



APRENDIZAGEM DE MÁQUINA

PROF. JOSENALDE OLIVEIRA

josenalde.oliveira@ufrn.br

<https://github.com/josenalde/machinelearning>

ANÁLISE E DESENVOLVIMENTO DE SISTEMAS - UFRN

MODELOS LINEARES DE REGRESSÃO

Machine Learning SUPERVISIONADA: estimar modelos que, embora simplificações da realidade, apresentem a melhor aderência possível entre os valores **reais (observados)** e os valores preditos

Esta é a equação geral da Regressão Linear Múltipla, com n features. Se considerarmos apenas uma, o objetivo é encontrar uma equação que apresente a relação entre uma variável dependente (target) e uma variável explicativa (feature)

$$\hat{y} = \theta_0 + \theta_1 x_1 + \theta_2 x_2 + \dots + \theta_n x_n, \text{ RLM}$$

$$\hat{y} = \theta_0 + \theta_1 x_1, \text{ RLS, onde } \theta_0 : \text{viés, intercepto, coef. linear}$$

x_i : é a i -ésima feature e θ_j é o j -ésimo parâmetro do modelo

Cada parâmetro do modelo (com exceção do intercepto) na verdade reflete o PESO de determinada feature, então vamos adotar uma notação mais comum em aprendizado de máquina, o peso como sendo a letra W

$$\hat{y} = w_0 + w_1 x_1 + w_2 x_2 + \dots + w_n x_n = \mathbf{W} \cdot \mathbf{x} = \mathbf{W}^T \mathbf{x}$$

MODELOS LINEARES DE REGRESSÃO

Notação

Medidas de desempenho comum em Regressão

Scikit-Learn.metrics – root_mean_squared_error (scikit-learn >= 1.4)

Raiz do Erro Médio Quadrático: $RMSE = \sqrt{\frac{1}{m} \sum_{i=1}^m (\mathbf{W}^T \mathbf{x}^{(i)} - y^{(i)})^2}$

Peso maior para grandes erros

Erro Médio Absoluto: $MAE = \frac{1}{m} \sum_{i=1}^m |\mathbf{W}^T \mathbf{x}^{(i)} - y^{(i)}|$

Indicado para bases com outliers...

m: número de instâncias no conjunto de dados

$\mathbf{x}^{(i)}$ é um vetor de todos os valores das features (excluindo o rótulo) da **i** – **esima** instância
 $y^{(i)}$ é seu rótulo, valor desejado da saída para aquela instância

$$\mathbf{x}^{(1)} = \begin{pmatrix} -118.20 \\ 33.91 \\ 1416 \\ 38372 \end{pmatrix} \quad y^{(1)} = 156400 \quad \mathbf{X} = \begin{pmatrix} (\mathbf{x}^{(1)})^T \\ (\mathbf{x}^{(2)})^T \\ \vdots \\ (\mathbf{x}^{(m)})^T \end{pmatrix} = \begin{pmatrix} -118.29 & 33.91 & 1416 & 38372 \\ \vdots & \vdots & \vdots & \vdots \end{pmatrix}$$

Exemplo: longitude, latitude, n. de habitantes, renda média anual com saída valor médio da moradia

MODELOS LINEARES DE REGRESSÃO

Notação

Medidas de desempenho comum em Regressão

O R-quadrado é uma medida estatística também conhecida como **coeficiente de determinação**, que serve para avaliar a existência de uma relação útil entre a variável dependente(Y) e as variáveis independentes (Xi) em um modelo de regressão.

Ele representa a proporção da variância de Y que foi explicada pelas variáveis independentes no modelo.

Ele fornece uma indicação da qualidade do ajuste e, portanto, uma medida de quão bem as amostras não vistas podem ser previstas pelo modelo, por meio da proporção da variância explicada.

O R-quadrado está sempre entre 0 e 1:

- * 0 indica que o modelo não explica nada da variabilidade dos dados de resposta ao redor de sua média.

- * 1 indica que o modelo explica toda a variabilidade dos dados de resposta ao redor de sua média.

Em geral, quanto maior o R-quadrado, melhor o modelo se ajusta aos seus dados.

$$R^2(y, \hat{y}) = 1 - \frac{\sum_{i=1}^n (y_i - \hat{y}_i)^2}{\sum_{i=1}^n (y_i - \bar{y})^2}$$

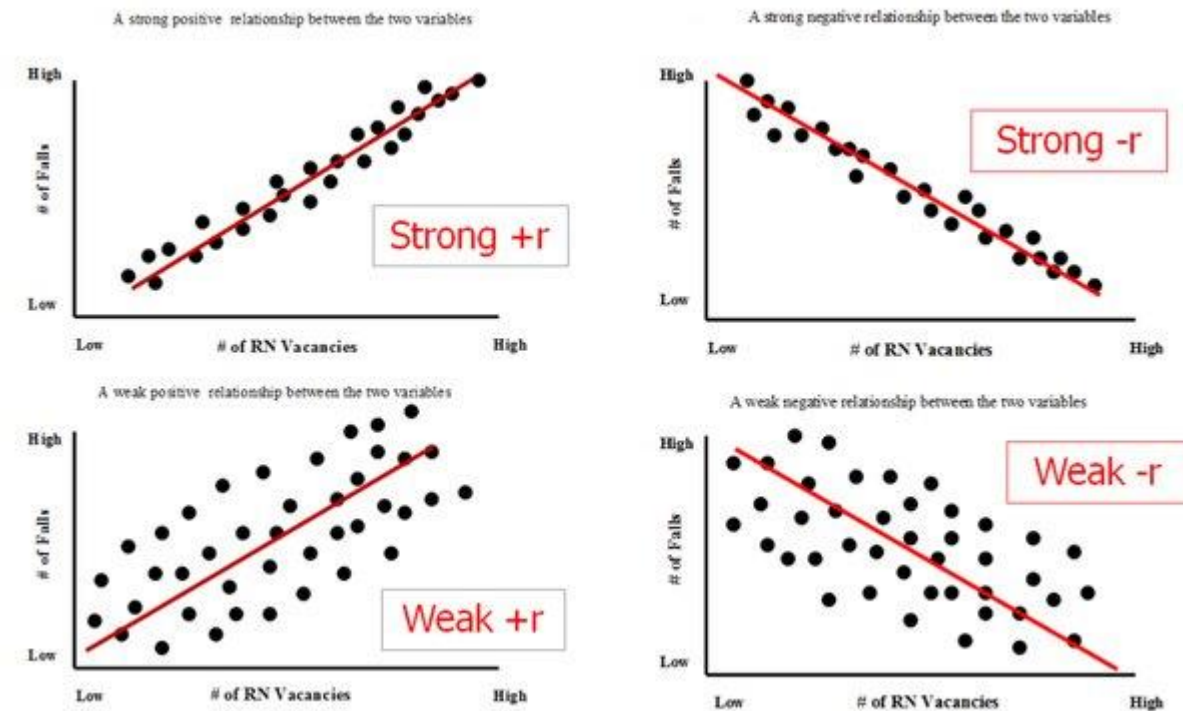
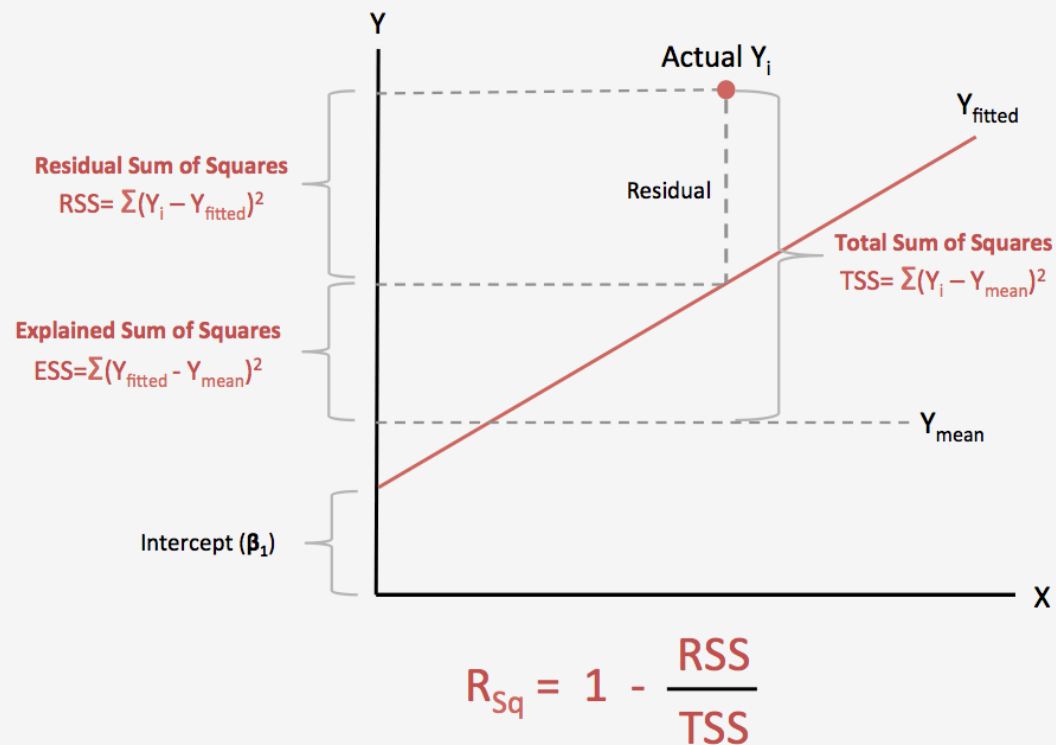
MODELOS LINEARES DE REGRESSÃO

Notação

Medidas de desempenho comum em Regressão

Intervalo de Confiança 95%

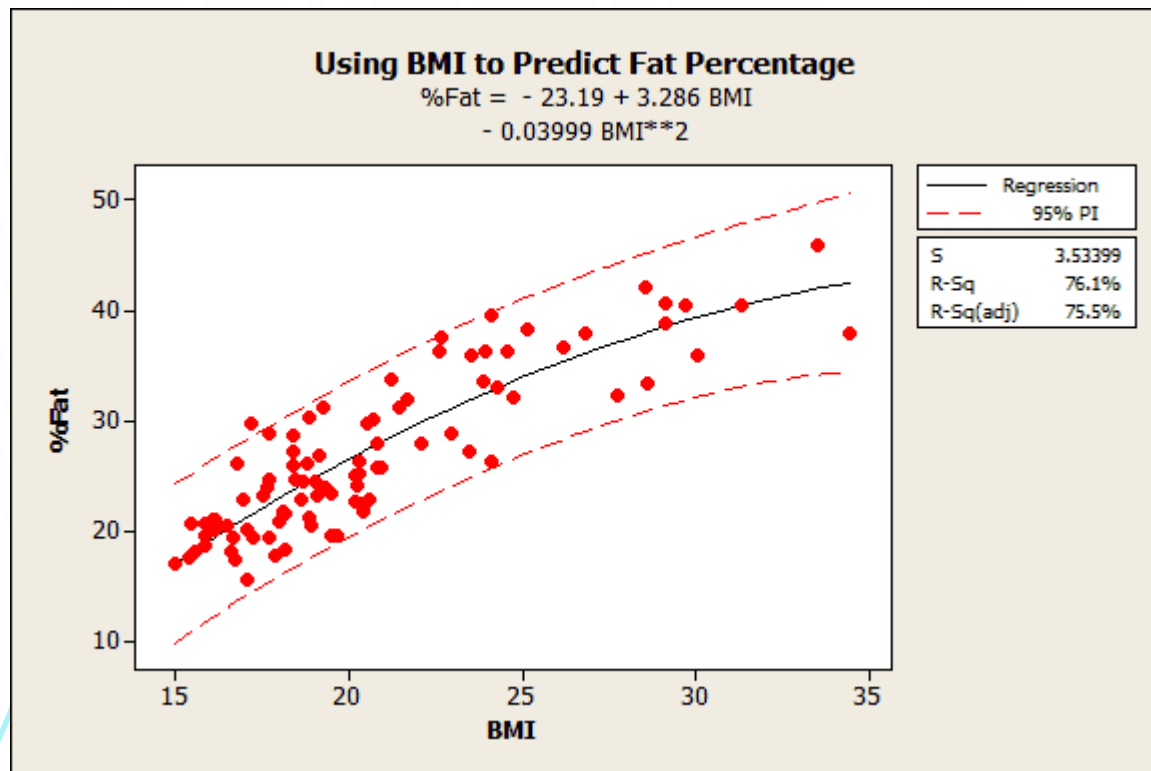
R-Squared Explanation



MODELOS LINEARES DE REGRESSÃO

Notação

Medidas de desempenho comum em Regressão



Existe uma **probabilidade de 95%** de que, no futuro, o verdadeiro valor do parâmetro da população (por exemplo, média) caia no intervalo **X** (limite inferior) e **Y** (limite superior).

O intervalo de confiança é importante para indicar a margem de incerteza (ou imprecisão) frente a um cálculo efetuado. Esse cálculo usa a amostra do estudo para estimar o tamanho real do resultado na população de origem.

O cálculo de um intervalo de confiança é uma estratégia que considera a amostragem de erro. A dimensão do resultado do seu estudo e seu intervalo de confiança caracterizam os valores presumíveis para a população original.

Quanto mais estreito for o intervalo de confiança, maior é a probabilidade da porcentagem da população de estudo representar o número real da população de origem dando maior certeza quanto ao resultado do objeto de estudo.

<https://www.significados.com.br/intervalo-de-confianca/>

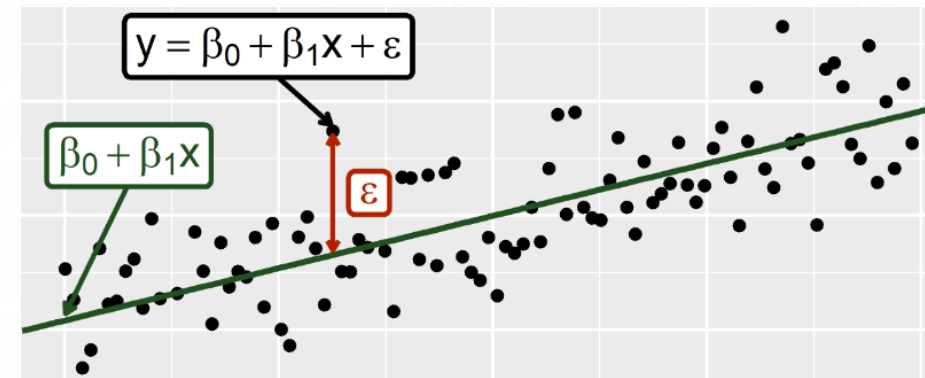
MODELOS LINEARES DE REGRESSÃO

A obtenção dos pesos do modelo (e do bias) consiste em resolver um problema de otimização de **minimização** da função de perda do erro quadrático médio (mean squared error – MSE) entre os valores observados e os valores preditos com a restrição de que este erro seja nulo

Ou seja, encontrar o vetor de pesos tal que $\min \frac{1}{m} \sum_{i=1}^m \left(\mathbf{W}^T \mathbf{x}^{(i)} - y^{(i)} \right)^2$, s.a. $\epsilon = 0$ onde $y = w_0 + w_1 x_1 + \epsilon$

Uma solução é a equação normal $\hat{\mathbf{W}} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X} \mathbf{y}$

Mas há uma inversa de produto de matrizes a resolver $(n+1) \times (n+1)$. A classe **LinearRegression** do scikit-learn é baseada na função `scipy.linalg.lstsq()` – mínimos quadrados, que usa uma versão mais efetiva com decomposição de valores singulares (SVD).



Complexidade eq. normal: $O(n^{2.4})$ a $O(n^3)$

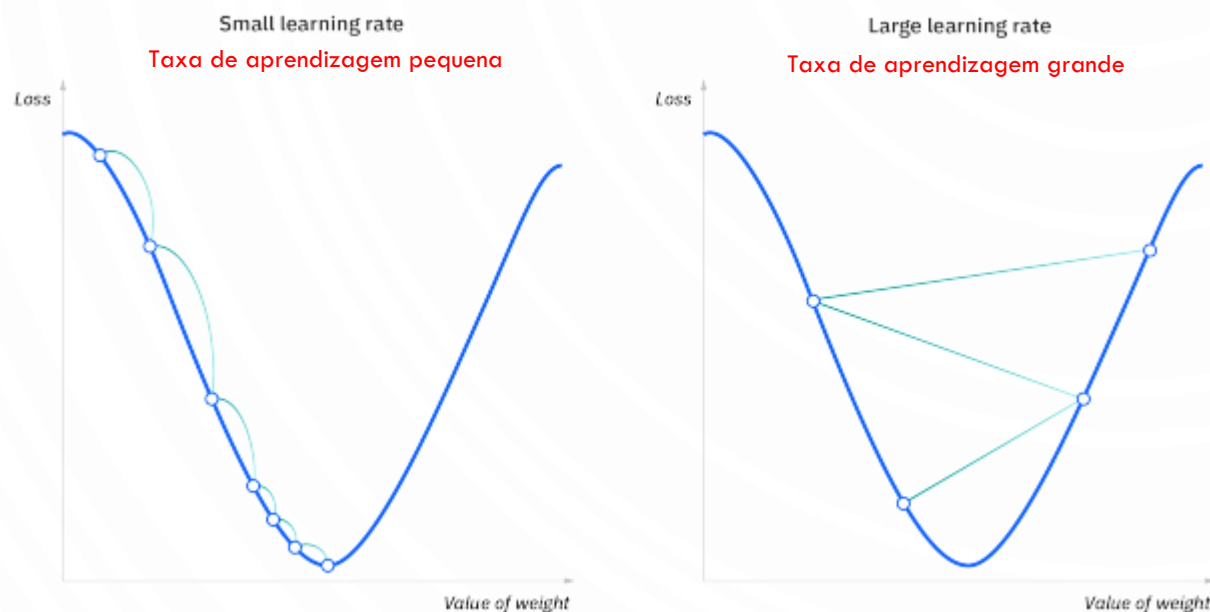
Complexidade SVD scikit-learn: $O(n^2)$

São contudo eficientes para lidar com grandes conjuntos de treinamento: $O(m)$

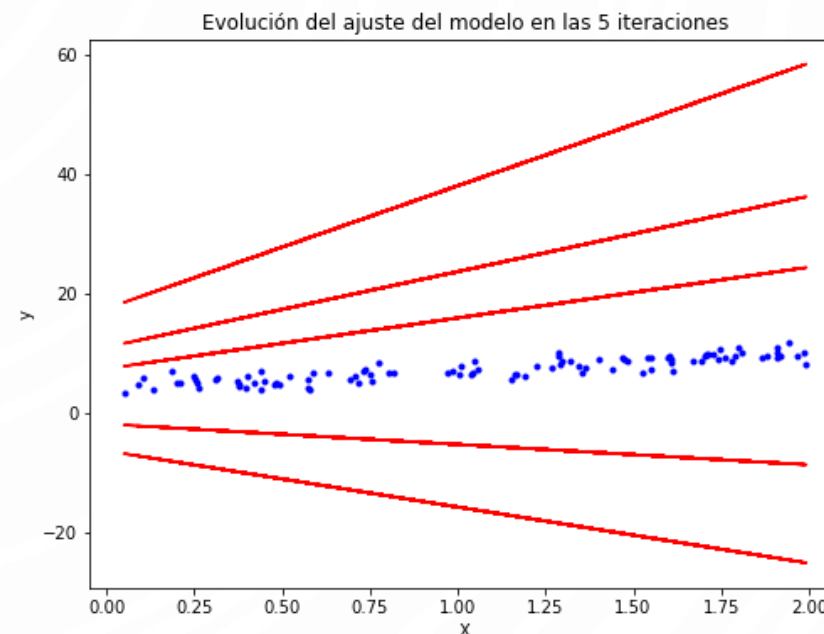
Para predição são lineares para número de instâncias e de características

GRADIENTE DESCENDENTE - GD

O gradiente descendente é um algoritmo de otimização que costuma ser usado para treinar modelos de aprendizado de máquina e redes neurais. Ele treina modelos de aprendizado de máquina minimizando os erros entre os resultados previstos e os reais, ou seja minimizando a função de perda associada.



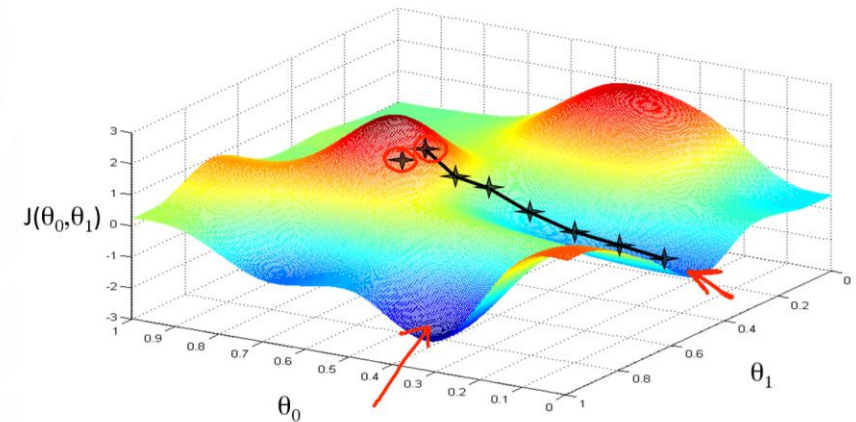
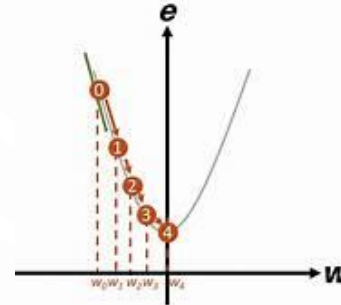
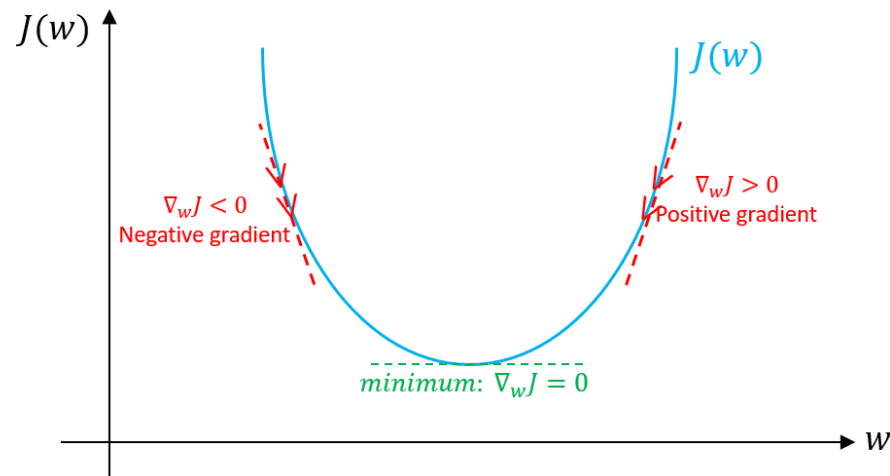
Fonte: [IBM](#)



Aprendizado

Parâmetros do modelo inicializados aleatoriamente. A cada passo calcula o gradiente (ou seja, a derivada parcial da função de perda em relação ao vetor de pesos W e segue em direção ao gradiente descendente, até o zero (convergir para o mínimo). O tamanho da etapa de aprendizado é proporcional à inclinação da função de perda, logo as etapas ficam gradualmente menores à medida que se os pesos se aproximam do mínimo

GRADIENTE DESCENDENTE EM LOTE (BATCH) - GD



É necessário calcular o gradiente da função de perda em relação à cada peso do modelo, ou seja, quanto a função mudará caso se modifique um pouco o peso. A cada etapa TODO o conjunto (lote de dados) é usado, sendo lento para conjunto de treinamento muito grandes. Contudo, escala bem o número de características, sendo muito mais rápido do que a equação normal ou SVD.

$$\frac{\partial}{\partial w_j} MSE(\mathbf{W}) = \frac{2}{m} \sum_{i=1}^m \left(\mathbf{W}^T \mathbf{x}^{(i)} - y^{(i)} \right) \mathbf{x}_j^{(i)}$$

$$\nabla_w MSE(\mathbf{W}) = \begin{pmatrix} \frac{\partial}{\partial w_0} MSE(\mathbf{W}) \\ \frac{\partial}{\partial w_1} MSE(\mathbf{W}) \\ \vdots \\ \frac{\partial}{\partial w_n} MSE(\mathbf{W}) \end{pmatrix} = \frac{2}{m} \mathbf{X}^T (\mathbf{XW} - \mathbf{y})$$

GRADIENTE DESCENDENTE EM LOTE (BATCH) - GD

Expressão de aprendizado – etapa do GD, supondo iterador k :

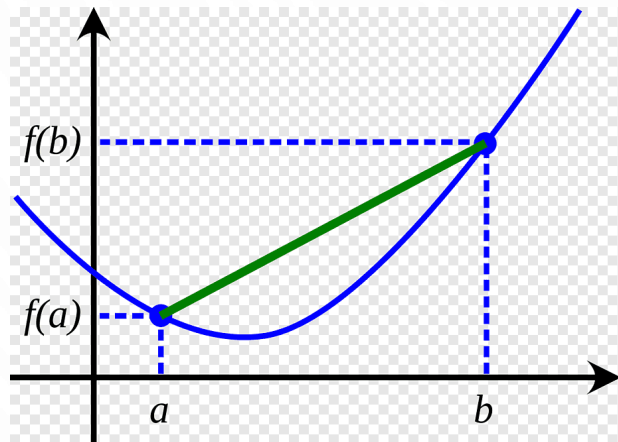
$$\mathbf{W}(k+1) = \mathbf{W}(k) - \eta \nabla_{\mathbf{w}} \text{MSE}(\mathbf{W}), \text{ onde } \eta: \text{ taxa de aprendizagem}$$

Taxa de aprendizagem baixa: muitas iterações para convergir, lento

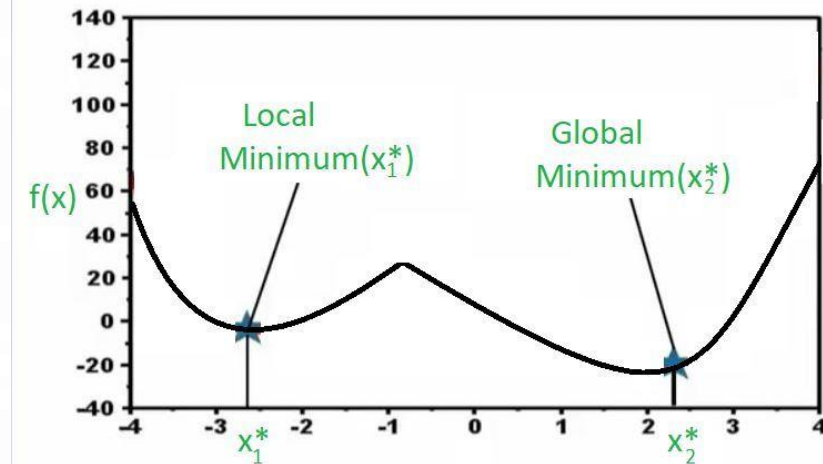
Taxa de aprendizagem alta: pode atravessar o vale e ir para o outro lado em lugar ainda mais alto do que onde estava, divergindo

Com funções convexas (como MSE) se pressupõe que não há mínimos locais, apenas um mínimo global e não muda abruptamente. A BUSCA é feita no espaço de pesos do modelo. **Quanto mais pesos, mais complexa a busca.**

Contudo é muito importante que as features do conjunto de treinamento tenham a mesma escala!



Dados dois pontos, o segmento de reta que os une nunca cruza a curva



GRADIENTE DESCENDENTE ESTOCÁSTICO (SGD)

Seleciona uma instância **aleatória** no conjunto de treinamento a cada etapa e calcula os gradientes com base apenas nesta instância única

Mais rápido, menos dados a cada iteração

Treinamento em grandes conjuntos de dados, pois só uma precisa estar na memória a cada iteração

Contudo, é aleatório, menos regular que o em batch. Não desce suave, mas fica oscilando, pode gerar solução subótima

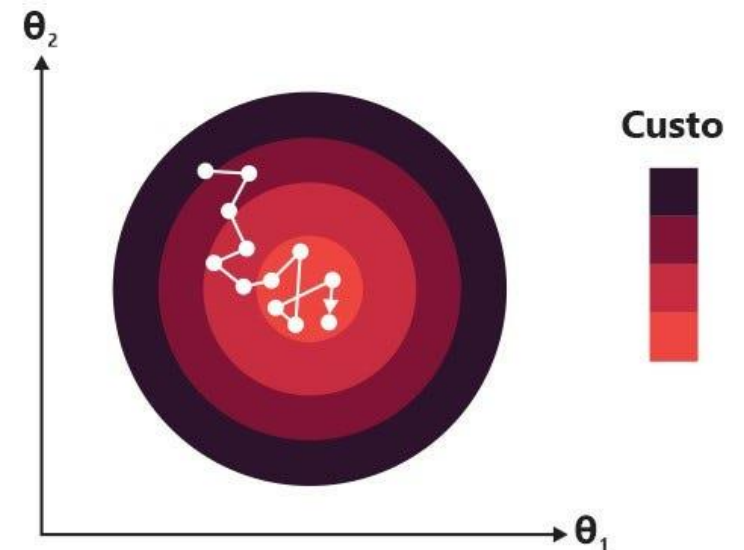
Para funções de perda irregulares (diferente da MSE) pode ajudar a fugir de mínimos locais e chegar ao global

Mas pode nunca se estabilizar no mínimo, e existem soluções de reduzir gradualmente a taxa de aprendizado

Implementar CRONOGRAMA DE APRENDIZADO

No Scikit-Learn, SGDRegressor

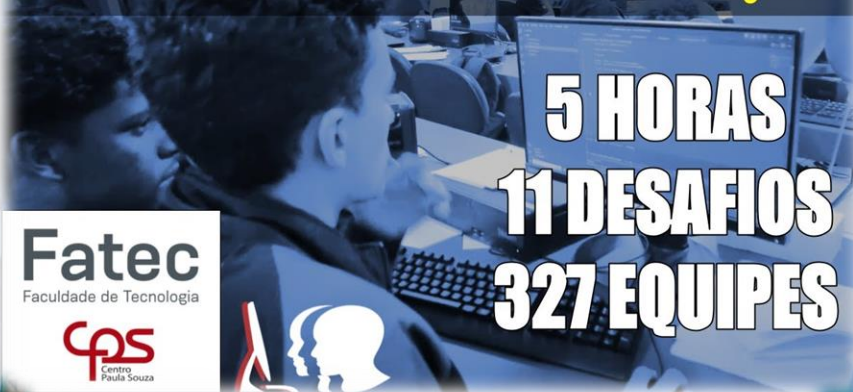
Embaralhar o conjunto de treinamento junto com os rótulos, e as instâncias precisam ser independentes



GRADIENTE DESCENDENTE ESTOCÁSTICO (SGD)

ÉPOCA DE TREINAMENTO: processo de FORWARD PASS (calcular as saídas) e ao processo de BACKWARD PASS (atualizar os pesos) em todo o conjunto de treinamento

MARATONA DE PROGRAMAÇÃO



Bob Burnquist - Skatista



Final Feminina – Carabina – Tiro esportivo

GRADIENTE DESCENDENTE MINI-BATCH

Calcula os gradientes com base em pequenos conjuntos aleatórios de instâncias

Permite obter ganho de desempenho com a otimização de hardware nas operações da matriz, especialmente se usar GPU

O progresso do algoritmo é menos irregular que o GD estocástico, principalmente com grandes mini-batches

Pode ser mais difícil escapar de mínimos locais

- O caminho do batch para no mínimo, SGD e mini-batch prosseguem, perto do mínimo.
- GD batch leva muito tempo para cada época, os outros também alcançariam com um bom cronograma de aprendizado

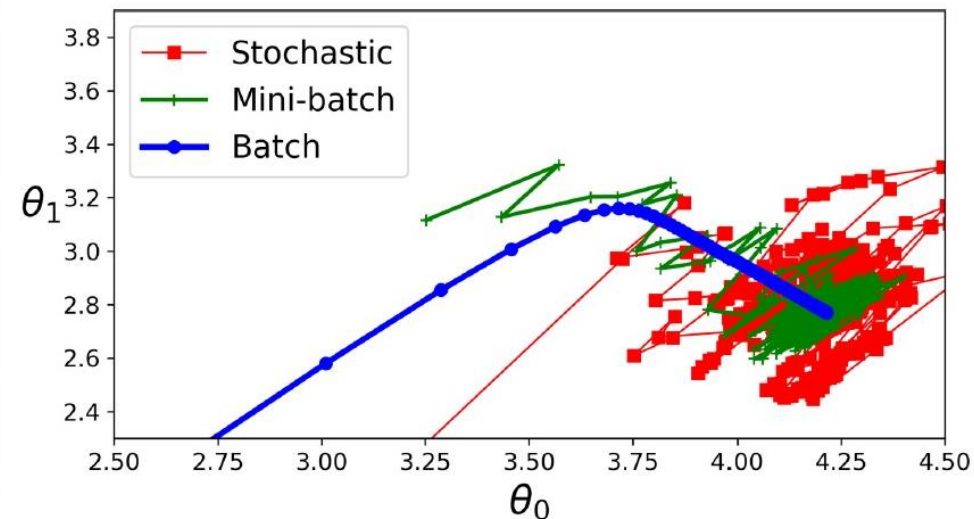


Figure 4-11. Gradient Descent paths in parameter space

(Géron, 2019)

REGRESSORES

Vários algoritmos na scikit-learn conhecidos como CLASSIFICADORES (Classifier), tem suas versões REGRESSORES. Exemplo:

KNeighborsRegressor (vizinhos mais próximos)

MLPRegressor (redes neurais)

DecisionTreeRegressor (árvore de decisão)

SGDRegressor

SVMRegressor (SVM)

RandomForestRegressor

XGBoost Regressor

etc.

Em geral se constrói uma classe para validação cruzada com Grid Search e comparação dos modelos

Várias empresas oferecem o serviço com modelos deste tipo. Exemplo: TOTVS

CUIDADOS AO USAR REGRESSÃO

Conforme lembrado pela AWS, analisar os resíduos do modelo obtido, para checar:

Relacionamento linear

Deve existir uma relação linear entre as variáveis independentes e dependentes. Para determinar essa relação, os cientistas de dados criam um gráfico de dispersão (uma coleção aleatória de valores x e y) para ver se eles se encaixam em uma linha reta. Se isso não acontecer, será possível aplicar funções não lineares, como raiz quadrada ou logarítmica, para criar matematicamente a relação linear entre as duas variáveis.

Correlação linear

Independência residual

ACF

Os cientistas de dados usam resíduos para medir a precisão da previsão. Um resíduo é a diferença entre os dados observados e o valor previsto. Resíduos não devem ter um padrão identificável entre eles. Por exemplo, os resíduos não devem crescer com o passar do tempo. Diferentes testes matemáticos, como o teste de Durbin-Watson, podem ser utilizados para determinar a independência residual. É possível usar dados fictícios para substituir qualquer variação de dados, como dados sazonais.

Testes Durbin-Watson...

Normalidade

Técnicas de representação gráfica, como gráficos Q-Q, determinam se os resíduos são normalmente distribuídos. Os resíduos devem se encaixar ao longo de uma linha diagonal no centro do gráfico. Se eles não forem normalizados, você poderá testar os dados quanto à presença de outliers aleatórios ou valores atípicos. Remover desses outliers ou realizar transformações não lineares pode corrigir o problema.

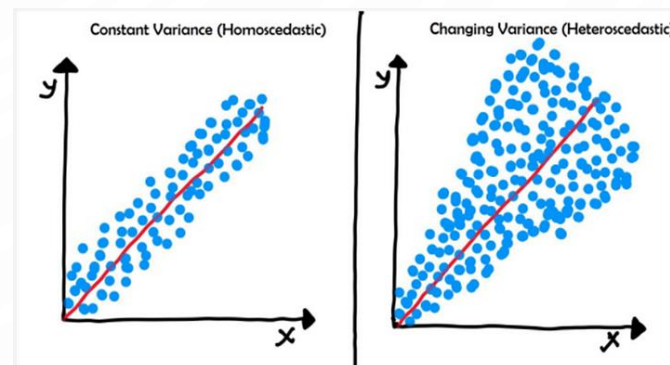
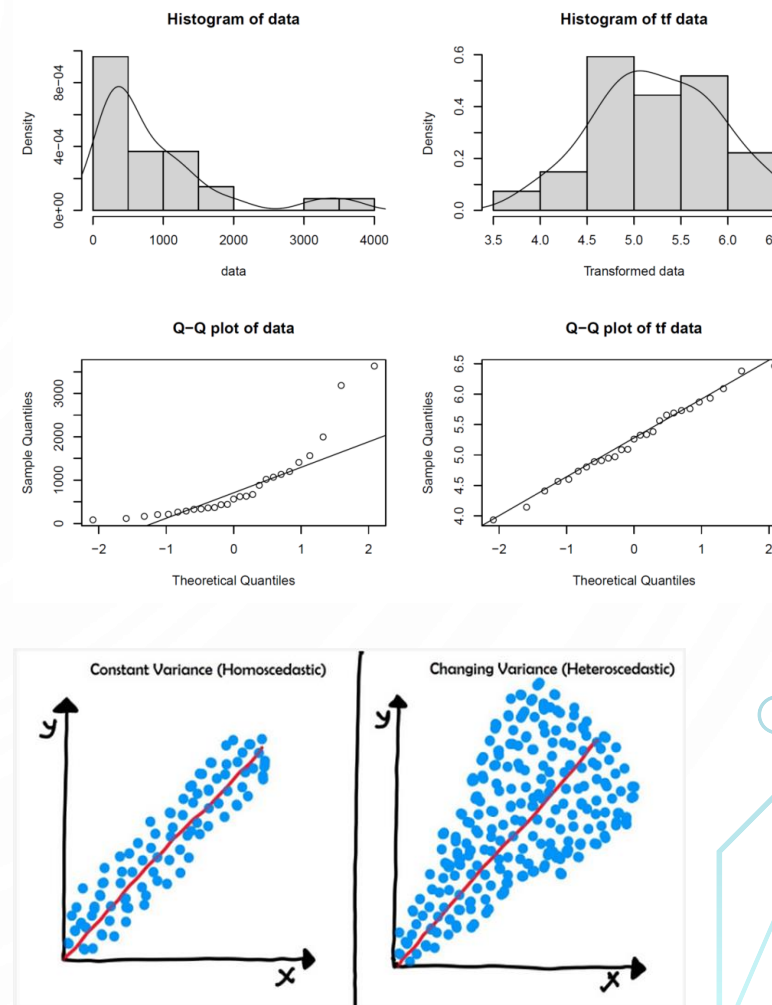
Box-cox

Testes Shapiro-wilk/francia, Jarque-berba...

Homocedasticidade

Testes Breusch-Pagan, White...

A homocedasticidade supõe que os resíduos tenham uma variância constante ou desvio padrão da média para cada valor de x . Se esse não for o caso, os resultados da análise podem não ser precisos. Se essa suposição não for atendida, talvez seja necessário alterar a variável dependente. Como a variância ocorre naturalmente em conjuntos de dados grandes, faz sentido alterar a escala da variável dependente. Por exemplo, em vez de usar o tamanho da população para prever o número de postos de bombeiros em uma cidade, é possível usar o tamanho da população para prever o número de postos de bombeiros por pessoa.



CUIDADOS AO USAR REGRESSÃO

Apesar de poderosa, a Análise de Regressão Linear exige fortes suposições para sua utilização, como normalidade, independência e homoscedasticidade dos erros. Foi a partir dessas limitações que surgiram os Modelos Lineares Generalizados.

Quando a correlação linear target-feature não é observada, existem os modelos lineares generalizados

Permite que a target tenha distribuição qualquer, e usa uma **função de ligação** que varia **linearmente** com as features, relacionando a média da variável target com a combinação linear das features

A target pode ter forma de contagem, contínuas simétricas e assimétricas, binárias e categóricas. É preciso portanto **conhecer a distribuição da variável target**, para escolha da função de ligação

Tipo de regressão	utilização	exemplo de uso
Poisson	usada para modelar dados de contagem.	número de mortes em determinada região ou o número de consumidores que entram em um estabelecimento comercial.
Bernoulli	utilizada na modelagem de fenômenos que podem ser resumidos em uma variável binária, ou seja, se ocorreu ou não um evento.	modelos de concessão de crédito ou em pesquisas clínicas que tem como objetivo verificar os fatores de influência na ocorrência ou não de uma determinada doença.
Gama	usada para modelar dados positivos e assimétricos. A regressão Gama modela variáveis contínuas.	estudo dos fatores que influenciam no valor de um imóvel ou ainda os fatores que influenciam na demanda de produtos em diferentes centros de distribuição.

CUIDADOS AO USAR REGRESSÃO

$\eta_i = \beta_0 + \beta_1 x_{1i} + \beta_2 x_{2i} + \dots + \beta_k x_{ki}$, funções ligam a η

Modelo de regressão	Característica target	Distribuição	Função de ligação
Linear	Quantitativa	Normal	\hat{y}
Box-cox	Quantitativa	Normal após transformação	$\frac{\hat{y}^\lambda - 1}{\lambda}$
Logística binária	Qualitativa com C=2	Bernoulli	$\ln\left(\frac{p}{1-p}\right)$
Logística multinomial	Qualitativa com C>2	Binomial	$\ln\left(\frac{p_c}{1-p_c}\right)$
Poisson	Contagem (int e > 0)	Poisson	$\ln(\lambda_{poisson})$
Binomial negativo	Contagem (int e > 0)	Poisson-Gama	$\ln(\lambda_{bneg})$

EXEMPLOS DE PROBLEMAS...

Contextualizando a Distribuição de Bernoulli, tem-se a seguinte situação: a experiência tem mostrado que até fevereiro o motorista que é parado em uma *blitz* tem 60% de chance de estar adimplente em relação ao Imposto sobre a Propriedade de Veículos Automotores (IPVA). Temos, portanto, uma probabilidade de sucesso (**o motorista não estar devendo o IPVA**) de 0,6 e uma probabilidade de estar devendo de 0,4 (**vem da diferença $q = 1 - 0,6$**).

O problema da loja de Roupas do Martin

Considere as decisões de compra dos próximos três clientes que entram na loja de roupas do Martin. Com base em experiências passadas, o gerente da loja estima que a probabilidade de que qualquer um dos clientes comprará é de 0,30. Qual é a probabilidade de que dois primeiros dos próximos três clientes realizarão a compra?

Exemplo: A análise dos dados dos últimos anos de uma empresa de energia elétrica forneceu o valor médio de um blecaute por ano. Pense na probabilidade de isso ocorrer no próximo ano:

- a) Nenhum blecaute.
- b) De 2 a 4 blecautes.
- c) No máximo 2 blecautes.

Observe que o exemplo afirma que a cada ano acontece em média um blecaute, ou seja, o **número de sucesso ocorrido em um intervalo contínuo**. Verificamos que a variável tem Distribuição Poisson:

$$P(X = x) = \frac{e^{-\mu} \mu^x}{x!}$$