

## Plotting for Exploratory Data analysis (EDA)

- EDA is a task of analysing data using simple tools from statistics, plotting tools and linear algebra.
- Vector is a n-dimensional numerical array.
- Dataset is whole given table
- Features / variable / input-variable / independent variable are sepal and petal length and width.
- Label / dependent variable / output-variable / class / Response label is the species to which particular flower belongs.

Dataset used

↳ Iris flower dataset

In [15]: # how many data-points and features?  
`print(iris.shape)`

(150, 5)

In [16]: # what are column names in our dataset?  
`print(iris.columns)`

`Index(['sepal-length', 'sepal-width', 'petal-length',  
 'petal-width', 'species'],  
 dtype='object')`

In[17]: # How many data points for each class are present  
# How many flowers for each species are present  
iris[ "species" ].value\_counts()

Out[17]: setosa 50  
virginica 50  
versicolor 50  
Name: species, dtype: int64

[ refer attached images ]

## (Pair-Plot)

In [11]: # pairwise scatter plot: Pair-Plot

# Dis-advantages:

## can be used when number of features are high

# cannot visualize higher dimensional patterns in 3D and 4-D

# only possible to view 2D patterns

```
plt.close();
```

```
sns.set_style("whitegrid");
```

```
sns.pairplot(iris, hue = "species", size=3);
```

```
plt.show()
```

[refer given pair-plot images]

### Observations

① petal-length and petal-width are most useful features to identify various flower types.

② while Setosa can be easily identified (linearly separable), Virginica and Versicolor have some overlap.

### Limitation of pair-plot

① Pair-plot are easy to understand when dimensionality of data is small but as soon as the dimensionality becomes greater than 6 pair-plot cannot help much.

## Histogram, PDF, CDF

In[175]: # What about 1-D scatter plot using one feature?  
# 1-D scatter plot of petal length  
import numpy as np  
iris\_setosa = iris.loc[iris["species"] == "setosa"];  
iris\_virginica = iris.loc[iris["species"] == "virginica"];  
iris\_versicolor = iris.loc[iris["species"] == "versicolor"];  
plt.plot(iris\_setosa["petal-length"], np.zeros\_like(iris\_setosa)  
 ['petal-length']);  
plt.plot(iris\_versicolor["petal-length"], np.zeros\_like(iris\_versicolor)  
 ['petal-length']);  
plt.plot(iris\_virginica["petal-length"], np.zeros\_like(iris\_virginica)  
 ['petal-length']);  
plt.show()

[refer to attached image]

In[172]: sns.FacetGrid(iris, hue = "species", size = 5)  
 .map(sns.distplot, "petal-length")  
 .add\_legend();  
plt.show();

[refer to attached image]

# Histogram on probability density Function (Univariate analysis)

CDF (Cumulative Distribution Function)

differentiating cdf  $\rightarrow$  pdf

integration of pdf  $\rightarrow$  cdf

In [12]: # How to construct and read cdf?

# plot CDF of petal length

counts, bin-edges = np.histogram(iris.setosa  
['petal-length'], bins=10, density=True)

pdf = counts / sum(counts)

print(pdf);

print(bin-edges)

compute CDF

cdf = np.cumsum(pdf)

plt.plot(bin-edges[1:], pdf)

plt.plot(bin-edges[1:], cdf)

plt.show()

(refer attached image)

## Mean, Variance and Std-dev

```
In[164]: # Mean, variance, std-dev  
print("Means:")  
print(np.mean(iris_setosa["petal_length"]))  
# Mean with outlier  
print(np.mean(np.append(iris_setosa["petal_length"],  
[150])));  
  
print(np.mean(iris_virginica["petal_length"]))  
print(np.mean(iris_versicolor["petal_length"]))  
  
print("\nstd-dev:");  
print(np.std(iris_setosa["petal_length"]))  
print(np.std(iris_virginica["petal_length"]))  
print(np.std(iris_versicolor["petal_length"]))
```

$$\text{Mean} = \frac{1}{n} \sum_{i=1}^n x_i \times \frac{1}{n}$$

$$\text{Var} = \frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2$$

$$\text{Std-dev} = \sqrt{\frac{1}{n} \sum_{i=1}^n (x_i - \mu)^2}$$

Means:

1. 464

2. 41568627451

3. 552

4. 26

Std-dev & spread

Std-dev:

0.171767284429

0.546347874527

0.465188133985

## Median, Percentile, Quantile, IQR, MAD

In[1167]: `pwint["\nMedians:\"]`

```
pwint(np.median(iris_setosa["petal_length"]))
```

# Median with outlier

```
pwint(np.median(np.append(iris_setosa["petal_length"], 50)))
```

```
pwint(np.median(iris_virginica["petal_length"]))
```

```
pwint(np.median(iris_versicolor["petal_length"]))
```

Medians:

1.5

1.5

5.55

4.35

### Median

$$n = \{1, 1.2, 1.1, 2.1, 1.8, 1.6, 1.2\}$$

① Sort them in order

$$\{1, 1.1, 1.2, 1.2, 1.6, 1.8, 2.1\}$$

1 2 3 4 5 6 7

② pick middle value  $\left(\frac{n+1}{2}\right)$  odd  $\left(\frac{2+1}{2}\right) = 4$

$\left(\frac{n}{2}\right)$  even

$\left[\frac{n}{2}, \frac{n}{2} + 1\right]$

Median of  $n = 1.2$

average

In[168]: print("90th Percentiles:")

```
print(np.percentile(iris['setosa']["petal length"], 90))  
print(np.percentile(iris['virginica']["petal length"], 90))  
print(np.percentile(iris['versicolor']["petal length"], 90))
```

90th Percentiles:

1.7

6.3)

4.8

### Percentile

$$n_s = \left\{ \begin{array}{c} 1 \ 2 \ 3 \ \dots \ 100 \\ \end{array} \right\}$$

$$n=100$$

$$\text{Median}(n) = \frac{n_{50} + n_{51}}{2} (\text{Mean}[n_{50}, n_{51}])$$

Median = 50<sup>th</sup> percentile value of  $n = n_s[50]$

10<sup>th</sup> percentile value of  $n = n_s[10]$

25<sup>th</sup> percentile, 50<sup>th</sup>, 75<sup>th</sup>, 100<sup>th</sup> = Quantile

```
In [167]: print("")\nQuantiles:\nprint(np.percentile(iris.setosa[\"petal-length\"],\n                    np.arange(0,100,25)))\nprint(np.percentile(iris-virginical[\"petal-length\"],\n                    np.arange(0,100,25)))\nprint(np.percentile(iris-versicolor[\"petal-length\"],\n                    np.arange(0,100,25)))
```

from statsmodels import robust

```
print(\"\\nMedian Absolute Deviation:\")
```

```
print(robust.mad(iris-setosa[\"petal-length\"]))
```

```
print(robust.mad(iris-virginica[\"petal-length\"]))
```

```
print(robust.mad(iris-versicolor[\"petal-length\"]))
```

Quantiles:

[1.	1.4	1.5	1.575 ]
4.5	5.1	5.55	5.875 ]
3.	4.	4.35	4.6 ]

Mean Absolute Deviation:

0.148260221851

0.667170998328

0.518910776477