



Insper

Machine Learning

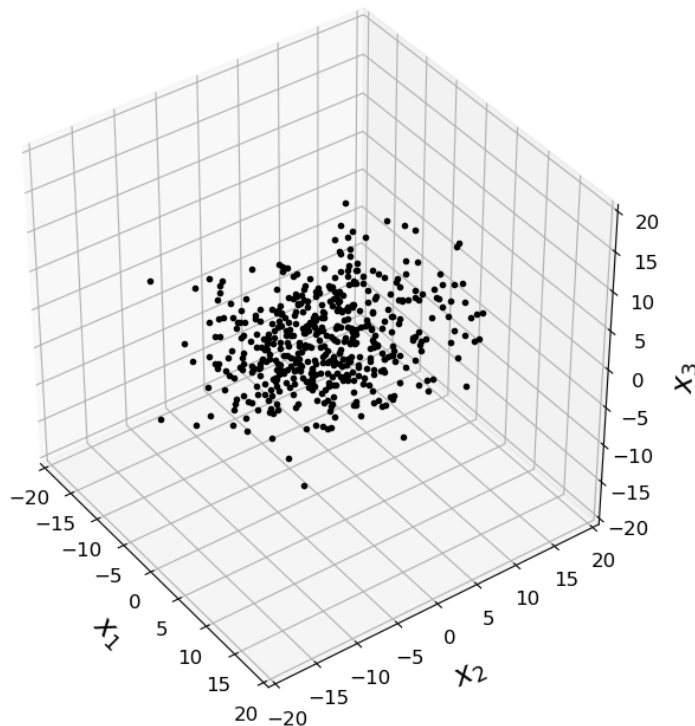
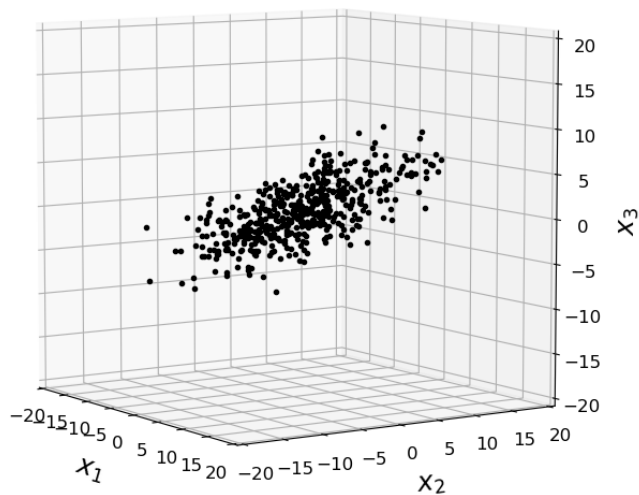
Aula 17 – Redução de dimensionalidade

2021 – Engenharia
Fábio Ayres <fabioja@insper.edu.br>

Redução de dimensionalidade

- Como visualizar um dataset de dimensão $n = 200$?
- Será que em um dataset de dimensão $n = 200$ realmente temos tudo isso de informação independente?
- Como “enxugar” um dataset?

Exemplo



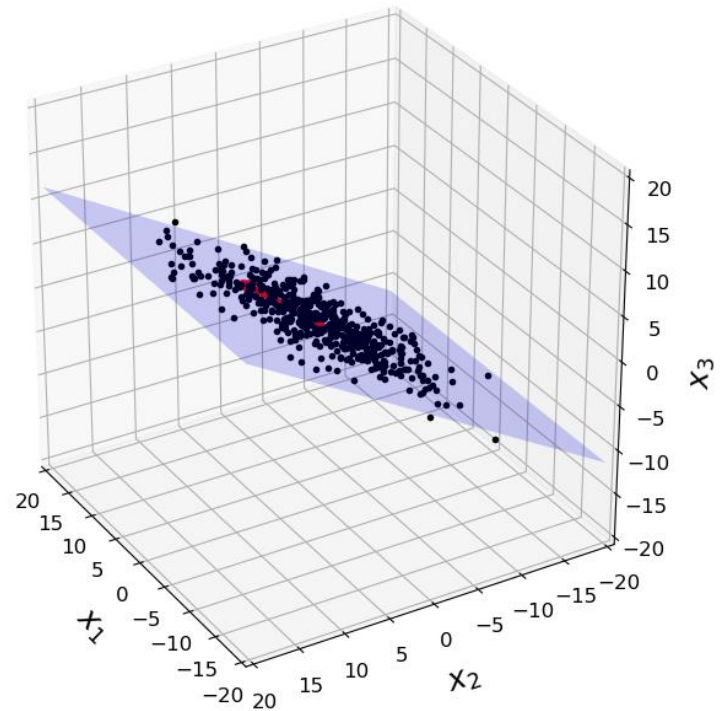
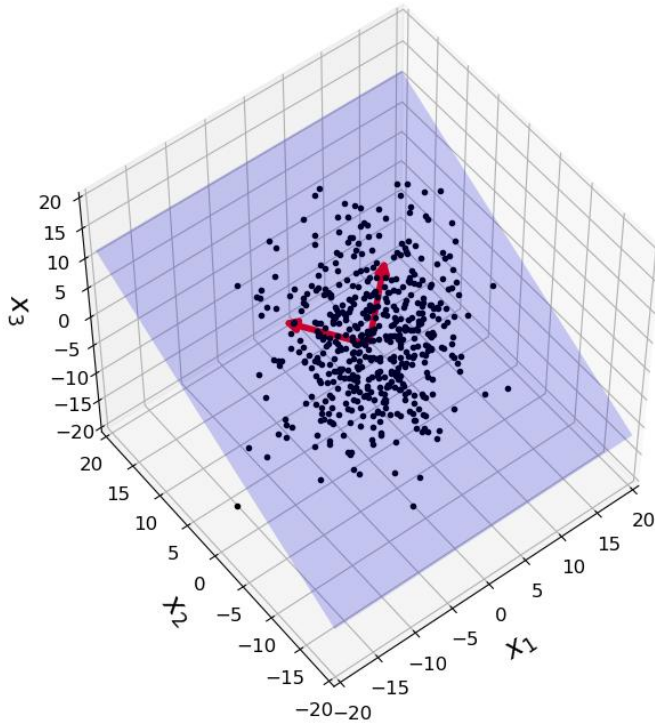
Ideia: projeção

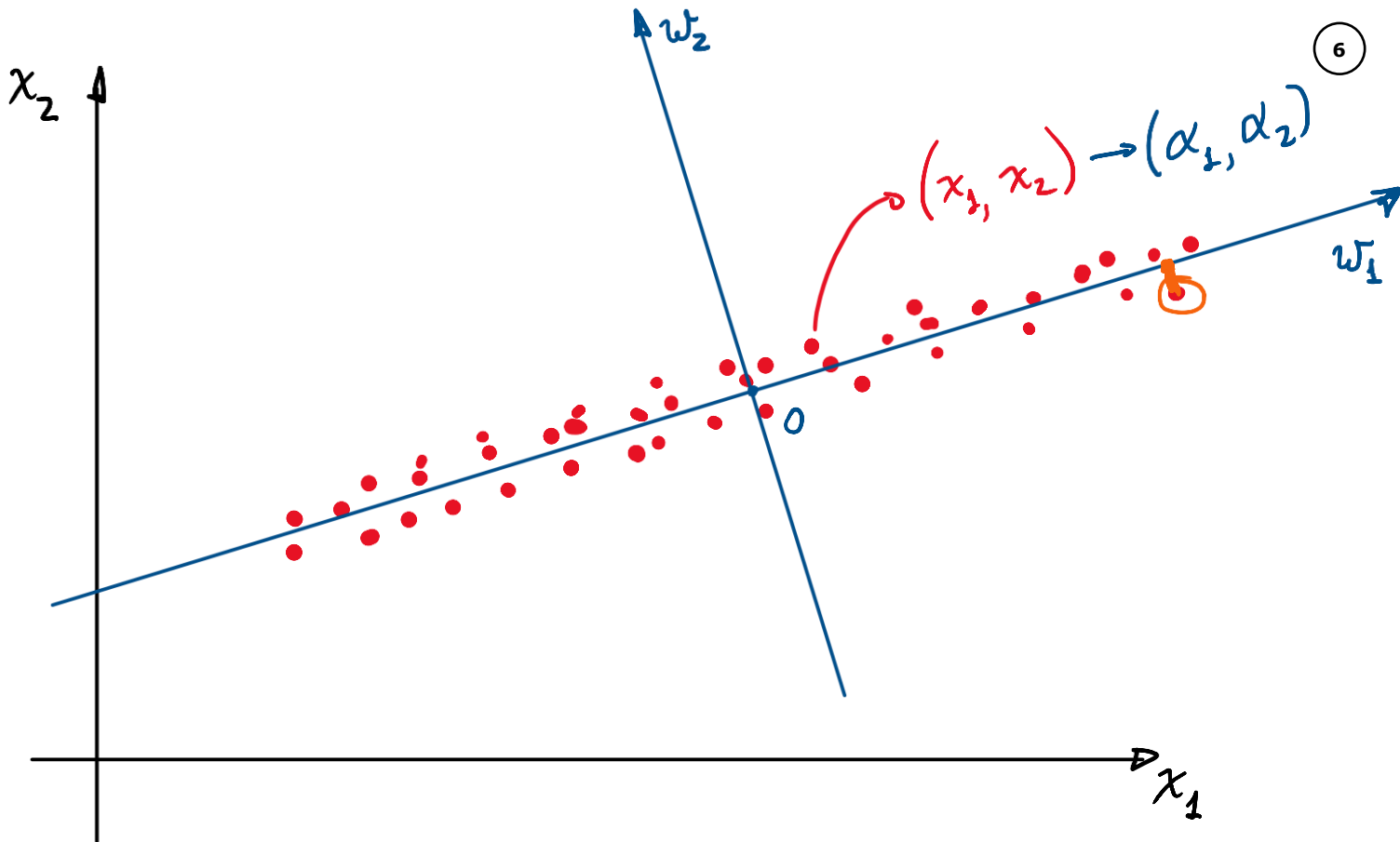
Podemos projetar esses dados em um hiperplano de menor dimensão!

Assim só precisamos guardar:

- A origem do plano
- Os vetores determinando o hiperplano
- As coordenadas de projeção por ponto: $d \ll n$

Exemplo





6

Aproximações sucessivas

- $d = 0$: não guardo nenhuma coordenada por ponto!
 - Guardo apenas o centroide da nuvem!
- $d > 0$: começo a gerar componentes de projeção por ponto.

i) guardar 1 vez por dataset

$(\alpha_1^{(i)}, \alpha_2^{(i)}, \dots, \alpha_d^{(i)}) \rightarrow d \ll n$ componentes

P, w_1, \dots, w_d

vectores no espaço original

$\hat{x}_i = p + \alpha_{i1}w_1 + \alpha_{i2}w_2 + \dots + \alpha_{id}w_d$

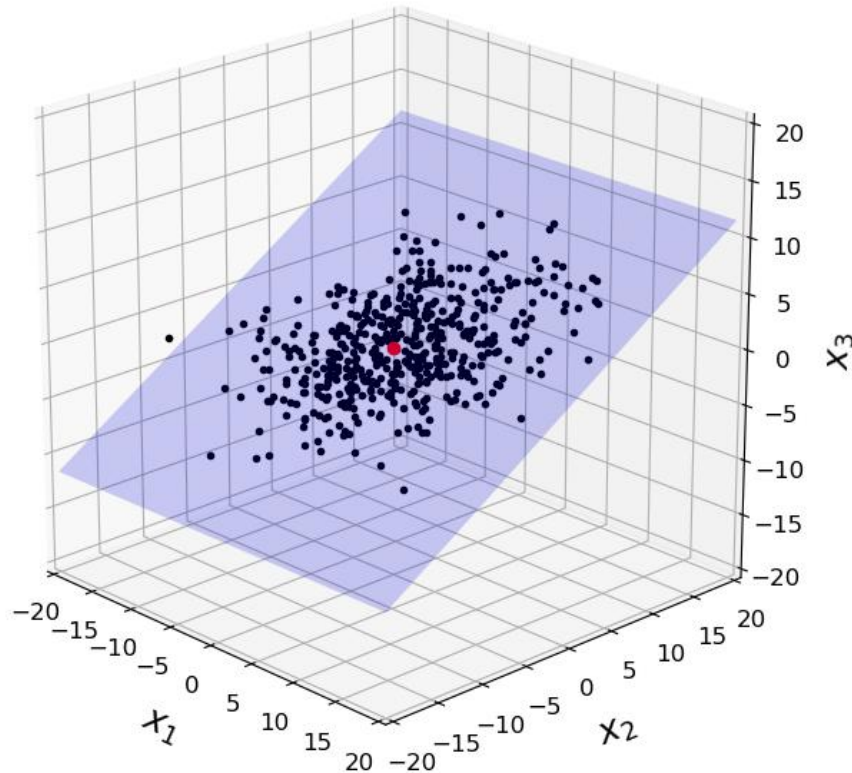
↑ origem ↑ direções.

aproximação de $x_i = (x_1^{(i)}, x_2^{(i)}, \dots, x_n^{(i)})$

n componentes

www.insper.edu.br

Zero-ésima aproximação



Zero-ésima aproximação

- Vamos chamar de \hat{x}_i a aproximação do ponto x_i
- Como temos uma aproximação de um ponto só, $\hat{x}_i = p$ para algum ponto p . Que ponto é esse?

Aproximar cada ponto x_i por $\hat{x}_i = p$
↳ todo mundo aproximado pelo mesmo ponto (!!!)

10

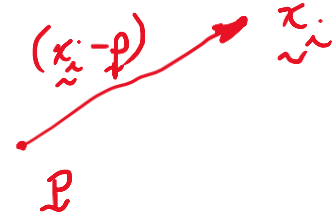
$$MSE = \frac{1}{m} \sum_{i=1}^m \| \underset{\sim}{x}_i - \underset{\sim}{\hat{x}}_i \|^2 = \frac{1}{m} \sum_{i=1}^m \| \underset{\sim}{x}_i - \underset{\sim}{p} \|^2$$

Qual o $\underset{\sim}{p}$ que minimiza o MSE?

$$\underset{\sim}{x}_i - \underset{\sim}{p} = \begin{bmatrix} x_{i1} \\ x_{i2} \\ \vdots \\ x_{in} \end{bmatrix} - \begin{bmatrix} p_1 \\ p_2 \\ \vdots \\ p_n \end{bmatrix}$$

$\underset{\sim}{x}_i$
(ponto original)

$\underset{\sim}{p}$
(aproximação)



$$\|x_i - p\| = [(x_{i1} - p_1)^2 + (x_{i2} - p_2)^2 + \dots + (x_{in} - p_n)^2]^{1/2}$$

$$\|x_i - p\|^2 = (x_{i1} - p_1)^2 + (x_{i2} - p_2)^2 + \dots + (x_{in} - p_n)^2$$

$$= (x_i - p)^T (x_i - p)$$

$$\begin{bmatrix} (x_{i1} - p_1), & (x_{i2} - p_2), & \dots, & (x_{in} - p_n) \end{bmatrix}$$

matriz linha
 $1 \times n$

$$\begin{bmatrix} (x_{i1} - p_1) \\ (x_{i2} - p_2) \\ \vdots \\ (x_{in} - p_n) \end{bmatrix}$$

matriz coluna
 $n \times 1$

$$\|x_i - p\|^2 = (x_i - p)^T (x_i - p) = x_i^T x_i - \underbrace{p^T x_i - x_i^T p}_{-2p^T x_i} + p^T p$$

$$\begin{aligned} \text{MSE} &= \frac{1}{m} \sum_{i=1}^m \|x_i - p\|^2 = \frac{1}{m} \sum_{i=1}^m (x_i^T x_i - 2p^T x_i + p^T p) \\ &= \frac{1}{m} \sum_{i=1}^m x_i^T x_i - 2p^T \cdot \frac{1}{m} \sum_{i=1}^m x_i + p^T p \end{aligned}$$

$$\frac{\partial \text{MSE}}{\partial p} = \frac{\partial}{\partial p} \left(\cancel{\frac{1}{m} \sum_{i=1}^m x_i^T x_i} - \underbrace{2p^T \cdot \frac{1}{m} \sum_{i=1}^m x_i}_{-2 \cdot \frac{1}{m} \sum_{i=1}^m x_i} + \underbrace{p^T p}_{2p} \right) = 0$$

$\xrightarrow{0}$
 $\underbrace{\frac{1}{m} \sum_{i=1}^m x_i}_{\bar{x}}$

$$\Rightarrow 2p - 2\bar{x} = 0 \Rightarrow p = \bar{x}$$

Zero-ésima aproximação

$$\begin{aligned}MSE &= \frac{1}{m} \sum_{i=1}^m \|x_i - p\|^2 = \frac{1}{m} \sum_{i=1}^m (x_i - p)^T (x_i - p) \\&= \frac{1}{m} \sum_{i=1}^m (x_i^T x_i - 2p^T x_i + p^T p)\end{aligned}$$

$$\frac{\partial}{\partial p} MSE = \frac{1}{m} \sum_{i=1}^m -2x_i + 2p = 0$$

$$\Rightarrow p = \frac{1}{m} \sum_{i=1}^m x_i$$

O ponto ótimo
é o centroide!

Próximas aproximações

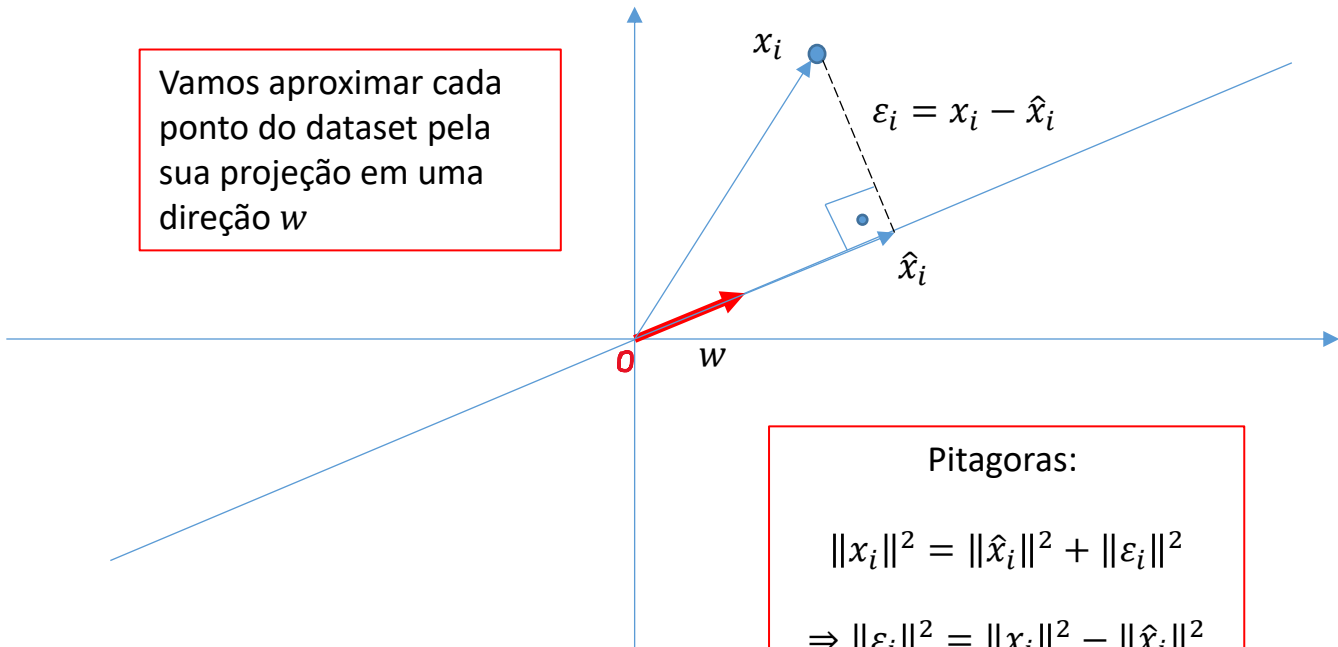
- Agora vamos descobrir a próxima aproximação:
 $d = 1$
- Vamos remover a zero-ésima aproximação do dataset:

$$x_i \leftarrow (x_i - p)$$

- Nosso dataset (remanescente) tem uma nuvem de pontos cujo centróide é a origem

Próximas aproximações

Vamos aproximar cada ponto do dataset pela sua projeção em uma direção w



Pitagoras:

$$\|x_i\|^2 = \|\hat{x}_i\|^2 + \|\epsilon_i\|^2$$

$$\Rightarrow \|\epsilon_i\|^2 = \|x_i\|^2 - \|\hat{x}_i\|^2$$

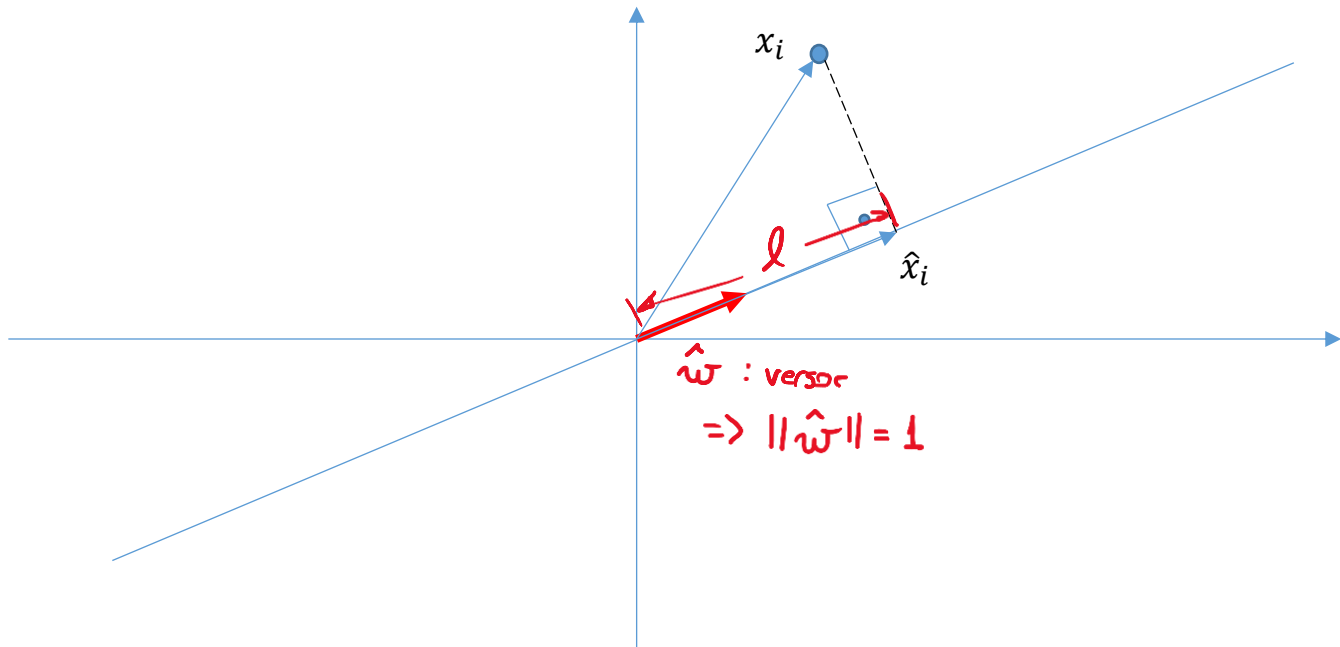
Próximas aproximações

$$\begin{aligned}
 MSE &= \frac{1}{m} \sum_{i=1}^m \|\varepsilon_i\|^2 = \frac{1}{m} \sum_{i=1}^m (\|x_i\|^2 - \|\hat{x}_i\|^2) \\
 &= \frac{1}{m} \sum_{i=1}^m \|x_i\|^2 - \frac{1}{m} \sum_{i=1}^m \|\hat{x}_i\|^2
 \end{aligned}$$

Portanto,

$$w_1 = \arg \min_w MSE = \arg \max_w \underbrace{\frac{1}{m} \sum_{i=1}^m \|\hat{x}_i\|^2}_{\mathfrak{M}(w)}$$

$$\begin{cases} l = x_i \cdot \hat{w} \\ \hat{x}_i = l \hat{w} \end{cases} \Rightarrow \boxed{\hat{x}_i = (x_i \cdot \hat{w}) \hat{w}}$$



$$\begin{cases} \hat{x}_i = (x_i \cdot \hat{w}) \hat{w} \\ \hat{w} = \frac{1}{\|w\|} w \end{cases}$$

$$\Rightarrow \hat{x}_i = \frac{(x_i \cdot w)}{\|w\|} \frac{w}{\|w\|}$$

$\|w\| \cdot \|w\| = \|w\|^2$

$$\begin{cases} (x_i \cdot w) = \underbrace{w^T}_{\text{matrizes}} x_i \\ \uparrow \\ \text{vetores} \end{cases}$$

$$\|w\|^2 = w^T w$$

$$\Rightarrow \hat{x}_i = \frac{w^T x_i}{w^T w} w$$

$$M(w) = \frac{1}{m} \sum_{i=1}^m \|\hat{x}_i\|^2 = \frac{1}{m} \sum_{i=1}^m \hat{x}_i^T x_i$$

$$\boxed{\hat{x}_i = \left(\frac{w^T x_i}{w^T w} \right) w} \Rightarrow M(w) = \frac{1}{m} \sum_{i=1}^m \left(\frac{w^T x_i}{w^T w} \right)^2 w^T w$$

\uparrow número \uparrow mat. coluna

$$= \frac{1}{m} \sum_{i=1}^m \frac{(w^T x_i)^2}{(w^T w)^2} \cdot \cancel{(w^T w)}$$

$$M(w) = \frac{1}{m} \sum_{i=1}^m \frac{(w^T x_i)^2}{w^T w}$$

$$M(w) = \frac{1}{m} \sum_{i=1}^m \frac{(w^T x_i)^2}{w^T w} = \frac{1}{m} \sum_{i=1}^m \frac{(w^T x_i)(x_i^T w)}{w^T w}$$

$$= \frac{1}{m} \frac{1}{w^T w} \sum_{i=1}^m w^T x_i x_i^T w$$

$n \times 1$
 $\begin{bmatrix} \end{bmatrix} \cdot \begin{matrix} 1 \times n \\ \begin{bmatrix} \end{bmatrix} \end{matrix} = \begin{matrix} n \times n \\ \begin{bmatrix} \end{bmatrix} \end{matrix}$

$$\Rightarrow M(w) = \frac{1}{w^T w} \cdot w^T \left(\frac{1}{m} \sum_{i=1}^m x_i x_i^T \right) w$$

$$C = \frac{1}{m} \sum_{i=1}^m x_i x_i^T$$

matriz de
covariância

pto original
menos
media \bar{x}

$$M(w) = \frac{w^T C w}{w^T w}$$

$$\text{maximizar } M(w) = \frac{w^T C w}{w^T w} = \underbrace{\frac{w^T}{\|w\|}}_{\hat{w}^T} \cdot C \cdot \underbrace{\frac{w}{\|w\|}}_{\hat{w}}$$

$$\Leftrightarrow \text{maximizar } \hat{w}^T C \hat{w} \quad \leftarrow \text{multiplicadores de Lagrange!}$$

sujeito a $\hat{w}^T \hat{w} = 1$

$$\Leftrightarrow \text{maximizar } L = \hat{w}^T C \hat{w} - \lambda (\hat{w}^T \hat{w} - 1)$$

$$\frac{\partial L}{\partial w} = 2C\hat{w} - 2\lambda\hat{w} = 0 \Rightarrow \boxed{C\hat{w} = \lambda\hat{w}} : \hat{w} \text{ é um autovetor de } C, \text{ com autovalor } \lambda$$

$$\begin{cases} L = \hat{\mathbf{w}}^T \mathbf{C} \hat{\mathbf{w}} - \lambda (\hat{\mathbf{w}}^T \hat{\mathbf{w}} - 1) \\ \mathbf{C} \hat{\mathbf{w}} = \lambda \hat{\mathbf{w}} \end{cases}$$

$$\begin{aligned} L &= \hat{\mathbf{w}}^T (\lambda \hat{\mathbf{w}}) - \lambda \hat{\mathbf{w}}^T \hat{\mathbf{w}} + \lambda \\ &= \lambda \hat{\mathbf{w}}^T \hat{\mathbf{w}} - \lambda \hat{\mathbf{w}}^T \hat{\mathbf{w}} + \lambda \Rightarrow L = \lambda \end{aligned}$$

$$\text{maximize } L = \hat{\mathbf{w}}^T \mathbf{C} \hat{\mathbf{w}} - \lambda (\hat{\mathbf{w}}^T \hat{\mathbf{w}} - 1)$$



$\hat{\mathbf{w}}$ é autovetor de \mathbf{C} (matriz de covariância dos dados)
associado ao maior autovalor de \mathbf{C}

Maximizando as projeções

\hat{x}_i = projeção de x_i na direção w

$$\Rightarrow \hat{x}_i = \left(\frac{w^T x_i}{w^T w} \right) w$$

$$M(w) = \frac{1}{m} \sum_{i=1}^m \|\hat{x}_i\|^2 = \frac{1}{m} \sum_{i=1}^m \hat{x}_i^T \hat{x}_i$$

$$= \frac{1}{m} \sum_{i=1}^m \left(\frac{w^T x_i}{w^T w} \right)^2 w^T w = \frac{1}{m} \sum_{i=1}^m \frac{(w^T x_i)^2}{w^T w}$$

Maximizando as projeções

$$M(w) = \frac{1}{m} \sum_{i=1}^m \frac{(w^T x_i)^2}{w^T w} = \frac{1}{m} \sum_{i=1}^m \frac{w^T x_i x_i^T w}{w^T w}$$

$$\Rightarrow M(w) = \frac{w^T C w}{w^T w}$$

onde

$$C = \frac{1}{m} \sum_{i=1}^m x_i x_i^T$$

Matriz de covariância

linhas

$$X = \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_m^T \end{bmatrix}$$

colunas

$$X^T = \begin{bmatrix} x_1 & x_2 & \cdots & x_m \end{bmatrix}$$

$$\frac{1}{m} X^T X = \frac{1}{m} \begin{bmatrix} x_1 & x_2 & \cdots & x_m \end{bmatrix} \begin{bmatrix} x_1^T \\ x_2^T \\ \vdots \\ x_m^T \end{bmatrix} = \frac{1}{m} \sum_{i=1}^m x_i x_i^T$$

Logo: $C = \frac{1}{m} X^T X$, que é a matriz de covariância dos dados!

Maximizando as projeções

Portanto, queremos achar w que maximiza

$$M(w) = \frac{w^T C w}{w^T w}$$

E agora?

Autovalores e autovetores

- Imagine o que acontece se w for um autovetor (dentre vários) da matriz C : $Cw = \lambda w$

- Portanto

$$M(w) = \frac{w^T C w}{w^T w} = \lambda \frac{w^T w}{w^T w} = \lambda$$

- Agora temos uma estratégia para maximizar $M(w)$:

$w =$ **autovetor de C** correspondente ao **maior autovalor**

Próximas aproximações?

Repita o processo:

- Remova a aproximação feita até o momento

$$X_k = X - \hat{X}_{k-1}$$

- Calcule a matriz de covariância atual

$$C_k = Cov[X_k] = \frac{1}{m} X_k^T X_k$$

- Calcule o autovetor do maior autovalor da matriz C_k

Principal Component Analysis

Ufa! Se você chegou até aqui, parabéns! Você acabou de derivar o **método das componentes principais**, um dos algoritmos mais importantes da estatística!

Principal Component Analysis

Calcule os parâmetros da PCA:

1. A média \bar{p} das amostras
2. A matriz de covariância C
3. Os d autovetores w_k correspondentes aos maiores autovalores de C

Principal Component Analysis

Agora calcule as **componentes principais** de cada ponto de dados:

1. Remova a média p de cada amostra x_i
2. Projete os pontos resultantes nas direções principais

$$\hat{x}_i = p + \alpha_{i1}w_1 + \alpha_{i2}w_2 + \cdots + \alpha_{id}w_d$$

Principal Component Analysis

- Medida de “desempenho”: soma dos autovalores
 - Representa o quanto da variabilidade original do dataset foi capturada no dataset reduzido
- Aplicações
 - Análise das direções: podem indicar os aspectos mais importantes do dataset
 - Compressão: representar aproximadamente o dataset com menos dados
 - Visualização: projetar em 2D ou 3D para observar como os dados se espalham
 - Redução de computação: treinar modelos mais rapidamente sem perder muita qualidade

Principal Component Analysis

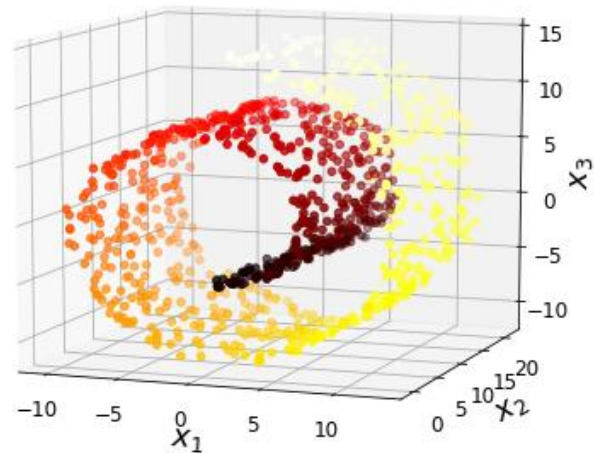
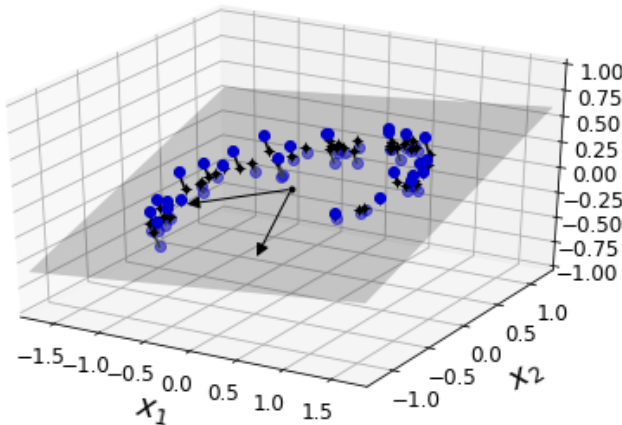
Implementações:

- Decomposição por Valores Singulares (SVD: Singular Value Decomposition)
- PCA incremental
- PCA aleatorizada

Estude mais sobre essas técnicas em seu livro texto

Limitações

- Só projeta os pontos em **hiperplanos**
- E se os pontos estão em outras superfícies?

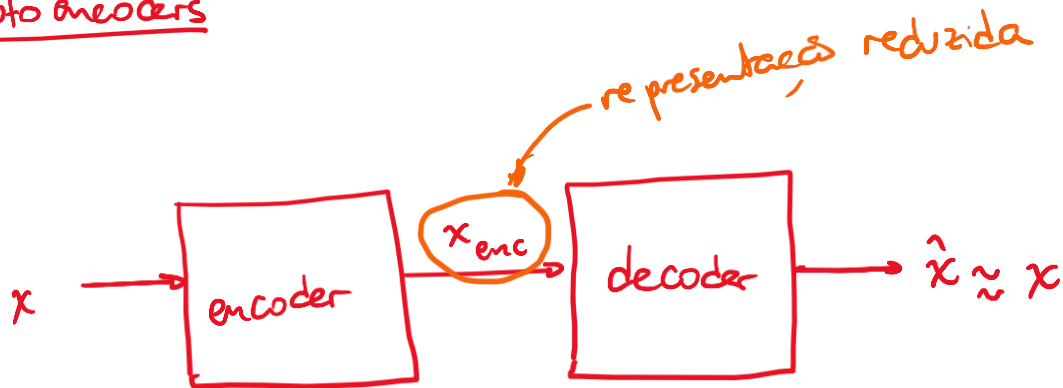


Outras técnicas

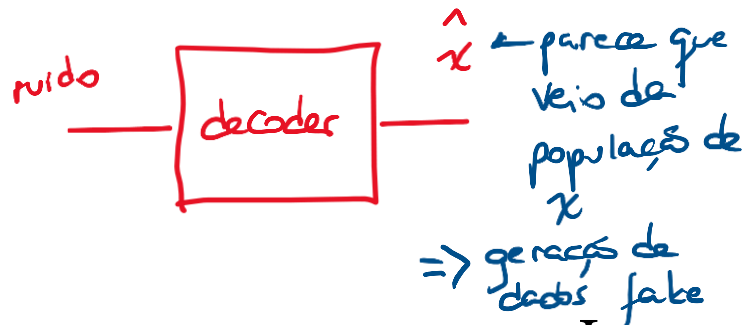
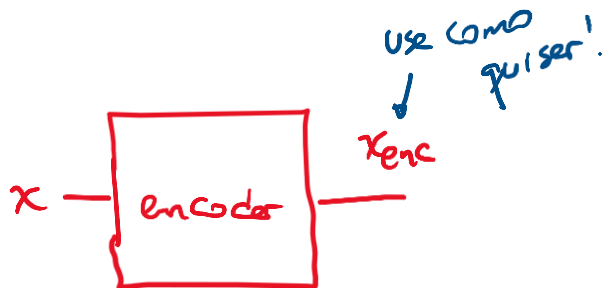
Existem técnicas mais sofisticadas para redução de dimensionalidade, que lidam com a não-linearidade:

- Kernel PCA
- Local Linear Embedding
- t-Distributed Stochastic Neighbor Embedding (t-SNE)

Auto encoders



"entra com x para prever x !!!"



The background of the slide is white and features several concentric, partial arcs in red and grey. These arcs are of varying thicknesses and are scattered across the frame, creating a dynamic, abstract pattern. Some arcs are solid, while others are thin outlines.

Insper