

Introduction & Linear Regression

Deep Dive: Mathematics and Implementation

Methods and Algorithms

MSc Data Science

Spring 2026

- 1 Problem
- 2 Method
- 3 Solution
- 4 Practice
- 5 Decision Framework
- 6 Summary

Finance Use Case: House Price Prediction

- Banks need accurate property valuations for mortgage decisions
- Insurers must estimate replacement costs
- Investors evaluate real estate portfolios

Why Linear Regression?

- Interpretable coefficients (how much does each feature matter?)
- Regulatory requirement for explainable models
- Fast, well-understood baseline

Linear regression: the workhorse of quantitative finance since the 1800s

Capital Asset Pricing Model – Linear Regression in Finance

$$R_i - R_f = \alpha_i + \beta_i(R_m - R_f) + \varepsilon_i \quad (1)$$

- R_i : Return of asset i
- R_f : Risk-free rate (e.g., T-bill)
- R_m : Market return (e.g., S&P 500)
- β_i : Systematic risk (market sensitivity)
- α_i : Abnormal return (should be zero if CAPM holds)

Interpretation: $\beta = 1.2$ means 10% market rise \Rightarrow 12% expected asset rise

CAPM: The original factor model – basis for portfolio management

The Model in Matrix Form

$$\mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (2)$$

- $\mathbf{y} \in \mathbb{R}^n$: Response vector
- $\mathbf{X} \in \mathbb{R}^{n \times (p+1)}$: Design matrix (with intercept column)
- $\boldsymbol{\beta} \in \mathbb{R}^{p+1}$: Coefficient vector
- $\boldsymbol{\varepsilon} \in \mathbb{R}^n$: Error vector

Matrix notation enables elegant derivations and efficient computation

The Design Matrix \mathbf{X}

$$\mathbf{X} = \begin{bmatrix} 1 & x_{11} & x_{12} & \cdots & x_{1p} \\ 1 & x_{21} & x_{22} & \cdots & x_{2p} \\ \vdots & \vdots & \vdots & \ddots & \vdots \\ 1 & x_{n1} & x_{n2} & \cdots & x_{np} \end{bmatrix} \quad (3)$$

- First column of 1s for intercept β_0
- Each row is one observation
- Each column (after first) is one feature

n observations, p features, $p + 1$ parameters

Classical Assumptions for Valid Inference

1. **Linearity:** $E[y|X] = X\beta$ (correct functional form)
2. **Exogeneity:** $E[\varepsilon|X] = 0$ (no omitted variables)
3. **Homoscedasticity:** $\text{Var}(\varepsilon|X) = \sigma^2 I$ (constant variance)
4. **No multicollinearity:** $\text{rank}(X) = p + 1$ (full rank)
5. **Normality** (for inference): $\varepsilon \sim N(0, \sigma^2 I)$

Violations? Robust standard errors, transformations, regularization

Assumptions 1-4 needed for unbiased estimates; 5 for t-tests and CIs

Why OLS is Special

Under assumptions 1-4 (linearity, exogeneity, homoscedasticity, no perfect multicollinearity):

OLS is BLUE – Best Linear Unbiased Estimator

What BLUE Means:

- **Best:** Lowest variance among all linear unbiased estimators
- **Linear:** Estimator is linear function of y
- **Unbiased:** $E[\hat{\beta}] = \beta$

Implication: You cannot find a better linear unbiased estimator than OLS

Gauss-Markov justifies why OLS is the default choice for linear regression

Sum of Squared Residuals (SSR)

$$L(\beta) = \sum_{i=1}^n (y_i - \hat{y}_i)^2 = (\mathbf{y} - \mathbf{X}\beta)^\top (\mathbf{y} - \mathbf{X}\beta) \quad (4)$$

Expanding:

$$L(\beta) = \mathbf{y}^\top \mathbf{y} - 2\beta^\top \mathbf{X}^\top \mathbf{y} + \beta^\top \mathbf{X}^\top \mathbf{X} \beta \quad (5)$$

Quadratic function in β – has unique minimum (if X full rank)

Deriving the Normal Equation

Taking the Derivative

$$\frac{\partial L}{\partial \beta} = -2\mathbf{X}^\top \mathbf{y} + 2\mathbf{X}^\top \mathbf{X} \beta \quad (6)$$

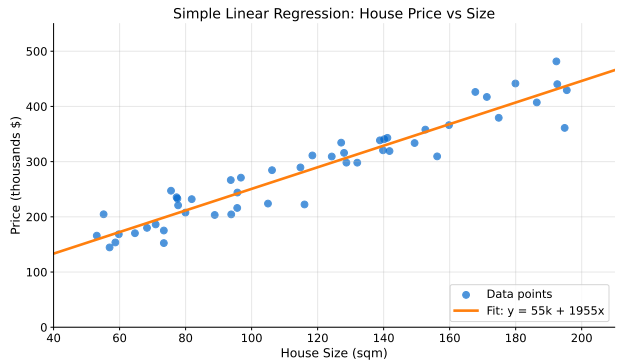
Setting to Zero:

$$\mathbf{X}^\top \mathbf{X} \hat{\beta} = \mathbf{X}^\top \mathbf{y} \quad (7)$$

$$\hat{\beta} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (8)$$

This is the closed-form OLS solution – the “normal equation”

Simple Regression Visualization

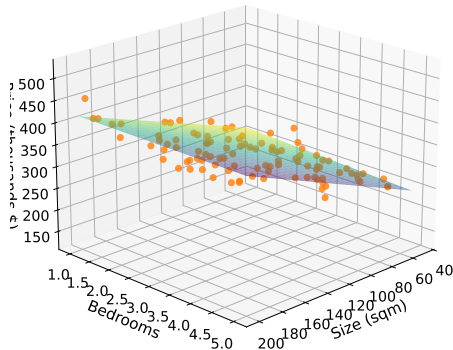


https://github.com/Digital-AI-Finance/methods-algorithms/tree/master/slides/01_Introduction_Linear_Regression/01_simple_regression

The fitted line minimizes vertical distances squared

Multiple Regression Surface

Multiple Regression: Price = $f(\text{Size}, \text{Bedrooms})$



https://github.com/Digital-AI-Finance/methods-algorithms/tree/master/slides/L01_introduction_Linear_Regression/02_multiple_regression_3d

With 2 features, we fit a plane; with p features, a hyperplane

Standardization: Zero Mean, Unit Variance

$$x_j^{\text{scaled}} = \frac{x_j - \bar{x}_j}{s_j} \quad (9)$$

Why Scale Features?

1. **Coefficient comparison:** After scaling, $|\beta_j|$ reflects relative importance
2. **Gradient descent:** Converges faster with similar feature scales
3. **Regularization:** Fair penalty across all features

Caution: Standardized coefficients lose original units (trade-off)

Always standardize for regularization; optional for OLS if only predicting

Why Gradient Descent?

Normal Equation Limitations

- Computing $(\mathbf{X}^\top \mathbf{X})^{-1}$ is $O(p^3)$
- Memory: Need to store $p \times p$ matrix
- For large p (millions of features): infeasible

Gradient Descent Advantages

- Memory efficient: process one sample at a time
- Scales to big data (SGD)
- Generalizes to non-linear models

For $p > 10,000$, gradient descent usually faster

Gradient of the Loss Function

$$\nabla L(\beta) = -2\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\beta) = -2\mathbf{X}^\top \mathbf{r} \quad (10)$$

where $\mathbf{r} = \mathbf{y} - \mathbf{X}\beta$ is the residual vector.

Intuition:

- Gradient points in direction of steepest ascent
- We move opposite to gradient (steepest descent)
- Scale by learning rate α

Gradient is a $p + 1$ dimensional vector

Update Rule

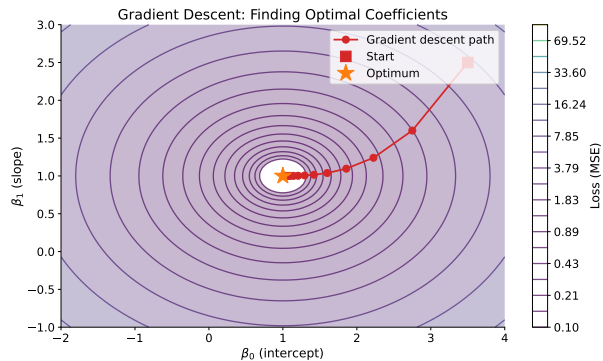
$$\beta^{(t+1)} = \beta^{(t)} - \alpha \nabla L(\beta^{(t)}) \quad (11)$$

Algorithm:

1. Initialize $\beta^{(0)}$ (often zeros or random)
2. Compute gradient $\nabla L(\beta^{(t)})$
3. Update: $\beta^{(t+1)} = \beta^{(t)} - \alpha \nabla L$
4. Repeat until convergence

Convergence: $\|\beta^{(t+1)} - \beta^{(t)}\| < \epsilon$ or max iterations

Gradient Descent Visualization



https://github.com/Digital-AI-Finance/methods-algorithms/tree/master/slides/01_introduction_linear_regression/04_gradient_descent

Contours show loss surface; path shows optimization trajectory

The Critical Hyperparameter

- **Too small:** Slow convergence, many iterations
- **Too large:** Divergence, oscillation
- **Just right:** Fast, stable convergence

Practical Approaches:

- Start with $\alpha = 0.01$ or 0.001
- Learning rate schedules (decay over time)
- Adaptive methods: Adam, AdaGrad, RMSprop

For OLS, optimal $\alpha = 1/\lambda_{\max}(X^T X)$

Mini-Batch Gradient Descent

Instead of full gradient:

$$\nabla L(\beta) = -\frac{2}{n} \mathbf{X}^\top \mathbf{r} \quad (12)$$

Use mini-batch of size m :

$$\nabla L_B(\beta) = -\frac{2}{m} \mathbf{X}_B^\top \mathbf{r}_B \quad (13)$$

- $m = 1$: Stochastic GD (noisy but fast)
- $m = n$: Batch GD (stable but slow)
- $m \in [32, 256]$: Mini-batch (good tradeoff)

SGD: Process data once per epoch, update many times

Coefficient of Determination

$$R^2 = 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}} = 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2} \quad (14)$$

Interpretation:

- Proportion of variance explained by model
- $R^2 = 0$: Model no better than mean
- $R^2 = 1$: Perfect fit
- $R^2 = 0.7$: 70% of variance explained

R^2 always increases with more features – use Adjusted R^2

Penalizing Model Complexity

$$R_{\text{adj}}^2 = 1 - \frac{(1 - R^2)(n - 1)}{n - p - 1} \quad (15)$$

Properties:

- Adjusts for number of predictors p
- Can decrease when adding irrelevant features
- Better for model comparison

Use R_{adj}^2 when comparing models with different p

Error Metrics in Original Units

$$\text{RMSE} = \sqrt{\frac{1}{n} \sum (y_i - \hat{y}_i)^2} \quad (16)$$

$$\text{MAE} = \frac{1}{n} \sum |y_i - \hat{y}_i| \quad (17)$$

Comparison:

- RMSE: Penalizes large errors more (sensitive to outliers)
- MAE: More robust, easier to interpret
- Units: Same as target variable (e.g., dollars)

Report both for comprehensive evaluation

Quantifying Uncertainty in Estimates

$$\text{Var}(\hat{\beta}) = \sigma^2(\mathbf{X}^\top \mathbf{X})^{-1} \quad (18)$$

Standard Error of $\hat{\beta}_j$:

$$\text{SE}(\hat{\beta}_j) = \hat{\sigma} \sqrt{[(\mathbf{X}^\top \mathbf{X})^{-1}]_{jj}} \quad (19)$$

where $\hat{\sigma}^2 = \frac{1}{n-p-1} \sum (y_i - \hat{y}_i)^2$ (unbiased variance estimate)

SE tells us how much $\hat{\beta}_j$ would vary across different samples

Is Feature j Significant?

- $H_0 : \beta_j = 0$ (feature has no effect)
- $H_1 : \beta_j \neq 0$ (feature matters)

t-Statistic:

$$t_j = \frac{\hat{\beta}_j - 0}{\text{SE}(\hat{\beta}_j)} \sim t_{n-p-1} \quad (20)$$

Decision Rule:

- p-value < 0.05 : Reject H_0 , coefficient is significant
- p-value ≥ 0.05 : Cannot reject H_0

Always check p-values before interpreting coefficients

95% CI for Coefficient β_j :

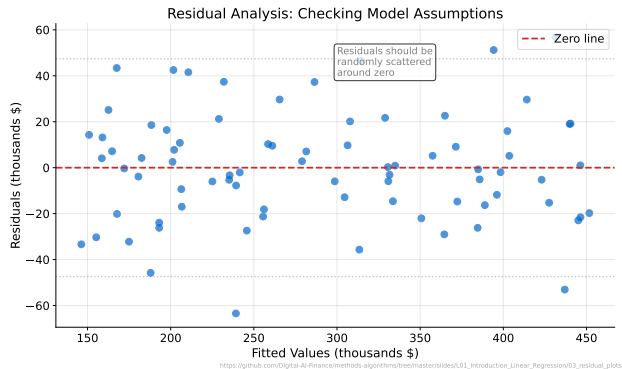
$$\hat{\beta}_j \pm t_{n-p-1, 0.975} \times \text{SE}(\hat{\beta}_j) \quad (21)$$

Interpretation: If we repeated the study many times, 95% of intervals would contain the true β_j

Prediction Intervals:

- **Confidence interval for mean:** $\hat{y} \pm t \cdot \text{SE}(\hat{y})$
- **Prediction interval for new observation:** Wider (includes σ^2)

CI for mean is narrower; prediction interval accounts for individual variability



Good: random scatter. **Bad:** patterns indicate model misspecification

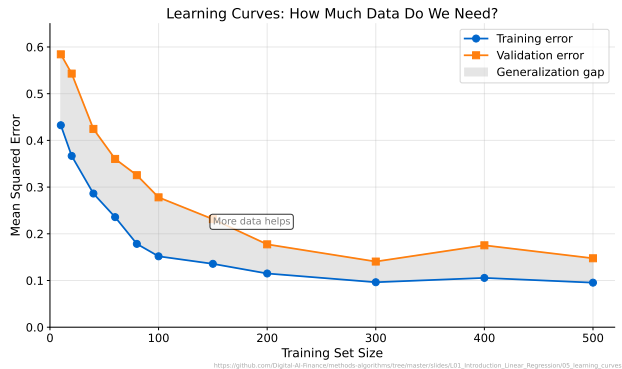
Evaluating Generalization

- Never evaluate on training data alone
- Split: 70-80% train, 20-30% test
- Report test set metrics

Cross-Validation (K-Fold):

- Split into K folds (typically $K = 5$ or 10)
- Train on $K - 1$ folds, validate on 1
- Repeat K times, average results

CV gives more reliable estimate with limited data



Gap between curves indicates overfitting; convergence shows saturation

The Overfitting Problem

When Models Memorize Instead of Learn

- High-dimensional data ($p \approx n$ or $p > n$)
- Coefficients become very large
- Perfect fit on training data, poor generalization

Solution: Add Penalty to Loss Function

$$L_{\text{reg}}(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda \cdot \text{penalty}(\beta) \quad (22)$$

λ controls strength of regularization

L2 Penalty: Sum of Squared Coefficients

$$L_{\text{ridge}}(\boldsymbol{\beta}) = \|\mathbf{y} - \mathbf{X}\boldsymbol{\beta}\|^2 + \lambda\|\boldsymbol{\beta}\|_2^2 \quad (23)$$

Closed-Form Solution:

$$\hat{\boldsymbol{\beta}}_{\text{ridge}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} \quad (24)$$

- Shrinks all coefficients toward zero
- Never sets coefficients exactly to zero
- Always invertible (even when $p > n$)

Ridge adds λ to diagonal – stabilizes inversion

L1 Penalty: Sum of Absolute Coefficients

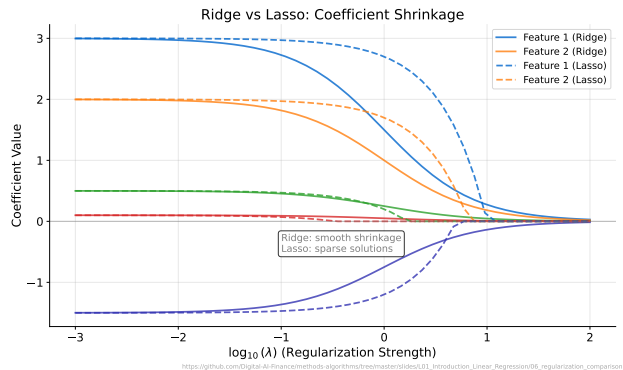
$$L_{\text{lasso}}(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda\|\beta\|_1 \quad (25)$$

Properties:

- Produces sparse solutions (some $\beta_j = 0$)
- Automatic feature selection
- No closed-form solution (use coordinate descent)

Lasso: Least Absolute Shrinkage and Selection Operator

Ridge vs Lasso Comparison



Ridge: smooth shrinkage. Lasso: sparse (feature selection)

Combining L1 and L2 Penalties

$$L_{\text{elastic}}(\beta) = \|\mathbf{y} - \mathbf{X}\beta\|^2 + \lambda_1 \|\beta\|_1 + \lambda_2 \|\beta\|_2^2 \quad (26)$$

Benefits:

- Sparsity from L1
- Stability from L2 (handles correlated features)
- Two hyperparameters to tune

Often best of both worlds for correlated features

Cross-Validation for Hyperparameter Tuning

1. Define grid of λ values (e.g., 10^{-4} to 10^4)
2. For each λ , perform K-fold CV
3. Select λ with lowest CV error
4. Refit on full training data

In Practice:

- `sklearn.linear_model.RidgeCV`
- `sklearn.linear_model.LassoCV`

Larger λ = more regularization = simpler model

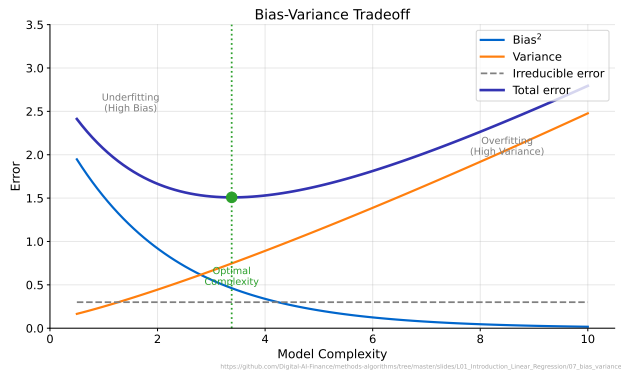
Expected Prediction Error

$$E[(y - \hat{f}(x))^2] = \text{Bias}^2(\hat{f}) + \text{Var}(\hat{f}) + \sigma^2 \quad (27)$$

- **Bias:** Error from wrong assumptions (underfitting)
- **Variance:** Error from sensitivity to training data (overfitting)
- σ^2 : Irreducible noise in data

We can't reduce irreducible error – focus on bias and variance

The Tradeoff Illustrated



Optimal complexity minimizes total error

How Regularization Helps

- Increasing λ : **increases bias, decreases variance**
- Decreasing λ : decreases bias, increases variance
- Optimal λ : minimizes total error

In Practice:

- Use CV to find optimal λ
- Regularization almost always helps when p is large

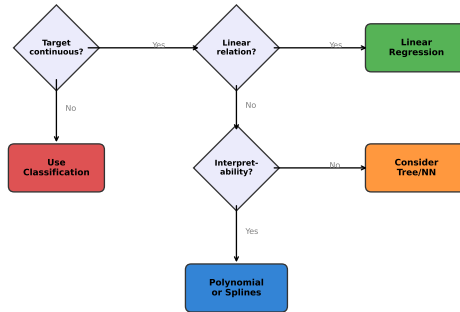
Regularization trades a little bias for a lot of variance reduction

Open the Colab Notebook

- Exercise 1: Implement OLS from scratch
- Exercise 2: Use scikit-learn LinearRegression
- Exercise 3: Compare with gradient descent

Link: See course materials on GitHub

Linear Regression Decision Guide



https://github.com/Digital-AI-finance/methods-algorithms/tree/master/slides/L01_Introduction_Linear_Regression/08_decision_flowchart

Use this framework when choosing regression methods

Use When:

- Continuous target variable
- Approximate linear relationships
- Interpretability is critical
- Inference on coefficients needed
- Fast prediction required

Avoid When:

- Target is categorical
- Strong non-linear patterns
- Many outliers present
- Features highly correlated
- Prediction accuracy paramount

When in doubt, linear regression is a strong baseline

Variance Inflation Factor (VIF)

$$\text{VIF}_j = \frac{1}{1 - R_j^2} \quad (28)$$

where R_j^2 is from regressing x_j on all other features.

Interpretation:

- $\text{VIF} = 1$: No correlation with other features
- $\text{VIF} > 5$: Moderate concern
- $\text{VIF} > 10$: Serious multicollinearity

Remedies:

- Remove highly correlated features
- Use Ridge regression ($\lambda > 0$ stabilizes)
- Apply PCA before regression

Always check VIF before trusting coefficient estimates

$$\text{Model: } \mathbf{y} = \mathbf{X}\boldsymbol{\beta} + \boldsymbol{\varepsilon} \quad (29)$$

$$\text{OLS Solution: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X})^{-1} \mathbf{X}^\top \mathbf{y} \quad (30)$$

$$\text{Gradient: } \nabla L = -2\mathbf{X}^\top (\mathbf{y} - \mathbf{X}\boldsymbol{\beta}) \quad (31)$$

$$\text{GD Update: } \boldsymbol{\beta}^{(t+1)} = \boldsymbol{\beta}^{(t)} - \alpha \nabla L \quad (32)$$

$$\text{Ridge: } \hat{\boldsymbol{\beta}} = (\mathbf{X}^\top \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^\top \mathbf{y} \quad (33)$$

$$R^2 : 1 - \frac{\sum (y_i - \hat{y}_i)^2}{\sum (y_i - \bar{y})^2} \quad (34)$$

Key Takeaways

1. Linear regression minimizes squared error – closed form or GD
2. Matrix notation enables efficient computation
3. Gradient descent scales to large datasets
4. Regularization (Ridge/Lasso) prevents overfitting
5. The bias-variance tradeoff guides model complexity
6. Always evaluate on held-out test data

Next Session: Logistic Regression for Classification

- James, Witten, Hastie, Tibshirani (2021). *Introduction to Statistical Learning*. Chapter 3.
- Hastie, Tibshirani, Friedman (2009). *Elements of Statistical Learning*. Chapter 3.
- Bishop (2006). *Pattern Recognition and Machine Learning*. Chapter 3.

Online Resources:

- scikit-learn: https://scikit-learn.org/stable/modules/linear_model.html
- Stanford CS229: <https://cs229.stanford.edu/>