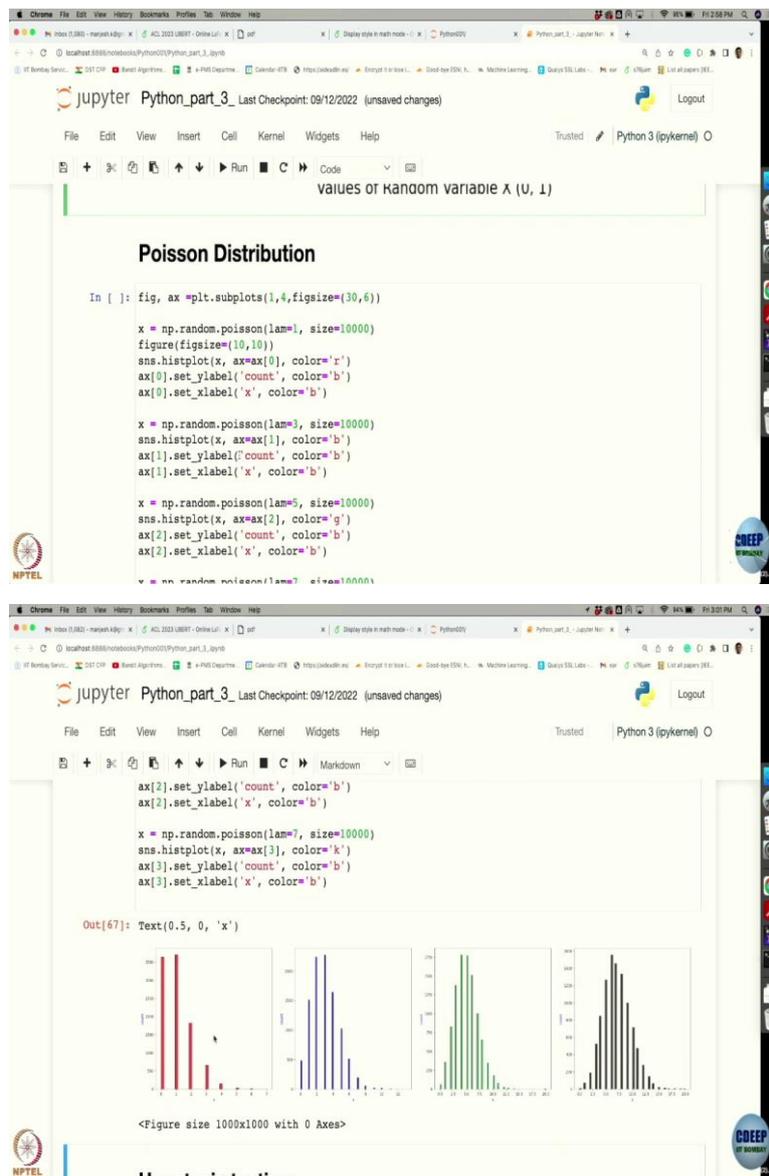


Engineering Statistics
Professor Manjesh Hanawal
Department of Industrial Engineering and Operational Research
Indian Institute of Technology, Bombay
Lecture 60
Generating Random Samples using Python,
Hypothesis Testing using Python (Contd.)

(Refer Slide Time: 0:15)



Next similarly, we can also do, Poisson distribution. But in the Poisson distribution, we will say like how to actually one figure we can use to index multiple figures, or like in one figure we can plot multiple figures as subplots. So, here I am creating subplots. And now I specifying 1 comma 4. That means what I am trying, asking Python to do is I want 4 figures in 1 row. It is like a 1 cross 4 matrix, and 4 figures I want which are at placed adjacent to each other. And the command to get the Poisson sample generated is you just need to invoke

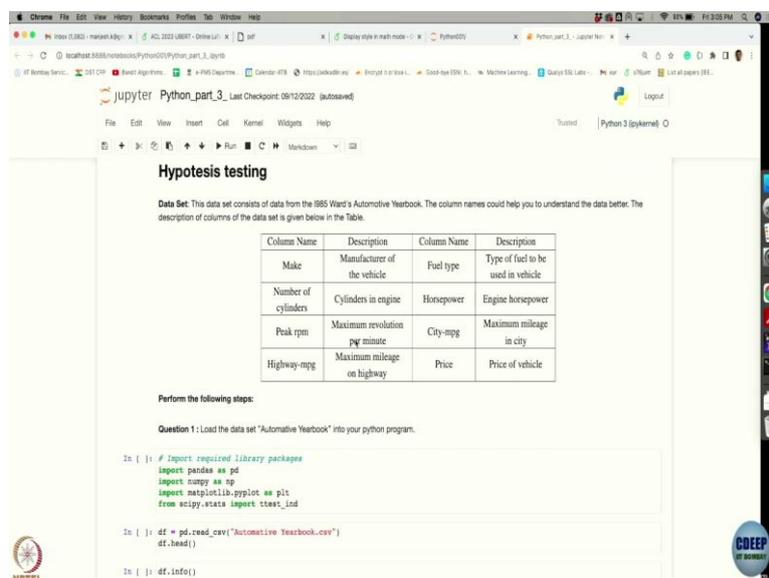
this function `random.poisson` by passing the lambda value, and how many samples you want, it will generate that many samples.

Now, I am going to decide what should be the size of my figure, and that figure size is mentioned here. And now I am going to draw the histogram of this data sample at the first element in my array of 1 cross 4, like I have declared to get 4 figures. And now the first figure I want to put my first figure in the first position. And for that, I am going to use this `x` equals `ax 0`, and then tell, plot the histogram of `x` to this function `histplot`. Let us first get that. So, and now you can set the label here, `y` label as count and your `x` label or `x`. And this one is the figure that corresponds to it.

And similarly, if you want to generate data according to another lambda parameters, and if you see its histogram, but you want it to be the second element in your array, you have to just use set `x` equals to `ax 1`, then you pass the argument to the `histplot`. Similarly, you can do, you can generate third and fourth plot. In that case, the only thing you want to do is like, whichever data generator you generate, and then specify that handler number of that figure, to a `histplot` to generate this plot. So, here, we wanted 4 figures placed adjacent to each other. And that is what we got here.

So, this is about basic distributions functions that are available in Python. I just talked about basic normal distribution, exponential distribution, gamma and beta distributions, so you can play with other distributions. Usually, this Python libraries are very well documented. So, just go and play around it.

(Refer Slide Time: 3:20)



The screenshot shows a Jupyter Notebook titled "Hypothesis testing". It contains a data set description for the 1985 Ward's Automotive Yearbook. The data set includes columns for Make, Number of cylinders, Peak rpm, Highway-mpg, Fuel type, Horsepower, City-mpg, and Price. Below the description, there is a code cell with the following Python code:

```
In [ ]: # Import required library packages
import pandas as pd
import numpy as np
import matplotlib.pyplot as plt
from scipy.stats import ttest_ind

In [ ]: df = pd.read_csv("Automotive Yearbook.csv")
df.head()

In [ ]: df.info()
```

Chrome File Edit View History Bookmarks Profiles Tab Window Help

Python_part_3_Last Checkpoint 09/12/2022 (unsaved changes)

```
In [63]: df = pd.read_csv('Automotive_Teasbook.csv')
df.head()
```

```
Out[63]:
```

	make	fuel-type	num-of-cylinders	horsepower	peak-mpg	city-mpg	highway-mpg	price
0	alfa-romeo	gas	four	111	5000	21	27	13495
1	alfa-romeo	gas	four	111	5000	21	27	16500
2	alfa-romeo	gas	six	154	5000	19	28	16500
3	audi	gas	four	102	5000	24	30	13650
4	audi	gas	five	115	5000	18	22	17450

```
In [70]: df.info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 205 entries, 0 to 204
Data columns (total 8 columns):
 #   column             non-null count  Dtype
---  --
 0   make                205 non-null   object
 1   fuel-type          205 non-null   object
 2   num-of-cylinders   205 non-null   object
 3   horsepower          205 non-null   object
 4   peak-mpg           205 non-null   object
 5   city-mpg            205 non-null   int64
 6   highway-mpg        205 non-null   int64
 7   price               205 non-null   object
dtypes: int64(2), object(6)
memory usage: 12.9+ KB
```

Added Information

Chrome File Edit View History Bookmarks Profiles Tab Window Help

Python_part_3_Last Checkpoint 09/12/2022 (autosaved)

```
In [71]: # To check missing value in a dataset.
print("Number of Question marks in each data column (NaN Values)")
for col in df.columns:
    print(f"{col}: ", df[col].value_counts().get("?"))
```

```
Number of Question marks in each data column (NaN Values)
make: None
fuel-type: None
num-of-cylinders: None
horsepower: 2
peak-mpg: 2
city-mpg: None
highway-mpg: None
price: 4
```

```
In [ ]: # Replace the missing value
for col in ["horsepower", "peak-mpg", "price"]:
    df[col] = df[col].str.replace("?", "", regex=False)
df[col] = pd.to_numeric(df[col])
```

```
In [ ]: df.describe()
```

Question 2: Create two data set of Highway-mpg for Diesel and Gas.

```
In [ ]: gas=df[df['fuel-type']=='gas']['highway-mpg']
```

Chrome File Edit View History Bookmarks Profiles Tab Window Help

Python_part_3_Last Checkpoint 09/12/2022 (unsaved changes)

```
gas[col] = gas[col].str.replace("?", "", regex=False)
df[col] = pd.to_numeric(df[col])
```

```
In [72]: df.describe()
```

```
Out[72]:
```

	city-mpg	highway-mpg
count	205.000000	205.000000
mean	25.219012	30.751220
std	6.542742	6.886443
min	13.000000	16.000000
25%	19.000000	25.000000
50%	24.000000	30.000000
75%	30.000000	34.000000
max	49.000000	54.000000

Question 2: Create two data set of Highway-mpg for Diesel and Gas.

```
In [ ]: gas=df[df['fuel-type']=='gas']['highway-mpg']
print(gas)
diesel=df[df['fuel-type']=='diesel']['highway-mpg']
print(diesel)
print(len(gas))
print(len(diesel))
```

```
In [ ]: dfa = df[df['fuel-type'] == 'gas']
dfb = df[df['fuel-type'] == 'diesel']
```

Question 3: Take 20 random samples from each data set.

Now, we will briefly touch upon how to apply or how to use Python in some of our other tests we did, we will going to talk about hypothesis testing here. If you recall, in hypothesis testing, we wanted to check whether a null hypothesis is true or not. And for that, we come up with a decision criteria. And based on the decision, we accept or reject the null hypothesis. Here I am interested in comparing performance of vehicles, when they use diesel or petrol. So, I think this petrol I do not know they use simply the word gas here. So, maybe we will just say diesel and gas here.

So, this data set so, for to understand this we are going to use a data set which is 1985 Ward's Automotive Yearbook, the column names here specifies what is the data type it is, for example, what make, so, the column are like a descriptions like for example, a column could be make which provides the manufacturer of the vehicle, and the column number of cylinder cues explanation about how number, how many cylinders are there in the engine, and the column peak RPM use what is the maximum revolution per minute like that. So, in this we are going to compare whether the performance of diesel and gas vehicles are going to be the same or not. And we are going to do that using this data here.

So, to play with data, we are going to use this Pandas and use NumPy to see some visualisation we have use matplotlib. And also, to some statistics, we use this scipy stats package and from that we use this ttest particularly. So, let us execute this and then let us read the file. And as we discussed yesterday when I read this file, I will see the first 5 rows with description about the headers and the data in it up to 5 rows, and this info will give me a brief summary of all the columns I have, there are looks like there are about 7 columns here sorry 8 here and they have no null values here.

So, they have all the about the discount 2 naught 5. So, there are 8 columns. And that is saying there are 2 naught 5 entries in here which are indexed from 0 to 2 naught 4. And all these columns have some values, that is why this is all of it is showing 2 naught 5 as a null, non-null values. Now, if you want to check there is any missing values in it.

Like for example, if there is no missing value, either that place would have been left empty or there would I am going to put a question mark. And since it is saying here do not find null values, maybe there is a question mark if there is data missing there. So, let us see this is, what I am trying to do here is I am trying to find in columns checking in each columns, if there is any question mark there.

So, for that, I am we define this variable column and taking it through all possible columns, and under each column I am checking if any of this value is a question marks. To check that we have used this function value count, and telling what is the value we are looking into. And this will give us how many times question mark is appearing in columns and that is printed here. So, notice that this looks like few lines of code, but this could be very compactly written like this in one line. So, in make column there are none, which are question mark, fuel type also there is none, and in horsepower there are 2 big question mark like that.

Now, you will want to remove these missing values. So, 2 missing values, first we identified the columns where the values are missing those are horsepower, peak time, and price. So, we have put the values, replace the question mark by false here. And now if you again can describe, and now everything is value here now we can count the statistics like mean, standard deviations, and all.

(Refer Slide Time: 10:15)

The screenshot shows a Jupyter Notebook interface with the following content:

Question 2 : Create two data set of Highway-mpg for Diesel and Gas.

```
In [73]: gas=df[df['fuel-type']=='gas']['highway-mpg']
print(gas)
diesel=df[df['fuel-type']=='diesel']['highway-mpg']
print(diesel)
print(len(gas))
print(len(diesel))
```

```
0    27
1    27
2    26
3    30
4    22
..
199  22
200  28
201  25
202  23
204  25
Name: highway-mpg, Length: 145, dtype: int64
```

```
In [ ]: dfa = df[df['fuel-type'] == "gas"]
dfb = df[df['fuel-type'] == "diesel"]
```

Question 3 : Take 20 random samples from each data set.

Python_part_3_Last Checkpoint 09/12/2022 (unsaved changes)

```

print(len(gas))
print(len(diesel))

69 25
70 25
90 50
188 33
110 25
112 33
114 25
116 33
158 36
159 47
174 33
182 46
184 46
187 42
192 38
203 27
Name: highway-mpg, dtype: int64
185
20

```

In []: `dfa = df[df['fuel-type'] == "gas"]`
`dfb = df[df['fuel-type'] == "diesel"]`

Question 3: Take 20 random samples from each data set.

In []: `sp.random.seed(42)`
`df1_20 = dfa.sample(n=20)`
`df2_20 = dfb.sample(n=20)`

Question 4: Give a summary of these two data sets. Compare the box-plot of each datasets and interpret your observations

Python_part_3_Last Checkpoint 09/12/2022 (unsaved changes)

```

182 46
184 46
187 42
192 38
203 27
Name: highway-mpg, dtype: int64
185
20

```

In []: `dfa = df[df['fuel-type'] == "gas"]`
`dfb = df[df['fuel-type'] == "diesel"]`

**Question 3 ** Take 20 random samples from each data set.

In [74]: `sp.random.seed(42)`
`df1_20 = dfa.sample(n=20)`
`df2_20 = dfb.sample(n=20)`

Question 4: Give a summary of these two data sets. Compare the box-plot of each datasets and interpret your observations

Summary of dataset GAS

In [75]: `df1_20.describe()`

```

Out[75]:
   city-mpg  highway-mpg
count  20.000000  20.000000
mean    23.500000  29.300000
std     7.189512   7.312629
min    14.000000  16.000000

```

Python_part_3_Last Checkpoint 09/12/2022 (unsaved changes)

In [75]: `df1_20.describe()`

```

Out[75]:
   city-mpg  highway-mpg
count  20.000000  20.000000
mean    23.500000  29.300000
std     7.189512   7.312629
min    14.000000  16.000000
25%    17.000000  23.750000
50%    23.000000  29.000000
75%    26.500000  34.000000
max    30.000000  43.000000

```

Summary of dataset DIESEL

In [43]: `df2_20.describe()`

```

Out[43]:
   horsepower  peak-rpm  city-mpg  highway-mpg  price
count  20.000000  20.000000  20.000000  20.000000  20.000000
mean    84.452000  4426.000000  30.200000  34.750000  15828.150000
std    25.958418  236.571967  6.610199  8.020201  7159.845713
min     52.000000  4190.000000  22.000000  25.000000  7996.000000
25%    62.000000  4187.500000  25.000000  25.000000  9120.000000
50%    84.000000  4425.000000  29.000000  33.000000  13862.000000
75%    97.000000  4700.000000  33.000000  37.000000  16770.000000
max   160.000000  5000.000000  40.000000  43.000000  24000.000000

```

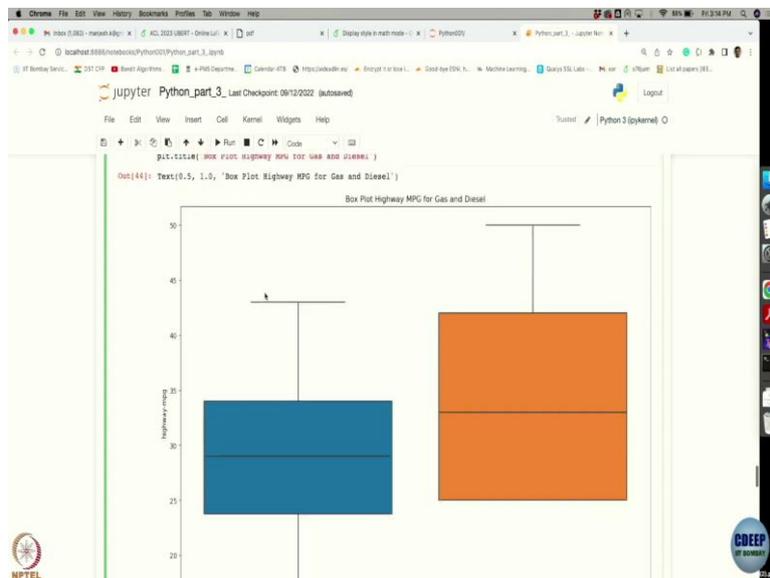
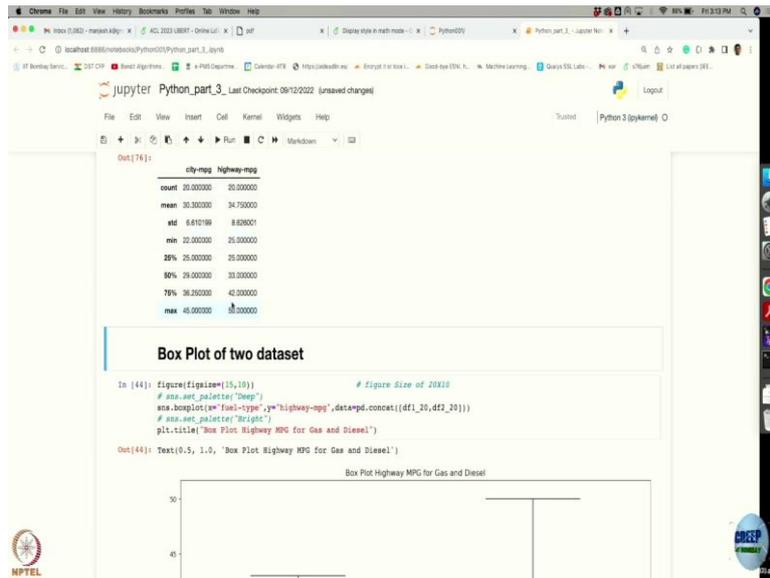
Now, what we want to do is we want to now look into their mileage. So, we want to look into their highway mileage. So, and I want to look into the vehicles which are gas and diesel type here. I have so I will look into all the columns which were all the vehicles which are running on gas and list their highway mileage. So, mpg stands for miles per gallon. And then again for fuel the vehicles which are running on diesel I again list their mileage, highway mileage and let us see, when I do this, I am going to get this, low the first maybe I am printing too much here, first I printed this gas.

So, there looks like a gas there are about 185 element vehicles which are running on gas and the second part corresponding to this print which is showing me that there are about, this just did not show me the length here, but when I look into the length of the diesel, it is showing me that is 20. So, here I have this diesel vehicles, sorry the vehicles running on gas are about 185, and the one which are running on diesel are about 20 here.

Now, I stored them separately in a new data frame, which the gas is stored in a data frame dfa, and the diesel one showed in dfp here. See, notice that storing getting extracting the column information is very easy in Python, all you have to do is first look into all the data where fuel type base gas. And when you pass it to only those indexes to give a data frame, which is storing all that data, it will only return those data's, and not others.

And similarly, here for the diesel case. Now, what I want to do is I want to take 20 sample from each one of these samples and use it to do my hypothesis testing whether diesel and gas have the similar mileage. Now, you will see that when I do the describe both before that, to achieve that 20, getting 20 samples, I am going to use a rand function, first the random function is seeded, so that you get the same value of this randomness, and in the first I will get 20 samples from the dfa, and another 20 samples from dfp. Now, when I look into the summary of this df1 20, you will see that there are 20 elements in that and we are seeing this mean, and median, and highway also has similar values populated here, this is you can also see.

(Refer Slide Time: 13:58)



```
Statistics

In [51]: import numpy as np
import random as rnd
from matplotlib.pyplot import figure
import matplotlib.pyplot as plt
from scipy.stats import norm
from scipy import stats
from scipy.stats import beta as beta_dist
from scipy.stats import beta

Normal Distribution

In [53]: # generate random number of normal distribution
mu, sigma = 0, 0.1 # mean and standard deviation
# = np.random.normal(mu, sigma, 1000)
# Generate 1000 random numbers from normal distribution

Out[53]: array([-7.7162832e-02,  2.9682822e-02,  1.00514257e-01, -2.52380043e-02,
  3.4860019e-02,  1.43518155e-01,  5.84402989e-02,  3.55503035e-02,
 -1.68470379e-01, -9.3372223e-02,  1.04695073e-01,  5.69503035e-02,
  7.21468274e-02,  6.35396932e-02,  1.28976653e-01, -2.61399251e-02,
  1.20499959e-01, -1.02867433e-01, -1.27802238e-01, -1.64374072e-01,
 -8.85183494e-02,  1.55261612e-01, -1.55272070e-01, -3.10535961e-02,
 -1.29314532e-01, -4.00973963e-02,  1.69693792e-02, -7.01158731e-02,
  1.08964874e-01,  1.53397812e-01, -2.24899573e-01, -6.82923202e-02,
  5.91421223e-03,  1.12611793e-01, -4.2424275e-02, -1.99353952e-03,
 -1.15744226e-01, -1.53279851e-01,  1.46538489e-01, -1.82406807e-01,
  2.32861901e-01,  7.25513399e-02,  6.79972076e-02,  5.75468667e-02,
```

Observation :

The box plot is comparatively short – see example for Gas. highway MPG for gas and diesel. Similarly, the box plot is comparatively tall – see examples of Diesel. This suggests highway MPG for gas and diesel hold quite different.

Here, we can see that if the null hypothesis is rejected, it indicates that the average mileage of the two data sets is equal, and if it is accepted, it indicates the same thing. Here, we choose a random sample from both sets of 20-person data, which can sometimes reject and sometimes accept the null hypothesis.

Question 5 : Test whether the average mileage of both the data sets are equal or not at level 5% of significance. Explain the procedure completely and interpret your results.

```
In [50]: stat, p = ttest_ind(df1_20['highway-mpg'], df2_20['highway-mpg'], equal_var=False)
print('t-statistic=%4.3f, p-value=%4.3f' % (stat, p))
if p > 0.05:
    print('Accept H0: average mileage of both data sets are equal')
else:
    print('Reject H0: average mileage of both data sets are not equal')
t-statistic=-2.155, p-value=0.038
Reject H0: average mileage of both data sets are not equal
```

Question 6 : Consider the whole data sets created in Step 2 (Here the no. of samples are not equal).

```
In [44]: dfa = df[df['fuel-type'] == "gas"]
dfb = df[df['fuel-type'] == "diesel"]
```

```
In [47]: dfa
```

```
Out[47]:
```

	make	fuel-type	num-of-cylinders	horsepower	peak-rpm	city-mpg	highway-mpg	price
0	alfa-romeo	gas	four	111	5000	21	27	13495
1	alfa-romeo	gas	four	111	6000	21	27	16000

Gas has a higher line in the middle of the box plot than diesel, indicating that it had a higher median value.

- In what ways does the distribution differ?
- Diesel has a significantly larger spread out than gas, as evidenced by the lengthier box plot for diesel.
- Are there any outliers?

Since neither dataset contained any obvious outliers, neither box plot features tiny circles that extend past the top or bottom whiskers.

Question 8 : Test at level 5% of significance that the average mileage of both the data sets are equal or not. Explain the procedure completely and interpret your results.

```
In [43]: stat, p = ttest_ind(dfa['highway-mpg'], dfb['highway-mpg'], equal_var=False)
print('stat=%4.3f, p=%4.3f' % (stat, p))
if p > 0.05:
    print("null hypothesis accepted")
    print("The average mileage of both the data sets are equal")
else:
    print("null hypothesis rejected")
    print("The average mileage of both the data sets are not equal")

stat=-2.228, p=0.037
null hypothesis rejected
The average mileage of both the data sets are not equal
```

Now, often these are like just like statistic count, maybe one can want to use it pictorially for that one can use the box plot here, the box plot to do this, we are going to use box plot function in this sns library, and the plotting, for that we are going to use this command so x is simply going to be the fuel type. We have to just say x is the string is fuel type, y is highway mpg, so highway mpg here, and x is fuel type here. And the data we are going to give you is this concatenated value of df1 and df2.

Now, when I plot this, you see this 1, 2, 3, 4, 5 end up us. The bottom one here corresponding to the lowest value that is the minimum value, which is 22 here, and the middle one, so, the second one here is going to correspond to the 25th quartile of the data and the middle one is the 50th quartile, which is also the median, and this is like a 75th quartile, and the last one here is going to represent the maximum value.

So, now we have named this figure as box plot highway mpg for gas and diesel. Now, based on this you see that if you compare this just take the mean, it gives a sense that diesel has a better efficiency, has better mileage. Now, we want to test this. Now, our hypothesis is we want to test the hypothesis our average mileage or both of them are same, and my alternate hypothesis is they are not.

Now, I want to test it, and I want to use the t test for it. So, if you recall in hypothesis testing, the accept or reject by comparing to a threshold. And now we want to first compute the statistics and the p value of that statistics. So, to compute the t states, I am going to use this t test ind, which is the t test we had already defined a t test, we did we define t test here before, no.

So, this t test is a function which takes the 2 data sets that we want to compare, and here equal variance tells whether I should be taking the variance of these 2 tests to be the same or not, here we have set it false. So, this will not take this data to be the same variance. And recall that when we are talking t tests, it is assumed that the data is following a quotient distribution.

Now, this function is going to return me 2 value stat and p, stat will be the value of the statistic that we compute and p is its probability that statistics is going to be larger than some number if I by default, that is like a set the significance level is set to be .05. So, we see that now p we compare against the significance level 0.05, if it is true, we say that null hypothesis accepted. And otherwise, we are going to say that it is rejected.

Now, it so happens that in this example, p value turns out to be minus 2.229 and whereas the p value is 0.037. And clearly 0.307 is going to be smaller than 0.05. So, we are going to see that the null hypothesis is going to be rejected, that is the performance of these 2 are not the same. Like this, you can do here I just showed you that, let me see why is this happening at some more things here.

(Refer Slide Time: 18:47)

```
Question 5: Test whether the average mileage of both the data sets are equal or not at level 5% of significance. Explain the procedure completely and interpret your results.
```

```
In [50]: stat, p = ttest_ind(df1_20['highway-mpg'], df2_20['highway-mpg'], equal_var=False)
print('t-statistic={:.3f}, p-value={:.3f}'.format(stat, p))
if p > 0.05:
    print('Accept H0: average mileage of both data sets are equal')
else:
    print('Reject H0: average mileage of both data sets are not equal')

t-statistic=-2.155, p-value=0.039
Reject H0: average mileage of both data sets are not equal
```

```
Question 6: Consider the whole data sets created in Step 2 (Here the no. of samples are not equal).
```

```
In [44]: dfa = df[df['fuel-type'] == "gas"]
dfb = df[df['fuel-type'] == "diesel"]

In [47]: dfa
Out[47]:
```

make	fuel-type	num-of-cylinders	horsepower	peak-rpm	city-mpg	highway-mpg	price
0	alfa-romeo	gas	four	111	5000	21	27 13495
1	alfa-romeo	gas	four	111	5000	21	27 16500
2	alfa-romeo	gas	six	154	5000	19	26 16500
3	audi	gas	four	102	5500	24	30 13860
4	audi	gas	five	115	5500	18	22 17450
...
199	volvo	gas	four	162	5100	17	22 18900
200	volvo	gas	four	114	5400	23	28 16845

```
well. The median values are compared in 1.
Gas has a higher line in the middle of the box plot than diesel, indicating that it had a higher median value.
2. In what ways does the distribution differ?
Diesel has a significantly larger spread out than gas, as evidenced by the lengthier box plot for diesel.
3. Are there any outliers?
Since neither dataset contained any obvious outliers, neither box plot features tiny circles that extend past the top or bottom whiskers.
```

```
Question 8: Test at level 5% of significance that the average mileage of both the data sets are equal or not. Explain the procedure completely and interpret your results.
```

```
In [43]: stat, p = ttest_ind(dfa['highway-mpg'], dfb['highway-mpg'], equal_var=False)
print('stat={:.3f}, p={:.3f}'.format(stat, p))
if p > 0.05:
    print('null hypothesis accepted')
    print('The average mileage of both the data sets are equal')
else:
    print('null hypothesis rejected')
    print('The average mileage of both the data sets are not equal')

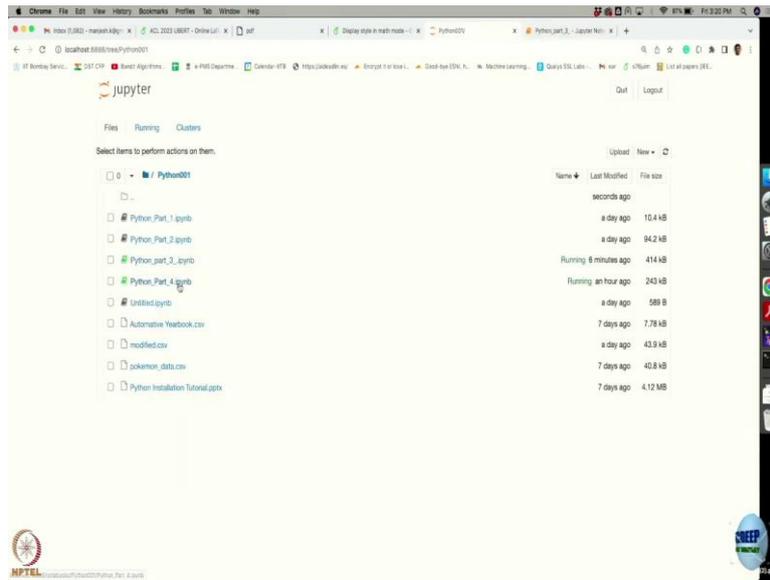
stat=-2.228, p=0.037
null hypothesis rejected
The average mileage of both the data sets are not equal
```

```
In [ ]:
```

So, here based on this, you can, here with the 20 random samples generated, you conclude that, yes, this your null hypothesis that both are the same performances is rejected. Now, you can do this again. But this time, instead of not necessarily taking 20 samples for both, maybe you can take all the possible samples available for each of these and compare.

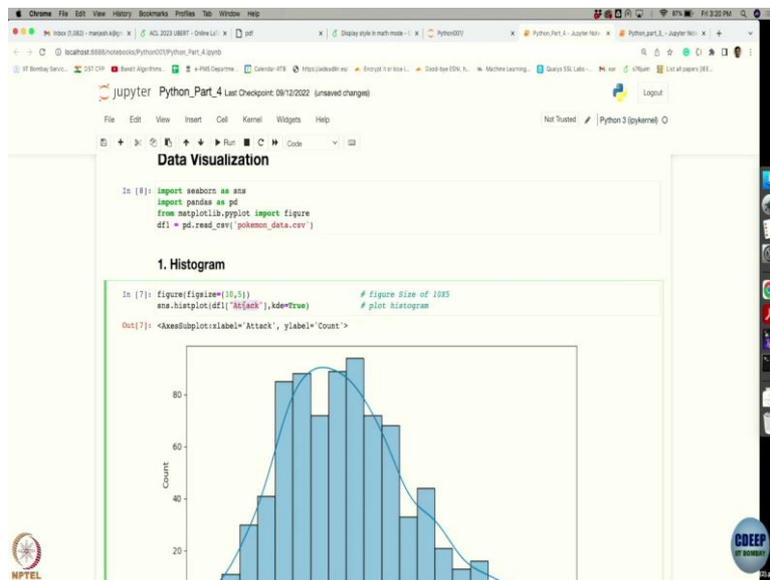
And you can do the same process here. I am not going to go into the details, but again, to get a sense here, you can first put the blocks plot. And again, here when I take all the data also like looks like again, diesel is performing better. And then you can try to test this and when you do this test, at significance level .05 again, the null hypothesis is rejected here. So, you see that like before we applying this hypothesis test, sometimes the visualisation gives more idea about how the data looks like.

(Refer Slide Time: 19:51)



The screenshot shows the JupyterLab interface with the file browser view. The file list includes:

Name	Last Modified	File size
..	seconds ago	
Python_Part_1.ipynb	a day ago	10.4 kB
Python_Part_2.ipynb	a day ago	94.2 kB
Python_part_3.ipynb	Running 8 minutes ago	414 kB
Python_Part_4.ipynb	Running an hour ago	243 kB
Untitled.ipynb	a day ago	589 B
Automotive Yearbook.csv	7 days ago	7.79 kB
modified.csv	a day ago	43.9 kB
pokemon_data.csv	7 days ago	40.8 kB
Python Installation Tutorial.pdf	7 days ago	4.12 MB



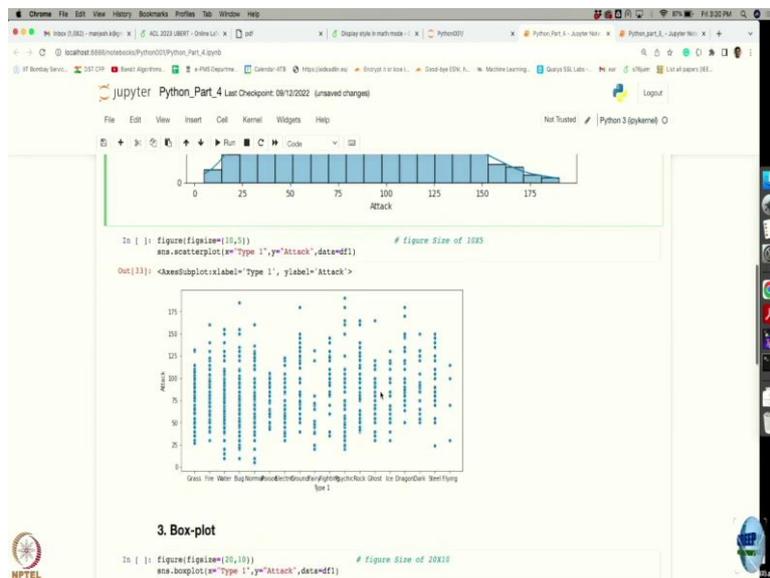
The screenshot shows the JupyterLab code editor with the following code and output:

```
In [8]: import seaborn as sns
import pandas as pd
from matplotlib.pyplot import figure
df1 = pd.read_csv('pokemon_data.csv')
```

1. Histogram

```
In [7]: figure(figsize=(10,5))
sns.histplot(df1['Attack'], kde=True)
Out[7]: <AxesSubplot:label='Attack', ylabel='Count'>
```

The output is a histogram showing the distribution of the 'Attack' attribute. The x-axis is labeled 'Attack' and the y-axis is labeled 'Count'. The histogram bars are blue, and a kernel density estimate (kde) is overlaid in a light blue line.



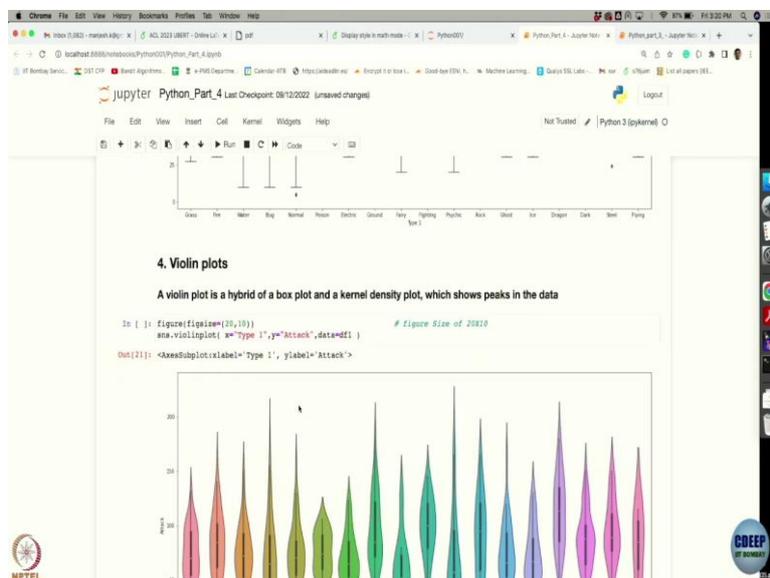
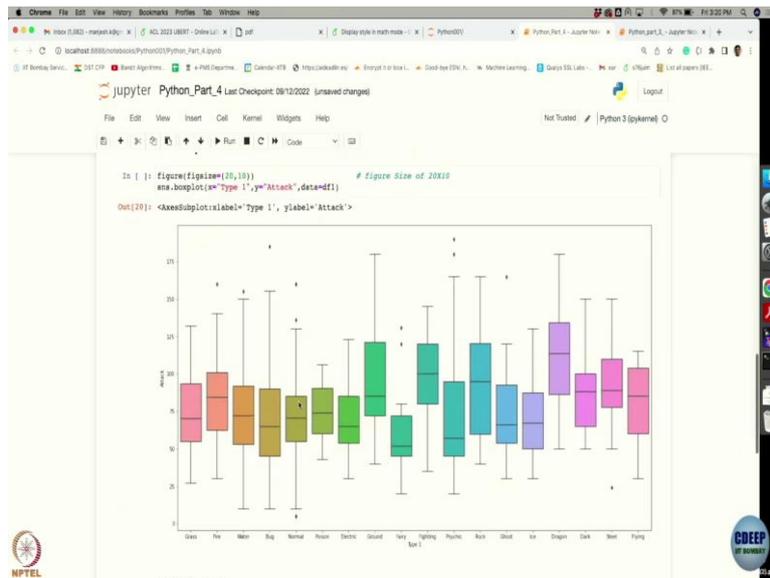
The screenshot shows the JupyterLab code editor with the following code and output:

```
In [1]: figure(figsize=(10,5))
sns.scatterplot(x='Type 1', y='Attack', data=df1)
Out[33]: <AxesSubplot:label='Type 1', ylabel='Attack'>
```

The output is a scatter plot showing the relationship between 'Type 1' and 'Attack'. The x-axis is labeled 'Type 1' and the y-axis is labeled 'Attack'. The data points are blue dots.

```
In [1]: figure(figsize=(10,10))
sns.boxplot(x='Type 1', y='Attack', data=df1)
```

The output is a box plot showing the distribution of 'Attack' for each 'Type 1'. The x-axis is labeled 'Type 1' and the y-axis is labeled 'Attack'. The box plots are blue.



So, I want to quickly spend one a minute on this various visualisation tools. So, that is why data visualisation is important and this seaborn tool provides some pretty good functions to do that. So, one function which we already used is histogram, and when you have a data, you can see that like, for example, earlier I do this Pokémon data. In the Pokémon data, let us say I want to see how this attack column looks like I can generate the histogram like this.

And also, there are other plots called scatter plots, box plot we just saw and there is another things violin plots, these are other good visualisation techniques. I encourage all of you to play with them. So, with that, we will stop. But this Python is a very comes with a very rich libraries. So, you can readily use it to do good statistical analysis. I hope whatever we discuss is going to be helpful to you.