

CS178 Homework 2

Due Wednesday, October 16th, 11:59pm

Instructions

This homework (and subsequent ones) will involve data analysis and reporting on methods and results using Python code. You will submit a **single PDF file** that contains everything to Gradescope. This includes any text you wish to include to describe your results, the complete code snippets of how you attempted each problem, any figures that were generated, and scans of any work on paper that you wish to include. It is important that you include enough detail that we know how you solved the problem, since otherwise we will be unable to grade it.

Your homeworks will be given to you as Jupyter notebooks containing the problem descriptions and some template code that will help you get started. You are encouraged to use these starter Jupyter notebooks to complete your assignment and to write your report. This will help you not only ensure that all of the code for the solutions is included, but also will provide an easy way to export your results to a PDF file (for example, doing *print preview* and *printing to pdf*). I recommend liberal use of Markdown cells to create headers for each problem and sub-problem, explaining your implementation/answers, and including any mathematical equations. For parts of the homework you do on paper, scan it in such that it is legible (there are a number of free Android/iOS scanning apps, if you do not have access to a scanner), and include it as an image in the Jupyter notebook.

Double check that all of your answers are legible on Gradescope, e.g. make sure any text you have written does not get cut off.

If you have any questions/concerns about using Jupyter notebooks, ask us on EdD. If you decide not to use Jupyter notebooks, but go with Microsoft Word or LaTeX to create your PDF file, make sure that all of the answers can be generated from the code snippets included in the document.

Summary of Assignment: 100 total points

- Problem 1: k-Nearest Neighbors (25 points)
 - Problem 1.1: Splitting data into training & test sets (10 points)
 - Problem 1.2: Plot predictions for different values of k (10 points)
 - Problem 1.3: Display performance as a function of k & select best (5 points)
- Problem 2: Linear Regression (25 points)
 - Problem 2.1: Train the model and plot the data along with its predictions (15 points)
 - Problem 2.2: Compute the MSE loss for the training and evaluation data (10 points)
- Problem 3: Feature transformations (25 points)
 - Problem 3.1: Train & display polynomial regression models using feature transforms (10 points)
 - Problem 3.2: Plot the training & evaluation error as a function of degree (10 points)
 - Problem 3.3: Select the best degree for these data (5 points)
- Problem 4: Cross-Validation (20 points)
 - Problem 4.1: Plot the five-fold cross validation error (10 points)
 - Problem 4.2: Select the best degree using cross-validation (5 points)
 - Problem 4.3: Compare cross-validation model selection to hold-out data (5 points)
- Statement of Collaboration (5 points)



```
In [2]: import numpy as np
import matplotlib.pyplot as plt
%matplotlib inline

from sklearn.model_selection import train_test_split
from sklearn.metrics import zero_one_loss
from sklearn.metrics import mean_squared_error as mse

from sklearn.neighbors import KNeighborsClassifier, KNeighborsRegressor
from sklearn.inspection import DecisionBoundaryDisplay

from sklearn.linear_model import LinearRegression    # Basic Linear Regression
```

```

from sklearn.linear_model import Ridge           # Linear Regression with L2 regularization
from sklearn.model_selection import KFold        # Cross-validation tools
from sklearn.preprocessing import PolynomialFeatures # Feature transformations
from sklearn.preprocessing import StandardScaler
from sklearn.pipeline import Pipeline           # Useful for sequences of transforms

import requests                                # reading data
from io import StringIO

seed = 1234

```

Training / Test Splits

As we've seen in lecture, it is difficult to tell how accurate our model is from only the data on which it has been trained. For this reason, we usually reserve some data for evaluation, often called "validation" or "test" data. We'll start by loading a one-dimensional regression data set to use in the rest of the homework. We will divide this data set into 75% training data, and 25% evaluation data:

```

In [3]: url = 'https://www.ics.uci.edu/~ihler/classes/cs178/data/curve80.txt'

with requests.get(url) as link: curve = np.genfromtxt(StringIO(link.text), delimiter=None)

X = curve[:,0:-1]      # extract features
Y = curve[:, -1]       # extract target values

# split into training and evaluation data
Xt, Xe, Yt, Ye = train_test_split(X, Y, test_size=0.25, random_state=seed)

```

P1: K-Nearest Neighbor Regression

P1.1: Visualizing the Data Splits

Plot the data for this regression problem, with the (scalar) feature x along the horizontal axis, and the real-valued target y as the vertical axis. Plot all the data, displaying the training data X_t in one color, and the evaluation data X_e in a different color.

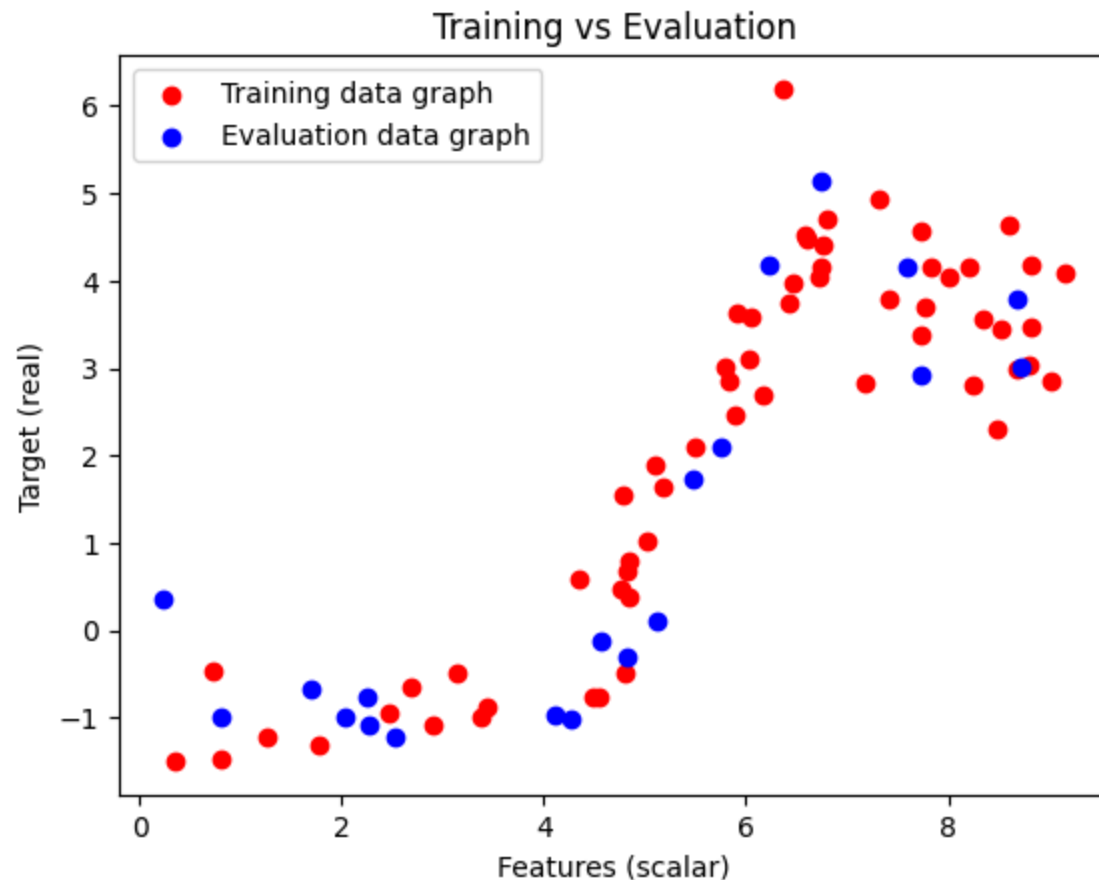
```
In [4]: plt.scatter(Xt, Yt, color="red", label="Training data graph")
plt.scatter(Xe, Ye, color="blue", label="Evaluation data graph")

plt.xlabel("Features (scalar)")
plt.ylabel("Target (real)")
plt.title("Training vs Evaluation")

plt.legend()

plt.show
```

```
Out[4]: <function matplotlib.pyplot.show(close=None, block=None)>
```



P1.2 Visualizing KNN Regression Predictions

Now use `sklearn`'s `KNeighborsRegressor` class to build a nearest neighbor regression model on your training data. Build three models, using $k = 1$, $k = 5$, and $k = 20$, and for each one display the training data, test data, and prediction function. (Note: you can evaluate the prediction function of your learner by predicting at a dense collection of locations `x_spaced` along the x-axis, and then predicting at these points and connecting them using `plot`.)

```
In [5]: # Create a figure with 1 row and 3 columns
fig, axes = plt.subplots(1, 3, figsize=(12, 3))

x_spaced = np.linspace(0,9,100).reshape(-1,1) # get a collection of x-locations at which to plot f(x)
```

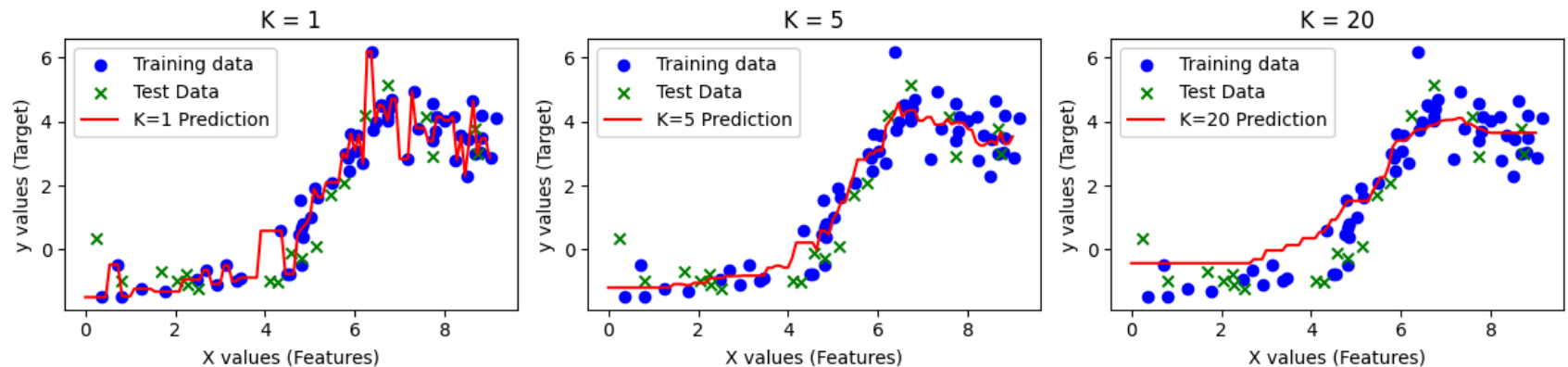
```
### YOUR CODE STARTS HERE ###
```

```
for idx, k in enumerate([1, 5, 20]):
    knn = KNeighborsRegressor(n_neighbors=k)
    knn.fit(Xt, Yt)
    y_pred = knn.predict(x_spaced)

    axes[idx].scatter(Xt, Yt, color="blue", label="Training data", marker="o")
    axes[idx].scatter(Xe, Ye, color='green', label='Test Data', marker='x')
    axes[idx].plot(x_spaced, y_pred, color='red', label=f'K={k} Prediction')

    axes[idx].set_title(f'K = {k}')
    axes[idx].set_xlabel('X values (Features)')
    axes[idx].set_ylabel('y values (Target)')
    axes[idx].legend()
### YOUR CODE ENDS HERE ###

fig.tight_layout()
```



P1.3: KNN Model Selection

Train a model for each k in $1 \leq k \leq 30$, and compute their training and validation MSE. Plot these values as a function of k . What is the best value of k for your model?

```
In [6]: k_values = list(range(1,30))
mse_train = []
mse_eval = []

for i,k in enumerate(k_values):

    ### YOUR CODE STARTS HERE ###
    knn = KNeighborsRegressor(n_neighbors=k)
    knn.fit(Xt, Yt)
    y_train_pred = knn.predict(Xt)
    y_eval_pred = knn.predict(Xe)

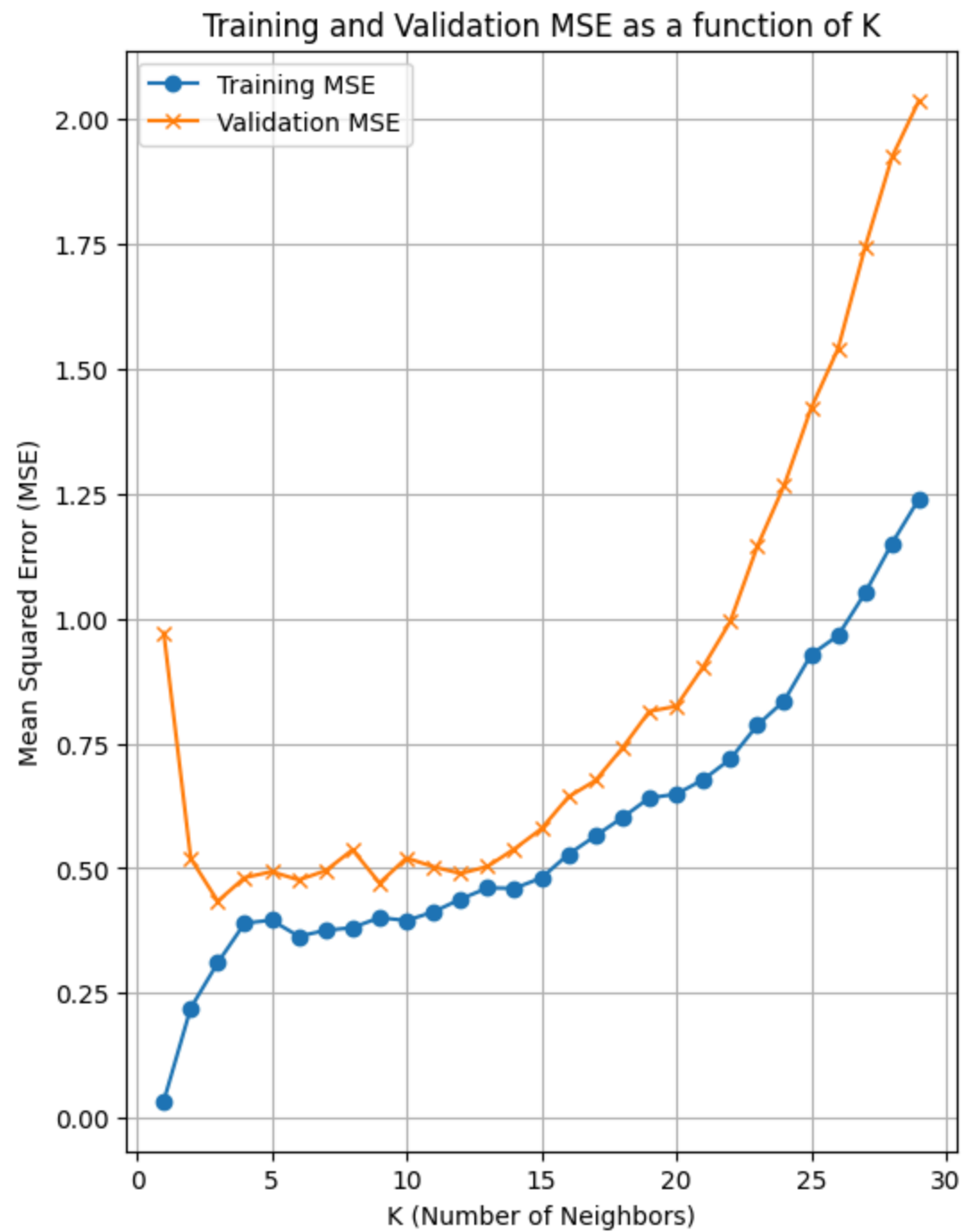
    train_mse = mse(Yt, y_train_pred)
    eval_mse = mse(Ye, y_eval_pred)

    mse_train.append(train_mse)
    mse_eval.append(eval_mse)

plt.figure(figsize=(6, 8))
plt.plot(k_values, mse_train, label='Training MSE', marker='o')
plt.plot(k_values, mse_eval, label='Validation MSE', marker='x')
plt.xlabel('K (Number of Neighbors)')
plt.ylabel('Mean Squared Error (MSE)')
plt.title('Training and Validation MSE as a function of K')
plt.legend()
plt.grid(True)
plt.show()

best_k = k_values[np.argmin(mse_eval)]
print(f"Best K value: {best_k} with Validation MSE: {min(mse_train)}")

### YOUR CODE ENDS HERE ###
```



Best K value: 3 with Validation MSE: 0.0337995723654002

P2: Linear Regression

P2.1: Train linear regression model

Now, let's train a simple linear regression model on the training data. After training the model, plot the training data (colored blue), evaluation data (colored red), and our linear fit (a line) together on a single plot. Also print out the coefficients (slope, `lr.coef_`, and intercept, `lr.intercept_`) of your model after fitting.

```
In [7]: plt.figure(figsize=(6,4))

        ### YOUR CODE STARTS HERE ###

lr = LinearRegression()

lr.fit(Xt, Yt)

y_train_pred = lr.predict(Xt)
y_test_pred = lr.predict(Xe)

plt.figure(figsize=(6, 4))
plt.scatter(Xt, Yt, color='blue', label='Training Data')
plt.scatter(Xe, Ye, color='red', label='Test Data')
plt.plot(np.linspace(1, 10, 100).reshape(-1, 1), lr.predict(np.linspace(1, 10, 100).reshape(-1, 1)), color='green')

plt.xlabel('X values (Features)')
plt.ylabel('y values (Target)')
plt.title('Linear Regression: Training Data, Test Data, and Linear Fit')
plt.legend()

x_spaced = np.linspace(0,10,200).reshape(-1,1)
yhat_spaced = lr.predict(x_spaced)

plt.plot(x_spaced, yhat_spaced, color='green', label='Prediction at dense points')
```

```
plt.show()
print(f"Model Coefficient (Slope): {lr.coef_[0]}")
print(f"Model Intercept: {lr.intercept_}")
```

```
### YOUR CODE ENDS HERE ###
```

<Figure size 600x400 with 0 Axes>



Model Coefficient (Slope): 0.7768472052927541

Model Intercept: -2.3463013180118275

P2.2 Evaluate your model's fit

Compute the mean squared error of your trained model on the training data (the data it was fit on) and the held-out evaluation data.

```
In [8]: train_mse = mse(Yt, y_train_pred)
        test_mse = mse(Ye, y_test_pred)

        print(f"Training MSE: {train_mse}")
        print(f"Test MSE: {test_mse}")
```

Training MSE: 1.270893125474928

Test MSE: 1.6723519225582435

Problem 3: Feature Transformations

Often we will want to transform our data (as we saw in class). A very simple version of this transformation is "normalizing" the data, in which we shift and scale the feature values to a desirable range; typically, zero mean and unit variance, for example. The `StandardScaler()` object in scikit-learn implements such a transformation.

Typically, a pre-processing transformation works in a similar way to training a model: we `fit` the object to our training data (in this case, computing the empirical mean and variance of the data), and save the parameters of the transformation (the shift and scale values) so that we can apply exactly the same transformation to subsequent data, for example when asked to predict on a new value of x .

So, for example:

```
In [9]: scale = StandardScaler().fit(Xt)
        X_transformed = scale.transform(Xt)

        # Now, we can train our model on X_transformed...
        # lr = LinearRegression()...

        # Before we predict, we also need to transform the test point's values:
        yhat_spaced = lr.predict(scale.transform(x_spaced))
```

If you like (and as described in the Discussion code), you can use `sklearn`'s `Pipeline` object to simplify the process of sequentially applying transformations before a predictor.

P3.1: Train polynomial regression models

As mentioned in the homework, you can create additional features manually, e.g.,

```
In [10]: m,n = Xt.shape           # rest of this cell assumes n=1 feature
Xt2 = np.zeros((m,2))
Xt2[:,0] = Xt[:,0]
Xt2[:,1] = Xt[:,0]**2
print (Xt.shape)
print (Xt2.shape)
print (Xt2[0:6,:])    # look at a few data points to check:

(60, 1)
(60, 2)
[[ 0.72580645  0.526795  ]
 [ 2.4769585   6.13532341]
 [ 7.7304147   59.75931143]
 [ 9.0207373   81.37370144]
 [ 8.6751152   75.25762373]
 [ 6.4631336   41.77209593]]
```

or, you can create them using SciKit's PolynomialFeatures transform object:

```
In [11]: Phi = PolynomialFeatures(degree=2,include_bias=False).fit(Xt)
Xt2 = Phi.transform(Xt)
print (Xt2[0:6,:])    # look at the same data points -- same values

[[ 0.72580645  0.526795  ]
 [ 2.4769585   6.13532341]
 [ 7.7304147   59.75931143]
 [ 9.0207373   81.37370144]
 [ 8.6751152   75.25762373]
 [ 6.4631336   41.77209593]]
```

Now, try fitting a linear regression model using different numbers of polynomial features of x .

For each degree $d \in \{0, 1, 3, 5, 7, 10, 15, 18\}$:

- Fit a linear regression model using features consisting of all powers of x up to degree d
 - Make sure you apply `StandardScaler` to the transformed data before training
- Plot the resulting prediction function $f(x)$, along with the training and validation data as before

```

In [12]: degrees = [0,1,3,5,7,10,15,18]
learners = [None]*len(degrees)

fig, ax = plt.subplots(2,4, figsize=(24,10))

for i,degree in enumerate(degrees):

    ### YOUR CODE STARTS HERE ###

    # Create a polynomial feature expansion of degree d
    poly = PolynomialFeatures(degree=degree, include_bias=(degree == 0))
    Xt_poly = poly.fit_transform(Xt)

    # Use StandardScaler to rescale the transformed data
    scaler = StandardScaler()
    Xt_poly_scaled = scaler.fit(Xt_poly)
    Xt_poly_scaled_transformed = Xt_poly_scaled.transform(Xt_poly)

    # Fit your linear regression and save it to "learners"
    lr = LinearRegression()
    lr.fit(Xt_poly_scaled_transformed, Yt)
    learners[i] = lr

    x_spaced_poly = poly.transform(x_spaced)
    x_spaced_poly_scaled = Xt_poly_scaled.transform(x_spaced_poly)
    yhat_spaced = lr.predict(x_spaced_poly_scaled)

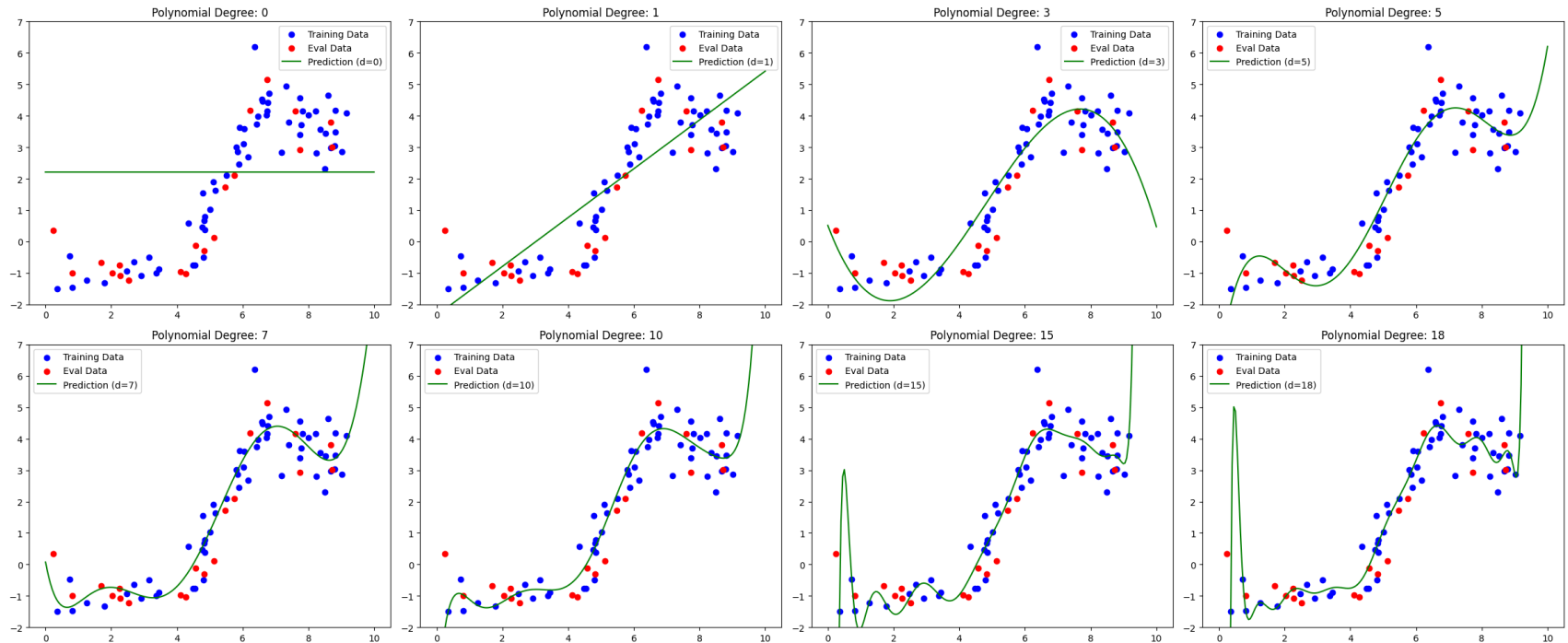
    axi = ax[i//4, i%4]
    axi.scatter(Xt, Yt, color='blue', label='Training Data')
    axi.scatter(Xe, Ye, color='red', label='Eval Data')
    axi.plot(x_spaced, yhat_spaced, color='green', label=f'Prediction (d={degree})')

    axi.set_ylim(-2, 7)
    axi.set_title(f'Polynomial Degree: {degree}')
    axi.legend()

    ### YOUR CODE ENDS HERE ###

plt.tight_layout()
plt.show()

```



P3.2 Model Performance

Compute the mean squared error (MSE) loss of each of your trained models on both the training data and the evaluation data. Plot these errors as a function of degree (so, degree along the horizontal axis, MSE loss as the vertical axis).

```
In [13]: mse_train = [0]*len(degrees)
mse_test = [0]*len(degrees)

for i, degree in enumerate(degrees):
    # Recompute the degree-d poly transform if you didn't save it!
    poly = PolynomialFeatures(degree=degree, include_bias=(degree == 0))
    Xt_poly = poly.fit_transform(Xt)
    Xe_poly = poly.transform(Xe)

    scaler = StandardScaler().fit(Xt_poly)
    Xt_poly_scaled = scaler.transform(Xt_poly)
    Xe_poly_scaled = scaler.transform(Xe_poly)
```

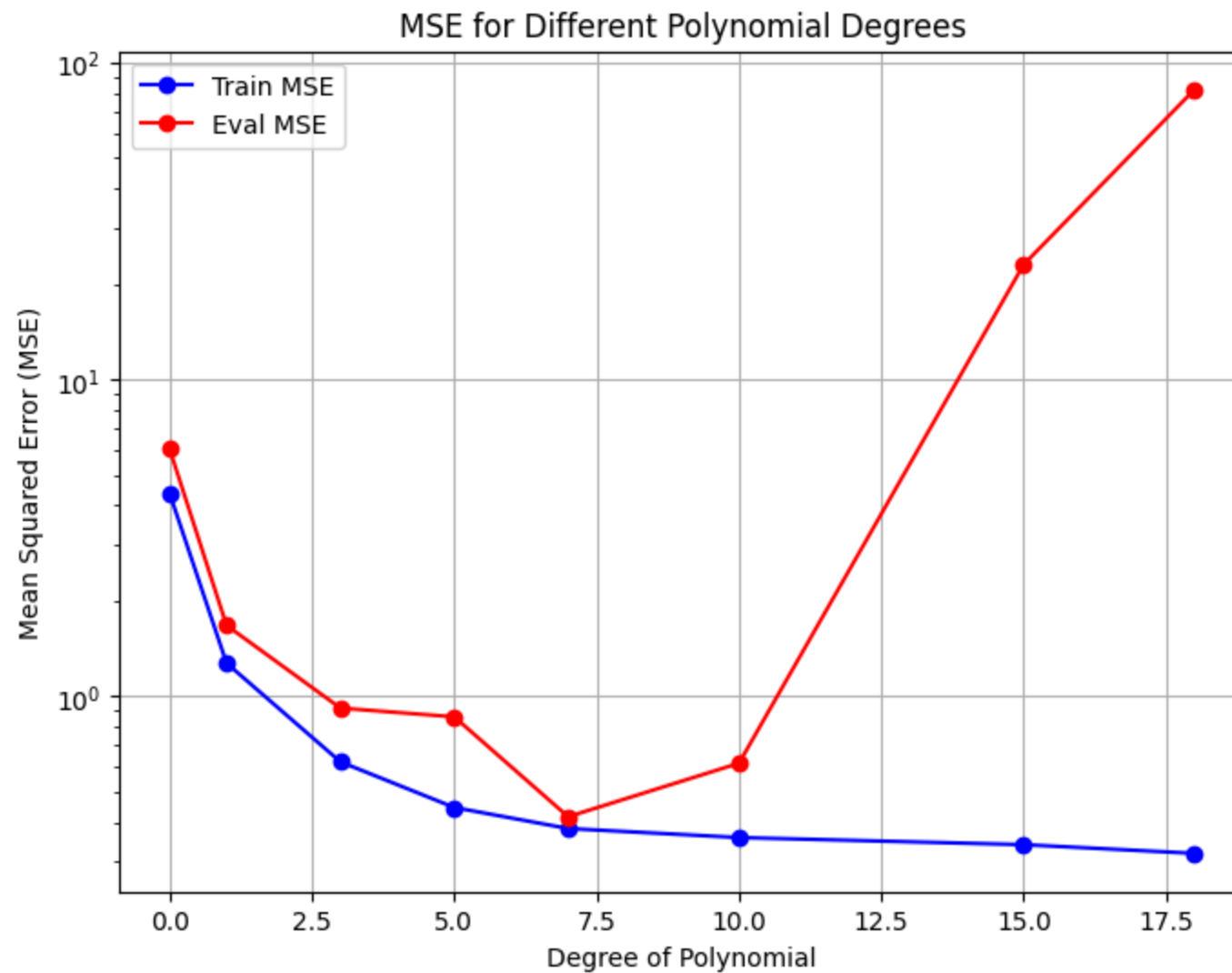
```
lr = learners[i]

Yt_pred = lr.predict(Xt_poly_scaled)
Ye_pred = lr.predict(Xe_poly_scaled)

mse_train[i] = mse(Yt, Yt_pred)
mse_test[i] = mse(Ye, Ye_pred)

plt.figure(figsize=(8, 6))
plt.semilogy(degrees, mse_train, label="Train MSE", marker='o', color='blue')
plt.semilogy(degrees, mse_test, label="Eval MSE", marker='o', color='red')

plt.xlabel('Degree of Polynomial')
plt.ylabel('Mean Squared Error (MSE)')
plt.title('MSE for Different Polynomial Degrees')
plt.legend()
plt.grid(True)
plt.show()
```



P3.3 Model Selection

Which degree would you select to use?

```
In [14]: best_degree = degrees[np.argmin(mse_test)]  
         print(f"The best degree based on evaluation MSE is: {best_degree}")
```


The best degree based on evaluation MSE is: 7

P4: Cross-validation

Cross validation is another method of model complexity assessment. We use it only to determine the correct setting of complexity-altering parameters ("hyperparameters"), such as how many and which features to use, or parameters like "k" in KNN, for which training error alone provides little information. In particular, cross validation will not produce a specific model (parameter values), only a setting of the hyperparameter values that cross-validation thinks will lead to a model (parameter values) with low test error.

P4.1: 5-Fold Cross-validation

In the previous problem, we decided what degree of polynomial fit to use based on the performance on a held-out set of test data. Now suppose that we do not have access to the target values of those data. How can we determine the best degree?

We could perform another split; but since this is reducing the number of data available, let us instead use cross-validation to evaluate the degrees. Cross-validation works by splitting the training data X_T multiple times, one for each of the K partitions (`n_splits` in the code), and repeat our entire training and evaluation procedure on each split:

```
In [15]: mse_xval = [ 0. ]*len(degrees)

for j,degree in enumerate(degrees):  # loop over desired degree values

    ### YOUR CODE STARTS HERE ###

    xval = KFold(n_splits = 5)        # split into k=5 splits
    mse_folds = []

    for train_index, val_index in xval.split(Xt):
        # Extract the ith cross-validation fold (training/validation split)
        Xti,Xvi,Yti,Yvi = Xt[train_index],Xt[val_index],Yt[train_index],Yt[val_index]

        # Now, build the model:
        # Create a polynomial feature expansion
        poly = PolynomialFeatures(degree=degree)
```

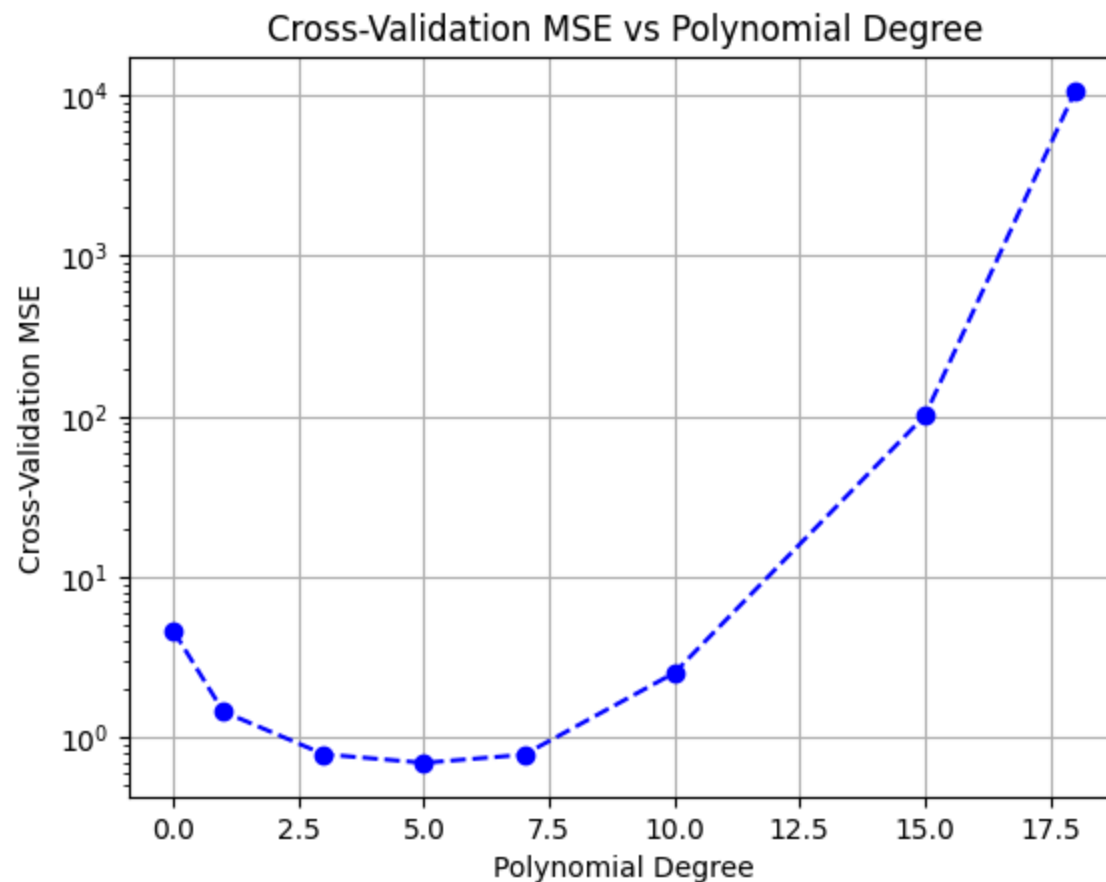
```
Xti_poly = poly.fit_transform(Xti)
Xvi_poly = poly.transform(Xvi)
# Create a StandardScaler
scaler = StandardScaler()
Xti_poly = scaler.fit_transform(Xti_poly)
Xvi_poly = scaler.transform(Xvi_poly)
# Fit the linear regression model on the training folds, Xti/Yti
model = LinearRegression()
model.fit(Xti_poly, Yti)
# Compute the MSE on the evaluation fold, Xvi/Yvi
Yvi_pred = model.predict(Xvi_poly)
mse_folds.append(mse(Yvi, Yvi_pred))

mse_xval[j] = np.mean(mse_folds)

# Evaluate the quality of this degree by averaging the MSE across the five folds

# Plot the estimated MSE from cross-validation as a function of the degree
plt.figure()
plt.semilogy(degrees, mse_xval, marker='o', linestyle='--', color='b')
plt.xlabel('Polynomial Degree')
plt.ylabel('Cross-Validation MSE')
plt.title('Cross-Validation MSE vs Polynomial Degree')
plt.grid(True)
plt.show()

### YOUR CODE ENDS HERE ###
```



P4.2: Cross-validation model selection

What degree would you choose based on the cross validation performance?

```
In [16]: best_degree = degrees[np.argmin(mse_xval)]  
         print(f"The best degree based on cross-validation performance is: {best_degree}")
```

The best degree based on cross-validation performance is: 5

P4.3 Comparison to test performance

How do the MSE estimates from 5-fold cross-validation compare to the estimated test performance you found from your held-out data, X_E ? Explain briefly.

The cross-validation process aims to find the degree that generalizes best across multiple splits, while the test set performance reflects how well the model fits one specific dataset. The degree 7 model may be slightly overfitting to the test set, whereas degree 5, chosen by cross-validation, likely offers better generalization across different datasets. Therefore, degree 5 is a more reliable choice based on cross-validation, even though degree 7 performs slightly better on the specific test set.

Statement of Collaboration (5 points)

It is **mandatory** to include a Statement of Collaboration in each submission, with respect to the guidelines below. Include the names of everyone involved in the discussions (especially in-person ones), and what was discussed.

All students are required to follow the academic honesty guidelines posted on the course website. For programming assignments, in particular, I encourage the students to organize (perhaps using EdD) to discuss the task descriptions, requirements, bugs in my code, and the relevant technical content before they start working on it. However, you should not discuss the specific solutions, and, as a guiding principle, you are not allowed to take anything written or drawn away from these discussions (i.e. no photographs of the blackboard, written notes, referring to EdD, etc.). Especially after you have started working on the assignment, try to restrict the discussion to EdD as much as possible, so that there is no doubt as to the extent of your collaboration.

N/A