# An introduction to PCA

Weekly AI pills

Fabio Brau.

2020-10-16

SSSA, Emerging Digital Technologies, Pisa.

ISTITUTO
DI TECNOLOGIE DELLA
COMUNICAZIONE,
DELL'INFORMAZIONE
E DELLA
PERCEZIONE

Scuola Superiore
Sant'Anna

etis
Real-Time Systems Laboratory

- The aim of Principal Component Analysis
- Derivation
  1. A Geometrical idea
  2. A statistical Derivation
  3. Singolar Value Decomposition
- PCA from Encoder Decoder NN
- Dummy examples

# Geometrical Introduction

## Geometrical Introduction

Let $X \in \mathbb{R}^{N \times n}$ be a dataset of $N$ observation within $n$ variables.

$$
X = \begin{bmatrix} x_1^T \\ \vdots \\ x_N^T \end{bmatrix} = \begin{bmatrix} x^{(1)} & | & \cdots & | & x^{(n)} \end{bmatrix} \tag{1}
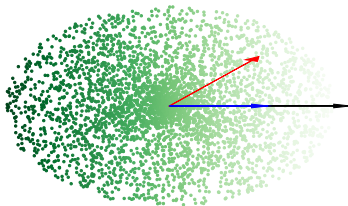$$

Notations:

- $x_i \in \mathbb{R}^n$ represents a single **observation**, i.e a **sample** in the feature space.
- $x^{(i)} \in \mathbb{R}^N$ represents the single **variable**, i.e a **column** of the dataset.
- The object $\mathbb{1}_n \in \mathbb{R}^n$ is the unitary columnar vector of length $n$ $\mathbb{1}_n = [1, \cdots, 1]$.

1. Scalar product measures the projection of $x_j$ along the direction $w$.

2. We are only interested on module.

3. Summation over samples to get the global projection's contribute.

4. Searching for $w$ which maximizes projection.

5. Adding constraint to avoid $w \to \infty$ solution.

$$w_1 \in \operatorname*{argmax}_{\|w\|=1} \sum_{j=1}^{N} \left( w \cdot x_j \right)^2$$

1. Scalar product measures the projection of $x_j$ along the direction $w$.

2. We are only interested on module.

3. Summation over samples to get the global projection's contribute.

4. Searching for $w$ which maximizes projection.

5. Adding constraint to avoid $w \to \infty$ solution.

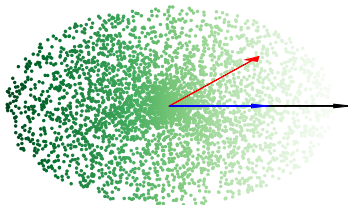$$w_1 \in \underset{\|w\|_2=1}{\arg\max} \sum_{j=1}^{N} \left( w \cdot x_j \right)^2$$

# Geometrical Introduction: Finding a principal direction.

1. Scalar product measures the projection of $x_j$ along the direction $w$.

2. We are only interested on module.

3. Summation over samples to get the global projection's contribute.

4. Searching for $w$ which maximizes projection.
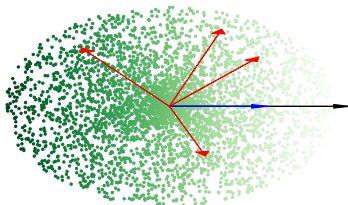
5. Adding constraint to avoid $w \to \infty$ solution.

$$w_1 \in \underset{\|w\|=1}{\mathrm{argmax}} \sum_{j=1}^{N} \left( w \cdot x_j \right)^2$$

1. Scalar product measures the projection of $x_j$ along the direction $w$.

2. We are only interested on module.

3. Summation over samples to get the global projection's contribute.

4. Searching for $w$ which maximizes projection.
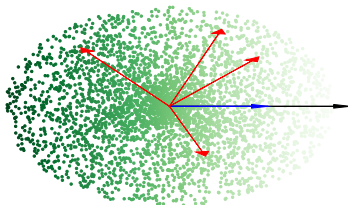
5. Adding constraint to avoid $w \to \infty$ solution.

$$w_1 \in \operatorname*{argmax}_{\|w\|_2 = 1} \sum_{j=1}^{N} \left( w \cdot x_j \right)^2$$

# Geometrical Introduction: Finding a principal direction.

1. Scalar product measures the projection of $x_j$ along the direction $w$.

2. We are only interested on module.

3. Summation over samples to get the global projection's contribute.

4. Searching for $w$ which maximizes projection.

5. Adding constraint to avoid $w \to \infty$ solution.

$$w_1 \in \operatorname*{argmax}_{\|w\|_2 = 1} \sum_{j=1}^{N} \left( w \cdot x_j \right)^2$$

## Geometrical Introduction: Finding other directions

We search for other orthogonal directions which maximize projections.

$$w_1 \in \underset{\|w\|_2=1}{\mathrm{argmax}} \sum_{j=1}^{N} (w \cdot x)^2$$

$$w_2 \in \underset{\|w\|_2=1}{\mathrm{argmax}} \sum_{j=1}^{N} (w \cdot x)^2 \quad \text{and} \quad w_2 \perp w_1$$

$$\vdots$$

$$w_n \in \underset{\|w\|_2=1}{\mathrm{argmax}} \sum_{j=1}^{N} (w \cdot x)^2 \quad \text{and} \quad w_2 \perp \{w_1, \ldots, w_{n-1}\}$$

Example

$$V(w) = \sum_j (w \cdot x_j)^2 \qquad \text{momentum along } w$$

If $w_1$, $w_2$, $w_3$ orthogonal that maximizes $V$ in the 3D example, then

1. $V(w_1) = 3181.20$ $\approx 82.5\%$
2. $V(w_2) = 646.25$ $\approx 17.0\%$
3. $V(w_3) = 19.23$ $\approx 0.5\%$

What if we forget the last direction?

Observation

- $x_j = \alpha_{1j}w_1 + \alpha_{2j}w_2 + \alpha_{3j}w_3$ (where $\alpha_{ij} = w_i \cdot x_j$).
- $\tilde{x}_j = \alpha_{1j}w_1 + \alpha_{2j}w_2$.

$$\frac{1}{N}\sum_j \|x_j - \tilde{x}_j\|^2 = \frac{V(w_3)}{N} \approx 4.8 \cdot 10^{-3} \qquad \text{(MSE)}$$

$$V(w) = \sum_j (w \cdot x_j)^2 \qquad \text{momentum along } w$$

If $w_1$, $w_2$, $w_3$ orthogonal that maximizes $V$ in the 3D example, then

1. $V(w_1) = 3181.20$ $\approx 82.5\%$
2. $V(w_2) = 646.25$ $\approx 17.0\%$
3. $V(w_3) = 19.23$ $\approx 0.5\ \%$

What if we forget the last direction?

Observation

- $x_j = \alpha_{1j} w_1 + \alpha_{2j} w_2 + \alpha_{3j} w_3$ (where $\alpha_{ij} = w_i \cdot x_j$).
- $\tilde{x}_j = \alpha_{1j} w_1 + \alpha_{2j} w_2$.

$$\frac{1}{N} \sum_j \|x_j - \tilde{x}_j\|^2 = \frac{V(w_3)}{N} \approx 4.8 \, 10^{-3} \qquad \text{(MSE)}$$

5

$$V(w) = \sum_j (w \cdot x_j)^2 \qquad \text{momentum along } w$$

If $w_1$, $w_2$, $w_3$ orthogonal that maximizes $V$ in the 3D example, then

1. $V(w_1) = 3181.20$           ≈82.5%
2. $V(w_2) = 646.25$           ≈17.0%
3. $V(w_3) = 19.23$           ≈0.5 %

What if we forget the last direction?

Observation

- $x_j = \alpha_{1j} w_1 + \alpha_{2j} w_2 + \alpha_{3j} w_3$ (where $\alpha_{ij} = w_i \cdot x_j$).
- $\tilde{x}_j = \alpha_{1j} w_1 + \alpha_{2j} w_2$.

$$\frac{1}{N} \sum_j \| x_j - \tilde{x}_j \|^2 = \frac{V(w_3)}{N} \approx 4.8 \cdot 10^{-3} \qquad \text{(MSE)}$$

5

$$V(w) = \sum_j (w \cdot x_j)^2 \qquad \text{momentum along } w$$

If $w_1$, $w_2$, $w_3$ orthogonal that maximizes $V$ in the 3D example, then

1. $V(w_1) = 3181.20$                                                 $\approx 82.5\%$
2. $V(w_2) = 646.25$                                               $\approx 17.0\%$
3. $V(w_3) = 19.23$                                               $\approx 0.5\ \%$

<div align="center">What if we forget the last direction?</div>

**Observation**

- $x_j = \alpha_{1j} w_1 + \alpha_{2j} w_2 + \alpha_{3j} w_3$ (where $\alpha_{ij} = w_i \cdot x_j$).

- $\tilde{x}_j = \alpha_{1j} w_1 + \alpha_{2j} w_2$.

$$\frac{1}{N} \sum_j \| x_j - \tilde{x}_j \|^2 = \frac{V(w_3)}{N} \approx 4.8 \, 10^{-3} \qquad \text{(MSE)}$$

$$V(w) = \sum_j (w \cdot x_j)^2 \qquad \text{momentum along } w$$

If $w_1$, $w_2$, $w_3$ orthogonal that maximizes $V$ in the 3D example, then

1. $V(w_1) = 3181.20$                                            $\approx 82.5\%$
2. $V(w_2) = 646.25$                                            $\approx 17.0\%$
3. $V(w_3) = 19.23$                                            $\approx 0.5\ \%$

<div align="center">What if we forget the last direction?</div>

### Observation

- $x_j = \alpha_{1j}w_1 + \alpha_{2j}w_2 + \alpha_{3j}w_3$ (where $\alpha_{ij} = w_i \cdot x_j$).
- $\tilde{x}_j = \alpha_{1j}w_1 + \alpha_{2j}w_2$.

$$\frac{1}{N} \sum_j \|x_j - \tilde{x}_j\|^2 = \frac{V(w_3)}{N} \approx 4.8\,10^{-3} \qquad \text{(MSE)}$$

- Given a set of data $X \in \mathbb{R}^{N \times n}$
- We can find $w_1, \cdots, w_n$ principal (orthonormal) directions the maximize their momentum.
- $V(w_1) > V(w_2) > \cdots > V(w_n)$
- Approximating X with $\tilde{X}$ by taking only the first $k$ directions we are getting an error that is $V(w_{k+1})/N$

What's the catch?

$$
\begin{aligned}
\max_{w \in \mathbb{R}^n} \quad & \sum_{j=1}^{N} (w \cdot x_j)^2 \\
\text{s.t} \quad & w_i \cdot w = 0, \ \forall i < k \\
& w \cdot w = 1
\end{aligned}
\tag{MP}
$$

6

- Given a set of data $X \in \mathbb{R}^{N \times n}$
- We can find $w_1, \cdots, w_n$ principal (orthonormal) directions the maximize their momentum.
- $V(w_1) > V(w_2) > \cdots > V(w_n)$
- Approximating X with $\tilde{X}$ by taking only the first $k$ directions we are getting an error that is $V(w_{k+1})/N$

### What's the catch?

$$\max_{w \in \mathbb{R}^n} \quad \sum_{j=1}^{N} (w \cdot x_j)^2$$
$$\text{s.t} \quad w_i \cdot w = 0, \; \forall i < k \tag{MP}$$
$$w \cdot w = 1$$

## Geometrical Introduction: Conclusion

- Given a set of data $X \in \mathbb{R}^{N \times n}$
- We can find $w_1, \cdots, w_n$ principal (orthonormal) directions the maximize their momentum.
- $V(w_1) > V(w_2) > \cdots > V(w_n)$
- Approximating X with $\tilde{X}$ by taking only the first $k$ directions we are getting an error that is $V(w_{k+1})/N$

### What's the catch?

$$
\begin{aligned}
\max_{w \in \mathbb{R}^n} \quad & \sum_{j=1}^{N} (w \cdot x_j)^2 \\
\text{s.t} \quad & w_i \cdot w = 0, \ \forall i < k \\
& w \cdot w = 1
\end{aligned}
\tag{MP}
$$

# Statistical Derivation