
Atividade 10 – Identificação de Lixo Eletrônico (SPAM).

Condições e Datas

O projeto deve ser realizado **individualmente** utilizando Python. Lembramos que o Python é livre e pode ser instalado, por exemplo, usando o ambiente Conda disponível em <https://conda.io>. Ele também pode ser acessado online usando o Google Colab através do link <https://research.google.com/colaboratory/>.

O projeto deve ser entregue no prazo especificado no Google Classroom. O arquivo deve descrever de forma clara os procedimentos adotados e as conclusões. Em particular, responda a(s) pergunta(s) abaixo de forma clara, objetiva e com fundamentos matemáticos. Recomenda-se que os códigos sejam anexados, mas **não serão aceitos trabalhos contendo apenas os códigos!** Pode-se submeter o arquivo .ipynb do Google Colab com os comandos e comentários.

O termo *lixo eletrônico* ou *spam* refere-se, de um modo geral, a mensagens (e-mails) indesejadas que são enviadas constantemente nos meios eletrônicos sem o consentimento do destinatário. Nesse projeto aplicaremos os conceitos de quadrados mínimos para auxiliar na identificação automática de *spams*. Especificamente, usaremos uma técnica de aprendizado de máquina que tem recebido bastante destaque nos últimos anos chamada *extreme learning machine* (ELM). Em termos gerais, uma ELM é sintetizada e avaliada com base num conjunto de mensagens que já foram identificadas como spam ou não-spam pelo(s) usuário(s). O conjunto usado para sintetizar a ELM é chamado *conjunto de treinamento* enquanto que o conjunto usado para avaliar o desempenho do modelo é chamado *conjunto de teste*. É importante destacar que o conjunto de teste não pode ser usado em nenhum momento para sintetizar a ELM. O aluno interessado em aprendizado de máquinas e nos detalhes da ELM pode consultar [1, 2].

Obtenção e Pre-Processamento dos Dados

Primeiramente, devemos baixar a base de dados *spambase* do repositório OpenML usando o scikit-learn (sklearn). Além de baixar os dados, iremos dividi-los em conjuntos de treinamento e teste e também aplicaremos uma sequência de pre-processamento. Especificamente, os seguintes passos incluem as bibliotecas e funções necessárias:

```
1 import numpy as np
2 from sklearn.datasets import fetch_openml
3 from sklearn.model_selection import train_test_split
4 from sklearn.preprocessing import StandardScaler
5 from sklearn.preprocessing import LabelEncoder
```

Obtemos os dados, dividimos em conjuntos de treinamento e teste e efetuamos o pre-processamento usando os comandos:

```
1 X, y = fetch_openml(data_id = 44, return_X_y=True)
2
3 Xtr, Xte, ytr, yte = train_test_split(X,y,random_state=1234)
4
5 scaler = StandardScaler()
6 Xtr = scaler.fit_transform(Xtr)
7 Xte = scaler.transform(Xte)
8
9 encoder = LabelEncoder()
10 ytr = 2*encoder.fit_transform(ytr)-1
11 yte = 2*encoder.transform(yte)-1
```

Dessa forma, obtemos uma matriz $\mathbf{X}_{\text{tr}} \in \mathbb{R}^{K \times N}$ e um vetor $\mathbf{y}_{\text{tr}} \in \{-1, 1\}^K$ em que $K = 3450$ e $N = 57$. A linha $\mathbf{X}_{\text{tr}}[k, :]$, ou equivalentemente $\mathbf{X}_{\text{tr}}[k]$, contém informações coletadas para identificação da k -ésima mensagem (e-mail). A componente $\mathbf{y}_{\text{tr}}[k]$ contém o valor $+1$ se a k -ésima mensagem foi identificada como *spam* e -1 se foi identificada como não-*spam*. O objetivo do projeto é construir um modelo capaz de classificar uma mensagem como *spam* ou não-*spam* utilizando somente o conjunto de treinamento, ou seja, \mathbf{X}_{tr} e \mathbf{y}_{tr} .

Funcionamento e Treinamento de uma Extreme Learning Machine

Uma *extreme learning machine* (ELM) é uma rede neural artificial de múltiplas camadas [1]. Nesse projeto, vamos considerar uma rede neural muito utilizada na literatura conhecida por *perceptron de múltiplas camadas*. Resumidamente, vamos assumir que a rede neural define uma função $\varphi : \mathbb{R}^N \rightarrow \mathbb{R}$ através da equação

$$\varphi(\mathbf{x}) = \alpha_1 g_1(\mathbf{x}) + \alpha_2 g_2(\mathbf{x}) + \dots + \alpha_M g_M(\mathbf{x}) = \sum_{i=1}^M \alpha_i g_i(\mathbf{x}), \quad (1)$$

em que $\alpha_1, \alpha_2, \dots, \alpha_M$ são parâmetros e as funções g_1, g_2, \dots, g_M são dadas por

$$g_i(\mathbf{x}) = \tanh \left(\sum_{j=1}^N w_{ij} x_j + \theta_i \right), \quad (2)$$

em que $\mathbf{w}_i = [w_{i1}, \dots, w_{iN}] \in \mathbb{R}^N$ e $\theta_i \in \mathbb{R}$ para todo $i = 1, \dots, M$ ¹. Em termos matriciais, podemos descrever a função $\varphi : \mathbb{R}^N \rightarrow \mathbb{R}$ como segue:

$$\varphi(\mathbf{x}) = \tanh(\mathbf{x}\mathbf{W} + \boldsymbol{\theta}) \cdot \boldsymbol{\alpha}, \quad (3)$$

em que $\boldsymbol{\alpha} = [\alpha_1, \dots, \alpha_M] \in \mathbb{R}^M$, $\mathbf{W} \in \mathbb{R}^{N \times M}$ é a matriz cujas linhas correspondem aos vetores \mathbf{w}_i e $\boldsymbol{\theta} = [\theta_1, \dots, \theta_M] \in \mathbb{R}^M$. O código que implementa a função φ descrita pela rede neural artificial em python é:

```
1 def EvalELM(alpha, W, theta, X):
2     return np.dot(np.tanh(X@W+theta), alpha)
```

Observação. Note que o comando

```
1 G = tanh(Xtr@W+theta)
```

fornece uma matriz $G \in \mathbb{R}^{K \times M}$ cujo elemento $G[k, i]$ corresponde à avaliação da i -ésima função base g_i calculada nos dados no vetor de características da k -ésima mensagem, isto é, $G[k, i] = g_i(\mathbf{X}_{\text{tr}}[k])$. Além disso, o produto $\mathbf{s} = G\boldsymbol{\alpha}$ fornece um vetor $\mathbf{s} \in \mathbb{R}^K$ contendo o valor de φ calculado em cada mensagem, ou seja, $\mathbf{s} = [s_1, \dots, s_K]$ em que $s_k = \varphi(\mathbf{X}_{\text{tr}}[k])$ para todo $k = 1, \dots, K$.

Numa ELM, os vetores $\mathbf{w}_i = [w_{i1}, \dots, w_{iN}] \in \mathbb{R}^N$ e o escalar $\theta_i \in \mathbb{R}$ que definem a função g_i são gerados aleatoriamente utilizando uma distribuição normal padrão, para todo $i = 1, \dots, M$. Usando a forma matricial explícita em (3), no python, utilizamos os seguintes comandos para construir \mathbf{W} e $\boldsymbol{\theta}$:

```
1 W = np.random.randn(N, M)
2 theta = np.random.randn(M)
```

Os parâmetros $\alpha_1, \dots, \alpha_M$ são determinados resolvendo o problema de quadrados mínimos

$$\varphi(\mathbf{X}_{\text{tr}}[k]) \approx \mathbf{y}_{\text{tr}}[k], \quad \forall k = 1, \dots, K, \quad (4)$$

¹Observe que K refere-se ao número de dados de treinamento enquanto que M corresponde ao número de parâmetros. Nesse projeto, temos $M = 1000$ e $K = 3450$ no conjunto de treinamento.

definido sobre o conjunto de treinamento. Em outras palavras, $\alpha_1, \dots, \alpha_M$ minimizam a soma dos quadrados dos desvios

$$J(\alpha_1, \dots, \alpha_n) = \sum_{k=1}^K (\alpha_1 g_1(\mathbf{x}_{tr}[k]) + \dots + \alpha_M g_M(\mathbf{x}_{tr}[k]) - y_{tr}[k])^2. \quad (5)$$

Finalmente, se $\mathbf{x} \in \mathbb{R}^N$ é o vetor contendo informações sobre uma mensagem, a identificação é efetuada como segue

$$\begin{cases} \text{A mensagem é um spam se } 0 \leq \varphi(\mathbf{x}), \\ \text{A mensagem não é spam caso contrário.} \end{cases} \quad (6)$$

Classificação de Novas Mensagens e Medidas de Desempenho

Conhecidos a função φ e o limiar L , podemos avaliar o desempenho do sistema usando um conjunto de dados que já foram identificados pelo(s) usuário(s). Por exemplo, podemos avaliar o desempenho do sistema no conjunto de teste, em que $\mathbf{x}_{te} \in \mathbb{R}^{1151 \times 57}$ e um vetor $\mathbf{y}_{te} \in \{-1, 1\}^{1151}$. Tal como nos dados de treinamento, $\mathbf{x}_{te}[k]$ e $\mathbf{y}_{te}[k]$ contém informações sobre o conteúdo da k -ésima mensagem (e-mail) e sua classificação como *spam* ou não-*spam*, respectivamente. O desempenho do sistema pode ser medido quantitativamente, por exemplo, calculando a *acurácia* (AC) ou a *taxa de falsos positivos* (TFP, também chamado “taxa de alerta falso”) definidos respectivamente pelas equações:

$$AC = \frac{\text{Número de mensagens identificadas corretamente pelo sistema}}{\text{Número total de mensagens}}, \quad (7)$$

e

$$TFP = \frac{\text{Número de mensagens identificadas como spam pelo sistema mas que não são spams}}{\text{Número de mensagens que não são spams}}. \quad (8)$$

Questão 1

Sintetize a aplicação φ resolvendo o problema de quadrados mínimos em (4) com respeito ao conjunto de treinamento considerando $M = 1000$.

Questão 2

Ainda usando o conjunto de treinamento, isto é, \mathbf{x}_{tr} e \mathbf{y}_{tr} , determine a acurácia e a taxa de falsos positivos.

Questão 3

Usando o conjunto de teste, isto é, \mathbf{x}_{te} e \mathbf{y}_{te} , calcule a acurácia e a taxa de falsos positivos. O desempenho no conjunto de teste é consistente com o esperado, isto é, a acurácia e a taxa de falsos positivos são semelhantes aos valores obtidos considerando o conjunto de treinamento?

Referências

- [1] HAYKIN, S. *Neural Networks and Learning Machines*, 3rd edition ed. Prentice-Hall, Upper Saddle River, NJ, 2009.
- [2] HUANG, G.-B., WANG, D., AND LAN, Y. Extreme learning machines: a survey. *Int. J. Machine Learning & Cybernetics* 2, 2 (2011), 107–122.