


```

Creating Embedding Matrix...
Embedding Matrix Created
number of null word embeddings: 350 in a total of 7907 words (4.43%)
words not found: 349
e.g. 10 words not found in the index : ['hydroxyferuloyl' 'superior' 'jingzaojing' 'dk18' 'pholidoto1'
'dhs201' '8821058' 'furtherstudy' 'metabolism' 'piceides' 'gp'
'danegry' 'npus2' 'ata38' 'tpicicrops1']
(7908, 200)
[[ 0. 0. 0. ... 0. 0.
 0. 0. 0. ... 0. 0.
 0. 0. 0. ... 0. 0.
 ...
 [-0.0848686 -0.0417668 -0.3703938 ... 0.07843895 0.08479623
 0.34163074]
 ...
 [-0.1518244 -0.13880545 0.16630234 ... -0.29269406 -0.08847184
 -0.05204243]
 [-0.1121414 -0.08872321 -0.33860227 ... -0.34313729 0.25840917
 0.20746794]
 [-0.26328463 0.19522475 -0.05939044 ... 0.18617576 -0.05572284
 -0.12646711]]

In [24]: if multiple_gpus:
        with strategy(scope()):
            model = HAM_opt(d1_config.embedding_matrix, d1_config.learning_rate=d1_config.learning_rate,
                             seed_value=d1_config.seed_value)
        else:
            model = HAM_opt(d1_config.embedding_matrix, d1_config.learning_rate=d1_config.learning_rate,
                             seed_value=d1_config.seed_value)

history = model.fit(X_train, y_train,
                    epochs=d1_config.epochs,
                    batch_size=d1_config.batch_size,
                    validation_data=(x_val, y_val),
                    callbacks=keras_callbacks)

if d1_config.keras_callbacks:
    model.load_weights(checkpoint_path)

train_loss, d1_config.train_acc = model.evaluate(X_train, y_train, verbose=0, batch_size = d1_config.batch_size)

print("Training Loss: %3f" % train_loss)
print("Training Accuracy: %3f" % (d1_config.train_acc))

plot_training_history(history_dict=history, d1_config=d1_config)

Model: "model_3"
Layer (type) Output Shape Param #
-----
input_4 (InputLayer) (None, 15, 501) 0
time_distributed_1 (TimeDist (None, 15, 256) 1984544
spatial_dropout1d_3 (Spatial (None, 15, 256) 0
bidirectional_3 (Bidirection (None, 15, 512) 1050624
attention_with_context_3 (At (None, 512) 263168
dense_1 (Dense) (None, 1) 513
Total params: 3,298,849
Trainable params: 1,717,249
Non-trainable params: 1,581,600

None
Epoch 1/50
10/10 [=====] - 59s 2s/step - loss: 0.6597 - accuracy: 0.6348 - val_loss: 0.6267 - val
accuracy: 0.709
Epoch 00001: val_loss improved from inf to 0.62668, saving model to models\WC_pro\WC_pro_0\checkpoint.hdf5
Epoch 2/50
10/10 [=====] - 19s 2s/step - loss: 0.6003 - accuracy: 0.6846 - val_loss: 0.5957 - val
accuracy: 0.708
Epoch 00002: val_loss improved from 0.62668 to 0.58568, saving model to models\WC_pro\WC_pro_0\checkpoint.hdf5
Epoch 3/50
10/10 [=====] - 15s 2s/step - loss: 0.5011 - accuracy: 0.7421 - val_loss: 0.6630 - val
accuracy: 0.724
Epoch 00003: val_loss did not improve from 0.58568
Epoch 4/50
10/10 [=====] - 14s 1s/step - loss: 0.3875 - accuracy: 0.8029 - val_loss: 0.5964 - val
accuracy: 0.755
Epoch 00004: val_loss did not improve from 0.58568
Epoch 5/50
10/10 [=====] - 15s 1s/step - loss: 0.4476 - accuracy: 0.7862 - val_loss: 0.6795 - val
accuracy: 0.716
Epoch 00005: val_loss did not improve from 0.58568
Epoch 6/50
10/10 [=====] - 15s 1s/step - loss: 0.4363 - accuracy: 0.8215 - val_loss: 0.9600 - val
accuracy: 0.732
Epoch 00006: val_loss did not improve from 0.58568
Epoch 7/50
10/10 [=====] - 16s 2s/step - loss: 0.4547 - accuracy: 0.7943 - val_loss: 0.6681 - val
accuracy: 0.740
Epoch 00007: val_loss did not improve from 0.58568
Epoch 00007: early stopping

Accuracy
Training acc
Validation acc
epoch
0.65
0.7
0.75
0.8
0.85
0.9
1
2
3
4
5
6
7

Loss
Training loss
Validation loss
epoch
0.4
0.5
0.6
0.7
0.8
0.9
1
2
3
4
5
6
7

In [25]: x_test, y_test = D1_preprocessing(X_test_df, y_test_df, d1_config, dataset="test")

yhat_probs = model.predict(x_test, verbose=0)
yhat_probs = yhat_probs[:, 0]

yhat_classes = np.where(yhat_probs > 0.5, 1, yhat_probs)
yhat_classes = np.where(yhat_classes < 0.5, 0, yhat_classes).astype(np.int64)

Index of Unknown Words: 1

In [26]: d1_config.test_roc_auc, d1_config.test_pr_auc = plot_roc_pr_curves(y_test, yhat_probs, d1_config = d1_config)
# ROC AUC
print("ROC AUC: %f" % d1_config.test_roc_auc)

# avg precision
d1_config.test_avg_prec = average_precision(y_test_df, yhat_probs)
print("Average Precision: %f" % d1_config.test_avg_prec)

# accuracy
d1_config.test_acc = accuracy_score(y_test, yhat_classes)
print("Accuracy: %f" % d1_config.test_acc)

# precision tp / (tp + fp)
d1_config.test_prec = precision_score(y_test, yhat_classes)
print("Precision: %f" % d1_config.test_prec)

# recall: tp / (tp + fn)
d1_config.test_recall = recall_score(y_test, yhat_classes)
print("Recall: %f" % d1_config.test_recall)

# F1: 2 tp / (2 tp + fp + fn)
d1_config.test_f1_score = f1_score(y_test, yhat_classes)
print("F1 score: %f" % d1_config.test_f1_score)

# confusion matrix
matrix = confusion_matrix(y_test, yhat_classes)
matrix = np.flip(matrix)
print("Confusion Matrix:\n %s\n" % matrix)

ROC AUC: 0.848009
Average Precision: 0.721673
Accuracy: 0.719780
Precision: 0.722222
Recall: 0.220339
F1 score: 0.337662
Confusion Matrix:
[[138 5]
 [ 46 131]]

ROC Curve
True Positive Rate
False Positive Rate
Random
Tuned Model

Precision-Recall Curve
Precision
Recall
"Naïve Skill"
Tuned Model

```

```

0    1    2    3    4    5    6    7    8    9    10
False Positive Rate   Recall
-- Panel Model --

In [27]: """Imports"""

from dl_models import Hierarchical_Attention_LSTM


In [38]: from nltk.corpus import stopwords
stopos = set(stopwords.words("english"))

dl_config = DLConfig(remove_punctuation="HTML", stop_words=stopos, lower=True, split_by_hyphen=True)

df = pandas.read_excel("C:/Users/z6 Freitas/Desktop/Mestrado/2ºSemestre/Projeto/code/datasets/dataset_final.xlsx")
index_col=df["index_col"]
idfinal = pandas_column_aslist(df, "Document")
docsfinal = pmidf_to_docs(idfinal, "pq42872&atunos.unimho.pt", dl_config)
docsfinal = docsfinal[0]
docsfinal(0).abstract_string
dataset_docs = docs_to_pandasdocs(docsfinal)

x_total = dataset_docs
print(x_total)

y_total = pandas.read_excel("C:/Users/z6 Freitas/Desktop/Mestrado/2ºSemestre/Projeto/code/datasets/dataset_final.xlsx")
index_col=y["Index_col"], usecols="B:F")

print(y_total)

relevance = pandas_column_aslist(y_total, "Relevance")
y_total = relevances_to_pandas(y_total, relevance)
print(y_total)

print(x_total["Document"][0].title_string)
print(x_total["Document"][0].abstract_string)

X_train_df, X_test_df, y_train_df, y_test_df = train_test_split(x_total, y_total, test_size=0.3, random_state=
                                                                    stratify=y_total)

print(X_train_df.shape)
print(X_test_df.shape)
print(y_train_df.shape)
print(y_test_df.shape)

Document
25873669 <data_structures.document.Document object at 0...
24269970 <data_structures.document.Document object at 0...
17653245 <data_structures.document.Document object at 0...
16663649 <data_structures.document.Document object at 0...
33907682 <data_structures.document.Document object at 0...
32079399 <data_structures.document.Document object at 0...
26976395 <data_structures.document.Document object at 0...
20138774 <data_structures.document.Document object at 0...
27936616 <data_structures.document.Document object at 0...
28090095 <data_structures.document.Document object at 0...

[606 rows x 1 columns] Document
25873669 <data_structures.document.Document object at 0...
24269970 <data_structures.document.Document object at 0...

```

```

17655245 <data_structures.document.Document object at 0...
16663489 <data_structures.document.Document object at 0...
33907482 <data_structures.document.Document object at 0...
...
310279939 <data_structures.document.Document object at 0...
...
26976595 <data_structures.document.Document object at 0...
20138774 <data_structures.document.Document object at 0...
23784616 <data_structures.document.Document object at 0...
28009095 <data_structures.document.Document object at 0...
...
[606 rows x 1 columns]
25873669 0
24269870 0
17655245 1
16663489 1
33907482 0
20138799 ''
26976595 0
20138774 1
23784616 1
28009095 1
Name: Label, Length: 606, dtype: int64
genetic diversity stilbene metabolism vitis sylvestris .
stilbenes , important secondary metabolites grapevine , represent central phytoalexins therefore constitute imp
ortant element basal immunity . study , potential genetic variation vitis vinifera sp. sylvestris , ancestor
cultivars grapevine , sought respect output stilbenes potential use resistance breeding , considerable variati
on stilbene inducibility identified v. vinifera sp. sylvestris . genotypic differences abundance profiles sti
lbenes reduced response to p < 0.05 show . two clusters stilbenes 'chemotera' emerged ; one cluster showed quic
k strong accumulation stilbenes , almost exclusively from non - glycosylated resveratrol viniferin , second c
luster accumulated fewer stilbenes relatively high proportions piceannanol glycosylated piceid . 86 genotypes ,
wide abundance stilbene pattern observed : piceid , resveratrol , piceannanol accumulated earlier , whereas vi
niferins found later . observed genotypes differences stilbene accumulation preceded differential accumulation
transcripts chalcone synthase ( chs ) stilbene - related genes : phenylalanine ammonia lyase ( pal ) , stilbene
dehydratase ( sty1 ) , resveratrol synthase ( rs ) . screen population respect susceptibility downy mildew grape
vine ( plasmagora viticola ) revealed considerable variability , subpopulation genotypes high stilbene induci
bility significantly less v. p < 0.05 show . two clusters genotypes , representative genotypes could show in
duced stilbene synthase up correlated inducibility pathogen .
(424, 1)
(424, 1)
(424, 1)
(182, 1)
In (29):
model_name = "BA_15TM"
dl_config = dlConfig(model_name=model_name, seed_value=seed_value)
dl_config.stop_words = stops
dl_config.lower = True
dl_config.remove_punctuation = False
dl_config.split_by_hyphen = True
dl_config.lemmatization = False
dl_config.stems = False

#Parameters
dl_config.padding = 'post' # 'pre' -> default; 'post' -> alternative
dl_config.truncating = 'post' # 'pre' -> default; 'post' -> alternative #####

```

```

dl_config.ovv_token = 'OOV'

dl_config.epochs = 50
dl_config.batch_size = 32 # e aumentar o batch
dl_config.learning_rate = 0.001 #experimentar diminuir

dl_config.max_sent_len = 50 #sentences will have a maximum of "max_sent_len" words #400/500
dl_config.max_nb_words = 100_000 #it will only be considered the top "max_nb_words" words in the dataset
dl_config.max_nb_sentences = 15 # set only for the hierarchical attention model!!!!

dl_config.embeddings = 'biowordvec'

# igual ao modelo anterior
dl_config.embedding_path = './embeddings/biowordvec'
dl_config.embedding_dim = 200
dl_config.embedding_format = 'word2vec'

dl_config.keras_callbacks = True

# compare as losses d strain e validation
if dl_config.keras_callbacks:
    dl_config.patience = 5 #early-stopping patience
    checkpoint_path = str(dl_config.model_id_path) + '\\checkpoint.hdf5'
    keras_callbacks = [
        EarlyStopping(monitor='val_loss', mode='min', verbose=1, patience=dl_config.patience),
        ModelCheckpoint(checkpoint_path, monitor='val_loss', mode='min', verbose=1, save_best_only=True) #
    ]
else:
    keras_callbacks = None

```

```
d_config.validation_percentage = 30 # talve aumentar
X_train, y_train, x_val, y_val = Df_preprocessing(X_train_df, y_train_df,
                                                    d_config, dataset='train',
                                                    validation_percentage = d_config.validation_percentage,
                                                    seed_value=d_config.seed_value)

X_train[0][0]

Found 7907 unique tokens.
Index of Unknown Words: 1
Training set with 487 samples
Validation set with 127 samples

Out[30]: array([[ 22.,  50., 20., 698., 1122., 3434., 178.,   20., 698., 431., 497.,
        353., 4794., 164., 85.,    1.,   0.,   0.,   0.,   0.,   0.,   0.,
         0.,   0.,   0.,   0.,   0.,   0.,   0.,   0.,   0.,   0.,
         0.,   0.,   0.,   0.,   0.,   0.]])

In [31]: d_config.embedding_matrix = compute_embedding_matrix(dl_config, embeddings_format = dl_config.embeddings_format
# quantas palavras encontram no vocabulario do embedding

print(dl_config.embedding_matrix.shape)
print(dl_config.embedding_matrix)

Creating Embedding Matrix...
Embedding Matrix Created

-----
number of null word-embeddings: 350 in total of 7907 words (4.43%)
words not found: 349
e.g., 1 words not found in the dictionary: "hidroxiferuloyloil""superior""jingzaojing""dk18""pholidotoil""
cha201""s8200w""furtherstady""metabolitism""picicled""he
danfeng""yiba2""etab36""spiroctropis"]
(7908, 200)
[[ 0.          0.          ...  0.          0.
 [ 0.          0.          ...  0.          0.
 [ 0.          0.          ...  0.          0.
 [ 0.          0.          ...  0.          0.
 [-0.0848686  -0.0417668 -0.3709338  -0.07843895  0.08479623
  0.34630074)
'''
[ 0.15182443 -0.13880545  0.16630234 ... -0.29269406 -0.08847184
  -0.05204243)
[-0.1121414  -0.08872321 -0.33860227 ... -0.34313729  0.25840917
  0.20746794)
[ 0.26528463  0.19522475 -0.05939044 ... 0.18617576 -0.05572284
```

```

-0.1264671]]]
In [32]: if multiple_gpus:
        with strategy.scope():
            model = Hierarchical_Attention_LSTM(dl_config.embedding_matrix, dl_config.learning_rate=dl_config.learning_rate,
                                                seed_value=dl_config.seed_value)
        else:
            model = Hierarchical_Attention_LSTM(dl_config.embedding_matrix, dl_config.learning_rate=dl_config.learning_rate,
                                                seed_value=dl_config.seed_value)

history = model.fit(X_train, y_train,
                    epochs=dl_config.epochs,
                    batch_size=dl_config.batch_size,
                    validation_data=(x_val, y_val),
                    callbacks=keras_callbacks)

if dl_config.keras_callbacks:
    model.load_weights(checkpoint_path)

Model: "model_4"

Layer (type)                 Output Shape                 Param #
-----
word_input (InputLayer)      [(None, 50)]                 0
word_embedding (Embedding)    (None, 50, 200)             1581600
word_gru (Bidirectional)     (None, 50, 100)             100400
word_dense (Dense)           (None, 50, 100)             10100
word_attention (AttentionLay  [(None, 100), (None, 50), 20400

Total params: 1,712,500
Trainable params: 130,800
Non-trainable params: 1,581,600

```

```

Non-trainable params: 1752,600

None
Model: "model_5"

Layer (type)                Output Shape                Param #
-----
sent_input (InputLayer)      [(None, 15, 50)]          0
sent_linking (TimeDistribute) (None, 15, 100)            1712500
sent_gru (Bidirectional)     (None, 15, 100)            60400
dropout (Dropout)           (None, 15, 100)            0
sent_dense (Dense)          (None, 15, 100)            10100
sent_attention (AttentionLay [(None, 100), (None, 15, 20400
sent_dropout (Dropout)      (None, 100)                 0
output (Dense)               (None, 1)                   101
-----
Total params: 1,603,501
Trainable params: 221,901
Non-trainable params: 1,581,600

None
Epoch 1/50
10/10 [=====] - 56s 1s/step - loss: 0.6769 - accuracy: 0.5541 - val_loss: 0.6099 - val
_accuracy: 0.708
=====

Epoch 00001: val_loss improved from inf to 0.60891, saving model to models\RA_LSTM\RA_LSTM_0\checkpoint.hdf5
Epoch 2/50
10/10 [=====] - 6s 608ms/step - loss: 0.6719 - accuracy: 0.6586 - val_loss: 0.6111 - v
al_accuracy: 0.708
=====

```

```
Epocho 00002: val_loss did not improve from 0.60891
Epoch 3/50
10/10 [=====] - 6s 61ms/step - loss: 0.6352 - accuracy: 0.6610 - val_loss: 0.6054 - v
al_accuracy: 0.708
7C#####

Epocho 00003: val_loss improved from 0.60891 to 0.60544, saving model to models\BA_LSTM\BA_LSTM_0\checkpoint.hdf
5
Epoch 4/50
10/10 [=====] - 7s 701ms/step - loss: 0.6026 - accuracy: 0.6876 - val_loss: 0.5739 - v
al_accuracy: 0.708
7C#####

Epocho 00004: val_loss improved from 0.60544 to 0.57395, saving model to models\BA_LSTM\BA_LSTM_0\checkpoint.hdf
5
Epoch 5/50
10/10 [=====] - 7s 678ms/step - loss: 0.5768 - accuracy: 0.6531 - val_loss: 0.5567 - v
al_accuracy: 0.700
8C#####

Epocho 00005: val_loss improved from 0.57395 to 0.55669, saving model to models\BA_LSTM\BA_LSTM_0\checkpoint.hdf
5
Epoch 6/50
10/10 [=====] - 7s 690ms/step - loss: 0.5097 - accuracy: 0.7159 - val_loss: 0.4829 - v
al_accuracy: 0.724
4C#####

Epocho 00006: val_loss did not improve from 0.55669
Epoch 7/50
10/10 [=====] - 6s 614ms/step - loss: 0.5229 - accuracy: 0.7595 - val_loss: 0.6538 - v
al_accuracy: 0.724
4C#####
```

```
Epoch 00007: val_loss did not improve from 0.55669
Epoch 8/50
10/10 [=====] - 7s 681ms/step - loss: 0.4294 - accuracy: 0.7785 - val_loss: 0.5995 - v
al_accuracy: 0.724
=====

Epoch 00008: val_loss did not improve from 0.55669
Epoch 9/50
10/10 [=====] - 6s 637ms/step - loss: 0.3833 - accuracy: 0.7948 - val_loss: 0.6423 - v
al_accuracy: 0.716
=====

Epoch 00009: val_loss did not improve from 0.55669
Epoch 10/50
10/10 [=====] - 7s 708ms/step - loss: 0.3938 - accuracy: 0.7814 - val_loss: 0.6620 - v
al_accuracy: 0.629
=====

Epoch 00010: val_loss did not improve from 0.55669
Epoch 00010: early stopping

In [33]: train_loss, dl_config, train_acc = model.evaluate(X_train, y_train, verbose=0, batch_size = dl_config.batch_size)
plot_training_history(history_dict=history, dl_config=dl_config)

X_test, y_test = DL_preprocessing(X_test_df, y_test_df, dl_config, dataset='test')

yhat_probs = model.predict(X_test, verbose=0)
yhat_probs = yhat_probs[:, 0]

yhat_classes = np.where(yhat_probs > 0.5, 1, yhat_probs)
yhat_classes = np.where(yhat_classes < 0, 0, yhat_classes)
accuracy = (yhat_classes == y_test).sum() / y_test.shape[0]
```

```

yhat_classes = np.where(yhat_classes < 0, 0, yhat_classes).astype(int).flatten()
dl_config.test_roc_auc, dl_config.test_auc = plot_roc_n_pr_curves(y_test, yhat_probs, dl_config)

# ROC AUC
print('ROC AUC: %f' % dl_config.test_roc_auc)

# avg precision
dl_config.test_avg_prec = average_precision(y_test_df, yhat_probs)
print('Average Precision: %f' % dl_config.test_avg_prec)

# accuracy
dl_config.test_acc = accuracy_score(y_test, yhat_classes)
print('Accuracy: %f' % dl_config.test_acc)

# precision tp / (tp + fp)
dl_config.test_prec = precision_score(y_test, yhat_classes)
print('Precision: %f' % dl_config.test_prec)

# recall: tp / (tp + fn)
dl_config.test_recall = recall_score(y_test, yhat_classes)
print('Recall: %f' % dl_config.test_recall)

# F1: 2 tp / (2 tp + fp + fn)
dl_config.test_f1_score = f1_score(y_test, yhat_classes)
print('F1 score: %f' % dl_config.test_f1_score)

# confusion matrix
matrix = confusion_matrix(y_test, yhat_classes)
matrix = np.flip(matrix)
print('Confusion Matrix\n %s\n' % matrix)

Index of Unknown Words: 1
ROC AUC: 0.857192
Average Precision: 0.748001
Accuracy: 0.686813
Recall: 0.750000

```

Recall: 0.050847
F1 score: 0.095238
Confusion Matrix:
[[122 1]
 [56 31]]

The figure displays four plots related to model performance over 10 epochs:

- Accuracy:** Training accuracy (blue line) starts at ~0.65, rises to ~0.75 by epoch 6, peaks at ~0.82 around epoch 8, and ends at ~0.81. Validation accuracy (orange line) starts at ~0.71, dips slightly at epoch 5, rises to ~0.72 by epoch 6, peaks at ~0.73 around epoch 7, and ends at ~0.63.
- Loss:** Training loss (blue line) starts at ~0.65, decreases to ~0.55 by epoch 6, dips to ~0.42 at epoch 7, and ends at ~0.32. Validation loss (orange line) starts at ~0.61, decreases to ~0.55 by epoch 5, spikes to ~0.92 at epoch 6, and ends at ~0.66.
- ROC Curve:** The training ROC curve (blue dashed line) is a diagonal line from (0,0) to (1,1). The validation ROC curve (orange solid line) is a step function, starting at (0,0), rising to ~0.9 at epoch 5, and reaching 1.0 by epoch 10.
- Precision-Recall Curve:** The training precision-recall curve (blue dashed line) is a horizontal line at precision ~0.9. The validation precision-recall curve (orange solid line) starts at precision ~0.6, rises to ~0.85 by epoch 5, and then fluctuates between 0.6 and 0.85 until epoch 10.

```
print(x_total)

relevance = pandas_column_satisfy(x_total, "Relevance")
y_total = relevance.to_pandas(y_total, relevance)
print(y_total)

print(x_total["Documentum"][0].title_string)
print(x_total["Documentum"][0].abstract_string)

X_train_df, X_test_df, y_train_df, y_test_df = train_test_split(x_total, y_total, test_size=0.3, random_state=
                                                                    stratify=y_total)

print(X_train_df.shape)
print(X_test_df.shape)
print(y_train_df.shape)
print(y_test_df.shape)

stilbenes, important secondary metabolites grapevine, represent central phytoalexins therefore constitute important element basic immunity, study potential genetic variation vitis vinifera ssp. sylvestris. ancestor cultivated grapevine, sought respect output stilbenes potential use resistance breeding. considerable variation on stilbene inducibility identified v. vinifera ssp. sylvestris. genotypic differences abundance profiles stilbenes induced response observed: two clusters stilbenes 'chemovars' emerged; one cluster showed quickly k strong accumulation stilbenes, almost exclusively form non-glycosylated resveratrol viniferin, second cluster ac accumulated fewer stilbenes relatively high proportions piceatannol glycosylated piceid. 96 genotypes, time dependence stilbene pattern observed: piceid, resveratrol, piceatannol accumulated earlier, whereas viniferins found later. observed genotypic differences stilbene accumulation preceded differential accumulation trans-epicatechin chalcone (ch) and piceatannol (p). genotypes: phenylalanine ammonia lyase (pal), stilbene synthase (st), resveratrol synthase (rs). screen population respect susceptibility down mildew grapevine (plasmopara viticola) revealed considerable variability. subpopulation genotypes high stilbene inducibility (st) similarly less susceptible compared low-stilbene genotypes. representative genotypes could show inducibility stilbene synthesis vs correlated inducibility pathogen.
```

```

320739991  <data_structures.document.Document object at 0...
320739992  <data_structures.document.Document object at 0...
20138774   <data_structures.document.Document object at 0...
23936616  <data_structures.document.Document object at 0...
28009095  <data_structures.document.Document object at 0...

[666 rows x 1 columns]
25973669  0
24269870  0
17655245  1
1663649   1
33907682  0

320739991  1
26978595  0
20138774  1
23936616  1
28009095  1
Name: Label, Length: 666, dtype: int64

genetic diversity stilbene metabolism vitis sylvestris .
stilbenes , important secondary metabolites grapevine , represent central phytoalexins therefore constitute im-
portant element basal immunity . study , potential genetic variation vitis vinifera ssp . sylvestris , ancestor
cultivated grapevine , sought respect output stilbenes potential use resistance breeding . considerable variation
in stilbene inducibility identified . vitifera ssp . sylvestris , genotypic differences abundance profiles atilbenes
induced response uv-pulse shown . two clusters stilbene 'chemovars' emerged : one cluster showed quick
response accumulation stilbenes , almost exclusively non-glycosylated resveratrolviniferin , second cluster
or accumulated fewer stilbenes relatively high proportions piceatannol glycosylated piceid . 86 genotypes , time
dependence stilbene pattern observed : piceid , resveratrol , piceatannol accumulated earlier , whereas viniferin
emerged later . observed genotypic differences stilbene accumulation preceded differential accumulation tran-
scripts chalcone synthase ( chs ) stilbene-related genes : phenylalanine ammonium lyase ( pal ) , stilbene syn-
thase ( sty ) , resveratrol synthase ( rs ) . screen population respect susceptibility downy mildew grapevine
vitifera vitis vinifera ssp . sylvestris , subpopulation genotypes high stilbene inducibility
significantly less susceptible compared low-stilbene genotypes , representative genotypes could shown inducibility
of stilbene synthase uv correlated inducibility pathogen .

(424 , 1)
(162 , 1)
(424 , 1)
(162 , 1)

```

```

in [37]:
model_name = "BURNS"
dl_config = DLConfig(model_name=model_name, seed_value=seed_value)
dl_config.stop_words = stops
dl_config.lower = True
dl_config.remove_punctuation = False
dl_config.split_by_hyphen = False
dl_config.immatination = False
dl_config.stems = True

dl_config.padding = 'post'           #'pre' -> default; 'post' -> alternative
dl_config.truncating = 'post'       #'pre' -> default; 'post' -> alternative      #####
dl_config.oov_token = "OOV"

dl_config.epochs = 50
dl_config.batch_size = 32           # e amentar o batch
dl_config.learning_rate = 0.001     #experimentar diminuir

dl_config.max_sent_len = 300         #sentences will have a maximum of "max_sent_len" words   #400/500
dl_config.max_nb_words = 100_000     #it will only be considered the top "max_nb_words" words in the dataset

dl_config.embeddings = 'biwordvec'

dl_config.embedding_path = '..\\embeddings\\biwordvec'
dl_config.embedding_dim = 200
dl_config.embedding_format = 'word2vec'

dl_config.keras_callbacks = True

# compare as losses d'otrain e validation
if dl_config.keras_callbacks:
    dl_config.patience = 5 #early-stopping patience
    checkpoint_path = str(dl_config.embedding_id_path) + '\\checkpoint.hdf5'
    keras_callbacks = [
        EarlyStopping(monitor='val_loss', mode='min', verbose=1, patience=dl_config.patience),

```

```
}  
    ModeCheckpoint(checkpoint_path, monitor='val_loss', mode='min', verbose=1, save_best_only=True) #  
  
else:  
    keras_callbacks = None  

```

In [38]:

```
d1_config.tokenizer = text.Tokenizer(num_words=d1_config.max_nb_words, oov_token=d1_config.oov_token)  
  
d1_config.validation_percentage = 30 # talvez aumentar  
X_train, y_train, x_val, y_val = dl_preprocessing(X_train_df, y_train_df,  
                                                  d1_config.dataset='train',  
                                                  validation_percentaged=d1_config.validation_percentage,  
                                                  seed_value=d1_config.seed_value)  
  
X_train[0][0]
```

Found 7907 unique tokens.
Training set with 297 samples
Validation set with 127 samples

Out[38]: 22

In [39]:

```
d1_config.embedding_matrix = compute_embedding_matrix(d1_config, embeddings_format = d1_config.embeddings_forma  
  
print(d1_config.embedding_matrix.shape)  
print(d1_config.embedding_matrix)
```

Creating Embedding Matrix...
Embedding Matrix Created

number of null word embeddings: 350 in a total of 7907 words (4.43%)
words not found: 349
e.g. 10 words not found in the index : ['hydroxyferuloyl' ''superior'' ''jingzaojiong'' 'dkil8' ''pholidotol'
''the201'' ''bz30m'' ''furtherastudy'' ''metabolism'' ''picieides'' ''6p'
''dafengse'' ''nyba2'' ''st36'' ''spirokropolis'']
(7908, 200)
[[0. 0. 0. 0.

```
[
  [ 0.         ]
  [ 0.         ] 0.         ... 0.         0.
  [ 0.         ]
  [-0.0848686  -0.0417668  -0.3709338  ... 0.07843895  0.08479623
    0.34163074]
  ...
  [ 0.15182443  -0.13880545  0.16630234  ... -0.29269406  -0.08847184
    -0.05204243]
  [-0.1121414  -0.08872321  -0.33860227  ... -0.34313729  0.25840917
    -0.20746794]
  [ 0.26528463  0.19522475  -0.05939044  ... 0.18617576  -0.05572284
    -0.12646711]]

In [40]: #if multiple gpus:
with tf.nn.nested_scope(''):
    model = Burns_CNNBiLSTM(dl_config.embedding_matrix, dl_config.learning_rate=dl_config.learning_rate,
                             seed_value=dl_config.seed_value)
else:
    model = Burns_CNNBiLSTM(dl_config.embedding_matrix, dl_config.learning_rate=dl_config.learning_rate,
                             seed_value=dl_config.seed_value)

history = model.fit(X_train, y_train,
                    epochs=dl_config.epochs,
                    batch_size=dl_config.batch_size,
                    validation_data=(x_val, y_val),
                    callbacks=keras_callbacks)

if dl_config.keras_callbacks:
    model.load_weights(checkpoint_path)

train_loss, dl_config.train_acc = model.evaluate(X_train, y_train, verbose=0)

print("Training Loss: %.3f" % (train_loss))
print("Training Accuracy: %.3f" % (dl_config.train_acc))
```

```
plot_training_history(history_dict == history, d1_config=mdl_config)

x_test, y_test = dl_preprocessing(X_test_df, y_test_df, d1_config, dataset='test')

yhat_probs = model_predict(x_test, verbose=0)
yhat_probs = yhat_probs[:, 0]

yhat_classes = np.where(yhat_probs > 0.5, 1, yhat_probs)
yhat_classes = np.where(yhat_classes < 0.5, 0, yhat_classes).astype(np.int64)
```



```

Model: "model_6"
Layer (type)                Output Shape          Param #          Connected to
input_5 (InputLayer)        (None, 3000)         0                (None, 3000)
embedding (Embedding)       (None, 300, 200)     1581600          input_5[0][0]
dropout_1 (Dropout)         (None, 300, 200)     0                embedding[0][0]
conv1d (Conv1D)             (None, 300, 16)      16016            dropout_1[0][0]
max_pooling1d (MaxPooling1D) (None, 75, 16)        0                conv1d[0][0]
conv1d_1 (Conv1D)           (None, 75, 16)       1286             max_pooling1d[0][0]
dropout_2 (Dropout)         (None, 75, 16)        0                conv1d_1[0][0]
bidirectional_4 (Bidirectional) (None, 128)         41472            dropout_2[0][0]
dropout_3 (Dropout)         (None, 128)          0                bidirectional_4[0][0]
dense_2 (Dense)             (None, 1)            129              dropout_3[0][0]
flatten (Flatten)           (None, 1)            0                dense_2[0][0]
activation (Activation)      (None, 1)            0                flatten[0][0]
repeat_vector (RepeatVector) (None, 128, 1)       0                activation[0][0]
permute (Permute)           (None, 1, 128)       0                repeat_vector[0][0]
multiply (Multiply)          (None, 1, 128)       0                dropout_3[0][0]
lambda (Lambda)             (None, 128)          0                permute[0][0]
dense_3 (Dense)             (None, 1)            129              multiply[0][0]
Total params: 1,640,642
Trainable params: 59,042
Non-trainable params: 1,581,600

Epoch 1/50
10/10 [=====] - 68s 5s/step - loss: 0.6748 - accuracy: 0.6406 - val_loss: 0.6078 - val
accuracy: 0.708
Epoch 00001: val_loss improved from inf to 0.60784, saving model to models\BURNS\BURNS_0\checkpoint hdf5
Epoch 3/50
10/10 [=====] - 2s 196ms/step - loss: 0.6622 - accuracy: 0.6586 - val_loss: 0.6003 - v
al accuracy: 0.708
Epoch 00003: val_loss improved from 0.60784 to 0.60034, saving model to models\BURNS\BURNS_0\checkpoint hdf5
Epoch 5/50
10/10 [=====] - 2s 196ms/step - loss: 0.6312 - accuracy: 0.6610 - val_loss: 0.5946 - v
al accuracy: 0.708
Epoch 00005: val_loss improved from 0.60034 to 0.59462, saving model to models\BURNS\BURNS_0\checkpoint hdf5
Epoch 7/50
10/10 [=====] - 2s 196ms/step - loss: 0.5870 - accuracy: 0.6909 - val_loss: 0.6461 - v
al accuracy: 0.629
Epoch 00007: val_loss did not improve from 0.59462
Epoch 9/50
10/10 [=====] - 2s 194ms/step - loss: 0.5876 - accuracy: 0.6710 - val_loss: 0.5403 - v
al accuracy: 0.724
Epoch 00009: val_loss improved from 0.59462 to 0.54026, saving model to models\BURNS\BURNS_0\checkpoint hdf5
Epoch 11/50
10/10 [=====] - 2s 191ms/step - loss: 0.5566 - accuracy: 0.7338 - val_loss: 0.5633 - v
al accuracy: 0.732
Epoch 00011: val_loss did not improve from 0.54026
Epoch 13/50
10/10 [=====] - 2s 194ms/step - loss: 0.5304 - accuracy: 0.7111 - val_loss: 0.5693 - v
al accuracy: 0.637
Epoch 00013: val_loss did not improve from 0.54026
Epoch 15/50
10/10 [=====] - 2s 196ms/step - loss: 0.4354 - accuracy: 0.7922 - val_loss: 0.5628 - v
al accuracy: 0.645
Epoch 00015: val_loss did not improve from 0.54026
Epoch 17/50
10/10 [=====] - 2s 198ms/step - loss: 0.4167 - accuracy: 0.8163 - val_loss: 0.6512 - v
al accuracy: 0.716
Epoch 00017: val_loss did not improve from 0.54026
Epoch 19/50
10/10 [=====] - 2s 195ms/step - loss: 0.4372 - accuracy: 0.8173 - val_loss: 0.6400 - v
al accuracy: 0.606
Epoch 00019: val_loss did not improve from 0.54026
Epoch 00020: early stopping
Training loss: 0.531
Training Accuracy: 0.717

Accuracy
Validation acc
Training acc
epochs
Loss
Training loss
Validation loss
epochs

d1_config.test_roc_auc, d1_config.test_pr_auc = plot_roc_n_pr_curves(y_test, y_hat_probs, d1_config\d1_config)

# ROC AUC
print("ROC AUC: %f" % d1_config.test_roc_auc)

# avg precision
d1_config.test_avg_prec = average_precision(y_test_df, y_hat_probs)
print("Average Precision: %f" % d1_config.test_avg_prec)

# accuracy
d1_config.test_acc = accuracy_score(y_test, y_hat_classes)
print("Accuracy: %f" % d1_config.test_acc)

# precision tp / (tp + fp)
d1_config.test_prec = precision_score(y_test, y_hat_classes)
print("Precision: %f" % d1_config.test_prec)

# recall: tp / (tp + fn)
d1_config.test_recall = recall_score(y_test, y_hat_classes)
print("Recall: %f" % d1_config.test_recall)

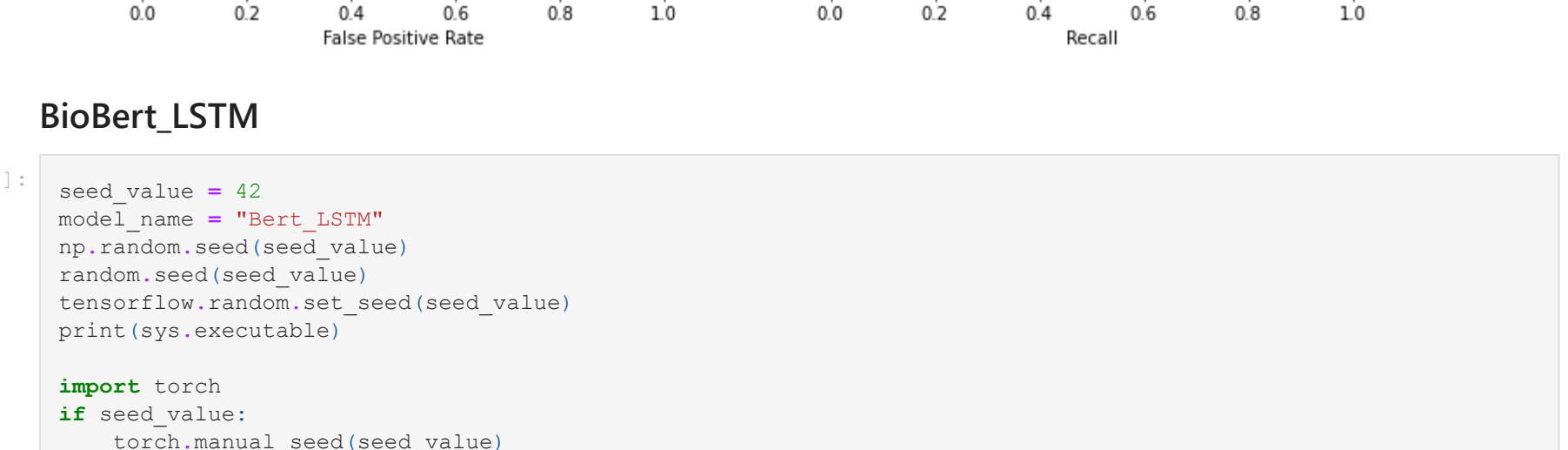
# f1: 2 tp / (2 tp + fp + fn)
d1_config.test_f1_score = f1_score(y_test, y_hat_classes)
print("F1 score: %f" % d1_config.test_f1_score)

# confusion matrix
matrix = confusion_matrix(y_test, y_hat_classes)
matrix = np.flip(matrix)
print("Confusion Matrix\n%s\n" % matrix)

ROC AUC: 0.853521
Average Precision: 0.696908
Accuracy: 0.723275
Precision: 0.736842
Recall: 0.237288
F1 score: 0.358974
Confusion Matrix:
[[128  3]
 [ 45 141]]

ROC Curve
Precision-Recall Curve
True Positive Rate
False Positive Rate
Precision
Recall
-- "No Skill"
-- Trained Model

```



```
torch.cuda.manual_seed_all(seed_value)
torch.backends.cudnn.deterministic = True
torch.backends.cudnn.benchmark = False
os.environ['PYTHONHASHSEED'] = str(seed_value)
os.environ['TF_DETERMINISTIC_OPS'] = '1'

if multiple_gpus:
    devices = []
    for gpu in multiple_gpus:
        devices.append('/gpu:' + str(gpu))
    strategy = tensorflow.distribute.MirroredStrategy(devices=devices)
    os.environ["CUDA_VISIBLE_DEVICES"] = ""

else:
    # Get the GPU device name.
    device_name = tensorflow.test.gpu_device_name()
    # The device name should look like the following:
    # /device:GPU:0
    print("Using GPU: {}".format(device_name))
    else:
        raise SystemError("GPU device not found")

os.environ["CUDA_VISIBLE_DEVICES"] = device_name
os.environ["CUDA_DEVICE_ORDER"] = "PCI_BUS_ID"
```

```

root@hpcia02:~/task0/DeviceGPU0:
INFO:tensorflow:Using MicroStrategy with devices (['/job:localhost/replica:0/task:0/DeviceGPU0:'])

***Imports***

from dl import Bert_preprocessing
from dl.models import Bert_LSTM
from transformers import BertTokenizer

os.environ["WANDB_API_KEY"] = "" # to silence warning
tf.nn.embedding_lookup.convert = "model_type bert"
__check_pointing "/Users/De Freitas/Desktop/Neatrodo/2Semestre/Projeto/code/hidbert_v1_l_pumbed/model_checkpoint" %C:/Users/De Freitas/Desktop/Neatrodo/2Semestre/Projeto/code/hidbert_v1_l_pumbed/bert_config.json"
__pytorch_checkpoint "/C:/Users/De Freitas/Desktop/Neatrodo/2Semestre/Projeto/code/hidbert_v1_l_pumbed.pytorch_model.bin"

Building Pytorch model from configuration: BertConfig {
  "attention_probs_dropout_prob": 0.1,
  "gradient_checkpointing": false,
  "hidden_act": "gelu",
  "hidden_dropout_prob": 0.1,
  "hidden_size": 768,
  "initializer_range": 0.02,
  "intermediate_size": 3072,
  "layer_norm_eps": 1e-12,
  "max_position_embeddings": 512,
  "model_type": "bert",
  "num_attention_heads": 12,
  "num_hidden_layers": 12,
  "pad_token_id": 0,
  "position_embedding_type": "absolute",
  "transformers_version": "4.6.1",
  "type_vocab_size": 2
  "use_cache": true,
  "vocab_size": 28996
}

Save Pytorch model to C:/Users/De Freitas/Desktop/Neatrodo/2Semestre/Projeto/code/hidbert_v1_l_pumbed/pytorch_model.bin

2021-06-02 19:00:48.749461: I tensorflow/stream_executor/platform/default/dso_loader.cc:60] Could not load dynamic
ml library 'cudart64_10.1.107.dll'; dlopen: cudart64_10.1.101 not found
2021-06-02 19:00:48.749591: I tensorflow/stream_executor/cuda/cuda_stub.cc:29] Ignore above cudart dlerror if you do not have a GPU set up on your machine.
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 0 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 1 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 2 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 3 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 4 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 5 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 6 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 7 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 8 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 9 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 10 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 11 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 12 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 13 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 14 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 15 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 16 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 17 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 18 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 19 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 20 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 21 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 22 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 23 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 24 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 25 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 26 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 27 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 28 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 29 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 30 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 31 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 32 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 33 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 34 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 35 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 36 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 37 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 38 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 39 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 40 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 41 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 42 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 43 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 44 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 45 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 46 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 47 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 48 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 49 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 50 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 51 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 52 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 53 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 54 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/cuda/cuda_device_reporter.cc:263] GeForce 1080Ti 55 capable GPU with 8GB memory
2021-06-02 19:00:48.750000: I tensorflow/stream_executor/c
```

[illegible][illegible]

```

y_total = dataset_data_docs
print(y_total)

y_total = pandas.read_excel('C:/Users/D/Downloads/Tratado/2ºSemestre/Projeto/code/datasets/dataset_final.xlsx', index_col=1, usecols="B:I")

print(y_total)

relevance = pandas.columns.tolist(y_total, "Relevance")
y_total = relevance.to_pandas(y_total, "relevance")
print(y_total)

print(y_total["Document"][0].title_string)
print(y_total["Document"][0].abstract_string)

X_train_df, X_test_df, y_train_df, y_test_df = train_test_split(x_total, y_total, test_size=0.3, random_state=y_total)

print(X_train_df.shape)
print(X_test_df.shape)
print(y_train_df.shape)
print(y_test_df.shape)

Stilbenes, as important secondary metabolites of grapevine, represent central phytoalexins and therefore constitute an important element of basal immunity. In this study, potential genetic variation in Vitis vinifera ssp. vinifera, the ancestor of cultivated grapevine, was sought with respect to their output of stilbenes and potential use for resistance breeding. Considerable variation in stilbene inducibility was identified in V. vinifera ssp. vinifera. Genotypic differences in abundance and profiles of stilbenes that are induced in response to grapevine pulp are shown. Two clusters of stilbene "chemovars" emerged; one cluster showed quick and strong accumulation of stilbenes, almost exclusively in the form of non-glycosylated resveratrol and viniferin while the second cluster accumulated fewer stilbenes and relatively high proportions of piceannol and the glycosylated piceid for all 86 genotypes, a time dependence of the stilbene pattern was observed: piceid, p-resveratrol, and piceannol accumulated earlier, whereas the viniferins were found later. It was further observed that the two clusters were preceded by differential accumulation of the transcripts for chalcone synthase (CHS) and stilbene-related genes: phenylalanine ammonium lyase (PAL), a stilbene synthase (StSy), and resveratrol synthase (RS). A screen of the population with respect to susceptibility to Vitis rotifoliorum mildew of grapevine (Vitis rotifoliorum) revealed considerable variability. The subpopulation of genotypes with high stilbene inducibility was significantly less susceptible as compared with low-stilbene genotypes, and for representative genotypes it could be shown that the inducibility of stilbene synthase B BY correlated with the inducibility by the pathogen.

Document
24268769 <data_structures.document.Document object at 0...
24268870 <data_structures.document.Document object at 0...
17655245 <data_structures.document.Document object at 0...
16663649 <data_structures.document.Document object at 0...
33907682 <data_structures.document.Document object at 0...
...
32079399 <data_structures.document.Document object at 0...
26976595 <data_structures.document.Document object at 0...
20138774 <data_structures.document.Document object at 0...
27936616 <data_structures.document.Document object at 0...
28080995 <data_structures.document.Document object at 0...

(606 rows x 1 columns)

Document
25873669 <data_structures.document.Document object at 0...
24268870 <data_structures.document.Document object at 0...
17655245 <data_structures.document.Document object at 0...
16663649 <data_structures.document.Document object at 0...
33907682 <data_structures.document.Document object at 0...
...
32079399 <data_structures.document.Document object at 0...
26976595 <data_structures.document.Document object at 0...
20138774 <data_structures.document.Document object at 0...
27936616 <data_structures.document.Document object at 0...
28080995 <data_structures.document.Document object at 0...

(606 rows x 1 columns)

25873669 1
24268870 1
17655245 1
16663649 1
33907682 0
32079399 1
20138774 1
27936616 1
28080995 1
Names: Label, Length: 606, dtype: int64

Genetic diversity of stilbene metabolism in Vitis vinifera ssp. vinifera. In this study, potential genetic variation in Vitis vinifera ssp. vinifera, the ancestor of cultivated grapevine, was sought with respect to their output of stilbenes and potential use for resistance breeding. Considerable variation in stilbene inducibility was identified in V. vinifera ssp. vinifera. Genotypic differences in abundance and profiles of stilbenes that are induced in response to grapevine pulp are shown. Two clusters of stilbene "chemovars" emerged; one cluster showed quick and strong accumulation of stilbenes, almost exclusively in the form of non-glycosylated resveratrol and viniferin while the second cluster accumulated fewer stilbenes and relatively high proportions of piceannol and the glycosylated piceid for all 86 genotypes, a time dependence of the stilbene pattern was observed: piceid, p-resveratrol, and piceannol accumulated earlier, whereas the viniferins were found later. It was further observed that the two clusters were preceded by differential accumulation of the transcripts for chalcone synthase (CHS) and stilbene-related genes: phenylalanine ammonium lyase (PAL), a stilbene synthase (StSy), and resveratrol synthase (RS). A screen of the population with respect to susceptibility to Vitis rotifoliorum mildew of grapevine (Vitis rotifoliorum) revealed considerable variability. The subpopulation of genotypes with high stilbene inducibility was significantly less susceptible as compared with low-stilbene genotypes, and for representative genotypes it could be shown that the inducibility of stilbene synthase B BY correlated with the inducibility by the pathogen.

```

```

viniferina esp., stylysteris). Genotypic differences in abundance and profiles of stilbenes that are induced in re-
sponse to a DVC-pulse were shown. Two clusters of stilbene ‘chemovars’ emerged: one cluster showed quick and
strong accumulation of stilbenes, almost exclusively in the form of non-glycosylated resveratrol and viniferin,
while the second cluster accumulated fewer stilbenes and relatively high proportions of piceatannol and the g
lycosylated piceid. For all 55 genotypes, a time dependence of the stilbene pattern was observed: piceid, r
esveratrol, and piceatannol accumulated earlier, whereas the viniferins were found later. It was further obs
erved that the genotypic differences in stilbene accumulation were preceded by differential accumulation of the
transcripts for chalcone synthase (CHS) and stilbene-related genes: phenylalanine ammonium lyase (PAL), a
stilbene synthase (STS), and resveratrol synthase (RS). A screen of the population with respect to suscep
tibility to downy mildew of grapevine (Plasmopara viticola) revealed considerable variability. The subsuppo
sition of genotypes with high stilbene inducibility was significantly less susceptible as compared with low-stilb
ene genotypes, and for representative genotypes it could be shown that the inducibility of stilbene synthase b
y DVC correlated with the inducibility by the pathogen.

(424, 1)
(182, 1)
(182, 1)
(424, 1)
(182, 1)

dl_config = DLConfig(model_name=model_name, seed_value=seed_value)

dl_config.stop_words = None
dl_config.lower = False #####
dl_config.remove_punctuation = False
dl_config.split_by_hyphen = False
dl_config.lemmatization = False #####
dl_config.stems = False

#parameters
dl_config.paddning = 'post'
dl_config.truncating = 'post'

dl_config.epochs = 2 # recommended number of epochs: 2, 3, 4
dl_config.batch_size = 16 # recommended batch-size: 16 or 32 # 8, 16, 32, 64, 128
dl_config.learning_rate = 3e-5 # recommended learning rate for Adam: 3e-5, 3e-5, 2e-5 # 3e-4, 1e-4,
dl_config.max_sent_len = 512 #sentences will have a maximum of "max_sent_len" words
dl_config.nm_sentences = 1 #1 or 2
dl_config.validation_percentage = 30

dl_config.keras_callbacks = True

if dl_config.keras_callbacks:
    checkpoint_path = str(dl_config.model_id_path) + "/model.{epoch:02d}.h5"
    keras_callbacks = [
        ModelCheckpoint(checkpoint_path,
            verbose=0,
            save_best_only=False,
            save_weights_only=True,
            monitor='epoch')
    ]

else:

```

```

keras_callbacks=None
):
    dl_config.tokenizer = BertTokenizer.from_pretrained('biobert_v1.1_pubmed', do_lower_case=False)

    x_train, y_train, x_val, y_val = Bert_preprocessing(X_train_df, y_train_df,
                                                         dl_config,
                                                         nmc_sentences = dl_config.nmc_sentences,
                                                         validation_percentage = dl_config.validation_percentage,
                                                         seed_valuedl_config.seed_value)

    Training set with 297 samples
    Validation set with 127 samples

biobert_path = './biobert_v1.1_pubmed'
if multiple_cpus:
    with strategy.scope():
        #model = Bert_FT(dl_config, learning_rate=dl_config.learning_rate, bert_name_or_path=biobert_path, bert_model = Bert_LSTM(dl_config, learning_rate=dl_config.learning_rate,static_bert=False, bert_name_or_path=biobert_path))
else:
    #model = Bert_FT(dl_config, learning_rate=dl_config.learning_rate, bert_name_or_path=biobert_path, bert_model = Bert_LSTM(dl_config, learning_rate=dl_config.learning_rate,static_bert=False, bert_name_or_path=biobert_path))

history = model.fit(x_train, y_train,
                    epochs=dl_config.epochs,
                    batch_size=dl_config.batch_size,
                    validation_data=(x_val, y_val),
                    callbacks=keras_callbacks)

```


Some weights of the PyTorch model were not used when initializing the TF 2.0 model TFBertModel: {'cls.predictions.transform.dense.bias', 'cls.predictions.transform.LayerNorm.bias', 'bert.embeddings.position_ids', 'cls.seq_relationship.weight', 'cls.predictions.decoder.weight', 'cls.predictions.transform.LayerNorm.weight', 'cls.seq_relationship.bias', 'cls.predictions.bias', 'cls.predictions.decoder.bias', 'cls.predictions.transform.dense.weight'}

- This is expected if you are initializing TFBertModel from a PyTorch model trained on another task or with another architecture (e.g. initializing a TFBertForSequenceClassification model from a BertForPreTraining model).

- This is NOT expected if you are initializing TFBertModel from a PyTorch model that you expect to be exactly identical (e.g. initializing a TFBertForSequenceClassification model from a BertForSequenceClassification model).

All the weights of TFBertModel were initialized from the Pytorch model.

If your task is similar to the task the model of the checkpoint was trained on, you can already use TFBertModel for predictions without further training.

WARNING:tensorflow: The parameters 'output_attentions', 'output_hidden_states' and 'use_cache' cannot be updated when calling a model. They have to be set to True/False in the config object (i.e.: 'config=XConfig.from_pretrained('name', output_attentions=True)').

WARNING:tensorflow: AutoGraph could not transform <bound method Socket.send of <mq.sugar.socket.Socket object at 0x0000237BB96A340> and will run it as-is.

Please report this to the TensorFlow team. When filing the bug, set the verbosity to 10 (on Linux, 'export AUTOGRAPH_VERBOSITY=10') and attach the full output.

Cause: module, class, method, function, traceback, frame, or code object was expected, got cython_function_or_method

To silence this warning, decorate the function with @tf.autograph.experimental.do_not_convert

WARNING: AutoGraph could not transform <bound method Socket.send of <mq.sugar.socket.Socket object at 0x0000237BB96A340> and will run it as-is.

Please report this to the TensorFlow team. When filing the bug, set the verbosity to 10 (on Linux, 'export AUTOGRAPH_VERBOSITY=10') and attach the full output.

Cause: module, class, method, function, traceback, frame, or code object was expected, got cython_function_or_method

To silence this warning, decorate the function with @tf.autograph.experimental.do_not_convert

WARNING:tensorflow: The parameter 'return_dict' cannot be set in graph mode and will always be set to 'True'. Model: "Model"

Layer (type)	Output Shape	Param #	Connected to
input_idx (InputLayer)	[(None, 512)]	0	
input_masks (InputLayer)	[(None, 512)]	0	
input_segments (InputLayer)	[(None, 512)]	0	
tf_bert_model (TFBertModel)	TFBaseModelOutputWith 108310272		input_idx[0][0] input_masks[0][0] input_segments[0][0]
bidirectional (Bidirectional)	(None, 512, 100)	327600	tf_bert_model[0][0]
dropout_37 (Dropout)	(None, 62)	0	bidirectional[0][0]
global_average_poolingid (Globa	(None, 100)	0	dropout_37[0][0]
dense (Dense)	(None, 62)	6262	global_average_poolingid[0][0]
dropout_38 (Dropout)	(None, 62)	0	dense[0][0]
dense_1 (Dense)	(None, 1)	63	dropout_38[0][0]

Epoch 1/2
WARNING:tensorflow: The parameters 'output_attentions', 'output_hidden_states' and 'use_cache' cannot be updated when calling a model. They have to be set to True/False in the Config object (i.e.: 'config=XConfig.from_pretrained('name', output_attentions=True)').
WARNING:tensorflow: The parameter 'return_dict' cannot be set in graph mode and will always be set to 'True'.
WARNING:tensorflow: AutoGraph could not transform <bound method Socket.send of <mq.sugar.socket.Socket object at 0x0000237BB96A340> and will run it as-is.
Please report this to the TensorFlow team. When filing the bug, set the verbosity to 10 (on Linux, 'export AUTOGRAPH_VERBOSITY=10') and attach the full output.
Cause: module, class, method, function, traceback, frame, or code object was expected, got cython_function_or_method
To silence this warning, decorate the function with @tf.autograph.experimental.do_not_convert
WARNING:tensorflow: The parameter 'return_dict' cannot be set in graph mode and will always be set to 'True'.
Model: "Model"

```
In [15]: plot_training_history(history_dict = history, dl_config=dl_config)

x_test, y_test = Bert_preprocessing(X_test_df, y_test_df, dl_config,
                                  nmr_sentences=dl_config.nmr_sentences)

yhat_probs = model.predict(x_test, verbose=0)
yhat_probs = yhat_probs[:, 0]

yhat_classes = np.where(yhat_probs > 0.5, 1, yhat_probs)
yhat_classes = np.where(yhat_classes < 0.5, 0, yhat_classes).astype(np.int64)
```

```
In [16]: dl_config.test_roc_auc, dl_config.test_pr_auc = plot_roc_pr_curves(y_test, yhat_probs, dl_config=dl_config)

# ROC AUC
print('ROC AUC: %f' % dl_config.test_roc_auc)

# avg precision
dl_config.test_avg_prec = average_precision(y_test_df, yhat_probs)
print('Average Precision: %f' % dl_config.test_avg_prec)

# accuracy
dl_config.test_acc = accuracy_score(y_test, yhat_classes)
print('Accuracy: %f' % dl_config.test_acc)

# precision tp / (tp + fp)
dl_config.test_prec = precision_score(y_test, yhat_classes)
print('Precision: %f' % dl_config.test_prec)

# recall: tp / (tp + fn)
dl_config.test_recall = recall_score(y_test, yhat_classes)
print('Recall: %f' % dl_config.test_recall)

# f1: 2 tp / (2 tp + fp + fn)
dl_config.test_f1_score = f1_score(y_test, yhat_classes)
print('F1 score: %f' % dl_config.test_f1_score)

# confusion matrix
matrix = confusion_matrix(y_test, yhat_classes)
matrix = np.flip(matrix)
print('Confusion Matrix:\n %s \n' % matrix)
```

Supervised Machine Learning

TF-IDF

```
In [7]: """Imports"""

from sklearn.feature_extraction.text import TfidfVectorizer
from sklearn.model_selection import train_test_split, cross_validate
import pandas as pd

from dl_config import DLConfig
from wrappers.pandas_wrapper import pandas_column_aslist, relevance_to_pandas, docs_to_pandasdocs
from web PubMed reader import pmids_to_docs
from nltk.corpus import stopwords
stops = set(stopwords.words("English"))

from sklearn.pipeline import make_pipeline
from sklearn.preprocessing import StandardScaler
from sklearn.svm import SVC
from sklearn.metrics import confusion_matrix, accuracy_score, precision_score, \
    precision_recall_curve, recall_score, average_precision_score, f1_score

vectorizer_w_stops = TfidfVectorizer(stop_words=stops)

dl_config = DLConfig()
```

```
In [2]: df = pd.read_excel("C:/Users/Z6 Freitas/Desktop/Mestrado/2ºSemestre/Projeto/code/datasets/dataset_final.xlsx",
                    index_col=0, usecols="A:F") # dataset que foi criado sem nenhum processamento
idfina1 = pandas_column_aslist(df, "Document")
docsfinal = pmids_to_docs(idfina1, "pq428728alunos.uminho.pt", dl_config)
docsfinal = docsfinal[0]
dataset_docs = docs_to_pandasdocs(docsfinal)
```

```
abst = pandas_column_aslist(df, "Abstract")
len(abst)
print(abst[0])
X = vectorizer_w_stops.fit_transform(abst)
print(X)
print(X.shape)
```

```
print(vectorizer_w_stops.get_feature_names())
print(len(vectorizer_w_stops.get_feature_names()))
analyzer = vectorizer_w_stops.build_analyzer()
print(analyzer(abst[0]))
```


[illegible]

