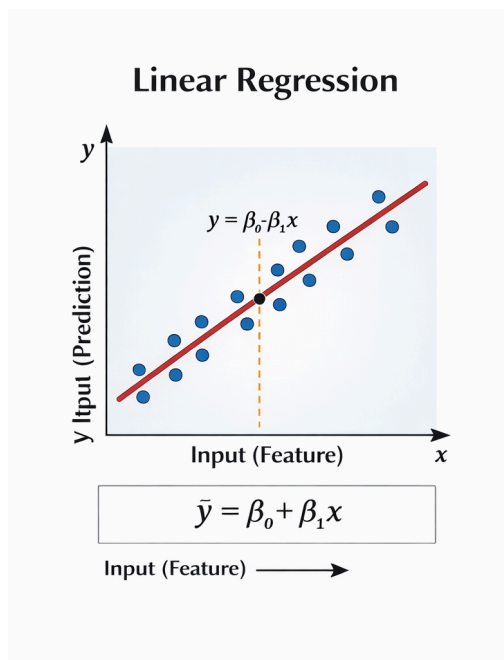


Documentation-Linear Regression



1 Problem Formulation

1.1 Input Space

The input space consists of feature vectors $\mathbf{X} \in \mathbb{R}^d$,

where each data point contains d numerical features.

1.2 Output Space

The output space is continuous and real-valued: $\mathbf{Y} \in \mathbb{R}$

1.3 Data Distribution

Assume data samples (\mathbf{x}_i, y_i) are drawn independently and identically distributed (i.i.d.) from an unknown joint distribution $\mathbf{P}(\mathbf{X}, \mathbf{Y})$.

1.4 Learning Objective (Expected Risk Minimization)

The goal is to learn a function $\mathbf{f}(\mathbf{x})$ that minimizes the expected prediction error:

2 Model Specification

2.1 Hypothesis Function

Linear Regression assumes a linear relationship between input features and output:

$$\mathbf{f}(\mathbf{x}) = \mathbf{w}^T \mathbf{x} + \mathbf{b}$$

where:

- $\mathbf{w} \in \mathbb{R}^d$ are model coefficients
- $\mathbf{b} \in \mathbb{R}$ is the bias term

2.2 Parameter Space

$$\Theta = \{ (\mathbf{w}, \mathbf{b}) \mid \mathbf{w} \in \mathbb{R}^d, \mathbf{b} \in \mathbb{R} \}$$

2.3 Structural Assumptions

- Linearity between features and target
- Additive noise model
- Independence of observations

3 Loss Function

3.1 Explicit Mathematical Form

The most common loss function is Mean Squared Error (MSE): $L(\mathbf{y}, \hat{\mathbf{y}}) = (\mathbf{y} - \hat{\mathbf{y}})^2$

3.2 Why This Loss? (Statistical Reasoning)

- Penalizes larger errors more heavily
- Differentiable and mathematically tractable
- Corresponds to Maximum Likelihood Estimation under Gaussian noise assumption

3.3 Convexity Properties

The squared error loss is convex with respect to model parameters, ensuring a unique global minimum.

4 Objective Function

4.1 Empirical Risk Expression

$$J(\mathbf{w}, \mathbf{b}) = (1/n) \sum_{i=1}^n (y_i - (\mathbf{w}^T \mathbf{x}_i + \mathbf{b}))^2$$

4.2 Regularized Formulation (if used)

- Ridge (L2) Regularization: $J(\mathbf{w}) = (1/n) \sum (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|^2$
- Lasso (L1) Regularization: $J(\mathbf{w}) = (1/n) \sum (y_i - \mathbf{w}^T \mathbf{x}_i)^2 + \lambda \|\mathbf{w}\|_1$

5 Optimization Method

5.1 Closed-form or iterative?

Closed-form solution via Normal Equation: $\mathbf{w} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$

5.2 Gradient expression

Iterative solution via Gradient Descent

$$\nabla J(\mathbf{w}) = -(2/n) \mathbf{X}^T (\mathbf{y} - \mathbf{X}\mathbf{w})$$

5.3 Convergence guarantees

Since the objective is convex, gradient descent converges to a global minimum with a proper learning rate.

5.4 Computational complexity

- Closed-form: $\mathcal{O}(d^3)$ due to matrix inversion
- Gradient Descent: $\mathcal{O}(nd)$ per iteration

6 Statistical Interpretation

6.1 MLE / MAP connection

Linear regression is equivalent to Maximum Likelihood Estimation (MLE) assuming:

$$\mathbf{Y} = \mathbf{w}^T \mathbf{X} + \epsilon$$

Where $\epsilon \sim N(0, \sigma^2)$

6.2 Noise model assumption

Errors are assumed to be normally distributed with constant variance.

6.3 Probabilistic meaning of outputs

The model predicts the conditional mean:

$$E[Y | X]$$

7 Regularization & Generalization

7.1 Bias–variance tradeoff

- High bias \rightarrow Underfitting
- High variance \rightarrow Overfitting
- Regularization balances both

7.2 Overfitting behavior

Occurs when model complexity is high or when multicollinearity exists.

7.3 Capacity control mechanism

Regularization limits coefficient magnitude to prevent overfitting.

8 Theoretical Properties

8.1 Convexity / Global optimality

The optimization problem is convex \rightarrow guarantees global optimum.

8.2 Consistency (if applicable)

Under standard assumptions, the estimator is statistically consistent as sample size increases.

8.3 Stability considerations

Sensitive to outliers due to squared loss.

9 Computational Complexity

9.1 Training time complexity

- Normal Equation: $O(d^3)$
- Gradient Descent: $O(nd \times \text{iterations})$

9.2 Inference time complexity

$O(d)$

9.3 Memory complexity

$O(nd)$

10 Limitations

10.1 When it fails

- Non-linear relationships
- High multicollinearity
- Heteroscedastic data

10.2 Assumption violations

- Non-normal residuals
- Correlated features
- Dependent observations

10.3 Sensitivity issues (outliers, scaling, etc.)

- Sensitive to outliers
- Requires feature scaling for numerical stability