# Week 7: Naive Bayes Classifiers in scikit-learn

In Week 7, we delved into the world of Naive Bayes classifiers, exploring their theoretical underpinnings and practical implementations within the scikit-learn library. We covered various types of Naive Bayes suitable for different data distributions and discussed their strengths and weaknesses.

## Naive Bayes: Core Concepts

The Naive Bayes classifier is a probabilistic classifier based on Bayes' theorem with a strong "naive" independence assumption. This assumption simplifies the calculations significantly, making it computationally efficient. The core idea is to estimate the probability of a data point belonging to each class and assign it to the class with the highest probability.

The Bayes' theorem is given by:

$$ P(y|x) = \frac{P(x|y)P(y)}{P(x)} $$

where:

- $P(y|x)$ is the posterior probability of class $y$ given the features $x$.
- $P(x|y)$ is the likelihood of observing features $x$ given class $y$.
- $P(y)$ is the prior probability of class $y$.
- $P(x)$ is the prior probability of features $x$.

The "naive" assumption is that the features are conditionally independent given the class:

$$ P(x|y) = \prod_{i=1}^{m} P(x_i|y) $$

This simplification significantly reduces the computational complexity.

## Types of Naive Bayes Classifiers in scikit-learn

Scikit-learn provides several implementations of Naive Bayes classifiers, each tailored to different data types:

- **GaussianNB:** Assumes features follow a Gaussian (normal) distribution. Suitable for continuous numerical data.

- **MultinomialNB:** Assumes features follow a multinomial distribution. Excellent for text classification (word counts).

- **BernoulliNB:** Assumes features are binary (0 or 1). Suitable for binary data or data that can be binarized.

- **CategoricalNB:** Handles categorical features with a categorical distribution. Suitable for discrete features with multiple categories.

- **ComplementNB:** A variation designed to handle imbalanced datasets. Often outperforms MultinomialNB,

especially in text classification.

# Implementation in scikit-learn

All the above classifiers share a similar implementation pattern in scikit-learn:

```python
from sklearn.naive_bayes import GaussianNB, MultinomialNB, BernoulliNB, CategoricalNB, ComplementNB

# Example using GaussianNB
gnb = GaussianNB()
gnb.fit(X_train, y_train)
y_pred = gnb.predict(X_test)

# Example using MultinomialNB
mnb = MultinomialNB()
mnb.fit(X_train, y_train)
y_pred = mnb.predict(X_test)

#Similar for BernoulliNB, CategoricalNB and ComplementNB
```

The `fit()` method estimates the parameters (e.g., means and variances for GaussianNB, probabilities for MultinomialNB) from the training data. The `predict()` method uses these estimated parameters to predict the class labels for new data points.

# Handling Imbalanced Datasets

Imbalanced datasets (where one class has significantly more samples than others) can affect the performance of Naive Bayes classifiers. The `ComplementNB` classifier is specifically designed to address this issue by focusing on the complement of the class probabilities.

# Computational Complexity and Performance

Naive Bayes classifiers are known for their low computational complexity. Training and prediction times are generally fast, making them suitable for large datasets. The time complexity is typically linear in the number of data points and features.

# Choosing the Right Naive Bayes Classifier

The choice of the appropriate Naive Bayes classifier depends heavily on the nature of the data:

| Data Type | Classifier |
|---|---|
| Continuous Numerical | GaussianNB |
| Discrete Counts | MultinomialNB |
| Binary | BernoulliNB |
| Categorical | CategoricalNB |
| Imbalanced Datasets | ComplementNB |

# Conclusion

In conclusion, Week 7 provided a thorough introduction to Naive Bayes classifiers and their practical application using scikit-learn. We learned about different types of Naive Bayes classifiers, their underlying assumptions, and how to choose the appropriate classifier for a given dataset. We also explored the implementation details and the handling of imbalanced datasets. The efficiency and simplicity of Naive Bayes make it a valuable tool in the machine learning practitioner's arsenal. Prof. Ashish's insights were invaluable in understanding the nuances of these powerful yet straightforward algorithms.