

67577) מבוא למערכות לומדות | תרגיל 2

שם: נמרוד בר גיורא | ת"ז: 207090622

חלק תאורטי

Let \mathbf{X} be the design matrix of a linear regression problem with m rows (samples) and d columns (variables/features). Let $\mathbf{y} \in \mathbb{R}^m$ be the response vector corresponding the samples in \mathbf{X} . Recall that for some vector space $V \subseteq \mathbb{R}^d$ the orthogonal complement of V is: $V^\perp := \{\mathbf{x} \in \mathbb{R}^d \mid \langle \mathbf{x}, \mathbf{v} \rangle = 0 \quad \forall \mathbf{v} \in V\}$

פתרונות למשוואות הנורמליות

שאלה 1

1. Prove that: $\text{Ker}(\mathbf{X}) = \text{Ker}(\mathbf{X}^\top \mathbf{X})$

הוכחה: נראה שיש הכלה דו-כיוונית בין הקבוצות וזה יוכיח את השוויון:

(\supseteq) יהי $v \in \text{ker}(\mathbf{X}^\top \mathbf{X})$ אזי $\mathbf{X}^\top \mathbf{X} v = 0$, ולכן:

$$0 = v^\top \cdot 0 = v^\top (\mathbf{X}^\top \mathbf{X} v) = (v^\top \mathbf{X}^\top) \mathbf{X} v = (\mathbf{X} v)^\top (\mathbf{X} v) = \|\mathbf{X} v\|^2$$

זו הנורמה האוקלידית

מכיוון שהנורמה היא חיובית בהחלט אז $\mathbf{X} v = 0$ בהכרח, כלומר $v \in \text{ker}(\mathbf{X})$, ובאופן כללי קיבלנו ש- $\text{ker}(\mathbf{X}) \supseteq \text{ker}(\mathbf{X}^\top \mathbf{X})$.

(\subseteq) יהי $v \in \text{ker}(\mathbf{X})$ אזי $\mathbf{X} v = 0$, ולכן:

$$\mathbf{X}^\top \mathbf{X} v = \mathbf{X}^\top \cdot 0 = 0$$

ולכן $v \in \text{ker}(\mathbf{X}^\top \mathbf{X})$ ובאופן כללי - $\text{ker}(\mathbf{X}) \subseteq \text{ker}(\mathbf{X}^\top \mathbf{X})$.

שאלה 2

2. Prove that for a square matrix A : $\text{Im}(A^\top) = \text{Ker}(A)^\perp$

נוכיח קודם 2 טענות עזר:

למה 1: לכל תמ"ו $V \subseteq \mathbb{R}^n$ מתקיים ש- $(V^\perp)^\perp = V$.

הוכחה: יהי $v \in V$, אז לכל $v^\perp \in V^\perp$ מתקיים ש- $\langle v, v^\perp \rangle = 0$, ולכן לפי הגדרת המרחב הניצב- $v \in (V^\perp)^\perp$. כלומר - $V \subseteq (V^\perp)^\perp$. בקורס ליניארית 2 מוכיחים שמרחב והמרחב הניצב לו הם סכום ישר של המרחב כולו. כלומר מתקיים ש-

$$\mathbb{R}^n = V \oplus V^\perp$$

ומצד שני מתקיים גם ש-

$$\mathbb{R}^n = V^\perp \oplus (V^\perp)^\perp$$

ולכן:

$$\dim V + \dim V^\perp = \dim V^\perp + \dim (V^\perp)^\perp \implies \dim V = \dim (V^\perp)^\perp$$

כעת יחד עם ההכלה שהראינו מתקיים ש- $V = (V^\perp)^\perp$. \square

למה 2: אם $U, V \subseteq \mathbb{R}^n$ תמ"וים כך ש- $U^\perp \subseteq V$ אז $V^\perp \subseteq U$.

הוכחה: יהי $v^\perp \in V^\perp$ אז לכל $v \in V$ מתקיים ש- $\langle v^\perp | v \rangle = 0$.

בפרט, מכיוון ש- $U^\perp \subseteq V$ אז לכל $u^\perp \in U^\perp$ מתקיים ש- $\langle v^\perp | u^\perp \rangle = 0$, ולכן

$$v^\perp \in (U^\perp)^\perp \stackrel{\text{לפי למה 1}}{=} U$$

כלומר מתקיים ש- $V^\perp \subseteq U$. \square

עכשיו נוכיח את הטענה הנדרשת:

הוכחה: נוכיח הכלה דו-כיוונית:

(\subseteq) יהי $v \in \text{Im}(A^\top)$ אז קיים $w \in \mathbb{R}^n$ כך ש- $A^\top w = v$. נשים לב שלכל $u \in \ker(A)$ מתקיים:

$$\langle v | u \rangle = v^\top u = (A^\top w)^\top u = w^\top A u = w^\top \cdot 0 = 0$$

ולכן $v \in (\ker(A))^\perp$. כלומר - $\text{Im}(A^\top) \subseteq (\ker(A))^\perp$.

(\supseteq) יהי $v \in (\ker(A))^\perp$ מכיוון שלכל $w \in \mathbb{R}^n$ מתקיים ש- $A^\top w \in \text{Im}(A^\top)$ אז לכל $w \in \mathbb{R}^n$ מתקיים:

$$0 = \langle v | A^\top w \rangle = v^\top A^\top w = (Av)^\top w$$

בפרט, זה נכון עבור $w = Av$, כלומר-

$$0 = (Av)^\top (Av) = \|Av\|^2$$

ולכן מחויבות בהחלט של הנורמה האוקלידית נובע ש- $Av = 0$, ולכן $v \in \ker(A)$.

כלומר הראינו ש- $\ker(A) \supseteq (\text{Im}(A^\top))^\perp$ ולכן לפי למה 2 - $\text{Im}(A^\top) \supseteq (\ker(A))^\perp$.

בסך הכל משתי ההכלות נובע ש- $\text{Im}(A^\top) = (\ker(A))^\perp$, כנדרש.

שאלה 3

3. Let $y = Xw$ be a non-homogeneous system of linear equations. Assume that X is square and not invertible. Show that the system has ∞ solutions $\Leftrightarrow y \perp \text{Ker}(X^\top)$.

הוכחה: (\Leftarrow) נניח ש- $y \perp \ker(X^\top)$. כלומר - $y \in (\ker(X^\top))^\perp$. לפי ההוכחה בשאלה הקודמת:

$$(\ker(X^\top))^\perp = \text{Im}((X^\top)^\top) = \text{Im}(X)$$

כלומר - $y \in \text{Im}(X)$, ולכן קיים למערכת המשוואות $Xw = y$ פתרון. נסמן אותו ב- p , כלומר - $Xp = y$.

מכיוון ש- X לא הפיכה, אז למערכת המשוואות ההומוגנית $Xw = 0$ קיימים ∞ פתרונות.

נשים לב שלכל פתרון v כזה, המקיים ש- $Xv = 0$, מתקיים ש-

$$X(p + v) = Xp + Xv = Xp + 0 = Xp = y$$

כלומר ניתן ליצור ∞ פתרונות למערכת $Xw = y$ בעזרת הפתרון p ואינסוף הפתרונות של המערכת ההומוגנית $Xw = 0$.
 (\Rightarrow) נניח שלמערכת המשוואות $Xw = y$ קיימים ∞ פתרונות. נובע מכך ש- $y \in \text{Im}(X)$, ולכן לפי ההוכחה בשאלה הקודמת:

$$y \in (\ker(X^\top))^\perp$$

כלומר $y \perp \ker(X^\top)$.

שאלה 4

4. Consider the (normal) linear system $X^\top Xw = X^\top y$. Using what you have proved above prove that the normal equations can only have a unique solution (if $X^\top X$ is invertible) or infinitely many solutions (otherwise).

הוכחה: נשים לב ש- $X^\top X$ היא מטריצה ריבועית (מסדר $d \times d$), ולכן אם היא הפיכה אז לכל $b \in \mathbb{R}^d$ קיים פתרון יחיד למערכת $X^\top Xw = b$.
 בפרט, עבור $b = X^\top y$ קיים פתרון יחיד, ואפשר לקבל אותו ע"י כפל במטריצה ההפוכה:

$$[X^\top X]^{-1} \cdot X^\top Xw = [X^\top X]^{-1} X^\top y \implies w = [X^\top X]^{-1} X^\top y$$

אם $X^\top X$ לא הפיכה אז לפי ההוכחה בשאלה הקודמת - קיימים ∞ פתרונות למערכת $X^\top Xw = X^\top y$ אם ורק אם מתקיים ש-

$$X^\top y \perp \ker(X^\top X)$$

אבל לפי ההוכחה בשאלה 1 - $\ker(X^\top X) = \ker(X)$, ומתקיים לכל $v \in \ker(X)$ ש-

$$\langle X^\top y | v \rangle = (X^\top y)^\top v = y^\top Xv \stackrel{v \in \ker(X)}{=} y^\top \cdot 0 = 0$$

כלומר $X^\top y \perp \ker(X)$ ולכן $X^\top y \perp \ker(X^\top X)$, ולכן אם $X^\top X$ לא הפיכה אז קיימים ∞ פתרונות למערכת.

מטריצות הטלה

שאלה 5

5. In this question you will prove some properties of orthogonal projection matrices seen in recitation 1. Let $V \subseteq \mathbb{R}^d$, $\dim(V) = k$ and let $\mathbf{v}_1, \dots, \mathbf{v}_k$ be an orthonormal basis of V . Define the orthogonal projection matrix $P = \sum_{i=1}^k \mathbf{v}_i \mathbf{v}_i^\top$ (notice this is an outer product)
- Show that P is symmetric.
 - Prove that the eigenvalues of P are 0 or 1 and that $\mathbf{v}_1, \dots, \mathbf{v}_k$ are the eigenvectors corresponding the eigenvalue 1.
 - Show that $\forall \mathbf{v} \in V \quad P\mathbf{v} = \mathbf{v}$.
 - Prove that $P^2 = P$.
 - Prove that $(I - P)P = 0$.

(a) **הוכחה:** נשים לב שמתכונות של פעולת השחלוף (שמוכיחים בקורס אלגברה ליניארית) נובע ש-

$$P^\top = \left(\sum_{i=1}^k v_i v_i^\top \right)^\top = \sum_{i=1}^k (v_i v_i^\top)^\top = \sum_{i=1}^k (v_i^\top)^\top v_i^\top = \sum_{i=1}^k v_i v_i^\top = P$$

ולכן P מטריצה סימטרית.

(b) **הוכחה:** כפי שמוכיחים בליניארית 2 - ניתן להשלים את הבסיס האורתונורמלי הנתון לבסיס של \mathbb{R}^d , ואז ע"י הליך גרס-שמידט לקבל בסיס $(u_1, \dots, u_k, u_{k+1}, \dots, u_d)$ אורתונורמלי של \mathbb{R}^d , כך שמשמרת התכונה ש- $V = \text{Span}(u_1, \dots, u_k)$, ומתקיים ש-

$$\mathbb{R}^d = \text{Span}(u_1, \dots, u_k) \oplus \text{Span}(u_{k+1}, \dots, u_d) = V \oplus V^\perp$$

יהי $v \in \mathbb{R}^d$, $v \neq 0$. אם $v \in V$ אז קיימים $\alpha_1, \dots, \alpha_k \in \mathbb{R}$ כך ש- $v = \sum_{j=1}^k \alpha_j v_j$, ולכן:

$$Pv = \left(\sum_{i=1}^k v_i v_i^\top \right) \cdot \sum_{j=1}^k \alpha_j v_j = \sum_{i=1}^k (v_i v_i^\top \sum_{j=1}^k \alpha_j v_j) = \sum_{i=1}^k \sum_{j=1}^k \alpha_j v_i v_i^\top v_j = \sum_{i=1}^k \sum_{j=1}^k \alpha_j v_i \langle v_i | v_j \rangle \stackrel{\text{⚡}}{=} \sum_{i=1}^k \alpha_i v_i = v = 1 \cdot v$$

$$\langle v_i | v_j \rangle = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases} \quad \text{⚡ כי } (v_1, \dots, v_k) \text{ בסיס א"נ ולכן לכל } 1 \leq i, j \leq k \text{ מתקיים ש-}$$

לכן 1 הוא ע"ע של P , וכל $v \in V$ הוא וקטור עצמי שמתאים לע"ע הזה, בפרט - v_1, \dots, v_k הם וקטורים עצמיים שמתאימים לו.

אם $v \in V^\perp$ אז $\langle v_i | v \rangle = 0$ לכל $i \in [k]$, ולכן:

$$Pv = \left(\sum_{i=1}^k v_i v_i^\top \right) \cdot v = \sum_{i=1}^k v_i v_i^\top v = \sum_{i=1}^k v_i \langle v_i | v \rangle = \sum_{i=1}^k v_i \cdot 0 = 0 = 0 \cdot v$$

לכן 0 הוא ע"ע של P , וכל $v \in V^\perp$ הוא וקטור עצמי שמתאים לע"ע הזה.

נשים לב שאלו הערכים העצמיים היחידים כי אם $v = u + w$ עם $u \in V$ ו- $w \in V^\perp$ אז לפי מה שהראינו:

$$Pv = P(u + w) = Pu + Pw = u + 0 = u$$

ולא קיים $\lambda \in \mathbb{R}$ כך ש- $Pv = u = \lambda v$, כי אז נקבל ש- $v \in V$ בסתירה לכך ש- $w \neq 0$. כלומר v אינו וקטור עצמי במקרה הזה. ■

(c) **הוכחה:** הראיתי את זה בסעיף הקודם. נשים את זה כאן שוב למקרה ש-:

יהי $v \in V$. אז קיימים $\alpha_1, \dots, \alpha_k \in \mathbb{R}$ כך ש- $v = \sum_{j=1}^k \alpha_j v_j$, ולכן:

$$Pv = \left(\sum_{i=1}^k v_i v_i^\top \right) \cdot \sum_{j=1}^k \alpha_j v_j = \sum_{i=1}^k (v_i v_i^\top \sum_{j=1}^k \alpha_j v_j) = \sum_{i=1}^k \sum_{j=1}^k \alpha_j v_i v_i^\top v_j = \sum_{i=1}^k \sum_{j=1}^k \alpha_j v_i \langle v_i | v_j \rangle \stackrel{\text{⚡}}{=} \sum_{i=1}^k \alpha_i v_i = v$$

$$\langle v_i | v_j \rangle = \begin{cases} 0 & i \neq j \\ 1 & i = j \end{cases} \quad \text{⚡ כי } (v_1, \dots, v_k) \text{ בסיס א"נ ולכן לכל } 1 \leq i, j \leq k \text{ מתקיים ש-}$$

(d) **הוכחה:** לכל $v \in V$ מתקיים ש-

$$\begin{aligned} P^2 v &= P \cdot Pv = P \cdot \sum_{j=1}^k v_j v_j^\top v = \left(\sum_{i=1}^k v_i v_i^\top \right) \cdot \sum_{j=1}^k v_j v_j^\top v = \sum_{i=1}^k (v_i v_i^\top \sum_{j=1}^k v_j v_j^\top v) = \\ &= \sum_{i=1}^k \sum_{j=1}^k v_i v_i^\top v_j v_j^\top v = \sum_{i=1}^k \sum_{j=1}^k v_i \langle v_i | v_j \rangle \underset{\substack{\uparrow \\ \text{כי הבסיס א"נ}}}{v_j^\top v} = \sum_{i=1}^k v_i v_i^\top v = Pv \end{aligned}$$

$$P^2 = P \quad \text{ולכן}$$

(e) **הוכחה:** לפי הסעיף הקודם $P^2 = P$ ולכן:

$$(I - P)P = P - P^2 = P - P = 0$$

6. Show that if $\mathbf{X}^\top \mathbf{X}$ is invertible, the general solution we derived in recitation equals to the solution you have seen in class. For this part, assume that $\mathbf{X}^\top \mathbf{X}$ is invertible.

נוכיח קודם טענת עזר:

למה: אם $O \in \mathbb{R}^{d \times r}$ מטריצה אורתוגונלית ו- $A \in \mathbb{R}^{r \times r}$ מטריצה הפיכה אז:

$$[OAO^\top]^{-1} = OA^{-1}O^\top$$

הוכחה: מכך ש- O אורתוגונלית נובע ש- $O \cdot O^\top = I$ ולכן:

$$(OA^{-1}O^\top)(OAO^\top) = OA^{-1}O^\top OAO^\top = OA^{-1}AO^\top = OO^\top = I$$

מכיוון ש- OAO^\top היא מטריצה ריבועית (מסדר $d \times d$) אז השיויון הזה מספיק כדי להסיק ש- $[OAO^\top]^{-1} = OA^{-1}O^\top$.
עכשיו נוכיח את הטענה הנדרשת:

הוכחה: פירוק ה-SVD של X הוא:

$$X = U\Sigma V^\top = \overbrace{\begin{bmatrix} | & & | \\ u_1 & \dots & u_m \\ | & & | \end{bmatrix}}^{m \times m} \overbrace{\begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_d \end{bmatrix}}^{m \times d} \overbrace{\begin{bmatrix} - & v_1^\top & - \\ & \vdots & \\ - & v_d^\top & - \end{bmatrix}}^{d \times d}$$

וראינו בתרגול שניתן להפוך את הפירוק הזה לפירוק compact SVD באופן הבא:
נסמן $r \leq d$ את מספר הערכים הסינגולרים השונים מ-0 של X , ונסמן ב- $\tilde{\Sigma} \in \mathbb{R}^{r \times r}$ את המטריצה האלכסונית שעל האלכסון שלה מופיעים r הערכים הסינגולרים הללו. נסמן ב- $\tilde{U} \in \mathbb{R}^{m \times r}$ את המטריצה המורכבת מ- r העמודות הראשונות של U , וב- $\tilde{V} \in \mathbb{R}^{d \times r}$ את המטריצה המורכבת מ- r העמודות הראשונות של V . נקבל ש-

$$X = U\Sigma V^\top = \begin{bmatrix} | & & | & & | \\ u_1 & \dots & u_r & u_{r+1} & \dots & u_m \\ | & & | & & | \end{bmatrix} \begin{bmatrix} \sigma_1 & & & 0 \\ & \ddots & & \vdots \\ & & \sigma_r & 0 \\ 0 & \dots & 0 & 0 \end{bmatrix} \begin{bmatrix} - & v_1^\top & - \\ & \vdots & \\ - & v_r^\top & - \\ & \vdots & \\ - & v_d^\top & - \end{bmatrix} = \begin{bmatrix} | & & | \\ u_1 & \dots & u_r \\ | & & | \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{bmatrix} \begin{bmatrix} - & v_1^\top & - \\ & \vdots & \\ - & v_r^\top & - \end{bmatrix} = \tilde{U}\tilde{\Sigma}\tilde{V}^\top$$

הפתרון הכללי למשוואות הנורמליות שראינו בתרגול הוא $\hat{w} = X^\dagger y$ כאשר $X^\dagger = V\Sigma^\dagger U^\top$, ו- $\Sigma^\dagger \in \mathbb{R}^{d \times m}$ אלכסונית המוגדרת ע"י:

$$\forall i \in [d]: \quad \Sigma_{i,i}^\dagger = \begin{cases} 1/\sigma_i & \sigma_i \neq 0 \\ 0 & \sigma_i = 0 \end{cases}$$

נשים לב שניתן לכתוב את X^\dagger באופן דומה ל-SVD compact:

$$X^\dagger = V\Sigma^\dagger U^\top = \begin{bmatrix} | & & | & & | \\ v_1 & \dots & v_r & v_{r+1} & \dots & v_m \\ | & & | & & | \end{bmatrix} \begin{bmatrix} \frac{1}{\sigma_1} & & & 0 \\ & \ddots & & \vdots \\ & & \frac{1}{\sigma_r} & 0 \\ 0 & \dots & 0 & 0 \end{bmatrix} \begin{bmatrix} - & u_1^\top & - \\ & \vdots & \\ - & u_r^\top & - \\ & \vdots & \\ - & u_d^\top & - \end{bmatrix} = \begin{bmatrix} | & & | \\ v_1 & \dots & v_r \\ | & & | \end{bmatrix} \begin{bmatrix} \frac{1}{\sigma_1} & & \\ & \ddots & \\ & & \frac{1}{\sigma_r} \end{bmatrix} \begin{bmatrix} - & u_1^\top & - \\ & \vdots & \\ - & u_r^\top & - \end{bmatrix} = \tilde{V}\tilde{\Sigma}^{-1}\tilde{U}^\top$$

עכשיו, מכך ש- $X^\top X$ הפיכה נובע כפי שראינו בכיתה שפתרון של המשוואות הנורמליות הוא $\bar{w} = [X^\top X]^{-1} X^\top y$. נראה ש- $\hat{w} = \bar{w}$ נציב את ה-SVD compact של X בנוסחה של \bar{w} ונקבל:

$$\bar{w} = [X^\top X]^{-1} X^\top y = \left[(\tilde{U} \tilde{\Sigma} \tilde{V}^\top)^\top \tilde{U} \tilde{\Sigma} \tilde{V}^\top \right]^{-1} (\tilde{U} \tilde{\Sigma} \tilde{V}^\top)^\top y = [\tilde{V} \tilde{\Sigma}^\top \tilde{U}^\top \tilde{U} \tilde{\Sigma} \tilde{V}^\top]^{-1} \tilde{V} \tilde{\Sigma}^\top \tilde{U}^\top y =$$

$$\tilde{U} \text{ אורתוגונלית} \rightarrow = [\tilde{V} \tilde{\Sigma}^2 \tilde{V}^\top]^{-1} \tilde{V} \tilde{\Sigma} \tilde{U}^\top y =$$

$$\tilde{V} \text{ אורתוגונלית} \rightarrow = \tilde{V} [\tilde{\Sigma}^2]^{-1} \tilde{V}^\top \tilde{V} \tilde{\Sigma} \tilde{U}^\top y =$$

$$\tilde{V} \text{ אורתוגונלית} \rightarrow = \tilde{V} [\tilde{\Sigma}^2]^{-1} \tilde{\Sigma} \tilde{U}^\top y = \star$$

נשים לב ש-

$$[\tilde{\Sigma}^2]^{-1} \tilde{\Sigma} = \begin{bmatrix} \sigma_1^2 & & \\ & \ddots & \\ & & \sigma_r^2 \end{bmatrix}^{-1} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{bmatrix} = \begin{bmatrix} \frac{1}{\sigma_1} & & \\ & \ddots & \\ & & \frac{1}{\sigma_r} \end{bmatrix} \begin{bmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{bmatrix} = \begin{bmatrix} \frac{1}{\sigma_1} & & \\ & \ddots & \\ & & \frac{1}{\sigma_r} \end{bmatrix} = \tilde{\Sigma}^{-1}$$

ולכן קיבלנו ש-

$$\star = \tilde{V} \tilde{\Sigma}^{-1} \tilde{U}^\top y = X^\dagger y = \hat{w}$$

כנדרש.

שאלה 7

7. Show that $\mathbf{X}^\top \mathbf{X}$ is invertible if and only if $\text{span}\{\mathbf{x}_1, \dots, \mathbf{x}_m\} = \mathbb{R}^d$.

הוכחה:

$$\begin{aligned} X^\top X \text{ הפיכה} &\iff \ker(X^\top X) = \{0\} \xLeftrightarrow{\text{לפי שאלה 1}} \ker(X) = \{0\} \iff \dim \ker(X) = 0 \iff \\ &\iff \dim \text{Span}\{x_1, \dots, x_m\} = \text{Rank}(X) = \dim \mathbb{R}^d \iff \text{Span}\{x_1, \dots, x_m\} = \mathbb{R}^d \end{aligned}$$

שאלה 8

8. Recall that if $\mathbf{X}^\top \mathbf{X}$ is not invertible then there are many solutions. Show that $\hat{\mathbf{w}} = \mathbf{X}^\dagger \mathbf{y}$ is the solution whose L_2 norm is minimal. That is, show that for any other solution $\bar{\mathbf{w}}$, $\|\hat{\mathbf{w}}\| \leq \|\bar{\mathbf{w}}\|$.

הוכחה: יהי $X = U \Sigma V^\top$ פירוק ה-SVD של X . נסמן ב- r את מספר הערכים הסינגולרים השונים מ-0 של X ונסמן:

$$U = \begin{bmatrix} U_1 & U_2 \end{bmatrix}, \quad V = \begin{bmatrix} V_1 & V_2 \end{bmatrix}, \quad \Sigma = \begin{bmatrix} \tilde{\Sigma} & 0 \\ 0 & 0 \end{bmatrix}$$

כאשר U_1, V_1 הן r העמודות הראשונות של U, V בהתאמה. אז כפי שראינו בסעיף 6 (ובתרגול) -

$$X = U_1 \tilde{\Sigma} V_1^\top$$

הוא פירוק ה-SVD compact של X , כאשר המימדים הם $U_1^\top \in \mathbb{R}^{m \times r}$, $\tilde{\Sigma} \in \mathbb{R}^{r \times r}$, $V_1^\top \in \mathbb{R}^{r \times d}$. נרממה מינימלית. עם $X^\top X w = X^\top y$ הוא פתרון של $\hat{w} = X^\dagger y$.

לכל פתרון $w \in \mathbb{R}^d$ נסמן $b = V^\top w$, וגם $b_1 = V_1^\top w$, $b_2 = V_2^\top w$. אז מתקיים ש-

$$b = V^\top w = \begin{bmatrix} V_1^\top \\ V_2^\top \end{bmatrix} w = \begin{bmatrix} V_1^\top w \\ V_2^\top w \end{bmatrix} = \begin{bmatrix} b_1 \\ b_2 \end{bmatrix}$$

ניזכר שפתרון של המשוואות הנורמליות הוא $\operatorname{argmin}_{w \in \mathbb{R}^n} \|Xw - y\|^2$.

ניזכר גם שמכך ש- $U^{-1}V$ אורתוגונליות נובע שהן משמרות נורמה. לכן \hat{w} הוא עם נורמה מינימלית אם ורק אם $\hat{b} = V^\top \hat{w}$ הוא עם נורמה מינימלית. בנוסף, ניעזר באורתוגונליות של U כדי לקבל לכל פתרון w ש-

$$\begin{aligned} \|Xw - y\|^2 &= \|U\Sigma V^\top w - y\|^2 = \|U(\Sigma b - U^\top y)\|^2 \stackrel{U \text{ אורת}}{=} \|\Sigma b - U^\top y\|^2 = \\ &= \left\| \begin{bmatrix} \tilde{\Sigma} & 0 \\ 0 & 0 \end{bmatrix} \begin{bmatrix} b_1 \\ b_2 \end{bmatrix} - \begin{bmatrix} U_1^\top y \\ U_2^\top y \end{bmatrix} \right\|^2 = \left\| \begin{bmatrix} \tilde{\Sigma} b_1 \\ 0 \end{bmatrix} - \begin{bmatrix} U_1^\top y \\ U_2^\top y \end{bmatrix} \right\|^2 = \left\| \begin{bmatrix} \tilde{\Sigma} b_1 - U_1^\top y \\ -U_2^\top y \end{bmatrix} \right\|^2 \stackrel{\star}{=} \\ &\stackrel{\star}{=} \|\tilde{\Sigma} b_1 - U_1^\top y\|^2 + \|U_2^\top y\|^2 \end{aligned}$$

☀: כי הנורמה בריבוע זה פשוט סכום הריבועים של הקואורדינטות של הוקטור $\begin{bmatrix} \tilde{\Sigma} b_1 - U_1^\top y \\ -U_2^\top y \end{bmatrix}$, ואפשר לחלק את הסכום אז לשני סכומים שנותנים את הנורמות של הוקטורים $\tilde{\Sigma} b_1 - U_1^\top y$ ו- $-U_2^\top y$ בנפרד.

מכיוון ש- $\|U_2^\top y\|^2$ קבוע אי-שליילי, אז $\bar{w} \in \mathbb{R}^d$ מקיים ש- $\bar{w} = \operatorname{argmin}_{w \in \mathbb{R}^n} \|Xw - y\|^2$ אם ורק אם $\bar{b}_1 = V_1^\top \bar{w}$ מקיים ש-

$$\bar{b}_1 = \operatorname{argmin}_{b_1 \in \mathbb{R}^n} \|\tilde{\Sigma} b_1 - U_1^\top y\|^2$$

כלומר אם ורק אם \bar{b}_1 הוא פתרון של המשוואות הנורמליות הבאות:

$$\begin{aligned} \tilde{\Sigma}^\top \tilde{\Sigma} \bar{b}_1 &= \tilde{\Sigma}^\top U_1^\top y \\ \Downarrow \text{כי } \tilde{\Sigma} \text{ אלכסונית} \\ \tilde{\Sigma}^2 \bar{b}_1 &= \tilde{\Sigma} U_1^\top y \\ \Downarrow \text{כפל ב-} \tilde{\Sigma}^{-2} \\ V_1^\top \bar{w} &= \bar{b}_1 = \tilde{\Sigma}^{-1} U_1^\top y \end{aligned}$$

כלומר $\bar{w} \in \mathbb{R}^d$ הוא פתרון של המשוואות הנורמליות $X^\top Xw = X^\top y$ אם ורק אם הוא מקיים ש- $\bar{b}_1 = V_1^\top \bar{w} = \tilde{\Sigma}^{-1} U_1^\top y$, וזה בכלל לא תלוי ברכיב השני - $\bar{b}_2 = V_2^\top \bar{w}$!

ניזכר עכשיו שפתרון \bar{w} הוא עם נורמה מינימלית אם ורק אם $V^\top \bar{w}$ הוא עם נורמה מינימלית, ולכן \bar{w} מקיים בהכרח ש- $V_2^\top \bar{w} = 0$, כי אחרת אם $V_2^\top \bar{w} \neq 0$ נוכל פשוט לקחת את הוקטור $\begin{bmatrix} V_1^\top \bar{w} \\ 0 \end{bmatrix}$ ולקבל לפי מה שראינו לעיל שהוא עדיין פתרון והוא גם מקיים-

$$\left\| \begin{bmatrix} V_1^\top \bar{w} \\ 0 \end{bmatrix} \right\| < \left\| \begin{bmatrix} V_1^\top \bar{w} \\ V_2^\top \bar{w} \end{bmatrix} \right\|$$

לסיום נראה שהפתרון $\hat{w} = X^\dagger y$ מקיים בדיוק את התכונה הזו:

לצורך כך נשתמש בצורת ה-SVD compact של X^\dagger שהראיתי בסעיף 6 (שם סימנתי $\tilde{U}, V_1 = \tilde{V}$). נקבל ש-

$$\begin{aligned} V_1^\top \hat{w} &= V_1^\top X^\dagger y = V_1^\top V_1 \tilde{\Sigma}^{-1} U_1^\top y \stackrel{(1)}{=} \tilde{\Sigma}^{-1} U_1^\top y \\ V_2^\top \hat{w} &= V_2^\top X^\dagger y = V_2^\top V_1 \tilde{\Sigma}^{-1} U_1^\top y \stackrel{(2)}{=} 0 \cdot \tilde{\Sigma}^{-1} U_1^\top y = 0 \end{aligned}$$

(1), (2): שני השיעיונות נובעים מכך ש- V אורתוגונלית - כי זה אומר שעמודותיה הן בסיס אורתונורמלי של \mathbb{R}^d . לכן במכפלה $V_1^\top V_1$ נקבל בדיוק את I_r , ובמכפלה $V_2^\top V_1$ נקבל רק מכפלות מהצורה $v_i^\top v_j$ עם $i \in [r]$ ו- $j \in [d] \setminus [r]$, השוות ל-0 מהאורתונורמליות (כש- v_i היא העמודה ה- i של V).

לכן \hat{w} הוא הפתרון עם הנורמה המינימלית למשוואות הנורמליות, כנדרש. ■

חלק מעשי

רגרסיה ליניארית

שאלה 1

- איזה פיצ'רים שמרתי ואיזה לא?
תשובה: מחקתי את הפיצ'רים id,date,lat,long (אחרי שעשיתי שימוש בנתונים שיש בעמודת date), שמרתי את שאר הפיצ'רים.
- איזה פיצ'רים מורכבים מערכים קטגוריים ואיך התמודדתי איתם?
תשובה: הפיצ'ר היחיד שהתייחסתי אליו כבעל ערכים קטגוריים הוא ה-zipcode. השתמשתי בשיטת one-hot encoding כדי להפוך אותו לערך של 0 או 1.
- איזה פיצ'רים נוספים הוספתי ומה היה ההיגיון מאחורי זה?
תשובה: הוספתי את הפיצ'רים הבאים:
– age: הוספתי לכל רשומה בדאטא ערך שמציין את ה"גיל" של הבית, ומוגדר ע"י:

$$\text{sale_year} - \max(\text{yr_built}, \text{yr_renovated})$$

עשיתי זאת כי לדעתי זו דרך טובה לתת ערך מספרי (עם יחס סדר) לתאריכים שיש לכל רשומה בדאטא - פחות ברור מה הקשר בין בית שנמכר ב-2018 ובית שנמכר ב-2015 מאשר הקשר בין בית שנמכר 20 שנים לאחר ששופץ לאחרונה ובין בית שנמכר 2 שנים לאחר ששופץ לאחרונה.

- sqft_living_ratio: פיצ'ר שמודד את הגודל היחסי של הבית ביחס לגודל הממוצע של 15 הבתים הקרובים אליו ביותר. מוגדר ע"י:

$$\frac{\text{sqft_living}_{15}}{\text{sqft_living}}$$

- sqft_lot_ratio: פיצ'ר שמודד את הגודל היחסי של הנכס ביחס לגודל הממוצע של 15 הנכסים הקרובים ביותר. מוגדר ע"י:

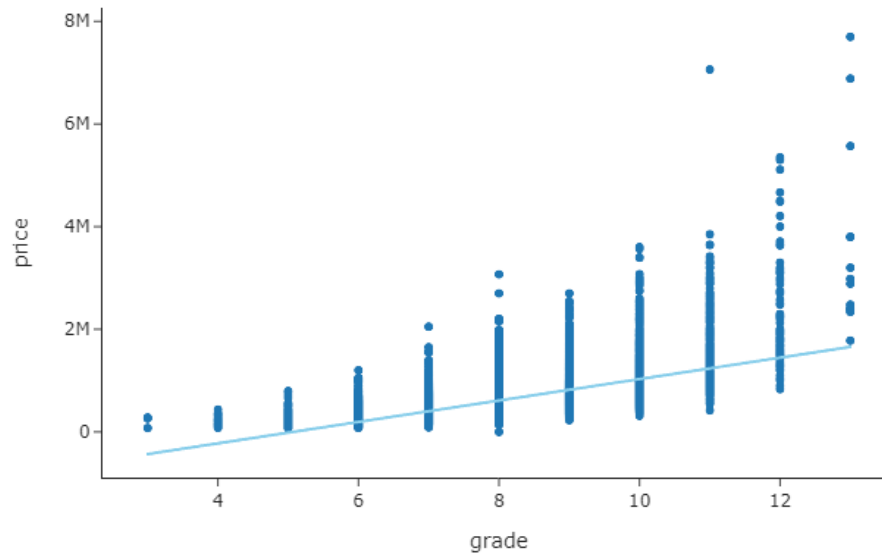
$$\frac{\text{sqft_lot}_{15}}{\text{sqft_lot}}$$

- איך התמודדתי עם ערכים חסרים או לא הגיוניים?
תשובה: מכיוון שלא היו הרבה ערכים חסרים או לא הגיוניים - מחקתי את כל הרשומות (השורות) שבהן היו ערכים כאלו. בסך הכל לא מחקתי יותר מ-20 רשומות.
- דברים נוספים שעשיתי:
בנוסף, מחקתי דגימות עם מחירי בתים גבוהים במיוחד (מעל 5 מיליון), כי ראיתי שבאופן יחסי אלו מקרי קצה והם מנפחים מאוד את ערך ה-MSE כפי שראינו בתרגול.

שאלה 2

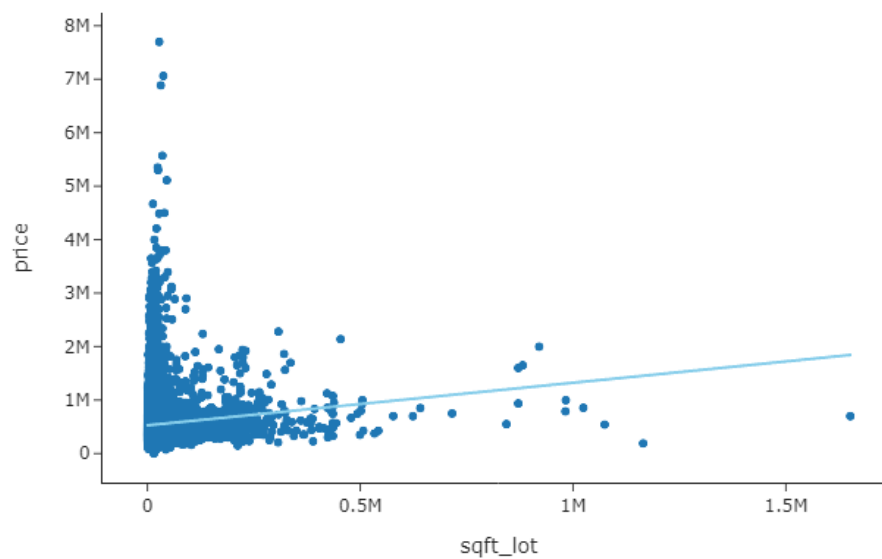
ערך שנראה שמועיל למודל הוא ערך ה-grade. הנה גרף שמציג אותו ביחס למחיר של כל בית בדאטא, ואת הקורלציה ביניהם:

Pearson Correlation of grade and price = 0.6676539796495431



ערך שנראה שפחות מועיל למודל הוא sqft_lot:

Pearson Correlation of sqft_lot and price = 0.08979401806069234



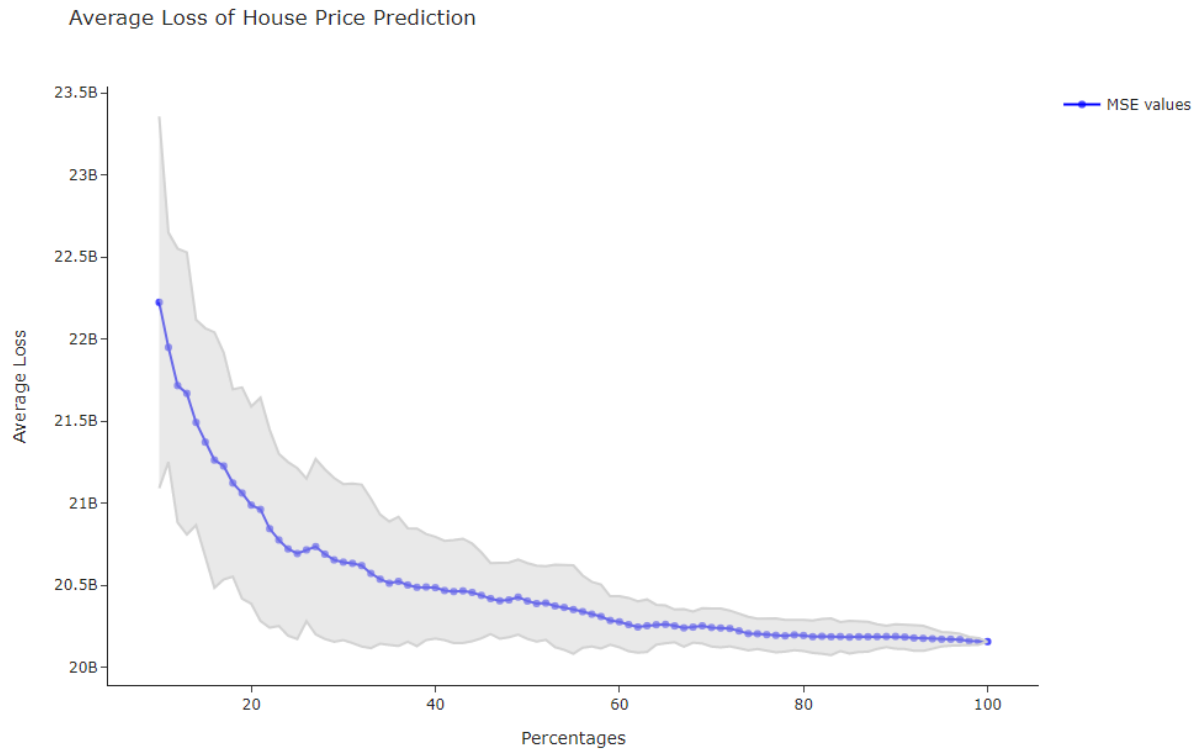
הסיבה שהסקתי שהראשון מועיל והשני פחות מועיל היא שהקורלציה בין ה-grade למחיר של הבית גבוהה בהרבה מהקורלציה בין ה-sqft_lot למחיר (שקרובה מאוד ל-0).

שאלה 3

בקוד.

שאלה 4

הגרף הנדרש:



הסבר:

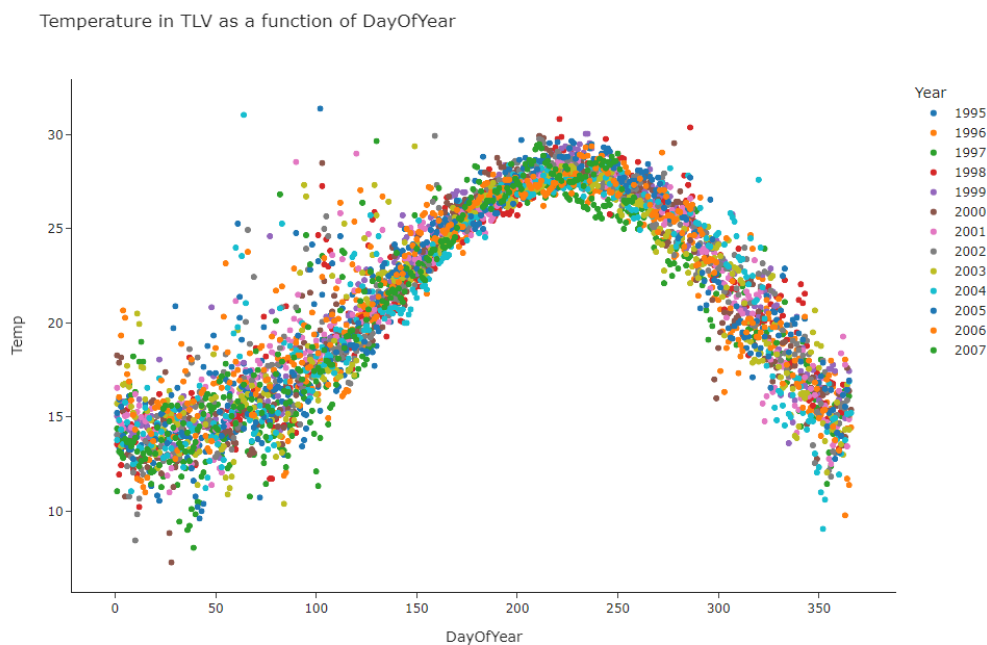
ככל שאחוז הדגימות מתוך סט האימון **גדול יותר**, ערך ה-Loss של המודל **קטן** (כאשר הוא מחושב כממוצע של האומדן עבור 10 דגימות אקראיות מסט ה-test), וה-Confidence interval הולך ו"מתהדק" - כלומר השונות של האומדן הולכת ו**קטנה** ככל שהמודל מתאמן על אחוז **גדול יותר** מתוך סט האימון. הירידה הזו בשגיאה מעידה על כך שהאומד שלנו משתפר ככל שמספר הדגימות שהוא מתאמן עליהן עולה, והוא הולך ונעשה פחות מוטה (biased) ועם שונות נמוכה יותר (variance).

הערה נוספת שראוי להוסיף על ערך ה-confidence interval, כלומר על השונות היא שכשמתקרבים ל-100% מתוך סט האימון ומריצים את המודל 10 פעמים, ההבדל בתוצאות נהיה קטן מאוד כי החפיפה בין השורות שנבחרות בכל הרצה הולכת וגדלה. זה יכול להסביר גם את הירידה בשונות.

Polynomial Fitting

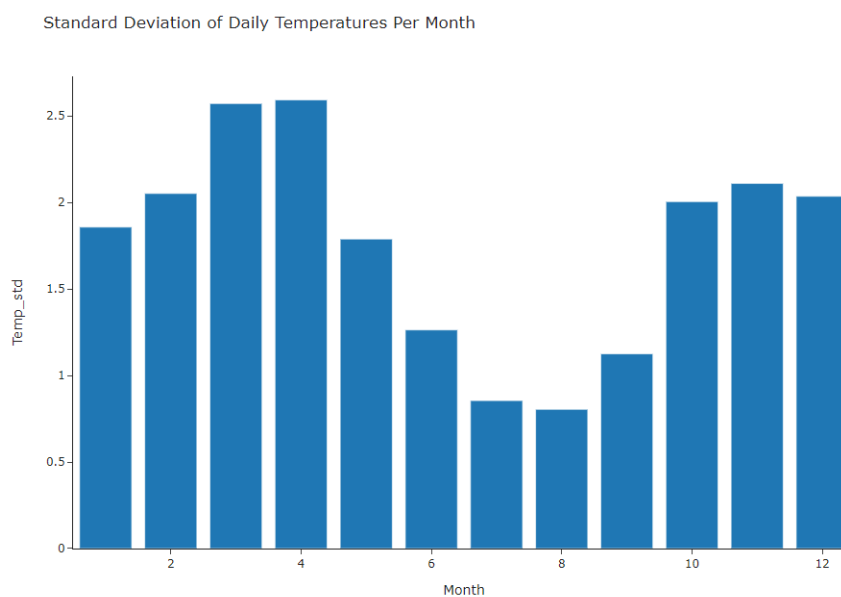
שאלה 2

- גרף שמציג את הטמפרטורה בישראל (בפועל בתל-אביב) בפונקציה של היום בשנה:



לדעתי פולינום מדרגה 3 או 4 יספיק כדי להתאים לנתונים האלו.

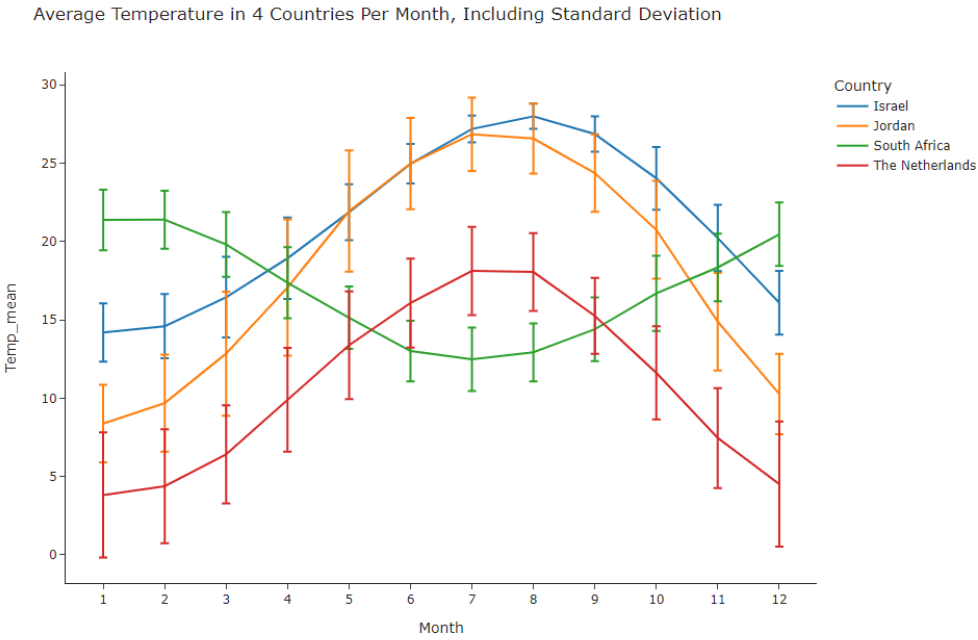
- גרף שמציג את סטיית התקן של הטמפרטורות היומיות בכל חודש:



ע"פ הגרף הזה אני מצפה שמודל פולינומיאלי יהיה מדויק יותר בחודשים 7, 8 שבהם סטיית התקן של הטמפרטורות נמוכה יותר.

שאלה 3

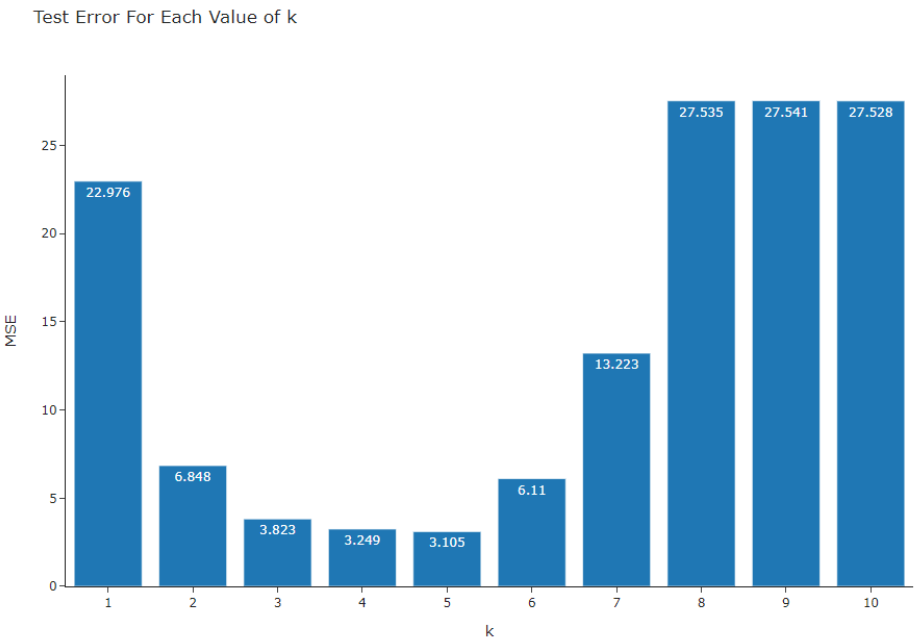
גרף המציג את הטמפרטורה הממוצעת בכל חודש בישראל, ירדן, דרום אמריקה והולנד, כולל סטיית התקן:



רואים בגרף לעיל שתבנית הטמפרטורה לפי חודשים של ישראל וירדן דיי דומה, והערכים דיי קרובים. אני חושב שמודל שלמד את הדאטא של ישראל בלבד יעבוד הכי טוב (באופן יחסי) על הדאטא של ירדן.

שאלה 4

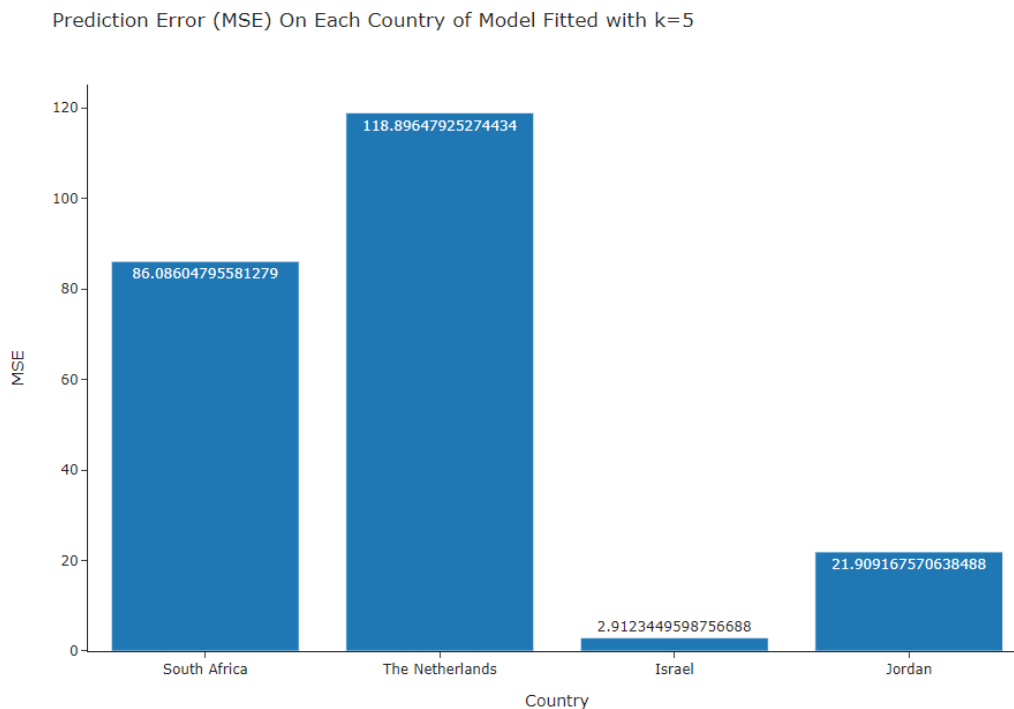
גרף המציג את הטעות (MSE) של המודל עבור כל ערך של k בין 1 ל-10:



לפי התוצאות האלו רואים שהערך של k עבורו המודל התאים ביותר לדאטא הוא 5. התוצאות שהתקבלו עבור $k = 4$ ו- $k = 3$ דומות מאוד, ולכן יתכן שעדיף לבחור באחד מהערכים הללו עבור k כי זה יתן לנו מודל פשוט יותר.

שאלה 5

גרף המציג את השגיאה (MSE) של מודל שלמד על הדאטא של ישראל בלבד, עם $k = 5$, על דאטא של כל אחת מ-4 המדינות שבדאטאסט כולו:



כצפוי, השגיאה הקטנה ביותר היא על הדאטא של ישראל, כי עליו המודל התאמן. כמו שצפיתי בתשובה לשאלה 3 - השגיאה שהתקבלה עבור הדאטא של ירדן היא הנמוכה מבין 3 המדינות הנוספות.