

UL = RANDOM PROJECTIONS

Random projections provide an alternative to PCA for dimension reduction.

The problem: given a learning sample $\mathcal{D}_n = \{x_1, \dots, x_n\}$ of size n , $x_i \in \mathbb{R}^D$, D supposed very large, project these points into a lower dimensional space, of dimension $d \ll D$, preserving some information.

↳ we are interested in preserving the geometry of x_1, \dots, x_n .

More precisely, our goal is to construct a mapping $T: \mathbb{R}^D \rightarrow \mathbb{R}^d$ (possibly random, and it will be)

such that

$$\mathbb{P}\left(\forall i \neq j : 1 - \varepsilon \leq \frac{\|T(x_i) - T(x_j)\|_2}{\|x_i - x_j\|_2} \leq 1 + \varepsilon\right) \geq 1 - \delta,$$

$\varepsilon > 0$ small
 $\delta \in (0, 1)$

With high probability, points that are close to each other stay close

different philosophy than PCA, which is projecting the data in directions of maximum variance.

We will see shortly that this is possible with a simple random linear transformation T . We recall first some background results about sub-Gaussian & sub-Gamma distributions.

I - PRELIMINARIES

2

I. 1. Sub-Gaussian distributions.

- Let $X \sim \mathcal{N}(\mu, \sigma^2)$. The Moment Generating Function (MGF) of X is given by $M_X(\theta) := \mathbb{E}\{e^{\theta X}\}$, with
 $(*) \quad \log[\mathbb{E}\{e^{\theta(X-\mu)}\}] = \frac{\theta^2 \sigma^2}{2}, \quad \forall \theta \in \mathbb{R}$.

Useful to obtain bounds on tail probabilities of a distribution:

CHERNOFF's BOUND: For a real-value RV X , $\forall t \in \mathbb{R}$, put $\Lambda_X(t) := \sup_{\theta \geq 0} (\theta t - \log M_X(\theta))$.

Then $\mathbb{P}(X > t) \leq \exp(-\Lambda_X(t))$

$$\begin{aligned} \text{Indeed, } \mathbb{P}(X > t) &= \mathbb{P}(e^{\theta X} > e^{\theta t}) \\ &\leq e^{-\theta t} \mathbb{E}(e^{\theta X}) \quad \text{Markov} \\ &= e^{-\theta t} M_X(\theta) \end{aligned}$$

For $X \sim \mathcal{N}(\mu, \sigma^2)$, we immediately get

$$\left. \begin{aligned} \mathbb{P}(X - \mu > t) \\ \mathbb{P}(X - \mu < -t) \end{aligned} \right\} \leq \exp\left(-\frac{t^2}{2\sigma^2}\right)$$

We introduce a class of distributions that have lighter tails than the Gaussian distribution:

A RV X is said to be SUB-GAUSSIAN if $\exists s^2 > 0$ s.t. $\forall \theta \in \mathbb{R}$, $\log[\mathbb{E}\{e^{\theta(X-\mathbb{E}X)}\}] \leq \frac{\theta^2 s^2}{2}$.

Whenever this holds, we write $X \in SG(s^2)$. (3)

The smallest $s^2 > 0$ for which $X \in SG(s^2)$ is called the VARIANCE PROXY of X , and is denoted $\|X\|_{vp}^2$.

It is given by

$$\|X\|_{vp}^2 = \sup_{\alpha \neq 0} \frac{2}{\alpha^2} \log \mathbb{E}\{e^{\alpha(X - \mathbb{E}X)}\}.$$

As the notation indicates, the variance proxy behaves almost like a norm: for two sub-Gaussian RVs X, Y , $\forall \alpha \in \mathbb{R}$,

- (i). $\|\alpha X\|_{vp} = |\alpha| \|X\|_{vp}$

$$(ii). \|X + Y\|_{vp} \leq \|X\|_{vp} + \|Y\|_{vp}$$

$$(iii). \|X\|_{vp} = 0 \Leftrightarrow X = c \text{ a.s.}$$

(i) follows from the definition of $\|X\|_{vp}^2$

(ii) $\forall 1 < p, q < +\infty$ s.t. $p^{-1} + q^{-1} = 1$, X, Y centered,

$$\begin{aligned} \mathbb{E}\{e^{\alpha(X+Y)}\} &= \mathbb{E}\{e^{\alpha X} e^{\alpha Y}\} \xrightarrow{\text{Hölder}} \\ &\leq (\mathbb{E}\{e^{p\alpha X}\})^{1/p} (\mathbb{E}\{e^{q\alpha Y}\})^{1/q} \\ &\leq \left(\exp\left\{\frac{p^2\alpha^2\|X\|_{vp}^2}{2}\right\}\right)^{1/p} \left(\exp\left\{\frac{q^2\alpha^2\|Y\|_{vp}^2}{2}\right\}\right)^{1/q} \\ &= \exp\left\{\frac{\alpha^2}{2}(p\|X\|_{vp}^2 + q\|Y\|_{vp}^2)\right\} \\ &= \exp\left\{\frac{\alpha^2}{2}(p\|X\|_{vp}^2 + \frac{p}{p-1}\|Y\|_{vp}^2)\right\} \end{aligned}$$

Holds for any $p > 1$.

$$\inf_{p>1} \left\{ p\|X\|_{vp}^2 + \frac{p}{p-1}\|Y\|_{vp}^2 \right\} = (\|X\|_{vp} + \|Y\|_{vp})^2 \text{ achieved for } p = 1 + \frac{\|Y\|_{vp}}{\|X\|_{vp}}$$

Thus $\mathbb{E}\{e^{\alpha(X+Y)}\} \leq \exp\left\{\frac{\alpha^2}{2}(\|X\|_{vp} + \|Y\|_{vp})^2\right\}$ (4)

(iii) Note that $\forall \alpha \in \mathbb{R}$ $1 = e^{\mathbb{E}(\alpha X)} \leq \mathbb{E}(e^{\alpha X})$

\uparrow
X assumed centered

In addition, $\|X\|_{vp} = 0 \Rightarrow \mathbb{E}(e^{\alpha X}) \leq 1$.
 \uparrow
X sub-gauss

$\Rightarrow e^{\mathbb{E}(\alpha X)} = \mathbb{E}(e^{\alpha X})$, which can happen only if $X = \text{constant a.s.}$ (here equal to 0)

Theorem: Suppose X sub-Gaussian. Then $\forall t \in \mathbb{R}$,

$$\left\{ \begin{array}{l} \mathbb{P}(X - \mathbb{E}X > t) \\ \mathbb{P}(X - \mathbb{E}X < -t) \end{array} \right\} \leq \exp\left(-\frac{t^2}{2\|X\|_{vp}^2}\right).$$

proof = obvious

\uparrow
X is sub gaussian
 $\Leftrightarrow -X$ is sub gaussian.

* Ex: Any bounded RV is sub-Gaussian:

let $X \in [a, b]$ a.s.

$$\text{Then } \forall \alpha \in \mathbb{R}, \log \left[\mathbb{E}\{e^{\alpha(X-\mathbb{E}X)}\} \right] \leq \frac{\alpha^2(b-a)^2}{8}$$

$$\Rightarrow \|X\|_{vp}^2 \leq \frac{(b-a)^2}{4}. \quad (\text{not obvious - proof required})$$

Remark = Any Gaussian RV is sub-Gaussian with variance proxy equal to its variance. However, in general, any sub-Gaussian RV X has variance smaller or equal to its variance proxy: $\text{Var } X \leq \|X\|_{vp}^2$.

I.2. Sub-gamma distributions.

(5)

Recall that for $X \sim \gamma(a, b)$, $a, b > 0$,

$$f_X(x) = \frac{b^a}{\Gamma(a)} x^{a-1} e^{-bx}, \quad x \geq 0.$$

$$\mathbb{E}\{e^{\theta(X - \mathbb{E}X)}\} = \begin{cases} \left(\frac{b}{b-\theta}\right)^a e^{-\frac{a\theta}{b}} & \text{if } \theta < b \\ +\infty & \text{if } \theta \geq b \end{cases}$$

$$\text{Using } -\log(1-u) - u \leq \frac{u^2}{2(1-u)}, \quad \forall u \in (0, 1),$$

we have $\forall \theta < b$,

$$\begin{aligned} \log \left[\mathbb{E}\{e^{\theta(X - \mathbb{E}X)}\} \right] &= -a \log \left(1 - \frac{\theta}{b}\right) - \frac{a\theta}{b} \\ &\leq \frac{\theta^2 \alpha^2}{2(1-\theta\beta)} \quad \text{with } \alpha = \frac{\sqrt{a}}{b} \\ \text{leads to } &\quad \beta = \frac{1}{b} \end{aligned}$$

Def = A real-valued RV X is said to be sub-gamma on the right tail if $\exists \alpha > 0, \beta > 0$ s.t.

$$\forall \theta < \frac{1}{\beta} \quad \log \left[\mathbb{E}\{e^{\theta(X - \mathbb{E}X)}\} \right] \leq \frac{\theta^2 \alpha^2}{2(1-\theta\beta)}$$

and we write $X \in \Gamma_+(\alpha, \beta)$

$$X \text{ is said to be sub-gamma in the left tail if } \forall \theta > -\frac{1}{\beta} \quad \log \left[\mathbb{E}\{e^{\theta(X - \mathbb{E}X)}\} \right] \leq \frac{\theta^2 \alpha^2}{2(1-\theta\beta)},$$

and we write $X \in \Gamma_-(\alpha, \beta)$

$$X \in \Gamma_+(\alpha, \beta) \Leftrightarrow -X \in \Gamma_-(\alpha, \beta).$$

(6)

We say that X is sub-gamma, and we write $X \in \Gamma(\alpha, \beta)$ if $X \in \Gamma_+(\alpha, \beta)$ and $-X \in \Gamma_-(\alpha, \beta)$.

From previous calculations, we see that for $a, b > 0$,

$$X \sim \gamma(a, b)$$

$$\Rightarrow X \in \Gamma_+\left(\frac{\sqrt{a}}{b}, \frac{1}{b}\right)$$

$$-X \in \Gamma\left(\frac{\sqrt{a}}{b}, \frac{1}{b}\right).$$

Remark: By definition, any sub-gaussian variable X is sub-gamma: $X \in \Gamma(\|X\|_{\text{vp}}, \beta)$ $\forall \beta > 0$. Moreover, if X is sub-gaussian, then $X^2 \in \Gamma(4\|X\|_{\text{vp}}^2, 4\|X\|_{\text{vp}}^2)$ (proof required)

We finish this section with a concentration inequality for sub-gamma RVs:

Theorem: Let X_1, \dots, X_n be independent sub-gamma RVs,

$X_i \in \Gamma(a_i, b_i)$. Then, $\forall \varepsilon \in (0, 1)$, $\forall t > 0$,

$$\begin{cases} \mathbb{P}\left(\sum_{i=1}^n (X_i - \mathbb{E}X_i) > t\right) \\ \mathbb{P}\left(\sum_{i=1}^n (X_i - \mathbb{E}X_i) < -t\right) \end{cases} \leq e^{-\frac{t}{2} \left(\frac{\varepsilon}{\max b_i} \wedge \frac{(1-\varepsilon)t}{\sum a_i^2} \right)}$$

Proof of this theorem + proofs of other results presented here can be found in the teaching material of Quentin Paris: www-qparis-math.com/teaching [High dimensional statistical methods]

II. JOHNSON - LINDENSTRAUSS THEOREM.

(7)

We are now in a position to solve the problem stated in the introduction, and to prove the following result, known as the Johnson-Lindenstrauss Theorem.

JOHNSON - LINDENSTRAUSS THEOREM

Let $x_1, \dots, x_n = n$ distinct (non-random) points in \mathbb{R}^D .

Fix $\varepsilon > 0$, $\delta \in (0, 1)$.

Let M be a $d \times D$ random matrix, whose rows r_1, \dots, r_d are \rightarrow independent

\rightarrow centered

\rightarrow isotropic $E(r_i r_i^t) = I_D$

\rightarrow sub-gaussian with variance proxy σ^2 .

Put $T := \frac{1}{\sqrt{d}} M$. $\hookrightarrow u^t r_i$ is sub-gaussian $\forall u \in \mathbb{R}^d$

Then, provided $d \geq \frac{\sigma^4}{\varepsilon^2} \log(\frac{n}{\delta})$,

constant times ...

(constant is explicit - see proof of theorem)

we have

$$P(\forall i \neq j : 1 - \varepsilon \leq \frac{\|T(x_i) - T(x_j)\|_2}{\|x_i - x_j\|_2} \leq 1 + \varepsilon) \geq 1 - \delta$$

\hookleftarrow T is very easy to implement.

Take for example a matrix $M \in \mathbb{R}^{d \times D}$ with iid entries $N(0, \sigma^2)$, or any distribution with bounded support.

proof = Denote $X := \{x_1, \dots, x_n\}$.

By linearity of T ,

$$\forall x, y \in X, \quad 1 - \varepsilon \leq \frac{\|T(x) - T(y)\|_2}{\|x - y\|_2} \leq 1 + \varepsilon$$

$$\Leftrightarrow 1 - \varepsilon \leq \|T(z)\|_2 \leq 1 + \varepsilon$$

$$\text{with } z := \frac{x - y}{\|x - y\|_2}$$

Since $(1 - \varepsilon)^2 \leq 1 - \varepsilon$, and $1 + \varepsilon \leq (1 + \varepsilon)^2$, the last written expression is equivalent to $1 - \varepsilon \leq \|T(z)\|_2^2 \leq 1 + \varepsilon$, and so it is enough to show that

$$P(\forall z \in \mathcal{Z} : 1 - \varepsilon \leq \|T(z)\|_2^2 \leq 1 + \varepsilon) \geq 1 - \delta, \quad (\#)$$

$$\text{where } \mathcal{Z} := \left\{ \frac{x - y}{\|x - y\|_2}, \quad x \neq y \in X \right\}.$$

$$\text{By construction, } T = \frac{1}{\sqrt{d}} M = \frac{1}{\sqrt{d}} \begin{pmatrix} -r_1 \\ -r_2 \\ \vdots \\ -r_d \end{pmatrix} \updownarrow^d,$$

$$\text{so that } \|T(z)\|_2^2 = \frac{1}{d} \sum_{i=1}^d r_i^t z \quad \updownarrow \text{Also denoted } \langle r_i, z \rangle.$$

Note that z is a unit vector.

+ r_i is sub-gaussian with variance proxy σ^2

def of the variance proxy of a random vector $\rightarrow \max_{\|u\|=1} \|u^t r_i\|_{\text{vp}}^2$

$$\Rightarrow \langle r_i, z \rangle \in SG(\sigma^2)$$

$$\& \langle r_i, z \rangle^2 \in \Gamma(4\sigma^2, 4\sigma^2)$$

In addition,

$$1 - \varepsilon \leq \|T(z)\|_2^2 \leq 1 + \varepsilon \Leftrightarrow \left| \frac{1}{d} \sum_{i=1}^d \langle r_i, z \rangle^2 - 1 \right| \leq \varepsilon \quad (9)$$

&

$$\begin{aligned} \mathbb{E} \langle r_i, z \rangle^2 &= \mathbb{E} ((r_i^t z)(r_i^t z)) \\ &= \mathbb{E} (z^t r_i r_i^t z) \\ &= z^t (\underbrace{\mathbb{E} r_i r_i^t}_\equiv) z = z^t z = 1. \end{aligned}$$

Thus,

$$\left| \frac{1}{d} \sum_{i=1}^d \langle r_i, z \rangle^2 - 1 \right| = \frac{1}{d} \left| \sum_{i=1}^d (\langle r_i, z \rangle^2 - \mathbb{E} \langle r_i, z \rangle^2) \right|,$$

and so it is enough to show that (see (*) page 8)

$$(\text{**}) \quad \mathbb{P} \left(\max_{z \in \mathcal{Z}} \left| \sum_{i=1}^d Y_i(z) - \mathbb{E} Y_i(z) \right| \geq d\varepsilon \right) \geq 1 - \delta,$$

where

$$Y_i(z) := \langle r_i, z \rangle^2 \in \Gamma(4e\sigma^2, 4e\sigma^2).$$

$$(\text{***}) \quad \mathbb{P} \left(\max_{z \in \mathcal{Z}} \left| \sum_{i=1}^d Y_i(z) - \mathbb{E} Y_i(z) \right| \geq d\varepsilon \right) \leq \delta.$$

Union-bound:

$$\mathbb{P} \left(\max_{z \in \mathcal{Z}} | \dots | \geq d\varepsilon \right) \leq |\mathcal{Z}| \max_{z \in \mathcal{Z}} \mathbb{P} (| \dots | \geq d\varepsilon)$$

This term can be bounded independently of z using concentration bounds for sub-gamma RVs (bottom of page 6).

$$\mathbb{P} (| \dots | \geq d\varepsilon) \leq 2 \exp \left\{ -\frac{1}{4} \left(\frac{d\varepsilon}{4e\sigma^2} \wedge \frac{d^2\varepsilon^2}{(4e\sigma^2)^2} \right) \right\} \quad (10)$$

absolute value

Take $\varepsilon = \frac{1}{2}$ in the theorem page 6, and using $t = d\varepsilon$.

In addition, $|\mathcal{Z}| \leq n^2$.

$$\Rightarrow \mathbb{P} \left(\max_{z \in \mathcal{Z}} | \dots | \geq d\varepsilon \right) \leq 2n^2 \exp \left\{ -\frac{d\varepsilon}{16e\sigma^2} \left(1 \wedge \frac{\varepsilon}{4e\sigma^2} \right) \right\}$$

We want this term to be less than $\delta \in (0, 1)$.

\Rightarrow Select d appropriately:

$$2n^2 \exp \left\{ -\frac{d\varepsilon}{16e\sigma^2} \left(1 \wedge \frac{\varepsilon}{4e\sigma^2} \right) \right\} \leq \delta$$

$$\Leftrightarrow -\frac{d\varepsilon}{16e\sigma^2} \left(1 \wedge \frac{\varepsilon}{4e\sigma^2} \right) \leq \log \left(\frac{\delta}{2n^2} \right)$$

$$\Leftrightarrow d \left(1 \wedge \frac{\varepsilon}{4e\sigma^2} \right) \geq \frac{16e\sigma^2}{\varepsilon} \log \left(\frac{2n^2}{\delta} \right)$$

$$\Leftrightarrow d \geq \frac{32e\sigma^2}{\varepsilon} \left(1 \vee \frac{4e\sigma^2}{\varepsilon} \right) \log \left(\sqrt{\frac{2}{\delta}} n \right),$$

which concludes the proof. \blacksquare

II - APPLICATION TO K-MEANS.

Random projections are particularly well-suited as a pre-processing step for the K-means algorithm, since it preserves the distance between points.

K-means revisited =

(11)

- Input: $x_1, \dots, x_n \in \mathbb{R}^D$, $T = \text{linear transformation}$.
- Compute: y_1, \dots, y_n ; $y_i = Tx_i \in \mathbb{R}^d$

→ Select $c_1^0, \dots, c_K^0 \in \mathbb{R}^d$

→ Construct clusters

$$\mathcal{C}_j^0 = \left\{ y_i \mid \|y_i - c_j^0\|_2 \leq \|y_i - c_l^0\|_2 \quad \forall l \neq j \right\}$$

→ Repeat $m = 0, 1, \dots$

$$c_j^{m+1} := \frac{1}{|\mathcal{C}_j^m|} \sum_{i \in \mathcal{C}_j^m} y_i$$

$$\mathcal{C}_j^{m+1} := \left\{ y_i \mid \|y_i - c_j^{m+1}\|_2 \leq \|y_i - c_l^{m+1}\|_2 \quad \forall l \neq j \right\}$$

- Output: $\mathcal{C}_j^{\text{final}} / c_j^{\text{final}}$.

and cluster x_1, \dots, x_n such that

$$\mathcal{C}_j := \left\{ x_i \mid y_i = Tx_i \in \mathcal{C}_j^{\text{final}} \right\}, \quad 1 \leq j \leq K$$



Performance bounds in Biau, Devroye, Lugosi (2008)

On the performance of clustering in Hilbert Spaces.

IEEE Trans. on Information Theory, Vol 54, Issue 2