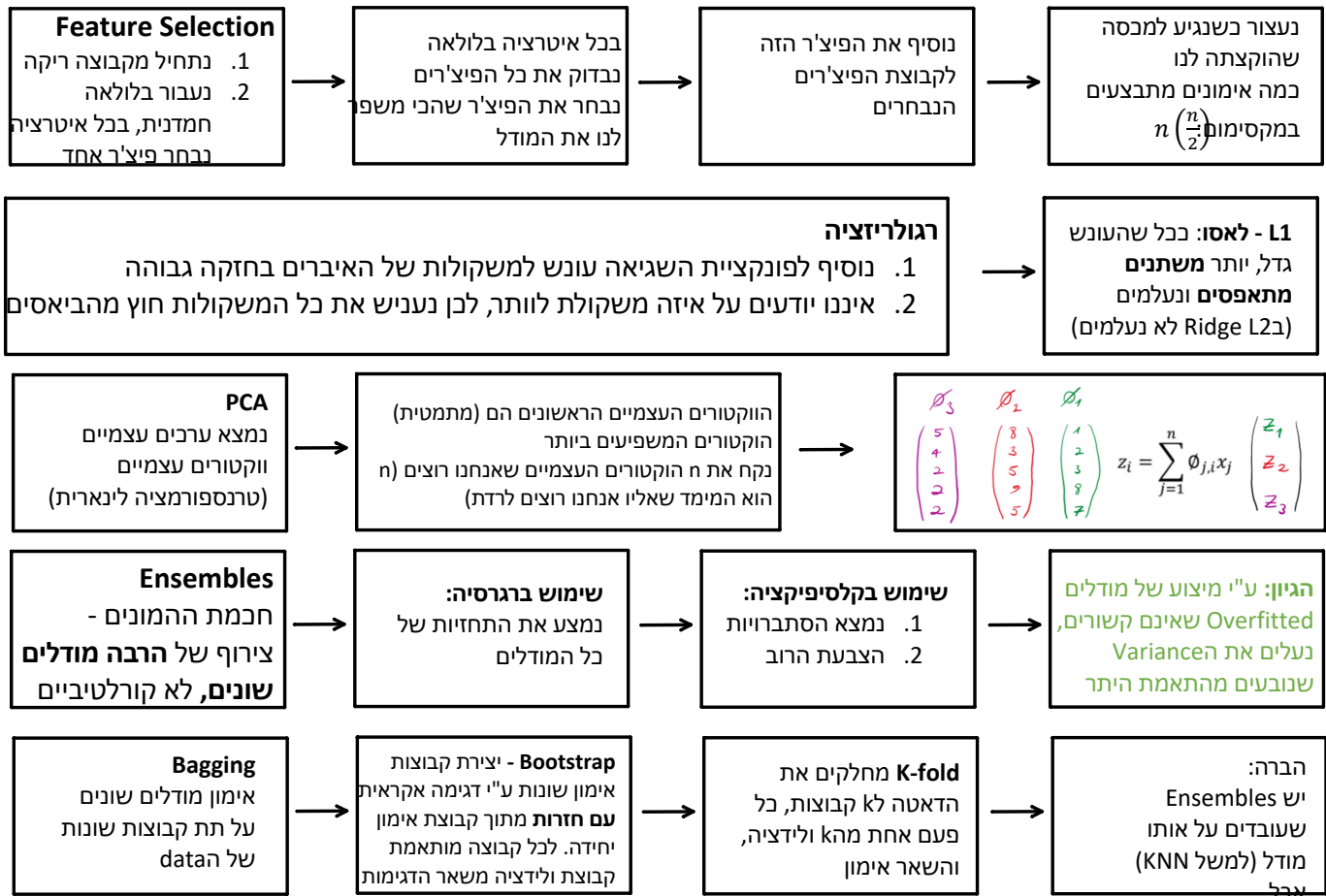


0:00:00

איך נפתור?



feature selection- Hybrid

1. Begin with Null model
2. $K=1\dots n$, until (stop criteria)
 - First add a feature as in Forward selection
 - Then, remove a feature that no longer provide improvement
 - Save M_k
3. Among the n models, Choose the best M_k using CV

שילוב 2 השיטות:
 התחל מקבוצת מאפיינים ריקה
 לכל $k=1\dots n$ בדוק מבין המאפיינים החסרים, איזה מאפיין
 להוסיף, הוסף את זה שמשפר הכי הרבה
 אח"כ, בדוק מבין המאפיינים שבקבוצה, איזה מאפיין כדאי
 להוריד (נוכחותו אינה הכרחית)
 בסוף התהליך בחר ע"י CV מבין k המודלים שניבנו

הבעיה בשיטה: עלולה להסיר feature שאינו משמעותי כשלעצמו, אבל בצירוף עם feature אחר הוא דווקא כן משמעותי להפחתת השגיאה

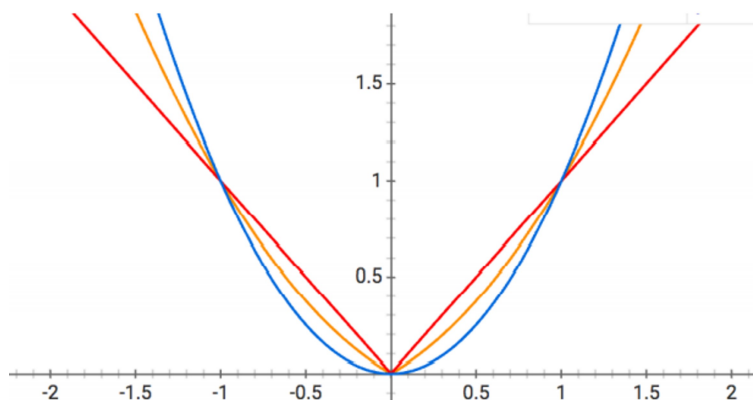
regularization

משתמשים בדר"ב לקלספיקציה

נוסיף לפונקציית ה-loss שלנו "עונש" למשקולות של האיברים בעלי חזקה גבוהה, נקבל פונקציית היפוטזה "חלקה" יותר.

נגדיר את γ (gamma) כקבוע הרגולריזציה, העונש למשקולות שערכה 1.

$$loss_{reg}(w) = loss(W) + \gamma ||W||_n$$



רגולריזציית Ridge – L2 – מרחק אוקלידי

$$\begin{aligned} ||W||_2 &= \frac{1}{2} \sqrt{\sum_{i=1}^n w_i^2} \\ w_i &= w_i - \lambda \frac{\partial loss}{\partial w_i} - (\gamma \cdot w_i) = \\ &= w_i(1 - \gamma) - \lambda \frac{\partial loss}{\partial w_i} \end{aligned}$$

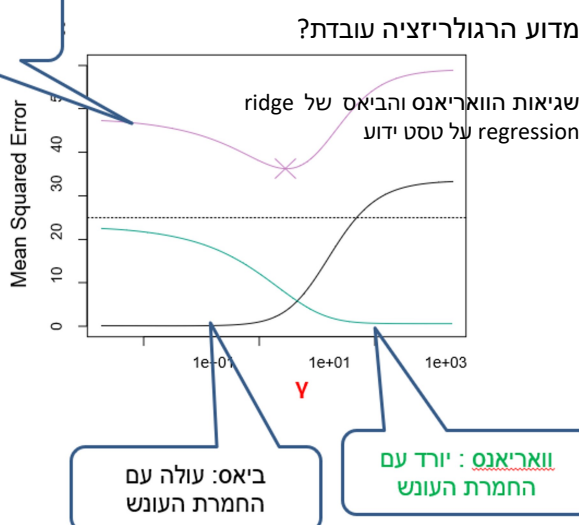
רגולריזציית Lasso – L1 – מרחק מנהטן

$$\begin{aligned} ||W||_1 &= \sum_{i=1}^n |w_i| \\ w_i &= w_i - \lambda \frac{\partial loss}{\partial w_i} - (\gamma \cdot \text{sign}(w_i)) \end{aligned}$$

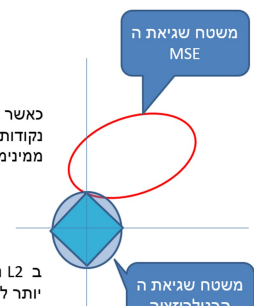
L1+L2 – Elastic Net

שילוב של Ridge ושל Lasso, מן המוצע משוקלל.
 במשקולות קטנות ה-Lasso ישפיע יותר.
 במשקולות גדולות ה-Ridge ישפיע יותר.

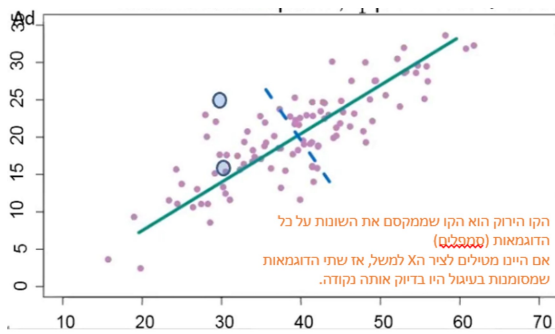
בגלל שהשגיאה תלויה ברבוע הביאס (פרבולה)
 העליה בשגיאת הביאס היא איטית בהתחלה ואח"כ
 מתדרדרת בקצב גדול יותר



כאשר הפינה של יהלום L1 פוגעת המשטח ה-MSE נקודות אחרות על ביהלום (כאלו בעלי אותו עונש) מתרחקות ממנימום ה-MSE עוד יותר.



PCA Principle Component Analysis



טרנספורמציות לינאריות ממימד n למימד $m < n$ המנסות לשמר כמה שיותר אינפורמציה באמצעות קומבינציות לינאריות של המאפיינים המקוריים

הוקטור העצמי הראשון שיבחר הוא הוקטור שממקסם את השונות של הנקודות (במקרה שלנו, הקו הירוק)

מכיוון שהוקטורים אורתוגונליים אחד לשני (ב90 מעלות \ בלתי תלויים לינארית)
הוקטור הבא חייב להיות הכחול (כי אנחנו ב2 ממדים, אין עוד דרגות חופש)

אינטואיציה:

1. אנחנו כל פעם מוצאים שמייצג את המידע בצורה הכי טובה (ימזער את ה MSE \ ימקסם את השונות)

2. אנחנו מטילים עליו את הנקודות

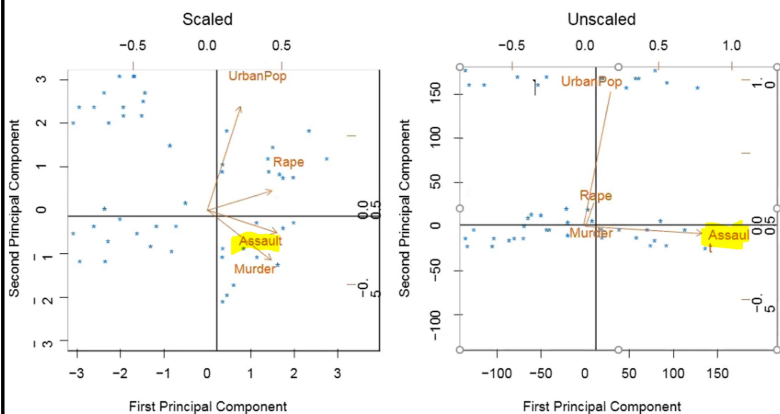
3. מתייחסים אליו ולנקודות עליו כציר חדש, ועוד פעם חוזרים ל1

הסתייגות:

אנחנו מוצאים את כל הוקטורים העצמיים בפעם אחת.

הפחתת מימדים

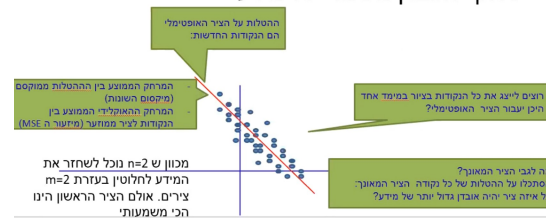
PCA נותן תוצאות שונות בתלות ב Scaling (אשר משפיע על השונות)



ל Assault יש שונות גדולה ולכן ישוקלל יותר מאחרים ב PC1

חלוקה של כל מאפיין בסטית התקן מונעת זאת לחילופין. אם היינו מודדים Assault פר אוכלוסיה של 100 אנשים, השונות הייתה קטנה משמעותית

בהינתן דוגמאות $D=\{x\}$ ממימד n
נסב דוגמאות אלו לנקודות במרחב ממימד $m < n$
נשאף לאובדן מינימלי של מידע



לדוגמה המרצה 5 ממדים ל3 ממדים:

$$z_i = \sum_{j=1}^n \phi_{j,i} x_j \begin{pmatrix} z_1 \\ z_2 \\ z_3 \end{pmatrix}$$

ϕ_3 ϕ_2 ϕ_1

$\begin{pmatrix} 5 \\ 4 \\ 2 \\ 2 \end{pmatrix}$ $\begin{pmatrix} 8 \\ 3 \\ 5 \\ 9 \end{pmatrix}$ $\begin{pmatrix} 1 \\ 2 \\ 3 \\ 8 \\ 7 \end{pmatrix}$

1. מספר הטרנספורמציות הן כגודל הממד שאליה אנחנו רוצים להגיע
2. כל טרנספורמציה כזו (של כל טרנספורמציה Φ_n היא וקטור מהמימד הישן של הספמלים n)

$$\Delta w_i = -\lambda \sum_{p \in D} \frac{\partial MSE_D(w)}{\partial y_{ip}} \frac{\partial y_{ip}}{\partial w_i}$$

$$= \frac{\lambda}{m} \sum_{p \in D} (t_p - y_p) \frac{\partial y_p}{\partial w_i}$$

$$= \frac{\lambda}{m} \sum_{p \in D} (t_p - y_p) x_{ip}$$

כלל עדכון המשקולות ב GD:

$$MSE_n(w) = \frac{1}{2m} \sum_{p \in D} (t_p - y_p)^2$$

$$y_p = h_w(x_p) = \sum_i w_i x_{ip} + b$$

$$\Delta w_i = -\lambda \frac{\partial loss_D(w)}{\partial w_i}$$