

DCE 88

TEMA 02 - APRENDIZADO DE MÁQUINA NÃO-SUPERVISIONADO

01 - INTRODUÇÃO

O APRENDIZADO DE MÁQUINA NÃO-SUPERVISIONADO, TAMBÉM CONHECIDO COMO MODELOS DESCRITIVOS, SÃO UMA TÉCNICA DE INTELIGÊNCIA ARTIFICIAL (IA) ONDE SUA PRINCIPAL CARACTERÍSTICA É TRABALHAR COM DADOS NÃO ROTULADOS. DIFERENTEMENTE DO APRENDIZADO SUPERVISIONADO, ONDE TODA A BASE DE DADOS ESTÁ DEVIDAMENTE ROTULADA, NO APRENDIZADO NÃO-SUPERVISIONADO ESSE PROCESSO DE MARCAÇÃO DE BASES PODE SER CARO, TANTO EM RECURSO HUMANO QUANTO EM TEMPO, NO CASO POR EXEMPLO DE MARCAÇÃO DE BASES DE DADOS DE IMAGENS, OU MUITAS VEZES IMPRATICÁVEL, COMO NO CASO DE RECOMENDAÇÃO DE NOVAS MÚSICAS BASEADO NO QUE UM USUÁRIO OUVI EM UM APLICATIVO DE STREAMING DE MÚSICAS.

O APRENDIZADO DE MÁQUINA NÃO-SUPERVISIONADO BUSCA, ATRAVÉS DE SUAS ^{TÉCNICAS} ALGORITMOS:

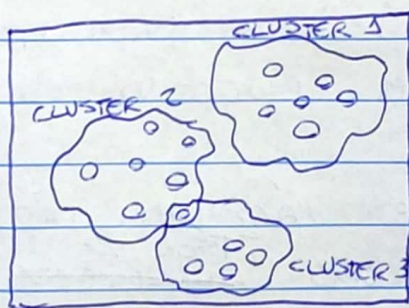
- 1- ENCONTRAR SIMILARIDADES ENTRE OS DADOS ATRAVÉS DE AGRUPAMENTO EM CLUSTERS, OU SEJA, AGRUPAR DADOS QUE POSSUAM CARACTERÍSTICAS SEMELHANTES;
- 2- ASSOCIAR CONJUNTOS DE DADOS QUE POSSUAM UMA FREQUÊNCIA DE USO EM CONJUNTO, POR EXEMPLO, QUEM COMPRO WHEY PROTEIN, FREQUENTEMENTE COMPRO CREATINA;
- 3- REDUZIR A DIMENSIONALIDADE DE UM GRANDE CONJUNTO DE DADOS, SIMPLIFICANDO-O PARA QUE O MESMO POSSA SER ANALISADO.

02 - TÉCNICAS ALGORITMOS DE APRENDIZADO NÃO-SUPERVISIONADO

AS TÉCNICAS ALGORITMOS DE APRENDIZADO NÃO-SUPERVISIONADO SÃO DIVIDIDOS EM TRÊS CLASSES: AGRUPAMENTO, ASSOCIAÇÃO E REDUÇÃO DE DIMENSIONALIDADE. NAS PRÓXIMAS SEÇÕES SERÃO APRESENTADOS OS DETALHES DE CADA CLASSE.

02.1 - TÉCNICAS ALGORITMOS DE AGRUPAMENTO

AS TÉCNICAS DE AGRUPAMENTO DE DADOS VISAM AGRUPAR EM CLUSTERS CONJUNTOS DE PONTOS QUE SÃO SIMILARES ENTRE SI E QUE DIFEREM DE OUTRO CONJUNTO DE PONTOS. NA IMAGEM ABAIXO É APRESENTADO UM EXEMPLO DE AGRUPAMENTO EM CLUSTER:



AS TÉCNICAS DE AGRUPAMENTO SÃO DIVIDIDAS EM AGRUPAMENTOS EXCLUSIVOS, SOBREPOSTOS, HIERÁRQUICOS E PROBABILÍSTICOS. NOS AGRUPAMENTOS EXCLUSIVOS, UM PONTO (DADO) PODE PERTENCER SÓ E SOMENTE SÓ A UM ÚNICO CLUSTER. O ALGORITMO FREQUENTEMENTE UTILIZADO É O K-MEANS, ONDE INICIALMENTE SÃO CRIADOS K CLUSTERS ALEATÓRIOS, O CENTROÍDE MÉDIO DE CADA CLUSTER É CALCULADO E OS CLUSTERS ENTÃO SÃO ATUALIZADOS ATRAVÉS DE UMA FUNÇÃO DE DISTÂNCIA, GERAL

MENTE A DISTÂNCIA EUCLIDIANA OU A DISTÂNCIA MANHATAN, ENTRE A DISTÂNCIA DE CADA PONTO PARA O ~~CLUSTRO~~ CENTROÍDE MÉDIO DE CADA CLUSTER. OS AGRUPAMENTOS SOBREPPOSTOS SE DIFEREM DOS EXCLUSIVOS POIS PERMITE QUE UM MESMO PONTO PERTENÇA A MAIS DE UM CLUSTER COM NÍVEIS DE PERTENCIMENTO DISTINTOS.

OS AGRUPAMENTOS HIERÁRQUICOS ATUALIZAM OS PONTOS ENTRE OS CLUSTERS ATRAVÉS DE REGRAS DE LIGAÇÃO ENTRE ^{CADA DOIS PONTOS} ~~OS PONTOS~~ DO CONJUNTO DE DADOS. AS LIGAÇÕES PRINCIPAIS SÃO:

- LIGAÇÃO MÉDIA: É A MÉDIA DA DISTÂNCIA ENTRE DOIS PONTOS
- LIGAÇÃO MÁXIMA: É A DISTÂNCIA MÁXIMA ENTRE DOIS PONTOS
- LIGAÇÃO MÍNIMA: É A DISTÂNCIA MÍNIMA ENTRE DOIS PONTOS

POR FIM, OS AGRUPAMENTOS PROBABILÍSTICOS ATUALIZAM OS PONTOS ENTRE OS CLUSTERS ATRAVÉS DA PROBABILIDADE DE UM DETERMINADO DADO PERTENCER A UM CLUSTER. NORMALMENTE A DISTRIBUIÇÃO GAUSSIANA É O MÉTODO UTILIZADO PARA NORMALIZAR OS DADOS.

02.2 - TÉCNICAS DE ASSOCIAÇÃO

AS TÉCNICAS DE ASSOCIAÇÃO LIGAM CONJUNTOS DE DADOS QUE ~~SÃO~~ ^{POSSUEM} UMA FREQUÊNCIA DE USO EM CONJUNTO. ESSAS TÉCNICAS SÃO PARTICULARMENTE VALIOSAS NO RAMO DE SUGESTÕES PARA USUÁRIOS TANTO EM SITES DE E-COMMERCE QUANTO STREAMING, OU SEJA, O PRÓPRIO USUÁRIO, ATRAVÉS DE SUAS COMPRAS OU REPRODUÇÕES AJUDAM A CONSTRUIR UM MODELO

DE IA QUE DARÁ RECOMENDAÇÕES PARA TODOS OS USUÁRIOS DA PLATAFORMA. SÃO EXEMPLOS DE RECOMENDAÇÕES O "QUEM COMPROU X TAMBÉM COMPRO Y" DA AMAZON E TAMBÉM AS PLAYLISTS "RECOMMENDED FOR YOU" DO SPOTIFY.

O ALGORITMO MAIS UTILIZADO NAS TÉCNICAS DE ASSOCIAÇÃO É O APRIORI. ESTE PARTE DO PRINCÍPIO QUE SE UM ITEM É FREQUENTE, TODOS OS SEUS SUBITEMS TAMBÉM SÃO FREQUENTES E ASSIM O MODELO É CONSTRUÍDO E É ATUALIZADO.

02.3 - TÉCNICAS DE REDUÇÃO DE DIMENSIONALIDADE

AS TÉCNICAS DE REDUÇÃO DE DIMENSIONALIDADE SÃO PARTICULARMENTE ÚTEIS EM CASOS ONDE A BASE DADOS É MUITO GRANDE OU REDUNDANTE. ESSAS TÉCNICAS SÃO FREQUENTEMENTE UTILIZADAS EM PRÉ-PROCESSAMENTO DE DADOS PARA UTILIZAÇÃO EM OUTROS MODELOS DE APRENDIZAGEM OU ATÉ MESMO OUTRAS APLICAÇÕES. UM EXEMPLO COTIDIANO ONDE SÃO APLICADAS TÉCNICAS DE REDUÇÃO DE DIMENSIONALIDADE SÃO OS "TRENDING TOPICS" DO SITE X (ANTIGO TWITTER). CADA TWEET É UM CONJUNTO DE DADOS ~~FA~~ JUNTAMENTE COM ~~TE~~ DE UMA VASTA BASE DE USUÁRIOS QUE, SEM A APLICAÇÃO DA REDUÇÃO DE DIMENSIONALIDADE, É IMPRATICÁVEL TRATAR OS DADOS. JUNTAMENTE COM TÉCNICAS DE PROCESSAMENTO DE LINGUAGEM NATURAL, AS PALAVRAS MAIS FALADAS DO MOMENTO CONSTITUEM O "TRENDING TOPICS".

O ALGORITMO MAIS CONHECIDA É O DE ANÁLISE DOS COMPONENTES PRINCIPAIS (PCA). NELE, CADA COMPONENTE

TE EXTRAÍDA DO CONJUNTO DE DADOS REPRESENTAM A VARIACÃO MÁXIMA DOS DADOS. CADA COMPONENTE É COMPLETAMENTE DIFERENTE UMA DAS OUTRAS.

03- Considerações Finais

NESTE TEXTO ~~FORAM~~ ^{FOI} APRESENTADO O APRENDIZADO DE MÁQUINA NÃO-SUPERVISIONADO E SUAS PRINCIPAIS TÉCNICAS. QUANDO A MARCAÇÃO DE UMA BASE DE DADOS É CUSTOSA, AS TÉCNICAS NÃO-SUPERVISIONADAS PODEM SER DE GRANDE UTILIDADE, INCLUSIVE SERVINDO DE PRÉ-PROCESSAMENTO DE DADOS PARA OUTRAS APLICAÇÕES DE INTELIGÊNCIA ARTIFICIAL, POR EXEMPLO, AS REDES NEURAIS.

AS TÉCNICAS NÃO-SUPERVISIONADAS SÃO DE GRANDE VALOR PARA EMPRESAS DE E-COMMERCE E STREAMING POIS AS PREFERÊNCIAS DE COMPRAS E GOSTOS DOS USUÁRIOS MUDAM CONSTANTEMENTE, DIFERENTE DE UMA BASE DE DADOS FIXAS, DE IMAGENS DE MAÇÃS E BANANAS.

É FUNDAMENTAL ~~TER~~ ^{SER} CRITERIOSO COM AS CONDIÇÕES DE PARADA DOS ALGORITMOS DE APRENDIZADO NÃO-SUPERVISIONADO, UMA VEZ QUE NESSE MODELO NÃO É POSSÍVEL APLICAR AS MÉTRICAS DE AVALIAÇÃO DE DADOS CONVENCIONAIS, COMO PRECISION, RECALL OU F1-SCORE.

O APRENDIZADO DE MÁQUINA NÃO-SUPERVISIONADO POSSUI UMA VASTA GAMA DE APLICAÇÕES. QUANDO BEM UTILIZADO É DE GRANDE VALIA NO ESPECTRO DA INTELIGÊNCIA ARTIFICIAL.