



ANÀLISI DE SENTIMENT DE PEL·LÍCULES

CONSTRUCCIÓ D'UN CLASSIFICADOR DE SENTIMENTS EN RESSENYES DE PEL·LÍCULES AMB SCIKIT-LEARN I PYTHON



TÍTOL
OFICIAL
DE PRECISIÓ
SOVIÈTICA

Data i hora
12 de desembre
MMXXIII
12:34:22

Agenda

INTRODUCCIÓ

DATASET

PREPROCESSAMENT

ANÀLISI

CONCLUSIONS

INTRO



INTRO



dataset molt polit a kaggle

INTRO



scikit-learn és guai

dataset molt polit a kaggle

INTRO

La Construcció d'un Classificador de Sentiments en Ressenyes de Pel·lícules.

Un model de **machine learning** amb la llibreria **scikit-learn** a Python per predir si una ressenya de pel·lícula és positiva o negativa a partir de la interpretació del text.

DATASET

<https://www.kaggle.com/code/lakshmi25npat/sentiment-analysis-of-imdb-movie-reviews>

Gràcies Kaggle!

S'ha fet servir un conjunt de dades de l'IMDB amb **50.000 ressenyes de pel·lícules** disponible a Kaggle.

Try Pitch

```
df = pd.read_csv('kaggle_movies_db.csv')
df
```

	review	sentiment
0	One of the other reviewers has mentioned that ...	positive
1	A wonderful little production. The...	positive
2	I thought this was a wonderful way to spend ti...	positive
3	Basically there's a family where a little boy ...	negative
4	Petter Mattei's "Love in the Time of Money" is...	positive
...
95	I thought this movie did a down right good job...	positive
96	Bad plot, bad dialogue, bad acting, idiotic di...	negative
97	I am a Catholic taught in parochial elementary...	negative
98	I'm going to have to disagree with the previou...	negative
99	No one expects the Star Trek movies to be high...	negative

100 rows × 2 columns

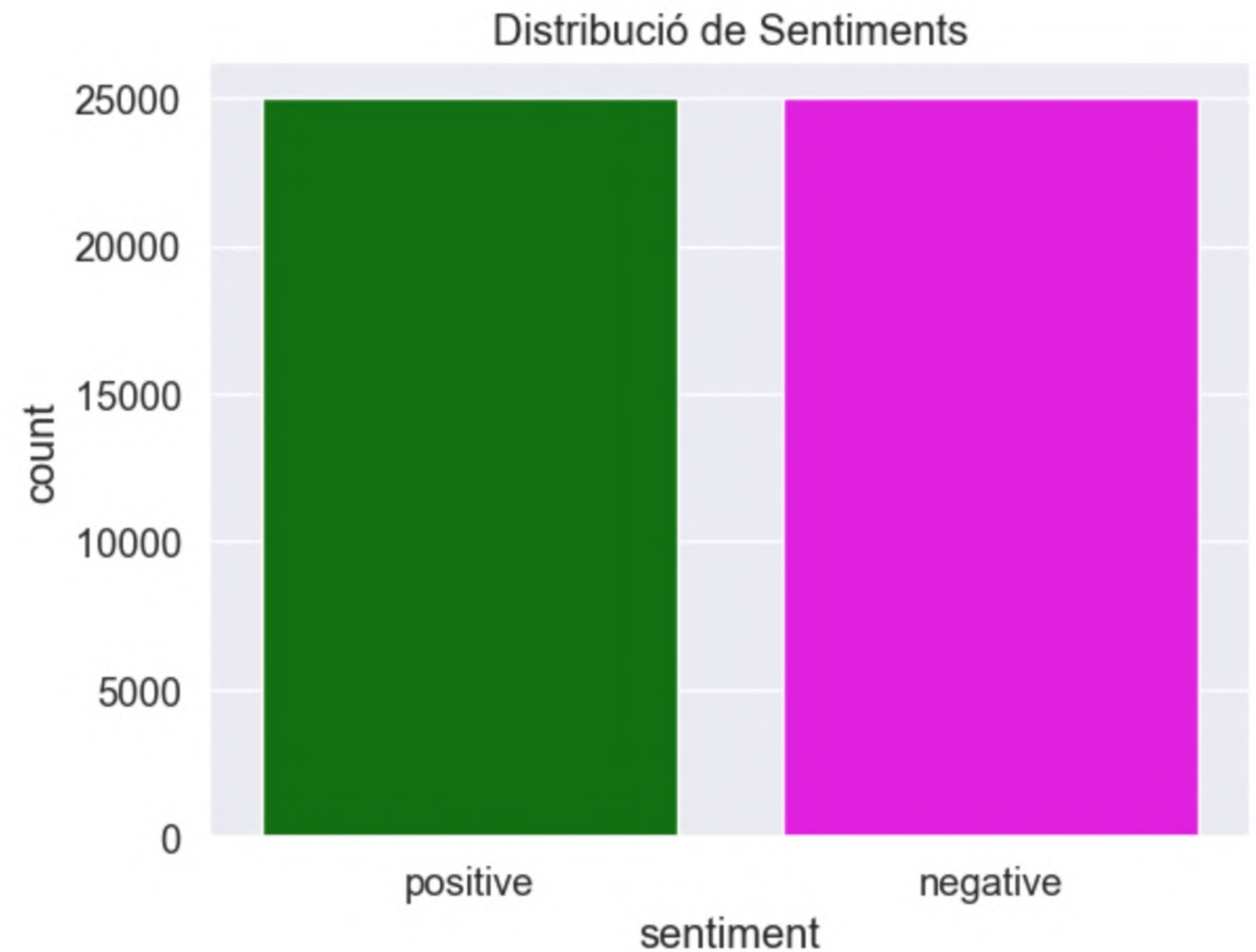
continuació, explorem les dades del dataset

```
df.describe()
```

	review	sentiment
count	50000	50000
unique	49582	2
top	Loved today's show!!! It was a variety and not...	positive
freq	5	25000

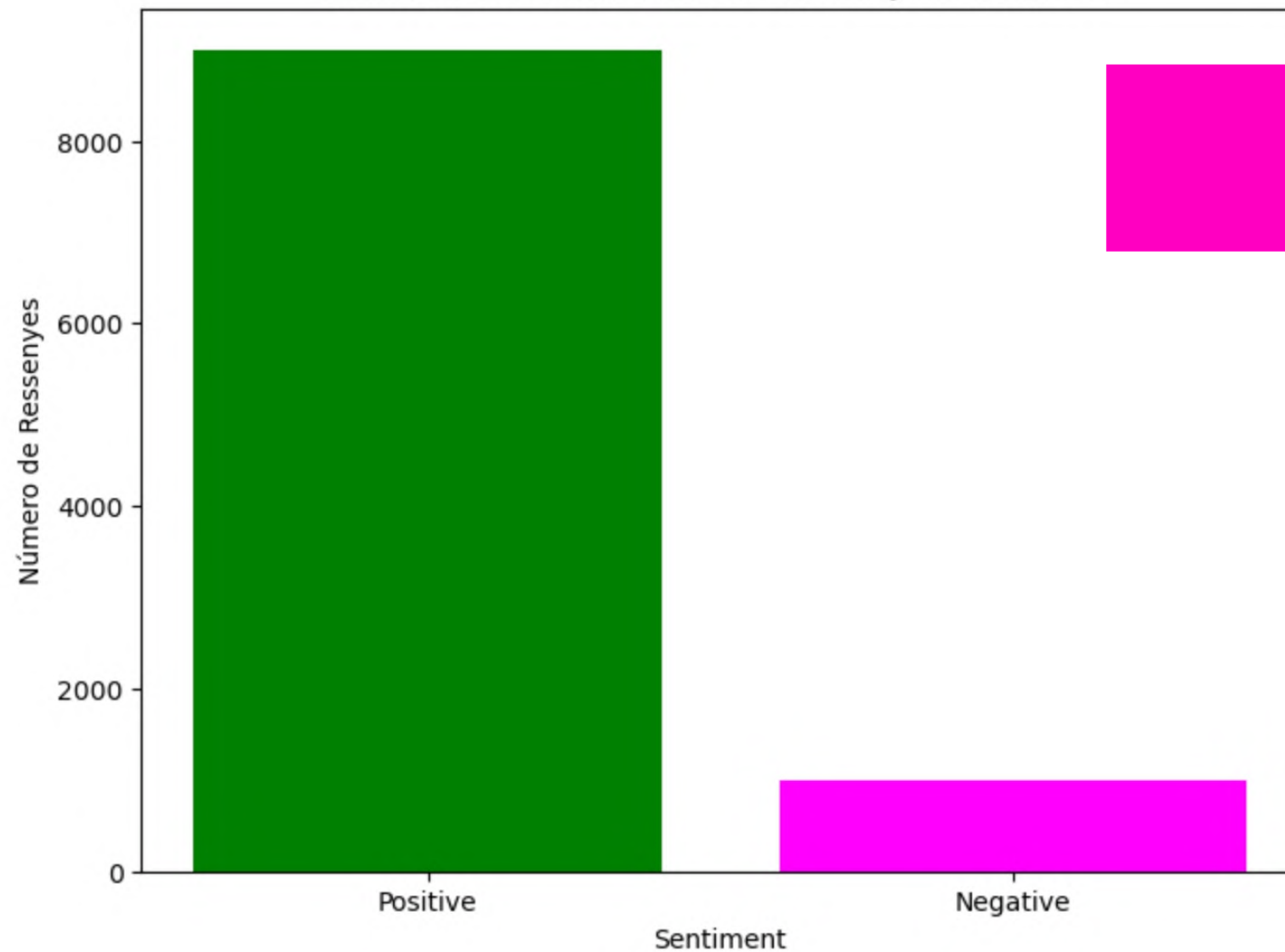
DATASET

<https://www.kaggle.com/code/lakshmi25npathi/sentiment-analysis-of-imdb-movie-reviews>



PREPROCESSAMENT

Distribució de Sentiments en el Conjunt de Dades



imbalanced-learn documentation

Date: Jul 08, 2023 Version: 0.11.0

Useful links: [Binary Installers](#) | [Source Repository](#) | [Issues & Ideas](#) | [Q&A Support](#)

Imbalanced-learn (imported as `imblearn`) is an open source, MIT-licensed library relying on scikit-learn (imported as `sklearn`) and provides tools when dealing with classification with imbalanced classes.

```
sentiment
negative      1000
positive      1000
dtype: int64
```


ANÀLISI: SCIKIT LEARN

```
from imblearn.under_sampling import RandomUnderSampler
from sklearn.model_selection import train_test_split

from sklearn.feature_extraction.text import CountVectorizer
from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.svm import SVC

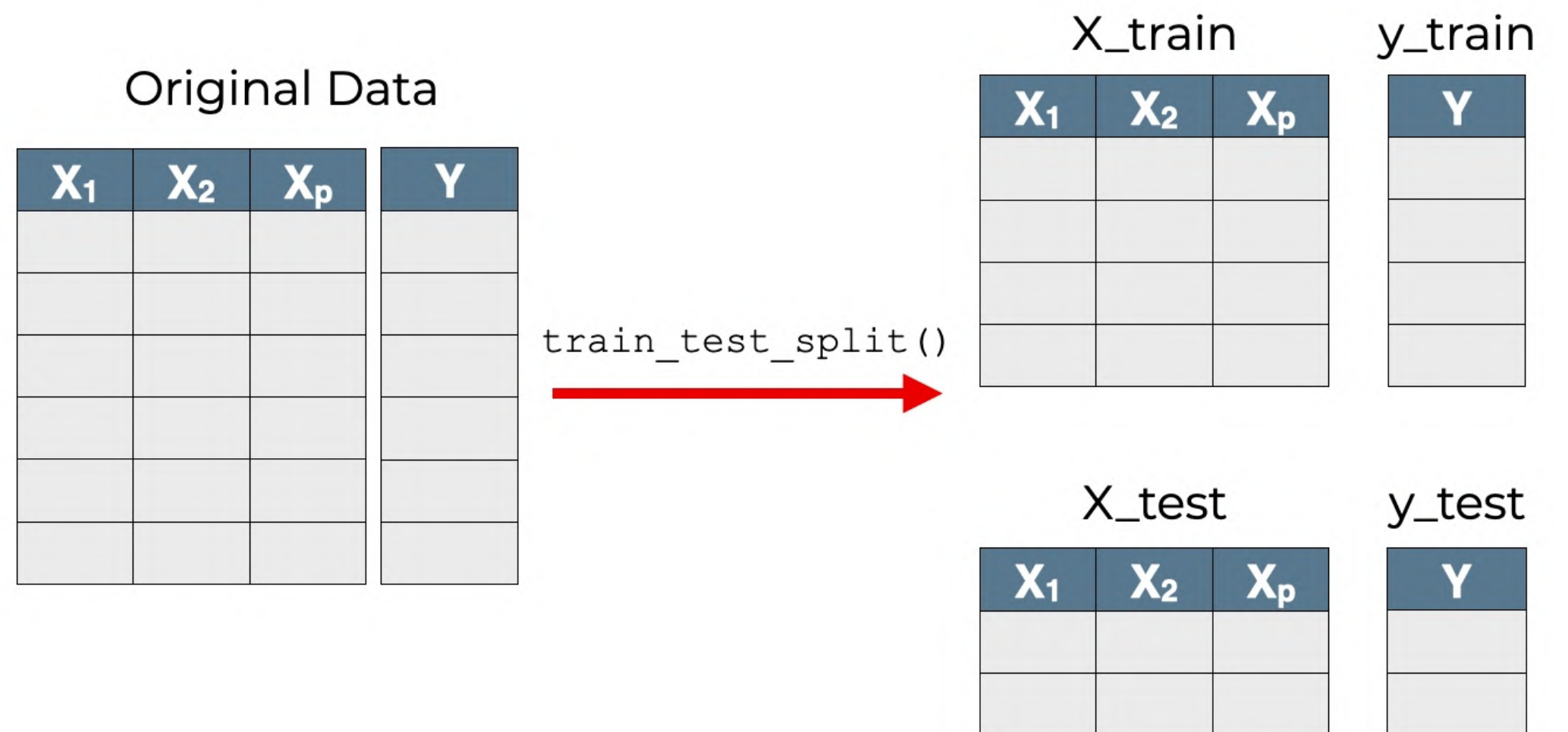
from sklearn.tree import DecisionTreeClassifier
from sklearn.naive_bayes import GaussianNB
from sklearn.linear_model import LogisticRegression

from sklearn.metrics import f1_score
from sklearn.metrics import classification_report
from sklearn.metrics import confusion_matrix

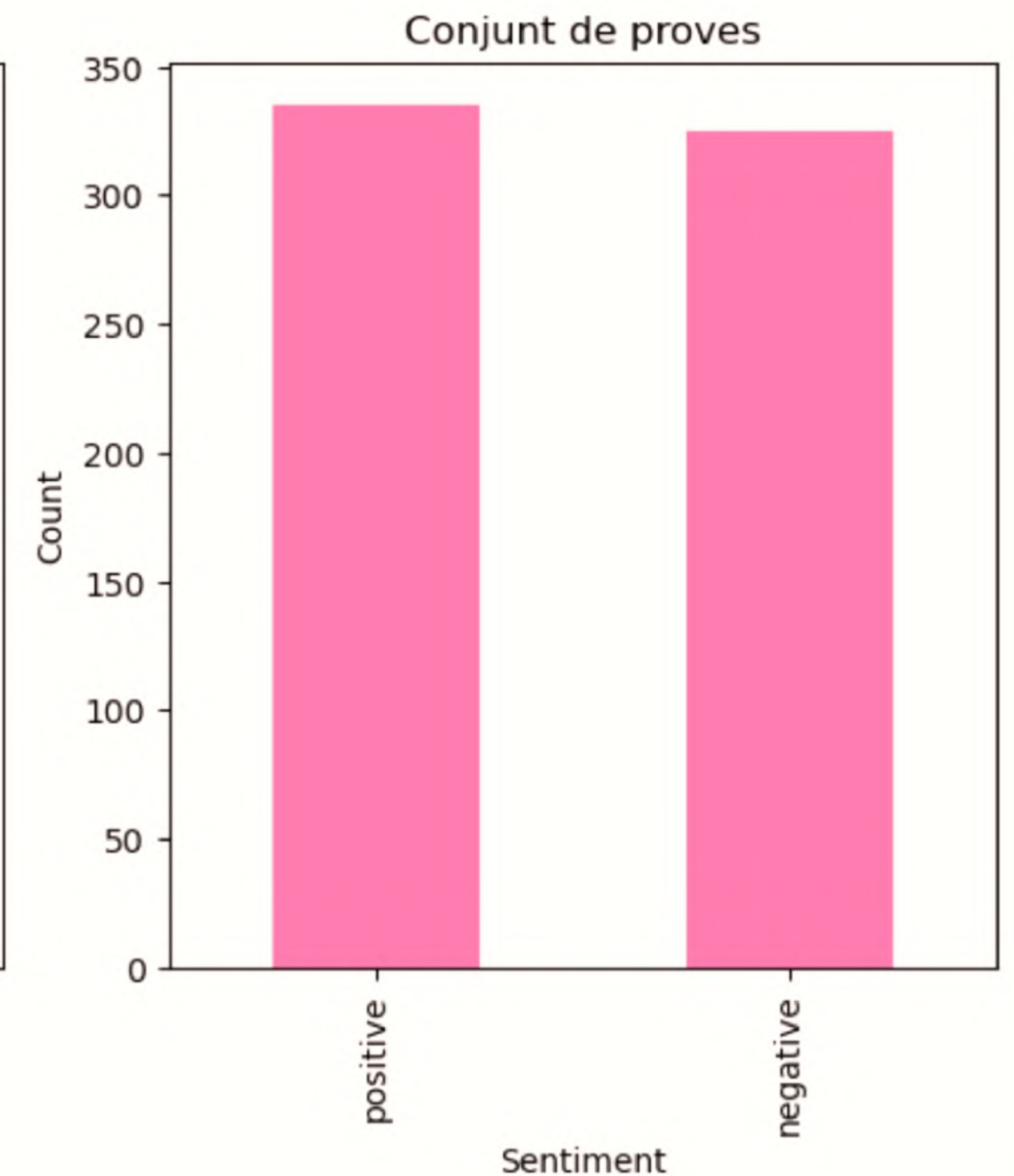
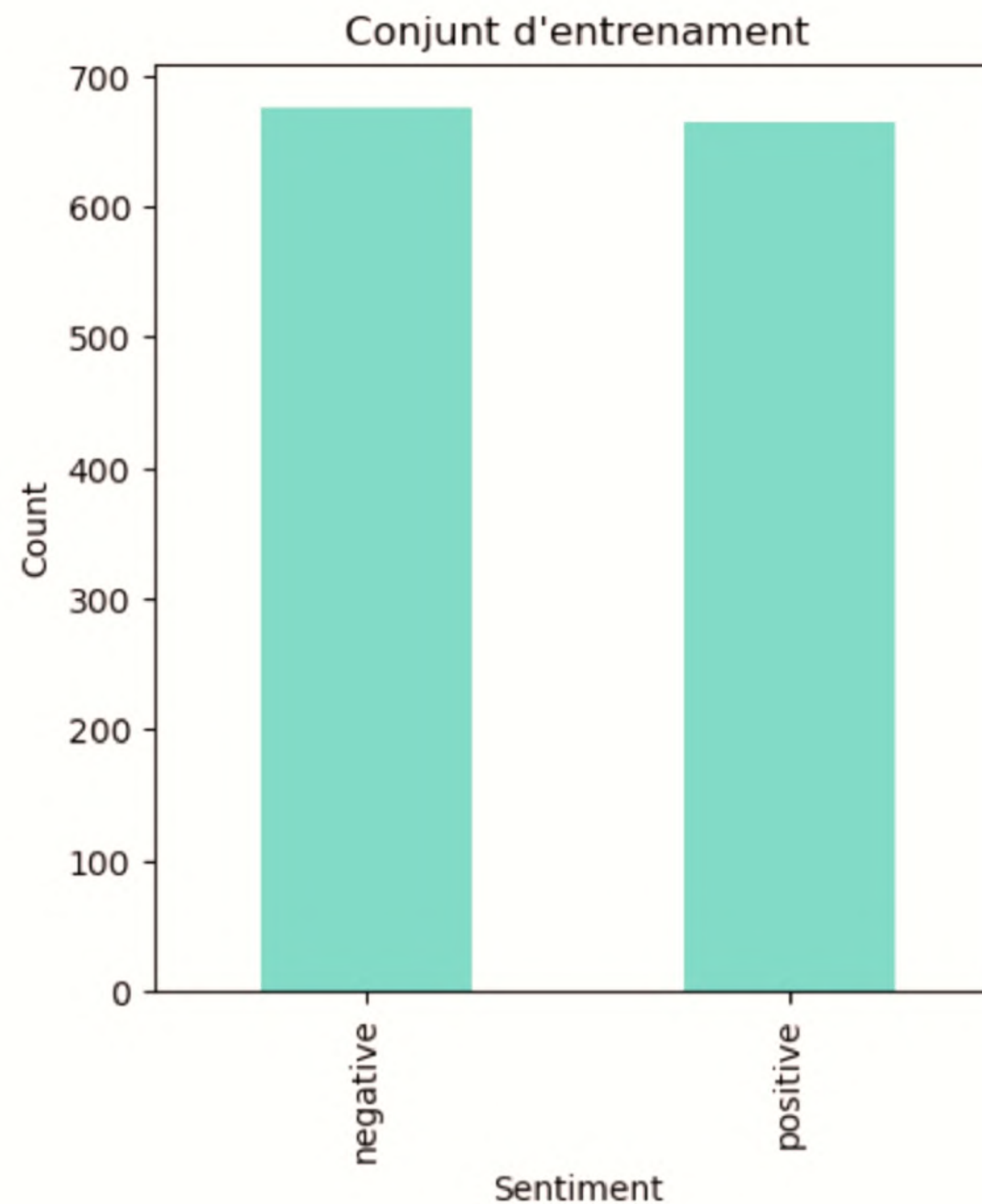
from sklearn.model_selection import GridSearchCV
```

ANÀlisi: TRAIN TEST SPLIT

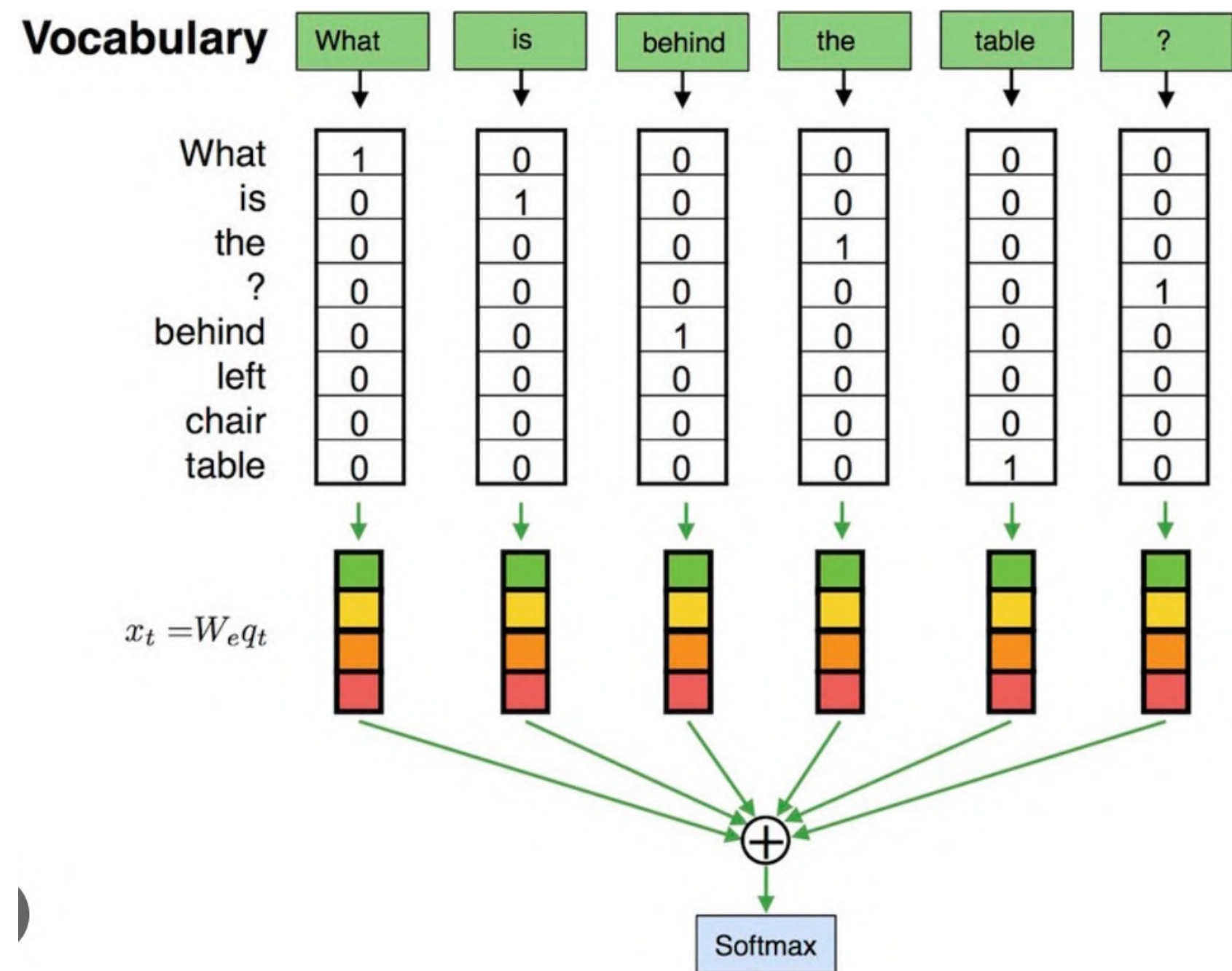
TRAIN_TEST_SPLIT SPLITS DATA INTO
TRAINING DATA AND TEST DATA



ANÀLISI: TRAIN TEST SPLIT



ANALISI: BOW (BAG OF WORDS)



ANÀlisi: APRENENTATGE SUPERVISAT

01

- ☐ **Support Vector Machines (SVM)**

02

- ☐ **Decision Tree**

03

- ☐ **Naive Bayes**

04

- ☐ **Logistic Regression**

05

- ☐ **Avaluació del Model**

ANÀLISI: AVALUACIÓ RENDIMENT

```
1 # Avaluu el rendiment del model SVM
2 svc_score = svc.score(test_x_tfidf, test_y)
3 print(f"Rendiment del model SVM: {round(svc_score, 2)}")
4
5 # Avaluu el rendiment del model Decision Tree
6 dec_tree_score = dec_tree.score(test_x_tfidf, test_y)
7 print(f"Rendiment del model Decision Tree: {round(dec_tree_score, 2)}")
8
9 # Avaluu el rendiment del model Naive Bayes
10 gnb_score = gnb.score(test_x_tfidf.toarray(), test_y)
11 print(f"Rendiment del model Naive Bayes: {round(gnb_score, 2)}")
12
13 # Avaluu el rendiment del model Logistic Regression
14 log_reg_score = log_reg.score(test_x_tfidf, test_y)
15 print(f"Rendiment del model Logistic Regression: {log_reg_score}")
16
```

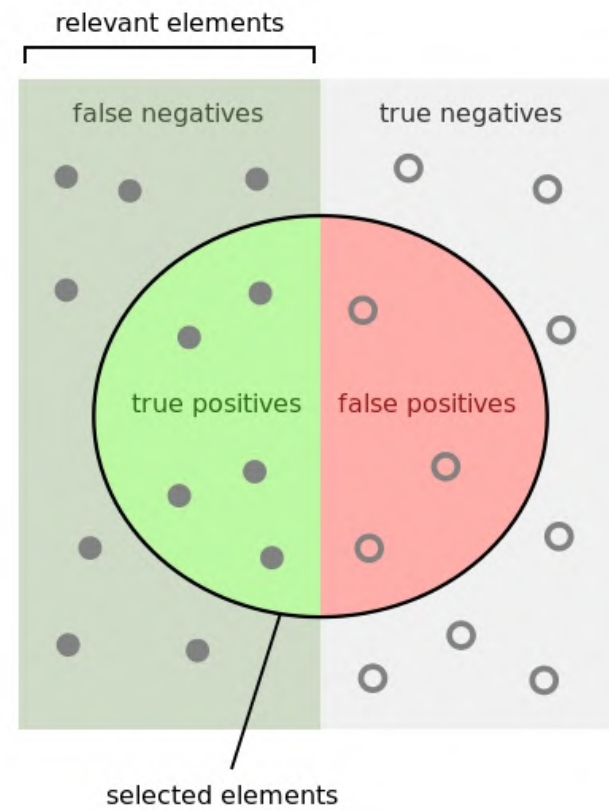
Rendiment del model SVM: 0.85

Rendiment del model Decision Tree: 0.51

Rendiment del model Naive Bayes: 0.63

Rendiment del model Logistic Regression: 0.85

F1 Score



How many selected items are relevant?

Precision =

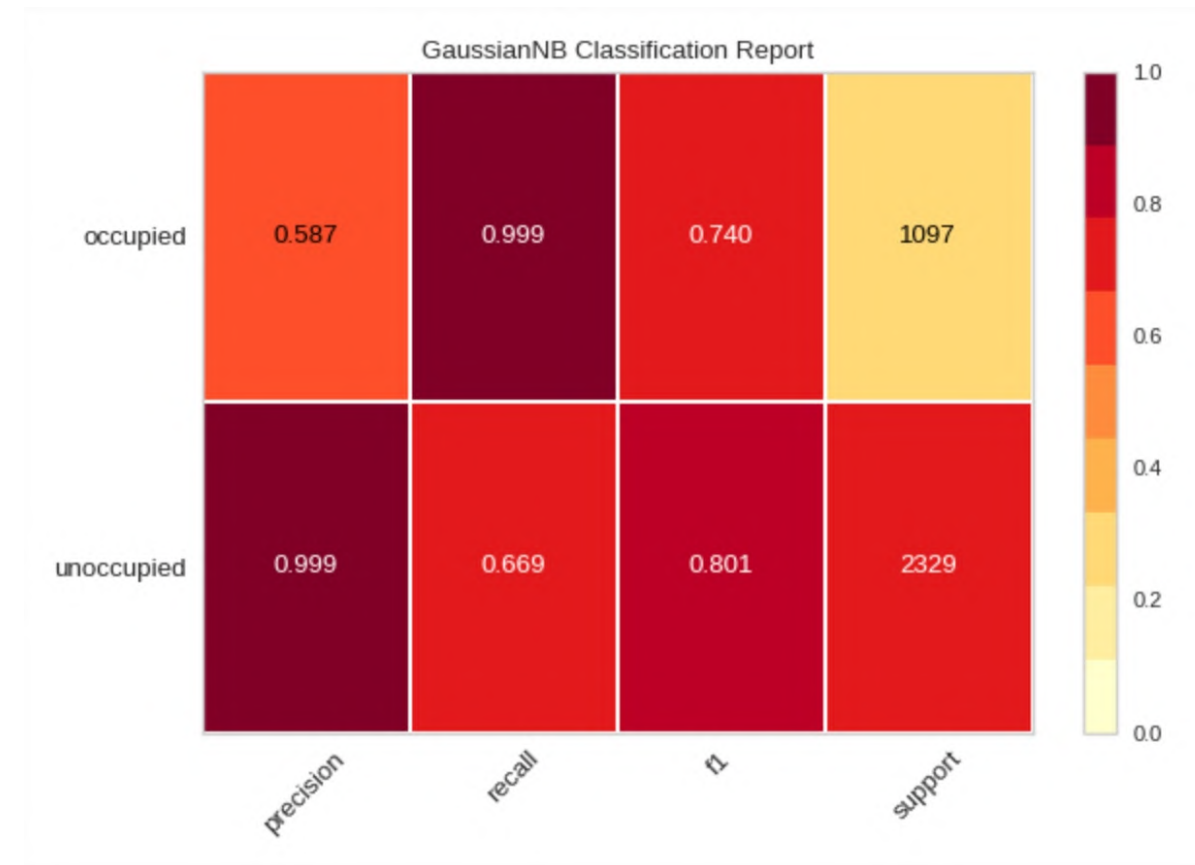


How many relevant items are selected?

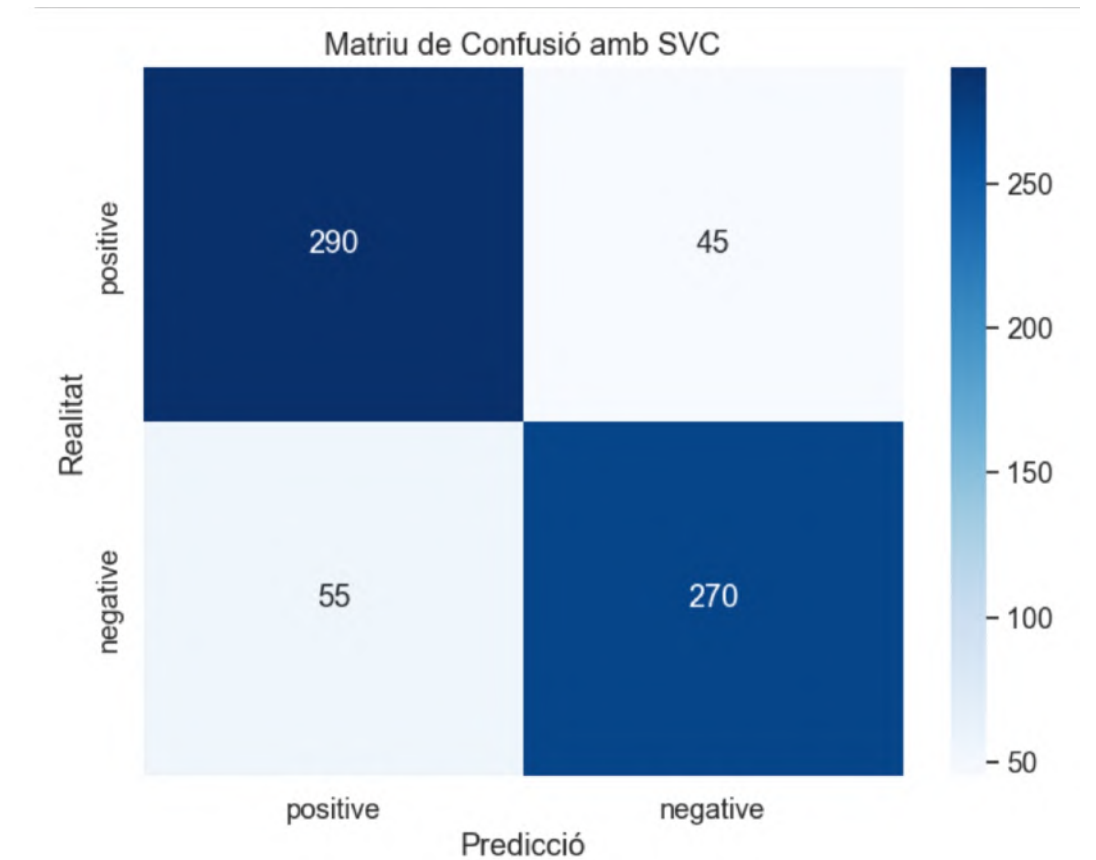
Recall =



Classification Report



Confusion Matrix



ANÀLISI: AVALUACIÓ DEL MODEL SEGONS SVM (SUPPORT VECTORS MACHINE)

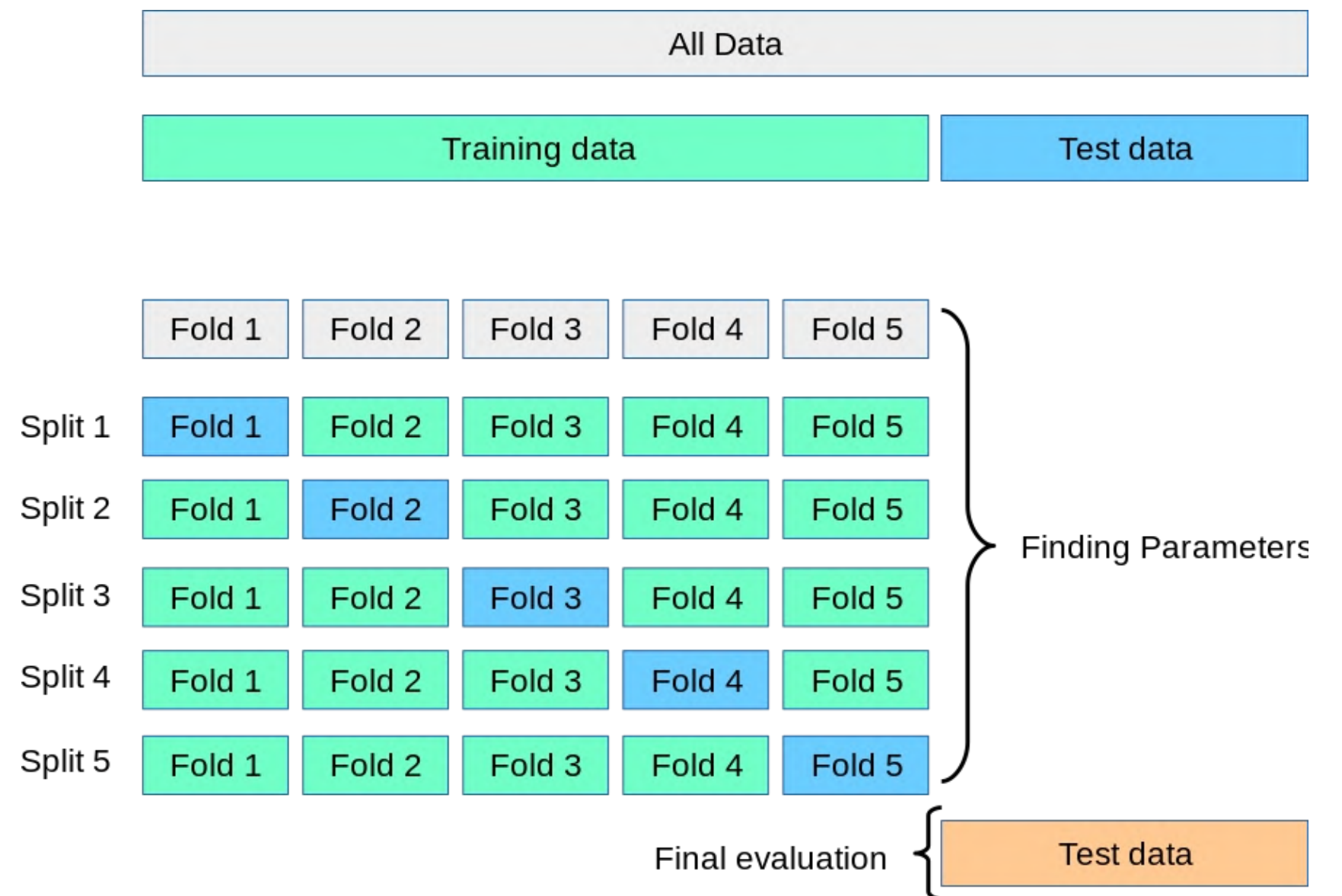
Classification Report

Informe de classificació amb SVC:

	precision	recall	f1-score	support
positive	0.84	0.87	0.85	335
negative	0.86	0.83	0.84	325
accuracy			0.85	660
macro avg	0.85	0.85	0.85	660
weighted avg	0.85	0.85	0.85	660

ANÀLISI: AVALUACIÓ DEL MODEL SEGONS SVM (SUPPORT VECTORS MACHINE)

AFINACIÓ DEL MODEL: GRID SEARCH CV



CONCLUSIÓ

Exactitud del model amb paràmetres
optimitzats: 0.8560606060606061

```
1 import time
2
3 # Inicia el rellotge
start_time = time.time()
4
5 # Resta el temps actual al temps d'inici per obtenir el temps d'inici
svc_grid = GridSearchCV(svc, parameters, cv=5)
svc_grid.fit(train_x_vector, train_y)
6
7 # Atura el rellotge
end_time = time.time()
8
9 # Calcula la diferència entre el temps d'inici i el temps final
10 execution_time = end_time - start_time
11
12 # Mostra els paràmetres òptims
13 print("Paràmetres òptims:", svc_grid.best_params_)
14
15 # Mostra el temps d'execució
16 print("Temps d'execució:", execution_time, "segons")
17
18
```

Paràmetres òptims: {'C': 4, 'kernel': 'rbf'}
Temps d'execució: 35.86082625389099 segons

```
1 # Millors paràmetres obtinguts de la cerca de paràmetres
2 best_params = {'C': 4, 'kernel': 'rbf'}
3
4 # Crea una nova instància de SVC amb els millors paràmetres
5 optimized_svc = SVC(**best_params)
6
7 # Entrena el model amb les dades d'entrenament
8 optimized_svc.fit(train_x_vector, train_y)
9
10 # Fes prediccions amb les dades de prova
11 predictions = optimized_svc.predict(test_x_vector)
12
13 # Avalua el rendiment del model
14 accuracy = optimized_svc.score(test_x_vector, test_y)
15 print(f"Exactitud del model amb paràmetres optimitzats: {accuracy}")
16
```

Exactitud del model amb paràmetres optimitzats: 0.8560606060606061

TENGO PRUEBAS!

L...

```
# Prediccions utilitzant el model SVM entrenat
prediction_1 = svc.predict(tfidf_vectorizer.transform(['A great film']))
prediction_2 = svc.predict(tfidf_vectorizer.transform(['A nice film but it could be better']))
prediction_3 = svc.predict(tfidf_vectorizer.transform(['A nice film']))
prediction_4 = svc.predict(tfidf_vectorizer.transform(['I was going to say something good, but I simply cannot be

# Mostrem les prediccions
print("Predicció per 'A great film':", prediction_1)
print("Predicció per 'A nice film but it could be better':", prediction_2)
print("Predicció per 'A nice film':", prediction_3)
print("Predicció per 'I was going to say something good, but I simply cannot because the film is so bad.':", prec
```

```
Predicció per 'A great film': ['positive']
Predicció per 'A nice film but it could be better': ['negative']
Predicció per 'A nice film': ['positive']
Predicció per 'I was going to say something good, but I simply cannot because the film is so bad.': ['negative']
```

Les conclusions del model són positives: fins ara, ha encertat tot.

CONCLUSIÓ

Desequilibri de les classes:

Ha permès l'aplicació de tècniques com l'undersampling per abordar aquesta qüestió.

Importància de la representativitat en la divisió de les dades:

La divisió aleatòria de les dades ajuda a evitar l'overfitting del model.

Aprenentatge supervisat:

Support Vector Machines (SVM) i Logistic Regression han presentat els millors resultats, acurats en un 85%.

Avaluació del model:

F1 Score, Classification Report i Confusion Matrix, i finalment la maximització amb GridSearchCV, han arribat a una precisió del un 85,6 %.

Reptes i millores:

Tècniques avançades de processament de llenguatge natural i cross-validation per avaluacions més robustes.

Contacting? Pones edaimon en google y ya te salgo.

THE END

FIN



Pitch

Want to make a presentation like this one?

Start with a fully customizable template, create a beautiful deck in minutes, then easily share it with anyone.

Create a presentation (It's free)

