

Part 2: Coding Task

EE214 - Assignment 1
(2025 Spring)

March 25, 2025

1 Overview

In Part 2 of the assignment, you are required to complete the task below, which is self-contained and focuses on exploring model performance with various basis functions. You are expected to use AI tools (e.g., ChatGPT, GitHub Copilot) only as an aid. You must document all AI assistance in your final report.

2 Comparative Study of Basis Functions on the Concrete Compressive Strength Dataset

2.1 Objective

Compare the performance of two different basis functions—polynomial and Gaussian—on a real-world regression task. Analyze how each approach affects the bias-variance tradeoff and overall prediction accuracy.

2.2 Data Source

Download the Concrete Compressive Strength dataset from the UCI Machine Learning Repository:

<https://archive.ics.uci.edu/ml/datasets/Concrete+Compressive+Strength>

Save the CSV file (e.g., `Concrete.Data.xls`) in your working directory.

2.3 Steps and Requirements

1. Data Loading and Preprocessing:

- Load the dataset using `pandas`.
- Handle missing values if present.
- Standardize the features.
- Split the data into training (80%) and test (20%) sets.

2. Modeling Approaches:

(a) Polynomial Regression Model:

For various degrees (e.g., 1, 3, 5, 7), transform the features using polynomial expansion and fit a linear regression model.

(b) Gaussian Basis Regression Model:

For various numbers of basis functions (e.g., 5, 7, 10, 15), create a design matrix using Gaussian basis functions. Use centers $\mu_j = j$ and $\sigma = 1.0$. Fit a linear regression model on the transformed data.

3. Evaluation:

For both approaches, compute the training and test MSE. Investigate how increasing model complexity (in polynomial regression, increasing degree; in Gaussian basis, varying the number of basis functions) affects performance.

4. Visualization:

Plot error curves (MSE vs. model complexity) for both approaches on the same graph for comparison.

5. Analysis Questions:

In your written report, answer:

- Which basis function (polynomial or Gaussian), and with which degree, provided better test performance and why?
- How did the training and test errors change as model complexity increased for each method? Based on this, identify examples of underfitting and overfitting.
- Discuss the bias-variance tradeoff observed in both approaches.
- What are the potential advantages and disadvantages of each basis function in real-world regression problems?
- Apply a **different basis function** from those used so far (polynomial or Gaussian), and analyze it in terms of the bias-variance tradeoff and the pros and cons in real-world regression problems. You may also try regularization, i.e., add a penalty term in the loss function.

2.4 Code Skeleton Example

```
1 import numpy as np
2 import pandas as pd
3 import matplotlib.pyplot as plt
4 from sklearn.linear_model import LinearRegression
5 from sklearn.metrics import mean_squared_error
6 from sklearn.model_selection import train_test_split
7 from sklearn.preprocessing import StandardScaler
8
9 % Load the dataset
10 data = pd.read_excel('Concrete_Data.xls')
11 print(data.head())
12
```

```

13 % Preprocess the data
14 X = data.drop('Concrete compressive strength(MPa, megapascals) ', axis=1).values
15 y = data['Concrete compressive strength(MPa, megapascals) '].values
16
17 # ----- Standardize features -----
18 # scaler = ...
19 # X_scaled = ...
20
21 # ----- Split the data -----
22 # X_train, X_test, y_train, y_test = ...
23
24 % (a) Polynomial Regression
25 from sklearn.preprocessing import PolynomialFeatures
26 degrees = [1, 3, 5, 7]
27 poly_train_errors = []
28 poly_test_errors = []
29
30 for degree in degrees:
31     # ----- Blank: Transform features into polynomial features -----
32     # Create a PolynomialFeatures object with the given degree (include bias term)
33     # Transform both the training and test data
34
35     # TODO: Fill in the code below
36     # poly = ...
37     # X_train_poly = ...
38     # X_test_poly = ...
39     # -----
40
41     # ----- Blank: Fit linear regression model and evaluate performance -----
42     # TODO:
43     # - Initialize and train LinearRegression model
44     # - Predict on both training and test sets
45     # - Compute train and test MSE
46     # - Append errors to respective lists
47
48     # model_poly = ...
49     # model_poly.fit(...)
50     # y_train_pred = ...
51     # y_test_pred = ...
52     # train_mse = ...
53     # test_mse = ...
54     # poly_train_errors.append(...)
55     # poly_test_errors.append(...)
56     # -----
57
58 % (b) Gaussian Basis Regression
59 def create_gaussian_design_matrix(x, degree, sigma=1.0):
60     """
61     Create a design matrix using Gaussian basis functions.
62     Args:
63         x: numpy array of shape (n_samples,)
64         degree: number of Gaussian basis functions
65         sigma: width of the Gaussian functions
66     Returns:
67         X: design matrix with shape (n_samples, degree + 1) including bias term
68     """

```

```

69     n_samples = x.shape[0]
70     X = np.ones((n_samples, degree + 1))
71     # ----- Blank: Fill in Gaussian basis transformation -----
72     # TODO: Fill in the loop to compute Gaussian basis functions
73     # for j in range(1, degree + 1):
74     #     mu = ...
75     #     X[:, j] = ...
76     # -----
77     return X
78
79 % For simplicity, use one feature (the first column) from the training set
80 x_train_feature = X_train[:, 0]
81 x_test_feature = X_test[:, 0]
82
83 gaussian_degrees = [5, 7, 10, 15]
84 gauss_train_errors = []
85 gauss_test_errors = []
86
87 for degree in gaussian_degrees:
88     X_train_gauss = create_gaussian_design_matrix(x_train_feature, degree, sigma
89 =1.0)
90     X_test_gauss = create_gaussian_design_matrix(x_test_feature, degree, sigma
91 =1.0)
92
93     # --- Blank: Fit Gaussian basis regression model and evaluate performance ---
94     # TODO:
95     # - Initialize and train LinearRegression model
96     # - Predict on both training and test sets
97     # - Compute train and test MSE
98     # - Append errors to respective lists
99
100     # model_gauss = ...
101     # model_gauss.fit(...)
102     # y_train_pred = ...
103     # y_test_pred = ...
104     # train_mse = ...
105     # test_mse = ...
106     # gauss_train_errors.append(...)
107     # gauss_test_errors.append(...)
108     # -----
109
110 % Plot comparison of test errors
111 plt.figure(figsize=(10,6))
112 plt.plot(degrees, poly_test_errors, 'ro-', label='Polynomial Regression Test MSE')
113 plt.plot(gaussian_degrees, gauss_test_errors, 'bo-', label='Gaussian Basis Test
114 MSE')
115 plt.xlabel('Model Complexity (Degree or Number of Basis Functions)')
116 plt.ylabel('Test Mean Squared Error')
117 plt.title('Comparison of Basis Functions on Concrete Data')
118 plt.legend()
119 plt.show()

```

Listing 1: Python Code for Part 2

3 Submission Guidelines for Part 2

- Submit a complete Jupyter Notebook containing your code, plots, and experimental results. Ensure the notebook runs end-to-end with a fixed random seed (where applicable) for reproducibility.
- Provide a written report that includes answers to the analysis questions, along with a thoughtful reflection on your work for Part 2.
- Clearly document every instance of AI assistance used (if any). Include which tool was used, what you asked, and how you integrated the response.

Good luck and happy coding!