# Naive Bayes

BY ANIL KUMAR APSSDC

# Bayes Theorem

Alex       Brenda

$P(\text{Alex}) = 0.5 \quad P(\text{Brenda}) = 0.5$

Prior      $P(\text{Alex}) = 0.5 \quad P(\text{Brenda}) = 0.5$

Posterior  $P(\text{Alex}) = 0.4 \quad P(\text{Brenda}) = 0.6$

Person had a red sweater

Alex wears red 2 times a week

Brenda wears red 3 times a week

# Principle of Naive Bayes Classifier

A Naive Bayes classifier is a probabilistic machine learning model that's used for classification task. The crux of the classifier is based on the Bayes theorem.

## Bayes Theorem

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)}$$

Using Bayes theorem, we can find the probability of **A** happening, given that **B** has occurred. Here, **B** is the evidence and **A** is the hypothesis. The assumption made here is that the predictors/features are independent. That is presence of one particular feature does not affect the other. Hence it is called naive.

# Types of Naive Bayes Classifier

1. **Multinomial Naive Bayes:**
This is mostly used for **document classification** problem, i.e. whether a document belongs to the category of sports, politics, technology etc. The features/predictors used by the classifier are the frequency of the words present in the document.

2. **Bernoulli Naive Bayes:**
This is similar to the multinomial naive bayes but the predictors are Boolean variables. The parameters that we use to predict the class variable take up only values yes or no, for example if a **word occurs in the text or not.**

3. **Gaussian Naive Bayes:**
When the predictors take up a **continuous value and are not discrete**, we assume that these values are sampled from a gaussian distribution.

$$P(x_i|y) = \frac{1}{\sqrt{2\pi\sigma_y^2}} exp\left(-\frac{(x_i-\mu_y)^2}{2\sigma_y^2}\right)$$

# Example

| | OUTLOOK | TEMPERATURE | HUMIDITY | WINDY | PLAY GOLF |
|---|---|---|---|---|---|
| 0 | Rainy | Hot | High | False | No |
| 1 | Rainy | Hot | High | True | No |
| 2 | Overcast | Hot | High | False | Yes |
| 3 | Sunny | Mild | High | False | Yes |
| 4 | Sunny | Cool | Normal | False | Yes |
| 5 | Sunny | Cool | Normal | True | No |
| 6 | Overcast | Cool | Normal | True | Yes |
| 7 | Rainy | Mild | High | False | No |
| 8 | Rainy | Cool | Normal | False | Yes |
| 9 | Sunny | Mild | Normal | False | Yes |
| 10 | Rainy | Mild | Normal | True | Yes |
| 11 | Overcast | Mild | High | True | Yes |
| 12 | Overcast | Hot | Normal | False | Yes |
| 13 | Sunny | Mild | High | True | No |

# Summarizing the Dataset

| | outlook | temp. | humidity | windy | play |
|-----|----------|------|----------|-------|------|
| D1 | sunny | hot | high | false | no |
| D2 | sunny | hot | high | true | no |
| D3 | overcast | hot | high | false | yes |
| D4 | rainy | mild | high | false | yes |
| D5 | rainy | cool | normal | false | yes |
| D6 | rainy | cool | normal | true | no |
| D7 | overcast | cool | normal | true | yes |
| D8 | sunny | mild | high | false | no |
| D9 | sunny | cool | normal | false | yes |
| D10 | rainy | mild | normal | false | yes |
| D11 | sunny | mild | normal | true | yes |
| D12 | overcast | mild | high | true | yes |
| D13 | overcast | hot | normal | false | yes |
| D14 | rainy | mild | high | true | no |

## Should I play Today?

- **Total Sample:** 14
- **No. of Attributes:** 4
- **No. of Class:** 1

**Total Yes:** 9  **P(Yes):** 9/14

**Total No:** 5  **P(No):** 5/14

# Frequency of Each Attribute

| | outlook | temp. | humidity | windy | play |
|---|---|---|---|---|---|
| D1 | sunny | hot | high | false | no |
| D2 | sunny | hot | high | true | no |
| D3 | overcast | hot | high | false | yes |
| D4 | rainy | mild | high | false | yes |
| D5 | rainy | cool | normal | false | yes |
| D6 | rainy | cool | normal | true | no |
| D7 | overcast | cool | normal | true | yes |
| D8 | sunny | mild | high | false | no |
| D9 | sunny | cool | normal | false | yes |
| D10 | rainy | mild | normal | false | yes |
| D11 | sunny | mild | normal | true | yes |
| D12 | overcast | mild | high | true | yes |
| D13 | overcast | hot | normal | false | yes |
| D14 | rainy | mild | high | true | no |

| Outlook | Yes | No |
|---|---|---|
| Sunny | 2 | 3 |
| Overcast | 4 | 0 |
| Rainy | 3 | 2 |

| Temperature | Yes | No |
|---|---|---|
| Hot | 2 | 2 |
| Mild | 4 | 2 |
| Cool | 3 | 1 |

| Humidity | Yes | No |
|---|---|---|
| High | 3 | 4 |
| Normal | 6 | 1 |

| Windy | Yes | No |
|---|---|---|
| False | 6 | 2 |
| True | 3 | 3 |

# Probability of Each Attribute P(A)

| | outlook | temp. | humidity | windy | play |
|----|---------|-------|----------|-------|------|
| D1 | sunny | hot | high | false | no |
| D2 | sunny | hot | high | true | no |
| D3 | overcast | hot | high | false | yes |
| D4 | rainy | mild | high | false | yes |
| D5 | rainy | cool | normal | false | yes |
| D6 | rainy | cool | normal | true | no |
| D7 | overcast | cool | normal | true | yes |
| D8 | sunny | mild | high | false | no |
| D9 | sunny | cool | normal | false | yes |
| D10 | rainy | mild | normal | false | yes |
| D11 | sunny | mild | normal | true | yes |
| D12 | overcast | mild | high | true | yes |
| D13 | overcast | hot | normal | false | yes |
| D14 | rainy | mild | high | true | no |

| Outlook | Yes | No | P(attrib) |
|---------|-----|-----|-----------|
| Sunny | 2 | 3 | 5/14 |
| Overcast | 4 | 0 | 4/14 |
| Rainy | 3 | 2 | 5/14 |
| Total | | | 100% |

| Temperature | Yes | No | P(attrib) |
|-------------|-----|-----|-----------|
| Hot | 2 | 2 | 4/14 |
| Mild | 4 | 2 | 6/14 |
| Cool | 3 | 1 | 4/14 |
| Total | | | 100% |

| Humidity | Yes | No | P(attrib) |
|----------|-----|-----|-----------|
| High | 3 | 4 | 7/14 |
| Normal | 6 | 1 | 7/14 |
| Total | | | 100% |

| Windy | Yes | No | P(attrib) |
|-------|-----|-----|-----------|
| False | 6 | 2 | 8/14 |
| True | 3 | 3 | 6/14 |
| Total | | | 100% |

# Probability of Each Attribute P(B)

| | outlook | temp. | humidity | windy | play |
|------|----------|-------|----------|-------|------|
| D1 | sunny | hot | high | false | no |
| D2 | sunny | hot | high | true | no |
| D3 | overcast | hot | high | false | yes |
| D4 | rainy | mild | high | false | yes |
| D5 | rainy | cool | normal | false | yes |
| D6 | rainy | cool | normal | true | no |
| D7 | overcast | cool | normal | true | yes |
| D8 | sunny | mild | high | false | no |
| D9 | sunny | cool | normal | false | yes |
| D10 | rainy | mild | normal | false | yes |
| D11 | sunny | mild | normal | true | yes |
| D12 | overcast | mild | high | true | yes |
| D13 | overcast | hot | normal | false | yes |
| D14 | rainy | mild | high | true | no |

| Play | | Probability |
|-------|-----|-------------|
| Yes | 9 | 9/14 |
| No | 5 | 5/14 |
| Total | 14 | 100% |

# Probability of Each Attribute P(A|B)

| Outlook | Yes | No | P(Yes) | P(No) |
|---|---|---|---|---|
| Sunny | 2 | 3 | 2/9 | 3/5 |
| Overcast | 4 | 0 | 4/9 | 0/5 |
| Rainy | 3 | 2 | 3/9 | 2/5 |
| Total | | | 100% | 100% |

| Temperature | Yes | No | P(Yes) | P(No) |
|---|---|---|---|---|
| Hot | 2 | 2 | 2/9 | 2/5 |
| Mild | 4 | 2 | 4/9 | 2/5 |
| Cool | 3 | 1 | 3/9 | 1/5 |
| Total | | | 100% | 100% |

| Humidity | Yes | No | P(Yes) | P(No) |
|---|---|---|---|---|
| High | 3 | 4 | 3/9 | 4/5 |
| Normal | 6 | 1 | 6/9 | 1/5 |
| Total | | | 100% | 100% |

| Windy | Yes | No | P(Yes) | P(No) |
|---|---|---|---|---|
| False | 6 | 2 | 6/9 | 2/5 |
| True | 3 | 3 | 3/9 | 3/5 |
| Total | | | 100% | 100% |

$$P(A|B) = \frac{P(A \bigcap B)}{P(B)} = \frac{P(A) \cdot P(B|A)}{P(B)}$$