

L02: Logistic Regression

Mathematical Foundations and Implementation

Methods and Algorithms

Spring 2026

The Classification Problem

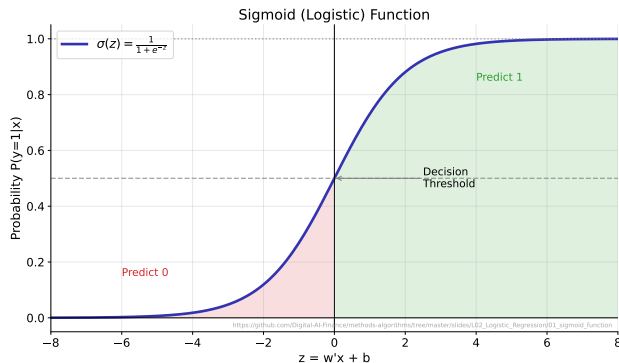
- Given features $\mathbf{x} \in \mathbb{R}^p$, predict $y \in \{0, 1\}$
- Linear regression: $\hat{y} = \mathbf{w}^T \mathbf{x} + b$ (unbounded)
- Need: $P(y = 1|\mathbf{x}) \in [0, 1]$

Solution: The Logistic Function

$$P(y = 1|\mathbf{x}) = \sigma(\mathbf{w}^T \mathbf{x} + b) = \frac{1}{1 + e^{-(\mathbf{w}^T \mathbf{x} + b)}}$$

The logistic function is also called the sigmoid function

The Sigmoid Function



Key Properties:

- Range: $(0, 1)$ – perfect for probabilities
- $\sigma(0) = 0.5$ – threshold for classification
- $\sigma'(z) = \sigma(z)(1 - \sigma(z))$ – simple gradient

Understanding the Model

- Odds: $\frac{P(y=1)}{P(y=0)} = \frac{p}{1-p}$
- Log-odds (logit): $\log\left(\frac{p}{1-p}\right) = \mathbf{w}^T \mathbf{x} + b$

Coefficient Interpretation

- w_j : change in log-odds per unit increase in x_j
- e^{w_j} : odds ratio – multiplicative effect on odds
- Example: $w_{\text{income}} = 0.5 \Rightarrow$ each unit increase in income multiplies odds by $e^{0.5} \approx 1.65$

Log-odds interpretation is key for regulatory compliance in banking

The Likelihood Function

For observations (x_i, y_i) , the likelihood is:

$$L(\mathbf{w}) = \prod_{i=1}^n p_i^{y_i} (1 - p_i)^{1-y_i}$$

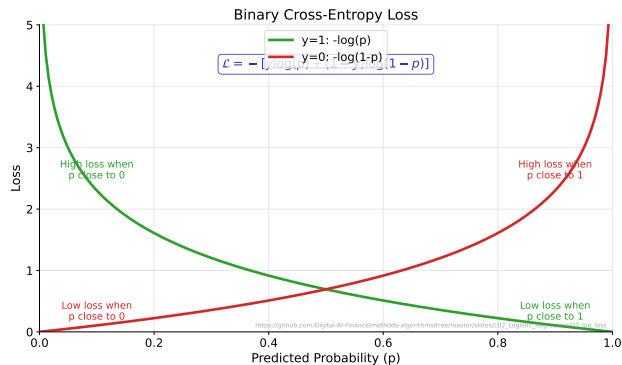
where $p_i = \sigma(\mathbf{w}^T \mathbf{x}_i + b)$

Log-Likelihood (easier to optimize)

$$\ell(\mathbf{w}) = \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$$

Maximize log-likelihood = minimize negative log-likelihood (cross-entropy)

Binary Cross-Entropy Loss



Loss Function

$$\mathcal{L} = -\frac{1}{n} \sum_{i=1}^n [y_i \log(p_i) + (1 - y_i) \log(1 - p_i)]$$

Cross-entropy loss is convex in the weights – guaranteed global optimum

Computing the Gradient

For a single sample:

$$\frac{\partial \mathcal{L}}{\partial w_j} = (p - y)x_j$$

In Matrix Form

$$\nabla_{\mathbf{w}} \mathcal{L} = \frac{1}{n} \mathbf{X}^T (\mathbf{p} - \mathbf{y})$$

where $\mathbf{p} = \sigma(\mathbf{X}\mathbf{w})$

Key Insight: Same form as linear regression gradient!

The elegance of logistic regression: gradient has the same form as linear regression

Update Rule

$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \eta \nabla_{\mathbf{w}} \mathcal{L}$$
$$\mathbf{w}^{(t+1)} = \mathbf{w}^{(t)} - \frac{\eta}{n} \mathbf{X}^T (\sigma(\mathbf{X} \mathbf{w}^{(t)}) - \mathbf{y})$$

Practical Considerations

- Feature scaling: standardize inputs for faster convergence
- Learning rate: start with $\eta = 0.01$, use line search or decay
- Convergence: monitor loss, check gradient norm ; tolerance

No closed-form solution like normal equation – must use iterative optimization

How Certain Are We About Each Coefficient?

Standard error measures uncertainty in our estimate $\hat{\beta}_j$:

$$SE(\hat{\beta}_j) = \sqrt{[\mathbf{H}^{-1}]_{jj}} \quad (1)$$

where \mathbf{H} is the Hessian matrix (measures how steep the loss landscape is around the optimum).

Intuition:

- Small SE: coefficient is precisely estimated
- Large SE: coefficient is uncertain (wide range of plausible values)
- More data \Rightarrow smaller SE \Rightarrow more certainty

SE tells us: "if we repeated this study, how much would $\hat{\beta}$ vary?"

Is This Feature Significant? (Wald Test)

The Question: Does feature j actually matter, or is its effect just noise?

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

Wald Statistic (z-score):

$$z = \frac{\hat{\beta}_j}{SE(\hat{\beta}_j)} \quad (2)$$

Decision Rule:

- $|z| > 1.96$ (or p-value < 0.05): coefficient is **significant**
- $|z| \leq 1.96$ (or p-value ≥ 0.05): cannot conclude it matters

Always check p-values before interpreting coefficients in business reports

95% CI for Coefficient:

$$\hat{\beta}_j \pm 1.96 \times \text{SE}(\hat{\beta}_j) \quad (3)$$

95% CI for Odds Ratio:

$$\exp\left(\hat{\beta}_j \pm 1.96 \times \text{SE}(\hat{\beta}_j)\right) \quad (4)$$

Example: Income coefficient $\hat{\beta} = 0.5$, $\text{SE} = 0.1$

- CI for β : [0.304, 0.696]
- CI for OR: [1.36, 2.01] – income increases odds by 36% to 101%

If CI for OR contains 1.0, the effect is not statistically significant

How Well Does the Model Fit the Data?

$$\text{Deviance} = -2 \times \text{Log-Likelihood} \quad (5)$$

Two Key Measures: (Think: how bad is each model?)

- **Null deviance:** Model with only intercept (baseline)
- **Residual deviance:** Model with all features

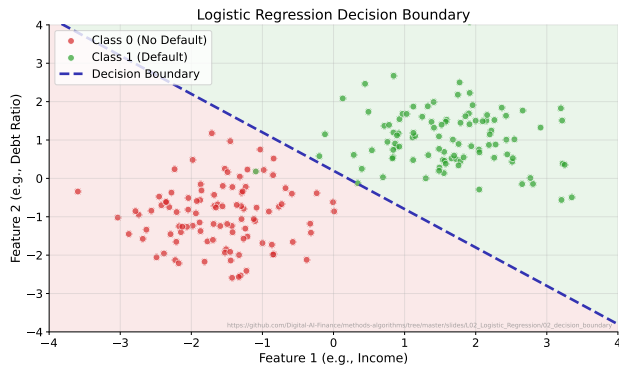
Pseudo R^2 (McFadden):

$$R^2_{\text{McFadden}} = 1 - \frac{\text{Residual Deviance}}{\text{Null Deviance}} \quad (6)$$

Interpretation: Values of 0.2–0.4 are considered good for logistic regression

Deviance drop = model improvement; compare models with Likelihood Ratio test

Linear Decision Boundary



Decision Rule: Predict $\hat{y} = 1$ if $\mathbf{w}^T \mathbf{x} + b \geq 0$

The decision boundary is always a hyperplane in the feature space

Default Threshold: 0.5

- Predict 1 if $P(y = 1|\mathbf{x}) \geq 0.5$
- Equivalent to: $\mathbf{w}^T \mathbf{x} + b \geq 0$

Custom Thresholds

- Lower threshold: more sensitive (higher recall)
- Higher threshold: more specific (higher precision)
- Choose based on business costs of FP vs FN

Example: Fraud Detection

- Cost of missing fraud (FN) \gg Cost of false alarm (FP)
- Use lower threshold, e.g., 0.3

Optimal threshold depends on the cost matrix of your application

Polynomial Features

- Original: $[x_1, x_2]$
- Expanded: $[x_1, x_2, x_1^2, x_2^2, x_1x_2]$
- Creates curved decision boundaries

Trade-offs

- More features: more flexible boundaries
- Risk: overfitting to training data
- Solution: regularization

Logistic regression is linear in parameters, but can model non-linear boundaries

One-vs-Rest (OvR)

- Train K binary classifiers
- Predict class with highest probability

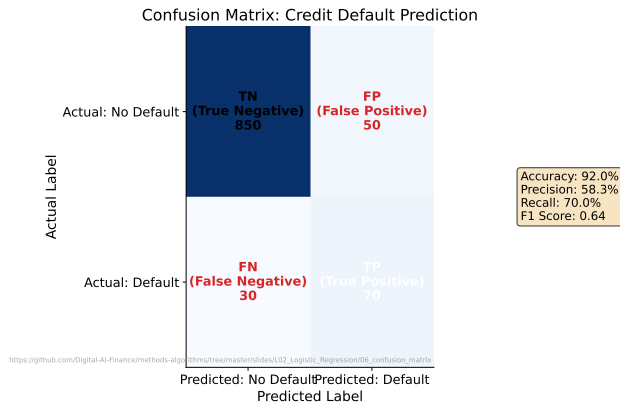
Multinomial (Softmax) Logistic Regression

$$P(y = k|\mathbf{x}) = \frac{e^{\mathbf{w}_k^T \mathbf{x}}}{\sum_{j=1}^K e^{\mathbf{w}_j^T \mathbf{x}}}$$

- Single model, probabilities sum to 1
- Loss: categorical cross-entropy

scikit-learn: `multi_class='multinomial'` for true softmax regression

Confusion Matrix



Always start evaluation by examining the confusion matrix

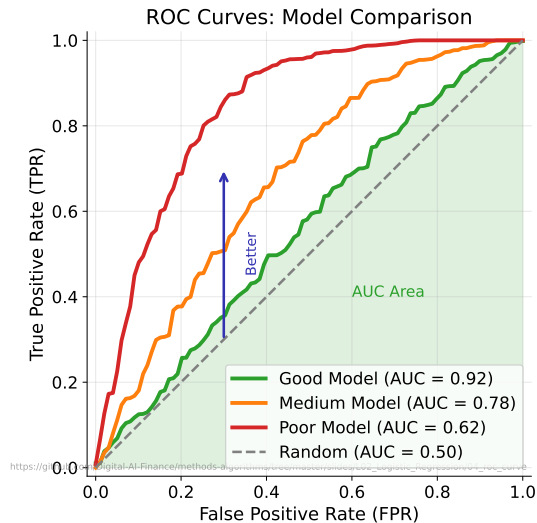
From the Confusion Matrix

- **Accuracy:** $\frac{TP+TN}{TP+TN+FP+FN}$ – overall correctness
- **Precision:** $\frac{TP}{TP+FP}$ – of predicted positives, how many correct?
- **Recall:** $\frac{TP}{TP+FN}$ – of actual positives, how many found?
- **F1 Score:** $\frac{2 \cdot \text{Precision} \cdot \text{Recall}}{\text{Precision} + \text{Recall}}$

When Accuracy Fails

- Imbalanced data: 99% negative class
- Predicting all negatives gives 99% accuracy!

Accuracy is misleading for imbalanced datasets



ROC = Receiver Operating Characteristic: X-axis FPR, Y-axis TPR

Interpretation

- $AUC = 0.5$: random guessing
- $AUC = 1.0$: perfect classifier
- $AUC = \text{probability}(\text{random positive} > \text{random negative})$

Guidelines

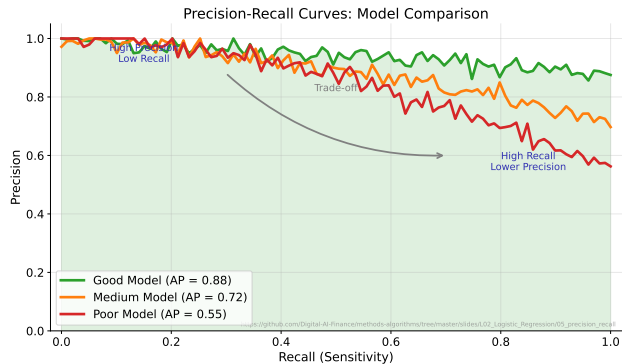
- 0.9–1.0: Excellent
- 0.8–0.9: Good
- 0.7–0.8: Fair
- 0.6–0.7: Poor
- 0.5–0.6: Fail

Finance Equivalent: $\text{Gini coefficient} = 2 \times AUC - 1$

- $AUC = 0.80$ means $\text{Gini} = 0.60$
- Banks often report Gini instead of AUC

AUC is threshold-independent – summarizes performance across all thresholds

Precision-Recall Curve



Use PR curves for imbalanced datasets where positive class is rare

When to Use ROC

- Balanced classes
- Care equally about both classes
- Comparing models at specific FPR

When to Use Precision-Recall

- Imbalanced classes (fraud, disease)
- Positive class is more important
- High precision required

ROC can be overly optimistic with imbalanced data

What is Calibration?

- Predicted 70% probability should mean 70% actually positive
- Well-calibrated: predicted probabilities match observed frequencies

Checking Calibration

- Reliability diagram (calibration plot)
- Brier score: $\frac{1}{n} \sum (p_i - y_i)^2$

Logistic Regression Advantage

- Naturally well-calibrated (MLE property)
- Unlike trees/random forests that may need calibration

Calibration is crucial when probabilities are used for decision-making

The Overfitting Problem

- Many features, limited data
- Model fits noise, not signal
- Perfect training accuracy, poor test performance

Solution: Penalize Large Coefficients

$$\mathcal{L}_{\text{regularized}} = \mathcal{L} + \lambda \cdot \text{penalty}(\mathbf{w})$$

- λ : regularization strength (hyperparameter)
- Larger λ = simpler model

Regularization trades bias for variance

L2 (Ridge)

$$\mathcal{L}_{\text{Ridge}} = \mathcal{L} + \lambda \sum_{j=1}^p w_j^2$$

- Shrinks coefficients toward zero
- Keeps all features, reduces magnitude

L1 (Lasso)

$$\mathcal{L}_{\text{Lasso}} = \mathcal{L} + \lambda \sum_{j=1}^p |w_j|$$

- Some coefficients exactly zero
- Automatic feature selection

L1 for sparse models, L2 when all features likely relevant

Best of Both Worlds

$$\mathcal{L}_{\text{ElasticNet}} = \mathcal{L} + \lambda_1 \sum |w_j| + \lambda_2 \sum w_j^2$$

Advantages

- Handles correlated features better than Lasso alone
- Can select groups of correlated features
- More stable feature selection

In scikit-learn

- `LogisticRegression(penalty='elasticnet', solver='saga', l1_ratio=0.5)`

Elastic Net: `l1_ratio = 1` is pure L1, `l1_ratio = 0` is pure L2

Cross-Validation

- Try grid of λ values: [0.001, 0.01, 0.1, 1, 10, 100]
- Use k-fold CV to estimate test performance
- Select λ with best CV score

scikit-learn Convenience

- `LogisticRegressionCV`: automatic λ search
- C s: inverse of λ (larger C = less regularization)

`LogisticRegressionCV` does cross-validation internally

Algorithm: Gradient Descent

```
1: Input:  $\mathbf{X}$ ,  $\mathbf{y}$ , learning rate  $\eta$ , max iterations  $T$ 
2: Initialize  $\mathbf{w} = \mathbf{0}$ 
3: for  $t = 1$  to  $T$  do
4:    $\mathbf{p} = \sigma(\mathbf{X}\mathbf{w})$ 
5:    $\nabla = \frac{1}{n}\mathbf{X}^T(\mathbf{p} - \mathbf{y})$ 
6:    $\mathbf{w} = \mathbf{w} - \eta\nabla$ 
7:   if  $\|\nabla\| < \epsilon$  then
8:     break
9:   end if
10: end for
11: return  $\mathbf{w}$ 
```

In practice, use quasi-Newton methods (L-BFGS) for faster convergence

Basic Usage

- `from sklearn.linear_model import LogisticRegression`
- `model = LogisticRegression()`
- `model.fit(X_train, y_train)`
- `y_pred = model.predict(X_test)`
- `y_proba = model.predict_proba(X_test)`

Key Parameters

- `C`: inverse regularization strength (default=1.0)
- `penalty`: 'l1', 'l2', 'elasticnet', 'none'
- `solver`: 'lbfgs', 'liblinear', 'saga'
- `class_weight`: 'balanced' for imbalanced data

`predict_proba` returns $[P(y=0), P(y=1)]$ – use `[:, 1]` for positive class

The Problem

- 99% negatives, 1% positives
- Model predicts all negatives: 99% accuracy!

Solutions

- **Class weights:** `class_weight='balanced'`
- **Oversampling:** SMOTE, random oversampling
- **Undersampling:** random undersampling
- **Threshold tuning:** optimize for F1 or business metric

Weighted Loss

$$\mathcal{L}_{\text{weighted}} = - \sum_i w_{y_i} [y_i \log p_i + (1 - y_i) \log(1 - p_i)]$$

`class_weight='balanced'` sets $w_k \propto 1/n_k$

For Logistic Regression

- **Standardization:** mean=0, std=1 for all features
- **Missing values:** impute or create indicator variable
- **Categorical:** one-hot encoding (drop one level)
- **Interactions:** $x_1 \times x_2$ if domain suggests
- **Non-linearity:** binning or polynomial features

Credit Scoring Example

- Age: may have non-linear effect (bin into groups)
- Debt-to-income ratio: interaction of two features
- Employment length: indicator for ≥ 2 years

Feature engineering often matters more than model selection

Coefficient Analysis

- Sign: direction of effect
- Magnitude: strength (after standardization)
- Odds ratio e^{w_j} : multiplicative effect

Example Interpretation

- $w_{\text{income}} = 0.5$: each \$1000 income increase multiplies odds of approval by $e^{0.5} = 1.65$
- $w_{\text{debt_ratio}} = -1.2$: each 0.1 increase in debt ratio multiplies odds by $e^{-0.12} = 0.89$ (11% decrease)

This interpretability makes logistic regression preferred in regulated industries

Available Solvers in scikit-learn

- `lbfgs`: default, works for L2 and no penalty
- `liblinear`: fast for small data, supports L1
- `saga`: supports all penalties, works for large data
- `newton-cg`: similar to `lbfgs`
- `sag`: stochastic, for very large data

Guidelines

- L1 penalty: use `liblinear` or `saga`
- Large data: `saga` or `sag`
- Default (L2): `lbfgs`

`solver='saga'` is the most versatile but may be slower for small datasets

Warning: “Convergence Warning”

- Model did not converge in `max_iter` iterations
- May mean poor solution

Solutions

- Increase `max_iter` (default=100)
- Standardize features
- Increase regularization (smaller C)
- Use different solver

Always check for convergence warnings in production code

Strengths of Logistic Regression

- Interpretable coefficients
- Well-calibrated probabilities
- Fast training and prediction
- Works well with few samples

Limitations

- Linear decision boundary
- May underfit complex patterns
- Sensitive to outliers (compared to trees)

When to Choose Alternatives

- Complex patterns: Random Forests, Gradient Boosting
- High-dimensional: SVM with RBF kernel
- Interpretability not required: Neural Networks

Start with logistic regression as baseline, then try more complex models

How Banks Actually Use Logistic Regression

- **PD (Probability of Default):** The probability a borrower won't repay
- Banks are **REQUIRED** to estimate PD under Basel regulations
- Logistic regression is the industry standard (interpretable, auditable)

Why Interpretability Matters:

- Regulators require explanation of every coefficient
- Must justify why income/debt ratio affects approval
- Black-box models (neural networks) often rejected by regulators

Basel II/III: International banking regulations requiring PD models

Credit Scorecards: Making Models Usable

Key Industry Metrics:

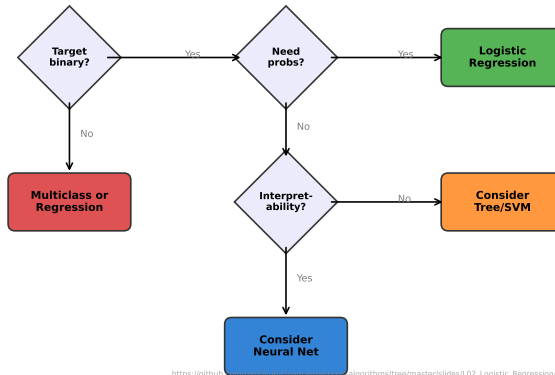
- **Gini Coefficient:** $\text{Gini} = 2 \times \text{AUC} - 1$
 - $\text{AUC} = 0.75 \Rightarrow \text{Gini} = 0.50$ (acceptable)
 - $\text{AUC} = 0.85 \Rightarrow \text{Gini} = 0.70$ (good)
- **KS Statistic:** Maximum separation between defaults and non-defaults

Scorecard Points:

- Convert log-odds to points: higher score = lower risk
- Example: "Each 20 points doubles the odds of being good"

Industry practice: Gini ≥ 0.40 is acceptable; Gini ≥ 0.60 is good

Logistic Regression Decision Guide



https://github.com/DataCamp/algorithms/tree/master/slides/L02_Logistic_Regression/07_decision_flowchart

Logistic regression: first choice for binary classification with interpretability

Interactive Notebook

- Open: `notebooks/L02_logistic_regression.ipynb`
- Dataset: Credit card fraud detection
- Tasks:
 - Train logistic regression with L2 regularization
 - Evaluate with confusion matrix, ROC, PR curves
 - Tune classification threshold
 - Handle class imbalance with class weights
 - Interpret coefficients

Google Colab

- Link: Colab Notebook
- Includes starter code and solutions

Practice exercises reinforce mathematical concepts with real-world implementation

Mathematical Foundation

- Sigmoid function maps linear combination to probability
- Maximum likelihood estimation via gradient descent
- Cross-entropy loss is convex, guaranteed global optimum

Evaluation

- Use confusion matrix, precision, recall, F1
- ROC/AUC for balanced data, PR curve for imbalanced
- Calibration matters when using probabilities

Practice

- Regularization prevents overfitting
- Class weights handle imbalance
- Coefficients are directly interpretable

Logistic regression: simple, fast, interpretable, and often competitive

Textbooks

- James et al. (2021). *ISLR*, Chapter 4: Classification
- Hastie et al. (2009). *ESL*, Chapter 4: Linear Methods

Documentation

- scikit-learn: LogisticRegression user guide
- statsmodels: Logit for statistical inference

Next Lecture

- L03: KNN and K-Means
- From parametric to non-parametric methods