

ACH2118 Introdução ao Processamento de Língua Natural

Semestre 2023-2 - Exercício prático 1: classificação de clareza das respostas na plataforma eSIC

1. O objetivo do trabalho é implementar um classificador de nível de clareza (c1,c234,c5) a partir de textos de respostas publicadas na plataforma eSIC. com base no conjunto de dados disponível na atividade.
2. O EP consiste de uma tarefa de classificação ternária indicando o grau de clareza do texto fornecido.
3. O trabalho é desenvolvido em duas partes: a primeira entrega consiste em desenvolver o classificador e reportar resultados médios de medida de **acurácia** usando validação cruzada de 10 partições sobre o conjunto de dados de treinamento completo, a ser avaliado de acordo com os critérios discutidos a seguir. Nesta primeira parte, programadores Python devem usar

```
result = cross_val_score(clf, x_train, y_train, scoring='accuracy', cv=10).mean()
```

4. Os dados de treinamento são fornecidos em um arquivo ep1_esic2023_clareza_TRAIN.xlsx que pode ser aberto em Excel, Google Sheets ou similar, ou diretamente em Python. Não é fornecido vocabulário adicional, e portanto os futuros dados de teste incluirão palavras desconhecidas em tempo de treinamento.
5. A segunda parte da avaliação consiste em gerar as predições e cada classificador para o conjunto de teste que será fornecido após a primeira entrega.
6. O que entregar:
 - **Entrega 1** (treinamento): uma pasta ZIP contendo (a) o código completo da implementação; e (b) o relatório descritivo dos classificadores considerados, deixando explícito qual deles é o classificador final que está sendo proposto, juntamente com os resultados de validação cruzada *apenas do classificador final* sobre os dados de treinamento. Ou seja, apresente apenas UM classificador final e apenas UM valor final de acurácia, e não valores para múltiplos classificadores ou *folds*.
 - **Entrega 2** (teste): uma pasta ZIP contendo (a) relatório atualizado com os resultados de teste, (b) slides para apresentação (15 min) e (c) arquivo XLSX com as predições sobre os dados de teste, conforme modelo a ser divulgado, contendo uma coluna de rótulos exatamente na mesma ordem dos dados do conjunto de teste, sem nenhuma linha a mais ou a menos (pois neste caso não haveria como avaliar).

A entrega deve ser feita ANTES do prazo estipulado, certificando-se de que o arquivo realmente foi carregado com sucesso, que é a versão correta, e que não está corrompido.

7. Avaliação

- **Critérios de avaliação:** acurácia de teste (70%), método inovador / bem elaborado (20%), relatório (5%), apresentação em aula (5%).
- **Penalidade por baixo desempenho:** EPs cuja acurácia média seja inferior à média de um baseline de classe majoritária recebem nota 3,0 e não são avaliados. Assim, por exemplo, se a classe majoritária obtém 60% de acurácia no teste, a acurácia do modelo não pode ser inferior a 60%, já que isso tornaria o aprendizado de máquina sem razão de ser.
- **Penalidade por formato inválido:** 5% de desconto na nota final por entrega de arquivo de predições em formato diferente do esperado.
- **Ranking da turma.** As notas serão ajustadas pelo melhor resultado (que receberá nota dez). É desejável portanto que você obtenha não apenas um bom resultado, mas melhor do que o de seus colegas. 😊

O EP pode ser desenvolvido individualmente ou em duplas, desde que estas sejam cadastradas no link a seguir antes do prazo a ser definido em aula.

https://docs.google.com/spreadsheets/d/1akG-zMcJEE-RN8oHs69pqkeHgOU_gR2Yrlnifur3Bw4/edit?usp=sharing