
IF867 - Introdução à Aprendizagem Profunda

1ª Lista Teórica

Discente:

- Gabriel D'assumpção de Carvalho - gdc2@cin.ufpe.br

Curso:

- Ciências Atuariais - 7º Período

04/12/2024

4. LIGTA 10

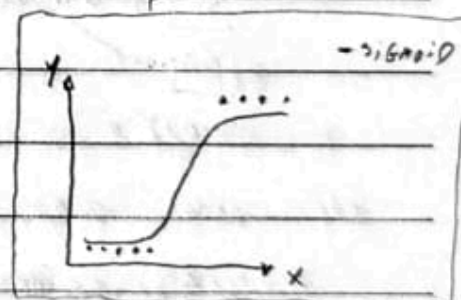
1) A REGRESSÃO LOGÍSTICA E LINEAR SÃO MODELOS ESTATÍSTICOS, QUE FAZEM ESTIMAÇÃO DE UMA VARIÁVEL DEPENDENTE (RESPOSTA Y) A PARTIR DE UM CONJUNTO DE VARIÁVEIS EXPLICATIVAS/INDEPENDENTES (\tilde{X}).

ESSES DOIS MODELOS SÃO MUITO UTILIZADOS PARA PREVISÃO DE RISCO DE INADIMPLÊNCIA, PREÇO DE IMÓVEIS, CUSTOS DE CLIENTES DE PLANOS DE SAÚDE, SE VAI CHOVER OU NÃO NO DIA.

POR MAIS QUE OS DOIS MODELOS SEJAM UTILIZADOS PARA FAZER PREVISÕES, A REGRESSÃO LOGÍSTICA É UTILIZADA APENAS QUANDO SE TEM A VARIÁVEL Y DICOTÔMICA, ASSUMINDO APENAS OS VALORES 0 (ZERO) E 1 (UM). COM BASE NISSO, ESSA REGRESSÃO VAI ESTIMAR A PROBABILIDADE DA OCORRÊNCIA DE UM EVENTO DADO AS VARIÁVEIS X E UM CONJUNTO DE PARÂMETRO θ .

$$P(Y=1 | \tilde{X}, \tilde{\theta}) = \frac{1}{1 + e^{-\tilde{\theta}^T \tilde{X}}}$$

↳ SIGMOID



NO MODELO DE REGRESSÃO LINEAR A VARIÁVEL Y DEVE SER CONTÍNUA ($Y \in \mathbb{R}$). PORTANTO, ESSE MÉTODO ESTIMA POTENCIALMENTE O VALOR DE Y A PARTIR DO CONJUNTO X E θ .

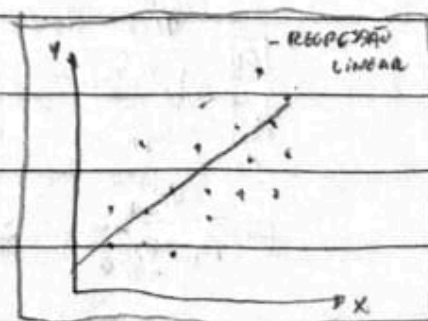
$$Y = \tilde{\theta}^T \tilde{X} + \epsilon$$

↳ RESÍDUO

↳ VETOR VARIÁVEIS EXPLICATIVAS

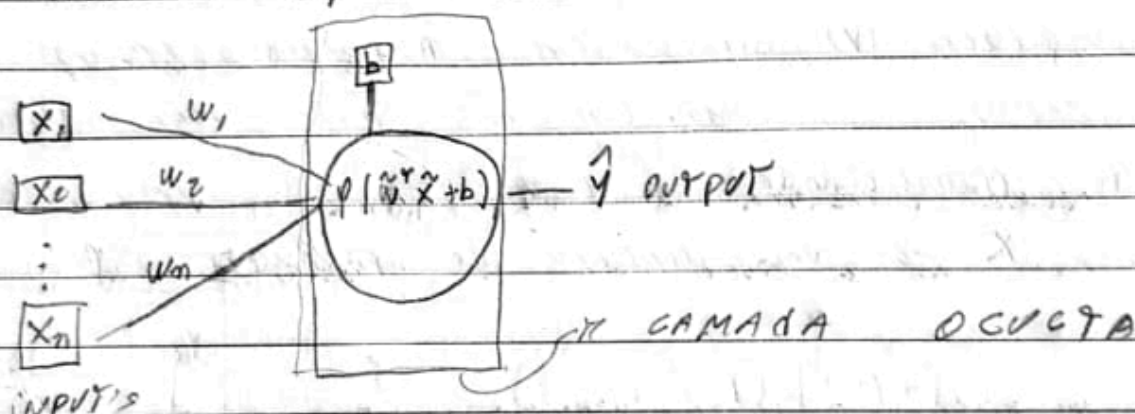
↳ VARIÁVEL DE PARÂMETROS

↳ VARIÁVEL RESPOSTA

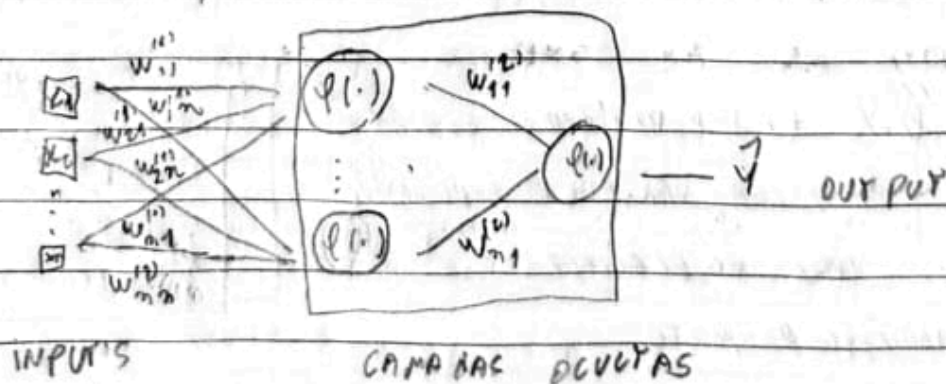


2) O modelo PERCEPTRON foi uma das primeiras estruturas de redes neurais artificiais propostas. Ele é inspirado no funcionamento de um único neurônio biológico, onde recebe informações de entrada (inputs) através de "dendritos", processa a informação em seu "núcleo" por meio de uma função de ativação, e transmite o resultado processado para "sinapses" como saída (output). Esse processo é representado matematicamente como uma combinação linear das entradas ponderada ($\sum w_i x_i + b$) seguida de uma função de ativação para prever o valor da saída.

• ESTRUTURA PERCEPTRON



Como pode ser visto acima, o modelo PERCEPTRON é bastante simples, sendo limitado a problemas linearmente separáveis. Alguns anos após o PERCEPTRON foi proposto o MULTILAYER PERCEPTRON, cuja estrutura é composta por mais camadas ocultas, a fim de resolver problemas não lineares.



2) TEOREMA DA APROXIMAÇÃO UNIVERSAL diz que um MLP com uma camada, M neurônios e o conjunto certo de pesos é possível modelar qualquer função contínua. Com o conjunto dos pesos, é importante encontrar os valores ótimos para os bias que é um parâmetro de deslocamento linear, sem o bias, todas as funções de ativação iriam passar pelo ponto $(0, 0)$, prejudicando o modelo a aprender padrões mais complexo. Portanto, sem o bias o teorema da aproximação universal não seria aplicado na prática.

3) A regularização em modelos estatísticos tem a finalidade de restringir o espaço paramétrico, para estimar um modelo mais simples e com uma melhor capacidade de generalização, diminuindo o overfitting. Essa restrição é feita com uma penalização na função de perda.

O regularizador L_1 (LASSO) penaliza a função de perda $L(\tilde{w}, \tilde{b})$ pelo produto da soma do módulo dos pesos $|\tilde{w}|$ por um hiperparâmetro (λ) que controla a força de regularização.

$$J(\tilde{w}, \tilde{b}) = L(\tilde{w}, \tilde{b}) + \sum_{i=1}^n |\tilde{w}_i| \lambda$$

Pelo fato da derivada do L_1 ser

$$\frac{\partial (|\tilde{w}_i| \lambda)}{\partial \tilde{w}_i} = \begin{cases} \lambda, & \text{se } \tilde{w}_i > 0 \\ -\lambda, & \text{se } \tilde{w}_i < 0 \end{cases}$$

quando for calculado o gradiente descendente o regularizador L_1 vai forçar os pesos para a origem, podendo zerar alguns

pesos. Portanto, o L_1 é bastante vantajoso quando queremos diminuir a dimensionalidade do modelo, observando as variáveis que são mais importantes.

Se o regularizador L_2 vai penalizar $L(\tilde{w}, \tilde{b})$ pelo produto da soma dos quadrados dos pesos por λ

$$J(\tilde{w}, \tilde{b}) = L(\tilde{w}, \tilde{b}) + \lambda \sum_{i=1}^n w_i^2$$

Derivada do L_2 , o L_2 puxa os pesos para a origem de forma contínua, pelo fato da sua derivada ser

$$\frac{\partial (w^2 \lambda)}{\partial w} = 2w\lambda$$

com isso o regularizador L_2 não zera absolutamente um w_i , mas podemos deixar ele bem próximo de zero. Então tanto esse é um bom regularizador para se utilizar quando se tem todas as variáveis importantes mas se pretende ter um modelo que consiga generalizar para novas observações.

Como visto anteriormente a escolha do hiperparâmetro λ impacta na força da regularização, controlando a complexidade do modelo e capacidade de generalização. Um λ baixo faz com que o erro do treinamento tenha pouca alteração, dando margem ao modelo se ajustar aos dados; com um λ alto o modelo fica forçado a ter uma taxa de erro mais alta, fazendo com que o modelo se torne menos complexo e evitando o overfitting mas podendo ter um underfitting. Portanto é importante escolher um λ nem tão alto e nem tão baixo, podendo utilizar técnicas de validação cruzada para encontrar λ .

4) NAS PRIMEIRAS REDES ARTIFICIAIS OS CIENTISTAS ESTAM INTERESSADOS EM MODELAR UM NEURÔNIO QUE DISPARA E NÃO DISPARA, UTILIZANDO UMA FUNÇÃO DE LIMITARIZAÇÃO. POR CONTA DO USO DO GRADIENTE A FUNÇÃO SIGMOIDE FOI UMA ESCOLHA VULGAR, SENDO UMA FUNÇÃO NÃO LINEAR, CONTÍNUA, SUAVA E DIMENSIONAL EM TERMOS DE PONTOS.

A FUNÇÃO SIGMOIDE É UTILIZADA QUANDO QUEREMOS MODELAR PROBABILIDADE, TENDO O VALOR DE SAÍDA PERTENCENDO AO INTERVALO $[0, 1]$. POR CONTA DISSO, ELA É BASTANTE UTILIZADA PARA PROBLEMAS DE CLASSIFICAÇÃO BINÁRIA.

COM O AVANÇO DA COMPUTAÇÃO, FOI POSSÍVEL CRIAR ESTRUTURAS DE NN COM MUITAS CAMADAS OCULTAS. ENTRETANTO, AO UTILIZAR A FUNÇÃO SIGMOIDE EM MUITAS CAMADAS A REDE SE TORNA INCAPAZ DE RETROPROPAGAR O GRADIENTE DA CAMADA DE SAÍDA DE NOVA PARA AS PRIMEIRAS CAMADAS. ISSO ACONTECE PORQUE O PONTO MÁXIMO DA DERIVADA DA SIGMOIDE É IGUAL A 0,25, APROXIMANDO O GRADIENTE PARA O PONTO ZERO À CADA CAMADA OCULTA NO BACKPROPAGATION, CAUSANDO A CHAMADA DISSIPACÃO DO GRADIENTE.

A FUNÇÃO RELU FOI IMPLEMENTADA EM REDES NEURAIS PROFUNDAS, PARA SUPRIMIR A LIMITAÇÃO DA SIGMOIDE NO TREINAMENTO E DEVIDO AO SEU BAIXO CUSTO COMPUTACIONAL. SENDO ELA

$$\text{ReLU}(x) = \max(0, x) \quad ; \quad \text{ReLU}'(x) = \begin{cases} 1, & \text{se } x > 0 \\ 0, & \text{se } x \leq 0 \end{cases}$$

51

② O modelo que melhor generaliza é o terceiro pois além de ter a 2ª maior acurácia de treinamento possui o melhor desempenho nos dados de teste.

③ O 1º modelo apresentou um possível overfitting, pelo fato da acurácia de teste estar em 80%, sendo 15% a baixo da acurácia de treinamento.

④ O modelo que apresentou um underfitting foi o segundo modelo, pois ele não foi capaz de aprender os padrões nos dados de treino e teste, ficando com uma acurácia de 70% e 68% respectivamente.

⑤ Existem diversas técnicas que podem ser utilizadas, dentre elas existem os regularizadores da função de perda, diminuir a profundidade da rede, early stop que faz uma parada antecipada ou usar um ensemble learning que consiste em treinar diversos modelos e fazer a previsão com base em uma média.

⑥ Diferente do overfitting, para evitar o underfitting pode tentar recolher novos dados e tratar de uma maneira melhor, aumentar a complexidade do modelo e aumentar a epoch time.

5) Quando o dropout é utilizado na matriz de pesos vai ser desligado em média $(1-p)$ sinapses de cada camada, sendo $(1-p)$ o número de sinapses e (p) a probabilidade de dropout, mas desligar apenas as sinapses o modelo ainda pode ficar dependente de alguns neurônios, fazendo com que essa técnica adicione pouco ruído para a rede. A utilização de dropout após a função de ativação é a prática mais comum, pois em média desliga $(1-p)$ neurônios da camada, sendo $(1-p)$ o número de neurônios. Isso faz com que o modelo não dependa de certos neurônios para fazer a sua estimativa, sendo uma técnica que adiciona mais ruído e melhora a capacidade de generalização do modelo.

6) O gradiente explosivo é um problema que se dá quando os parâmetros da rede recebem grandes atualizações no backpropagation fazendo com que os parâmetros sempre fiquem com valores altos e não vão para o mínimo. Isso pode acontecer quando os valores iniciais dos parâmetros são altos ou a derivada da função de ativação seja maior que um, fazendo com que o gradiente sempre aumente a cada camada oculta.

O vanishing gradient é o inverso do explosivo, que seria a atualização insignificante dos parâmetros na etapa de backpropagation. Esse problema pode acontecer se a derivada da função de ativação tiver grandes pontos de saturação, como a função sigmoide.

Quando escolhemos um grande valor para a taxa de aprendizado o ajuste dos parâmetros sempre vão dar um salto

TO EM TORNO DO MÍNIMO LOCAL OU GLOBAL, EVUJINDO DA CONVERGÊNCIA. ESSA SITUAÇÃO É DENOMINADA DE OVER Shooting.

3) QUANDO TEMOS UM PROBLEMA DE REGRESSÃO ESTAMOS INTERESSADOS EM FAZER PREVISÃO PARA VALORES CONTÍNUOS, PARA FAZER ESSAS PREVISÕES DEVEMOS MINIMIZAR O NOSSO ERRO QUE SERIA A DISTÂNCIA DO VALOR PREVISTO (\hat{y}) PARA O VALOR OBSERVADO (y), COM ISSO PODAMOS USAR AS SEGUINTE FUNÇÕES DE PERDA:

$$MSE = n^{-1} \cdot \sum_{i=1}^n (\hat{y}_i - y_i)^2 \quad \text{OU} \quad MAE = n^{-1} \cdot \sum_{i=1}^n |\hat{y}_i - y_i|$$

COMO PODEMOS VER A LINA, A FUNÇÃO MSE ELEVA A DISTÂNCIA DE (\hat{y}) E (y) AO QUADRADO, FAZENDO COM QUE ERROS PEQUENOS SEJAM MAIS INSIGNIFICANTES MAS PUNALIZA MAIS OS ERROS GRANDES. SE QUERMOS GRANDES QUANTIDADES DE OUTLIERS É INTERESSANTE UTILIZAR A FUNÇÃO MAE POR SER MAIS ROBUSTA, PELA FAZTO DE CALCULAR A DISTÂNCIA ENTRE (\hat{y}) E (y) DE FORMA LINEAR.

QUANDO TEMOS PROBLEMAS DE CLASSIFICAÇÃO FICAMOS INTERESSADOS EM MODULAR A PROBABILIDADE DE CADA CLASSE (y) DADO UM CONJUNTO DE CARACTERÍSTICAS X E UM CONJUNTO DE PARÂMETRO θ . UMA DAS FUNÇÕES DE PERDA MAIS UTILIZADA PARA CLASSIFICAÇÃO É A CROSS-ENTROPY, QUE É DADA POR

$$L_{CE} = n^{-1} \cdot \sum_{i=1}^n [y_i \cdot \log(\hat{y}_i) + (1 - y_i) \cdot \log(1 - \hat{y}_i)]$$

SENDO $\hat{y}_i = P(y_i | X; \theta)$

8)

① O gradiente descendente é um algoritmo de otimização utilizado para minimizar a função de perda, sendo utilizado na fase de treinamento do modelo. Ele funciona ajustando iterativamente os parâmetros θ na direção oposta ao gradiente da função de perda em relação aos parâmetros.

$$\theta_{t+1} = \theta_t - \eta \frac{\partial L(\theta)}{\partial \theta}$$

• θ_t : PARÂMETRO NA EPOCH T

• θ_{t+1} : " " " " $T+1$

• η : TAXA DE APRENDIZADO

• $L(\theta)$: FUNÇÃO DE PERDA

Esse algoritmo é bastante utilizado, pois é o método mais eficiente para encontrar os θ^* , sendo capaz de ser utilizados em modelos de regressões e redes neurais e tem a capacidade de funcionar para grandes quantidades de dados.

②

• BATCH: esse é um método que consiste em utilizar todas as observações do conjunto de treinamento para calcular a média da função de perda para depois fazer a atualização dos parâmetros. Portanto, a sua convergência para o mínimo é suave mas só é recomendada para uma quantidade pequena de observações, caso contrário levaria muitas épocas até sua convergência.

• **Stochastic**: É um método que utiliza apenas uma observação aleatória da base de treinamento, fazendo com que os parâmetros sejam frequentemente modificados. Interpretando, a sua natureza estocástica faz com que ele possa estar capar de mínimos locais mas oscila em torno do mínimo global, sendo mais instável, o que pode dificultar a convergência.

• **Mini-batch**: É o método mais comum, pois é uma junção de batch com o stochastic. Ele funciona selecionando uma amostra de tamanho K da base de treinamento, fazendo com que não demore tanto tempo para atualizar os parâmetros e tem possibilidade de fugir de mínimos locais.

9)

① A taxa de aprendizado é um hiperparâmetro que controla o tamanho do passo que o algoritmo de otimização dá em direção ao mínimo da função de custo. Ele controla a velocidade da convergência, afetando a estabilidade e eficiência durante o treinamento.

② Uma taxa de aprendizagem muito alta vai fazer com que os parâmetros de grandes saltos, se essa taxa for constante ao longo do treinamento pode ser que o modelo não converja para o mínimo, ficando sempre em torno dele.

Ao contrário, quando se tem uma taxa muito baixa

O modelo leva muitas épocas para convergir, podendo ficar preso em mínimos locais.

③ As estratégias adaptativas ajustam automaticamente a taxa de aprendizado com base no gradiente, tornando o treinamento mais estável e eficiente.

Um dos otimizadores bastante utilizado é o Adam, sendo de

$$SGD \Rightarrow \theta_{t+1} = \theta_t - \eta \cdot (\sqrt{G_t})^{-1} \cdot \frac{\partial L(\theta)}{\partial \theta}$$

• θ_t : parâmetro no tempo t

• θ_{t+1} : " " " " $t+1$

• η : taxa de aprendizado

• G_t : soma dos quadrados dos gradientes

• $\frac{\partial L(\theta)}{\partial \theta}$: gradiente da função de perda em relação a θ

A taxa de aprendizado é ajustado para cada parâmetro de acordo com o gradiente da função de perda em relação a θ . Portanto se um parâmetro tem gradientes grandes a taxa diminui, caso contrário, a taxa aumenta.

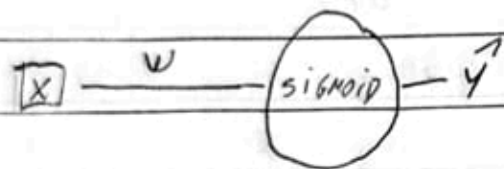
101 O backpropagation é um algoritmo de otimização que utiliza o gradiente descendente e a regra da cadeia para ajustar os valores dos parâmetros da última camada até a primeira.

O forward é o primeiro passo no ciclo de treinamento. Funcionando com o recebimento de um conjunto X de variáveis, sendo passado para frente o resultado da função $\psi(w, x, b)$. Até a última camada, que faz a estimativa

depois do FORWARD PASS é feito o BACKWARD, que consiste em utilizar o BACKPROPAGATION PARA AJUSTAR OS VALORES DOS PARÂMETROS DA ÚLTIMA CAMADA ATÉ A PRIMEIRA.

* QUESTÃO 60005

ESTRUTURA DA REDE:



CONO PODE SER OBSERVADO A REDE PROPOSTA É UM PERCEPTRON QUE RECEBE APENAS 1 FEATURES E SUA FUNÇÃO DE ATIVAÇÃO É UMA SIGMOID, ISSO SUGERE QUE O PROBLEMA A SER RESOLVIDO É DE CLASSIFICAÇÃO BINÁRIA.

①

$$z = 0,9 \cdot 3 = 2,7$$

$$\hat{y} = \frac{1}{1 + e^{-z}} \approx 0,9315$$

$$E \approx (0,9 - 0,9315)^2 \approx 0,001$$

②

$$\frac{\partial (y - \hat{y})^2}{\partial \hat{y}} = 2(y - \hat{y})$$

$$\frac{\partial (1 + e^{-z})^{-1}}{\partial z} = -e^{-z} (1 + e^{-z})^{-2} = \frac{-1}{(1 + e^{-z})} \cdot \frac{e^{-z}}{(1 + e^{-z})} = \frac{-1}{(1 + e^{-z})} \cdot \left[1 - \frac{1}{(1 + e^{-z})} \right] = -\hat{y} \cdot (1 - \hat{y})$$

$$\frac{\partial (w \cdot x)}{\partial w} = x$$

$$\textcircled{3} \frac{\partial E}{\partial w} = 2(y - \hat{y}) \cdot \hat{y} \cdot (1 - \hat{y}) \cdot x$$

④

$$w_{t+1} = 0,4 - 0,1 \cdot 2 \cdot (0,2315 - 0,9) \cdot 0,2315 \cdot (1 - 0,2315) \cdot 3 = 0,4714$$

$$\frac{\partial E}{\partial w} = -0,7136$$