

Comprehensive Notes: Bayes' Theorem, Naive Bayes & Laplace Smoothing

1. Bayes' Theorem Fundamentals

Definition and Formula

Bayes' Theorem describes how to update probabilities based on new evidence. It answers the question: "What is the probability of event A occurring given that event B has occurred?"

Formula:

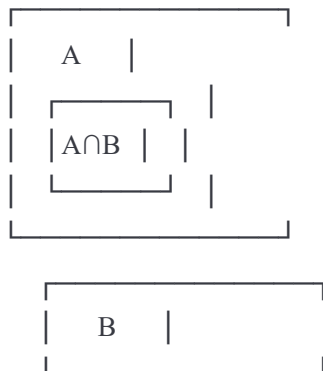
$$P(A|B) = [P(B|A) \times P(A)] / P(B)$$

Components:

- **P(A|B)**: Posterior probability (probability of A given B has occurred)
- **P(B|A)**: Likelihood (probability of observing B given A is true)
- **P(A)**: Prior probability (initial probability of A before any evidence)
- **P(B)**: Evidence or marginal probability (total probability of observing B)

Understanding with Venn Diagrams

Imagine a Venn diagram with two circles representing events A and B:



- The **overlap** ($A \cap B$) represents instances where both events occur
- **P(A|B)** measures what fraction of circle B overlaps with circle A
- **P(B|A)** measures what fraction of circle A overlaps with circle B

Mathematical interpretation:

$$P(A|B) = P(A \cap B) / P(B)$$

$$P(B|A) = P(A \cap B) / P(A)$$

The intersection $P(A \cap B)$ is the same in both cases, but we divide by different denominators depending on what we're conditioning on.

Real-World Example: Medical Diagnosis

Suppose a disease affects 1% of the population, and a test is 99% accurate:

- $P(\text{Disease}) = 0.01$ (prior: 1% of people have disease)
- $P(\text{Positive}|\text{Disease}) = 0.99$ (likelihood: test is 99% accurate if you have disease)
- $P(\text{Positive}) = P(\text{Positive}|\text{Disease}) \times P(\text{Disease}) + P(\text{Positive}|\neg\text{Disease}) \times P(\neg\text{Disease})$
- $P(\text{Positive}) = 0.99 \times 0.01 + 0.01 \times 0.99 = 0.0198$

What's the probability you actually have the disease given a positive test?

$$\begin{aligned} P(\text{Disease}|\text{Positive}) &= [P(\text{Positive}|\text{Disease}) \times P(\text{Disease})] / P(\text{Positive}) \\ &= [0.99 \times 0.01] / 0.0198 \\ &= 0.0099 / 0.0198 \\ &= 50\% \end{aligned}$$

Despite the test being 99% accurate, a positive result only gives you a 50% chance of having the disease! This is because the disease is rare.

2. Mutually Exclusive vs. Independent Events

Mutually Exclusive Events

Events A and B are **mutually exclusive** if they cannot occur simultaneously.

Characteristics:

- $A \cap B = \emptyset$ (empty set)
- $P(A \cap B) = 0$
- $P(A|B) = 0$
- $P(B|A) = 0$

Example: Rolling a die

- $A = \text{"rolling a 3"}$
- $B = \text{"rolling a 5"}$
- These cannot happen on the same roll, so $P(A \cap B) = 0$

Visual with Venn diagram:



(no overlap)

Independent Events

Events A and B are **independent** if the occurrence of one doesn't affect the probability of the other.

Characteristics:

- $P(A \cap B) = P(A) \times P(B)$
- $P(A|B) = P(A)$
- $P(B|A) = P(B)$

Example: Two coin flips

- $A = \text{"first flip is heads"}$
- $B = \text{"second flip is heads"}$
- These are independent: knowing the first is heads doesn't change the probability the second is heads

Key Difference

Mutually Exclusive \neq Independent

In fact, if two events are mutually exclusive ($P(A \cap B) = 0$) and both have non-zero probabilities, they are **negatively dependent** (knowing one occurred makes the other impossible).

3. The Zero Frequency Problem in Naive Bayes

What is Naive Bayes?

Naive Bayes is a probabilistic classification algorithm that uses Bayes' theorem with the simplifying assumption that features are conditionally independent given the class label.

For a text classification example:

$$P(\text{Class}|\text{Document}) \propto P(\text{Class}) \times P(\text{word1}|\text{Class}) \times P(\text{word2}|\text{Class}) \times \dots \times P(\text{wordN}|\text{Class})$$

The Zero Frequency Problem

Suppose we're classifying emails as "Spam" or "Ham" and we encounter the word "cryptocurrency" only in spam emails (never in ham).

During training:

- $\text{count}(\text{cryptocurrency}, \text{Spam}) = 50$
- $\text{count}(\text{cryptocurrency}, \text{Ham}) = 0$
- $P(\text{cryptocurrency}|\text{Ham}) = 0/100 = 0$

During prediction:

- If a new email contains "cryptocurrency" but also contains other words, the entire posterior becomes:

$$\begin{aligned} P(\text{Ham}|\text{email}) &\propto P(\text{Ham}) \times P(\text{cryptocurrency}|\text{Ham}) \times P(\text{other_words}|\text{Ham}) \\ &= P(\text{Ham}) \times 0 \times P(\text{other_words}|\text{Ham}) \\ &= 0 \end{aligned}$$

This means the email is classified as spam with 100% certainty, even if it contains many typical ham words.

This is problematic because:

1. We've never seen "cryptocurrency" in ham emails during training
2. But that doesn't mean it's impossible for it to appear in ham emails in the future
3. Treating unseen events as **impossible** (zero probability) is overly harsh

4. Laplace Smoothing: The Solution

What is Laplace Smoothing?

Laplace smoothing (also called add-one smoothing) is a technique that prevents zero probabilities by adding a small constant to the count of each feature-class combination.

Formula:

$$P(x = v \mid y) = [\text{count}(x=v, y) + 1] / [\text{count}(y) + K]$$

Where:

- **count(x=v, y):** Number of times feature x has value v in class y
- **count(y):** Total count of instances in class y
- **K:** Number of possible feature values for feature x

Why K in the Denominator?

When we add 1 to the numerator, we must also add K to the denominator to maintain valid probability distributions. If we only added 1, probabilities wouldn't sum to 1.

Example: Email Classification with Laplace Smoothing

Scenario: 100 ham emails, word "cryptocurrency" appears 0 times

Without Laplace smoothing:

$$P(\text{cryptocurrency}|\text{Ham}) = 0 / 100 = 0$$

With Laplace smoothing (assuming 50,000 possible words in vocabulary):

$$\begin{aligned} P(\text{cryptocurrency}|\text{Ham}) &= (0 + 1) / (100 + 50000) \\ &= 1 / 50100 \\ &\approx 0.00002 \end{aligned}$$

This small probability is much better than zero! It says: "We haven't seen this word in ham before, but it's still slightly possible."

Another Example: Spam Classification

Training data:

- Spam emails: 200 total
- Word "free" appears in 150 spam emails
- Word "crypto" appears in 0 spam emails
- Vocabulary size K = 5000 words

Probability calculations:

Without smoothing:

$$P(\text{free}|\text{Spam}) = 150/200 = 0.75$$

$$P(\text{crypto}|\text{Spam}) = 0/200 = 0$$

With Laplace smoothing:

$$P(\text{free}|\text{Spam}) = (150 + 1) / (200 + 5000) = 151/5200 \approx 0.029$$

$$P(\text{crypto}|\text{Spam}) = (0 + 1) / (200 + 5000) = 1/5200 \approx 0.0002$$

Notice how Laplace smoothing **shrinks high probabilities toward the uniform distribution** ($1/K$) while **lifting zero probabilities above zero**.

5. Intuition Behind Laplace Smoothing

Conceptual Understanding

Think of Laplace smoothing as adding a "pseudo-count" or "virtual instance" for each feature-class combination. It's like saying:

┆ "Before seeing any real data, assume each feature value could appear once in each class."

This represents a **uniform prior belief**: in the absence of evidence, all feature values are equally likely.

The Connection to Mutual Exclusivity

Without smoothing, we create **false mutual exclusivity**. We're saying: "Because we didn't see X in class Y during training, X and Y are mutually exclusive (impossible together)."

Laplace smoothing prevents this by adding a small overlap in the Venn diagram:

Without Laplace smoothing:

Features seen in Class A Features seen in Class B



(completely separate - false mutual exclusivity)

With Laplace smoothing:

Features seen in Class A Features seen in Class B



(allows for unseen feature-class combinations)

6. Practical Considerations

When to Use Laplace Smoothing

- **Always use it in Naive Bayes** for text classification and similar tasks
- When you have rare features or small datasets
- When vocabulary or feature space is large relative to training data size

Effect on Results

- **With small datasets:** Laplace smoothing has a larger impact on probabilities
- **With large datasets:** The effect diminishes (adding 1 to large counts barely matters)
- **With large vocabularies:** The impact is more pronounced because K is larger

Alternative Smoothing Techniques

1. **Add-k smoothing:** Instead of adding 1, add a smaller constant k ($0 < k < 1$)

$$P(x = v \mid y) = [\text{count}(x=v, y) + k] / [\text{count}(y) + K \times k]$$

2. **Good-Turing smoothing:** Uses frequency of frequencies to estimate unseen event probabilities
3. **Backoff and interpolation:** Combine probabilities from different feature levels

7. Interview Summary

One-Liner Answer

"Laplace smoothing prevents zero probabilities in Naive Bayes caused by unseen feature-class combinations. It adds a pseudo-count of 1 to each feature value, ensuring that even unseen combinations have a small non-zero probability instead of being treated as impossible."

Key Points to Remember

1. **Problem:** Zero frequency of a feature-class combination makes the entire posterior probability zero
2. **Solution:** Add 1 to the numerator and K to the denominator
3. **Effect:** Balances between empirical probabilities and a uniform prior
4. **Why it works:** Prevents overfitting to training data and handles unseen events gracefully
5. **Trade-off:** Slightly biases probabilities toward uniform distribution but prevents catastrophic zero probabilities