

Dimensionality Reduction & Subspace Embeddings

Monday, March 2, 2020 1:11 PM

Goal:

$$\text{Linear regression} \quad \min_x \|Ax - b\|_2$$

$$\nabla f = 0 \iff A^T A x - A^T b = 0.$$

$$\text{So } x = (A^T A)^{-1} A^T b$$

Computational complexity:

$$O(nd^2 + d^{\omega}).$$



Q. Can we do it faster?

Iterative methods.

Richardson Iteration:

$$x^{(k+1)} = x^{(k)} - (A^T A x^{(k)} - A^T b) = (I - A^T A)^{(k)} x + A^T b$$

$$\text{Im. } k = \frac{\lambda_{\max}(A^T A)}{\lambda_{\min}(A^T A)} \quad \text{with} \quad A^T A \prec I$$

$$\dots \rightarrow \|x^{(k)} - x^*\|$$

$$\text{Then } \|x^{(k+1)} - x^*\|_2 \leq \left(1 - \frac{1}{k}\right) \|x^{(k)} - x^*\|_2.$$

$$\begin{aligned} & \tilde{A}^T b + (I - \tilde{A}^T \tilde{A}) A^T b + (\quad)^2 A^T b + \dots \\ \rightarrow & (\tilde{A}^T \tilde{A})^{-1} = \frac{1}{(I - (I - \tilde{A}^T \tilde{A}))} = (I + (I - \tilde{A}^T \tilde{A}) + \dots) \end{aligned}$$

We will prove a more general theorem using a more general algorithm.

(above could be very slow if k is large.)

Now Suppose we know M s.t.

$$A^T A \succcurlyeq M \succcurlyeq k \cdot A^T A$$

$$\text{let } x^{(k+1)} = x^{(k)} - M^{-1} (A^T A x^{(k)} - A^T b)$$

$$\text{Then } \|x^{(k+1)} - x^*\|_M \leq \left(1 - \frac{1}{k}\right) \|x^{(k)} - x^*\|_M.$$

Note $\|y\|_M^2 = y^T M y$.

Pf. $x^{(k+1)} - x^* = x^k - x^* - M^{-1}(A^T A x^{(k)} - A^T A x^*)$
 $= (I - M^{-1} A^T A)(x^{(k)} - x^*)$

$$\begin{aligned}\|x^{(k+1)} - x^*\|_M^2 &= (x^{(k)} - x^*)^T (I - A^T A M^{-1}) M (I - M^{-1} A^T A) (x^{(k)} - x^*) \\ &= (x^{(k)} - x^*)^T M^{\frac{1}{2}} (I - M^{\frac{1}{2}} A^T A M^{\frac{1}{2}}) (I - M^{\frac{1}{2}} A^T A M^{\frac{1}{2}})^T M^{\frac{1}{2}} (x^{(k)} - x^*) \\ &= (x^{(k)} - x^*)^T M^{\frac{1}{2}} (I - H)^2 M^{\frac{1}{2}} (x^{(k)} - x^*)\end{aligned}$$

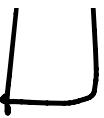
$$\frac{1}{k} I \preccurlyeq H \preccurlyeq I \Rightarrow I - H \preccurlyeq \left(1 - \frac{1}{k}\right) \cdot I$$

$$\leq \left(1 - \frac{1}{k}\right)^2 \|x^{(k)} - x^*\|_M^2.$$

What M to choose?

Goal is to approximate $A^T A$.
st. $\|Ax\|^2 \approx \|Bx\|^2$, $M = B^T B$.

At. $\|Ax\| \approx \|Dx\|$, ...



There is a perfect M .

$$A = U\Sigma V^T \quad \forall y \in \{Ax\}$$

$$\|U^Ty\|^2 = \|U^TV\Sigma V^Tx\|^2 = x^TA^TAx = \|Ax\|^2$$

But finding this V needs SVD. Typically more expensive.

How about random Π ? ΠA

$$M = (\Pi A)^T (\Pi A)$$

Is this any good? What size of Π ?

Low-distortion embedding

Π is a ε -low-dist. emb. for a set S of vectors in \mathbb{R}^n

if $\forall y \in S$

$$(1-\varepsilon) \|y\|^2 \leq \|\Pi y\|^2 \leq (1+\varepsilon) \|y\|^2$$

dimension of Π = # rows?

Dimension of $U = \# rows$.

S for w_0 is the subspace $\{Ax\}$.

We know $\Pi = U$ is perfect ($\epsilon = 0$).

Oblivious Subspace Embedding.

Random matrix Π is a (d, ϵ, δ) -OSE for a fixed d -dim subspace S if it preserves $\|y\|^2$ to within $(1 \pm \epsilon)$ w.p. δ with prob. at least $1 - \delta$.

Alternatively. If $U \in \mathbb{R}^{n \times d}$

$$\Pr\left(\|U^T \Pi^T \Pi U - I_{d \times d}\|_{op} \geq \epsilon\right) \leq \delta$$

Pf. $S = \{Uz\}$

$$(1-\epsilon)\|y\|^2 \leq \|\Pi y\|^2 \leq (1+\epsilon)\|y\|^2$$

$$\Leftrightarrow (1-\epsilon)U^T U \preceq U^T \Pi^T \Pi U \preceq (1+\epsilon)U^T U$$

$$U^T U = I$$

$$\Leftrightarrow \|U^T \Pi^T \Pi U - I\|_{op} \leq \epsilon.$$

$$\Leftrightarrow \|\mathbf{U}^T \boldsymbol{\Pi} \mathbf{U}^{-1}\|_{\text{op}} = \varepsilon.$$

Thm [Johnson-Lindenstrauss] $\Pi_{ij} \sim N(0, \frac{1}{\sqrt{m}})$ with
 $m = O\left(\frac{1}{\varepsilon^2} \log \frac{1}{\delta}\right)$ rows is a $(1, \varepsilon, \delta)$ -OSE.
 i.e. $\Pr(|\|\boldsymbol{\Pi}x\|^2 - 1| \geq \varepsilon) \leq \delta.$

Thm A $(1, \varepsilon, \delta)$ -OSE is a $(d, 4\varepsilon, 5^d \delta)$ -OSE.

(so it suffices to handle $d=1$).

Lemma [ε -net]. $\exists N \subseteq S^{n-1}$ st. $\forall x \in B^n$, $\exists x_i \in N$
 st. $\|x - x_i\| \leq \varepsilon$ and $|N| \leq \left(1 + \frac{2}{\varepsilon}\right)^n$.

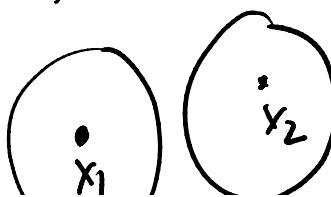
Pf. start with any $x \in S^{n-1}$.

[while $\exists x$ st. $\nexists x_i \in N$ $\|x - x_i\| > \varepsilon$

add x to N

At the end $\forall x \in S^n$, $\exists x_i \in N$ st. $\|x - x_i\| \leq \varepsilon$.

$x_i + \frac{\varepsilon}{2} B_n$ are disjoint



$$\dots \subset \subset (1 + \varepsilon) B_n$$

$$\left(\begin{array}{c} \bullet \\ x_1 \end{array} \right) \cup \left(\begin{array}{c} \bullet \\ x_2 \end{array} \right) \bigcup_i x_i + \frac{\varepsilon}{2} B_n \subseteq \left(1 + \frac{\varepsilon}{2}\right) B_n$$

$$\therefore |N| \leq \frac{\text{Vol}((1 + \frac{\varepsilon}{2}) B_n)}{\text{Vol}(\frac{\varepsilon}{2} B_n)} = \left(\frac{1 + \frac{\varepsilon}{2}}{\frac{\varepsilon}{2}}\right)^n = \left(1 + \frac{2}{\varepsilon}\right)^n.$$

Lemma 2. $\forall x \in B_n, \exists t_1, \dots, t_i \quad t_i \leq \frac{1}{2^i}$ s.t.

$$x = \sum_i t_i x_i \quad x_i \in N.$$

Pf. Take N with $\varepsilon = \frac{1}{2}$.

$\forall x, \exists x_1$ s.t. $\|x - x_1\| \leq \frac{1}{2}$

$\|x\| = 1 \quad \therefore \exists x_2, t_2, t_2 \leq \frac{1}{2} \quad \text{s.t.} \quad \|x - x_1 - t_2 x_2\| \leq \frac{1}{4}$

(applied to $\frac{1}{2} B^n$) continue to get conclusion.

Pf. (of Thm OSE):

$$x^\top (U^\top \Pi^\top \Pi U - I) x = \sum_{i,j} t_i t_j x_i (U^\top \Pi^\top \Pi U - I) x_j$$

$$\leq \sum_{i,j} t_i t_j \max_{x_i, x_j} x_i (U^\top \Pi^\top \Pi U - I) x_j$$

$$= 1. \max x^\top (U^\top \Pi^\top \Pi U - I) x$$

$$\begin{aligned}
 &\leq 4 \cdot \max_{x \in \mathcal{N}} x^T (U^T \Pi^T \Pi U - I) x \\
 &= 4 \max_{x \in \mathcal{U}\mathcal{N}} x^T |\Pi^T \Pi - I| x \\
 &= 4 \max_{x \in \mathcal{U}\mathcal{N}} |\|\Pi x\|^2 - 1|
 \end{aligned}$$

Since Π is an $(1, \varepsilon, \delta)$ -OSE

$$\Pr(|\|\Pi x\|^2 - 1| \geq \varepsilon) \leq \delta. \text{ for any single } x$$

And for all the $|N| \leq 5^d$ x^s in $\mathcal{U}\mathcal{N}$,

$$\Pr(\forall x \in \mathcal{U}\mathcal{N} | \|\Pi x\|^2 - 1 | \geq \varepsilon) \leq 5^d \cdot \delta.$$

$\therefore \Pi$ is a $(d, 4\varepsilon, 5^d \delta)$ -OSE.

\therefore Using $m = O\left(\frac{1}{\varepsilon^2}(d + \log \frac{1}{\delta})\right)$ rows

suffices to get a (d, ε, δ) -OSE.

When $\Pi_{ij} \sim N(0, \frac{1}{m})$ or $\Pi_{ij} = \pm \frac{1}{\sqrt{m}}$.

$\varepsilon = \Theta(1)$ suffices for linear regression

2. t. r. r. $\Pi \Delta$ takes $O(nd^2)$

but. computing $\pi\pi^T$ takes $O(nd^2)$

So no saving on $\pi\pi^T$.

How about a sparse random matrix?

$\pi_{ij} = \pm \frac{1}{\sqrt{d}}$ w.p. $\frac{1}{m}$ and 0 o.w. π has m rows.

Thm. π as above is a (d, ε, δ) -OSE for

$$S = O\left(\frac{1}{\varepsilon^2} \log^2 \frac{d}{\delta}\right) \text{ and } m = O\left(\frac{d \log \frac{d}{\delta}}{\varepsilon^2}\right).$$

$$U^\top \pi \pi^\top \pi U = \sum_{r=1}^m (\pi U)_r^T (\pi U)_r$$

We will use 2 more lemmas to analyze this sum.

Thm 1 (Matrix Chernoff) $M_1, M_2, \dots, M_T \in \mathbb{R}^{n \times n}$, $M_i \neq 0$, $\|M_i\| = 1$
 $M_i \not\approx R \cdot I$

$$(1 - O(\varepsilon))I \leq \frac{1}{T} \sum M_j \leq (1 + O(\varepsilon))I$$

$$\text{where } T \geq \frac{R}{\varepsilon^2} \log \frac{n}{\delta}.$$

Thm 2. (Hausman-Wright) σ_i iid $E\sigma_i = 0$ $|\sigma_i| \leq 1$.

Then $|\sigma^T A \sigma - E\sigma^T A \sigma| \leq C \cdot (||A||_F \sqrt{\log \frac{1}{\delta}} + ||A||_{op} \log \frac{1}{\delta})$

w.p. 1- δ .

Lemma. $M_r = m \ U^T \pi_r \pi_r^T U$

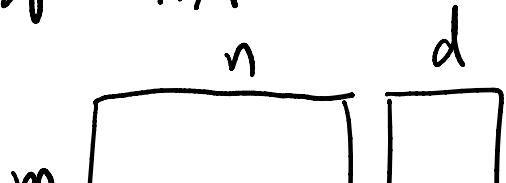
$$M_r \succcurlyeq 0 \quad E M_r = I$$

$$M_r \lesssim m \cdot \pi_r^T U U^T \pi_r \cdot I$$

For $A \geq \frac{m}{n} \log \frac{1}{\delta}, \frac{1}{\varepsilon^2} \log^2 \frac{1}{\delta}$, $m \geq \frac{d}{\varepsilon^2} \log \frac{1}{\delta}$.

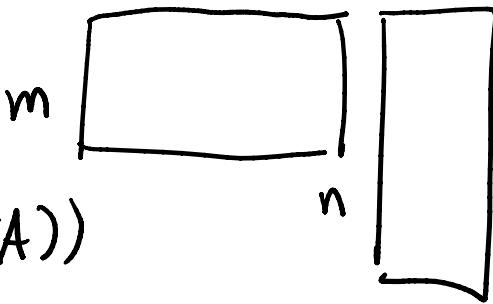
$$\pi_r^T U U^T \pi_r \leq \frac{\varepsilon^2}{\log \frac{n}{\delta}} \quad \text{w.p. } 1-\delta.$$

Now the running time for πA is



$$m \cdot n \cdot \frac{p}{m}$$

$\mathcal{O}(nnz(A))$



to get an $\tilde{\mathcal{O}}(d) \times d$ matrix.

$$\tilde{\mathcal{O}}(nnz(A) + d^w)$$
