



Processamento de linguagem natural com Python

Coordenador
Fernando Ferreira



Institucional

Reitor

Prof. Eduardo Ramos

Eduardo.ramos@Infnet.edu.br

Coordenação do Curso

Prof. Fernando Ferreira

Fernando.gFerreira@prof.infnet.edu.br

Gerente Acadêmica

Ana Curi

ana.curi@Infnet.edu.br

Coordenador



Prof. Fernando Ferreira

fernando.gferreira@infnet.edu.br

<http://twitter.com/nandoferreira/>

<https://www.linkedin.com/in/nandoferreira/>

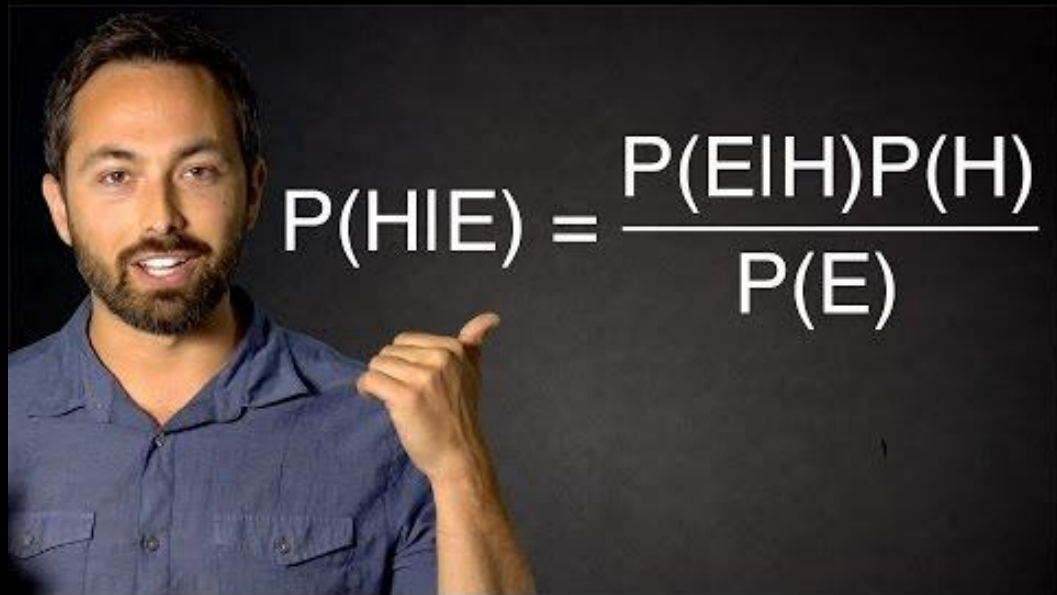
<https://www.facebook.com/fernandogferreira/>

Mini-bio

- Doutor pelo Programa de Engenharia Elétrica da Universidade Federal do Rio de Janeiro (Coppe-UFRJ) - ênfase em Inteligência Computacional.
- Sócio-fundador da empresa Twist (<https://twist.systems/>) - residente do Parque Tecnológico da UFRJ.
- Coordenador do MBA de Data Science e do MIT de Inteligência Artificial do Instituto Infnet.
- Professor de disciplinas de desenvolvimento ágil de software na graduação da Infnet. Professor do MBA de Business Analytics e Big Data da Fundação Getúlio Vargas.
- Membro da colaboração internacional CERN/UFRJ.
- [Twist] Empresa membro-nato do CDIA.Rio - Rede de Ciência de Dados & Inteligência Artificial do Rio de Janeiro

Exemplo de NLP: Classificação de Sentimentos



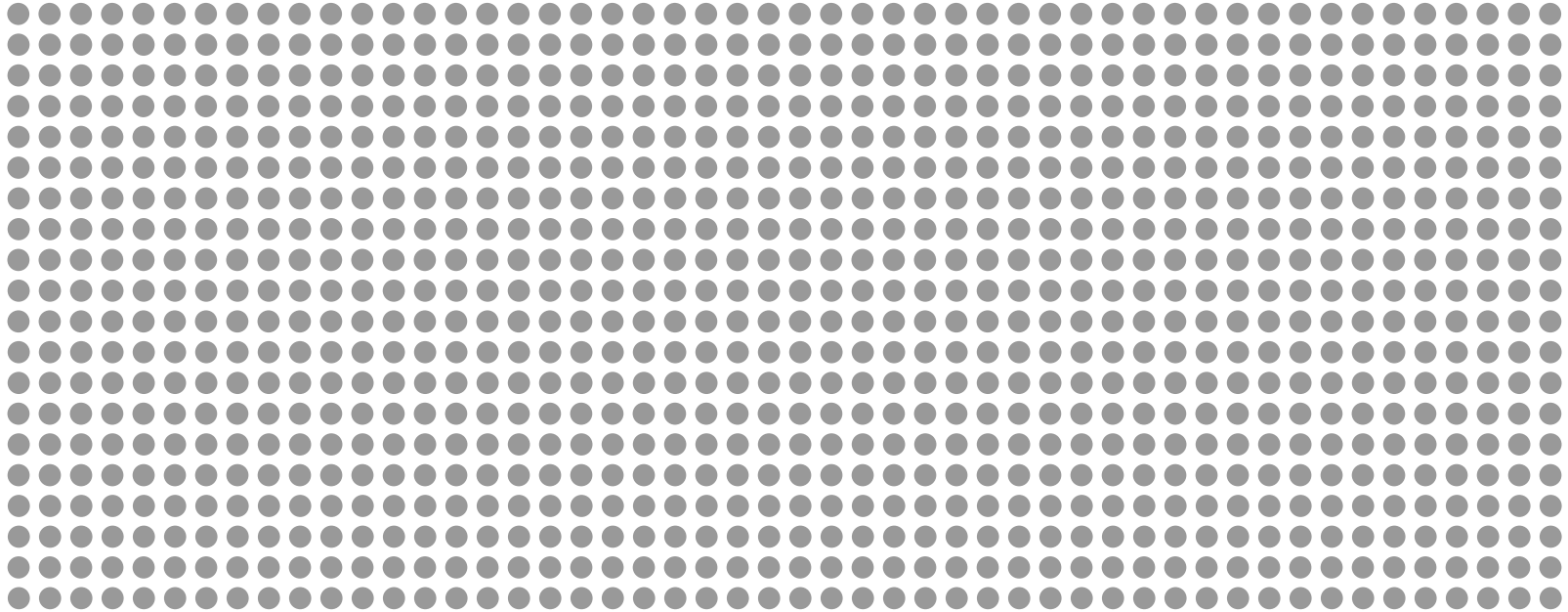


Teorema de Bayes - Caso

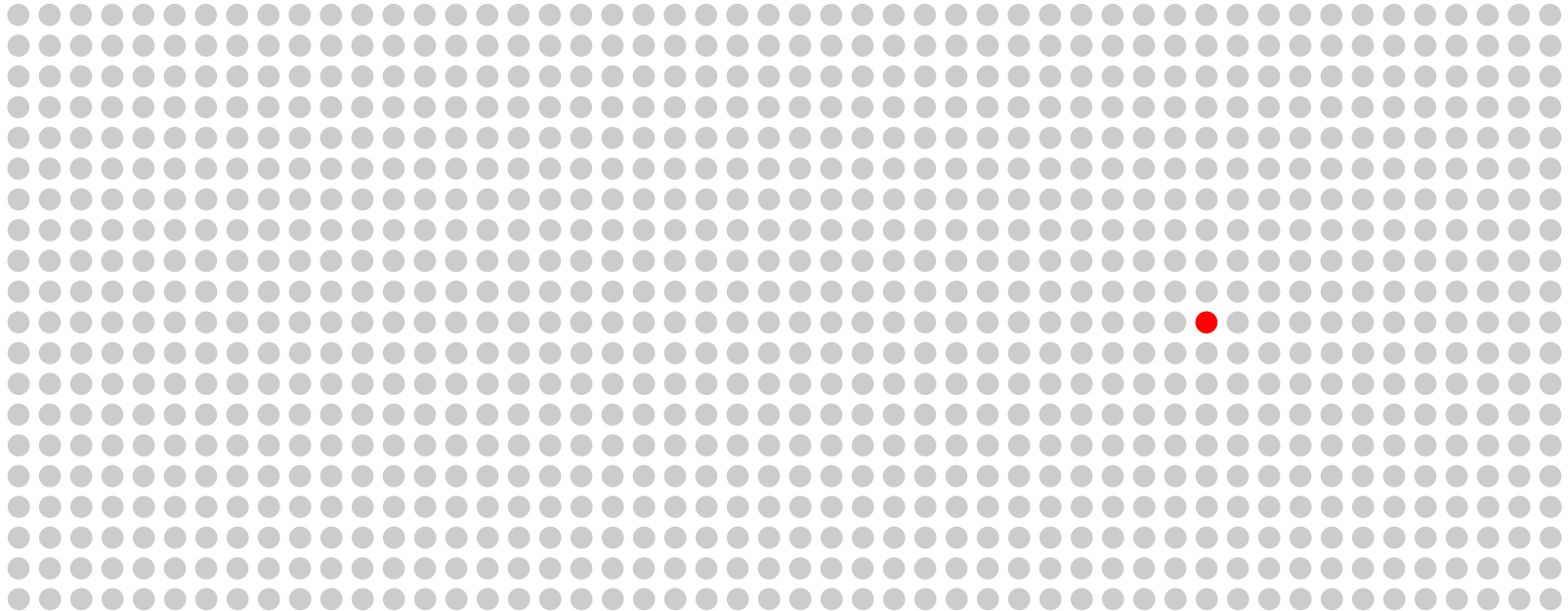
- Um teste diagnostica uma doença muito rara, que afeta 0.1% da população
- O teste identifica corretamente 99% das pessoas que possuem a doença
- O teste incorretamente acusa 1% das pessoas que não possuem a doença

Qual a chance da pessoa realmente ter a doença se teste der positivo?

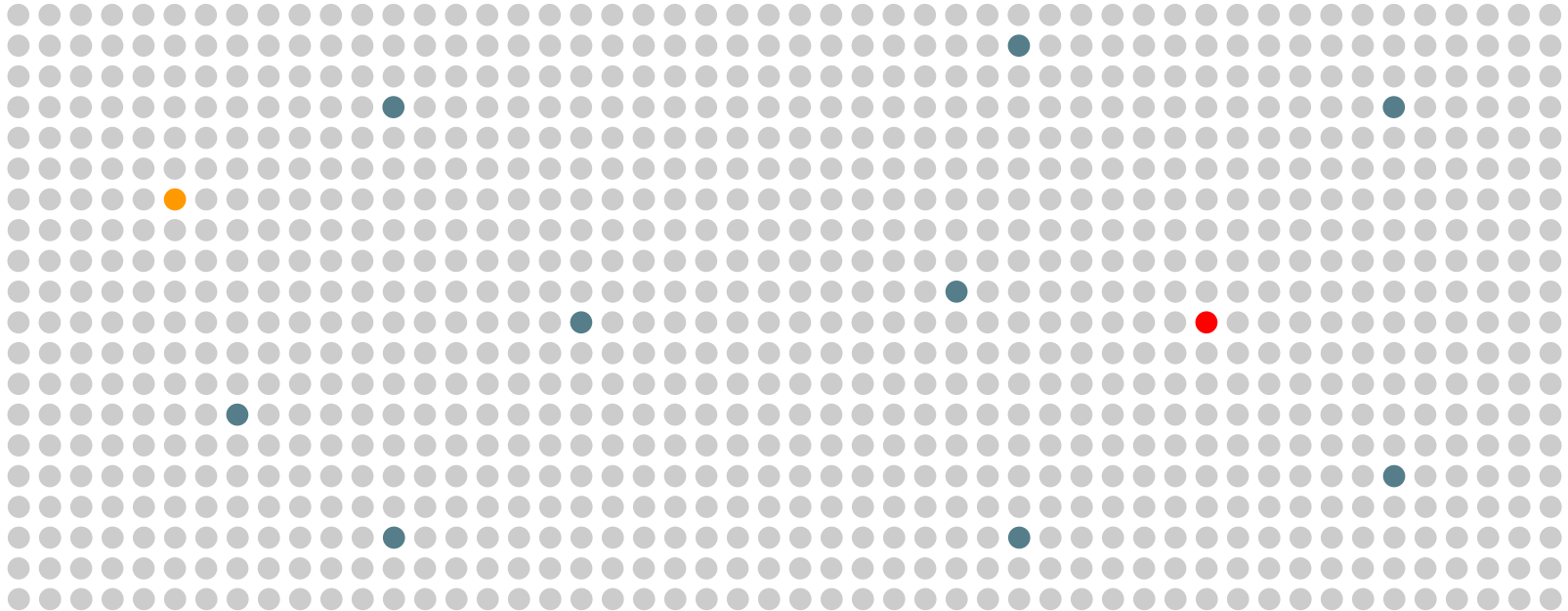
- 0.1%?
- 99%?
- 9%?



- A doença afeta 0.1 % da população



- A doença afeta 0.1% da população
- O teste incorretamente acusa 1% das pessoas que **não possuem** a doença



- A doença afeta 0.1% da população
- O teste incorretamente acusa 1% das pessoas que **não possuem** a doença
- De 11 pessoas que testaram positivo, **só uma** realmente tem a doença



Probabilidade da pessoa ter a doença, dado que o teste deu positivo é **9%**

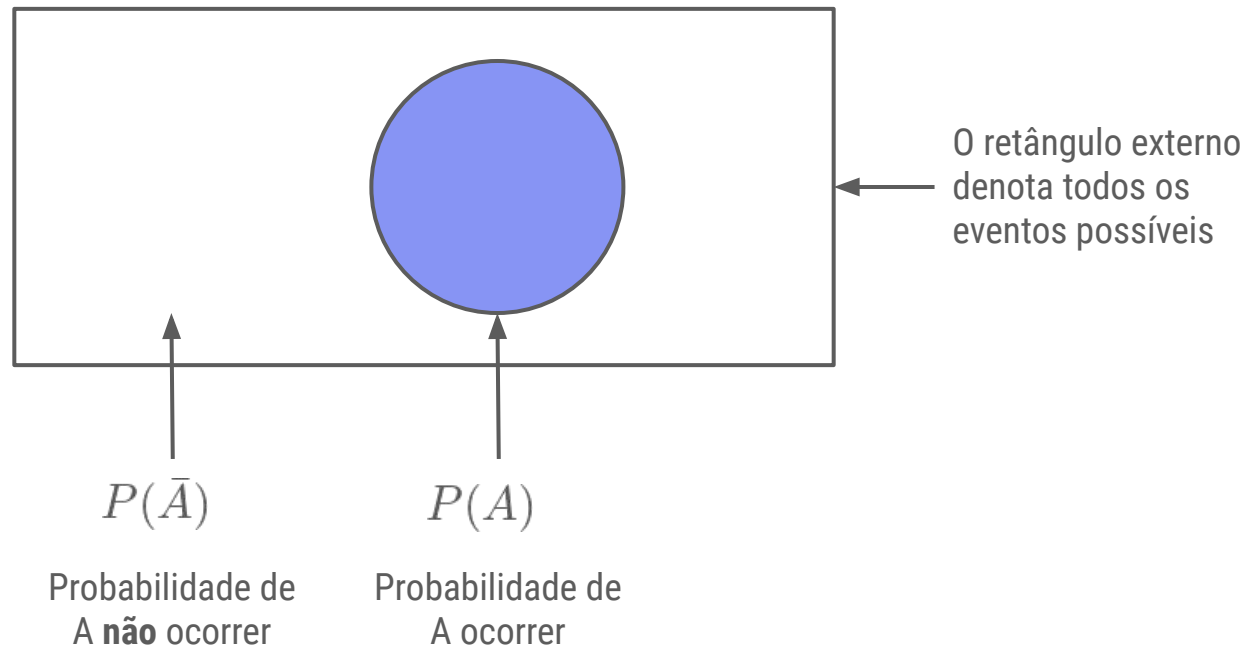
Teorema de Bayes - Caso

- Um teste diagnostica uma doença muito rara, que afeta **0.1% da população**
- O teste identifica corretamente **99%** da pessoas que **possuem** a doença
- O teste incorretamente acusa 1% das pessoas que **não possuem** a doença

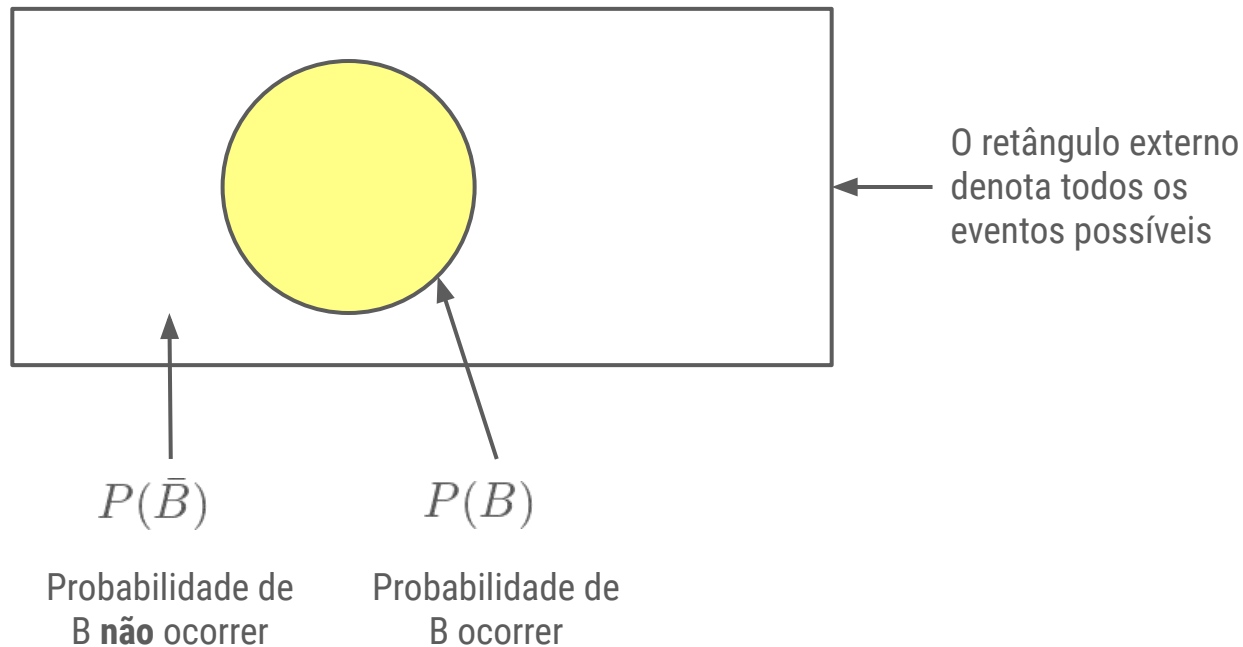
Qual a chance da pessoa realmente ter a doença?

- 0.1% - Pode ser usada como probabilidade *a priori*
- 99%?
- 9%

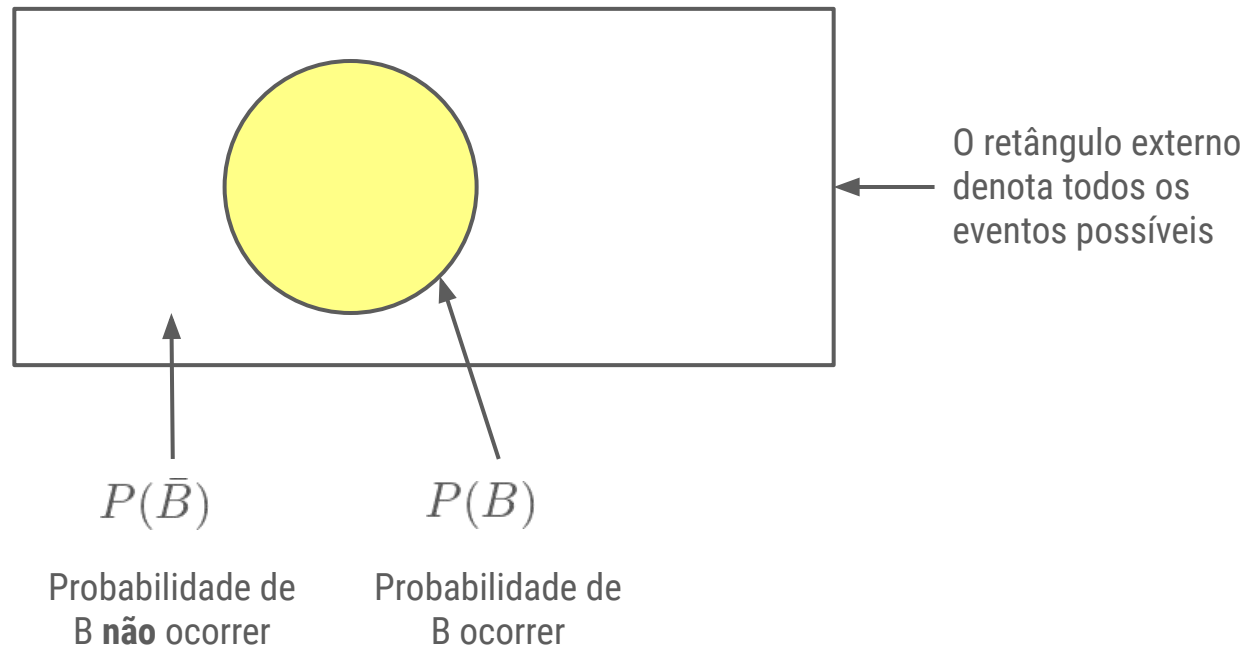
Probabilidades



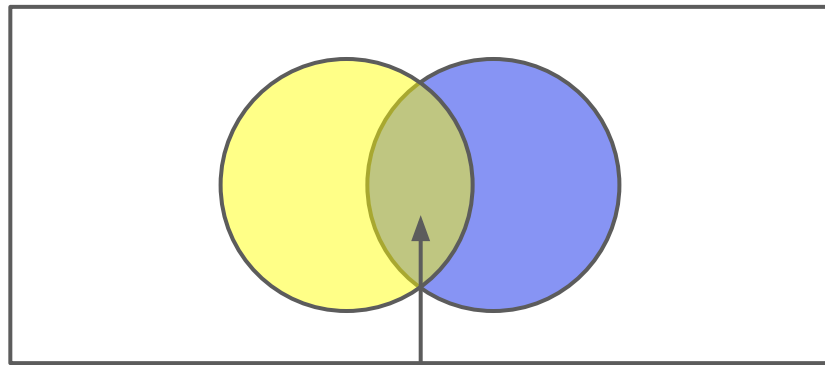
Probabilidades



Probabilidades



Probabilidade Conjunta



$$P(A \wedge B)$$

Probabilidade de
A e B ocorrerem ao
mesmo tempo

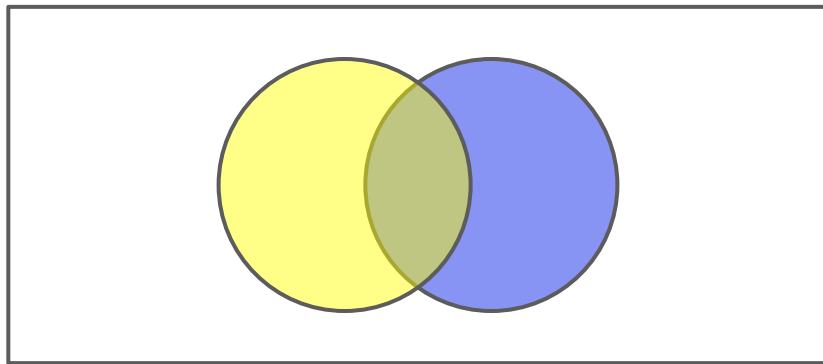
Se A e B forem **independentes**, temos

$$P(A \wedge B) = P(A)P(B)$$

Probabilidade Condicional



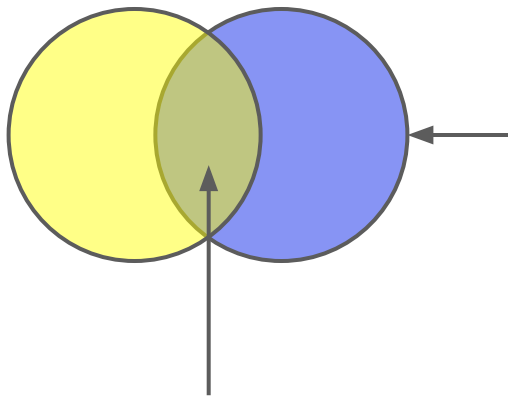
E se tivermos **evidência de que A ocorreu**, qual a probabilidade de B ocorrer?



Probabilidade Condicional

E se tivermos **evidência de que A ocorreu**, qual a probabilidade de B ocorrer?

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



Dado que A ocorreu, sabemos que o universo de eventos possíveis é o círculo azul

$P(B|A)$

Probabilidade de B ocorrer,
dado que A ocorreu

Teorema de Bayes



Probabilidade **a priori** para A, ou seja, o
que achávamos antes de observar B



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

Probabilidade **a posteriori** de A
ocorrer, dado que temos a
evidência B

Teorema de Bayes - Caso

- Um teste diagnostica uma doença muito rara, que afeta **0.1%** da população
- O teste identifica corretamente **99%** da pessoas que **possuem** a doença
- O teste incorretamente acusa 1% das pessoas que **não possuem** a doença

Qual a chance da pessoa realmente ter a doença?

$$P(\text{doente}|+) = \frac{P(+|\text{doente})P(\text{doente})}{P(+)}$$

Teorema de Bayes - Caso

- Um teste diagnostica uma doença muito rara, que afeta 0.1% da população
- O teste identifica corretamente 99% da pessoas que possuem a doença
- O teste incorretamente acusa 1% das pessoas que não possuem a doença

Qual a chance da pessoa realmente ter a doença?

$$P(\text{doente}|+) = \frac{P(+|\text{doente})P(\text{doente})}{P(+)}$$

$$P(\text{doente}|+) = \frac{P(+|\text{doente})P(\text{doente})}{P(\text{doente})P(+|\text{doente}) + P(\text{nao doente})P(+|\text{nao doente})}$$

$$P(\text{doente}|+) = \frac{0.99 \times 0.001}{0.001 \times 0.99 + 0.999 \times 0.01} = 9\%$$

Teorema de Bayes - Caso

- Um teste diagnostica uma doença muito rara, que afeta 0.1% da população
- O teste identifica corretamente 99% da pessoas que possuem a doença
- O teste incorretamente acusa 1% das pessoas que não possuem a doença

E se um segundo teste também der positivo?

$$P(\text{doente}|+) = \frac{P(+|\text{doente})P(\text{doente})}{P(+)}$$

$$P(\text{doente}|+) = \frac{P(+|\text{doente})P(\text{doente})}{P(\text{doente})P(+|\text{doente}) + P(\text{nao doente})P(+|\text{nao doente})}$$

$$P(\text{doente}|+) = \frac{0.99 \times 0.09}{0.09 \times 0.99 + 0.91 \times 0.01} = 90.7\%$$

Classificação de texto com teorema de Bayes



- Possíveis classificações
 - SPAM vs não SPAM
 - Sentimento positivo ou negativo
 - Relevante ou não relevante
- Evidências: presença de palavras (ou tokens)
- Ordem das palavras importa!
- Palavras **não são independentes!**

$$P(\text{SPAM} | w_1 \wedge w_2 \wedge \dots \wedge w_n) = \frac{P(w_1 \wedge w_2 \wedge \dots \wedge w_n | \text{SPAM}) P(\text{SPAM})}{P(w_1 \wedge w_2 \wedge \dots \wedge w_n)}$$

Classificador Naïve Bayes -

considerar que evidências são independentes



- Possíveis classificações
 - SPAM vs não SPAM
 - Sentimento positivo ou negativo
 - Relevante ou não relevante
- Evidências: presença de palavras (ou tokens)
- ~~• Ordem das palavras importa!~~
- ~~• Palavras não são independentes~~

$$P(\text{SPAM} | w_1 \wedge w_2 \wedge \dots \wedge w_n) = \frac{P(w_1 \wedge w_2 \wedge \dots \wedge w_n | \text{SPAM}) P(\text{SPAM})}{P(w_1 \wedge w_2 \wedge \dots \wedge w_n)}$$

$$P(\text{SPAM} | w_1 \wedge w_2 \wedge \dots \wedge w_n) = \frac{1}{Z} P(\text{SPAM}) \prod_i P(w_i | \text{SPAM})$$

Classificador Naïve Bayes - exemplo



Compra o remédio X, o segredo para emagrecer.	SPAM
Seu médico não quer que você saiba sobre X.	SPAM
Vamos confirmar a reunião de manhã.	Não-SPAM
Você vai agendar a visita ao médico?	Não-SPAM
Confirmei sua visita ao médico amanhã.	?

$$P(\text{SPAM}|w_1 \wedge w_2 \wedge \dots \wedge w_n) = \frac{1}{Z} P(\text{SPAM}) \prod_i P(w_i|\text{SPAM})$$