

PRIMEIRA AVALIAÇÃO DE APRENDIZADO ESTATÍSTICO DE MÁQUINA

Miguel Zanchettin de Oliveira*

Maio de 2024

EXERCÍCIO 01 (ML)

Diferencie os seguintes termos e utilize um exemplo.

01.A FUNÇÃO DE PERDA E ERRO

A resposta do exercício está em: Resposta - Exercício 01.A

01.B VALIDAÇÃO E DATA SPLITTING

A resposta do exercício está em: Resposta - Exercício 01.B

01.C OVERFITTING E COMPLEXIDADE

A resposta do exercício está em: Resposta - Exercício 01.C

01.D RISCO E RISCO EMPÍRICO

A resposta do exercício está em: Resposta - Exercício 01.D

01.E DIMENSÃO VC E COEF. DE SHATTERING

A resposta do exercício está em: Resposta - Exercício 01.E

EXERCÍCIO 02

Escolha dois métodos de aprendizado de máquina diferentes de SVM ou KNN. Pesquise sobre a Dimensão VC de cada um. Crie um texto explicativo sobre sua pesquisa e vincule com o que foi visto na disciplina.

Qualquer pesquisa relacionada ao tema revelará, de forma rápida, que são pouco frequentes aquelas capazes de concluir algo. Principalmente, que têm forma fechada e simplificada de resultado. Isto porque, denotar a possibilidade

*Mestrando do Programa de Pós-Graduação em Métodos Numéricos Aplicados à Engenharia - PPGMNE, da Universidade Federal do Paraná - UFPR.

de dimensões VC para quaisquer modelos demonstra-se de relativa dificuldade quando por métodos dedutivos. Embora sejam difíceis, tais teorias não são impossíveis. Para esta resposta, serão expostas as conclusões para redes neurais feedforward e as redes neurais com funções de ativação Piecewise polinomial, dois métodos de classificação com modelos similares, porém distintos para a classificação.

As redes neurais feedforward, com saídas binárias, possuem dimensão VC $O(w \cdot \log w)$ [2]. Sendo w a quantidade de pesos e O a utilização da Big-O notation, medida de complexidade de um algoritmo. A definição formal realizada pelo autor está em seu Teorema 2.1, citado abaixo.

Seja \mathcal{N} uma rede neural feedforward arbitrária com w pesos que consiste em portões de limiar linear. Então, $VCdim(\mathcal{N}) = O(w \cdot \log w)$ [3](Tradução pelo autor).

Quanto às redes neurais com funções de ativação Piecewise polinomial, sua dimensão VC é descrita como $wl \log w + wl^2$, sendo w a quantidade de pesos e l , de camadas.

Suponha que \mathcal{N} é uma rede feed-forward com w pesos, l camadas, e todos os portões não de saída tendo uma função de ativação polinomial em partes fixas com um número fixo de partes. Então, $VCdim(\mathcal{N}) = O(wl \log w + wl^2)$. [1].

Assim, percebe-se que, pelas dimensões VC demonstradas pelos autores, seria possível identificar, dada a dimensão VC, um mínimo global para o risco empírico nas funções, porém, este problema recai sobre uma estrutura ótima de formulação das redes neurais em si e das matrizes de pesos utilizadas. Em função disso, ambos os modelos não têm convergência garantida para o mínimo global, uma vez que, dadas quantidades de w ou l , diferentes mínimos locais são passíveis de aparecimento.

EXERCÍCIO 03 (ML)

Leia artigo RAPER, Simon. Leo Breiman's "two cultures". Significance, v. 17, n. 1, p. 34-37, 2020. Acesso aqui. Argumente relacionando o artigo com os conteúdos de aprendizado de máquina vistos na disciplina.

A resposta do exercício, elaborada à mão, está nos anexos (Resposta - Exercício 03).

EXERCÍCIO 04

Pesquise duas métricas de capacidade preditiva não vistas em sala de aula para o contexto de classificação. Implemente a versão micro das métricas escolhidas. Utilize os dados Wall-Following Robot Navigation e o classificador KNN.

O código utilizado para responder o exercício está em: Código: Exercício 4. Os métodos utilizados foram: Recall e F-Beta Score. Embora o Recall seja bem comum, não foi apresentado em aula, já o F-Beta Score é, de fato, menos convencional.

EXERCÍCIO 05

Escolha um conjunto de dados para classificação binária e de seu interesse. Aplique os métodos de validação holdout, holdout repetido, k-fold, k-fold repetido e leave-one-out. Compare os resultados.

Os resultados estão disponíveis na tabela 1. Para validar um modelo de regressão em uma base de dados binária gerada aleatoriamente na biblioteca Sickit-Learn do Python, foi utilizado o Matthews Correlation Coefficient (MCC), pois costuma ser um dos métodos mais utilizados em pesquisas de Machine Learning, conforme informado pelo professor em classe.

Metodologia de validação	Reg. Log.
Houldout Simples	0.986681
Houldout Repetido (20 vezes)	0.986681
KFold	0.984071
KFold Repetido (20 vezes)	0.985376
Leave One Out	0.993333

Table 1: Resultados calculados

Os códigos utilizados, por sua vez, estão disponíveis em Código: Exercício 5.

EXERCÍCIO 06 (ML)

Considere $n = 3$ observações, $Y = -1, +1$ e $\mathcal{X} = \mathbf{R}$ e $G = I_{[a,b]} | a < b \in \mathbf{R}$. Calcule o coeficiente de shattering $\mathcal{M}(G, n)$.

A resposta do exercício, elaborada à mão, está nos anexos (Resposta - Exercício 06 e Resposta - Exercício 06 (Continuação)).

EXERCÍCIO 07 (ML)

Mostre matematicamente que o princípio de minimização do risco empírico não é consistente para o coef. de shattering com crescimento exponencial. Interprete o resultado.

A resposta do exercício, elaborada à mão, está nos anexos (Resposta - Exercício 07 e Resposta - Exercício 07 (Continuação)).

EXERCÍCIO 08 (ML)

Justifique o Método de Máquina de Vetores de Suporte com base na Teoria do Aprendizado Estatístico.

A resposta do exercício, elaborada à mão, está nos anexos (Resposta - Exercício 08).

ANEXOS

RESOLUÇÃO À MÃO LIVRE - EXERCICIO 01.A

Questão 01 - ①: referência:

a. função de perda e erro

O erro é a medida de quanto o resultado predito está distante, ou próximo, do resultado. Isto para cada observação. As funções de perda, por sua vez, são formas de tornar o vetor de erros em um escalar.

Ex:

$$\begin{matrix} \vec{y} & & \vec{\hat{y}} \\ \begin{bmatrix} y_1 \\ \vdots \\ y_n \end{bmatrix} & - & \begin{bmatrix} \hat{y}_1 \\ \vdots \\ \hat{y}_n \end{bmatrix} \end{matrix} \text{ erro} \quad \underbrace{\left\{ \sum_{i=1}^n |\hat{y}_i - y_i|, \sum_{i=1}^n (\hat{y}_i - y_i)^2, \dots \right\}}_{\text{funções de perda}}$$

y = vetor observado
 \hat{y} = vetor predito

Figure 1: Resposta - Exercício 01.A

RESOLUÇÃO À MÃO LIVRE - EXERCICIO 01.B

b. Validação e Data splitting.

Data Splitting consiste em uma técnica de simular dados não vistos pelo modelo, a fim de estimar qual seria seu resultado frente a novas observações. Possibilitando, assim, treinar modelos em um conjunto comum e avaliá-los em outro conjunto comum, tornando-os comparáveis.

A validação, por sua vez, consiste em formas de melhorar a capacidade de aprendizado do modelo em uma mesma base de testes (ou validação), com a intenção de maximizar a capacidade de generalização do modelo.

Ex.

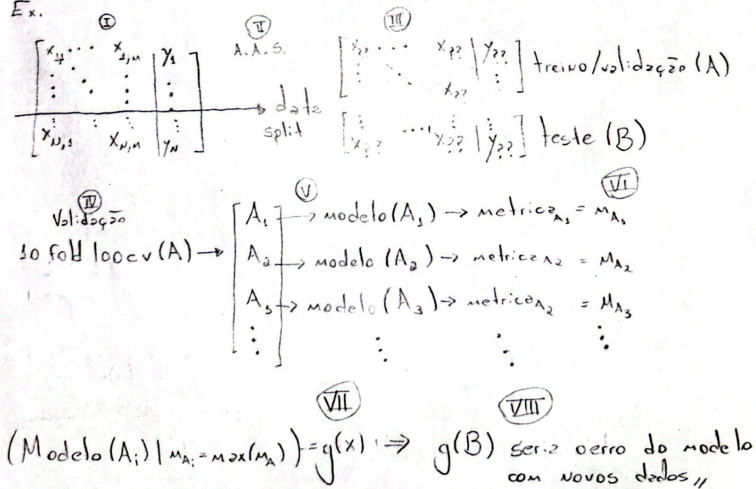


Figure 2: Resposta - Exercício 01.B

RESOLUÇÃO À MÃO LIVRE - EXERCÍCIO 01.C

Questão 01-c. Overfitting e complexidade

Embora sejam conceitos relacionados, não são equivalentes.

Overfitting, "sobreajuste", refere-se a um modelo incapaz de gerar uma generalização a dados não vistos.

A complexidade, por outro lado, refere-se a quantidade de "mudanças" no formato da função. Por exemplo, $f(x) = 1+x$ é uma função mais simples que $f(x) = 1+x^2$, uma vez que a primeira possui apenas um valor possível de $f(x)$ para cada x , enquanto a segunda, mais de um.

Ambos os conceitos estão relacionados. Espera-se, no Aprendizado supervisionado, que maior complexidade do modelo esteja associada a um melhor desempenho na base de treino, porém uma deterioração no performance na base de testes.

ex:

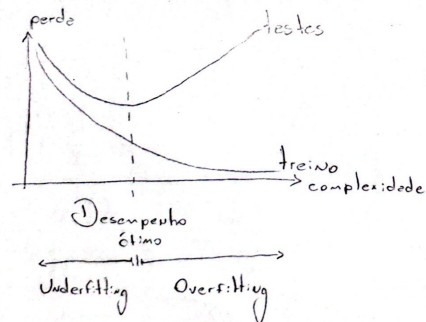


Figure 3: Resposta - Exercício 01.C

RESOLUÇÃO À MÃO LIVRE - EXERCÍCIO 01.D

Questão 1.d - Risco e risco empírico

O risco refere-se à natureza, a aleatoriedade desconhecida. Por outro lado, o risco empírico é a observação de uma amostra do risco. Assim, espera-se que o risco empírico (R_{emp}) convirja para o risco (R) quando a quantidade de observações da amostra (n) tende ao infinito:

$$R_{emp} \xrightarrow{n \rightarrow \infty} R$$

O exemplo clássico da estatística, um dado, representa bem este conceito. Quanto mais lançamentos de um dado realizador, menor será o risco empírico ao prever qual o próximo resultado em um modelo que aprende a cada nova observação. Contudo, mesmo que este experimento se repita infinitas vezes, a predição ainda estará sujeita a aleatoriedade da natureza, isto é

$$R_{emp} \leq R \quad \forall n \in \mathbb{N}^*$$

Figure 4: Resposta - Exercício 01.D

RESOLUÇÃO À MÃO LIVRE - EXERCÍCIO 01.E

Questão 1.E - ① Dimensão VC e Coeficiente de shattering.

A dimensão VC refere-se a maior quantidade de observações (n) de uma amostra tal que uma classe de funções pode realizar corretamente a classificação de todas as observações.

O coeficiente de shattering, por sua vez, representa a maior quantidade de funções capazes de classificar corretamente uma quantidade n de observações.

Por exemplo,
Seja $X = \mathbb{R}^d$ e $\mathcal{F}_P :=$ hiperplanos lineares $\Rightarrow VC(\mathcal{F}) = d+1$ ^{dimensão VC}

Isto leva ao Sauer-Shelah, que Vapnik e Chervonovskis utilizaram para relacionar os dois conceitos:

$$N(\mathcal{F}, n) \leq \sum_{i=0}^d \binom{n}{i}$$

^{coef. de shattering} \leftarrow $d \rightarrow$ VC dimension

Isto faz com que, para $n \geq d$, $N(\mathcal{F}, n) \leq \left(\frac{en}{d}\right)^d$

Portanto,

$$\begin{array}{c} \text{VC} = \infty \quad n = d \quad \text{VC} \neq \infty \\ \hline \underbrace{\hspace{10em}}_{N(\mathcal{F}, n) \approx 2^n} \quad \underbrace{\hspace{10em}}_{N(\mathcal{F}, n) = n^d} \quad \xrightarrow{\quad} n \text{ (observações na amostra)} \end{array}$$

\hookrightarrow possibilita que $R_{\text{emp}} \xrightarrow{n \rightarrow \infty} R$ (convergência),

Figure 5: Resposta - Exercício 01.E

RESOLUÇÃO À MÃO LIVRE - EXERCÍCIO 03

Exercício 03 - Avaliação 01

O artigo de Réper, publicado em 2020, exemplifica que os problemas mencionados por Breiman, em seu relato original, ainda persistem, mesmo depois de 20 anos.

Breiman, ao retornar de um trabalho de consultoria para o meio acadêmico, deparou-se com uma academia obsoleta, com métodos que já não mais representavam as necessidades práticas da sociedade, deixando com que outras ciências desenvolvessem pesquisas mais relevantes.

É evidente, pelo conteúdo ministrado nesta disciplina, que a Teoria do Aprendizado Estatístico seria motivo de orgulho à luz de Breiman. Incorporando a ciência estatística ao processo de aprendizado de máquina, Vapnik e colaboradores foram capazes de desenvolver uma teoria com relevantes implicações para a sociedade, como o modelo de SVM, ou mesmo o Princípio de Minimização do Risco Empírico.

Contudo, Réper finaliza mencionando desafios criados pela cultura incentivada por Breiman. O autor menciona que, desta vez, mesmo a cultura de modelagem algorítmica tem cedido aos problemas outrora da modelagem de dados, como buscar soluções de machine learning sem mesmo antes estudar o comportamento dos dados. Sendo, assim, os desafios futuros da Teoria do Aprendizado Estatístico.

Figure 6: Resposta - Exercício 03

RESOLUÇÃO À MÃO LIVRE - EXERCÍCIO 06

Questão 6-

$n=3$ observações $Y = \{-1, +1\} = \mathcal{Y} = \mathcal{R}$

Seja $G = \{I_{[a,b]} \mid a < b \in \mathcal{R}\}$, calcular coeficiente de shattering $\mathcal{M}(G, n)$

Dada uma função de perda $L(Y, g(X))$, pode-se definir o risco como sua esperança.

$$R(g) = E_{X,Y} [L(Y, g(X))]$$

Porém, $R(g)$ representa o erro, produzido pela natureza, sendo somente observável sua parte empírica ($R_{\text{emp}}(g)$).

$$R_{\text{emp}}(g) = \frac{1}{n} \sum_{i=1}^n L(Y_i, g(X_i))$$

Assim, segundo o Princípio da Minimização do Risco Empírico na Teoria do Aprendizado estatístico, tem-se que:

$$P(|R_{\text{emp}}(g) - R(g)| > \epsilon) \leq 2e^{-2n\epsilon^2}, \text{ quando } n \rightarrow \infty$$

Considerando o G enunciado como um conjunto de possíveis funções g tal que $g_i \in G$, então

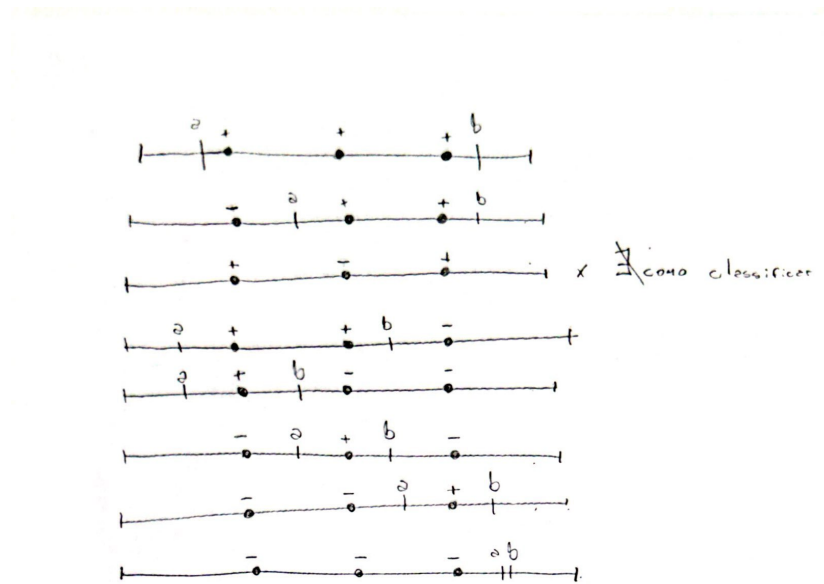
$$P(\sup_{g \in G} |R_{\text{emp}}(g) - R(g)| > \epsilon) \leq 2\mathcal{M}(G, n)e^{-2n\epsilon^2}$$

A definição de $\mathcal{M}(G, n)$ é a quantidade de funções $g \in G$ capazes de classificar n observações.

Assumindo que nenhuma das observações possui valor exatamente igual a a ou b , que $I_{[a,b]}$ é positiva no intervalo e negativa fora dele, pode-se representar os dados apenas nas 8 combinações abaixo:

Figure 7: Resposta - Exercício 06

RESOLUÇÃO À MÃO LIVRE - EXERCÍCIO 06 (CONTINUAÇÃO)



Dado que $a < b$, a única forma em que $I_{[a,b]}$ não seria capaz de prever corretamente os dados, seria quando o valor negativo estivesse entre a e b e positivos fora dele. Assim, das 8 combinações, apenas 7 são classificáveis corretamente. Portanto $M(G, 3) = 7$.

Figure 8: Resposta - Exercício 06 (Continuação)

RESOLUÇÃO À MÃO LIVRE - EXERCÍCIO 07

Exercício 7 - Mostre, matematicamente, que o coeficiente de shattering exponencial não é consistente com o princípio de minimização do risco empírico. Interprete o resultado.

A desigualdade proposta por Vapnik (1) representa o Princípio de Minimização do Risco Empírico:

$$(1) \quad P(\sup_{g \in G} |R_{\text{emp}}(g) - R(g)| > \epsilon) \leq 2me^{-2n\epsilon^2}$$

Seja:

g uma função de predição, G um conjunto de possíveis funções g e g^* a função mais próxima possível da natureza, isto é, a função predição ótima, $R(g^*)$ sendo a esperança de sua função de perda, caso todos os dados fossem conhecidos. $R_{\text{emp}}(g)$, por sua vez, representa a esperança de função de perda de uma função $g \in G$ calculada em uma amostra com n observações. E representa um termo de erro. m , por fim, represente a maior quantidade de funções g capazes de classificar n observações corretamente.

O comportamento de m será fator determinante para garantir que a equação (1) convirja para 0 quando $n \rightarrow \infty$.

A forma do coeficiente de shattering exponencial, como se sabe, é dada pela equação (2).

$$(2) \quad m(n) \approx 2^n$$

Figure 9: Resposta - Exercício 07

RESOLUÇÃO À MÃO LIVRE - EXERCÍCIO 07 (CONTINUAÇÃO)

É possível afirmar que o limite da função (2) tende ao infinito quando $n \rightarrow \infty$ (3).

$$(3) \quad \lim_{n \rightarrow \infty} u(n) \approx \lim_{n \rightarrow \infty} 2^n = \infty$$

Tautologicamente, se o resultado (3) é verdadeiro, pode-se dizer que (1) não converge para 0 quando $n \rightarrow \infty$ (4).

$$(4) \quad P(\sup_{g \in G} |R_{np}(g) - R(g)| > \varepsilon) \leq \infty$$

Isto é o mesmo que dizer que, caso o coeficiente de shattering tende ao ∞ quando $n \rightarrow \infty$, então não existe uma função em G equivalente a q^* . De forma simplificada, não há uma solução ótima para o problema de classificação, para as premissas assumidas.

Figure 10: Resposta - Exercício 07 (Continuação)

RESOLUÇÃO À MÃO LIVRE - EXERCÍCIO 08

Exercício 8 - Justifique o Método de Máquinas de Vetores de Suporte com base na Teoria do Aprendizado Estatístico.

A desigualdade proposta por Vapnik (1) demonstra a possibilidade de existência de uma combinação de coeficiente de shattering n de funções g tal que, mesmo no pior dos casos (supremo) da diferença absoluta entre o risco empírico (R_{emp}) e o risco caso todos os dados fossem conhecidos (R), converja para 0 quando o número de observações n tende ao infinito.

$$(1) P(\sup_{g \in G} |R_{emp}(g) - R(g)| > \epsilon) \leq 2me^{-2n\epsilon^2}$$

Destaca-se que a existência de um grupo de funções G pode ser pensado de forma exclusiva a fim de gerar essa convergência. Assim, surge a ideia do método de Máquinas de Vetores de Suporte (SVM).

Figure 11: Resposta - Exercício 08

CÓDIGOS

CÓDIGO: EXERCÍCIO 4

```
#-----  
# Libraries  
#  
  
import numpy as np  
import pandas as pd  
  
from sklearn.metrics import recall_score  
from sklearn.metrics import fbeta_score  
  
from sklearn.model_selection import train_test_split  
from sklearn.neighbors import KNeighborsClassifier  
  
from rich import print as pprint  
  
#-----  
# Data Handling  
#  
  
# Extract  
url = ('https://raw.githubusercontent.com/'  
       'andersonara/datasets/master/'  
       'wall-robot-navigation.csv')  
  
df = pd.read_csv(url, delimiter=';')  
  
# Transform data into numpy matrices  
X = df[['X1', 'X2']].to_numpy()  
y = np.ravel(df[['Y']])  
  
# Data split  
X_train, X_test, \  
y_train, y_test = train_test_split(X, y,  
                                    train_size=0.7,  
                                    shuffle=True,  
                                    random_state=2002)
```

```

#-----
# Model evaluation
#

def get_best_knn_by_metric(metric,
                           maxiter: int = 10,
                           **kwargs):

    def train_knn(k):
        model = KNeighborsClassifier(n_neighbors=k)
        model.fit(X_train, y_train)

        y_pred = model.predict(X_test)

        error = metric(y_test, y_pred, **kwargs)

        return (model, round(error, 5))

    def k_fit(maxiter):
        metrics = {}

        for k in range(1, maxiter + 1):
            _, metrics[k] = train_knn(k)

        best_k = max(metrics, key=metrics.get)

        return best_k

    def get_result(metric, maxiter):
        r = {}
        r['k'] = k_fit(maxiter)
        r['model'], r['metric_value'] = train_knn(r['k'])
        r['metric'] = metric

        return r

    return get_result(metric, maxiter)

```

```

#-----
# Show results
#

metricas = {
    'Recall Micro': dict(metric=recall_score, average='micro'),
    'Recall Macro': dict(metric=recall_score, average='macro'),
    'F-Beta Micro': dict(metric=fbeta_score, average='micro',
        beta=2),
    'F-Beta Macro': dict(metric=fbeta_score, average='macro',
        beta=2)
}

resultados = {}
for nome, p in metricas.items():
    resultados[nome] = get_best_knn_by_metric(**p)['
        metric_value']

pprint(resultados)

# {
#   'Recall Micro': 0.99023,
#   'Recall Macro': 0.98358,
#   'F-Beta Micro': 0.99023,
#   'F-Beta Macro': 0.98508
# }

```

CÓDIGO: EXERCÍCIO 5

```
import numpy as np
import pandas as pd

from sklearn.linear_model import LogisticRegression
from sklearn.datasets import make_classification

from sklearn.model_selection import \
    cross_val_score, \
    train_test_split, \
    KFold, \
    LeaveOneOut

from sklearn.metrics import \
    matthews_corrcoef, \
    f1_score, \
    make_scorer

X, y = make_classification(
    n_samples=10000,
    n_features=2,
    n_informative=2,
    n_redundant=0,
    n_clusters_per_class=1,
    random_state=2002,
)

X_train, X_test, \
y_train, y_test = train_test_split(X, y,
                                    train_size=0.7,
                                    shuffle=True,
                                    random_state=2002)

model = LogisticRegression(penalty='l2')
model.fit(X_train, y_train)
```

```

mcc = make_scorer(matthews_corrcoef)

r = {}

# Houldout simples 50%
_, HS_X, _, HS_y = train_test_split(X_test, y_test,
                                     test_size=0.5,
                                     random_state=2002)

r['Houldout Simples'] = matthews_corrcoef(HS_y,
                                           model.predict(HS_X))

# Houldout repetido
HR = []
for _ in range(20):
    _, HS_X, _, HS_y = train_test_split(X_test, y_test,
                                         test_size=0.5,
                                         random_state=2002)

    HR.append(matthews_corrcoef(HS_y, model.predict(HS_X)))

r['Houldout Repetido (20 vezes)'] = np.mean(HR)

# KFold
cv = KFold(n_splits=10, shuffle=True, random_state=2002)
r['KFold'] = np.mean(cross_val_score(model,
                                     X_test, y_test,
                                     cv=cv, scoring=mcc))

# KFold Repetido
KR = []
for _ in range(20):
    cv = KFold(n_splits=10, shuffle=True, random_state=2002)
    HR.append(np.mean(cross_val_score(model,
                                     X_test, y_test,
                                     cv=cv, scoring=mcc)))

r['KFold Repetido (20 vezes)'] = np.mean(HR)

# LOOCV
loo = LeaveOneOut()
total_error = 0
for train_index, test_index in loo.split(X_test):
    X_tst = X_test[test_index, ]
    y_tst = y_test[test_index]
    # Evaluate the model on the test data
    error = model.score(X_tst, y_tst)
    total_error += error

r['Leave One Out'] = np.mean(total_error / len(X_test))

df = pd.DataFrame.from_dict(r, orient='index')
df.columns = ['Reg. Log.']
df.head()

```

REFERENCES

- [1] Peter Bartlett, Vitaly Maierov, and Ron Meir. Almost linear vc dimension bounds for piecewise polynomial networks. *Advances in neural information processing systems*, 11, 1998. Apud [2].
- [2] Peter L Bartlett and Wolfgang Maass. Vapnik-chervonenkis dimension of neural nets. *The handbook of brain theory and neural networks*, pages 1188–1192, 2003.
- [3] Eric Baum and David Haussler. What size net gives valid generalization? *Advances in neural information processing systems*, 1, 1988. Apud [2].