

CS-EJ3211 Machine Learning with Python

Session 6 - Feature Learning

Shamsi Abdurakhmanova

Aalto University
FITech

24.03.22

Feature Learning

- The efficiency of ML methods depends on the choice of features
- Features have to be relevant (useful properties of data points)
- Too many unnecessary features result in:
 - increase in computational resources used
 - overfitting

Feature Learning

Feature Learning - automate the choice of finding good features.

Feature Learning can be:

- supervised - dictionary learning, ANN
- unsupervised - PCA

Learn a hypothesis map that reads in some representation of a datapoint (e.g features) and transforms it to a set of (new) features:

$$\mathbf{z} = (z_1, \dots, z_D) \rightarrow \mathbf{x} = (x_1, \dots, x_n) ,$$

Dimensionality reduction

Transform high dimensional (many features) dataset Z to a dataset X with lower dimensionality.



Dimensionality reduction

Why?

- reduce the use of computational resources
- prevent overfitting
- data visualization

How?

- How to map datapoint to a space with lower dimensionality?
- How to quantify the information loss after dataset transformation?

Dimensionality reduction - Compression

The goal is to find (learn) a compression map:

$$h(\cdot) : \mathbb{R}^D \rightarrow \mathbb{R}^n,$$

that transforms a long feature vector $\mathbf{z} \in \mathbb{R}^D$ to a short feature vector $h(\mathbf{z}) = \mathbf{x} \in \mathbb{R}^n$ ($D \gg n$).

The new feature vector $\mathbf{x} = h(\mathbf{z})$ is a compressed representation (code) of the original feature vector \mathbf{z} .

Dimensionality reduction - Reconstruction

We can reconstruct original vector using a reconstruction map:

$$r(\cdot) : \mathbb{R}^n \rightarrow \mathbb{R}^D$$

Reconstructed original feature vector: $\hat{\mathbf{z}} = r(\mathbf{x})$

The difference (reconstruction error): $\hat{\mathbf{z}} - \mathbf{z}$

Need to learn a compression map such that: $\hat{\mathbf{z}} \approx \mathbf{z}$, i.e. reduce dimensionality with minimal information loss.

Principal Component Analyses

PCA is a dimensionality reduction method where compression map $h(\cdot)$ is a **linear** map.

PCA is a statistical procedure that uses an orthogonal transformation to convert a set of observations of possibly correlated variables (entities each of which takes on various numerical values) into a set of values of linearly uncorrelated variables called principal components.

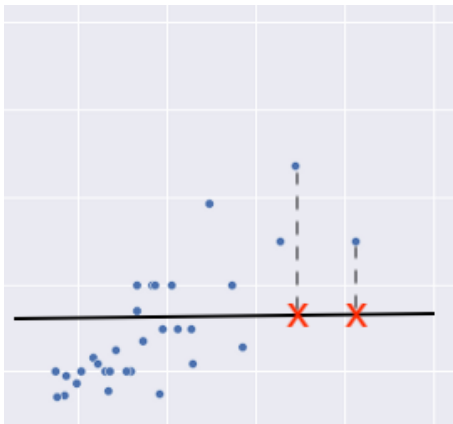
Linear Transformation



3Blue1Brown

Principal Component Analyses

How to map datapoint to a space with lower dimensionality?



PCA, gif

Principal Component Analyses

Linear transformation: $h(\cdot) : \mathbb{R}^2 \rightarrow \mathbb{R}^1$:

$\mathbf{z} = (z_1, z_2) \rightarrow \mathbf{x} = (x_1)$

$\mathbf{x} = W\mathbf{z}$

Reconstruction: $\hat{\mathbf{z}} = W^T \mathbf{x}$

How to quantify the information loss after dataset transformation?

Information loss is measured with **reconstruction error**:

$$(1/m) \sum_{i=1}^m \|\mathbf{z}^{(i)} - \hat{\mathbf{z}}^{(i)}\|_2^2$$

Reconstruction error, gif.
PC1, gif.
Source.

Principal Component Analyses

How to interpret PC components?

Cocktail recipe ($PC1 = x$ amount of feature 1 + y amount of feature 2),
i.e. linear combination of features.

Check out PCA by StatQuest for more explanations.

Principal Component Analyses

$$\begin{array}{ccc} \mathbf{Z} & & \mathbf{X} & & \hat{\mathbf{Z}} \\ \begin{bmatrix} z_1^{(1)} & z_2^{(1)} & \dots & z_D^{(1)} \\ z_1^{(2)} & z_2^{(2)} & \dots & z_D^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ z_1^{(m)} & z_2^{(m)} & \dots & z_D^{(m)} \end{bmatrix} & \xrightarrow[\text{.transform(Z)}]{\mathbf{X} = (\mathbf{Z} - \bar{\mathbf{Z}})\mathbf{W}^T} & \begin{bmatrix} x_1^{(1)} & x_2^{(1)} & \dots & x_n^{(1)} \\ x_1^{(2)} & x_2^{(2)} & \dots & x_n^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ x_1^{(m)} & x_2^{(m)} & \dots & x_n^{(m)} \end{bmatrix} & \xrightarrow[\text{.inverse_transform(X)}]{\hat{\mathbf{Z}} = \mathbf{X}\mathbf{W} + \bar{\mathbf{Z}}} & \begin{bmatrix} \hat{z}_1^{(1)} & \hat{z}_2^{(1)} & \dots & \hat{z}_D^{(1)} \\ \hat{z}_1^{(2)} & \hat{z}_2^{(2)} & \dots & \hat{z}_D^{(2)} \\ \vdots & \vdots & \ddots & \vdots \\ \hat{z}_1^{(m)} & \hat{z}_2^{(m)} & \dots & \hat{z}_D^{(m)} \end{bmatrix} \\ (m \times D) & & (m \times n) & & (m \times D) \end{array}$$