

Remise à niveau - Statistique

Réduction de dimension

Claire Boyer

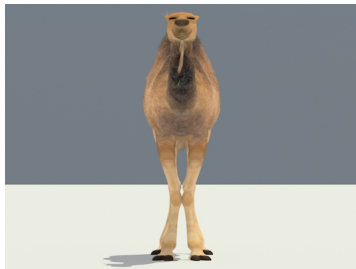
High Dimension Geometry Curse

- ▶ Folks theorem: In high dimension, everyone is alone.
- ▶ Theorem: If X_1, \dots, X_n in the hypercube of dimension d such that their coordinates are i.i.d then

$$d^{-1/p} (\max \|X_i - X_j\|_p - \min \|X_i - X_j\|_p) = 0 + O\left(\sqrt{\frac{\log n}{d}}\right)$$
$$\frac{\max \|X_i - X_j\|_p}{\min \|X_i - X_j\|_p} = 1 + O\left(\sqrt{\frac{\log n}{d}}\right).$$

- ▶ When d is large, all the points are almost equidistant...
- ▶ Nearest neighbors are meaningless!

- ▶ How to view a high-dimensional dataset?
- ▶ High-dimension: dimension larger than 2!
- ▶ *Projection* in a 2D space.



- ▶ How to view a high-dimensional dataset?
- ▶ High-dimension: dimension larger than 2!
- ▶ *Projection* in a 2D space.



- ▶ How to view a high-dimensional dataset?
- ▶ High-dimension: dimension larger than 2!
- ▶ *Projection* in a 2D space.



- ▶ How to view a high-dimensional dataset?
- ▶ High-dimension: dimension larger than 2!
- ▶ *Projection* in a 2D space.

Dimension reduction

- ▶ Training data $\mathcal{D} = \{x_1, \dots, x_n\} \in \mathcal{X}^n$ (i.i.d. $\sim \mathbb{P}$)
- ▶ Space \mathcal{X} of possibly high dimension.

Dimension Reduction Map

Construct a map Φ from the space \mathcal{X} into a space \mathcal{X}' of smaller dimension:

$$\begin{aligned}\Phi : \mathcal{X} &\rightarrow \mathcal{X}' \\ x &\mapsto \Phi(x)\end{aligned}$$

Criterion

- ▶ Reconstruction error **focus here in these slides**
- ▶ Distance preservation

- ▶ Construct a map Φ from the space \mathcal{X} into a space \mathcal{X}' of smaller dimension:

$$\begin{aligned}\Phi : \mathcal{X} &\rightarrow \mathcal{X}' \\ x &\mapsto \Phi(x)\end{aligned}$$

- ▶ Construct $\tilde{\Phi}$ from \mathcal{X}' to \mathcal{X}
- ▶ Control the error between x and its reconstruction $\tilde{\Phi}(\Phi(x))$
- ▶ Canonical example for $x \in \mathbb{R}^d$: find Φ and $\tilde{\Phi}$ in a parametric family that minimize

$$\frac{1}{n} \sum_{i=1}^n \|x_i - \tilde{\Phi}(\Phi(x_i))\|^2$$

- ▶ $x_1, \dots, x_n \in \mathbb{R}^d$
- ▶ $m = \frac{1}{n} \sum_{i=1}^n x_i$

Two views on inertia

- ▶ Inertia:

$$\begin{aligned} I &= \frac{1}{n} \sum_{i=1}^n \|x_i - m\|^2 \\ &= \frac{1}{2n^2} \sum_{i,j} \|x_i - x_j\|^2 \end{aligned}$$

- ▶ 2 times the mean squared distance to the mean = Mean squared distance between individual
- ▶ Heuristic: a good representation is a representation with a large inertia
- ▶ Large dispersion \sim Large average separation!

- ▶ What if we replace x by its projection $\tilde{x} = P(x - m) + m$?

Two views on inertia

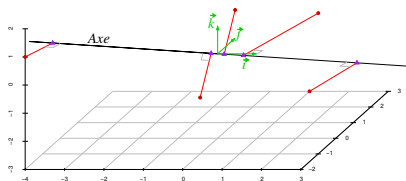
- ▶ Inertia:

$$\begin{aligned}\tilde{I} &= \frac{1}{n} \sum_{i=1}^n \|\tilde{x}_i - m\|^2 \\ &= \frac{1}{2n^2} \sum_{i,j} \|\tilde{x}_i - \tilde{x}_j\|^2\end{aligned}$$

- ▶ Inertia:

$$\begin{aligned}\tilde{I} &= I - \frac{1}{n} \sum_{i=1}^n \|\tilde{x}_i - x_i\|^2 \\ &= I - \frac{1}{2n^2} \sum_{i,j} (\|x_i - x_j\|^2 - \|\tilde{x}_i - \tilde{x}_j\|^2)\end{aligned}$$

- ▶ Four different way to obtain a large inertia!



- ▶ 1D case: $\tilde{x} = m + a^\top (x - m)a$ with $\|a\| = 1$
- ▶ Inertia:
$$\tilde{l} = \frac{1}{n} \sum_{i=1}^n a^\top (x_i - m)(x_i - m)^\top a$$

Principal Component Analysis : optimization of the projection

- ▶ Maximization of

$$\tilde{l} = \frac{1}{n} \sum_{i=1}^n a^\top (x_i - m)(x_i - m)^\top a = a^\top \Sigma a$$

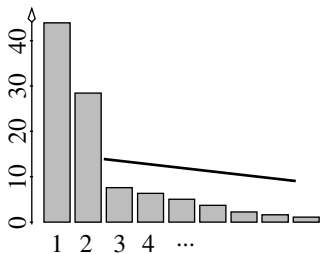
with $\Sigma = \frac{1}{n} \sum_{i=1}^n (x_i - m)(x_i - m)^\top$ the **empirical covariance matrix**.

- ▶ Explicit optimal choice given by the eigenvector of the largest eigenvalue of Σ .

Principal Component Analysis : optimization of the projection

- ▶ **Explicit** optimal solution obtain by the projection on the eigenvectors of the largest eigenvalues of Σ .
- ▶ Projected inertia given by the **sum of those eigenvalues**.

% d'inertie



- ▶ Often fast decay of the eigenvalues: some dimensions are much more important than other.
- ▶ Not exactly the curse of dimensionality setting...
- ▶ Yet a lot of *small* dimension can drive the distance

Take-home message

- ▶ Principal components = Eigenvectors of the empirical covariance matrix
- ▶ Principal components = "New" variables obtained as linear combinations of the initial variables
 - ▶ less interpretable
 - ▶ but capture the dataset variance better

- ▶ $\mathcal{X} \in \mathbb{R}^d$ and $\mathcal{X}' = \mathbb{R}^{d'}$
- ▶ Linear map with V orthonormal

$$\Phi(x) = V^\top(x - m) \quad \text{and} \quad \tilde{\Phi}(x') = m + Vx'$$

- ▶ Reconstruction error criterion:

$$\frac{1}{n} \sum_{i=1}^n \|x_i - (m + VV^\top(x_i - m))\|^2$$

- ▶ *Explicit solution:*
 - ▶ m is the empirical mean
 - ▶ V is any orthonormal basis of the space spanned by the d' first eigenvectors (the one with largest eigenvalues) of the empirical covariance matrix $\frac{1}{n} \sum_{i=1}^n (x_i - m)(x_i - m)^\top$.

PCA Algorithm

- ▶ Compute the empirical mean $m = \frac{1}{n} \sum_{i=1}^n x_i$
- ▶ Compute the empirical covariance matrix $\frac{1}{n} \sum_{i=1}^n (x_i - m)(x_i - m)^\top$.
- ▶ Compute the d' first eigenvectors of this matrix: $V^{(1)}, \dots, V^{(d')}$
- ▶ Set $\Phi(x) = V^\top(x - m)$

Remarks

- ▶ Complexity: $O(n(1 + d^2) + d'd^2)$
- ▶ Interpretation:
 $\Phi(x) = V^\top(x - m)$: coordinates in the restricted space
 $V^{(i)}$: influence of each original coordinates in the i th new one.
- ▶ Not invariant to a scaling of the variables
 \rightsquigarrow It is custom to normalize the variables before applying PCA.

- ▶ PCA assumes $\mathcal{X} = \mathbb{R}^d$
- ▶ How to deal with categorical values?
- ▶ MFA = PCA with clever coding strategy for categorical values.

Categorical value code for a single variable

- ▶ Classical redundant dummy coding:

$$x \in \{1, \dots, C\} \mapsto P(x) = (1_{x=1}, \dots, 1_{x=C})^t$$

- ▶ Compute the mean (i.e. the empirical proportions):

$$\bar{P} = \frac{1}{n} \sum_{i=1}^n P(x_i)$$

- ▶ Renormalize $P(x)$ by $1/\sqrt{(C-1)\bar{P}}$:

$$P(x) \mapsto P^r(x)$$

$$(1_{x=1}, \dots, 1_{x=C}) \mapsto \left(\frac{1_{x=1}}{\sqrt{(C-1)\bar{P}_1}}, \dots, \frac{1_{x=C}}{\sqrt{(C-1)\bar{P}_v}} \right)$$

- ▶ PCA becomes the minimization of

$$\begin{aligned} & \frac{1}{n} \sum_{i=1}^n \|P^r(x_i) - (m + VV^t(P^r(x_i) - m))\|^2 \\ &= \frac{1}{n} \sum_{i=1}^n \sum_{v=1}^V \frac{\left| 1_{x_i=v} - (m' + \sum_{l=1}^{d'} V^{(l)t}(P(x_i) - m')V^{(l,v)}) \right|^2}{(V-1)\bar{P}_v} \end{aligned}$$

- ▶ Interpretation:

- ▶ $m' = \bar{P}$
- ▶ $\Phi(x) = V^t(P^r x - m)$: coordinates in the restricted space.
- ▶ $V^{(l)}$ can be interpreted as a probability profile.

- ▶ Complexity: $O(n(V + V^2) + d'V^2)$
- ▶ Link with Correspondence Analysis (CA)

MFA Algorithm

- ▶ Redundant dummy coding of each categorical variable.
 - ▶ Renormalization of each block of dummy variable.
 - ▶ Classical PCA algorithm on the resulting variables
-
- ▶ Interpretation as a reconstruction error with a rescaled/ χ^2 metric.
 - ▶ Interpretation:
 - ▶ $\Phi(x) = V^t(P^r(x) - m)$: coordinates in the restricted space.
 - ▶ $V^{(l)}$: influence of each modality/variable in the l th new coordinates.
 - ▶ **Scaling:** This method is not invariant to a scaling of the continuous variables
↪ It is custom to normalize the variables (at least within groups) before applying PCA.

PCA Model

- ▶ PCA: Linear model assumption

$$\mathbf{x} \simeq \mathbf{m} + \sum_{l=1}^{d'} x'^{(l)} \mathbf{V}^{(l)} = \mathbf{m} + \mathbf{V} \mathbf{x}'$$

- ▶ with
 - ▶ $\mathbf{V}^{(l)}$ orthonormal
 - ▶ $x'^{(l)}$ without constrains.
- ▶ Two directions of extension:
 - ▶ Other constrains on \mathbf{V} (or the coordinates in the restricted space):
ICA, NMF, Dictionary approach
 - ▶ PCA on a non linear image of \mathbf{x} : kernel-PCA
- ▶ Much more complex algorithm!

ICA (Independent Component Analysis)

- ▶ Linear model assumption

$$\mathbf{x} \simeq \mathbf{m} + \sum_{l=1}^{d'} \mathbf{x}'^{(l)} \mathbf{V}^{(l)} = \mathbf{m} + \mathbf{V} \mathbf{x}'$$

- ▶ with
 - ▶ $\mathbf{V}^{(l)}$ without constrains.
 - ▶ $\mathbf{x}'^{(l)}$ independent

NMF (Non Negative Matrix Factorization)

- ▶ (Linear) Model assumption

$$\mathbf{x} \simeq \mathbf{m} + \sum_{l=1}^{d'} \mathbf{x}'^{(l)} \mathbf{V}^{(l)} = \mathbf{m} + \mathbf{V} \mathbf{x}'$$

- ▶ with
 - ▶ $\mathbf{V}^{(l)}$ non negative
 - ▶ $\mathbf{x}'^{(l)}$ non negative.

Dictionary

- ▶ (Linear) Model assumption

$$\mathbf{x} \simeq \mathbf{m} + \sum_{l=1}^{d'} \mathbf{x}'^{(l)} \mathbf{V}^{(l)} = \mathbf{m} + \mathbf{V} \mathbf{x}'$$

- ▶ with
 - ▶ $\mathbf{V}^{(l)}$ without constraints
 - ▶ \mathbf{x}' sparse (with a lot of 0)

kernel PCA

- ▶ Linear model assumption

$$\Psi(\mathbf{x} - \mathbf{m}) \simeq \sum_{l=1}^{d'} \mathbf{x}'^{(l)} \mathbf{V}^{(l)} = \mathbf{V} \mathbf{x}'$$

- ▶ with
 - ▶ $\mathbf{V}^{(l)}$ orthonormal
 - ▶ \mathbf{x}'_l without constraints.

Deep Auto Encoder

- ▶ Construct a map Φ with a NN from the space \mathcal{X} into a space \mathcal{X}' of smaller dimension:

$$\begin{aligned}\Phi : \mathcal{X} &\rightarrow \mathcal{X}' \\ x &\mapsto \Phi(x)\end{aligned}$$

- ▶ Construct $\tilde{\Phi}$ with a NN from \mathcal{X}' to \mathcal{X}
- ▶ Control the error between x and its reconstruction $\tilde{\Phi}(\Phi(x))$:

$$\frac{1}{n} \sum_{i=1}^n \|x_i - \tilde{\Phi}(\Phi(x_i))\|^2$$

- ▶ Optimization by gradient descent.
- ▶ NN can be replaced by another parametric function...

Pairwise distances

- ▶ Different point of view: input is distances between feature vectors
- ▶ Use only distances $d(x_i, x_j)$.

Distance Preservation. Construct a map Φ from the space \mathcal{X} into a space \mathcal{X}' of smaller dimension:

$$\begin{aligned}\Phi : \mathcal{X} &\rightarrow \mathcal{X}' \\ x &\mapsto \Phi(x) = x' \\ d(x_i, x_j) &\approx d'(x'_i, x'_j)\end{aligned}$$

such that

A natural criterion is

$$\frac{1}{n^2} \sum_{i=1}^n \sum_{j=1}^n |d(x_i, x_j) - d'(x'_i, x'_j)|^2$$

- ▶ Distance preserving transform
- ▶ Random transform
- ▶ Data-independent!
- ▶ Dimension-independent!

Theorem (Johnson-Lindenstrauss Lemma)

Let $\mathcal{S} \subseteq \mathbb{R}^d$ be a finite set of vectors with cardinality $n \geq 2$ and $W \in \mathbb{R}^{p \times d}$ be a random matrix such that its entries $\{W_{k\ell}\}_{1 \leq k \leq p, 1 \leq \ell \leq d}$ are i.i.d. and distributed according to $\mathcal{N}\left(0, \frac{1}{p}\right)$. For any $(\varepsilon, \delta) \in (0, 1)^2$, if

$$p \geq 16\varepsilon^{-2} \log(n/\sqrt{\delta}),$$

then with probability at least $1 - \delta$ on the random matrix W ,

$$\forall (x_i, x_j) \in \mathcal{S}^2: \quad (1 - \varepsilon) \|x_i - x_j\|^2 \leq \|Wx_i - Wx_j\|^2 \leq (1 + \varepsilon) \|x_i - x_j\|^2.$$

The mapping $x \in \mathbb{R}^d \mapsto Wx \in \mathbb{R}^p$ is called an ε -isometry on \mathcal{S} .

The underlying idea of this approach is that the reduction mapping $x \in \mathbb{R}^d \mapsto Wx \in \mathbb{R}^p$ is an exact isometry “in expectation”:

$$\forall x \in \mathbb{R}^d : \quad \mathbb{E} \left(\|Wx\|^2 \right) = \|x\|^2 .$$

Indeed, since for all $x \in \mathbb{R}^d$ such that $x \neq 0$, $\frac{p\|Wx\|^2}{\|x\|^2} \sim \chi_p^2$, one has

$$\mathbb{E} \left(\|Wx\|^2 \right) = \frac{\|x\|^2}{p} \mathbb{E} \left(\frac{p\|Wx\|^2}{\|x\|^2} \right) = \frac{\|x\|^2}{p} p = \|x\|^2 .$$

Let us remark that it is enough for $\{W_{ij}\}_{\substack{1 \leq i \leq p \\ 1 \leq j \leq d}}$ to be independent with $\mathbb{E}W_{ij} = 0$ and $\mathbb{V}(W_{ij}) = \frac{1}{p}$ (for all $i \in [p]$ and $j \in [d]$) in order to get an exact isometry “in expectation”:

$$\mathbb{E} \left(\|Wx\|^2 \right) = \sum_{i=1}^p \mathbb{E} \left[\left(\sum_{j=1}^d W_{ij} x_j \right)^2 \right] = \sum_{i=1}^p \left(\sum_{j=1}^d x_j^2 \mathbb{V}(W_{ij}) + \left(\sum_{j=1}^d x_j \mathbb{E}W_{ij} \right)^2 \right)$$

Remark. Note that $p \geq 8\varepsilon^{-2} \log(n^2/\delta)$ **does not depend on the original dimension d** This means that we could consider data in infinite-dimensional Hilbert space...!

Remark. In practice, Multi-Dimensional Scaling (MDS) techniques such as Isomap are more used.