

# Linear Regression From Scratch: Mathematics and Implementation

(Your Name)

October 3, 2025

## 1 Feature Matrix (Stick to Notes)

We define the feature matrix  $X \in \mathbb{R}^{m \times n}$  as

$$X = \begin{bmatrix} x_{11} & x_{12} & \dots & x_{1n} \\ x_{21} & x_{22} & \dots & x_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ x_{m1} & x_{m2} & \dots & x_{mn} \end{bmatrix}, \quad \vec{x}^{(i)} = (x_{i1}, x_{i2}, \dots, x_{in}).$$

Each row of  $X$  is a feature vector  $\vec{x}^{(i)}$ .

## 2 Prediction Rule

For a sample  $\vec{x}^{(i)}$ , the prediction is

$$\hat{y}^{(i)} = f(\vec{x}^{(i)}) = \vec{w} \cdot \vec{x}^{(i)} + b,$$

so vectorized:

$$\hat{\vec{y}} = X\vec{w} + b\mathbf{1}.$$

## 3 Error Vector

For each sample,

$$e^{(i)} = \hat{y}^{(i)} - y^{(i)},$$

and the stacked error vector is

$$\vec{e} = \begin{bmatrix} \hat{y}^{(1)} - y^{(1)} \\ \hat{y}^{(2)} - y^{(2)} \\ \vdots \\ \hat{y}^{(m)} - y^{(m)} \end{bmatrix} \in \mathbb{R}^m.$$

## 4 Gradient of the Multivariable Cost

We use the mean squared error with the  $\frac{1}{2}$  convention:

$$J(\vec{w}, b) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 = \frac{1}{2m} \|\vec{e}\|_2^2.$$

Its gradients are

$$\nabla_{\vec{w}} J = \frac{1}{m} X^\top \vec{e}, \quad \frac{\partial J}{\partial b} = \frac{1}{m} \sum_{i=1}^m e^{(i)}.$$

## 5 Parameter Updates (Stick to Notes)

With learning rate  $\alpha$ :

$$\vec{w} \leftarrow \vec{w} - \alpha \nabla_{\vec{w}} J = \vec{w} - \alpha \cdot \frac{1}{m} X^\top \vec{e}, \quad b \leftarrow b - \alpha \frac{\partial J}{\partial b} = b - \frac{\alpha}{m} \sum_{i=1}^m e^{(i)}.$$

## 6 Normalization (z-Score)

Before regression we normalize each feature using training statistics:

$$x'_j = \frac{x_j^{(i)} - \mu_j}{\sigma_j}, \quad \mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}, \quad \sigma_j = \sqrt{\frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2},$$

and if  $\sigma_j = 0$  we set  $\sigma_j \leftarrow 1$ .

## 7 Dataset Split (Stick to Notes)

We split the dataset as:

Train: 80%, Validation: 10%, Test: 10%.

## 8 Implementation Snippets (Exact Code)

Loading Data and Building  $X, \vec{y}$

```
1 import numpy as np
2 import pandas as pd
3
4 df = pd.read_csv("housing.csv", delim_whitespace=True, header=
    None)
5 X = df.iloc[:, :-1].to_numpy(dtype=float)
6 y = df.iloc[:, -1].to_numpy(dtype=float)
```

## z-Score Normalization (with zero-std guard)

```
1 m, n = X.shape
2 mu = X.mean(axis=0)
3 sigma = X.std(axis=0, ddof=0)
4 sigma[sigma == 0] = 1.0
5 X = (X - mu) / sigma
```

## Parameters and Prediction Function

```
1 w = np.zeros(n, dtype=float)
2 b = 0.0
3
4 def f(X, w, b):
5     # Vectorized: \hat{y} = X w + b 1
6     return X @ w + b
```

## Gradient Step (Implements $X^T e$ and bias average)

```
1 def step(X, y, w, b, alpha):
2     m = len(y)
3     e = f(X, w, b) - y           # e = \hat{y} - y
4     grad_w = (X.T @ e) / m       # \nabla_w J = (1/m) X^T e
5     grad_b = e.mean()            # dJ/db = (1/m) sum e_i
6     return w - alpha * grad_w, b - alpha * grad_b
```

## Cost Function $J(\vec{w}, b) = \frac{1}{2m} \|e\|^2$

```
1 def cost(X, y, w, b):
2     e = f(X, w, b) - y
3     return 0.5 * np.mean(e ** 2)
```

## Batch Gradient Descent with Tolerance

```
1 alpha = 0.05
2 max_epochs = 5000
3 tol = 1e-8
4 prev_J = None
5
6 for epoch in range(1, max_epochs + 1):
7     w, b = step(X, y, w, b, alpha)
8     J = cost(X, y, w, b)
9     if prev_J is not None and abs(prev_J - J) < tol:
10         break
11     prev_J = J
```

## Predictions and Metrics

```
1 y_pred = f(X, w, b)
2 mse = np.mean((y - y_pred) ** 2)
3
4 print(f"final_cost J(w,b): {J:.6f}")
5 print(f"MSE: {mse:.6f}")
6 print("bias b:", round(b, 6))
7 print("weights w:", np.round(w, 6))
```