

## Fitting A Linear Regression Model

1. ראשית נרצה לבנות את הפונקציה `loadData()` אשר תבצע תהליך קדם-עיבוד למידע שקיבלנו עבור הבתים. אציין באילו פיצ'רים השתמשתי, על איזו החלטתי לוותר ואיזה עיבוד בוצע.

- **שדות שלא נעשה בהם שימוש:** `id`, `Date` אלו שדות שנועדו ליצר סדר בתוך בסיס הנתונים שקיבלנו – לא רק שאינם קטגוריים הם אינם נותנים לנו כל מידע כלל על הנכס ולכן אינם רלוונטיים כלל למודל הלינארי שלנו.

כמו כן, החלטתי כמו כן להתעלם מעמודות **שקשורות למיקום** – בעולם האמיתי למקום הנכנס יש חשיבות גבוהה מאוד ולכן אלו שדות שנרצה מאוד שישפיעו על המודל שלנו, אך במקרה זה בו אין אנו יודעים את המיקום או הקשר הספציפי שלו לשווי הנכס אי הנכנסתו לשקלול לא עלולה להשפיע רבות על יכולת החיזוי של המודל שלנו (ניתן לדוגמא למצע את המיקום ולהסתכל על המרחק בין הנקודות, אך שוב אין זה נותן אינדיקציה).

שיקול דומה הפעלתי עבור שדה ה-`zipcode`, מהניתוח עולה כי ישנה רק דגימת אחת עם קוד 0, אך זו דגימה לא תקינה (20671) ולכן אתמודד עם מקרים אלו בעיבוד שנבצע בהמשך.

שדות `sqft_living15` ו-`sqft_lot15`, ממעבר על המידע שהתקבל, נראה כי שדות אלו רגישים לגודל המרתף או השטח העליון אם זה קיים. לעיתים מתקיים קשר ישיר ולפעמים יש תנודות במידע שלא ניתן להבין מהו הקשר. מכיוון שאנו משתמשים בכל שאר השדות, החלטתי לוותר על חישוב זה ולהתחשב במדידה המקורית, כאשר הגדלים האחרים שמסופקים לנו ישפיעו על המודל באופן זהה.

- **שדות שנלקחו כפי שהם:** מספר החדרים, קומות, מספר חדרי השירותים דירוגים כאלו ואחרים שמופיעים (כמו מצב, ציון ונוף) אלו שדות שכנראה יתארו יחסית בצורה טובה את הקשר הלינארי בניהם לבין ה-`response` שלנו – המחיר. לכן, לא נעבד אותם ונשתמש בהם כפי שקיבלנו במודל.

- **שדות שדרושים עיבוד:** נשים לב כי ישנם שדות, כגון שנת שיפוץ וגודל מרתף, אשר מתארים מאורע אם קרה ואם לא מופיע 0. נרצה להימנע ממקרה זה מכיוון שהנ"ל ייצר לנו רעש גבוהה בדגימות שלנו מה שישפיע על יכולת חיזוי המודל. לכן, נדרש לחשוב על כללי עיבוד שיעזרו לנו להחליט מה לעשות בכל שדה:

i. **שנת השיפוץ** – נייצר משתנה דמה `renovated` אשר יקבל ערכים  $\{0,1\}$ , כאשר 1 מצביע על כך שהיה שיפוץ ו-0 על כך שלא היה.

ii. **גודל המרתף** – כדי להקטין את הפיזור ניתן שוב להגדיר משתנה מסווג עם קיומו של מטרף, אך נוכל לחלק את הגודל ב-1000, ולהקטין משמעותית את הפזור בעת חישוב המודל כך נוכל בכל זאת להתחשב במקרים בהם יש מרתף ממש גדול. חשבתי לבצע מהלך זה מכיוון שלאחר שבצעתי ממוצע לעמודה, מתקבל כי הממוצע הינו 291.482 רגל מרובע, ולכן חלוקה ב-1000 תשמר אותנו יחסית קרובים ל-0 אך בכל זאת תייצר את ההבדל שהקטן שאנו רוצים עבור קיום המרתף. **נשים לב** כי זהו לא המצב בגודל השטח העליון ולכן לא נשנה שדה זה, כאשר ה-0 היחיד שהם הינו טעות במידע.

iii. **שנת בנייה** – נרצה להמיר את שנת הבנייה לגיל הנכס כאשר ככל שהנכס יותר

חדש כך נקבל ערך גבוהה יותר, קרי  $\frac{1}{2022 - age}$  תוך וידאו שאין אנו מחלקים מ-0

(לא בוצעה בדיקתי כי בדקתי במידע כי זה לא ייתכן).

2. נתבונן בשני פיצ'רים מתוך מטריצת הדגימות שלנו. נבחר את עמודות sqft\_above ו waterfront.

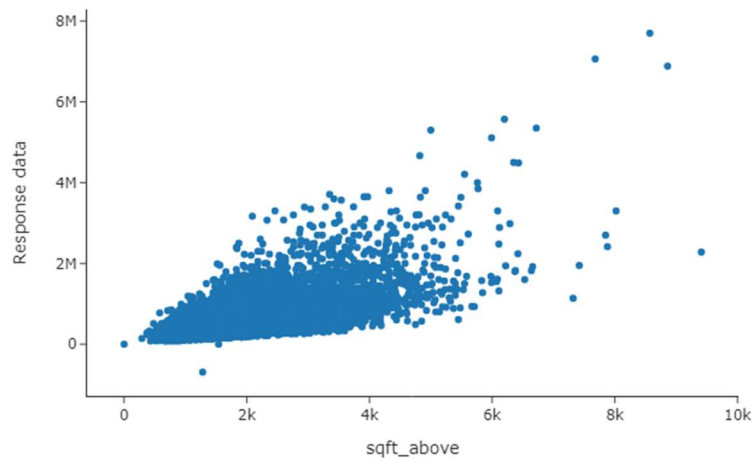
מכיוון שהמודל שלנו מחפש קשר ליניארי, ניתן לראות כי שדות בהם יש קשר לינארי כלשהו (כמו גודל מול מחיר) הינם שדות שיועילו יותר למודל שלנו מאשר שדות בינארים, כגון יש/אין מרתף.

מתאם פירסון שחישבנו נותן לנו את עוצמת הקשר הליניארי בין  $X, Y$  (חיובי, שלילי או לא קיים).

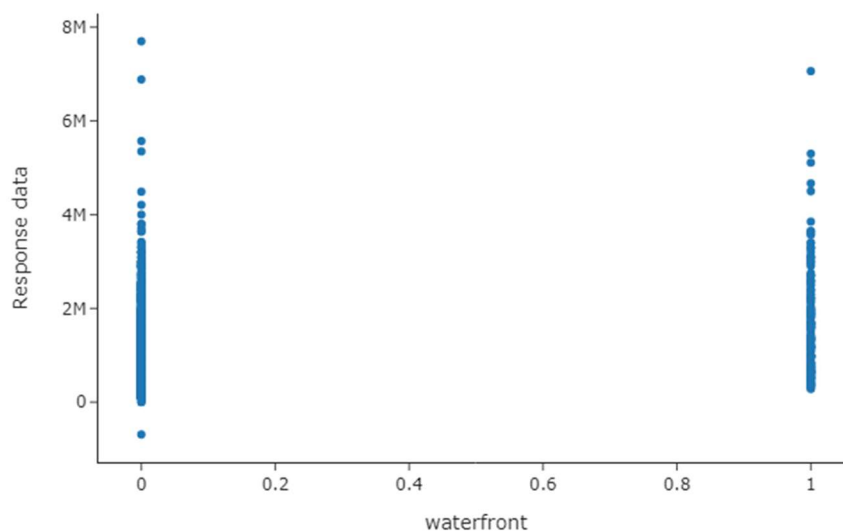
ניתן לראות זאת בברור שנסתכל על הפלטים שהוצאנו בחישוב של המתאם לכל פיצ'ר.

עבור הפיצ'רים שנבחרו מתקבל:

sqft\_above - Pearson Correlation is 0.6056177332802692



waterfront - Pearson Correlation is 0.26629523409982625

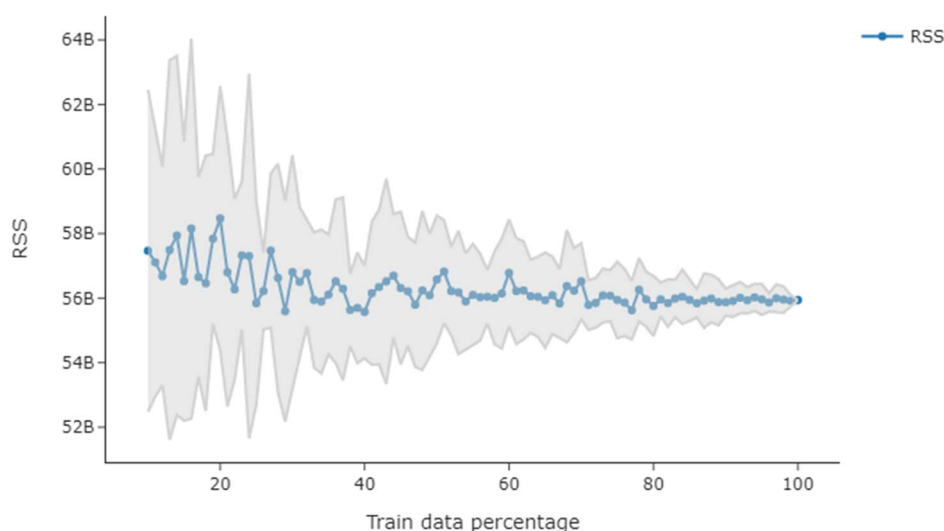


נשים לב כי ההערכה שלנו נכונה ובאמת המתאם עומד מול מה שציפינו, עבור פיצר גודל השטח העליון, נקבל כי המתאם גובהה יותר – זאת גם בהתאם למגמה הכללית שאנו רואים בגרף, ביחס מתאם נמוך עבור פיצר בינארי שלא מגלם כל קשר לינארי. לכן, ניתן לראות כי באמת שדה waterfront פחות יועיל למודל שלנו מכיוון שהוא לא תורם למציאת הקשר הלינארי, זאת ביחס לקשר שברור נוכח בשדה השטח.

3. הפונקציה המתאימה מומשה בקבצים

4. מצב גרף הפלט עבור שאלה זו.

RSS per train data percentage - Linear regression



נתבונן בגרף הפלט, באופן ישיר אנו רואים ככל שאנו מאמנים את המודל שלנו על יותר מידע, כך השגיאה מתכנסת וקטנה לדיוק המודל, שגיאה ריבועית של כ-56 מיליארד, שאילו נוציא שורש נקבל כי השגיאה  $MSE$  הינה  $\sim 230$  אלף דולר בהערכת מחיר בית עבור מידע שעוד לא נראה לפני.

מכיוון שלא הגדרנו מהי הטעות שאנו מאפשרים למודל להשיג, מבחינתו זה נתון נוסף שניתן ללמוד ממנו ובהמשך לשנות אולי את פרמטרי עיבוד המידע שלנו, או לנסות לשנות פיצרים נוספים ע"מ לקבל אומד טוב יותר, אם זה אפשרי – תוך ניסיון להימנע מלבצע overfit, כאשר כפי שציפינו שימוש בכל סט האימון הקטין את השגיאה וצמצם את מרחב הבחירה של המודל.

כמו כן, מבחינת confidence interval – ראינו בהרצאה כי זה נתון שמאפשר לנו להמחיש את השונות, כאשר אנו יכולים ללמוד מן הפלט כי תחום זה הולך ומצטמצם ככל שאחוז עולה, קרי מתקיים כי היכולת של המודל לחזות עבור  $\hat{y}$  ערכים גדולים ורחוקים יותר מצטמצם משמעותית, מה שמאפשר לנו לשפר את הדיוק, תחת סט האימון שקיבלנו, אך במחיר של גמישות במודל (אותו תשלום שראינו בין ה-bias ל-variance).

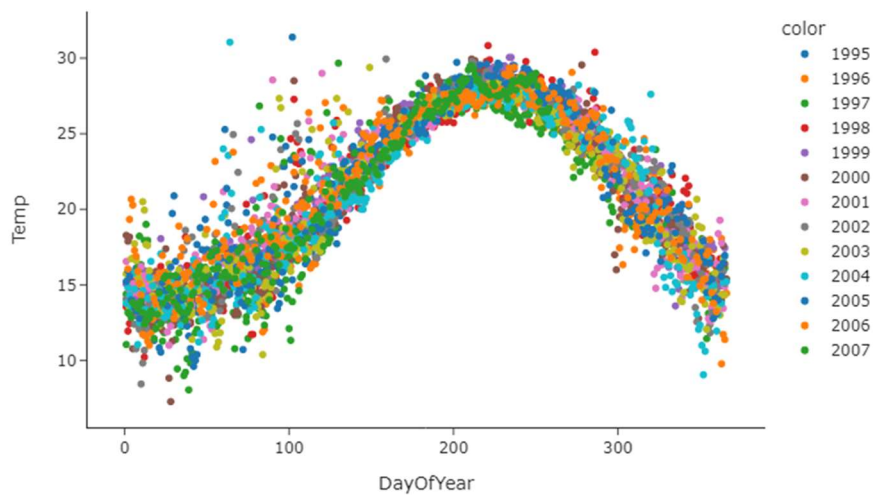
## Polynomial Fitting

1. מומש בקבצים.

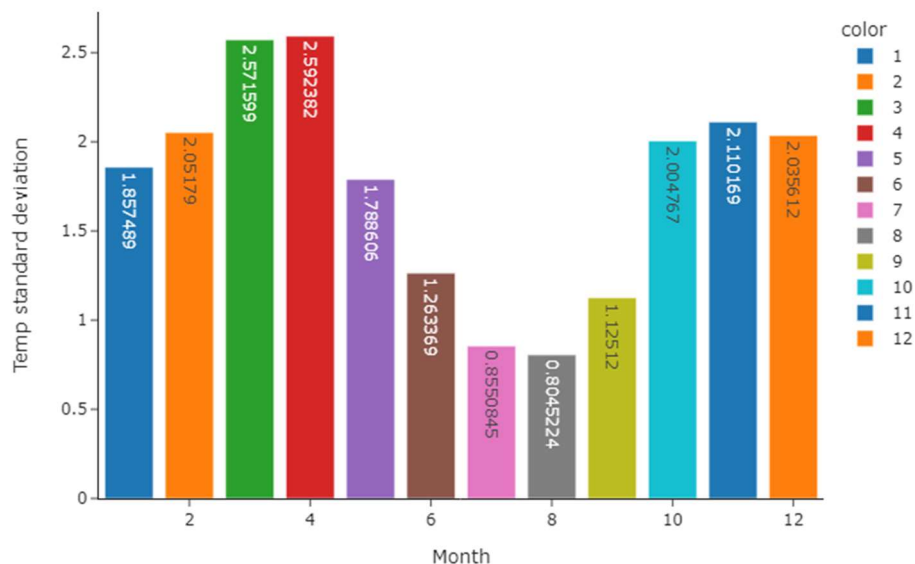
2.

a. מצייב גרף הפלט של הסעיף המבוקש. ניתן לראות כי הגרף מייצג מגמה כללית של פולינום מדרגה 3 או 4, על בסיס העיקול וירידה בגודל השיפוע כאשר הערכים מתרכבים ל-0.

Q2.1 - Temp to DayOfYear Israel



b. להלן גרף המייצג את סטיית התקן ביחס לחודשים כפי שהתבקשנו לעשות

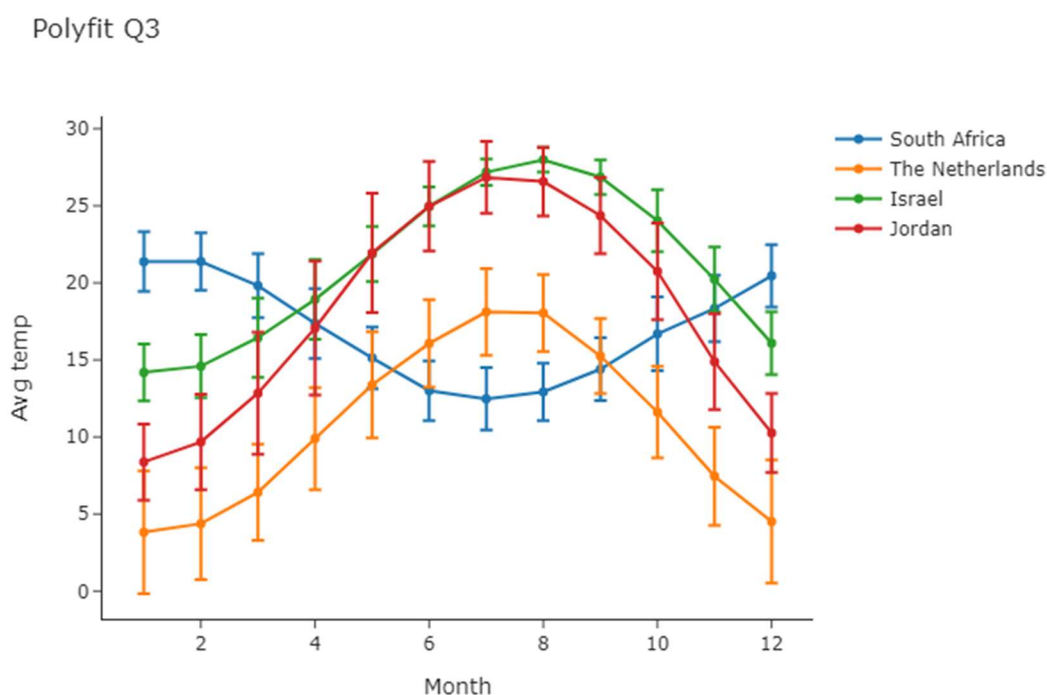


כפי שנתן לראות בגרף, אנו יכולים להעריך כי המודל שלנו לא יוכל לבצע פרדיקציה זהה על כל החודשים, מכיוון שסטיות התקן משתנות. ראינו כי סטיית התקן והשונויות הן המדד לנו להבין את פזור המידע ולכן כ"כ שאלו גבוהים יותר אנו למדים כי המידע שיקבלנו יותר

מפוזר ולכן יהיה יותר קשה לקבל פרדיקציה טובה עבור מודלים שמתבססים על רגרסיה ליניארית.

לכן, לפי תובנה זו ניתן לראות כי המודל שלנו יתקשה לבצע חיזוי מדויק בחודשי החורף ביחס לחודשי הקיץ – תובנה הגיונית עבור המעלות שבישראל, כאשר ישנן שנים עם שינויים גדולים בחורף ביחס לקיץ שהמעלות בו יחסית זהות.

3. להלן הפלט הנדרש:

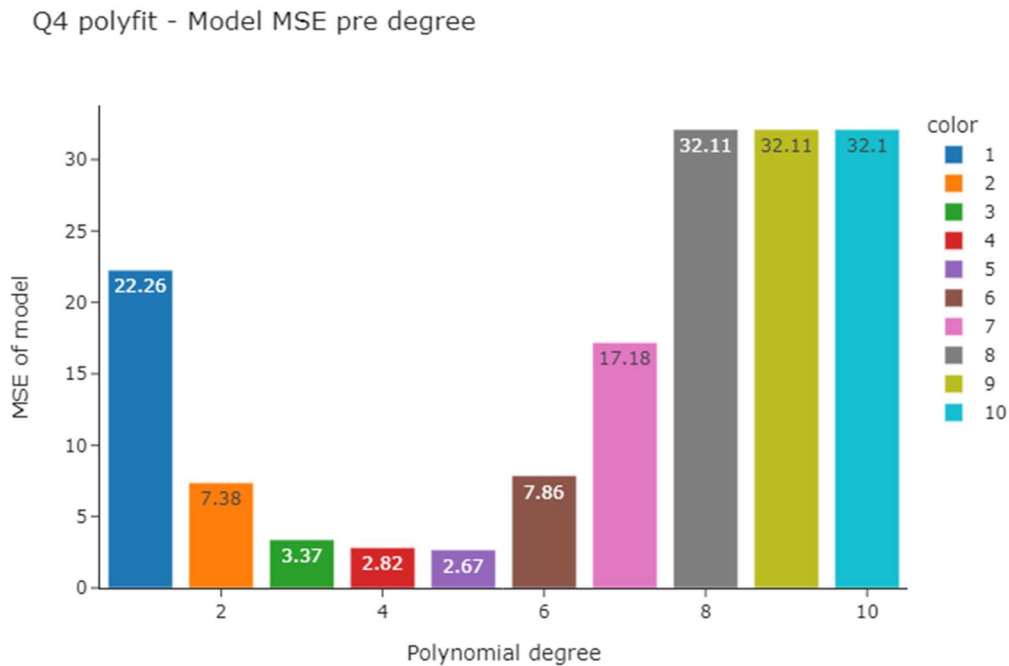


בהתבסס על גרף זה ניתן לראות כי לא לכל המדינות הדפוס זהה, לדוגמא דרום אפריקה הפוכה מישראל, ירדן או הולנד (דבר הגיוני שאנו רוצים לראות מכיוון שזו בצד השני של קו המשווה). לכן, ניתן לראות כי אילו נתאים מודל עבור ישראל בלבד, בהסתברות גבוהה מאוד שהמודל יספק פרדיקציה טובה מאוד עבור ירדן – ניתן לראות כי תוך התחשבות בסטיות התקן, קיימת חפיפה רבה בין ממוצעים ולכן הגיוני כי המודל יאפשר לקבל התאמה טובה גם במקרה זה.

מנגד, המודל לא יעבוד טוב כלל עבור דרום אפריקה מכיוון שסט האימון עבורו לא רלוונטי כלל. הנ"ל מתנהג בצורה שונה ועל כן לא נוכל להתבסס – שזו סה"כ תובנה שציפינו לה.

עבור הולנד נשים לב כי הפונקציה בצורתה דומה מאוד לפונקציה שמיוצרת ע"י הממוצעים של ישראל אך קיים פער בין התוצאות. לכן מודל של ישראל כנראה ישגה להשיג פרדיקציה טובה עבור הולנד, אך אילו נצליח לנרמל את הקבוע שמפריד בין הפונקציות נוכל לבצע התאמה ולקבל מודל שחווה בצורה טובה.

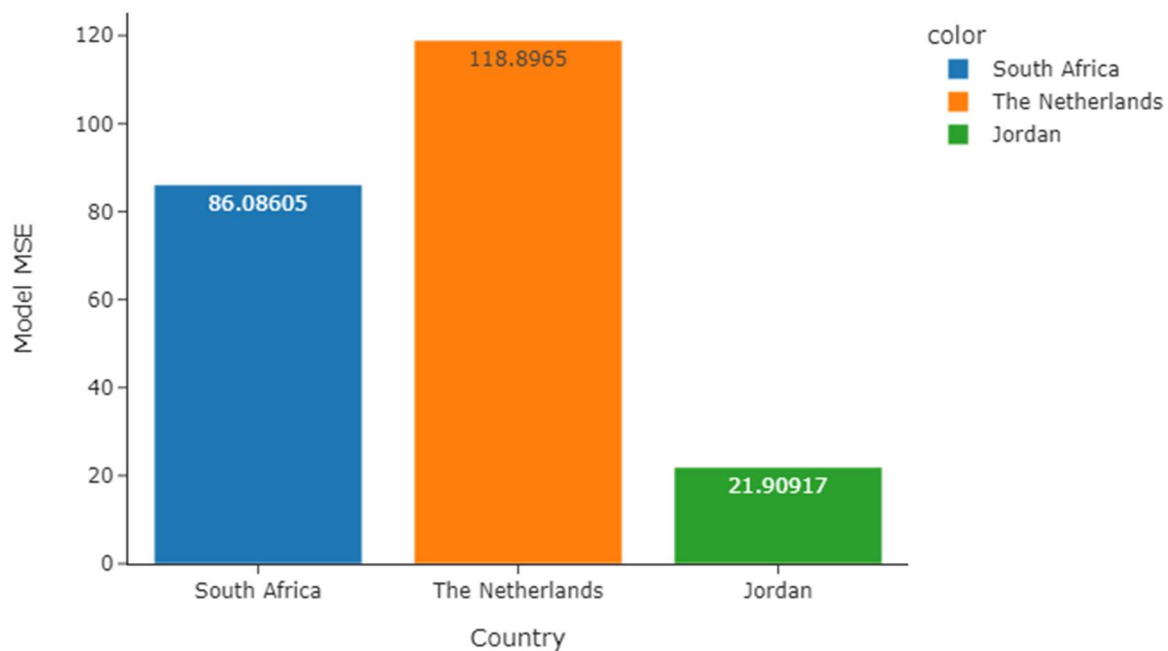
4. להלן גרף הפלט שמסכם את שגיאות המודל תחת כל דרגה



ניתן ראות כי פולינום מדרגה 5 השיג את השגיאה הקטנה ביותר, כאשר דרגה 4 מאוד קרובה בשגיאתה. פולינום מדרגה מאוד גבוהה כבר לא משיגים קירוב מתאים. נשים לב כי לא נרצה להמשיך ולהגדיל את דרגות הפולינום מכיוון שאנו נצבע overfit לסט האימון שלנו – ולכן נשאף להשתמש בפולינום מדרגה 4 או 5 כפונקציית הפרדיקציה שלנו עבור המקרה הנייל.

5. להלן פלט שגיאות המודל, כאשר התאמנו פולינום מדרגה 5 על כל המידע שהגיע מישראל

Q5 polyfit - Israel\_Model MSE for other countries, k=5



ניתן לראות כי מה שצפינו בשאלה 3 בהחלט מתקיים, המודל הצליח לחזות יחסית בצורה טובה את הטמפרטורה עבור ירדן, אך השגיאה הריבועית עלתה בצורה משמעותית כאשר ניסינו לחזות את עבור דרום אפריקה והולנד.

ניתן גם להסביר מדוע השגיאה מדרום אפריקה אפילו קטנה יותר מזו של הולנד למרות שראינו שהגרפים של ישראל והולנד היו דומים – שגיאה זו עלתה בהפרשי הטמפרטורה, אומנם עונות השנה הפוכות בין ישראל לדרום אפריקה, אך באופן ממוצע ראינו כי הטמפרטורות יחסית קרובות, ביחס להולנד. לכן המודל השיג תוצאות שבאופן כללי קרובות יותר למידע עליו הוא אומן – ולכן השגיאה במקרה זה קטנה יותר.

אילו היינו רוצים, ניתן היה להוסיף התערבות נוספת ולנרמל את תוצאת המודל במרחק הממוצע בין הטמפרטורות ולקבל חיזוי שככה"נ משמעותית טוב יותר עבור הולנד מן המודל הקיים.