

מבוא ללמידת מכונה - תרגיל 2

מגיש: יואב לוי 314963257

Solutions of the Normal Equations

שאלה 1

יש להוכיח $\text{Ker}(X^T) = \text{Ker}(XX^T)$, אראה זאת בעזרת הכלה דו-כיוונית.

$$\bullet \text{Ker}(X^T) \subseteq \text{Ker}(XX^T)$$

יהי וקטור $w \in \text{Ker}(X^T)$ מתקיים,

$$XX^T w = X\vec{0} = \vec{0} \implies w \in \text{Ker}(XX^T)$$

$$\bullet \text{Ker}(X^T) \supseteq \text{Ker}(XX^T)$$

יהי וקטור $w \in \text{Ker}(XX^T)$ (נסמן ב- \odot), נשים לב כי בעזרת פירוק ה-SVD נוכל לרשום את המטריצות הבאות כך:

קיימות מטריצות U, V אורתוגונליות ומטריצה Σ ריבועית אלכסונית

$$X = U\Sigma V^T$$

$$X^T = V\Sigma^T U^T \stackrel{\Sigma \text{ is squared mat}}{=} V\Sigma U^T$$

$$XX^T = U\Sigma^2 U^T$$

כעת, נשים לב שמתקיים

$$XX^T w = U\Sigma^2 U^T w \stackrel{\odot}{=} \vec{0}$$

למה 0.1 - בהינתן מטריצה אורתוגונלית $U \in M_{n \times n}(\mathbb{R})$ הפיכה (תכונה של מטריצות אורתוגונליות)

$$\text{אז } \dim(\text{Ker}(U^T)) = 0$$

הוכחת למה 0.1:

ראשית מהגדרת ההפיכות נובע כי $\dim(\text{Ker}(U)) = 0$, מלינארית אנו יודעים שמימד מרחב השורות שווה למימד מרחב העמודות של מטריצה.

בנוסף אנחנו יודעים כי מימד מרחב השורות של U הוא $\dim(\text{Im}(U))$ והוא כאמור שווה למימד מרחב השורות של U^T . נשים לב שלפי משפט המימדים נקבל:

$$\dim(\text{Im}(U)) + \dim(\text{Ker}(U)) = n$$

ולכן גם:

$$\dim(\text{Im}(U^T)) + \dim(\text{Ker}(U^T)) = n \iff \dim(\text{Ker}(U^T)) = n - \dim(\text{Im}(U)) = \dim(\text{Ker}(U)) = 0$$

■

כעת נחזור למשוואה $XX^Tw = U\Sigma^2U^Tw = \vec{0}$ ונשים לב שמהלמה 0.1 נקבל ש- $U^Tw \neq 0$ ואם $\Sigma^2U^Tw \neq 0$ אז שוב מלמה 0.1 $U\Sigma^2U^Tw \neq \vec{0}$ בסתירה למשוואה שהגענו אליה, ולכן $\Sigma^2U^Tw = \vec{0}$ בהכרח. מכאן שגם $\Sigma U^Tw = \vec{0}$, זאת מכיוון שהמטריצה Σ ריבועית ואלכסונית ולכן להכפיל ווקטור בה זה כמו להכפיל כל קורדינטה של הווקטור. פי סקלר (הסקלרים שעל האלכסון), אם כך שלכפול פי עוד סיגמה זה כמו לכפול פי הסקלר בריבוע פשוט, לכן אם הסקלר אינו אפס מלכתחילה הוא לא יאפס את הווקטור (אלא אם הווקטור הוא אפס מלכתחילה), ומכאן נוכל לקבל בקלות שמתקיים:

$$X^Tw \stackrel{S.V.D}{=} V\Sigma U^Tw = \vec{0}$$

כלומר $w \in \text{Ker}(X^T)$ כנדרש.

■

שאלה 2

נוכיח את הטענה בעזרת הכלה דו-כיוונית.

$$\bullet \text{ } \text{Im}(A^T) \subseteq \text{Ker}(A)^\perp$$

יהי $w \in \text{Im}(A^T)$ כלומר קיים ווקטור $v \in \mathbb{R}^n$ כך שמתקיים $A^Tv = w$. אזי לכל ווקטור $\mu \in \text{Ker}(A)$ מתקיים:

$$\langle w, \mu \rangle = \langle A^Tv, \mu \rangle \stackrel{\text{Hermitian adjoint over } \mathbb{R}}{=} \langle v, A\mu \rangle \stackrel{\mu \in \text{Ker}(A)}{=} \langle v, 0 \rangle = 0$$

ולכן $w \in \text{Ker}(A)$

$$\bullet \text{ } \text{Im}(A^T) \supseteq \text{Ker}(A)^\perp$$

נניח בשלילה שעבור ווקטור $w \notin \text{Im}(A^T)$ ועלינו להראות כי $w \notin \text{Ker}(A)^\perp$. כלומר עלינו להראות כי קיים $v \in \text{Ker}(A)$ כך שמתקיים $\langle w, v \rangle \neq 0$. נגדיר את $v \in \text{Im}(A^T)^\perp$ ונראה כי מכיוון ש- $w \notin \text{Im}(A^T)$ אזי $w \in \text{Im}(A^T)^\perp$ כלומר $\langle w, v \rangle \neq 0$ וגם

$$\left\langle v, \overbrace{A^T Av}^{\in \text{Im}(A^T)} \right\rangle = 0$$

ולכן,

$$\|Av\|^2 = \langle Av, Av \rangle = \langle v, A^T Av \rangle = 0$$

כלומר $Av = 0$ ז"א $v \in \text{Ker}(A)$ ולכן $\text{Im}(A^T) \supseteq \text{Ker}(A)^\perp$ כנדרש.

■

שאלה 3

בהינתן שהמטריצה X^T סינגולרית צריך להוכיח $X^T w = y$ has ∞ solutions $\iff y \perp \text{Ker}(X)$ (⇐) אנחנו יודעים שיהיו למערכת המשוואות $X^T w = y$ או 0 או ∞ פתרונות. כעת, אם נראה ש- $y \in \text{Im}(X^T)$ אז זה אומר שקיים פתרון, ולכן יש ∞ פתרונות. נשים לב כי אנחנו מניחים שמתקיים $y \perp \text{Ker}(X)$, כלומר $y \in \text{Ker}(A)^\perp$. כעת נעזר בשאלה הקודמת ובה הראינו את השקילות $y \in \text{Ker}(A)^\perp = \text{Im}(X^T)$, ולכן y אכן בתמונה ולמערכת יש ∞ פתרונות. (⇒) כעת נניח שלמערכת יש ∞ פתרונות. נסמן פתרון אחד ב- w ונשים לב כי $y = X^T w$, כעת לכל וקטור $v \in \text{Ker}(X)$

$$\langle y, v \rangle = \langle X^T w, v \rangle = \langle w, Xv \rangle \stackrel{v \in Ker(X)}{=} \langle w, 0 \rangle = 0 \implies y \perp Ker(X)$$

כנדרש.

שאלה 4

- אם XX^T היא הפיכה (Invertible) אזי מתקיים,

$$XX^T w = Xy \iff (XX^T)^{-1} XX^T w = (XX^T)^{-1} Xy \iff w = (X^T)^{-1} y$$

- אחרת אם XX^T סינגולרית (אינה הפיכה) אזי או שיש 0 פתרונות או שיש $-\infty$, אם אצליח להראות כי $XY \in Im(XX^T)$

הרי שיש ∞ פתרונות למערכת המשוואות.

נשים לב שניתן להמיר את הבעיה לכך שצריך להראות כי $Xy \in Ker(XX^T)^\perp$ מתוך $xy \in Ker(XX^T)$ (הוכחנו בשאלה השנייה).

כעת, מהלמה שהוכחנו בשאלה הראשונה ניתן להמיר שוב את הבעיה לכך שצריך להראות כי $Xy \in Ker(X^T)^\perp$ ואת זה קל לנו כבר לפתור, יהי $v \in Ker(X^T)$ נשים לב כי מתקיים

$$\langle Xy, v \rangle = \langle y, X^T v \rangle = \langle y, 0 \rangle = 0$$

לכן $Xy \in Im(XX^T)$ כלומר $Xy \in Ker(X^T)^\perp = Ker(XX^T)^\perp = Im(XX^T)$ כנדרש.

ולכן באם המטריצה XX^T סינגולרית אז למערכת המשוואות יש ∞ פתרונות.

Projection Matrices

שאלה 5

(a) מטריצת ההטלה על $V \subseteq \mathbb{R}^d$ מוגדרת ע"י המכפלה החיצונית הבאה, (כאשר $P = \sum_{i=1}^k v_i \cdot v_i^T$)
 בסיס אורתונורמלי.

מכאן שכל איבר בסכום הוא מטריצה בגודל $d \times d$, אם אראה שכל מטריצה כזו היא סימטרית הרי ש- P עצמה סימטרית, כתוצאה של חיבור של מטריצות סימטריות.

לכן, נתבונן במחובר כללי בסכום (למשל כאשר $i = j$)

$$\vec{v}_j := \begin{bmatrix} v_{j,1} \\ v_{j,2} \\ \vdots \\ v_{j,d} \end{bmatrix} \in \mathbb{R}^d$$

$$\vec{v}_j \cdot \vec{v}_j^T = \begin{bmatrix} v_{j,1} \\ v_{j,2} \\ \vdots \\ v_{j,d} \end{bmatrix} \cdot \begin{bmatrix} v_{j,1} & v_{j,2} & \cdots & v_{j,d} \end{bmatrix} = \begin{pmatrix} (v_{j,1})^2 & (v_{j,1} \cdot v_{j,2}) & \cdots & (v_{j,1} \cdot v_{j,d}) \\ (v_{j,2} \cdot v_{j,1}) & (v_{j,2})^2 & \ddots & \vdots \\ \vdots & \ddots & \ddots & \vdots \\ (v_{j,d} \cdot v_{j,1}) & (v_{j,d} \cdot v_{j,2}) & \cdots & (v_{j,d})^2 \end{pmatrix}$$

נסמן את המטריצה שלעיל ב- A וכך נוכל לכתוב בכתוב קומפקטי:

$$A_{x,y} = (v_{j,x} \cdot v_{j,y})$$

לכן מכיוון שהכפל קומוטטיבי אז המטריצה הזאת היא סימטרית, וזה גורר שהמטריצה המקורית P סימטרית גם כן.

(b) ראשית אוכיח שהערכים העצמיים היחידים של מטריצת ההטלה הם 0,1. נשים לב שמתקיים $P^2 = P$ כי

$$P^2 = \left(\sum_{i=1}^k v_i \cdot v_i^T \right)^2 = \left((v_1 \cdot v_1^T) + \cdots + (v_k \cdot v_k^T) \right)^2 = \sum_{i=j} (v_i \cdot v_i^T) \cdot (v_j \cdot v_j^T) + 2 \sum_{i \neq j} (v_i \cdot v_i^T) (v_j \cdot v_j^T) =$$

$$\stackrel{(v_1, \dots, v_k \text{ orthonormal basis})}{=} \sum_{i=1}^k v_i \cdot v_i^T = P$$

כעת, עבור ווקטור עצמי $v \in V$ $v \neq 0$ נקבל,

$$\exists \lambda \in \mathbb{R} \text{ s.t. } Pv = \lambda v$$

$$P^2 v = Pv = \lambda v$$

$$P^2 v = P(\lambda v) = \lambda Pv = \lambda^2 v$$

\Downarrow

$$\lambda^2 v = \lambda v \iff \lambda^2 v - \lambda v = 0 \iff \lambda v (\lambda - 1) = 0$$

כלומר $\lambda = 0 \vee \lambda = 1$, כנדרש.
 כעת, נשים לב שעבורים הווקטורים v_1, \dots, v_k הערכים העצמיים המשוויכים להם הם 1:

$$\forall j \in [k]: P(v_j) = \sum_{i=1}^k v_i \cdot v_i^T \cdot (v_j) \stackrel{(v_1, \dots, v_k \text{ orthonormal basis})}{=} \sum_{i=1}^k v_i \cdot \overbrace{\delta_{i,j}}^{\text{Kronecker delta}} = 1 \cdot v_j$$

כנדרש.

(c) יהי $v \in V$, כלומר ניתן לייצגו כסופרפוזיציה של הבסיס האורתונורמלי של V הנתון לנו $B = (v_1, \dots, v_k)$ בהינתן סקלרים $c_1, \dots, c_k \in \mathbb{R}$ נקבל:

$$v = c_1 \cdot v_1 + c_2 \cdot v_2 + \dots + c_k \cdot v_k$$

כעת,

$$Pv = P(c_1 \cdot v_1 + c_2 \cdot v_2 + \dots + c_k \cdot v_k) = \left(\sum_{i=1}^k v_i \cdot v_i^T \right) \cdot \left(\sum_{i=1}^k c_i \cdot v_i \right) =$$

$$= \sum_{j=1}^k \sum_{i=1}^k v_j \cdot \overbrace{v_j^T \cdot v_i}^{\delta_{i,j}} \cdot c_i = \sum_{i=1}^k c_i \cdot v_i = v$$

כנדרש.

(d) הוכח כבר בסעיף ב' בשאלה זו (🧐)

$$(I - P)P = P - P^2 \stackrel{(d)}{=} P - P = \mathbf{0} \quad (e)$$

Least Squares

שאלה 6

•

$$(XX^T)^{-1} = \left(U \underbrace{\Sigma \underbrace{V^T V}_{=0 \text{ because orthonormal matrix}} \Sigma^T}_{=0 \text{ because orthonormal matrix}} U^T \right)^{-1} = (UDU^T)^{-1} = UD^{-1}U^T$$

כאשר השייוויון האחרון נובע מהעובדה שההופכי של מטריצות אורתוגונליות הוא המטריצה המשוחלפת.

•

$$(XX^T)^{-1}X = UD^{-1}U^T X = UD^{-1}U^T U \Sigma V^T = UD^{-1} \Sigma V^T = \left(\overbrace{V \Sigma^T \Sigma^T (\Sigma \Sigma^T)^{-1} U^T}^{(\ominus)} \right)^T$$

$$([\Sigma]_{i,j})^\dagger := \begin{cases} \frac{1}{[\Sigma]_{i,j}} & \text{if } [\Sigma]_{i,j} \neq 0 \\ 0 & \text{otherwise} \end{cases} \quad \text{כעת נעבור לנוטציה אינדקסים, ואגדיר את האופרטור } \dagger \text{ על סקלרים באופן הבא:}$$

ונשים לב שמתקיים: (\ominus)

$$\Sigma^T \Sigma^T (\Sigma \Sigma^T)^{-1} \xrightarrow{\text{Index Notation}} [\Sigma]_{j,i} \cdot \left([\Sigma]_{i,j} \cdot [\Sigma]_{j,i} \right)^{-1} = \left([\Sigma]_{i,j} \right)^\dagger$$

ולכן,

$$(X X^T)^{-1} X = X^T{}^\dagger$$

כנדרש.

שאלה 7

ראשית נשים לב שהדרגה של המטריצה $X^T \in M_{m \times d}(\mathbb{R})$ היא המימד הנפרש ע"י הווקטורים x_1, \dots, x_m . כעת, מהשאלה הראשונה בתרגול נובע כי $\text{Ker}(X^T) = \text{Ker}(X X^T)$ ומכיוון ש- $X X^T$ הפיכה אז $\text{Ker}(X^T) = \text{Ker}(X X^T) = 0$ ולפי משפט המימדים נובע,

$$\text{Im}(X^T) + \text{Ker}(X^T) = d \iff \text{Im}(X^T) = d$$

כלומר, מימד מרחב השורות של X^T שווה ל- d כלומר $\text{span}\{x_1, \dots, x_m\} = \mathbb{R}^d$, וזה קורה אם ורק אם המטריצה $X X^T$ הפיכה.

שאלה 8

מכיוון שכל קורדינטה ב- \hat{w} היא ייחודית לכל ערך סינגולרי של X אז לכל פתרון \bar{w} אחר נקבל, $\bar{w}_i = \hat{w}_i$ עבור כל $i = 1, \dots, r$, אבל $\hat{w}_i = 0$ $\forall i \in r+1, \dots, d$ ולכן:

$$\|\bar{w}\|^2 = \sum_{i=1}^d \bar{w}_i^2 = \sum_{i=1}^r \bar{w}_i^2 + \sum_{i=r+1}^d \bar{w}_i^2 \leq \sum_{i=1}^r \hat{w}_i^2 = \sum_{i=1}^d \hat{w}_i^2 = \|\hat{w}\|^2$$

שאלה 9+10+11+12

אלו שאלות קוד שמומשו בקובץ "linear_model.py".

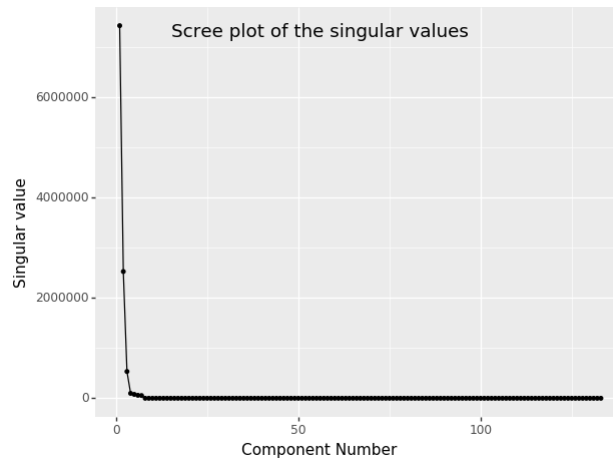
שאלה 13

מצאתי שה-features הבאים הם `categoricals`:

- `zipcode` - המיקוד, כמו שנוסח בשאלה, אין קורלציה ישירה בין הגדלים של זיפ קוד אחד לאחר אך סביר להניח שהמיקוד יכול להועיל לרגרסיה הלינארית, מאחר וסביר שתהיה קורלציה בין המחיר של הבית למיקום שלו (שבא לידי ביטוי במיקוד).
- `date` - זהו התאריך של מכירת הבית, נשים לב כי קיבלנו אותו בפורמט קצת מוזר `yyymmddT000000`, התמודדתי איתו על ידי עריכה של העמודה ופיצולה ל-3 עמודות שונות (שנה, חודש, יום) ועבור כל אחד מהעמודות האלה, יצרתי `dummy variables`.

שאלה 14

הגרף שמתאר את הערכים הסינגולרים:

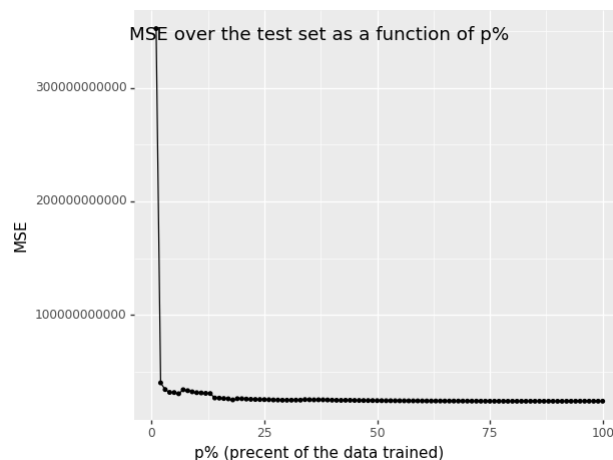


שאלה 15

מתוך ההרצאה "LectureHandout_2_Linear_Regression.pdf": (תרגומי לעברית)
 "לעיתים XX^T היא פורמלית הפיכה אך קרובה ללא הפיכה (=סינגולרית). הדבר קורה כאשר העמודות של X הן כמעט תלויות לינארית, או כאשר אחת מהעמודות של X^T הוא כמעט נפרש על-ידי שאר העמודות. במקרה כזה, מספר הערכים הסינגולרים של X שהם לא אפסים יהיה קטן".

כעת, בגרף שמופיע בשאלה 14 יש את התוצאה שלי, נשים לב כי רוב הערכים הסינגולרים הם דווקא אפס, כלומר מספר הערכים הסינגולרים של X שהם לא אפס קטן, ולכן סביר להניח שהיא **אכן** קרובה להיות סינגולרית.

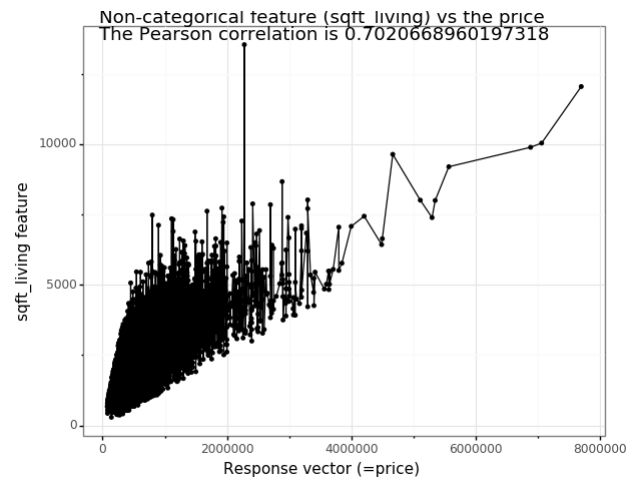
שאלה 16



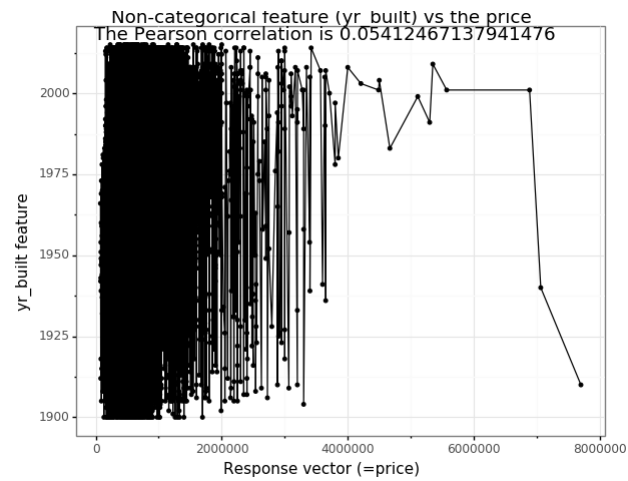
נשים לב כי אנחנו מתכנסים מהר מאוד ל- $MSE = 24,000,000,000$ זאת מכיוון שאנחנו מאמנים על עוד נתונים. בנוסף שמתי לב שכשאני לא מעבד את הדאטא שלי ולא מחליף משתנים קטגוריים אלא מוחק אותם לוקח לי יותר זמן להתכנס, פה אני מתכנס די מהר.

שאלה 17

בגרף הבא, ניתן לראות קורלציה גבוהה (יחסית) בין מחיר הדירה לשטחה ב-square feet. ניתן לראות זאת בכך שה-Pearson correlation קרוב ל-1, בנוסף בגרף ניתן לראות התנהגות יחסית לינארית ככל שהמחיר עולה. לכן ניתן להסיק שהשטח של הדירה יכול להועיל לגרסיה הלינארית. (חוץ מזה אינטואיטיבית לא מפתיע שהמחיר של הדירה קשור בקשר הדוק לגודל הדירה)



בגרף הבא, ניתן לראות קורלציה נמוכה (יחסית) בין מחיר הדירה לשנה שבה היא נבנתה. ניתן לראות זאת בכך שה-Pearson correlation קרוב ל-0, בנוסף בגרף ניתן לראות התנהגות שאינה ליניארית לכל אורך הגרף. לכן ניתן להסיק שהשטח של הדירה יכול להועיל לרגרסיה הליניארית.

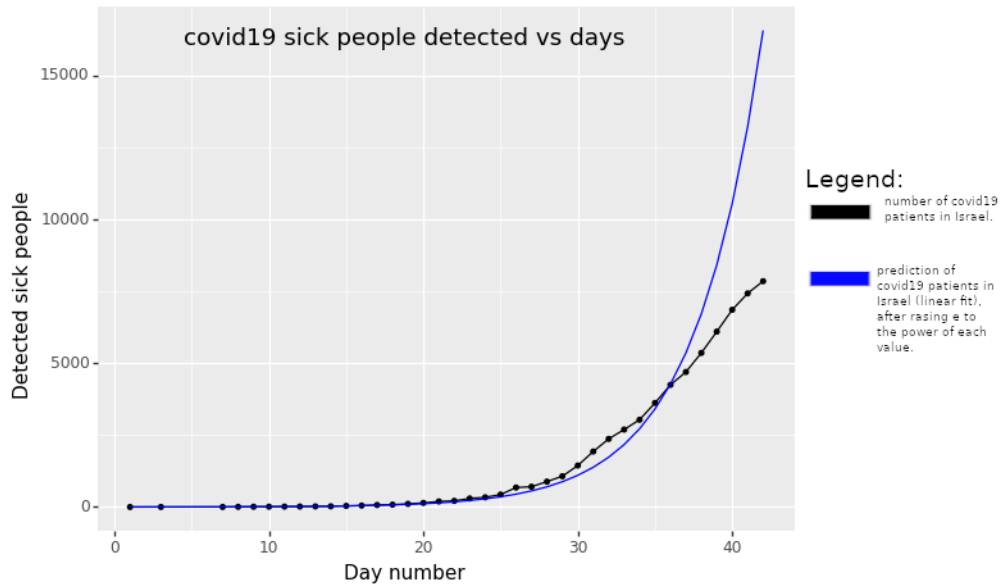
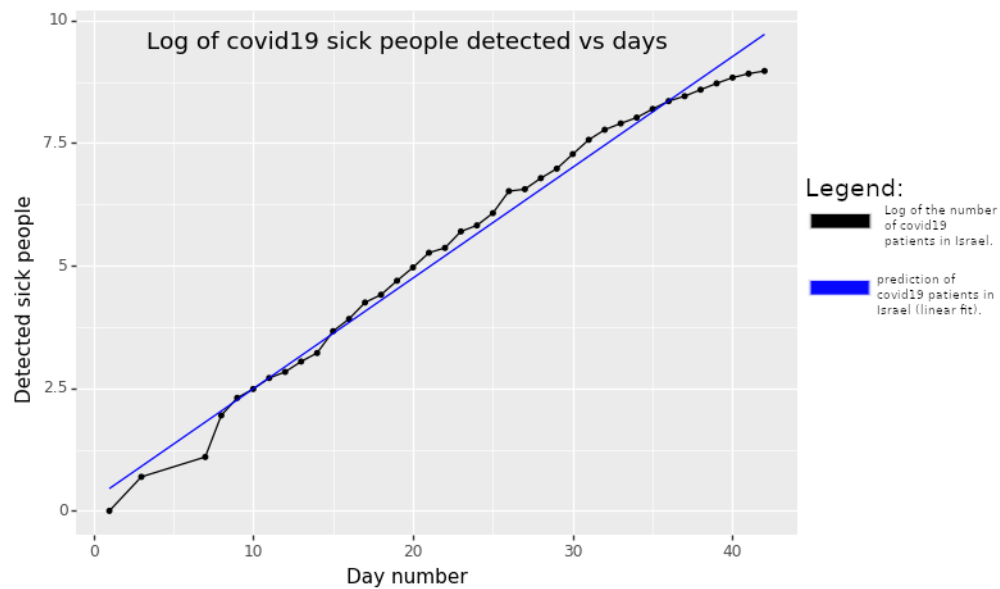


שאלות 18+19+20

אלו שאלות קוד הממומשות בקובץ "covid19.py".

שאלה 21

הגרפים הבאים מראים את המספר חולי הקורונה בישראל אל מול מספר הימים מאז שהתגלה החולה הראשון. הגרף הימני מתאר את מספר החולים באופן לוגריתמי והגרף השני לא. בנוסף בשני הגרפים יש את הקו הכחול שהוא השיערוך שלי, כאשר בגרף השמאלי זה השיערוך של החולים מועלים בחזקת e .



שאלה 22

- אם נרצה להתאים את אותו מודל לינארי אך כאשר נשתמש ב- y המקורי בפונקציית המחר, ולא ב- $\log(y)$, אז עבור (x, y) דגימה כלשהי, ומשקל $w \in \mathbb{R}$ אנחנו נקבל

$$Loss(f_w, (x, y)) = (\langle w, x \rangle - y)$$

נשים לב שמכיוון ש- y אינו מתנהג באופן לינארי, ואנחנו מנסים להתאים לה פונקציה לינארית, אז נקבל פונקציית לוס בעלת ערכים גדולים יחסית.

- אם היינו משתמשים בפונקציית ה-least square loss, אז הייתי מנסה לעשות Gradient Descent שזו שיטת אופטימיזציה איטרטיבית מסדר ראשון למציאת מינימום מקומי של פונקציה. בשיטה זו, נעשה צעד נגדי לגרדיאנט ביחס לנקודה הנוכחית.