

למידה עמוקה 046211

גליון יבש 2

315681593

208731422

גל גרנות

ניר טבת

שאלה 1

הגדרנו את ה-Bayes Risk ואת הסיכון של w_μ בתור:

$$\bar{R}(w) = E_{\epsilon, w_{true}}(R), \quad R(w_\mu) = \|w_\mu - w_{true}\|^2 = \|(H_\mu^{-1}H - I)w_{true} + H_\mu^{-1}X^T\epsilon\|^2$$

נוכיח כי:

$$\bar{R}(w_\mu) = \sum_{i=1}^d \frac{(\sigma_w^2/d)\mu^2 + \sigma_\epsilon^2\lambda_i}{(\lambda_i + \mu)^2}$$

עבור:

$$1. \quad \epsilon \sim N(0, \sigma_\epsilon^2 I)$$

$$2. \quad w_{true} \sim N\left(0, \frac{\sigma_w^2}{d} I\right)$$

$$3. \quad H_\mu = \mu I + X^T X$$

$$4. \quad H = X^T X$$

$$5. \quad \lambda_i \text{ הע"ע של המטריצה } H$$

$$6. \quad w_\mu = H_\mu^{-1} X^T y = (\mu I + X^T X)^{-1} X^T y$$

הוכחה:

$$\begin{aligned}
 & \| (H_\mu^{-1}H - I)w_{true} + H_\mu^{-1}X^T\epsilon \|^2 \\
 &= \left((H_\mu^{-1}H - I)w_{true} + H_\mu^{-1}X^T\epsilon \right)^T \left((H_\mu^{-1}H - I)w_{true} + H_\mu^{-1}X^T\epsilon \right) \\
 &= \left(\left((H_\mu^{-1}H - I)w_{true} \right)^T + \left(H_\mu^{-1}X^T\epsilon \right)^T \right) \left((H_\mu^{-1}H - I)w_{true} + H_\mu^{-1}X^T\epsilon \right) \\
 &= \| (H_\mu^{-1}H - I)w_{true} \|^2 + \| H_\mu^{-1}X^T\epsilon \|^2 + \text{cross terms} = *
 \end{aligned}$$

מכיוון ש- ϵ ו- w_{true} חסרי קורלציה ($E(\epsilon w_{true}) = 0$) איברי ה-cross יתאפסו תחת התוחלת ולכן נתעלם מהם. נחשב בנפרד את התוחלת עבור הנורמות של שני האיברים:

$$\begin{aligned}
 &= \| (H_\mu^{-1}H - I)w_{true} \|^2 = \left((H_\mu^{-1}H - I)w_{true} \right)^T (H_\mu^{-1}H - I)w_{true} \\
 &= w_{true}^T (H_\mu^{-1}H - I)^T (H_\mu^{-1}H - I)w_{true}
 \end{aligned}$$

הביטוי המתקבל הוא סקלר, או מטריצה מסדר 1×1 ולכן שווה לעקבה של עצמו. נוכיח זהות קצרה בשימוש בזהות הציקלית של העקבה:

$$a^T b = \text{Tr}(ba^T) \Rightarrow a^T M a = \text{Tr}(Maa^T)$$

נשתמש בזהות עם $a = w_{true}$ ו- $M = (H_\mu^{-1}H - I)^T (H_\mu^{-1}H - I)$:

$$w_{true}^T (H_\mu^{-1}H - I)^T (H_\mu^{-1}H - I)w_{true} = \text{Tr} \left((H_\mu^{-1}H - I)^T (H_\mu^{-1}H - I)w_{true}w_{true}^T \right)$$

נפעיל את התוחלת ונשתמש בלינאריות העקבה והתוחלת, כאשר $E(w_{true}w_{true}^T)$ היא מטריצת הקוואריאנס של הוקטור האקראי w_{true} שלפי הנתון היא $\frac{\sigma_w^2}{d}I$:

$$E \left(\text{Tr} \left((H_\mu^{-1}H - I)^T (H_\mu^{-1}H - I)w_{true}w_{true}^T \right) \right) = \text{Tr} \left(E \left((H_\mu^{-1}H - I)^T (H_\mu^{-1}H - I)w_{true}w_{true}^T \right) \right)$$

$$\text{Tr} \left((H_\mu^{-1}H - I)^T (H_\mu^{-1}H - I)E(w_{true}w_{true}^T) \right) = \text{Tr} \left((H_\mu^{-1}H - I)^T (H_\mu^{-1}H - I) \frac{\sigma_w^2}{d} I \right)$$

$$\frac{\sigma_w^2}{d} \text{Tr} \left((H_\mu^{-1}H - I)^T (H_\mu^{-1}H - I) \right)$$

נחזור לאיבר השני ב-*, ונשתמש בעובדה ש- H_μ^{-1} סימטרית ולכן שווה לשחלוף שלה:

$$\| H_\mu^{-1}X^T\epsilon \|^2 = (H_\mu^{-1}X^T\epsilon)^T H_\mu^{-1}X^T\epsilon = \epsilon^T XH_\mu^{-1}H_\mu^{-1}X^T\epsilon$$

$$E \left(\| H_\mu^{-1}X^T\epsilon \|^2 \right) = E(\epsilon^T XH_\mu^{-1}H_\mu^{-1}X^T\epsilon) = \sum_{i=1}^N \sum_{j=1}^N E(\epsilon_i\epsilon_j) (XH_\mu^{-1})_i (H_\mu^{-1}X^T)_j$$

$$= \sum_{i=1}^N \sum_{j=1}^N \delta_{ij} \sigma_\epsilon^2 (XH_\mu^{-1})_i (H_\mu^{-1}X^T)_j = \sigma_\epsilon^2 \sum_{i=1}^N (XH_\mu^{-1})_i (H_\mu^{-1}X^T)_i = \sigma_\epsilon^2 \text{Tr}(XH_\mu^{-2}X^T)$$

נאחד את התוצאות שקיבלנו:

$$\bar{R}(w_\mu) = \frac{\sigma_w^2}{d} \text{Tr} \left((H_\mu^{-1}H - I)^T (H_\mu^{-1}H - I) \right) + \sigma_\epsilon^2 \text{Tr}(XH_\mu^{-2}X^T)$$

H היא מטריצה סימטרית ולכן קיים לה לכסון אורתוגונלי $U\Lambda U^T$ עבור U אורתוגונלית ו- Λ אלכסונית. נקבל:

$$\begin{aligned} &= \frac{\sigma_w^2}{d} \text{Tr} \left((H_\mu^{-1}H - I)^2 \right) + \sigma_\epsilon^2 \text{Tr}(XX^T H_\mu^{-2}) \\ &= \frac{\sigma_w^2}{d} \text{Tr}(((H + \mu I)^{-1}H - I)^2) + \sigma_\epsilon^2 \text{Tr}(H(H + \mu I)^{-2}) \\ &= \frac{\sigma_w^2}{d} \text{Tr}((U\Lambda U^T + \mu I)^{-1}U\Lambda U^T - I)^2 + \sigma_\epsilon^2 \text{Tr}(U\Lambda U^T (U\Lambda U^T + \mu I)^{-2}) \end{aligned}$$

נשתמש בעובדה כי $UU^T = I$ ובחילוף לכפל של מטריצת היחידה μI , נכפול משמאל וימין בהתאמה:

$$\begin{aligned} &= \frac{\sigma_w^2}{d} \text{Tr}((U\Lambda U^T + U\mu I U^T)^{-1}U\Lambda U^T - I)^2 + \sigma_\epsilon^2 \text{Tr}(U\Lambda U^T (U\Lambda U^T + U\mu I U^T)^{-2}) \\ &= \frac{\sigma_w^2}{d} \text{Tr}((U(\Lambda + \mu I)U^T)^{-1}U\Lambda U^T - I)^2 + \sigma_\epsilon^2 \text{Tr}(U\Lambda U^T (U\Lambda U^T + U\mu I U^T)^{-2}) \\ &= \frac{\sigma_w^2}{d} \text{Tr} \left(U((\Lambda + \mu I))^{-1}U^T U\Lambda U^T - I \right)^2 + \sigma_\epsilon^2 \text{Tr}(U\Lambda U^T U(\Lambda + \mu I)^{-2}U^T) \\ &= \frac{\sigma_w^2}{d} \text{Tr}(U(\Lambda + \mu I)^{-1}U^T - I)^2 + \sigma_\epsilon^2 \text{Tr}(U\Lambda(\Lambda + \mu I)^{-2}U^T) \\ &= \frac{\sigma_w^2}{d} \text{Tr}(U(\Lambda + \mu I)^{-1}\Lambda U^T - I)^2 + \sigma_\epsilon^2 \text{Tr}(U\Lambda(\Lambda + \mu I)^{-2}U^T) \end{aligned}$$

המטריצה $\Lambda + \mu I$ אלכסונית ולכן ההופכית שלה היא אלכסונית עם כל האלמנטים ההפוכים:

$$\begin{aligned} \Lambda + \mu I &= \text{diag}(\lambda_i) + \text{diag}(\mu) = \text{diag}(\lambda_i + \mu) \\ &= \frac{\sigma_w^2}{d} \text{Tr} \left(U \cdot \text{diag} \left(\frac{1}{\lambda_i + \mu} \right) \cdot \text{diag}(\lambda_i) U^T - I \right)^2 + \sigma_\epsilon^2 \text{Tr} \left(U \cdot \text{diag} \left(\frac{\lambda_i}{(\lambda_i + \mu)^2} \right) U^T \right) \\ &= \frac{\sigma_w^2}{d} \text{Tr} \left(\text{diag} \left(\frac{\lambda_i}{\lambda_i + \mu} - 1 \right) \right)^2 + \sigma_\epsilon^2 \text{Tr} \left(\text{diag} \left(\frac{\lambda_i}{(\lambda_i + \mu)^2} \right) \right) \\ &= \frac{\sigma_w^2}{d} \sum_{i=1}^d \left(\frac{\lambda_i - \lambda_i - \mu}{\lambda_i + \mu} \right)^2 + \sigma_\epsilon^2 \sum_{i=1}^d \frac{\lambda_i}{(\lambda_i + \mu)^2} = \sum_{i=1}^d \left(\frac{\frac{\sigma_w^2}{d} \mu^2}{(\lambda_i + \mu)^2} + \frac{\sigma_\epsilon^2}{(\lambda_i + \mu)^2} \right) \\ &\quad \sum_{i=1}^d \frac{\frac{\sigma_w^2}{d} \mu^2 + \sigma_\epsilon^2 \lambda_i}{(\lambda_i + \mu)^2} \end{aligned}$$

כנדרש.

שאלה 2

1. נתון ש- w נדגם בצורה אחידה מתוך Q אשר מכיל q ערכים. וכן נתון שכל רכיב בוקטור w נדגם באופן i.i.d. לשאר הרכיבים.

$$P_{w \sim P_w}(w = w_*) = \frac{1}{q}$$

אם w היה וקטור באורך 1, קל היה לראות שמתקבל: $\frac{1}{q}$

$$P_{w \sim P_w}(w = w_*) = \left(\frac{1}{q}\right)^k$$

2. כדי שהביטוי יתקיים נצטרך שכל המשקולות יתאימו בין הרשתות. נשים לב שיש במטריצה $d_0 * d_1, W_1$ משקולות, וכן במטריצה W_2 יש d_1 משקולות. כמו שראינו בסעיף 1, לכל משקולת יש הסתברות של $\frac{1}{q}$ להתאים בין שתי הרשתות. כמו כן אנחנו יודעים שיש רק $d_1 < d_{1*}$ משקולות שאינן אפס בשכבה 1. סך הכל נקבל:

$$p_* \geq \left(\frac{1}{q}\right)^{d_0 d_{1*} + d_1} = q^{-d_1 - d_0 d_{1*}}$$

3. T הוא זמן העצירה של הרשת. ניתן להסתכל על T כעל רצף של ניסויי ברנולי עם פרמטר הצלחה p_* שהוגדר בסעיף 2. כלומר T מתפלג גאומטרית. לכן ניתן לומר:

$$P(t > T) = (1 - p_*)^T$$

נפעיל לוג על שני האגפים:

$$\log(P(t > T)) = T \cdot \log(1 - p_*)$$

$$T = \frac{\log(P(t > T))}{\log(1 - p_*)} \rightarrow [T] \leq \frac{\log(P(t > T))}{\log(1 - p_*)}$$

4. משילוב של שני הקירובים שנתונים ברמז והצבה שלהם בתוצאה של סעיף 3 נקבל:

$$T = \frac{\log(P(t > T))}{-p_*}$$

נציב את התוצאה מסעיף 1 ונקבל:

$$T \leq -(\log(P(t > T))) q^{d_1 + d_0 d_{1*}}$$

$$\eta = \log(P(t > T))$$

כמו כן במקרה שלנו:

$$\log(|\mathcal{F}|) = \log T \leq \log\left(\frac{\log \eta}{-p_*}\right) = \log \log \eta + \log p_* \leq \log \log \left(\frac{1}{\eta}\right) + (d_0 d_{1*} + d_1) \log q$$

נעת נציב הכל במשפט 2 ונקבל:

$$\epsilon < \frac{\log \log \left(\frac{1}{\eta}\right) + (d_0 d_{1*} + d_1) \log q + \log \left(\frac{1}{\delta}\right)}{N}$$

כנדרש.

5. ביטוי 4 עבור משקולות שיכולות לקבל q ערכים שונים נקבל:

$$\log |\mathcal{F}| = k \log q = (d_0 d_1 + d_1) \log q$$

נקבל ביטוי שגדול יותר מהביטוי בסעיף 3 ולכן ההכללה של ביטוי 4 תהיה פחות טובה.

שאלה 3

1. בכל הגרפים הנקודה הקריטית היא הנקודה שבה ה- test loss מפסיק לעלות ומתחיל לרדת שוב.
2.
 - a. הגרף הכחול והירוק הם שני מודלים שונים, הכחול הוא מודל שמתרגם מגרמנית לאנגלית והירוק הוא מודל שמתרגם מאנגלית לצרפתית. אנו יכולים לראות כי ה- test loss לאחר הגדלת גודל המודל מתחיל לעלות בשלב מסוים ולאחר גודל מסוים מתחיל לרדת, בנוסף ניתן לראות כי ה- train loss מונוטוני יורד עם הגדלת גודל המודל ולכן זהו model-wise double decent מכיוון שהאזור הקריטי מופיע ב- test loss .
 - b. אנו יכולים לראות כי זהו $\text{model-wise double decent}$ ככל שכמות הפרמטרים במודל גדלה כך השגיאה קטנה, עבור מודלים גדולים, בניגוד לשגיאה הקלאסית. הנקודה הקריטית מתקבלת בנק' מקסימום בה הגרף מתחיל לרדת לאחר העלייה (בערך ברוחב 10 פרמטרים)
 - c. אנו יכולים לראות כי $\text{epoch-wise double decent}$ מכיוון שעבור המודל הגדול (האדום) אנחנו רואים שעלייה בכמות ה- EPOCHS גורמת לירידה בשגיאה.

שאלה 4

1. ההתפלגות של W סימטרית <- ההתפלגות של U סימטרית:

$$E[\varphi^2(u_{l-1})] = E[\max(0, u_{l-1})^2] = \frac{1}{2} E[u_{l-1}^2] = \frac{1}{2} \sigma_{u_{l-1}}^2$$

ניתן לראות כי הביטוי המחושב זהה לחישוב השונות של u אך החלק השלילי מאופס. מכיוון שההתפלגות של u סימטרית, החלק השלילי תרומה שווה לתרומה של החלק החיובי ונקבל חצי מהשונות של u .

$$\sigma_l = \frac{1}{\sqrt{\sum_j E[\varphi^2(u_{l-1}[j])]}} = \frac{1}{\sqrt{\sum_j \frac{1}{2} \sigma_{u_{l-1}[j]}^2}} = \sqrt{\frac{2}{d_{l-1}}}$$

תחת ההנחה שהשונות של u_{l-1} היא 1.

2. לפי תאוריית הגבול המרכזי:

נדיר משתנה Z שמתפלג נורמלי עם תוחלת 0 ושונות 1.

$$\varphi(u_{l-1}) = \max(0, u_{l-1}) = \sigma_{u_{l-1}} \max(0, z)$$

$$E[\varphi^2(u_{l-1})] = E[(\sigma_{u_{l-1}} \max(0, z))^2] = \sigma_{u_{l-1}}^2 E[\max(0, z)^2] = \frac{1}{2} \sigma_{u_{l-1}}^2 \sigma_z = \frac{1}{2} \sigma_{u_{l-1}}^2$$

וההמשך זהה לסעיף 1.

שאלה 5

צד ראשון: אם מתקיים $\forall \tau \in H, W[i, j] = W[\tau(i), \tau(j)]$ אז f_W היא $equivariant$:

$$f_W(\tau x) = \phi(W(\tau x)) = \phi(\tau(Wx))$$

כאשר השוויון האחרון נובע מהנתון ומכך ש ϕ פועלת על כל אלמנט בנפרד.

מ- $equivariance$ של f ל- τ נקבל:

$$\phi(\tau(Wx)) = \tau\phi(Wx) = \tau \cdot f_W(x)$$

ובסך הכול קיבלנו $equivariance$ של f ל- H .

צד שני: אם f_W היא $equivariant$ אז מתקיים $\forall \tau \in H, W[i, j] = W[\tau(i), \tau(j)]$

$$f_W(\tau x) = \tau f_W(x) \rightarrow \phi(W\tau x) = \tau \cdot \phi(Wx)$$

מחד-חד ערכיות של ϕ ניתן לומר שחייב להתקיים שוויון בארגומנטים:

$$W\tau x = \tau(Wx) \rightarrow \tau W = W\tau$$

ומכאן ניתן להסיק שלכל איבר W מתקיים:

$$W[i, j] = W[\tau(i), \tau(j)]$$

שאלה 6

Layer	Output_dim	Number of parametes
INPUT	224X224X3	0
CONV3-64	224X224X64	$(3 \times 3 \times 3 + 1) \times 64 = 1792$
ReLU	224X224X64	0
POOL2	112X112X64	0
CONV3-128	112X112X128	$(3 \times 3 \times 64 + 1) \times 128 = 73856$
ReLU	112X112X128	0
POOL2	56X56X128	0
CONV3-256	56X56X256	$(3 \times 3 \times 128 + 1) \times 256 = 295168$
ReLU	56X56X256	0
CONV3-256	56X56X256	$(3 \times 3 \times 256 + 1) \times 256 = 590080$
ReLU	56X56X256	0
POOL2	28X28X256	0
CONV3-512	28X28X512	$(3 \times 3 \times 256 + 1) \times 512 = 2359808$
ReLU	28X28X512	0
CONV3-512	28X28X512	$(3 \times 3 \times 512 + 1) \times 512 = 4718592$
ReLU	28X28X512	0
POOL2	14X14X512	0
CONV3-512	14X14X512	$(3 \times 3 \times 512 + 1) \times 512 = 4718592$
ReLU	14X14X512	0
CONV3-512	14X14X512	$(3 \times 3 \times 512 + 1) \times 512 = 4718592$
ReLU	14X14X512	0
POOL2	7X7X512	0
FC-4096	4096X1	$(7 \times 7 \times 512 + 1) \times 4096 = 102764544$
FC-4096	4096X1	$(4096 + 1) \times 4096 = 16781312$
FC-1000	1000X1	$(4096 + 1) \times 1000 = 4097000$
SOFTMAX	1000X1	0

1. מספר הפרמטרים הכולל: 114,119,336

2. החלק היחסי של הפרמטרים של שכבות ה-FC:

$$\frac{123,642,856}{141,119,336} \cdot 100\% = 87.62\%$$