

TP 1 – AOS1

PCA

Corrigé

1 Python warm up: PCA by hand

- ① Generate a dataset with the following instruction

```
| X = np.random.multivariate_normal([1, 3], [[2, 1], [1, 2]], 100)
```

How many samples are generated? How many features? What is the underlying distribution of samples in X ?

```
In [1]: | import matplotlib.pyplot as plt
         | import numpy as np
         | import scipy.linalg as linalg
         | X = np.random.multivariate_normal([1, 3], [[2, 1], [1, 2]], 100)
         | X.shape
Out [1]: | (100, 2)
```

There are 100 samples and 2 features. Samples are drawn according to $\mathcal{N}\left(\begin{pmatrix} 1 \\ 3 \end{pmatrix}, \begin{pmatrix} 2 & 1 \\ 1 & 2 \end{pmatrix}\right)$.

- ② Verify the relation that exists between singular values and eigenvalues using a matrix X . To use the functions provided by the `scipy` library, use the following command:

```
| import scipy.linalg as linalg
```

and look at the functions `linalg.eig`, `linalg.eigh`, `linalg.eigvals`, `linalg.eigvalsh`, `linalg.svd` `linalg.svdvals`

```
In [2]: | X = np.random.normal(size=(6, 2))
         | print(linalg.eigvalsh(X.T @ X))
Out [2]: | [3.02532233 4.56956673]
In [3]: | print(linalg.svdvals(X)**2)
Out [3]: | [4.56956673 3.02532233]
```

Nonzero eigenvalues of $X^T X$ (or XX^T) are squared singular values of X .

- ③ Compute the principal directions and principal components by hand using the unbiased variance–covariance estimator. Verify that they coincide with the ones computed by `scikit-learn`.

```
In [4]: n, p = 100, 15
        X = np.random.normal(size=(n, p))
        X0 = X - X.mean(axis=0)
        V = 1/(n-1) * X0.T @ X0
        vp, U = linalg.eigh(V)
        print(vp)

Out [4]: [0.40500148 0.51331787 0.58970774 0.66158489 0.68821847 0.74879712
         0.79081476 0.89547448 1.00224725 1.02939506 1.07853708 1.32517166
         1.40961574 1.52915075 1.81708517]

In [5]: Xpca = X0 @ U
        print(n / (n-1) * Xpca.std(axis=0)**2)

Out [5]: [0.40500148 0.51331787 0.58970774 0.66158489 0.68821847 0.74879712
         0.79081476 0.89547448 1.00224725 1.02939506 1.07853708 1.32517166
         1.40961574 1.52915075 1.81708517]

In [6]: from sklearn.decomposition import PCA
        pca = PCA()
        pca.fit(X)

Out [6]: PCA(copy=True, iterated_power='auto', n_components=None, random_state=None,
         svd_solver='auto', tol=0.0, whiten=False)

In [7]: print(pca.explained_variance_)

Out [7]: [1.81708517 1.52915075 1.40961574 1.32517166 1.07853708 1.02939506
         1.00224725 0.89547448 0.79081476 0.74879712 0.68821847 0.66158489
         0.58970774 0.51331787 0.40500148]
```

We have the same eigenvalues (in reverse order).

2 PCA for dimension reduction

In this section, we use the `boston` regression dataset. To load it use

```
from sklearn.datasets import load_boston
boston = load_boston()
```

- ④ Perform a PCA on this dataset and study how many number of principal components should be retained from the two empirical methods seen in class.

```

In [8]: from sklearn.datasets import load_boston
        boston = load_boston()
        X = boston.data
        y = boston.target
        pca = PCA()
        pca.fit(X)

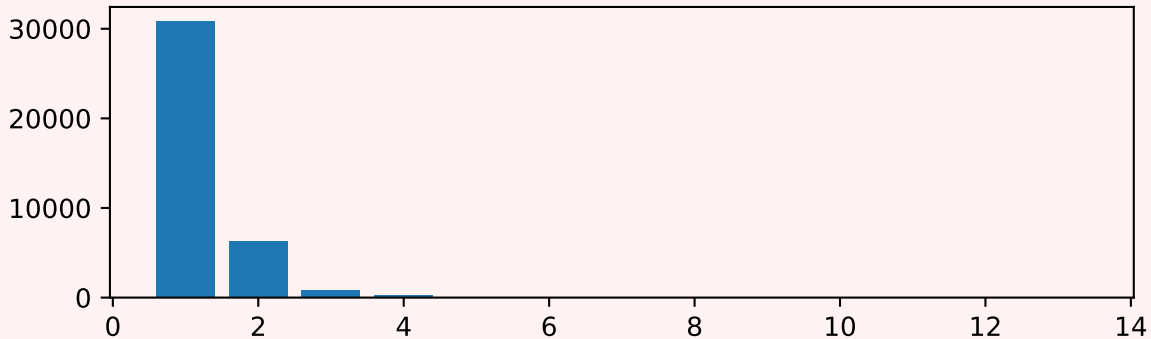
Out [8]: PCA(copy=True, iterated_power='auto', n_components=None, random_state=None,
        svd_solver='auto', tol=0.0, whiten=False)

In [9]: plt.figure()
        plt.bar(range(1, X.shape[1]+1), pca.explained_variance_)

Out [9]: <BarContainer object of 13 artists>

In [10]: plt.show()

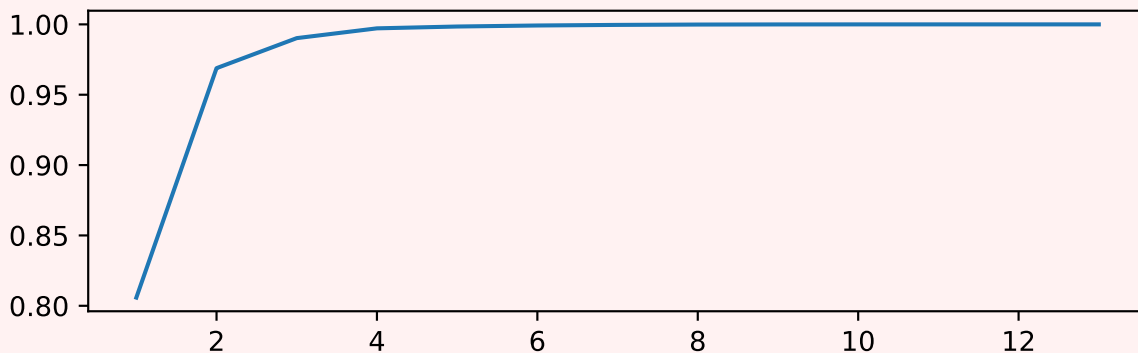
```



```

In [11]: plt.figure()
        plt.plot(range(1, X.shape[1]+1), np.cumsum(pca.explained_variance_ratio_))
        plt.show()

```



The scree plot suggest to choose 2 components. The cumulative sum of explained variance ratio is telling us that to retain 80% of the explained variance, only one component is enough.

- ⑤ Describe the following code. What is it supposed to be doing? Adapt it to determine the optimal number of principal components for the regression task at hand.

```

from sklearn.linear_model import LinearRegression
from sklearn.decomposition import PCA
from sklearn.model_selection import GridSearchCV
from sklearn.pipeline import Pipeline

pca = PCA()
lin = LinearRegression()

```

```
pca_lin = Pipeline([("pca", pca), ("lin", lin)])
clf = GridSearchCV(
    estimator=pca_lin,
    scoring="neg_mean_squared_error",
    cv=10,
    iid=False,
    param_grid=dict(pca__n_components=range(1, X.shape[1] + 1)),
)
clf.fit(X, y)
```

This snippet of code is defining a pipeline consisting of a PCA followed by a linear regression on the boston data. The optimal number of components is then computed by cross-validation with 10 folds.

```
In [12]: from sklearn.linear_model import LinearRegression
         from sklearn.decomposition import PCA
         from sklearn.model_selection import GridSearchCV
         from sklearn.pipeline import Pipeline

         pca = PCA()
         lin = LinearRegression()
         pca_lin = Pipeline([("pca", pca), ("lin", lin)])
         clf = GridSearchCV(
             estimator=pca_lin,
             scoring="neg_mean_squared_error",
             cv=10,
             iid=True,
             param_grid=dict(pca__n_components=range(1, X.shape[1] + 1)),
         )
         clf.fit(X, y)

Out [12]: GridSearchCV(cv=10, error_score=nan,
                      estimator=Pipeline(memory=None,
                                         steps=[('pca',
                                                  PCA(copy=True, iterated_power='auto',
                                                       n_components=None,
                                                       random_state=None,
                                                       svd_solver='auto', tol=0.0,
                                                       whiten=False)),
                                                  ('lin',
                                                  LinearRegression(copy_X=True,
                                                                    fit_intercept=True,
                                                                    n_jobs=None,
                                                                    normalize=False))],
                                         verbose=False),
                      iid=True, n_jobs=None,
                      param_grid={'pca__n_components': range(1, 14)},
                      pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
                      scoring='neg_mean_squared_error', verbose=0)

/home/sylvain/.local/lib/python3.8/site-packages/sklearn/model_selection/_search.py:823:
↪ FutureWarning: The parameter 'iid' is deprecated in 0.22 and will be removed in 0.24.
  warnings.warn(

In [13]: print(clf.best_params_)
Out [13]: {'pca__n_components': 10}
```

The best number of principal components using the mean square error and a 10-CV is then 10 (it might vary)

⑥ Is standardizing data improving the optimal number of principal components?

This snippet of code is defining a pipeline consisting of a PCA followed by a linear regression on the boston data. The optimal number of components is then computed by cross-validation with 10 folds.

```
In [14]: from sklearn.linear_model import LinearRegression
         from sklearn.decomposition import PCA
         from sklearn.model_selection import GridSearchCV
         from sklearn.pipeline import Pipeline

         pca = PCA()
         lin = LinearRegression()
         pca_lin = Pipeline([("pca", pca), ("lin", lin)])
         clf = GridSearchCV(
             estimator=pca_lin,
             scoring="neg_mean_squared_error",
             cv=10,
             iid=True,
             param_grid=dict(pca__n_components=range(1, X.shape[1] + 1)),
         )
         Y = X / X.std(axis=0)
         clf.fit(Y, y)

Out [14]: GridSearchCV(cv=10, error_score=nan,
                      estimator=Pipeline(memory=None,
                                         steps=[('pca',
                                                  PCA(copy=True, iterated_power='auto',
                                                       n_components=None,
                                                       random_state=None,
                                                       svd_solver='auto', tol=0.0,
                                                       whiten=False)),
                                                  ('lin',
                                                   LinearRegression(copy_X=True,
                                                                    fit_intercept=True,
                                                                    n_jobs=None,
                                                                    normalize=False))],
                                         verbose=False),
                      iid=True, n_jobs=None,
                      param_grid={'pca__n_components': range(1, 14)},
                      pre_dispatch='2*n_jobs', refit=True, return_train_score=False,
                      scoring='neg_mean_squared_error', verbose=0)

/home/sylvain/.local/lib/python3.8/site-packages/sklearn/model_selection/_search.py:823:
↪ FutureWarning: The parameter 'iid' is deprecated in 0.22 and will be removed in 0.24.
  warnings.warn(

In [15]: print(clf.best_params_)
Out [15]: {'pca__n_components': 9}
```

The best number of principal components using the mean square error and a 10-CV is then 9 (it might vary)

3 Problem: band reduction in multispectral images

A multispectral image is an image that has several components. For example, a color image has 3 components: red, green and blue and each pixel can be viewed as a vector in \mathbb{R}^3 . More generally a multispectral image of size $N \times M$ with P spectral bands can be stored as a $N \times M \times P$ array. There are $N \times M$ pixels living in \mathbb{R}^P .

When the number of spectral bands P is too large, it is desirable to somehow reduce that number ultimately to 3 for viewing purposes. This process is called band reduction.

Propose a method using the PCA performing a band reduction to 3 bands and use it on the provided multispectral image.

Some multispectral images are available on the internet to test your band reduction

algorithm. See for example the following website

- <http://lesun.weebly.com/hyperspectral-data-set.html>

Most of them are available as a Matlab data file (.mat files). It can be loaded with `scipy` with the following function

```
| scipy.io.loadmat
```

You will probably have to reshape arrays. It can be done with the `reshape` method. For example, an array of size $6 \times 6 \times 3$ can be “linearized” using `reshape`

```
| X_lin = X.reshape((-1, 3))
```

the `-1` is automatically inferred from the number of elements in the array. The array is then reshaped into an array of size 36×3 .

It might be handy to be able to rescale the data when it has to belong to some specific range. `scikit-learn` has several rescalers available. For example

```
| from sklearn.preprocessing import MinMaxScaler
```

rescales the data between 0 and 1.

`matplotlib` can display images with the function

```
| plt.imshow
```

Beware of the type of the array (float or integers)!