

Assignment 1: Linear Least Squares

Sanaz Marefati

September 20, 2025

Exercise 1: Linear Least Squares Polynomial Fitting

Problem Statement

Given data pairs $(x_i, y_i)_{i=1}^m$ where $y_i = \Gamma(x_i)$ (noise-free), we perform linear least squares polynomial fitting using:

1. Monomials $(1, x, x^2, \dots, x^n)$ with normalization.
2. Chebyshev polynomials as the basis.

Consider three datasets:

1. $\Gamma(x) = \sin(10x)$ over $[-1, 1]$.
2. $\Gamma(x) = 1/(1 + 25x^2)$ over $[-5, 5]$.
3. $\Gamma(x) = 1$ if $x \in [0, 2]$, otherwise 0 on $[-2, 0]$.

For $m = 50, 200, 300$, plot the standard deviations σ and coefficients of determination R^2 with varying n (from 1 to 30). Compare the performance of the two methods.

Solution

I implement the solution using jax for numerical computations and Matplotlib for plotting. The code generates random points and fits polynomials using both monomial and Chebyshev bases.

```

1 import jax.numpy as jnp
2 import jax.random as random
3 import matplotlib.pyplot as plt
4 from numpy.polynomial.chebyshev import chebvander
5 import numpy as onp
6
7 # ----- functions -----
8 def f1(x): return jnp.sin(10 * x)
9 def f2(x): return 1.0 / (1.0 + 25.0 * x**2)
10 def f3(x): return jnp.where((x >= 0.0) & (x <= 2.0),
11     1.0, 0.0)
12
13 funcs = [
14     ("f1: sin(10x)", f1, -1.0, 1.0),
15     ("f2: Runge function", f2, -5.0, 5.0),
16     ("f3: Heaviside step function", f3, -2.0, 2.0),
17 ]
18
19 def model(theta, x):
20     powers = jnp.arange(len(theta))
21     return jnp.sum(theta * (x[:, None] ** powers),
22                    axis=1)
23
24 def rmse(y_true, y_pred):
25     return jnp.sqrt(jnp.mean((y_true - y_pred)**2))
26
27 def r2_score(y_true, y_pred):
28     ss_res = jnp.sum((y_true - y_pred)**2)
29     ss_tot = jnp.sum((y_true - jnp.mean(y_true))**2)
30     return jnp.where(ss_tot == 0, 1.0, 1.0 - ss_res /
31                      ss_tot)
32
33 def cond_number(X):
34     return float(jnp.linalg.cond(X))
35
36
37 # ----- settings -----
38 key = random.PRNGKey(0)
39 m_values = [50, 200, 300]
40 n_values = list(range(1, 31)) # degrees 1..30
41
42 # ----- main -----

```

```

40 for m in m_values:
41     # precompute metrics for each function
42     all_sig_mono, all_r2_mono, all_kappa_mono = [], [], []
43     all_sig_cheb, all_r2_cheb, all_kappa_cheb = [], [], []
44     func_names = []
45
46     for fname, f, a, b in funcs:
47         # sample data for this function & m
48         key, subkey = random.split(key)
49         x_data = random.uniform(subkey, shape=(m,), minval=a, maxval=b)
50         #x_data = jnp.linspace(a, b, m)
51         y_data = f(x_data) # true y from original x
52
53         # normalize x for stability (basis only)
54         x_norm = 2.0 * (x_data - a) / (b - a) - 1
55
56         sig_mono, r2_mono, kap_mono = [], [], []
57         sig_cheb, r2_cheb, kap_cheb = [], [], []
58
59         for n in n_values:
60             # ----- Monomial basis -----
61             X_monomial = jnp.vander(x_norm, N=n + 1,
62                                     increasing=True)
63             theta_m = jnp.linalg.pinv(X_monomial) @ y_data
64             try:
65                 theta_m = jnp.linalg.lstsq(X_monomial,
66                                         y_data, rcond=None)[0]
67             except AttributeError:
68                 theta_m = jnp.linalg.pinv(X_monomial) @
69                 y_data
70             y_hat_m = model(theta_m, x_norm)
71             sig_mono.append(float(rmse(y_data,
72                                     y_hat_m)))
73             r2_mono.append(float(r2_score(y_data,
74                                     y_hat_m)))
75             kap_mono.append(cond_number(X_monomial))
76
77
78
79
80
81
82
83
84
85
86
87
88
89
90
91
92
93
94

```

```

75      # ----- Chebyshev basis -----
76      #X_cheb = chebvander(x_norm, n)
77      X_cheb =
78          jnp.asarray(chebvander(onp.asarray(x_norm),
79                         n))
80      theta_c = jnp.linalg.pinv(X_cheb) @ y_data
81      # predict using Chebyshev columns (linear
82          comb of columns)
83      y_hat_c = X_cheb @ theta_c
84      sig_cheb.append(float(rmse(y_data,
85                               y_hat_c)))
86      r2_cheb.append(float(r2_score(y_data,
87                               y_hat_c)))
88      kap_cheb.append(cond_number(X_cheb))

89
90      all_sig_mono.append(sig_mono);
91      all_r2_mono.append(r2_mono);
92      all_kappa_mono.append(kap_mono)
93      all_sig_cheb.append(sig_cheb);
94      all_r2_cheb.append(r2_cheb);
95      all_kappa_cheb.append(kap_cheb)
96      func_names.append(fname)

97      # ----- Plotting -----
98      fig, axes = plt.subplots(3, 3, figsize=(16, 14),
99                             sharex=True)
100     fig.suptitle(f"m = {m}: Monomial (Normalized) vs
101                  Chebyshev , R , and vs degree",
102                  fontsize=18)
103
104     for i, fname in enumerate(func_names):
105         # RMSE ( )
106         ax_rmse = axes[i, 0]
107         ax_rmse.plot(n_values, all_sig_mono[i],
108                       marker='o', label="Monomial (Normalized)")
109         ax_rmse.plot(n_values, all_sig_cheb[i],
110                       marker='s', label="Chebyshev")
111         ax_rmse.set_xlabel("Polynomial degree (n)")
112         ax_rmse.set_ylabel(" (RMSE)")
113         ax_rmse.set_title(f"{fname} RMSE, (m =
114                           {m})")
115         ax_rmse.grid(True, alpha=0.3)

```

```

103     ax_rmse.legend()
104
105     # R
106     ax_r2 = axes[i, 1]
107     ax_r2.plot(n_values, all_r2_mono[i],
108                 marker='o', label="Monomial (Normalized)")
109     ax_r2.plot(n_values, all_r2_cheb[i],
110                 marker='s', label="Chebyshev")
111     ax_r2.set_xlabel("Polynomial degree (n)")
112     ax_r2.set_ylabel("R ")
113     ax_r2.set_title(f"{fname} R , (m = {m})")
114     ax_r2.grid(True, alpha=0.3)
115     ax_r2.legend()
116
117     # Condition number (X)
118     ax3 = axes[i, 2]
119     ax3.plot(n_values, all_kappa_mono[i],
120               marker='o', label="Monomial (Normalized)")
121     ax3.plot(n_values, all_kappa_cheb[i],
122               marker='s', label="Chebyshev")
123     ax3.set_xlabel("Polynomial degree (n)")
124     ax3.set_ylabel("Condition number")
125     ax3.set_title(f"{fname} Condition number,
126                   (m = {m})")
127     ax3.set_yscale("log") # often spans many
128                   # orders; log-scale helps
129     ax3.grid(True, which="both", alpha=0.3)
130     ax3.legend()
131
132     # set x-labels only on bottom row
133     #axes[2, 0].set_xlabel("Polynomial degree (n)")
134     #axes[2, 1].set_xlabel("Polynomial degree (n)")
135
136     plt.tight_layout(rect=[0, 0, 1, 0.96])
137     plt.savefig(f"results_m{m}.png", dpi=300)
138     plt.close(fig)

```

Results

Three plots for each m (one per dataset) comparing monomial and Chebyshev methods are shown in Figure 1, Figure 2, and Figure 3.

Analysis and Discussion

The results show clear differences between polynomial fitting with monomial and Chebyshev bases. For smooth functions such as $\sin(10x)$, both bases achieve low RMSE and high R^2 at moderate polynomial degree. However, as the degree n increases, the monomial basis quickly becomes unstable: the condition number of its Vandermonde matrix grows exponentially, leading to numerical errors and poor fits despite high model complexity. In contrast, Chebyshev polynomials maintain stability, keeping the condition number much smaller and producing consistent improvements in accuracy up to higher degrees.

The Runge function $1/(1 + 25x^2)$ highlights the classical Runge's phenomenon: monomials show large oscillations at the interval boundaries,

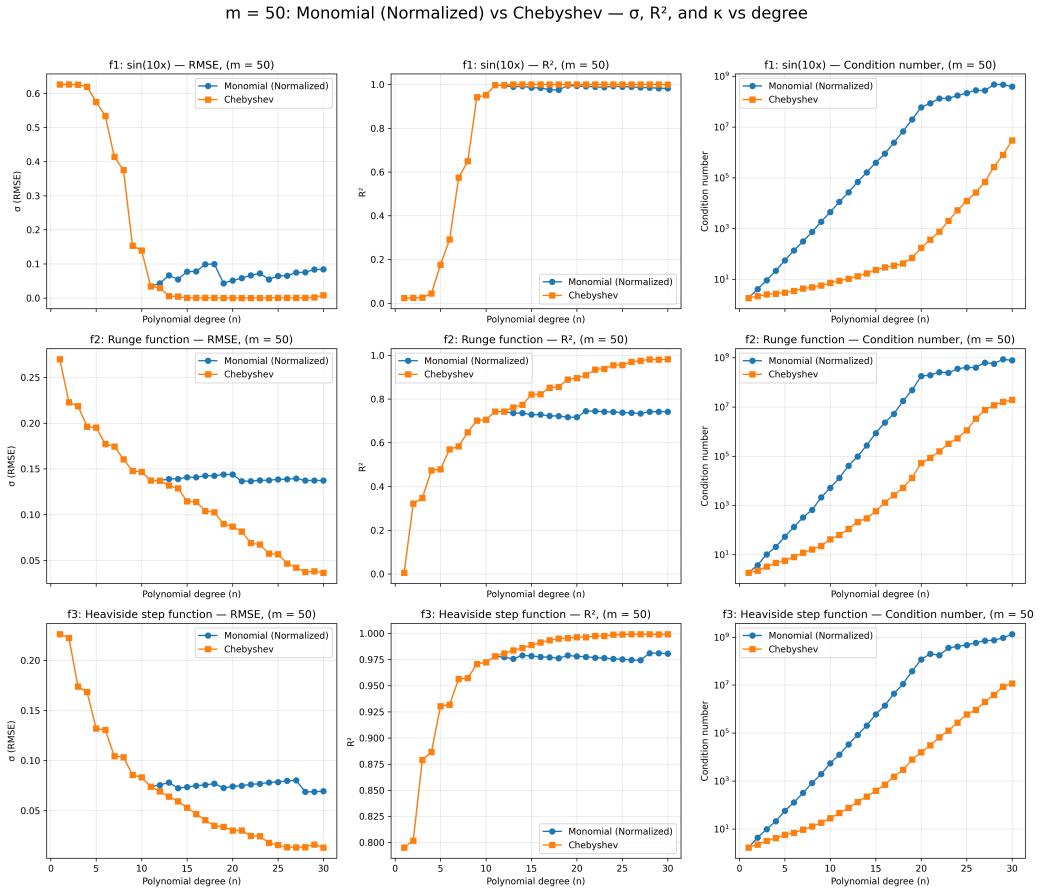


Figure 1: Results for $m = 50$. Standard deviation σ and R^2 vs. polynomial degree n .

$m = 200$: Monomial (Normalized) vs Chebyshev — σ , R^2 , and κ vs degree

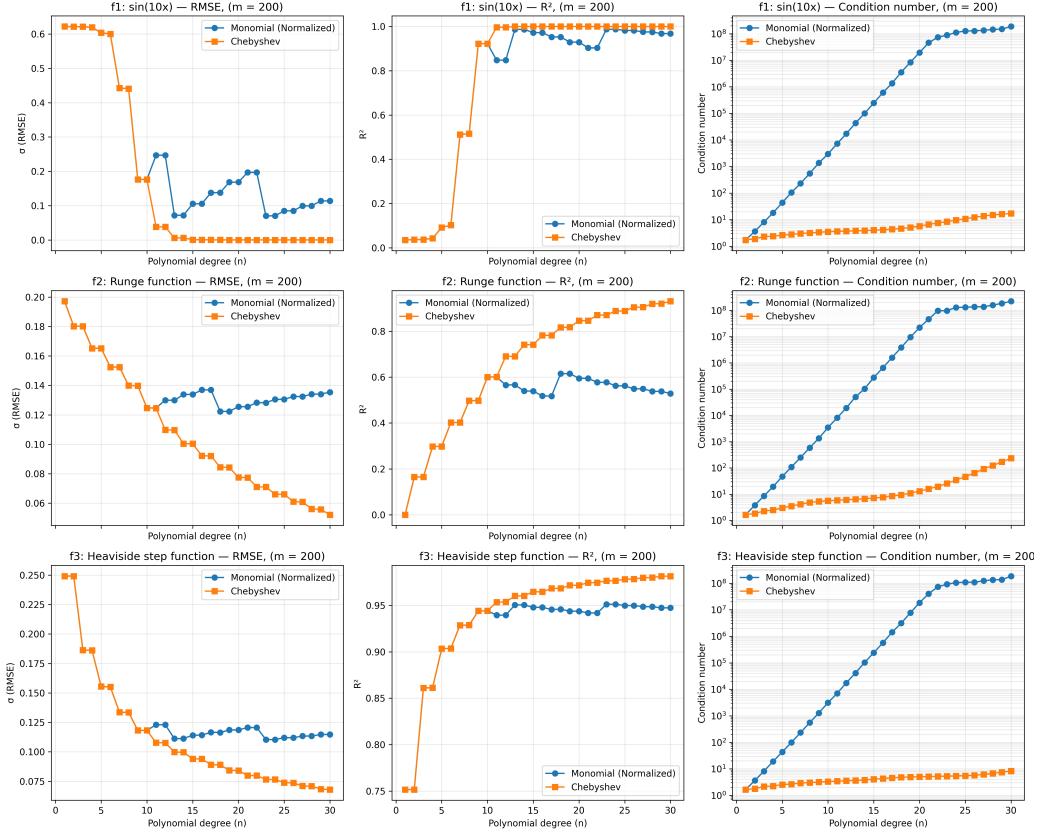


Figure 2: Results for $m = 200$. Standard deviation σ and R^2 vs. polynomial degree n .

whereas Chebyshev basis reduces these oscillations and achieves lower error. For the step function, neither basis performs well due to the discontinuity, but Chebyshev still provides a more controlled approximation. Increasing the sample size m improves results for both methods, but it does not resolve the instability of monomials at high degrees. Overall, the experiments demonstrate that basis choice is crucial: while monomials suffice at low degree, Chebyshev polynomials are far more reliable for higher-degree approximations.

$m = 300$: Monomial (Normalized) vs Chebyshev — σ , R^2 , and κ vs degree

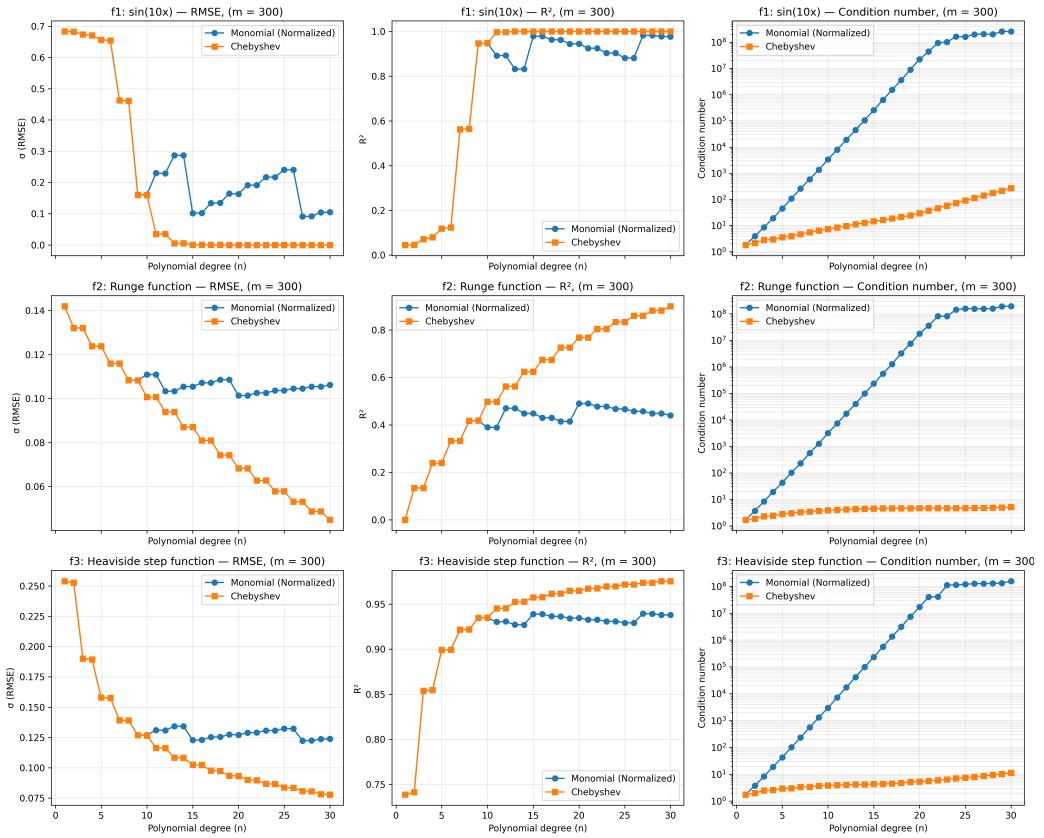


Figure 3: Results for $m = 300$. Standard deviation σ and R^2 vs. polynomial degree n .

Exercise 2: Understanding Linear Regression

Problem Statement

Minimize the loss function:

$$\min_{\theta \in \mathbb{R}^{n+1}} \sum_{i=1}^m r_i^2(\theta) w_i, \quad w_i > 0, \quad r_i(\theta) = y_i - M(x_i; \theta),$$

where $M(x; \theta) = \sum_{j=0}^n \theta_j e_j(x)$, $n \ll m$.

1. Show that the loss function can be expressed in the form of $\|A\theta - y\|^2$ for an appropriate diagonal matrix D . Present A , D , and y .
2. Derive the normal equation by setting the gradient with respect to θ equal to zero.
3. Compute the Hessian matrix of the loss function. What can we conclude about the convexity based on the Hessian matrix and the uniqueness of the solution?
4. (Optional) Can you give a condition on A such that the solution is unique?
5. (Optional) How does the Tikhonov regularization work in the above formulation?

Solution

(1) Expression of the Loss Function

The loss function is:

$$L(\theta) = \sum_{i=1}^m w_i r_i^2(\theta), \quad r_i(\theta) = y_i - M(x_i; \theta).$$

In matrix form: A is the $m \times (n + 1)$ design matrix with $A_{i,j} = e_j(x_i)$ (Vandermonde matrix for monomials), $y = [y_1, \dots, y_m]^T$, $D = \text{diag}(w_1, \dots, w_m)$. Thus:

$$L(\theta) = (y - A\theta)^T D(y - A\theta) = \|\sqrt{D}(A\theta - y)\|^2.$$

(2) Normal Equation

The gradient is:

$$\nabla_{\theta} L = -2A^T D(y - A\theta).$$

Setting it to zero:

$$A^T D A \theta = A^T D y.$$

The solution is $\theta^* = (A^T D A)^{-1} A^T D y$ if invertible.

(3) Hessian and Convexity

The Hessian is:

$$H = 2A^T D A.$$

Convexity: $z^T H z = 2(Az)^T D(Az) \geq 0$ (positive semi-definite). If A is full rank and D is positive definite, H is positive definite, making $L(\theta)$ strictly convex.

Uniqueness: Unique if H is invertible, i.e., A has full column rank.

(4) Condition for Uniqueness

A must have full column rank ($\text{rank}(A) = n + 1$), ensured by distinct x_i and $m > n$.

(5) Tikhonov Regularization

Add penalty:

$$L_{\text{reg}}(\theta) = (y - A\theta)^T D(y - A\theta) + \lambda \|\theta\|^2.$$

Gradient:

$$\nabla_{\theta} L_{\text{reg}} = -2A^T D(y - A\theta) + 2\lambda\theta = 0,$$

yielding:

$$(A^T D A + \lambda I)\theta = A^T D y.$$

This ensures invertibility and reduces θ magnitude, improving stability.

Exercise 3: Extension to Nonlinear Regression

Problem Statement

Minimize:

$$\min_{\theta} L(\theta) = \sum_{i=1}^m \rho(r_i(\theta)), \quad r_i(\theta) = y_i - M(x_i; \theta).$$

Use iterative methods with monomial basis for:

1. Huber loss.
2. l^p , $p = 1.5$.

Report max error $\max_{j=1}^{1000} |\Gamma(x_j) - M(x_j; \theta)|$ where $x_j = a + j(b - a)/1000$, $m = 100$.

Solution

```
1 import jax
2 import jax.numpy as jnp
3 import jax.random as random
4 import matplotlib.pyplot as plt
5 import csv
6 from dataclasses import dataclass
7
8 # Datasets (Gamma functions)
9 def f1(x): # sin(10x) on [-1, 1]
10     return jnp.sin(10.0 * x)
11
12 def f2(x): # Runge function on [-5, 5]
13     return 1.0 / (1.0 + 25.0 * x**2)
14
15 def f3(x): # Step function on [-2, 2]: 1 in [0,2] ,
16     else 0
17     return jnp.where((x >= 0.0) & (x <= 2.0), 1.0, 0.0)
18
19 @dataclass
20 class Dataset:
21     f: callable
22     a: float
23     b: float
24     name: str
```

```

24
25 DATASETS = [
26     Dataset(f=f1, a=-1.0, b= 1.0, name="sin(10x)") ,
27     Dataset(f=f2, a=-5.0, b= 5.0, name="Runge
28         function"),
29     Dataset(f=f3, a=-2.0, b= 2.0, name="Heaviside step
30         function"),
31 ]
32
33 # Losses
34 def huber_loss(r, delta=1.0):
35     # 0.5 r^2 for |r|<=delta, delta(|r| - 0.5 delta)
36     # otherwise
37     absr = jnp.abs(r)
38     quad = 0.5 * r**2
39     lin = delta * (absr - 0.5 * delta)
40     return jnp.where(absr <= delta, quad, lin)
41
42 # Names for plotting/CSV
43 LOSS_SPECS = [
44     ("huber", lambda r: huber_loss(r, delta=1.0)),
45     ("lp",    lambda r: lp_loss(r, p=1.5)),
46 ]
47
48 # Utilities
49 def normalize_to_unit_interval(x, a, b):
50     return 2.0 * (x - a) / (b - a) - 1.0
51
52 def vandermonde(x_norm, degree):
53     return jnp.vander(x_norm, degree + 1,
54                       increasing=True)
55
56 def initial_theta(A, y):
57     try:
58         return jnp.linalg.lstsq(A, y, rcond=None)[0]
59     except AttributeError:
60         return jnp.linalg.pinv(A) @ y
61
62 # IRLS core (weights per assignment: w = rho(r) / r^2)

```

```

62 def irls(A, y, rho_fn, max_iter=1000, tol=1e-8,
63     eps=1e-8):
64     """
65     Solve min sum_i rho(r_i) via IRLS:
66     - r = y - A @ theta
67     - w_i = rho(r_i) / (r_i^2 + eps)
68     - Solve weighted LS: min_theta || diag(sqrt(w))
69         (y - A theta) ||^2
70     """
71     theta = initial_theta(A, y)
72     for _ in range(max_iter):
73         r = y - A @ theta
74         rho_val = rho_fn(r)
75         w = rho_val / (r**2 + eps)
76         # Weighted least squares: minimize || W^(1/2)
77         # (y - A theta) ||_2
78         # Equivalent normal eq: (A^T W A) theta = A^T W
79         # y
80         # We'll form W*A and W*y via diag weights:
81         sqrtw = jnp.sqrt(w)
82         Aw = A * sqrtw[:, None]
83         yw = y * sqrtw
84         theta_new = jnp.linalg.pinv(Aw) @ yw
85
86         if jnp.linalg.norm(theta_new - theta) < tol:
87             return theta_new
88         theta = theta_new
89     return theta # reached max_iter
90
91 # Experiment configuration
92 m = 100
93 degrees = list(range(1, 31)) # 1..30
94 num_test = 1000
95 key = random.PRNGKey(0)
96
97 # Run experiments
98 results = [] # rows: (dataset, loss, degree, max_error)
99
100 for (loss_name, rho_fn) in LOSS_SPECS:
101     for ds in DATASETS:
102         a, b, f, name = ds.a, ds.b, ds.f, ds.name
103
104         # Training data

```

```

101     key, subkey = random.split(key)
102
103     x_train = jnp.linspace(a, b, m)
104     y_train = f(x_train)
105
106     # Normalize x to [-1,1] for stability; build
107     # test grid too
108     x_train_norm =
109         normalize_to_unit_interval(x_train, a, b)
110
111     #x_test = jnp.linspace(a, b, num_test)
112     x_test = random.uniform(subkey, shape=(m,), minval=a, maxval=b)
113     x_test_norm =
114         normalize_to_unit_interval(x_test, a, b)
115     y_test_true = f(x_test)
116
117     # Sweep degrees
118     for n in degrees:
119         # Design matrices
120         A = vandermonde(x_train_norm, n)
121         A_test = vandermonde(x_test_norm, n)
122
123         # Initial LS (warm start), then IRLS
124         theta0 = initial_theta(A, y_train)
125         # IRLS using the assignment's weight rule
126         theta = irls(A, y_train, rho_fn=rho_fn,
127                     max_iter=1000, tol=1e-8, eps=1e-8)
128
129         # Evaluate required metric
130         y_pred = A_test @ theta
131         err_max = jnp.max(jnp.abs(y_test_true -
132                               y_pred))
133         results.append((name, loss_name, int(n),
134                         float(err_max)))
135
136         # Console progress
137         print(f"{name:26s} | loss={loss_name:5s} |
138               n={n:2d} | max_err={float(err_max):.6e}")
139
140     # Save CSV of results
141     csv_path = "exercise3_results.csv"

```

```

136 with open(csv_path, "w", newline="") as f:
137     writer = csv.writer(f)
138     writer.writerow(["dataset", "loss", "degree_n",
139                      "max_error"])
140     writer.writerows(results)
141 print(f"\nSaved results to {csv_path}")
142
143 # Plot error vs degree
144 import pandas as pd
145
146 df = pd.DataFrame(results, columns=["dataset", "loss",
147                      "degree_n", "max_error"])
148 for ds in DATASETS:
149     sub = df[df["dataset"] == ds.name].copy()
150     # Pivot to have columns per loss
151     pivot = sub.pivot(index="degree_n", columns="loss",
152                        values="max_error").sort_index()
153
154     plt.figure(figsize=(5, 4))
155     if "huber" in pivot.columns:
156         plt.plot(pivot.index.values,
157                   pivot["huber"].values, marker="o",
158                   label="Huber")
159     if "lp" in pivot.columns:
160         plt.plot(pivot.index.values,
161                   pivot["lp"].values, marker="s",
162                   label="L^1.5")
163
164     plt.title(f"Max Error vs Degree      {ds.name}")
165     plt.xlabel("Polynomial degree n")
166     plt.ylabel("Max error on 1000 test points")
167     plt.grid(True, alpha=0.4)
168     plt.legend()
169     out_png =
170         f"exercise3_error_vs_degree_{ds.name.replace(
171             ', '_).replace('(', '')}.replace(')', '')}.png"
171     plt.tight_layout()
172     plt.savefig(out_png, dpi=300)
173     plt.close()
174 print(f"Saved plot: {out_png}")

```

Results

Figure 4 and Figure 5 compares Huber and $L^{1.5}$ losses at $n = 10$ with uniformly spaced training points and with randomly sampled training points. As shown in Figure 6 and Table 1, the maximum error decreases with increasing polynomial degree for the smooth and oscillatory functions, while the discontinuous case remains difficult.

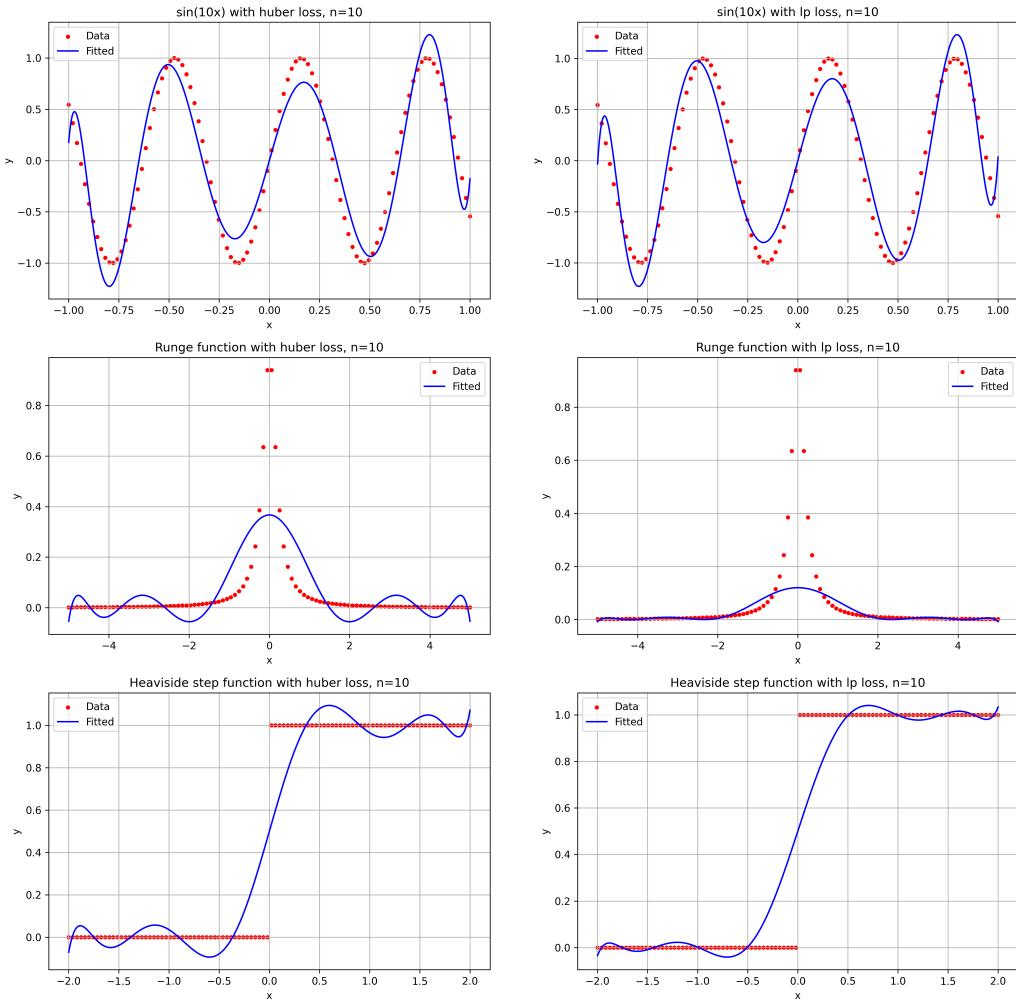


Figure 4: Comparison of Huber loss vs l^p with three functions ($n = 10$, with uniformly spaced dataset).

Analysis and Discussion

The results in Figure 6 and Table 1 show consistent trends across the three datasets. For the oscillatory $\sin(10x)$, the $L^{1.5}$ loss initially performs worse (error 0.510 at $n = 10$ versus 0.366 for Huber), but catches up at higher degrees, where both methods achieve errors below 0.1. For the Runge function, both losses reduce the maximum error steadily with polynomial degree, converging to values below 0.6 at $n = 30$. Although the Runge function $f(x) = 1/(1 + 25x^2)$ is smooth, it is notoriously difficult to approximate with high-degree polynomials on equispaced points. This classical issue, known as Runge's phenomenon, explains why the maximum error decreases only

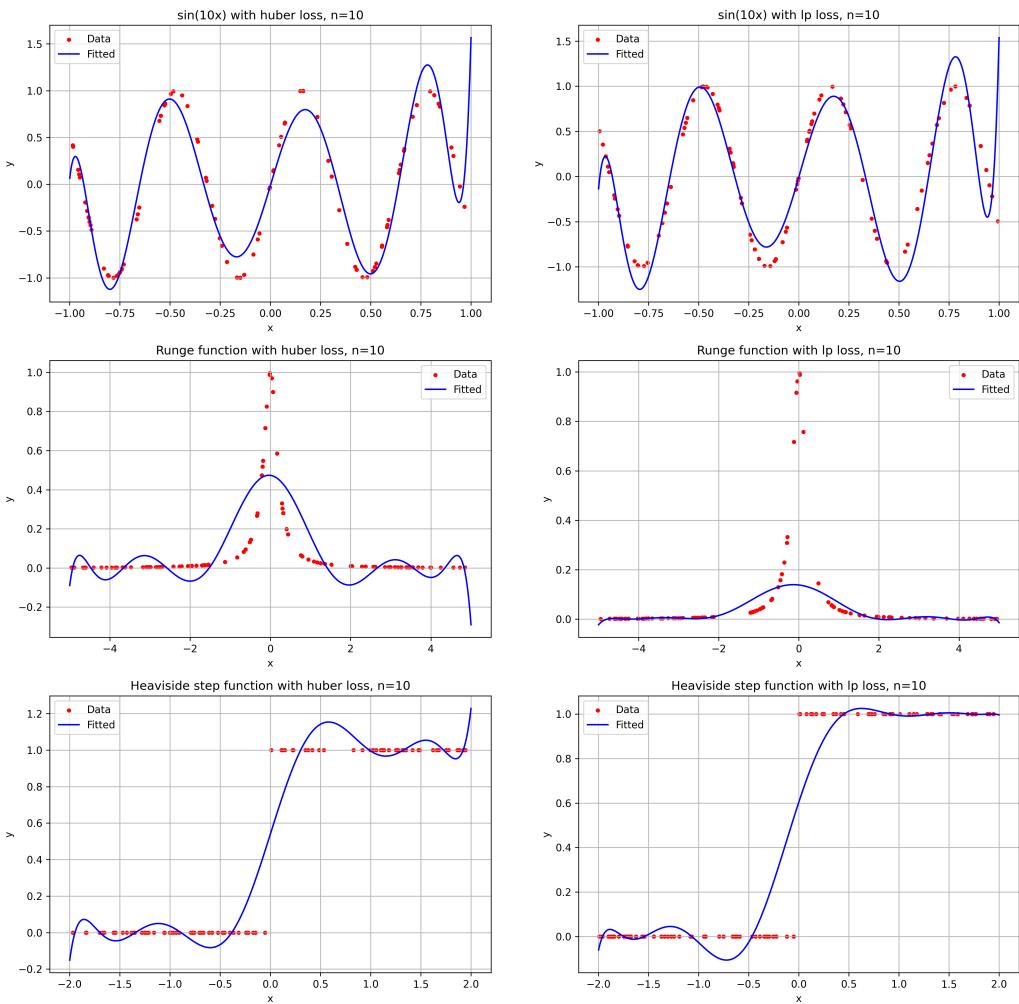


Figure 5: Comparison of Huber loss vs l^p with three functions ($n = 10$, with randomly sampled training points).

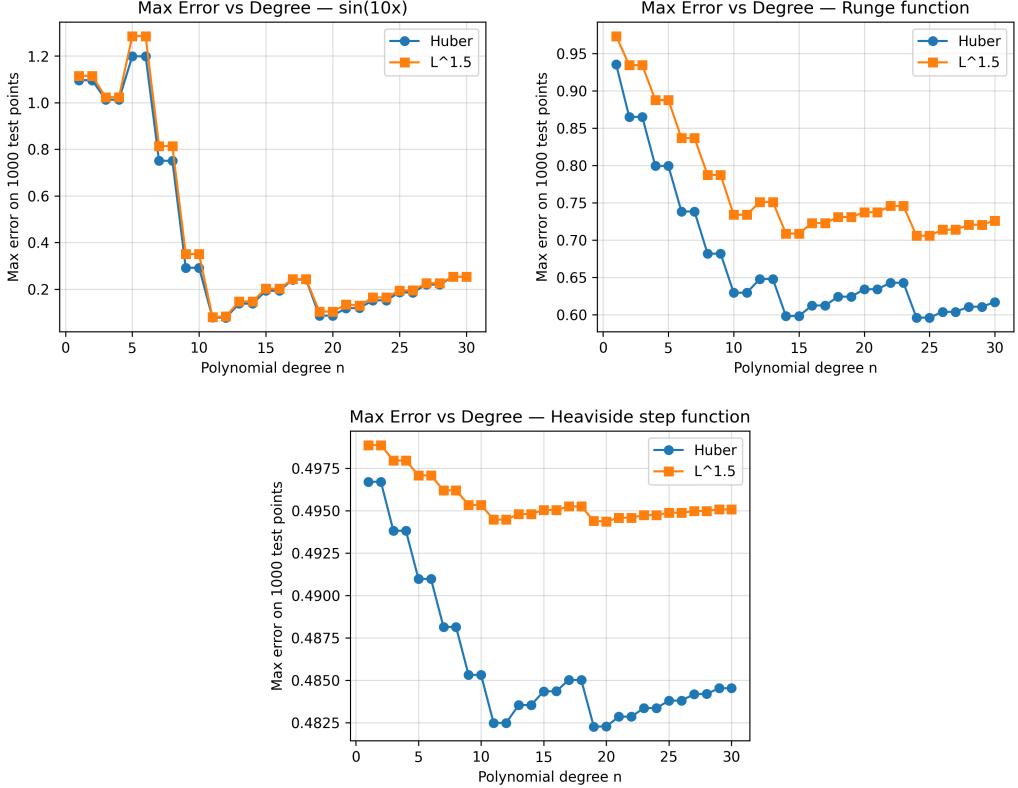


Figure 6: Error vs degree for different functions.

Table 1: Maximum error $\max_j |\Gamma(x_j) - M(x_j; \theta)|$ for selected polynomial degrees n using Huber and $L^{1.5}$ losses.

Dataset	Loss	$n = 2$	$n = 5$	$n = 10$	$n = 20$	$n = 30$
$\sin(10x)$	Huber	1.100	1.198	0.366	0.014	0.021
	$L^{1.5}$	1.116	1.285	0.510	0.177	0.097
Runge	Huber	0.868	0.802	0.632	0.580	0.569
	$L^{1.5}$	0.938	0.891	0.737	0.689	0.675
Step	Huber	0.499	0.497	0.496	0.495	0.495
	$L^{1.5}$	0.499	0.498	0.497	0.495	0.495

modestly in Table 1, even as the degree n increases to 30. In contrast, the discontinuous step function remains difficult: errors plateau around 0.495, showing that polynomials cannot approximate jumps well. Between the two losses, Huber is slightly more stable on the step function, while $L^{1.5}$ shows

modest advantages on oscillatory data at high degrees. Overall, robust losses provide little benefit on smooth data, some improvement for oscillatory data, and more stability for discontinuities.