Creative Data Mining

Using ML algorithms in python

Artem Chirkin Dr. Daniel Zünd Danielle Griego

> Lecture 7 10.04.2017

Chair of Information Architecture



DARCH

R P P P 1/1

What we will cover today

ETT

Outline



ETH



At first, we need to set up our environment.

We'll use several package, most of which you have already installed

pip install numpy pandas matplotlib sklearn

1# We need pandas to import and use , csv files 2 import pandas as pd 4 # numpy is used for various numeric operations 5 import numpy as np 7# matplotlib lets us draw nice plots 8 import matplotlib.pyplot as plt 10 # We will use several modules of sklearn package 11# that is a "swiss knife" ML toolset for python 12 from sklearn import cluster, metrics, decomposition



Example data for clustering

The Iris flower data set is a multivariate data set introduced by the British statistician and biologist Ronal Fisher in 1936.

You can read more about it on Wikipedia's page https://en.wikipedia.org/wiki/Iris_flower_data_ set

You can find the .csv file easily; I got it from https://raw.githubusercontent.com/vincentarelbundock/ Rdatasets/master/csv/datasets/iris.csv

DARCH

Chair of Information

Let's download it and use in our python script





Example data for clustering

Now, we are ready to load the .csv file into the interpreter using pandas package.

DARCH

Chair of Informatior



= 3°° - 4/1





Eidgenössische Technische Hochschule Zürich Swiss Federal Institute of Technology Zurich

Trying useful functions from pandas

The first useful command is head()

It returns first several lines of a dataframe. Very useful to get an idea of what kind of data you have!

1 ourData.head()

					XX 19H \4	1°	S & 2 /~
#		Unnamed: C	Sepal.Length	Sepal.Width	Petal.Length	Petal.Width	Species
#	0	⊘ 1	. 5.1	3.5	1.4	200.2	setosa
#	1	2	2 4.9	3.0	1.4	0,2	setosa
#	2	3	3 4.7	3.2	1.3	0.2	setosa
#	3	4	4.6	3.1	1.5	0.2	setosa
#	4	⇒ √5	5,0	3.6	1.4	0.2	setosa
	F F	9 9 10 10 10 10 10 10 10 10 10 10 10 10 10					
× T							

5 DQC

-5117





Chair of Information

Trying useful functions from pandas

Second command, info(), can give a more precise information _1 ourD on what data types we operate on.

- ► The dataset has 150 entries
- Sepal and petal parameters are floating-point numbers
- Species is recognized as object type in fact, this is a text

DARCH





Trying useful functions from pandas

The last command is more advanced: ourData.describe().

This command gives some statistical information about numeric values in your dataset.

DARCH

It is useful to understand what is the range of your values.

1 ourData . describe ()

#		Unnamed: 0	Sepal, Length	Sepal.Width	Petal.Length	Petal.Width
#	count	150.000000	150.000000	150.000000	150.000000	150.000000
#	mean	75.500000	5.843333	3.057333	3.758000	1,199333
#	std	43.445368	0.828066	0.435866	1.765298	0.762238
#	min	1,000000	4.300000	2.000000	1,000000	0.100000
⊘#,	25%	38.250000	5.100000	2.800000	1.600000	0.300000
∕#	50%	75.500000	5,800000	3.000000	4.350000	1.300000
#	75%	112.750000	6.400000	3.300000	5.100000	1.800000
#	max 🖉	150.000000	7.900000	4.400000	6.900000	2.500000
			- VN - 7 6 A VA - 3 A - 7 Y			





Trying useful functions from pandas

E

Eidgenössische Technische Hochschule Zürich Swiss Federal Institute of Technology Zurich

But what about the last column Species?

We've seen these are textual values. Now, let us see what kinds of values does this column have?

DARCH

Chair of Information



ETH



Draw some plots Using matlplotlib

Let's try to plot something!

ΞT

Eidgenössische Technische Hochschule Zürich Swiss Federal Institute of Technology Zurich

In this example I map names of the species onto colors. Play with this code a little bit changing the columns to plot. This gives you an insight how the data looks like.

DARCH



9 #

use either

1# plot different species in different colours

, ourData['Sepal.Width']
, c=ourData['Species']

.apply(lambda x: species_dict[x])

2 colors = ['g', 'r', 'b', 'c', 'm', 'k'] 3 species_dict = dict(zip(species, colors))

4 plt.scatter(ourData['Sepal Length]

Draw some plots

Using matlplotlib

...And here is how the result looks like





DARCH

ETT

Eidgenössische Technische Hochschule Zürich Swiss Federal Institute of Technology Zurich

Copy the data to a new variable ourDataNoLabels. We will play like we don't know the labels and want to estimate them.

DARCH

Chair of Information



Principal Component analysis

Using sklearn.decomposition

This code performs PCA (lines 3-7). (Most of the code is just to draw three plots).

Principal component transform is a procedure that rotates point space in such a way that points variance is the biggest along X-axis, the second is along Y-axis, and so on.

The method is simple but very powerful for data exploration. If data have many dimensions (i.e. 50) it is very convenient to look only at 2-5 most significant dimensions. Otherwise visual inspection of the data would be too difficult.

DARCH

Chair of Informatio

1# Run Principal Component Analysis 2 # to get more understanding how our data looks like 3 ourDataReduced = decomposition .PCA(n_components=3) . fit_transform (ourDataNoLabels) 6 7 plt.scatter(ourDataReduced[:,0] , ourDataReduced [:,1] 8 , c=ourData['Species'] .apply(lambda x: species_dict[x])) 10 11 plt.savefig("real-pca-1-2.png") 12 plt.scatter(ourDataReduced[:,0] , ourDataReduced [:,2] 13 , c=ourData['Species'] 14 .apply(lambda x: species_dict[x])) 15 16 plt.savefig("real-pca-1-3.png") 17 plt.scatter(/ourDataReduced[:,1] , ourDataReduced [:,2] 18 _, c=ourData['Species'] 19 apply(lambda x: species_dict[x])) 20 21 plt.savefig("real-pca-2-3.png")

P. S. C. P. 2017

Eidgenössische Technische Hochschule Züric Swiss Federal Institute of Technology Zurich

Draw some plots

Using matlplotlib

...And here are the plots



ETTH Eidgenössische Technische Hochschule Zürich Swiss Federal Institute of Technology Zurich DARCH

ETH



K-means clustering

Using sklearn.cluster

Finally, we are ready for some clustering!

```
1# We know that there must be 3 different clusters
2 #
                                     - species of flowers
3# So let's give this hint to the algorithm
4 \text{ kmeans} = \text{cluster} \cdot \text{KMeans}(3)
                    . fit (np.array (ourDataNoLabels))
5
6 foundLabels
    = pd.DataFrame( kmeans.labels_
                    , columns=['K-means clusters'])
9
10
11 plt.scatter( ourData['Sepal.Length']
                ourData ['Sepal. Width']
12
               , c=foundLabels['K-means clusters']
13
                  .apply(lambda x: colors[x])
16 plt. savefig ("predicted -labels.png")
```





Chair of Information

K-means clustering

Using sklearn.cluster

Let's compare plots...



-15/17

...would you say this is a good result?







ETTH Eidgenössische Technische Hochschule Zürich Swiss Federal Institute of Technology Zurich

Assess clustering quality

Using pandas

A very good way to assess performance of an unsupervised clustering algorithm is to look at co-occurrence tables.

Package pandas provides a special function crosstab() that calculates how many times a value from one column occurs together with a value from another column.

DARCH

Chair of Information

So, what would be a conclusion now?





Using sklearn.cluster

This code does everything the same way as KMeans clustering example.

Try it! And compare the results. Which algorithm performs better?











-17 / 17

What we have covered today

Thanks for your attention!

ETH

