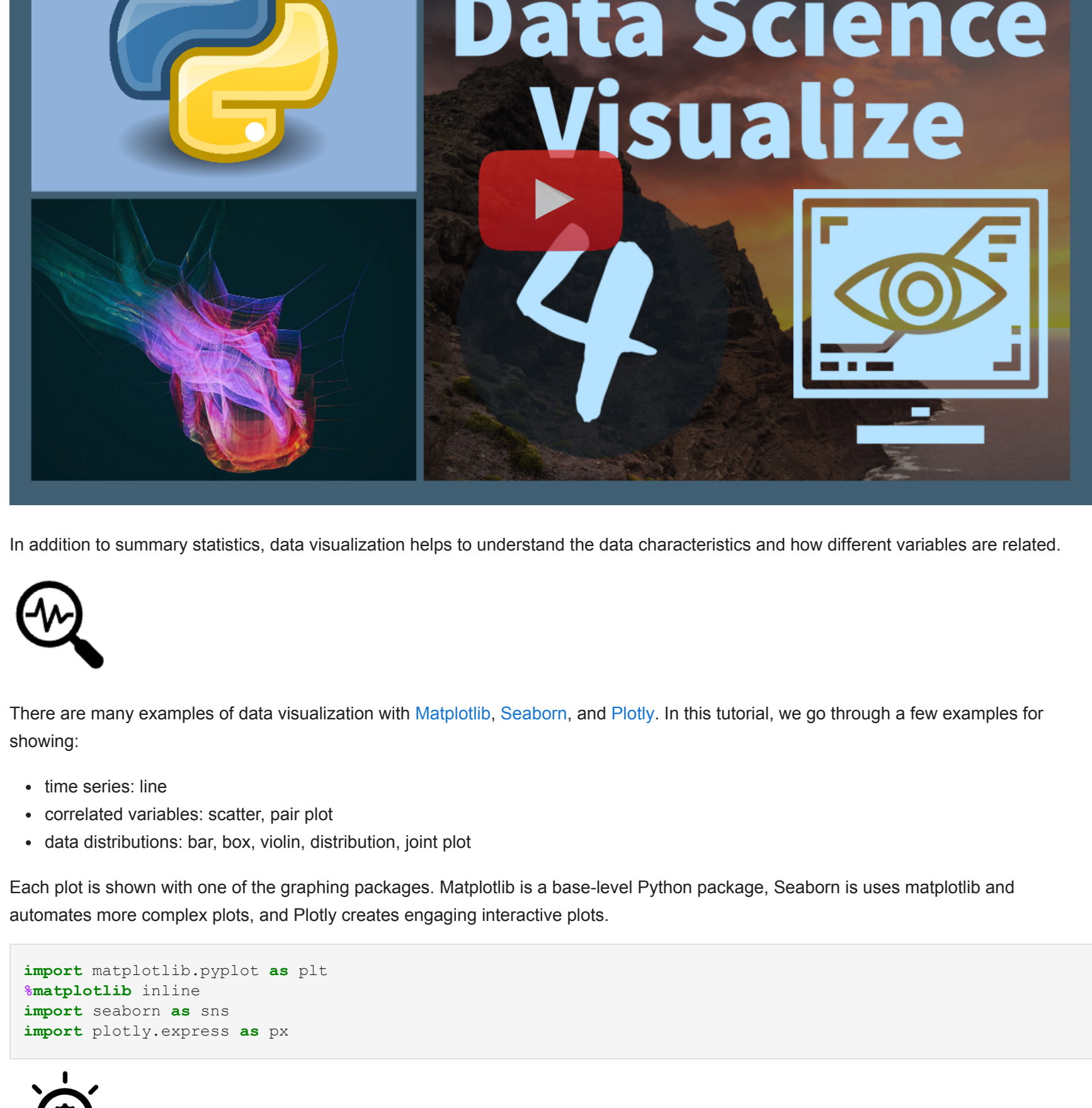
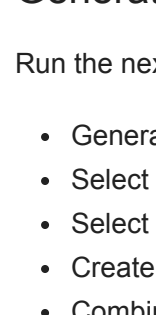


4. Visualize Data

Data Science Playlist on YouTube



In addition to summary statistics, data visualization helps to understand the data characteristics and how different variables are related.



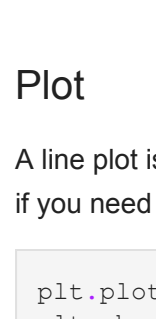
There are many examples of data visualization with **Matplotlib**, **Seaborn**, and **Plotly**. In this tutorial, we go through a few examples for showing:

- time series: line
- correlated variables: scatter, pair plot
- data distributions: bar, box, violin, distribution, joint plot

Each plot is shown with one of the graphing packages. Matplotlib is a base-level Python package, Seaborn is uses matplotlib and automates more complex plots, and Plotly creates engaging interactive plots.

In [2]:

```
import matplotlib.pyplot as plt
import pandas as pd
import numpy as np
import random
import time
import sys
import os
import math
import datetime
import json
import csv
import pickle
import pickle5
import pickletools
import picklecloud
import pickles
import pickles3
import pickles4
import pickles5
import pickles6
import pickles7
import pickles8
import pickles9
import pickles10
import pickles11
import pickles12
import pickles13
import pickles14
import pickles15
import pickles16
import pickles17
import pickles18
import pickles19
import pickles20
import pickles21
import pickles22
import pickles23
import pickles24
import pickles25
import pickles26
import pickles27
import pickles28
import pickles29
import pickles30
import pickles31
import pickles32
import pickles33
import pickles34
import pickles35
import pickles36
import pickles37
import pickles38
import pickles39
import pickles40
import pickles41
import pickles42
import pickles43
import pickles44
import pickles45
import pickles46
import pickles47
import pickles48
import pickles49
import pickles50
import pickles51
import pickles52
import pickles53
import pickles54
import pickles55
import pickles56
import pickles57
import pickles58
import pickles59
import pickles60
import pickles61
import pickles62
import pickles63
import pickles64
import pickles65
import pickles66
import pickles67
import pickles68
import pickles69
import pickles70
import pickles71
import pickles72
import pickles73
import pickles74
import pickles75
import pickles76
import pickles77
import pickles78
import pickles79
import pickles80
import pickles81
import pickles82
import pickles83
import pickles84
import pickles85
import pickles86
import pickles87
import pickles88
import pickles89
import pickles90
import pickles91
import pickles92
import pickles93
import pickles94
import pickles95
import pickles96
import pickles97
import pickles98
import pickles99
import pickles100
```



Generate Data

Run the next cell to:

- Generate n linearly spaced values between 0 and $n-1$ with `np.linspace(start,end,count)`
- Select random samples from a uniform distribution between 0 and 1 with `np.random.rand(count)`
- Select random samples from a normal (Gaussian) distribution with `np.random.normal(mean,std,count)`
- Create a time series that changes based on `y[i]*0.1` staying within the range -3 to 3
- Combine `tt`, `x`, `y`, and `z` with a vertical stack `np.vstack` and transpose `.T` for column oriented data
- Create pandas DataFrame with columns `tt`, `x`, `y`, and `z`

```
import numpy as np
import pandas as pd
np.random.seed(0) # change seed for different answer
n = 1000
tt = np.linspace(0,n-1,n)
x = np.random.rand(n)*tt/500
y = np.random.normal(0,x,n)
z = [0]
for i in range(1,n):
    z.append(min(max(-3,z[i-1]*y[i]*0.1),3))
data = pd.DataFrame(np.vstack((tt,x,y,z)).T,
                    columns=['time','x','y','z'])
data['w'] = 0.499
for i in range(int(n/2),n):
    data.at[i,'w'] = '500-999'
data.head()
```

```
Out[3]:
   time    x         y         z    w
0    0.0  0.00014 -0.05813  0.000000  0.499
1    1.0  0.00189  0.01327  0.001383  0.499
2    2.0  0.006763  0.12284  0.13609  0.499
3    3.0  0.050883 -0.117891  0.101811  0.499
4    4.0  0.431655 -0.215403  0.080271  0.499
```

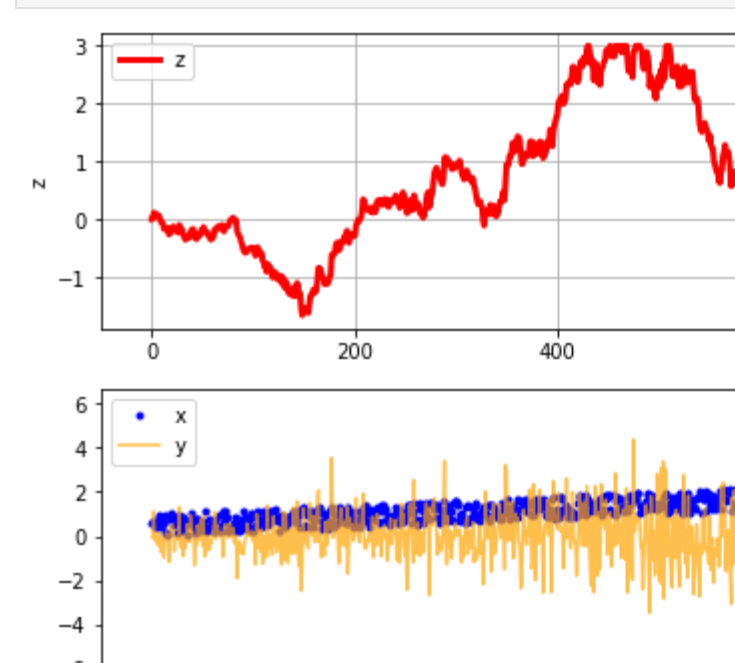


Plot

A line plot is the most basic type. There is an introductory tutorial on plots in the [Begin Python Course, Lesson 12](#). Visit that course module if you need additional information on basic plots such as `plt.plot()`

In [4]:

```
plt.plot(tt,z)
plt.show()
```



The line plot can also be improved with customized trend styles. Below is an example with common options.

cColors

```
=====
character      color
-----
'b'           blue
'g'           green
'r'           red
'y'           yellow
'k'           black
=====
```

mMarkers

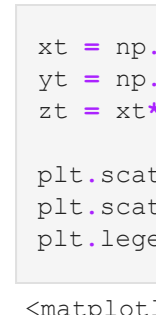
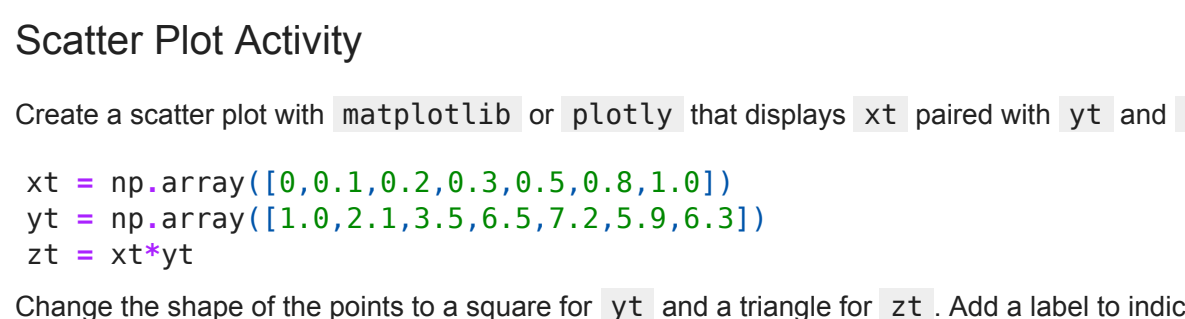
```
=====
character      description
-----
'.'            point marker
'o'            circle marker
's'            square marker
'^'            triangle marker
'*'            star marker
=====
```

lLine Styles

```
=====
character      description
-----
'-'            solid line style
'--'           dashed line style
'-.'           dash-dot line style
':'           dotted line style
=====
```

In [5]:

```
plt.figure(figsize=(10,6))
ax=plt.subplot(2,1,1)
plt.plot(tt,z,'-',linewidth=3,label='z')
ax.grid()
plt.ylabel('z'); plt.legend()
plt.subplot(2,1,2)
plt.plot(tt,x,'-',label='x')
plt.plot(tt,y,color='orange',label='y',alpha=0.7)
plt.savefig('testfig.png',transparent=True,dpi=600)
plt.show()
```



Plot Activity

Create a plot that displays the data:

```
xt = [0,0.1,0.2,0.3,0.5,0.8,1.0]
yt = [1,0.2,1,3,5,6,5,7,2,5,9,6,3]
```

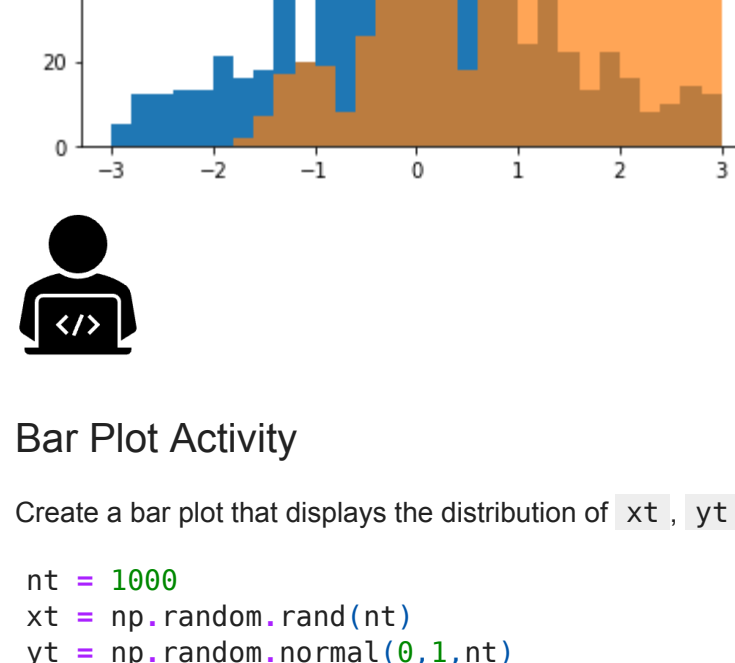


Scatter Plot

Scatter plots are similar to regular plots but they show individuals points instead of values connected in series. Matplotlib and Plotly are used in this example. Matplotlib is fast and simple while Plotly has features for interactive plots.

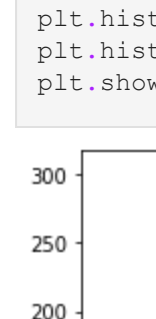
In [6]:

```
# matplotlib
plt.scatter(xt,yt)
plt.show()
```



In [7]:

```
# plotly
fig = px.scatter(data,x='xt',y='yt',color='w',size='x',hover_data=['w'])
fig.show()
```



Scatter Plot Activity

Create a scatter plot with `matplotlib` or `plotly` that displays `xt` paired with `yt` and `zt`:

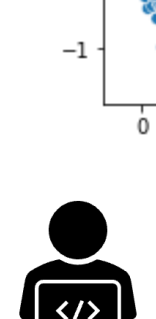
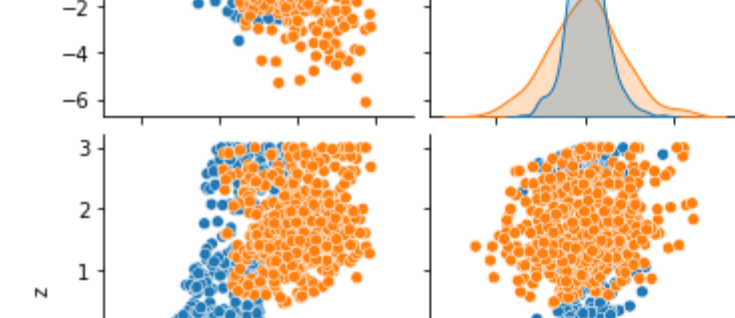
```
xt = np.array([0,0.1,0.2,0.3,0.5,0.8,1.0])
yt = np.array([1,0.2,1,3,5,6,5,7,2,5,9,6,3])
zt = xt*yt
```

Change the shape of the points to a square for `yt` and a triangle for `zt`. Add a label to indicate which points are `yt` and `zt`.

In [16]:

```
xt = np.array([0,0.1,0.2,0.3,0.5,0.8,1.0])
yt = np.array([1,0.2,1,3,5,6,5,7,2,5,9,6,3])
zt = xt*yt
plt.scatter(xt,yt,marker='s',label='y')
plt.scatter(xt,zt,marker='^',label='z')
plt.legend()
```

Out[16]:

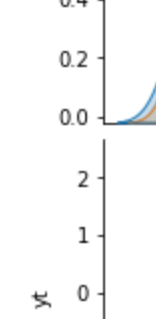
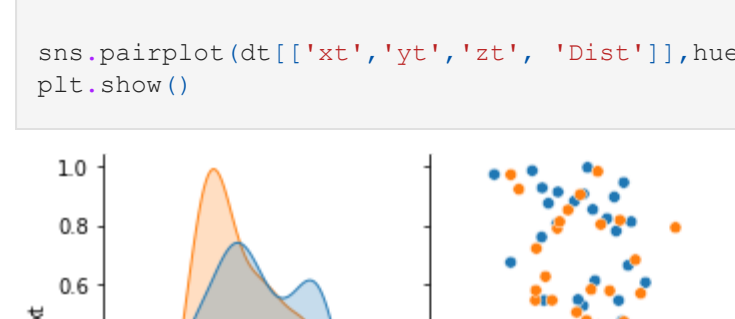


Bar Chart

Bar charts show a histogram distribution of count in a bin range. The `alpha` option is the transparency between 0 and 1. A value of 0.7 is a good value to use to show the overlying and underlying data.

In [17]:

```
bins = np.linspace(-3,3,31)
plt.hist(yt,bins,label='y')
plt.hist(zt,bins,alpha=0.7,label='z')
plt.legend()
plt.show()
```



Bar Plot Activity

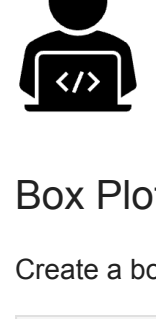
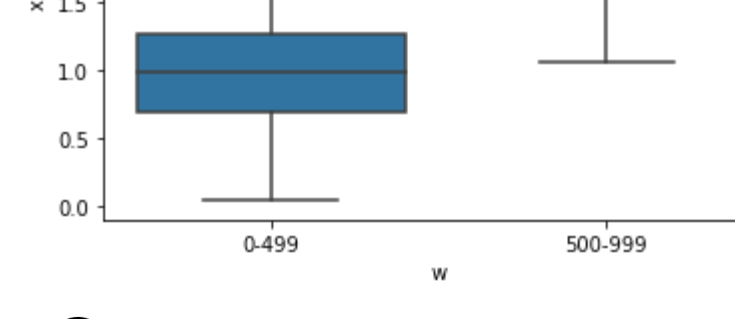
Create a bar plot that displays the distribution of `xt`, `yt`, and `zt`:

```
nt = 1000
xt = np.random.rand(nt)
yt = np.random.normal(0,1,nt)
zt = xt*yt
```

Use `bins = np.linspace(-3,3,31)` to create the histogram distribution.

In [21]:

```
nt = 1000
xt = np.random.rand(nt)
yt = np.random.normal(0,1,nt)
zt = xt*yt
bins = np.linspace(-3,3,31)
plt.hist(xt,bins)
plt.hist(yt,bins)
plt.hist(zt,bins)
plt.show()
```

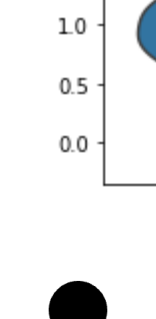
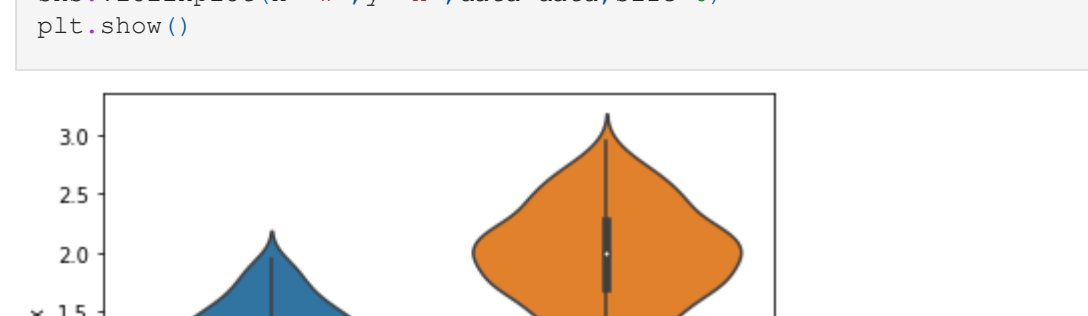


Pair Plot

A pair plot shows the correlation between variables. It has bar distributions on the diagonal and scatter plots on the off-diagonal. A pair plot also shows a different color (hue) by category `w`. Pair plots show correlations between pairs of variables that may be related and gives a good indication of features (explanatory inputs) that are used for classification or regression.

In [22]:

```
sns.pairplot(data[['xt','yt','zt','w']],hue='w')
plt.show()
```



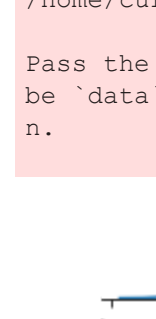
Pair Plot Activity

Create a pair plot that displays the correlation between `xt`, `yt`, and `zt` between the first 500 and second 500 random numbers that are categorized as `Dist`. Create a `pandas` dataframe with:

```
nt = 100
xt = np.random.rand(nt)
yt = np.random.normal(0,1,nt)
zt = xt*yt
dt = pd.DataFrame(np.column_stack([xt,yt,zt]),columns=['xt','yt','zt'])
dt['Dist'] = 'First'
for i in range(int(nt/2),nt):
    dt.at[i,'Dist'] = 'Second'
```

In [27]:

```
nt = 100
xt = np.random.rand(nt)
yt = np.random.normal(0,1,nt)
zt = xt*yt
dt = pd.DataFrame(np.column_stack([xt,yt,zt]),columns=['xt','yt','zt'])
dt['Dist'] = 'First'
for i in range(int(nt/2),nt):
    dt.at[i,'Dist'] = 'Second'
sns.pairplot(dt[['xt','yt','zt','Dist']],hue='Dist')
plt.show()
```

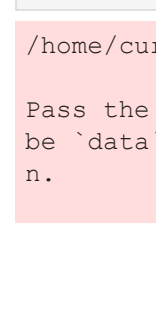


Box Plot

A box plot shows data quartiles. In this case, we are comparing the first 500 points with the last 500 points.

In [28]:

```
sns.boxplot(x='w',y='xt',data=data)
plt.show()
```

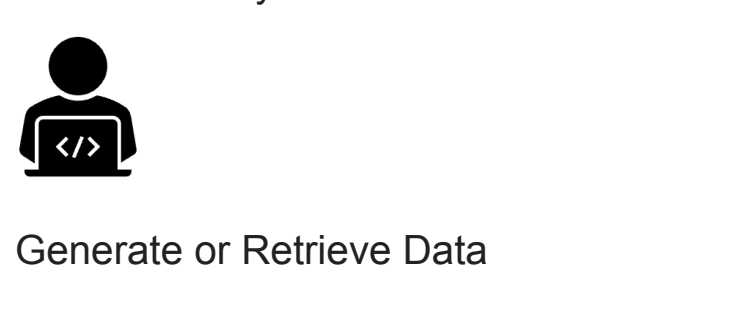


Box Plot Activity

Create a box plot that shows the quartiles of `yt` by first and second sets as indicated in `Dist`.

In [29]:

```
sns.boxplot(x='Dist',y='yt',data=dt)
plt.show()
```

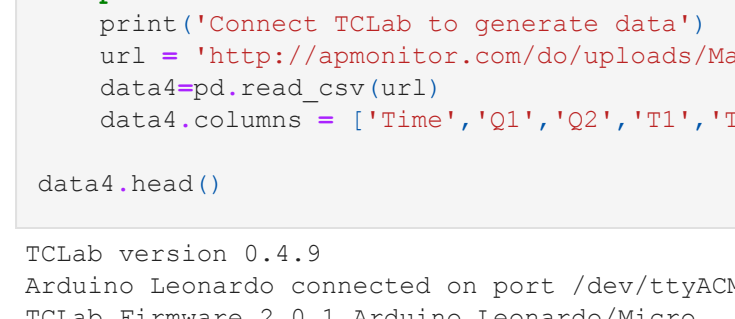


Violin Plot

A violin plot combines the box plot quartiles with the distribution.

In [30]:

```
sns.violinplot(x='w',y='xt',data=data,size=6)
plt.show()
```

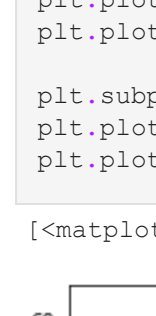
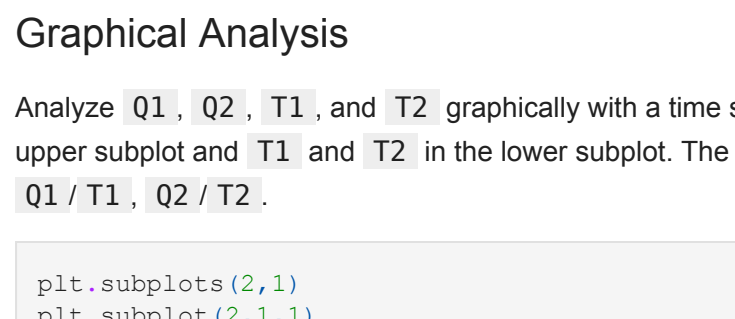


Violin Plot Activity

Create a violin plot that shows the quartiles and distribution of `zt` by first and second sets as indicated in `Dist` in the DataFrame `dt`.

In [31]:

```
sns.violinplot(x='Dist',y='zt',data=dt,size=6)
plt.show()
```



Joint Plot

A joint plot shows two variables, with the univariate and joint distributions. Try `kind='reg'`, `'kde'`, and `'hex'` to see different joint plot styles.

In [32]:

```
sns.jointplot('xt','zt',data=data,kind='kde')
plt.show()
```

/home/curtis/.local/lib/python3.8/site-packages/seaborn/_decorators.py:36: FutureWarning:
Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.

Joint Plot Activity

Create a joint plot that shows the joint distribution of `yt` and `zt` in the DataFrame `dt`.

In [33]:

```
sns.jointplot('yt','zt',data=dt,kind='kde')
plt.show()
```

Pass the following variables as keyword args: x, y. From version 0.12, the only valid positional argument will be 'data', and passing other arguments without an explicit keyword will result in an error or misinterpretation.

TCLab Activity

Generate or Retrieve Data

A sample data file loads if you do not have a TCLab connected. Otherwise, generate a file from the TCLab data with seconds (`T`), heater levels (`Q1` and `Q2`), and temperatures (`T1` and `T2`). Record data every second for 120 seconds and change the heater levels every 30 seconds to a random number between 0 and 80 with `np.random.randint()`. There is no need to change this program, only run it to collect the data over 2 minutes.

In [34]:

```
import tclab, time, csv
import numpy as np
try:
    n = 120
    wth=open('04-tclab.csv',mode='w',newline='') as f:
        cw = csv.writer(f)
        cw.writerow(['time','Q1','Q2','T1','T2'])
        with tclab.TCLab() as lab:
            for i in range(n):
                if i%30==0:
                    Q1 = np.random.randint(0,81)
                    Q2 = np.random.randint(0,81)
                    lab.Q1(Q1); lab.Q2(Q2)
                    cw.writerow([i,Q1,Q2,lab.T1,lab.T2])
                    if i%5==0:
                        print(i,Q1,Q2,lab.T1,lab.T2)
                time.sleep(1)
            data4=pd.read_csv('04-tclab.csv')
except:
    print('Connect TCLab to generate data')
    url = 'http://apmonitor.com/do/uploads/Main/tclab_dyn_data2.txt'
    data4pd=pd.read_csv(url)
    data4.columns = ['Time','Q1','Q2','T1','T2']
data4.head()
```

TCLab version 0.4.9
Arduino Leonardo connected on port /dev/ttyACM1 at 115200 baud.
TCLab firmware 2.0.1: Arduino Leonardo/Micro.
t Q1 Q2 T1 T2
0 48 53 23.477 22.893
1 48 53 23.541 23.09
2 48 53 24.121 23.509
3 48 53 24.093 23.696
4 48 53 24.121 23.509
5 48 53 26.377 24.508
6 48 53 26.344 25.249
7 48 53 26.377 24.508
8 48 53 26.344 25.249
9 48 53 26.377 24.508
10 48 53 26.344 25.249
11 48 53 26.377 24.508
12 48 53 26.344 25.249
13 48 53 26.377 24.508
14 48 53 26.344 25.249
15 48 53 26.377 24.508
16 48 53 26.344 25.249
17 48 53 26.377 24.508
18 48 53 26.344 25.249
19 48 53 26.377 24.508
20 48 53 26.344 25.249
21 48 53 26.377 24.508
22 48 53 26.344 25.249
23 48 53 26.377 24.508
24 48 53 26.344 25.249
25 48 53 26.377 24.508
26 48 53 26.344 25.249
27 48 53 26.377 24.508
28 48 53 26.344 25.249
29 48 53 26.377 24.508
30 48 53 26.344 25.249
31 48 53 26.377 24.508
32 48 53 26.344 25.249
33 48 53 26.377 24.508
34 48 53 26.344 25.249
35 48 53 26.377 24.508
36 48 53 26.344 25.249
37 48 53 26.377 24.508
38 48 53 26.344 25.249
39 48 53 26.377 24.508
40 48 53 26.344 25.249
41 48 53 26.377 24.508
42 48 53 26.344 25.249
43 48 53 26.377 24.508
44 48 53 26.344 25.249
45 48 53 26.377 24.508
46 48 53 26.344 25.249
47 48 53 26.377 24.508
48 48 53 26.344 25.249
49 48 53 26.377 24.508
50 48 53 26.344 25.249
51 48 53 26.377 24.508
52 48 53 26.344 25.249
53 48 53 26.377 24.508
54 48 53 26.344 25.249
55 48 53 26.377 24.508
56 48 53 26.344 25.249
57 48 53 26.377 24.508
58 48 53 26.344 25.249
59 48 53 26.377 24.508
60 48 53 26.344 25.249
61 48 53 26.377 24.508
62 48 53 26.344 25.249
63 48 53 26.377 24.508
64 48 53 26.344 25.249
65 48 53 26.377 24.508
66 48 53 26.344 25.249
67 48 53 26.377 24.508
68 48 53 26.344 25.249
69 48 53 26.377 24.508
70 48 53 26.344 25.249
71 48 53 26.377 24.508
72 48 53 26.344 25.249
73 48 53 26.377 24.508
74 48 53 26.344 25.249
75 48 53 26.377 24.508
76 48 53 26.344 25.249
77 48 53 26.377 24.508
78 48 53 26.344 25.249
79 48 53 26.377 24.508
80 48 53 26.344 25.249
81 48 53 26.377 24.508
82 48 53 26.344 25.249
83 48 53 26.377 24.508
84 48 53 26.344 25.249
85 48 53 26.377 24.508
86 48 53 26.344 25.249
87 48 53 26.377 24.508
88 48 53 26.344 25.249
89 48 53 26.377 24.508
90 48 53 26.344 25.249
91 48 53 26.377 24.508
92 48 53 26.344 25.249
93 48 53 26.377 24.508
94 48 53 26.344 25.249
95 48 53 26.377 24.508
96 48 53 26.344 25.249
97 48 53 26.377 24.508
98 48 53 26.344 25.249
99 48 53 26.377 24.508
100 48 53 26.344 25.249
101 48 53 26.377 24.508
102 48 53 26.344 25.249
103 48 53 26.377 24.508
104 48 53 26.344 25.249
105 48 53 26.377 24.508
106 48 53 26.344 25.249
107 48 53 26.377 24.508
108 48 53 26.344 25.249
109 48 53 26.377 24.508
110 48 53 26.344 25.249
111 48 53 26.377 24.508
112 48 53 26.344 25.249
113 48 53 26.377 24.508
114 48 53 26.344 25.249
115 48 53 26.377 24.508
116 48 53 26.344 25.249
117 48 53 26.377 24.508
118 48 53 26.344 25.249
119 48 53 26.377 24.508
120 48 53 26.344 25.249
TCLab disconnected successfully.

```
Out[34]:
   Time  Q1  Q2  T1  T2
0    0  48  53  23.477 22.896
1    1  48  53  23.477 22.832
2    2  48  53  23.444 23.025
3    3  48  53  23.444 22.735
4    4  48  53  23.509 23.025
```

Graphical Analysis

Analyze `Q1`, `Q2`, `T1`, and `T2` graphically with a time series plot and a pair plot. The time series plot should show `Q1` and `Q`

