



g 1

Dimension  
Reduction

日期: Lecture 1: Singular Value Decomposition (SVD).

Consider data as vectors  $v = \begin{pmatrix} v_1 \\ \vdots \\ v_d \end{pmatrix} \in \mathbb{R}^d$ .

then,  $A \in \mathbb{R}^{n \times d}$ , 有两种看法: ① Matrix.

$$= \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix}.$$

②  $n$  points in  $d$ -dimension space.

Fact: SVD finds the "best-fitting"  $k$ -dim subspace of  $A$ .

Aim: "best"  $\Rightarrow$  ①  $\min_{S: k\text{-dim space}} \sum_{i=1}^n \text{dist}^2(a_i, S)$

$\Leftrightarrow$  ②  $\max_{S: k\text{-dim space}} \sum_{i=1}^n \text{length}^2(\text{proj}(a_i, S))$ .

范数:

$$\|x\|_p = \left( \sum_{i=1}^n |x_i|^p \right)^{\frac{1}{p}}$$

$$\|x\|_\infty = \max_{i=1}^n |x_i|$$

$$\|x\| := \|x\|_2.$$

Algorithm:

Consider  $k=1$ :  $\max_{\|v\|=1} \sum_{i=1}^n \text{length}^2(\text{proj}(a_i, \text{span}(v)))$ .

$$= \max_{\|v\|=1} \sum_{i=1}^n |k a_i, v|^2$$

$$= \max_{\|v\|=1} \|Av\|^2.$$

Define  $\sigma_1(A) = \max_{\|v\|=1} \|Av\|$

$v_1$  is "the first" singular vector if  $\|Av_1\| = \sigma_1(A)$ . (注: 不一定唯一).

Consider  $k=2$ :

在  $k=1$  基础上, 剔除  $v_1$  影响后继续  $\max_{\|v\|=1} \|Av\|$ .

$$\begin{aligned} &\|v\|=1, \\ &v \perp v_1 \end{aligned}$$

Define: the second singular value:  $\sigma_2(A) = \max_{\substack{\|v\|=1 \\ v \perp v_1}} \|Av\|$ .

"the second" singular vector:  $v_2 = \arg \max_{\substack{\|v\|=1 \\ v \perp v_1}} \|Av\|$ .

日期:

/

... (这里指从求 "best-fitting" 的空间维数).

Until:  $k=r+1$ . (或  $k=\overline{k+1}$ )

find  $\sigma_{r+1}(A)=0 \Rightarrow$  算法停止.

Problem Remained:

Definition

1. How to compute  $\arg \max_{\substack{\|V\|=1 \\ V \perp V_1, \dots, V_k}} \|AV\|$ .

2. Greedy's correctness.

Thm 1:  $A \in \mathbb{R}^{n \times d}$ . singular vectors  $v_1, \dots, v_r$ .

$\forall 1 \leq k \leq r$ , let  $V_k = \text{span}(v_1, \dots, v_k)$ .  $\Rightarrow$  Correctness

Then,  $V_k$ .  $V_k$  is the "best-fitting"  $k$ -dim subspace of  $A$ .

Pf:  $k=1$  时, 由前已证.

by induction: Suppose  $k=k-1$  成立.

$k=k$  时. Suppose  $W$  is the "best-fitting"  $k$ -dim space.

总可选出一组 orthonormal basis  $w_1, \dots, w_k$ . st.  $v_i \perp w_k, \dots, v_{k-1} \perp w_k$ .

由归纳,  $\text{span}(v_1, \dots, v_{k-1})$  is "best-fitting".

$$\|Aw_1\|^2 + \dots + \|Aw_k\|^2 \leq \|Av_1\|^2 + \dots + \|Av_{k-1}\|^2.$$

又由  $v_k$  选择,  $\|Aw_k\|^2 \leq \|Av_k\|^2$ .

$$\therefore \|Aw_1\|^2 + \dots + \|Aw_k\|^2 \leq \|Av_1\|^2 + \dots + \|Av_{k-1}\|^2. \quad \square$$

日期:

(Rmk: 将  $A$  作用在  $\text{span}(v_1, \dots, v_r)$  上, 在旋转之外, 相当于将  $v_1, \dots, v_r$  伸长  $\sigma_i(A)$  倍)  
(且  $A \perp (\text{span}(v_1, \dots, v_r))^{\perp} \Rightarrow A = \text{Span}(v_1, \dots, v_r)$ )  
 $\hookrightarrow \dim A = r$ .

Def: (Frobenius Norm / F-Norm).

$$\|A\|_F := \sqrt{\sum a_{ij}^2} \quad (= \sqrt{\sum_{i=1}^n \|\alpha_i\|^2})$$

Rmk:  $\|UAV^T\|_F = \|A\|_F$  (即旋转不影响 F-Norm).

Lem:  $\underbrace{\sum \sigma_i^2(A)}_{\sim} = \|A\|_F^2$ .

(Geometric Exp.:  $\|A\|_F^2 = \sum_{i=1}^n \|\alpha_i\|^2$ .  $\alpha_i \in \mathbb{R}^d$ . (标准基下的模长计算))

换 basis 不改变模长.

而  $\sigma_i(A) = \|AV_i\|$  即  $a_1, \dots, a_n$  在  $v_i$  方向上投影的平方和.

而  $\text{span}(v_1, \dots, v_r) = \text{span}(a_1, \dots, a_n)$ .

$\therefore \text{LHS}$  即在  $\{v_1, \dots, v_r\}$  基下计算  $a_1, \dots, a_n$  模长平方和.

$\therefore \text{LHS} = \text{RHS. } \square. )$

Def:  $A \in \mathbb{R}^{n \times d}$ . define the left singular vector as:

$$u_i = \frac{1}{\sigma_i(A)} AV_i$$

Lem:  $\{u_i\}_{i=1, \dots, r}$  两两正交.  $\{A$ 's singular value

Pf: 反证. let  $i$  be the smallest index. s.t.

$$\langle u_i, u_j \rangle \neq 0 \text{ for some } j > i.$$

日期:

WLOG suppose  $\langle u_i, u_j \rangle = \delta > 0$ .  
(不失一般性)

for some  $\varepsilon > 0$ , define  $v_i' = \frac{u_i + \varepsilon v_j}{\|u_i + \varepsilon v_j\|}$ . Then  $\|v_i'\| = 1$ .

①  $\|Av_i'\| \geq u_i^T A v_i'$ .

$$> \frac{u_i^T G_i u_i + u_i^T \varepsilon g_j u_j}{\sqrt{1+\varepsilon^2}}$$

$$\begin{aligned} &= \frac{g_i + g_j \cdot \varepsilon \cdot \delta}{\sqrt{1+\varepsilon^2}} \\ &\xrightarrow{\frac{1}{\sqrt{1+\varepsilon^2}} \geq 1 - \frac{\varepsilon^2}{2}} \geq (g_i + g_j \cdot \varepsilon \cdot \delta) (1 - \frac{\varepsilon^2}{2}) > g_i. \end{aligned}$$

②  $i=1$  时,  $\|Av_1'\| \stackrel{\text{①}}{>} g_1 = \max_{\|v\|=1} \|Av\|$ . 矛盾.

$i > 1$  时,

by suppose  
⇒ 由  $v_i' = \frac{u_i + \varepsilon v_j}{\|u_i + \varepsilon v_j\|}$  is orthogonal to  $v_1, \dots, v_{i-1}$ .

$$\therefore \|Av_i'\| \leq g_i. \text{ 矛盾. } \square.$$

Thm:  $A = \underbrace{\sum_{i=1}^r g_i u_i v_i^T}_{\text{(SVD)}}$

Pf: for  $\forall i$ ,  $Av_i = g_i u_i$ .

(Rmk: for  $\forall v \in R^d$ ,  $Av = Bv$ , then  $A = B$ .)

∴ for  $\forall v \in R^d$ , suppose  $v = \alpha_1 v_1 + \dots + \alpha_r v_r + \dots + \alpha_d v_d$ .

$$\therefore Av = \sum_{i=1}^r \alpha_i Av_i = \sum_{i=1}^r \alpha_i g_i u_i.$$

$$\text{而 } (\sum_{i=1}^r g_i u_i v_i^T)v = \sum_{i=1}^r g_i u_i v_i^T v_i \cdot \alpha_i = \sum_{i=1}^r \alpha_i g_i u_i.$$

$$\therefore A = \sum_{i=1}^r g_i u_i v_i^T. \quad \square.$$

日期:

Cor1:  $U = \underbrace{(u_1, \dots, u_r)}_{\in \mathbb{R}^{n \times r}}, V = \underbrace{(v_1, \dots, v_r)}_{\in \mathbb{R}^{d \times r}}, D = \underbrace{\begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r \end{pmatrix}}_{\in \mathbb{R}^{r \times r}}$ .

$$\Rightarrow A = UDV^T$$

用V<sub>d</sub>替换D<sub>r</sub>，  
Cor2:  $U = \underbrace{(u_1, \dots, u_r, \underbrace{u_{r+1}, \dots, u_n}_{\in \mathbb{R}^{n \times n}})}_{\in \mathbb{R}^{n \times n}}, V = \underbrace{(v_1, \dots, v_r, \underbrace{v_{r+1}, \dots, v_d}_{\in \mathbb{R}^{d \times d}})}_{\in \mathbb{R}^{d \times d}}, D = \underbrace{\begin{pmatrix} \sigma_1 & & \\ & \ddots & \\ & & \sigma_r & 0 \end{pmatrix}}_{\in \mathbb{R}^{n \times d}}$

$$\Rightarrow A = UDV^T$$

## I. Singular vectors & eigenvectors.

Suppose  $A \in \mathbb{R}^{n \times d}$ .

① If A has SVD:  $A = \sum_{i=1}^r \sigma_i u_i v_i^T$

then  $\underbrace{Av_i}_{\sigma_i u_i}, \underbrace{A^T u_i}_{\sigma_i v_i}$

Pf: 由  $A^T = \sum_{i=1}^r \sigma_i v_i u_i^T \therefore A^T u_i = \sigma_i v_i \quad \square$

( Rmk: 类似于 eigenvectors  $\alpha$  与对应 eigenvalue  $\lambda$  的性质

$$A\alpha = \lambda \alpha \text{ (for square } A\text{). }$$

② Let  $B = A^T A$ , then  $v_i$  is eigenvectors of B,

对应 eigenvalue is  $\sigma_i^2$ .

同理,  $AA^T$ , then  $u_i$  is eigenvectors of B,

对应 eigenvalue is  $\sigma_i^2$ .

Pf:  $Bu_i = (\sum \sigma_j u_j u_j^T)(\sum \sigma_k v_k v_k^T)u_i = (\sum \sigma_j^2 u_j u_j^T)u_i = \sigma_i^2 u_i \quad \square$

日期:

## II. low rank approximation.

$$A_{n \times d} \leftarrow A'_{n \times d} \text{ (low rank).}$$

e.g. customer-movie.

hypothesis:  $\exists k \uparrow$  underlying factors for  $a_{ij}$ .

$\Rightarrow$  low rank.

$$A \approx \left( \begin{matrix} x_{ij} \\ \vdots \\ x_{ik} \end{matrix} \right) \left( \begin{matrix} y_{ij} \\ \vdots \\ y_{ik} \end{matrix} \right)_{k \times d}.$$

customer-factors    factors-movie  
(客户矩阵)    (属性矩阵)

Consider SVD:  $A = \underbrace{U}_{X} \underbrace{V^T}_{Y^T}$  (rank: 可能出现负元 (但可忍受))

Goal: find  $B$  with rank  $k$ , s.t.  $B$  is close to  $A$

(体现为优化问题: e.g. s.t.  $\|A-B\|_F$  or  $\|A-B\|_2^2$  is minimal.).

Def: def  $A_k$  as: rows of  $A_k$  are projections ( $r(A_k)=k$ )

of rows of  $A$  onto the subspace  $V_k = \text{span}\{v_1, \dots, v_k\}$

$$\text{即. } \underline{A_k} = \left( \underline{a_i} \sum_{j=1}^k \underline{v_j} \underline{v_j^T} \right) \xrightarrow{\sum \langle a_i, v_j \rangle v_j^T \cdot (\text{单行的 proj.})}.$$

$$= A \sum_{j=1}^k v_j v_j^T$$

$$= (\sum_{i=1}^n g_i u_i u_i^T) (\sum_{j=1}^k v_j v_j^T)$$

$$= \underbrace{\sum_{i=1}^n g_i u_i u_i^T}_{\sim}$$

日期:

By above,  $\|A - A_k\|_F^2 = \sum_{i=1}^n \text{dist}^2(a_i, V_k) = \min_{\substack{S: \text{space.} \\ \dim(S)=k}} \sum_{i=1}^n \text{dist}^2(a_i, S)$ .

Thm1:  $A_k$  is the best rank approximation wrt. F-Norm.  $\rightarrow$  with respect to.

即.  $\forall B$  with  $\text{rank}(B) \leq k$ ,  $\|A - A_k\|_F \leq \|A - B\|_F$ .

Pf: Suppose  $B \in \mathbb{R}^{n \times d}$  s.t.  $\|A - B\|_F^2 \min$ .  $\rightarrow$  geometric.

设  $V$  为  $\text{span}(\text{rows of } B)$ .

then  $\|A - B\|_F^2 = \sum_{i=1}^n \text{dist}^2(a_i, V) \geq \|A - A_k\|_F^2$ .  $\square$ .

Def: (2-Norm):  $\|A\|_2 = \max_{\|x\| \leq 1} \|Ax\|$

(Note that 1°  $\|A\|_2 = \sigma_1(A)$  2°  $\|A\|_2$  与 向量  $v$  的  $\|v\|_2$  定义相容.)

Thm2:  $A_k$  is the best rank approximation wrt. 2-Norm.

即.  $\forall B$  with rank  $k$ ,  $\|A - A_k\|_2 \leq \|A - B\|_2$

Lem:  $\|A - A_k\|_2^2 = \sigma_{k+1}^2$

Note that,  $A - A_k = \sum_{i=k+1}^r \sigma_i u_i v_i^T$ .

For  $\forall v \in \mathbb{R}^d$  不妨设  $v = \sum_{j=1}^r \alpha_j v_j$ . ( $\alpha_j \in \mathbb{R}$ )

then.  $(A - A_k)v = \sum_{i=k+1}^r \alpha_i \sigma_i u_i$

$\therefore \|(A - A_k)v\| = \left\| \sum_{i=k+1}^r \alpha_i \sigma_i u_i \right\| = \sqrt{\sum_{i=k+1}^r \alpha_i^2 \sigma_i^2}$

日期:

let  $v^*$  is the unit vector maximizing above.

$$\therefore \sum_{i=1}^r \alpha_i^2 = 1.$$

$$\lambda: 6_{k+1} \geq \dots \geq 6_r$$

$\therefore v^*$  satisfy  $\alpha_{k+1} = 1$ ,  $\alpha_i = 0$  ( $i \neq k+1$ )

$$\therefore \|A - A_k\|_2 = \|(A - A_k)v^*\| = 6_{k+1}. \quad \square$$

Thm-Pf: If  $\text{rank}(A) \leq k$ . then  $\|A - A_k\|_2 = 0$ .

Assume  $\text{rank}(A) > k$ .  $\therefore \|A - A_k\|_2^2 = 6_{k+1}^2$

Consider the null space of  $B$ , i.e.  $\text{Null}(B) = \{v : Bv = 0\}$ .

Note that  $\dim(\text{Null}(B)) \geq d - k$ . ( $\because \text{rank}(B) \leq k$ ).

then  $\exists$  vector  $z \neq 0$  such that  $z \in \text{Null}(B) \cap \text{Span}(v_1, \dots, v_{k+1})$

$$\therefore \|z\| = 1.$$

Note that  $\|A - B\|_2^2 \geq \|(A - B)z\|_2^2 \downarrow z \in \text{Null}(B)$ .

$$= \|Az\|_2^2$$

$$\stackrel{z \in \text{Span}(v_1, \dots, v_k)}{\downarrow} = \left\| \sum_{i=1}^{k+1} 6_i (v_i^T z) v_i \right\|_2^2$$

$$= \sum_{i=1}^{k+1} 6_i^2 (v_i^T z)^2$$

$$\geq 6_{k+1}^2 \underbrace{\sum_{i=1}^{k+1} (v_i^T z)^2}_{= \|z\|^2}$$

$$= 6_{k+1}^2. \quad \square.$$

(Rmk: If  $A_k$  is best approximation over 2-form & F-Form.)

日期:

## III. Compute SVD:

1. Exact Computation:  $O(nd^2)$ .  $\rightarrow$  太大, 考虑近似的.

2. Approximate SVD:

1. Start with  $v_1$ : (Power Method).  $\rightarrow$  通过 power method 得到影响.

$$\textcircled{1} \quad B = A^T A = (\sum \sigma_i v_i u_i^T) (\sum \sigma_j u_j v_j^T) = \sum \sigma_i^2 v_i u_i^T.$$

$$B^2 = (\sum \sigma_i^2 v_i u_i^T) (\sum \sigma_i^2 v_i u_i^T) = \sum \sigma_i^4 v_i u_i^T.$$

$$\dots \\ B^k = \sum \sigma_i^{2k} v_i u_i^T.$$

If  $\sigma_1 > \sigma_2$ , &  $k$  large enough, then,  $B^k \rightarrow \sigma_1^{2k} v_1 u_1^T$ .

$\Rightarrow B^k$ 's 1st 列  $b \approx c v_1$ . ( $c$  is Const.).

$\Rightarrow$  Output  $\frac{b}{\|b\|}$  as  $v_1$ .

② Compute  $B^k$  too expensive!

Idea: matrix  $\cdot$  matrix  $O(d^3) \rightarrow$  matrix  $\cdot$  vector  $O(d^2)$ .

Choose a random vector  $x \in \mathbb{R}^d$ .

then compute  $B^k x = B^{k-1} (Bx)$ .

Then,  $B^k x \approx c v_1$

[直观-Pf: Assume  $x = \sum_{i=1}^d \alpha_i v_i$ .

$$\Rightarrow B^k x = (\sum \sigma_i^{2k} v_i u_i^T) (\sum \alpha_i v_i) = \sum \sigma_i^{2k} \alpha_i v_i$$

$$b \rightarrow \sigma_1^{2k} \alpha_1 v_1 = c' v_1.$$

$\sigma_1 > \sigma_2$ ,  
 $k \leq j \leq k$ .

日期: /

严格地有 Thm:

→ P2.

Let  $A \in \mathbb{R}^{n \times d}$ .  $x$  is a vector s.t.  $\|x\|=1$  &  $|x^T v_i| \geq \delta$ .

Let  $V$  be the space spanned by right singular vectors of  $A$ .

→  $\text{proj}_V = \text{span}(v_i | \sigma(v_i) \rightarrow \sigma(v_i))$ .

corresponding to singular values  $\sigma_i$  s.t.  $\sigma_i \geq (1-\epsilon)\sigma_1$ .

Let  $w$  be normalized vector after  $k = \frac{\ln(1/\epsilon)}{2\epsilon}$  iterations

of power method. i.e.  $w = \frac{B^k x}{\|B^k x\|}$ .

→ 4.7-4c

Then  $w$  has a component of at most  $\epsilon$  perpendicular to  $V$ ,

垂直于  $V$

i.e.  $\|w - \text{proj}(w, V)\| \leq \epsilon$ .

(Rmk: 即若  $\epsilon \rightarrow 0$ , 则  $w \rightarrow V$ , 而  $V \approx \text{span}(v_i | \sigma(v_i) \rightarrow \sigma(v_i))$ ).

(Rmk: PF method 不是 spectral cluster 等中经常使用.)

PF: let  $A = \sum_{i=1}^r \sigma_i u_i v_i^T = \sum_{i=1}^d \sigma_i u_i v_i^T$ ,  $\# \sigma_{r+1} = \dots = \sigma_d = 0$ .

$$x = \sum_{i=1}^d \alpha_i v_i$$

$$\text{then } B^k x = \sum_{i=1}^d \sigma_i^{2k} \alpha_i v_i.$$

Suppose  $\sigma_1, \dots, \sigma_m \geq (1-\epsilon)\sigma_1$ ,  $\sigma_{m+1} \leq (1-\epsilon)\sigma_1$ .

$$\text{Consider: } \|B^k x\|^2 = \sum_{i=1}^d \sigma_i^{4k} \alpha_i^2$$

$$\geq \sigma_1^{4k} \alpha_1^2 = \sigma_1^{4k} \cdot \sigma^2.$$

$$\leq \sum_{i=1}^d \alpha_i^2 = \|x\|^2 = 1$$

$$\text{then, } \|w - \text{proj}(w, V)\|^2 = \frac{\left\| \sum_{i=m+1}^d \sigma_i^{2k} \alpha_i v_i \right\|^2}{\|B^k x\|^2} \leq \frac{(1-\epsilon)^4 \sigma_1^{4k} \sum_{i=m+1}^d \alpha_i^2}{\sigma_1^{4k} \sigma^2} = \frac{(1-\epsilon)^4 k}{\sigma^2}.$$

日期:

we want  $\frac{(1-\varepsilon)^{4k}}{\delta^2} \leq \varepsilon^2$ .

$$\Rightarrow k \geq \frac{\ln(1/\varepsilon\delta)}{2\varepsilon}. \quad \square.$$

Rmk: 取过 larger  $\delta$  s.t.  $k$  不过 for some  $\varepsilon$ .

↳ How to find proper  $x$ :

1) Take  $k$  samples  $\sim N(0, I)$ :  $x_1, \dots, x_d \in \mathbb{R}^d$ .

2) Normalize  $x$ . s.t.  $\|x\|=1$ .  $\Rightarrow$  random projection.

(以较大概率 st.  $|x^\top v_i| \geq \delta$ .)

2. Approximate first  $k$  singular vectors:

Method 1:

1) select random vector  $r$ .

2) define  $u_1, \dots, u_k, v_1, \dots, v_k$ . as:

$$u_1 \leftarrow r, u_2 \leftarrow Br, \dots, u_k \leftarrow B^{k-1}r.$$

$v_1, \dots, v_k \leftarrow u_1, \dots, u_k$  施密特正交化.

3) Repeat.

update  $u_1, \dots, u_k, v_1, \dots, v_k$ . as:

$$u_1 \leftarrow v_1, u_2 \leftarrow Bv_2, \dots, u_k \leftarrow Bv_k.$$

$v_1, \dots, v_k \leftarrow u_1, \dots, u_k$  施密特正交化.

日期:

### Method 2:

1) select  $d \times k$  matrix  $V$  st.  $V^T V = I_k$ . define  $Z \leftarrow 0$ .

2) Repeat:

update  $Z \leftarrow BV$ .

find  $Z$  的 QR decomposition. ( $Q$ : 正交阵,  $R$ : 上三角).

update  $V \leftarrow Q$ .

### Method 3:

1) select  $w_1$  as above.

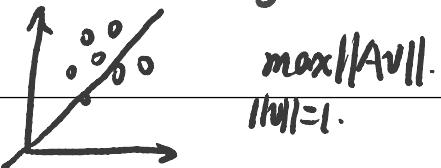
2) Consider  $A - \tilde{\sigma}_1 \tilde{u}_1 w_1$ :

select  $w_2$  as above.

Repeat.

## IV. Principal Component Analysis (PCA)

之前考慮的是



best-fitting  $k$ -dim subspace.

Now Consider:

best-fitting  $k$ -dim affine subspace.

$$\text{即 } \max_{\|v\|=1} \| [A - \begin{pmatrix} \cdot & \cdot \\ \cdot & \cdot \end{pmatrix} \cdot \vec{r}] v \|.$$

Fact: best-fitting  $k$ -dim affine subspace 通过它们的中心

日期:

## PCA Alg:

1) Input  $A \in \mathbb{R}^{n \times d}$ . ( $i \geq A = (a_{ij})$ .)

2)  $\bar{a}_{ij} = \frac{1}{n} \sum_{i=1}^n a_{ij}$ .  $A' = A - \begin{pmatrix} 1 \\ \vdots \\ 1 \end{pmatrix} (\bar{a}_1, \dots, \bar{a}_d)$ . (中心化)

3) Compute  $A'$ . SVD.

↓  
中心化

↓  
低秩近似(SVD)

## V. Clustering mixtures of Gaussian distributions.

↳ 基本分布的带权求和，本身也是一种分布。

$$f = w_1 p_1 + \dots + w_k p_k. \quad w_i > 0, \sum w_i = 1.$$

Problem: find  $w_i$ ,  $\mu(p_i)$  ( $i=1, \dots, k$ ),  $\sigma$  to fit  $n$  i.i.d. (f). samples.

(恢复类  $\rightarrow$  每类分别恢复 Gaussian).

Assume  $p_1, \dots, p_k$  spherical Gaussian (但  $\neq \mu, \sigma^2$ )

TASK: ① Cluster the set of samples into  $k$  clusters  $C_1, \dots, C_k$ .

② for  $\forall$  clusters  $C_i$  求  $p_i$ .

Alg: ① for  $\forall$  sample pair  $\|x-y\|^2$ .

② put  $x, y$  in same cluster if  $\|x-y\|^2 \leq T$ .

## 2 Case:

i. if 2 samples  $x, y$  from same spherical Gaussian with  $\sigma^2$ .

日期:

$$\text{then } \sup \|x-y\| \approx 2(\sqrt{d} \pm o(1))^2 \cdot \sigma^2 \quad (=A)$$

ii. if 2 samples  $x, y$  from diff. spherical Gaussian with  $\sigma^2, \Delta$ .

$$\text{then } \inf \|x-y\| \approx 2(\sqrt{d} \pm o(1))^2 \cdot \sigma^2 + \Delta^2 \quad (=B).$$

$\Rightarrow$  要求  $B > A$ . (否则无法有效分类).

$$\Rightarrow \Delta > C \cdot d^{1/4} \cdot \sigma. \quad (C \text{ is some Const}).$$

$\downarrow$   
与  $d$  有关, 可能太强]. (考虑降维, s.t. 只要  $\Delta > C \cdot \sigma$  即可).

An SVD-based alg.: (Input:  $n$  samples w.p.d. f.).

① Compute  $A$ 's SVD to find first  $k$  singular vectors,  
project  $n$  samples to  $\text{span}(v_1, \dots, v_k)$ .

② Cluster them in  $\text{span}(v_1, \dots, v_k)$ .

Thm: The above alg. can separate  $k$  different Gaussians  
with  $\Delta > ck^{1/4}\sigma$ , whp. (不证)

一些直观: ① the set of samples can fit prob. density f.

② for the prob. density f, the best-fitting subspace contains  
the  $k$  centers. ?

(Rmk:  $v_k = \operatorname{argmax}_{V: k\text{-dim}} E \|\operatorname{proj}(x, V)\|$ ).

(Cor:  $\mu$  在投影下不变).

日期:

③ the projection of a spherical Gaussian remains the variance  $\sigma^2$ . ?

$\Rightarrow \mu_i$  与  $\sigma^2$  都不变.

$\Rightarrow$  从而只要  $\Delta > C \cdot k^{1/4} \cdot \sigma$ .

## Lecture 2. Random Projection & JL Lemma.

Thm 1: JL Lemma:

$\forall 0 < \epsilon < 1$ .  $n > 0 \in \mathbb{Z}$ . let  $k = O\left(\frac{C \cdot \ln n}{\epsilon^2}\right)$  (即  $k$  不太小). for some Const.  $C$ .

$\forall$  set of  $n$  points  $\{a_1, \dots, a_n\}$  in  $\mathbb{R}^d$ . Define random projection:

$$f: \mathbb{R}^d \rightarrow \mathbb{R}^k \quad f(a) = \frac{1}{\sqrt{k}}(u_1^T a, \dots, u_k^T a)^T. \quad (= \frac{1}{\sqrt{k}} F \cdot a, F = \begin{pmatrix} u_1^T \\ \vdots \\ u_k^T \end{pmatrix} \in \mathbb{R}^{k \times d})$$

where  $u_i$  is a  $d$ -dim vectors.  $\sim$  Gaussian  $(0, 1)$ . at each coordinate.

Then  $wp. \geq 1 - \frac{c}{n}$ . (for some fixed const  $c$ )  $\forall (a_i, a_j)$  (注意: 整体 Prob.).

$$(1-\epsilon) \|a_i - a_j\| \leq \|f(a_i) - f(a_j)\| \leq (1+\epsilon) \|a_i - a_j\|. \quad \text{whp. dist}(a_i, a_j) \text{ 变化不大.}$$

Thm 2: Random-Projection Thm

$\exists$  Const  $C > 0$ . s.t. for  $\epsilon \in (0, 1)$ .  $\forall a$ ,

$$P[\|F(a)\| - \sqrt{k} \|a\| \geq \epsilon \sqrt{k} \|a\|] \leq 3e^{-ck\epsilon^2}. \quad \text{whp. } \|a\| \text{ 变化不大.}$$

日期:

$\text{Thm 2} \Rightarrow \text{Thm 1}$ : 取  $k = \frac{3(n\epsilon)}{C\epsilon^2}$ ,  $a = a_i - a_j$  代入 Thm 2:

$$\text{s.t. } \Pr \left[ \left| \frac{1}{\sqrt{k}} \|F(a_i - a_j)\| - \|a_i - a_j\| \right| \geq \epsilon \|a_i - a_j\| \right] \leq \frac{3}{n^3}$$

by union bound.

w.p.  $1 - \frac{3}{2n}$  有 Thm 1 结论.

$\text{Thm 3: Gaussian Annulus Thm.}$

For d-dim spherical Gaussian with variance 1 各分量.  
( $\beta$  不)

for  $\forall \beta \leq \sqrt{d}$ . 至少  $1 - 3e^{-c\beta^2}$  of the prob. mass lie in the annulus

$\sqrt{d} - \beta \leq \|x\| \leq \sqrt{d} + \beta$ . for fixed const.  $c > 0$ .  
 $\approx \frac{1}{96}$ ?

(Rmk: 1° for d-dim spherical Gaussian, 大部分 mass 集中于中心. ( $\sqrt{d}$  两边).)

2° for d-dim ball,  $d \rightarrow \infty$ , mass 集中于表面.) whp. 在中心.

$$\forall \beta \geq 1, \Pr[\|x\| \leq \sqrt{d} + \beta] \geq 1 - 3e^{-c\beta^2}.$$

$$\text{s.t. } 1 - \beta \leq \|x\| \leq 1 + \beta.$$

fix  $\beta$ .  $\sqrt{d} \rightarrow \infty$ . whp.  $\|x\| = 1 \leq \beta$ . 即 mass 集中于表面.

$\text{Thm 3} \Rightarrow \text{Thm 2}$ : WLOG assume  $\|a\|=1$ .

$$\text{由 } a \cdot u_i = \sum_{j=1}^d a_j u_{ij}, \quad u_{ij} \sim N(0, 1).$$

$$\therefore a \cdot u_i \sim N(0, \sum a_j^2 \text{Var}(u_{ij})) = N(0, \sum a_j^2) = N(0, 1).$$

$\Rightarrow F(a) \sim k\text{-dim spherical Gaussian with variance 1 各分量.}$

by Thm 3:  $\sum \beta = \sum \epsilon_j \sqrt{k}$ . then.

$$(1 - \epsilon) \sqrt{k} \leq \|F(a)\| \leq (1 + \epsilon) \sqrt{k} \quad \text{w.p.} \leq 3e^{-c k \epsilon^2}. \quad (\text{注意: } \|a\|=1).$$

日期:

Thm 3 - Pf: (when  $\beta \leq \frac{\sqrt{d}}{5}$ ).

Assume  $X = (X_1, \dots, X_n)$ .  $X_i \sim N(0, 1)$ . let  $r = \|X\|$ .

$$\sqrt{d} - \beta \leq \|X\| \leq \sqrt{d} + \beta \Leftrightarrow |\sqrt{d} - r| \leq \beta.$$

$$\therefore \text{Thm 3} \Leftrightarrow P[|\sqrt{d} - r| \geq \beta] \leq 3e^{-C\beta^2}.$$

$$= P[|\sqrt{d} - r| \geq \beta(r + \sqrt{d})]$$

↓ 放缩.

$$\text{要证 } P[|\sqrt{d} - r| \geq \beta\sqrt{d}] \leq 3e^{-C\beta^2}.$$

$$\text{由 } r^2 = \|X\|^2 = \sum_{i=1}^d X_i^2 \sim \chi^2(d).$$

$$(\text{Fact: } g_{\chi^2(d)}(t) = (1-2t)^{-\frac{d}{2}}. \quad \boxed{g_{\chi^2(d)} = (1-2t)^{-\frac{1}{2}}}).$$

$$P[r \geq d + \beta\sqrt{d}] = P[e^{\lambda r^2} \geq e^{\lambda(d + \beta\sqrt{d})}] \stackrel{\text{Markov}}{\leq} \frac{E[e^{\lambda X^2}]}{e^{\lambda(d + \beta\sqrt{d})}}$$

↓ positive.  $(\lambda > 0)$ .

~~(趁手数上用 Markov)~~  $= \frac{\prod_{i=1}^d E[e^{\lambda X_i^2}]}{e^{\lambda(d + \beta\sqrt{d})}}.$

$$\begin{aligned} &= \frac{(1-2\lambda)^{-\frac{d}{2}}}{e^{\lambda(d + \beta\sqrt{d})}} \\ (\ln(1-\lambda) &\geq x - \frac{x^2}{2} - \frac{x^3}{2}, \forall 0 < x < \frac{1}{2}) \\ &= e^{-\frac{d}{2}\ln(1-2\lambda) - \lambda(d + \beta\sqrt{d})} \\ &\leq e^{d(\frac{1}{2} + 2\lambda^2) - \lambda\beta\sqrt{d}}. \quad (\text{if } \lambda > 0). \end{aligned}$$

$$\text{取 } \lambda = \frac{\beta}{2\sqrt{d}}, \text{ 即有 } P[X \geq d + \beta\sqrt{d}] \leq e^{-\frac{\beta^2}{5}}$$

$P[X \leq d - \beta\sqrt{d}]$  类似可证.  $\square$ .

(或者用 Master Tail Bound 证).

日期:

Rmk On JL Lemma:

① JL Lemma 保证所有 pairs of  $a_i, a_j$ . (w.h.p.)

② 可降维至  $k = O(\frac{\ln n}{\varepsilon^2})$ . (与  $d$  无关).

(data set 另会影响正确率).

③ Random Projection function  $f$  与 data set  $(x_i)$  无关.

④ Running Time:  $O(n \cdot d \cdot k) \approx O(n \cdot d \cdot \frac{\ln n}{\varepsilon^2})$ .

⑤  $f: \mathbb{R}^d \rightarrow \mathbb{R}^k$ . bijection. w.h.p.

More Extension:

① [Ailon, Chazelle]. gives a new projection  $f$ . s.t.

each projection  $f(x)$  takes  $O(d \cdot \log d + \frac{(\ln n)^{1+o(1)}}{\varepsilon^2})$ .

② [Achlioptas] Gaussian  $\rightarrow \pm 1$  r.v.

即令  $f_{ij} = \begin{cases} 1 & \text{w.p. } \frac{1}{2} \\ -1 & \text{w.p. } \frac{1}{2} \end{cases}$  或者  $f_{ij} = \begin{cases} 1 & \text{w.p. } \frac{1}{6} \\ 0 & \text{w.p. } \frac{2}{3} \\ -1 & \text{w.p. } \frac{1}{6} \end{cases}$

(rmk: 1/3-化因子须对应修改).

③ [Larsen, Nelson].

$O(\frac{\ln n}{\varepsilon^2})$  是  $k$  的 lower bound. (即不可再改进).

④ [Nelson] Strong JL Lemma:

for all  $x \in A$ , for all  $y \in R^d$ , (JL: for all  $x, y \in A$ ).

$(1-\varepsilon) \|x-y\| \leq \|f(x)-f(y)\| \leq (1+\varepsilon) \|x-y\|$ .

日期:

## Lecture 3. Nearest Neighbor Search (NNS).

Input: a.  $n$  points ( $\in \mathbb{D}^d$ ) database.  
 $\hookrightarrow$  (domain:  $\mathbb{R}/\mathbb{Z}/[r]/\dots$ ).

Query:  $x$ .

Problem: find nearest point  $a^*$  to  $x$ . i.e.  $a^* = \underset{a \in \text{database}}{\operatorname{argmin}} d(x, a)$ .

(Rmk: Hamming distance: (for 0-1 bits).

the number of positions with different values. ).

$$(\text{dist.} \Leftrightarrow \|x-y\|_1 = \sum |x_i - y_i|)$$

Solution:

① Preprocessing: process database and build a data structure.

② Query: upon query  $x$ , answer  $a^*$ .

Goal: ① query time: 尽量快 (ideally  $\text{poly}(\log n)$ ).

② used space for data structure: 尽量小.

③ 预处理时间可接受 ( $\text{poly}(n, d)$ ).

Exact Solution:

Curse of dimensionality:

if beats  $O(d \cdot n)$  query time requires  $\Omega(2^d)$  space.

(Rmk: Exact solution isn't proper!).

日期:

Approximate Solution:

Problem 1: C-Approximate Nearest Neighbor Search. (C-ANN)

Goal: Output  $a^*$  s.t.  $d(x, a^*) \leq C \cdot \min_i d(x, a_i)$ . ( $C > 1$ )

Problem 2: (C, r)-ANN.

Promise:  $\exists a_0$  s.t.  $d(x, a_0) \leq r$ . then.

Goal: Output  $a_i$  s.t.  $d(x, a_i) \leq C \cdot r$ . (否则 output no).

等价搜索(c).

From (C, r)-ANN. to C-ANN:

若 no, 则直接 output  $a_i$ .

此时,  $d(x, a_i) \geq \frac{1}{C-1} D_{\max}$ .  $\forall i \Rightarrow (C-1)d \geq D_{\max}$ .

有  $d(x, a_i) \leq d(x, a_j) + d(a_i, a_j) \leq c \cdot d$ .

即  $D_{\max} \leq d$ .

① Preprocess:  $D_{\max} = \max(a_i, a_j)$ .  $D_{\min} = \min(a_i, a_j)$ .  $R = \frac{D_{\max}}{D_{\min}} \cdot \frac{1}{C-1}$ .

② Solve (C, r)-ANN for  $r_0 = D_{\min}$ ,  $r_1 = C \cdot D_{\min}$ , ...,  $r_k = \frac{1}{C-1} \cdot D_{\max}$ . (G.P. Search).

Method 1: (利用 Hamming distance:  $\text{Ham}(x, y)$ ).

Thm:  $\exists$  a data structure uses space  $n^{O(1)}$ , query time  $O(d \cdot \ln n)$ .

solves (c, r)-ANN w.h.p.

Alg: (Use dimension reduction).

1) Parameters:  $k = \frac{\log n}{(\frac{1}{8} - 2^{-C-2})^2}$ .  $p = \frac{1}{2} - (\frac{1}{2})^{1+\frac{1}{r}}$ .  $S = (\frac{3}{8} - 2^{-(C+2)}) \cdot k$ .

2) Sample a  $k \times d$  matrix  $U \in \mathbb{F}_2^{k \times d}$ :

日期:

四路: 在低秩上单步枚举 query. ( $2^k$ ).

$$U_{ij} = \begin{cases} 1 & \text{w.p. } p \\ 0 & \text{w.p. } 1-p \end{cases} \quad \text{i.i.d.} \quad \Rightarrow \text{JL lemma's variant}$$

3) for each  $i \in n$ , compute:

$$z_i = U \cdot a_i \in \{0,1\}^k.$$

4) store all  $s$ -balls:  $B_s(u) = \{a_i \mid \text{Ham}(u, z_i) \leq s\}$ . for  $\forall u \in \{0,1\}^k$ .  
 $(\# 2^{k+1})$ .

5) Upon a query  $x \in \{0,1\}^d$ :

compute  $Ux$ :

$\begin{cases} \text{if } B_s(Ux) \neq \emptyset, \text{return any } a_i \in B_s(Ux). \\ \text{if } B_s(Ux) = \emptyset, \text{return no.} \end{cases}$

Analysis:

$$\Pr[Ux \neq Uy] = \Pr[x-y \neq 0].$$

① Space:  $2^k$  balls. 每个  $n \times d$ :  $O(nd \cdot 2^k) = n^{O(1)}$ .

② Time:  $O(k \cdot d) = O(d \cdot \ln n)$

③ Correctness:

Lemma:

$$\forall x, y \in \{0,1\}^d. \begin{cases} \text{Ham}(x, y) \leq r \Rightarrow \Pr[\text{Ham}(Ux, Uy) \geq s] < e^{-\Omega(rk)}. \\ \text{Ham}(x, y) \geq c \cdot r \Rightarrow \Pr[\text{Ham}(Ux, Uy) < s] < e^{-\Omega(rk)}. \end{cases}$$

Pf:  $\text{Ham}(Ux, Uy) = \sum_{i=1}^k x_i$ . where  $x_i = \begin{cases} 1, & (Ux)_i \neq (Uy)_i \\ 0, & (Ux)_i = (Uy)_i \end{cases}$ .

Lem:  $\Pr[(Ux)_i \neq (Uy)_i] = \frac{1}{2} (1 - (1-2p)^{\text{Ham}(x, y)}) = \frac{1}{2} (1 - (\frac{1}{2})^{\text{Ham}(x, y)/r}).$   
(由递推)

日期: /

Apply Chernoff Bound:

$$P[\text{Ham}(U \times U) \geq s] = P\left[\sum_{i=1}^k X_i \geq s\right]$$

$$\leq P\left[|\sum_{i=1}^k X_i - pk| \geq s - pk\right]$$

$$\leq 3e$$

对称 Method: 障碍有限 ( $\text{clog}n$ ). 精度 wh.p. 一次到位 (或找  $2^k$  query) w.h.p.  $dn^{O(1)}$  space  
日期 Method 2: / 障碍任意. 精度常数. 4 次到位 (或建  $n^{\frac{1}{2}}$  的 hash 表). w.p.  $\geq 0.6$   $n^{1+\frac{1}{d}}$  space

Me 路径: LSH + 查询至多  $L$  次, +LT.

Method 2. (基于 LSH: Locality Sensitive Hashing). (针对 Hamming).

Goal:  $\exists$  a randomized data structure solves  $(c, r)$ -ANN w.p.  $\geq 0.6$ .

Space:  $O(n^{1+\frac{1}{d}} + nd)$  query time:  $O(dn^{1/c})$ .

Def: (LSH): random function:  $h: \{0,1\}^d \rightarrow \{0,1\}^k$  is  $(p_1, p_2)$ -LSH. if:

1) if  $\text{Ham}(x, y) \leq r$ . then  $P[h(x)=h(y)] \geq p_1$  ( $p_1$  not so small).

2) if  $\text{Ham}(x, y) \geq c \cdot r$ . then  $P[h(x)=h(y)] \leq p_2$ . ( $p_2$  small).

Def: (Hash family  $\mathcal{H}$ ) ( $g \in \mathcal{H}$  is called LSH)

$\mathcal{H} = \{g: \{0,1\}^d \rightarrow \{0,1\}^k\}$ . for  $p \in \{0,1\}^d$ .  $g(p) = \{p_{i1}, \dots, p_{ik}\}$ .

where  $i_j$  is independently u.a.r. chosen from  $[d]$ .

def  $h_j$ :  $h_j(p) = p_{ij}$ .  $i_j$  u.a.r. chosen from  $[d]$ . then  $g(p) = (h_1(p), \dots, h_k(p))$ .

def: parameter  $\rho = \frac{\log(1/p_1)}{\log(1/p_2)}$ .

Lemma:  $\rho_g = \rho_h \leq \frac{1}{c}$ . (Rmk:  $P_{1,g} = P_{1,h}$ ,  $P_{2,g} = P_{2,h}$ ).

Pf:  $P_{1,g} = \prod_{i=1}^k P[h_i(p) = h_i(q)] = P_{1,h}^k$ . if  $\text{Ham}(p, q) \leq r$ .

$P_{2,g} = \prod_{i=1}^k P[h_i(p) \neq h_i(q)] = P_{2,h}^k$ . if  $\text{Ham}(p, q) \geq c \cdot r$ .

$$\therefore \rho_g = \frac{\log(1/P_{1,h}^k)}{\log(1/P_{2,h}^k)} = \frac{\log(1/P_{1,h})}{\log(1/P_{2,h})} = \rho_h.$$

$\forall p, q \in \{0,1\}^d \therefore P[h(p) = h(q)] = \frac{d - \text{Ham}(p, q)}{d}$ .

$$\therefore \text{Ham}(p, q) \leq r \Rightarrow P[h(p) = h(q)] \geq 1 - \frac{r}{d} = P_{1,h}.$$

日期:

$$\text{Ham}(p, q) \geq cr \Rightarrow P[h(p) = h(q)] \leq 1 - \frac{cr}{d} = p_{2,h}$$

$$\therefore p_h \leq \frac{-\log(1-r/d)}{-\log(1-cr/d)} \leq \frac{1}{c}. \quad (\text{as: } (1-\frac{r}{d})^c \geq 1 - \frac{cr}{d}, \forall c > 1).$$

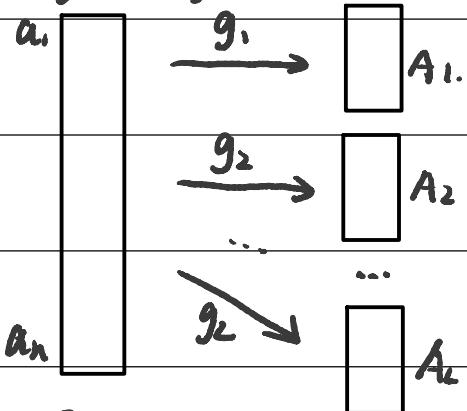
Alg.

Preprocessing:

lemma



- 1) set  $k = \log n / \log(1/p_{2,h})$ ,  $L = 2n^\rho$ . where  $\rho = \rho_g (\leq \frac{1}{c})$ .
- 2) Allocate  $L$  hash tables  $A_1, \dots, A_L$ , each a Hamming LSH  $g_i$  ( $i \in [L]$ ).
- 3) Hash  $\{a_1, \dots, a_n\}$  into tables  $A_1, \dots, A_L$ .



Query: (Input  $x$ ).

顺序遍历  $A_1[g_1(x)], \dots, A_L[g_L(x)]$ , return first  $a$  s.t.  $d(x, a) \leq cr$ .

if no such  $a$  is found before  $4L = 8n^\rho$  elements, output no.

Thm: The alg solves  $(c, r)$ -ANN w.p.  $\geq 0.6$ .

Pf: ① if  $a^*$  is a point s.t.  $\text{Ham}(x, a^*) \leq r$ .

若有  $\leq r$  个  $a^*$ ,

则  $P[a^* \notin A_1[g_1(x)] \cup \dots \cup A_L[g_L(x)]]$

w.p.  $< \frac{1}{7}$  不在其中.

$$= \prod_{i=1}^L P[a^* \notin A_i[g_i(x)]]$$

日期:

$$= \prod_{i=1}^L P[g_i(x) \neq g_i(a^*)]$$

$$\leq (1 - p_{1,g})^L.$$

$\ell$ 's origin.

$$\text{Consider, } P_{1,g} = P_{1,h} = P_{1,h}^{\frac{\log n / \log(\ell P_{2,h})}{\log(1/\rho_h) / \log(1/P_{2,h})}} = n^{-\log(1/\rho_h) / \log(1/P_{2,h})} = n^{-\rho}.$$

$$\text{从而 } P[\alpha^* \notin A_1[g_1(x)] \cup \dots \cup A_L[g_L(x)]] \leq (1 - n^{-\rho})^L = \left(1 - \frac{1}{n^\rho}\right)^{2n^\rho} \leq \left(\frac{1}{e}\right)^2 < \frac{1}{7}.$$

② let  $Y$  be the number of pairs  $(a, j)$  s.t.

$$d(a, x) > cr \text{ 且 } g_j(a) = g_j(x).$$

即坏点数 w.p.  $\leq \frac{1}{4}$ .  $\geq 4L$ .

(Rmk: 即放在  $L$  tables 中的 distance far by points 的 num).

$$\text{Pf: } E[Y] \leq L \cdot n \cdot P[g_j(a) = g_j(x) \text{ for some fixed } j].$$

$\uparrow$   
union.

$$= n \cdot L \cdot P_{2,g}.$$

Parameters:  $\Rightarrow P_{2,g} \Rightarrow \rho \Rightarrow L$ .

$= L$ .  $\leftarrow P_{2,g}$ 's value's origin. (erase  $n$ 's impact).

$$\text{By Markov, } P[Y \geq 4L] = P[Y \geq 4E[Y]] \leq \frac{1}{4}.$$

$$P_{2,g} = \frac{1}{n} \cdot k$$

$$k = \frac{\log(1/n)}{\log(1/\ell P_{2,h})}.$$

即 w.p.  $\leq \frac{1}{4}$ .  $Y \geq 4L$ .

③ Assume  $\exists a^*$  s.t.  $\text{Ham}(a^*, x) \leq r$ .

by ①, ② w.p.  $\geq 1 - \frac{1}{7} - \frac{1}{4} > 0.6$ .

alg finds a s.t.  $d(x, a) \leq r$ .

(Rmk: 无法保证 find. a s.t.  $d(x, a) \leq r$ ).

Thm: Space:  $O(n^{1+\frac{r}{d}} + nd)$ . query time:  $O(d \cdot n^{1/d})$ .  $L = nr$ .



# Streaming Algorithm

日期:

/

## Lecture 4: Counting.

(Link: the important point is space).

Problem: A sequence of events, at any given time, output the number of events.

Goal: find a data structure maintaining a value  $n$  and supporting:

① update    ② query:

Aim: Use as small space as possible.

Alg:

1. Exact Solutions:

maintain  $n$ . Space  $\Theta(\log n)$ .

2. Approximate Solutions:

Goal: Maintain  $\tilde{n}$  st.

$$P[|n - \tilde{n}| > \varepsilon \cdot n] < \delta. \quad \varepsilon, \delta \in (0, 1). \quad \Rightarrow \quad \text{w.p. } 1-\delta, \quad (1-\varepsilon)n \leq \tilde{n} \leq (1+\varepsilon)n.$$

① Morris Alg:

1)  $X \leftarrow 0$ .

Space:  $O(\log \log n)$ .

2) For each update, increment  $X$  by 1 w.p.  $\frac{1}{2^X}$ .

↑  
# of bits.

3) For query, output  $\tilde{n} = 2^X - 1$ .

Jensen Ineq:

$$2^{E[X_n]} \leq E[2^{X_n}] = n + 1.$$

日期:

$$\Rightarrow E[X_n] = O(\log n).$$

$$\Rightarrow E[\log X_n] = O(\log \log n).$$

Correctness Analysis:

let  $X_n$  be  $x$  after  $n$  update.

(Aim:  $E[X_n]$ ,  $\text{Var}[X_n]$ .)

Lemma 1:  $E[2^{X_n}] = n+1$ .  $\Rightarrow E[\hat{n}] = n$ .

Pf: By Induction.

$$E[2^{X_{n+1}}] = \sum P[X_n=j] \cdot E[2^{X_{n+1}} | X_n=j].$$
$$= \sum P[X_n=j] \cdot [2^j \cdot (1 - \frac{1}{2^j}) + \frac{1}{2^j} \cdot 2^{j+1}].$$

$$= \sum P[X_n=j] 2^j + \sum P[X_n=j].$$

$$= E[2^{X_n}] + 1 \stackrel{\uparrow}{=} n+2. \quad \square.$$

induction.

Lemma 2:  $\text{Var}[\hat{n}]$ .

$$\text{Var}[\hat{n}] = E[(\hat{n} - n)^2]. \quad (E[2^{X_n}] = \frac{3}{2}n^2 + \frac{3}{2}n + 1).$$
$$= \frac{n^2}{2} - \frac{n}{2} < \frac{n^2}{2}$$

By Chebyshev's Ineq:

$$P[|\hat{n} - n| > \epsilon n] \leq \frac{\text{Var}[\hat{n}]}{\epsilon^2 n^2} < \frac{1}{2\epsilon^2}. \quad \begin{array}{l} \text{太大不能 s.t. } \epsilon, \delta \in (0,1) \text{ 且可尽量小.} \\ (\text{即O最好与 } \epsilon \text{ 无关).} \end{array}$$

$\Rightarrow$  Bound 稍多. 改进:

② Morris + alg:

日期:

① Independently run  $s$  Morris alg. ( $s \geq \frac{1}{2\delta\epsilon^2}$ ).

② Output  $\hat{n} = \frac{1}{s} \sum_{i=1}^s \tilde{n}_i$ . Space:  $O(\frac{1}{2\delta\epsilon^2} \log n)$ .

Correctness Analysis:

$$E[\hat{n}] = E\left[\frac{1}{s} \sum_{i=1}^s \tilde{n}_i\right] = n.$$

$$\text{w.p. } \frac{1}{(1+\epsilon)^s} \uparrow \quad \hat{n} = \frac{(1+\epsilon)^s - 1}{\epsilon}$$

$$\text{Var}(\hat{n}) = \frac{1}{s^2} \sum_{i=1}^s \text{Var}(\tilde{n}_i) < \frac{1}{s^2} \cdot s \cdot \frac{n^2}{2} < \frac{n^2}{2s}$$

$$\frac{2n^2}{2s}$$

By Chebyshev's Ineq:

$$P[|\hat{n} - n| > \epsilon n] \leq \frac{\text{Var}(\hat{n})}{\epsilon^2 n^2} < \frac{1}{2\epsilon^2 s}.$$

$$\frac{2}{2\epsilon^2 s}$$

let  $s \geq \frac{1}{2\delta\epsilon^2}$ , then  $P[|\hat{n} - n| > \epsilon n] < \delta$ .  $\square$ .

③ Morris ++ alg:

$$s = \frac{32}{2\epsilon^2}$$

① let  $s = \frac{1}{2 \cdot \frac{1}{3} \epsilon^2}$ . Independently run  $t$  Morris + alg. ( $t \geq 18 \ln(1/\delta)$ ).

② Output the median of  $\tilde{n}_1, \dots, \tilde{n}_t$ .

Analysis:

$$\text{def } Y_i = \begin{cases} 1, & (1-\epsilon)n \leq \tilde{n}_i \leq (1+\epsilon)n \\ 0, & \text{otherwise.} \end{cases} \Rightarrow E[Y_i] = P[Y_i] \geq \frac{2}{3} \quad \text{示性变量}$$

$$\text{def } Y = \sum_{i=1}^t Y_i$$

$$E[Y] = \sum_{i=1}^t E[Y_i] \geq \frac{2}{3}t$$

计算  $E[\sum Y_i]$ .  
 $\downarrow$

$$P[\tilde{n} \text{ bad}] \leq P[Y \leq \frac{t}{2}] \leq P[Y - E[Y] < -\frac{t}{6}]$$

$$\text{Chernoff } \leq e^{-2(\frac{t}{6})^2 t}$$

求中位数.  
 $\downarrow$   
转化为 Chernoff.

日期:

/

let  $t \geq 18 \ln(1/\delta) \Rightarrow P[\tilde{n} \text{ bad}] \leq \delta$ .  $\square$ .

Space: ④ (s.t. space( $\tilde{x}$ )) = ④ ( $\frac{\ln(1/\delta)}{\epsilon^2} \text{ space}(x)$ ).

If  $\tilde{x}$  reaches  $\log(stn/\delta)$ , then at any time,  $\tilde{x}$  increase wp.  $\frac{1}{2^{\tilde{x}}} \leq \frac{\delta}{stn}$ .

then  $P[\tilde{x} \text{ increase in the next } n \text{ updates}] \leq \frac{\delta}{stn} \cdot n = \frac{\delta}{st}$ .

By Union Bound, wp.  $\geq 1 - \delta$ ,

$\tilde{x}_i \leq \log(stn/\delta), i=1, \dots, st. \quad (stn/\delta = \frac{n}{\epsilon^2 \delta} \log(\frac{1}{\delta}) \leq \frac{n^2}{\epsilon^2 \delta^2})$   
 $\Rightarrow \text{Space} \leq O(st \cdot \log \log(stn/\delta)) = O(\frac{1}{\epsilon^2} \log(\frac{1}{\delta}) \log \log(\frac{n}{\epsilon \delta})). \quad \text{wp.} \geq 1 - \delta$

其他估计:

$$[2st \log \log(2st/\delta)]$$

Flajolet 85:  $O(\log(1/\epsilon) + \log \log n + \log(1/\delta))$ .

Nelson 2020: ④ ( $\log(\frac{1}{\epsilon}) + \log \log n + \log \log(1/\delta)$ ). (需要对 Morris++ 稍作改进)

## Lecture 5: Reservoir Sampling.

Input: A Sequence of numbers  $a_1, a_2, \dots, a_m$  (m<sup>2</sup> 只关心) revealed over time.

Problem: Output a uniform sample  $s$  from all the revealed numbers.

Algorithm: ① 维护一个数字  $s$ . initially let  $s=0$ .

② on seeing the m-th number  $a_m$ , then set  $s=a_m$  wp.  $\frac{1}{m}$ .

日期:

/

else, maintain  $S = S.$

Correctness:  $P[S = a_i] = \frac{1}{l} \cdot \left(1 - \frac{1}{i+1}\right) \cdots \left(1 - \frac{1}{m}\right) = \frac{1}{m} \quad \forall i.$

Space: (If  $a_i \in n, \forall i$ ). then  $O(\log n) + O(\log m)$  bits.  
↓ store  $n_i$ .      ↗ store  $m$ .

(Rmk: ①  $k$  samples  $\Rightarrow$  Independently run  $k \uparrow$  Alg.)

②  $\tilde{m}$  (Morris Alg 优化).

Lecture 6: Estimating the number of distinct elements.

Input: a stream of integers  $i_1, \dots, i_m \in [n] = \{1, \dots, n\}$ .

Output: the number of distinct elements seen in the stream.

(Rmk: denoted as 数据流的一范数问题).

Exact Solutions:

① Store an array  $[n]$ :  $O(n)$

② Store all the elements.  $m \cdot O(n)$ .

Approximate:

Maintain a  $\hat{t}$  st.  $P[|\hat{t} - t| > \varepsilon t] < \delta$ . where  $t$  is the number of distinct elements in stream.  $\varepsilon, \delta \in (0, 1)$ .

Alg 1: Idealed FM alg:

日期:

① Pick a random function:  $h: [n] \rightarrow [0,1]$ .

② Maintain a counter  $X = \min_{i \in \text{Stream}} h(i)$ .

③ Output  $\hat{t} = 1/X - 1$ .

Correctness Pf:

Claim1:  $E[X] = \frac{1}{t+1}$ .

Lemma:  $E[X] = \int_0^\infty f(y) y dy$

$$= \int_0^\infty \int_0^y f(y) d\lambda dy$$

$$= \int_0^\infty \int_\lambda^\infty f(y) dy d\lambda$$

$$= \int_0^\infty P[X \geq \lambda] d\lambda.$$

由 Lemma,  $E[X] = \int_0^\infty P[X \geq \lambda] d\lambda$

$$= \int_0^\infty P[\forall a_i \in \text{stream}, h(a_i) \geq \lambda] d\lambda$$

$$= \int_0^1 (1-\lambda)^t d\lambda$$

$$= \frac{1}{t+1}.$$

Claim2:  $\text{Var}(X) = \frac{t}{(t+1)^2(t+2)} < \frac{1}{(t+1)^2}$ .

Pf:  $E[X^2] = \int_0^\infty P[X^2 \geq \lambda] d\lambda$ .

$$= \int_0^1 (1-\lambda)^t d\lambda.$$

日期:

?

$$= 2 \int_0^1 u^t (1-u) du = \frac{2}{(t+1)(t+2)}.$$

$$\therefore \text{Var}[X^2] = E[X^2] - (E[X])^2 = \frac{t}{(t+1)^2(t+2)} < \frac{1}{(t+1)^2}.$$

Alg2: FM + Alg. (Mean Trick).  $\rightarrow$  降低 Var t, 获得 Correctness.

① Independently run  $S = \frac{25}{\epsilon^2 S}$ . copies of FM Alg. and obtain counters

$x_1, \dots, x_S$ .

② Output  $\hat{t} = \frac{1}{S} - 1$ , where  $Z = \frac{1}{S} \sum_{i=1}^S x_i$ .

Correctness Pf:

$$\text{Lem1: } P\left[|Z - \frac{1}{t+1}| > \frac{\epsilon}{5(t+1)}\right] < \delta.$$

$$\text{Pf: } E[Z] = \frac{1}{t+1}. \quad \text{Var}[Z] = \frac{1}{S^2} \sum_{i=1}^S \text{Var}[x_i] < \frac{1}{S(t+1)^2}.$$

By Chebyshev Ineq.

$$P\left[|Z - \frac{1}{t+1}| > \frac{\epsilon}{5(t+1)}\right] < \frac{\text{Var}[Z]}{\left(\frac{\epsilon}{5(t+1)}\right)^2} < \frac{25}{\epsilon^2 S} = \delta.$$

Claim:  $P[|\hat{t} - t| > \epsilon t] < \delta$

Pf: if  $t=0$  ✓.

if  $t \geq 1$ . by Lem1, w.p.  $\geq 1 - \delta$ .

$$\frac{1-\epsilon/5}{(t+1)} < Z < \frac{1+\epsilon/5}{(t+1)}$$

$$\Rightarrow \frac{t-\epsilon/5}{1+\epsilon/5} < \hat{t} = \frac{1}{S} - 1 < \frac{t+\epsilon/5}{1+\epsilon/5}$$

$$\frac{t-\epsilon/5}{1+\epsilon/5} > (t-\epsilon/5)(1-\epsilon/5) > t - 2\epsilon t/5.$$

日期:

另一边同理,  $\Rightarrow P(|\hat{t} - t| > \varepsilon t) < \delta$ .  $\square$ .

Alg 3. FM++ Alg. (Median Trick)  $\rightarrow$  降低 space complexity.

① Independently run  $q = 8 \ln(1/\delta)$  copies FM+Alg. with  $s' = \frac{1}{3}$ . Obtain  $\hat{t}_1, \dots, \hat{t}_q$ .

② Output median of  $\hat{t}_1, \dots, \hat{t}_q$  as  $\hat{t}$ .

Correctness Pf:

Let  $Y = \sum Y_i$  where  $Y_i = \begin{cases} 1, & |\hat{t}_i - t| < \varepsilon t \\ 0, & \text{o.w.} \end{cases}$

Note  $EY \geq \frac{2}{3}q$ .

$\Leftarrow$  Median Trick 对于连续变化的量 ( $t$ ) 均可行

$\therefore P(|\hat{t} - t| > \varepsilon t) \leq P[Y \leq \frac{q}{2}]$

Chernoff  $\Downarrow$

$\leq e^{-2 \frac{1}{6} q^2} \leq \delta$ .  $\square$ .

Space: #counters =  $O(q \cdot s) = O(\frac{1}{\varepsilon^2} \ln(1/\delta))$ .

----- way to get hash functions.

I.  $k$ -wise independent hash families:

Def: a family  $\mathcal{H}$  of functions mapping  $[a] \rightarrow [b]$  is  $k$ -wise independent if:

$\forall j_1, \dots, j_k \in [b]$ , distinct  $i_1, \dots, i_k \in [a]$  有

$$P_{h \in \mathcal{H}}[h(i_1) = j_1 \wedge \dots \wedge h(i_k) = j_k] = \frac{1}{b^k}.$$

(Rmk:  $P[h(i) = j] = \frac{1}{b}$ :  $\forall i, j$ ,  $k$ -wise independent  $\Rightarrow k \uparrow$  增强).

e.g.:  $\mathcal{H} = \{\text{all functions: } [a] \rightarrow [b]\}$  is  $k$ -wise independent.

日期:

存储  $h \in H_{\text{poly}}$  need space  $O(\alpha \log b)$ .

e.g.: let  $a = b = q$ , where  $q = p^r$ ,  $p$  prime.

$H_{\text{poly}} = \{ \text{polynomials } \in F_q[X] \mid \text{degree} \leq k-1 \}$ .

$H_{\text{poly}}$  is  $k$ -wise independent.  $\downarrow$   
 $k$  位数字.

存储  $h \in H_{\text{poly}}$  need space  $O(k \log q)$ .

(Rmk: 用小空间存储 hash 函数).

II. KMV alg (BJKST). ( $k$ -minimum values). Space:  $O(\frac{\lg n}{\varepsilon^2})$

① 2-wise independent hash  $h: [n] \rightarrow [M]$ .  $M = n^3$ .

② Maintain  $k$  smallest hash values of numbers in stream.

③ if  $< k$  distinct hash values. output # distinct hash values.

o.w. output  $\tilde{t} = k \cdot M / X$ . ( $X$  is  $k$ -th smallest value)

Analysis:

If  $\frac{1}{\sqrt{n}} < \varepsilon < \frac{1}{2}$ , then w.p.  $\geq \frac{2}{3}$ .  $(1-\varepsilon)t < \tilde{t} < (1+\varepsilon)t$ .

(Rmk: use median trick, run  $O(\lg \frac{1}{\delta})$  times, to boost prob. to  $1-\delta$ .)

Pf: ①  $P[\tilde{t} > (1+\varepsilon)t]$

$$= P[X < \frac{km}{(1+\varepsilon)t}]$$

$$= P[Y \geq k].$$

def 离性变量  $Y_i = \begin{cases} 1, & h(a_i) < \frac{km}{(1+\varepsilon)t} \\ 0, & h(a_i) \geq \frac{km}{(1+\varepsilon)t}. \end{cases}$

$$Y = \sum_{i=1}^t Y_i$$

$$\text{Note: } EY_i < \frac{k}{1+\varepsilon t} \Rightarrow EY < \frac{k}{1+\varepsilon t}$$

$$\text{Var} Y < \frac{k}{1+\varepsilon t}$$

日期:

$$\leq P[Y - E[Y] \geq k - \frac{k}{4\epsilon}]$$

Chebyshev

$$< \frac{k}{4\epsilon} \cdot \frac{(4\epsilon)^2}{k^2 \epsilon^2} = \frac{16\epsilon}{\epsilon^2 k} < \frac{1}{6}$$

$$② P[\hat{t} < (1-\epsilon)t]$$

$$= P[X > \frac{km}{(1-\epsilon)t}]$$

$$= P[Z < k].$$

...

$$< \frac{1}{6}$$

□.

$$\text{def 离性变量 } Z_i = \begin{cases} 1, & h(a_i) < \frac{km}{(1-\epsilon)t} \\ 0, & h(a_i) \geq \frac{km}{(1-\epsilon)t} \end{cases}$$
$$Z = \sum_{i=1}^t Z_i.$$

$$\text{Note: } \mathbb{P}[Z_i = 1] > \frac{k}{(1-\epsilon)t} - \frac{1}{m}$$

$$\Rightarrow EZ_i > \frac{k}{(1-\epsilon)t} - \frac{1}{m}$$

$$\Rightarrow \frac{(1+\epsilon)k}{t} - \frac{ek}{4t}. \quad (M = n^3 > \frac{4k}{\epsilon}).$$

## Lecture 7: Frequent items

### §1: Finding the majority:

Input: a stream.  $a_1, \dots, a_n \in [m]$ .

Problem: Check if  $\exists i \in [n]$  st. appears  $> \frac{m}{2}$  times. if yes, output it.

Misra-Gries alg:

Space:  $O(\log m + \log n)$ .

Maintain an ID and a counter C. (initial ID=1, C=0).

if  $a = ID \Rightarrow C \leftarrow C + 1$ .

else if  $c = 0 \Rightarrow ID \leftarrow a, C \leftarrow 1$ .

else  $\Rightarrow C \leftarrow C - 1$ .

(Rmk: 若不存在 majority, 会输出 any number.)

日期:

Way to follow: 再跑一遍 stream (以验证.)

§2: Finding frequent items.

$O(k(\log n + \log m))$

(Addition Stream)

Input: a stream.  $a_1, \dots, a_n \in [m]$ ,  $k \geq 1$ .

Problem: Check if  $\exists i \in [n]$  st. appears  $> \frac{m}{k+1}$  times. if yes, output all.  
(最多  $k\hat{n}$ ).

Misra-Gries alg: (general).

Maintain  $k$  ID and  $k$  counter  $C$ . (initial  $ID_i = L$ ,  $c_i = 0$ ).

if  $\exists i$  st.  $a = ID_i \Rightarrow C_i \leftarrow C_i + 1$

else if  $\exists i$  st.  $C_i = 0 \Rightarrow ID_i \leftarrow a$ .  $C_i \leftarrow 1$

else  $C_j \leftarrow C_j - 1$ .  $\forall j$ . (decrement-all (DA) step).

(Rmk: the output 不一定都 s.t. condition.  $\Rightarrow$  通过 second pass).

Proof:

Denote  $f_i$  as the times  $i$  appears in the stream.

$\tilde{f}_i$  as  $\begin{cases} C_j & , \text{if } i = ID_j \\ 0 & , \text{o.w.} \end{cases}$

Thm: for  $\forall i \in [n]$ ,  $f_i - \frac{m}{k+1} \leq \tilde{f}_i \leq f_i$ . 即要求至少  $\frac{m}{2(k+1)}$  DA.

Pf: the alg  $\Leftrightarrow$  Maintain Counter  $C_i$  for  $\forall i \in [n]$ .

日期:

input  $i$ :

if  $c_i \neq 0$ ,  $\Rightarrow c_i \leftarrow c_{i+1}$

else if  $< k$  nonzero counters.  $\Rightarrow c_i \leftarrow 1$

else  $\Rightarrow c_j \leftarrow c_j - 1$  for  $\forall c_j \neq 0$ .

$(\tilde{f}_i)$  对于  $c_i$ .  $\tilde{f}_i - f_i$ .

input  $i$ . input  $j$

+1 0 0

+1 0 0

0 -1 -1

由上易见,  $\tilde{f}_i \leq f_i$ . 要证  $f_i - \frac{m}{k+1} \leq \tilde{f}_i$ . 只要证最多发生  $\frac{m}{k+1}$  次 DA steps.

Pf: 设发生  $t$  次. 则  $(m-t) - kt \geq 0 \Rightarrow t \leq \frac{m}{k+1}$ .  $\square$ .

势  $\varphi(m) = \sum \tilde{f}_i^m$  势.

Cor:  $HH_{k+1}(S) \subseteq H$ .

Extension:

denote  $HH_{k+1}(S) \stackrel{\text{stream}}{=} \{j \in [n] : f_j > \frac{1}{k+1}m\}$  Alg output  $H$ .

for Misra-Gries alg,  $HH_{k+1}(S) \subseteq H$ .

Here's an alg, s.t.  $HH_{k+1}(S) \subseteq H \subseteq HH_{2(k+1)}(S)$ .

(aim:  $f_i - \frac{m}{k+1} \leq \tilde{f}_i \leq f_i - \frac{m}{2(k+1)}$ )

§3: turnstile streaming model:

General Model:

Input  $i_1, \dots, i_m \in [n]$ .

let  $x = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$ . denote each  $c_j$  in stream as  $(i_j, \Delta_j)$ .  $\Delta_j \in R$ .

update:  $x_{i_j} \leftarrow x_{i_j} + \Delta_j$ .

日期:

Cash Register Model:  $\Delta_j > 0$ . (especially  $\Delta_j = 1$ ).

Strict Turnstile Streaming Model:  $\Delta_j$  arbitrary, but  $x \geq 0$ . at all time.

Problems: (for General Model).

I:  $(k, L_1)$ -point query problem:

query i

↓ return  $\tilde{x}_i$  s.t.  $|x_i - \tilde{x}_i| \leq \frac{1}{k} \|x\|_1$ . w.h.p.

II:  $(k, L_1)$ -heavy hitters:

query.

return a set  $L$  s.t. 1)  $L = O(k)$ . w.h.p.

2) if  $|x_i| > \frac{1}{k} \|x\|_1$ , then  $i \in L$ .

Lemma: If  $A$  is an alg. for  $(3k, L_1)$ -point query problem

w.p.  $\geq 1 - \frac{\delta}{n}$  & w.s.  $s$  bits. Then there's an alg.  $A'$  solve  $(k, L_1)$ -heavy hitters. w.p.  $1 - \delta$  & w.s.  $s + O(k \log n)$  bits.

Proof: Define  $A'$  as:

for all  $i \in [n]$ . use  $A$  to answer query (i).

remember  $3k$  indices with the  $3k$  largest point query values.

Correctness: w.p.  $\geq 1 - \frac{\delta}{n} \cdot n = 1 - \delta$ .

日期:

$\forall$  queries( $i$ ) return  $\tilde{x}_i$ , st.  $|\tilde{x}_i - x_i| \leq \frac{\|x_i\|_1}{3k}$ .

then for  $\forall i$  st.  $x_i > \frac{\|x_i\|_1}{k}$ , 有:  $\tilde{x}_i > \frac{\|x_i\|_1}{k} - \frac{\|x_i\|_1}{3k} = \frac{2\|x_i\|_1}{3k}$ .

$\forall i$  st.  $x_i \leq \frac{\|x_i\|_1}{3k}$ , 有  $\tilde{x}_i \leq \frac{\|x_i\|_1}{3k} + \frac{\|x_i\|_1}{3k} = \frac{2\|x_i\|_1}{3k}$ .

Note that,  $\#\{i : x_i > \frac{\|x_i\|_1}{3k}\} \leq 3k$ .

$\Rightarrow \#\{i : \tilde{x}_i > \frac{2\|x_i\|_1}{3k}\} \leq \#\{i : x_i > \frac{\|x_i\|_1}{3k}\} \leq 3k$ .

而  $\{i : x_i > \frac{\|x_i\|_1}{k}\} \subseteq \{i : \tilde{x}_i > \frac{2\|x_i\|_1}{3k}\}$ .  $\square$ .

Count Min Sketch: (for:  $(k, \epsilon)$ -point query).

Let  $w \geq k$ ,  $d = \Omega(\log(1/\delta))$  be parameters.

① Choose 2-wise independent hash functions.  $h_1, \dots, h_d : [n] \rightarrow [w]$ .

let  $C[l, s] = 0$ .  $\begin{cases} l \in [d], \\ 1 \leq s \leq w \end{cases}$ .

② for each item  $e_t = (i_t, \Delta_t)$  in stream.

$$C[l, h_l(i_t)] \leftarrow C[l, h_l(i_t)] + \Delta_t \quad \forall l.$$

③ for  $\forall i \in [n]$ , set  $\tilde{x}_i = \min_{l \in [d]} C[l, h_l(i)]$  MIN.

④ for a query  $i$ , output  $\tilde{x}_i$ .

Lemma (for strict turnstile model):

For  $\forall$  fixed  $i$ ,  $x_i \in \tilde{x}_i$  and  $P[\tilde{x}_i \geq x_i + \|x_i\|_1/k] \leq \delta$ .

Pf: ①  $\tilde{x}_i = \min_{j \in d} \tilde{x}_{ji} = \min_{j \in d} (x_i + \sum_{i \in I_j} x_{is}) \geq x_i$ .  
 $\sum_{i \in I_j} x_{is} \geq 0$

日期:

$$\begin{aligned} \textcircled{2} \quad \forall i, \quad & E[\underbrace{C[l, h_c(i)]}_{:= Z_L}] = x_i + \sum_{j \neq i} \Pr[h_c(j) = h_c(i)] \cdot x_j \geq x_i \\ & = x_i + \frac{1}{w} \sum_{j \neq i} x_j \\ & < x_i + \frac{\|x\|_1}{2k}. \end{aligned}$$

$$\Rightarrow E[Z_L - x_i] < \frac{\|x\|_1}{2k}.$$

$$\begin{aligned} \text{Markov} \Rightarrow \Pr[Z_L - x_i > \frac{\|x\|_1}{k}] & < \frac{\|x\|_1 / 2k}{\|x\|_1 / k} = \frac{1}{2}. \\ Z_L - x_i > 0 \end{aligned}$$

$$\Rightarrow \Pr[\tilde{x}_i - x_i > \frac{\|x\|_1}{k}] < (\frac{1}{2})^d < \delta.$$

Rank: let  $d = \Theta(\ln n)$ ,  $w = 2k$ , then w.p.  $\geq 1 - \frac{1}{n}$  st.  $\tilde{x}_i \leq x_i + \frac{\|x\|_1}{k}$  for some  $x_i$ .

Space Complexity:  $d w$  counters =  $O(k \log n \cdot \log m)$ . ( $\|x\|_1 \leq m$ ).

Count Sketch (for  $(k, L_2)$ -point query).

Let  $w \geq 3k^2$ ,  $d = 4(\log(1/\delta))$ , be parameters.

① Choose 2-wise independent hash functions.  $h_1, \dots, h_d: [n] \rightarrow [w]$ .

let  $C[l, s] = 0$ .  $\begin{cases} 1 \leq l \leq d, \\ 1 \leq s \leq w \end{cases}$ . hash functions  $g_1, \dots, g_d: [n] \rightarrow \{\pm 1\}$

② for each item  $t$ :  $(i_t, \Delta_t)$  in stream.

$$C[l, h_l(i_t)] \leftarrow C[l, h_l(i_t)] + \Delta_t \cdot g_l(i_t).$$

③ for  $\forall i \in [n]$ , set  $\tilde{x}_i = \underset{l \leq d}{\text{median}}[C[l, h_l(i)] \cdot g_l(i)]$ . Median.

④ for a query  $i$ , output  $\tilde{x}_i$ .

日期:

Lemma: for  $\forall$  fixed  $i$ ,  $\mathbb{E}[\hat{x}_i] = x_i$ .

$$\textcircled{2} \quad \Pr[|\hat{x}_i - x_i| \leq \frac{\|x\|_2}{k}] \geq 1 - \delta.$$

Pf: fix  $i$ .  $\textcircled{1}$ : let  $Z_i = \mathbb{E}[h_i(i)]$ .

$$\begin{aligned} \Rightarrow EZ_i &= x_i g(i) + \sum_{j \neq i} \Pr[h_i(j) = h_i(i)] \cdot x_j g(j) \\ &= x_i g(i) + \sum_{j \neq i} \frac{1}{w} (\frac{1}{2} x_j - \frac{1}{2} x_j) \\ &= x_i g(i). \end{aligned}$$

$$\Rightarrow \mathbb{E}[\hat{x}_i] = x_i.$$

$$\textcircled{2} \quad EZ_i^2 = x_i^2 + \sum_{j \neq i} \Pr[h_i(j) = h_i(i)] \cdot x_j^2 + 2 \sum \Pr[h_i(j) = h_i(i)] \cdot x_i g(i) \underbrace{\mathbb{E}[x_j g(j)]}_{=0} \\ < x_i^2 + \frac{\|x\|_2^2}{w}.$$

$$\Rightarrow \text{Var} Z_i = EZ_i^2 - (EZ_i)^2 < \frac{\|x\|_2^2}{w}.$$

$$\Pr[|Z_i - x_i| \geq \frac{\|x\|_2}{k}] \leq \frac{\frac{\|x\|_2^2}{w}}{\frac{\|x\|_2^2}{w}/k^2} = \frac{k^2}{w} = \frac{1}{3}$$

由 Median Trick,  $\Pr[|\hat{x}_i - x_i| \geq \frac{\|x\|_2}{k}] \leq \delta$ .  
(-j+± Chernoff).

(Rmk: i.e Count Sketch ≈ g, 则 g=linear. 即  $g(x + \Delta_j \vec{1}_j) = g(x) + \Delta_j \cdot g(\vec{1}_j)$ .)

Applications:

I.  $\Rightarrow$  Point query.  $\Rightarrow$  heavy hitters.

II. Range Queries:

日期:

## Count Min Sketch + 线段树

given  $i, j \in [n]$ , estimate  $\sum_{l \in [i, j]} x_l$ .

Ideally  $O(c \log(nm))$ .

(Rmk: 不能直接  $\tilde{x}_{ij} = \sum_{l \in [i, j]} \tilde{x}_l$ ).

Def: Dyadic interval:

$[i, j]$  is a dyadic interval/range if:

- ①  $j-i+1 = 2^k$
- ②  $2^k$  divides  $i-1$  for some  $k \geq 0$ .

Alg: (线段树).

① 对线段树每层用 Count Min Sketch. (with  $w = 2k$ ,  $d = \log n \log(1/\delta)$ )

② upon a query  $[i, j]$ , decompose it to  $[i, j], \dots, [i_s, j]$ ,  $s \leq \log n$ .

③ 对  $s$  个 dyadic intervals, maintain  $\tilde{x}_{ij}, \dots, \tilde{x}_{jsj}$ .

④ Output  $\tilde{x}_{ij} = \tilde{x}_{ij} + \dots + \tilde{x}_{jsj}$ .

Space:  $O(k \log^2 n \log m)$ . update  $O(1)$ . (时间复杂度). query:  $O(\log n)$ .

## III. Sparse Recovery.

Given a vector  $x \in \mathbb{R}^n$  (Streaming model) find  $z$  s.t.  $z$  has  $\leq k$  non-zeros.

and  $\|x-z\|_p$  is minimized for some  $p \geq 1$ .

(Rmk: Offline Model:  $z$  picks the largest  $k$  entries of  $x$ ).

Aim: turnstile streams model.  $p=2$ .  $\tilde{O}(ck) = O(k \cdot \log^c k)$  space.

Define:  $\text{err}_2^k(x) = \min_{z: \|z\|_0 \leq k} \|x-z\|_2^2$ .

$$z: \|z\|_0 \leq k$$

日期:

Thm: The algorithm use  $O(\frac{k}{\varepsilon^2} \text{poly}(\log n))$  spaces and return  $\tilde{z}$ , s.t.

$$\|z\|_0 \leq k \text{ w.p. } \geq 1 - \frac{1}{n}, \|x - z\|_2 \leq (1 + \varepsilon) \text{err}_2^k(x).$$

(Rmk: in the algorithm, if  $x$  is  $k$ -sparse, then  $x = z$ )

Alg: Count Sketch + 选择  $k$  个.

→  $k$  len. queue  $\tilde{x}_1, \dots, \tilde{x}_n$ .

① let  $w = 3k/\varepsilon^2$ .  $d = \Omega(\log n)$ . use Count Sketch. (Maintain  $\tilde{x}_1, \dots, \tilde{x}_n$ ).

② Output the  $k$  coordinate of the  $k$  maximal estimate.

Rmk: space:  $O(k/\varepsilon^2 \cdot \log n \log m)$ . update:  $O(n \cdot \log k)$ .

Lemma1: Count Sketch with  $w = 3k/\varepsilon^2$ .  $d = \Omega(\log n)$ . ensures.

$$\forall i, |\tilde{x}_i - x_i| \leq \frac{\varepsilon}{\sqrt{k}} \text{err}_2^k(x). \text{ w.p. } \geq 1 - \frac{1}{n}.$$

Lemma2:

let  $x, y$  be vectors st.  $\|x - y\|_\infty \leq \frac{\varepsilon}{\sqrt{k}} \text{err}_2^k(x)$ .

Then let  $T$  be the set of  $k$  largest indices of  $y$ .

let  $z_i = y_i$ , if  $i \in T$ , then  $\|z - x\|_2 \leq (1 + \varepsilon) \text{err}_2^k(x)$ .

Pf1: denote  $T_{\text{big}} = \{i \mid x_i \text{ is } k \text{ largest in } x\}$ .  $T_{\text{small}} = [n] \setminus T_{\text{big}}$ .

Fix  $i \in [n]$ . Consider  $l$  s.t.  $l \leq l \leq d$

① denote  $A_l$  as event:  $h_l(i) = h_l(i')$  for some  $i' \in T_{\text{big}}, (i' \neq i)$ .

To prove  $\Pr[A_l] \leq \frac{\varepsilon^2}{3}$ .

日期:

② denote  $Z_l = g_{l(i)} + C[l, h_l(i)]$ .

then  $\hat{x}_i = x_i + \underbrace{\sum_{j \in T_{\text{big}} \setminus i} g_{l(j)} g_{l(j)} y_j x_j}_{:= Z_l''} + \underbrace{\sum_{j \in T \setminus i} g_{l(j)} g_{l(j)} y_j x_j}_{\substack{\text{Small} \\ := Z_l'}}$

To prove  $\Pr[|Z_l'| \geq \frac{\epsilon}{\sqrt{k}} \text{err}_2^k(x)] \leq \frac{1}{3}$ .

③ Assuming ①② holds, then.

$$|Z_l - x_i| = |Z_l'' + Z_l'|$$

$$= |Z_l'| \quad \text{w.p. } \geq 1 - \frac{\epsilon^2}{3} \Rightarrow \text{w.p. } \leq \frac{2}{3}$$

$$< \frac{\epsilon}{\sqrt{k}} \text{err}_2^k(x) \quad \text{w.p. } \geq \frac{2}{3}$$

Then by median Tick,  $|\hat{x}_i - x_i| < \frac{\epsilon}{\sqrt{k}} \text{err}_2^k(x)$  w.p.  $1 - \frac{1}{n}$ . ( $d = \Omega(\log n)$ ).

Prove of ①: 设  $Y_j = \Pr[h_l(j) = h_l(i)] \Rightarrow \Pr[Y_j = 1] = \frac{1}{n} \leq \frac{\epsilon^2}{3k}$ .

$$\xrightarrow{\text{Markov}} \Pr[\sum Y_j \geq 1] \leq k \cdot \frac{\epsilon^2}{3k} = \frac{\epsilon^2}{3}$$

Prove of ②:  $E[Z_l'] = 0$ .  $\text{Var}[Z_l'] \leq \frac{\epsilon^2}{3k} \sum_{j \in T \setminus i} x_j^2 = \frac{\epsilon^2}{3k} \text{err}_2^k(x)$ .

$$\Rightarrow \Pr[|Z_l'| \geq \frac{\epsilon}{\sqrt{k}} \text{err}_2^k(x)] \leq \frac{1}{3}$$

Pf2:  $\|x - z\|_2^2 = \sum_{i \in T} |x_i - z_i|^2 + \sum_{i \in S \setminus T} |x_i - z_i|^2 + \sum_{i \in \{x\} \setminus S \setminus T} |x_i - z_i|^2 \quad (S \triangleq T_{\text{big}})$

$$\stackrel{(i)}{\leq} k \cdot \left(\frac{\epsilon}{\sqrt{k}} \text{err}_2^k(x)\right)^2 \stackrel{(ii)}{\leq} k \cdot \left(\frac{\epsilon}{\sqrt{k}} \text{err}_2^k(x)\right)^2 \stackrel{(iii)}{\leq} \text{err}_2^k(x).$$

IV. Matrix Sketch (The Frequent Directions Problem).

Input:  $n$  rows of  $A \in \mathbb{R}^{n \times d}$  one by one (Bp in streaming model).

日期:

(Streaming Truncated SVD)

Output:  $B \in \mathbb{R}^{k \times d}$ . where  $k < n$ . s.t.  $A^T A \approx B^T B$ .

Precisely,  $\forall \|x\|=1$ ,  $0 \leq \|Ax\|^2 - \|Bx\|^2 \leq 2\|A\|_F^2/k$ .

i.e.  $\|A^T A - B^T B\| \leq 2\|A\|_F^2/k$ . 且  $A^T A - B^T B \geq 0$ .

Rmk: ① Frequent Items Problem's推广.

② If  $x$  is a direction ( $\|x\|=1$ ) that  $\|Ax\|^2 \geq \epsilon \|A\|^2$ . ( $\|A\| := \max_{\|x\|=1} \|Ax\|$ ).

then the Alg can recover  $x$  by choosing  $k \geq \frac{2}{\epsilon} \cdot \frac{\|A\|_F^2}{\|A\|^2}$ .

Alg. (需要整理)

① Initialize  $B \in \mathbb{R}^{k \times d}$  to be all-zero matrix.

② For each row  $A_i$  ( $i \in [n]$ ).

1° Insert  $A_i$  into a zero valued row of  $B$ .

2° If  $B$  has no zero valued row, then,

1) Compute SVD of  $B$ :  $B = UDV^T$

2) let  $C = PV^T$ ,  $\delta = 6k^2/2$ ,  $\tilde{D} = \sqrt{\max(D^2 - I_k \cdot \delta, 0)}$   
(至少一半以上变>0)

3) let  $B = \tilde{D}V^T$

↳不要U, 因为保证下面行为0.

③ Output  $B$ .

Correctness:

Claim: for output  $B$ ,  $0 \leq B^T B \leq A^T A$ .

日期:

Pf:  $0 \leq B^T B$  显然.

Denote  $B^{(i)}, C^{(i)}$  as  $B$  and  $C$  after  $i$  loops.

Note that,  $B^{(0)} = 0$ ,  $B^{(n)} = B$ .

$$\text{then } \|Ax\|^2 - \|Bx\|^2 = \sum_{i=1}^n (\langle A_i x \rangle^2 + \|B^{(i-1)}x\|^2 - \|B^{(i)}x\|^2).$$

$$\stackrel{(*)}{=} \sum_{i=1}^n \|C^{(i)}x\|^2 - \|B^{(i)}x\|^2 \\ \stackrel{(\star)}{\geq} 0.$$

Why (\*) :  $B^{(i+1)} + A_i \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix} = U^T C^{(i)}$

Why (\star) : Case 1:  $\|C^{(i)}x\|^2 - \|B^{(i)}x\|^2 = 0$ .

Case 2:  $x^T C^{(i)T} C^{(i)} x - x^T B^{(i)T} B^{(i)} x$

$$= (V^T x)^T D^{(i)T} D^{(i)} (V^T x) - (V^T x)^T \tilde{D}^{(i)T} \tilde{D}^{(i)} (V^T x)$$

$$= (V^T x)^T (D^{(i)} - \tilde{D}^{(i)}) (D^{(i)} - \tilde{D}^{(i)})^T (V^T x)$$

$$\geq 0 \quad (D^{(i)} - \tilde{D}^{(i)}) \geq 0.$$

Claim 2:  $\|A^T A - B^T B\| \leq 2 \|A\|_2^F / k$ .

Pf: let  $x$  be the unit eigenvector for the max eigenvalue of  $A^T A - B^T B$ .

Note that,  $\|A^T A - B^T B\| = \max_{\|x\|=1} \|(A^T A - B^T B)x\| = \|(A^T A - B^T B)x\| = \lambda$ .

$$x^T (A^T A - B^T B) x = x^T \lambda x = \lambda.$$

日期:

$$\therefore \|\mathbf{A}^T \mathbf{A} - \mathbf{B}^T \mathbf{B}\| = \|\mathbf{A}\mathbf{x}\|^2 - \|\mathbf{B}\mathbf{x}\|^2$$

$$\stackrel{\text{Claim!}}{=} \sum_{i=1}^n \|\mathbf{C}^{(i)} \mathbf{x}\|^2 - \|\mathbf{B}^{(i)} \mathbf{x}\|^2$$

$$\stackrel{(*)}{\leq} \sum_{i=1}^n \|\mathbf{C}^{(i)T} \mathbf{C}^{(i)} - \mathbf{B}^{(i)T} \mathbf{B}^{(i)}\|.$$

$$= \sum_{i=1}^n \|(\tilde{\mathbf{D}}^{(i)})^2 - (\mathbf{D}^{(i)})^2\|.$$

$$= \sum_{i=1}^n \left\| \begin{pmatrix} \delta_{11}^{(i)2} & & \\ & \ddots & \\ & & \delta_{nn}^{(i)2} \end{pmatrix} \right\|$$

$$= \sum_{i=1}^n \delta_{ii}^{(i)2} \cdot 6n^2 + 1$$

$$\text{Why (*)? } \sum_{i=1}^n \|\mathbf{C}^{(i)} \mathbf{x}\|^2 - \|\mathbf{B}^{(i)} \mathbf{x}\|^2$$

$$= \langle \mathbf{x}, (\mathbf{C}^{(i)T} \mathbf{C}^{(i)} - \mathbf{B}^{(i)T} \mathbf{B}^{(i)}) \mathbf{x} \rangle$$

$$\stackrel{\text{Cauchy Ineq}}{\leq} \|\mathbf{x}\| \cdot \|(\mathbf{C}^{(i)T} \mathbf{C}^{(i)} - \mathbf{B}^{(i)T} \mathbf{B}^{(i)}) \mathbf{x}\|.$$

$$\leq \|(\mathbf{C}^{(i)T} \mathbf{C}^{(i)} - \mathbf{B}^{(i)T} \mathbf{B}^{(i)})\|.$$

$$\text{Lem: } \|\mathbf{A}\|_F^2 = \text{Tr}(\mathbf{A}^T \mathbf{A}).$$

$$\|\mathbf{B}\|_F^2 = \sum_{i=1}^n \|\mathbf{B}^{(i)}\|_F^2 - \|\mathbf{B}^{(i)}\|_F^2$$

$$= \sum_{i=1}^n (\|\mathbf{C}^{(i)}\|_F^2 - \|\mathbf{B}^{(i)}\|_F^2) - (\|\mathbf{C}^{(i)}\|_F^2 - \|\mathbf{B}^{(i)}\|_F^2).$$

$$\stackrel{(*)}{=} \sum_{i=1}^n \|A_{ii}\|^2 - \text{Tr}(\mathbf{C}^{(i)T} \mathbf{C}^{(i)} - \mathbf{B}^{(i)T} \mathbf{B}^{(i)}).$$

$$= \|\mathbf{A}\|_F^2 - \text{Tr}((\mathbf{D}^{(i)})^2 - (\tilde{\mathbf{D}}^{(i)})^2)$$

$$\leq \|\mathbf{A}\|_F^2 - k/2 \cdot \delta_{ii}^{(i)2}$$

日期:

$$\Rightarrow \underline{g^{(i)} \leq 2(\|A\|_F^2 - \|B\|_F^2)/k} \leq \underline{2\|A\|_F^2/k}$$

Space:  $O(d \cdot k)$ . ?

Time:  $O(ndk)$

s 3

Machine  
Learning

日期:

## Lecture 8. Classification PAC learning.

(Model, Optimization, Generalization).

Space  $X \subseteq \mathbb{R}^d$

Goal: Find a "simple" rule performs well on training data  
and as long as training data are representative to future data,  
then it also performs well on future data.

I: The perception algorithm:

Problem: Labels:  $l_1, \dots, l_n$  (Assume hyperplane 通过原点).  
examples:  $x_1, \dots, x_n \in X$ .

Goal: to find weight vector  $w$ , and a threshold  $t$ , s.t.

$$w^T x_i > t \text{ if } l_i = 1.$$

⇒ 要求线性可分性.

$$w^T x_i < t \text{ if } l_i = -1.$$

Solutions: (Assuming  $\exists w, t$  ).

(1) linear programming algorithm:

$$(\min c^T x).$$

solvable but time-consuming.

$$\text{s.t. } Ax \geq b.$$

(2) The Perception alg:

$$\text{Define } \hat{x}_i = l_i \cdot \begin{pmatrix} x_i \\ 1 \end{pmatrix}. \quad \hat{w} = \begin{pmatrix} w \\ -t \end{pmatrix}$$

日期:

$$\text{then, } \begin{cases} w^T x_i > t & \text{if } l_i = 1 \\ w^T x_i < t & \text{if } l_i = -1. \end{cases} \Leftrightarrow \hat{w}^T \hat{x}_i > 0.$$

Rmk:  $\hat{l} = \text{sign}(\hat{w}^T \hat{x})$ . 在训练集上总有  $\hat{l}_i = l_i$ .

Alg:

(1) Initialize  $\hat{w} \leftarrow 0$ .

(2) if  $\exists i \text{ st. } \hat{w}^T \hat{x}_i \leq 0$ . update  $\hat{w} \leftarrow \hat{w} + l_i (\hat{x}_i)$

Thm: If  $\exists w^* \text{ st. } w^{*T} x_i l_i \geq 1 \quad \forall i$   $\xrightarrow{\text{即相隔一定距离}}$

then the alg. finds a solution  $w$  st  $w^T x_i l_i > 0 \quad \forall i$ .

and at most  $r^2 \|w^*\|^2$  updates where  $r = \max_i \|x_i\|$ .

(Rmk: ① Assume hyperplane 过原点)

$$② \text{margin} = \min_i \text{dist}(x_i, w^{*T} x = 0) = \min_i \frac{|w^{*T} x_i|}{\|w^*\|} \geq \min_i \frac{1}{\|w^*\|}$$

③  $r = \max_i \|x_i\|$  是 radius of data ball.

$$\|w^*\| = \frac{1}{\text{margin}}.$$

是 radius 的  $\frac{1}{\text{margin}}$ , margin 越大, alg 越快.

④ the output  $w$  is  $\{x_i\}$ 's linear combination. i.e.  $w = \sum c_i x_i$ .

Pf: Each update, consider  $w^T w^* \leq \|w\|^2$ :

$$① (w + x_i l_i)^T w^* = w^T w^* + (x_i l_i)^T w^* \geq w^T w^* + 1.$$

日期:

$$\begin{aligned} \textcircled{2} \quad \|w + x_i l_i\|^2 &= \|w\|^2 + 2\cancel{w^T x_i l_i} + \|x_i l_i\|^2 \\ &\leq \|w\|^2 + \|x_i l_i\|^2 \stackrel{\leq 0}{=} \\ &\leq \|w\|^2 + r^2. \end{aligned}$$

Let  $m$  be numbers of updates. then,

$$ww^* \geq m. \quad \|w\|^2 \leq mr^2.$$

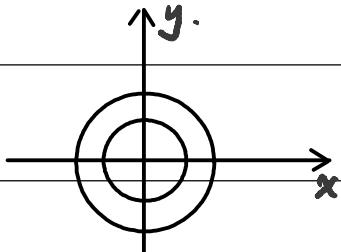
$$\therefore \text{由 Cauchy Ineq. } m \leq |ww^*|^2 \leq \|w\|^2 \|w^*\|^2 \leq mr^2 \|w^*\|^2.$$

$$\Rightarrow m \leq r^2 \|w^*\|^2.$$

(注: 可用反证法, 得最终  $w^T x_i l_i > 0 \forall i$ .)

Improvement:

Alg can't work when  $\exists w^*$ . e.g.



Solutions:

① use nonlinear function to partition.

② 在高维空间  $\nexists$  may have  $\exists w^*$ .

e.g.  $(x, y) \xrightarrow{\varphi} (x, y, x^2 + y^2)$ . then  $\exists w^*$ .

缺点:  $\varphi$  难以确定 (dim?, entries of new dim?).

③ Kernel Method.

(Perception Alg. 需要有  $k(x_i, x_j) := \varphi(x_i)^T \varphi(x_j)$ . 即可.)

日期:

Reason:

dual perception Alg.

① Initialize  $\alpha = 0 \in \mathbb{R}^d$ .  $w = \sum \alpha_i l_i x_i$ .  $k(x_i, x_j)$ .

② Iterate:  $\forall i$ , let  $\hat{y}_i = \text{sgn}(w^T x_i) = \text{sgn}(\sum \alpha_i l_i x_i^T x_i)$   
if  $\exists i$  s.t.  $\hat{y}_i \neq l_i$ , then  $\alpha_i \leftarrow \alpha_i + 1$ .

else return  $\alpha$  and  $w$ .

For Prediction of  $x$ : output  $\text{sgn}(\sum \alpha_i l_i x_i^T x)$ .

kernel Matrix: (for finite data-set).

Define  $\Phi: \mathbb{R}^d \rightarrow \mathbb{R}^{d'} (d' > d)$ . define  $k(x, y) = \Phi(x)^T \Phi(y)$ .

Kernel Matrix  $K = (K_{ij})_{\substack{i=1 \\ j=1}}^{m \times n}$ .  $K_{ij} = k(x_i, x_j)$ .

Lemma: A matrix  $K$  is a Kernel Matrix iff  $K$  is PSD

Pf: (FACT: PSD  $K \Leftrightarrow K = B^T B$  for some  $B$ .)

Define  $\Phi(x_i)$  is  $i$ -th column of  $B$ .

then  $K_{ij} = (B^T B)_{ij} = \Phi(x_i)^T \Phi(x_j)$ .

反之亦然.

日期:

Thm: If  $\mathcal{K} = \{\text{kernel functions } k(x, y)\}$ .

①  $\mathcal{K}$  is a ring. ? geometric.

②  $\forall$  scalar function  $f$ ,  $\forall k \in \mathcal{K}$ .

$$k'(x, y) = f(x)f(y)k(x, y)$$

常用 e.g.: 1) Polynomial Kernel:  $k(x, y) = (1 + x^T y)^t$ . for  $t \geq 0$ .

2) Gaussian Kernel:  $e^{-c\|x - y\|^2}$  ( $= f(x)f(y)e^{-2c x^T y}$ ).  
 $\downarrow e^{-c\|x\|^2}$  Taylor off.  
scalar function.

## II. SVM (支持向量机). (optimum hyperplane)

由 Perceptron Alg. Assume  $(x_i, l_i)$  is linear separable (have margin).

We can find  $w, t$ . s.t.

$$(w^T x_i + t)l_i \geq 1, \forall i \quad (<= (w^T x_i + t)l_i > 0, \forall i).$$

SVM with hard margin:

Assume linear separable. find  $(w, t)$  s.t.

$$\begin{cases} \text{maximize margin } \frac{1}{\|w\|}. (\Leftrightarrow \text{minimize } \frac{1}{2} w^T w.) \\ (w^T x_i + t)l_i \geq 1, \forall i \end{cases}$$

(Rmk: It's a convex programming. 有 poly. time · solution.).

日期:

Def: support vector:  $\frac{1}{2} \nparallel (w^T x_i + b) l_i = 1$  线上的 Vectors.

SVM with soft margin:

No linear separable assumption.

Try: ① Kernel method.

② minimize<sub>w,t</sub> # mistakes under  $(w^T x_j + b) l_j \geq 1 \quad \forall j$ .

③ minimize<sub>w,t</sub> ( $\frac{1}{2} w^T w + C \cdot \# \text{mistakes}$ )  $\Rightarrow$  NP-hard Problem.

考慮到  $(w^T x_j + b) l_j \geq 1 \quad \forall j$ .

$$\text{loss}_j = \max(0, 1 - (w^T x_j + b) l_j).$$

$$\Rightarrow \underset{w,t}{\text{minimize}} \frac{1}{2} w^T w + C \cdot \sum_j \max(0, 1 - (w^T x_j + b) l_j).$$

$$\left\{ \begin{array}{l} \underset{w,t}{\text{minimize}} \frac{1}{2} w^T w + C \cdot \sum \xi_j \\ (w^T x_j + b) l_j \geq 1 - \xi_j \quad \forall j \end{array} \right.$$

$$\left. \begin{array}{l} \xi_j \geq 0 \end{array} \right.$$

$$L(w, t, \xi, \alpha, \beta) = \frac{1}{2} w^T w + C \sum \xi_j - \sum \alpha_j [(w^T x_j + b) l_j - 1 + \xi_j] - \sum \beta_j \xi_j.$$

↓ KKT 条件

$$\left\{ \begin{array}{l} \nabla_w L = w - \sum \alpha_j l_j x_j = 0. \quad (1) \\ \nabla_t L = -\sum \alpha_j l_j = 0. \quad (2) \end{array} \right.$$

$$\left\{ \begin{array}{l} \nabla_{\xi_j} L = C - \alpha_j - \beta_j = 0. \quad \forall j. \quad (3) \\ \alpha_j ((w^T x_j + b) l_j - 1 + \xi_j) = 0. \quad \forall j. \quad (4) \end{array} \right.$$

$$\left\{ \begin{array}{l} \beta_j \xi_j = 0. \quad \forall j. \quad (5) \end{array} \right.$$

(\*1)

$$\text{对 } \alpha_i \in (0, C) \Rightarrow \beta_i > 0 \Rightarrow \xi_i = 0 \Rightarrow t = l_i - w^T x_i. \quad (6)$$

$$\text{日期: } \sum_j \alpha_j = \sum (\beta_j + \alpha_j) \stackrel{3^{\circ}}{=} \sum \alpha_j \stackrel{5^{\circ}}{=} \sum \alpha_j (1 - b_j(w^T x_j + t)) \stackrel{4^{\circ}}{=} \sum \alpha_j - w^T v - \sum \alpha_j b_j t \stackrel{=0 \quad (2)}{=} \sum \alpha_j - w^T w. \quad (7)$$

$\Rightarrow$  Dual Problem:

$$K(x_i, x_j).$$

$$\begin{cases} \text{maximize}_{\alpha} L(w, t, \beta, \alpha, \beta) = \sum \alpha_i - \frac{1}{2} \sum \alpha_i \alpha_j b_j K(x_i, x_j) & (7) \\ \sum \alpha_i b_i = 0. & (2) \end{cases}$$

(\*)

求解 (\*) 后, 由 (1).  $w = \sum \alpha_j b_j x_j$ . (1)  
 (凸优化在多项式时间内可解。  
 多用椭球法等方法).  $t = b_{j_0} - w^T x_{j_0}$  (for some  $j_0$  st  $\alpha_{j_0} \in (0, c)$ ) (6)

$\Rightarrow$  classification function:

$$\begin{aligned} f(x) &= \text{sgn}(w^T x + t) \\ &= \text{sgn}(\underbrace{\sum \alpha_j b_j x_j^T x}_{K(x_j, x)} + b_{j_0} - \underbrace{\sum \alpha_i b_i x_i^T x_{j_0}}_{K(x_i, x_{j_0})}). \end{aligned}$$

(Rmk: SVM 同样可用 Kernel Method 优化).

iii PAC Learning (Probably Approximately Correct) (适用于离散假设集)

Want to generalize the rule learned from training data.

$\Rightarrow$  need data to be "representative".

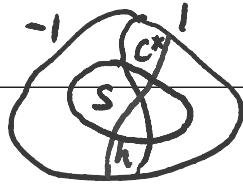
推广.  
 $\Rightarrow$  均匀抽样  $\Rightarrow$  随意分布抽样

Assumption:  $\exists$  a prob. distribution  $D$  over instance space  $X$ .

st. 1) training set  $S$  consists of  $n$  points drawn independently  
 at random from  $D$ .

2) Goal is to predict well on new points drawn from  $D$ .

日期:



Def: 1).  $C^*$ : target concept:

a subset of  $X$  corresponding to positive class in  $X$ .

2)  $h$ : hypothesis

3) true/generalization error of  $h$ :

$$\text{err}_D(h) = \Pr_{x \sim D} [x \in h \Delta C^*].$$

4) training/empirical error:

$$\text{err}_S(h) = \frac{\# S \cap (h \Delta C^*)}{\# S}.$$

5) Hypothesis class  $\mathcal{H}$  over  $X$ . ( $\forall h \in \mathcal{H}$  is a hypothesis.).

Goal: to produce  $h$  with low. true error.

Intuition: If  $S$  is large enough, then w.h.p.  $\forall h \in \mathcal{H}$  have:

$$\text{err}_S(h) \approx \text{err}_D(h).$$

Thm (PAC Learning).

Let  $\mathcal{H}$  be a hypothesis class,  $\epsilon, \delta \in (0,1)$ .

If training set  $S$  of size  $n \geq \frac{1}{\epsilon^2}(\ln|\mathcal{H}| + \ln(1/\delta))$  is drawn from distribution  $D$ ,

then w.p.  $\geq 1 - \delta$  every  $h \in \mathcal{H}$  with training error 0 has true error  $\leq \epsilon$ .

(Rmk:  $h \in \mathcal{H}$  with 0 training error is Probably Approximately Correct (PAC))

Pf: 逆否命題: w.p.  $\geq 1 - \delta$ .  $\forall h \in \mathcal{H}$  with  $\text{err}_S(h) \geq \epsilon$ , has  $\text{err}_D(h) > 0$ .

日期:

fix  $h$  with  $\text{err}_S(h) \geq \varepsilon$ . then.

$$\Pr_{S \sim D} [\text{err}_S(h) = 0] \leq (1 - \text{err}_D(h))^{|S|} \leq (1 - \varepsilon)^n.$$

then  $\Pr[\exists h \in \mathcal{H} \text{ with } \text{err}_D(h) \geq \varepsilon \text{ has } \text{err}_S(h) = 0]$

$$\leq |\mathcal{H}| \cdot \Pr_{S \sim D} [\text{err}_S(h) = 0]$$

$$= |\mathcal{H}| \cdot (1 - \varepsilon)^n$$

$$\leq |\mathcal{H}| \cdot e^{-\varepsilon n}$$

$$\leq |\mathcal{H}| \cdot e^{-(\ln |\mathcal{H}| + \ln(1/\delta))} = \delta.$$

□.

Thm (PAC Learning 2):

Let  $\mathcal{H}$  be a hypothesis class,  $\varepsilon, \delta \in (0, 1)$ .

If training set  $S$  of size  $n \geq \frac{1}{2\varepsilon^2}(\ln |\mathcal{H}| + \ln(2/\delta))$  is drawn from distribution  $D$ ,

then w.p.  $\geq 1 - \delta$ .  $\forall h \in \mathcal{H}$  有  $|\text{err}_S(h) - \text{err}_D(h)| \leq \varepsilon$ .

Pf: (思路同上: fix  $h \Rightarrow P[h \text{ opposite}] \leq *$   $\Rightarrow$  Union bound).

fix  $h \in \mathcal{H}$ . for  $j$ -th element in  $S$ . define

$$X_j = \begin{cases} 1 & \text{if } h \text{ disagree with } c^* \text{ on } j\text{-th element} \\ 0 & \text{o.w.} \end{cases}$$

$$\text{then. } P[X_j = 1] = \text{err}_D(h). \quad \text{err}_S(h) = \frac{\sum_j X_j}{|S|}.$$

$$E[\text{err}_S(h)] = \frac{1}{|S|} \sum_j E[X_j] = \text{err}_D(h).$$

By Chernoff-Hoeffding,

日期:

$$\Pr[|\text{err}_{\mathcal{S}}(h) - \text{err}_D(h)| > \varepsilon] \leq 2 e^{-2\varepsilon^2 n}.$$

Then,  $\Pr[\exists h \in \mathcal{H}, \text{st. } |\text{err}_{\mathcal{S}}(h) - \text{err}_D(h)| > \varepsilon]$

$$\leq 2|\mathcal{H}| \cdot e^{-2n \cdot \varepsilon^2}$$
$$\leq 8.$$

□

#### IV. Occam's Razor:

Thought: Simple explanations are better than complicated ones.

#### Thm (Occam's Razor)

w.p.  $\geq 1-\delta$ ,  $\forall h$  with  $\text{err}_{\mathcal{S}}(h)=0$  that can be described using  $< b$  bits.

has  $\text{err}_D(h) \leq \varepsilon$ . if.  $|S| \geq \frac{1}{\varepsilon}(\ln 2 + \ln(1/\delta))$ .

(Rmk: 即 w.p.  $\geq 1-\delta$ ,  $\forall h$  with  $\text{err}_{\mathcal{S}}(h)=0$  that can be described

using  $< b$  bits. has  $\text{err}_D(h) \leq \frac{\ln 2 + \ln(1/\delta)}{|S|}$

Pf: let  $\mathcal{H} = \{h : \forall h \text{ can be described using } < b \text{ bits}\}$ .

then  $|\mathcal{H}| \leq 2^b$ . plugging it into Thm (PAC). □

#### V. VC dimension $\star$ (适用于连续假设集)

Def: System Set  $(X, \mathcal{H})$ .

$X$ : a set, instance space.

$\mathcal{H}$ : a class of subsets of  $X$ . hypothesis class.

Def:  $(X, \mathcal{H})$  shatters a set  $A \subseteq X$ :

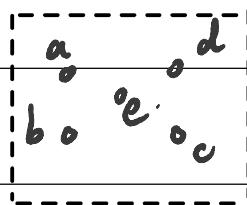
日期:

$\forall$  subset of  $A$  can be expressed as  $Anh$  for some  $h \in \mathcal{H}$ .

Def: VC dimension of  $(X, \mathcal{H})$ :

$$\max_{(A, \mathcal{H}) \text{ shatters } A} \# A.$$

eg1:  $X = \mathbb{R}^2$ .  $\mathcal{H} = \{\text{all rectangles}\}$ .



$$VC(X, \mathcal{H}) = 4.$$

当  $A = \{a, b, c, d, e\}$ . 必有一个在中间, 不妨为  $e$ . 则  $\{a, b, c, d\} \neq Anh$  for some  $h$

eg2:  $X = \mathbb{R}^2$ .  $\mathcal{H} = \{[a, b] \mid a \leq b\}$ .

$$VC(X, \mathcal{H}) = 2.$$

eg3: Consider  $X = \mathbb{R}^d$ .

linear separators:  $\mathcal{H} = \{w^T x + b \geq 0, w \in \mathbb{R}^d, b \in \mathbb{R}\}$



(Rmk:  $\#\mathcal{H} = +\infty$ , 不能直接用 PAC Learning)

Thm:  $VC(X, \mathcal{H}) = d+1$ .

Pf: ① 取  $A = \{e_1, e_2, \dots, e_d, 0_d\}$ .  $e_i = \begin{pmatrix} 0 \\ \vdots \\ i \\ 0 \end{pmatrix}$   $0_d = \begin{pmatrix} 0 \\ \vdots \\ 0 \end{pmatrix}$ .

1° for  $\forall$  subset  $A' = \{e_{i_1}, \dots, e_{i_j}, 0_d\}$ .

$$\exists w' = 1_d - \sum_{k=1}^j e_{ik} \quad b=0.$$

则  $A' = A \cap \{w'^T x + b \leq 0\}$ .

日期:

$2^d$  for  $\forall$  subset  $A' = \{e_{i_1}, \dots, e_{i_d}\}$ .

$$\text{tr } w' = 1_d - 2 \sum_{k=1}^d e_{ik}, \quad b=1.$$

$\therefore A' = A \cap \{w'^T x + b \leq 0\}. \quad \therefore \text{VC}(X, H) \geq d+1.$

② By Radon Thm,  $\text{VC}(X, H) \leq d+1$ .  $\square$ .

Def:  $S$  is a set of  $s$  vertices, then

convex hull of  $S$ :  $\text{convex}(S) = \left\{ \sum_{i=1}^s \alpha_i a_i \mid \alpha_i \geq 0, \sum \alpha_i = 1, a_1, \dots, a_s \in S \right\}$ .

Thm(Radon): Any set  $S \subseteq \mathbb{R}^d$  with  $|S| \geq d+2$  can be partitioned into two disjoint subsets  $S_1, S_2$  s.t.  $\text{convex}(S_1) \cap \text{convex}(S_2) \neq \emptyset$ .

Pf: Assume  $S = \{a_1, \dots, a_{d+2}\}$ .

Let  $A = (a_1, \dots, a_{d+2}) \in \mathbb{R}^{d \times (d+2)}$ .

$B = \begin{pmatrix} A \\ 1^T \end{pmatrix} \in \mathbb{R}^{(d+1) \times (d+2)}$

Consider  $Bx = 0$ .  $\because \text{rank}(B) \leq d+1$ .  $x \in \mathbb{R}^{d+2}$ .

$\Rightarrow$  has a non-zero solution  $x = (x_1, \dots, x_{d+2})$ .

Reorder columns s.t.  $x_1, \dots, x_s \geq 0$ .  $x_{s+1}, \dots, x_{d+2} < 0$ .

$\Rightarrow$  normalize  $x$  s.t.  $\sum_{i=1}^s |x_i| = 1$ .

Let  $b_i = \begin{pmatrix} x_i \\ 1 \end{pmatrix}$ . then  $B = (b_1, \dots, b_{d+2})$ . then  $\sum_{i=1}^{d+2} b_i x_i = 0$

$\Rightarrow \sum_{i=1}^s b_i / |x_i| = \sum_{i=s+1}^{d+2} b_i / |x_i|$ .

日期:

$$b_i = \binom{a_i}{1} \Rightarrow \sum_{i=1}^s a_i |x_i| = \sum_{i=s+1}^{d+2} a_i |x_i| := a.$$

$$\sum_{i=1}^s |x_i| = \sum_{i=s+1}^{d+2} |x_i| (=1).$$

2)  $a \in \text{convex}(a_1, \dots, a_s) \cap \text{convex}(a_{s+1}, \dots, a_{d+2})$ .  $\square$ .

eg.4: spheres in  $d$ -dimension:

$$X = \mathbb{R}^d. \quad \mathcal{H} = \{x \mid \|x - x_0\| \leq r\}. \quad \text{VC}(X, \mathcal{H}) = d+1. \quad (\text{不证}).$$

eg.5.  $X = \mathbb{R}$ .  $\mathcal{H} = \{h \mid h \text{ contains finite number of real numbers}\}$ .

$$\underline{\text{VC}(X, \mathcal{H}) = \infty}.$$

Rmk: VC-dimension 可用于描述  $(X, \mathcal{H})$ 's complexity.

Shatter function (extension of VC dimension).

Def: Given  $(X, \mathcal{H})$ , the shatter function

$$\pi_{\mathcal{H}}(n) := \max_{|\mathcal{A}|=n} \#\{A' \subseteq A \mid A' = \bigcup_{h \in \mathcal{H}} h\}.$$

(Rmk: if  $\text{VC}(X, \mathcal{H}) = n_0$ , then  $\pi_{\mathcal{H}}(n_0) = 2^{n_0}$ .  $\pi_{\mathcal{H}}(n_0+1) < 2^{n_0+1}$ ).

Lemma (Sauer) (用 VC dimension 估算  $\pi_{\mathcal{H}}(n)$ ).

$$\text{VC}(X, \mathcal{H}) = d. \text{ then } \pi_{\mathcal{H}}(n) \leq \binom{n}{\leq d} \triangleq \binom{n}{0} + \binom{n}{1} + \dots + \binom{n}{d} \leq n^d + 1.$$

Thm:  $(X, \mathcal{H})$ .  $X \sim D$ . let  $n$  s.t.  $n \geq \frac{8}{\epsilon}$  且  $n \geq \frac{2}{\epsilon} [\log_2(2\pi_{\mathcal{H}}(2n)) + \log_2(4/\delta)]$ .  
 (Key Thm).  $\frac{4(1-\delta)}{\epsilon} \rightarrow 2\pi_{\mathcal{H}}(2n)$ .

Let  $S \sim D$  且  $\#S = n$ . w.p.  $\geq 1-\delta$   $\forall h \in \mathcal{H}$  w.p.  $\geq 1-\epsilon$  intersects  $S$ .  
 (针对  $D$ )

日期: Double Sampling.

Pf:  $S_1 \sim D$ ,  $S_2 \sim D$ .

def A:  $\exists h \in H$  w.p.  $\geq \varepsilon$  disjoint from  $S_1$

B:  $\exists h \in H$  w.p.  $\geq \varepsilon$  st.

$h \cap S_1 = \emptyset$  且  $\#h \cap S_2 \geq \frac{\varepsilon n}{2}$ .

Lemma 1:  $P_r[B|A] \geq \frac{1}{2}$

Lemma 2:  $P_r[B] \leq \frac{\delta}{2}$ .

Assume the above correct, then

$$\frac{\delta}{2} \geq P_r[B] \geq P_r[AB] = P_r[B|A]P_r[A] \geq \frac{1}{2}P_r[A]. \quad \square$$

Lemma 2 Pf:

st.  $S_1, S_2$  产生 connection.  
fix  $S_1, S_2$  random  $h \rightarrow$  fix  $h$ , random  $S_1, S_2$ .

$S_1 \sim D$ ,  $S_2 \sim D$

依序抽样  $\rightarrow$  混合抽样 + 划分

$\Leftrightarrow$  double sampling:

转换随机对象.

1°  $S_3 \sim D$  st.  $\#S_3 = 2n$ .

2° randomly partition  $S_3$  into  $S_1, S_2$ . st.  $\#S_1 = \#S_2$

于是, 有  $\#\{S_3 \cap h \mid h \in H\} \leq \pi_H(2n)$ .

Lemma 3: Given  $h \in H$  st  $h \cap S_3 \neq \emptyset$ . i.e.  $h = h \cap S_3$ .

randomly partition  $S_3$  into  $S_1, S_2$ .

then  $\#S_1 \cap h' = 0$ ,  $\#S_2 \cap h' \geq \frac{\varepsilon n}{2}$  w.p.  $\leq \frac{\delta}{2\pi_H(2n)}$

日期:

Assume correct, then  $\Pr[B] \leq \frac{\delta}{2\pi_{\mathcal{H}}(2n)} \cdot \pi_{\mathcal{H}}(2n) = \frac{\delta}{2}$ .

Pf of Lemma 3:

case I:  $\# h \cap S_3 < \frac{\varepsilon n}{2}$

$$\Pr[\# h' \cap S_1 = 0, \# S_2 \cap h' \geq \frac{\varepsilon n}{2}] = 0.$$

case II:  $\# h \cap S_3 \geq \frac{\varepsilon n}{2}$

$$\Pr[\# h' \cap S_1 = 0, \# S_2 \cap h' \geq \frac{\varepsilon n}{2}]$$

$$= \Pr[\# h \cap S_1 = 0].$$

$$\leq (\frac{1}{2})^{\frac{\varepsilon n}{2}} \leq \frac{\delta}{2\pi_{\mathcal{H}}(2n)}. \quad \text{By Union Bound. } \square.$$

Pf of Lemma 1:

let  $h$  be a set w.p.  $\geq \varepsilon$   $S_1 \cap h = \emptyset$ .

let  $X_i = \begin{cases} 1, & \text{the } i\text{-th point in } S_2 \text{ belong to } h. \\ 0, & \text{o.w.} \end{cases}$

$$\text{let } X = \sum X_i.$$

$$\therefore \eta := E[X_i] \geq \varepsilon. \quad \text{Var}[X_i] = \eta - \eta^2 < \eta.$$

$$\Rightarrow E[X] = n\eta. \quad \text{Var}[X] < n\eta.$$

$$\downarrow n \geq \frac{8}{\varepsilon} \geq \frac{8}{\eta}$$

$$\Rightarrow \Pr[|X - E[X]| \geq \frac{n\eta}{2}] \leq \frac{n\eta}{(n\eta/2)^2} \leq \frac{4}{n\eta} \leq \frac{1}{2}$$

$$\Rightarrow \Pr[B | A] = \Pr[X \geq \frac{\varepsilon n}{2}] \geq \Pr[X \geq \frac{\eta n}{2}] \geq \Pr[|X - \eta n| \leq \frac{\eta n}{2}]$$

$$= \Pr[|X - E[X]| \leq \frac{n\eta}{2}] \geq \frac{1}{2}.$$

日期:

即预测错误的部分.

Prop: A target concept  $c^*$ . let  $\mathcal{H}' = \{h \Delta c^* \mid h \in \mathcal{H}\}$

Fact:  $VCC(X, \mathcal{H}') = VCC(X, \mathcal{H})$

and  $(X, \mathcal{H}')$ ,  $(X, \mathcal{H})$  have same shutter function.

Thm1:  $(X, \mathcal{H})$ .  $X \sim D$ . let  $n$  s.t.  $n \geq \frac{8}{\epsilon^2} \text{ 且 } n \geq \frac{2}{\epsilon} [\log_2(\bar{\alpha}_{\mathcal{H}}(2n)) + \log_2(1/\delta)]$ .

$S \stackrel{D}{\sim} X$ . Then w.p.  $\geq 1-\delta$ .  $\forall h \in \mathcal{H}$  with  $\text{err}_S(h) = 0$ . has  $\text{err}_D(h) \leq \epsilon$ .

Pf:  $\forall h$  反证: Assume  $\text{err}_D(h) > \epsilon$ .

Apply Key Thm to  $\mathcal{H}' = \{h \Delta c^* \mid h \in \mathcal{H}\}$ .

Then, w.p.  $\geq 1-\delta$ ,  $\text{err}_S(h) = \Pr_{x \sim D}[x \in h \Delta c^*] > \epsilon$ . has

$|S \cap h \Delta c^*| > 0$ . 这与  $\text{err}_S(h) = 0$  矛盾.  $\square$ .

Thm2:  $(X, \mathcal{H})$ .  $X \sim D$ . let  $n \geq \frac{8}{\epsilon^2} \text{ 且 } n \geq \frac{8}{\epsilon^2} [\log_2(\bar{\alpha}_{\mathcal{H}}(2n)) + \log_2(1/\delta)]$ .

$S \stackrel{D}{\sim} X$ . Then w.p.  $\geq 1-\delta$ .  $\forall h \in \mathcal{H}$ ,  $|\text{err}_S(h) - \text{err}_D(h)| \leq \epsilon$ .

Rmk: 需要将 key Thm 修改为  $|S_1 \cap h| - |S_2 \cap h| \geq \frac{\epsilon n}{2}$  ?

并用 Chernoff Bound 做改进.

Rmk: ① By Sauer Bound. denote  $d = VCC(X, X)$ . Improvement:

$$\begin{aligned} \text{Thm1 : } n &\geq \frac{c}{\epsilon^2} (d \log n + \log(\frac{1}{\delta})). \\ &n \geq \frac{c'}{\epsilon^2} (d \log \frac{1}{\epsilon} + d \log \frac{1}{\delta}) \end{aligned}$$

$$\text{Thm2 : } n \geq \frac{c'}{\epsilon^2} (d \log \frac{1}{\epsilon} + d \log \frac{1}{\delta}).$$

②  $|\mathcal{H}| < \infty$  时,  $\bar{\alpha}_{\mathcal{H}}(2n)$  不符 Thm2 表现很可能不如 Thm1 版本

(除非  $\bar{\alpha}_{\mathcal{H}}(2n) < \log |\mathcal{H}|$ ).

sy

Clustering

日期:

Input:  $n$  points  $\in \mathbb{R}^d / \{0,1\}^d$ .

Goal:  $k$  clusters.

{ centerbased clustering. eg. k-means, k-median, k-centers.  
{ high-density clustering. eg. PBSCAN.

I. Centerbased clustering.

Form in Euclidean space:

$$A \subseteq M = \mathbb{R}^d. D(x,y) = \|x-y\|_2 \quad \forall x,y \in \mathbb{R}^d$$

Find  $C = \{c_1, \dots, c_k\}$  s.t.  $\sum_i D(a_i, \{c_1, \dots, c_k\})^q$  for some  $q$ . is minimized.

then have  $c_i = \{a \in A \mid \forall j, D(a, c_i)^q \leq D(a, c_j)^q\}$ .

ties broken arbitrary. ( $\Rightarrow \exists D(a, c_i) = D(a, c_j) = \min$ , randomly choose one)

k-centers:  $q = +\infty \Rightarrow \min_{\{c_1, \dots, c_k\}} \max_i D(a_i, \{c_1, \dots, c_k\})$ . (+∞)

k-median:  $q = 1 \Rightarrow \min_{\{c_1, \dots, c_k\}} \sum_i D(a_i, \{c_1, \dots, c_k\})$ .

k-means:  $q = 2 \Rightarrow \min_{\{c_1, \dots, c_k\}} \sum_i D(a_i, \{c_1, \dots, c_k\})^2$ .

Form in general space:

Def: (Metric Space).

$M = \mathbb{R}^d / \{0,1\}^d / \text{other set}$

$D: M \times M \rightarrow \mathbb{R}$ . for  $\forall x, y, z \in M$ . st.

日期:

1°  $D(x, y) = 0 \Rightarrow x = y$ .

2°  $D(x, y) = D(y, x)$ .

3°  $D(x, y) \leq D(x, z) + D(y, z)$ .

then  $(M, D)$  is a Metric Space.

e.g.  $D_{L_2}(x, y) = \|x - y\|_2 \rightarrow$  Euclidean space.

$D_{L_p}(x, y) = \|x - y\|_p \rightarrow L_p$  space.

(Rmk: form in general space  $\rightarrow$  replace  $(M, D)$  by a Metric Space).

The exact opt problem is NP-hard.

Lloyd's algorithm (for k-means clustering. in Euclidean space).

$$A \subseteq M = \mathbb{R}^d, k \geq 1.$$

① choose  $k$  centers  $c_1, \dots, c_k \in \mathbb{R}^d$ . arbitrarily.

② Repeat: (Until convergence).  
     $\xrightarrow{\text{迭代}} \text{eg. loss 变化不大.}$

1° find clusters  $C_1, \dots, C_k \rightarrow$  (assignment)

2° compute centroid  $c(C_i) = \frac{1}{|C_i|} \sum a$  } update.

3°  $c_i \leftarrow c(C_i)$ .

③ Return  $c_1, \dots, c_k, C_1, \dots, C_k$ .

日期:



Lemma 1:  $\forall x \in \mathbb{R}^d$

$$\sum_{a \in A'} \|a - x\|^2 = \sum_{a \in A'} \|a - c(A')\|^2 + |A'| \cdot \|c(A') - x\|^2.$$

Pf: LHS =  $\sum_{a \in A'} \|a - c(A') + c(A') - x\|^2$

$$= \sum_{a \in A'} (\|a - c(A')\|^2 + \|c(A') - x\|^2 + 2 \cdot (c(A') - x)^T \cdot (a - c(A'))).$$

$$= \sum_{a \in A'} \|a - c(A')\|^2 + |A'| \cdot \|c(A') - x\|^2 + 2(c(A') - x)^T \sum_{a \in A'} (a - c(A')).$$

$$= RHS \quad \square.$$

Lemma 2: halt after finite steps. the running time is  $O(nkdR^*)$ .

$R^*$  is the number of iterations until convergence.

(Rmk: 有人证明],  $R^* \leq 2^{\frac{nkd}{n}}$ ).

Pf: Loss  $\downarrow$  each iteration. 因此.

(Rmk: 某些实例中 k-means 的结果与最优解相差较远, 但一般效果不错)

(Rmk: 注意到主要是初始点集较差时效果差, 尝试改变抽样方式).

补充:

Thm: k-means Problem is NP-hard.

Thm: If P  $\neq$  NP, then no poly. time alg. can approximate k-means

Problem within 1.0013 multiplicative factor.

日期:

$D^2$ -sampling + Lloyd's.

## II. k-means++ algorithm.

Input  $A, k$ .

Alg: ① Run  $D^2$ -sampling on  $A, k$  to obtain  $C = \{c_1, \dots, c_k\}$ .

② Run k-means on  $A, k$  with  $C$  initially.

$D^2$  distribution: (针对  $S$ ).

for  $\forall S \subseteq \mathbb{R}^d$ ,  $|S| < \infty$ . denote the distribution on  $A$  by:

$$P_S(a) = \frac{D(a, S)}{\sum_{b \in A} D(b, S)} = \frac{\min_{s \in S} D(a, s)}{\sum_{b \in A} \min_{s \in S} D(b, s)} \quad (\text{设 } D := D_{(2)}^2).$$

$D^2$ -sampling:

① choose point  $c_1 \in A$  u.a.r. set  $C' \leftarrow C_1$ .

② loop: ( $2 \sim k$ ):

draw point  $c_i \in A$  according to  $D^2$ -distribution  $P_{C'^i}(\cdot)$ .

set  $C' \leftarrow C'^i \cup \{c_i\}$ .

Thm: let  $C$  be output of  $D^2$ -sampling,  $C^*$  be the opt. solution for

input  $A, k$ . Then,

$$\Leftrightarrow E[\text{cost}] \leq 8(\ln k + 2) \cdot \text{OPT}.$$

$$E[D(A, C)] \leq 8(\ln k + 2) D(A, C^*).$$

(Rmk: This implies output  $C'$  of k-means also st. above ineq.)

Analysis:

日期:

Lem 1  $\star$ : let  $S \subseteq \mathbb{R}^d$  be a point set,  $c \in S$  a point chosen

w.a.r. from  $S$ . Then,

$$E[D(S, \{c\})] = 2 \cdot D(S, c(S))$$

where  $c(S)$  is centroid of  $S$ .

Pf:  $E[D(S, \{c\})] = \sum_{b \in S} \frac{1}{|S|} D(S, \{b\}).$

$$= \frac{1}{|S|} \sum_{b \in S} \sum_{a \in S} \|a - b\|^2$$

$$= \frac{1}{|S|} \sum_{b \in S} \left( \sum_{a \in S} \|a - c(S)\|^2 + |S| \cdot \|c(S) - b\|^2 \right)$$

$$= \sum_{b \in S} \|b - c(S)\|^2 + \sum_{a \in S} \|a - c(S)\|^2$$

$$= 2D(S, c(S)). \quad \square.$$

Cor: ①  $\sum_{a \in S} \|a - c(S)\|^2 = \frac{1}{2|S|} \cdot \sum_{a \in S} \sum_{b \in S} \|a - b\|^2.$

② let  $S$  be a cluster from optimal clustering  $C^* = \{C_1^*, \dots, C_k^*\}$ .

If  $c \stackrel{\text{w.a.r}}{\sim} A$ , then

$$E[D(S, \{c\}) | c \in S] = 2 \cdot D(S, c(S)).$$

Lem 2:  $S \subseteq \mathbb{R}^d$  point set.  $C \subseteq \mathbb{R}^d$ . Both finite set.

$x$  is chosen from  $S$  by  $P_C(\cdot)$ . Then,

$$E[D(S, C \cup \{x\})] \leq 8 D(S, c(S)).$$

看论文

日期:

Rmk: Time  $\geq \Omega(d)$ .

## II. Dimension Reduction of k-means:

Another Version:  $A = \begin{pmatrix} a_1 \\ \vdots \\ a_n \end{pmatrix} \in \mathbb{R}^{n \times d}$ . centers set  $\mathcal{C} \subseteq \mathbb{R}^d$ .  $\#\mathcal{C} = k$ .

denote  $S(\mathcal{C}) = \begin{bmatrix} \underset{c \in \mathcal{C}}{\operatorname{arg\min}} \|a_1 - c\|_2 \\ \vdots \\ \underset{c \in \mathcal{C}}{\operatorname{arg\min}} \|a_n - c\|_2 \end{bmatrix} \in \mathbb{R}^{n \times d}$ .

$$\arg \min_{B \in \mathbb{R}^{k \times d}} \|A - B\|_F^2 = A_k$$

$$\text{then cost } D(A, \mathcal{C}) = \sum_{c \in \mathcal{C}} \min_{a_i \in A} \|a_i - c\|^2 \\ = \|A - S(\mathcal{C})\|_F^2.$$

## Partition-based k-means:

find a partition  $A_1, \dots, A_k$  of  $A$

to minimize  $D((A_i)_{i=1}^k) := \sum_{i=1}^k D(A_i - c(A_i))$   
↳ centroid.

Rmk: Partition-based Problem's solutions  $\supseteq$  Center-based Problem's solutions

(OPT-par. 不一定是 OPT-cen.).

且若 Center-based Problem find OPT.

then Partition-based Problem find OPT.

Apply JL-lemma:

$|A| = n$ . let  $\varepsilon \in (0, 1)$ .  $d \geq \Omega\left(\frac{\log n}{\varepsilon^2}\right)$ . f random projection.

then for any partition  $A_1, \dots, A_k$  of  $A$ ,

$$|D((A_i)_{i=1}^k) - D(f(A_i)_{i=1}^k)| \leq \varepsilon \cdot D((A_i)_{i=1}^k)$$

(由 JL-lemma 3 证)

日期:

JL-DimRed-k-means ( $A, k, \varepsilon$ ):

1°  $\varepsilon' = \frac{\varepsilon}{6d}$ .  $d' \geq \Omega(\frac{\log n}{\varepsilon'^2})$ .  $f$  from JL-lemma.

2° Use a  $\alpha$ -approximation alg. for k-means on  $f(A)$ . to obtain a partition  $(Af_i)_{i=1}^k$ .

3° Translate  $(Af_i)_{i=1}^k$  to  $(A_i)_{i=1}^k$ . Return  $(A_i)_{i=1}^k$ .

⇒  $(2+\varepsilon)$ -approximation w.h.p.

#### IV. Other Problems.

##### Local Search.

LocalSearch k-means ( $A \subseteq \mathbb{R}^d, k \in \mathbb{N}$ ).

① choose  $k$  centers arbitrarily  $T = \{c_1, \dots, c_k\}$ .

② while  $\exists c \in T, c' \in A$  st.  $D(A, T + \{c'\} - \{c\}) < D(A, T)$ .

$$C \leftarrow C + \{c'\} - \{c\}.$$

③ Return  $C$ .

Rank: 50-approximation.

##### k-median:

$A \subseteq \mathbb{R}^d$ .  $d \geq 1$ . find  $\ell \subseteq \mathbb{R}^d$ .  $|\ell| = k$ .  $D := D_{\ell_2}$ . minimize.

$$D(A, \ell) = \sum_{a \in A} \min_{c \in \ell} D(a, c).$$

日期:

Rmk: ① NP-hard.

② Linear-programming  $\Rightarrow$  constant-app.

③ Local search  $\Rightarrow$  b-app.

k-center:

$A \subseteq \mathbb{R}^d$ .  $d \geq 1$ . find  $C \subseteq A$ .  $|C|=k$ .  $D := D_{C_2}$ . minimize.

$$\bar{D}(A, C) = \sum_{c \in C} \max_{a \in A} D(a, c).$$

Rmk: NP-hard.

Gonzalez's alg: (furthest-first-traversal)

① choose  $C_1$  arbitrarily.

② For  $i=2, \dots, d$ . do:

choose  $c_i$  maximize  $\bar{D}(x, C^{i-1})$ .

$$\text{i.e. } c_i = \arg \max_{x \in A} D(x, C^{i-1}).$$

$$C^i \leftarrow C^{i-1} \cup \{c_i\}.$$

Rmk: Time:  $O(|A| \cdot k) \cdot O(d)$   $\xrightarrow{\text{compute distance}}$

Thm: alg above outputs  $C$  s.t.  $\bar{D}(A, C) \leq 2 \bar{D}(A, C^*)$ .

V. Coreset by k-means clustering:

Goal of Coreset:  $n$  points  $\Rightarrow O(\log n)/O(\frac{k}{\epsilon})$  points.

日期:

(Pink: If dimension reduction:  $A \in \mathbb{R}^{n \times d} \rightarrow A' \in \mathbb{R}^{n \times d'} (d' < d)$ ).

Def: Let  $A \subseteq \mathbb{R}^d$  n points.  $\epsilon \in (0, 1)$ .  $k \geq 1$ .

$S \subseteq \mathbb{R}^d$  with weight function  $w: S \rightarrow \mathbb{R}^+$  is  $(k, \epsilon)$ -coreset

if for  $\forall C \subseteq \mathbb{R}^d$ :  $|C|=k$ .  $|D(A, C) - D(S, w, C)| \leq \epsilon D(A, C)$ .  
↳ Ap cost function 差别不大.

Construction Alg:

Tools:

①  $\epsilon$ -ball-cover:

$\epsilon$ -ball-cover of unit sphere is a point set  $B$

s.t.  $\forall p \in$  unit sphere,  $d(p, B) \leq \epsilon$ .

lemma:  $U \subseteq \mathbb{R}^d$  unit sphere. Then,  $\forall \epsilon \in (0, 1) \exists \epsilon$ -ball-cover  $B$

of size  $(1 + \frac{2}{\epsilon})^d$  s.t.  $\forall p, \min_{b \in B} \|p - b\| \leq \epsilon$ .

Cor: ①  $\exists$  alg. construct a cover of size  $\epsilon^{-O(d)}$ .

② if  $U$  has radius  $r$ . then ...  $\forall p, \min_{b \in B} \|p - b\| \leq \epsilon \cdot r$ .

③ App-Alg: Thm:  $\exists$  a 6.357-app. alg for k-means problems.

④ (Generalized Triangle Inequality)

$a, b, c \in \mathbb{R}^d$ . for  $\forall \epsilon \in (0, 1)$ . has.

$$|\|a - c\|^2 - \|b - c\|^2| \leq \frac{12}{\epsilon} \|a - b\|^2 + 2\epsilon \|a - c\|^2.$$

日期:

i.e.  $C \subseteq \mathbb{R}^d$ .  $|D(a,C) - D(b,C)| \leq \frac{12}{\varepsilon} D(a,b) + 2\varepsilon D(a,c)$ .

Alg: ① Input  $A, k$ .

② compute 10-app. k-means, maintain  $C_1, \dots, C_k$   $c_1, \dots, c_k$ .

③ for  $j=1, \dots, k$ :

1° let  $F = C_j$ .

2° let  $B^i$  be ball with  $r_i = \sqrt{\frac{2^i}{n} \sum_{x \in F} \|x - c_j\|^2}$ , centered  $c_j$ .

and  $S^i$  be  $\varepsilon/192$ -ball-cover of  $B^i$ .  $i=1, 2, \dots, \log(n)$ .

let  $S_j = \bigcup_{i=0}^{\log(n)} S^i$ .

3°  $\forall x \in C_j$  let  $B(x)$  be the nearest point in  $S_j$ .

$\forall y \in S_j$ , let  $w_{xy}$  be  $\#\{x : B(x)=y\}$

④  $S = \bigcup B(x)$ . Return  $(S, w)$ .

$x \in A$ :  $\sqrt{k^i \text{cluster}} \rightarrow \varepsilon\text{-ball-cover}$ .  
(Rmk: size of  $S$ :  $\underline{k \cdot \log(n) \cdot \varepsilon^{-O(d)}} \rightarrow \log(n) + 3k$ ).

Analysis:

Lemma:  $F \subseteq \mathbb{R}^d$   $n$  points.  $B^i$  is ball with  $r_i = \sqrt{\frac{2^i}{n} \sum_{x \in F} \|x - c_j\|^2}$ , centered  $c_j$ .

let  $S^i$  be  $\varepsilon/3$ -ball-cover of  $B^i$ .  $S = \bigcup_{i=0}^{\log(n)} S^i$ . Then,

$$\sum_{\substack{x \in F \\ x \in S}} \min_{s \in S} \|x - s\|^2 \leq \varepsilon^2 \sum_{x \in F} \|x - c_j\|^2$$

日期:

/ 远近分别放缩.

Pf: denote  $F_{close} := S^0 = \{y \in F \mid \|y\|^2 \leq \frac{1}{n} \sum_{j \in F} \|x - c_j\|^2\}$ .  $F_{far} = F \setminus F_{close}$ .

$$\textcircled{1} \quad \sum_{x \in F_{close}} \min_{s \in S^0} \|x - s\|^2 \leq |F_{close}| \cdot \frac{1}{n} \sum_{j \in F} \|x - c_j\|^2 \cdot (\frac{\varepsilon}{3})^2 \leq \frac{\varepsilon^2}{9} \sum_{x \in F} \|x - c_j\|^2.$$

$$\textcircled{2} \quad \min_{s \in S^2} \|x - s\|^2 \leq \frac{\varepsilon^2}{9} \cdot r_i^2 \leq \frac{4\varepsilon^2}{9} r_i^2 \leq \frac{4\varepsilon^2}{9} \|x - c_j\|^2. \quad (i \geq 1).$$

$$\Rightarrow \sum_{x \in F_{far}} \min_{s \in S} \|x - s\|^2 \leq \frac{4\varepsilon^2}{9} \sum_{x \in F_{far}} \|x - c_j\|^2 \leq \frac{4\varepsilon^2}{9} \sum_{x \in F} \|x - c_j\|^2$$

$$\begin{aligned} \Rightarrow \sum_{x \in F} \min_{s \in S} \|x - s\|^2 &\leq \sum_{x \in F} \min_{s \in S_{close}} \|x - s\|^2 + \sum_{x \in F_{far}} \min_{s \in S} \|x - s\|^2 \\ &\leq \frac{\varepsilon^2}{9} \sum_{x \in F} \|x - c_j\|^2 + \frac{4\varepsilon^2}{9} \sum_{x \in F} \|x - c_j\|^2 \\ &\leq \varepsilon^2 \sum_{x \in F} \|x - c_j\|^2 \end{aligned}$$

Thm:

Pf:  $\forall C \subseteq \mathbb{R}^d$ .  $|C| = k$ .

$$|D(A, C) - D(S, w, C)| = \left| \sum_{x \in A} \min_{c \in C} \|x - c\|^2 - \sum_{x \in A} \min_{c \in C} \|B(x) - c\|^2 \right|.$$

$$\leq \sum_{x \in A} \min_{c \in C} \|x - c\|^2 - \sum_{x \in A} \min_{c \in C} \|B(x) - c\|^2.$$
$$\sum_{x \in A} \min_{c \in C} \|x - c\|^2 \leq \frac{12}{\varepsilon} \sum_{x \in A} \|x - B(x)\|^2 \leq \frac{12}{\varepsilon} \cdot \varepsilon^2 \sum_{x \in A} \|x - c_j\|^2 \leq 12\varepsilon \cdot 6.357 \sum_{x \in A} \|x - c_j\|^2.$$

$$\leq \sum_{x \in A} \min_{c \in C} \left( \frac{12}{\varepsilon} \|x - B(x)\|^2 + 2\varepsilon \|x - c\|^2 \right) \quad (\text{triangle ineq.})$$

$$\leq O(\varepsilon) \cdot D(A, C). \quad (\varepsilon\text{-ball-cover, lemma, k-means})$$

(Rmk: 除了“ $\varepsilon$ -ball-cover”，还常用“Importance Sampling”构造coreset.)

即.  $\forall a \in A$  def importance (e.g. by the prop. to be center).

sampling by the importance. )

日期:

## VI. Coreset in Streaming:

Lemma (Composability (组合性)).

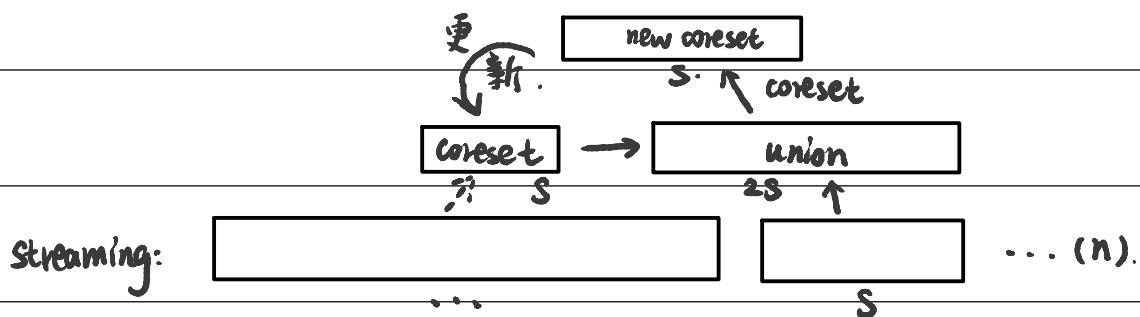
$A_1, A_2 \subseteq \mathbb{R}^d$  不相交. 设  $(S_1, w_1)$  是  $A_1$ 's coreset.

$(S_2, w_2)$  是  $A_2$ 's coreset.

则  $(S_1 \cup S_2, w_1 + w_2)$  是  $A_1 \cup A_2$ 's coreset.

Merge & Reduce.

每流进  $s \uparrow$  data, 求 coreset, 与之前的 coreset 合并 (按对的方式)



(Rmk: need to store  $\leq 2$  coresets each level 同时.)

coreset size =  $S = O(k \log n \varepsilon^{-d}) = O(k \log n)^{d+1} \varepsilon^{-d}$ .

$$\varepsilon' = \frac{\varepsilon}{\log n} \quad (\text{因为误差产生了累加}).$$

## VII. Hierarchical Clustering:

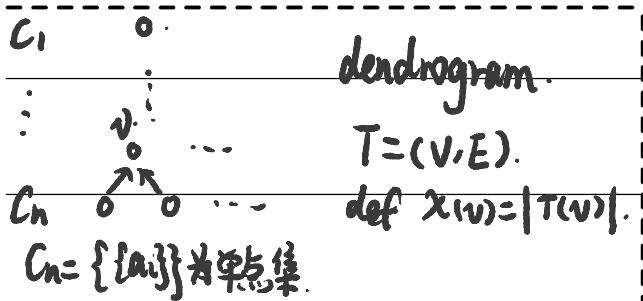
a sequence of clustering:  $C_1, \dots, C_n$ . with  $|C_k| = c_k$ .

is called hierarchical clustering of  $A$  if  $\forall H \in C_k$ :

日期:

1°  $H \in C_{k+1}$

or 2°  $\exists B, C \in C_{k+1}, H = B \cup C, C_k = (C_{k+1} \setminus \{B\} \cup \{C\}) \cup \{H\}$ .



Construct 1: 自顶向下.

Main Idea: 每层对其中一类做二分类, 直到划分为单点集.

Construct 2: 自底向上.

Main Idea:

①  $C_n = \{C^1, \dots, C^n\}, C^i = \{a_i\}$ .

② each step, replace  $C^i, C^j$  (which are "closest") by  $C^i \cup C^j$ .

until get  $C_1$ .

③ return  $C_1, \dots, C_n$ .

How to describe "closest"?

1. Agg CL: (基于 complete linkage).

用  $D_{CL}(C^i, C^j) = \max_{\substack{x \in C^i \\ y \in C^j}} D(x, y)$  衡量 "closest".

日期:

Def: 直径:  $\text{diam}^D(S) = \max_{x,y \in S} D(x,y)$

CL 在优化

直径代价  $\text{cost}_{\text{diam}}^D(C_k) = \max_{1 \leq i \leq k} \text{diam}^D(C_i)$ .

$\Rightarrow$

直径代价

Problem:  $\text{opt}_k^{\text{diam}}(A) = \min_{|C|=k} \text{cost}_{\text{diam}}^D(C)$ .

Thm: Agg CL computes a  $k$ -clustering  $C_k$  with:

$\text{cost}_{\text{diam}}^D(C_k) \leq O_d(\log k) \text{opt}_k^{\text{diam}}(A)$ .

$\hookrightarrow$  n无关.

for  $\forall k \leq |A|$ .

Rmk: ①  $O_d(\cdot)$  表示 hide  $\exists d$  的依赖. 即  $O_d(\cdot) = g(d)\cdots$ , 一般表示可将  $d$  视作常数.

② 对  $\forall k$  可给出近似解.

2. Agg SL: (基于 single linkage).

用  $D_{SL}(C^1, C^2) = \min_{\substack{x \in C^1 \\ y \in C^2}} D(x,y)$  衡量 "closest".

Rmk: SL 在优化 Problem: maximize:  $\min_{\substack{x \in C_i \\ y \in C_j \\ i \neq j}} D(x,y)$ .

3. Agg ML: (基于 mean linkage).

用  $D_{SL}(C^1, C^2) = \text{mean}_{\substack{x \in C^1 \\ y \in C^2}} D(x,y)$  衡量 "closest".

Construct 3: (Gonzalez's alg-based)

日期:

full-furthest-first-traversal:

① choose  $c_1$  arbitrarily.

② For  $i=2, \dots, |A|$ , do:

choose  $c_i$  maximize  $R_i = \max_{x \in A} D(x, c^{i-1})$ .

$C^i \leftarrow C^{i-1} \cup \{c_i\}$ .

③ maintain  $c_1, \dots, c_{|A|}, R_2, \dots, R_{|A|}$ . (rank:  $\{c_1, \dots, c_{|A|}\} = A$ ).

④ define  $L_0 = \{c_1\}$ .  $R = R_2$  (largest)  $L_i = \{c_j \mid R_j \in (\frac{R}{2^i}, \frac{R}{2^{i-1}}]\}$

define  $\text{parent}(c_i) = \arg \min_{y \in \bigcup_{j>0} L_j} D(c_i, y)$ .

⑤  $i$  from  $n$  to 1: merge  $c_i$  with  $\text{parent}(c_i)$ . as  $\text{parent}(c_i)$ .