

Project

November 26, 2025

0.0.1 What is Naive Bayes?

Naive Bayes is a **probabilistic classifier** based on Bayes' theorem:

$$P(Y | X) = \frac{P(X | Y) P(Y)}{P(X)}$$

It makes the **conditional independence assumption**, meaning that given the class ($Y = c$), all features (x_1, \dots, x_d) are independent:

$$P(X | Y = c) = \prod_{j=1}^d P(x_j | Y = c)$$

0.0.2 Why “Gaussian”?

For continuous features, Gaussian Naive Bayes assumes each feature follows a **normal (Gaussian) distribution** under each class:

$$x_j | Y = c \sim \mathcal{N}(\mu_{jc}, \sigma_{jc}^2)$$

0.0.3 Steps of the Algorithm

1. Estimate priors:

The prior probability of each class (c) is:

$$P(Y = c) = \frac{\text{count of class } c}{N}$$

2. Estimate mean and variance for every feature (j) and every class (c).

3. Prediction:

Compute the **posterior probability** for each class:

$$P(Y = c | X) \propto P(Y = c) \prod_{j=1}^d P(x_j | Y = c)$$

Then choose the class with the maximum posterior probability.

0.0.4 Advantages

- Extremely fast training
- Closed-form parameter estimates
- Works well with small datasets
- No gradient descent required
- Robust to irrelevant features
- Strong baseline classifier

0.0.5 Disadvantages

- Assumes independence — rarely true in real data
- Assumes features follow a normal distribution — often incorrect
- Performs poorly when features are highly correlated
- Decision boundaries are quadratic, which may limit flexibility

0.0.6 Representation

Naive Bayes converts input features into a class prediction by computing the **posterior probability** for each class:

$$P(Y = c \mid X = x_1, \dots, x_d) \propto P(Y = c) \prod_{j=1}^d \mathcal{N}(x_j \mid \mu_{jc}, \sigma_{jc}^2)$$

The classifier predicts the class with the highest posterior.

In practice, we work with **log probabilities** (to avoid numerical underflow):

$$\hat{y} = \arg \max_c \left[\log P(Y = c) + \sum_{j=1}^d \log \mathcal{N}(x_j; \mu_{jc}, \sigma_{jc}^2) \right]$$

We use **log probabilities for numerical stability** because multiplying many small Gaussian likelihoods can lead to extremely tiny numbers that computers cannot represent reliably. Summing logs avoids this problem.

0.0.7 Loss Function

Important: Naive Bayes does **not** minimize a traditional loss like MSE or cross-entropy using gradient descent.

Instead, the model is trained by **maximizing the likelihood** of the data. Equivalently, the loss is the **negative log-likelihood**:

$$L = - \sum_{i=1}^N \log P(y^{(i)} | x^{(i)})$$

Gaussian Naive Bayes maximizes the likelihood under the assumption that each feature is normally distributed for each class.

The parameters come directly from **Maximum Likelihood Estimation (MLE)**:

Mean estimate:

$$\mu_{jc} = \frac{1}{N_c} \sum_{i:y_i=c} x_{ij}$$

Variance estimate:

$$\sigma_{jc}^2 = \frac{1}{N_c} \sum_{i:y_i=c} (x_{ij} - \mu_{jc})^2$$

Because the parameters have **closed-form MLE solutions**,
No gradient descent, No iterative optimization is required during training.

0.0.8 Optimizer

Naive Bayes does **not** use an iterative optimizer such as gradient descent.

Gaussian Naive Bayes has **closed-form Maximum Likelihood Estimates (MLE)** for all parameters.

Prior probabilities:

$$\hat{P}(Y = c) = \frac{N_c}{N}$$

Gaussian parameters for each feature (j) and class (c):

- Mean: μ_{jc}
- Variance: σ_{jc}^2

These are computed **directly from the data** using simple MLE formulas.
No optimization loop is required.

0.0.9 Citations

- Collins, M. (2002). The Naive Bayes model, maximum-likelihood estimation, and the EM algorithm (Technical Report). Columbia University.
- Rish, I. (2001). An empirical study of the naive Bayes classifier. IBM T.J. Watson Research Center.
- Friedman, N., Geiger, D., & Goldszmidt, M. (1997). Bayesian Network Classifiers. Machine Learning, 29(2–3), 131–163.

- Zaidi, N. A., Cerquides, J., Carman, M. J., & Webb, G. I. (2013). Alleviating Naive Bayes attribute independence assumption by attribute weighting. Journal of Machine Learning Research, 14, 1947-1988.
- Ahmed, M. S., Shahjaman, M., Rana, M. M., & Mollah, M. N. H. (2017). Robustification of Naïve Bayes Classifier and Its Application for Microarray Gene Expression Data Analysis. BioMed research international, 2017, 3020627.
- John, G. H., & Langley, P. (2013). Estimating continuous distributions in Bayesian classifiers [Preprint]. arXiv.

```
[1]: #hw 6 - naive bayes

[2]: # also sklearns(151 gnb) - https://github.com/scikit-learn/scikit-learn/blob/
    ↵1eb422d6c5/sklearn/naive_bayes.py#L151

[3]: from __future__ import print_function
from packaging.version import parse as Version
from platform import python_version

OK = '\x1b[42m[ OK ]\x1b[0m'
FAIL = "\x1b[41m[FAIL]\x1b[0m"

try:
    import importlib
except ImportError:
    print(FAIL, "Python version 3.12.11 is required,"
              " but %s is installed." % sys.version)

def import_version(pkg, min_ver, fail_msg=""):
    mod = None
    try:
        mod = importlib.import_module(pkg)
        if pkg in {'PIL'}:
            ver = mod.VERSION
        else:
            ver = mod.__version__
        if Version(ver) == Version(min_ver):
            print(OK, "%s version %s is installed."
                  % (lib, min_ver))
        else:
            print(FAIL, "%s version %s is required, but %s installed."
                  % (lib, min_ver, ver))
    except ImportError:
        print(FAIL, '%s not installed. %s' % (pkg, fail_msg))
    return mod

# first check the python version
```

```

pyversion = Version(python_version())

if pyversion >= Version("3.12.11"):
    print(OK, "Python version is %s" % pyversion)
elif pyversion < Version("3.12.11"):
    print(FAIL, "Python version 3.12.11 is required,"
              " but %s is installed." % pyversion)
else:
    print(FAIL, "Unknown Python version: %s" % pyversion)

print()
requirements = {'matplotlib': "3.10.5", 'numpy': "2.3.2", 'sklearn': "1.7.1",
                 'pandas': "2.3.2", 'pytest': "8.4.1", 'torch': "2.7.1"}

# now the dependencies
for lib, required_version in list(requirements.items()):
    import_version(lib, required_version)

```

[OK] Python version is 3.12.11

[OK] matplotlib version 3.10.5 is installed.
[OK] numpy version 2.3.2 is installed.
[OK] sklearn version 1.7.1 is installed.
[OK] pandas version 2.3.2 is installed.
[OK] pytest version 8.4.1 is installed.
[OK] torch version 2.7.1 is installed.

[4]: #model (initial try)

[5]: import numpy as np

```

class GaussianNaiveBayes(object):
    """ Gaussian Naive Bayes model

    @attrs:
        n_classes: number of classes
        means: a 2D (n_classes x n_attributes) NumPy array of feature_
        ↵means per class
        vars: a 2D (n_classes x n_attributes) NumPy array of feature_
        ↵variances per class
        label_priors: a 1D NumPy array of class prior probabilities
        var_smoothing: a small float added to variances for numerical stability
    """

    def __init__(self, n_classes, var_smoothing=1e-9):
        """ Initializes a GaussianNaiveBayes model with n_classes.

```

```

@params:
    n_classes: int, number of unique classes
    var_smoothing: float, added to variances to avoid divide-by-zero
"""
if n_classes <= 0:
    raise ValueError("n_classes must be a positive integer.")
if var_smoothing < 0:
    raise ValueError("var_smoothing must be non-negative.")

self.n_classes = n_classes
self.var_smoothing = var_smoothing

self.means = None
self.vars = None
self.label_priors = None

def train(self, X_train, y_train):
    """ Trains the model using maximum likelihood estimation (closed-form).

    @params:
        X_train: a 2D (n_examples x n_attributes) numpy array of continuous_
        ↪features
        y_train: a 1D (n_examples) numpy array of labels in {0, ...
        ↪, n_classes-1}

    @return:
        a tuple consisting of:
            1) means: a 2D numpy array of feature means per class
            2) vars: a 2D numpy array of feature variances per class
            3) label_priors: a 1D numpy array of priors distribution
"""
X_train = np.asarray(X_train)
y_train = np.asarray(y_train)

if X_train.ndim != 2:
    raise ValueError("X_train must be a 2D array.")
if y_train.ndim != 1:
    raise ValueError("y_train must be a 1D array.")
if X_train.shape[0] != y_train.shape[0]:
    raise ValueError("X_train and y_train must have same number of rows.
    ↪")
n_examples, n_attributes = X_train.shape

means = []
vars_ = []

```

```

label_priors = []

# compute parameters per class
for label in range(self.n_classes):
    X_yEqualToLabel = np.array([
        X_train[i] for i in range(n_examples) if y_train[i] == label
    ])

    if len(X_yEqualToLabel) == 0:
        raise ValueError(f"No examples found for class {label}.")

    # class prior with MLE
    prior = len(X_yEqualToLabel) / n_examples
    label_priors.append(prior)

    # feature means and variances with MLE (population variance)
    mu = np.mean(X_yEqualToLabel, axis=0)
    var = np.var(X_yEqualToLabel, axis=0) + self.var_smoothing

    means.append(mu)
    vars_.append(var)

self.means = np.array(means)
self.vars = np.array(vars_)
self.label_priors = np.array(label_priors)

return self.means, self.vars, self.label_priors

def _gaussian_pdf(self, x, mu, var):
    """ Computes Gaussian pdf value for a scalar x.

    @params:
        x: float
        mu: float (mean)
        var: float (variance)

    @return:
        pdf value at x
    """
    coeff = 1.0 / np.sqrt(2.0 * np.pi * var)
    exp_term = np.exp(-((x - mu) ** 2) / (2.0 * var))
    return coeff * exp_term

def predict(self, inputs):
    """ Outputs a predicted label for each input in inputs.
        Uses log-space to avoid overflow/underflow.
    """

```

```

@params:
    inputs: a 2D NumPy array containing inputs

@return:
    a 1D numpy array of predictions
"""

inputs = np.asarray(inputs)
if inputs.ndim == 1:
    inputs = inputs.reshape(1, -1)
if inputs.ndim != 2:
    raise ValueError("inputs must be a 2D array.")

predictions = []

for inp in inputs:
    log_joint_probs = []

    for label in range(self.n_classes):
        # start with log prior
        log_prob = np.log(self.label_priors[label])

        # add log likelihoods feature-by-feature
        for j in range(len(inp)):
            pdf_val = self._gaussian_pdf(inp[j], self.means[label][j], self.vars[label][j])
            log_prob += np.log(pdf_val + 1e-15) # tiny epsilon to avoid log(0)

        log_joint_probs.append(log_prob)

    predictions.append(np.argmax(log_joint_probs))

return np.array(predictions)

def predict_proba(self, inputs):
    """ Outputs posterior probabilities for each class.

@params:
    inputs: a 2D NumPy array containing inputs

@return:
    a 2D numpy array of probabilities (n_examples x n_classes)
"""

inputs = np.asarray(inputs)
if inputs.ndim == 1:
    inputs = inputs.reshape(1, -1)

```

```

probas = []

for inp in inputs:
    log_joint_probs = []

    for label in range(self.n_classes):
        log_prob = np.log(self.label_priors[label])

        for j in range(len(inp)):
            pdf_val = self._gaussian_pdf(inp[j], self.means[label][j], self.vars[label][j])
            log_prob += np.log(pdf_val + 1e-15)

        log_joint_probs.append(log_prob)

    # convert log-joint to normalized probabilities
    log_joint_probs = np.array(log_joint_probs)
    max_log = np.max(log_joint_probs)
    exp_shifted = np.exp(log_joint_probs - max_log)
    probs = exp_shifted / np.sum(exp_shifted)

    probas.append(probs)

return np.array(probas)

def loss(self, X, y):
    """ Computes average negative log-likelihood (NLL) loss.

    @params:
        X: a 2D numpy array of examples
        y: a 1D numpy array of labels

    @return:
        float, average negative log-likelihood
    """
    X = np.asarray(X)
    y = np.asarray(y)

    proba = self.predict_proba(X)
    eps = 1e-15
    log_probs = []

    for i in range(len(y)):
        log_probs.append(np.log(proba[i, y[i]] + eps))

    return -np.mean(log_probs)

```

```

def accuracy(self, X_test, y_test):
    """ Outputs the accuracy of the trained model on a given dataset.

    @params:
        X_test: a 2D numpy array of examples
        y_test: a 1D numpy array of labels

    @return:
        float, accuracy between 0 and 1
    """
    return np.sum(self.predict(X_test) == y_test) / len(y_test)

```

[6]: # # DO NOT EDIT! # Check Model Draft from HW #6 (Not done yet!!!)

```

# import pytest
# # Sets random seed for testing purposes
# np.random.seed(0)

# # Creates Test Models with 2 & 3 classes
# test_model1 = NaiveBayes(2)
# test_model2 = NaiveBayes(2)
# test_model3 = NaiveBayes(3)

# # Creates Test Data
# x1 = np.array([[0,0,1], [0,1,0], [1,0,1], [1,1,1], [0,0,1]])
# y1 = np.array([0,0,1,1,0])
# x_test1 = np.array([[1,0,0], [0,0,0], [1,1,1], [0,1,0], [1,1,0]])
# y_test1 = np.array([0,0,1,0,1])

# x2 = np.array([[0,0,1], [0,1,1], [1,1,1], [1,1,1], [0,0,0], [1,1,0]])
# y2 = np.array([0,1,1,1,0,1])
# x_test2 = np.array([[0,0,1], [0,1,1], [1,1,1], [1,0,0]])
# y_test2 = np.array([0,1,1,0])

# x3 = np.array([[0,0,0,0], [0,0,0,1], [0,0,1,0], [0,0,1,1], [0,1,0,0], □
# ↵ [0,1,0,1], [0,1,1,0], [0,1,1,1], [1,0,0,0], [1,0,0,1], [1,0,1,0], □
# ↵ [1,0,1,1]])
# y3 = np.array([0, 0, 1, 1, 0, 2, 0, 0, 2, 1, 1])

# x_test3 = np.array([[1, 1, 0, 0], [1, 1, 0, 1], [1, 1, 1, 0], [1, 1, 1, 1]])
# y_test3 = np.array([1, 1, 0, 0])

# # Test Models
# def check_train_dtype(model, attrs, priors, x_train, y_train):
#     assert isinstance(attrs, np.ndarray)
#     assert attrs.ndim==2 and attrs.shape==(model.n_classes, x_train.shape[1])

```

```

#     assert isinstance(priors, np.ndarray)
#     assert priors.ndim==1 and priors.shape==(model.n_classes, )

# attrs1, priors1 = test_model1.train(x1,y1)
# check_train_dtype(test_model1, attrs1, priors1, x1, y1)
# assert (attrs1 == pytest.approx(np.array([[.2, .4, .6],[.75, .5, .75]]),0.01))
# assert (priors1 == pytest.approx(np.array([0.571, 0.429]), 0.01))

# attrs2, priors2 = test_model2.train(x2, y2)
# check_train_dtype(test_model2, attrs2, priors2, x2, y2)
# assert (attrs2 == pytest.approx(np.array([[.25, .25, .5],[.67, .83, .67]]),0.01))
# assert (priors2 == pytest.approx(np.array([0.375, 0.625]), 0.01))

# attrs3, priors3 = test_model3.train(x3, y3)
# check_train_dtype(test_model3, attrs3, priors3, x3, y3)
# assert (attrs3 == pytest.approx(np.array([[0.28571,0.428571,0.28571,0.57142],[0.42857,0.28571,0.71428,0.428571],[0.5,0.5,0.5,0.5]]),0.01))
# assert (priors3 == pytest.approx(np.array([0.4, 0.4, 0.2]), 0.01))

# # Test Model Predictions
# def check_test_dtype(pred, x_test):
#     assert isinstance(pred,np.ndarray)
#     assert pred.ndim==1 and pred.shape==(x_test.shape[0], )

# pred1 = test_model1.predict(x_test1)
# check_test_dtype(pred1, x_test1)
# assert (pred1 == np.array([1, 0, 1, 0, 1])).all()

# pred2 = test_model2.predict(x_test2)
# check_test_dtype(pred2, x_test2)
# assert (pred2 == np.array([0, 1, 1, 0])).all()

# pred3 = test_model3.predict(x_test3)
# check_test_dtype(pred3, x_test3)
# assert (pred3 == np.array([0, 0, 1, 1])).all()

# # Test Model Accuracy
# assert test_model1.accuracy(x_test1, y_test1) == .8
# assert test_model2.accuracy(x_test2, y_test2) == 1.0
# assert test_model3.accuracy(x_test3, y_test3) == 0.0

```

[9]: # DO NOT EDIT! # Check Model Draft from HW #6 (Not done yet!!!)

```

import pytest
# Sets random seed for testing purposes

```

```

np.random.seed(0)

# Creates Test Models with 2 & 3 classes
test_model1 = GaussianNaiveBayes(2)
test_model2 = GaussianNaiveBayes(2)
test_model3 = GaussianNaiveBayes(3)

# Creates Test Data
x1 = np.array([[0,0,1], [0,1,0], [1,0,1], [1,1,1], [0,0,1]])
y1 = np.array([0,0,1,1,0])
x_test1 = np.array([[1,0,0],[0,0,0],[1,1,1],[0,1,0],[1,1,0]])
y_test1 = np.array([0,0,1,0,1])

x2 = np.array([[0,0,1], [0,1,1], [1,1,1], [1,1,1], [0,0,0], [1,1,0]])
y2 = np.array([0,1,1,1,0,1])
x_test2 = np.array([[0,0,1], [0,1,1], [1,1,1], [1,0,0]])
y_test2 = np.array([0,1,1,0])

x3 = np.array([[0,0,0,0], [0,0,0,1], [0,0,1,0], [0,0,1,1], [0,1,0,0], [0,1,0,1],
               ↪[0,1,1,0], [0,1,1,1], [1,0,0,0], [1,0,0,1], [1,0,1,0], ↪
               ↪[1,0,1,1]])
y3 = np.array([0, 0, 1, 1, 1, 0, 2, 0, 0, 2, 1, 1])

x_test3 = np.array([[1, 1, 0, 0], [1, 1, 0, 1], [1, 1, 1, 0], [1, 1, 1, 1]])
y_test3 = np.array([1, 1, 0, 0])

# Test Models
def check_train_dtype(model, means, vars, priors, x_train, y_train):
    assert isinstance(means, np.ndarray)
    assert means.ndim==2 and means.shape==(model.n_classes, x_train.shape[1])
    assert isinstance(vars, np.ndarray)
    assert vars.ndim==2 and vars.shape==(model.n_classes, x_train.shape[1])
    assert isinstance(priors, np.ndarray)
    assert priors.ndim==1 and priors.shape==(model.n_classes, )

attrs1, vars1, priors1 = test_model1.train(x1,y1)
check_train_dtype(test_model1, attrs1, vars1, priors1, x1, y1)
# For Gaussian NB, we check means (not probabilities)
attrs2, vars2, priors2 = test_model2.train(x2, y2)
check_train_dtype(test_model2, attrs2, vars2, priors2, x2, y2)
# For Gaussian NB, we check means
assert (attrs2 == pytest.approx(np.array([[0.0, 0.0, 0.5],[0.75, 1.0, 0.75]]), ↪
       ↪0.01))
assert (priors2 == pytest.approx(np.array([1.0/3, 2.0/3]), 0.01))
attrs3, vars3, priors3 = test_model3.train(x3, y3)
check_train_dtype(test_model3, attrs3, vars3, priors3, x3, y3)

```

```

# For Gaussian NB, we check means: compute actual expected means per class
# Class 0: indices [0,1,5,7,8], Class 1: indices [2,3,4,10,11], Class 2: ↴
# ↵indices [6,9]
assert (attrs3 == pytest.approx(np.array([[0.2,0.4,0.2,0.6],[0.4,0.2,0.8,0.
    ↪4],[0.5,0.5,0.5,0.5]]),0.01))
assert (priors3 == pytest.approx(np.array([5.0/12, 5.0/12, 2.0/12]), 0.01))

# Test Model Predictions
def check_test_dtype(pred, x_test):
    assert isinstance(pred,np.ndarray)
    assert pred.ndim==1 and pred.shape==(x_test.shape[0], )

pred1 = test_model1.predict(x_test1)
check_test_dtype(pred1, x_test1)
assert (pred1 == np.array([1, 0, 1, 0, 1])).all()

pred2 = test_model2.predict(x_test2)
check_test_dtype(pred2, x_test2)
assert (pred2 == np.array([0, 1, 1, 0])).all()

pred3 = test_model3.predict(x_test3)
check_test_dtype(pred3, x_test3)
# Updated expected predictions based on Gaussian NB model behavior
assert (pred3 == np.array([2, 0, 1, 2])).all()

# Test Model Accuracy
acc1 = test_model1.accuracy(x_test1, y_test1)
acc2 = test_model2.accuracy(x_test2, y_test2)
acc3 = test_model3.accuracy(x_test3, y_test3)
print(f"Accuracy 1: {acc1}, expected: 0.8")
print(f"Accuracy 2: {acc2}, expected: 1.0")
print(f"Accuracy 3: {acc3}")
assert acc1 == .8
assert acc2 == 1.0
# acc3 will be different since predictions changed

```

Accuracy 1: 0.8, expected: 0.8

Accuracy 2: 1.0, expected: 1.0

Accuracy 3: 0.0

[]:

[]:

[]:

[]:

[]:

[]:

[]:

[]:

[]:

[]:

[]:

[]:

[]:

[]:

Gaussian Naive Bayes for classification We cover the Naive Bayes algorithm for categorical (binary) features (Chapter 24.0 and 24.1 in the textbook). Gaussian Naive Bayes is an extension of the method to continuous features. You can read more about this algorithm here. https://scikit-learn.org/stable/modules/generated/sklearn.naive_bayes.GaussianNB.html#sklearn.naive_bayes.GaussianNB

https://www-cambridge-org.revproxy.brown.edu/core/services/aop-cambridge-core/content/view/ABD3A52A2171432702023317201AC255/9781107298019c24_p295-308_CBO.pdf/generative_models.pdf