

Linear Regression

1 Problem Formulation

Input space:

$$X \subseteq \mathbb{R}^d \quad \mathcal{X} \subseteq \mathbb{R}^d$$

Feature vector $x = (x_1, x_2, \dots, x_d)$

Output space:

$$Y \subseteq \mathbb{R} \quad \mathcal{Y} \subseteq \mathbb{R}$$

Continuous target variable

Data distribution:

Samples $(x_i, y_i) \in \mathcal{X} \times \mathcal{Y}$ drawn i.i.d. from unknown distribution $P(X, Y)$

Learning Objective (Expected Risk Minimization):

$$R(w) = E(X, Y)[(Y - f_w(X))^2] \quad R(w) = \mathbb{E}_{(X, Y)}[(Y - f_w(X))^2]$$

Goal:

$$\min_{\{f_0\}} w R(w) \quad \min_w R(w)$$

2 Model Specification

Hypothesis function:

$$f_w(x) = w^T x + b \quad f_w(x) = w^T x + b$$

Parameter space:

$$w \in \mathbb{R}^d, b \in \mathbb{R} \quad w \in \mathbb{R}^d, b \in \mathbb{R}$$

Structural assumptions:

- Linearity in parameters

- Additive noise
 - Homoscedasticity
 - No multicollinearity (ideal case)
-

3 Loss Function

Explicit form (Squared Error Loss):

$$L(y, \hat{y}) = (y - \hat{y})^2$$

Why this loss?

- Corresponds to Gaussian noise assumption
- Maximum Likelihood Estimation under normal errors

Convexity:

- Convex
 - Differentiable
 - Global minimum exists
-

4 Objective Function

Empirical Risk:

$$J(w) = \frac{1}{n} \sum_{i=1}^n (y_i - w^T x_i)^2$$

Regularized (Ridge example):

$$J(w) = \frac{1}{n} \sum_{i=1}^n (y_i - w^T x_i)^2 + \lambda \|w\|^2$$

5 Optimization Method

Closed-form solution (Normal Equation):

$$w = (X^T X)^{-1} X^T y$$

Gradient:

$$\nabla J(w) = -2nXT(y - Xw) \quad J(w) = \frac{1}{2} \|y - Xw\|^2$$

Convergence:

- Convex \Rightarrow global optimum guaranteed

Computational complexity:

- Closed-form: $O(d^3)$
 - Gradient Descent: $O(nd)$ per iteration
-

6 Statistical Interpretation

MLE connection:

Assume:

$$Y = w^T X + \epsilon, \epsilon \sim N(0, \sigma^2) \quad Y = w^T X + \epsilon, \epsilon \sim \mathcal{N}(0, \sigma^2)$$

Then minimizing MSE = maximizing likelihood.

Noise model: Gaussian

Probabilistic meaning:

Predicts conditional mean:

$$E[Y|X] = w^T X$$

7 Regularization & Generalization

Bias–Variance Tradeoff:

- No regularization \rightarrow low bias, high variance
- With Ridge \rightarrow slightly higher bias, lower variance

Overfitting behavior:

- High-degree polynomial features can overfit

Capacity control:

- L2 (Ridge)
 - L1 (Lasso)
-

8 Theoretical Properties

- Convex optimization
 - Unique global minimum (if $X^T X X^T X$ invertible)
 - Consistent estimator (under assumptions)
-

9 Computational Complexity

Training:

- $O(nd^2)O(nd^2)O(nd^2)$ or $O(d^3)O(d^3)O(d^3)$

Inference:

- $O(d)O(d)O(d)$ per prediction

Memory:

- $O(nd)O(nd)O(nd)$
-

10 Limitations

- Sensitive to outliers
- Assumes linear relationship
- Requires homoscedasticity
- Multicollinearity issues