



Dimension Reduction Algorithms

Carlos J. Costa





Learning Goals

- Know concept of Dimension Reduction Algorithms
- Distinguish between Feature Extraction and Feature Selection
- Distinguish between main algorithms
- Apply algorithms by using python libraries



Summary

- Dimension reduction
- Feature Extraction vs Feature Selection
- Feature Extraction (PCA, LDA, NMF, TSVD)
- Feature Selection



Dimension Reduction Algorithms

dimensionality reduction

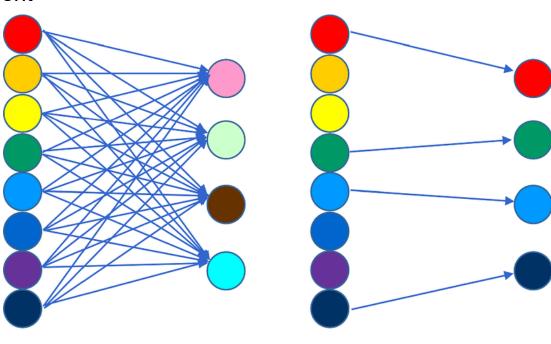
seek and exploit the inherent

structure in the data

Unsupervised learning

Feature extraction

Feature selection



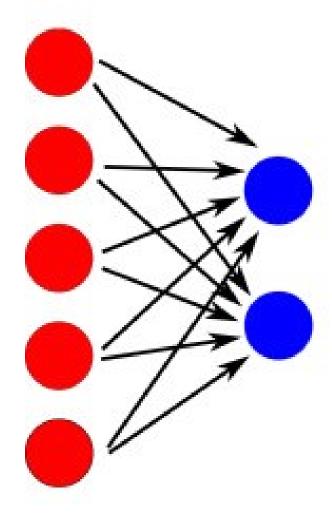
feature extraction

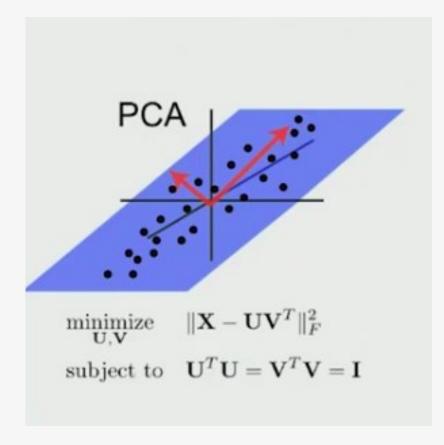
feature selection

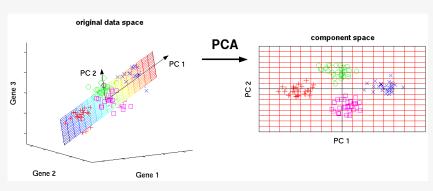
Dimension Reduction Algorithms

- Feature Extraction
 - PCA (principal Components analysis)
 - LDA (Linear Discriminant Analysis)
 - NMF (Non-negative Matrix Factorization)
 - TSVD (Truncate Singular Value Decomposition)

Feature Extraction







PCA

- Principal Component Analysis
- Data is first centered around its mean
- then finding the eigenvectors and eigenvalues of the covariance matrix.
- The eigenvectors represent the directions of maximum variance, while the eigenvalues represent the amount of variance explained by each eigenvector.
- The eigenvectors are then used to project the data onto a lower-dimensional space.
- The number of principal components to keep is determined by the amount of variance we want to retain.

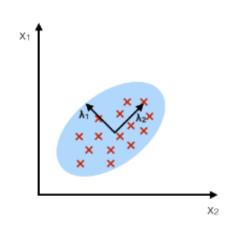


PCA

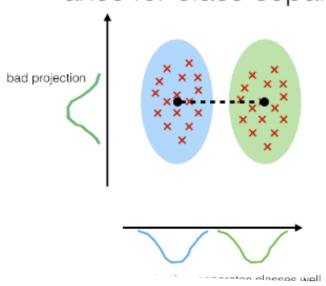
```
# Load libraries
 2 from sklearn import datasets
 3 from sklearn.decomposition import PCA
 1 # Load the Iris flower dataset:
 2 iris = datasets.load iris()
 3 X = iris.data
 4 y = iris.target
 1 # Create an PCA that will reduce the data down to 2 feature
 2 PCAModel = PCA(n components=2)
 3
 4 # run an PCA and use it to transform the features
 5 XPCA = PCAModel.fit(X).transform(X)
 1 # Print the number of features
 2 print('Original number of features:', X.shape[1])
 3 print('Reduced number of features:', XPCA.shape[1])
 Original number of features: 4
 Reduced number of features: 2
 1 ## View the ratio of explained variance
 2 PCAModel.explained variance ratio
array([0.92461621, 0.05301557])
```



PCA: component axes that maximize the variance



LDA: maximizing the component axes for class-separation



- Linear Discriminant Analysis (LDA)
- Supervised Learning

LDA

- Is a linear transformation techniques that is commonly used for dimensionality reduction (like PCA)
- Reducing features by maximizing class separation



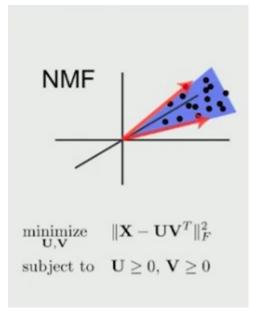
LDA

```
1 # Load libraries
 2 from sklearn import datasets
 3 from sklearn.discriminant analysis import LinearDiscriminantAnalysis
 1 # Load the Tris flower dataset:
 2 iris = datasets.load iris()
 3 X = iris.data
 4 y = iris.target
 1 # Create an LDA that will reduce the data down to 1 feature
 2 | ldaModel = LinearDiscriminantAnalysis(n components=2)
 4 # run an LDA and use it to transform the features
 5 | XLda = ldaModel.fit(X, y).transform(X)
 1 # Print the number of features
 2 print('Original number of features:', X.shape[1])
 3 print('Reduced number of features:', XLda.shape[1])
 Original number of features: 4
 Reduced number of features: 2
 1 ## View the ratio of explained variance
 2 | ldaModel.explained variance ratio
array([0.99147248, 0.00852752])
```



NMF

- Non-negative Matrix Factorization
- where a matrix is factorized into (usually) two matrices, with the property that all three matrices have no negative elements.
- Performs matrix factorization
- It can be applied for:
 - Recommender Systems,
 - Collaborative Filtering
 - topic modelling
 - · dimensionality reduction.
- Does not provide the explained variance



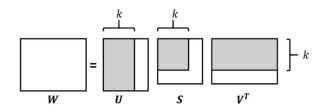
NMF

```
# Load libraries
 from sklearn import datasets
 from sklearn.decomposition import NMF
  # Load the Iris flower dataset:
2 iris = datasets.load iris()
 X = iris.data
  # Create an NMF that will reduce the data down to 2 feature
  NMFModel = NMF(n components=2)
  # run an LDA and use it to transform the features
  XNMF = NMFModel.fit(X).transform(X)
 # Print the number of features
  print('Original number of features:', X.shape[1])
  print('Reduced number of features:', XNMF.shape[1])
```



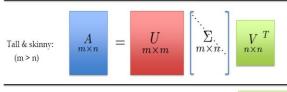
TSVD

- Truncated Singular Value Decomposition
- Used in sparce feature matrix
- Contrary to PCA, this estimator does not center the data before computing the singular value decomposition.



Any real rectangular matrix A can be factored into the form

$$A_{m \times n} = U_{m \times m} \sum_{m \times n} V_{n \times n}^{T}$$



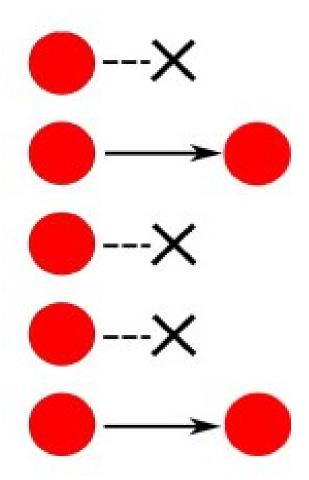
TSVD

```
1 # Load libraries
 2 from sklearn.preprocessing import StandardScaler
 3 from sklearn.decomposition import TruncatedSVD
 4 from scipy.sparse import csr matrix
 5 from sklearn import datasets
 6 import numpy as np
 1 # Load the data
 2 digits = datasets.load digits()
 3 # Standardize the feature matrix
 4 X = StandardScaler().fit transform(digits.data)
    # Make sparse matrix
   X sparse = csr matrix(X)
    # Create a TSVD
    tsvdModel = TruncatedSVD(n components=10)
    # Conduct TSVD on sparse matrix
    X sparse tsvd = tsvdModel.fit(X sparse).transform(X
   # Show results
 2 print('Original number of features:', X sparse.shape
 3 print('Reduced number of features:', X sparse tsvd.s
 Original number of features: 64
 Reduced number of features: 10
 1 # Sum of first three components' explained variance
 2 tsvdModel.explained variance ratio [0:3].sum()
0.3003938538627934
```

Dimension Reduction Algorithms

- Feature Selection
 - Thresholding numerical features variance
 - Thresholding binary features variance
 - Handling high correlated features
 - Removing irrelevant features for Classification
 - RFE (Recursive Feature Elimination)

Feature Selection



Thresholding numerical features variance

The dataset has set of numerical features

Approach:

- Remove those with the low variance
- Low variance likely contains little information



Thresholding numerical features variance

```
1 from sklearn import datasets
 2 from sklearn.feature selection import VarianceThreshold
 1 # Load iris data
 2 iris = datasets.load iris()
 4 # Create features and target
 5 X = iris.data
 6 y = iris.target
 1 # Create VarianceThreshold object with a variance with a
 2 #threshold of 0.5
 3 thresholder = VarianceThreshold(threshold=.5)
 5 # Conduct variance thresholding
 6 XHighVariance = thresholder.fit transform(X)
 1 # View first five rows with features with variances above
 2 # threshold
 3 XHighVariance[0:5]
array([[5.1, 1.4, 0.2],
      [4.9, 1.4, 0.2],
      [4.7, 1.3, 0.2],
      [4.6, 1.5, 0.2],
      [5., 1.4, 0.2]])
```



Handling high correlated features

```
1 # Load libraries
 2 import pandas as pd
 3 import numpy as np
1 # Create feature matrix with two highly correlated features
2 X = np.array([[6, 12, 1],
                [5, 10, 0],
                 [4, 8, 1],
                [3, 3, 0],
                [2, 5, 1],
                [1, 2, 0],
                 [3, 6, 1],
                [5, 10, 0],
                 [9, 19, 1]])
12 # Convert feature matrix into DataFrame
13 df = pd.DataFrame(X)
1 # Create correlation matrix
2 corr matrix = df.corr().abs()
3 # Select upper triangle of correlation matrix
 4 upper = corr matrix.where(np.triu(np.ones(corr matrix.shape), k=1).astype(np.bool))
 5 # Find index of feature columns with correlation greater than 0.95
 6 to drop = [column for column in upper.columns if any(upper[column] > 0.95)]
1 # Drop features
2 df.drop(df[to drop], axis=1)
```



Removing irrelevant features for Classification

Categorical features:

 Calculate Chi-square statistic between each feature and target

Quantitative features:

 Calculate ANOVA F-Value between each feature and target



Recursive Feature Elimination

```
1 # Load libraries
2 from sklearn.datasets import make regression
3 from sklearn.feature selection import RFECV
4 from sklearn import datasets, linear model
5 import warnings
7 # Suppress an annoying but harmless warning
8 | warnings.filterwarnings(action="ignore", module="scipy", message="^internal gelsd")
1 # Generate features matrix, target vector, and the true coefficients
2 X, y = make regression(n samples = 10000,
                         n features = 100,
                         n informative = 2,
                         random state = 1)
1 # Create a linear regression
2 olsModel = linear model.LinearRegression()
1 # Create recursive feature eliminator that scores features by mean squared errors
 rfecvModel = RFECV(estimator=olsModel, step=1, scoring='neg mean squared error')
 # Fit recursive feature eliminator
 rfecvModel.fit(X, y)
7 # Recursive feature elimination
8 rfecvModel.transform(X)
1 # Number of best features
2 rfecvModel.n features
```



Conclusion

- Dimensionality reduction
- Feature Extraction and Feature Selection
- Feature Extraction (PCA, LDA, NMF, TSVD)
- Feature Selection

