
TENSORFLOW 2.0: APRENDIZAGEM PROFUNDA E INTELIGÊNCIA ARTIFICIAL

ESCRITO POR: NAILSON CHAGAS

Index

1	O que é Aprendizado de Máquina?	2
1.1	Aprendizado de Máquina é apenas um problema geométrico	2
2	Teoria de Classificação Linear	3
3	Teoria de Regressão Linear	6
4	O Neurônio	8
4.1	Como isso é relacionado com um neurônio?	8
4.2	Rede neural artificial	9
4.3	Como um modelo aprende?	9
5	Feedforward Neural Network	10
5.1	Arquitetura do Modelo	10
5.2	A Geometria	11
5.3	Funções de Ativação	11
5.3.1	Sigmóide	11
5.3.2	Tangente Hiperbólica	11
5.3.3	ReLU	12
5.3.4	ELU	12
5.3.5	LReLU	12
5.3.6	Softplus	12
5.4	Classificação Multiclasse	12
5.5	O que Redes Neurais podem fazer?	12
5.6	Imagens	12
5.6.1	Como imagens são armazenadas?	12
5.6.2	Como representar um dataset de imagens?	13
5.7	Código	13

Sobre essas anotações

Minhas anotações do curso: <https://udemy.com/course/deep-learning-tensorflow-2>

Exemplos práticos: <https://github.com/NailsonChagas/Tensorflow-2.0-Deep-Learning-and-Artificial-Intelligence>

O que é Aprendizado de Máquina?

Aprendizagem de máquina é um subcampo da Engenharia e da ciência da computação que evoluiu do estudo de reconhecimento de padrões e da teoria do aprendizado computacional em inteligência artificial. Em 1959, Arthur Samuel definiu aprendizado de máquina como o "campo de estudo que dá aos computadores a habilidade de aprender sem serem explicitamente programados"

1.1 Aprendizado de Máquina é apenas um problema geométrico

- Um estatístico diria: Aprendizado de Máquina é apenas um ajuste de curva glorificado.
- Tanto regressão quanto classificação são exemplos de aprendizado supervisionado .

Teoria de Classificação Linear

Obs: Para facilitar o exemplo lida só com classificação linear binária.

A linha que separa um conjunto de pontos em um plano cartesiano pode ser escrita através da equação: (1) $y = m * x + b$

Também podendo ser escrita como sendo: (2) $w_1 * x_1 + w_2 * x_2 + b = 0$, com x_1 sendo o eixo horizontal e x_2 sendo o eixo vertical.

Como podemos usar essa linha para classificar pontos?

$$(3) a = w_1 * x_1 + w_2 * x_2 + b$$

Usando a equação (3), podemos tomar uma decisão:

- Se $a \geq 0 \implies$ Classe 1
- Se $a < 0 \implies$ Classe 0

Matematicamente podemos encapsular esse processo de decisão em uma função de Heaviside (função degrau: assume valor 0 ou 1):

$$(4) \hat{y} = u(a), a = w_1 * x_1 + w_2 * x_2 + b$$

Como em Deep Learn nós preferimos utilizar funções mais suaves e diferenciáveis, é utilizada uma função Sigmoid (função logística: assume valores entre 0 e 1):

$$(5) \hat{y} = \sigma(a), a = w_1 * x_1 + w_2 * x_2 + b$$

Com \hat{y} da equação (5) podendo ser interpretado como a probabilidade de que:

$$p(y = 1|x) = \sigma(w_1 * x_1 + w_2 * x_2 + b), \text{ ou seja, } y = 1 \text{ dado } x$$

Com o resultado da probabilidade podemos:

- Se $p(y = 1|x) \geq 0 \implies$ prever 1, se não 0

Quando aplicamos uma função Sigmoidal em cima de uma função linear, nós chamamos o modelo de Regressão Logística, e o argumento de uma função Sigmoidal é chamado de ativação.

Usando o que sabemos até agora, pode-se notar um problema de notação: e se tivermos mais de duas entradas ($w_1 * x_1, w_2 * x_2, \dots, w_n * x_n$)? R: Sem problemas, considerando w como um vetor de pesos e x como um vetor de características :

$$(6) \quad p(y = 1|x) = \sigma(w^T * x + b) = \sigma\left(\sum_{d=1}^D w_d * x_d + b\right)$$

Mas como podemos implementar, o que vimos até agora no Tensorflow? A expressão $w^T * x + b$ é implementada através da função:

```
tf.keras.layers.Dense(output_size)
```

Para podermos implementar o modelo de Regressão Logística utilizaremos dois layers:

```
"""
    Input(shape=(D,)), basicamente falando para o keras o tamanho do vetor de entrada x
    Dense(1, activation='sigmoid') especifica o tamanho de saida e função de ativação
"""

model = tf.keras.models.Sequential([
    tf.keras.layers.Input(shape=(D,)),
    tf.keras.layers.Dense(1, activation='sigmoid')
])

#configura o modelo para treino
model.compile(
    optimizer='adam',
    loss='binary_crossentropy',
    metrics=['accuracy']
)

r = model.fit(
    X_train, y_train,
    validation_data=(X_test, y_test),
    epochs=100 #numero de iterações para calculo dos pesos (vetor w)
)

"""
    Visualiar a função de perda e outras métricas ao passar das iterações
    Com isso podemos avaliar se precisamos de mais epochs
    ou alterar outros hiper parametros
"""

plt.plot(r.history['loss'], label='loss')
plt.plot(r.history['val_loss'], label='val_loss')
```

Para olhar o código funcionando, usar o script classification na pasta ML and Neurons.

Teoria de Regressão Linear

Obs: Para facilitar o exemplo lida só com regressão linear.

Em estatística regressão linear é uma equação para se estimar a condicional (valor esperado) de uma variável y , dados os valores de algumas outras variáveis x . A regressão, em geral, tem como objetivo tratar de um valor que não se consegue estimar inicialmente. A regressão linear é chamada "linear" porque se considera que a relação da resposta às variáveis é uma função linear de alguns parâmetros.

- Equação da reta: $\hat{y} = \sum_{d=1}^D w_d * x_d + b = w^T * x + b$

Diferente da classificação, a regressão, devido a sua saída poder assumir "qualquer" valor, não necessita de uma função de ativação (ex: sigmoide)

```
model = tf.keras.models.Sequential([
    tf.keras.layers.Input(shape=(D,)),
    tf.keras.layers.Dense(1, activation=None)
])
```

```
model.compile(
    #escolhido no lugar do adam pois esse performou melhor nesse caso
    optimizer = tf.keras.optimizers.SGD(0.001, 0.9),
    #como a saída não é binária (mse = mean square erro)
    loss = 'mse'
)
```

Em regressão o uso da acurácia não faz sentido pois ela quase sempre seria 0. No caso da regressão faz mais sentido utilizar a medida estatística R^2 (coeficiente de determinação)

O Neurônio

Chamamos regressão logística de neurônio, com um neurônio sendo a base fundamental de uma rede neural.

$$\hat{y} = \sigma\left(\sum_{d=1}^D w_d * x_d + b\right)$$

Com cada x_d sendo uma característica, é possível afirmar que seu peso w_d diz o quanto importante essa característica é para predizer a saída.

4.1 Como isso é relacionado com um neurônio?

Os neurônios têm um papel essencial na determinação do funcionamento e comportamento do corpo humano e do raciocínio. Eles são formados pelos dendritos, que são um conjunto de terminais de entrada, pelo corpo central, e pelos axônios que são longos terminais de saída. Neurônios se comunicam através de sinapses. Sinapse é a região onde dois neurônios entram em contato e através da qual os impulsos nervosos são transmitidos entre eles. Os impulsos recebidos por um neurônio A, em um determinado momento, são processados, e atingindo um dado limiar de ação, o neurônio A dispara, produzindo uma substância neuro transmissora que flui do corpo celular para o axônio, que pode estar conectado a um dendrito de um outro neurônio B. O neuro transmissor pode diminuir ou aumentar a polaridade da membrana pós-sináptica, inibindo ou excitando a geração dos pulsos no neurônio B. Este processo depende de vários fatores, como a geometria da sinapse e o tipo de neuro transmissor. Tendo essas informações, é possível fazer algumas suposições básicas:

- A atividade de um neurônio é um processo tudo ou nada (binário).

- Um certo número fixo (>1) de entradas devem ser excitadas dentro de um período de adição latente para excitar um neurônio.
- Único atraso significativo é o atraso sináptico.
- A atividade de qualquer sinapse inibitória previne absolutamente a excitação do neurônio.
- A estrutura das interconexões não muda com o tempo.

Através dessas suposições, pode-se representar um neurônio artificial através do algoritmo de Regressão Logística:

$$\hat{y} = \sigma \left(\sum_{d=1}^D w_d * x_d + b \right)$$

Com os n elementos do vetor x sendo as entradas (dendritos) do neurônio, \hat{y} sendo o terminal de saída (axônio), os elementos do vetor w descrevendo o comportamento das sinapses, $w_d * x_d$ o efeito da sinapse e b representando o bias (Tendência).

4.2 Rede neural artificial

Uma rede neural artificial é composta por várias unidades de processamento (neurônios), cujo funcionamento é bastante simples. Essas unidades, geralmente são conectadas por canais de comunicação que estão associados a determinado peso. As unidades fazem operações apenas sobre seus dados locais, que são entradas recebidas pelas suas conexões. O comportamento inteligente de uma Rede Neural Artificial vem das interações entre as unidades de processamento da rede. A maioria dos modelos de redes neurais possui alguma regra de treinamento, onde os pesos de suas conexões são ajustados de acordo com os padrões apresentados. Em outras palavras, elas aprendem através de exemplos. Arquiteturas neurais são tipicamente organizadas em camadas, com unidades que podem estar conectadas às unidades da camada posterior. Usualmente as camadas são classificadas em três grupos:

- Camada de Entrada: onde os padrões são apresentados à rede.
- Camadas Intermediárias ou Escondidas: Onde é feita a maior parte do processamento, através das conexões ponderadas; podem ser consideradas como extratoras de características.
- Camada de Saída: onde o resultado final é concluído e apresentado.

4.3 Como um modelo aprende?

Uma rede neural aprende calculando os pesos de cada característica, em cada neurônio. O Tensorflow ira fazer todas as diferenciações para calcular os pesos com menores custos, sendo preciso selecionar apenas a taxa de aprendizado. Obs: Para escolher a taxa de aprendizado, normalmente, tentar potencias de 10 (0.1, 0.01, 0.001, ...)

Feedforward Neural Network

Feedforward Neural Network (FFNN) ou Rede Neural de Alimentação Direta, é o tipo mais básico de rede neural e forma a base de outros tipos de redes neurais como:

1. Convolutional Neural Network (CNN)
2. Recurrent Neural Network (RNN)

Em uma FFNN o sinal apenas vai da esquerda para a direita, ou seja, os neurônios da camada anterior sempre passam as informações para os neurônios da próxima camada sem que esses tenham como enviar sinais à camada anterior.

5.1 Arquitetura do Modelo

A rede neural é dividida em n camadas, com cada camada possuindo n neurônios, com todos os neurônios possuindo estrutura uniforme (ou seja, todos neurônios são iguais em todas as camadas). A mesma entrada pode ser usada em múltiplos neurônios de uma mesma camada, com cada um calculado algo diferente. Pode-se considerar que uma rede neural possui duas dimensões, sua largura (quantos neurônios há por camada) e profundidade (quantas camadas existem na rede). Chamando a saída de um neurônio de z_j em uma camada com M neurônios, pode-se representar uma camada de uma rede neural através da expressão:

$$z_j = \sigma(w_j^T * x + b_j), \text{ para } j = 1 \text{ até } M$$

Ou vetorizando:

$$z = \sigma(W^T * x + b)$$

Onde z é um vetor de tamanho M representando as saídas de uma camada, x é um vetor de tamanho D representando as entradas dos neurônios na camada, W é uma matriz de tamanho $D \times M$ representando os pesos de cada entrada em cada neurônio e b um vetor de tamanho M representando o bias. Tendo a representação vetorizada, cada as camadas de uma rede com L camadas, podem ser escritas como:

$$\begin{aligned}z^{(1)} &= \sigma(W^{(1)T} * x + b^{(1)}) \\z^{(2)} &= \sigma(W^{(2)T} * z^{(1)} + b^{(2)}) \\z^{(3)} &= \sigma(W^{(3)T} * z^{(2)} + b^{(3)}) \\&\dots \\p(y = 1|x) &= \sigma(W^{(L)T} * z^{(L-1)} + b^{(L)})\end{aligned}$$

Para uma rede neural usada para regressão, a única coisa que mudaria é que na última camada não haverá a função sigmóide:

$$\hat{y} = W^{(L)T} * z^{(L-1)} + b^{(L)}$$

5.2 A Geometria

Por que usar redes neurais ao invés de usar Regressão Logística e Linear? Nem todos os problemas do mundo podem ser representados bidimensionalmente e/ou linearmente, as vezes algo que queremos pode ser apenas representado através de curvas extremamente complexas.

5.3 Funções de Ativação

5.3.1 Sigmóide

$$\sigma(a) = \frac{1}{1 + \exp(-a)}$$

- Mapeia a saída para valores entre 0 e 1 \rightarrow Imita um neurônio biológico
- Faz com que a rede neural não seja linear
- Problema: A sigmóide foi usada em todos os exemplos até aqui, porém como é preferido que as entradas dos neurônios sejam centradas ao redor de 0 e tenham aproximadamente a mesma distância do centro (exemplo de como não deve ser: uma entrada tem valor de 1 até 5 milhões e outra de 0 até 0.0000001), a sigmóide não é mais tão utilizada (seu centro é 0.5)

5.3.2 Tangente Hiperbólica

$$\tanh(a) = \frac{\exp(2a) - 1}{\exp(2a) + 1}$$

- Mapeia a saída para valores entre -1 e 1
- Possui o mesmo formato de uma sigmóide, porém seu centro é em 0
- Problema de padronização: Apesar de ser um pouco melhor que a sigmóide, ela ainda é problemática (ver The Vanishing Gradient Problem, sigmóide tem o mesmo problema).

5.3.3 ReLU

- Escolha padrão
- Não possui o problema apresentado por The Vanishing Gradient Problem

5.3.4 ELU

$$f(x) = \begin{cases} x & \text{se } x \geq 0 \\ \alpha(\exp(x) - 1) & \text{se } x < 0 \end{cases}$$

Autores dizem que acelera o aprendizado e leva a melhor precisão

5.3.5 LReLU

5.3.6 Softplus

$$f(x) = \log(1 + e^x)$$

5.4 Classificação Multiclasse

Para classificação binária é utilizado ReLUs nas camadas intermediárias (hidden layers) e a sigmóide na saída. Para a classificação de K classes é utilizada a função Softmax na saída.

```
tf.keras.layers.Dense(K, activation='softmax')
```

5.5 O que Redes Neurais podem fazer?

Tarefa	Função de ativação comumente usada na camada de saída
Regressão	Nenhuma
Classificação Binária	Sigmóide
Classificação Multiclasse	Softmax

5.6 Imagens

5.6.1 Como imagens são armazenadas?

- Imagens são armazenadas como matrizes com dimensões: largura x altura x canais-de-cores (Ex: RGB tem 3 canais \rightarrow red, green e blue)
- Cada canal de cor de uma célula (pixel) da imagem assume valores entre 0 e 255 (8bits)
- Sem compressão, uma imagem RGB de 500x500 ira ocupar 732KB
(Motivo: $\frac{500px * 500px * 3 \text{ canais} * 8bits}{8bits} = 750000bytes$)
- Diferentes formatos de imagem \rightarrow diferentes métodos de comprimir a imagem

- Para usar em redes neurais é conveniente escalar os valores de 0 a 255 para 0 a 1 (uma exceção ao padrão de centralizar valores ao redor de 0, devido a representação de 0 a 1 poder ser interpretada como probabilidades)

5.6.2 Como representar um dataset de imagens?

Entrada de uma rede neural deve ter formato N (numero de amostras) x D (numero de características)

1. Converter a matriz de cada imagem para um vetor 1 x D (com D sendo o numero total de pixels, isso levando em conta que a imagem é preto e branca, caso não for também terá os canais de cores)
2. Colocar essas imagens em um vetor de tamanho N (numero de imagens), assim teremos uma matriz de dimensões N x D

5.7 Código

Irá utilizar o dataset MNIST (numeros escritos no papel, 0 a 9, sendo de tamanho 28 x 28 em preto e branco), já incluso no tensorflow.

```
model = tf.keras.models.Sequential([
    tf.keras.layers.Flatten(input_shape=(28, 28)),
    tf.keras.layers.Dense(128, activation='relu'),
    tf.keras.layers.Dropout(0.2),
    tf.keras.layers.Dense(10, activation='softmax')
])

model.compile(
    optimizer='adam',
    loss='sparse_categorical_crossentropy',
    metrics=['accuracy']
)

model.fit(x_train, y_train)
```