,i} + y_i \cdot \log(\lambda_{k,i}) - \log(y_i!) \end{aligned}$$

נרצה לחשב את θ הממקסמת את הנראות, לכן

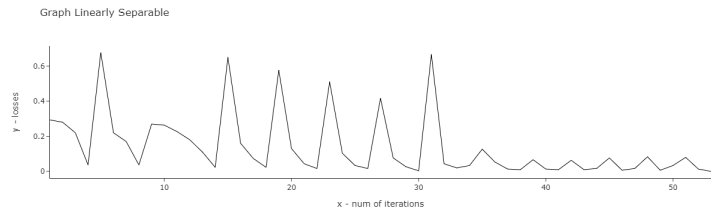
$$\begin{aligned} \operatorname{argmax}_{\theta} \ell(\theta, S) &= \operatorname{argmax}_{\theta} \sum_{i=1}^m \log(\pi_{y_i}) - \lambda_{k,i} + y_i \cdot \log(\lambda_{k,i}) - \log(y_i!) \\ &= \operatorname{argmax}_{\theta} \sum_{k=1}^d n_k \cdot \log(\pi_k) + \sum_{i:y_i=k} -\lambda_{k,i} + y_i \cdot \log(\lambda_{k,i}) - \log(y_i!) \end{aligned}$$

נקבל מהגזירה אותו דבר כמו החישוב בתרגול (כי יש את אותו ביטוי עם π_k) לכן $\hat{\pi}_k = \frac{n_k}{m}$ $\lambda_{k,i}$

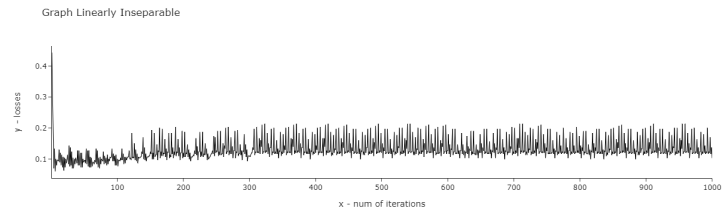
$$\begin{aligned} \frac{\partial \ell}{\partial \lambda_{k,i}} &= \sum_{i:y_i=k} -1 + y_i \cdot \frac{1}{\lambda_{k,i}} = 0 \\ -n_k + n_k \cdot y_i \cdot \frac{1}{\lambda_{k,i}} &= 0 \\ \lambda_{k,i} &= y_i \end{aligned}$$

חלק פרקטי

Perceptron Classifier

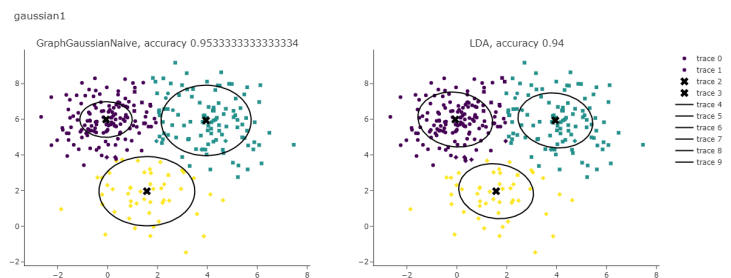


אנו יכולים להסיק מהגרף כי עבור המקרה בו הנתונים יכולים להיחלק לינארית, האלגוריתם יכול להגיע לתוצאה טובה (להפריד את הנתונים בצורה מושלמת) אחרי מספר איטרציות הקטן ממספר האיטרציות המקסימלי.

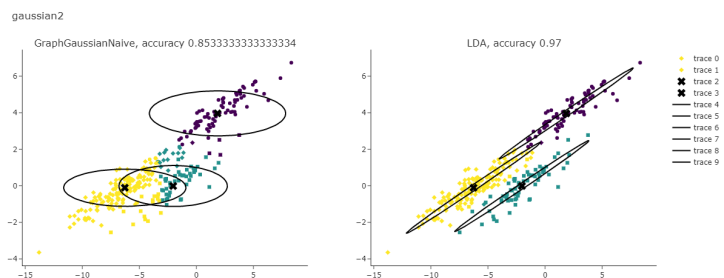


אנו יכולים להסביר את ההבדל בין שני הגרפים בכך שהגרף הזה עובד על נתונים שאינם יכולים להיפרד ע"י חלוקה לינארית. לכן האלגוריתם perceptron אינו יכול להגיע לתשובה טובה אף פעם. לכן הוא פעל עד מספר האיטרציות הסופי - שבמקרה שלנו 1000.

Bayes Classifiers



אנו יכולים להסיק שעבור הדאטא הראשון, לשני המסווגים תוצאות דומות.



ניתן להסיק מהגרף הזה כי לdataset הנוכחי כדאי להשתמש בLDA משום שהוא סיווג יותר טוב, עם $\text{accuracy} = 0.97$. ניתן להבין מכך שהsamples בדאטא היו תלויים לינארית. לכן החיזוי של Gaussian Naive היה פחות טוב כי הוא עובד עם מידע בלתי תלוי לינארית.