

Naïve Bayes: Theory and Practice

September 9, 2025

1 Introduction

Naïve Bayes (NB) is a family of simple yet effective probabilistic classifiers. Under the *conditional independence* assumption, the posterior is

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

where y is the class and $\mathbf{x} = (x_1, \dots, x_d)$ are features. Despite the strong assumption, NB can work surprisingly well in many domains, especially with high-dimensional sparse inputs.

2 Theory and Formulas

Consider the Gaussian NB model for continuous features. For class $c \in \{1, \dots, C\}$, assume

$$x_j \mid y = c \sim \mathcal{N}(\mu_{c,j}, \sigma_{c,j}^2), \quad j = 1, \dots, d. \quad (2)$$

Then the class-conditional density factorizes: $p(\mathbf{x} \mid y = c) = \prod_j \mathcal{N}(x_j; \mu_{c,j}, \sigma_{c,j}^2)$. With prior $p(y = c)$, the (unnormalized) log-posterior is

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

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

The predicted class is $\hat{y} = \arg \max_c \log p(y = c \mid \mathbf{x})$. Estimation is straightforward via sample means and variances within each class.

Notes. NB variants include Gaussian NB for continuous features and Multinomial/Bernoulli NB for count/binary features with Laplace (additive) smoothing. Calibration may be needed if probabilities are used downstream.

3 Applications and Tips

- **When it works:** high-dimensional sparse text features (bag-of-words), simple sensor data, baseline models.

- **Preprocessing:** standardize continuous features for Gaussian NB; for text, TF-IDF or raw counts for Multinomial NB.
- **Class priors:** either empirical (class frequencies) or domain-informed.
- **Independence assumption:** correlations between features may harm performance; use as a baseline and compare.
- **Evaluation:** compare with logistic regression/SVMs; use cross-validation.

4 Python Practice

The script below generates figures for Gaussian NB decision boundaries and simple diagnostics. Run it in the chapter folder; it saves images under **figures/**.

Listing 1: Generate Naive Bayes figures

```
1 # Terminal
2 python gen_naive_bayes_figures.py
```

Listing 2: gen_naive_bayes_figures.py

```
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(Z).
58                 size)
59     # Training points
60     scatter = ax.scatter(X[:, 0], X[:, 1], c=y, cmap=cmap_bold, edgecolors="k",
61                          s=25)
62     ax.set_xlabel("Feature 1")
63     ax.set_ylabel("Feature 2")
64     ax.set_title(title)
65     return scatter
66
67 def fig_gnb_decision_boundary_2class(out_dir: str) -> str:
68     np.random.seed(42)
69     X, y = make_blobs(n_samples=400, centers=2, cluster_std=[1.2, 1.2],
70                      random_state=42)
71
72     clf = GaussianNB()
73     clf.fit(X, y)
74
75     cmap_light = ListedColormap(["#FFEEEE", "#EEEEFF"])
76     cmap_bold = ListedColormap(["#E74C3C", "#3498DB"])
77
78     fig, ax = plt.subplots(figsize=(5.5, 4.5), dpi=150)
79     _plot_decision_boundary(ax, clf, X, y, "Gaussian Naive Bayes (2-class)",
80                            cmap_light, cmap_bold)
81     out_path = os.path.join(out_dir, "gnb_decision_boundary_2class.png")
82     fig.tight_layout()

```

```

80     fig.savefig(out_path)
81     plt.close(fig)
82     return out_path
83
84
85 def fig_gnb_decision_boundary_3class(out_dir: str) -> str:
86     np.random.seed(7)
87     X, y = make_blobs(
88         n_samples=600,
89         centers=3,
90         cluster_std=[1.1, 1.0, 1.2],
91         random_state=7,
92     )
93
94     clf = GaussianNB()
95     clf.fit(X, y)
96
97     cmap_light = ListedColormap(["#FFEEEE", "#EEFFEE", "#EEEEFF"])
98     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-class)",
102        cmap_light, cmap_bold)
103    out_path = os.path.join(out_dir, "gnb_decision_boundary_3class.png")
104    fig.tight_layout()
105    fig.savefig(out_path)
106    plt.close(fig)
107    return out_path
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",
132         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 assumption
145     # 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", color="black")
164
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="correlation")
171     out_path = os.path.join(out_dir, "feature_independence_heatmap.png")
172     fig.tight_layout()
173     fig.savefig(out_path)
174     plt.close(fig)
175     return out_path
176
177
178 def fig_gnb_vs_logreg_boundary(out_dir: str) -> str:
179     # Dataset with partially overlapping Gaussians
180     np.random.seed(0)
181     X, y = make_blobs(n_samples=500, centers=[(-2, -2), (2.5, 2.0)],
182                       cluster_std=[1.6, 1.2], random_state=0)
183
184     scaler = StandardScaler()
185     Xs = scaler.fit_transform(X)

```

```

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

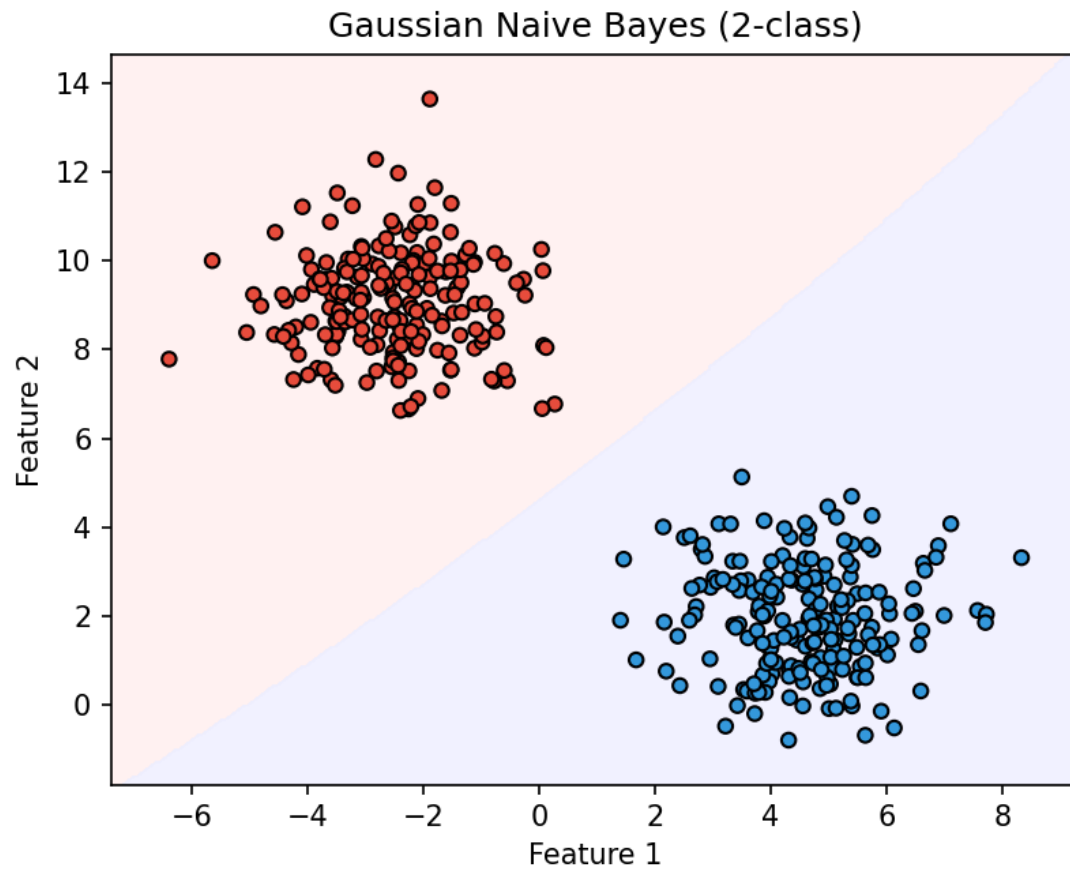


Figure 1: Gaussian NB decision boundary (2-class).

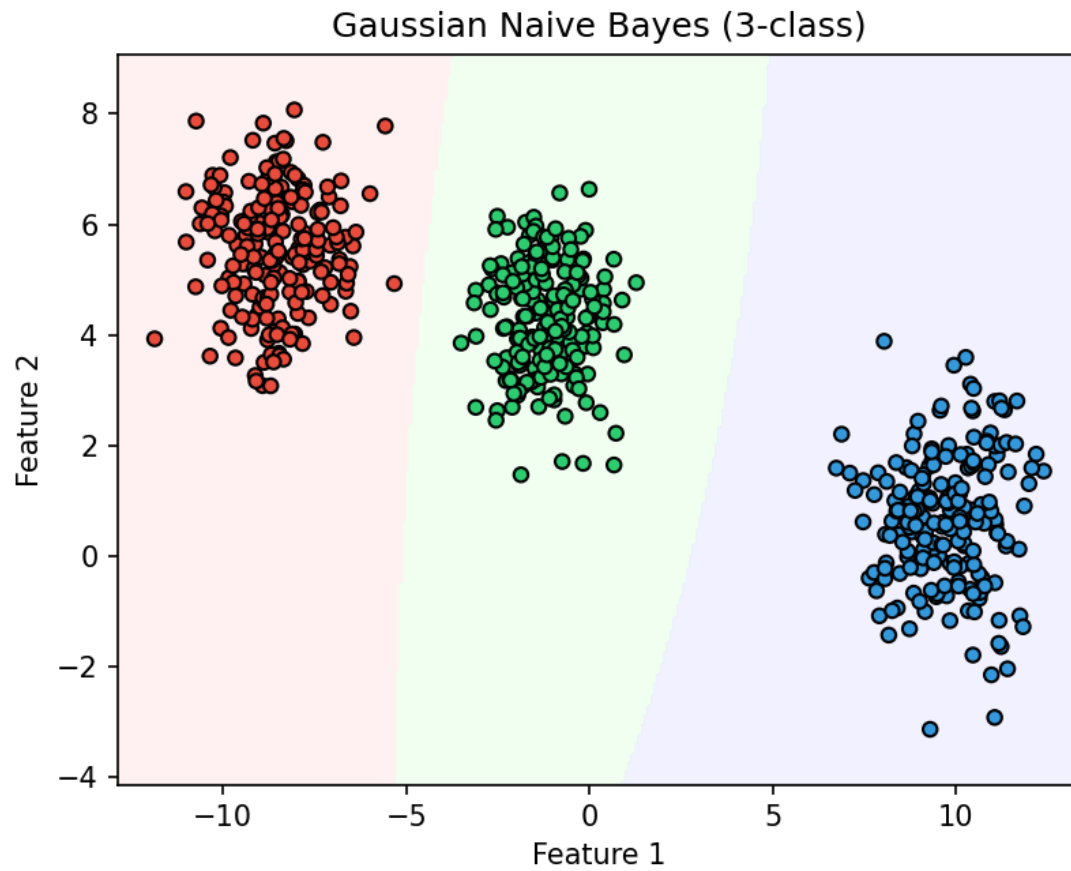


Figure 2: Gaussian NB decision regions (3-class).

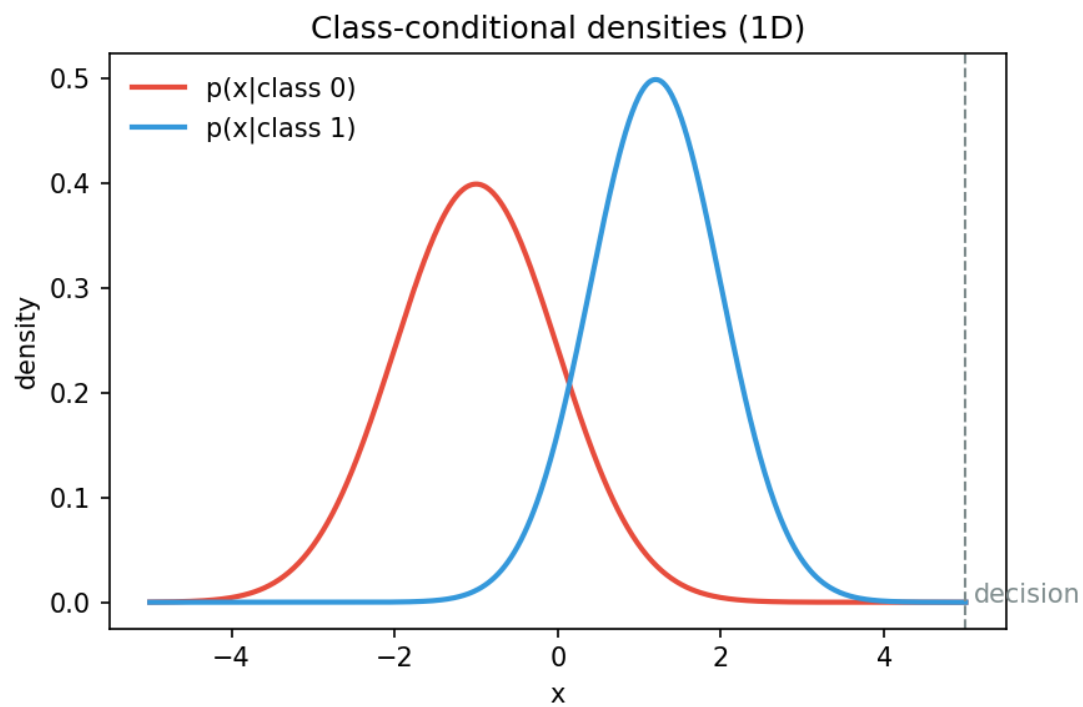


Figure 3: Class-conditional densities in 1D and a decision threshold.

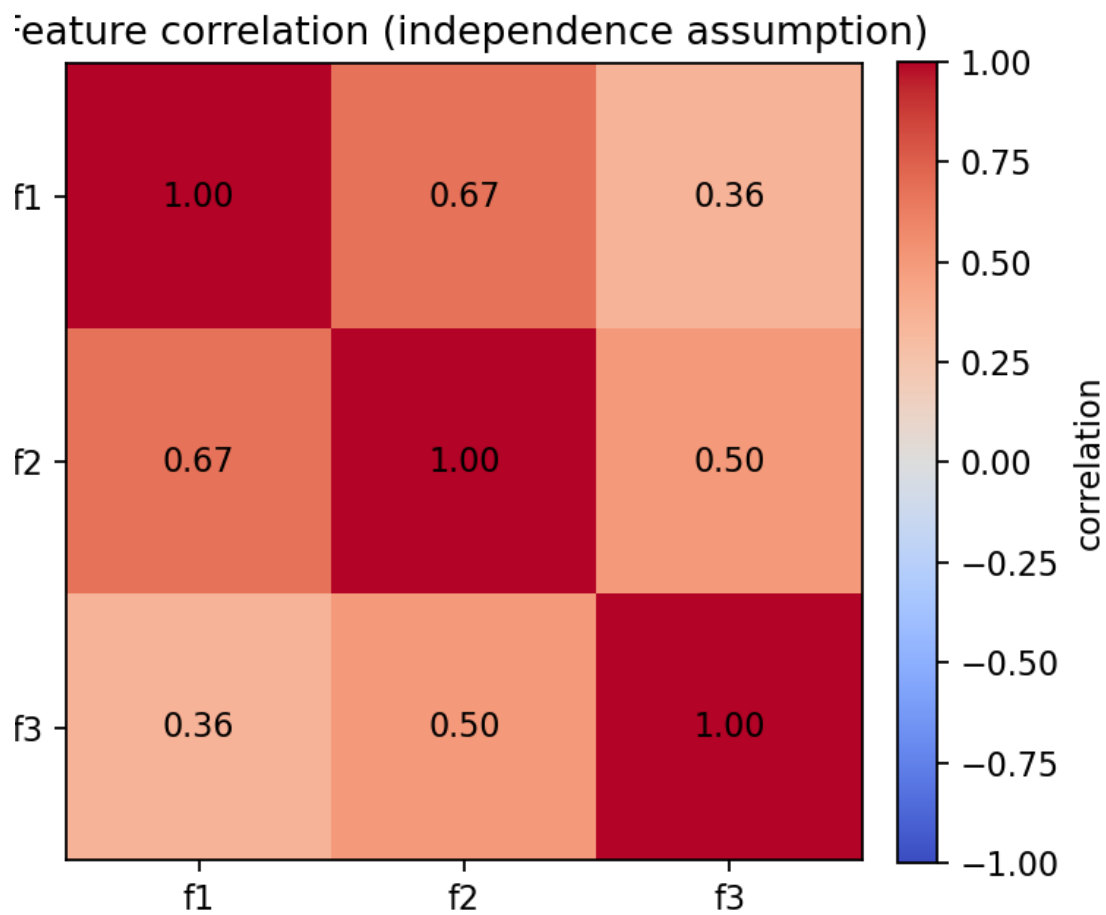


Figure 4: Feature correlation heatmap (independence assumption illustration).

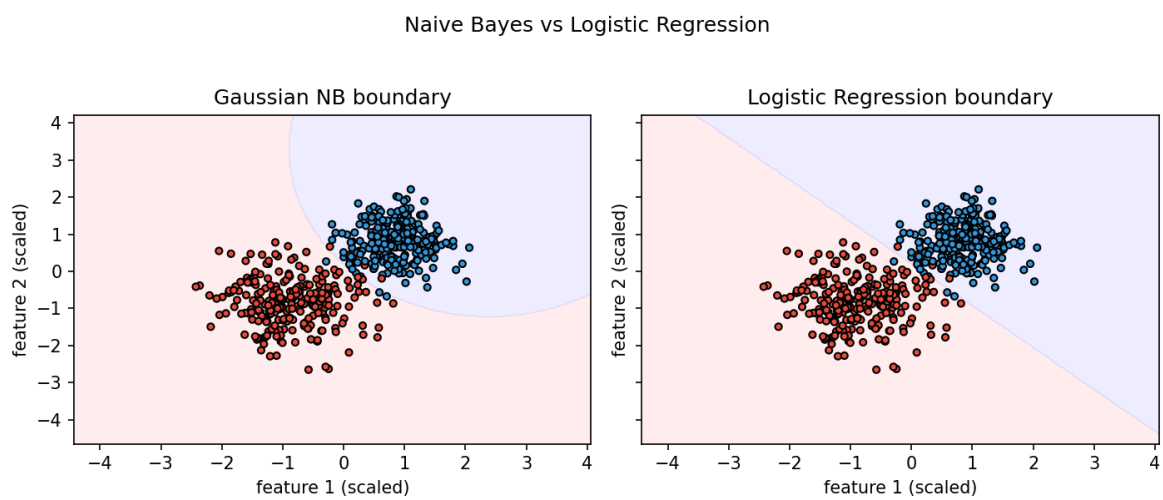


Figure 5: Decision boundary comparison: Gaussian NB vs Logistic Regression.

6 Summary

Naïve Bayes offers a fast, interpretable baseline. Its core idea is simple—combine class priors with per-feature likelihoods under conditional independence. While the assumption is often violated, NB remains competitive on certain problems and serves as a strong baseline against more flexible discriminative models.