

## Sumário

	Pág
1. Introdução	1
2. Aprendizado de máquina não-supervisionado	2
2.1 Tarefas	2
2.1.1 Agrupamento	3
2.1.2 Redução de dimensionalidade	4
2.1.3 Associação	5
3. Ciclo de vida do processo de aprendizagem de máquina não-supervisionado	6
4. Desafios e limitações	7
5. Tendências	8
6. Conclusão	8

Aprendizado de máquina não supervisionado

## 1. Introdução

Aprendizado de máquina é o ramo da inteligência artificial que permite que os sistemas aprendam sem que sejam explicitamente programados. Ela envolve a construção e treinamento de modelos (algoritmos), capazes de resolver tarefas complexas, sem que as regras de resolução das tarefas, sejam explicitamente colocadas nos códigos destes modelos.

A classificação geral do aprendizado de máquina distingue duas categorias: aprendizado de máquina supervisionado e o aprendizado de máquina não-supervisionado. Outras classificações incluem o aprendizado semi-supervisionado, aprendizado por reforço. O aprendizado de máquina supervisionado distingue-se do não-supervisionado pelos rótulos (ou respostas) dos dados de treinamento, enquanto o primeiro (supervisionado) possui rótulos, o segundo (não supervisionado) não tem.

Para exemplificar, supomos que os dados de treinamento são referente a imóveis, tais como: área ( $m^2$ ), nº de cômodos,

localização. No caso do aprendizado supervisionado poderíamos ter como rótulos o tipo de imóvel: popular ou luxo. No caso do aprendizado não supervisionado não existem rótulos.

O aprendizado não supervisionado resolve muitas tarefas da vida real, já que a realidade da maioria das situações é um grande conjunto de dados sem rótulos.

## 2. Aprendizado de máquina não supervisionado

O aprendizado não supervisionado busca encontrar padrões ocultos nos dados. Possui três componentes básicos: entradas, métricas de distância ou semelhança e algoritmos. As entradas também conhecidas pelo termo inglês "features", referem-se aos dados, que trazem atributos e características, que serão usados para treinar os modelos. As métricas de distância/semelhança serão utilizadas para avaliar a qualidade dos padrões encontrados. Algoritmos são usados para treinar os modelos e dependem do tipo de problema ou tarefa a ser executada.

### 2.1 Tarefas

Há três principais tarefas no

aprendizado não supervisionado: agrupamento (clusterização), redução de dimensionalidade e associação.

### 2.1.1 Agrupamento

Agrupamento é uma das tarefas mais populares do aprendizado não supervisionado e consiste em reunir (agrupar) as entradas de acordo com critérios de distância e/ou semelhança entre os pontos dos dados. É importante lembrar que as máquinas precisam que os dados sejam numéricos para conseguirem processá-los. Assim, após esta conversão os dados de entrada podem ser representados por pontos num espaço de dimensão " $n$ ".

O algoritmo mais conhecido para a tarefa de agrupamento é o K-means. O algoritmo cria os centróides, que são os pontos centrais dos pontos próximos e forma um número " $k$ " de grupos ou clusters. No caso do K-means, o centróide é a média aritmética dos pontos contidos no grupo. Outros algoritmos populares são o DBSCAN e o agrupamento hierárquico. Duas métricas importantes para o agrupamento são coesão e distância entre clusters. Coesão diz respeito a quão próximos uns dos outros os pontos contidos num cluster estão. A dis-

PCA3

4

ância entre clusters mede o quão distantes os grupos estão entre si. O número com seus pontos próximos dentro de cada cluster (coeso) e com seus diferentes clusters distantes entre si, indicam um bom agrupamento. Portanto, a ideia seria combinar ambos coesão e distância entre clusters para avaliar a qualidade do modelo. Assim existem métricas que se preocupam com coesão e distância entre grupos como: coeficiente de silhouette, coeficiente Dunn e coeficiente Davis Bouldin.

O agrupamento é aplicado em segmentação de clientes, estratégias de marketing para encontrar perfis de usuários, entre outras aplicações.

## 2.1.8. Redução de Dimensionalidade

Em muitos casos, quando os dados não convertidos em pontos num espaço de dimensão "n", este "n" é muito alto, o que além de dificultar a compreensão, torna-se computacionalmente custoso processá-lo. A redução de dimensionalidade visa reduzir este espaço dimensional mantendo o máximo de informações e variedade de dos dados de entrada. Os algoritmos mais conhecidos PCA (Análise de Componentes Principais) e o t-SNE. Um desafio é saber qual o número ideal de

DCE33

Componentes os dados terão, de modo a manter suas características e diversidade. Este número ideal depende do algoritmo utilizado, mas há métricas e ferramentas que nos ajudam a encontrá-lo, por exemplo a variância explicada e o erro. Com a variância explicada podemos saber o percentual de variância dos dados que foi mantida com um número "x" de componentes. Traçando o gráfico da variância explicada acumulada podemos verificar quantos componentes não são necessários para atingir o percentual de variância desejado. Já o erro, calcula o quanto o componente corresponde ao valor original (diferença entre o ponto gerado e o ponto original). De modo similar, podemos encontrar o nº de componentes que produza um erro até um certo limite, por exemplo,  $erro_{max} = 0.01$ .

A redução de dimensionalidade é bastante usada no pré-processamento de texto e imagens, que possuem alta dimensão para seus dados.

### 2.1.3. Associação

Um pouco menos popular a associação é uma tarefa, que busca encontrar padrões que juntos indiquem uma tendência. Por exemplo, baseado nos itens de

## DCE 33

compras de um consumidor, o sistema pode sugerir novos itens. Um algoritmo utilizado é o Apriari, que busca encontrar padrões baseado na variância dos dados.

### 3. Ciclo de vida do processo de aprendizado não supervisionado (ANS)

O ciclo de vida do ANS inclui como principais etapas: a coleta de dados, o pré-processamento, o treinamento, a avaliação do desempenho e a implantação.

A coleta de dados refere-se a etapa onde os dados são reunidos para serem utilizados. Usa-se conjuntos de dados públicos, API, técnicas de raspagem de dados da internet ou utiliza-se dados gerados internamente.

O pré-processamento dos dados é uma etapa crucial para garantir dados de qualidade, pois não há como gerar bons modelos se a matéria-prima (dados de entrada) não é de boa qualidade. O pré-processamento inclui a limpeza dos dados (duplicidades, erros, etc), o tratamento de dados faltantes (exclusão de amostras, imputação de dados), normalização/ padronização, redução de dimensionalidade, entre outras.

# DEE33

1/1/

O treinamento do modelo inclui utilizar o algoritmo escolhido, ajustando os parâmetros para reduzir o erro. Técnicas como Gridsearch e validação cruzada não usadas em alguns treinamentos para evitar o superajuste dos dados (overfitting) ou um subajuste dos dados (underfitting).

A avaliação de desempenho é a etapa onde usamos as métricas para avaliar o quão bom nosso modelo está. É importante usar métricas apropriadas para a tarefa e seu objetivo e impacto dos resultados.

Implantação do modelo refere-se a etapa em que o modelo é colocado em um sistema num ambiente de produção para ser usado pelos usuários finais. Inclui monitorar, manter e escalar o modelo.

## 4. Desafios e limitações

Os modelos de aprendizado de máquina não supervisionado podem ser difíceis de explicar e compreender os padrões, sendo muitas vezes necessária a atuação de especialistas do domínio que está sendo trabalhado. Outro problema é quando trabalhamos dados de altas dimensões ou com um volume grande





~~DEE33~~

de dados, tornando-se bastante custoso computacionalmente. A presença de dados anômalos (outliers) pode distorcer os resultados dos modelos.

### 5. Tendências

A integração do aprendizado de máquina não supervisionado com as novas tecnologias como "Internet das Coisas" e Big Data, podem trazer uma grande evolução. Assim como a combinação deste tipo de aprendizado com outros tipos. Fala-se em auto-aprendizagem criando modelos que criam e treinam outros modelos e no aprendizado com redes neurais profundas, como forte tendência associada ao aprendizado não-supervisionado.

### 6. Conclusão

O aprendizado de máquina não supervisionado torna-se um campo de bastante interesse, tendo em vista a produção exponencial de dados não rotulados. Importantes algoritmos como K-means, DBSCAN e PCA, ainda são bastante usados e tem respondido com eficiência em muitos problemas reais. Porém, é evidente que pesquisas neste campo de



DCE33

\_ / \_ / \_

todos são necessários para adaptar ou desenvolver novas técnicas para lidar com tecnologias como Internet das Coisas e Big Data.