

Equipe 16

- Ana Beatriz Kindinger
- Daniel Victor Andrade
- Igor Buess Atala Y Mansour
- Marlon Mateus Prudente de Oliveira
- Ronaldo Santana da Silva Moco

Identificador automático de idioma

Problema: Dados um texto de entrada, é possível identificar em qual língua o texto está escrito?

Entrada: "texto qualquer"

Saída: português ou inglês ou francês ou italiano ou...

O processo de Reconhecimento de Padrões

O objetivo desse trabalho é demonstrar o processo de "construção de atributos" e como ele é fundamental para o **Reconhecimento de Padrões (RP)**.

Primeiro um conjunto de "amostras" previamente conhecido (classificado)

```
1  #
2  # amostras de texto em diferentes línguas
3  #
4  ingles = [
5  "Hello, how are you?",
6  "I love to read books.",
7  "The weather is nice today.",
8  "Where is the nearest restaurant?",
9  "What time is it?",
10 "I enjoy playing soccer.",
11 "Can you help me with this?",
12 "I'm going to the movies tonight.",
13 "This is a beautiful place.",
14 "I like listening to music.",
15 "Do you speak English?",
16 "What is your favorite color?",
17 "I'm learning to play the guitar.",
18 "Have a great day!",
19 "I need to buy some groceries.",
20 "Let's go for a walk.",
21 "How was your weekend?",
22 "I'm excited for the concert.",
23 "Could you pass me the salt, please?",
24 "I have a meeting at 2 PM.",
25 "I'm planning a vacation.",
26 "She sings beautifully.",
27 "The cat is sleeping.",
28 "I want to learn French.",
29 "I enjoy going to the beach.",
30 "Where can I find a taxi?",
31 "I'm sorry for the inconvenience.",
32 "I'm studying for my exams.",
33 "I like to cook dinner at home.",
34 "Do you have any recommendations for restaurants?",
35 ]
36
37 espanhol = [
38 "Hola, ¿cómo estás?",
39 "Me encanta leer libros.",
40 "El clima está agradable hoy.",
41 "¿Dónde está el restaurante más cercano?",
42 "¿Qué hora es?",
43 "Voy al parque todos los días.",
44 "¿Puedes ayudarme con esto?",
45 "Me gustaría ir de vacaciones.",
46 "Este es mi libro favorito.",
47 "Me gusta bailar salsa.",
48 "¿Hablas español?",
49 "¿Cuál es tu comida favorita?",
50 "Estoy aprendiendo a tocar el piano.",
51 "¡Que tengas un buen día!",
52 "Necesito comprar algunas frutas.",
53 "Vamos a dar un paseo.",
54 "¿Cómo estuvo tu fin de semana?",
55 "Estoy emocionado por el concierto.",
56 "¿Me pasas la sal, por favor?",
57 "Tengo una reunión a las 2 PM.",
58 "Estoy planeando unas vacaciones.",
59 "Ella canta hermosamente.",
60 "El perro está jugando."
```



```

58 pu_espanhol = re.findall(padrao_espanhol, texto)
59
60 return len(pd_espanhol)
61
62 def caracteristicas_portugues(texto):
63     padrao_portugues = r'\w(?:nh|lh|rr|ss|am)'
64     pd_portugues = re.findall(padrao_portugues, texto)
65
66     return len(pd_portugues)
67
68 def extraiCaracteristicas(frase):
69     # frase é um vetor [ 'texto', 'lingua' ]
70     texto = frase[0]
71     pattern_regex = re.compile('[^\\w+]', re.UNICODE)
72     texto = re.sub(pattern_regex, ' ', texto)
73     #print(texto)
74     caracteristica1=tamanhoMedioFrases(texto)
75     caracteristica2=contar_consoantes(texto)
76     caracteristica3=contar_acentuados(texto)
77     caracteristica4=caracteristicas_ingles(texto)
78     caracteristica5=caracteristicas_espanhol(texto)
79     caracteristica6=caracteristicas_portugues(texto)
80     # acrescente as suas funcoes no vetor padrao
81     padrao = [caracteristica1, caracteristica2, caracteristica3, caracteristica4,
82               caracteristica5,caracteristica6, frase[1] ]
83     return padrao
84
85 def geraPadroes(frases):
86     padroes = []
87     for frase in frases:
88         padrao = extraiCaracteristicas(frase)
89         padroes.append(padrao)
90     return padroes
91
92 # converte o formato [frase classe] em
93 # [caracteristica_1, caracteristica_2,... caracteristica n, classe]
94 padroes = geraPadroes(pre_padroes)
95
96 #
97 # apenas para visualizacao
98 print(padroes)
99
100 dados = pd.DataFrame(padroes)
101 dados

```

[illegible]

Com os conjuntos separados, podemos "treinar" o modelo usando a SVM.

```
1 from sklearn import svm
2 from sklearn.metrics import confusion_matrix
3 from sklearn.metrics import classification_report
4
5 treinador = svm.SVC() #algoritmo escolhido
6 modelo = treinador.fit(X_train, y_train)
7
8 #
9 # score com os dados de treinamento
10 acuracia = modelo.score(X_train, y_train)
11 print("Acurácia nos dados de treinamento: {:.2f}%".format(acuracia * 100))
12
13 #
14 # melhor avaliar com a matriz de confusão
15 y_pred = modelo.predict(X_train)
16 cm = confusion_matrix(y_train, y_pred)
17 print(cm)
18 print(classification_report(y_train, y_pred))
19
20 #
21 # com dados de teste que não foram usados no treinamento
22 print('métricas mais confiáveis')
23 y_pred2 = modelo.predict(X_test)
24 cm = confusion_matrix(y_test, y_pred2)
25 print(cm)
26 print(classification_report(y_test, y_pred2))
27
```

→ Acurácia nos dados de treinamento: 71.01%

```
[[19  2  2]
 [ 5 14  3]
 [ 8  0 16]]
```

	precision	recall	f1-score	support
espanhol	0.59	0.83	0.69	23
inglês	0.88	0.64	0.74	22
português	0.76	0.67	0.71	24
accuracy			0.71	69
macro avg	0.74	0.71	0.71	69
weighted avg	0.74	0.71	0.71	69

métricas mais confiáveis

```
[[4 1 2]
 [3 4 1]
 [3 0 5]]
```

	precision	recall	f1-score	support
espanhol	0.40	0.57	0.47	7
inglês	0.80	0.50	0.62	8
português	0.62	0.62	0.62	8
accuracy			0.57	23
macro avg	0.61	0.57	0.57	23
weighted avg	0.62	0.57	0.57	23