

Machine Learning Notes

Amit Thakur

January 3, 2025

Contents

Table of Contents	2
1 Linear Regression	3
1.1 Linear Regression with Single Feature	3
1.1.1 Linear Regression Model	3
1.1.2 Cost Function	3
1.1.3 Minimizing Cost	4
1.1.4 Gradient Descent Algorithm to find optimal θ_0 and θ_1	5
1.2 Linear Regression with Multiple Features	6
1.2.1 General Form	6
1.2.2 Cost Function	7
1.2.3 Gradient of the Cost Function	8
1.2.4 Analytical Solution	8
1.2.5 Gradient Descent Algorithm	8
1.3 Linear Models for Regression	9
1.3.1 Polynomial Curve Fitting	9
1.3.2 Analytical Solution	9
1.3.3 Linear Models with Basis Functions	9
1.4 Probabilistic Interpretation of Linear Regression	11
1.4.1 Equivalence of least square error and maximum likelihood estimation	11

Chapter 1

Linear Regression

Linear Regression is a statistical method used for modeling the relationship between a dependent variable (target or output) and one or more independent variables (predictors or features).

1.1 Linear Regression with Single Feature

This regression deals with one independent variable (x).

1.1.1 Linear Regression Model

$$y = \theta_0^t + \theta_1^t x + \epsilon \quad (1.1)$$

$$\hat{y} = h_\theta(x) = \theta_0 + \theta_1 x \quad (1.2)$$

where:

- y : the dependent variable (target)
- x : the independent variable or input feature used to predict y .
- θ_1^t : the slope or coefficient of x at iteration t
- θ_0^t : intercept at iteration t
- ϵ : the error term or residual. It captures the noise or other unmodeled effects.
- \hat{y} : the output of the linear regression model (\cdot) for a given input (x)
- $h_\theta(x)$: the hypothesis function for linear regression.

1.1.2 Cost Function

The cost function for linear regression measures the average squared error between predicted and actual values. It is defined as:

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (\hat{y}^{(i)} - y^{(i)})^2 = \frac{1}{2m} \sum_{i=1}^m ((\theta_0 + \theta_1 x^{(i)}) - y^{(i)})^2 \quad (1.3)$$

where:

- m : Number of input and output datapoints.
- $\frac{1}{2}$: A factor kept for convenience, as it simplifies the derivative calculations during gradient descent.

1.1.3 Minimizing Cost

The error for each data point can be written as:

$$e^{(i)} = y^{(i)} - (\theta_0 + \theta_1 x^{(i)}) \quad (1.4)$$

The cost function becomes:

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m e^{(i)2} \quad (1.5)$$

We need to find θ_0 and θ_1 that minimize J . This requires setting the partial derivatives of J with respect to θ_0 and θ_1 to zero as the cost function is quadratic and there's just one critical point.

Solving for θ_0 :

$$\frac{\partial J}{\partial \theta_0} = -\frac{1}{m} \sum_{i=1}^m (y^{(i)} - (\theta_0 + \theta_1 x^{(i)})) = 0$$

$$\begin{aligned} \sum_{i=1}^m (y^{(i)} - (\theta_0 + \theta_1 x^{(i)})) &= 0 \\ \sum_{i=1}^m y^{(i)} &= m\theta_0 + \theta_1 \sum_{i=1}^m x^{(i)} \end{aligned}$$

Divide through by m :

$$\bar{y} = \theta_0 + \theta_1 \bar{x}$$

where \bar{y} and \bar{x} are the means of $y^{(i)}$ and $x^{(i)}$, respectively.

Rearranging:

$$\theta_0 = \bar{y} - \theta_1 \bar{x} \quad (1.6)$$

Solving for θ_1 :

$$\frac{\partial J}{\partial \theta_1} = -\frac{1}{m} \sum_{i=1}^m x^{(i)} (y^{(i)} - (\theta_0 + \theta_1 x^{(i)})) = 0$$

Expanding:

$$\sum_{i=1}^m x^{(i)} y^{(i)} = \theta_0 \sum_{i=1}^m x^{(i)} + \theta_1 \sum_{i=1}^m (x^{(i)})^2$$

Substitute $\theta_0 = \bar{y} - \theta_1 \bar{x}$:

$$\sum_{i=1}^m x^{(i)} y^{(i)} = (\bar{y} - \theta_1 \bar{x}) \sum_{i=1}^m x^{(i)} + \theta_1 \sum_{i=1}^m (x^{(i)})^2$$

Simplify:

$$\sum_{i=1}^m x^{(i)} y^{(i)} = \bar{y} \sum_{i=1}^m x^{(i)} - \theta_1 \bar{x} \sum_{i=1}^m x^{(i)} + \theta_1 \sum_{i=1}^m (x^{(i)})^2$$

Reorganize terms:

$$\theta_1 \left(\sum_{i=1}^m (x^{(i)})^2 - \frac{1}{m} \left(\sum_{i=1}^m x^{(i)} \right)^2 \right) = \sum_{i=1}^m x^{(i)} y^{(i)} - \frac{1}{m} \sum_{i=1}^m x^{(i)} \sum_{i=1}^m y^{(i)}$$

Using simplified notation:

- $\bar{x} = \frac{1}{m} \sum_{i=1}^m x^{(i)}$ (mean of x)
- $\bar{y} = \frac{1}{m} \sum_{i=1}^m y^{(i)}$ (mean of y)

$$\theta_1 = \frac{\sum_{i=1}^m (x^{(i)} - \bar{x}) (y^{(i)} - \bar{y})}{\sum_{i=1}^m (x^{(i)} - \bar{x})^2} \quad (1.7)$$

1.1.4 Gradient Descent Algorithm to find optimal θ_0 and θ_1

1. Gradient descent moves the parameters in the direction of the negative gradient (steepest descent) of the cost function.
2. The learning rate α determines the size of each step. If α is too large, the algorithm may overshoot the minimum. If it is too small, convergence may take too long.
3. The iterative process ensures gradual improvement in the model parameters until the cost function is minimized.

Algorithm 1 Gradient Descent for Linear Regression

-
- 1: **Input:** Learning rate α , initial values for θ_0 and θ_1 , and maximum iterations or convergence threshold ϵ .
 - 2: **Output:** Optimized parameters θ_0 and θ_1 .
 - 3: Set initial values for θ_0 and θ_1 .
 - 4: **repeat**
 - 5: Compute updates for parameters:

$$\begin{aligned}\text{temp0} &:= \theta_0 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) \\ \text{temp1} &:= \theta_1 - \alpha \cdot \frac{1}{m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)}) x^{(i)}\end{aligned}$$

- 6: Update parameters simultaneously:

$$\begin{aligned}\theta_0 &:= \text{temp0} \\ \theta_1 &:= \text{temp1}\end{aligned}$$

- 7: Compute cost function:

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (\theta_0 + \theta_1 x^{(i)} - y^{(i)})^2$$

- 8: **until** maximum iterations reached **or** change in $J(\theta_0, \theta_1)$ is below threshold ϵ .
 - 9: **return** θ_0, θ_1
-

1.2 Linear Regression with Multiple Features

$$y = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \cdots + \theta_n x_n + e \quad (1.8)$$

where

- y : The dependent variable (or target/output variable) that the model predicts
- θ_0 : The intercept term, representing the value of y when all $x_i = 0$
- $\theta_1, \theta_2, \dots, \theta_n$: The coefficients or weights for the independent variables x_1, x_2, \dots, x_n .

1.2.1 General Form

$$y = \boldsymbol{\theta}^T \mathbf{x} + e \quad (1.9)$$

where $\boldsymbol{\theta} \in \mathbb{R}^{(n+1) \times 1}$ is the parameter vector (column vector of coefficients), defined as:

$$\boldsymbol{\theta} = \begin{bmatrix} \theta_0 \\ \theta_1 \\ \vdots \\ \theta_n \end{bmatrix} \quad (1.10)$$

And the input feature vector $\mathbf{x} \in \mathbb{R}^{(n+1) \times 1}$ defined as:

$$\mathbf{x} = \begin{bmatrix} 1 \\ x_1 \\ \vdots \\ x_n \end{bmatrix} \quad (1.11)$$

The prediction function is:

$$\hat{y} = h_{\theta}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x} \quad (1.12)$$

1.2.2 Cost Function

The residuals (E) represent the difference between the actual target values (y) and the predicted values ($X \cdot \boldsymbol{\theta}$). For m data points:

$$E = \begin{bmatrix} e^{(1)} \\ e^{(2)} \\ \vdots \\ e^{(m)} \end{bmatrix} = X \cdot \boldsymbol{\theta} - y \quad (1.13)$$

The cost function $J(\boldsymbol{\theta})$ is:

$$J(\boldsymbol{\theta}) = \frac{1}{2m} \sum_{i=1}^m (e^{(i)})^2 \quad (1.14)$$

Using $E = X \cdot \boldsymbol{\theta} - y$, we can write:

$$J(\boldsymbol{\theta}) = \frac{1}{2m} \|E\|^2$$

Expanding this using the transpose:

$$J(\boldsymbol{\theta}) = \frac{1}{2m} E^T E$$

Substituting $E = X \cdot \boldsymbol{\theta} - y$:

$$J(\boldsymbol{\theta}) = \frac{1}{2m} (X \cdot \boldsymbol{\theta} - y)^T (X \cdot \boldsymbol{\theta} - y) \quad (1.15)$$

To minimize $J(\boldsymbol{\theta})$, we compute the gradient (partial derivative) with respect to each parameter θ_j :

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x_j^{(i)} \quad (1.16)$$

We now compute the partial derivative of $J(\boldsymbol{\theta})$ with respect to θ_j . The chain rule will be used.

$$J(\boldsymbol{\theta}) = \frac{1}{2m} [(X \cdot \boldsymbol{\theta})^T (X \cdot \boldsymbol{\theta}) - 2y^T (X \cdot \boldsymbol{\theta}) + y^T y]$$

$$J(\boldsymbol{\theta}) = \frac{1}{2m} [\boldsymbol{\theta}^T X^T X \boldsymbol{\theta} - 2y^T X \boldsymbol{\theta}] + \text{constant terms independent of } \boldsymbol{\theta}$$

Compute the derivative with respect to $\boldsymbol{\theta}$ (matrix calculus):

- Derivative of $\boldsymbol{\theta}^T X^T X \boldsymbol{\theta}$:

$$\frac{\partial}{\partial \boldsymbol{\theta}} (\boldsymbol{\theta}^T X^T X \boldsymbol{\theta}) = 2X^T X \boldsymbol{\theta}$$

- Derivative of $-2y^T X \boldsymbol{\theta}$:

$$\frac{\partial}{\partial \boldsymbol{\theta}} (-2y^T X \boldsymbol{\theta}) = -2X^T y$$

Substituting above values, we get:

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \frac{1}{m} (X^T X \boldsymbol{\theta} - X^T y) \quad (1.17)$$

From the general gradient, we can isolate the derivative with respect to a single parameter θ_j . The gradient is:

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \boldsymbol{\theta}} = \frac{1}{m} X^T (X \cdot \boldsymbol{\theta} - y)$$

The j -th element of the gradient corresponds to:

$$\frac{\partial J(\boldsymbol{\theta})}{\partial \theta_j} = \frac{1}{m} (X_{\text{column},j})^T (X \cdot \boldsymbol{\theta} - y)$$

1.2.3 Gradient of the Cost Function

The gradient of the cost function with respect to all parameters $\boldsymbol{\theta}$ is a vector:

$$\text{grad} = \begin{bmatrix} \frac{\partial J(\boldsymbol{\theta})}{\partial \theta_0} \\ \frac{\partial J(\boldsymbol{\theta})}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\boldsymbol{\theta})}{\partial \theta_n} \end{bmatrix} \quad (1.18)$$

$$\text{grad} = \frac{1}{m} X^T (X \cdot \boldsymbol{\theta} - y) \quad (1.19)$$

1.2.4 Analytical Solution

To find the optimal $\boldsymbol{\theta}$, set the gradient grad to zero:

$$\frac{1}{m} X^T (X \cdot \boldsymbol{\theta} - y) = 0$$

Simplify:

$$X^T (X \cdot \boldsymbol{\theta}) = X^T y$$

Rearranging:

$$\boldsymbol{\theta} = (X^T X)^{-1} X^T y \quad (1.20)$$

This is known as the normal equation, which gives the closed-form solution for $\boldsymbol{\theta}$ in linear regression.

1.2.5 Gradient Descent Algorithm

Algorithm 2 Gradient Descent for Linear Regression for multiple features

-
- 1: **Input:** Learning rate α , initial values for $\boldsymbol{\theta}$, and maximum iterations or convergence threshold ϵ .
 - 2: **Output:** Optimized parameter vector, $\boldsymbol{\theta}$.
 - 3: **repeat**
 - 4: Compute gradient vector: $\mathbf{grad} \in \mathbb{R}^{(n+1) \times 1}$:

$$\mathbf{grad} = X^T(X \cdot \boldsymbol{\theta} - y)$$

- 5: Update the parameter vector, $\boldsymbol{\theta}$:

$$\boldsymbol{\theta} := \boldsymbol{\theta} - \alpha \cdot \mathbf{grad}$$

- 6: Compute cost function (just for monitoring):

$$J(\boldsymbol{\theta}) = \frac{1}{2m}(X \cdot \boldsymbol{\theta} - y)^T(X \cdot \boldsymbol{\theta} - y)$$

- 7: **until** maximum iterations reached **or** change in $J(\boldsymbol{\theta})$ is below threshold ϵ .
 - 8: **return** $\boldsymbol{\theta}$
-

1.3 Linear Models for Regression

1.3.1 Polynomial Curve Fitting

$$\hat{y}(x, \boldsymbol{\theta}) = \theta_0 + \theta_1 x + \theta_2 x^2 + \cdots + \theta_n x^n \quad (1.21)$$

$$\text{Let } X = \begin{bmatrix} 1 & x^{(1)} & (x^{(1)})^2 & (x^{(1)})^3 & \cdots & (x^{(1)})^n \\ 1 & x^{(2)} & (x^{(2)})^2 & (x^{(2)})^3 & \cdots & (x^{(2)})^n \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x^{(m)} & (x^{(m)})^2 & (x^{(m)})^3 & \cdots & (x^{(m)})^n \end{bmatrix} \quad (1.22)$$

and the target vector:

$$y = \begin{bmatrix} y^{(1)} \\ y^{(2)} \\ \vdots \\ y^{(m)} \end{bmatrix} \quad (1.23)$$

1.3.2 Analytical Solution

Just like previous multi-feature solution:

$$\boldsymbol{\theta}^* = (X^T X)^{-1} X^T \mathbf{y} \quad (1.24)$$

1.3.3 Linear Models with Basis Functions

$$\hat{y} = h_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^T \boldsymbol{\Phi}(\mathbf{x}) = \sum_{j=0}^n \theta_j \phi_j(\mathbf{x}) \quad (1.25)$$

where $\Phi(\mathbf{x})$ is a matrix of basis functions:

$$\Phi = \begin{bmatrix} \phi_0(x^{(1)}) & \phi_1(x^{(1)}) & \phi_2(x^{(1)}) & \cdots & \phi_n(x^{(1)}) \\ \phi_0(x^{(2)}) & \phi_1(x^{(2)}) & \phi_2(x^{(2)}) & \cdots & \phi_n(x^{(2)}) \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ \phi_0(x^{(m)}) & \phi_1(x^{(m)}) & \phi_2(x^{(m)}) & \cdots & \phi_n(x^{(m)}) \end{bmatrix} \quad (1.26)$$

The cost function is:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (e^{(i)})^2$$

Using the matrix form, the residuals for all data points can be written as $E = \Phi \cdot \theta - y$, and the cost function becomes:

$$J(\theta) = \frac{1}{2m} \mathbf{E}^T \mathbf{E} = \frac{1}{2m} (\Phi \cdot \theta - y)^T (\Phi \cdot \theta - y)$$

The gradient of $J(\theta)$ with respect to the parameters θ is:

$$\text{grad} = \frac{\partial J(\theta)}{\partial \theta} = \begin{bmatrix} \frac{\partial J(\theta)}{\partial \theta_0} \\ \frac{\partial J(\theta)}{\partial \theta_1} \\ \vdots \\ \frac{\partial J(\theta)}{\partial \theta_n} \end{bmatrix} = \frac{1}{m} \Phi^T \cdot (\Phi \cdot \theta - y)$$

From the cost function:

$$J(\theta) = \frac{1}{2m} (\Phi \cdot \theta - y)^T (\Phi \cdot \theta - y)$$

Taking the derivative with respect to θ :

$$\text{grad} = \frac{1}{m} \Phi^T \cdot (\Phi \cdot \theta - y)$$

To find the optimal θ^* , we set the gradient to zero:

$$\frac{\partial J(\theta)}{\partial \theta} = 0$$

This gives the normal equation:

$$\Phi^T \cdot \Phi \cdot \theta = \Phi^T \cdot y$$

Rearranging:

$$\theta^* = (\Phi^T \cdot \Phi)^{-1} \cdot \Phi^T \cdot y$$

Gradient Descent Update Rule:

$$\theta := \theta - \alpha \cdot \text{grad}$$

Substituting the gradient:

$$\theta := \theta - \alpha \cdot \frac{1}{m} \Phi^T \cdot (\Phi \cdot \theta - y)$$

Here the α is the learning rate, controlling the step size.

1.4 Probabilistic Interpretation of Linear Regression

Treating the target variable y as a random variable conditioned on the input features \mathbf{x} , with the assumption that the observed y values are drawn from a probability distribution centered around the model's prediction.

1.4.1 Equivalence of least square error and maximum likelihood estimation

$$y = \boldsymbol{\theta}^T \mathbf{x} + e$$

Here we have:

- y : The observed output (dependent variable or target)
- $\boldsymbol{\theta}$: is the parameter vector
- \mathbf{x} : is the feature/input vector
- e : The random noise term, representing unmodeled or unexplained effects

For i^{th} data point:

$$y^{(i)} = \boldsymbol{\theta}^T \mathbf{x}^{(i)} + e^{(i)}$$

Model Assumptions:

- The relationship between the features (\mathbf{x}) and the target (y) is linear:

$$\mu_y = h_{\boldsymbol{\theta}}(\mathbf{x}) = \boldsymbol{\theta}^T \mathbf{x},$$

where μ_y is the mean of y , predicted by the linear model.

- Observations y deviate from the linear prediction due to noise or unmodeled effects, captured by an additive random error term e :

$$y = \boldsymbol{\theta}^T \mathbf{x} + e.$$

The noise e is assumed to follow a Gaussian distribution:

$$e \sim \mathcal{N}(0, \sigma^2),$$

where σ^2 is the variance of the noise.

- Given \mathbf{x} , y is conditionally normally distributed:

$$y \mid \mathbf{x} \sim \mathcal{N}(\boldsymbol{\theta}^T \mathbf{x}, \sigma^2).$$

Probabilistic Model:

The probability density function (PDF) of y given \mathbf{x} is:

$$p(y \mid \mathbf{x}, \boldsymbol{\theta}, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{(y - \boldsymbol{\theta}^T \mathbf{x})^2}{2\sigma^2}\right). \quad (1.27)$$

This represents the likelihood of observing a particular y value for a given input \mathbf{x} . Now we can estimate the parameters $\boldsymbol{\theta}$ and σ^2 by maximizing the likelihood of the observed data. For a dataset $\{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^m$, the likelihood is the joint probability of all m observations:

$$\mathcal{L}(\boldsymbol{\theta}, \sigma^2) = \prod_{i=1}^m p(y^{(i)} | \mathbf{x}^{(i)}, \boldsymbol{\theta}, \sigma^2) = \prod_{i=1}^m \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(y^{(i)} - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}))^2}{2\sigma^2}\right) \quad (1.28)$$

Taking \ln both sides:

$$\ell(\boldsymbol{\theta}, \sigma^2) = \ln \mathcal{L}(\boldsymbol{\theta}, \sigma^2) = -\frac{m}{2} \ln(2\pi\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^m (y^{(i)} - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}))^2 \quad (1.29)$$

Here, maximize the log-likelihood with respect to $\boldsymbol{\theta}$ is equivalent to minimizing the sum of squared errors, under the assumption of Gaussian distribution for the error:

$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} \sum_{i=1}^m (y^{(i)} - \boldsymbol{\theta}^T \mathbf{x}^{(i)})^2. \quad (1.30)$$

This leads to the familiar solution of linear regression (normal equation):

$$\boldsymbol{\theta}^* = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{y} \quad (1.31)$$

Thus, the optimal value of σ^2 can be calculated as:

$$\sigma_{ML}^2 = \frac{1}{m} \sum_{i=1}^m (y^{(i)} - h_{\boldsymbol{\theta}_{ML}}(\mathbf{x}^{(i)}))^2 \quad (1.32)$$

which is equivalent to minimized cost function.

For any new input \mathbf{x} , the predicted value \hat{y} is computed using $h_{\boldsymbol{\theta}_{ML}}(\mathbf{x})$.

Bibliography

- [1] Weidong Kuang. *Fundamentals of deep learning: a step-by-step guide*. Self, 2024. <https://faculty.utrgv.edu/weidong.kuang/book/book.html> (visited 2024-12-01).
- [2] Simon J.D. Prince. *Understanding Deep Learning*. The MIT Press, 2023.
- [3] Amit Thakur. *Maths for Machine Learning Notes*. Self, 2025. <http://amitthakur.org/learnings/maths/maths-for-ml>.