

M6-L1 Problem 3

SciKit-Learn only offers a few built-in preprocessors, such as MinMax and Standard scaling. However, it also offers the ability to create custom data transformation functions, which can be integrated into your pipeline. In this problem, you will implement a log transform and observe how using it changes a regression result.

Start by running this cell to import modules and load data:

```
In [1]: import numpy as np
import matplotlib.pyplot as plt
from sklearn.model_selection import train_test_split
from sklearn.metrics import mean_squared_error
from sklearn.pipeline import Pipeline
from sklearn.preprocessing import FunctionTransformer
from sklearn.svm import SVR

def plot(X_train, X_test, y_train, y_test, model=None, log = False):
    plt.figure(figsize=(5, 5), dpi=200)
    if model is not None:
        X_fit = np.linspace(min(X_train)-0.15, max(X_train)+0.2)
        y_fit = model.predict(X_fit)
        plt.plot(np.log(X_fit+1) if log else X_fit, y_fit, c="red", label="Prediction")
    if log:
        X_train = np.log(X_train+1)
        X_test = np.log(X_test+1)

    plt.scatter(X_train, y_train, s=30, c="powderblue", edgecolors="navy", linewidths=.5, label="Train")
    plt.scatter(X_test, y_test, s=30, c="orange", edgecolors="red", linewidths=.5, label="Test")
    plt.legend()
    plt.xlabel("log(x+1)" if log else "x")
    plt.ylabel("y")
    plt.show()

x = np.array([ 5.83603919,  1.49205924,  2.66109578,  9.40172515,  6.47247125,  0.376331])
X = x.reshape(-1, 1)
y = np.array([ 4.32538472e+00, -5.59312420e+00, -4.57455876e+00,  4.23667057e+01,  1.04111111e+01])

X_train, X_test, y_train, y_test = train_test_split(X, y, random_state=1, train_size=0.8)
```

Distribution of x data

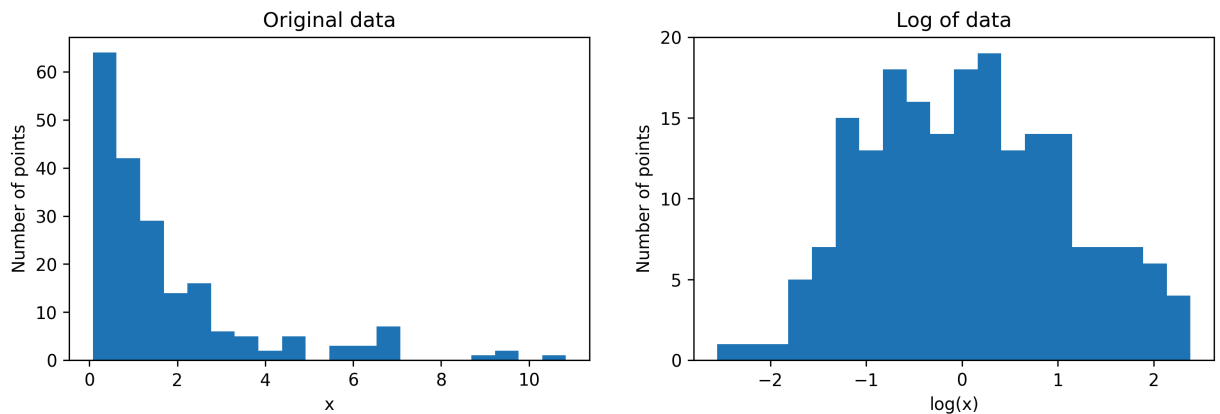
Let's visualize how the original input feature is distributed, alongside the log of the data -- notice that performing this log transformation makes the data much closer to normally distributed.

```
In [2]: plt.figure(figsize=(12, 3.4), dpi=300)
plt.subplot(1, 2, 1)
plt.hist(x, bins=20)
plt.xlabel("x")
plt.ylabel("Number of points")
plt.title("Original data")

plt.subplot(1, 2, 2)
plt.hist(np.log(x), bins=20)
plt.xlabel("log(x)")
plt.ylabel("Number of points")
```

```
plt.title("Log of data")
plt.ylim(0, 20)
plt.yticks([0, 5, 10, 15, 20])

plt.show()
```



No log transform

First, we do support vector regression on the untransformed inputs. The code to do this has been provided below.

```
In [3]: model = SVR(C=100)

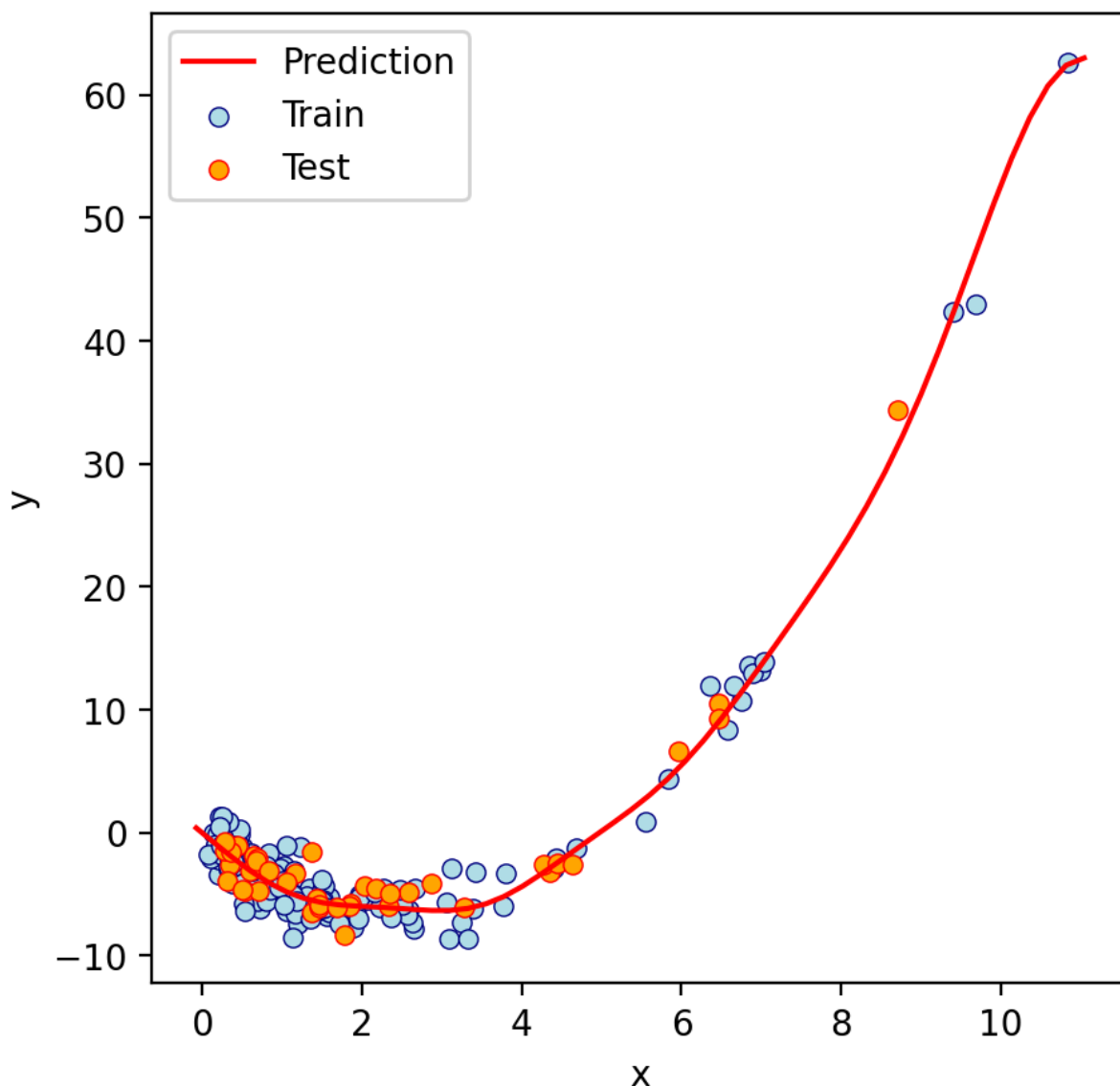
pipeline = Pipeline([("SVR", model)])

# Fit the pipeline to the training data
pipeline.fit(X_train, y_train)

# Make predictions with the pipeline
pred_train = pipeline.predict(X_train)
pred_test = pipeline.predict(X_test)
print("Training MSE:", mean_squared_error(y_train, pred_train), "    Testing MSE:", mea

# Plot the predictions
plot(X_train, X_test, y_train, y_test, pipeline)
```

Training MSE: 2.071006155233616 Testing MSE: 1.9453578716771704



With log transform

Notice that the data are not spread uniformly across the x axis. Instead, most input data points have low values -- this is a roughly "log normal" distribution. If we take the log of the input, we saw it was more normally distributed, which can improve machine learning model results in some cases. The transform function has been given below. Add this to a new pipeline, train the pipeline, and compute the train MSE and test MSE. Show a plot as above. Note the subtle change in behavior of the fitting curve.

Also, make another plot setting the `log` argument to `True`. This will show the scaling of the x-axis used by the model.

```
In [6]: def log_transform(x):
        return np.log(x + 1.)

        transform = FunctionTransformer(log_transform)
        model = SVR(C=100)

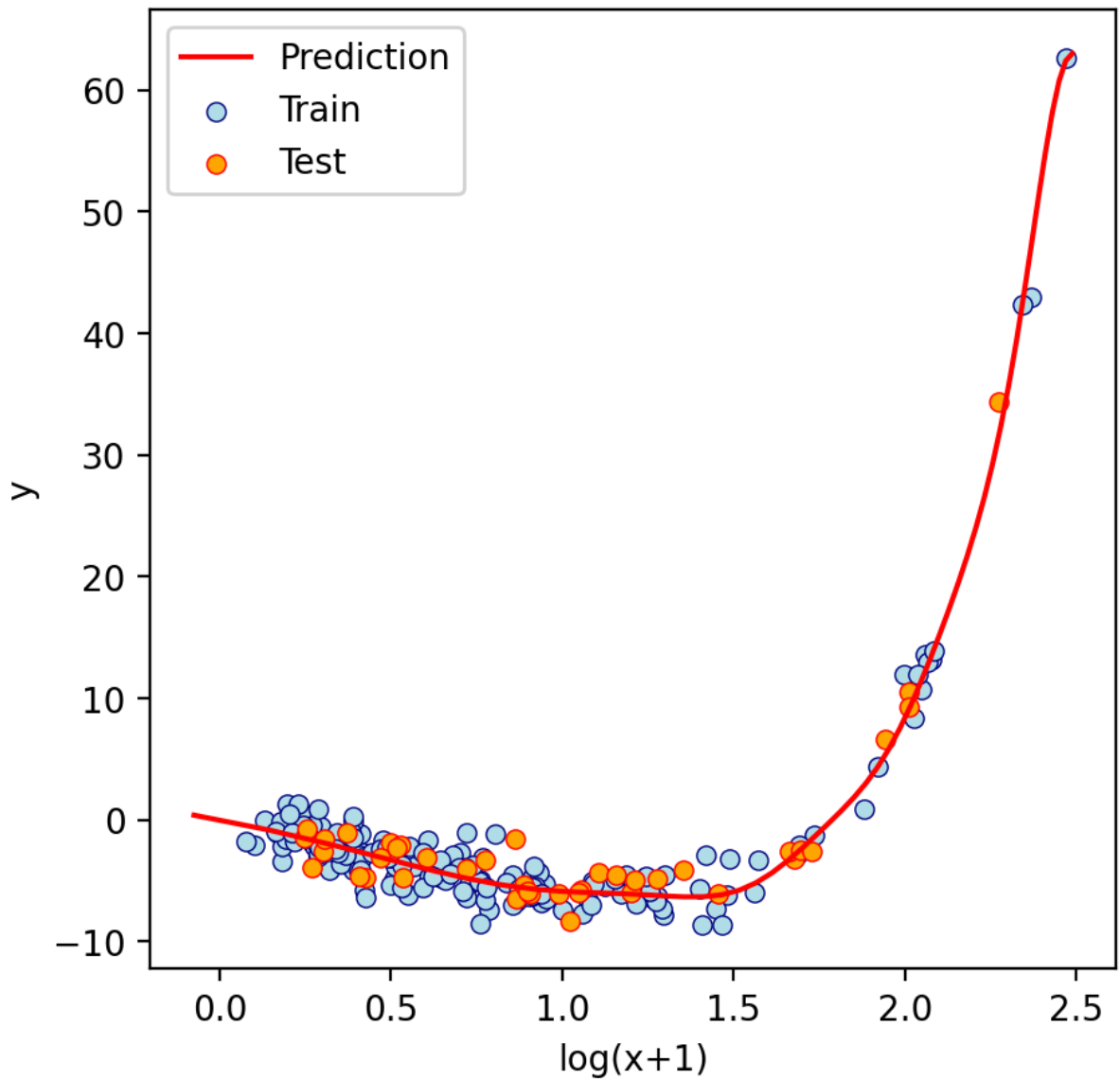
        # YOUR CODE GOES HERE
        pipeline = Pipeline([("SVR", model)])
```

```
# Fit the pipeline to the training data
pipeline.fit(X_train, y_train)

# Make predictions with the pipeline
pred_train = pipeline.predict(X_train)
pred_test = pipeline.predict(X_test)
print("Training MSE:", mean_squared_error(y_train, pred_train), "    Testing MSE:", mean_squared_error(y_test, pred_test))

# Plot the predictions
plot(X_train, X_test, y_train, y_test, pipeline, log = True)
```

Training MSE: 2.071006155233616 Testing MSE: 1.9453578716771704



In []: