

Lesson 25 Summary: Logistic Regression

Data Science with Python – Key Concepts

Data Science Program

Logistic Regression

Sigmoid Function

Maps to (0,1)
Probability output
S-shaped curve

Decision Boundary

Linear in feature space
Threshold at $p=0.5$
Adjustable threshold

Interpretation

Odds ratios
Log-odds (logit)
 $\text{Exp}(\text{coef}) = \text{OR}$

Multiclass Extension

Softmax for K classes
One-vs-Rest (OvR) or multinomial

Regularization (C parameter)

$C=1/\lambda$ (inverse strength)
Smaller C = more regularization

`LogisticRegression(C=1.0).fit(X, y) | model.predict_proba(X)`

Logistic regression is the go-to model for binary classification

The Sigmoid Function

Transforms linear combination to probability:

- **Formula:** $\sigma(z) = \frac{1}{1+e^{-z}}$
- **Output:** Always between 0 and 1
- **Input:** $z = w^T x + b$ (linear combination)

Sigmoid maps any real number to a probability

Separating classes:

- **Threshold:** Default at probability = 0.5
- **Geometry:** Linear boundary in feature space
- **Flexibility:** Can adjust threshold based on costs

If $P(y=1|x) > \text{threshold}$, predict class 1

The decision boundary is where $w^T x + b = 0$

Odds ratio interpretation:

- **Coefficient:** Change in log-odds per unit x
- **Odds ratio:** $e^{\text{coefficient}}$
- **OR > 1:** Feature increases positive class odds

Odds ratios are multiplicative effect on odds

Getting class probabilities:

- **predict():** Returns class labels (0 or 1)
- **predict_proba():** Returns probability array
- **Use case:** Ranking, threshold tuning

Probabilities are more useful than binary predictions

Regularization (C parameter)

Controlling complexity:

- **C**: Inverse of regularization strength
- **High C**: Less regularization, may overfit
- **Low C**: More regularization, simpler model

Default $C=1.0$ works well for most cases

Extending beyond binary:

- **One-vs-Rest (OvR):** K binary classifiers
- **Multinomial:** Softmax for K classes simultaneously
- **sklearn default:** OvR for small data, multinomial for large

Set `multi_class` parameter explicitly for control

Important for convergence:

- **Problem:** Different scales slow optimization
- **Solution:** StandardScaler before fitting
- **Benefit:** Faster training, better convergence

Always scale features for logistic regression

Credit default prediction:

- **Target:** Default (1) vs No default (0)
- **Features:** Debt ratio, income, history
- **Output:** Probability of default score

Probability calibration is crucial for risk scoring

Essential Commands:

Task	Code
Fit model	<code>LogisticRegression().fit(X, y)</code>
Predict class	<code>model.predict(X)</code>
Get probabilities	<code>model.predict_proba(X)</code>
Coefficients	<code>model.coef_, model.intercept_</code>
Regularization	<code>LogisticRegression(C=0.1)</code>

Logistic regression is interpretable and efficient