

A.I. Linear regression

Dataset

x	x
27	442
26	442
20	457
31	431
14	471
16	463
25	434
7	487
6	480
10	480
15	470
31	430
6	489
28	436
26	444
16	464
28	447

24	430
11	475
24	445

Z-score Standardization

We used z-score standardization to normalize our data. Z-score standardization puts all features to the same scale by setting their mean to zero and standard deviation to one. By doing so, it preserves the original shape of the data, as the transformation only shifts and rescales the values without changing their distribution and internal relationships. This standardization is useful because many models—such as k-NN, linear regression and logistic regression—assume that the input data are centered and scaled. Moreover, it is easy to interpret: a z-score of 1 just indicates that a value lies one standard deviation above the mean. Z-score is defined as:

$$z_i = \frac{x_i - \mu}{\sigma}, \quad (1)$$

where

$$\mu = \frac{1}{n} \sum_{i=1}^n x_i, \quad (2)$$

$$\sigma = \sqrt{\frac{1}{n-1} \sum_{i=1}^n (x_i - \mu)^2} \quad (3)$$

We comparted the data through pictures:

-0.4167070940175797	0.4049359749530624
0.9918802660418444	-0.4397157519428983
0.5223511460220364	-1.284367478838859
-1.0036184940423398	0.9514753276504487
0.5223511460220364	-0.5390865433424231

Least Square

For one demintional feature, we can simply use:

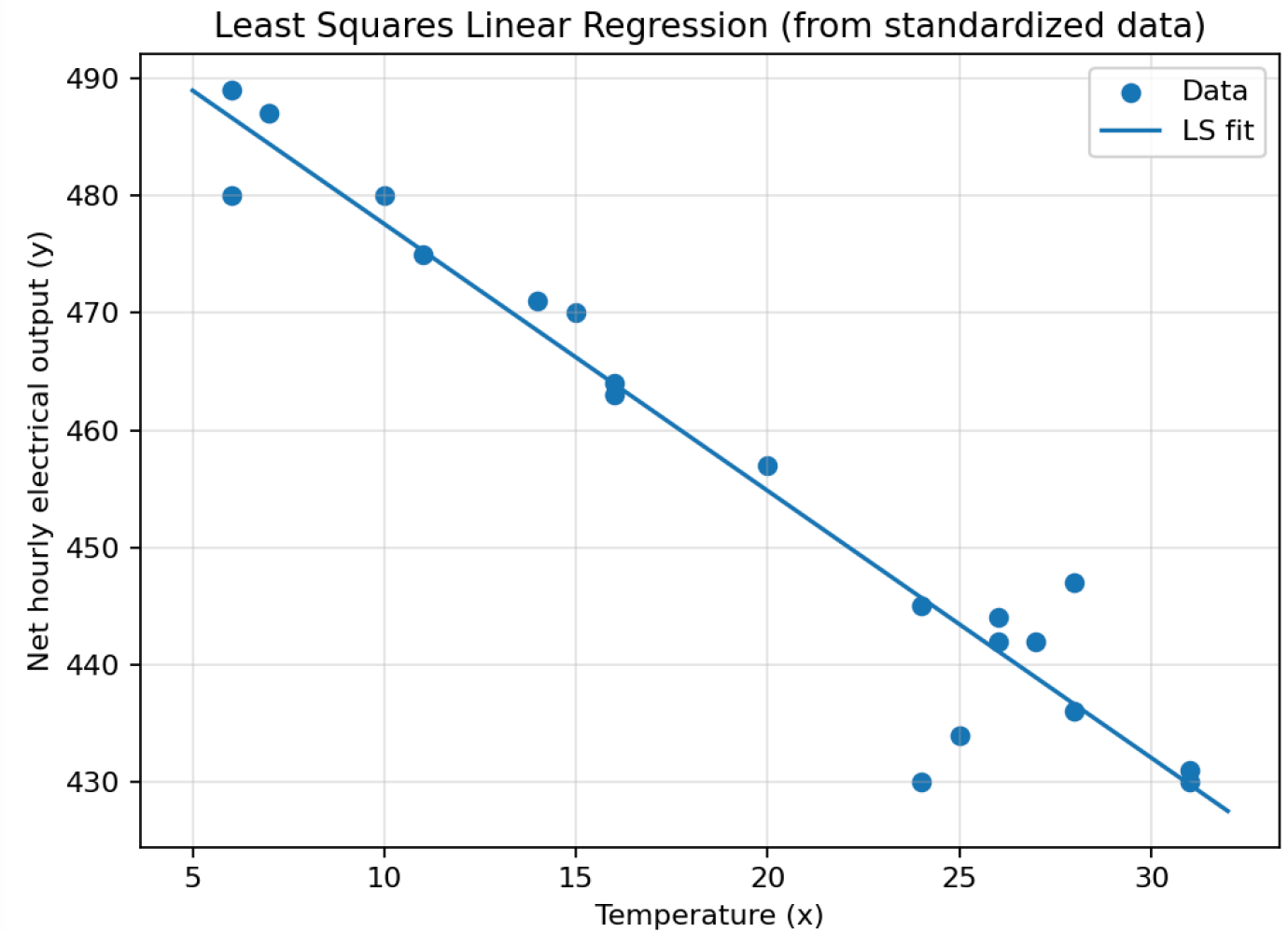
$$m = \frac{\sum_{i=1}^n (x_i - \bar{x})(y_i - \bar{y})}{\sum_{i=1}^n (x_i - \bar{x})^2} \quad (4)$$

to caculate its slope, and then using it to calculate the intercept:

$$b = \bar{y} - m\bar{x}. \quad (5)$$

Finially, we can get the equation of the line:

$$y = -2.275898x + 500.343812. \quad (6)$$



Linear regression with Gradient Descent

We define the hypothesis function as

$$h_{\theta}(x_i) = \theta_0 + \theta_1 x_i, \quad i = 1, \dots, m \quad (7)$$

where m is the numbers of training samples.

We choose the **Mean Squared Error** as the cost function:

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)^2 \quad (8)$$

Reason

MSE is smooth and differentiable, allowing gradient-based optimization and it measures the average squared difference between predicted and true values, and the square term penalizes large errors more strongly, helping the algorithm converge more stably.

Partial Derivative

Partial derivative with respect to θ_0 :

$$\begin{aligned}\frac{\partial J}{\partial \theta_0} &= \frac{1}{2m} \sum_{i=1}^m 2(h_{\theta}(x_i) - y_i) \frac{\partial h_{\theta}(x_i)}{\partial \theta_0} \\ &= \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)\end{aligned}\tag{9}$$

Partial derivative with respect to θ_1 :

$$\begin{aligned}\frac{\partial J}{\partial \theta_1} &= \frac{1}{2m} \sum_{i=1}^m 2(h_{\theta}(x_i) - y_i) \frac{\partial h_{\theta}(x_i)}{\partial \theta_1} \\ &= \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)x_i\end{aligned}\tag{10}$$

The gradient descent update rules are given by:

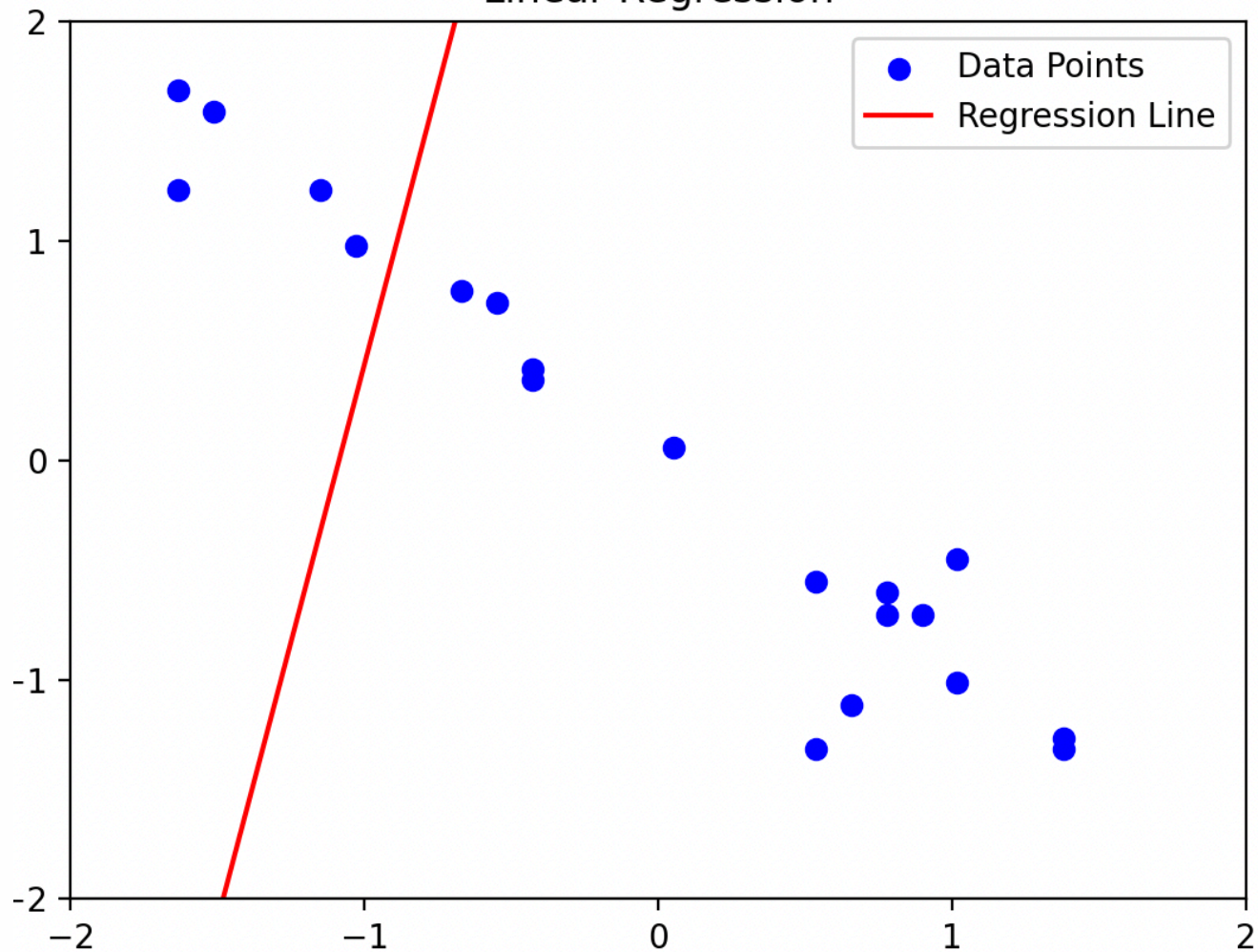
$$\theta_0 \leftarrow \theta_0 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)\tag{11}$$

$$\theta_1 \leftarrow \theta_1 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x_i) - y_i)x_i\tag{12}$$

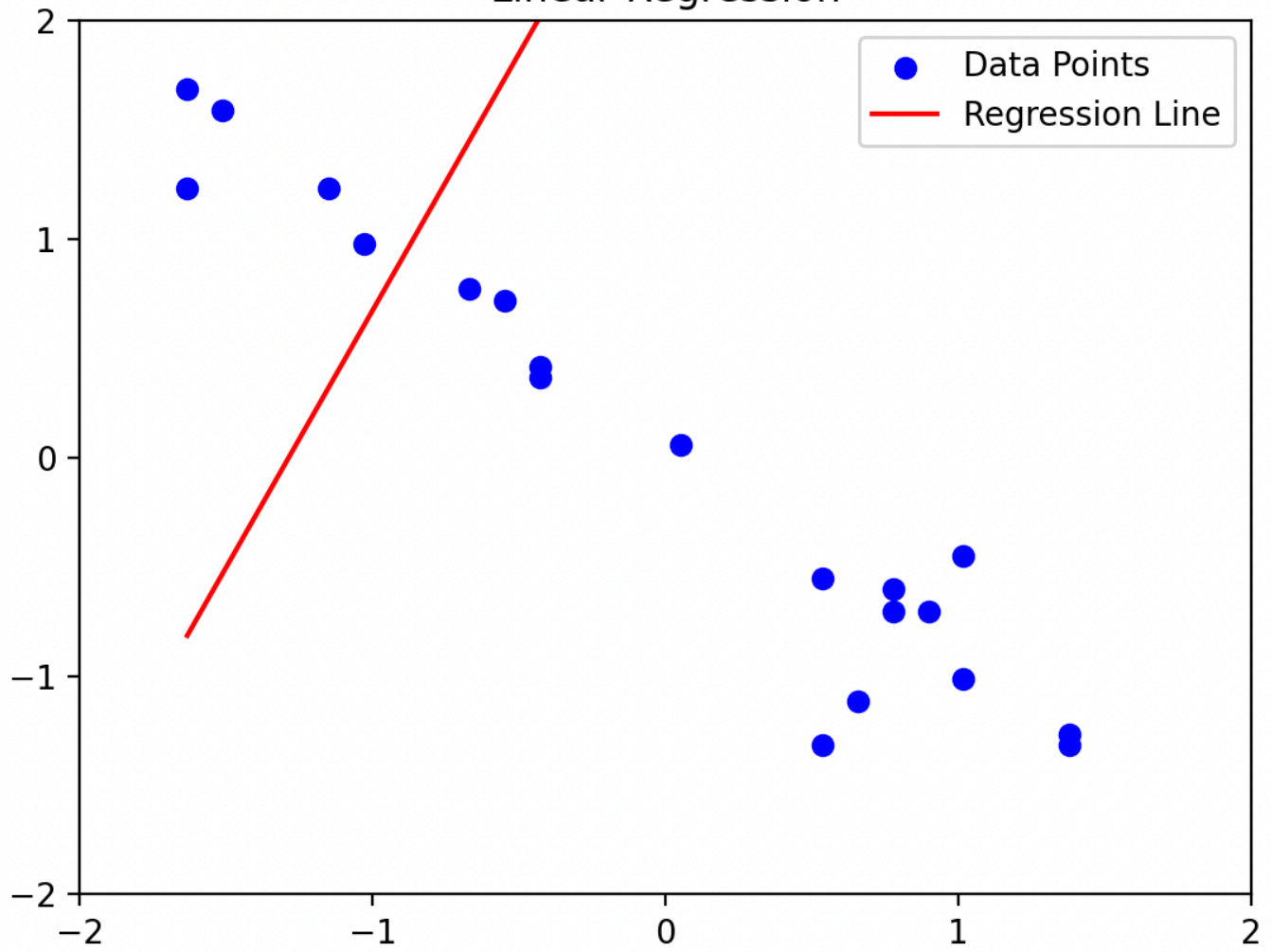
We started with $(\theta_0, \theta_1) = (10, 10)$ and a learning rate of 0.9; after the first, second, and third iterations, the parameters were updated to $(5.5, 5.07)$, $(3.03, 2.35)$, and $(1.66, 0.86)$ respectively, and finally converged to about $(0, -0.96)$.

Here are the regression lines after the first, second, and third iterations, as well as the final converged line:

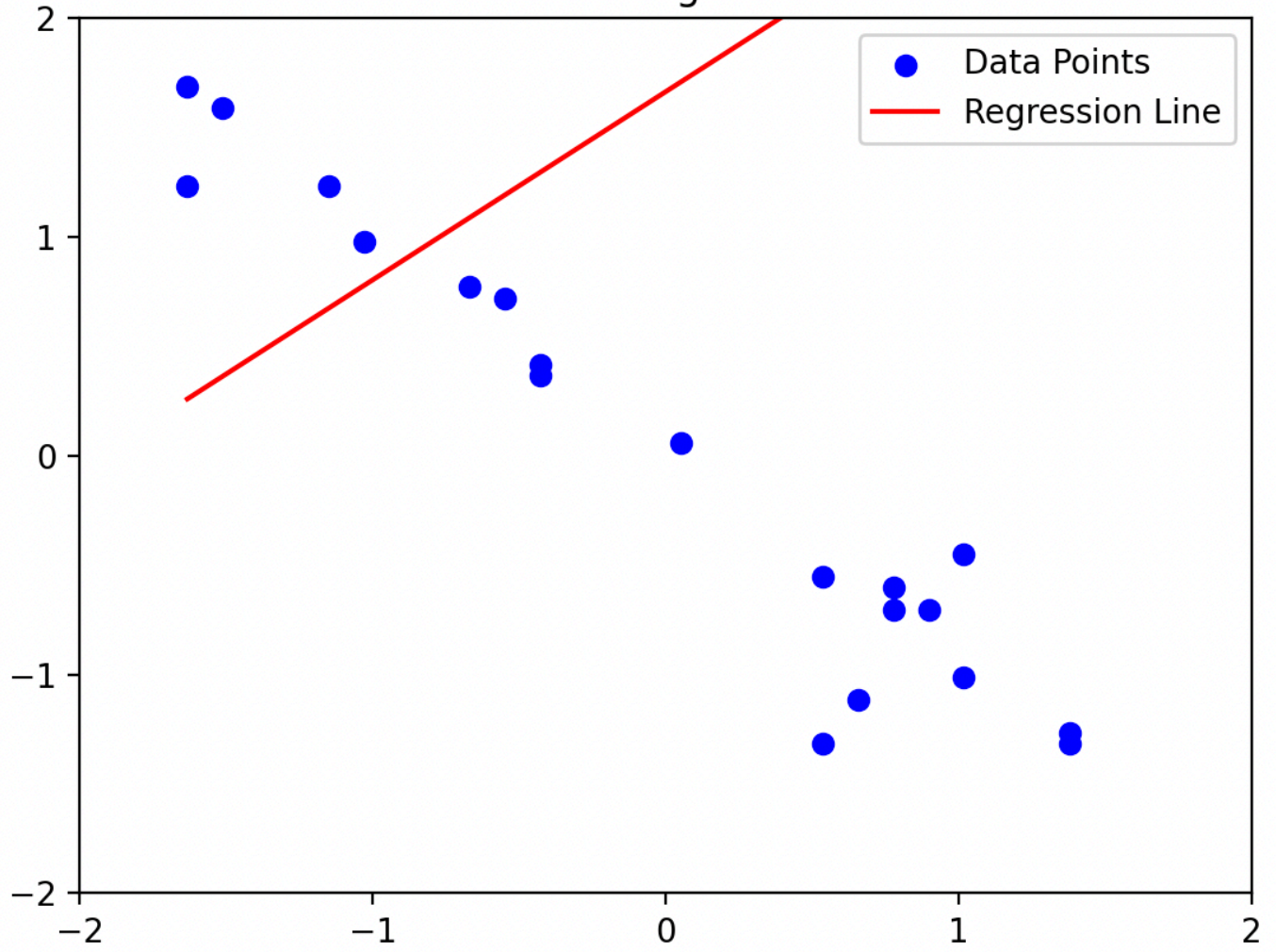
Linear Regression

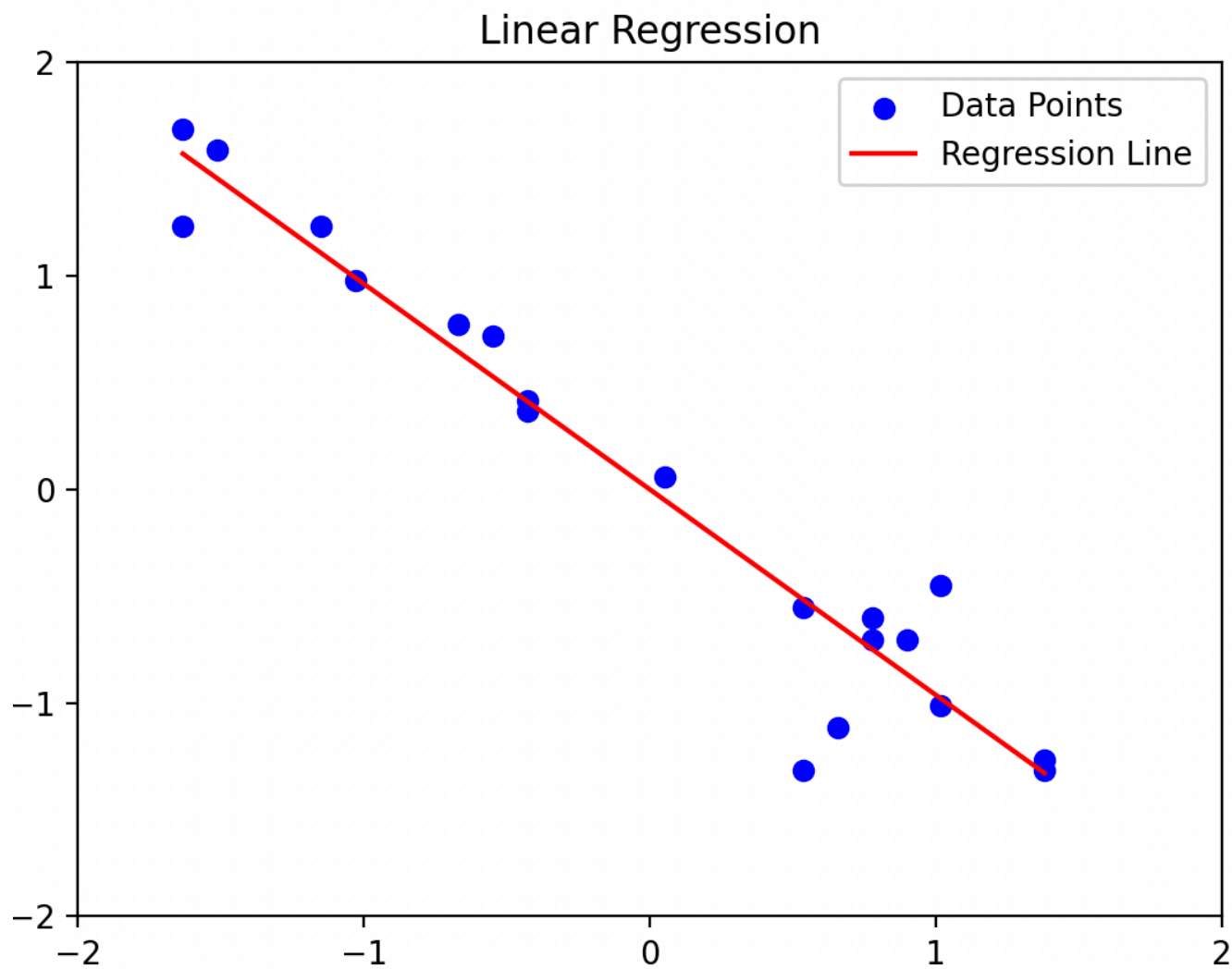


Linear Regression

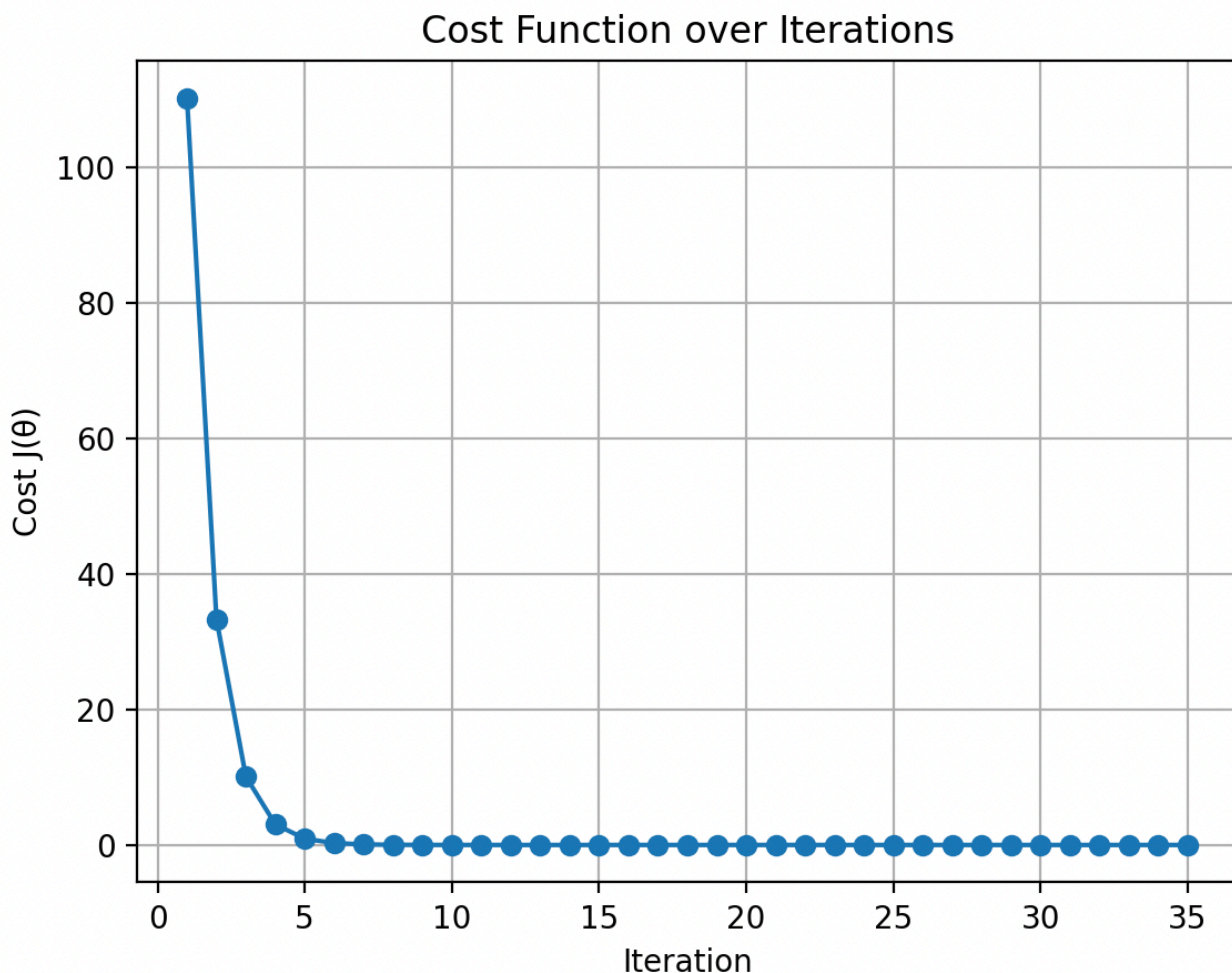


Linear Regression





Here is the picture for the cost:



Comparison

Similarities:

Both methods aim to find the optimal parameters for a linear regression model.

Differences:

- Least Squares: Computes the analytical (closed-form) solution in one step without iteration. However, when the feature dimension is large, calculating the inverse matrix can be very slow. It generally has low bias but high variance.
- Gradient Descent: An iterative method that gradually approaches the optimal solution. It is more suitable for large-scale datasets or high-dimensional features. The results are affected by the choice of learning rate and initial values. It typically has low variance but high bias.

Performance:

Least Squares performs well on small to medium-sized datasets with relatively few features. It provide an exact solution efficiently. By contrast, Gradient Descent performs better on large-scale or high-dimensional datasets, in which computing the matrix inverse in Least Squares may be computationally expensive. However, Gradient Descent may require more time to converge and its performance depends on proper hyperparameter tuning.

A.II. Poly regression

(1) Partial Derivatives of the Cost Function

Given the cost function:

$$J(\theta) = \frac{1}{4n} \sum_{i=1}^n \left(y^{(i)} - h_{\theta}(x^{(i)}) \right)^4 \quad (13)$$

where the hypothesis function is:

$$h_{\theta}(x^{(i)}) = \theta_2(x^{(i)})^2 + \theta_1 x^{(i)} + \theta_0 \quad (14)$$

We define the error term as:

$$\varepsilon^{(i)} = y^{(i)} - h_{\theta}(x^{(i)}) \quad (15)$$

Thus, the cost function can be rewritten as:

$$J(\theta) = \frac{1}{4n} \sum_{i=1}^n (\varepsilon^{(i)})^4 \quad (16)$$

Derivatives with respect to each parameter:

With respect to θ_2 :

$$\frac{\partial J}{\partial \theta_2} = -\frac{1}{n} \sum_{i=1}^n (\varepsilon^{(i)})^3 (x^{(i)})^2 \quad (17)$$

With respect to θ_1 :

$$\frac{\partial J}{\partial \theta_1} = -\frac{1}{n} \sum_{i=1}^n (\varepsilon^{(i)})^3 x^{(i)} \quad (18)$$

With respect to θ_0 :

$$\frac{\partial J}{\partial \theta_0} = -\frac{1}{n} \sum_{i=1}^n (\varepsilon^{(i)})^3 \quad (19)$$

(2) Parameter Update Rules

We use gradient descent to update the parameters.

The general update rule is:

$$\theta_j := \theta_j - \alpha \frac{\partial J}{\partial \theta_j} \quad (20)$$

where:

- α is the learning rate,
- $\frac{\partial J}{\partial \theta_j}$ is the gradient for parameter θ_j .

Specifically, in our problem, the update rules for each parameter are:

$$\theta_2 := \theta_2 - \alpha \left(-\frac{1}{n} \sum_{i=1}^n (\varepsilon^{(i)})^3 (x^{(i)})^2 \right) \quad (21)$$

$$\theta_1 := \theta_1 - \alpha \left(-\frac{1}{n} \sum_{i=1}^n (\varepsilon^{(i)})^3 x^{(i)} \right) \quad (22)$$

$$\theta_0 := \theta_0 - \alpha \left(-\frac{1}{n} \sum_{i=1}^n (\varepsilon^{(i)})^3 \right) \quad (23)$$

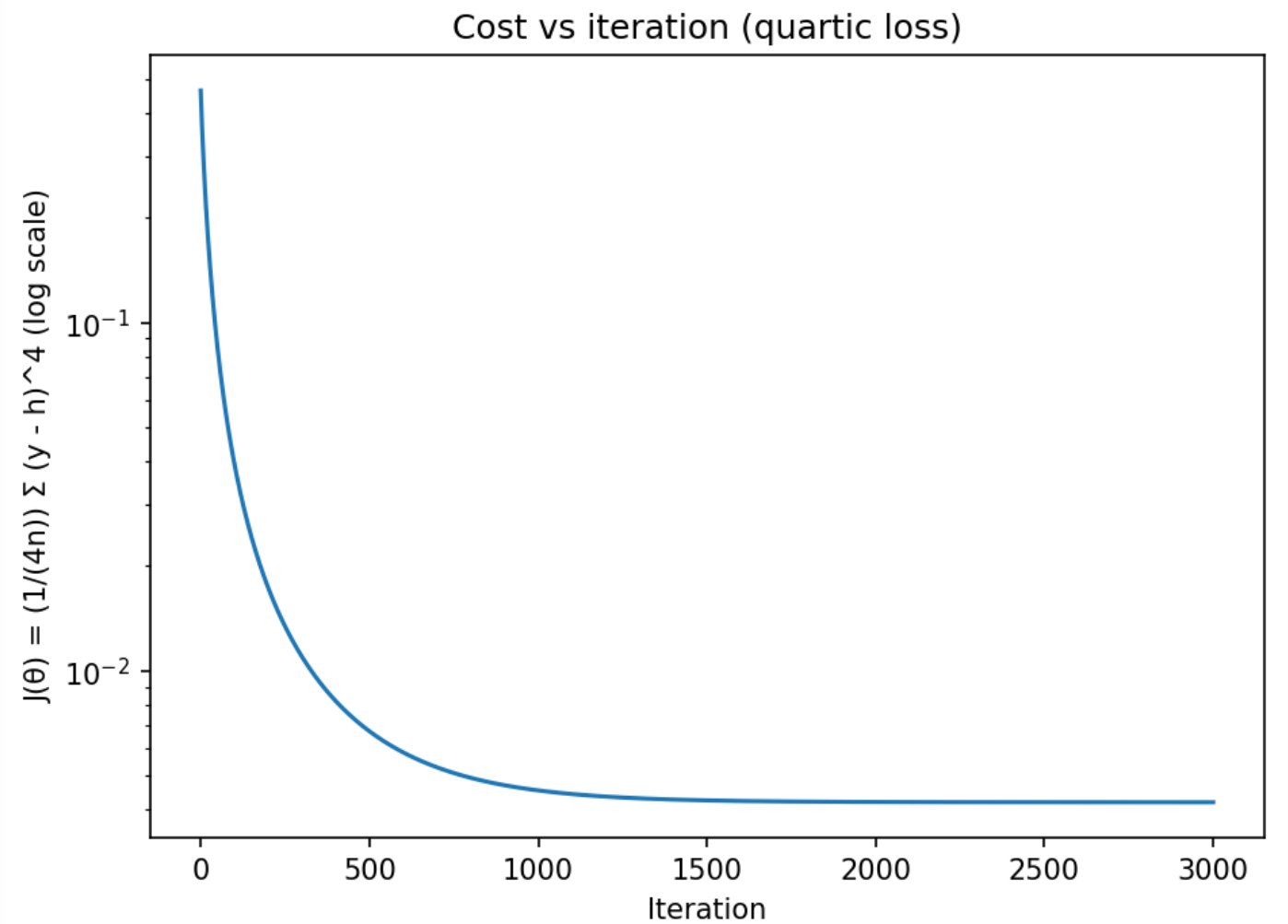
(3) Iterations

Parameter evolution and Cost during training (learning rate $\alpha = 0.01$):

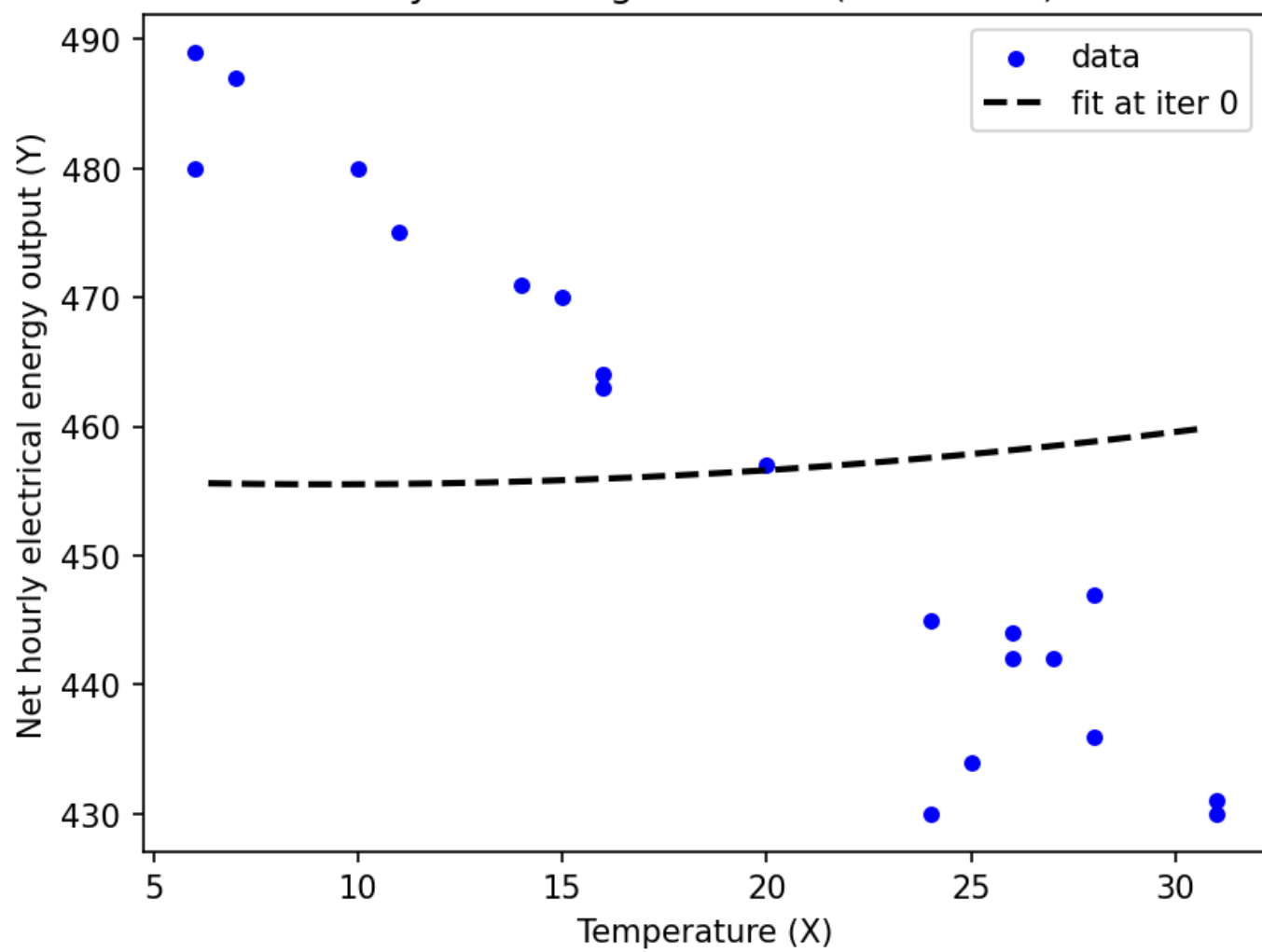
Iteration	θ_0	θ_1	θ_2	Cost (J)
0	0.03455842	0.08216181	0.03304371	0.4643122
100	-0.06609397	-0.4526044	0.09791056	0.03924249
200	-0.1092821	-0.5906548	0.1053334	0.01732243
500	-0.1683796	-0.7493620	0.1199282	0.006708472
1000	-0.2047500	-0.8408934	0.1345493	0.004537533

2000	-0.2221575	-0.8875152	0.1457154	0.004204270
3000	-0.2249141	-0.8943075	0.1478593	0.004196765

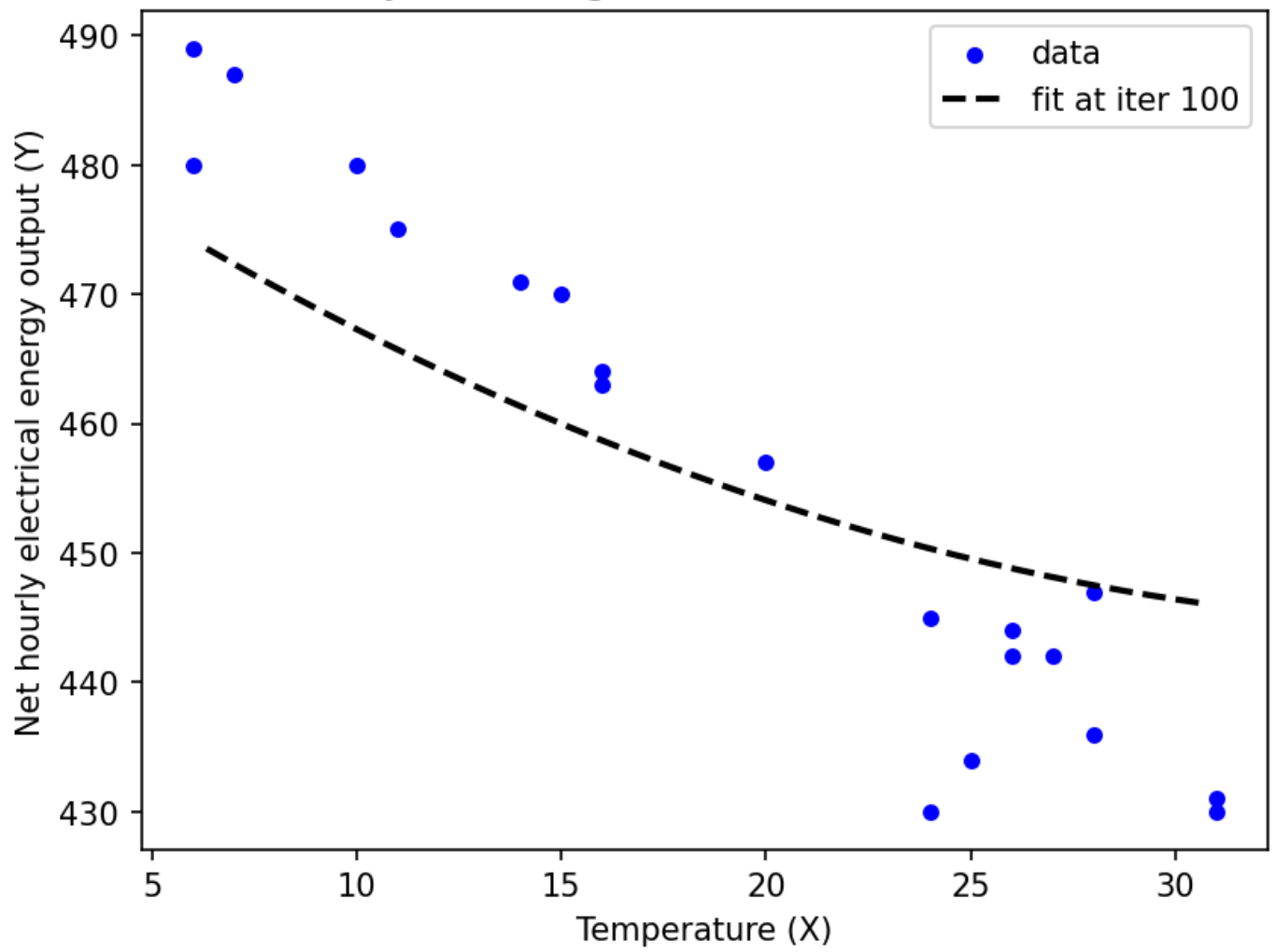
The first image illustrates the decrease of the cost function over iterations, while the others display the fitting results at different stages. It is obvious that the model becomes stable and shows good convergence after about 1,000 iterations.



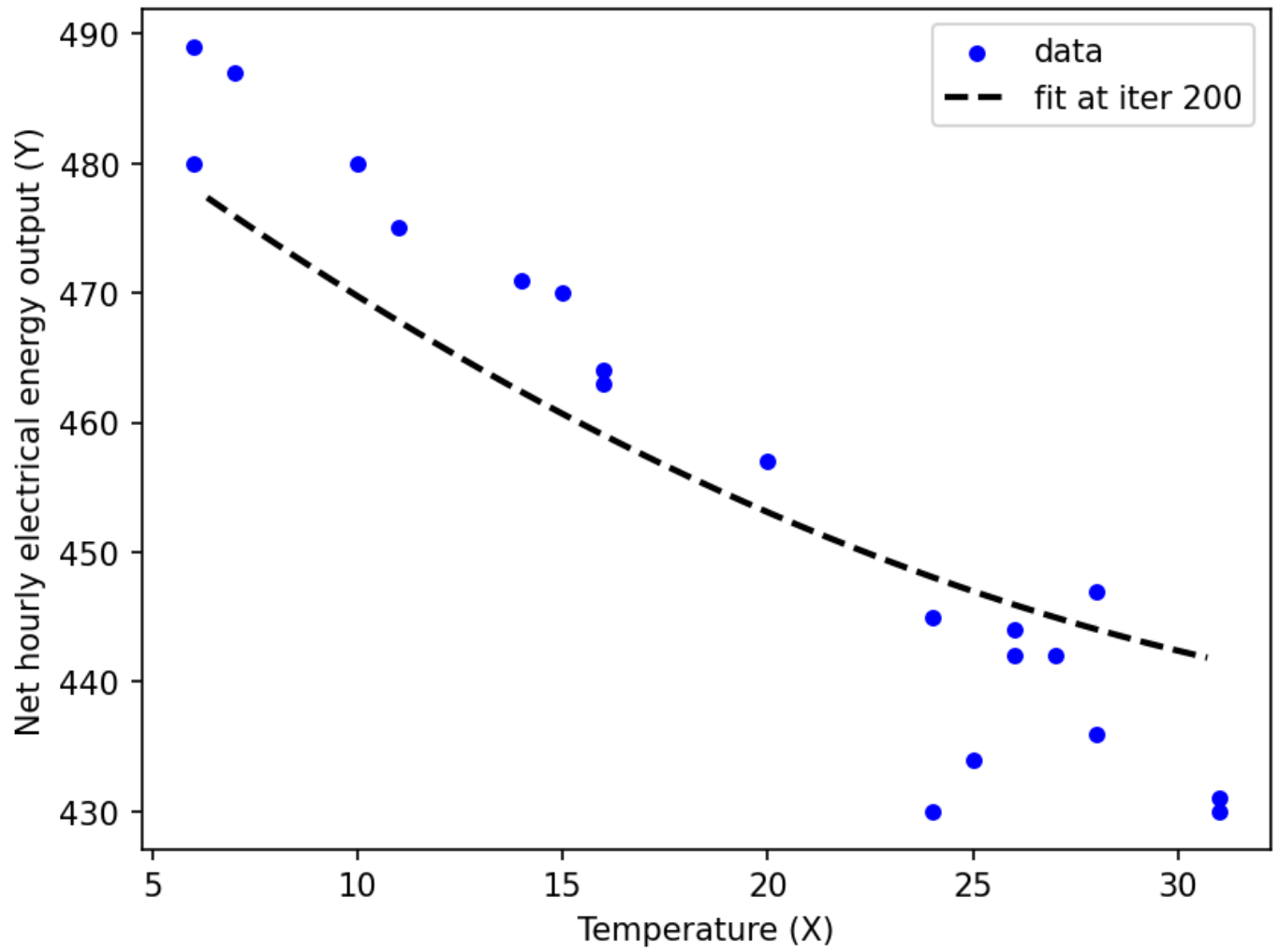
Polynomial regression fit (iteration 0)



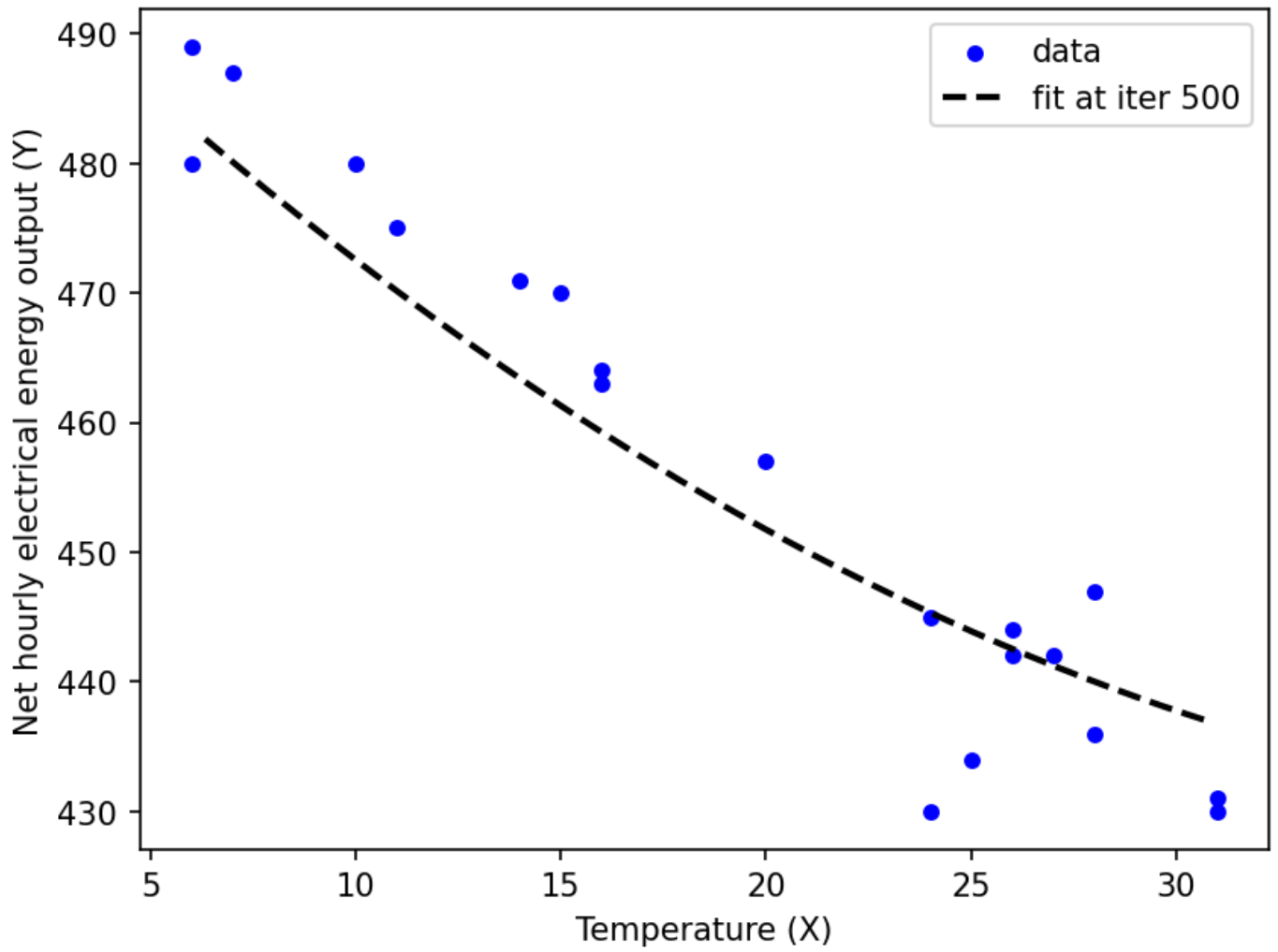
Polynomial regression fit (iteration 100)



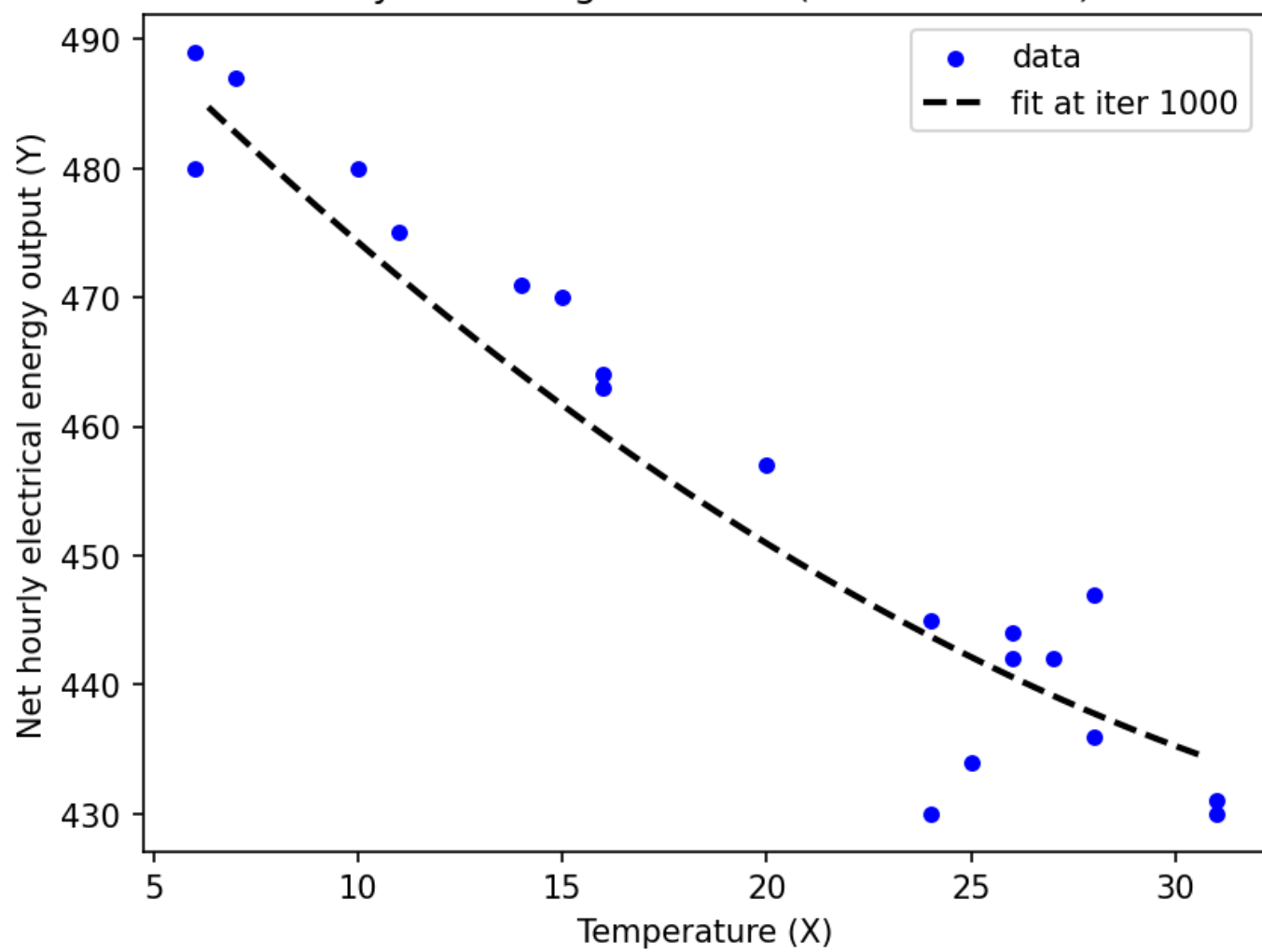
Polynomial regression fit (iteration 200)



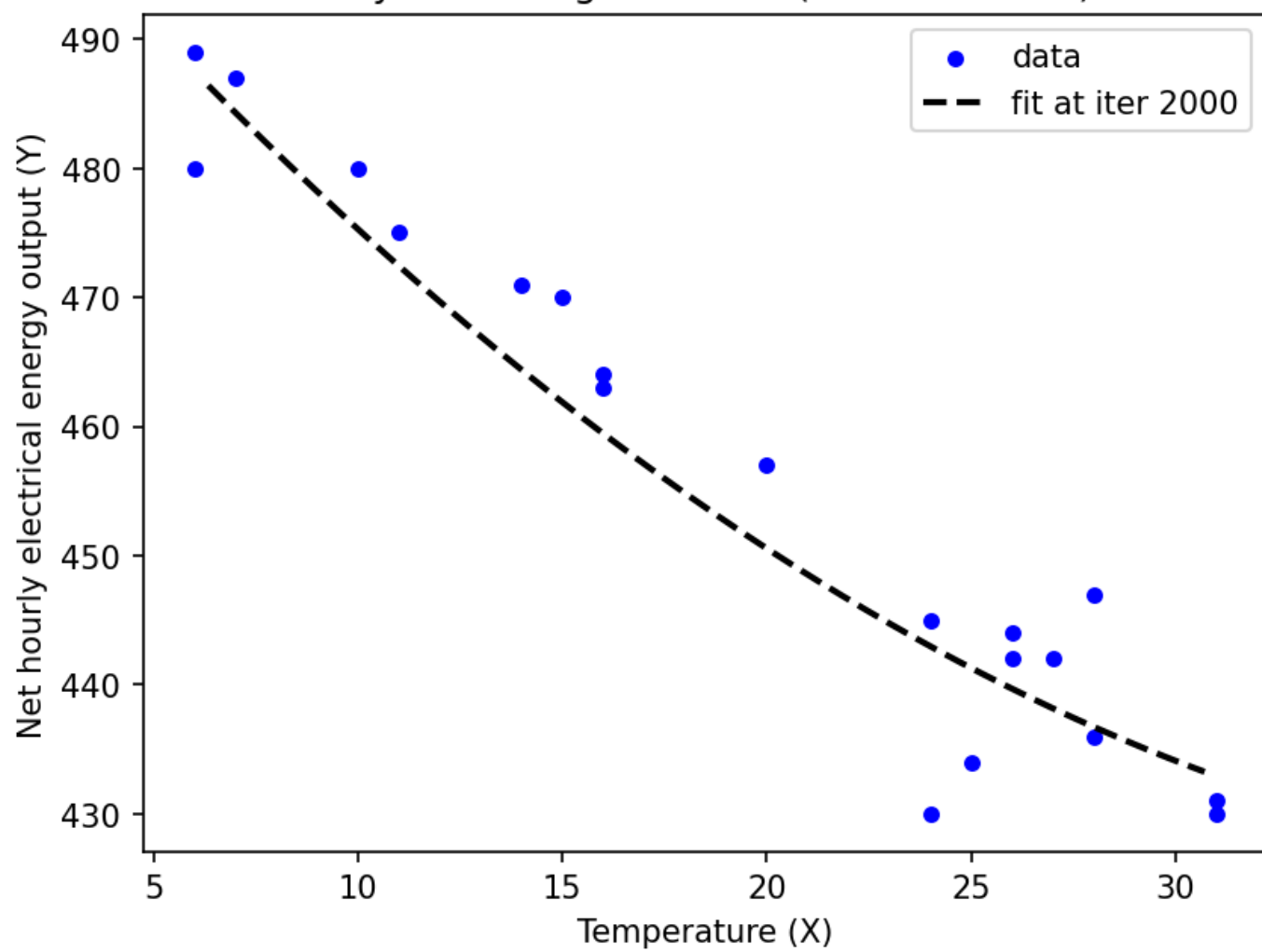
Polynomial regression fit (iteration 500)

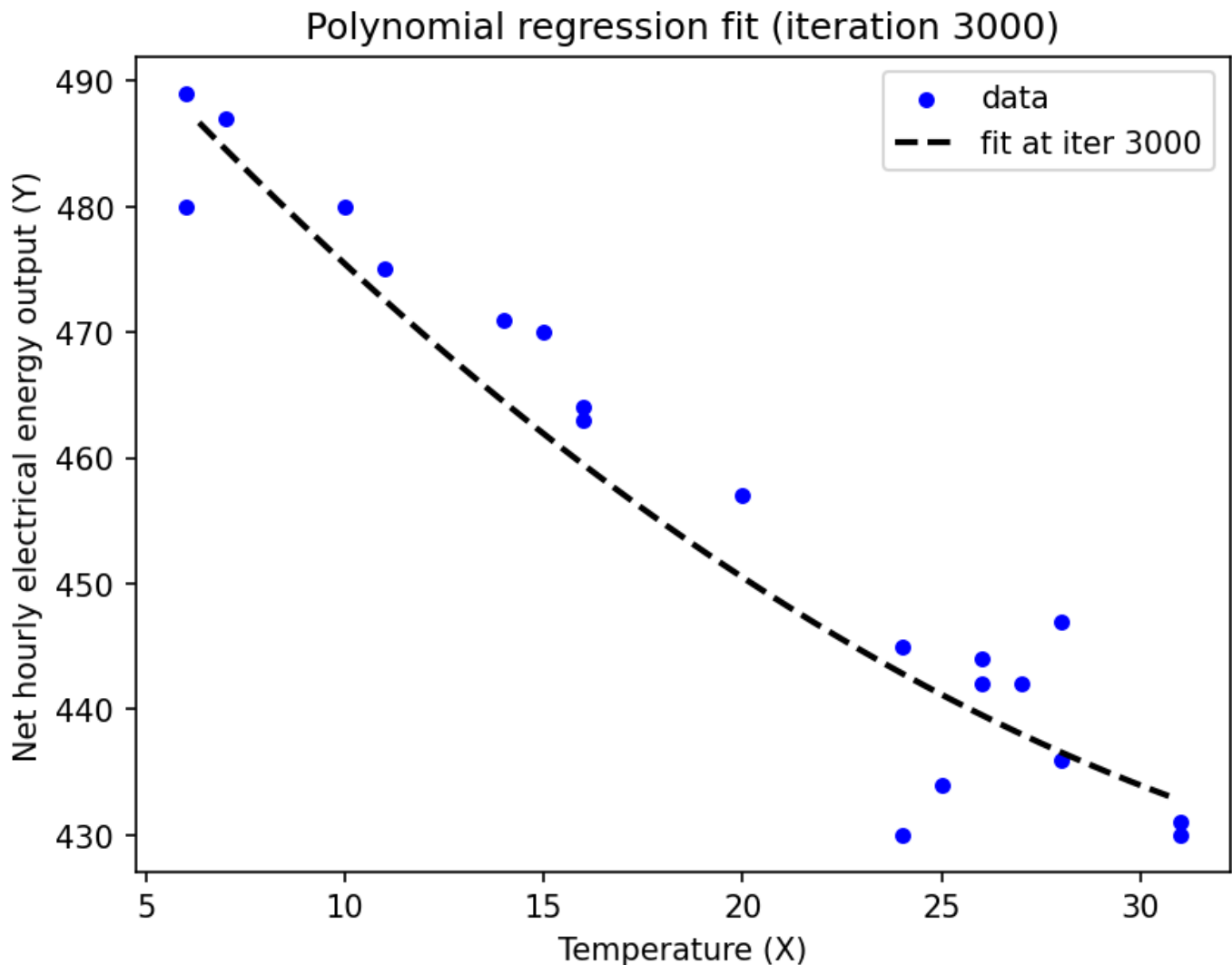


Polynomial regression fit (iteration 1000)



Polynomial regression fit (iteration 2000)





Similarities

The two are very consistent in their core principles and objectives: their goal is the same: regression algorithms within supervised learning. They aim to find a model (a line or a curve) that minimizes error on the training data. This error is typically measured by Mean Squared Error or the Sum of Squared Errors .

They have the same mathematical principle :The term “linear” in “linear regression” refers to linearity with respect to the model parameters (coefficients β), not necessarily linearity with respect to the features.

Univariate linear regression model:

$$y = \beta_0 + \beta_1 x \quad (24)$$

This model is linear in the parameters β_0 and β_1 .

For Second-order polynomial regression model:

$$y = \beta_0 + \beta_1x + \beta_2x^2 \tag{25}$$

The model remains linear in the parameters $\beta_0, \beta_1, \beta_2$.

In addition, they have the same approach to get $h(\vec{x})$: Least Square and Gradient Descent.

Differences

Views	Linear Regression	Polynomial Regression
Assumed Relationship	Assumes a linear relationship between features (X) and target (Y).	Assumes a non-linear relationship between features (X) and target (Y).
Model Form	$y = \beta_0 + \beta_1x + \dots \beta_nx_n$	$y = \beta_0 + \beta_1x + \beta_2x^2 + \dots + \beta_nx^n$
Graphical Representation	A straight line.	A curve.
Model Complexity	Low. The model is simple.	High (depends on the degree n).
Bias	High Bias. It will underfit if the relationship is non-linear.	Low Bias. Can fit the training data better.
Variance	Low Variance. Not sensitive to small changes in training data.	High Variance. Prone to overfitting, especially with a high degree n .