

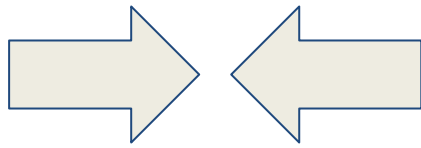
Métodos Matriciais e Análise de Clusters

Similaridade e Distância

Laura de Oliveira Fernandes Moraes

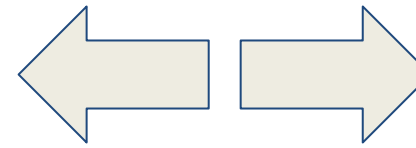
Por que usar?

- Sistemas de recomendação (Amazon, Netflix)
- Agrupamentos
- Classificação e regressão
- Detecção de anomalias/outliers
- Descoberta de casos parecidos (diagnóstico médico, precedentes legais)



Similaridade

- Aumenta se os objetos são **mais** parecidos.
- Normalmente entre **0 e 1**, onde 1 representa objetos idênticos.



Distância

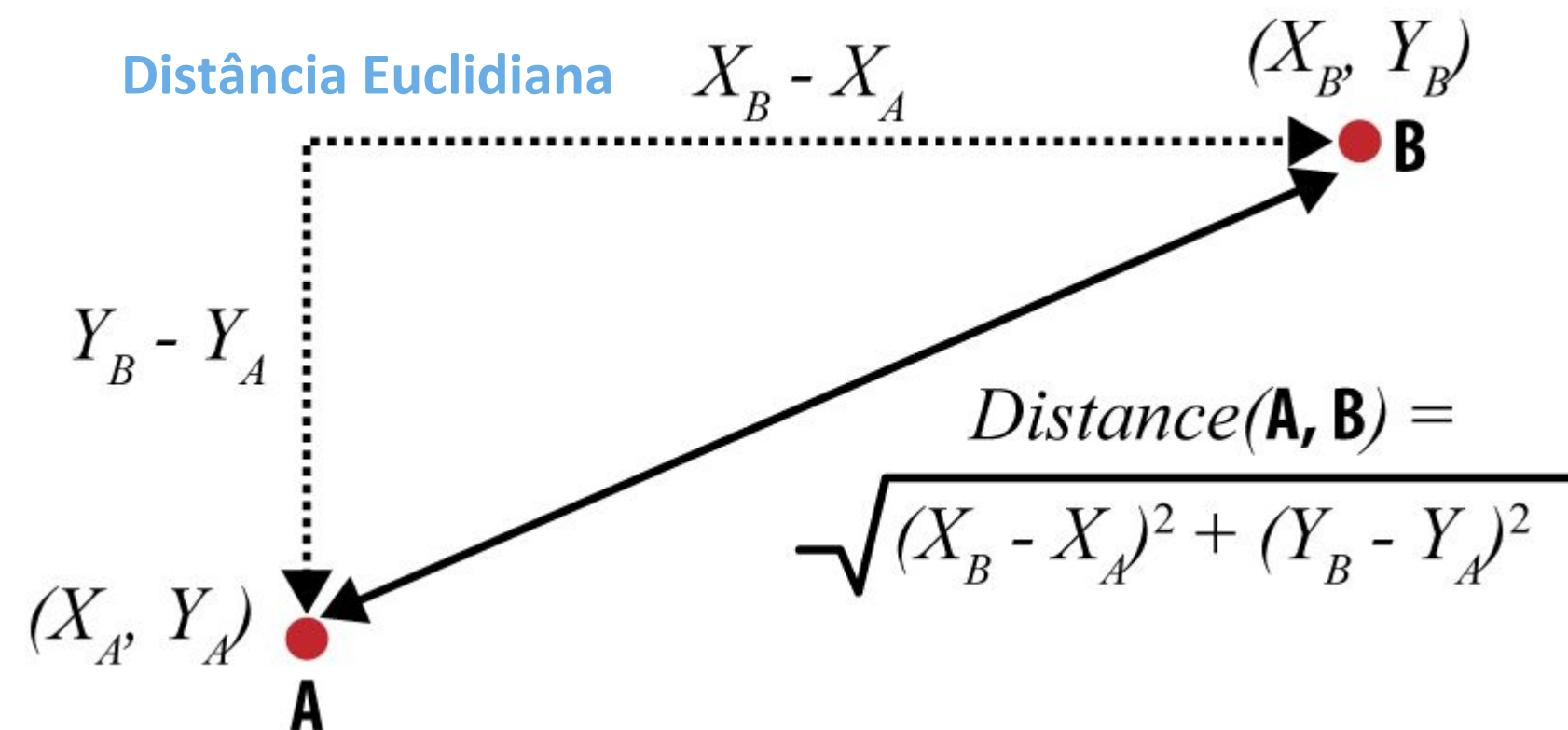
- Aumenta se os objetos são **menos** parecidos.
- Valor sempre será **maior ou igual a 0**, onde 0 representa objetos idênticos.

- Se dois objetos são representados por vetores, é possível calcular a distância entre eles:

Atributo	Pessoa A	Pessoa B
Idade	23	40
Anos no endereço atual	2	10
Estado residencial (1=Proprietário, 2=Inquilino, 3=Outro)	2	1

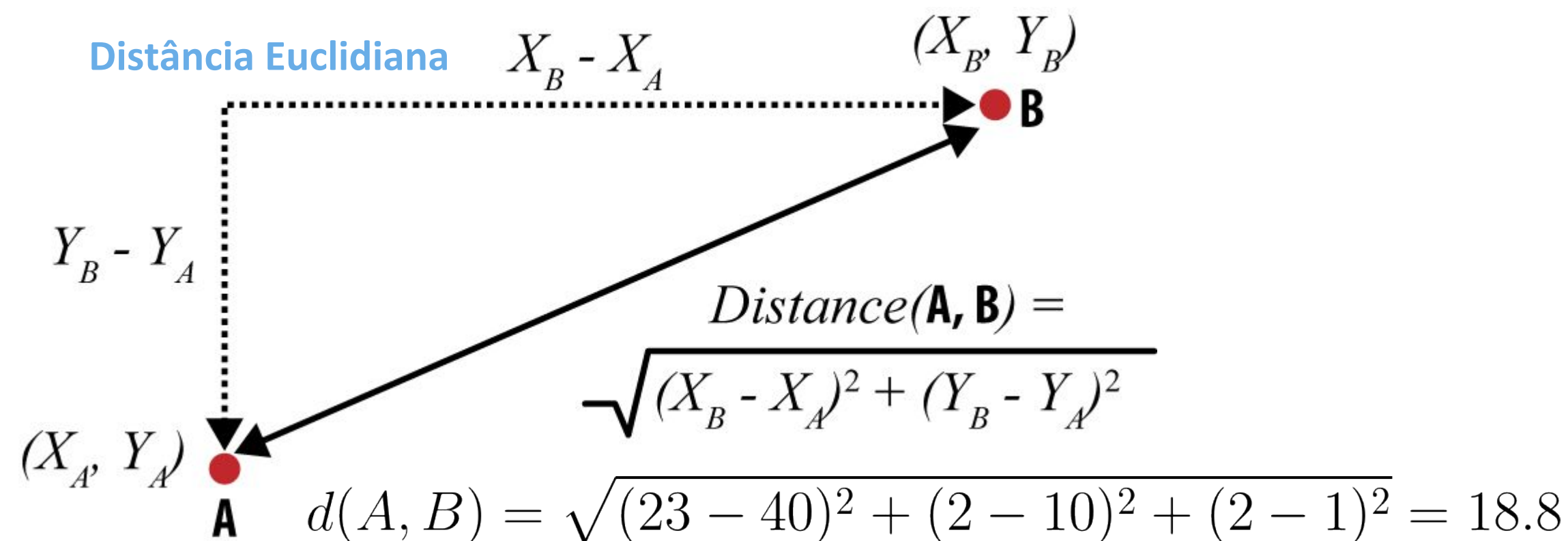
- Se dois objetos são representados por vetores, é possível calcular a distância entre eles:

Atributo	Pessoa A	Pessoa B
Idade	23	40
Anos no endereço atual	2	10
Estado residencial (1=Proprietário, 2=Inquilino, 3=Outro)	2	1



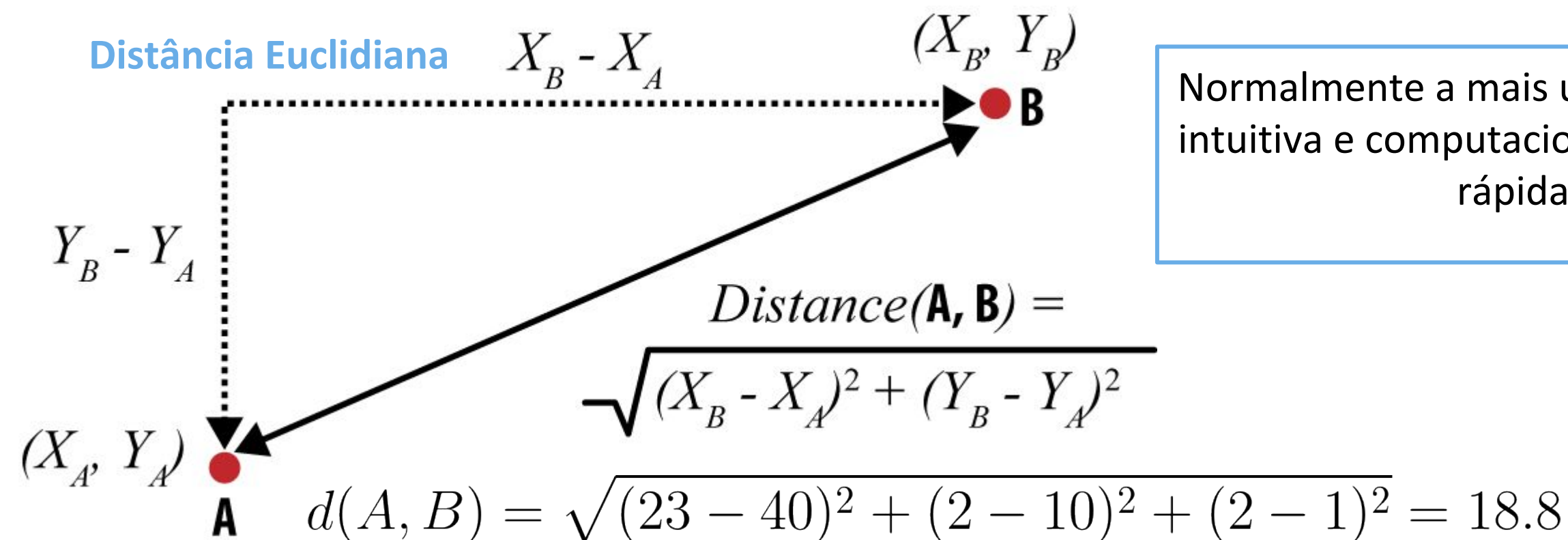
- Se dois objetos são representados por vetores, é possível calcular a distância entre eles:

Atributo	Pessoa A	Pessoa B
Idade	23	40
Anos no endereço atual	2	10
Estado residencial (1=Proprietário, 2=Inquilino, 3=Outro)	2	1



- Se dois objetos são representados por vetores, é possível calcular a distância entre eles:

Atributo	Pessoa A	Pessoa B
Idade	23	40
Anos no endereço atual	2	10
Estado residencial (1=Proprietário, 2=Inquilino, 3=Outro)	2	1



Normalmente a mais utilizada. É geral, intuitiva e computacionalmente muito rápida.

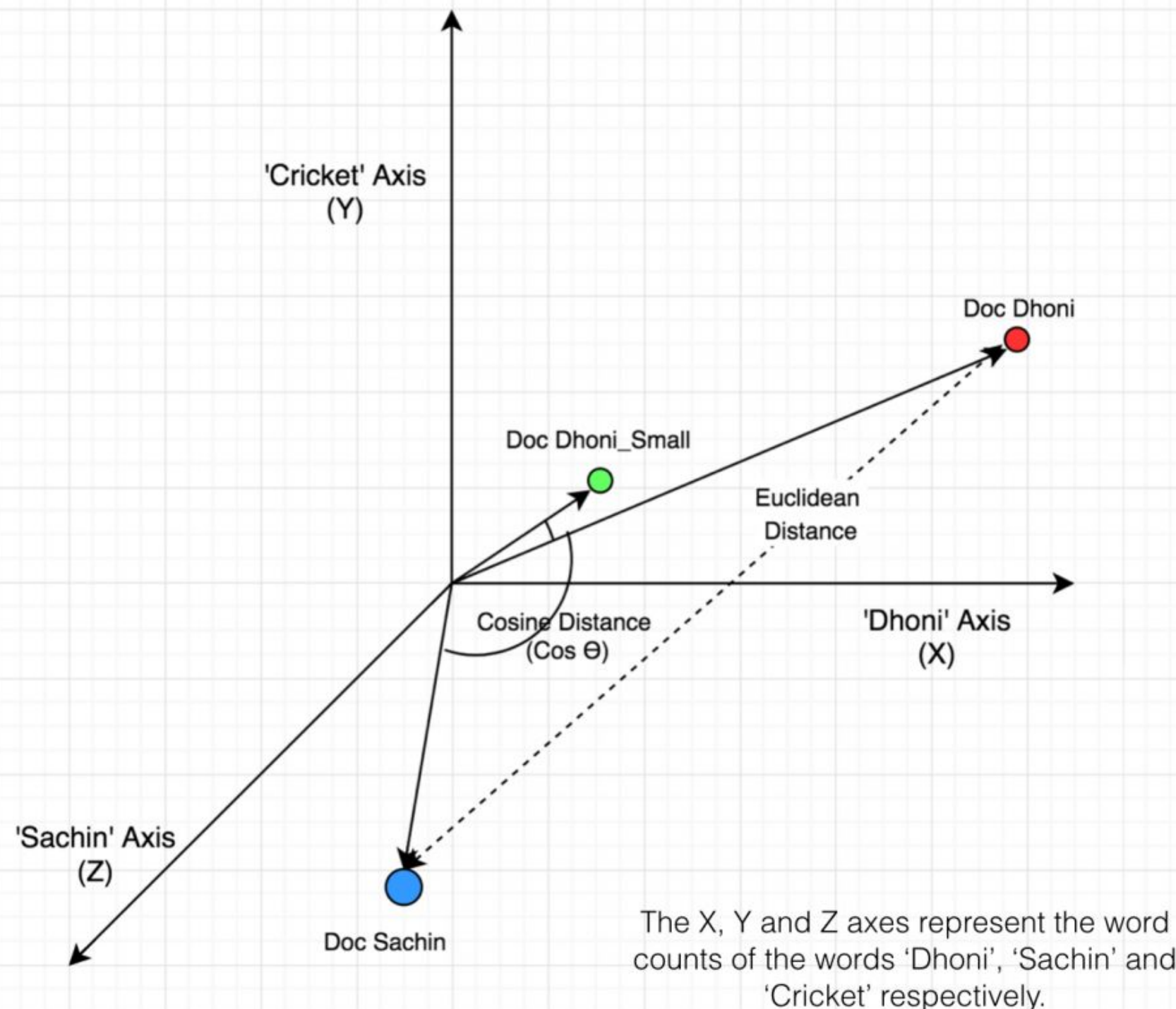
Outras funções de distância

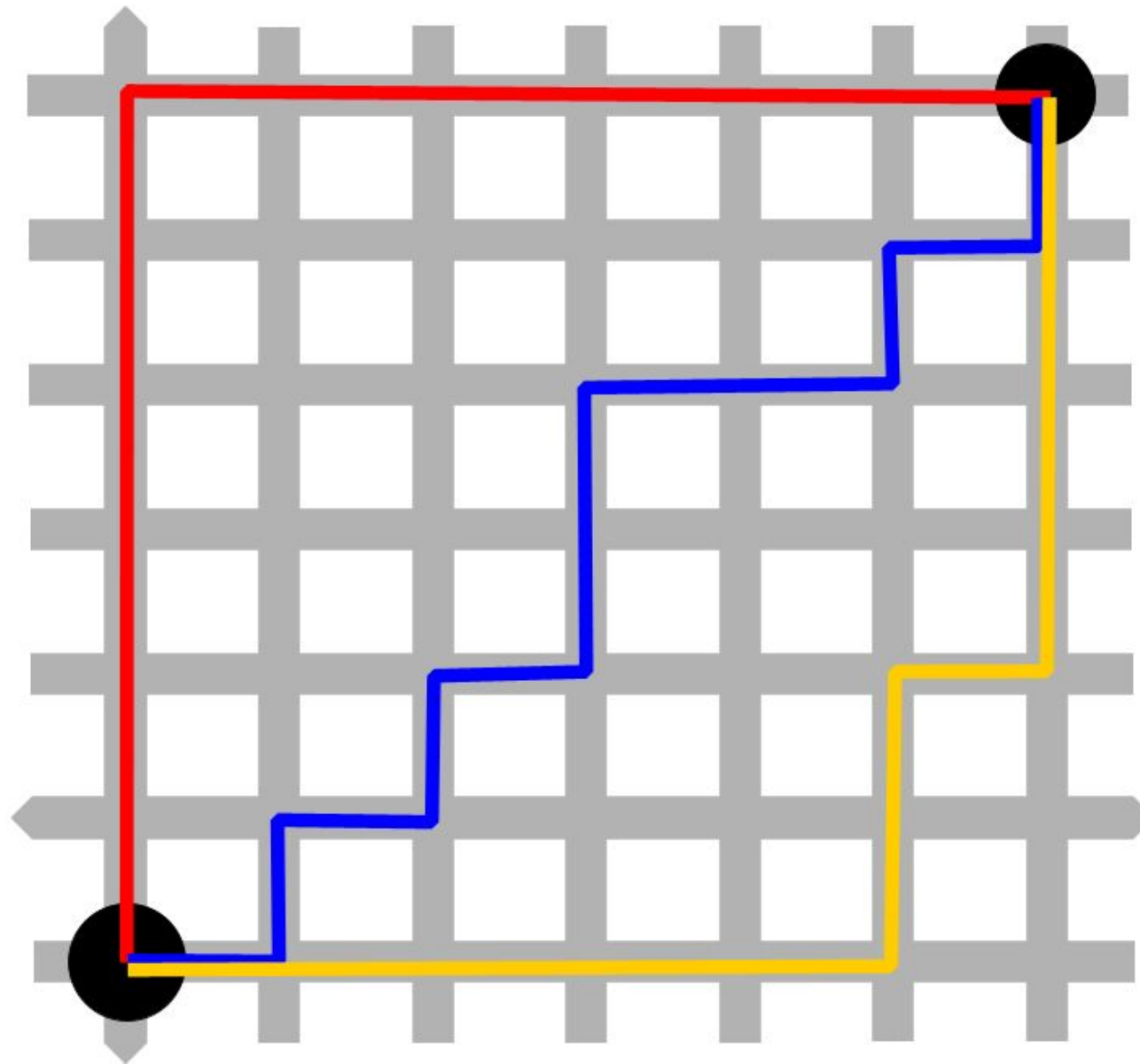
Projection of Documents in 3D Space

Distância de cosseno

$$d(A, B) = 1 - \frac{A \cdot B}{\|A\|_2 \cdot \|B\|_2}$$

Útil quando se quer ignorar diferenças de escala, como o tamanho dos textos





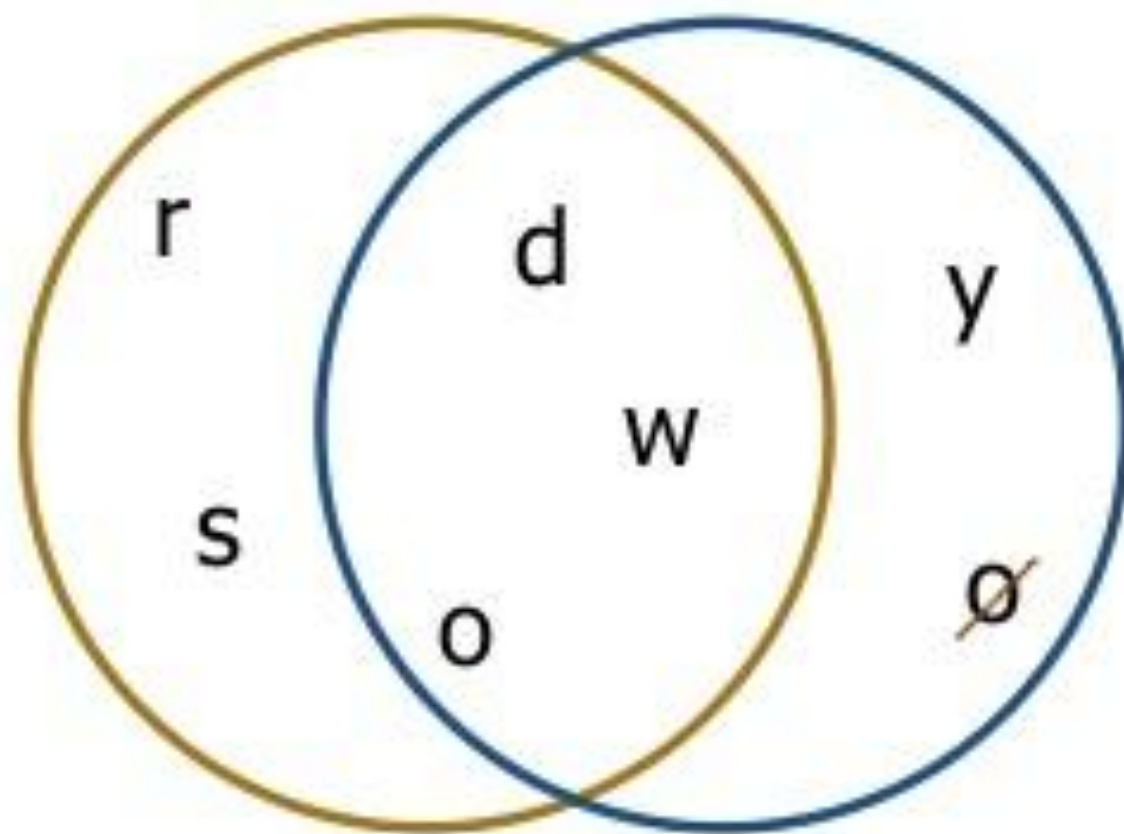
Distância de Manhattan

$$d(A, B) = \|A - B\|_1 = |x_1 - y_1| + |x_2 - y_2| + \dots$$

Outras funções de distância

Distância de Jaccard

"words" vs "woody"



Trata dois objetos como **conjuntos** de características. Adequado para problemas em que a posse de uma característica comum é importante, mas a ausência não.

$$J(A, B) = \frac{|A \cap B|}{|A \cup B|} = \frac{3}{7} = 0.43$$

Handwritten corrections: The denominator 7 is crossed out and replaced with 6. The result 0.43 is crossed out and replaced with 0.5.