

# Dimensionality Reduction

Master 2 Maths en Action

---

Aurélien Garivier

2018-2019



# Table of contents

1. Dimension reduction: PCA
2. Dimension reduction: random projections

# Dimensionality reduction

- Data:  $X = \begin{pmatrix} x_1^T \\ \vdots \\ x_n^T \end{pmatrix} \in \mathcal{M}_{n,p}(\mathbb{R})$ ,  $p \gg 1$ .
- Dimensionality reduction: replace  $x_i$  with  $y_i = Wx_i$ , where  $W \in \mathcal{M}_{d,p}(\mathbb{R})$ ,  $d \ll p$ .
- Hopefully, we do not lose too much by replacing  $x_i$  by the  $y_i$ .  
2 approaches:
  - Quasi-invertibility: there exists a recovering matrix  $U \in \mathcal{M}_{p,d}(\mathbb{R})$  such that for all  $i \in \{1, \dots, n\}$ ,

$$\tilde{x}_i = Uy_i \approx x_i .$$

- More modest goal: distance-preserving property

$$\forall 1 \leq i, j \leq n, \quad \frac{\|Wx_1 - Wx_2\|}{\|x_1 - x_2\|} \approx 1 .$$

## **Dimension reduction: PCA**

---

PCA aims at finding the compression matrix  $W$  and the recovering matrix  $U$  such that the total squared distance between the original and the recovered vectors is minimal:

$$\arg \min_{W \in \mathcal{M}_{d,p}(\mathbb{R}), U \in \mathcal{U} \subset \mathcal{M}_{p,d}(\mathbb{R})} \sum_{i=1}^n \|x_i - UWx_i\|^2.$$

**Property.** A solution  $(W, U)$  is such that  $U^T U = I_d$  and  $W = U^T$ .

**Proof.** Let  $W \in \mathcal{M}_{n,p}(\mathbb{R})$ ,  $U \in \mathcal{U} \subset \mathcal{M}_{p,d}(\mathbb{R})$ , and let  $R = \{UWx : x \in \mathbb{R}^p\}$ .  $\dim(R) \leq d$ , and we can assume that  $\dim(R) = d$ . Let  $V = (v_1 \mid \dots \mid v_d) \in \mathcal{M}_{p,d}(\mathbb{R})$  be an orthogonal basis of  $R$ , hence  $V^T V = I_d$  and for every  $\tilde{x} \in R^p$  there exists  $y \in \mathbb{R}^d$  such that  $\tilde{x} = Vy$ . But for every  $x \in \mathbb{R}^p$ ,

$$\arg \min_{\tilde{x} \in R} \|x - \tilde{x}\|^2 = V \cdot \arg \min_{y \in \mathbb{R}^d} \|x - Vy\|^2 = V \cdot \arg \min_{y \in \mathbb{R}^d} \|x\|^2 + \|y\|^2 - 2y^T (V^T x) = VV^T x$$

(as can be seen easily by differentiation in  $y$ ), and hence

$$\sum_{i=1}^n \|x_i - UWx_i\|^2 \geq \sum_{i=1}^n \|x_i - VV^T x_i\|^2.$$

# The PCA solution

Corollary: the optimization problem can be rewritten

$$\arg \min_{U \in \mathcal{U} \in \mathcal{M}_{p,d}(\mathbb{R}): U^T U = I_d} \sum_{i=1}^n \|x_i - UU^T x_i\|^2.$$

Since  $\|x_i - UU^T x_i\|^2 = \|x_i\|^2 - \text{Tr}(U^T x_i x_i^T U)$ , this is equivalent to

$$\arg \max_{U \in \mathcal{U} \in \mathcal{M}_{p,d}(\mathbb{R}): U^T U = I_d} \text{Tr} \left( U^T \sum_{i=1}^n x_i x_i^T U \right).$$

Let  $A = \sum_{i=1}^n x_i x_i^T$ , and let  $A = VDV^T$  be its spectral decomposition:  $D$  is diagonal, with  $D_{1,1} \geq \dots \geq D_{p,p} \geq 0$  and  $V^T V = VV^T = I_p$ .

# Solving PCA by SVD

**Theorem** Let  $A = \sum_{i=1}^n x_i x_i^T$ , and let  $u_1, \dots, u_d$  be the eigenvectors of  $A$  corresponding to the  $d$  largest eigenvalues of  $A$ . Then the solution to the PCA optimization problem is  $U = \left( u_1 \mid \dots \mid u_d \right)$ , and  $W = U^T$ .

**Proof.** Let  $U \in \mathcal{M}_{p,d}(\mathbb{R})$  be such that  $U^T U = I_d$ , and let  $B = V^T U$ . Then  $VB = U$ , and  $U^T A U = B^T V^T V D V^T V B = B^T D B$ , hence

$$\text{Tr}(U^T A U) = \sum_{j=1}^p D_{j,j} \sum_{i=1}^d B_{j,i}^2.$$

Since  $B^T B = U^T V V^T U = I_d$ , the columns of  $B$  are orthonormal and  $\sum_{j=1}^p \sum_{i=1}^d B_{j,i}^2 = d$ .

In addition, completing the columns of  $B$  to an orthonormal basis of  $\mathbb{R}^p$  one gets  $\tilde{B}$  such that  $\tilde{B}^T \tilde{B} = I_p$ , and for every  $j$  one has  $\sum_{i=1}^p \tilde{B}_{j,i}^2 = 1$ , hence  $\sum_{i=1}^d B_{j,i}^2 \leq 1$ .

Thus,

$$\text{Tr}(U^T A U) \leq \max_{\beta \in [0,1]^p: \|\beta\|_1 \leq d} \sum_{j=1}^p D_{j,j} \beta_j = \sum_{j=1}^d D_{j,j},$$

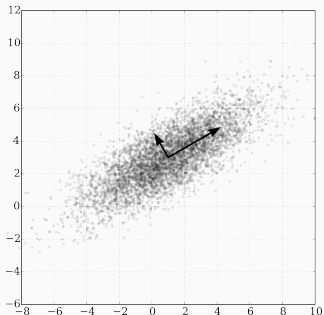
which can be reached if  $U$  is made of the  $d$  leading eigenvectors of  $A$ .

# PCA: comments

Interpretation: PCA aims at maximizing the projected variance.

Often, the quality of the result is measured by the proportion of the variance explained by the  $d$  principal components:

$$\frac{\sum_{i=1}^d D_{i,i}}{\sum_{i=1}^p D_{i,i}}.$$



[Src: wikipedia.org]

In practice: sometimes cheaper to compute svp of  $B = X^T X \in \mathcal{M}_n(\mathbb{R})$ , since if  $u$  is such that  $Bu = \lambda u$  then for  $v = X^T u / \|X^T u\|$  one has  $Av = \lambda v$ .



# Computing the PCA: iteration method

Let  $\lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n$  be the eigenvalues of  $A$ , and let  $v^1$  be such that  $\|v^1\| = 1$  and  $Av^1 = \lambda_1 v^1$ . Goal: approximate  $v^1$ .

Algorithm:  $u_0 = \left[ \frac{\epsilon_1}{\sqrt{n}}, \dots, \frac{\epsilon_n}{\sqrt{n}} \right]$  where  $\epsilon_i \stackrel{iid}{\sim} \mathcal{U}(\{-1, 1\})$ , then  $\|u_0\|^2 = 1$ .  
 $u_{k+1} = \frac{Au_k}{\|Au_k\|}$ .

## Theorem

With probability at least  $3/16$ ,

$$|\langle u_t, v^1 \rangle| \geq 1 - 2n \left( \frac{\lambda_2}{\lambda_1} \right)^{2t}.$$

Thus, it takes at most  $t = \frac{\log \frac{2n}{\epsilon}}{2 \log \frac{\lambda_1}{\lambda_2}}$  iterations to ensure that

$$|\langle u_t, v^1 \rangle| \geq 1 - \epsilon.$$

Remark: one can similarly show that with non-vanishing probability

$$\langle u_t, Au_t \rangle \geq \lambda_1 \times \frac{1-\epsilon}{1+4n(1-\epsilon)^{2t}}. \quad \text{http://theory.stanford.edu/~trevisan/expander-online/lecture03.pdf.}$$

## The complexity of the iteration method 1/2

Observe that  $\langle u_0, v^1 \rangle$  has expectation 0 and variance  $\sum_{i=1}^n (v_i^1)^2/n = 1/n$ . Hence,  $Z = \langle u_0, v^1 \rangle^2$  has expectation  $1/n$  and variance such that

$$\begin{aligned} n^2 \text{Var}[Z] &= \mathbb{E} \left[ \sum_{1 \leq i, j, k, l \leq d} \epsilon_i \epsilon_j \epsilon_k \epsilon_l \right] = \sum_{1 \leq j \leq d} (v_j^1)^4 + 6 \sum_{1 \leq j < k \leq d} (v_j^1)^2 (v_k^1)^2 \\ &= 3 \left( \|v\|^2 \right)^2 - 2 \sum_{1 \leq j \leq d} (v_j^1)^4 \leq 3. \end{aligned}$$

By the Cauchy-Schwartz inequality, for every  $\delta \in (0, 1)$

$$\mathbb{E}[Z] = \mathbb{E}[Z \mathbb{1}\{Z < \delta \mathbb{E}[Z]\}] + \mathbb{E}[Z \mathbb{1}\{Z \geq \delta \mathbb{E}[Z]\}] \leq \delta \mathbb{E}[Z] + \mathbb{E}[Z^2] \mathbb{P}(Z \geq \delta \mathbb{E}[Z]).$$

and hence, for  $\delta = 1/4$ :

$$\mathbb{P}(Z \geq \delta \mathbb{E}[Z]) \geq (1 - \delta)^2 \frac{\mathbb{E}[Z]^2}{\mathbb{E}[Z^2]} \geq \left(\frac{3}{4}\right)^2 \frac{1/n^2}{3/n^2} = \frac{9}{16} \times \frac{1}{3} \geq \frac{3}{16}.$$

## The complexity of the iteration method 2/2

But whenever  $\langle u_0, v^1 \rangle^2 > \frac{1}{4n}$ :

$$\begin{aligned} |\langle u_t, v^1 \rangle| &= \frac{|\langle u_0, v^1 \rangle| \lambda_1^t}{\sqrt{\sum_{i=1}^n \langle u_0, v^i \rangle^2 \lambda_j^{2t}}} = \frac{1}{\sqrt{1 + \frac{1}{\langle u_0, v^1 \rangle^2} \sum_{i=2}^n \langle u_0, v^i \rangle^2 \left(\frac{\lambda_j}{\lambda_1}\right)^{2t}}} \\ &\geq \frac{1}{\sqrt{1 + 4n \sum_{i=2}^n \langle u_0, v^i \rangle^2 \left(\frac{\lambda_2}{\lambda_1}\right)^{2t}}} \\ &\geq 1 - 2n \left(\frac{\lambda_2}{\lambda_1}\right)^{2t}. \end{aligned}$$

# Dimension reduction: random projections

---

# Johnson-Lindenstrauss Lemma

## Theorem

Let  $x_1, \dots, x_n \in \mathbb{R}^p$ , and let  $\epsilon > 0$ . Then, for every

$d \geq \frac{4 \log(n)}{-\log(1 - 2\epsilon) - 2\epsilon}$ , there exists a matrix  $W \in \mathcal{M}_{d,p}(\mathbb{R})$  such that

$$\forall 1 \leq i \leq j, \quad (1 - \epsilon) \|x_i - x_j\|^2 \leq \|Wx_i - Wx_j\|^2 \leq (1 + \epsilon) \|x_i - x_j\|^2.$$

**Remark 1: on the dependence on  $\epsilon$ .**

$$\frac{4 \log(n)}{-\log(1 - 2\epsilon) - 2\epsilon} \leq \frac{8 \log(n)}{\epsilon^2} \left(1 + \frac{\epsilon}{3}\right)^2.$$

**Remark 2: how to find such a matrix  $W$ .**

For every  $d \geq \frac{4 \log(n) + 2 \log(1/\delta)}{-\log(1 - 2\epsilon) - 2\epsilon}$ , the probability that a random matrix with entries  $W_{i,j} \stackrel{iid}{\sim} \mathcal{N}(0, \frac{1}{d})$  satisfies the lemma is larger than  $1 - \delta$ .

# Proof of the Johnson-Lindenstrauss Lemma

Method: (constructive) probabilistic method. We choose  $W_{i,j} \stackrel{iid}{\sim} \mathcal{N}(0, \frac{1}{d})$ . Let  $y \in \mathbb{R}^p$  and  $Y = Wy$ . Then, for all  $1 \leq i \leq d$ ,

$Y_i = \sum_{j=1}^p y_j W_{i,j} \sim \mathcal{N}(0, \frac{\|y\|^2}{d})$ . Hence  $\mathbb{E}[\|Y\|^2] = \|y\|^2$ . Besides, by the deviation bound for the  $\chi^2$  distribution presented below,

$$\mathbb{P}\left(\|Y\|^2 \geq (1+\epsilon)\|y\|^2\right) = \mathbb{P}\left(\sum_{i=1}^d \left(\frac{\sqrt{d}Y_i}{\|y\|}\right)^2 \geq d(1+\epsilon)\right) \leq \exp(-d\phi^*(\epsilon)) \leq \frac{1}{n^2}$$

and similarly  $\mathbb{P}\left(\|Y\|^2 \leq (1-\epsilon)\|y\|^2\right) \leq \exp(-d\phi^*(\epsilon)) \leq \frac{1}{n^2}$ .

Applying this result to all  $y_{i,j} = x_i - x_j$ ,  $1 \leq i < j \leq n$ , we obtain the conclusion by the union bound:

$$\begin{aligned} \mathbb{P}\left(\bigcup_{1 \leq i < j \leq n} \left(\|W(x_i - x_j)\| \geq (1+\epsilon)\|x_i - x_j\| \cup \|W(x_i - x_j)\| \leq (1-\epsilon)\|x_i - x_j\|\right)\right) \\ \leq \frac{n(n-1)}{n^2} < 1, \end{aligned}$$

and hence there exists at least a matrix  $W$  for which the lemma holds.

# Deviations of the $\chi^2$ distribution: rate function

## Lemma

If  $U \sim \mathcal{N}(0, 1)$  and  $X = U^2 - 1$ , then

$$\phi^*(x) = \sup_{\lambda} \lambda x - \log \mathbb{E} [e^{\lambda X}] = \frac{x - \log(1+x)}{2} \geq \frac{x^2}{4 \left(1 + \frac{x}{3}\right)^2}.$$

**Proof:** For every  $\lambda < 1/2$ ,

$$\mathbb{E} [e^{\lambda X}] = \frac{1}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{\lambda(u^2-1)} e^{-\frac{u^2}{2}} du = \frac{e^{-\lambda}}{\sqrt{2\pi}} \int_{\mathbb{R}} e^{-\frac{(1-2\lambda)u^2}{2}} du = e^{-\lambda} \frac{1}{\sqrt{1-2\lambda}}.$$

Hence  $\phi(\lambda) = \log \mathbb{E} [e^{\lambda X}] = -\frac{1}{2} \log(1-2\lambda) - \lambda$ . The concave function  $\lambda \mapsto \lambda x - \phi(\lambda)$  is maximized at  $\lambda^*$  s.t.  $0 = \phi'(\lambda^*) = \frac{1}{1-2\lambda^*} - 1 - x$ , that is at  $\lambda^* = \frac{1}{2} \left(1 - \frac{1}{1+x}\right) = \frac{x}{2(1+x)}$ . Hence

$$\phi^*(x) = \lambda^* x - \phi(\lambda^*) = \frac{x - \log(1+x)}{2}.$$

The last inequality is obtained by "Pollard's trick" applied to  $g(x) = x - \log(1+x)$ : since  $g(0) = g'(0) = 0$  and since  $g''(x) = 1/(1+x)^2$  is convex, by Jensen's inequality

$$\frac{x - \log(1+x)}{x^2/2} = \int_0^1 g''(sx) 2(1-s) ds \geq g'' \left( \int_0^1 sx 2(1-s) ds \right) = g'' \left( \frac{x}{3} \right).$$

# Deviations of the $\chi^2(d)$ distribution

By Chernoff's method, if  $Z \sim \chi^2(d) \stackrel{\text{dist}}{=} U_1^2 + \dots + U_d^2$  where  $U_i \stackrel{iid}{\sim} \mathcal{N}(0, 1)$ :

$$\mathbb{P}(Z \geq d(1 + \epsilon)) \leq \exp(-d\phi^*(\epsilon)) \leq \exp\left(-\frac{d\epsilon^2}{4\left(1 + \frac{\epsilon}{3}\right)^2}\right).$$

Note: the Laurent-Massart inequality states that for every  $u > 0$ ,

$$\mathbb{P}(Z \geq d + 2\sqrt{du} + 2u) \leq \exp(-u).$$

It can be deduced from the previous bound by noting that for every  $u > 0$

$$\begin{aligned}\phi^*(2\sqrt{u} + 2u) &= u + \frac{1}{2} \left( 2\sqrt{u} - \log \left( 1 + 2\sqrt{u} + \frac{(2\sqrt{u})^2}{2} \right) \right) \\ &\geq u + \frac{1}{2} \left( 2\sqrt{u} - \log(\exp(2\sqrt{u})) \right) \geq u.\end{aligned}$$

The proof of Laurent and Massart (which takes elements from Birgé and Massart 1998) is a bit different: they note that

$$\phi(\lambda) = -\frac{1}{2} \log(1 - 2\lambda) - \lambda = \sum_{k=2}^{\infty} \frac{(2\lambda)^k}{2k} = \lambda^2 \sum_{\ell=0}^{\infty} \frac{4(2\lambda)^\ell}{2(\ell+2)} \leq \lambda^2 \sum_{\ell=0}^{\infty} (2\lambda)^\ell = \frac{\lambda^2}{1 - 2\lambda},$$
 and deduce that

$$\phi^*(x) \geq \psi^*(x) = \sup_{\lambda} \lambda x - \frac{\lambda^2}{1 - 2\lambda} = \frac{x + 1 - \sqrt{2x + 1}}{2},$$
 while  $x > 0$  and  $\psi^*(x) = u$  implies  $x = 2\sqrt{u} + 2u$ . Also note in

$$\text{passing that by Pollard's trick } \phi^*(x) \geq \psi^*(x) \geq \frac{x^2}{4\left(1 + \frac{2x}{3}\right)^{3/2}}.$$

Moreover, since  $\phi^*(-\epsilon) = -(\epsilon + \log(1 - \epsilon))/2 \geq \epsilon^2/4$ ,

$$\mathbb{P}(Z \leq d(1 - \epsilon)) \leq \exp\left(-\frac{d\epsilon^2}{4}\right).$$