

# Contents

<a href="#">1</a>	<a href="#">covid19.py</a>	2
<a href="#">2</a>	<a href="#">ex 2 Amit Baskin.pdf</a>	4
<a href="#">3</a>	<a href="#">linear model.py</a>	18

# 1 covid19.py

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 import os
5
6
7 def fit_linear_regression(design_mat, response_vec):
8     singular_values = np.linalg.svd(design_mat, compute_uv=False)
9     coefficients_vec = np.linalg.pinv(design_mat.T) @ response_vec
10    return coefficients_vec, singular_values
11
12
13 def load_data(path):
14     os.chdir(os.path.dirname(path))
15     return pd.read_csv(os.path.basename(path))
16
17
18 def add_log_detected_col(path):
19     data_df = load_data(path)
20     data_df['log_detected'] = np.log(data_df['detected'])
21     return data_df
22
23
24 def fit_covid(data_df):
25     column_length = len(data_df)
26     ones_col = np.ones(column_length)
27     data_df.insert(loc=0, column='ones_col', value=ones_col)
28
29     design_mat = data_df[['ones_col', 'day_num']].values.T
30     response_vec = data_df['log_detected'].T
31     coefficients_vec = fit_linear_regression(design_mat, response_vec)[0]
32     day_num_vec = np.array(data_df['day_num'])
33     linear_prediction_vec = design_mat.T @ coefficients_vec
34
35     plt.plot(day_num_vec, linear_prediction_vec)
36     plt.scatter(day_num_vec, response_vec)
37     plt.title('log of the number of detected cases\n'
38             'with the corona virus in Israel\n'
39             'as a function of the number of days\n'
40             'since the first case detected')
41     plt.xlabel('number of days')
42     plt.ylabel('log of the number of detected')
43     plt.legend(['prediction for the log of the number of detected',
44             'actual log of the number of detected'])
45     plt.show()
46
47     zero_coef = coefficients_vec[0]
48     first_coef = coefficients_vec[1]
49     exponential_prediction_vec = np.exp(first_coef * day_num_vec) * np.exp(
50         zero_coef)
51
52     plt.plot(day_num_vec, exponential_prediction_vec)
53     plt.scatter(day_num_vec, np.exp(response_vec))
54     plt.title('the number of detected cases\n'
55             'with the corona virus in Israel\n'
56             'as a function of the number of days\n'
57             'since the first case detected')
58     plt.xlabel('number of days')
59     plt.ylabel('number of detected')
```

```
60     plt.legend(['prediction for the number of detected',  
61                'actual number of detected'])  
62     plt.show()
```

## מערכות לומדות תרגיל 2

עמית בסקין 312259013

1. צ.להוכיח ש-  $\ker(X^t) = \ker(XX^t)$ . הוכחה:

כיוון ראשון:  $\ker(X^t) \subseteq \ker(XX^t)$ : יהא  $v \in \mathbb{R}^p$  כך ש-  $X^t v = \mathbf{0}_{\mathbb{R}^n}$ , אז בפרט מתקיים ש-  
 $XX^t v = X\mathbf{0}_{\mathbb{R}^n} = \mathbf{0}_{\mathbb{R}^p}$ .

כיוון שני:  $\ker(X^t) \supseteq \ker(XX^t)$ : יהא  $v \in \mathbb{R}^p$  כך ש-  $XX^t v = \mathbf{0}_{\mathbb{R}^p}$ . אז:

$$v^t XX^t v = 0$$

$$(X^t v)^t (X^t v) = 0$$

$$\langle X^t v, X^t v \rangle = 0$$

קרי:

$$X^t v = \mathbf{0}_{\mathbb{R}^n}$$

$$v \in \ker(X^t)$$

מש"ל.

2. צ.להוכיח עבור מטריצה ריבועית  $A$  ש-  $\operatorname{Im}(A^t) = (\ker(A))^\perp$ . הוכחה:

כיוון ראשון:  $\operatorname{Im}(A^t) \subseteq (\ker(A))^\perp$ : יהא  $v \in \operatorname{Im}(A^t)$  אז קיים  $u$  כך ש-  $A^t u = v$ . יהא  $x \in \ker(A)$  קרי  $Ax = 0$ , ונראה ש-  $\langle v, x \rangle = 0$ . אכן:

$$A^t u = v$$

$$u^t A = v^t$$

$$u^t A x = v^t x$$

$$0 = v^t x$$

$$(v, x) = 0$$

כיוון שני:  $\text{Im}(A^t) \supseteq (\ker(A))^\perp$  יהא  $v \notin (\ker(A))^\perp$  ונראה ש-  $v \notin \text{Im}(A^t)$ . אז קיים  $x \in \ker(A)$  כך ש-  $\langle v, x \rangle \neq 0$ . נניח בשלילה שיש  $u$  כך ש-  $A^t u = v$ , קרי:

$$u^t A = v^t$$

$$u^t A x = v^t x$$

$$0 = v^t x$$

$$\langle v, x \rangle = 0$$

סתירה.

מש"ל.

3. יהא  $y = X^t w$  מערכת לינארית אי־הומוגנית. נניח ש-  $X^t$  ריבועית ולא הפיכה. צ.להראות שלמערכת יש אינסוף פתרונות אם  $y$  מאונך ל-  $\ker(X)$ . הוכחה:  
כיוון ראשון: נניח שלמערכת יש אינסוף פתרונות. נניח בשלילה שיש  $x \in \ker(X)$  כך ש-  $\langle y, x \rangle \neq 0$ . אז:

$$y^t = w^t X$$

$$y^t x = w^t X x$$

$$y^t x = 0$$

$$\langle y, x \rangle = 0$$

סתירה.

כיוון שני: נניח ש-  $y$  מאונך ל-  $\ker(X)$ . אז מהסעיף הקודם מתקיים ש-  $y \in \operatorname{Im}(X^t)$ , קרי יש  $w$  כך ש-  $X^t w = y$ . נניח בשלילה שזה הפיתרון היחיד ונקבל ש-  $X$  מדרגה מלאה ולכן הפיכה בסתירה לנתון. מש"ל

4. נתבונן במערכת הנורמלית הלינארית  $XX^t w = Xy$ . צ.להוכיח שלמערכת הנורמלית יש פיתרון יחיד אם  $XX^t$  הפיכה ואינסוף פתרונות אחרת. הוכחה: מצד אחד אם  $XX^t$  הפיכה אז  $w$  חייב לקיים:

$$w = (XX^t)^{-1} Xy$$

קרי יש פיתרון יחיד.

מצד שני נניח ש-  $XX^t$  לא הפיכה. לפי סעיף 3 מספיק להראות ש-  $Xy$  מאונך ל-  $\ker(XX^t)$  ונקבל שיש אינסוף פתרונות.

כמו כן הראנו בסעיף 1 שלהיות מאונך ל-  $\ker(XX^t)$  זה שקול ללהיות מאונך ל-  $\ker(X^t)$ , ולפי סעיף 2 זה שקול ללהיות ב-  $\operatorname{Im}(X)$ , ואכן,  $Xy \in \operatorname{Im}(X)$ , כנדרש.

מש"ל

5.

א. צ.להוכיח ש-  $P$  סימטרית. הוכחה: לכל  $i$  נסמן:

$$v_i = (v_{i1}, \dots, v_{in})$$

ונקבל:

$$P = \sum_{i=1}^k v_i \otimes v_i^t = \sum_{i=1}^k \begin{bmatrix} v_{i1}^2 & v_{i1}v_{i2} & \cdots & v_{i1}v_{ik} \\ v_{i2}v_{i1} & v_{i2}^2 & \cdots & v_{i2}v_{ik} \\ \vdots & \vdots & \ddots & \vdots \\ v_{ik}v_{i1} & v_{ik}v_{i2} & \cdots & v_{ik}^2 \end{bmatrix} =$$

$$= \begin{bmatrix} \sum_{i=1}^k v_{i1}^2 & \sum_{i=1}^k v_{i1}v_{i2} & \cdots & \sum_{i=1}^k v_{i1}v_{ik} \\ \sum_{i=1}^k v_{i2}v_{i1} & \sum_{i=1}^k v_{i2}^2 & \cdots & \sum_{i=1}^k v_{i2}v_{ik} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^k v_{ik}v_{i1} & \sum_{i=1}^k v_{ik}v_{i2} & \cdots & \sum_{i=1}^k v_{ik}^2 \end{bmatrix}$$

ואכן עבור אינדקס  $j$  כלשהו, השורה ה- $j$  היא:

$$\left( \sum_{i=1}^k v_{ij}v_{i1} \cdots \sum_{i=1}^k v_{ij}^2 \cdots \sum_{i=1}^k v_{ij}v_{ik} \right)$$

והעמודה ה- $j$ :

$$\begin{pmatrix} \sum_{i=1}^k v_{i1}v_{ij} \\ \vdots \\ \sum_{i=1}^k v_{ij}^2 \\ \vdots \\ \sum_{i=1}^k v_{ik}v_{ij} \end{pmatrix}$$

אז אם נסמן את וקטור העמודה ה- $j$  ב- $c_j$  ואת וקטור השורה ה- $j$  ב- $r_j$ , נקבל שמתקיים:

$$c_j = r_j^t$$

כלומר,  $P$  סימטרית.

מש"ל

ב. צלוחיך שהוקטורים העצמיים של  $P$  הם 0 או 1 וש- $v_1, \dots, v_k$  הם הוקטורים העצמיים המתאימים

לערכך העצמי 1. הוכחה: נתבונן במטריצה הסימטרית האורתונורמלית:

$$Q = \begin{bmatrix} v_{11} & v_{21} & \cdots & v_{k1} \\ v_{12} & v_{22} & \cdots & v_{k2} \\ \vdots & \vdots & \ddots & \vdots \\ v_{1k} & v_{2k} & \cdots & v_{kk} \end{bmatrix}$$

ונשים לב כי:

$$QIQ^t =$$

$$= \begin{bmatrix} v_{11} & v_{21} & \cdots & v_{k1} \\ v_{12} & v_{22} & \cdots & v_{k2} \\ \vdots & \vdots & \ddots & \vdots \\ v_{1k} & v_{2k} & \cdots & v_{kk} \end{bmatrix} \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix} \begin{bmatrix} v_{11} & v_{12} & \cdots & v_{1k} \\ v_{21} & v_{22} & \cdots & v_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ v_{k1} & v_{k2} & \cdots & v_{kk} \end{bmatrix} =$$

$$\begin{aligned}
&= \begin{bmatrix} v_{11} & v_{21} & \cdots & v_{k1} \\ v_{12} & v_{22} & \cdots & v_{k2} \\ \vdots & \vdots & \ddots & \vdots \\ v_{1k} & v_{2k} & \cdots & v_{kk} \end{bmatrix} \begin{bmatrix} v_{11} & v_{12} & \cdots & v_{1k} \\ v_{21} & v_{22} & \cdots & v_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ v_{k1} & v_{k2} & \cdots & v_{kk} \end{bmatrix} = \\
&= \begin{bmatrix} \sum_{i=1}^k v_{i1}^2 & \sum_{i=1}^k v_{i1}v_{i2} & \cdots & \sum_{i=1}^k v_{i1}v_{ik} \\ \sum_{i=1}^k v_{i2}v_{i1} & \sum_{i=1}^k v_{i2}^2 & \cdots & \sum_{i=1}^k v_{i2}v_{ik} \\ \vdots & \vdots & \ddots & \vdots \\ \sum_{i=1}^k v_{ik}v_{i1} & \sum_{i=1}^k v_{ik}v_{i2} & \cdots & \sum_{i=1}^k v_{ik}^2 \end{bmatrix} = P
\end{aligned}$$

כלומר, מדובר בפירוק ספקטרלי של  $P$  ולכן האיברים באלכסון של  $I$  הם הערכים העצמיים, קרי הערך העצמי 1 וכן הוקטורים העצמיים הם העמודות של  $Q$ , קרי וקטורי הבסיס האורתונורמלי של המרחב  $V$ , כנדרש.

מש"ל

ג. ישירות מהסעיף הקודם מתקבל, לכל  $v \in V$  מתקיים  $Pv = v$ , כי  $v$  הוא וקטור עצמי עם ערך עצמי 1.

מש"ל

ד.

$$P^2 = QQ^tQQ^t$$

אבל:

$$Q^tQ =$$

$$\begin{aligned}
&= \begin{bmatrix} v_{11} & v_{12} & \cdots & v_{1k} \\ v_{21} & v_{22} & \cdots & v_{2k} \\ \vdots & \vdots & \ddots & \vdots \\ v_{k1} & v_{k2} & \cdots & v_{kk} \end{bmatrix} \begin{bmatrix} v_{11} & v_{21} & \cdots & v_{k1} \\ v_{12} & v_{22} & \cdots & v_{k2} \\ \vdots & \vdots & \ddots & \vdots \\ v_{1k} & v_{2k} & \cdots & v_{kk} \end{bmatrix} = \\
&= \begin{bmatrix} v_1^t v_1 & v_1^t v_2 & \cdots & v_1^t v_k \\ v_2^t v_1 & v_2^t v_2 & \cdots & v_2^t v_k \\ \vdots & \vdots & \ddots & \vdots \\ v_k^t v_1 & v_k^t v_2 & \cdots & v_k^t v_k \end{bmatrix}
\end{aligned}$$



ומכך שמדובר בוקטורים של בסיס אורתונורמלי:

$$= \begin{bmatrix} 1 & 0 & \cdots & 0 \\ 0 & 1 & \cdots & 0 \\ \vdots & \vdots & \ddots & \vdots \\ 0 & 0 & \cdots & 1 \end{bmatrix} = I$$

ולכן:

$$P^2 = QQ^tQQ^t = P^2 = QIQ^t = QQ^t = P$$

מש"ל.

ד. ישירות מהסעיף הקודם:

$$P^2 - P = 0$$

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

מש"ל.

6.

א. יהא  $X = U\Sigma V^t$  פירוק SVD של  $X$ . אז:

$$XX^t = U\Sigma V^t V \Sigma^t U^t$$

ומכך שההפכית של מטריצה אורתוגונלית היא השחלוף שלה:

$$= U\Sigma I \Sigma^t U^t = U\Sigma \Sigma^t U^t$$

ולפי הסימון בתרגיל:

$$= UDU^t$$

ומאחר שאנחנו מניחים ש-  $XX^t$  הפיכה, אזי שגם כל אחת מ-  $UDU^t$  הפיכה כי הדטרמיננטה של  $XX^t$

היא מכפלת הדטרמיננטות של המכפלה הנ"ל. מכאן שקיימת  $D^{-1}$ , ובפרט:

$$(UDU^t)^{-1} = (U^t)^{-1} D^{-1} U^{-1}$$

ושוב מאחר שההפכית של אורתונורמלית היא השחלוף:

$$= U D^{-1} U^t$$

מש"ל.

ב. מצד אחד:

$$(XX^t)^{-1}X = UD^{-1}U^tU\Sigma V^t =$$

$$= UD^{-1}\Sigma V^t = U(\Sigma\Sigma^t)^{-1}\Sigma V^t$$

אבל  $\Sigma\Sigma^t$  אלכסונית וסימטרית ולכן:

$$(\Sigma\Sigma^t)^{-1} = (\Sigma\Sigma^t)^\dagger = \Sigma^{t\dagger}\Sigma^\dagger$$

קרי:

$$U(\Sigma\Sigma^t)^{-1}\Sigma V^t = U\Sigma^{t\dagger}\Sigma^\dagger\Sigma V^t =$$

$$= U\Sigma^{t\dagger}IV^t = U\Sigma^{t\dagger}V^t$$

ומצד שני:

$$X^{t\dagger} = (V\Sigma^tU^t)^\dagger = U\Sigma^{t\dagger}V^t$$

מש"ל.

7. בשאלה 1 ראינו ש-  $\ker(X^t) = \ker(XX^t)$ . עתה  $XX^t$  הפיכה אסס  $\ker(XX^t) = \{0_{\mathbb{R}^d}\}$  אסס  $\ker(X^t) = \{0_{\mathbb{R}^m}\}$ . אבל  $X: \mathbb{R}^m \rightarrow \mathbb{R}^d$ , אז לפי משפט המימדים השני, הנ"ל מתקיים אסס מימד מרחב העמודות שהוא מימד התמונה של  $X$ , שווה ל-  $d$ , אבל מרחב העמודות הוא בדיוק  $\text{span}\{x_1, \dots, x_m\}$ , קרי הנ"ל קורה אסס  $\text{span}\{x_1, \dots, x_m\} = \mathbb{R}^d$  כנדרש.

מש"ל

8. נסמן  $\hat{w} = X^{t\dagger}y$ . צ. להראות שלכל פיתרון  $\bar{w}$  של  $XX^tw = Xy$ , מתקיים ש-  $\|\hat{w}\|_2 \leq \|\bar{w}\|_2$ .

הוכחה: יהא  $X = U\Sigma V^t$  פירוק SVD של  $X$ . ניזכר כי:

$$XX^t = U\Sigma V^tV\Sigma^tU^t = U\Sigma\Sigma^tU^t$$

ונתבונן במשוואה שלנו:

$$XX^tw = Xy$$

כלומר:

$$U\Sigma\Sigma^tU^tw = U\Sigma V^ty$$

נסמן ב-  $\Sigma_1$  את הבלוק של  $\Sigma$  שאינו אפסים ונקבל:

$$U \begin{bmatrix} \Sigma_1^2 & O \\ O & O \end{bmatrix} U^tw = U \begin{bmatrix} \Sigma_1 & O \\ O & O \end{bmatrix} V^ty$$

נכפיל ב-  $U^t$  משמאל:

$$\begin{bmatrix} \Sigma_1^2 & O \\ O & O \end{bmatrix} U^tw = \begin{bmatrix} \Sigma_1 & O \\ O & O \end{bmatrix} V^ty$$

נחלק לבלוקים את  $U^tw$  ו-  $V^ty$  בהתאם לחלוקה הקודמת:

$$\begin{bmatrix} \Sigma_1^2 & O \\ O & O \end{bmatrix} \begin{bmatrix} U_1^tw \\ U_2^tw \end{bmatrix} = \begin{bmatrix} \Sigma_1 & O \\ O & O \end{bmatrix} \begin{bmatrix} V_1^ty \\ V_1^ty \end{bmatrix}$$

$$: \begin{bmatrix} (\Sigma_1^\dagger)^2 & O \\ O & O \end{bmatrix} \text{ נכפיל משמאל ב-}$$

$$\begin{bmatrix} (\Sigma_1^\dagger)^2 & O \\ O & O \end{bmatrix} \begin{bmatrix} \Sigma_1^2 & O \\ O & O \end{bmatrix} \begin{bmatrix} U_1^tw \\ U_2^tw \end{bmatrix} = \begin{bmatrix} (\Sigma_1^\dagger)^2 & O \\ O & O \end{bmatrix} \begin{bmatrix} \Sigma_1 & O \\ O & O \end{bmatrix} \begin{bmatrix} V_1^ty \\ V_1^ty \end{bmatrix}$$

ונקבל:

$$\begin{bmatrix} U_1^tw \\ O \end{bmatrix} = \begin{bmatrix} I_1 & O \\ O & O \end{bmatrix} \begin{bmatrix} U_1^tw \\ U_2^tw \end{bmatrix} = \begin{bmatrix} \Sigma_1^\dagger & O \\ O & O \end{bmatrix} \begin{bmatrix} V_1^ty \\ V_2^ty \end{bmatrix} = \begin{bmatrix} \Sigma_1^{-1}V_1^ty \\ \mathbf{0}_{d-r} \end{bmatrix}$$

ובה"כ קיבלנו:

$$U_1^tw = \Sigma_1^\dagger V_1^ty$$

קרי:

$$w_1 = U_1 \Sigma_1^\dagger V_1^ty$$

ו-  $U_2^tw$  חופשי. אז נסמן:

$$U_2^tw = z$$

קרי:

$$w_2 = U_2 z$$

ובסה"כ:

$$w = \begin{bmatrix} U_1 \Sigma_1^\dagger V_1^t y \\ U_2 z \end{bmatrix}$$

וזאת עבור פיתרון כללי כאשר דרגת החופש היא ב- $z$ .

אז:

$$\begin{aligned} \|w\|_2^2 &= \left\| \begin{bmatrix} U_1 \Sigma_1^\dagger V_1^t y \\ U_2 z \end{bmatrix} \right\|_2^2 = \\ &= \|U_1 \Sigma_1^\dagger V_1^t y\|_2^2 + \|U_2 z\|_2^2 \end{aligned}$$

ונורמה מינימלית מתקבלת עבור  $z = 0$ , קרי:

$$\|U_1 \Sigma_1^\dagger V_1^t y\|_2^2 = \|U \Sigma^\dagger V^t y\|_2^2$$

אבל:

$$\hat{w} = X^{t\dagger} y = U \Sigma^{t\dagger} V^t y$$

כלומר קיבלנו את הנורמה של  $\hat{w}$ , כנדרש.

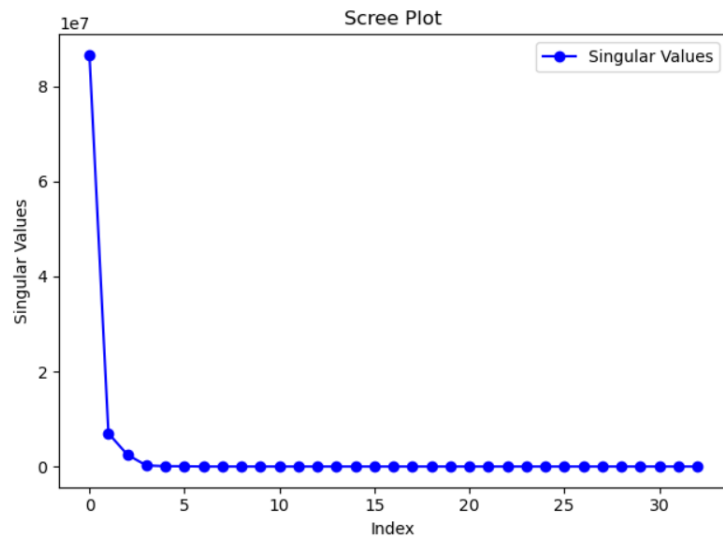
מש"ל

13.

- באשר לעמודה של התאריך: החלטתי לפצל לפי שנה וחודש, על מנת לקבל אינדיקציה על השפעה של החודש והשנה שבהם בית נמכר, וזאת בניגוד לתאריכים ספיציפיים מדי, לפי ימים, ללא משמעות כמותית שאותה ניתן לכלול ברגרסיה הלינארית.
- כמו כן, החלטתי להתבונן במחיר הממוצע פר מיקוד במקום על המיקוד עצמו, על מנת לתת איפיון נוסף לטיב המיקום של הבית, שכן, המיקוד לבדו לא חושף מידע רלוונטי עבור הרגרסיה הלינארית שהרי אין לו משמעות כמותית בפני עצמו.
- את העמודה של תעודת הזהות מצאתי כלא רלוונטית ללמידה שנפעיל על המידע ולכן הסרתי את העמודה.

- באשר לעמודות שמכילות טווח ערכים מצומצם שאיננו מנורמל, טווח שנקודות הקצה שלו קשורים במשמעות של המידע, כמו קו רוחב, קו גובה, שנת בנייה, שנת שיפוץ, את כל אלה החלטתי לנרמל, שהרי המידע החשוב הוא מיקום הערך ביחס לקצוות של הטווח, וכך ניתן לראות את המידע בצורה ברורה יותר ביחס למה שמעניין אותנו לגביו, וכן המספרים יותר נוחים לעיבוד.

15.



- המטריצה  $X$  היא בקירוב מסדר  $30 \times 20,000$ , קרי יש לנו כ- 20,000 דגימות עם כ- 30 תכונות לכל דגימה. אז  $X$  לא ולכן לא הפיכה, ובפרט מימד הגרעין של  $X$  הינו גדול מאפס.

- נתבונן ב-  $X$  כהעתקה מ-  $\mathbb{R}^{20,000}$  ל-  $\mathbb{R}^{30}$ . לפי משפט המימדים השני, מתקיים כי:

$$\dim(\text{Im}(X)) + \dim(\ker(X)) = \dim(\mathbb{R}^{20,000})$$

ובכן, הדרגה של המטריצה היא לכל היותר המינימום שבין מימד מרחב העמודות למימד מרחב השורות, קרי 30.

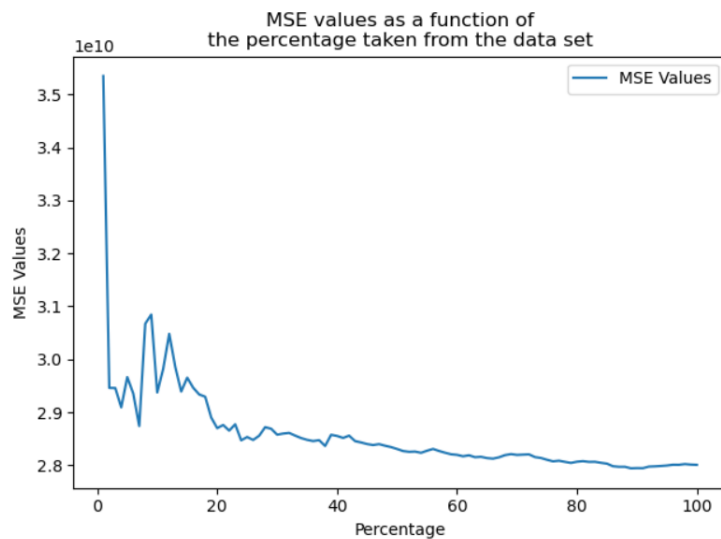
- אבל מימד התמונה הוא מימד מרחב העמודות. לכן:  $\dim(\text{Im}(X)) \leq 30$ , ומכאן ש-  $\dim(\ker(X)) \geq 19,970$ .

- נתבונן בפירוק  $SVD$  של  $X$ , קרי  $X = U\Sigma V^t$ . הערכים על האלכסון של  $\Sigma$  הם הערכים הסינגולריים של  $X$  שמקיימים  $Xv = \sigma u$  עבור  $u, v$  וקטורי העמודות של  $U, V$  בהתאמה. אז לכל  $v$  עם ערך סינגולרי  $\sigma = 0$ , מתקיים ש-  $v \in \ker(X)$ . ונשים לב שכל שורה של  $X$  היא שורה של פיצ'ר. אז אם  $\sigma = 0$ , זאת אומרת שלכל שורה  $x_i$  ב-  $X$  מתקיים ש-  $x_i v = 0$ , קרי יש תלות לינארית בין הדגימות ביחס לפיצ'ר ה-  $i$ , וזאת לכל  $i$ . אז

ככל שיש לנו יותר ערכים סינגולריים אפסיים כך התלות בין עמודות המטריצה גדולה יותר (מה שצפוי שיקרה מאחר שיש לנו סדר גודל של 30 משתנים לעומת כ-20,000 משוואות).

- אם כן, הנ"ל מסתדר עם הגרף המוצג: העובדה שרוב הערכים הסינגולריים הם אפס או מאוד קרובים לאפס, עולה בקנה אחד עם העובדה שמימד הגרעין של  $X$  גדול.
- בשאלה 1 ראינו ש-  $\ker(X) = \ker(XX^t)$ . מכאן שגם  $XX^t$  לא הפיכה, שהרי לכל מטריצה מתקיים שהיא הפיכה אםס מימד הגרעין שלה הוא אפס.
- כמו כן, בשאלה 4 ראינו שלמשוואה הנורמלית יש פיתרון יחיד אםס  $XX^t$  הפיכה, ואינסוף פתרונות אחרת.
- ואכן, יש לנו הרבה מאוד דגימות ביחס למספר המשתנים, ולאור הניתוח הנ"ל, ניתן להיווכח בתרומה של מספר גדול של דגימות, למציאת פתרון לבעיית הרגרסיה הלינארית שהמטריצה  $X$  היא מטריצת העיצוב שלה.
- אז לסיכום, ניתן לומר שגרף הערכים הסינגולריים מאשרר את הסינגולריות של  $X$ , ואכן, כאמור בניתוח הנ"ל,  $X$  רחוקה מאוד מלהיות הפיכה.

16.



הסבר:

$w$  הוא וקטור המקדמים שמתאר את הקשר הלינארי בין התכונות של כל דגימה ובין המחיר של אותו בית שעליו נעשתה הדגימה.

$\hat{y}$  הוא וקטור המחירים שאנו מנחשים באמצעות  $w$  ו-  $y$  הוא וקטורים המחירים האמיתיים.

ככל שלוקחים יותר דגימות, כך  $w$  נעשה מדויק יותר, ולכן, כלומר  $\hat{y}$  קרוב יותר ל-  $y$ , ולכן ממוצע סכום ההפרשים ביניהם דועך ככל שמספר הדגימות גדל.

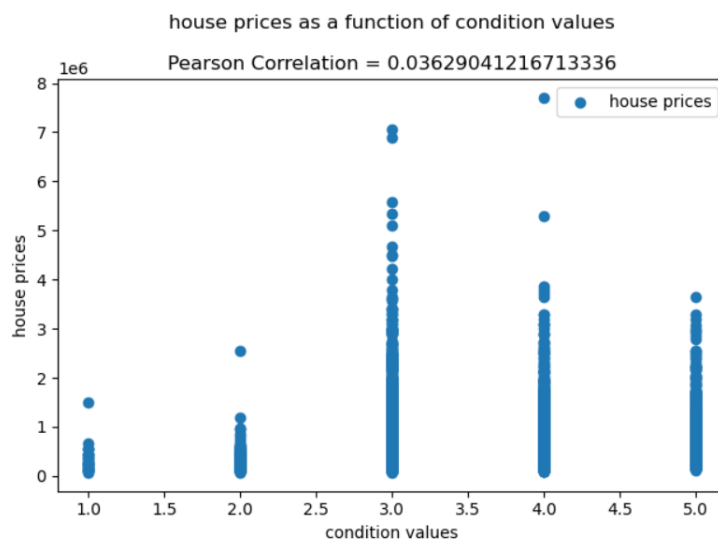
- תכונה שתורמת למודל: **שטח הבית**



אכן, ניתן לראות בגרף שיש קשר בין השטח של הבית ובין המחיר שבו הוא נמכר: ככל שהשטח גדול יותר כך מחיר הבית נוטה להיות גבוה יותר, ואכן הקורולציה המחושבת היא הגבוהה ביותר מבין הקורולציות המחושבות.

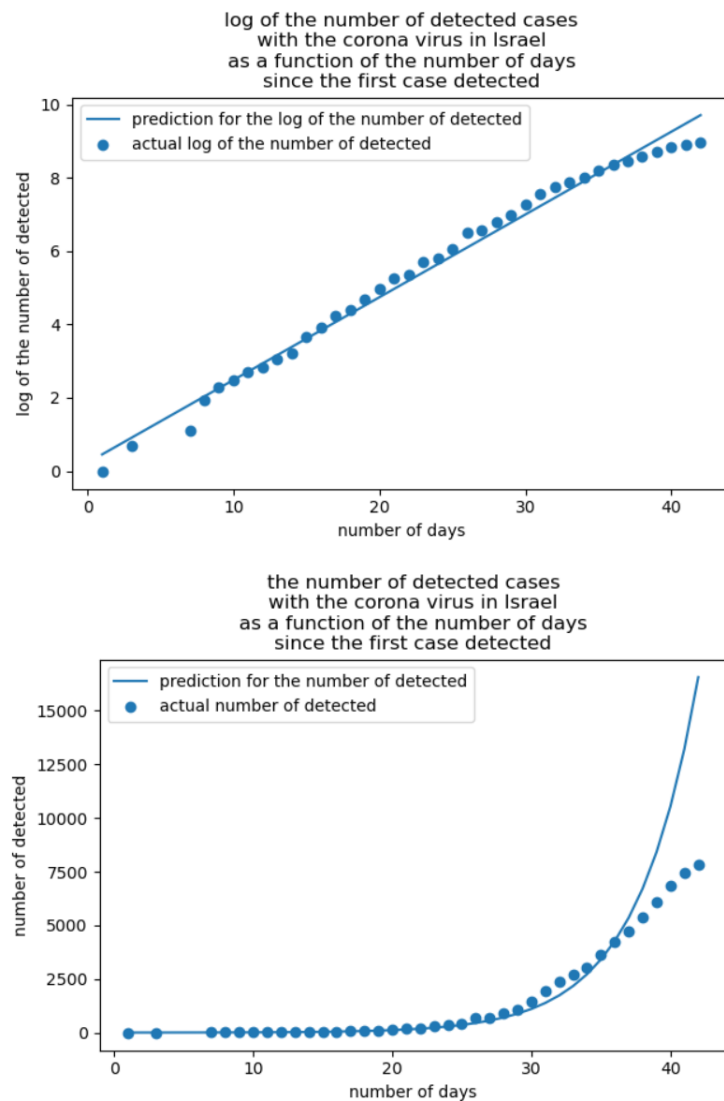
מכאן שזוהי תכונה שעוזרת למודל להתאים בין התכונות של הבית ובין המחיר שלו, ולכן היא תורמת.

- תכונה שאינה תורמת למודל: **מצב הבית**



ניתן לראות שדווקא כאשר מצב הבית הינו באמצע בין הגרוע ביותר לטוב ביותר, יש כמות גבוהה יותר של מחירים גבוהים. כמו כן, בכל אחד מהדירוגים יש בתים במחירים נמוכים. אמנם המחירים מטפסים קצת יותר כאשר מצב הבית טוב יותר, אבל עצם זה שהרוב מרוכז באמצע, מותיר עמימות בפני המודל, שכן זה לא נותן אינדיקציה לינארית לגבי הקשר בין מצב הבית למחיר. אכן, מדד הקורולציה הוא גם הנמוך ביותר. מכאן שתכונה זו לא תורמת למודל.

21.



אופן החישוב של הגרף האקספוננציאלי:

נסמן את המשתנה של prediction\_detected ב-  $\hat{y}$  ואת המשתנה של day\_num ב-  $x$ . נניח שוקטור



המקדמים שמצאנו הוא  $(\beta_0, \beta_1)$  כאשר  $\beta_1$  הוא המקדם של  $x$  ו- $\beta_0$  הוא הגורם החופשי. אז:

$$\log \hat{y} = \beta_1 x + \beta_0$$

$$e^{\log \hat{y}} = e^{\beta_1 x + \beta_0}$$

$$\hat{y} = e^{\beta_1 x} e^{\beta_0}$$

22. בהתבסס על הנימוק שסיפקתי בשאלה הקודמת, במקרה האקספוננציאלי, עבור  $x = (x_0, x_1, \dots, x_n)$  ו- $w = (w_0, w_1, \dots, w_n)$  כאשר  $x_0 = 1$ , פונקציית ההפסד צריכה להיות:

$$\left( y - \prod_{i=1}^n e^{w_i x_i} \right)^2 = \left( y - e^{x^t w} \right)^2$$

ואז כדי למצוא את  $w$ , הפיתרון ל- $ERM$ , נפעל בדומה להרצאה: נרצה למזער את:

$$\sum_{i=1}^m \left( y_i - e^{x_i^t w} \right)^2$$

אז נגזור לפי  $w$  ונשווה לאפס. נפתור את המשוואה שמתקבלת ונקבל את  $w$ .

## 3 linear model.py

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from dateutil.parser import parse
5 import re
6 import os
7
8
9 def fit_linear_regression(design_mat, response_vec):
10     singular_values = np.linalg.svd(design_mat, compute_uv=False)
11     coefficients_vec = np.linalg.pinv(design_mat.T) @ response_vec
12     return coefficients_vec, singular_values
13
14
15 def predict(design_mat, coefficients_vec):
16     return np.matmul(design_mat.T, coefficients_vec)
17
18
19 def mse(prediction_vec, response_vec):
20     diff = prediction_vec - response_vec
21     vec_len = len(response_vec)
22     return np.sum(diff ** 2) / vec_len
23
24
25 def is_int_and_not_negative(num):
26     if isinstance(num, int) or (isinstance(num, float) and num % 1 == float(0)):
27         if num >= 0:
28             return True
29     return False
30
31
32 def is_float_and_not_negative(num):
33     if isinstance(num, float) or isinstance(num, int):
34         if num >= 0:
35             return True
36     return False
37
38
39 def load_data(path):
40     data_df = load_valid_data(path)
41     data_df = data_df.iloc[:, 1:21]
42
43     data_price_by_zip_code = data_df[['price', 'zipcode']]
44     mean_price = data_price_by_zip_code.groupby('zipcode').mean()
45     mean_price = mean_price.rename(
46         columns={'zipdoe': 'zipcode', 'price': 'mean_price_by_zipcode'})
47
48     data_df = data_df.merge(mean_price, on='zipcode')
49     data_df = data_df.drop(['zipcode'], axis=1)
50
51     data_df[['year', 'month']] = data_df.date.str.split(expand=True)
52     data_df = pd.get_dummies(data_df, columns=['year', 'month'])
53     data_df = data_df.drop(['date'], axis=1)
54
55     cols = data_df.columns.tolist()
56
57     def convert_to_int(string):
58         return int(string) if string.isdigit() else string
59
```

```

60     def keys(string):
61         return [convert_to_int(component)
62                 for component in re.split('\d+', string)]
63
64     cols.sort(key=keys)
65     data_df = data_df[cols]
66
67     return data_df
68
69
70 def load_valid_data(path):
71     os.chdir(os.path.dirname(path))
72     data_df = pd.read_csv(os.path.basename(path))
73     for row in data_df.iterrows():
74         if not is_id(row):
75             data_df = data_df.drop([row[0]], axis=0)
76             continue
77
78         if not is_date(row):
79             data_df = data_df.drop([row[0]], axis=0)
80             continue
81         date = row[1]['date']
82         actual_date = parse(date, fuzzy=True)
83         year = actual_date.year
84         month = actual_date.month
85         data_df.at[row[0], 'date'] = str(year) + ' ' + str(month)
86
87         if not is_price(row):
88             data_df = data_df.drop([row[0]], axis=0)
89             continue
90
91         if not is_bedrooms_num(row):
92             data_df = data_df.drop([row[0]], axis=0)
93             continue
94
95         if not is_bathrooms_num(row):
96             data_df = data_df.drop([row[0]], axis=0)
97             continue
98
99         if not is_sqft_living(row):
100             data_df = data_df.drop([row[0]], axis=0)
101             continue
102
103         if not is_sqft_lot(row):
104             data_df = data_df.drop([row[0]], axis=0)
105             continue
106
107         if not is_floors(row):
108             data_df = data_df.drop([row[0]], axis=0)
109             continue
110
111         if not is_waterfront(row):
112             data_df = data_df.drop([row[0]], axis=0)
113             continue
114
115         if not is_view(row):
116             data_df = data_df.drop([row[0]], axis=0)
117             continue
118
119         if not is_condition(row):
120             data_df = data_df.drop([row[0]], axis=0)
121             continue
122
123         if not is_grade(row):
124             data_df = data_df.drop([row[0]], axis=0)
125             continue
126
127         if not is_sqft_above(row):

```

```

128         data_df = data_df.drop([row[0]], axis=0)
129         continue
130
131     if not is_sqft_basement(row):
132         data_df = data_df.drop([row[0]], axis=0)
133         continue
134
135     if not is_yr_built(row):
136         data_df = data_df.drop([row[0]], axis=0)
137         continue
138     data_df.at[row[0], 'yr_built'] /= 2015
139
140     if not is_yr_renovated(row):
141         data_df = data_df.drop([row[0]], axis=0)
142         continue
143     data_df.at[row[0], 'yr_renovated'] /= 2015
144
145     if not is_zipcode(row):
146         data_df = data_df.drop([row[0]], axis=0)
147         continue
148
149     if not is_lat(row):
150         data_df = data_df.drop([row[0]], axis=0)
151         continue
152     data_df.at[row[0], 'lat'] /= 47.8
153
154     if not is_long(row):
155         data_df = data_df.drop([row[0]], axis=0)
156         continue
157     data_df.at[row[0], 'long'] /= -123
158
159     if not is_sqft_living15(row):
160         data_df = data_df.drop([row[0]], axis=0)
161         continue
162
163     if not is_sqft_lot15(row):
164         data_df = data_df.drop([row[0]], axis=0)
165         continue
166
167     return data_df
168
169
170 def is_id(row):
171     id_num = row[1]['id']
172     if is_int_and_not_negative(id_num):
173         return 10 ** 6 < id_num < 10 ** 10
174
175
176 def is_date(row):
177     date = row[1]['date']
178     if isinstance(date, str):
179         nums = list(map(int, re.findall(r'\d+', date)))
180     else:
181         return False
182     is_date = None
183     date_str = ''
184     if len(nums) > 0:
185         date_str = str(nums[0])
186     try:
187         actual_date = parse(date, fuzzy=True)
188         year = actual_date.year
189         if year == 2014 or year == 2015:
190             is_date = True
191
192     except ValueError:
193         is_date = False
194     if is_date:
195         date_format = date_str + 'T000000'

```

```

196         if date_format._eq_(date):
197             return True
198         return False
199
200
201 def is_price(row):
202     price = row[1]['price']
203     if is_float_and_not_negative(price):
204         return 75 * 10 ** 3 <= price <= 8 * 10 ** 6
205
206
207 def is_bedrooms_num(row):
208     num = row[1]['bedrooms']
209     if is_int_and_not_negative(num):
210         return num <= 33
211
212
213 def is_bathrooms_num(row):
214     num = row[1]['bathrooms']
215     if is_float_and_not_negative(num):
216         return num <= 8
217
218
219 def is_sqft_living(row):
220     num = row[1]['sqft_living']
221     if is_float_and_not_negative(num):
222         return 200 <= num <= 14 * 10 ** 3
223
224
225 def is_sqft_lot(row):
226     num = row[1]['sqft_lot']
227     if is_float_and_not_negative(num):
228         return 520 <= num <= 1.7 * 10 ** 6
229
230
231 def is_floors(row):
232     num = row[1]['floors']
233     if is_float_and_not_negative(num):
234         return 1 <= num <= 5
235
236
237 def is_waterfront(row):
238     num = row[1]['waterfront']
239     if is_int_and_not_negative(num):
240         return num <= 1
241
242
243 def is_view(row):
244     num = row[1]['view']
245     if is_int_and_not_negative(num):
246         return num <= 4
247
248
249 def is_condition(row):
250     num = row[1]['condition']
251     if is_float_and_not_negative(num):
252         return num <= 5
253
254
255 def is_grade(row):
256     num = row[1]['grade']
257     if is_float_and_not_negative(num):
258         return 1 <= num <= 13
259
260
261 def is_sqft_above(row):
262     num = row[1]['sqft_above']
263     if is_float_and_not_negative(num):

```

```

264         return 200 <= num <= 10 ** 4
265
266
267 def is_sqft_basement(row):
268     num = row[1]['sqft_basement']
269     if is_float_and_not_negative(num):
270         return num <= 4820
271
272
273 def is_yr_built(row):
274     num = row[1]['yr_built']
275     if is_int_and_not_negative(num):
276         return 1900 <= num <= 2015
277
278
279 def is_yr_renovated(row):
280     num = row[1]['yr_renovated']
281     if is_int_and_not_negative(num):
282         return num == 0 or 1900 <= num <= 2015
283
284
285 def is_zipcode(row):
286     num = row[1]['zipcode']
287     if is_int_and_not_negative(num):
288         return 98 * 10 * 3 <= num <= 98.2 * 10 ** 3
289
290
291 def is_lat(row):
292     num = row[1]['lat']
293     if is_float_and_not_negative(num):
294         to_return = 47 <= num <= 48
295         return to_return
296
297
298 def is_long(row):
299     num = row[1]['long']
300     return -123 <= num <= -121
301
302
303 def is_sqft_living15(row):
304     num = row[1]['sqft_living15']
305     if is_float_and_not_negative(num):
306         return 300 <= num <= 7 * 10 ** 3
307
308
309 def is_sqft_lot15(row):
310     num = row[1]['sqft_lot15']
311     if is_float_and_not_negative(num):
312         return 600 <= num <= 880 * 10 ** 3
313
314
315 def plot_singular_values(singular_values):
316     singular_values.sort()
317     singular_values = singular_values[::-1]
318     plt.plot(np.arange(singular_values.shape[0]), singular_values, 'bo-')
319     plt.title('Scree Plot')
320     plt.xlabel('Index')
321     plt.ylabel('Singular Values')
322     plt.legend(['Singular Values'])
323     plt.show()
324
325
326 def add_ones_col(data_df):
327     column_length = len(data_df)
328     ones_col = np.ones(column_length)
329     data_df.insert(loc=0, column='ones_col', value=ones_col)
330
331

```

```

332 def q15(path):
333     data_df = load_data(path)
334     add_ones_col(data_df)
335     data_df = data_df.drop('price', axis=1)
336     data_mat = data_df.values.T
337     singular_values = np.linalg.svd(data_mat, compute_uv=False)
338     plot_singular_values(np.array(singular_values))
339
340
341 def q17(path):
342     data_df = load_data(path)
343     add_ones_col(data_df)
344     data_df = data_df.sample(frac=1)
345     response_vec = np.array(data_df['price'])
346     data_df = data_df.drop('price', axis=1)
347     desin_mat = data_df.values
348
349     rows_amount = desin_mat.shape[0]
350     test_mat_size = int(0.25 * rows_amount)
351     test_mat = desin_mat[:test_mat_size]
352     test_response_vec = response_vec[:test_mat_size]
353     train_mat = desin_mat[test_mat_size:]
354     train_response_vec = response_vec[test_mat_size:]
355
356     coefficients_vecs_lst = []
357     train_mat_rows_amount = train_mat.shape[0]
358     for i in range(1, 101):
359         current_rows_amount = int(i / 100 * train_mat_rows_amount)
360         current_train_mat = train_mat[:current_rows_amount]
361         current_response_vec = train_response_vec[:current_rows_amount]
362         current_coefficients_vec = \
363             fit_linear_regression(current_train_mat.T, current_response_vec)[0]
364         coefficients_vecs_lst.append(current_coefficients_vec)
365
366     mse_lst = []
367     for i in range(0, 100):
368         current_coefficients_vec = coefficients_vecs_lst[i]
369         current_prediction_vec = predict(test_mat.T, current_coefficients_vec)
370         mse_lst.append(mse(current_prediction_vec, test_response_vec))
371
372     mse_array = np.array(mse_lst)
373     percentages_array = np.array(range(1, 101))
374     plt.plot(percentages_array, mse_array)
375     plt.title('MSE values as a function of\nthe percentage taken from the '
376             'data set')
377     plt.xlabel('Percentage')
378     plt.ylabel('MSE Values')
379     plt.legend(['MSE Values'])
380     plt.show()
381
382
383 def feature_evaluation(path):
384     data_df = load_data(path)
385     response_vec = np.array(data_df['price'])
386     response_vec_std = np.std(response_vec)
387     non_categorical_features = ['lat', 'long', 'price', 'year_2014',
388                               'year_2015', 'yr_built', 'yr_renovated']
389     months_lst = ['month_' + str(i) for i in range(1, 13)]
390     features_to_drop = non_categorical_features + months_lst
391     data_df = data_df.drop(features_to_drop, axis=1)
392     features = data_df.columns.to_list()
393     for feature in features:
394         current_feature = np.array(data_df[feature])
395         current_feature_std = np.std(current_feature)
396         feature_mean = np.mean(current_feature)
397         response_mean = np.mean(response_vec)
398         feature_from_mean = current_feature - feature_mean
399         response_from_mean = response_vec - response_mean

```

```

400     mul = feature_from_mean * response_from_mean
401     cor_numerator = np.mean(mul)
402     cor_denominator = current_feature_std * response_vec_std
403     current_pearson_correlation = cor_numerator / cor_denominator
404
405     plt.scatter(current_feature, response_vec)
406     plt.title('house prices as a function of ' + feature + ' values\n\n'
407             + 'Pearson Correlation = ' + str(current_pearson_correlation))
408     plt.xlabel(feature + ' values')
409     plt.ylabel('house prices')
410     plt.legend(['house prices'])
411     plt.show()

```