# Dimensionality Reduction & Subspace Embeddings

Goal:

Linear regression $\quad \min_x \|Ax - b\|_2$

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

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

Computational complexity:

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

(with diagram labeled $d$, $n$)

Q. Can we do it faster?

---

Iterative methods.

Richardson Iteration:

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

Thm. $\kappa = \dfrac{\lambda_{max}(A^T A)}{\lambda_{min}(A^T A)} \quad \text{with} \quad A^T A \preceq I$

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

---

$$A^T b + (I - A^T A)A^T b + \left( \qquad \right)^2 A^T b + \cdots$$

$$\longrightarrow (A^T A)^{-1} = \underline{\quad 1 \quad} = \left(I + (I - A^T A) + \cdots \right)_{A^T}$$

$$\rightarrow \quad (A^T A)^{-1} = \frac{1}{(I - (I - A^T A))} = \left( I + (I - A^T A) + \cdots \right) A^T b.$$

---

We will prove a more general theorem using a more general algorithm.
(above could be very slow if $k$ is large.)

$\boxed{Thm}$
Suppose we know $M$ s.t.

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

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

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:
$$x^{(k+1)} - x^* = x^k - x^* - M^{-1}\left( A^T A x^{(k)} - A^T A x^* \right)$$

$$= \left( I - M^{-1} A^T A \right)\left( x^{(k)} - x^* \right)$$

$$\| x^{(k+1)} - x^* \|_M^2 = \left( x^{(k)} - x^* \right)^T \left( I - A^T A M^{-1} \right) M \left( I - M^{-1} A^T A \right)\left( x^{(k)} - x^* \right)$$

$$\|x^{(k+1)} - x^*\|_M^2 = (x^{(k)} - x^*)^T (I - A^T A M^{-1}) M (I - M^{-1} A^T A)(x^{(k)} - x^*)$$

$$= (x^{(k)} - x^*)^T M^{\frac{1}{2}} (I - M^{-\frac{1}{2}} A^T A M^{-\frac{1}{2}})(I - M^{-\frac{1}{2}} A^T A M^{-\frac{1}{2}}) M^{\frac{1}{2}} (x^{(k)} - x^*)$$

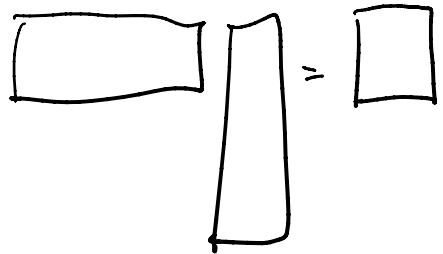$$= (x^{(k)} - x^*)^T M^{\frac{1}{2}} (I - H)^2 M^{\frac{1}{2}} (x^{(k)} - x^*)$$

$$\frac{1}{k} I \preceq H \preceq I \qquad \Rightarrow \qquad I - H \preceq (1 - \frac{1}{k}) \cdot I$$

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

---

what $M$ to choose?

Goal is to approximate $A^T A$.

s.t. $\|Ax\|^2 \approx \|Bx\|^2$, $M = B^T B$..



There is a perfect $M$.

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

$$\|U^T y\|^2 = \|U^T U \Sigma V^T x\|^2 = x^T A^T A x = \|Ax\|^2$$

But finding this $U$ needs SVD. Typically more expensive.

---

How about **random** $\Pi$? $\qquad \Pi A$

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

Is this any good? What size of $\Pi$?

Is this any good? what size of $\Pi$ ?

## Low-distortion embedding

$\Pi$ is a $\varepsilon$ 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 $\Pi$ = # rows.

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

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

## Oblivious Subspace Embedding.

Random matrix $\Pi$ is a $(d, \varepsilon, \delta)$-OSE for a fixed $d$-dim subspace $S$ if it preserves $\|\|^2$ to within $(1\pm\varepsilon)$ $\forall y \in S$ with prob. at least $1-\delta$.

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

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

Pf. $S = \{Uz\}$
$$(1-\varepsilon)\|y\|^2 \leq \|\Pi y\|^2 \leq (1+\varepsilon)\|y\|^2$$

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

$$\Longleftrightarrow \quad (1-\varepsilon)U^T U \preceq U^T \Pi^T \Pi U \preceq (1+\varepsilon)U^T U$$

$$U^T U = I$$

$$\Longleftrightarrow \quad \| U^T \Pi^T \Pi U - I \|_{op} \le \varepsilon.$$
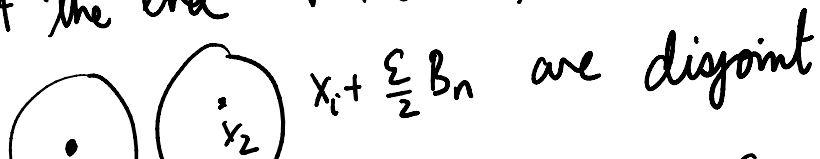
---

**Thm** [Johnson-Lindenstrauss] $\Pi_{ij} \sim N(0, \frac{1}{\sqrt{m}})$ with $m = O(\frac{1}{\varepsilon^2} \log \frac{1}{\delta})$ rows is a $(1, \varepsilon, \delta)$-OSE.

i.e. $\quad \Pr\left( \left| \|\Pi x\|^2 - 1 \right| \ge \varepsilon \right) \le \delta.$
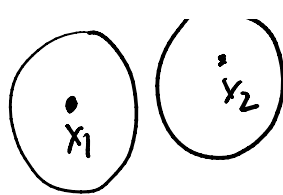
---

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

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

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

**Pf.**

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

$\left[\begin{array}{l} \text{while } \exists x \text{ s.t. } \forall x_i \in N \ \|x - x_i\| > \varepsilon \\ \text{add } x \text{ to } N \end{array}\right.$

At the end $\forall x \in S^{n-1}$, $\exists x_i \in N$ s.t. $\|x - x_i\| \le \varepsilon$.

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

 $X_i + \frac{2}{2} B_n$ are disjoint

$$\bigcup_i X_i + \frac{\varepsilon}{2} B_n \subseteq \left(1 + \frac{\varepsilon}{2}\right) B_n$$

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

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

$$X = \sum_i t_i X_i \qquad X_i \in N.$$

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

$\forall X, \quad \exists X_1, \text{ s.t. } \|X - X_1\| \leq \frac{1}{2}$
$\|X\| = 1$

$\therefore \exists X_2, t_2, \quad t_2 \leq \frac{1}{2}$ s.t. $\|X - X_1 - t X_2\| \leq \frac{1}{4}$

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


**Pf. (of Thm OSE):**

$$X^T (U^T \Pi^T \Pi U - I) x = \sum_{i,j} t_i t_j X_i (U^T \Pi^T \Pi U - I) X_j$$

$$\leq \sum_{i,j} t_i t_j \max_{X_i, X_j} X_i (U^T \Pi^T \Pi U - I) X_j$$

$$\leq 4 \cdot \max_{X \in N} X^T (U^T \Pi^T \Pi U - I) x$$

$$= 4 \max_{X \in UN} X^T | \Pi^T \Pi - I | x$$

$$= 4 \max_{x \in UN} \left| \| \Pi x \|^2 - 1 \right|$$

Since $\Pi$ is an $(1, \epsilon, \delta)$-OSE

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

And for all the $|N| \leq 5^d$ $x^s$ in $UN$,

$$\Pr\left( \forall x \in UN \left| \| \Pi x \|^2 - 1 \right| \geq \epsilon \right) \leq 5^d \cdot \delta.$$

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

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

suffices to get a $(d, \epsilon, \delta)$-OSE.

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

$\epsilon = \theta(1)$ suffices for linear regression.

But computing $\Pi A$ takes $O(nd^2)$

So no saving on $A^T A$.

How about a _sparse_ random matrix?

$\Pi_{ij} = \pm \frac{1}{\sqrt{s}}$ w.p. $\frac{s}{m}$ and $0$ o.w. $\Pi$ has $m$ rows.

**Thm.** $\Pi$ as above is a $(d, \epsilon, \delta)$-OSE for

__Thm__. $\Pi$ as above is a $(d, \varepsilon, \delta)$-OSE $\quad$ if

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

---

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

We will use 2 thms + 1 Lemma to analyze this sum.

__Thm 1__ (Matrix Chernoff) $M_1, M_2, \ldots \in \mathbb{R}^{n \times n}$, $M_i \succeq 0$, $\mathbb{E} M_i = I$
$\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad M_i \preceq R \cdot I$

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

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

---

__Thm 2.__ (Hanson-Wright) $\sigma_i$ iid $\quad \mathbb{E} \sigma_i = 0 \quad |\sigma_i| \leq 1$.

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

w.p. $1 - \delta$.

---

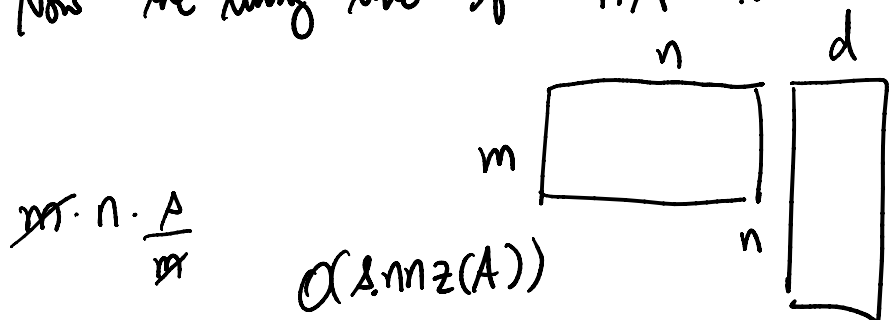__Lemma.__ $\qquad M_r = m \, U^T \Pi_r \Pi_r^T U$

$\qquad M_r \succeq 0 \qquad \mathbb{E} M_r = I$

$\qquad\qquad M_r \preceq m \cdot \Pi_r^T U U^T \Pi_r \cdot I$

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

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

---

Now the running time for $\pi A$ is



$\cancel{m} \cdot n \cdot \frac{\Delta}{\cancel{m}}$     $O(\Delta \cdot nnz(A))$

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

$\tilde{O}\left(nnz(A) + d^w\right)$ .

---

Pf of Lemma.    $\pi_r$ has $\frac{\Delta n}{m}$ non zeros in expectation

and at most $\frac{2\Delta n}{m}$ non zeros

(CHERNOFF: $P_r\left(\sum X_i > \mathbb{E}(\sum X_i)(1+\varepsilon)\right) \leq e^{-c\varepsilon^2 \mathbb{E}(X_i)}$ ) .

$\quad |X_i| \leq 1$

with prob $1-\delta$ assuming $\mathbb{E}(X) = \frac{\Delta n}{m} > c \cdot \log\frac{1}{\delta}$ .

let $\sigma = \pi_r|_I$     $I = \{i : (\pi_r)_i \neq 0\}$. $|I| \leq \frac{2\Delta n}{m}$ .

$\qquad \qquad \qquad D = (UU^T)$

$$\sigma_i = \pm \frac{1}{\sqrt{\Delta}} \qquad\qquad P = (UU^T)_{I \times I}$$

$$\Pi_r^T U U^T \Pi_r = \sigma^T P \sigma \qquad\qquad P \succeq 0.$$

By the H-W inequality

$$\left| \sigma^T P \sigma - \mathbb{E}(\sigma^T P \sigma) \right| \le \frac{c}{\Delta}\left( \|P\|_F \sqrt{\log \frac{1}{\delta}} + \|P\| \log \frac{1}{\delta} \right)$$

$$\|P\|_{op} \le \|UU^T\|_{op} \le 1.$$

$$\|P\|_F \le \sqrt{\text{tr } P} \qquad\qquad \begin{array}{l} P \text{ is a } 2\frac{\Delta n}{m} \text{ diagonal block of } UU^T \\[4pt] \text{tr}(UU^T) = d \\[4pt] \mathbb{E}(\text{tr}(P)) = \dfrac{2\frac{\Delta n}{m}}{n} \cdot d = \dfrac{2\,\Delta d}{m}. \end{array}$$

W.p. $1-\delta$ $\quad \text{tr } P \le 4\dfrac{\Delta d}{m}$ (Chernoff again)

Also $\quad \mathbb{E}(\sigma^T P \sigma) = \dfrac{1}{\Delta} \text{tr}(P)$

So $\quad \Pi_r^T U U^T \Pi_r \le \dfrac{c}{\Delta}\left( \text{tr}(P) + \sqrt{\text{tr } P}\sqrt{\log \frac{1}{\delta}} + \log \frac{1}{\delta} \right)$

$$\le c \cdot \left( \frac{d}{m} + \sqrt{\frac{d}{m\Delta} \log \frac{1}{\delta}} + \log \frac{1}{\delta} \right).$$

$$m \ge c \cdot \frac{d \log \frac{d}{\delta}}{\varepsilon^2} \quad , \quad \Delta \ge c \cdot \frac{\log^2(\frac{d}{\delta})}{\varepsilon^2} \qquad\qquad \searrow \le \frac{\varepsilon^2}{\log \frac{d}{\delta}}.$$

_____

To prove the theorem, we can now apply matrix-chernoff.

Another proof (classical) of JL $(1, \varepsilon, \delta)$-OSE.

**Lema.** $\Pi_{ij} \sim N\left(0, \frac{1}{m}\right)$. Then for any fixed $x \in \mathbb{R}^n$

$$P_r\left(\left|\|\Pi x\|^2 - \|x\|^2\right| > \varepsilon \|x\|^2\right) \leq 2e^{-m \cdot \frac{(\varepsilon^2 - \varepsilon^3)}{4}}.$$

**Pf.** $\|\Pi x\|^2 = \sum_{r=1}^{m} (\Pi_r^T x)^2$ $\qquad Y_r \sim \Pi_r^T x \sim N(0, \sigma^2)$

$$\sigma^2 = \sum_{i=1}^{n} x_i^2 \cdot \frac{1}{m} = \frac{\|x\|^2}{m}.$$

$Y = \sum_{r=1}^{m} Y_r^2$

$\mathbb{E}(Y) = \|x\|^2$ $\qquad$ Y has a chi-squared distribution.

$$P_r(Y > t\mathbb{E}(Y)) = P_r\left(e^{\alpha Y} > e^{\alpha t \mathbb{E}(Y)}\right) \leq \frac{\mathbb{E}(e^{\alpha Y})}{e^{\alpha t \mathbb{E}(Y)}}.$$

$$\mathbb{E}(e^{\alpha Y}) = \prod_{r=1}^{m} \mathbb{E}(e^{\alpha Y_r})$$

$$\mathbb{E}(e^{\alpha Y_r}) = \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} e^{\alpha y^2} \cdot e^{-\frac{y^2}{2\sigma^2}} dy$$

$$= \frac{1}{\sqrt{2\pi}\sigma} \int_{-\infty}^{\infty} e^{-\frac{y^2}{2}\left(\frac{1}{\sigma^2} - 2\alpha\right)} dy$$

$(2\alpha\sigma^2 < 1)$ $\qquad = \frac{1}{\sqrt{(1-2\alpha\sigma^2)}} \int_{-\infty}^{\infty} e^{-\frac{y^2}{2\sigma^2/1-2\alpha\sigma^2}} dy$

$$(2\alpha\sigma^2 < 1) \qquad = \frac{1}{\sqrt{1-2\alpha\sigma^2}} \frac{\sqrt{(1-2\alpha\sigma^2)}}{\sqrt{2\pi}\,\sigma} \int_{-\infty}^{\infty} e^{-\frac{y^2/2\sigma^2}{1-2\alpha\sigma^2}} \cdot dy.$$

$$= \frac{1}{\sqrt{1-2\alpha\sigma^2}} \cdot$$

$$\therefore P_\lambda\left(Y > t\,\mathbb{E}(Y)\right) \le \frac{e^{-\alpha t}}{(1-2\alpha\sigma^2)^{\frac{m}{2}}} = \frac{e^{-\alpha(1+\varepsilon)}}{(1-2\frac{\alpha}{m})^m}$$

$$= \left( \frac{e^{-2\alpha(1+\varepsilon)}}{(1-2\alpha)} \right)^{\frac{m}{2}}$$

we can choose $\alpha$.

$$-2(1+\varepsilon)\frac{1}{(1-2\alpha)} + \frac{2}{(1-2\alpha)^2} = 0$$

$$(1-2\alpha) = \frac{1}{1+\varepsilon} \implies \alpha = \frac{\varepsilon}{2(1+\varepsilon)}$$

$$= \left( e^{-\varepsilon}(1+\varepsilon) \right)^m$$

$$\le e^{-\frac{(\varepsilon^2 - \varepsilon^3)}{4} \cdot m}.$$

Other direction is the same, except

$$P_\lambda\left(Y < t\,\mathbb{E}(Y)\right) = P_\lambda\left( e^{-\alpha Y} > e^{-\alpha t \mathbb{E}(Y)} \right) \text{ for } \alpha \ge 0.$$

---

A.t. related to subspace embedding : Now

Another approach to subspace embedding: row sampling.