

# 朴素贝叶斯：理论与实践

2025 年 9 月 9 日

## 1 引言

朴素贝叶斯（Naïve Bayes, NB）是在条件独立假设下建立的概率分类器族：

$$p(y \mid \mathbf{x}) \propto p(y) \prod_{j=1}^d p(x_j \mid y), \quad (1)$$

其中  $y$  为类别， $\mathbf{x} = (x_1, \dots, x_d)$  为特征。尽管独立性假设较强，NB 在高维稀疏特征（如文本）等任务中常表现稳健，且训练、预测效率较高。

## 2 原理与公式

以高斯朴素贝叶斯（Gaussian NB）为例，对于连续特征，假设对每个类别  $c \in \{1, \dots, C\}$  与每个特征  $j$  有：

$$x_j \mid y = c \sim \mathcal{N}(\mu_{c,j}, \sigma_{c,j}^2). \quad (2)$$

则条件似然分解为  $p(\mathbf{x} \mid y = c) = \prod_j \mathcal{N}(x_j; \mu_{c,j}, \sigma_{c,j}^2)$ 。结合先验  $p(y = c)$ ，（未归一化的）对数后验为：

$$\log p(y = c \mid \mathbf{x}) \propto \log p(y = c) + \sum_{j=1}^d \log \mathcal{N}(x_j; \mu_{c,j}, \sigma_{c,j}^2) \quad (3)$$

$$\propto \log p(y = c) - \sum_{j=1}^d \left[ \frac{1}{2} \log(2\pi\sigma_{c,j}^2) + \frac{(x_j - \mu_{c,j})^2}{2\sigma_{c,j}^2} \right]. \quad (4)$$

预测类别为  $\hat{y} = \arg \max_c \log p(y = c \mid \mathbf{x})$ 。参数估计可由各类别内样本的均值与方差直接得到。

**备注** 变体包括：连续特征的 Gaussian NB；计数/二值特征的 Multinomial/Bernoulli NB（常配合拉普拉斯平滑）。若需使用概率值做后续决策，建议做概率校准。

## 3 应用场景与要点

- **适用场景：**高维稀疏文本（BOW/TF-IDF）、简单传感器数据、作为强基线。
- **预处理：**Gaussian NB 建议对连续特征做标准化；文本常用计数或 TF-IDF（Multinomial NB）。
- **类别先验：**可用经验频率或领域知识设定。
- **独立性假设：**特征强相关时性能可能下降；建议与逻辑回归/线性 SVM 等对比。
- **评估：**采用交叉验证对比不同模型与超参数。

## 4 Python 实战

在章节目录内运行下述脚本，图片将保存到本目录下的 `figures/`：

Listing 1: 生成朴素贝叶斯配图

```
1 # 在 4_Naive Bayes 目录中执行：  
2 python gen_naive_bayes_figures.py
```

Listing 2: `gen_naive_bayes_figures.py` 源码

```
1 """  
2 Figure generator for the Naive Bayes chapter.  
3  
4 Generates illustrative figures and saves them into the local 'figures/'  
   folder.  
5  
6 Requirements:  
7 - Python 3.8+  
8 - numpy, matplotlib, scikit-learn  
9  
10 Install (if needed):  
11     pip install numpy matplotlib scikit-learn  
12  
13 This script avoids newer or experimental APIs to stay compatible with  
   older  
14 versions of the dependencies.  
15 """  
16 from __future__ import annotations  
17  
18 import os  
19 import math
```

```
20 import numpy as np
21 import matplotlib.pyplot as plt
22 from matplotlib.colors import ListedColormap
23
24 try:
25     from sklearn.datasets import make_blobs
26     from sklearn.naive_bayes import GaussianNB
27     from sklearn.linear_model import LogisticRegression
28     from sklearn.preprocessing import StandardScaler
29 except Exception as e:
30     raise SystemExit(
31         "Missing scikit-learn dependency. Please install with: pip
           install scikit-learn"
32     )
33
34
35 def _ensure_figures_dir(path: str | None = None) -> str:
36     """Create figures directory under this chapter regardless of CWD.
37
38     If `path` is None, resolve to `<this_file_dir>/figures`.
39     """
40     if path is None:
41         base = os.path.dirname(os.path.abspath(__file__))
42         path = os.path.join(base, "figures")
43     os.makedirs(path, exist_ok=True)
44     return path
45
46
47 def _plot_decision_boundary(ax, clf, X, y, title: str, cmap_light,
                             cmap_bold):
48     # Create a mesh grid for decision surface
49     x_min, x_max = X[:, 0].min() - 1.0, X[:, 0].max() + 1.0
50     y_min, y_max = X[:, 1].min() - 1.0, X[:, 1].max() + 1.0
51     xx, yy = np.meshgrid(
52         np.linspace(x_min, x_max, 300), np.linspace(y_min, y_max, 300)
53     )
54     Z = clf.predict(np.c_[xx.ravel(), yy.ravel()])
55     Z = Z.reshape(xx.shape)
56
57     ax.contourf(xx, yy, Z, cmap=cmap_light, alpha=0.8, levels=np.unique
68                 (Z).size)
59
60     # Training points
61     scatter = ax.scatter(X[:, 0], X[:, 1], c=y, cmap=cmap_bold,
```

```

        edgecolors="k", s=25)
60 ax.set_xlabel("Feature 1")
61 ax.set_ylabel("Feature 2")
62 ax.set_title(title)
63 return scatter
64
65
66 def fig_gnb_decision_boundary_2class(out_dir: str) -> str:
67     np.random.seed(42)
68     X, y = make_blobs(n_samples=400, centers=2, cluster_std=[1.2, 1.2],
69                       random_state=42)
69
70     clf = GaussianNB()
71     clf.fit(X, y)
72
73     cmap_light = ListedColormap(["#FFEEEE", "#EEEEFF"])
74     cmap_bold = ListedColormap(["#E74C3C", "#3498DB"])
75
76     fig, ax = plt.subplots(figsize=(5.5, 4.5), dpi=150)
77     _plot_decision_boundary(ax, clf, X, y, "Gaussian Naive Bayes (2-
78                             class)", cmap_light, cmap_bold)
79     out_path = os.path.join(out_dir, "gnb_decision_boundary_2class.png"
80                             )
81     fig.tight_layout()
82     fig.savefig(out_path)
83     plt.close(fig)
84     return out_path
85
86 def fig_gnb_decision_boundary_3class(out_dir: str) -> str:
87     np.random.seed(7)
88     X, y = make_blobs(
89         n_samples=600,
90         centers=3,
91         cluster_std=[1.1, 1.0, 1.2],
92         random_state=7,
93     )
94
95     clf = GaussianNB()
96     clf.fit(X, y)
97
98     cmap_light = ListedColormap(["#FFEEEE", "#EEFFEE", "#EEEEFF"])
99     cmap_bold = ListedColormap(["#E74C3C", "#2ECC71", "#3498DB"])

```

```

99
100 fig, ax = plt.subplots(figsize=(5.5, 4.5), dpi=150)
101 _plot_decision_boundary(ax, clf, X, y, "Gaussian Naive Bayes (3-
      cmap_light, cmap_bold)
102 out_path = os.path.join(out_dir, "gnb_decision_boundary_3class.png"
      )
103 fig.tight_layout()
104 fig.savefig(out_path)
105 plt.close(fig)
106 return out_path
107
108
109 def _gaussian_pdf(x: np.ndarray, mu: float, sigma: float) -> np.ndarray
      :
110     coef = 1.0 / (math.sqrt(2.0 * math.pi) * sigma)
111     return coef * np.exp(-0.5 * ((x - mu) / sigma) ** 2)
112
113
114 def fig_class_conditional_densities_1d(out_dir: str) -> str:
115     # Two 1D Gaussians with equal priors
116     mu0, sigma0 = -1.0, 1.0
117     mu1, sigma1 = 1.2, 0.8
118     xs = np.linspace(-5, 5, 500)
119     p_x_c0 = _gaussian_pdf(xs, mu0, sigma0)
120     p_x_c1 = _gaussian_pdf(xs, mu1, sigma1)
121
122     # Decision threshold where  $p(x|c_0) = p(x|c_1)$ 
123     # For illustration, compute numerically
124     idx = np.argmin(np.abs(p_x_c0 - p_x_c1))
125     x_star = xs[idx]
126
127     fig, ax = plt.subplots(figsize=(6, 4), dpi=150)
128     ax.plot(xs, p_x_c0, label="p(x|class 0)", color="#E74C3C", lw=2)
129     ax.plot(xs, p_x_c1, label="p(x|class 1)", color="#3498DB", lw=2)
130     ax.axvline(x_star, color="#7F8C8D", ls="--", lw=1)
131     ax.text(x_star + 0.1, max(p_x_c0[idx], p_x_c1[idx]) * 0.9, "
          decision", color="#7F8C8D")
132     ax.set_xlabel("x")
133     ax.set_ylabel("density")
134     ax.set_title("Class-conditional densities (1D)")
135     ax.legend(frameon=False)
136     out_path = os.path.join(out_dir, "class_conditional_densities_1d.
          png")

```

```
137     fig.tight_layout()
138     fig.savefig(out_path)
139     plt.close(fig)
140     return out_path
141
142
143 def fig_feature_independence_heatmap(out_dir: str) -> str:
144     # Create 3 correlated features to illustrate independence
145     # assumption violation
146     np.random.seed(123)
147     mean = np.array([0.0, 0.0, 0.0])
148     cov = np.array(
149         [
150             [1.0, 0.7, 0.4],
151             [0.7, 1.0, 0.5],
152             [0.4, 0.5, 1.0],
153         ]
154     )
155     X = np.random.multivariate_normal(mean, cov, size=1000)
156     # Empirical correlation matrix
157     C = np.corrcoef(X, rowvar=False)
158
159     fig, ax = plt.subplots(figsize=(4.8, 4.2), dpi=160)
160     im = ax.imshow(C, cmap="coolwarm", vmin=-1, vmax=1)
161     for i in range(C.shape[0]):
162         for j in range(C.shape[1]):
163             ax.text(j, i, f"{C[i, j]:.2f}", ha="center", va="center",
164                     color="black")
165     ax.set_xticks([0, 1, 2])
166     ax.set_yticks([0, 1, 2])
167     ax.set_xticklabels(["f1", "f2", "f3"])
168     ax.set_yticklabels(["f1", "f2", "f3"])
169     ax.set_title("Feature correlation (independence assumption)")
170     fig.colorbar(im, ax=ax, fraction=0.046, pad=0.04, label="
171         correlation")
172     out_path = os.path.join(out_dir, "feature_independence_heatmap.png")
173
174     fig.tight_layout()
175     fig.savefig(out_path)
176     plt.close(fig)
177     return out_path
```

```

176 def fig_gnb_vs_logreg_boundary(out_dir: str) -> str:
177     # Dataset with partially overlapping Gaussians
178     np.random.seed(0)
179     X, y = make_blobs(n_samples=500, centers=[(-2, -2), (2.5, 2.0)],
180                       cluster_std=[1.6, 1.2], random_state=0)
181
182     scaler = StandardScaler()
183     Xs = scaler.fit_transform(X)
184
185     gnb = GaussianNB().fit(Xs, y)
186     # Use lbfgs which supports multinomial/binary and is widely
187     # available
188     lr = LogisticRegression(solver="lbfgs", max_iter=1000).fit(Xs, y)
189
190     x_min, x_max = Xs[:, 0].min() - 2.0, Xs[:, 0].max() + 2.0
191     y_min, y_max = Xs[:, 1].min() - 2.0, Xs[:, 1].max() + 2.0
192     xx, yy = np.meshgrid(np.linspace(x_min, x_max, 300), np.linspace(
193         y_min, y_max, 300))
194     grid = np.c_[xx.ravel(), yy.ravel()]
195     Z_gnb = gnb.predict(grid).reshape(xx.shape)
196     Z_lr = lr.predict(grid).reshape(xx.shape)
197
198     fig, axes = plt.subplots(1, 2, figsize=(9.5, 4.2), dpi=150, sharex=
199         True, sharey=True)
200     for ax, Z, title in [
201         (axes[0], Z_gnb, "Gaussian NB boundary"),
202         (axes[1], Z_lr, "Logistic Regression boundary"),
203     ]:
204         ax.contourf(xx, yy, Z, alpha=0.25, levels=np.unique(y).size,
205                    cmap=ListedColormap(["#FFBBBB", "#BBBBFF"]))
206         ax.scatter(Xs[:, 0], Xs[:, 1], c=y, s=15, cmap=ListedColormap([
207             "#E74C3C", "#3498DB"]), edgecolors="k")
208         ax.set_title(title)
209         ax.set_xlabel("feature 1 (scaled)")
210         ax.set_ylabel("feature 2 (scaled)")
211     fig.suptitle("Naive Bayes vs Logistic Regression")
212     out_path = os.path.join(out_dir, "gnb_vs_logreg_boundary.png")
213     fig.tight_layout(rect=[0, 0.03, 1, 0.95])
214     fig.savefig(out_path)
215     plt.close(fig)
216     return out_path

```

```
213 def main():
214     # Always save figures inside the current chapter directory
215     out_dir = _ensure_figures_dir(None)
216     generators = [
217         fig_gnb_decision_boundary_2class,
218         fig_gnb_decision_boundary_3class,
219         fig_class_conditional_densities_1d,
220         fig_feature_independence_heatmap,
221         fig_gnb_vs_logreg_boundary,
222     ]
223
224     print("Generating figures into:", os.path.abspath(out_dir))
225     for gen in generators:
226         try:
227             path = gen(out_dir)
228             print("Saved:", path)
229         except Exception as e:
230             print("Failed generating", gen.__name__, ":", e)
231
232
233 if __name__ == "__main__":
234     main()
```



## 5 结果

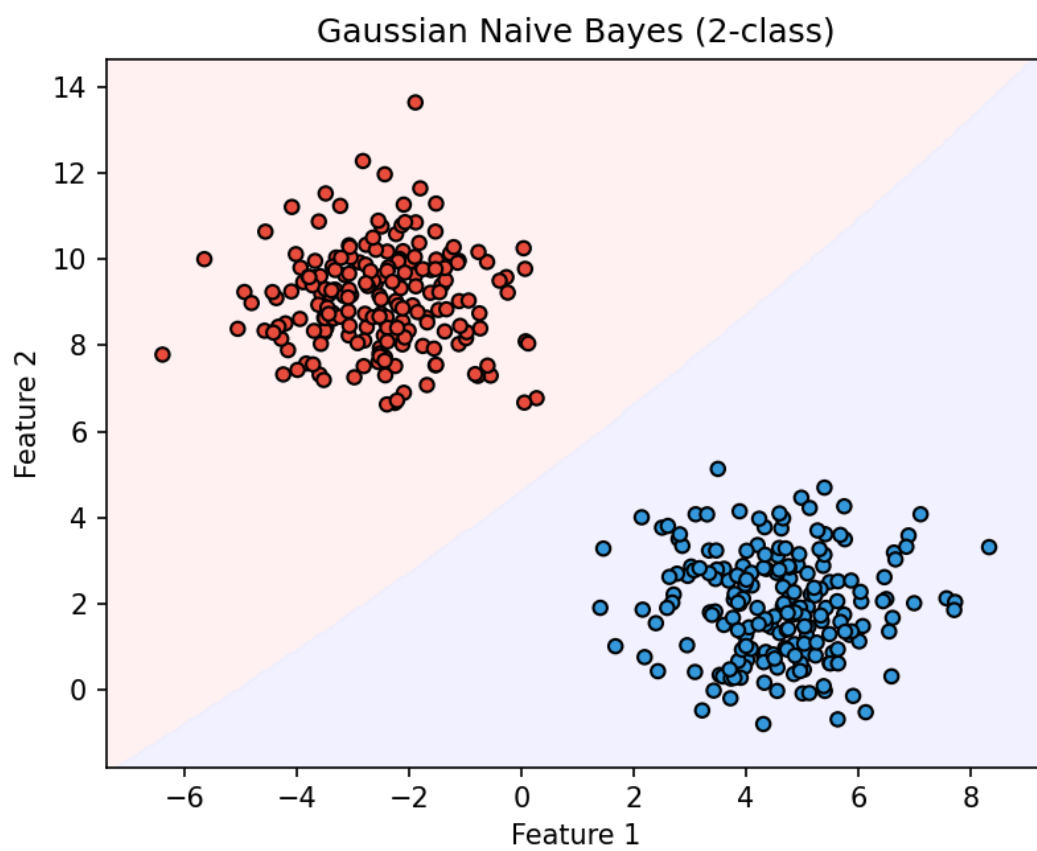


图 1: Gaussian NB 分类边界 (两类)。

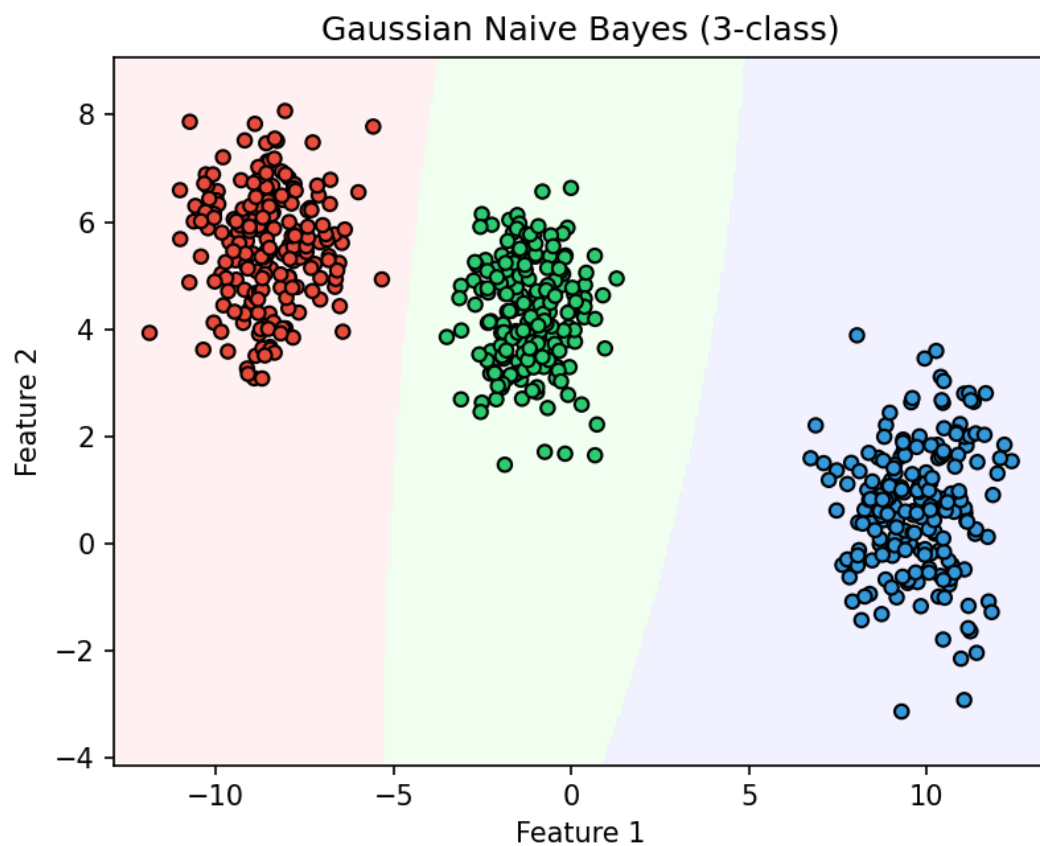


图 2: Gaussian NB 决策区域 (三类)。

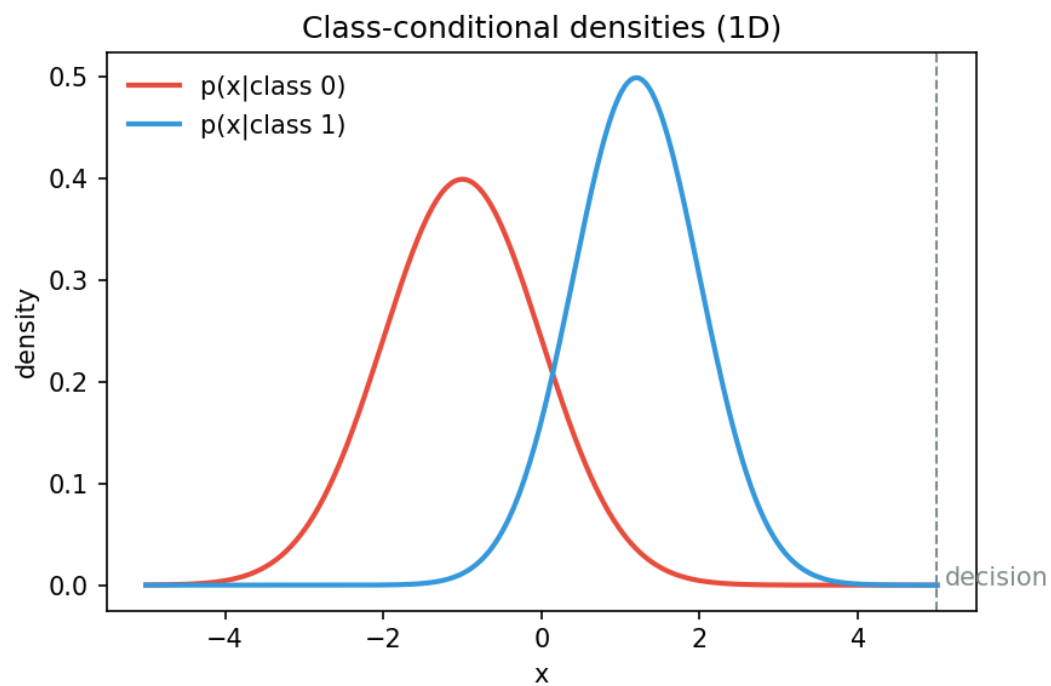


图 3: 一维类别条件密度与决策阈值。

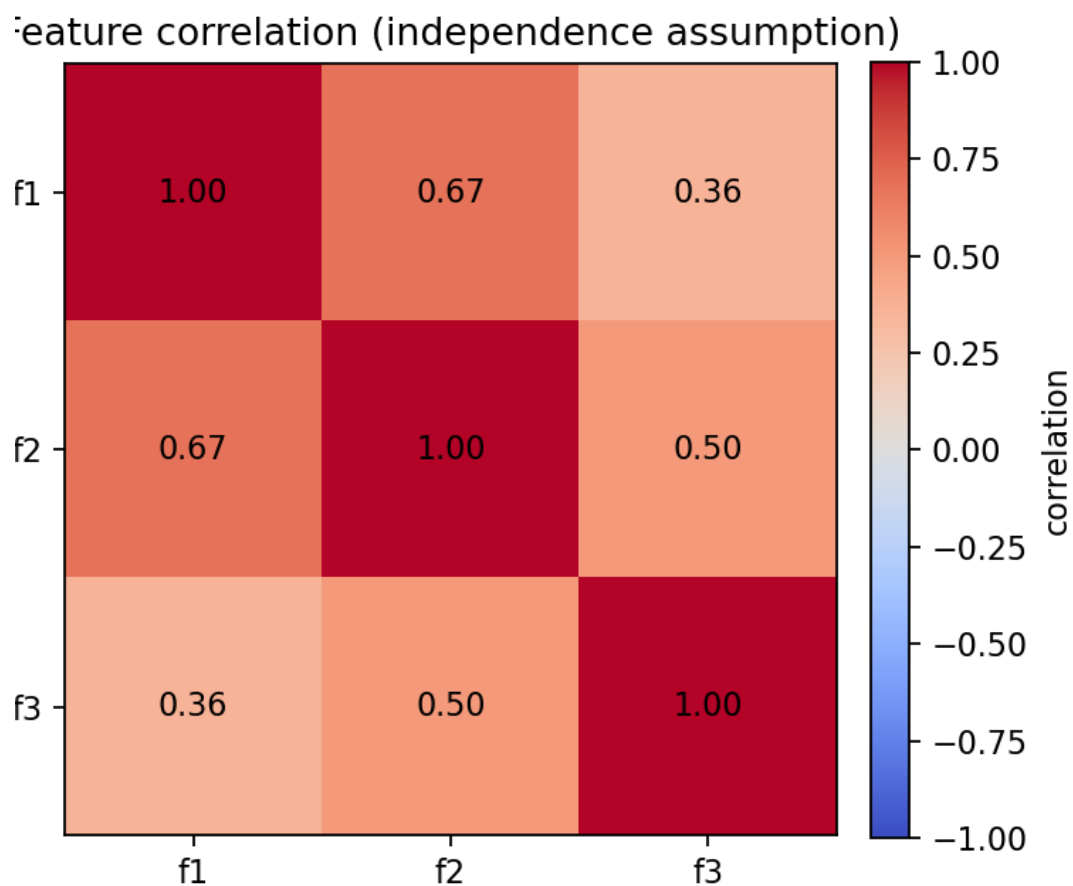


图 4: 特征相关性热力图（独立性假设示意）。

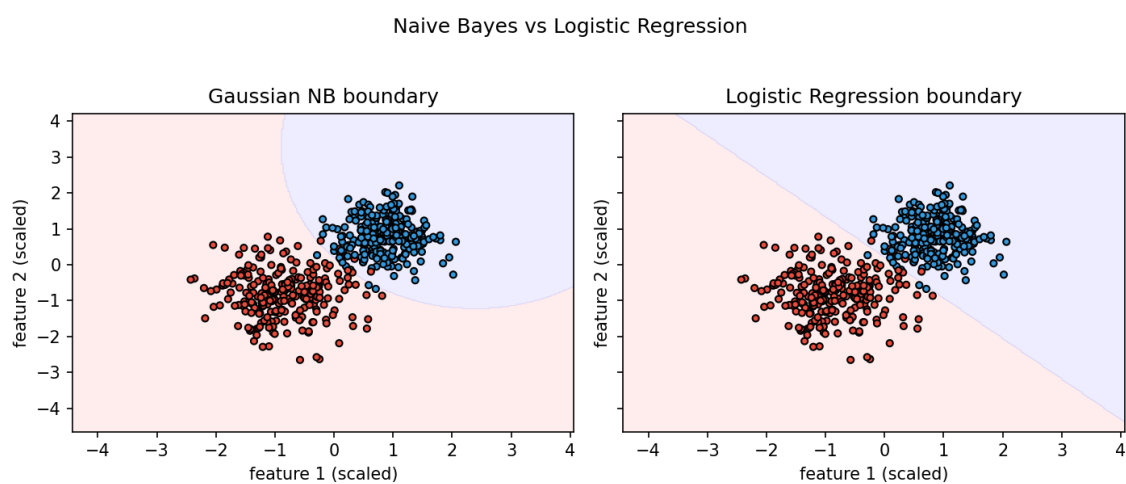


图 5: Gaussian NB 与逻辑回归的决策边界对比。

## 6 总结

朴素贝叶斯以简洁可解释、训练与预测高效为特点：核心是先验与逐特征似然的乘积（条件独立）。虽然假设并非总成立，它依然是可靠的基线模型，常用于与更强的判别式模型进行对比。