

UNIVERSITÀ DEGLI STUDI DI NAPOLI

FEDERICO II



Confronto tra SVM e KNN
applicando i test di significatività statistica
discussi in
*“The Hitchhiker’s Guide to Testing Statistical Significance in
Natural Language Processing”*

Prof.ssa Anna Corazza

De Lucia Gianluca (N97000311)

A.A. 2019/2020

Sommario

1. Introduzione.....3

2. Esperimento5

3. Riferimenti.....7

4. Appendice.....8

1. Introduzione

Con questo esperimento vogliamo mettere a confronto due algoritmi di classificazione quali Knn e SVM.

Svm (Space Vector Machine): divide uno spazio vettoriale con un iperpiano, definendo lo spazio delle classi.

Knn (K nearest neighbor): tramite k vicini e vedendo la maggioranza di essi a quale classe appartengono, riesce a classificare un nuovo documento.

Per valutare e confrontare due algoritmi come questi utilizziamo il test di significatività statistica per la NPL [1].

Ma come si può conoscere la distribuzione statistica dei test? Una possibilità è quella di applicare prove progettate per valutare la distribuzione di un campione di osservazioni.

Esistono vari test, che si dividono in parametrici e non parametrici.

PARAMETRICI:

- Ad esempio, il test di Shapiro-Wilk (Shapiro e Wilk, 1965) verifica l'ipotesi nulla che un campione provenga da una popolazione normalmente distribuita;
- il test di Kolmogorov-Smirnov quantifica la distanza tra la funzione di distribuzione empirica del campione e la funzione di distribuzione cumulativa della distribuzione di riferimento;
- il test Anderson-Darling (Anderson e Darling, 1954) verifica se un dato campione di dati viene estratto da una data distribuzione di probabilità;
- il t-test Questo test valuta se la popolazione media di due serie di misurazioni differisce l'una dall'altra e si basa sul presupposto che entrambi i campioni provengano da una distribuzione normale (Fisher, 1937). calcola il numero medio di previsioni corrette per esempio di input. i punteggi sono significativamente più grandi di zero, il che può servire da indicazione che un parser è migliore dell'altro.

NON PARAMETRICI:

Quando la distribuzione statistica del test è sconosciuta, è necessario utilizzare test di significatività non parametrici.

Ci sono due famiglie dei non parametrici:

- La prima famiglia è composta da test che non tengono conto dei valori effettivi delle misure di valutazione.
- La seconda famiglia considera i valori delle misure: verifica ripetutamente il campione dai dati del test e stima il valore p in base ai valori statistici del test nei campioni.

Ci riferiamo alla prima famiglia come famiglia di test senza campionamento e alla seconda come famiglia di test basati su campionamento. Sono i seguenti test:

- McNemar: Questo test è progettato per osservazioni nominali accoppiate (etichette binarie). Il test viene applicato a una tabella di contingenza 2×2 , che tabula i risultati di due algoritmi su un campione di n esempi. L'ipotesi nulla per questo test afferma che la probabilità marginale per ciascun risultato (etichetta uno o etichetta due) è la stessa per entrambi gli algoritmi. Cioè, quando applichiamo entrambi gli algoritmi sugli stessi dati, ci aspetteremmo che siano corretti / errati sulla stessa proporzione di elementi;
- Wilcoxon: Questo test viene utilizzato quando si confrontano due campioni corrispondenti. La sua ipotesi nulla è che le differenze seguano una distribuzione simmetrica attorno allo zero;
- Permutation: Questo test stima la distribuzione delle statistiche di test sotto l'ipotesi nulla calcolando i valori di questa statistica sotto tutte le possibili etichettature (permutazioni) del set di test. Il valore p (bilaterale) del test viene calcolato come la proporzione di queste permutazioni in cui la differenza assoluta era maggiore o uguale al valore assoluto della differenza nell'output dell'algoritmo;
- Bootstrap: Questo test è molto simile alla randomizzazione approssimativa del test di permutazione, con la differenza che il campionamento viene effettuato con sostituzioni (ad esempio, un esempio dai dati del test originale può apparire più di una volta in un campione). L'idea di bootstrap è di usare i campioni come popolazioni surrogate, allo scopo di approssimare la distribuzione campionaria della statistica. Il valore p viene calcolato in modo simile al test di permutazione. Il test è meno efficace per piccoli set di test, poiché presuppone che la distribuzione del set di test non si discosti troppo dalla distribuzione della popolazione.

Quando il test risulta significativo possiamo affermare con certezza la superiorità di un algoritmo rispetto all'altro.

2. Esperimento

Ho raccolto in due file i valori di knn di svm applicati a 4 dataset differenti [2].

Secondo i test svolti:

DataSet	Accuratezza KNN (%)	Accuratezza SVM (%)
Data 1	92.5925	98.1481
Data 2	87.0370	90.7407
Data 3	81.8182	84.0909

Notiamo subito che SVM risulta essere più accurato di KNN.

Ho applicato il test di significatività su questi valori di entrambi gli algoritmi utilizzando uno script python allegato all'articolo scientifico in considerazione [2].

Per valutare la prestazione di due algoritmi, A e B, su un set di dati X, si confrontano usando una misura di valutazione M. Indichiamo $M(\text{ALG}, X)$ come valore della misura di valutazione M quando l'algoritmo ALG viene applicato al set di dati X.

Senza perdita di generalità, assumiamo che valori più alti della misura siano migliori.

Definiamo la differenza nelle prestazioni tra i due algoritmi secondo la misura M sul set di dati X come:

$$\delta(X) = M(A, X) - M(B, X).$$

Consideriamo $\delta(X)$ come la nostra statistica di prova.

Usando questa notazione formuliamo il seguente problema di verifica dell'ipotesi statistica:

$$H_0: \delta(X) \leq 0$$

$$H_1: \delta(X) > 0$$

Per decidere se rifiutare o meno l'ipotesi nulla, che sta raggiungendo la conclusione che $\delta(X)$ è effettivamente maggiore di 0, di solito calcoliamo un *p-value* per il test. Il valore p è definito come la probabilità, sotto l'ipotesi nulla H_0 , di ottenere un risultato uguale o più estremo di quanto effettivamente osservato.

Indicando con α il valore di soglia da utilizzare per considerare la significatività e quindi il livello di significatività e con *p-value* il valore:

$$P r(\delta(X) \geq \delta_{\text{observed}} | H_0)$$

Quando il valore $p < \alpha$ allora l'ipotesi nulla viene respinta.

Applicazione del Test di significatività statistica con $\alpha = 0.5$:

Tipo di Test	Significativo (Si/No)	P-Value
T-test	SI	0.028013616053
Wilcoxon	SI	0.108809430041
Permutation	SI	0.124487551245
Bootstrap	SI	0.0

Applicazione del Test di significatività statistica con $\alpha = 0.8$:

Tipo di Test	Significativo (Si/No)	P-Value
T-test	SI	0.028013616053
Wilcoxon	SI	0.108809430041
Permutation	SI	0.122587741226
Bootstrap	SI	0.0

È cambiato solo il p-value del test di Permutation.

McNemar non passa il test poiché accetta solo etichettature binarie.

Bootstrap per funzionare al meglio dovrebbe essere eseguiti su molti più valori per poter essere accertato.

Per quanto riguarda i test di Shapiro-Wilk, di Kolmogorov-Smirnov e di Anderson-Darling, non avendo una distribuzione normale dei dati, il test risulta essere non significativo in tutti i casi.

A conclusione di tutto possiamo affermare che l'algoritmo SVM è migliore rispetto a KNN in maniera significativa in accordo con il t-test, Wilcoxon test, Permutation test e Bootstrap test.

3. Riferimenti

[1]: “The Hitchhiker’s Guide to Testing Statistical Significance in Natural Language Processing” Rotem Dror Gili Baumer Segev Shlomov Roi Reichart Faculty of Industrial Engineering and Management, Technion, IIT

[2]: <http://archive.ics.uci.edu/ml/datasets.php>

[3]: <https://github.com/rtmdrr/testSignificanceNLP.git>

[4]: <https://gitlab.com/gigernau/cvs-classification.git>

4. Appendice

Il codice utilizzato per valutare l'accuratezza dei due algoritmi è il seguente (python 3):

```
1 # Importing the libraries
2 import numpy as np
3 import matplotlib.pyplot as plt
4 import pandas as pd
5 from sklearn.metrics import accuracy_score
6
7 # Importing the dataset
8 dataset = pd.read_csv('wine_data.csv')
9 X = dataset.iloc[:, 1:13].values
10 y = dataset.iloc[:, 0].values
11
12 # Splitting the dataset into the Training set and Test set
13 from sklearn.model_selection import train_test_split
14 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.30, random_state = 0)
15
16 # Feature Scaling
17 from sklearn.preprocessing import StandardScaler
18 sc = StandardScaler()
19 X_train = sc.fit_transform(X_train)
20 X_test = sc.transform(X_test)
21
22 """#linear discriminant analysis
23 from sklearn.discriminant_analysis import LinearDiscriminantAnalysis
24 lda = LinearDiscriminantAnalysis(n_components=5)
25 X_train = lda.fit_transform(X_train,y_train)
26 X_test = lda.fit_transform(X_test,y_test)
27 #explained_variance=pca.explained_variance_ratio_
28 or
29
30 from sklearn.decomposition import PCA
31 pca = PCA(n_components=2)
32 X_train = pca.fit_transform(X_train)
33 X_test = pca.fit_transform(X_test)
34 #explained_variance=pca.explained_variance_ratio_"""
35
36
37 # Fitting KNN to the Training set
38
39 from sklearn.neighbors import KNeighborsClassifier
40
41
42 classifier = KNeighborsClassifier(n_neighbors=5)
43 trained_model=classifier.fit(X_train,y_train)
44 trained_model.fit(X_train,y_train )
45
46 # Predicting the Test set results
47
48 y_pred = classifier.predict(X_test)
49
```



```

50 # Making the Confusion Matrix
51
52 from sklearn.metrics import confusion_matrix
53
54 cm_KNN = confusion_matrix(y_test, y_pred)
55 print(cm_KNN)
56 print("Accuracy score of train KNN")
57 print(accuracy_score(y_train, trained_model.predict(X_train))*100)
58
59 print("Accuracy score of test KNN")
60 print(accuracy_score(y_test, y_pred)*100)
61 TotalAmount = accuracy_score(y_test, y_pred)*100
62
63
64
65
66 with open("resA.txt", "a+") as text_file:
67     print(f"{TotalAmount}", file=text_file)
68
69
70
71 # # Fitting SVM to the Training set
72 from sklearn.svm import SVC
73
74 classifier = SVC(kernel = 'linear', random_state = 0)
75 trained_model=classifier.fit(X_train,y_train)
76 trained_model.fit(X_train,y_train )
77
78
79 # Predicting the Test set results
80 y_pred = classifier.predict(X_test)
81
82 # Making the Confusion Matrix
83 from sklearn.metrics import confusion_matrix
84 cm_SVM = confusion_matrix(y_test, y_pred)
85 print(cm_SVM)
86 print("Accuracy score of train SVM")
87 print(accuracy_score(y_train, trained_model.predict(X_train))*100)
88
89 print("Accuracy score of test SVM")
90 print(accuracy_score(y_test, y_pred)*100)
91 TotalAmount2 = accuracy_score(y_test, y_pred)*100
92
93 with open("resB.txt", "a+") as text_file2:
94     print(f"{TotalAmount2}", file=text_file2)
95

```