

02 - Aprendizado de Maq. Não-Supervisionada

Os modelos de regressão linear ou redes neurais utilizam um conjunto de parâmetros para prever ~~uma~~ a saída de uma determinada hipótese. Tais parâmetros são ajustados na etapa de treinamento de acordo com o que é esperado na saída de modo a obter-se um menor erro possível entre o valor predito e o valor real. Neste caso, há uma supervisão na etapa de treinamento ~~visando~~ visando os ajustes dos parâmetros. Portanto, temos um modelo paramétrico e de aprendizado supervisionado.

Quando se tem milhões ou bilhões de dados ~~de~~ no conjunto de treinamento (Dataset), uma abordagem que pode trazer excelentes resultados seria construir um modelo onde os dados "falem" por si e não por ~~um~~ intermédio de vetores parametrizados que muitas vezes são relativamente pequenos para o problema. Além disso, vale ressaltar que para modelos paramétricos o número de parâmetros não muda com o tamanho do dataset de treinamento. Se for considerado uma hipótese que faz uso de todos os dados de treinamento para prever uma saída, tem-se um modelo não paramétrico e de aprendizagem não-supervisionada. Nesse caso, não há ajuste de parâmetros e, claramente, não há nenhuma supervisão, o treinamento considera o data set sem nenhum ajuste orientado a um valor esperado.

Um exemplo bem simples, e com viés didático, seria inserir todos os dados de treinamento em uma tabela. Quando for necessário fazer a previsão da hipótese ($h(x)$)

o valor x pode ser consultado na tabela e retornar o y correspondente. Entretanto, tal proposta demonstra claramente uma dificuldade de generalização, pois caso x não exista na tabela, um valor padrão é retornado. Note que no modelo proposto não houve ajuste de parâmetros de forma supervisionada.

O aprendizado de máquina não-supervisionado possui características extremamente relevantes no contexto de lidar com muitos dados de entrada, pois a etapa de treinamento desses modelos requer menos tempo e menor esforço computacional. Além disso, tais modelos possuem maior facilidade na definição dos ~~tipos de~~ hiperparâmetros pelo projetista devido a ser um modelo menos complexo em geral. Por outro lado, tais modelos possuem dificuldades de lidar com generalizações e problemas complexos.

Um clássico exemplo de modelos de aprendizagem não-supervisionada é o modelo K -means. Para caso de classificação, onde o objetivo é atribuir a(s) entrada(s), ou hipótese ($h(x)$), a um certo grupo (saída), o modelo K -means faz a análise dos K pontos mais próximos de x , usando distância Euclidiana. Neste caso, é como se cada ponto selecionado votasse que x pertence ao seu ~~o~~ grupo (grupo de pontos selecionado). Por fim, x será classificado ao grupo que receber mais votos.

Para casos de regressão, onde a saída pode ser valores do domínio real, é possível definir regressões lineares a partir dos centroides definidos ~~para~~ dos grupos definidos pelo modelo

~~Clustering~~ K-means.