

# 逻辑回归（Logistic Regression）：原理、公式、应用与实战

2025 年 9 月 7 日

## 目录

### 1 引言

逻辑回归通过 S 形函数（sigmoid）将线性组合的输出映射到  $[0, 1]$ ，从而建模类别为 1 的条件概率。它具有良好的可解释性与概率输出，常用于风险评估、医疗诊断、CTR 预测等任务。

### 2 原理与公式

设  $\mathbf{x} \in \mathbb{R}^d$ ,  $y \in \{0, 1\}$ , 模型为：

$$p(y = 1 | \mathbf{x}) = \sigma(z), \quad z = w_0 + \mathbf{w}^\top \mathbf{x}, \quad \sigma(t) = \frac{1}{1 + e^{-t}}. \quad (1)$$

对数几率（logit）线性： $\log \frac{p}{1-p} = w_0 + \mathbf{w}^\top \mathbf{x}$ 。

给定样本  $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$ ，二元交叉熵（负对数似然）为：

$$\mathcal{L}(\mathbf{w}, w_0) = - \sum_{i=1}^n [y_i \log p_i + (1 - y_i) \log(1 - p_i)], \quad p_i = \sigma(w_0 + \mathbf{w}^\top \mathbf{x}_i). \quad (2)$$

梯度：

$$\nabla_{\mathbf{w}} \mathcal{L} = \sum_{i=1}^n (p_i - y_i) \mathbf{x}_i, \quad \frac{\partial \mathcal{L}}{\partial w_0} = \sum_{i=1}^n (p_i - y_i). \quad (3)$$

加入  $\ell_2$  正则可缓解过拟合： $\frac{\lambda}{2} \|\mathbf{w}\|^2$ ； $\ell_1$  则有助于稀疏化与特征选择。

阈值取 0.5 时的判别边界满足  $\sigma(z) = 0.5 \iff z = 0$ ，即超平面  $w_0 + \mathbf{w}^\top \mathbf{x} = 0$ 。

### 3 应用场景与要点

- 特征缩放：有助于优化收敛与系数可解释性；
- 类别不平衡：可调阈值、设定类权重或重采样；
- 正则化： $\ell_2$  稳定系数， $\ell_1$  促稀疏，缓解多重共线性；
- 概率输出：便于排序与代价敏感决策；
- 系数解释：关注胜算比  $e^{w_j}$  的含义。

## 4 Python 实战

运行配套脚本以生成本章使用的图片。脚本仅依赖 NumPy 与 Matplotlib，并内置简易的逻辑回归实现，避免版本兼容问题。

Listing 1: gen\_logistic\_regression\_figures.py

```

1  """
2  Generate figures for the Logistic Regression chapter.
3
4  Figure list (saved under ./figures/):
5      - sigmoid_curve.png           : Sigmoid function curve
6      - logistic_loss_curves.png    : Per-sample logistic losses for y=0
7                                     and y=1 vs logit z
8      - decision_boundary.png       : 2D synthetic data with learned
9                                     decision boundary
10     - probability_contours.png     : Predicted probability contours over
11                                     a grid
12     - confusion_matrix.png         : Confusion matrix heatmap on a held-
13                                     out split
14
15 Dependencies:
16     - numpy, matplotlib
17
18 Notes on compatibility:
19     - Avoids optional or newer Matplotlib parameters; uses standard
20       pyplot API.
21     - Implements a simple Logistic Regression via gradient descent to
22       avoid external deps.
23
24 Usage:
25     python gen_logistic_regression_figures.py

```

```
20 """
21
22 from __future__ import annotations
23
24 import os
25 import numpy as np
26 import matplotlib.pyplot as plt
27
28
29 def sigmoid(z: np.ndarray) -> np.ndarray:
30     """Numerically stable sigmoid."""
31     # For large negative z, exp(-z) can overflow; using np.clip is a
32     # simple safeguard.
33     z = np.clip(z, -50, 50)
34     return 1.0 / (1.0 + np.exp(-z))
35
36 def binary_cross_entropy(z: np.ndarray, y: np.ndarray) -> np.ndarray:
37     """Per-sample logistic loss as a function of logit z and label y in
38     {0,1}."""
39     p = sigmoid(z)
40     # Clip for numerical stability in log
41     eps = 1e-12
42     p = np.clip(p, eps, 1.0 - eps)
43     return -(y * np.log(p) + (1 - y) * np.log(1 - p))
44
45 def make_gaussian_2class(n_per_class: int = 200, seed: int = 42):
46     """Generate a linearly separable-ish 2D dataset of two Gaussian
47     blobs."""
48     rng = np.random.RandomState(seed)
49
50     mean0 = np.array([-1.0, -1.0])
51     mean1 = np.array([+1.2, +1.2])
52     cov = np.array([[0.6, 0.2], [0.2, 0.6]])
53
54     X0 = rng.multivariate_normal(mean0, cov, size=n_per_class)
55     X1 = rng.multivariate_normal(mean1, cov, size=n_per_class)
56     y0 = np.zeros(n_per_class, dtype=int)
57     y1 = np.ones(n_per_class, dtype=int)
58
59     X = np.vstack([X0, X1])
60     y = np.concatenate([y0, y1])
```

```
60
61     # Shuffle
62     idx = rng.permutation(X.shape[0])
63     X, y = X[idx], y[idx]
64     return X, y
65
66
67 def train_logreg_gd(X: np.ndarray, y: np.ndarray, lr: float = 0.1,
68     n_iter: int = 1000, reg_l2: float = 0.0, seed: int = 42):
69     """Train a simple logistic regression via batch gradient descent.
70
71     Parameters
72     -----
73     X : (n_samples, n_features)
74     y : (n_samples,) in {0,1}
75     lr : learning rate
76     n_iter : number of iterations
77     reg_l2 : L2 regularization strength (applied to weights, not bias)
78     seed : random seed for initialization
79
80     Returns
81     -----
82     w0 : bias (float)
83     w : weights (n_features,)
84     history : dict with loss per iteration
85     """
86     rng = np.random.RandomState(seed)
87     n, d = X.shape
88
89     # Initialize small random weights for symmetry breaking
90     w = rng.normal(scale=0.01, size=d)
91     w0 = 0.0
92
93     hist_loss = []
94     for _ in range(n_iter):
95         z = w0 + X.dot(w)
96         p = sigmoid(z)
97         # Gradients
98         err = (p - y)
99         grad_w = X.T.dot(err) / n + reg_l2 * w
100         grad_b = err.mean()
101         # Update
102         w -= lr * grad_w
```

```
102     w0 -= lr * grad_b
103
104     # Track loss
105     loss = binary_cross_entropy(z, y).mean() + 0.5 * reg_l2 * np.
        dot(w, w)
106     hist_loss.append(loss)
107
108     return w0, w, {"loss": np.array(hist_loss)}
109
110
111 def plot_sigmoid(out_path: str):
112     t = np.linspace(-10, 10, 500)
113     s = sigmoid(t)
114     plt.figure(figsize=(6, 4))
115     plt.plot(t, s, color="tab:blue", lw=2)
116     plt.axhline(0.5, color="gray", lw=1, ls="--")
117     plt.axvline(0.0, color="gray", lw=1, ls="--")
118     plt.title("Sigmoid Function")
119     plt.xlabel("t")
120     plt.ylabel("sigma(t)")
121     plt.grid(alpha=0.3)
122     plt.tight_layout()
123     plt.savefig(out_path, dpi=300, bbox_inches="tight")
124     plt.close()
125
126
127 def plot_logistic_losses(out_path: str):
128     z = np.linspace(-10, 10, 500)
129     loss_y1 = binary_cross_entropy(z, np.ones_like(z))
130     loss_y0 = binary_cross_entropy(z, np.zeros_like(z))
131
132     plt.figure(figsize=(6.5, 4.2))
133     plt.plot(z, loss_y1, label="y=1", color="tab:blue", lw=2)
134     plt.plot(z, loss_y0, label="y=0", color="tab:orange", lw=2)
135     plt.title("Logistic Loss vs Logit z")
136     plt.xlabel("z")
137     plt.ylabel("Per-sample loss")
138     plt.legend(frameon=False)
139     plt.grid(alpha=0.3)
140     plt.tight_layout()
141     plt.savefig(out_path, dpi=300, bbox_inches="tight")
142     plt.close()
143
```

```

144
145 def plot_decision_boundary_and_data(X: np.ndarray, y: np.ndarray, w0:
    float, w: np.ndarray, out_path: str):
146     plt.figure(figsize=(6.8, 5.2))
147
148     # Scatter points
149     m0 = y == 0
150     m1 = y == 1
151     plt.scatter(X[m0, 0], X[m0, 1], s=20, c="tab:orange", alpha=0.8,
        label="Class 0")
152     plt.scatter(X[m1, 0], X[m1, 1], s=20, c="tab:blue", alpha=0.8,
        label="Class 1")
153
154     # Decision boundary  $w_0 + w_1 x + w_2 y = 0$ 
155     if abs(w[1]) > 1e-12:
156         xs = np.linspace(X[:, 0].min() - 0.5, X[:, 0].max() + 0.5, 200)
157         ys = -(w0 + w[0] * xs) / w[1]
158         plt.plot(xs, ys, color="k", lw=2, label="Decision boundary (z
            =0)")
159     else:
160         # Vertical boundary
161         x_b = -w0 / (w[0] + 1e-12)
162         plt.axvline(x_b, color="k", lw=2, label="Decision boundary (z
            =0)")
163
164     plt.title("Logistic Regression Decision Boundary")
165     plt.xlabel("x1")
166     plt.ylabel("x2")
167     plt.legend(frameon=False)
168     plt.grid(alpha=0.25)
169     plt.tight_layout()
170     plt.savefig(out_path, dpi=300, bbox_inches="tight")
171     plt.close()
172
173
174 def plot_probability_contours(X: np.ndarray, w0: float, w: np.ndarray,
    out_path: str):
175     # Grid covering the data extent
176     x_min, x_max = X[:, 0].min() - 0.8, X[:, 0].max() + 0.8
177     y_min, y_max = X[:, 1].min() - 0.8, X[:, 1].max() + 0.8
178     xx, yy = np.meshgrid(
179         np.linspace(x_min, x_max, 200),
180         np.linspace(y_min, y_max, 200),

```

```

181     )
182     grid = np.c_[xx.ravel(), yy.ravel()]
183     z = w0 + grid.dot(w)
184     p = sigmoid(z).reshape(xx.shape)
185
186     plt.figure(figsize=(6.8, 5.2))
187     cs = plt.contourf(xx, yy, p, levels=21, cmap="RdBu_r", alpha=0.8)
188     cbar = plt.colorbar(cs)
189     cbar.set_label("p(y=1|x)")
190     # Decision contour at p=0.5 (z=0)
191     plt.contour(xx, yy, p, levels=[0.5], colors=["k"], linewidths=2)
192     plt.title("Predicted Probability Contours")
193     plt.xlabel("x1")
194     plt.ylabel("x2")
195     plt.tight_layout()
196     plt.savefig(out_path, dpi=300, bbox_inches="tight")
197     plt.close()
198
199
200 def plot_confusion_matrix(y_true: np.ndarray, y_prob: np.ndarray,
201     threshold: float, out_path: str):
202     y_pred = (y_prob >= threshold).astype(int)
203     # Compute confusion matrix counts
204     tp = int(((y_true == 1) & (y_pred == 1)).sum())
205     tn = int(((y_true == 0) & (y_pred == 0)).sum())
206     fp = int(((y_true == 0) & (y_pred == 1)).sum())
207     fn = int(((y_true == 1) & (y_pred == 0)).sum())
208     cm = np.array([[tn, fp], [fn, tp]], dtype=float)
209
210     plt.figure(figsize=(4.8, 4.2))
211     im = plt.imshow(cm, interpolation="nearest", cmap="Blues")
212     plt.title("Confusion Matrix (thr=%.2f)" % threshold)
213     plt.colorbar(im, fraction=0.046, pad=0.04)
214     tick_marks = np.arange(2)
215     plt.xticks(tick_marks, ["Pred 0", "Pred 1"])
216     plt.yticks(tick_marks, ["True 0", "True 1"])
217
218     # Annotate counts
219     for i in range(2):
220         for j in range(2):
221             plt.text(j, i, "%d" % cm[i, j], ha="center", va="center",

```

```
222     plt.tight_layout()
223     plt.ylabel("True label")
224     plt.xlabel("Predicted label")
225     plt.savefig(out_path, dpi=300, bbox_inches="tight")
226     plt.close()
227
228
229 def main():
230     # Ensure output directory exists
231     out_dir = os.path.join(os.path.dirname(__file__), "figures")
232     if not os.path.isdir(out_dir):
233         os.makedirs(out_dir)
234
235     # 1) Sigmoid curve
236     plot_sigmoid(os.path.join(out_dir, "sigmoid_curve.png"))
237
238     # 2) Logistic losses
239     plot_logistic_losses(os.path.join(out_dir, "logistic_loss_curves.
240                             png"))
241
242     # 3) Synthetic data + split
243     X, y = make_gaussian_2class(n_per_class=250, seed=42)
244     # Simple train/test split
245     n = X.shape[0]
246     split = int(0.7 * n)
247     X_train, y_train = X[:split], y[:split]
248     X_test, y_test = X[split:], y[split:]
249
250     # 4) Train simple logistic regression
251     w0, w, hist = train_logreg_gd(X_train, y_train, lr=0.15, n_iter
252                                     =800, reg_l2=0.01, seed=42)
253
254     # 5) Decision boundary with training data
255     plot_decision_boundary_and_data(X_train, y_train, w0, w, os.path.
256                                     join(out_dir, "decision_boundary.png"))
257
258     # 6) Probability contours over full extent
259     plot_probability_contours(X, w0, w, os.path.join(out_dir, "
260                                     probability_contours.png"))
261
262     # 7) Confusion matrix on test set
263     z_test = w0 + X_test.dot(w)
264     p_test = sigmoid(z_test)
```



```

261     plot_confusion_matrix(y_test, p_test, threshold=0.5, out_path=os.
262         path.join(out_dir, "confusion_matrix.png"))
263
264     print("Figures written to:", out_dir)
265
266 if __name__ == "__main__":
267     main()

```

## 5 运行效果

核心插图如下所示。

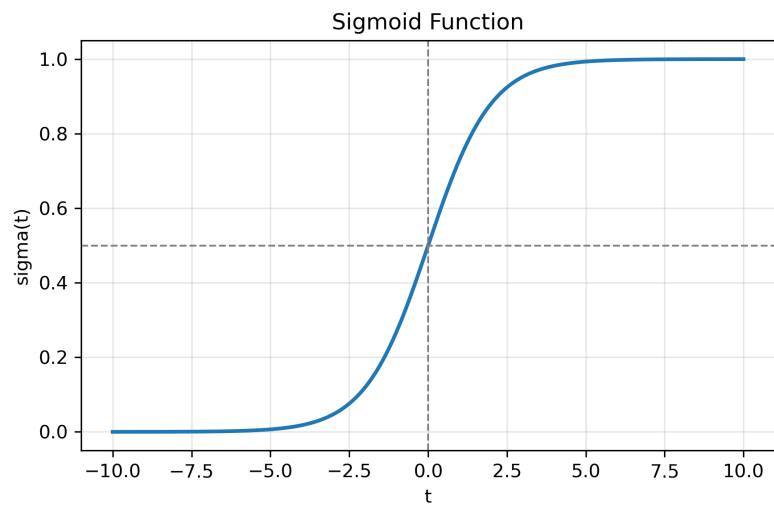


图 1: S 形函数  $\sigma(t) = 1/(1 + e^{-t})$

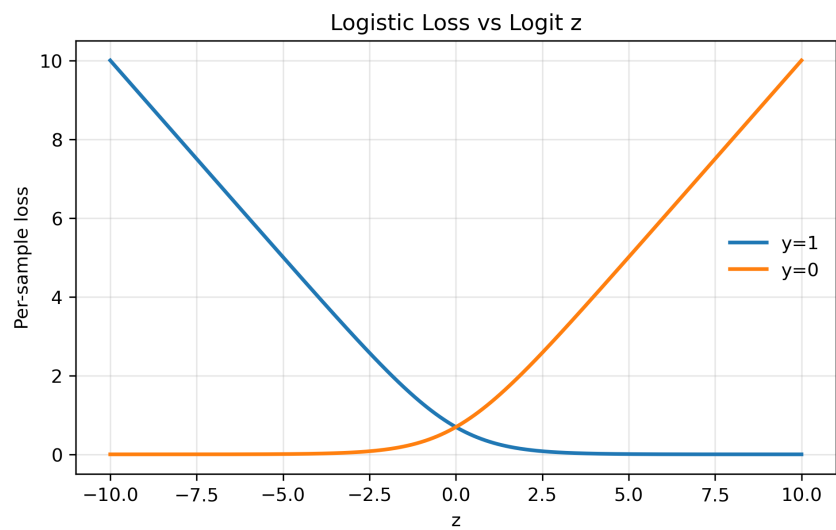


图 2: 单样本二元交叉熵随对数几率  $z$  的变化曲线 ( $y=0$  与  $y=1$ )

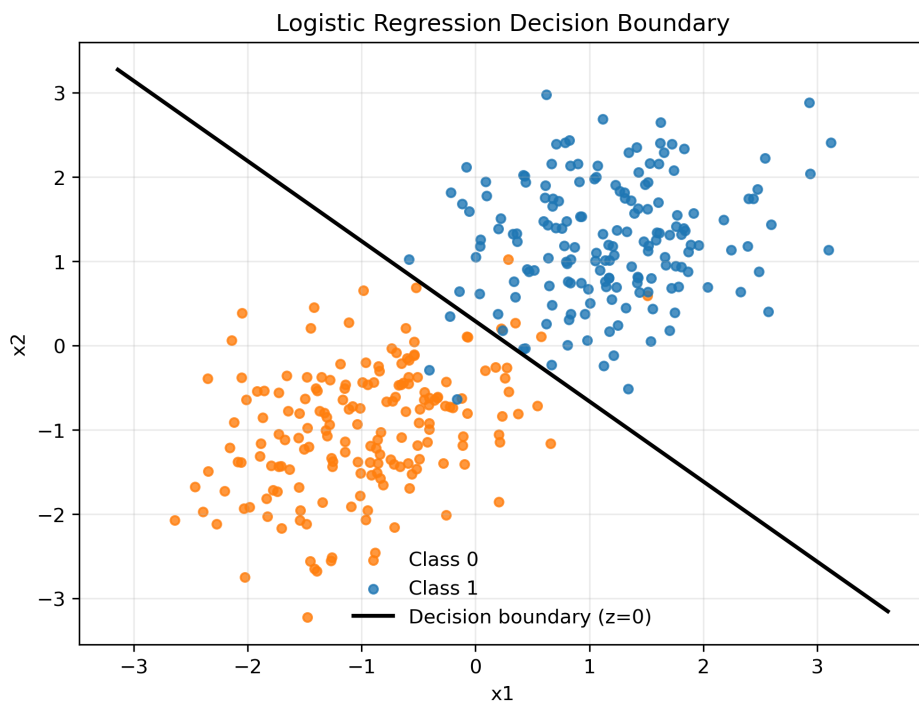
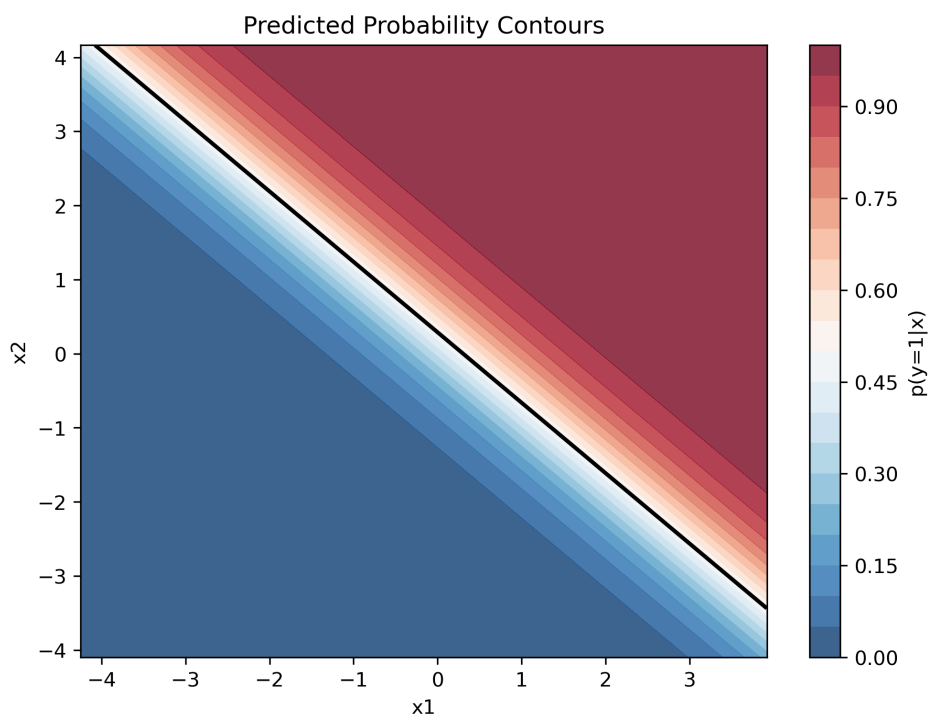


图 3: 二维合成数据与学习到的逻辑回归判别边界

图 4: 网格上的预测概率等高线  $p(y=1|\mathbf{x})$

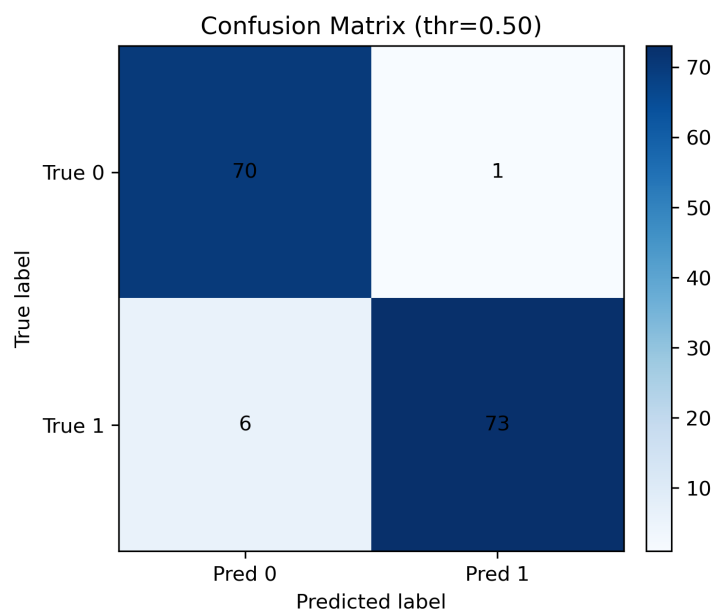


图 5: 在留出集上（阈值 0.5）的混淆矩阵

## 6 小结

逻辑回归以简洁、可解释和高效著称，是分类任务中实用的基线模型。其概率化建模、凸优化训练和线性判别边界使之在工程与研究中长期占据重要地位。