

Estudos em Regressão Linear

Matheus Mortatti Diamantino
RA 156740
matheusmortatti@gmail.com

José Renato Vicente
RA 155984
joserrenatovi@gmail.com

I. INTRODUCTION

Este projeto teve como intuito o estudo prático do método de regressão linear em Machine Learning. Foram utilizados os algoritmos de *Gradient Descent* conhecidos como *Stochastic*, *Batch*, *Mini Batch* e *Equação Normal*, de modo a compará-los em termos de complexidade e acurácia.

Foi feito um estudo de predição de preço de diamantes, utilizando uma base de dados com 54000 exemplos, em que são apresentados seus preços e nove features como tamanho, cor e número de quilates.

II. ACTIVITIES

A. Regressão Linear

Regressão Linear é um método muito conhecido de Machine Learning, utilizado para prever o valor de uma variável dependente baseado em valores de variáveis independentes. Essa regressão é chamada linear porque se considera que a relação da resposta às variáveis é uma função linear de alguns parâmetros. Desta forma, dado um vetor Theta de tamanho igual ao número de features, cujo valor queremos determinar, temos que:

$$\text{PreçoAlvoEsperado} = \sum_{i=1}^m \theta_i X_i = h_{\theta}(x) \quad (1)$$

Em que X é um vetor com os valores das features para um dado diamante, cujo preço queremos determinar. Para encontrar esse valor de Theta, utilizaremos alguns algoritmos e compararemos os resultados obtidos com cada um.

Cada algoritmo utilizado é baseado no método de *Descida de Gradiente* (ou *Gradient Descent*). Este é um método utilizado para achar o ponto mínimo de uma função, aproximando gradativamente seu valor até um ponto quando não é possível ser diminuído mais (i.e. a derivada da função neste ponto é zero). Este método consegue apenas achar mínimos locais e, com isso, não é garantido que o resultado obtido é o melhor para o dado problema.

Para medirmos a eficácia do algoritmo, utilizamos uma *Função de Custo* que nos diz o quão perto do resultado desejado estamos, dado um conjunto de dados. Esta função é definida por:

$$J(\theta) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i)^2 \quad (2)$$

Como queremos minimizar a função de custo, queremos que cada passo da nossa descida de gradiente se aproxime mais

do mínimo local. Para extrairmos a direção que temos que ir, utilizamos a derivada da função de custo:

$$\frac{\partial J}{\partial \theta_j} = \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i) x^i \quad (3)$$

Logo, para aproximarmos os valores de θ de modo a nos aproximar do mínimo local, utilizamos a fórmula

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i) x^i \quad (4)$$

, onde $0 \leq j \leq m$, $x_0 = 1$ e α é o que chamamos de *Learning Rate* que define o quão agressivamente tentaremos nos aproximar do mínimo. Este método de descida de gradiente é chamada de *Batch Gradient Descent*. A seguir, veremos três outras variações deste algoritmo

a) *Stochastic Gradient Descent*: Neste método, utiliza-se apenas um exemplo de treino para cada passo da descida de gradiente. Deste modo, a equação 4 se transforma em

$$\theta_j := \theta_j - \alpha (h_{\theta}(x^i) - y^i) x^i \quad (5)$$

Utiliza-se este método quando procura-se rapidez de execução. Contudo, um ponto negativo deste método é que, como usamos como amostra apenas um exemplo de treino, um passo pode nos levar a um custo mais alto. Como o número de iterações é alta, porém, o método converge ao mínimo e mais eficientemente do que o *Batch* até um certo limite.

b) *Mini Batch Gradient Descent*: Para obtermos um resultado balanceado, utiliza-se uma mistura dos métodos *Batch* e *Stochastic*, em que define-se um tamanho para o lote de exemplos de treino que serão utilizados para atualizar θ .

$$\theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=k}^{k+b-1} (h_{\theta}(x^i) - y^i) x^i \quad (6)$$

, onde b é o tamanho do lote, $k = 0, b, 2b, \dots, m - 1$.

c) *Equação Normal*: Para a regressão linear, é possível derivar uma fórmula direta para o ponto de mínimo local que desejamos. Para isso, utilizamos manipulações matriciais na forma

$$\theta = (X^T X)^{-1} X^T y \quad (7)$$

, onde X é a matriz de features, y é a matriz dos dados que queremos prever e θ é a matriz dos coeficientes de $h(X)$.

B. Normalização de Features

Como cada feature tem seu valor em uma escala diferente (i.e. algumas estão na ordem de milhares e outras na ordem de centenas), o processo de descida do gradiente poderá acontecer de forma lenta. Isso se dá pelo fato de que a atualização dos θ s não ocorrerá de forma uniforme entre as features, já que a distância de um dado θ_i a seu valor esperado pode ser maior do que de outro $\theta_j, j \neq i$. Assim, realizamos o que é chamado de *Normalização de Features*, onde colocamos todas as features x_i em um valor entre $0.5 \leq x_i \leq 0.5$. Isso é feito através da fórmula:

$$x_i = \frac{x_i - \frac{\text{size}(x_i)}{2}}{\text{size}(x_i)} \quad (8)$$

C. Transformação de Features

Em alguns casos, podemos ter features que não possuem um valor no domínio dos números reais. Por exemplo, uma feature pode representar a cor de um dado elemento. Deste modo, é necessário realizar uma transformação de tais features para o domínio dos reais. Para realizar tal transformação, criamos uma nova feature para cada possível valor da feature que queremos transformar, de modo que, se para o exemplo de dado e_i temos que esta feature possui um valor x_j , a nova feature correspondente a x_j terá o valor 1 e as demais features criadas terão o valor 0. Tomemos como exemplo uma feature de Cor que pode receber os valores *Azul*, *Amarelo*, *Vermelho* e *Verde*. Assim, o resultado da transformação acontecerá da seguinte forma:

Feature x_i	Azul	Amarelo	Vermelho	Verde
Azul	1	0	0	0
Amarelo	0	1	0	0
Vermelho	0	0	1	0
Verde	0	0	0	1

Onde cada coluna representa a nova feature criada e cada linha representa o valor original da feature x_i .

D. Regularização

Regularização é um método utilizado para *evitar* overfitting dos dados. É realizado de forma a penalizar os valores de θ para que estes mantenham valores pequenos, sendo menos propenso ao overfitting. Isto é feito na forma de um novo parâmetro λ que definirá o quanto cada θ será penalizado:

$$\theta_j := \theta_j * (1 - \lambda * \frac{\alpha}{m}) - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^i) - y^i) x^i \quad (9)$$

Este valor deve ser modificado de forma a melhorar o resultado final. Se o valor for muito alto, a influência da regularização pode ser negativa, criando um aumento no erro final. Se muito baixo, não terá o efeito desejado. Vale mencionar que a regularização na equação normal se dá pela forma

$$\theta = (X^T X + \lambda I)^{-1} X^T y \quad (10)$$

, onde I é a matriz identidade.

III. PROPOSED SOLUTIONS

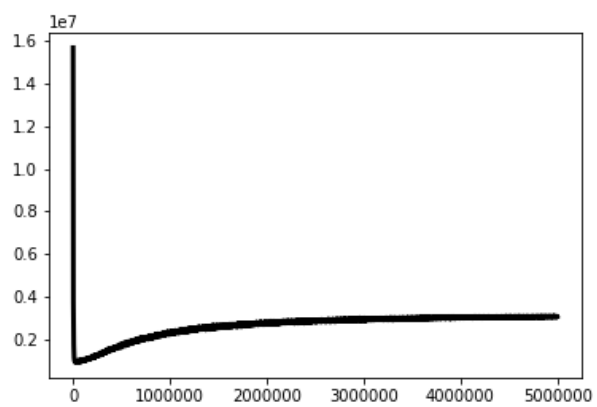
IV. EXPERIMENTS AND DISCUSSION

V. CONCLUSIONS AND FUTURE WORK

The main conclusions of the work as well as some future directions for other people interested in continuing this work. [1]

REFERENCES

- [1] C. M. Bishop, "Pattern recognition and machine learning," 2006.



(a) $\lambda = 0.0000005$, learning rate = 0.01