

**פרק 9: מימד, קשר ותבנית – יסודות ליניאריים
לאינטליגנציה מלאכותית**

**Chapter 9: Dimension, Correlation, and Pattern – Linear
Foundations for Artificial Intelligence**

ד"ר יורם סגל

© Dr. Segal Yoram - כל הזכויות שמורות

ספטמבר 2025

תוכן העניינים

4	1 מבוא: עולמו הנסתר של הווקטור
4	1.1 מבט ראשון: מה רואים כשמביטים בתמונה?
4	1.2 דקארט והמצאת השפה המספרית של המרחב
4	1.3 הווקטור: לא רק חץ במרחב
5	1.4 מתמונת חתול לווקטור: תהליך הוקטוריזציה
5	1.5 מכפלה סקלרית: מדד הדמיון האולטימטיבי
6	1.6 הוכחה: שתי ההגדרות זהות
7	1.7 קשר ל-IA: Cosine Similarity ומנועי חיפוש
8	1.8 Word Embeddings: המהפכה של Word2Vec
9	1.9 תרגיל תכנות עצמי 1.1 – חיפוש מילים דומות
10	1.10 אזהרה: הסכנה בווקטורים – Bias ודעות קדומות
11	1.11 מרחבים וקטוריים: הפורמליזם המתמטי
12	1.12 המעבר למימד גבוה: ברכה או קללה?
13	1.13 תרגיל תכנות עצמי 1.2 – השוואת מרחקים בממדים שונים
14	1.14 סיכום ומבט קדימה
16	2 הדיכוטומיה של המידע: קללת המימדיות
16	2.1 הפרדוקס: מתי "יותר" זה "פחות"?
16	2.2 בלמן והולדת המושג
17	2.3 הוכחה מתמטית: התרחקות מהמרכז
18	2.4 הוכחה: התכנסות המרחקים
19	2.5 השפעה על אלגוריתמים: K-Nearest Neighbors
20	2.6 יחס פיצ'ר-דגימה: כמה נתונים צריך?
21	2.7 הפתרון: רשתות CNN ומבנה מרחבי
21	2.8 פתרונות נוספים: הפחתת מימדיות
23	2.9 תרגיל תכנות עצמי 2.1 – הפחתת מימדיות עם PCA
25	2.10 רגולריזציה: מניעת Overfitting במימד גבוה
26	2.11 מחקר עדכני: אקסטרפולציה מול אינטרפולציה
27	2.12 סיכום ומבט קדימה
29	3 הערכת מודלים: מקדם הקביעה R^2
29	3.1 השאלה המרכזית: מה זה מודל "טוב"?
29	3.2 היסטוריה: מי המציא את R^2 ?
29	3.3 הגדרה: מהו R^2 ?
30	3.4 משמעות אינטואיטיבית: מהו "הסבר שונות"?
31	3.5 הוכחה: R^2 תמיד בין 0 ל-1
32	3.6 ייצוג גרפי: חלוקת השונות

32	קשר למקדם הקורלציה	3.7
32	סכנה 1: R^2 עולה תמיד כשמוסיפים תכונות	3.8
34	הפתרון: Adjusted R^2	3.9
36	סכנה 2: R^2 לא מתאים לסיווג	3.10
36	קשר לאלגברה ליניארית: הקרנות	3.11
36	תרגיל תכנות עצמי 3.1 – חישוב R^2 ידני	3.12
38	תרגיל תכנות עצמי 3.2 – השוואת R^2 ו- R^2_{detsujdA}	3.13
39	מקרי קצה: מתי R^2 מטעה?	3.14
40	סיכום ומבט קדימה	3.15

4 English References

42

1 מבוא: עולמו הנסתר של הווקטור

Introduction: The Hidden World of the Vector

בפרק זה נחקור את הבסיס הפילוסופי והמתמטי של הייצוג הווקטורי בבינה מלאכותית. נראה כיצד כל יחידת מידע – תמונה, מילה, או חולה במרפאה – הופכת לנקודה במרחב רב-ממדי, וכיצד מושג זה מהווה את היסוד של כל AI מודרני.

1.1 מבט ראשון: מה רואים כשמביטים בתמונה?

כשאתם מביטים בתמונה של חתול, מה המוח שלכם באמת רואה? האם הוא רואה יצור פרוותי חמוד עם שפם? או שמא הוא רואה משהו עמוק יותר – רשת מורכבת של 1000000 נקודות אור, כל אחת מהן נושאת מספר המייצג עוצמת בהירות? זוהי השאלה המרכזית שמניעה את כל תחום הבינה המלאכותית המודרני: **כיצד ניתן להפוך מידע לווקטור, ואיך ווקטורים הופכים למשמעות?**

1.2 דקארט והמצאת השפה המספרית של המרחב

ב-1637, כשרנה דקארט (René Descartes) פרסם את ספרו "La Géométrie", הוא לא יכול היה לדמיין שהרעיון הפשוט שלו – לייצג נקודה במרחב באמצעות מספרים – יהפוך ליסוד של מהפכה טכנולוגית שתגיע 400 שנה מאוחר יותר [1].

דקארט המציא את מערכת הצירים הקרטזית, ועמה את הרעיון המהפכני שכל נקודה במרחב ניתנת לתיאור מספרי. אך רק במאה ה-20, כשמדעני המחשב החלו לחשוב על דרכים לייצג מידע, הבינו שהרעיון של דקארט הוא לא רק כלי גיאומטרי – **הוא שפה אוניברסלית לתיאור המציאות.**

שאלה למחשבה: אם דקארט יכול היה לייצג נקודה במרחב תלת-ממדי עם 3 מספרים, כמה מספרים נדרשים כדי לייצג תמונת חתול?

1.3 הווקטור: לא רק חץ במרחב

בקורסים מתמטיים מסורתיים, וקטור מוגדר לרוב כ"בהרמב קח" – יש לו כיוון ואורך. אך בעולם מדעי הנתונים ולמידת המכונה, הווקטור הוא משהו עמוק יותר בהרבה: **הוא ייצוג של ישות במרחב מופשט [1].**

דוגמה מהחיים – חולה במרפאה:

אם נתון לכם חולה במרפאה, איך תתארו אותו כווקטור? התשובה: כל תכונה נמדדת של החולה – גיל, משקל, לחץ דם, רמת סוכר, טמפרטורה – הופכת ל- $t_{enopmoc}$ אחד בווקטור. אם יש לנו 50 תכונות נמדדות, החולה מיוצג כנקודה במרחב בעל 50 ממדים:

$$\vec{p}_{\text{patient}} = \begin{bmatrix} \text{age} \\ \text{weight} \\ \text{blood_pressure} \\ \text{glucose} \\ \vdots \\ \text{feature}_{50} \end{bmatrix} \in \mathbb{R}^{50}$$

הסימון \mathbb{R}^{50} משמעותו "מידמם 50 לש ידילקוא בחרמ" – מרחב שבו כל נקודה מוגדרת על ידי 50 מספרים ממשיים. זהו המרחב שבו כל החולים שלנו חיים, במובן מתמטי.

1.4 מתמונת חתול לווקטור: תהליך הוקטוריזציה

כפי שהדגיש ד"ר יורם סגל בהרצאתו, תמונה של 1000×1000 פיקסלים בצבע RGB מכילה 3×10^6 ערכים מספריים (כל פיקסל מיוצג על ידי 3 ערכים: אדום, ירוק, כחול).

המתמטיקה מאחורי הוקטוריזציה:

כדי להפוך אותה לייצוג וקטורי הניתן לניתוח על ידי מודל למידת מכונה, אנו מבצעים תהליך של "Flattening" (השטחה). המטריצה המבנית מפורקת לרצף של מספרים, היוצר וקטור באורך של שלושה מיליון.

אם $\mathbf{I} \in \mathbb{R}^{h \times w \times c}$ היא תמונה בגובה h , רוחב w ו- c ערוצי צבע, אזי הוקטוריזציה מוגדרת כ:

$$(1) \quad \text{vec}(\mathbf{I}) = \mathbf{v} \in \mathbb{R}^{h \cdot w \cdot c}$$

כאשר $\mathbf{v} = [I_{1,1,1}, I_{1,1,2}, \dots, I_{h,w,c}]^T$ – רצף של כל הפיקסלים בזה אחר זה. **משמעות עמוקה:** תמונת החתול הפכה לנקודה אחת ויחידה במרחב בעל שלושה מיליון ממדים. כל תמונה שונה במעט תהיה נקודה אחרת באותו מרחב עצום.

1.5 מכפלה סקלרית: מדד הדמיון האולטימטיבי

אחת הפעולות החשובות ביותר על וקטורים היא **מכפלה סקלרית** (Dot Product, Inner Product). היא מוגדרת כך:

הגדרה 1.1 – מכפלה סקלרית (אלגברית):

$$(2) \quad \langle \vec{u}, \vec{v} \rangle = \sum_{i=1}^n u_i v_i = u_1 v_1 + u_2 v_2 + \dots + u_n v_n$$

שאלה מהותית שהעלה המרצה: מה משמעות המכפלה הסקלרית? למה היא כל כך חשובה?

התשובה מגיעה מנוסחה חלופית, גיאומטרית, שגילה המתמטיקאי הגרמני Hermann von Helmholtz במאה ה-19:

$$(3) \quad \langle \vec{u}, \vec{v} \rangle = \|\vec{u}\| \cdot \|\vec{v}\| \cdot \cos(\theta)$$

כאשר θ היא הזווית בין שני הווקטורים, ו- $\|\vec{u}\|$ היא הנורמה (Norm) או אורך הווקטור:

$$(4) \quad \|\vec{u}\| = \sqrt{u_1^2 + u_2^2 + \dots + u_n^2} = \sqrt{\sum_{i=1}^n u_i^2}$$

1.6 הוכחה: שתי ההגדרות זהות

משפט 1.1: ההגדרה האלגברית (משוואה 1.2) וההגדרה הגיאומטרית (משוואה 1.3) של מכפלה סקלרית שוות.

הוכחה:

צעד 1: משפט הקוסינוסים.

במשולש עם צלעות $\|\vec{u}\|$, $\|\vec{v}\|$ וזווית θ ביניהן, הצלע השלישית היא $\|\vec{u} - \vec{v}\|$. לפי משפט הקוסינוסים:

$$\|\vec{u} - \vec{v}\|^2 = \|\vec{u}\|^2 + \|\vec{v}\|^2 - 2\|\vec{u}\|\|\vec{v}\|\cos(\theta)$$

צעד 2: פיתוח אלגברי של האגף השמאלי.

נפתח את $\|\vec{u} - \vec{v}\|^2$ באמצעות ההגדרה האלגברית:

$$\begin{aligned} \|\vec{u} - \vec{v}\|^2 &= \sum_{i=1}^n (u_i - v_i)^2 \\ &= \sum_{i=1}^n (u_i^2 - 2u_i v_i + v_i^2) \\ &= \sum_{i=1}^n u_i^2 + \sum_{i=1}^n v_i^2 - 2 \sum_{i=1}^n u_i v_i \\ &= \|\vec{u}\|^2 + \|\vec{v}\|^2 - 2\langle \vec{u}, \vec{v} \rangle \end{aligned} \quad (5)$$

צעד 3: השוואת הביטויים.

$$\text{מצעד 1: } \|\vec{u} - \vec{v}\|^2 = \|\vec{u}\|^2 + \|\vec{v}\|^2 - 2\|\vec{u}\|\|\vec{v}\|\cos(\theta)$$

$$\text{מצעד 2: } \|\vec{u} - \vec{v}\|^2 = \|\vec{u}\|^2 + \|\vec{v}\|^2 - 2\langle \vec{u}, \vec{v} \rangle$$

ביטול $\|\vec{u}\|^2 + \|\vec{v}\|^2$ משני האגפים מוביל ל:

$$2\langle \vec{u}, \vec{v} \rangle = 2\|\vec{u}\| \cdot \|\vec{v}\| \cdot \cos(\theta)$$

חלוקה ב-2 נותנת:

$$\langle \vec{u}, \vec{v} \rangle = \|\vec{u}\| \cdot \|\vec{v}\| \cdot \cos(\theta) \quad \blacksquare$$

משמעות עמוקה: מכפלה סקלרית היא מדד דמיון גיאומטרי.

- אם $\theta = 0^\circ$ (וקטורים באותו כיוון): $\cos(\theta) = 1$ והמכפלה מקסימלית \square דמיון מלא
- אם $\theta = 90^\circ$ (וקטורים אורתוגונליים): $\cos(\theta) = 0$ והמכפלה מתאפסת \square אין דמיון כלל
- אם $\theta = 180^\circ$ (וקטורים מנוגדים): $\cos(\theta) = -1$ והמכפלה מינימלית \square ניגוד מלא

1.7 קשר ל-Cosine Similarity ומנועי חיפוש

מנועי חיפוש מודרניים כמו Google ומערכות המלצה כמו Netflix משתמשים בווריאציה על מכפלה סקלרית שנקראת Cosine Similarity – דמיון קוסינוס:

$$(6) \quad \text{similarity}(\vec{u}, \vec{v}) = \frac{\langle \vec{u}, \vec{v} \rangle}{\|\vec{u}\| \cdot \|\vec{v}\|} = \cos(\theta)$$

חישוב זה מנרמל את המכפלה הסקלרית כך שהתוצאה תמיד בין -1 ל- $+1$, ללא תלות באורכי הווקטורים.

למה נורמליזציה חשובה?

בלי נורמליזציה, וקטור ארוך יקבל מכפלה סקלרית גדולה יותר רק בגלל האורך, לא בגלל הדמיון. Cosine Similarity פותר זאת על ידי מדידת הזווית בלבד.

דוגמה מהחיים – מנוע חיפוש:

כשאתם מקלידים שאילתה "machine learning tutorial", מנוע החיפוש:

1. ממיר את השאילתה לווקטור \vec{q} : כל מילה מקבלת משקל (למשל, באמצעות TF-IDF או Word2Vec [2])

2. משווה את \vec{q} לכל מסמך \vec{d}_i במאגר באמצעות Cosine Similarity

3. מחזיר את המסמכים עם הדמיון הגבוה ביותר (ציון קוסינוס הכי קרוב ל-1)

פונקציות NumPy למכפלה סקלרית ונורמה:

טבלה 1: פונקציות NumPy למכפלה סקלרית, נורמה ו-ytiralmiS enisoC

Function	תפקיד	שימוש והסבר
np.dot(u, v)	מכפלה סקלרית	מחשבת $\sum u_i v_i$ – הליבה של כל חישוב דמיון
np.linalg.norm(u)	נורמה (L2)	מחשבת $\sqrt{\sum u_i^2}$ – אורך הווקטור
@ (רוטרפוא)	מכפלת מטריצות	תחביר קצר: $u @ v$ זהה ל- $\text{np.dot}(v, u)$

הערה חשובה: פונקציות אלה הן הבסיס של כל אלגוריתם AI. נשתמש בהן שוב ושוב לאורך הספר.

פסאודו-קוד - Cosine Similarity

Cosine Similarity - מימוש ממוקד

```
import numpy as np

def cosine_similarity(u, v):
    """
    סימילאריות קוסין בין שני וקטורים.
    Cosine Similarity מתבטא כ:
    Args:
        u, v: וקטורים מסוג NumPy
    Returns:
        float: ערך בין -1 ל-1
    """
    # הנומ - תירלקסהלפכח
    dot_product = np.dot(u, v)

    # הנכח - תומרון
    norm_u = np.linalg.norm(u)
    norm_v = np.linalg.norm(v)

    # Cosine Similarity
    return dot_product / (norm_u * norm_v)

# (סימילאריות TF-IDF) וקטורים מסוג TF-IDF
doc1 = np.array([1, 2, 3, 0, 0])
doc2 = np.array([1, 1, 0, 4, 5])

print(f"Similarity: {cosine_similarity(doc1, doc2):.3f}")
# האצות: 0.277 - שלחן ומד
```

1.8 Word Embeddings: המהפכה של Word2Vec

אחת ההצלחות המרשימות ביותר של ייצוג וקטורי היא Word2Vec, שפותחה על ידי Tomas Mikolov ועמיתיו ב-elgooG ב-2013 [2].

הרעיון המהפכני: לייצג כל מילה בשפה כווקטור בן 300 ממדים, כך שמילים דומות במשמעות יהיו קרובות במרחב הווקטורי.

הפלא המתמטי של Word2Vec:

אחד הממצאים המפורסמים ביותר הוא שחישובים אריתמטיים על וקטורי מילים מניבים תוצאות משמעותיות סמנטית:

$$\vec{\text{king}} - \vec{\text{man}} + \vec{\text{woman}} \approx \vec{\text{queen}}$$

שאלה שמעלה פלא: איך זה אפשרי? איך חיבור וחיסור של מספרים יכול לבטא יחסים סמנטיים?

התשובה: הווקטורים לומדים יחסים. הווקטור $\vec{\text{king}} - \vec{\text{man}}$ מייצג את המושג המופשט "תירכו תוכלי", וכשמוסיפים $\vec{\text{woman}}$ מקבלים "תישנ תוכלי" = $\vec{\text{queen}}$.

זה עובד בגלל הנחת הדיסטריבוציה (Distributal Hypothesis) שניסח הבלשן Zellig Harris ב-1954 [3]:

"Words that occur in similar contexts tend to have similar meanings"

"מילים שמופיעות בהקשרים דומים נוטות לשאת משמעות דומה"

אלגוריתם Word2Vec לומד וקטורים שמנבאים את הקונטקסט של מילה, ובכך מקודדים משמעות סמנטית כמרחק גיאומטרי.

1.9 תרגיל תכנות עצמי 1.1 – חיפוש מילים דומות

מטרה: להבין איך Cosine Similarity משמש למציאת מילים דומות במודל Word2Vec. **משימה:**

1. טענו וקטורי Word2Vec מוכנים (למשל, מ-misne או מודל GloVe)

2. בחרו מילה (למשל, "computer")

3. חשבו Cosine Similarity בינה לבין כל המילים האחרות במאגר

4. הציגו את 10 המילים הדומות ביותר

פסאודו-קוד:

תוצאה צפויה:

- laptop: 0.82

- software: 0.78

- technology: 0.75

- hardware: 0.72

- PC: 0.70

```

from gensim.models import KeyedVectors

# (טנרטיאחמ) דירוהלןכומ Word2Vec לדומתניעט
model = KeyedVectors.load_word2vec_format(
    'GoogleNews-vectors-negative300.bin',
    binary=True
)

# Cosine Similarity - בתשמתשמהיצקנופה - תומודסילימתאיצח
similar_words = model.most_similar('computer', topn=10)

# הספדה
print("-'computer':")
for word, similarity in similar_words:
    print(f"{word}: {similarity:.3f}")

```

1.10 אזהרה: הסכנה בווקטורים – Bias ודעות קדומות

ב-2016, חוקרים מאוניברסיטת Princeton גילו תופעה מדאיגה: מודלי Word2Vec שאומנו על טקסטים מהאינטרנט הטמיעו **דעות קדומות חברתיות** [4].
דוגמאות להטיות שנמצאו:

man : programmer כמו woman : homemaker -

man : doctor כמו woman : nurse -

- קשרים חזקים יותר בין שמות אירופאים למילים חיוביות לעומת שמות אפרו-אמריקאים

המודל למד שגברים קשורים למקצועות טכנולוגיים, ונשים למקצועות ביתיים – לא משום שזו אמת, אלא כי כך התייחסו הטקסטים שממנו למד.
שאלה מוסרית שהעלה ד"ר סגל: אם AI לומד מהעולם, והעולם מוטה – האם AI מחזק את המוטיות?

התשובה המחקרית: כן, אלא אם כן מתערבים באופן פעיל. זו אחת המשימות המרכזיות של תחום **Fairness in AI**. חוקרים כמו Tolga Bolukbasi ועמיתיו פיתחו שיטות לניטרול הטיות (Debiasing) בווקטורי מילים [5].

שיטת הניטרול:

השיטה מבוססת על זיהוי **כיווני הטיה** במרחב הווקטורי. למשל, הכיוון $\vec{d}_{gender} = \vec{he} - \vec{she}$ מייצג את ממד המגדר. לאחר מכן, מסירים את המרכיב הזה ממילים שאמורות להיות ניטרליות:

$$\vec{\text{doctor}}_{\text{debiased}} = \vec{\text{doctor}} - \text{proj}_{\vec{d}_{\text{gender}}}(\vec{\text{doctor}})$$

כאשר $\text{proj}_{\vec{d}}(\vec{v})$ היא ההטלה של \vec{v} על הכיוון \vec{d} .

מסקנה: ייצוג וקטורי הוא כלי עוצמתי, אך הוא משקף את הנתונים שממנו למד. **אחריות החוקר היא להבטיח שהמודל לא מנציח אפליה.**

1.11 מרחבים וקטוריים: הפורמליזם המתמטי

עד כה דיברנו על וקטורים בצורה אינטואיטיבית. אך מה באמת הופך אוסף של מספרים לווקטור? התשובה מגיעה מתורת המרחבים הווקטוריים (Vector Space Theory), שפותחה במאה ה-19 על ידי מתמטיקאים כמו Hermann Grassmann (1844) ו-onaeP eppesuiG (1888).

1.2 הגדרה – מרחב וקטורי:

מרחב וקטורי V מעל שדה \mathbb{F} (בדרך כלל \mathbb{R} או \mathbb{C}) הוא קבוצה עם שתי פעולות:

1. **חיבור וקטורים:** $\vec{u}, \vec{v} \in V$ לכל $\vec{u} + \vec{v} \in V$

2. **כפל בסקלר:** $\alpha \vec{u} \in V$ לכל $\alpha \in \mathbb{F}$ ו- $\vec{u} \in V$

שמקיימות 8 אקסיומות:

טבלה 2: אקסיומות מרחב וקטורי

תכונה	אקסיומה
קומוטטיביות: $\vec{u} + \vec{v} = \vec{v} + \vec{u}$	1
אסוציאטיביות: $(\vec{u} + \vec{v}) + \vec{w} = \vec{u} + (\vec{v} + \vec{w})$	2
קיום איבר אפס: קיים $\vec{0}$ כך ש- $\vec{u} + \vec{0} = \vec{u}$	3
קיום איבר נגדי: לכל \vec{u} קיים $-\vec{u}$ כך ש- $\vec{u} + (-\vec{u}) = \vec{0}$	4
כפל באחד: $1 \cdot \vec{u} = \vec{u}$	5
דיסטריבוטיביות: $\alpha(\vec{u} + \vec{v}) = \alpha\vec{u} + \alpha\vec{v}$	6
דיסטריבוטיביות: $(\alpha + \beta)\vec{u} = \alpha\vec{u} + \beta\vec{u}$	7
אסוציאטיביות כפל: $(\alpha\beta)\vec{u} = \alpha(\beta\vec{u})$	8

למה זה חשוב ל-IA?

ברגע שהוכחנו שאוסף של ישויות (תמונות, מסמכים, חולים) יוצר מרחב וקטורי, אנחנו יכולים להשתמש בכל הכלים של אלגברה ליניארית:

- פירוק SVD (Singular Value Decomposition)

- ערכים עצמיים (Eigenvalues) ווקטורים עצמיים (Eigenvectors)

- הטלות (Projections) והטרנספורמציות ליניאריות

- PCA (Principal Component Analysis) להפחתת ממדיות

כל אלו הם כלים שנשתמש בהם לאורך הספר לניתוח ועיבוד נתונים.

1.12 המעבר למימד גבוה: ברכה או קללה?

ב-1961, הסטטיסטיקאי והמתמטיקאי Richard Bellman טבע את המונח המפורסם "Curse of Dimensionality" – קללת המימדיות [6].

הבעיה המתמטית:

Bellman הבין שככל שמוסיפים תכונות (ממדים) לנתונים, נפח המרחב גדל באופן אקספוננציאלי, והנתונים הופכים ספרסיים (Sparse) יותר ויותר.

דוגמה מספרית שהדגים המרצה:

נניח שאנחנו רוצים לדגום 10 נקודות לאורך כל ממד כדי לכסות את המרחב בצפיפות סבירה:

- ב-1 ממד: $10^1 = 10$ נקודות

- ב-2 ממדים: $10^2 = 100$ נקודות

- ב-3 ממדים: $10^3 = 1000$ נקודות

- ב-10 ממדים: $10^{10} = 10000000000$ נקודות (עשרה מיליארד!)

- ב-100 ממדים: 10^{100} – מספר גדול יותר מכמות האטומים ביקום!

שאלת המפתח שהעלה ד"ר סגל:

אם יש לנו תמונת חתול 1000×1000 פיקסלים (מיליון תכונות), כמה דוגמאות אימון נדרשות כדי לאמן מודל Brute Force?

התשובה המתמטית:

לפי כלל האצבע שנידון בהרצאה:

- למידת מכונה קלאסית: מינימום 30 דוגמאות לכל תכונה

- למידה עמוקה: מינימום 100 דוגמאות לכל תכונה

עבור מיליון תכונות: $100 \times 10^6 = 100000000$ דוגמאות (מאה מיליון!)

הפתרון – רשתות CNN:

זו הסיבה שרשתות Convolutional Neural Networks (CNN) חיוניות [7], [8]. הן מפחיתות באופן דרמטי את מספר הפרמטרים באמצעות:

1. **שיתוף משקלים** (Weight Sharing): אותו פילטר משמש בכל התמונה

2. **קישוריות מקומית** (Local Connectivity): כל נוירון מחובר רק לאזור קטן

3. **הירארכיה של תכונות:** למידה הדרגתית מתכונות פשוטות (קצוות) למורכבות (פנים)

הישג מפורסם: רשת AlexNet של Krizhevsky, Sutskever ו-notniH (2012) אימנה על 1.2 מיליון תמונות עם 60 ~ מיליון פרמטרים, והשיגה פריצת דרך בסיווג תמונות [8]. זו תהיה נקודת המוצא לפרק הבא, שבו נצלול לעומקה של קללת המימדיות ונראה כיצד היא משפיעה על כל אלגוריתם למידה.

1.13 תרגיל תכנות עצמי 1.2 – השוואת מרחקים בממדים שונים

מטרה: להמחיש את קללת המימדיות באופן מעשי.
משימה:

1. צרו 100 נקודות אקראיות במרחבים בממדים שונים: 2, 10, 100, 1000
2. חשבו את המרחק האוקלידי בין כל זוג נקודות
3. חשבו את ממוצע המרחקים ואת סטיית התקן
4. הציגו גרף: ציר $X =$ מימד, ציר $Y =$ יחס סטיית תקן/ממוצע

פסאודו-קוד:

תוצאה צפויה: היחס שואף ל-0 ככל שהממד גדל – כל הנקודות נעשות "באותו מרחק" זו מזו, ומדדי מרחק מאבדים משמעות.

1.14 סיכום ומבט קדימה

מה למדנו בפרק זה?

1. **וקטור הוא ייצוג מופשט** – לא רק חץ, אלא כל ישות הניתנת לתיאור מספרי (תמונה, מילה, חולה)
2. **מכפלה סקלרית היא מדד דמיון גיאומטרי** – הבסיס לחיפוש, המלצות, ו-PLN
3. **Cosine Similarity הוא הכלי המרכזי** – מנרמל מרחקים ומודד זווית בלבד
4. **Word2Vec הוא דוגמה מבריקה** – מילים הופכות לווקטורים והסמנטיקה הופכת לגיאומטריה
5. **הטיות ב-IA הן בעיה אמיתית** – מודלים לומדים מהעולם, כולל דעות קדומות
6. **מרחבים וקטוריים הם הפורמליזם** – מבטיחים שנוכל להשתמש באלגברה ליניארית
7. **ממד גבוה = אתגר עצום** – קללת המימדיות דורשת ארכיטקטורות חכמות כמו CNN

מבט קדימה – פרק 2:

בפרק הבא נצלול לעומק הדיכוטומיה של המידע – מדוע יותר תכונות לא תמיד טוב יותר, ואיך להתמודד עם מרחבים בעלי ממד גבוה. נראה:

```

import numpy as np
import matplotlib.pyplot as plt

def curse_of_dimensionality_demo(dims, n_points=100):
    """
    סיקורמטבושיחידים לעטות ידמימהתללקטאש יחממטטטט.
    """
    results = []

    for d in dims:
        # תויארקאטודוקנתריצי
        points = np.random.rand(n_points, d)

        # סיקורמהלכבושיח
        distances = []
        for i in range(n_points):
            for j in range(i+1, n_points):
                dist = np.linalg.norm(points[i] - points[j])
                distances.append(dist)

        # תוקיטטיטטט
        mean_dist = np.mean(distances)
        std_dist = np.std(distances)
        ratio = std_dist / mean_dist # -0, לפאושהזשכ

        results.append(ratio)
        print(f"Dim={d}: Mean={mean_dist:.3f}, Std={std_dist:.3f}, Ratio={ratio:.3f}")

    return results

# הצרה
dimensions = [2, 5, 10, 50, 100, 500, 1000]
ratios = curse_of_dimensionality_demo(dimensions)

# פרג
plt.plot(dimensions, ratios, marker='o')
plt.xlabel('Number of Dimensions')
plt.ylabel('Std/Mean Ratio')
plt.title('Curse of Dimensionality')
plt.grid(True)
plt.show()

```

- הוכחה מתמטית מלאה של קללת המימדיות
- השפעה על אלגוריתמים: K-Nearest Neighbors, SVM, Decision Trees
- יחס פיצ'ר-דגימה: כמה נתונים באמת צריך?
- פתרונות: Feature Selection, PCA, Regularization

שאלת מחשבה לסיום:

אם תמונת חתול 1000×1000 היא נקודה במרחב מיליון-ממדי, וכמעט כל "הזזה" קטנה היא אקסטרפולציה – איך מודלים גנרטיביים כמו DALL-E מצליחים ליצור תמונות חדשות שנראות הגיוניות?
התשובה תחכה לפרק 5 על קונוולוציה ורשתות CNN.

מטלות וקריאה מורחבת

- תרגיל 1.1:** מצאו מילים דומות באמצעות Word2Vec (ראו פסאודו-קוד).
תרגיל 1.2: המחישו את קללת המימדיות (ראו פסאודו-קוד).
תרגיל 1.3: חשבו ידנית מכפלה סקלרית, נורמה, ו- ytiralimiS enisoC בין:

$$\vec{u} = [1, 2, 3], \quad \vec{v} = [4, 5, 6]$$

$$\text{פתרון: } \langle \vec{u}, \vec{v} \rangle = 32, \|\vec{u}\| = \sqrt{14}, \|\vec{v}\| = \sqrt{77}, \text{ דמיון } = 0.975$$

קריאה מורחבת:

- [1] Linear Algebra and Learning from Data, פרקים 1-2: יסודות אלגברה ליניארית
- [2] – המאמר המקורי על Word2Vec: "Efficient Estimation of Word Representations in Vector Space"
- [4] – מחקר על הטיות במודלי שפה: "Semantics Derived Automatically from Language Corpora Contain Human-like Biases"
- [9] Deep Learning, פרק 5: למידת מכונה בסיסית

שאלות להעמקה:

1. מדוע Cosine Similarity עדיף על מרחק אוקלידי במרחבים רב-ממדיים?
2. האם ניתן להסיר לחלוטין הטיות ממודלי Word2Vec? מה המחיר?
3. איך רשתות CNN "מרמות" את קללת המימדיות?

סיום פרק 1

2 הדיכוטומיה של המידע: קללת המימדיות

The Information Dichotomy: Curse of Dimensionality

בפרק זה נחקור את אחד האתגרים המרכזיים של למידת מכונה: הפרדוקס שבו הוספת מידע (תכונות) עלולה להחליש את המודל במקום לחזק אותו. נבחן את היסודות המתמטיים של קללת המימדיות, נוכיח את השפעותיה, ונראה כיצד היא משפיעה על אלגוריתמים שונים.

2.1 הפרדוקס: מתי "יותר" זה "פחות"?

דמיינו רופא שמנסה לאבחן מחלה. בתחילה, יש לו שלוש תכונות: טמפרטורה, דופק, ולחץ דם. המודל שלו עובד היטב. לפתע, הוא מקבל גישה למעבדה מתקדמת שמודדת 1000 תכונות ביוכימיות. אינטואיטיבית, המידע הנוסף אמור לשפר את האבחנה. אך במציאות, דיוק המודל יורד.

למה זה קורה?

התשובה טמונה בתופעה מתמטית מפתיעה שגילה Richard Bellman ב-1957 [6]: ככל שמספר הממדים גדל, המרחב "מתנפח" באופן אקספוננציאלי, והנתונים הופכים דלילים (Sparse) – כמו כוכבים ביקום המתרחב.

2.2 בלמן והולדת המושג

Richard Ernest Bellman (1920–1984) היה מתמטיקאי אמריקאי שתרם תרומות מהפכניות לתחומי תכנות דינמי, תורת הבקרה, ולמידת מכונה. בספרו "Dynamic Programming" (1957), הוא זיהה בעיה מהותית: **מורכבות חישובית גדלה באופן אקספוננציאלי עם מספר המשתנים.**

התובנה המרכזית של בלמן:

במרחב חד-ממדי, אם רוצים לכסות קטע באורך 1 ברשת של נקודות במרווח ϵ , נדרשות $\sim 1/\epsilon$ נקודות. אך במרחב d -ממדי, נדרשות $\sim (1/\epsilon)^d$ נקודות – גידול אקספוננציאלי [6].

דוגמה מספרית:

טבלה 3: מספר הנקודות הנדרש לכיסוי מרחב $[0,1]$ בצפיפות $\epsilon = 0.1$

סדר גודל	נקודות נדרשות	ממד
עשרות	$10^1 = 10$	1
מאות	$10^2 = 100$	2
אלפים	$10^3 = 1000$	3
עשרה מיליארד	10^{10}	10
גוגול (יותר מאטומי היקום)	10^{100}	100
בלתי נתפס	$10^{1000000}$	1000000

מסקנה: תמונת חתול 1000×1000 פיקסלים (מיליון ממדים) דורשת מספר דגימות שהוא מעבר ליכולת חישובית של כל מחשבי העולם ביחד.

2.3 הוכחה מתמטית: התרחקות מהמרכז

נוכיח תופעה מפתיעה: במרחבים רב-ממדיים, כמעט כל הנפח מרוכז בקליפה החיצונית, רחוק מהמרכז.

משפט 2.1 – ריכוז הנפח בקליפה:

יהי כדור יחידה $B_d = \{\vec{x} \in \mathbb{R}^d : \|\vec{x}\| \leq 1\}$ בממד d . נפח הכדור הוא:

$$(7) \quad V_d(r) = \frac{\pi^{d/2}}{\Gamma(d/2 + 1)} r^d$$

כאשר Γ היא פונקציית הגמא, ו- r הוא הרדיוס.

הוכחה – יחס הנפחים:

נחשב את יחס הנפח של הכדור הפנימי $B_d(0.9)$ (רדיוס 0.9) לכדור המלא $B_d(1)$:

$$\begin{aligned} \frac{V_d(0.9)}{V_d(1)} &= \frac{(0.9)^d \cdot \pi^{d/2} / \Gamma(d/2 + 1)}{1^d \cdot \pi^{d/2} / \Gamma(d/2 + 1)} \\ &= (0.9)^d \end{aligned} \quad (8)$$

נחשב עבור ממדים שונים:

טבלה 4: יחס הנפח הפנימי (90% מהרדיוס) לנפח המלא

ממד d	$(0.9)^d$	אחוז מהנפח
2	0.81	81%
3	0.729	73%
10	0.349	35%
100	2.66×10^{-5}	0.003%
1000	$\sim 10^{-46}$	כמעט אפס

מסקנה מדהימה: ב-100 ממדים, רק 0.003% מהנפח נמצא ב-90% המרכזיים! כמעט כל הנפח מרוכז בקליפה הדקה בין $r = 0.9$ ל- $r = 1$.

משמעות ל-IA:

אם דגימות האימון שלנו הן "נקודות במרכז", המודל שלנו יצטרך לעשות אקסטרפולציה עצומה כדי לחזות במרחב האמיתי (הקליפה החיצונית) שבו נמצאות רוב ה"תמונות האפשריות".

2.4 הוכחה: התכנסות המרחקים

תופעה נוספת של קללת המימדיות: במרחבים רב-ממדיים, כל המרחקים בין נקודות הופכים להיות דומים.

משפט 2.2 – התכנסות יחס המרחקים:

יהיו $\{\vec{x}_i\}_{i=1}^n$ נקודות בלתי-תלויות במרחב \mathbb{R}^d . יהיו d_{\max} ו- d_{\min} המרחקים המקסימלי והמינימלי מנקודת מבחן \vec{q} . אזי:

$$(9) \quad \lim_{d \rightarrow \infty} \frac{d_{\max} - d_{\min}}{d_{\min}} \rightarrow 0$$

הוכחה (סקיצה):

צעד 1: מרחק אוקלידי במימד d בין \vec{q} ל- \vec{x}_i :

$$d_i = \|\vec{q} - \vec{x}_i\| = \sqrt{\sum_{j=1}^d (q_j - x_{ij})^2}$$

צעד 2: אם המרכיבים $(q_j - x_{ij})$ הם משתנים אקראיים בלתי-תלויים עם תוחלת μ ושונות σ^2 , אזי לפי חוק המספרים הגדולים:

$$\frac{1}{d} \sum_{j=1}^d (q_j - x_{ij})^2 \xrightarrow{P} \mu^2 + \sigma^2$$

צעד 3: לכן:

$$d_i = \sqrt{d \cdot (\mu^2 + \sigma^2 + o(1))} = \sqrt{d} \cdot \sqrt{\mu^2 + \sigma^2} \cdot (1 + o(1))$$

צעד 4: כל המרחקים d_i גדלים כ- \sqrt{d} , אך ההבדלים ביניהם גדלים לאט יותר (רק כ- $\sqrt{d \cdot \text{Var}}$). לכן:

$$\frac{d_{\max} - d_{\min}}{d_{\min}} \sim \frac{O(\sqrt{d})}{O(\sqrt{d})} \cdot \frac{\text{Var}}{\text{Mean}^2} \rightarrow 0 \quad \text{כ-} d \rightarrow \infty \text{ רשאכ}$$

■

משמעות מעשית:

במרחב בעל 1000 ממדים, אם הנקודה הקרובה ביותר נמצאת במרחק 100, הנקודה הרחוקה ביותר נמצאת במרחק $100.01 \sim$ – כמעט אותו דבר. מדדי מרחק מאבדים יכולת הבחנה [10], [11].

2.5 השפעה על אלגוריתמים: K-Nearest Neighbors

אלגוריתם K-Nearest Neighbors (KNN), שפותח על ידי Evelyn Fix ו-segdoH hpesoJ ב-1951 [12], הוא אחד האלגוריתמים הפשוטים והאינטואיטיביים ביותר בלמידת מכונה.

הרעיון: כדי לסווג נקודה חדשה \vec{q} , מצא את k הנקודות הקרובות ביותר במרחב האימון, והצבע לפי רוב.

למה KNN קורס במימד גבוה?

מכיוון שכל המרחקים הופכים דומים (משפט 2.2), אין משמעות ל"קרוב ביותר". הנקודה ה-1 הקרובה ביותר והנקודה ה-1000 הקרובה ביותר נמצאות כמעט באותו מרחק!

ניסוי מספרי - סימולציה:

תוצאה צפויה: הדיוק יורד מ-95% (במימד 2) ל-55% (במימד 1000) - כמעט אקראי.

2.6 יחס פיצ'ר-דגימה: כמה נתונים צריך?

כפי שהדגיש ד"ר יורם סגל בהרצאתו, קיימים כללי אצבע לקביעת מספר הדגימות הנדרש:

כלל 1 - למידת מכונה קלאסית:

$$(10) \quad n_{\text{samples}} \geq 30 \times n_{\text{features}}$$

כלל 2 - למידה עמוקה:

$$(11) \quad n_{\text{samples}} \geq 100 \times n_{\text{features}}$$

הצדקה מתמטית:

כללים אלו נובעים מתורת הלמידה הסטטיסטית. לפי **גבול VC (Vapnik-Chervonenkis)** [13], שגיאת ההכללה של מודל עם VC-dimension d מוגבלת ע"י:

$$(12) \quad \epsilon \leq \sqrt{\frac{d \log(n/d) + \log(1/\delta)}{n}}$$

כאשר n הוא מספר הדגימות, δ רמת הביטחון, ו- ϵ שגיאת ההכללה.

פירוש: כדי לקבל שגיאה קטנה ϵ , נדרש $n \sim O(d/\epsilon^2)$ - מספר הדגימות צריך לגדול ליניארית עם המימד.

דוגמה מהחיים - תמונת חתול:

תמונה 1000×1000 פיקסלים BGR:

- מספר תכונות: $1000 \times 1000 \times 3 = 3000000$

- דגימות נדרשות (כלל 2): $100 \times 3000000 = 300000000$ (שלוש מאות מיליון תמונות!)

למה זה בלתי אפשרי?

- ImageNet, אחד ממאגרי התמונות הגדולים בעולם, מכיל רק ~ 1.2 מיליון תמונות אימון [14]

```

import numpy as np
from sklearn.neighbors import KNeighborsClassifier
from sklearn.datasets import make_classification
from sklearn.model_selection import train_test_split

def knn_curse_demo(dims, n_samples=1000):
    """
    מידממה הרפסמלש הציקנופכ KNN יעוציב ון ווב.
    """
    results = []

    for d in dims:
        # ייטניס ינותנתיצי
        X, y = make_classification(
            n_samples=n_samples,
            n_features=d,
            n_informative=min(d, 10), # קר 10 תויביתמרופניאתונוכת
            n_redundant=0,
            random_state=42
        )

        X_train, X_test, y_train, y_test = train_test_split(
            X, y, test_size=0.3, random_state=42
        )

        # וניא KNN
        knn = KNeighborsClassifier(n_neighbors=5)
        knn.fit(X_train, y_train)

        # קויד
        accuracy = knn.score(X_test, y_test)
        results.append(accuracy)

    print(f"Dim={d}: Accuracy={accuracy:.3f}")

    return results

# הצרה
dimensions = [2, 5, 10, 50, 100, 500, 1000]
accuracies = knn_curse_demo(dimensions)

```

- גם Google לא צילמה 300 מיליון תמונות חתולים

אז איך Deep Learning מצליח?
התשובה: שיתוף משקלים והנחות אינדוקטיביות [7], [9].

2.7 הפתרון: רשתות CNN ומבנה מרחבי

רשתות Convolutional Neural Networks (CNN), שפותחו על ידי Yann LeCun ב-1989 [15], מתגברות על קללת המימדיות באמצעות שלושה עקרונות:

עיקרון 1 – קישוריות מקומית (Local Connectivity):

במקום לחבר כל פיקסל לכל נוירון (שיצטרך $10^{12} = 1000^2 \times 1000^2$ משקלים), כל נוירון מחובר רק לאזור קטן (למשל, 3×3 פיקסלים).

עיקרון 2 – שיתוף משקלים (Weight Sharing):

אותו פילטר (קרנל) משמש בכל מקום בתמונה. במקום 10^{12} משקלים, נדרשים רק $576 = 3 \times 3 \times 64$ משקלים לשכבה!

עיקרון 3 – הירארכיה של תכונות:

השכבות הראשונות לומדות תכונות פשוטות (קצוות, צבעים), והשכבות העמוקות לומדות תכונות מורכבות (עיניים, אוזניים, פנים).

השוואת מספר הפרמטרים:

טבלה 5: השוואת מספר פרמטרים: Fully Connected מול CNN

דגימות נדרשות (כלל 2)	מספר פרמטרים	ארכיטקטורה
10^{14} (בלתי אפשרי)	$\sim 10^{12}$	Fully Connected
6×10^9	$\sim 60 \times 10^6$	CNN (AlexNet)
2.5×10^9	$\sim 25 \times 10^6$	CNN (ResNet-50)

הצלחה מעשית:

רשת AlexNet [8] אומנה על 1.2 מיליון תמונות בלבד (הרבה פחות מהנדרש תיאורטית), והשיגה פריצת דרך בתחרות ImageNet ב-2012, עם הפחתה של 10% בשגיאה לעומת השיטות הקודמות.

2.8 פתרונות נוספים: הפחתת ממדיות

מלבד CNN, קיימות שיטות נוספות להתמודדות עם קללת המימדיות:

שיטה 1 – Feature Selection (בחירת תכונות):

בחירת תת-קבוצה של התכונות החשובות ביותר. שיטות נפוצות:

- **Filter Methods:** מיון לפי קורלציה, mutual information, chi-square

- **Wrapper Methods:** בחירה לפי ביצועי המודל (forward selection, backward elimination)

- **Embedded Methods**: רגולריזציה L1 (Lasso) שמאלצת משקלים לאפס

שיטה 2 – PCA (Principal Component Analysis)

המצאה של Karl Pearson (1901) [16] ו-H. Hotelling (1933) [17]. PCA מוצא את הכיוונים בעלי השונות המקסימלית ומקרין עליהם.

רעיון PCA:

אם $\mathbf{X} \in \mathbb{R}^{n \times d}$ היא מטריצת הנתונים, PCA מוצא את הווקטורים העצמיים של מטריצת הקוריאנס $\mathbf{C} = \frac{1}{n} \mathbf{X}^T \mathbf{X}$.

הווקטורים העצמיים $\{\vec{v}_1, \dots, \vec{v}_d\}$ עם הערכים העצמיים הגדולים ביותר $\{\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_d\}$ מהווים כיוונים של שונות מקסימלית.

הקרנה על k מרכיבים עיקריים:

$$(13) \quad \mathbf{Z} = \mathbf{XV}_k$$

כאשר $\mathbf{V}_k = [\vec{v}_1, \dots, \vec{v}_k] \in \mathbb{R}^{n \times k}$ ו- $(k \ll d)$. הוא הייצוג המופחת.

כמה מרכיבים לבחור?

כלל אצבע: בחר k כך שהשונות המוסברת היא לפחות 95%:

$$(14) \quad \frac{\sum_{i=1}^k \lambda_i}{\sum_{i=1}^d \lambda_i} \geq 0.95$$

פסאודו-קוד PCA – yPmuN:

הערה על SVD: Singular Value Decomposition הוא שיטה יעילה יותר לחישוב PCA מאשר חישוב ישיר של ערכים עצמיים, במיוחד כאשר $n \ll d$ [18].

2.9 תרגיל תכנות עצמי 2.1 – הפחתת ממדיות עם PCA

מטרה: להמחיש את יכולת PCA לשמר מידע תוך הפחתת ממדים.
משימה:

1. טענו את מערך הנתונים MNIST (תמונות ספרות כתובות ביד, 28×28 פיקסלים = 784 תכונות)

2. הפחיתו ל-50 מרכיבים עיקריים באמצעות PCA

3. שחזרו תמונות מהייצוג המופחת

4. השוו את התמונות המקוריות למשוחזרות

5. חשבו את אחוז השונות המוסברת

פסאודו-קוד מורחב:

תוצאה צפויה: עם 50 מרכיבים (במקום 784), ניתן לשמר ~85% מהמידע, והתמונות המשוחזרות יהיו קריאות למדי.

```

import numpy as np

def pca_manual(X, n_components=2):
    """
    SVD תועצמאב ינדי PCA עצובת.

    Args:
        X: (n_samples, n_features) סינות צירטת.
        n_components: סיירקיע סיביכרמ רפטת.

    Returns:
        X_reduced: (n_samples, n_components) סייתח פומ סינות.
        components: סימצעה סיירוטקווה.
        explained_var_ratio: תרבסומה תונושה הרועיש.
    """
    # (עצומח) רוסיח סינותנה זוכיר
    X_centered = X - np.mean(X, axis=0)

    # SVD - סיירלוגניסיכרעלקוריפ
    # U: eigenvectors לש  $XX^T$ 
    # S: (לש  $eigenvalues$ ) שרוש סיירלוגניסיכרע
    # Vt: eigenvectors לש  $X^TX$  סיביכרמ
    U, S, Vt = np.linalg.svd(X_centered, full_matrices=False)

    # סיירקיע סיביכרמ k תריחב
    components = Vt[:n_components]

    # הנרקה
    X_reduced = X_centered @ components.T

    # תרבסומתונושבושיח
    explained_variance = (S**2) / (len(X) - 1)
    explained_var_ratio = explained_variance[:n_components] / np.sum(
        explained_variance)

    print(f"Explained variance ratio: {explained_var_ratio}")
    print(f"Total: {np.sum(explained_var_ratio):.3f}")

    return X_reduced, components, explained_var_ratio

# המגוד
X = np.random.rand(1000, 100) # תומיגד 1000, תונוכת 100
X_reduced, components, var_ratio = pca_manual(X, n_components=10)
print(f"Original shape: {X.shape}, Reduced shape: {X_reduced.shape}")

```

```

from sklearn.datasets import fetch_openml
from sklearn.decomposition import PCA
import matplotlib.pyplot as plt
import numpy as np

# תניעט MNIST
mnist = fetch_openml('mnist_784', version=1, parser='auto')
X = mnist.data.to_numpy()[ :1000] # תונושאר תומיגד 1000
y = mnist.target.to_numpy()[ :1000]

# בושחלומרנ (!)
X = X / 255.0

print(f"Original_shape: {X.shape}") # (1000, 784)

# סיביכרמ -50 להתחפהל PCA
pca = PCA(n_components=50)
X_reduced = pca.fit_transform(X)

print(f"Reduced_shape: {X_reduced.shape}") # (1000, 50)
print(f"Explained_variance: {np.sum(pca.explained_variance_ratio_):.3f}"
)

# תונומתרוזחש
X_reconstructed = pca.inverse_transform(X_reduced)

# הגצה
fig, axes = plt.subplots(2, 5, figsize=(12, 5))
for i in range(5):
    # תירוקמהנומת
    axes[0, i].imshow(X[i].reshape(28, 28), cmap='gray')
    axes[0, i].set_title(f'Original: {y[i]}')
    axes[0, i].axis('off')

    # תרוזחושמהנומת
    axes[1, i].imshow(X_reconstructed[i].reshape(28, 28), cmap='gray')
    axes[1, i].set_title(f'Reconstructed')
    axes[1, i].axis('off')

plt.tight_layout()
plt.show()

```


2.10 רגולריזציה: מניעת Overfitting במימד גבוה

כאשר מספר התכונות גדול ביחס למספר הדגימות, המודל נוטה להתאמת יתר (Overfitting) – הוא "לומד בעל-פה" את נתוני האימון ונכשל בנתונים חדשים.

רגולריזציה (Regularization) היא טכניקה שמוסיפה "קנס" על מורכבות המודל, ובכך מעודדת פשטות.

שתי שיטות עיקריות:

רגולריזציה L2 (Ridge):

$$(15) \quad \mathcal{L}_{\text{Ridge}} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^d w_j^2$$

קנס על גודל המשקלים – מעודד משקלים קטנים, אך לא מאפס אותם.

רגולריזציה L1 (Lasso):

$$(16) \quad \mathcal{L}_{\text{Lasso}} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 + \lambda \sum_{j=1}^d |w_j|$$

קנס על סכום ערכי המשקלים – מאלץ חלק מהמשקלים להיות בדיוק אפס, ובכך מבצע בחירת תכונות אוטומטית.

פרמטר הקנס λ :

- $\lambda = 0$: אין רגולריזציה (סיכון ל-gnittifrevO)

- λ גדול מאוד: המודל פשוט מדי (Underfitting)

- λ אופטימלי: נקבע באמצעות Cross-Validation

מי המציא?

Ridge: Arthur Hoerl ו-dranneK treboR (1970) [19]

Lasso: Robert Tibshirani (1996) [20]

פסאודו-קוד – השוואת Ridge ו-Lasso:

תוצאה צפויה: Lasso יאפס ~ 490 מתוך 500 המשקלים, וישמור רק את התכונות האינפורמטיביות.

2.11 מחקר עדכני: אקסטרפולציה מול אינטרפולציה

מחקר פורץ דרך של Anadi Chaman ו-cinamkoD navI (2021) [21] גילה תוצאה מדהימה: במרחבים עם יותר מ-100 ממדים, דגימה חדשה כמעט אף פעם אינה אינטרפולציה – היא כמעט תמיד אקסטרפולציה.

הגדרות:

- **אינטרפולציה:** דגימה חדשה נמצאת בתוך המעטפת הקמורה (Convex Hull) של נתוני האימון

```

from sklearn.linear_model import Ridge, Lasso
from sklearn.datasets import make_regression
from sklearn.model_selection import train_test_split
import numpy as np

# (!תוידמימה) תללק  $d < n$  : מייטתניסנינותנתריצי
X, y = make_regression(
    n_samples=100,      # דבלבתומיגד 100
    n_features=500,     # תונוכת 500!
    n_informative=10,   # תויטנוולר 10 קר
    noise=10,
    random_state=42
)

X_train, X_test, y_train, y_test = train_test_split(
    X, y, test_size=0.3, random_state=42
)

# Ridge
ridge = Ridge(alpha=1.0) #  $\alpha = \lambda$ 
ridge.fit(X_train, y_train)
ridge_score = ridge.score(X_test, y_test)

# Lasso
lasso = Lasso(alpha=1.0)
lasso.fit(X_train, y_train)
lasso_score = lasso.score(X_test, y_test)

print(f"Ridge  $R^2$ : {ridge_score:.3f}")
print(f"Lasso  $R^2$ : {lasso_score:.3f}")

# ספא Lasso מילקשמהחכ
zero_weights = np.sum(np.abs(lasso.coef_) < 1e-5)
print(f"Lasso zeroed out {zero_weights}/{len(lasso.coef_)} features")

```

- **אקסטרפולציה:** דגימה חדשה נמצאת מחוץ למעטפת הקמורה

המעטפת הקמורה (Convex Hull):

$$(17) \quad \text{Conv}(\mathcal{X}) = \left\{ \sum_{i=1}^n \alpha_i \vec{x}_i \mid \vec{x}_i \in \mathcal{X}, \alpha_i \geq 0, \sum_{i=1}^n \alpha_i = 1 \right\}$$

התוצאה המפתיעה:

בניסויים על מערכי נתונים אמיתיים (MNIST, CIFAR-10), Chaman ו-cinamko מצאו ש-99.9% מדגימות הבדיקה הן אקסטרפולציה במרחב המקורי!

משמעות ל-IA:

כל פעם שמודל Deep Learning מסווג תמונה חדשה, הוא למעשה **משער מעבר לנתונים שראה**. ההצלחה של מודלים מודרניים נובעת מההנחות **האינדוקטיביות** שלהם (כמו מבנה ה-NNC) ולא מכיסוי מלא של המרחב.

שאלת מחשבה שהעלה המרצה:

"אם כל תחזית היא אקסטרפולציה, האם AI באמת 'מבין' את העולם, או שהוא רק מנחש בצורה חכמה?"
זוהי אחת השאלות הפילוסופיות המרכזיות של AI מודרני.

2.12 סיכום ומבט קדימה

מה למדנו בפרק זה?

1. **קללת המימדיות היא אמיתית ומוכחת** – נפח המרחב גדל אקספוננציאלית, והנתונים הופכים דלילים
2. **מרחקים מתכנסים** – במימד גבוה, כל הנקודות נעשות "באותו מרחק", ומדדים כמו KNN קורסים
3. **יחס פיצ'ר-דגימה הוא קריטי** – נדרשות לפחות 30-100 דגימות לכל תכונה
4. **רשתות CNN מתגברות על הקללה** – שיתוף משקלים וקישוריות מקומית מפחיתים דרמטית את מספר הפרמטרים
5. **הפחתת ממדיות עובדת** – PCA, Feature Selection, Regularization מאפשרים עבודה יעילה
6. **רוב התחזיות הן אקסטרפולציה** – AI מצליח בזכות הנחות אינדוקטיביות, לא כיסוי מלא

מבט קדימה – פרק 3:

בפרק הבא נעבור מהבעיות לפתרונות. נחקור את **מקדם הקביעה** R^2 – הכלי המרכזי להערכת מודלי רגרסיה. נראה:

- הגדרה מתמטית מלאה של R^2
- הוכחה: למה R^2 תמיד בין 0 ל-1?
- סכנות: מתי R^2 מטעה?
- Adjusted R^2 – הגרסה המתוקנת
- קשר לקורלציה ולמכפלה סקלרית

שאלת מחשבה לסיום:

אם מודל רגרסיה ליניארית השיג $R^2 = 0.95$ על נתוני האימון, אך רק $R^2 = 0.60$ על נתוני הבדיקה – מה קרה? האם זו קללת המימדיות, Overfitting, או משהו אחר?

מטלות וקריאה מורחבת

- תרגיל 2.1:** הפחיתו את MNIST עם PCA (ראו פסאודו-קוד).
- תרגיל 2.2:** השוו Ridge ו-Lasso על נתונים עם $n < d$ (ראו פסאודו-קוד).
- תרגיל 2.3:** הוכיחו בעצמכם את משפט 2.1 (ריכוז הנפח בקליפה) עבור $d = 10$.
- רמז:** חשבו $V_d(r) = \frac{\pi^{d/2}}{\Gamma(d/2+1)} r^d$ עבור $r = 0.9$ ו- $r = 1$.
- קריאה מורחבת:**

- [6] – Dynamic Programming: המקור המקורי של "קללת המימדיות"
- [10] – "When Is 'Nearest Neighbor' Meaningful?": ניתוח מתמטי של קריסת KNN
- [21] – "Truly Generative or Just Extrapolation?": המחקר על אינטרפולציה מול אקסטרפולציה
- [9] – Deep Learning, פרק 7: רגולריזציה והפחתת Overfitting
- [22] – The Elements of Statistical Learning, פרק 3: רגרסיה ליניארית ורגולריזציה

שאלות להעמקה:

1. מדוע Lasso מאפס משקלים אך Ridge לא? רמז: חשבו על הגיאומטריה של הקנס (L_1 מול L_2)
2. האם קללת המימדיות משפיעה גם על רשתות Transformer ב-PLN?
3. אם 99.9% מהתחזיות הן אקסטרפולציה, איך מסבירים את ההצלחה של GPT-4?

סיום פרק 2

3 הערכת מודלים: מקדם הקביעה R^2

Model Evaluation: The Coefficient of Determination R^2

בפרק זה נחקור את אחד הכלים המרכזיים להערכת מודלי רגרסיה: **מקדם הקביעה R^2** (R-squared). נבין את המשמעות המתמטית שלו, נוכיח מדוע הוא תמיד בין 0 ל-1, ונראה מתי הוא עלול להטעות אותנו.

3.1 השאלה המרכזית: מה זה מודל "טוב"?

דמיינו שבניתם מודל רגרסיה ליניארית לחיזוי מחיר דירה על סמך שטח, מיקום וגיל. המודל נותן תחזיות, אך **כיצד נדע אם הוא טוב?** שאלה זו מעסיקה סטטיסטיקאים ומדעני נתונים מאז ראשית המדע. התשובה המודרנית, שהתגבשה במאה ה-20, היא **מקדם הקביעה** – מספר יחיד שמסכם את "טיב ההתאמה" של המודל לנתונים.

3.2 היסטוריה: מי המציא את R^2 ?

המושג פותח בהדרגה על ידי מספר סטטיסטיקאים:
Francis Galton (1822–1911) היה הראשון לחקור קורלציה ורגרסיה. בעבודתו "Re- gression towards Mediocrity in Hereditary Stature" (1886) [23], הוא גילה שילדים של הורים גבוהים נוטים להיות נמוכים יותר מהוריהם (רגרסיה לממוצע).
Karl Pearson (1857–1936), תלמידו של Galton, פיתח את **מקדם הקורלציה r** (1896) [24], שהוא השורש של R^2 במקרה של רגרסיה פשוטה.
Ronald Fisher (1890–1962) הרחיב את המושג ל**רגרסיה מרובה** (Multiple Regression) וניתח את חלוקת השונות [25].
Sewall Wright (1889–1988) טבע את הסימון R^2 ב-1921 [26].

3.3 הגדרה: מהו R^2 ?

הגדרה 3.1 – מקדם הקביעה:

R^2 הוא מדד סטטיסטי המצביע על **שיעור השונות במשתנה התלוי Y שמוסברת על ידי המשתנים הבלתי תלויים X** במודל רגרסיה.

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}} \quad (18)$$

כאשר:

SS_{res} – **סכום ריבועי השאריות** (Residual Sum of Squares)

SS_{tot} – **סכום הריבועים הכולל** (Total Sum of Squares)

פירוט המשתנים:

$$SS_{\text{res}} = \sum_{i=1}^n (y_i - \hat{y}_i)^2 \quad (19)$$

$$SS_{\text{tot}} = \sum_{i=1}^n (y_i - \bar{y})^2 \quad (20)$$

כאשר:

y_i - הערך האמיתי של התצפית ה- i

\hat{y}_i - הערך החזוי על ידי המודל עבור התצפית ה- i

$\bar{y} = \frac{1}{n} \sum_{i=1}^n y_i$ - ממוצע הערכים האמיתיים

3.4 משמעות אינטואיטיבית: מהו "הסבר שונות"?

כדי להבין את R^2 , נתחיל מ**מודל הבסיס** - המודל הפשוט ביותר האפשרי.

מודל הבסיס (ledoM enilesaB): תחזה תמיד את הממוצע \bar{y} .

זהו המודל הטיפש ביותר - הוא מתעלם מכל התכונות X ופשוט אומר "התשובה היא

הממוצע". שגיאת מודל זה היא SS_{tot} - **השונות הכוללת** בנתונים.

המודל שלנו: משתמש ב- X כדי לחזות \hat{y}_i .

שגיאת המודל שלנו היא SS_{res} - **השונות שנותרה** אחרי שהמודל עשה את עבודתו.

הפרשנות:

$$R^2 = 1 - \frac{\text{שגיאת המודל שלנו}}{\text{שגיאת מודל הבסיס}} = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}} \quad (21)$$

- אם $R^2 = 0$: המודל שלנו לא טוב יותר מהממוצע (חסר תועלת!)

- אם $R^2 = 1$: המודל מושלם - אין שגיאות כלל

- אם $R^2 = 0.8$: המודל מסביר 80% מהשונות

דוגמה מספרית:

נניח שאנחנו מנבאים מחירי דירות. הממוצע הוא 500000 ש"ח.

טבלה 6: דוגמה: חישוב R^2 למחירי דירות

תצפית	y_i (אמיתי)	\hat{y}_i (חזוי)	$(y_i - \bar{y})^2$	$(y_i - \hat{y}_i)^2$
1	600000	580000	10000000000	400000000
2	450000	470000	2500000000	400000000
3	700000	680000	40000000000	400000000
סה"כ			$SS_{\text{tot}} = 52500000000$	$SS_{\text{res}} = 1200000000$

חישוב:

$$R^2 = 1 - \frac{1200000000}{52500000000} = 1 - 0.023 = 0.977$$

פרשנות: המודל מסביר 97.7% מהשונות במחירים – מודל מצוין!

3.5 הוכחה: R^2 תמיד בין 0 ל-1

משפט 3.1 – תחום R^2 :

עבור רגרסיה ליניארית עם איבר חופשי, מתקיים: $0 \leq R^2 \leq 1$.

הוכחה:

חלק א': $R^2 \leq 1$

נוכיח ש- $SS_{\text{res}} \leq SS_{\text{tot}}$.

צעד 1: מודל הרגרסיה הליניארית ממזער את SS_{res} על ידי פתרון:

$$\min_{\vec{w}} \sum_{i=1}^n (y_i - \vec{w}^T \vec{x}_i)^2$$

צעד 2: מודל הבסיס (חיזוי הממוצע) הוא מקרה פרטי של רגרסיה ליניארית עם $\vec{w} = \vec{0}$ (למעט איבר חופשי \bar{y}).

צעד 3: מכיוון שהרגרסיה ממזערת, היא לא יכולה להיות גרועה יותר ממודל הבסיס:

$$SS_{\text{res}} \leq SS_{\text{tot}}$$

צעד 4: לכן:

$$\frac{SS_{\text{res}}}{SS_{\text{tot}}} \leq 1 \Rightarrow 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}} \leq 1 \Rightarrow R^2 \leq 1 \quad \blacksquare$$

חלק ב': $R^2 \geq 0$

זה נכון רק לרגרסיה ליניארית עם איבר חופשי. במקרים מיוחדים (כמו רגרסיה דרך הראשית), R^2 יכול להיות שלילי!

עבור רגרסיה רגילה, המודל תמיד יכול לחזות את הממוצע (במקרה הגרוע ביותר), כך ש- $SS_{\text{res}} \leq SS_{\text{tot}}$ ולכן $R^2 \geq 0$.

■

3.6 ייצוג גרפי: חלוקת השונות

ניתן לחשוב על השונות הכוללת כ"פאי" שמתחלק לשניים:

$$(22) \quad \underbrace{SS_{\text{tot}}}_{\text{שונות כוללת}} = \underbrace{SS_{\text{reg}}}_{\text{שונות מוסברת}} + \underbrace{SS_{\text{res}}}_{\text{שונות לא מוסברת}}$$

כאשר $SS_{\text{reg}} = \sum_{i=1}^n (\hat{y}_i - \bar{y})^2$ - השונות שהמודל "לכד".
קשר ל- R^2 :

$$(23) \quad R^2 = 1 - \frac{\text{שגיאת המודל שלנו}}{\text{שגיאת מודל הבסיס}} = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}$$

זו הגרסה החיובית של ההגדרה - במקום "אחד פחות שאריות", זה "שיעור המוסבר".

3.7 קשר למקדם הקורלציה

במקרה של רגרסיה ליניארית פשוטה (משתנה בלתי תלוי אחד), מתקיים:

$$(24) \quad R^2 = r^2$$

כאשר r הוא מקדם הקורלציה של פירסון בין X ו- Y :

$$(25) \quad r = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum_{i=1}^n (x_i - \bar{x})^2} \sqrt{\sum_{i=1}^n (y_i - \bar{y})^2}}$$

הוכחה (סקיצה):

ברגרסיה פשוטה, $\hat{y}_i = a + bx_i$, כאשר:

$$b = \frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sum (x_i - \bar{x})^2} = r \cdot \frac{\sigma_y}{\sigma_x}$$

ניתן להוכיח (באלגברה מייגעת אך ישירה) ש- $SS_{\text{reg}} = r^2 \cdot SS_{\text{tot}}$, ולכן $R^2 = r^2$. ■

משמעות: ברגרסיה פשוטה, R^2 הוא ריבוע הקורלציה. אם $r = 0.9$ (קורלציה גבוהה), אזי $R^2 = 0.81$.

3.8 סכנה 1: R^2 עולה תמיד כשמוסיפים תכונות

הבעיה הקריטית של R^2 : הוא אף פעם לא יורד כשמוסיפים תכונות, גם אם הן חסרות תועלת!

משפט 3.2 - מונוטוניות R^2 :

אם מוסיפים תכונה למודל רגרסיה, R^2 לא יורד (ובדרך כלל עולה).

הוכחה:

צעד 1: נניח מודל עם p תכונות מושג R_p^2 . נוסיף תכונה נוספת (סה"כ $p+1$).
 צעד 2: המודל החדש יכול תמיד לבחור משקל $w_{p+1} = 0$ לתכונה החדשה, ובכך להשיג את אותה שגיאה כמו המודל הישן.
 צעד 3: מכיוון שהרגרסיה מיזער את SS_{res} , היא תמצא משקלים טובים או שווים:

$$SS_{\text{res}}^{(p+1)} \leq SS_{\text{res}}^{(p)}$$

צעד 4: לכן:

$$R_{p+1}^2 = 1 - \frac{SS_{\text{res}}^{(p+1)}}{SS_{\text{tot}}} \geq 1 - \frac{SS_{\text{res}}^{(p)}}{SS_{\text{tot}}} = R_p^2 \quad \blacksquare$$

המשמעות המסוכנת:

אם נוסיף 1000 תכונות אקראיות (רעש טהור), R^2 יעלה! זה יוביל אותנו לחשוב שהמודל השתפר, אך באמת הוא פשוט **התאים יתר** (Overfitted).

דוגמה מספרית – סימולציה:

תוצאה צפויה:

טבלה 7: עליית R^2 עם הוספת תכונות אקראיות

R^2 (אימון)	מספר תכונות
0.850	1
0.852	2
0.870	6
0.895	11
0.930	21
0.975	51

R^2 עלה מ-0.85 ל-0.975 למרות שהתכונות הנוספות היו רעש טהור!

3.9 הפתרון: Adjusted R^2

Adjusted R^2 (מקדם הקביעה המתוקן) פותר את הבעיה על ידי הטלת "קנס" על הוספת תכונות.

הגדרה 3.2 – Adjusted R^2 :

$$(26) \quad R_{\text{adj}}^2 = 1 - \frac{SS_{\text{res}}/(n-p-1)}{SS_{\text{tot}}/(n-1)}$$

כאשר:

n – מספר התצפיות

```
import numpy as np
from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score

# שער  $Y = 2X +$  מיטושפמינות
np.random.seed(42)
n = 100
X_real = np.random.rand(n, 1)
y = 2 * X_real.flatten() + np.random.randn(n) * 0.5

# תוירקאתוננוכטרפטמלשהיצקנופכ  $R^2$  קודבנ
r2_scores = []

for num_random_features in [0, 1, 5, 10, 20, 50]:
    # שערותוירקאתוננוכתתפטוה (!)
    X_random = np.random.randn(n, num_random_features)
    X_combined = np.hstack([X_real, X_random])

    # לומיאל
    model = LinearRegression()
    model.fit(X_combined, y)

    # לומיאהינוותנלע  $R^2$ 
    y_pred = model.predict(X_combined)
    r2 = r2_score(y, y_pred)
    r2_scores.append(r2)

    print(f"Num_features: {1+num_random_features},  $R^2$ : {r2:.4f}")

# שערתיאסמג , הנוכתלכסעהלוע  $R^2$ : האצות
```

- p - מספר התכונות (לא כולל איבר חופשי)

ניסוח חלופי:

$$(27) \quad R_{\text{adj}}^2 = 1 - (1 - R^2) \cdot \frac{n - 1}{n - p - 1}$$

משמעות:

המונה והמכנה מחולקים במספר **דרגות החופש** (Degrees of Freedom). ככל שיש יותר תכונות p , דרגות החופש קטנות יותר, והקנס גדל.

תכונות של R^2 Adjusted:

- **יכול לרדת** כשמוסיפים תכונה חסרת תועלת

- **יכול להיות שלילי** (אם המודל גרוע מאוד)

$$R_{\text{adj}}^2 \leq R^2 \quad \text{תמיד}$$

מתי להשתמש?

- R^2 : להערכת התאמה על נתוני אימון בודדים

- R^2 Adjusted: להשוואה בין מודלים עם מספר תכונות שונה

חזרה על הסימולציה עם R^2 Adjusted:

תוצאה צפויה:

טבלה 8: השוואה: R^2 מול R^2 Adjusted

תכונות	R^2	R^2 Adjusted
1	0.850	0.849
2	0.852	0.849
6	0.870	0.862
11	0.895	0.880
21	0.930	0.900
51	0.975	0.895

R^2 Adjusted יורד כשמוסיפים תכונות חסרות תועלת - בדיוק כמו שצריך!

3.10 סכנה 2: R^2 לא מתאים לסיווג

R^2 מיועד ל**רגרסיה** (תחזית של ערך מספרי). עבור **סיווג** (תחזית של קטגוריה), הוא לא מתאים.

למה?

בסיווג, $y_i \in \{0, 1\}$ אך \hat{y}_i יכול להיות כל מספר (למשל, 0.73 = הסתברות לקטגוריה 1). חישוב SS_{res} לא משקף את דיוק הסיווג.

```
def adjusted_r2(r2, n, p):
    """
    מחשב Adjusted  $R^2$ .

    Args:
        r2:  $R^2$  ליגור
        תומיגד תרפסמ: ת
        p: תרפסמ (אל תונוכת תרפסמ: ת)
    """
    return 1 - (1 - r2) * (n - 1) / (n - p - 1)

# שדחמבו שיח
for i, num_random in enumerate([0, 1, 5, 10, 20, 50]):
    r2 = r2_scores[i]
    n = 100
    p = 1 + num_random
    adj_r2 = adjusted_r2(r2, n, p)

    print(f"Features: {p},  $R^2$ : {r2:.4f}, Adj  $R^2$ : {adj_r2:.4f}")
```

מדדים נכונים לסיווג:

טבלה 9: מדדי הערכה למשימות שונות

משימה	מדד	הסבר
רגרסיה	R^2 , RMSE, MAE	R^2 מודד שיעור שונות מוסברת
סיווג	Accuracy, Precision, Recall, F1	מודדים אחוז סיווגים נכונים
סיווג (התפלגות)	Log-Loss, AUC-ROC	מודדים איכות הסתברויות

3.11 קשר לאלגברה ליניארית: הקרנות

ניתן לראות רגרסיה ליניארית כהטלה של \vec{y} על המרחב שנפרש על ידי עמודות המטריצה X .

נוסחת הרגרסיה במטריצות:

$$(28) \quad \vec{\hat{y}} = X(X^T X)^{-1} X^T \vec{y} = P \vec{y}$$

כאשר $P = X(X^T X)^{-1} X^T$ היא מטריצת ההטלה (Projection Matrix).
משמעות גיאומטרית של R^2 :

$$R^2 = \frac{\|\vec{\hat{y}} - \vec{\bar{y}1}\|^2}{\|\vec{y} - \vec{\bar{y}1}\|^2} = \cos^2(\theta)$$

כאשר θ היא הזווית בין $\vec{y} - \vec{\bar{y}1}$ (הנתונים הממורכזים) ל- $\vec{\hat{y}} - \vec{\bar{y}1}$ (החיזוי הממורכז).
פרשנות: R^2 הוא ריבוע הקוסינוס של הזווית בין הנתונים לחיזוי במרחב הממורכז – מדד גיאומטרי של התאמה.

3.12 תרגיל תכנות עצמי 3.1 – חישוב R^2 ידני

מטרה: להבין את המשמעות של R^2 על ידי חישוב ידני.
משימה:

1. צרו נתונים סינתטיים: שער $y = 3x + 5$

2. אמנו מודל רגרסיה ליניארית

3. חשבו את R^2 בשלוש דרכים:

- נוסחה ישירה: $1 - SS_{\text{res}}/SS_{\text{tot}}$

- באמצעות קורלציה: r^2

- באמצעות sklearn

4. השוו את התוצאות (צריכות להיות זהות!)

פסאודו-קוד מלא:

תוצאה צפויה: כל שלוש השיטות יתנו $R^2 \approx 0.95$ (תלוי ברעש).

3.13 תרגיל תכנות עצמי 3.2 – השוואת R^2 ו- R^2_{detsujdA}

מטרה: להמחיש את ההבדל בין R^2 ל- R^2_{detsujdA} כאשר מוסיפים תכונות.
משימה:

1. צרו נתונים עם תכונה אחת רלוונטית

2. הוסיפו בהדרגה תכונות אקראיות (רעש)

3. עקבו אחרי R^2 ו- R^2_{detsujdA} כפונקציה של מספר התכונות

4. הציגו גרף: ציר X = מספר תכונות, ציר Y = ציון

תוצאה צפויה: R^2 עולה בהתמדה, אך R^2_{Adjusted} יורד אחרי כמה תכונות.

```

import numpy as np
from sklearn.linear_model import LinearRegression
from sklearn.metrics import r2_score

# 1. סינונתרצי
np.random.seed(42)
n = 100
X = np.random.rand(n, 1) * 10 # X ל־0 ו־10
y = 3 * X.flatten() + 5 + np.random.randn(n) * 2 # y = 3x + 5 + שער

# 2. לדומןוחיא
model = LinearRegression()
model.fit(X, y)
y_pred = model.predict(X)

print(f"דמלנשםלדומה: y={model.coef_[0]:.3f}x+{model.intercept_:.3f}")

# 3. הרשיהחסונ 1: דרד
y_mean = np.mean(y)
SS_tot = np.sum((y - y_mean)**2)
SS_res = np.sum((y - y_pred)**2)
r2_manual = 1 - (SS_res / SS_tot)

print(f"\nחסונ1מדרד(הרשיה):")
print(f"SS_tot={SS_tot:.2f}")
print(f"SS_res={SS_res:.2f}")
print(f"R²={r2_manual:.6f}")

# 4. הרשיהחסונ 2: דרד (! הטושפהיסרגל) קרעובירבהיצלרוק
correlation = np.corrcoef(X.flatten(), y)[0, 1]
r2_from_corr = correlation**2

print(f"\nהיצלרוק2מדרד():")
print(f"r={correlation:.6f}")
print(f"R²=r²={r2_from_corr:.6f}")

# 5. דרד 3: sklearn
r2_sklearn = r2_score(y, y_pred)

print(f"\nחסונ3מדרד(sklearn):")
print(f"R²={r2_sklearn:.6f}")

# 6. האוושה
print(f"ההזמלכה? {np.allclose([r2_manual, r2_from_corr, r2_sklearn], r2_sklearn)}")

```

```

import matplotlib.pyplot as plt

def adjusted_r2(r2, n, p):
    return 1 - (1 - r2) * (n - 1) / (n - p - 1)

# סינון
np.random.seed(42)
n = 100
X_real = np.random.rand(n, 1)
y = 2 * X_real.flatten() + np.random.randn(n) * 0.5

# בקנה
num_features_list = []
r2_list = []
adj_r2_list = []

for num_random in range(0, 51):
    # תו'אקאוןוכתפטסה
    if num_random == 0:
        X_combined = X_real
    else:
        X_random = np.random.randn(n, num_random)
        X_combined = np.hstack([X_real, X_random])

    # לוחיא
    model = LinearRegression()
    model.fit(X_combined, y)
    y_pred = model.predict(X_combined)

    # בושיו
    r2 = r2_score(y, y_pred)
    p = X_combined.shape[1]
    adj_r2 = adjusted_r2(r2, n, p)

    num_features_list.append(p)
    r2_list.append(r2)
    adj_r2_list.append(adj_r2)

# גרף
plt.figure(figsize=(10, 6))
plt.plot(num_features_list, r2_list, label='R²', marker='o', markersize=3)
plt.plot(num_features_list, adj_r2_list, label='Adjusted R²', marker='s', markersize=3)
plt.xlabel('Number of Features')
plt.ylabel('Score')
plt.title('R² vs Adjusted R²: Adding Random Features')
plt.legend()
plt.grid(True, alpha=0.3)
plt.show()

```

3.14 מקרי קצה: מתי R^2 מטעה?

מקרה 1: יחסים לא-ליניאריים

אם הקשר בין X ל- Y הוא לא-ליניארי (למשל, ריבועי), רגרסיה ליניארית תיכשל ו- R^2 יהיה נמוך – אך זה לא אומר שאין קשר!

דוגמה: $y = x^2$

רגרסיה ליניארית תיתן R^2 נמוך, אך הקשר מושלם (רק לא ליניארי).

פתרון: השתמש בפולינומים או מודלים לא-ליניאריים (כמו רשתות נוירונים).

מקרה 2: Outliers (ערכי חריגים)

ערך חריג אחד יכול להוריד דרמטית את R^2 , גם אם המודל טוב ל-99% מהנתונים.

פתרון: זהה וטפל ב-sreiltuO לפני אימון (באמצעות Z-score, IQR, או שיטות עמידות כמו RANSAC).

מקרה 3: מתאם אינו סיבתיות

R^2 גבוה לא אומר ש- X גורם ל- Y . ייתכן שיש משתנה שלישי Z שגורם לשניהם.

דוגמה קלאסית: צריכת גלידה וטביעות מים מתואמות – לא כי גלידה גורמת לטביעות, אלא כי קיץ גורם לשניהם.

3.15 סיכום ומבט קדימה

מה למדנו בפרק זה?

1. R^2 מודד שיעור שונות מוסברת – מספר יחיד שמסכם את איכות המודל

2. תחום: 0 עד 1 – הוכחנו מתמטית למה זה תמיד נכון (לרגרסיה עם intercept)

3. קשר לקורלציה – ברגרסיה פשוטה, $R^2 = r^2$

4. בעיה: עולה תמיד – הוספת תכונות מעלה את R^2 גם אם הן חסרות תועלת

5. פתרון: Adjusted R^2 – מנרמל לפי מספר התכונות ומונע Overfitting

6. לא לסיווג! – R^2 מיועד לרגרסיה, לא לסיווג

7. מגבלות – מתאם אינו סיבתיות, רגיש ל-sreiltuO, מניח ליניאריות

מבט קדימה – פרק 4:

בפרק הבא נעבור מהערכת מודלים למדידת קשרים בין משתנים. נחקור:

- קו-ווריאנס (Covariance) – מדד גולמי של קשר ליניארי

- קורלציה (Correlation) – הגרסה המנורמלת של קו-ווריאנס

- ייצוג אלגברי ליניארי – קו-ווריאנס כמכפלה סקלרית, קורלציה כ- ytiralimiS enisoC

- סכנות – "Correlation is not Causation", קורלציות מזויפות

- מטריצת קורלציה - כלי לזיהוי Multicollinearity

שאלת מחשבה לסיום:

אם שני משתנים X ו- Y בעלי קורלציה $r = 0$ (אורתוגונליים), האם זה אומר שהם בלתי תלויים סטטיסטית?
רמז: התשובה היא לא! קורלציה אפס משמעה רק אין קשר ליניארי, אך עדיין יכול להיות קשר לא-ליניארי (למשל, $Y = X^2$).

מטלות וקריאה מורחבת

- תרגיל 3.1:** חשבו R^2 ידנית בשלוש דרכים (ראו פסאודו-קוד).
תרגיל 3.2: השוו R^2 ו- detsujdA עם תכונות אקראיות (ראו פסאודו-קוד).
תרגיל 3.3: הוכיחו בעצמכם שברגרסיה פשוטה $R^2 = r^2$.
רמז: התחילו מהביטוי $SS_{\text{reg}} = \sum (\hat{y}_i - \bar{y})^2$ והשתמשו ב- $\hat{y}_i = a + bx_i$.
קריאה מורחבת:

- [23] - "Regression towards Mediocrity in Hereditary Stature": המאמר המקורי של Galton
- [24] - "Mathematical Contributions to the Theory of Evolution": Pearson על קורלציה
- [25] - "Statistical Methods for Research Workers": Fisher על ניתוח שונות
- [26] - "Correlation and Causation": Wright טובע את R^2
- [22] - The Elements of Statistical Learning, פרק 3: רגרסיה ליניארית ברזולוציה גבוהה

שאלות להעמקה:

1. למה $\text{Adjusted } R^2$ יכול להיות שלילי, אך R^2 לא?
2. מה יקרה ל- R^2 אם ננרמל את X ו- Y (נחסיר ממוצע ונחלק בסטיית תקן)?
3. איך R^2 מתייחס למודלים לא-ליניאריים (כמו רשתות נוירונים)?

סיום פרק 3

4 English References

- 1 G. Strang, *Linear Algebra and Learning from Data*. Wellesley, MA, USA: Wellesley-Cambridge Press, 2019.
- 2 T. Mikolov, K. Chen, G. Corrado, and J. Dean, “Efficient estimation of word representations in vector space,” in *Proceedings of the International Conference on Learning Representations (ICLR)*, 2013.
- 3 Z. S. Harris, “Distributional structure,” *Word*, vol. 10, no. 2-3, 146–162, 1954.
- 4 A. Caliskan, J. J. Bryson, and A. Narayanan, “Semantics derived automatically from language corpora contain human-like biases,” 6334, 356, 2017, 183–186.
- 5 T. Bolukbasi, K.-W. Chang, J. Zou, V. Saligrama, and A. Kalai, “Man is to computer programmer as woman is to homemaker? debiasing word embeddings,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 29, 2016.
- 6 R. E. Bellman, *Dynamic Programming*. Princeton, NJ, USA: Princeton University Press, 1957.
- 7 Y. LeCun, L. Bottou, Y. Bengio, and P. Haffner, “Gradient-based learning applied to document recognition,” 11, 86, 1998, 2278–2324.
- 8 A. Krizhevsky, I. Sutskever, and G. E. Hinton, “Imagenet classification with deep convolutional neural networks,” in *Advances in Neural Information Processing Systems (NeurIPS)*, 25, 2012.
- 9 I. Goodfellow, Y. Bengio, and A. Courville, *Deep Learning*. Cambridge, MA, USA: MIT Press, 2016.
- 10 K. Beyer, J. Goldstein, R. Ramakrishnan, and U. Shaft, “When is “nearest neighbor” meaningful?” *Lecture Notes in Computer Science*, vol. 1540, 217–235, 1999.
- 11 C. C. Aggarwal, A. Hinneburg, and D. A. Keim, “On the surprising behavior of distance metrics in high dimensional space,” *Lecture Notes in Computer Science*, vol. 1973, 420–434, 2001.
- 12 E. Fix and J. L. Hodges, “Discriminatory analysis: Nonparametric discrimination: Consistency properties,” *Technical Report, USAF School of Aviation Medicine*, 1951.
- 13 V. N. Vapnik and A. Y. Chervonenkis, “On the uniform convergence of relative frequencies of events to their probabilities,” *Theory of Probability and Its Applications*, vol. 16, no. 2, 264–280, 1971.
- 14 J. Deng, W. Dong, R. Socher, L.-J. Li, K. Li, and L. Fei-Fei, “Imagenet: A large-scale hierarchical image database,” in *Proceedings of the IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, 2009, 248–255.

- 15 Y. LeCun et al., "Backpropagation applied to handwritten zip code recognition," *Neural Computation*, vol. 1, no. 4, 541–551, 1989.
- 16 K. Pearson, "On lines and planes of closest fit to systems of points in space," *Philosophical Magazine*, vol. 2, no. 11, 559–572, 1901.
- 17 H. Hotelling, "Analysis of a complex of statistical variables into principal components," *Journal of Educational Psychology*, vol. 24, no. 6, 417–441, 1933.
- 18 G. H. Golub and C. Reinsch, "Singular value decomposition and least squares solutions," *Numerische Mathematik*, vol. 14, 403–420, 1970.
- 19 A. E. Hoerl and R. W. Kennard, "Ridge regression: Biased estimation for nonorthogonal problems," *Technometrics*, vol. 12, no. 1, 55–67, 1970.
- 20 R. Tibshirani, "Regression shrinkage and selection via the lasso," *Journal of the Royal Statistical Society: Series B*, vol. 58, no. 1, 267–288, 1996.
- 21 A. Chaman and I. Dokmanic, "Truly generative or just extrapolation? on the ability of generative models to create truly novel content," *arXiv preprint arXiv:2109.09018*, 2021.
- 22 T. Hastie, R. Tibshirani, and J. Friedman, *The Elements of Statistical Learning: Data Mining, Inference, and Prediction*, 2nd. New York, NY, USA: Springer, 2009.
- 23 F. Galton, "Regression towards mediocrity in hereditary stature," *Journal of the Anthropological Institute of Great Britain and Ireland*, vol. 15, 246–263, 1886.
- 24 K. Pearson, "Mathematical contributions to the theory of evolution. iii. regression, heredity, and panmixia," *Philosophical Transactions of the Royal Society of London. Series A*, vol. 187, 253–318, 1896.
- 25 R. A. Fisher, *Statistical Methods for Research Workers*. Edinburgh, UK: Oliver and Boyd, 1925.
- 26 S. Wright, "Correlation and causation," *Journal of Agricultural Research*, vol. 20, no. 7, 557–585, 1921.