# Dimensionality Reduction Algorithms

**CARLOS J. COSTA**

# Dimensional Reduction Algorithms

- dimensionality reduction seek and exploit the inherent structure in the data,

- unsupervised learning

# Dimensional Reduction Algorithms

- Feature Extraction
- Feature Selection

# Dimensional Reduction Algorithms

- Feature Extraction

  - PCA (principal Components analysis)

  - LDA (Linear Discriminant Analysis )

  - NMF (Non-negative Matrix Factorization)

  - TSVD (Truncate Singular Value Decomposition)

# PCA

- Principal Component Analysis, or PCA, is a dimensionality-reduction method

- It is often used to reduce the dimensionality of large data sets

- The purpose is transforming a large set of variables into a smaller

- Containing most of the information in the large set.

- When data is linearly inseparable using PCA extension using kernels

# PCA

- Standardization

- Covariance Matrix computation.

- Compute the eigenvectors and eigenvalues of the covariance matrix to identify the principal components

- Feature Vector

- Recast the Data Along the Principal Components Axes

```python
1  # Load libraries
2  from sklearn import datasets
3  from sklearn.decomposition import PCA
```

```python
1  # Load the Iris flower dataset:
2  iris = datasets.load_iris()
3  X = iris.data
4  y = iris.target
```

```python
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)
```

```python
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
```

```python
1  ## View the ratio of explained variance
2  PCAModel.explained_variance_ratio_
```

```
array([0.92461621, 0.05301557])
```

# LDA

- Linear Discriminant Analysis (LDA)

- Is a linear transformation techniques that is commonly used for dimensionality reduction (like PCA)

- Reducing features by maximizing class separation

```python
1   # Load libraries
2   from sklearn import datasets
3   from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
```

```python
1   # Load the Iris flower dataset:
2   iris = datasets.load_iris()
3   X = iris.data
4   y = iris.target
```

```python
1   # Create an LDA that will reduce the data down to 1 feature
2   ldaModel = LinearDiscriminantAnalysis(n_components=2)
3
4   # run an LDA and use it to transform the features
5   XLda = ldaModel.fit(X, y).transform(X)
```

```python
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
```

```python
1   ## View the ratio of explained variance
2   ldaModel.explained_variance_ratio_
```

```
array([0.99147248, 0.00852752])
```

# NMF

- Non-negative Matrix Factorization

- Performs matrix factorization

- It can be applied for:
  - Recommender Systems,
  - Collaborative Filtering
  - topic modelling
  - dimensionality reduction.

- Does not provides the explained variance

# NMF

```python
# Load libraries
from sklearn import datasets
from sklearn.decomposition import NMF
```

```python
# Load the Iris flower dataset:
iris = datasets.load_iris()
X = iris.data
```

```python
# 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

- Truncate Singular Value Decomposition

- Used in sparce feature matrix

# 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)
5   # Make sparse matrix
6   X_sparse = csr_matrix(X)
```

```
1   # Create a TSVD
2   tsvdModel = TruncatedSVD(n_components=10)
```

```
1   # Conduct TSVD on sparse matrix
2   X_sparse_tsvd = tsvdModel.fit(X_sparse).transform(X_sparse)
```

```
1   # Show results
2   print('Original number of features:', X_sparse.shape[1])
3   print('Reduced number of features:', X_sparse_tsvd.shape[1])
```

```
Original number of features: 64
Reduced number of features: 10
```

```
1   # Sum of first three components' explained variance ratios
2   tsvdModel.explained_variance_ratio_[0:3].sum()
```

```
0.3003938538627934
```

# Dimensional Reduction Algorithms

- Feature Selection

  – Thresholding numerical features variance

  – Thresholding binary features variance

  – Handling high correlated features

  – Removing irrelevant features for Classification

  – RFEC

# 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()
3
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)
4
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

```python
1  # Load libraries
2  import pandas as pd
3  import numpy as np
```

```python
1   # Create feature matrix with two highly correlated features
2   X = np.array([[6, 12, 1],
3                 [5, 10, 0],
4                 [4, 8, 1],
5                 [3, 3, 0],
6                 [2, 5, 1],
7                 [1, 2, 0],
8                 [3, 6, 1],
9                 [5, 10, 0],
10                [9, 19, 1]])
11
12  # Convert feature matrix into DataFrame
13  df = pd.DataFrame(X)
```

```python
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)]
```

```python
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 Eliminating Feature

```python
# Load libraries
from sklearn.datasets import make_regression
from sklearn.feature_selection import RFECV
from sklearn import datasets, linear_model
import warnings

# Suppress an annoying but harmless warning
warnings.filterwarnings(action="ignore", module="scipy", message="^internal gelsd")
```

```python
# Generate features matrix, target vector, and the true coefficients
X, y = make_regression(n_samples = 10000,
                       n_features = 100,
                       n_informative = 2,
                       random_state = 1)
```

```python
# Create a linear regression
olsModel = linear_model.LinearRegression()
```

```python
# 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)

# Recursive feature elimination
rfecvModel.transform(X)
```

```python
# Number of best features
rfecvModel.n_features_
```