

Plotting for Exploratory Data Analysis:-

→ Analyzing data using plotting tools, statistics, linear algebra etc.

→ column to be predicted = class label / dependent variable

(Data already there = data points / vectors)

→ columns are called features / input variable / independent variable

→ 1D vector = scalar.

$[a], [1], [8]$

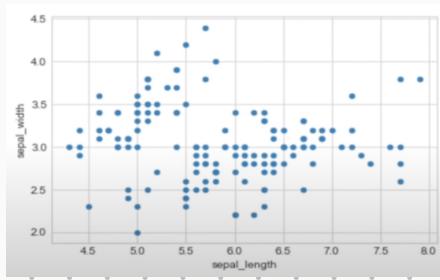
$\downarrow \quad \downarrow \quad \downarrow$

'a' 1 8

→ Balanced Dataset → Each class has equal number of datapoints.

→ When reading a plot, always read axis labels & values - It doesn't always start at 0.

→ iris.plot (kind='scatter', x="sepal_length", y="petal_width");

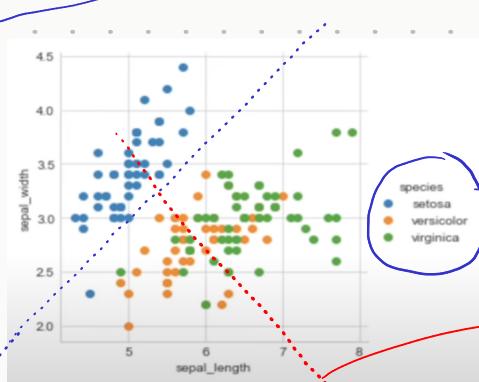


→ import seaborn as sns

sns.set_style("whitegrid") → white grid structure

sns.FacetGrid((iris), hue="species", size=4).map(plt.scatter, "sepal_length", "sepal_width").add_legend();

dataset
by which
column should
the dataset be
colored



type of plot
x-Axis
y-Axis

FacetGrid

(sns.scatterplot() only plots one plot)
(sns.FacetGrid() can plot multiple plots.)

legend
cont separate as too many outliers

line could separate setosa from versicolor & virginica (this is called linearly separable)

Observations:- ① S-length & S-width can separate setosa flowers from others.
② Separating versicolor & virginica is harder.

QUICK SEABORN INTRODUCTION:-

→ simpler way to plot attractive plots.

→ high level interface to matplotlib.

→ some features include:

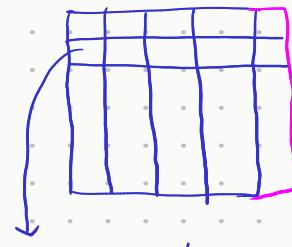
- default aesthetic themes.
- custom color palettes

* Visualizing info from matrices & dataframes -

* attractive statistical plots.

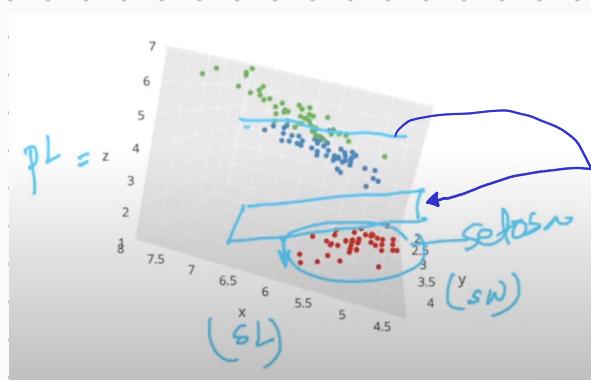
* easily displaying distributions

→ Seaborn is a complement, not a substitute of matplotlib



1 row = 1 array/vector.

3D Scatter Plot :-



when 3D plot is plotted, we use a plane to separate instead of a line.

Plotting petal length helps notice that plotting a line can separate them which we were not able to do in a 2D plot.

Drawbacks :- → while being interactive as a plot, it can't be viewed on 2D surface like paper or in an ipynb.

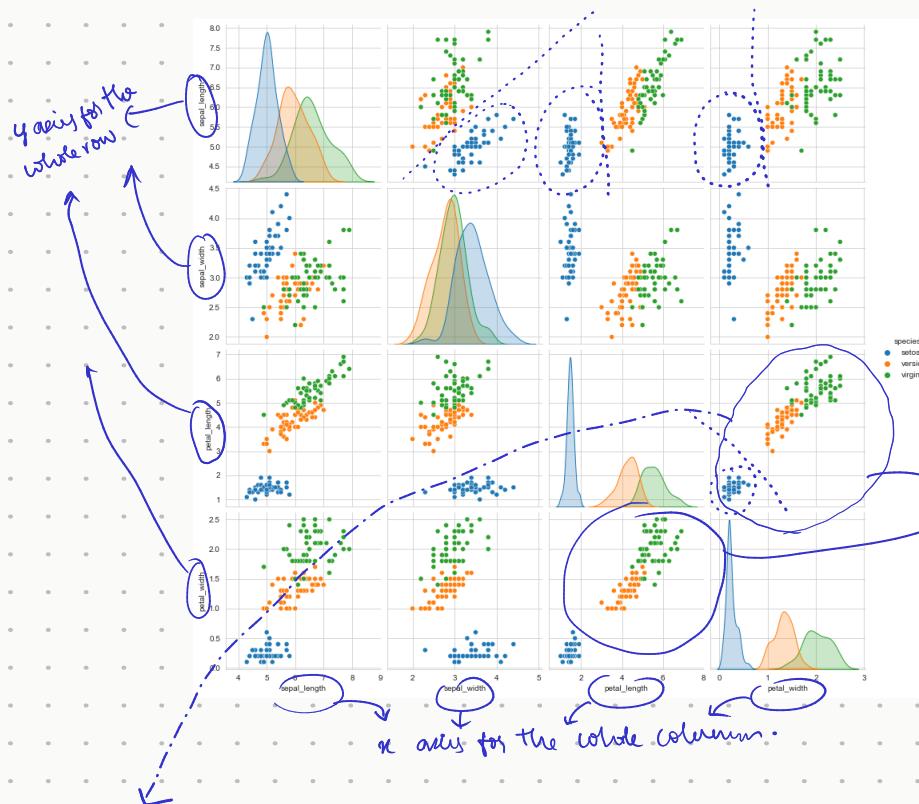
→ Since we can't visualize 4D datasets, we use pair plots.

↓
there are 4 variable $\Rightarrow 4 \times 3$ pairs are possible.

$$6 \rightarrow \begin{cases} (SL, SW), (SW, PL) \\ (SL, PL), (PW, PL) \\ (SL, PW), (SW, PW) \end{cases}$$

Since we can't visualize 4D,
we try visualizing these 6.

< sns.pairplot(iris, hue="species", size=3) >

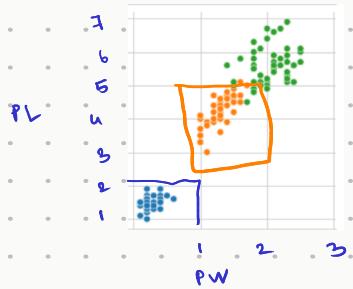


(The diagonal plots are PDFs of each feature).

moreover - As both are b/w PL, PW.

Similarly there are 3 pairs in total.

So we form only 6 plots.



\Rightarrow if $PL \leq 2$ & $PW \leq 1$
then flower type = setosa

| if $(PW \leq 2 \& PW \geq 1) \& (PL \leq 5 \& PL \geq 2.5)$
then versicolor

→ There will be some error but that's OK.

Drawback of pairplots :-

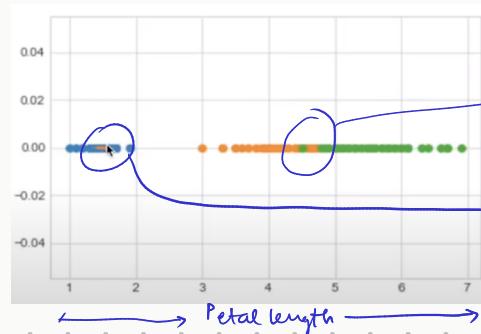
→ When dimensionality is 4, we get 4C_2 plots.

If the dimensionality is 7, we will get 7C_2 plots which is not that useful.

So we use techniques like PCA, t-SNE that will help with this.

In real world, we limit pairplots to only important features & not apply it on all features.

1D Scatter Plots :-



→ harder to read.
→ green overlaps orange.

→ Number of points here is not clear.

To overcome this, instead of using a scatterplot, we can use a histogram.

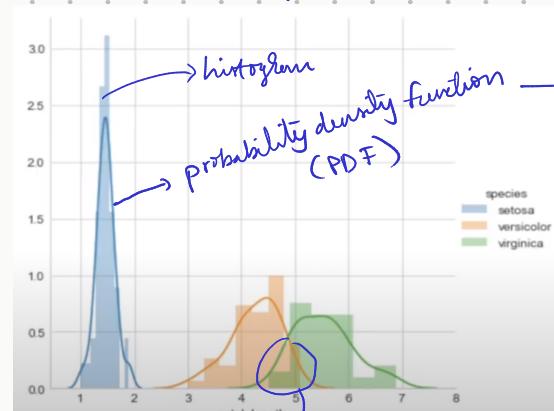
`sns.FacetGrid(iris, hue="species", size=5).`

`map(sns.distplot, "petal_length").add_kde()`

distribution plot aka histogram aka density plot
y-axis: count.

if petal length < 2
then setosa.
if petal length > 2 &
petal length < 4.7
then versicolor
if petal length > 4.7
then virginica.

But there will be overlap



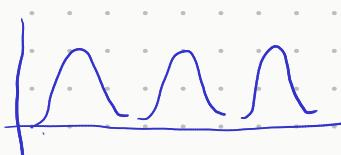
overlapping region

Smoothed histogram.
Smoothed using kernel Density Estimation. (KDE)

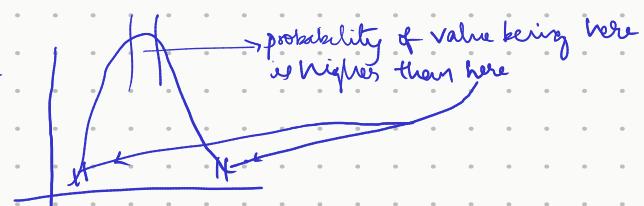
Univariate Analysis :-

→ Analysis of one variable (Uni + variate)

→ The farther the distributions are, the better.

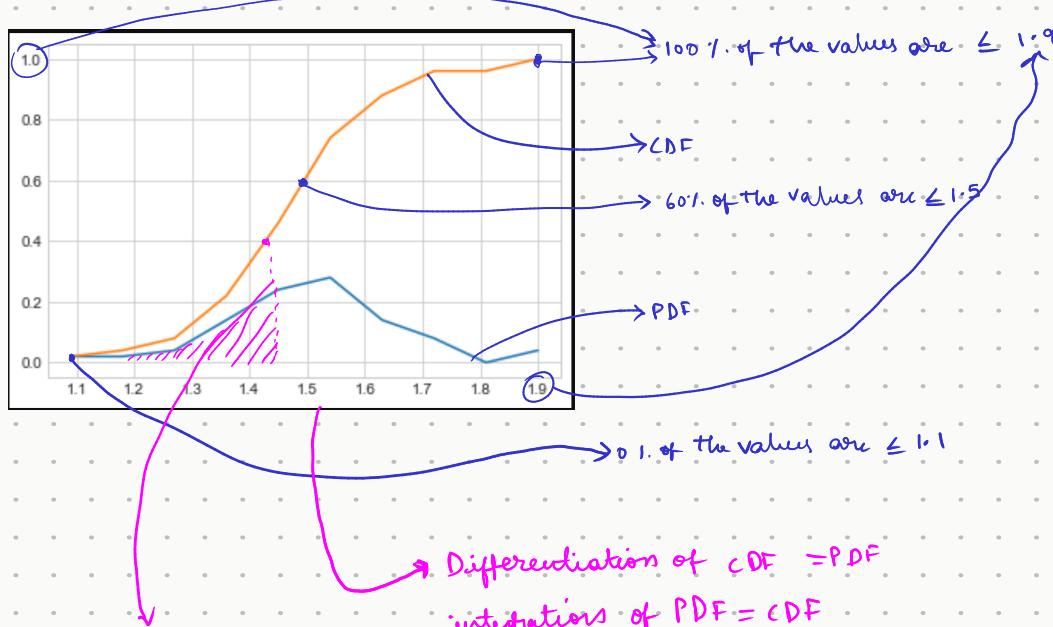


PDF → probability of variable to be in the range of values. Ex:-



Cumulative Data Function :- (CDF)

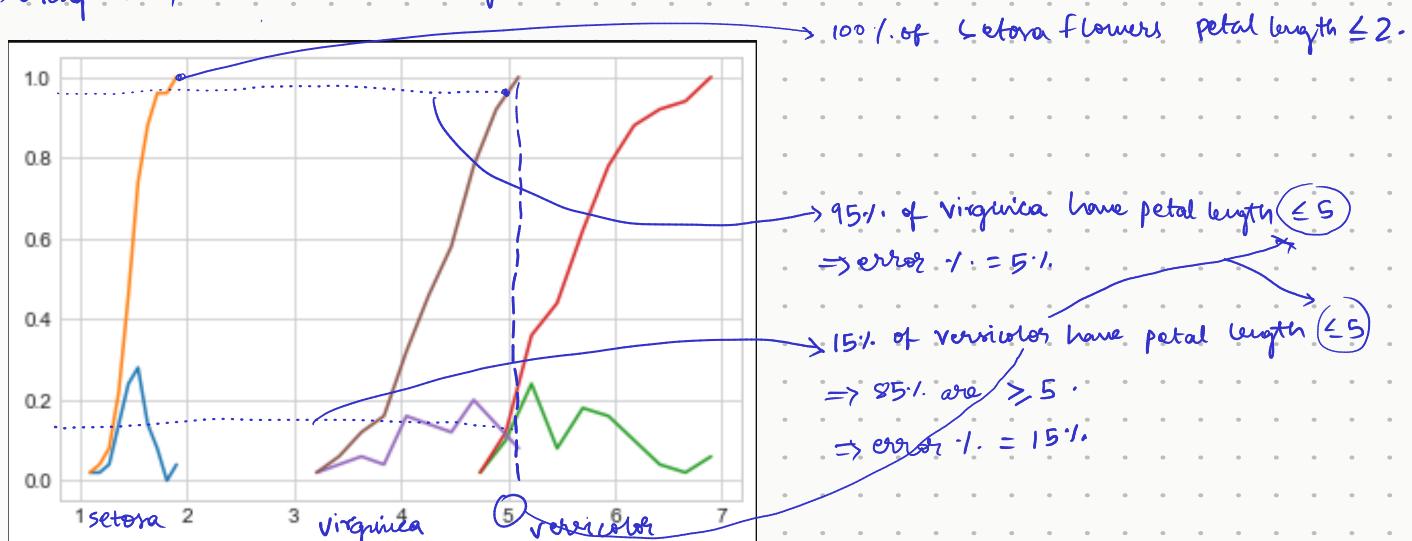
→ what % of values are \leq the value at that point.



The CDF

Value at this point = sum(values) till that point = cumulative sum.

→ Using CDF, we can make assumptions such as



bin = number of bars in a region in a histogram.

→ np.histogram returns counts & bin.edges.



How to compute CDF ? :-

```

counts, bin_edges = np.histogram(iris['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)

```

number of bins

if false, it will just return the number of samples in each bin (aka normal hist). But if its true - returns PDF

cumulative sum.

→ Another way to plot CDF :-

`sns.kdeplot(data, cumulative=True)`

→ But can't adjust bins in this.

Mean, Variance & Standard deviation :-

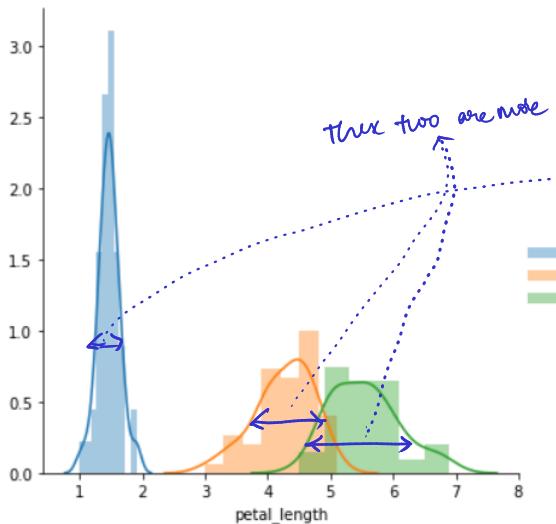
$$\text{Mean} = \frac{1}{n} \times \sum_{i=1}^n (x_i)$$

Gives an average value. But outliers affect this value. For ex, mean of setosa = 1.464.

Adding one extra value 50 will change the mean to 2.41.

How does mean help ?

Spread = how spreaded is the graph.



These two are more widespread than setosa.

We calculate standard deviation which gives us a range. Useful when we can't plot the data.

$$\text{std-dev} = \sqrt{\text{variance}}$$
$$\text{variance} = \frac{1}{n} \sum_{i=1}^n (m - x_i)^2$$

mean of dataset
Square because diff could be -ve.
unit²

`np.mean()` → mean

`np.std()` → standard deviation.

Now we can use this to estimate type of iris.

$$m \text{ of setosa} = 1.4$$

$$\text{std-dev of setosa} = 0.171$$

Now if flower is in the range $(1.4 - 0.171)$ to $(1.4 + 0.171)$ it would be a setosa.

Median :-

→ To overcome mean outliers problem, we use median.

→ Median = sort all elements & pick middle element.

If number of elements = even, mean of middle two elements = median.

→ Median is not corrupted as long as number of corrupted elements $\leq 50\%$ of total elements.

`np.median()` to get median.

Percentile & Quantile :-

Let x be an array containing some values (100 values)

$$x = [\dots \dots \dots]$$

Let x_s be the sorted version of x .

$$x_s = [\dots \dots \dots] \rightarrow \text{indices}$$

The mean of elements at positions 50, 51 = medians of x .

↳ This is also known as the 50th percentile.

↳ 50% of the values are \leq this value.

Similarly $x[10] = 10^{\text{th}}$ percentile.

↳ 10% of the values are \leq this value.

→ 25th, 50th, 75th, 100th percentiles are called Quantiles.

How to get percentiles?

`np.percentile(dataarr, 90)` → will give 90th percentile

`np.percentile(dataarr, np.arange(0, 100, 25))` → will give

↳ takes care of sorting:

We don't want to pass a sorted array.

0th percentile,
25th percentile
50th percentile &
75th percentile

→ 90th & 95th percentile in commerce helps make sense out of data.

e.g.: In Amazon, 90th percentile of delivery dates tell us 90% of orders are delivered in \leq that day & they can work on improving that by adding more trucks, personnel etc.,

→ Duplicate values shift percentile.

```
a1=[1,2,3,4,5]  
a2=[1,2,2,2,2,2,2,2,3,4,5]  
print(np.percentile(a1,90)) -- 4.6  
print(np.percentile(a2,90)) -- 3.9
```

Median Absolute Deviations :-

→ Same notion as standard deviation.

$$\text{median} \left(\underbrace{|x_i - \text{median}|}_{\text{absolute deviation}} \right)_{i=1}^n$$

How to calculate?

→ Not readily available in numpy.

```
from statsmodels import robust  
robust.mad(dataarr)
```

mean → median
Standard deviations → Median Absolute Deviations

don't have problems with outliers.

IQR :- Inter Quartile Range = 75th percentile - 25th percentile.

↳ 50% of values lie in this range.

👤 Aman Sharma

I have seen many students are confused about why we are using SD rather than MAD. Don't worry I have an example.

e.g.

1,2,3,4,5

if you calculate the standard deviation of this number it will be 1.4142

now calculate MAD it will come around 1

now add or append 50

1,2,3,4,5,50

the standard deviation will be 17.1884

and MAD will be 1.5

so this is the power of MAD>SD. Try to explain all this concept with an example.

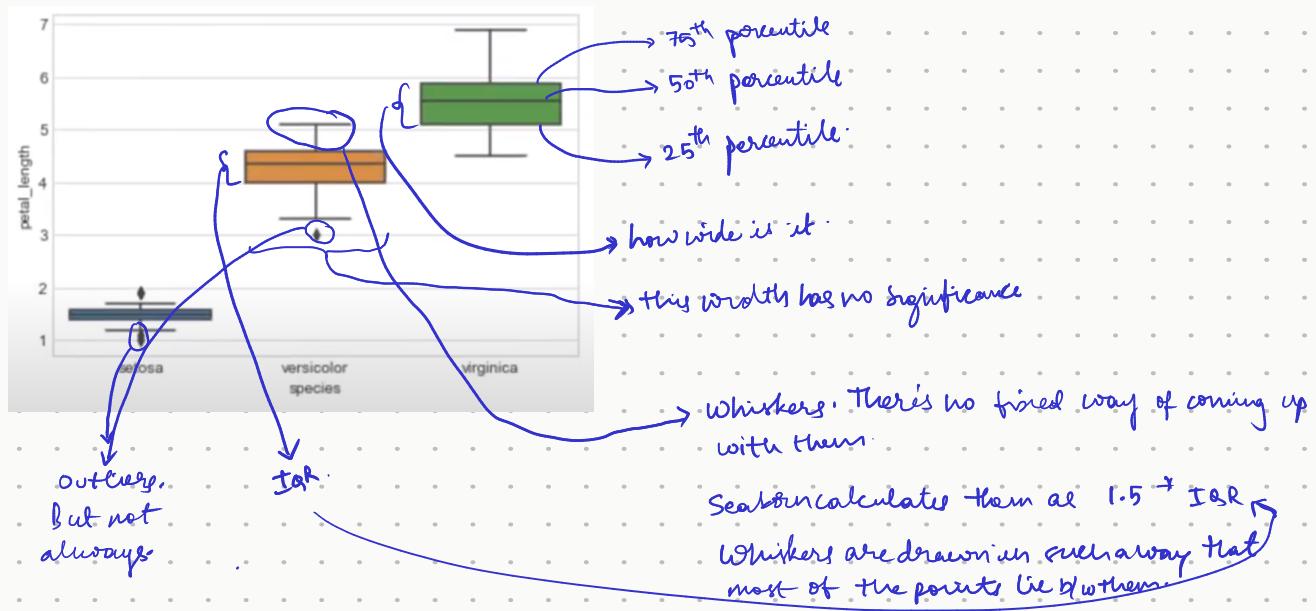
Boxplot & whiskers :-

→ histograms are not good at telling where the percentile lies.

→ We use boxplots for this.

→ ex:- for iris dataset:-

```
sns.boxplot(x="species", y="petal_length", data=iris)
```



→ Box plot also tells us the skewness of the data. If box is cut in half, it means data isn't too skewed. Otherwise it is skewed. Sometimes it may also mean spread of data is high, like in the above graph.

Deleting outliers :-

if the point lies within $[Q_1 - (1.5 \times IQR), Q_3 + (1.5 \times IQR)]$
at \Rightarrow not an outlier. Otherwise it is.

Whiskers are extended based on the following conditions:

The value/Extent of whiskers are decided based on minimum value of the following two conditions:

(1) + or - 1.5 * IQR

(2) Maximum / Minimum value present in the data set.

>The whiskers are extended-above up to

minimum value of (Max value present in data set AND Q3+1.5 IQR)

>The whiskers are extended-below up to

minimum value of (Min value present in data set AND Q1-1.5 IQR)

JUSTIFICATION:

** This rule is mentioned clearly in the Data Mining: Concepts and Techniques

Book by Jiawei Han Page No. 50

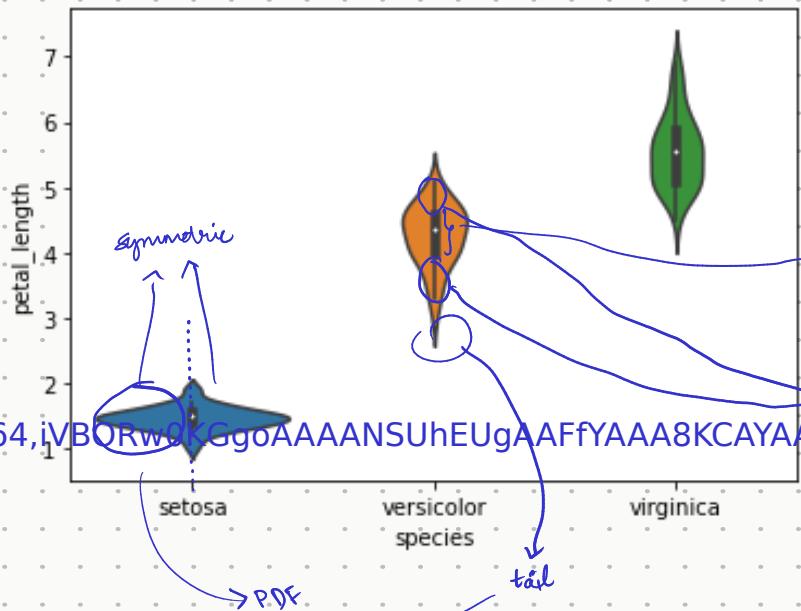
Book PDF Source: <https://www.pdfdrive.com/data-mining-concepts-and-techniques-3rd-2011-e34881113.html>

Violin Plot :-

→ Best of both histogram & box plot.

How to plot:-

```
sns.violinplot(x="species", y="petal_length", data=iris, height = 8)
```



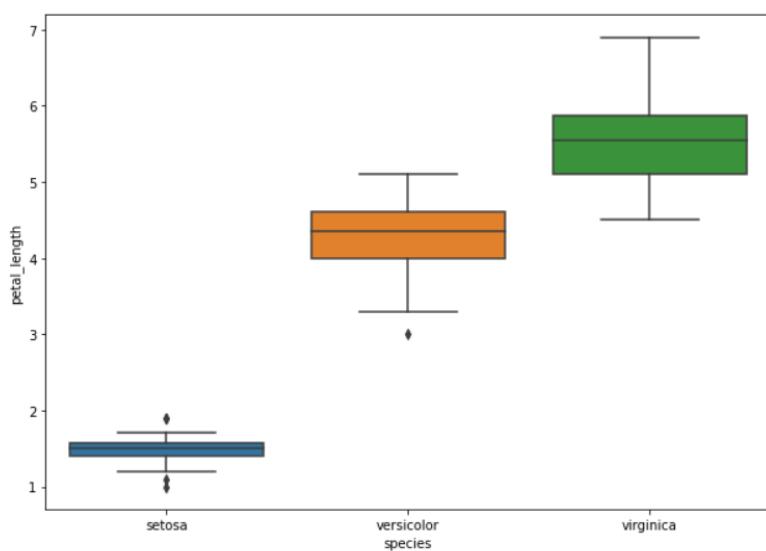
Box plot:
 white dot = 50th percentile
 Top/bottom edge = 75th & 25th.

whiskers.

Violin plots have longer tail if there are outliers.

How to have plots of various sizes :-

```
plt.figure(figsize=(10,7)) → size of plot
ax=plt.subplot(111) → subplot 111
sns.boxplot(x='species',y='petal_length', data=iris,ax=ax)
plt.show()
```



Multivariate probability Density K contour plot :-

→ Univariate → histogram, CDF, PDF

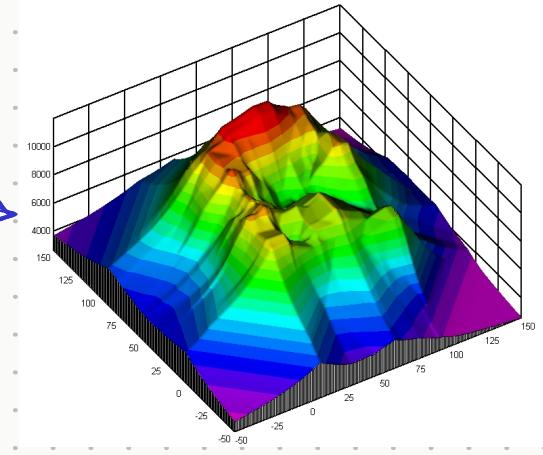
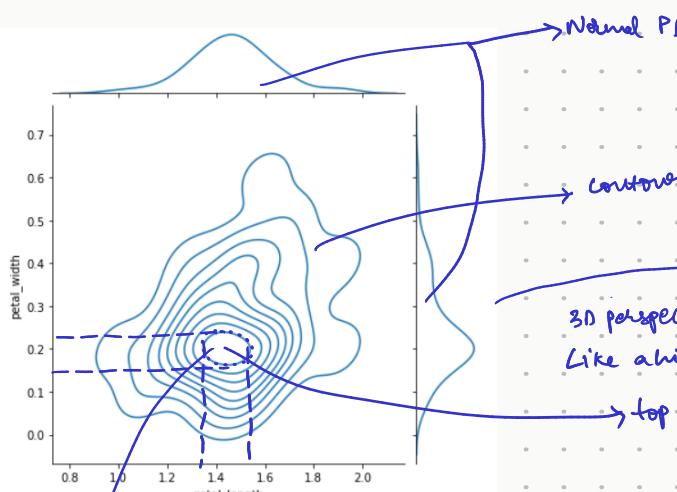
→ Bivariate → scatter plot, pair plot

→ For multivariate we use scatter plot

How to plot

```
| sns.jointplot(x="petal length", y="petal width", data=iris_setosa, kind="kde")
| plt.show()
```

→ The third dimension for this contour plot = density.



This area has the most density (most points)