

Analyzing out-of-this world data

Using data collected from the Open Exoplanet Catalogue database:

https://github.com/OpenExoplanetCatalogue/open_exoplanet_catalogue/
(https://github.com/OpenExoplanetCatalogue/open_exoplanet_catalogue/).

Data License

Copyright (C) 2012 Hanno Rein

Permission is hereby granted, free of charge, to any person obtaining a copy of this database and associated scripts (the "Database"), to deal in the Database without restriction, including without limitation the rights to use, copy, modify, merge, publish, distribute, sublicense, and/or sell copies of the Database, and to permit persons to whom the Database is furnished to do so, subject to the following conditions:

The above copyright notice and this permission notice shall be included in all copies or substantial portions of the Database. A reference to the Database shall be included in all scientific publications that make use of the Database.

THE DATABASE IS PROVIDED "AS IS", WITHOUT WARRANTY OF ANY KIND, EXPRESS OR IMPLIED, INCLUDING BUT NOT LIMITED TO THE WARRANTIES OF MERCHANTABILITY, FITNESS FOR A PARTICULAR PURPOSE AND NONINFRINGEMENT. IN NO EVENT SHALL THE AUTHORS OR COPYRIGHT HOLDERS BE LIABLE FOR ANY CLAIM, DAMAGES OR OTHER LIABILITY, WHETHER IN AN ACTION OF CONTRACT, TORT OR OTHERWISE, ARISING FROM, OUT OF OR IN CONNECTION WITH THE DATABASE OR THE USE OR OTHER DEALINGS IN THE DATABASE.

Setup

```
In [1]: %matplotlib inline

import matplotlib.pyplot as plt
import numpy as np
import pandas as pd
import seaborn as sns
```

EDA

```
In [2]: planets = pd.read_csv('data/planets.csv')
planets.head()
```

Out[2]:

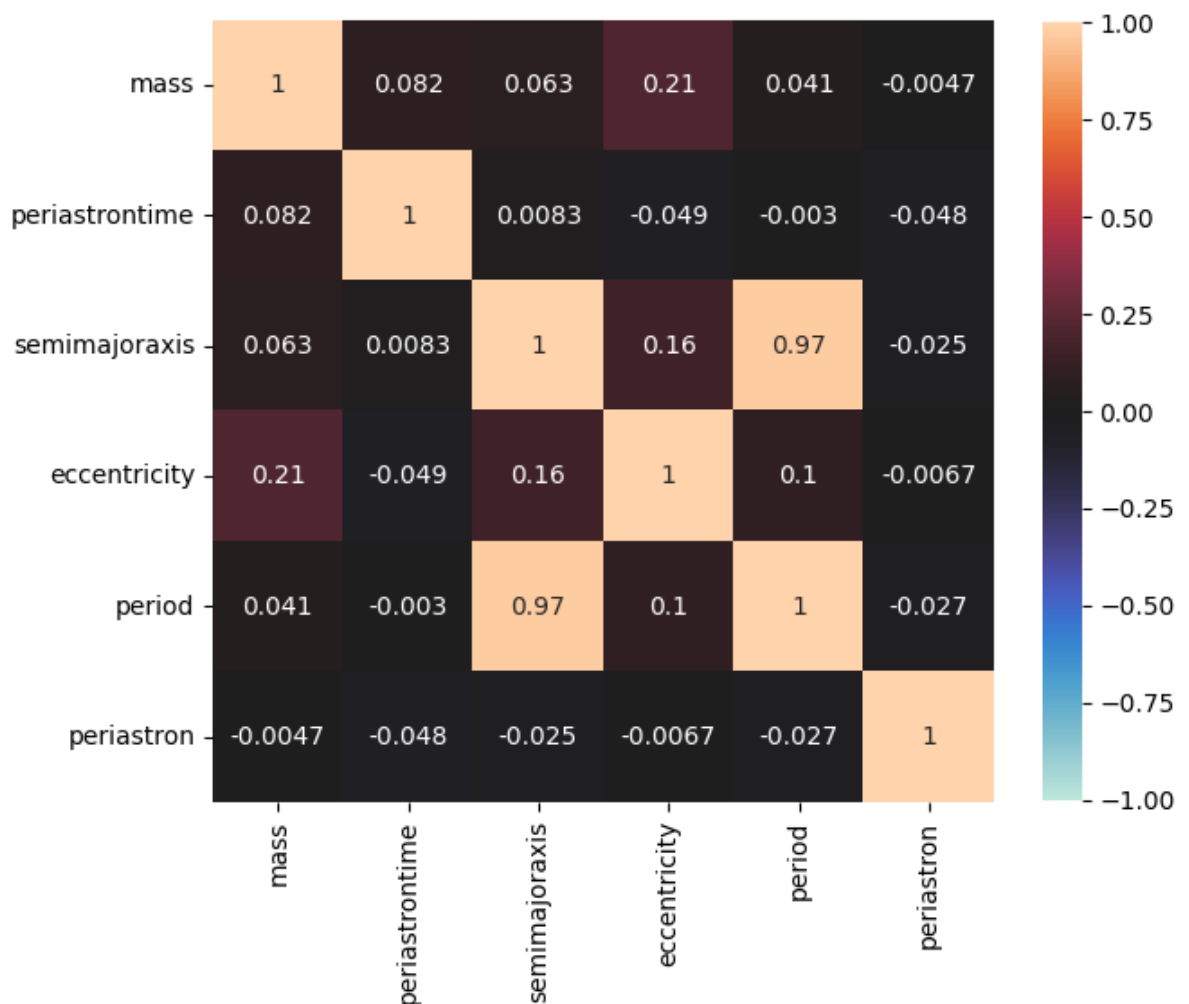
	mass	description	periastrontime	semimajoraxis	discoveryyear	list	eccentricity	
0	19.400	11 Com b is a brown dwarf-mass companion to th...	2452899.60	1.290	2008.0	Confirmed planets	0.231	
1	11.200	11 Ursae Minoris is a star located in the cons...	2452861.04	1.540	2009.0	Confirmed planets	0.080	
2	4.800	14 Andromedae is an evolved star in the conste...	2452861.40	0.830	2008.0	Confirmed planets	0.000	
3	4.975	The star 14 Herculis is only 59 light years aw...	NaN	2.864	2002.0	Confirmed planets	0.359	1
4	7.679	14 Her c is the second companion in the system...	NaN	9.037	2006.0	Controversial	0.184	9

Looking for correlated features

It's important to perform an in-depth exploration of the data before modeling. This includes consulting domain experts, looking for correlations between variables, examining distributions, etc. The visualizations covered in chapters 5 and 6 will prove indispensable for this process. One such visualization is the heatmap which we can use to look for correlated features:

```
In [3]: fig = plt.figure(figsize=(7, 7))
sns.heatmap(
    planets.drop(columns='discoveryyear').corr(),
    center=0, vmin=-1, vmax=1, square=True, annot=True,
    cbar_kws={'shrink': 0.8}
)
```

Out[3]: <AxesSubplot:>



Looking at Orbit shape

Eccentricity	Orbit Shape
0	Circular
(0, 1)	Elliptical
1	Parabolic
> 1	Hyperbolic

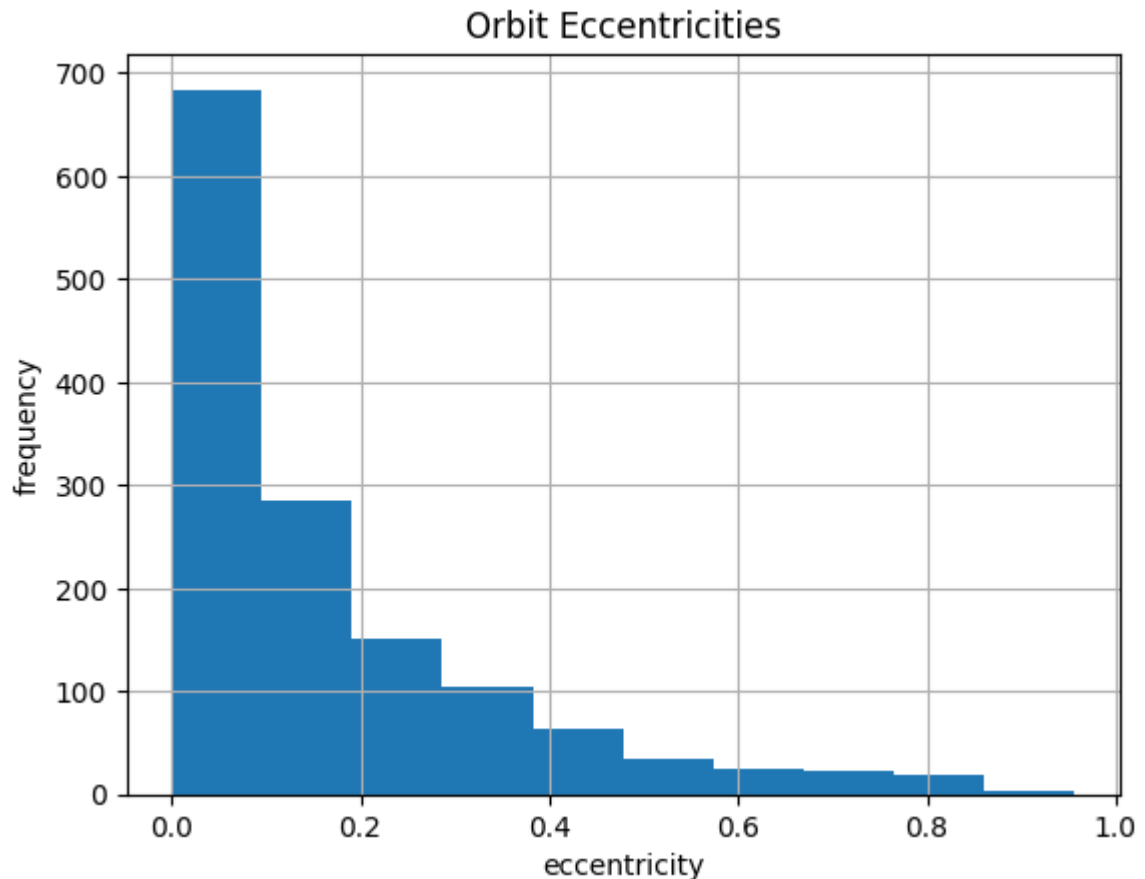
```
In [4]: planets.eccentricity.min(), planets.eccentricity.max()
```

Out[4]: (0.0, 0.956)

All of the planets in the data have circular or elliptical orbits. Let's see the distribution:

```
In [5]: planets.eccentricity.hist()  
plt.xlabel('eccentricity')  
plt.ylabel('frequency')  
plt.title('Orbit Eccentricities')
```

```
Out[5]: Text(0.5, 1.0, 'Orbit Eccentricities')
```

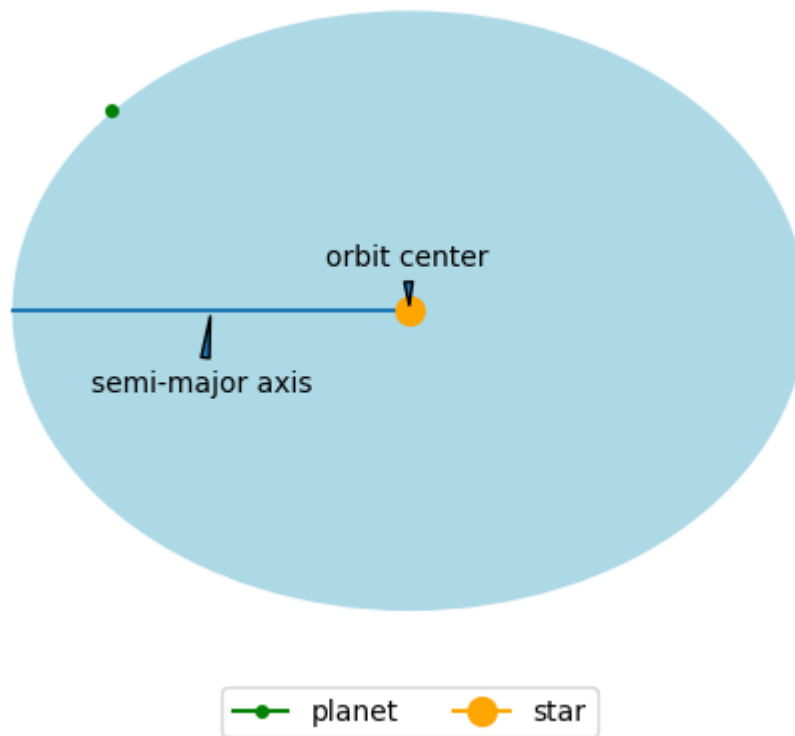


Understanding the semi-major axis

An ellipse, being an elongated circle, has 2 axes: **major** and **minor** for the longest and smallest ones, respectively. The *semi*-major axis is half the major axis. When compared to a circle, the axes are like the diameter crossing the entire shape and the semis are akin to the radius being half the diameter.

```
In [6]: from visual_aids import misc_viz
misc_viz.elliptical_orbit()
```

Out[6]: <AxesSubplot:>



Checking data values

With just the variables of interest, we have a lot of missing data:

```
In [7]: planets[['period', 'eccentricity', 'semimajoraxis', 'mass']].info()
```

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 4094 entries, 0 to 4093
Data columns (total 4 columns):
#   Column          Non-Null Count  Dtype
---  -
0   period          3930 non-null   float64
1   eccentricity    1388 non-null   float64
2   semimajoraxis   1704 non-null   float64
3   mass            1659 non-null   float64
dtypes: float64(4)
memory usage: 128.1 KB
```

If we drop it, we are left with about 30% of it:

```
In [8]: planets[['period', 'eccentricity', 'semimajoraxis', 'mass']].dropna().shape
```

Out[8]: (1222, 4)

We use `describe()` to get a summary of the variables of interest:

```
In [9]: planets[['period', 'eccentricity', 'semimajoraxis', 'mass']].describe()
```

Out[9]:

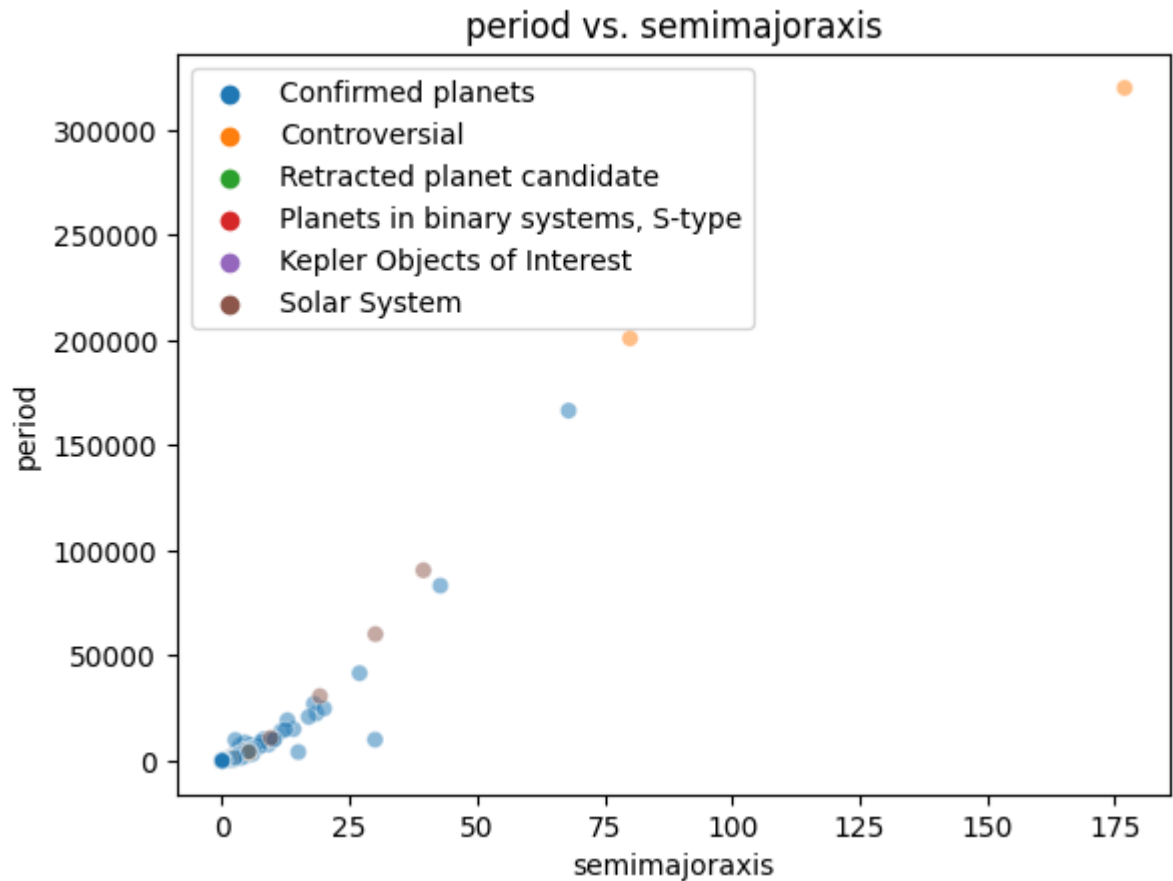
	period	eccentricity	semimajoraxis	mass
count	3930.000000	1388.000000	1704.000000	1659.000000
mean	524.084969	0.159016	5.837964	2.702061
std	7087.428665	0.185041	110.668743	8.526177
min	0.090706	0.000000	0.004420	0.000008
25%	4.552475	0.013000	0.051575	0.085000
50%	12.364638	0.100000	0.140900	0.830000
75%	46.793136	0.230000	1.190000	2.440000
max	320000.000000	0.956000	3500.000000	263.000000

Visualizing Year and Orbit Length

We have information on the planet list each planet belongs to. We may be wondering: are these planets are controversial because they are so far away?

```
In [10]: sns.scatterplot(
          x=planets.semimajoraxis, y=planets.period,
          hue=planets.list, alpha=0.5
        )
plt.title('period vs. semimajoraxis')
plt.legend(title='')
```

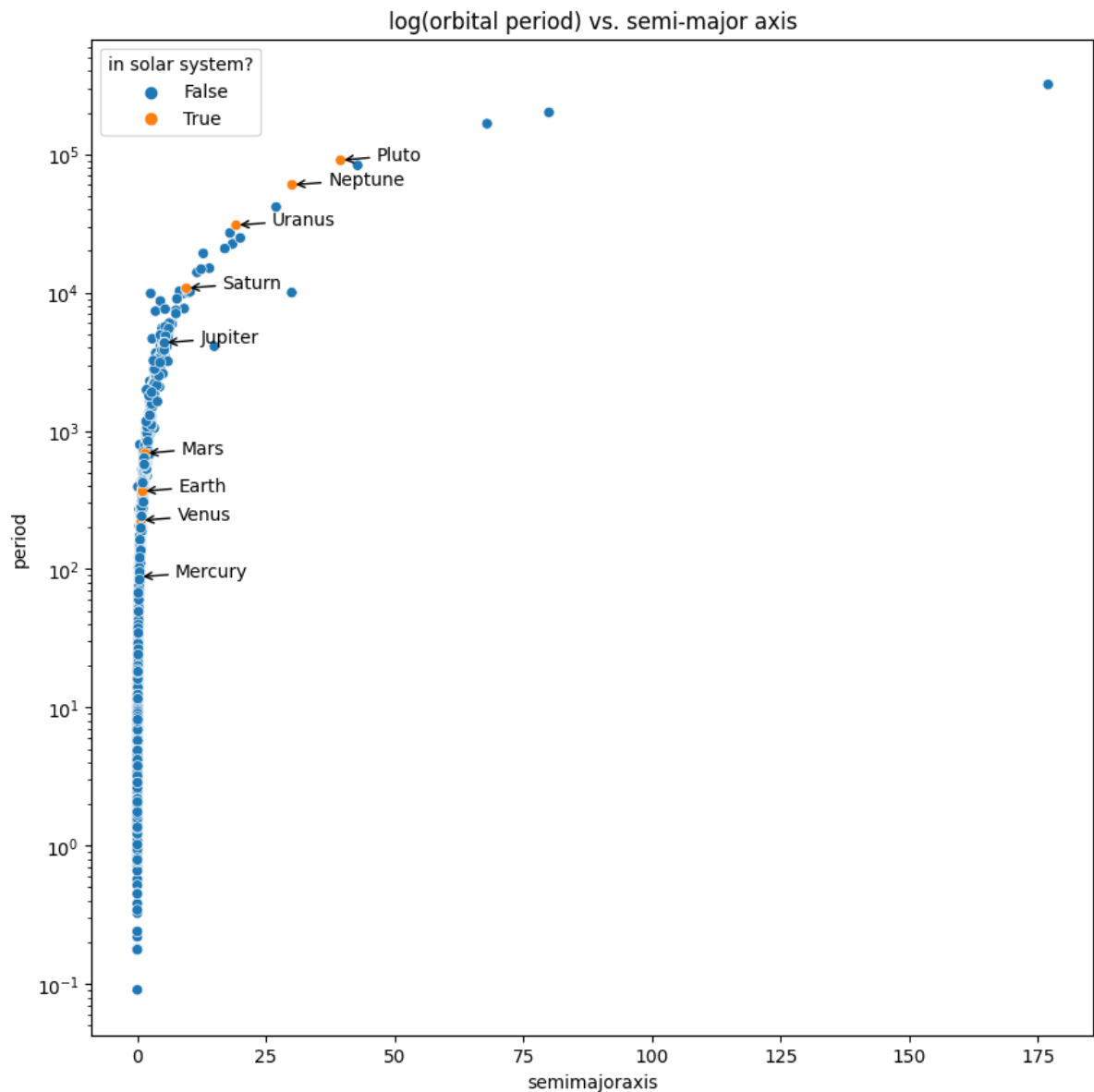
Out[10]: <matplotlib.legend.Legend at 0x1a45c236d60>



Since semi-major axis is highly correlated with period, let's see how the planets compare and label those in our solar system:

```
In [11]: fig, ax = plt.subplots(1, 1, figsize=(10, 10))
in_solar_system = (planets.list == 'Solar System').rename('in solar system?')
sns.scatterplot(
    x=planets.semimajoraxis,
    y=planets.period,
    hue=in_solar_system,
    ax=ax
)
ax.set_yscale('log')
solar_system = planets[planets.list == 'Solar System']
for planet in solar_system.name:
    data = solar_system.query(f'name == "{planet}"')
    ax.annotate(
        planet,
        (data.semimajoraxis, data.period),
        (7 + data.semimajoraxis, data.period),
        arrowprops=dict(arrowstyle='->')
    )
ax.set_title('log(orbital period) vs. semi-major axis')
```


Out[11]: Text(0.5, 1.0, 'log(orbital period) vs. semi-major axis')



Finding Similar Planets with k-Means Clustering

Since we want to perform clustering to learn more about the data, we will build our pipeline standardizing the data before running k-means and fit it on the all the data:

```
In [12]: from sklearn.cluster import KMeans
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import StandardScaler

kmeans_pipeline = Pipeline([
    ('scale', StandardScaler()),
    ('kmeans', KMeans(8, random_state=0))
])
```

Grab the data and fit the model:

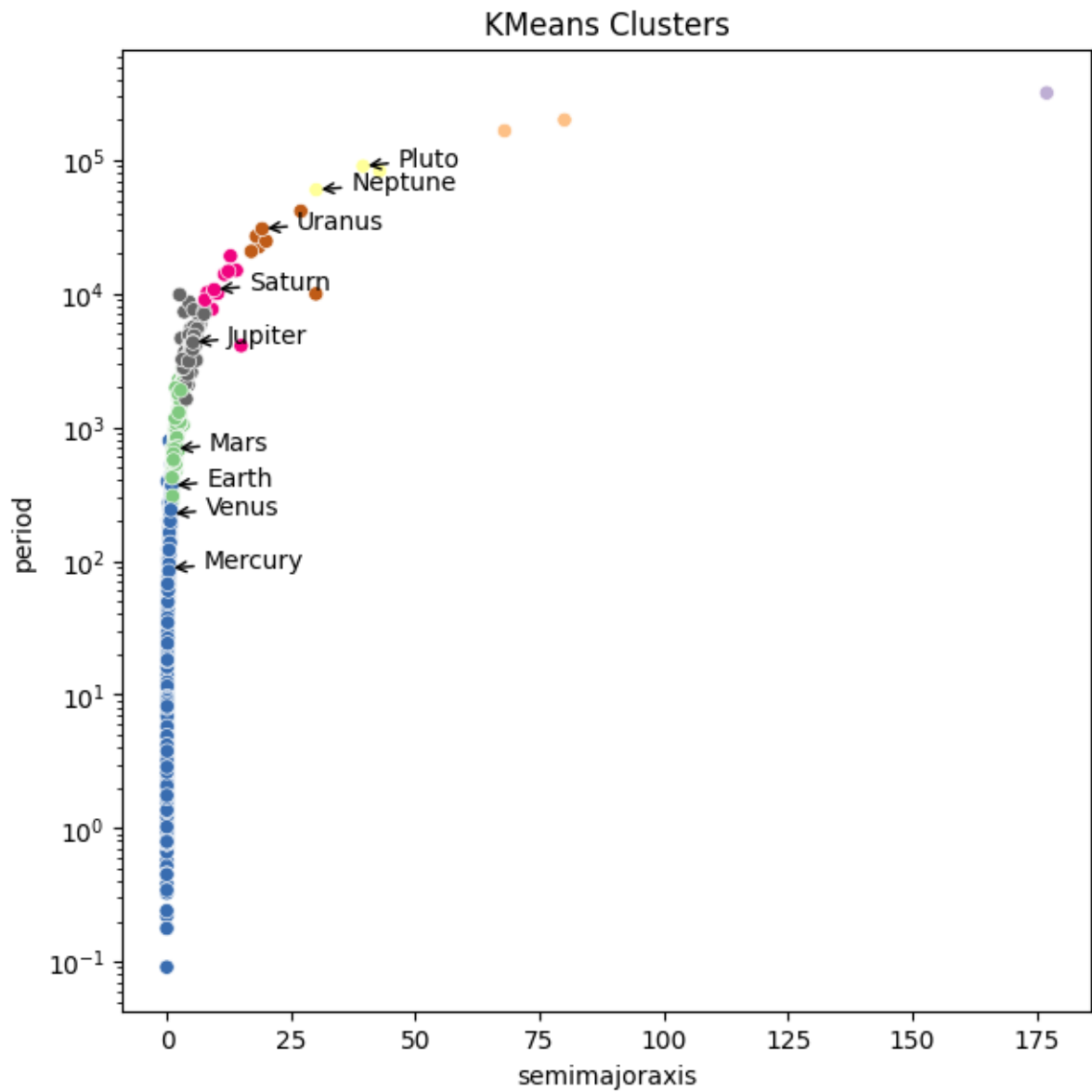
```
In [13]: kmeans_data = planets[['semimajoraxis', 'period']].dropna()  
kmeans_pipeline.fit(kmeans_data)
```

```
Out[13]: Pipeline(steps=[('scale', StandardScaler()),  
                           ('kmeans', KMeans(random_state=0))])
```

We can recreate our plot from before and this time, color by the cluster k-means put each planet in:

```
In [14]: fig, ax = plt.subplots(1, 1, figsize=(7, 7))
sns.scatterplot(
    x=kmeans_data.semimajoraxis,
    y=kmeans_data.period,
    hue=kmeans_pipeline.predict(kmeans_data),
    ax=ax, palette='Accent'
)
ax.set_yscale('log')
solar_system = planets[planets.list == 'Solar System']
for planet in solar_system.name:
    data = solar_system.query(f'name == "{planet}"')
    ax.annotate(
        planet,
        (data.semimajoraxis, data.period),
        (7 + data.semimajoraxis, data.period),
        arrowprops=dict(arrowstyle='->')
    )
ax.get_legend().remove()
ax.set_title('KMeans Clusters')
```

```
Out[14]: Text(0.5, 1.0, 'KMeans Clusters')
```

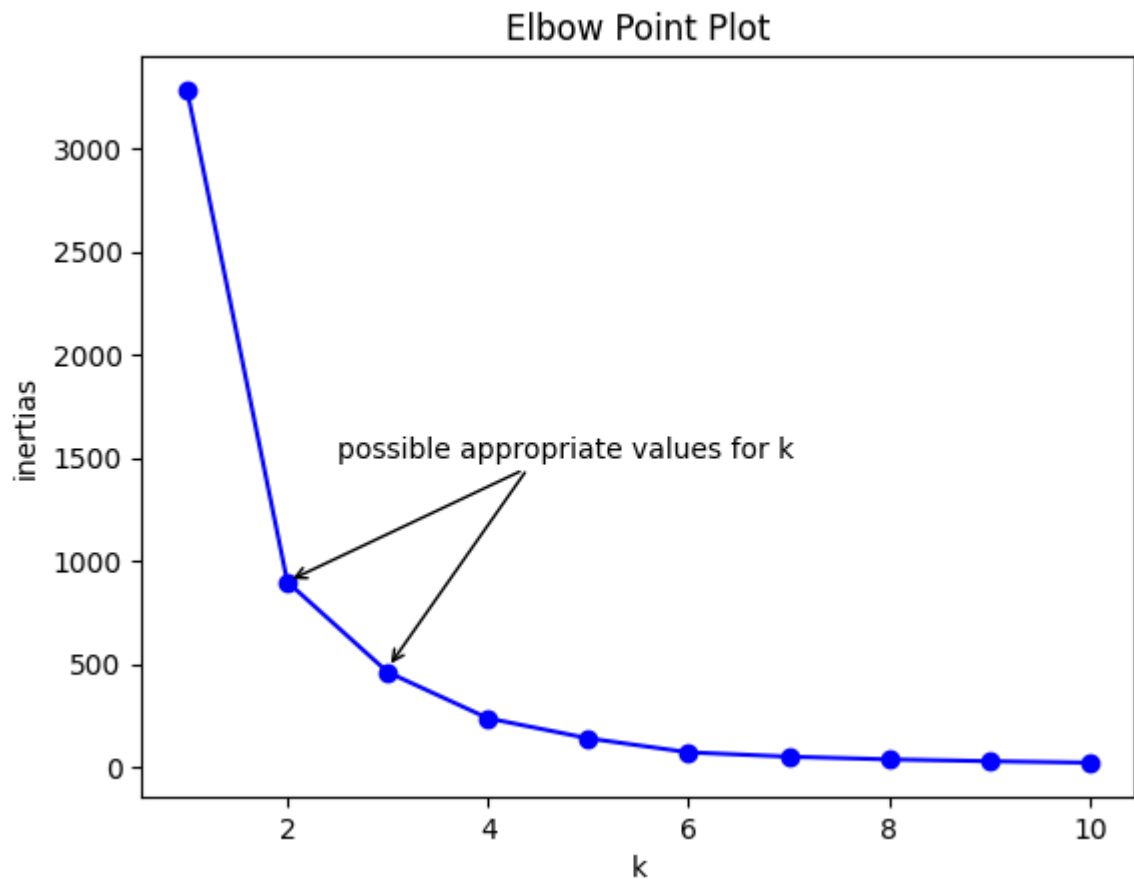


The elbow point method can be used to pick a good value for k . This value will be where we begin to see diminishing returns in the reduction of the value of the objective function:

```
In [15]: from ml_utils.elbow_point import elbow_point

ax = elbow_point(
    kmeans_data,
    Pipeline([
        ('scale', StandardScaler()),
        ('kmeans', KMeans(random_state=0))
    ])
)
ax.annotate(
    'possible appropriate values for k', xy=(2, 900), xytext=(2.5, 1500),
    arrowprops=dict(arrowstyle='->')
)
ax.annotate(
    '', xy=(3, 480), xytext=(4.4, 1450), arrowprops=dict(arrowstyle='->')
)
```

Out[15]: Text(4.4, 1450, '')

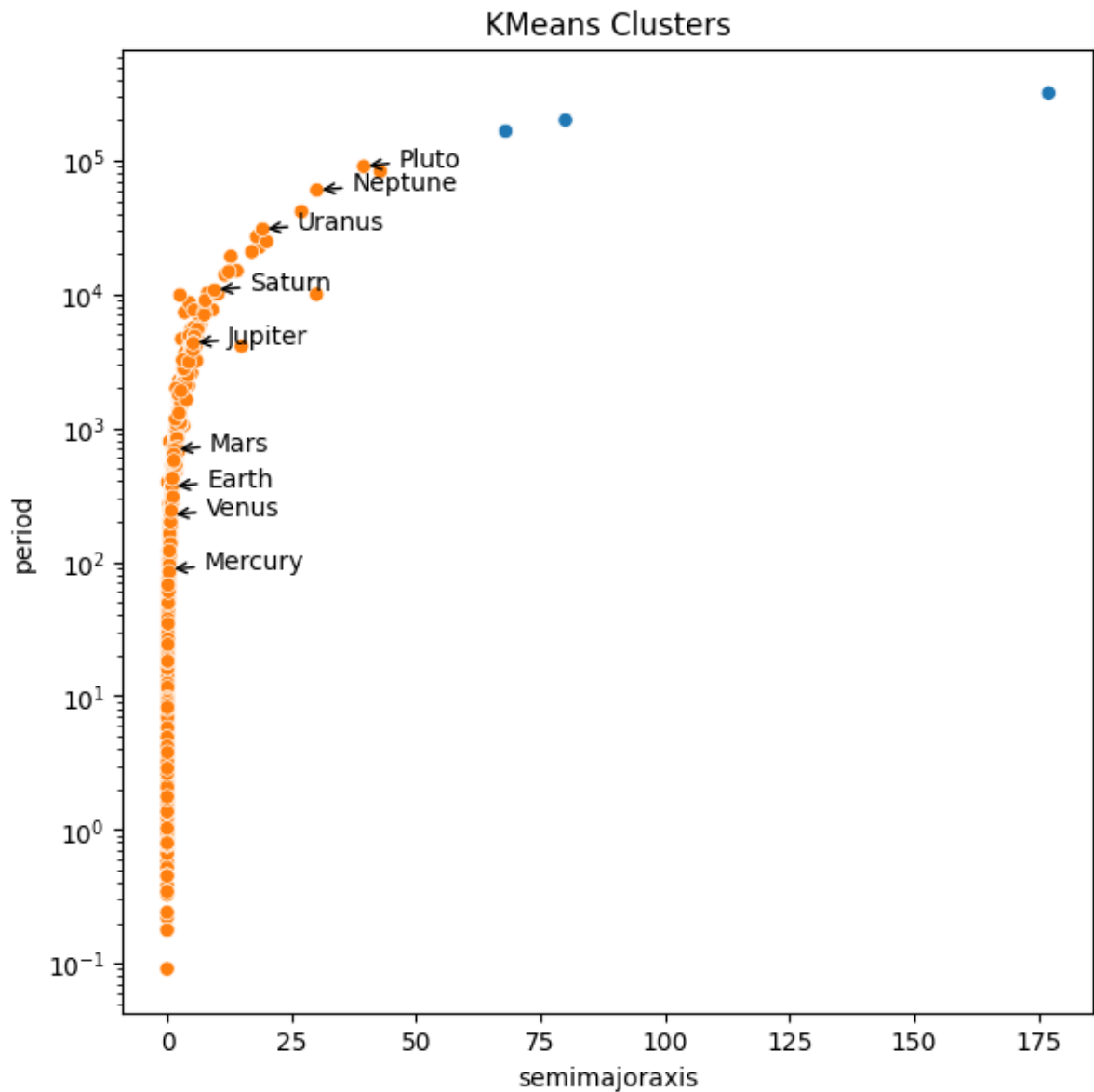


k-means with the "optimal" k of 2

```
In [44]: kmeans_pipeline_2 = Pipeline([
    ('scale', StandardScaler()),
    ('kmeans', KMeans(2, random_state=0))
]).fit(kmeans_data)

fig, ax = plt.subplots(1, 1, figsize=(7, 7))
sns.scatterplot(
    x=kmeans_data.semimajoraxis,
    y=kmeans_data.period,
    hue=kmeans_pipeline_2.predict(kmeans_data),
    ax=ax
)
ax.set_yscale('log')
solar_system = planets[planets.list == 'Solar System']
for planet in solar_system.name:
    data = solar_system.query(f'name == "{planet}"')
    ax.annotate(
        planet,
        (data.semimajoraxis, data.period),
        (7 + data.semimajoraxis, data.period),
        arrowprops=dict(arrowstyle='->')
    )
ax.get_legend().remove()
ax.set_title('KMeans Clusters')
```

Out[44]: Text(0.5, 1.0, 'KMeans Clusters')



Visualizing the cluster space

Since we standardized the data, looking at the centers tells us the second cluster contains "outliers" for period and semi-major axis:

```
In [17]: kmeans_pipeline_2.named_steps['kmeans'].cluster_centers_
```

```
Out[17]: array([[18.9113263 , 20.86736585],
                [-0.03463613, -0.03821862]])
```

We can also visualize the clusters:

```

In [18]: # set up layout
fig = plt.figure(figsize=(8, 6))
outside = fig.add_axes([0.1, 0.1, 0.9, 0.9])
inside = fig.add_axes([0.6, 0.2, 0.35, 0.35])

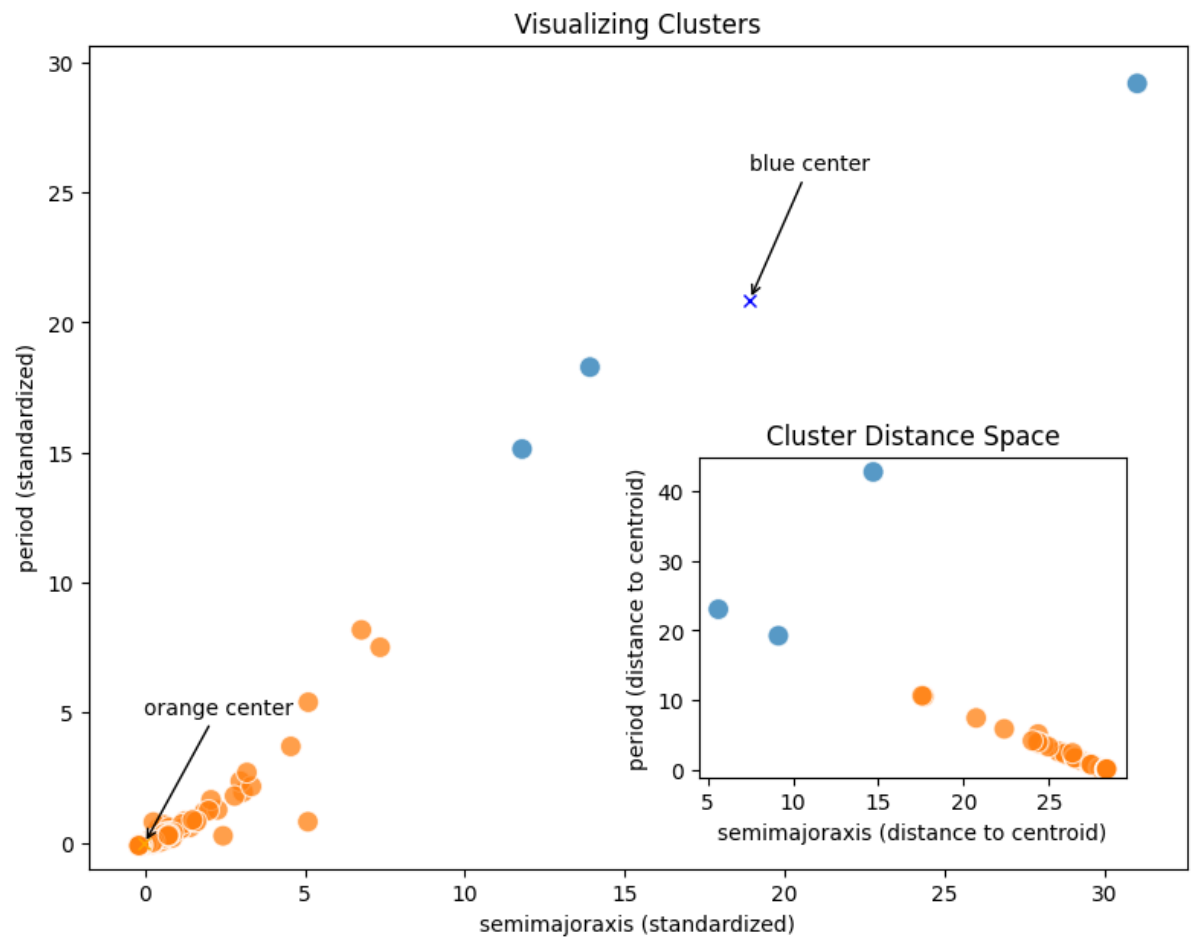
# scaled data and cluster distance data
scaled = kmeans_pipeline_2.named_steps['scale']\
    .fit_transform(kmeans_data)
cluster_distances = kmeans_pipeline_2\
    .fit_transform(kmeans_data)

for ax, data, title, axes_labels in zip(
    [outside, inside], [scaled, cluster_distances],
    ['Visualizing Clusters', 'Cluster Distance Space'],
    ['standardized', 'distance to centroid']
):
    sns.scatterplot(
        x=data[:,0], y=data[:,1], ax=ax, alpha=0.75, s=100,
        hue=kmeans_pipeline_2.named_steps['kmeans'].labels_
    )

    ax.get_legend().remove()
    ax.set_title(title)
    ax.set_xlabel(f'semimajoraxis ({axes_labels})')
    ax.set_ylabel(f'period ({axes_labels})')
    ax.set_ylim(-1, None)

# add the centroids to the outside plot
cluster_centers = kmeans_pipeline_2.named_steps['kmeans'].cluster_centers_
for color, centroid in zip(['blue', 'orange'], cluster_centers):
    outside.plot(*centroid, color=color, marker='x')
    outside.annotate(
        f'{color} center', xy=centroid, xytext=centroid + [0, 5],
        arrowprops=dict(arrowstyle='->')
    )

```

Notes on the `scikit-learn` API

Method	Action	Used when...
<code>fit()</code>	Train the model or preprocessor	Modeling, preprocessing
<code>transform()</code>	Transform the data into the new space	Clustering, preprocessing
<code>fit_transform()</code>	Run <code>fit()</code> , followed by <code>transform()</code>	Clustering, preprocessing
<code>score()</code>	Evaluate the model using the default scoring method	Modeling
<code>predict()</code>	Use model to predict output values for given inputs	Modeling
<code>fit_predict()</code>	Run <code>fit()</code> , followed by <code>predict()</code>	Modeling
<code>predict_proba()</code>	Like <code>predict()</code> , but returns the probability of belonging to each class	Classification

Evaluation of model

There are many metrics to choose from, but since we don't know the true labels of our data, we can only use unsupervised ones. We will use a few different metrics to get a more well-rounded view of our performance:

Silhouette Score

- true labels not known
- higher = better defined (more separated) clusters
- -1 is worst, 1 is best, near 0 indicates overlapping clusters

```
In [19]: from sklearn.metrics import silhouette_score
          silhouette_score(kmeans_data, kmeans_pipeline.predict(kmeans_data))
```

```
Out[19]: 0.7579771626036677
```

Davies-Bouldin Score

- true labels not known
- ratio of within-cluster distances to between-cluster distances
- zero is the best partition

```
In [20]: from sklearn.metrics import davies_bouldin_score
          davies_bouldin_score(kmeans_data, kmeans_pipeline.predict(kmeans_data))
```

```
Out[20]: 0.4632311032231894
```

Calinski and Harabasz Score

- true labels not known
- higher = better defined (more separated) clusters

```
In [21]: from sklearn.metrics import calinski_harabasz_score  
calinski_harabasz_score(kmeans_data, kmeans_pipeline.predict(kmeans_data))
```

Out[21]: 21207.276781867335

Predicting Length of Year in Earth Days (Period)

1. separate x and y data, dropping nulls
2. create the training and testing sets
3. train a linear regression model (no preprocessing since we want to interpret the coefficients)
4. isolate the coefficients from the model
5. evaluate the model

Step 1:

```
In [22]: data = planets[  
        ['semimajoraxis', 'period', 'mass', 'eccentricity']  
        ].dropna()  
X = data[['semimajoraxis', 'mass', 'eccentricity']]  
y = data.period
```

Step 2:

```
In [23]: from sklearn.model_selection import train_test_split  
  
X_train, X_test, y_train, y_test = train_test_split(  
    X, y, test_size=0.25, random_state=0  
)
```

Linear Regression

Step 3:

```
In [24]: from sklearn.linear_model import LinearRegression  
lm = LinearRegression().fit(X_train, y_train)
```

Get equation

Step 4:

```
In [25]: # get intercept  
lm.intercept_
```

Out[25]: -622.9909910671802

```
In [26]: # get coefficients
[(col, coef) for col, coef in zip(X_train.columns, lm.coef_)]
```

```
Out[26]: [('semimajoraxis', 1880.4365990440922),
          ('mass', -90.18675916509216),
          ('eccentricity', -3201.0780593330896)]
```

Evaluation of model

Step 5

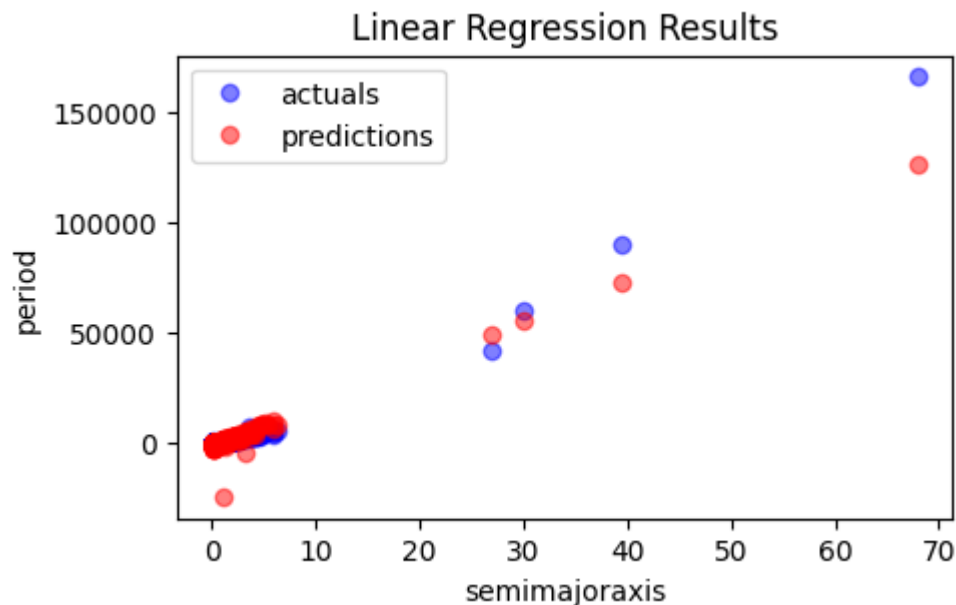
In order to evaluate our model's predictions against the actual values, we need to make predictions for the test set:

```
In [27]: preds = lm.predict(X_test)
```

We can then plot the predictions and actual values:

```
In [28]: fig, axes = plt.subplots(1, 1, figsize=(5, 3))
axes.plot(X_test.semimajoraxis, y_test, 'ob', label='actuals', alpha=0.5)
axes.plot(X_test.semimajoraxis, preds, 'or', label='predictions', alpha=0.5)
axes.set(xlabel='semimajoraxis', ylabel='period')
axes.legend()
axes.set_title('Linear Regression Results')
```

```
Out[28]: Text(0.5, 1.0, 'Linear Regression Results')
```



The correlation between the predictions and the actual values tells us they trend together, but we need to look at other metrics to quantify the errors our model makes:

```
In [29]: np.corrcoef(y_test, preds)[0][1]
```

```
Out[29]: 0.9692104355988056
```

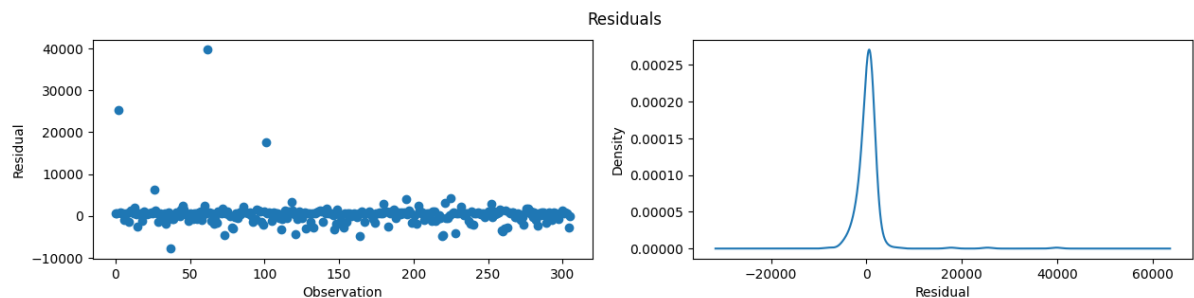
Residuals

Our residuals have no pattern (left subplot); however, the distribution has some negative skew, and the residuals aren't quite centered around zero (right subplot):

```
In [30]: from ml_utils.regression import plot_residuals
```

```
plot_residuals(y_test, preds)
```

```
Out[30]: array([<AxesSubplot:xlabel='Observation', ylabel='Residual'>,
                <AxesSubplot:xlabel='Residual', ylabel='Density'>], dtype=object)
```



R^2

By default, the `score()` method of the `LinearRegression` object will give us the R^2 :

```
In [31]: lm.score(X_test, y_test)
```

```
Out[31]: 0.9209013475842682
```

If not, we can use the `r2_score()` function from `sklearn.metrics`:

```
In [32]: from sklearn.metrics import r2_score
r2_score(y_test, preds)
```

```
Out[32]: 0.9209013475842682
```

Adjusted R^2

R^2 increases when we add regressors whether or not they actually improve the model. Adjusted R^2 penalizes additional regressors to address this:

```
In [33]: from ml_utils.regression import adjusted_r2
adjusted_r2(lm, X_test, y_test)
```

Out[33]: 0.9201155993814629

Problems with R^2

R^2 doesn't tell us about the prediction errors or if we specified the model correctly. Consider Anscombe's quartet from chapter 1:

Anscombe's Quartet

All four data sets have the same summary statistics (mean, standard deviation, correlation coefficient), despite having different data:

```
In [34]: anscombe = sns.load_dataset('anscombe').groupby('dataset')
anscombe.describe()
```

Out[34]:

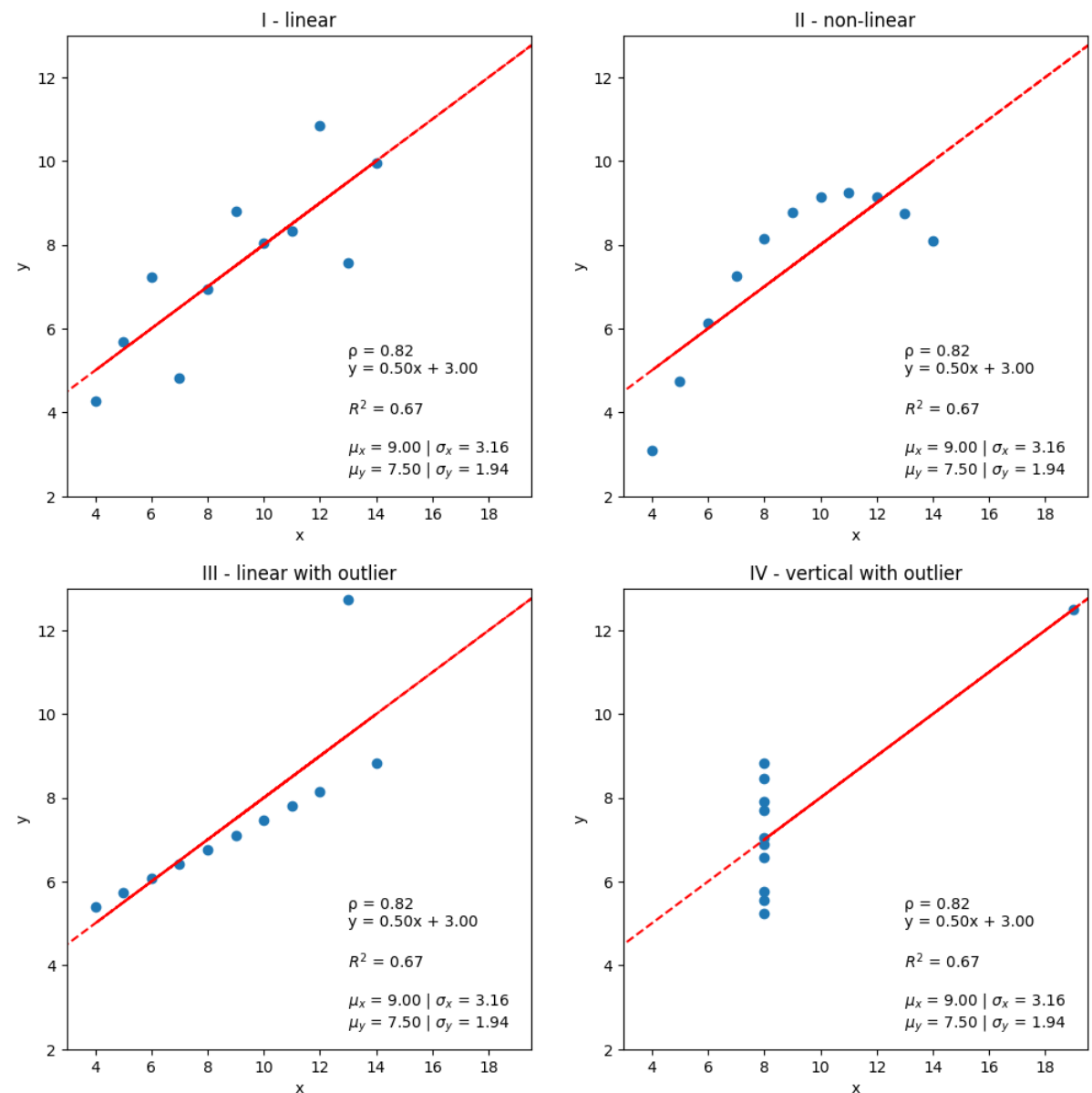
	x								y						
	count	mean	std	min	25%	50%	75%	max	count	mean	std	min	25%		
dataset															
I	11.0	9.0	3.316625	4.0	6.5	9.0	11.5	14.0	11.0	7.500909	2.031568	4.26	6.5		
II	11.0	9.0	3.316625	4.0	6.5	9.0	11.5	14.0	11.0	7.500909	2.031657	3.10	6.5		
III	11.0	9.0	3.316625	4.0	6.5	9.0	11.5	14.0	11.0	7.500000	2.030424	5.39	6.5		
IV	11.0	9.0	3.316625	8.0	8.0	8.0	8.0	19.0	11.0	7.500909	2.030579	5.25	6.5		

When fitted with a regression line, they all have the same R^2 despite some of them not indicating a linear relationship between x and y:

```
In [35]: from visual_aids import stats_viz
stats_viz.anscombes_quartet(r_squared=True)
```

```
Out[35]: array([<AxesSubplot:title={'center':'I - linear'}, xlabel='x', ylabel='y'>,
<AxesSubplot:title={'center':'II - non-linear'}, xlabel='x', ylabel='y'>,
<AxesSubplot:title={'center':'III - linear with outlier'}, xlabel='x',
ylabel='y'>,
<AxesSubplot:title={'center':'IV - vertical with outlier'}, xlabel='x',
ylabel='y'>],
dtype=object)
```

Anscombe's Quartet



Explained Variance

The percentage of the variance in the data is explained by our model:

```
In [36]: from sklearn.metrics import explained_variance_score
         explained_variance_score(y_test, preds)
```

```
Out[36]: 0.9220144218429369
```

Mean Absolute Error (MAE)

This gives us an idea of how far off our predictions are on average (in Earth days):

```
In [37]: from sklearn.metrics import mean_absolute_error
         mean_absolute_error(y_test, preds)
```

```
Out[37]: 1369.4418170735328
```

Root Mean Squared Error (RMSE)

We can use this to punish large errors more:

```
In [40]: from sklearn.metrics import mean_squared_error
         np.sqrt(mean_squared_error(y_test, preds))
```

```
Out[40]: 3248.4999619283776
```

Median Absolute Error

We can also look at the median absolute error to ignore any outliers in prediction errors and get a better picture of our error:

```
In [41]: from sklearn.metrics import median_absolute_error
         median_absolute_error(y_test, preds)
```

```
Out[41]: 759.861335833544
```

[← Chapter 8](#) ([../ch_08/anomaly_detection.ipynb](#)) [Planet Data Collection](#) ([../planet_data_collection.ipynb](#))
[Preprocessing](#) ([../preprocessing.ipynb](#)) [Red Wine](#) ([../red_wine.ipynb](#)) [Red + White Wine](#) ([../wine.ipynb](#))
</div> [Solutions](#) ([../solutions/ch_09/exercise_1.ipynb](#)) [Chapter 10 →](#) ([../ch_10/red_wine.ipynb](#))
