

איפה נלמדת מידע? (אנליזה), ארגון של המידע, preprocessing...

מנקים: המרחב הן $x \in \mathbb{R}^p$, יש n דוגמאות (nmp)

p הוא מימד המרחב. אנו לא יכולים לראות מספר דוגמאות

אקסטרנסיבי p - p כפי שהיה. אנו הורחבנו p (ה')

מימדים: ambient dimensionality - כמה פיזיקלי יש ישירות בדאטא?

למשל במחשבה: מס' הפיקסלים RGB

intrinsic dimensionality - כמה אינפורמציה סתומה יש בדאטא?

כי תמונה סופית ממש לא סורטלית

\mathbb{R}^{30000} (גודל תמונה סטנדרטית 3)

נקטל לא האינפורמציה שבדאטא אפילו

כנראה איננו מסוגלים להבחין קטן יותר

של פיזיקליים.

אפשר לחשוב על
מרחב מידע של המידע
של הדאטא.

אז כמו שהיה בעמים מיוחדים תנונו שהדאטא כמו sparsity,

אנחנו נעסוק בדאטא שאפשר לחקור שיש לו מידע אינטרינטי. נחקר.

① דאטא שבו יש מידע אינטינטי. מידע נחקר.

נרצה להבין את הדאטא למרחב הנמוך יותר, ונבין לא

נאסר אינפורמציה (אנליזה) בעזרת אנחנו בין המידע.

Δ אפשר לעשות את זה אם יש דאטא בעל מרחב קטן המרחב

הוא חסוך מרחב.

אין נגלה מידע המידע של הדאטא? (כמה המידע של המרחב כוונתו שכן)

- נשים את הדאטא במטריצה קצת, נבדוק את המטריצה,

ונבדוק מה הדגה של המטריצה המורגת - זו הדגה לחישוב.

- למעשה, נלמדת SVD על אגרה מטריצה קצת, ומספר הערכים

הסינגולריים (אנליזה האנליזה של U - V U^T) שאנחנו, זו הדגה.

משהו שהיה נראה יותר העבר הדאטא למידע אחר זה לא

המידע בין אנליזה בדאטא.

הוכחה מחדש של PCA - תהליך

PCA - 1.1

מכנסים את מטרית השונות האמפירית של הנתונים (p x p)

$$S = \frac{1}{n-1} \sum_{i=1}^n (x_i - \bar{x})(x_i - \bar{x})^T$$

↑ Empirical Covariance Matrix

המטרית S היא כזו של קונסטיטנטים (population covariance)

היא גם PSD, p.f. $S = U \Lambda U^T$ ו- Λ אלכסונית.

כל (עצם) הערכים של המטרית (הערך האולן) הם:

$$y_i = x_i \cdot U_d$$

data in new dimension d

data (as row vecs)

first d columns of U (d x p)

MDS - 2.1

השיטה היא אלמנטרית וזוהי ה אחת מהשיטות הפשוטות ביותר

$$\Delta_{ij} = \|x_i - x_j\|^2$$

ועדיין נקרא מרחק מניאלי

מה זה יקרה לנו? אם נשתמש בשיטה הזו / נתון בין

אנשים על שטח חלקי משתנים קטנים / כמות...

אם נלמד מטרית Similarity (n x n):

הפרדה של המטרית

$$S = -\frac{1}{2} H \cdot \Delta \cdot H$$

עבור data centering $H = I - \frac{1}{n} \mathbf{1} \mathbf{1}^T$ מטרית סימטרית לחלוטין

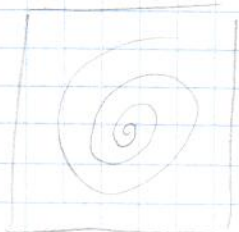
אז מטרית ה- S היא מטרית סימטרית

לכן אם S היא מטרית סימטרית $S = U \Lambda U^T$ ו- Λ אלכסונית

הערות: u_1, u_2, \dots, u_n (הערות האולן של A)

כל y_i הוא הערך ה- i של המטרית הזו.

② אסא להויל אפסיל לוקלי:



Swiss Roll Data

המקרה העליוני מה קויל לא ענייני.

מסנה למה, המרחק האקלידיס d ambient space



המרחק בין נקודות

המרחק "לוקלי" זה הכולל ולא להשתמש ישירות

מרחקים

APML
למה

Manifolds: מרחב שבו יש מרחב חלק (smooth) ונקודה

ק' למקרה מסוימת אפשר להחזיר אותו לפי המרחק הקרוב

manifold מרחב ק' ניתן לקרוב הקרוב בצורה מרחב לוקלי מרחב ק'.

← אסא לוקלי - Manifold

נקודות x_1, x_2, \dots, x_n הנמצאות שטח מרחב קרוב לקואורדינטות

נקודות $y_1, \dots, y_n \in \mathbb{R}^d$ כאשר הנחנו שהאסא יונק מה Manifold מרחב ד.

נקודות אסא
המרחקים הקטנים
לוקלי - $\|x_i - x_j\|$
if $\|x_i - x_j\| < \epsilon$

Diffusion maps \gg אקורדיות: LLE

Locally Linear Embedding: LLE

המרחב: נעשה חידוש מרחב לוקליס קרוב, ואח"כ נשלב.

→ זה היינו מחפשים למרחב לוקליס - Manifold מקרוב נראה לוקלי.

(כמו שפירוש האלף נראה לנו שטוח ...)

האלגוריתם 1. לכל נקודה x_i , נבחר neighbourhood של נקודה לקרוב

אלה במרחק אקלידיס \mathbb{R}^d - מרחב KNN

2. נקודה חודשה מרחב לוקליס אל x_i הנקודה

3. נחבר אלמנט האלמן חלק פשוט.

חידוש שלם 2: נחשב מרחב W שמראה למרחב אסא $\text{residual sum of squares}$

מרחב reconstruction של x_i מהשלים שלו. → זה מרחב הקרוב לוקליס של x_i אל השלים שלו.

W נקרא מרחב Sparse ואפשר לחשוב אלה כמרחב מרחב x_i

אל השלים שלו (המרחב הקרוב לוקליס - קרוב ו-1 מרחב KNN)

$$\|x_i - \sum_{j=1}^n W_{ij} x_j\|^2$$

פירוט שלם כי אין נחמד את הנקודות אחרי תורגלו הנושא?

נמצא קואורדינטות Y שמקלים למינימום את:

$$\Phi(Y) = \sum_{i=1}^n \left\| x_i - \sum_{j=1}^n W_{ij} y_j \right\|^2$$

$$\text{כך } Y^T Y = I \text{ ו- } \sum_{j=1}^n y_j = 0$$

↑
לומר הפיצולים אורגניזציה אחת עשני.

קרייזיין

* בחינוך א שנים הני קרובים, אסטרטגיה זקוקה ϵ פשוט ולקחת

את ϵ הטנגים להפחית הרעיון ϵ ספיט הנקודה.

ה-KNN אלוהי זהירות לנו נקודה רחוקה מסעיג ובכבוד ϵ אנחנו

לא יודעים כמה נקודה יהיו. זה כי שקול קדאסא למה מפורסם בורה

אחריה...

* איך עושים חיפוש KNN? מה הסיכויים?

המקרה הקשה זה $(n, 0)$ $\left(\begin{smallmatrix} n \\ 0 \end{smallmatrix} \right)$.

הפער המשמש ב-k-d tree (אם יש לנו את x_1, \dots, x_n וליקח מחקים)

ויל גם אלקוריתמים לקיום של KNN, שהם הרבה יותר יעילים.

רעיון נוסף: random-projection של הדיסטנס L_2

↑
רוח הגמון זה טוב כאלפין מפתח וזה הסקופיה
הרבה יותר טובה (פצרה FFT ממש)

חיפוש k-NN במרחב המינימלי

ואפשר לעבור אל המרחב הזה כמה פעמים ולמצא את

המקומות הקרובים חכמה בשלטי.

* אחריו W זה קסמים מוצאה של רגעים זינאריה על כל נקודה,

ומעשה מקלים קומקוציה אפניג לכל x .

* אפשר גם לדעת מזה כי לא מקלים את הדיסטנס אלא רק מרחקים

וסגין וולא אמר הלקח מקרי בין הנקודות—

המשך: Manifold Learning

Diffusion Maps

מבקשים את ההצגה של המרחב המרחק - המרחק "המרחק" בין
האיברים בקבוצה אחת (רשת).

diffusion distance

$$\Delta_{ij}^t = \sqrt{\sum_{k=1}^n \frac{1}{d_k} (A_{ik}^t - A_{jk}^t)^2}$$

זה המרחק בין שני ענני ההתפלגות בסיים i וסיים j בזמן t .
בצורה מופשטת יותר, קנייה באלמנטים אלו.

(אם A הוא המרחב A : Φ_1, \dots, Φ_n : המרחב A : $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n = 1$).

והקנייה: $\Phi_t: x_i \mapsto (\lambda_1^t \Phi_1(i), \dots, \lambda_n^t \Phi_n(i))$

diffusion map

עכשיו נקרא מרחקים המרחקים!

$$\|\Phi_t(x_i) - \Phi_t(x_j)\| = \Delta_{ij}^t$$

זה המרחק
המרחק המרחק n

ואפשר לתאר את
המרחק Δ_{ij}^t כ-
שקילה של המרחק
המרחק המרחק

