

הגבול של  $\epsilon$  kernel הוא  $G \subseteq H$  גיה  $\epsilon$ -kernel אם  $\epsilon$  לכל  
 $(1-\epsilon)\epsilon_H(v) \leq \epsilon_G \leq \epsilon_H(v)$ ,  $\forall v \in V$

הוא, במקום של מילונים, אפשר לקרוא פונקציות אלה  
במקום אחרים.  $\epsilon$ -kernel פונקציות, ויש ענייני צורה אחרים  
אליהם בנויה על מילונים.  
אפשר יהיה לשים מילונים, ולקחת את  $\epsilon$ -kernel, וה- $\epsilon$ -kernel  
יוכל ע"י צורה המון צורה ש"מילונים".

הרצאה

היום א"מ נחמק על LSH - Locality Sensitive Hashing

עד 2 נחמק מילונים Nearest Neighbor בקלות. זמן לזכור לבנות  
ל- $\log$  query. עם המילונים זה נהיה אקספוננציאלי.

ה-LSH זה הישגה הישגה כיום לאנליזה אקספוננציאלית, אבל זה  
ה-LSH קצר יותר יקר. במקום query ע"י זה אולי, פשוט  
ה-brute force.

ההסבר הוא להשתמש ב-hashing, והוא יותר פשוט ופשוט  
נראה אולי זה מילונים.

יש 2 הסברים ל-LSH בספר של  $\epsilon$ .  $\epsilon$  Charikar  
או  $\epsilon$  Indyk & Motwani ...

אם  $\epsilon$ , ההפרדה: (צורה)  
למה  $H$  של פונקציות היא Locality Sensitive כמו  $\epsilon$ -kernel, similarity  
 $0 \leq sim(p,q) \leq 1$  אם  
 $Pr[h(p)=h(q)] = sim(p,q)$   
במילים אחרות, אפשר לומר שיש קשר בין  $\epsilon$ -kernel ל-similarity.

סומת, ככל שכל נק' יוגר "קומת", יוגר "1", ההסתברות שילכו  
סומת bucket יוגר זקוף, וזאת אלוה ע- sim שלהם.

Hamming Similarity

אם שני הנק' של שני strings באורך m בימים מןן לנקיים  
פונק' קומת של

$$sim(p,q) = 1 - \frac{ham(p,q)}{m}$$

כאשר ham(p,q) זהו מספר הבימים ש- p ו- q אינם  
אז המספר H הוא קי למקבל

$$H = \sum h_i(p) = p \text{ זהו מספר}$$

Jaccard?  
עוד קומת Jaccard - מספר המקיים p ו- q קומת אלקי  
מספר המקיים יש סה"כ.

$$\begin{matrix} 101110 \rightarrow 1,3,4,5 \\ 101010 \rightarrow 1,3,5 \end{matrix} \rightarrow \frac{3}{4}$$

אם כן bit שונה יכול עייצ term שלומא שלומא ב- p ו- q.  
פונקציה H-ה כן: מספרים של פונקציה אקראית ומספרים של פונקציה

$$H = \sum h_r(p) = \min_{i \in \pi} \text{of the items in } p$$

מספרים אלו... של מספר פונקציה יש אזה המספר זהו ראשון.

אבל זו מספר עצומה. מים אחר שלג באמ צחים אג ס הפונקציה

הן צחים ליום נכח 2 פונקציה על אזה יהיה המספר

שזה זהו הראשון, אחר הןן כיוגר. אז יש כמ מני' מומיים על

באמת הפונקציה...

אם קומת:

יש נק' ב- R והמספר היא כמ העל מיליונים. כמ h של  
על מיליון אמרה נק' ע-1 אם הנק' מעל.

המספר שונה בין כל נק' פונקציה לנק' אחר, אחר המספר  
שזה מיליון בקיוק עזרי כנייה. אם כן, ילכו ע-2 buckets

$$Pr[h_r(p) = h_r(q)] = 1 - \frac{\theta}{\pi}$$



שניים.

מה הבטיחה? קל"ם (buckets) זה לא מספיק. אולי עוקבים  
 ספרתיקה בטל הרבה עם מילונים אקראיים אולי פונקצי  
 ה- hash היא למעשה שרשרת של כמות; מה השאלה.  
 אולי נקרא עברו כל שבת ומוסססס אולי זה כולו עשן  
 ונקרצ'ה עם hamming distance ואולי מילים אחרות עם  
 האם מקבלים מדיוק אולי מרוחק, אולי מתחילים עם במחשבו  
 "קראבוג".

לשלי:

אם ק"מ לשל H עבר פונק' similarity, אז  $d(p,q) = 1 - sim$  מקיים  
 או אי שיוון השלש.

לכ"ו:

$$X(p,q) = \begin{cases} 1 & \text{אם שייך ל bucket} \\ 0 & \text{אחרת} \end{cases}$$

$$X(p,q) + X(q,r) \geq X(p,r)$$

אולי יש כמה מקרים של אם  $p+q$  במילון או  $p, r$  לא אז  $q, r$  לא  
 וכן הלאה. זה כמות מקיים בכלום ואולי זה בסדר גמור.

$$\Rightarrow 1 - sim(p,q) + 1 - sim(q,r) \geq 1 - sim(p,r)$$

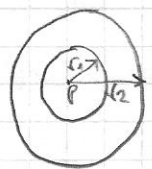
$$d(p,q) + d(q,r) \geq d(p,r)$$

באמת מקיים אי שיוון השלש וזאת יש לה מסוים למתיקה.

וכן ההפסדה האמת  $\delta$  - LSH:

לשלה H של פונק' היא  $(r_1 < r_2, p_1 > p_2)$ -sensitive אז  
 $d(p,q) \leq r_1 \rightarrow Pr[H(p)=H(q)] \geq p_1$   
 $d(p,q) \geq r_2 \rightarrow Pr[H(p)=H(q)] \leq p_2$

שומר, נק' קבלת העלם בהוס' גבוהה שזוהי אולי מוקד, אז  
 בהוס' נאכה הוס' שזוהי אולי.



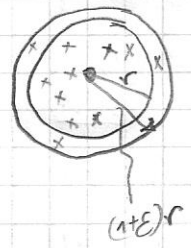
אולי שומר או הקשר להפסדה הקבלת, בין המוק' ו  
 research.



LS,  $p_1 = 1 - \frac{r_1}{m}$ ,  $p_2 = 1 - \frac{r_2}{m}$ .  
אם  $p_1 = 1 - \frac{r_1}{m}$ ,  $p_2 = 1 - \frac{r_2}{m}$ .

אם לא יוקעים כלומר איך עסקי את בע"מ זה  $N$ , אבל נבדוק  
מה קרה.  $\approx$  עושים וק' קציה נכס"ג זה  $AW$  (Approximate)  
עב"ה זה  $(r, \epsilon)$ -neighbor problem

אם קיים סמן  $p$  כך ש-  $d(p, q) \leq r$  אז נרצה להוכיח  $d(p, q) \leq (1+\epsilon)r$



אז הבעיה נכרעה עסקי.

נבנה מבנה נתונים הסגור (גורף פונק' זה hash הקטנים), כמות  
שנבנה בהסגר קבועה, אוז נאם להכניס אתה ע"י הפעולה כמו בעמ"ס.

בסקי עסקי  $\sqrt{m}$  זה  $(r, \epsilon)$ -Neighbor (נכ"ק א שיהא  $LS$  ע  
הפונקציות:  $(p_1, p_2, r, (1+\epsilon)r)$  והביצועים שלנו יהיו גלויים  
זה  $Gap$  בין  $p_1$  ו-  $p_2$ .  
זה - Hamming dist. וקעים נהים:  $p_1 = 1 - \frac{r}{m}, p_2 = 1 - \frac{(1+\epsilon)r}{m}$ .

אנחנו רוצים הסגר הפעולה קבועה, גיבור, שלא גרעה גלוייה זה  $(1+\epsilon)r$   
וזמן הריצה יהיה גלוי זה  $p_1$  ו-  $p_2$ .  
כדי לגבוס נק' צמיק עמנו עם הניסוי  $\frac{1}{m}$  בעמ"ס... כי הסיכוי לטלן הוא  $1-p$   
ואם נמני עם זה היתה בעמ"ס נקרא  $\frac{1}{m}$  או זהו...

(טענה נקראת טענה)  
אם יש המון נק' עם שבר הכנור הנדדיוס  $(1+\epsilon)r$  אז יש סנה של  $\frac{1}{p_1} \cdot p_2$   
אם נמני עם הניסוי  $\frac{1}{m}$  בעמ"ס (באחר נקרא היתה פונקציה) ... וזה ע"ה  
כי זה עסקי...

מה עושים? זה פונק' hash הולך להיות מספר קטני  $k$  פונק' hash  
ואנחנו  $k$  bucket יורב זה  $k$  ערכים שמיצבים אתנו.  
במס עסקי boosting עכסגרות וקילנו משנה  $(p_1 < r_2, p_2 > p_1)$

אם כן קי לגמול נ' בייק  $\frac{1}{p_1^k}$  בוקצ'ים, וגם לקב'  $n(\frac{p_2}{p_1})^k$  נ' פן false positives.

אם זה האל' ?

מג'ע query, נסמן ב-q. מוסי' אג q ע'י ב' אומ' מנה -  $\frac{1}{p_1^k}$  בוקצ'ים. אוסי' אג הנ' ע'י ש'  $n(\frac{p_2}{p_1})^k$ , ומג'יים אג ה'וכ' קרובה. בהסג'ות  $\frac{1}{4}$  נג'י נ' ש'פ'א false positive, כ' י'ל ל'נו ז'א המס' ע' ג'מ'ל' ס'מ'ת - false positives. א'ס, מקב'ים ש'כ'ן קרוב בהסג'ה  $\leq 1 - \epsilon$ .

נ'ל'א ר'ק ע'מ'ור אג א, ו'מ' ש'מ'נ'ן א'מ'ן ז'ה ה- querytime או space, ס'מ'ת ה'ז'כ'ר'ן.

כ'מ'ן ה'ר'יב'ה מ'מ'ס'כ ל'פ'י א ב'ג'ים ע'ס א ב'ונ'ק' ? ע'וסי' אג ה'א'ל'מ'י' כ'ז'י'ה ו'ר'מ'ים מ'נו ה-א ה'א'ל'מ'י' (ב'ס'ק'ה)

$$\Rightarrow k = \log_{1/p_2}(n) - \Theta(\log \log n)$$

נ'ז'י'ב ו'ר'א'ה ש'ה'מ'ן ש' ה- query י'ז'כ'א:  $n^p = n^{\frac{\log(1/p_1)}{\log(1/p_2)}}$

מ'מ'ת' ע'ז'ר א'ז'ו ב'ונ'ק' hash ה'ש'ל'ה ש'מ'ת'ה ב'וקב'ים מ'מ'ים ל'פ'ק' כ'י י'ש מ'מ'ו ב'וקב'ים ר'וק'ים ו'ל'ו ר'ז'י'ב ה'מ'ן ב'וקב'ים ס'ל'י'ים. כ'ק, מקב'ים (פ'מ'ת'ס) מ'ק'א. א'ג ה'ש'ל'ה ה'נו'ס'ט ה'י'א ס'א'ג'ס פ'ר'ס'ו'ר'ז'י'ו'נ'י' ע'-n, (מ'ג'מ'ס ע'מ'ת' ה- buckets ה'מ'מ'ים).

ע'ס, ע'ס ה- ב'וקב'י' ש'ל (hamming) מקב'ים query  $n^{\frac{1}{1+\epsilon}}$  space  $n^{\frac{1}{1+\epsilon}}$

א'ם-ע' ק'ן ז'ה ק'י נ'א'ו...

ו'ש מ'מ'ר ע'-ANM מ- (ע'ה)...

כ'מ' א'י'ים מ'מ'ה ש'ני'סו ל'מ'ש'ל' ב'ב'ר'ים ע'ס א ר'ק ע'ס hamming. ע'כ'ד'ו ע'ס  $L_2$  ו'ש כ'מ'ה ש'ק'ב'ים ע'ס ז'ה, Andoni ו'מ'י'ק' ש'ל indy ס'י'מ'ס

א'ג ז'ה ב- 2005, א'ס מ'י'ם ו'מ'י'ה ה'ש'מ'ל'ו כ'ז'ה... ו'מ'ב'ע'ה ק'ק ← ELLSH

הראו ש  $\frac{1}{4\epsilon}$  הוא סף ונתתייג סף  $\epsilon - \rho$  שהתקבל היה בין  $\frac{1}{4\epsilon}$  ל  $\frac{1}{(4+\epsilon)^2}$  והייתה שורה שמה היום אנחנו עושים  $\frac{1}{(4+\epsilon)^2}$ .

הגלגלה הגדלה ככן, ע"י זריק (או כמו זריקים) של כדורים של מסכים את החתום. וכל פונק' היא אוסף (סדרה) של ~~זריקים~~ זריקים אקסיויים של הפונקציה. זה שכן מסתמך על החתום, וזה פונק' hash יחידה. והכל של הפק' זה היוצא י"י אנחנו אלמנט location point ב-  $\epsilon$  נמקדים את הפונקציות והזריקים הם רצף זריקים וזריק ע"י  $\epsilon$  הוא כדור הפק'. אוכלנו  $33$  פונק' hash בסביבת  $\epsilon$  בעל נפח הפונקציות, ולכן צולטים הסוגר  $NN$  כן מתקנים  $\epsilon$  (אם פועלים ב-  $\epsilon$  ו  $\epsilon$ ).

מורה פאייט!

עכשיו א"כ רוצה עומר סוף כמה נתיים על coresets כחלק משאר העבודה.

כפבור, במספר הפתיחות היו לנו  $n$  נק' ב-  $\epsilon$ , ונתנו  $\epsilon$ -kernel  $Q \subseteq \mathbb{R}^d$ , במקרה  $\sim \epsilon^{\frac{d+1}{2}}$ .

במספר הקואורנט -  $H$ : קב' של  $n$  על נילונים  $h_1, \dots, h_n$  פונקציות עיקריות  $(x_1, \dots, x_n) \rightarrow x_d = h_d(x_1, \dots, x_n)$ .

$(x_1, \dots, x_n)$  הוא כולן קואור' אינטגרים על הוויסטה בין המעטת העיונית והמטונה ולפי קואור' זה extent:

$$\max_i h_i(x_1, \dots, x_n) - \min_i h_i(\dots)$$

הוא שכלו היא קרבה מ רחבה  $G \subseteq H$  כקב' -

$$(1-\epsilon) \text{extent}_H(\vec{x}) \leq \text{extent}_G(\vec{x}) \leq \text{extent}_H(\vec{x}) \quad \vec{x} \in \mathbb{R}^d$$

וקב' (גל) שנה הואת הקואור' עצובי.

במקרה (בפתיחות), עובי כוללן  $u$ :

$$\max_{PEP} p \cdot u - \min_{PEP} p \cdot u$$

↓

$$\sum_{i=1}^n p_i u_i$$



אלו הנק'  $p = (p_1, \dots, p_d)$  גמולי עמ'ה העל מ'טלר  $\xi$  :

$$X_d = \sum_{i=1}^{d-1} p_i X_i + p_d$$

שזה בדיוק התפלגה  $U_i$  : אוקיי גמולת, עם קצת נורמליזציה. כמה?  
כ"כ  $(x_1, \dots, x_{d-1})$  זה לא בדיוק טיוון עם קצת הנורמליזציה אולי  
אפשר לנרמל את זה יהיה, ועכ"ן ההפסד יהיה אולי דבר כמו  
קודם עם כד' בקואור. כלי. אולי זהו צדדיו נש"י מנ' ומקס'  
יהיה אולי דבר.

אכן, מה שמש'ן בהימ'אולי. אפשר לעשות עם כן וההפך.

נש"י אפשר ע"י כן מנ' סלקים נמנדים. עש"ש, הענ"ן יעבור עם  
עבורה ע"י אולי. עש"ש פונק' ריבועיות:

$$h_i(x, y) = a_i x^2 + b_i xy + c_i y^2 + \dots$$

נש"י פנימ' שפונק' ין עינאית ע"י עינאית צביה. זה מוס'ד מנדיים  
כאופן הבאו:

$$\Rightarrow X_4 = a_1 X_1 + b_1 X_2 + c_1 X_3 + \dots$$

$\downarrow$   $\downarrow$   $\downarrow$   
 $x^2$   $xy$   $y^2$

אכן אפשר לקבל קירוב  $\epsilon$ - extent עם כד'  $(1-\epsilon)$  גם עבור  
העל עינאית.

אפשר לעשות דברים עם הרבה יותר מורכבים.  
(2-ד מנדיים)  
עש"ש, נניח שיש ע"י נק' שמש' במחמה כפונק' של הענ"ן, ונניח  
שאולי נמצאים במקום מסוים והוציג עש"ש בכל הענ"ן והנ' הכי  
קרובה והכי רחוקה. "גן שרצה לבנות Coreset, מנ' נק' 'פולחן'  
שצדקה אמרוק וכן יהיו המאמ'ר בכל הענ"ן עש"ש גמול המנ' ונש"י  
אם הנק' עצמן כפונק' של הענ"ן, וע"י עינאית צביה נוס'ד מנדיים  
אולי עש"ש ...

$$\leftarrow p_1 X_1 + p_2 X_2 + \dots$$
$$(a_1 + b_1 t) X_1 + (a_2 + b_2 t) X_2 + \dots$$

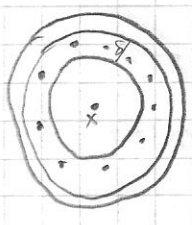
אולי היקואולי בקצה פי 2.

וקינול, עזר קינולן:  
p = n נק' ב-2 ממדים

חוצים לעם קינול בלולא עניינול שמסילא אל p.  
עזובי עניינול

בלולא, ממסיס לקרר אל p ע"ו בקור.

בהנחה שאין הנקודות הם בערך מ'צנול בקור  
בלצם חוצים למצול אל ממכז הבלול, וקמיליסי  
אלולא עקבילא עניינולסי אל:



$$\max_{p \in P} \|x-p\| - \min_{p \in P} \|x-p\|$$

אה-מה-מה, יש כולן שילרס ריבול' וזה ממכז. אלולא יולול עזמל  
פה גנול עינול!

מה עושים? מעלים בריבול':  
$$f_p = \|x-p\|^2 = \|x\|^2 - 2x \cdot p + \|p\|^2$$

בשיל' עינול ריבול', ממק' ריבול' נקבל  $x_{d+2}$  (ול קולול' גליל' הערק)  
וכל אל מילולסי ע- $R^{d+2}$  ונלסי עב שבלול  $f_p$  הילא עינולסי, וממילסי

$$\max_p \sqrt{f_p(x)} - \min_p \sqrt{f_p(x)}$$

וקר' חוצים עהיכלר מהמילסי. יש לול  $\frac{\epsilon \sqrt{H}}{n}$  עס מילולסי  $x_{d+2} = f_p(x)$   
נבנול G שהולו E-kernel (ול Coreset) עזר H.  
קולול' ע-QCP.

$$\Rightarrow \begin{aligned} A &= \max_{p \in P} f_p & \text{ב- } \vec{x} \text{ צנול' ב' } \\ B &= \min_{p \in P} f_p \\ C &= \max_{p \in Q} f_p \\ D &= \min_{p \in Q} f_p \end{aligned}$$

וביל' שולו Coreset קיבלול  
$$(1-\epsilon)(A-B) \leq C-D \leq A-B$$

אלול מה קולול עס המילולסי? וכלן  $\sqrt{|A|} - \sqrt{|B|} \leq \sqrt{|C|} - \sqrt{|D|}$  וליק נכלן. צצז הלול:

$$\sqrt{|C|} - \sqrt{|D|} = \frac{C-D}{\sqrt{|C|} + \sqrt{|D|}}, \quad \sqrt{|A|} - \sqrt{|B|} = \frac{A-B}{\sqrt{|A|} + \sqrt{|B|}} \Rightarrow \frac{\sqrt{|C|} - \sqrt{|D|}}{\sqrt{|A|} - \sqrt{|B|}} = \frac{C-D}{A-B}$$



$$\frac{BD}{CA} \approx \frac{\sqrt{C} + \sqrt{D}}{\sqrt{A} + \sqrt{B}} - \epsilon$$

כל עוד  $\epsilon$  קטן מספיק, נקבל  $(1 - \epsilon)(\sqrt{A} - \sqrt{B}) \leq \sqrt{C} - \sqrt{D}$

כלומר,  $\dots$