

Polynomial Regression and Visualization Using LAPACK and Cairo in C with Python Data Generation

Amaan Rahman

Introduction

In this project, a set of n sample points (x_i, y_i) are stored in a text file. The objective is to compute an optimal n -degree polynomial that best fits the given data in the least-squares sense. The computation of polynomial coefficients is performed in C using the high-performance numerical routines provided by **BLAS** and **LAPACK**, ensuring both efficiency and numerical stability.

The dataset is generated in Python, while the regression computation and visualization are done in C. The **Cairo Graphics Library** is used to plot the sample points along with the fitted polynomial curve, providing a clear visual representation of the model's accuracy.

Objective

To compute an n -degree polynomial regression that best fits a noisy dataset, using optimized numerical routines, and visualize the dataset along with the true polynomial curve.

Mathematical Formulation

For polynomial regression of degree n , the model is:

$$y = c_0x^n + c_1x^{n-1} + \cdots + c_{n-1}x + c_n.$$

In matrix form:

$$Y = X\beta + E,$$

where X is the Vandermonde matrix and E is the error vector.

Minimizing the squared error gives the Normal Equation:

$$X^T X \beta = X^T Y.$$

LAPACK solves this efficiently without computing the inverse.

OLS Formulation Using LAPACK

We represent the system in matrix form:

$$Y = X\beta + E, \quad E = Y - X\beta$$

Minimizing the sum of squared errors gives:

$$E^T E = (Y - X\beta)^T (Y - X\beta)$$

Differentiating and setting the gradient to zero:

$$X^T X \beta = X^T Y$$

Hence, the optimal parameter vector is:

$$\beta = (X^T X)^{-1} X^T Y$$

Rather than explicitly computing the inverse, LAPACK provides an efficient and numerically stable way to solve this system.

Implementation via LAPACK

Using LAPACK and BLAS, the steps are:

1. Compute $X^T X$ using `dgemm`.
2. Compute $X^T Y$ using `dgemv`.
3. Solve $(X^T X)\beta = X^T Y$ using `dposv` or `dgesv`.

Advantages of LAPACK Implementation

- **Efficiency:** Optimized for modern CPUs and parallel computation.
- **Numerical Stability:** Reduces rounding errors common in manual inversion.
- **Ease of Integration:** Directly linkable in C for high-performance computation.

Python Implementation for Generating Random Points

Part 1: Polynomial & Random Point Generation

```
1 import numpy as np
2 import random
3 import sympy as sp
4 import matplotlib.pyplot as plt
5 from scipy.optimize import curve_fit
6
7 def generate_points_and_polynomial():
8     n = int(input("Enter the polynomial order n: "))
9     x_vals = np.array([random.uniform(-10, 10) for _ in range(n)])
10    y_vals = np.array([random.uniform(-10, 10) for _ in range(n)])
11    coeffs = np.polyfit(x_vals, y_vals, n)
12    x = sp.Symbol('x')
13    poly_expr = sum(sp.Float(coeffs[i]) * x**(n - i) for i in range(n + 1))
14
15    print("\nGenerated Points (x, y):")
16    for xi, yi in zip(x_vals, y_vals):
17        print(f"({xi:.4f}, {yi:.4f})")
18
19    print("\nPolynomial coefficients:")
20    print(coeffs)
21
22    print("\nPolynomial (Symbolic):")
23    print(sp.expand(poly_expr))
24
25    return n, coeffs, poly_expr
```

Part 2: Synthetic Data

```
1 def generate_synthetic_data(poly_expr, x_range=(-5, 5), n_points=100,
2   noise_std=1.0):
3     x = np.linspace(x_range[0], x_range[1], n_points)
4     y_true = np.array([poly_expr.subs('x', val) for val in x], dtype=float)
5     y_noisy = y_true + np.random.normal(0, noise_std, size=len(x))
6     return x, y_noisy, y_true
7
8 def build_curve_fit_model(order):
9     def model(x, *params):
10         return sum(params[i] * x**(order - i) for i in range(order + 1))
11     return model
12
13 def main():
14     n, coeffs, poly_expr = generate_points_and_polynomial()
15     x_data, y_noisy, y_true = generate_synthetic_data(poly_expr, n_points
16       =200, noise_std=2.0)
17     model = build_curve_fit_model(n)
18     initial_guess = np.ones(n + 1)
19     popt, pcov = curve_fit(model, x_data, y_noisy, p0=initial_guess)
20     y_fit = model(x_data, *popt)
21
22     print("\nFitted Parameters:")
23     print(popt)
24
25     filename = "synthetic_polynomial_data.txt"
26     with open(filename, "w") as f:
27         f.write(f"Polynomial order: {n}\n")
28         f.write(f"Number of data points: {len(x_data)}\n")
29         f.write(f"x\ty_noisy\ty_true\n")
30         for xi, yi_noisy, yi_true in zip(x_data, y_noisy, y_true):
31             f.write(f"{xi:.6f}\t{yi_noisy:.6f}\t{yi_true:.6f}\n")
32
33     print(f"\nData saved to '{filename}'.")
34
35     plt.figure(figsize=(10, 6))
36     plt.scatter(x_data, y_noisy, label="Noisy data", s=15)
37     plt.plot(x_data, y_true, label="True polynomial", linewidth=2)
38     plt.plot(x_data, y_fit, label="Fitted curve", linestyle='--',
39       linewidth=2)
40     plt.title(f"Polynomial Order {n}: True vs Fitted Curve")
41     plt.xlabel("x")
42     plt.ylabel("y")
43     plt.legend()
44     plt.grid(True)
45     plt.show()
46
47 main()
```

C Implementation for Solving Polynomial Coefficients Using LAPACK

C Program (Part 1): Reading Data and Memory Setup

```
1  #include <stdio.h>
2  #include <stdlib.h>
3  #include <string.h>
4  #include <math.h>
5  #include <Accelerate/Accelerate.h>
6
7  /* Safe malloc */
8  void *xmalloc(size_t s) {
9      void *p = malloc(s);
10     if (!p) { printf("Out of memory!\n"); exit(1); }
11     return p;
12 }
13
14 /* Read data file generated by Python */
15 int read_data(const char *filename, int *order, int *N,
16              double **x, double **y_noisy, double **y_true)
17 {
18     FILE *f = fopen(filename, "r");
19     if (!f) { printf("Error: cannot open file %s\n", filename); return -1; }
20     char line[512];
21
22     if (!fgets(line, sizeof(line), f)) return -1;
23     sscanf(line, "Polynomial order: %d", order);
24     if (!fgets(line, sizeof(line), f)) return -1;
25     sscanf(line, "Number of data points: %d", N);
26     fgets(line, sizeof(line), f); /* Skip header row */
27
28     *x = xmalloc(sizeof(double) * (*N));
29     *y_noisy = xmalloc(sizeof(double) * (*N));
30     *y_true = xmalloc(sizeof(double) * (*N));
31
32     for (int i = 0; i < *N; i++) {
33         double xv, yn, yt;
34         if (!fgets(line, sizeof(line), f)) { printf("Unexpected EOF at\nline %d\n", i); return -1; }
35         if (sscanf(line, "%lf %lf %lf", &xv, &yn, &yt) != 3) { printf("Bad\n data line: %s\n", line); return -1; }
36         (*x)[i] = xv; (*y_noisy)[i] = yn; (*y_true)[i] = yt;
37     }
38     fclose(f);
39     return 0;
40 }
```

C Program (Part 2): Vandermonde Matrix Construction

```
1  /* Build Vandermonde matrix (column-major for BLAS) */
2  void build_vandermonde(int N, int deg, const double *x, double *X)
3  {
4      int P = deg + 1;
5      for (int j = 0; j < P; j++) {
6          int power = deg - j;
7          for (int i = 0; i < N; i++) {
8              X[j*N + i] = pow(x[i], power);
9          }
10     }
11 }
```

C Program (Part 3): LAPACK Solution and Coefficient Extraction

```

1  int main(int argc, char **argv)
2  {
3      const char *infile = "synthetic_polynomial_data.txt";
4      if (argc > 1) infile = argv[1];
5
6      int N, deg;
7      double *x, *y_noisy, *y_true;
8
9      printf("Reading: %s\n", infile);
10
11     if (read_data(infile, &deg, &N, &x, &y_noisy, &y_true) != 0) {
12         printf("File read error.\n");
13         return 1;
14     }
15
16     printf("Loaded %d points, degree = %d\n", N, deg);
17
18     int P = deg + 1;
19     double *X = xmalloc(sizeof(double) * N * P);
20     double *XtX = xmalloc(sizeof(double) * P * P);
21     double *Xty = xmalloc(sizeof(double) * P);
22
23     build_vandermonde(N, deg, x, X);
24
25     /* XtX = X^T X */
26     cblas_dgemm(CblasColMajor, CblasTrans, CblasNoTrans,
27                P, P, N, 1.0,
28                X, N, X, N,
29                0.0, XtX, P);
30
31     /* Xty = X^T y */
32     cblas_dgemv(CblasColMajor, CblasTrans,
33                N, P, 1.0,
34                X, N, y_noisy, 1,
35                0.0, Xty, 1);
36
37     /* Solve (XtX) beta = Xty with LAPACK */
38     int n = P, nrhs = 1, lda = P, ldb = P, info;
39     char uplo = 'U';
40
41     dposv_(&uplo, &n, &nrhs, XtX, &lda, Xty, &ldb, &info);
42
43     if (info != 0) {
44         printf("Error: dposv failed with info=%d\n", info);
45         return 1;
46     }
47
48     printf("\nOptimal Polynomial Coefficients:\n");
49     for (int i = 0; i < P; i++) {
50         printf("c[%d] = %.12f\n", i, Xty[i]);
51     }
52
53     FILE *out = fopen("optimal_coeffs.txt", "w");
54     for (int i = 0; i < P; i++)
55         fprintf(out, "%.12f\n", Xty[i]);
56     fclose(out);
57     printf("Saved coefficients to optimal_coeffs.txt\n");
58
59     free(x); free(y_noisy); free(y_true);

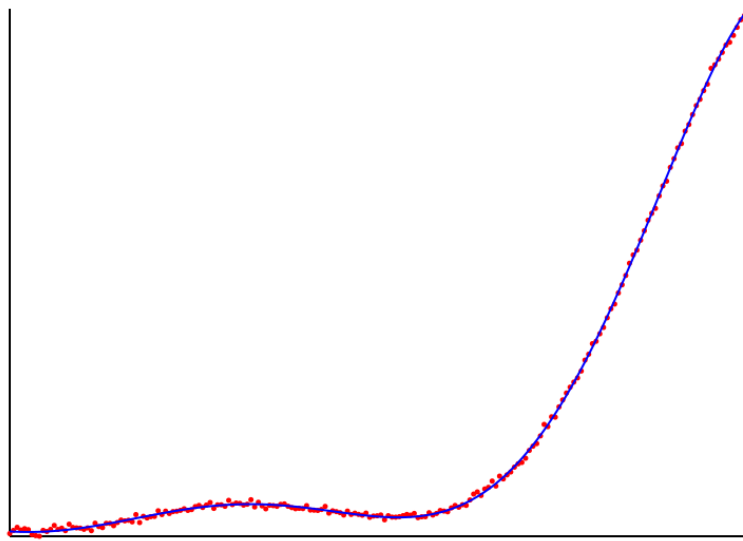
```

```
00     free(X); free(XtX); free(Xty);  
01     return 0;  
02 }
```


Visualization using Cairo Library in C

After computing the regression line, visualization is done using the Cairo graphics library. Steps include:

- Initialize a Cairo surface and context using `cairo_create()`.
- Plot all n sample points as small circles/dots using `cairo_arc()` and `cairo_fill()`.
- Draw the polynomial curve using `cairo_move_to()` and `cairo_line_to()`.
- Label endpoints and annotate minimal distance using `cairo_show_text()`.
- Save the image to PNG using `cairo_surface_write_to_png()`.



Example plot for an 8th order polynomial using Cairo in C

Input/Output Format

Input:

- Text file containing n sample points (x_i, y_i) .

Output:

- Optimal coefficients β minimizing the squared error.
- Visualization of data points and best-fit polynomial curve.