

Sistemas Inteligentes Aplicados: Pré-processamento dos dados Parte 2

Prof. Arnaldo Candido Junior
UTFPR – Medianeira

Similaridade e dissimilaridade

- Similaridade: mede o quanto dois atributos são semelhantes
 - Quanto mais parecidos, maior o valor
 - Geralmente valor entre $[0, 1]$
- Dissimilaridade: mede o quanto dois atributos são diferentes
 - Quanto mais diferentes, maior o valor
 - Caso especial: medidas de distância

Similaridade e dissimilaridade ⁽⁴⁾

- Exemplos de similaridade s e dissimilaridade d medidas para dois valores v_1 e v_2 .
- Nominal:
 - $s = 1$ se $v_1 = v_2$
 0 se $v_1 \neq v_2$
 - $d = 0$ se $v_1 = v_2$
 1 se $v_1 \neq v_2$

Similaridade e dissimilaridade ₍₅₎

- Ordinal:

- $s = 1 - (|v_1 - v_2| / (n - 1))$

- $d = |v_1 - v_2| / (n - 1)$

- $n = \text{número de valores ordinais}$

Similaridade e dissimilaridade ⁽⁶⁾

- Intervalar ou razão:
 - $s_1 = 1 / (1 + |v_1 - v_2|)$ ou
 - $s_2 = 1 - d_2$
 - $d_1 = |v_1 - v_2|$
 - $d_2 = |v_1 - v_2| / \max(|v_i - v_j|)$
- Distâncias (a seguir) são boas medidas de dissimilaridade

Distância Euclidiana

- Medida clássica de distância
- Para duas instâncias \hat{u} e \hat{v} no R^n , a distância é dada por:

$$d(\hat{u}, \hat{v}) = \sqrt{\sum_{i=0}^n (u_i - v_i)^2}$$

Distância Euclidiana ⁽²⁾

- Exemplo:

- $\hat{u} = (2, 1)$

- $\hat{v} = (3, 4)$

- $$\begin{aligned} d &= \sqrt{(2 - 3)^2 + (1 - 4)^2} \\ &= \sqrt{(-1)^2 + (3)^2} \\ &= \sqrt{10} \end{aligned}$$

Distâncias de Minkowski

- Generalização da distância euclidiana
- r é um parâmetro escolhido pelo usuário

$$m_r(\hat{u}, \hat{v}) = \sqrt[r]{\sum_{i=0}^n (|u_i - v_i|)^r}$$

Distância de Minkowski

- Representa diferentes distâncias
 - m_1 : distância bloco cidade (Manhattan), distância de hamming (valores binários e strings)
 - m_2 : distância euclidiana

Distância de Minkowski

- m_{∞} : distância suprema (distância do eixo mais distante; distância dos atributos mais distantes)
 - Também chamada de distância quadrática ou distância de Chebyshev
 - $m_{\infty} = \max(|u_i - v_i|)$

Distância de Minkowski ⁽²⁾

- Exemplos para $\hat{u} = (1, 2)$ e $\hat{v} = (3, 5)$
 - $m_1 = 5,000$
 - $m_2 = 3,606$
 - $m_3 = 3,271$
 - $m_4 = 3,138$
 - $m_5 = 3,075$
 - ...
 - $m_\infty = 3,000 = \max(|u_i - v_i|)$

Avançado: Distância de Mahalanobis

- Generalização da distância euclidiana para atributos correlacionados que tenham distribuição normal
 - Leva em conta escala dos atributos
 - Leva em conta distribuição estatística das instâncias ao calcular distâncias
 - S : matriz de covariância (medida estatística obtida através da análise das instâncias)

Avançado: Distância de Mahalanobis ⁽²⁾

- Cálculo da esperança $E(X)$ para cada valor x_i com probabilidade p_i que o atributo X assume

$$E[X] = x_1p_1 + x_2p_2 + \cdots + x_kp_k .$$

- Cálculo da Matrix de covariância S para dois atributos X_i e X_j

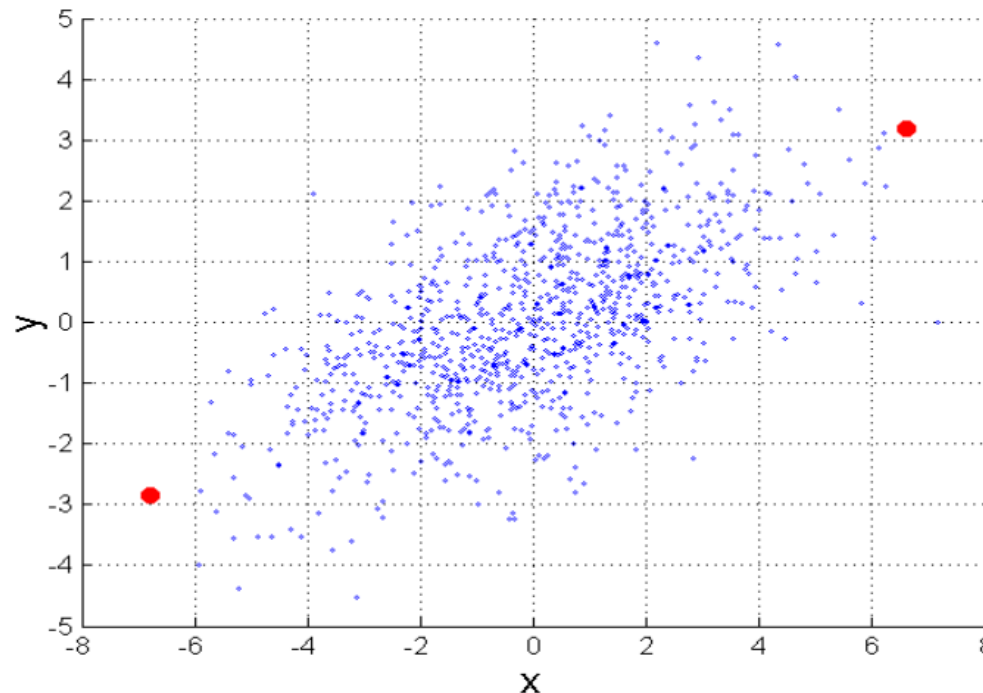
$$S_{i,j} = E[(x_i - E[x_i])(x_j - E[x_j])]$$

- Cálculo da distância de Mahalanobis

$$d(\hat{u}, \hat{v}) = \sqrt{(\hat{u} - \hat{v})^T S^{-1} (\hat{u} - \hat{v})}$$

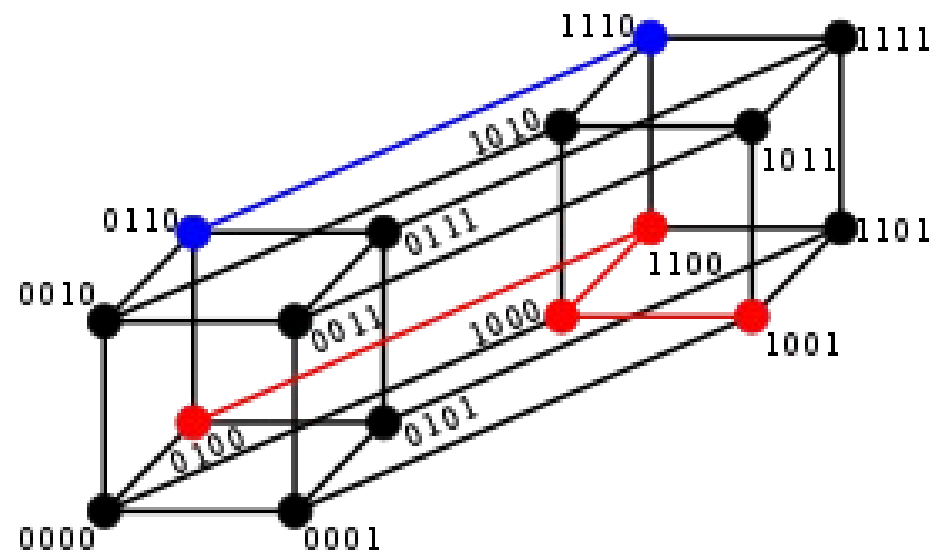
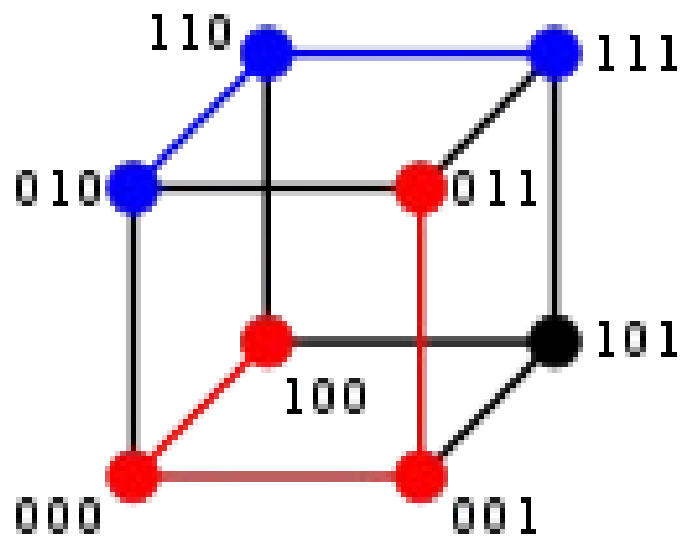
Avançado: Distância de Mahalanobis ⁽³⁾

- Para pontos vermelhos: distância euclidiana = 14.7; Mahalanobis = 6



Distância de Hamming

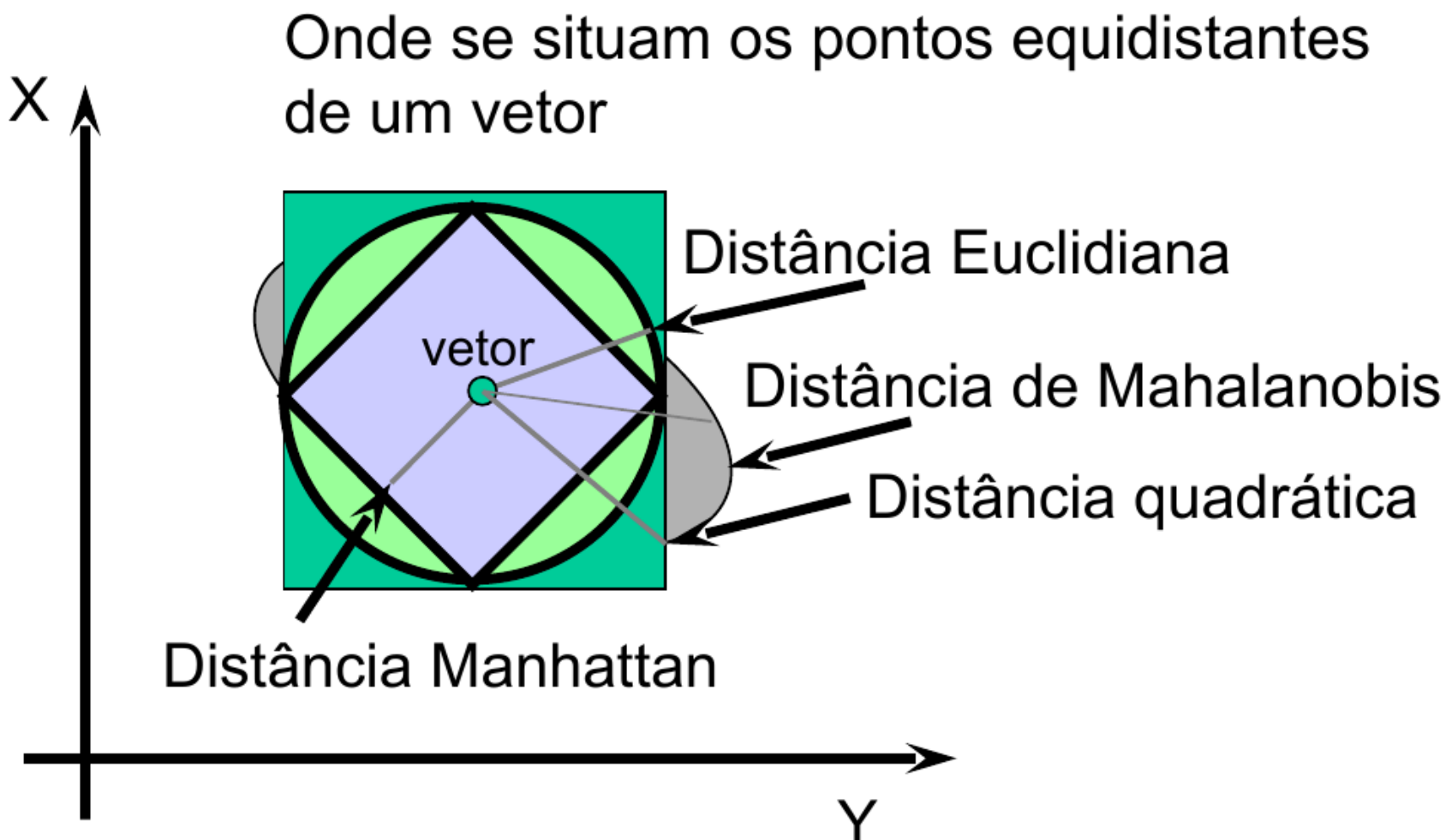
- Caso especial da distância Manhattan para valores binários (e também para Strings)
 - O cálculo é o mesmo



Distâncias

- Exercício
 - Para $\hat{u} = (1, 2, -3, 2)$ e $\hat{v} = (0, 6, 2, -1)$
 - Calcule as distâncias:
 - Manhattan
 - Euclidiana
 - Suprema

Propriedades das distâncias



Propriedades das distâncias ₍₂₎

- Medidas que satisfazem essas propriedades são chamadas **métricas**
 - $d(\hat{u}, \hat{v}) \geq 0 \quad \forall \hat{u}, \hat{v}$
 - $d(\hat{u}, \hat{v}) = 0 \leftrightarrow \hat{u} = \hat{v}$
 - $d(\hat{u}, \hat{v}) = d(\hat{v}, \hat{u})$
(simetria)
 - $d(\hat{u}, \hat{w}) \leq d(\hat{u}, \hat{v}) + d(\hat{v}, \hat{w})$
(desigualdade triangular)

Propriedades das distâncias ⁽³⁾

- Medidas de similaridade também possuem propriedades bem definidas
 - Seja $s(\hat{u}, \hat{v})$ a similaridade entre duas instâncias
 - $s(\hat{u}, \hat{v}) = 1 \leftrightarrow \hat{u} = \hat{v}$
 - $s(\hat{u}, \hat{v}) = s(\hat{v}, \hat{u})$

Conjuntos e vetores binários: (dis)similaridades

- Instâncias com apenas valores binários não são uma ocorrência incomum
- Conjuntos (exemplo, bag of words) podem ser mapeados para vetores binários
 - $X = \{A, B, C, D\}$
 - $Y = \{B, E, F\}$
 - $X \cup Y = \{A, B, C, D, E, F\}$
 - $X' = (1, 1, 1, 1, 0, 0)$
 - $Y' = (0, 1, 0, 0, 1, 1)$

Similaridade entre vetores binários

- Considerar duas instâncias originais \hat{u} e \hat{v}
 - $m_{0,0}$ número de atributos que possuem valor 0 tanto em \hat{u} quanto \hat{v}
 - $m_{1,0}$ número de atributos que possuem valor 1 em \hat{u} e valor 0 em \hat{v}
 - $m_{0,1}$ número de atributos que possuem valor 0 em \hat{u} e valor 1 em \hat{v}
 - $m_{1,1}$ número de atributos que possuem valor 1 tanto em \hat{u} quanto \hat{v}

Similaridade entre vetores binários ⁽²⁾

- **Coeficiente de casamento simples**

- $$CS = \frac{(m_{0,0} + m_{1,1})}{(m_{0,0} + m_{0,1} + m_{1,0} + m_{1,1})}$$

- **Coeficiente de Jaccard** (recomendado para vetores muito esparsos)

- $$j = m_{1,1} / (m_{0,1} + m_{1,0} + m_{1,1})$$

Similaridade entre vetores binários ⁽³⁾

- Exemplo

- $A = (1, 1, 1, 1, 0, 0)$

- $B = (0, 1, 0, 0, 1, 0)$

- $m_{0,0} = 1$

- $m_{0,1} = 1$

- $m_{1,0} = 3$

- $m_{1,1} = 1$

- $CS = (1 + 1) / (1 + 1 + 3 + 1) \cong 0.333$

- $j = 1 / (1 + 3 + 1) = 0.20$

Similaridade entre vetores binários ⁽⁴⁾

- Exercício
 - Calcular similaridade usando casamento simples e coeficiente de Jaccard
 - $A = 1 \ 0 \ 0 \ 1 \ 1 \ 0 \ 1$
 - $B = 0 \ 1 \ 0 \ 0 \ 1 \ 1 \ 0$

Similaridade do cosseno

- Bastante usada em mineração de textos
 - Atributos assimétricos e esparsos
 - Trata instâncias como vetores no espaço
 - Extrair cosseno
 - Ideia geral: instâncias cujos vetores apontam na mesma direção devem ser similares
 - $\cos(\hat{u}, \hat{v}) = (\hat{u} \cdot \hat{v}) / (||\hat{u}|| * ||\hat{v}||)$

Similaridade do cosseno ⁽²⁾

- Exemplo 1

- $\hat{u} = (2, 3)$

- $\hat{v} = (4, 3)$

- $|\hat{u} \cdot \hat{v}| = 2 * 4 + 3 * 3 = 17$

- $||\hat{u}|| = \sqrt{2^2 + 3^2} = \sqrt{13} \cong 3.6$

- $||\hat{v}|| = \sqrt{4^2 + 3^2} = \sqrt{25} = 5.0$

- $\cos(\hat{u}, \hat{v}) \cong 17 / (3.6 * 5.0) \cong 17/18$
 $\cong 0.94$

Similaridade do cosseno ⁽³⁾

- Exemplo 2

- $\hat{u} = (1, 1)$

- $\hat{v} = (3, 3)$

- $|\hat{u} \cdot \hat{v}| = 3 + 3 = 6$

- $||\hat{u}|| = \sqrt{1^2 + 1^2} = \sqrt{2}$

- $||\hat{v}|| = \sqrt{3^2 + 3^2} = \sqrt{18}$

- $\cos(\hat{u}, \hat{v}) = 6 / (\sqrt{2} * \sqrt{18}) = 6 / \sqrt{36}$
 $= 6 / 6 = 1$

Similaridade do cosseno ₍₄₎

- Exercício: calcular o cosseno para
 - $\hat{u} = (1, 0, 0, 4)$
 - $\hat{v} = (0, 5, 0, 2)$

Correlação

- Medida se dois atributos x e y estão linearmente relacionadas
 - Valores no intervalo [-1, 1]

$$\frac{\sum (x_i - \bar{x})(y_i - \bar{y})}{\sqrt{\sum (x_i - \bar{x})^2 \cdot \sum (y_i - \bar{y})^2}}$$

Correlação ₍₂₎

- Caso particular: podemos comparar a similaridade entre duas instâncias R^n :
- Isto é, ver se são proporcionais, semelhante a cosseno)
- Fazemos isso quebrando as instâncias originais em vetores do R^2

Correlação ₍₂₎

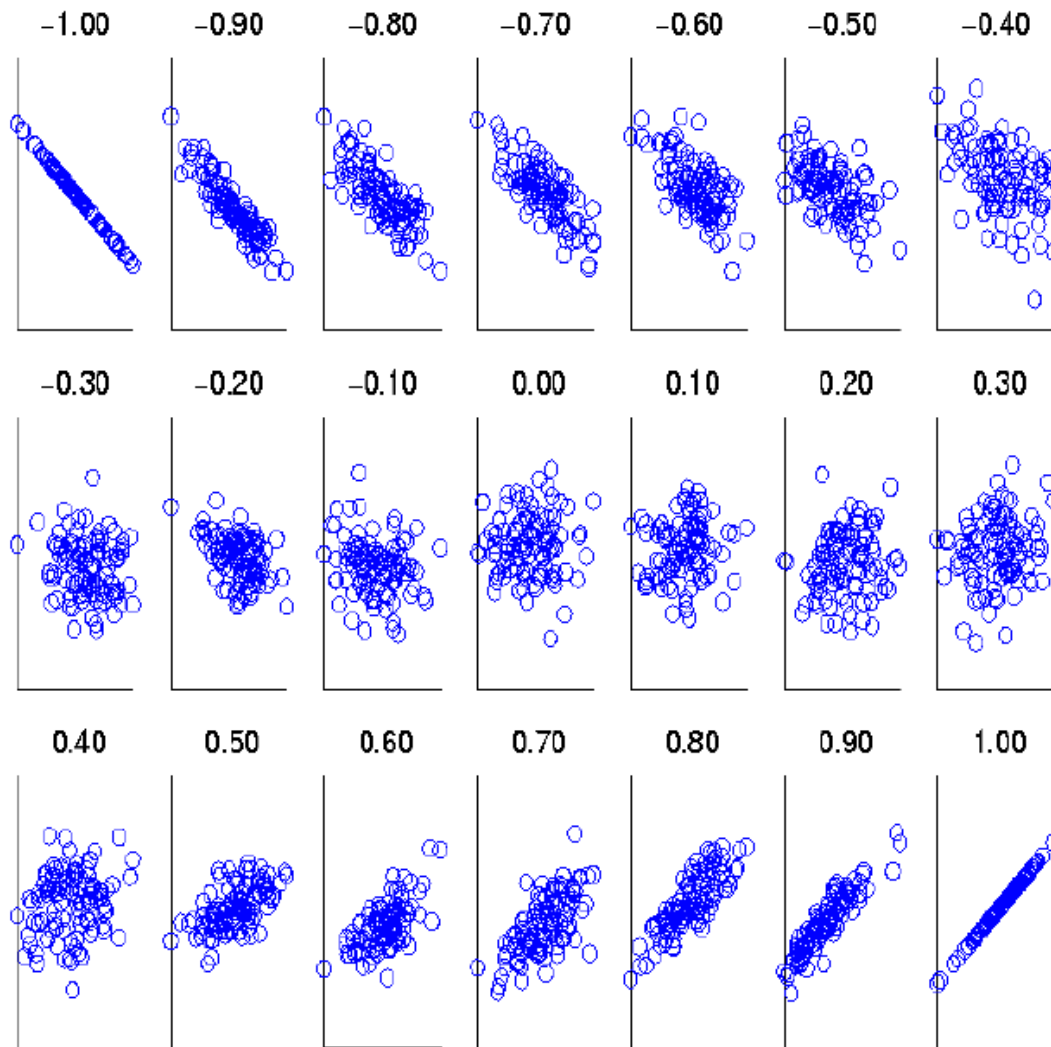
- Exemplo: $\hat{\mathbf{u}} = (0, 3, 4, -3)$; $\hat{\mathbf{v}} = (1, 1, 0, 0)$
 - $\hat{w}_1 = (0, 1)$; $\hat{w}_2 = (3, 1)$; $\hat{w}_3 = (4, 0)$; $\hat{w}_4 = (-3, 0)$;
 - Calcule a correlação:

$$\frac{\sum (\mathbf{x}_i - \bar{\mathbf{X}}) \cdot (\mathbf{y}_i - \bar{\mathbf{Y}})}{\sqrt{\sum (\mathbf{x}_i - \bar{\mathbf{X}})^2 \cdot \sum (\mathbf{y}_i - \bar{\mathbf{Y}})^2}}$$

Correlação ₍₃₎

- $c(\hat{u}, \hat{v}) = 1$
 - Relacionamento linear positivo perfeito
- $c(\hat{u}, \hat{v}) = -1$
 - Relacionamento linear negativo perfeito
- $c(\hat{u}, \hat{v}) = 0$:
 - Não existe relacionamento
- Relacionamento linear: $u_i = av_i + b$

Correlação ₍₄₎



- Comparando instâncias \hat{u} e \hat{v} (eixos)
- Pontos: valor de um atributo em um \hat{u} e em \hat{v}

Exercícios vistos em aula

- 1. Para $\hat{u} = (1, 2, -3, 2)$ e $\hat{v} = (0, 6, 2, -1)$, calcule as distâncias: euclidiana; Manhattan; suprema
- 2. Calcular similaridade usando casamento simples e coeficiente de Jaccard
$$A = \begin{matrix} & 1 & 0 & 0 & 1 & 1 & 0 & 1 \end{matrix}$$
$$B = \begin{matrix} & 0 & 1 & 0 & 0 & 1 & 1 & 0 \end{matrix}$$
- 3. Calcular o cosseno para
 $u = (1, 0, 0, 4)$
 $v = (0, 5, 0, 2)$

Exercício

- 4. Correlação (visto em sala):
 $\hat{u} = (0, 3, 4, -3)$; $\hat{v} = (1, 1, 0, 0)$
- 5. No weka, rodar KNN para base Iris

Pontos chaves

- Escalas
- Similaridade: cosseno, casamento simples, Jaccard
- Dissimilaridade (todas as distâncias)

Agradecimentos/referências

- Notas de aula do Prof. André de Carvalho (USP)