

INF-393/INF-578 MACHINE LEARNING. TAREA 2.

MÉTODOS LINEALES PARA CLASIFICACIÓN

Prof. Ricardo Ñanculef - jnancu@inf.utfsm.cl
Prof. Carlos Valle - cvalle@inf.utfsm.cl
Pablo Ibarra S - pablo.ibarras@alumnos.usm.cl
Francisco Mena - francisco.mena.13@sansano.usm.cl

Temas

- Implementación y evaluación de clasificadores en *sklearn*.
- Clasificadores bayesianos ingenuos (Bernoulli).
- LDA versus QDA. Reducción de dimensionalidad para clasificación.
- Regresión Logística. Selección de hiper-parámetros estructurales.
- SVMs Lineales. Selección de hiper-parámetros estructurales.
- Clasificadores k -NN. Selección de hiper-parámetros estructurales.
- Representación vectorial de texto.
- Procesamiento de lenguaje natural.
- Clases desbalanceadas

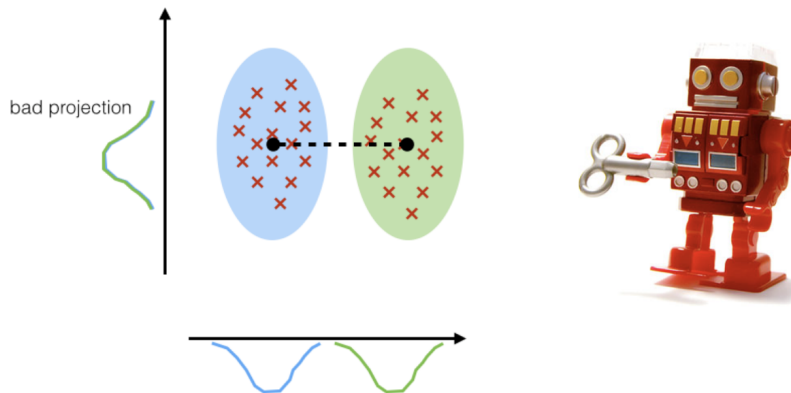
Formalidades

- Equipos de trabajo de: 2 personas.*.
- Se debe preparar un (breve) Jupyter/IPython notebook que explique la actividad realizada y las conclusiones del trabajo.
- En la clase del Viernes 17 Noviembre, el grupo (todos los integrantes) debe venir preparado para presentar cada punto de la tarea. Puede utilizar Jupiter, diapositivas, o cualquier material de apoyo.
- Se debe mantener un respaldo de cualquier tipo de código utilizado, informe y presentación en Github.
- Deadline: Viernes 17 Noviembre a las 22:00hrs.
- Formato de entrega: envío de link Github al correo electrónico del ayudante incluyendo a todos los profesores en copia y especificando asunto: [TallerXXX-CursoXXX-Semestre-Año].

*La modalidad de trabajo en equipo nos parece importante, porque en base a nuestra experiencia enriquece significativamente la experiencia de aprendizaje. Sin embargo, esperamos que esto no se transforme en una cruda división de tareas. Cada miembro del equipo debe estar en condiciones de realizar una presentación y discutir sobre cada punto del trabajo realizado.

1 Reducción de Dimensionalidad para Clasificación

Como hemos discutido en clases, la reducción de dimensionalidad es extremadamente importante en análisis de datos, no sólo porque nos permite visualizar y por lo tanto explorar los datos a nuestra disposición, si no que también permite reducir el costo computacional de procesarlos, sino porque reduce de modo significativo el riesgo de *overfitting*. En problemas de clasificación, la técnica que hemos discutido en el capítulo 1 (PCA), puede exhibir resultados diametralmente opuestos a aquello que se busca: preservar la información correspondiente a la clase a la que pertenece un ítem.



Para experimentar con este problema, en esta sección trabajaremos con una colección de sonidos fonéticos que deben ser identificados con vocales del inglés británico. Los datos han sido representados en un espacio de $d = 10$ características. Existen 528 datos de entrenamiento y 462 de pruebas, que pueden ser descargados desde el sitio web mantenido por los autores de nuestro texto guía. El mejor desempeño reportado por los autores corresponde a un 56% de *accuracy*, y es alcanzado por un modelo de vecinos más cercanos y una red neuronal artificial de radio basal.

- (a) Construya un dataframe con los datos a analizar descargando los datos desde la URL. Determine cuántos registros contiene el conjunto de entrenamiento y cuántos el conjunto de pruebas. Determine además el número promedio de palabras por ítem en cada clase.

```
1 import urllib
2 import pandas as pd
3 train_data_url = "http://statweb.stanford.edu/~hastie/ElemStatLearn/datasets/vowel.train"
4 test_data_url = "http://statweb.stanford.edu/~hastie/ElemStatLearn/datasets/vowel.test"
5 train_data_f = urllib.urlretrieve(train_data_url, "train_data.csv")
6 test_data_f = urllib.urlretrieve(test_data_url, "test_data.csv")
7 train_df = pd.DataFrame.from_csv('train_data.csv', header=0, index_col=0)
8 test_df = pd.DataFrame.from_csv('test_data.csv', header=0, index_col=0)
9 train_df.head()
10 test_df.tail()
```

- (b) Construya matrices X e y que contengan las características y las etiquetas correspondientes a los datos de entrenamiento y pruebas. Normalice apropiadamente los datos antes de empezar a trabajar.

```
1 from sklearn.preprocessing import StandardScaler
2 X = train_df.loc[:, 'x.1': 'x.10'].values
3 y = train_df.loc[:, 'y'].values
4 Scaler = StandardScaler().fit(X)
5 X_std = Scaler.transform(X)
```

- (c) Utilizando PCA genere una representación en 2 dimensiones de la data original (10 dimensiones) identificando cada clase con un color distinto (elija una paleta apropiada).

```

1  from sklearn.decomposition import PCA
2  from matplotlib import pyplot as plt
3  import seaborn as sns
4  import numpy as np
5  sklearn_pca = PCA(n_components=2)
6  Xred_pca = sklearn_pca.fit_transform(X_std)
7  cmap = plt.cm.get_cmap('ChooseAnAppropriatePalette')
8  mclasses=(1,2,3,4,5,6,7,8,9)
9  mcolors = [cmap(i) for i in np.linspace(0,1,10)]
10 plt.figure(figsize=(12, 8))
11 for lab, col in zip(mclasses,mcolors):
12     plt.scatter(Xred_pca[y==lab, 0],Xred_pca[y==lab, 1],label=lab,c=col)
13 plt.xlabel('Principal Component 1')
14 plt.ylabel('Principal Component 2')
15 leg = plt.legend(loc='upper right', fancybox=True)
16 plt.show()

```

- (d) Utilizando LDA genere una representación en 2 dimensiones de la data original (10 dimensiones) identificando cada clase con un color distinto (elija una paleta apropiada).

```

1  from sklearn lda import LDA
2  sklearn_lda = LDA(n_components=2)
3  Xred_lda = sklearn_lda.fit_transform(X_std,y)
4  cmap = plt.cm.get_cmap('ChooseAnAppropriatePalette')
5  mclasses=(1,2,3,4,5,6,7,8,9)
6  mcolors = [cmap(i) for i in np.linspace(0,1,10)]
7  plt.figure(figsize=(12, 8))
8  for lab, col in zip(mclasses,mcolors):
9      plt.scatter(Xred_lda[y==lab, 0],Xred_lda[y==lab, 1],label=lab,
10                  c=col)
11 plt.xlabel('LDA/Fisher Direction 1')
12 plt.ylabel('LDA/Fisher Direction 2')
13 leg = plt.legend(loc='upper right', fancybox=True)
14 plt.show()

```

- (e) Compare cualitativamente los resultados obtenidos en c y d.
- (f) Construya un clasificador que determine la clase de un dato x aleatoriamente sin considerar las características sino que solamente la probabilidad *a-priori* de cada clase. Por ejemplo, si la clase $y = 0$ ocurre el 25% de las veces, su clasificador debe predecir esta clase para un determinado x con probabilidad 0.25, independiente de los atributos de x .
- (g) Compare el desempeño de LDA, QDA y un modelo de *Vecinos Más Cercanos* (k -NN)[†] sin reducir dimensionalidad. ¿Qué técnica se comporta mejor sobre el conjunto de entrenamiento? ¿Sobre el conjunto de pruebas? Describa, utilizando un gráfico, el efecto de cambiar el parámetro de k en el tercer modelo.

```

1  from sklearn lda import LDA
2  from sklearn qda import QDA
3  from sklearn.neighbors import KNeighborsClassifier
4  import matplotlib.pyplot as plt

```

[†]Este clasificador busca los k datos de entrenamiento más similares al patrón que se quiere clasificar y predice la clase más popular entre éstos

```

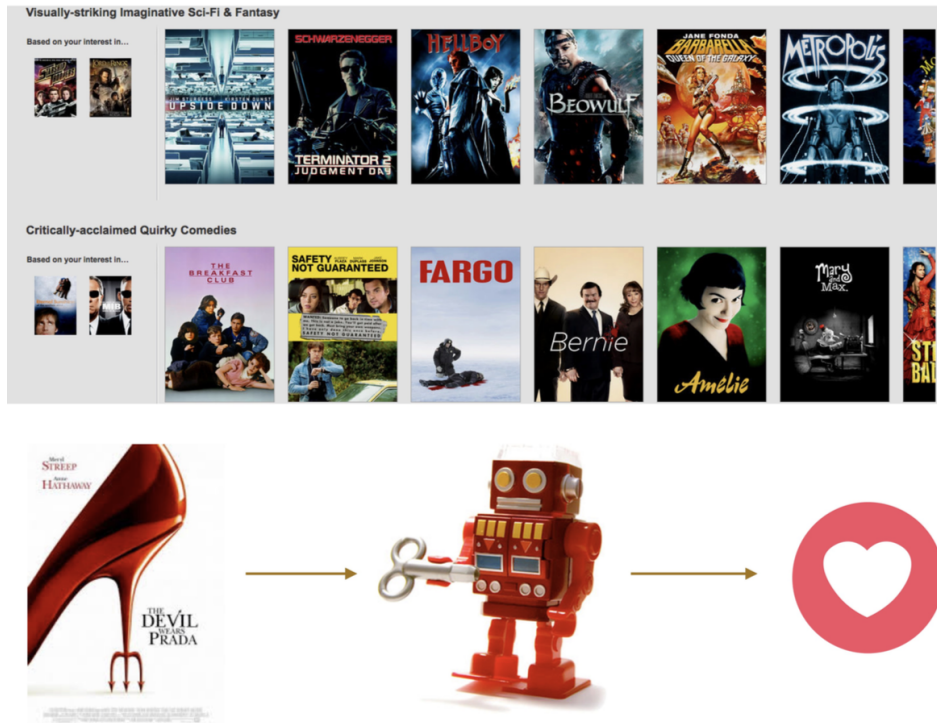
5  Xtest = test_df.loc[:, 'x.1': 'x.10'].values
6  ytest = test_df.loc[:, 'y'].values
7  X_std_test = Scaler.transform(Xtest)
8  lda_model = LDA()
9  lda_model.fit(X_std, y)
10 print "-----LDA-----"
11 print lda_model.score(X_std, y)
12 print lda_model.score(X_std_test, ytest)
13 qda_model = QDA()
14 qda_model.fit(X_std, y)
15 print "-----QDA-----"
16 print qda_model.score(X_std, y)
17 print qda_model.score(X_std_test, ytest)
18 print "-----KNN-----"
19 knn_scores_training = []
20 knn_scores_testing = []
21 k = range(1, 51)
22 for c in k:
23     knn_model = KNeighborsClassifier(n_neighbors=c)
24     knn_model.fit(X_std, y)
25     knn_scores_training.append(knn_model.score(X_std, y))
26     knn_scores_testing.append(knn_model.score(X_std_test, ytest))
27 plt.figure(figsize=(10, 10))
28 plt.subplot(211)
29 plt.xlabel("k")
30 plt.ylabel("Score")
31 plt.title("Training/Testing")
32 plt.plot(k, knn_scores_training, '-o', k, knn_scores_testing, '-o')
33 plt.legend(("Training", "Testing"), loc = "upper right")
34 plt.show()

```

- (h) Utilice PCA para generar una representación de la data en $d' = 1, 2, 3, \dots, 10$ dimensiones. Para cada caso entrene un modelo LDA, QDA y de k-NN. Construya un gráfico que muestre cómo evoluciona el error de entrenamiento versus d' . Sobreponga a este gráfico el error de pruebas versus d' . Concluya.
- (i) Utilice LDA para generar una representación de la data en $d' = 1, 2, 3, \dots, 10$ dimensiones. Para cada caso entrene un modelo LDA, QDA y de k-NN. Construya un gráfico que muestre cómo evoluciona el error de entrenamiento versus d' . Sobreponga a este gráfico el error de pruebas versus d' . Concluya.

2 Análisis de Opiniones sobre Películas

El análisis de sentimiento (o minería de opiniones) se refiere al proceso de extraer información acerca de la actitud que una persona (o grupo de ellas) manifiesta, en un determinado medio o formato digital, con respecto a un tópico o contexto de comunicación. Uno de los casos más estudiados corresponde a determinar la *polaridad* de un trozo de texto, es decir, clasificar una determinada evaluación escrita (ó *review*), en que una persona manifiesta una opinión, como *positiva*, *negativa* o *neutral*. La dificultad de este problema radica en el carácter altamente ambiguo e informal del lenguaje que utilizan naturalmente las personas así como el manejo de negaciones, sarcasmo y abreviaciones en una frase.



Los datos que usaremos para esta actividad corresponden a un subconjunto de los datos publicados en *Kaggle*, en el contexto de una competencia organizada por la Universidad de Stanford [?]. Cada registro disponible corresponderá a una opinión sobre una película, registrada sobre el sitio *Rotten Tomatoes*. Para empezar nos limitaremos a estudiar textos anotados como positivos o negativos, clases que codificaremos como +1 y 0 respectivamente. Para construir un clasificador que determine automáticamente la polaridad de un trozo de texto, vamos a necesitar representar los textos $\{d_i\}_{i=1}^n$ disponibles como vectores de características (features). El tipo de características más utilizado consiste en contar cuántas veces aparecen ciertos términos/palabras en el texto. Para esto, necesitaremos un *vocabulario* que, para esta actividad, construiremos mediante la unión de todas las palabras que observemos en los textos que tenemos a disposición. Para aumentar la eficacia de las características extraídas es conveniente ejecutar algunas técnicas de preprocesamiento básicas como: pasar todo el texto a minúsculas (lower-casing), eliminar signos de puntuación y eliminar palabras sin significado como artículos, pronombres y preposiciones (stop word removal [4]). Otra técnica que suele ser útil para obtener buenas características (features) es la lematización [6], es decir la reducción de todas las palabras a su tronco léxico base. Una técnica similar y más utilizada en la práctica es el *stemming* [5].

- Construya un dataframe con los datos a analizar descargando los datos desde la URL local. Determine cuántos registros de cada clase contiene el conjunto de entrenamiento y cuántos el conjunto de pruebas.

```

1 import urllib
2 import pandas as pd
3 train_data_url = "http://www.inf.utfsml.cl/~jnancu/stanford-subset/polarity.train"
4 test_data_url = "http://www.inf.utfsml.cl/~jnancu/stanford-subset/polarity.dev"
5 train_data_f = urllib.urlretrieve(train_data_url, "train_data.csv")
6 test_data_f = urllib.urlretrieve(test_data_url, "test_data.csv")
7 ftr = open("train_data.csv", "r")
8 fts = open("test_data.csv", "r")
9 rows = [line.split(" ",1) for line in ftr.readlines()]
10 train_df = pd.DataFrame(rows, columns=['Sentiment','Text'])
11 train_df['Sentiment'] = pd.to_numeric(train_df['Sentiment'])
12 rows = [line.split(" ",1) for line in fts.readlines()]
13 test_df = pd.DataFrame(rows, columns=['Sentiment','Text'])
14 test_df['Sentiment'] = pd.to_numeric(test_df['Sentiment'])
15 print train_df.shape
16 print test_df.shape

```

- (b) Construya una función, denominada *word_extractor*, que devuelva una lista de las palabras contenidas en un determinado un trozo de texto. Incorpore en su función las operaciones de lower-casing y stemming. Pruebe la función con las frases sugeridas en el código, invente otras similares y comente. Compare con los resultados obtenidos si no se hace stemming.

```

1 import re, time
2 from nltk.corpus import stopwords
3 from nltk import WordNetLemmatizer, word_tokenize
4 from nltk.stem.porter import PorterStemmer
5 def word_extractor(text):
6     ps = PorterStemmer()
7     commonwords = stopwords.words('english')
8     commonwords.remove("not")
9     text = re.sub(r'([a-z])\1+', r'\1\1',text)
10    words = ""
11    wordtokens = [ ps.stem(word.lower()) \
12                  for word in word_tokenize(text.decode('utf-8', 'ignore')) ]
13    for word in wordtokens:
14        if word not in commonwords:
15            words+=" "+word
16    return words
17 print word_extractor("I love to eat cake")
18 print word_extractor("I love eating cake")
19 print word_extractor("I loved eating the cake")
20 print word_extractor("I do not love eating cake")
21 print word_extractor("I don't love eating cake")

```

- (c) Construya una función, denominada *word_extractor2*, análoga a la función anterior, pero que lematice las palabras en vez de hacer stemming. Pruebe la función con las frases sugeridas en el código anterior y discuta las diferencias que observa.

```

1 def word_extractor2(text):
2     wordlemmatizer = WordNetLemmatizer()
3     commonwords = stopwords.words('english')
4     commonwords.remove("not")
5     text = re.sub(r'([a-z])\1+', r'\1\1',text)
6     words = ""
7     wordtokens = [ wordlemmatizer.lemmatize(word.lower()) \
8                   for word in word_tokenize(text.decode('utf-8','ignore')) ]

```

```

9         for word in wordtokens:
10             if word not in commonwords:
11                 words+=" "+word
12     return words

```

- (d) Utilizando la función *CountVectorizer* de la librería *sklearn* y de acuerdo a las directrices mencionadas en la introducción, genere una representación vectorial del texto de entrenamiento y del conjunto que usaremos para realizar pruebas. Explore el vocabulario utilizado y determine cuáles son las palabras más frecuentes en el conjunto de entrenamiento y pruebas.

```

1  import numpy as np
2  from sklearn.feature_extraction.text import CountVectorizer
3  texts_train = [word_extractor2(text) for text in train_df.Text]
4  texts_test = [word_extractor2(text) for text in test_df.Text]
5  vectorizer = CountVectorizer(ngram_range=(1, 1), binary='False')
6  vectorizer.fit(np.asarray(texts_train))
7  features_train = vectorizer.transform(texts_train)
8  features_test = vectorizer.transform(texts_test)
9  labels_train = np.asarray((train_df.Sentiment.astype(float)+1)/2.0)
10 labels_test = np.asarray((test_df.Sentiment.astype(float)+1)/2.0)
11 vocab = vectorizer.get_feature_names()
12 dist=list(np.array(features_train.sum(axis=0)).reshape(-1,))
13 for tag, count in zip(vocab, dist):
14     print count, tag

```

- (e) Construya una función que evalúe el desempeño obtenido por un clasificador genérico en el conjunto de entrenamiento y en el conjunto de pruebas. Utilice y explique las métricas que calcula la función *classification_report* de la librería *sklearn*.

```

1  from sklearn.metrics import classification_report
2  def score_the_model(model,x,y,xt,yt,text):
3      acc_tr = model.score(x,y)
4      acc_test = model.score(xt,yt)
5      print "Training Accuracy %s: %f"%(text,acc_tr)
6      print "Test Accuracy %s: %f"%(text,acc_test)
7      print "Detailed Analysis Testing Results ..."
8      print(classification_report(yt, model.predict(xt), target_names=['+', '-']))

```

- (f) Construya una función que entrene/ajuste un clasificador *Bayesiano Ingenuo (Binario)* (las características no nulas serán tratadas como 1) y mida el error de predicción obtenido sobre los datos de entrenamiento y pruebas. Utilice esta función con las características extraídas en el punto (d). Mida el efecto de filtrar *stopwords* y de eliminar este paso de pre-procesamiento típico. Determine además, qué representación obtiene un mejor resultado: si aquella obtenida vía lematización o aquella obtenida vía stemming. Finalmente, tome un subconjunto aleatorio de los textos de prueba y analice las predicciones del modelo (explore las predicciones, así como las probabilidades que el clasificador asigna a cada clase).

```

1  from sklearn.naive_bayes import BernoulliNB
2  import random
3  def do_NAIVE_BAYES(x,y,xt,yt):
4      model = BernoulliNB()
5      model = model.fit(x, y)
6      score_the_model(model,x,y,xt,yt,"BernoulliNB")
7      return model
8  model=do_NAIVE_BAYES(features_train,labels_train,features_test,labels_test)
9  test_pred = model.predict_proba(features_test)
10 spl = random.sample(xrange(len(test_pred)), 15)

```

```

11 for text, sentiment in zip(test_df.Text[spl], test_pred[spl]):
12     print sentiment, text

```

- (g) Construya una función que entrene/ajuste una *Máquina de Vectores de Soporte (SVM) Lineal* y mida el error de predicción obtenido sobre los datos de entrenamiento y pruebas. Incluya en su función la exploración de diferentes valores del parámetro de regularización C . Discuta el significado y efecto esperado de este parámetro. Utilice la función construida con los atributos extraídos en el punto (d). Mida el efecto de filtrar *stopwords* y de eliminar este paso de pre-procesamiento típico. Determine además, qué representación obtiene un mejor resultado: si aquella obtenida vía lematización o aquella obtenida vía stemming. Finalmente, tome un subconjunto aleatorio de los textos de prueba y analice las predicciones del modelo (explore las predicciones, así como las probabilidades que el clasificador asigna a cada clase).

```

1 from sklearn.svm import LinearSVC
2 def do_SVM(x,y,xt,yt):
3     Cs = [0.01,0.1,10,100,1000]
4     for C in Cs:
5         print "C Value: %f"%C
6         model = LinearSVC(C=C)
7         model = model.fit(x, y)
8         score_the_model(model,x,y,xt,yt,"SVM")
9 do_SVM(features_train,labels_train,features_test,labels_test)

```

3 Fraude en transacciones bancarias

El mundo está más globalizado, servicios que antes solo se podían hacer de manera presencial hoy en día están disponibles de forma electrónica, el más claro ejemplo de esto son las transacciones bancarias. La ventaja de este cambio es que cualquier persona en el mundo puede enviar dinero de un lado a otro en cosa de segundos, por consecuencia, la cantidad de transacciones aumentó mucho. La desventaja es que es más fácil cometer fraudes, debido a que las únicas barreras de protección que tenemos ante estas malas prácticas son códigos de seguridad (Clave secreta, matriz de números, etc.) que si caen en manos de gente deshonesto lo más probable es que se cometa algún tipo fraude.

En esta sección utilizaremos un dataset de *Kaggle* [3] que contiene transacciones bancarias, el desafío está en que sus clases están desbalanceadas, por lo tanto, utilizaremos 3 técnicas para hacer frente a este problema. El *dataset* contiene transacciones hechas por tarjetas de créditos europeas durante un periodo de 2 días, donde 492 de 284807 fueron fraudulentas. Cada dato tiene como atributo 28 componentes principales obtenidos por PCA (Esto se hace para enmascarar datos sensibles), una variable tiempo y el monto de la transacción.

- (a) Descargue y cargue el archivo, genere un gráfico que compare la cantidad de elementos que hay por clase.

```

1 import pandas as pd
2 import seaborn as sns
3 import matplotlib.pyplot as plt
4 import numpy as np
5 dt = pd.read_csv("creditcard.csv",header = 0)
6 sns.countplot("Class",data=dt)
7 plt.show()

```

- (b) Estandarice el monto de la transacción. ¿Por qué solo estandarizamos un atributo?

```

1 from sklearn.preprocessing import StandardScaler
2 dt["nAmount"] = StandardScaler().fit_transform(dt['Amount'].values.reshape(-1, 1))
3 dt.drop(["Time", "Amount"],axis=1,inplace=True)
4 dt.head()

```


- (c) Defina una función que reciba como *input* un modelo, los datos de entrenamiento y datos de testeo. Esta función deberá mostrar luego de generar el modelo y haber probado la data de testing una matriz de confusión, junto a sus respectivas métricas. Deberá además mostrar una curva ROC junto al valor del área bajo la curva de esta. Finalmente, responda las siguientes preguntas ¿Qué es una matriz de confusión? ¿Que métricas se pueden calcular de esta y cuales nos serán útiles en este experimento? ¿Que es una curva ROC y qué relación tiene con la matriz de confusión?

```

1 def model(model,features_train,features_test,labels_train,labels_test):
2     clf = model
3     clf.fit(features_train,labels_train.values.ravel())
4     pred=clf.predict(features_test)
5     cnf_matrix=confusion_matrix(labels_test,pred)
6     print "The recall for this model is: %f"%(float(cnf_matrix[1,1])/(cnf_matrix[1,1]+
7                                                                                   cnf_matrix[1,0]))
8
9     fig= plt.figure(figsize=(6,3))
10    print "TP: %d"%cnf_matrix[1,1,]
11    print "TN: %d"%cnf_matrix[0,0]
12    print "FP: %d"%cnf_matrix[0,1]
13    print "FN: %d"%cnf_matrix[1,0]
14    sns.heatmap(cnf_matrix,cmap="coolwarm_r",annot=True,linewidths=0.5)
15    plt.title("Confusion_matrix")
16    plt.xlabel("Predicted_class")
17    plt.ylabel("Real class")
18    plt.show()
19    print "\n-----Classification Report-----"
20    print classification_report(labels_test,pred)
21
22    ''' You have to complete this function, ROC Curve is missing'''

```

- (d) Defina una función que genere un *Training* y *Testing Set* de un *Dataset* cualquiera. Indique que está haciendo la función.

```

1 from sklearn.model_selection import train_test_split
2 def data_preparation(x):
3     x_features= x.iloc[:,x.columns != "Class"]
4     x_labels=x.iloc[:,x.columns=="Class"]
5     x_train,x_test,y_train,y_test=train_test_split(x_features,x_labels,test_size=0.3)
6     print "Length of training data: %d"%len(x_train)
7     print "Length of test data: %d"%len(x_test)
8     return(x_train,x_test,y_train,y_test)

```

- (e) Defina una función que realice *Undersample* [8] . ¿Explique en que consiste esta técnica?

```

def undersample(data,times):
    fraud_indices= np.array(data[data.Class==1].index)
    normal_indices = np.array(data[data.Class==0].index)
    Count_Normal_transacation = len(data[data["Class"]==0])
    Count_Fraud_transacation = len(data[data["Class"]==1])
    Normal_indices_undersample = np.array(np.random.choice(normal_indices,
        (times*Count_Fraud_transacation),replace=False))
    undersample_data= np.concatenate([fraud_indices,Normal_indices_undersample])
    undersample_data = data.iloc[undersample_data,:]
    normal = (float(len(undersample_data[undersample_data.Class==0]))
        /len(undersample_data["Class"]))
    print "The normal transacation proportion is : %f"%normal

```

```

fraud = (float(len(undersample_data[undersample_data.Class==1]))
/len(undersample_data["Class"]))
print "The fraud transaction proportion is : %f"%fraud
return(undersample_data)

```

- (f) Genere un modelo de *Logistic Regression* y otro a su elección utilizando *UnderSample* [8] . Concluya a partir de los resultados.

```

from sklearn.linear_model import LogisticRegression
from sklearn.metrics import confusion_matrix,recall_score,precision_recall_curve, auc, roc_curve,
roc_auc_score, classification_report
x_tr,x_ts,y_tr,y_ts = data_preparation(dt)
x_tr["Class"] = y_tr["Class"]
x_tr = x_tr.reset_index(drop=True)
for i in range(1,4):
    print "The undersample data for " + str(i) + " proportion\n"
    Undersample_data = undersample(x_tr,i)
    print "\n-----Validation Set-----"
    print "\nThe model classification for " + str(i) + " proportion\n"
    under_x,under_xt,under_yx,under_yxt=data_preparation(Undersample_data)
    clf=LogisticRegression()
    model(clf,under_x,under_xt,under_yx,under_yxt)
    print "-----Testing Set-----"
    model(clf,under_x,x_ts,under_yx,y_ts)
    print "-----"

```

- (g) Vuelva a realizar los puntos (e) y (f), esta vez utilizando *OverSampling* [7] y *SMOTE* [8] (Hint: Para *SMOTE* usted puede apoyarse de la libreria *imblearn*)

```

from imblearn.over_sampling import SMOTE
os = SMOTE(random_state=0)
os_data_X,os_data_y=os.fit_sample(x_tr,y_tr)

```

References

- [1] Hastie, T.; Tibshirani, R., Friedman, J. (2009), The Elements of Statistical Learning, Second Edition. Springer New York Inc.
- [2] Joshi, M., Das, D., Gimpel, K., Smith, N. A. (2010). Movie reviews and revenues: An experiment in text regression. In the 2010 Annual Conference of the North American Chapter of the Association for Computational Linguistics (pp. 293-296). Association for Computational Linguistics.
- [3] <https://www.kaggle.com/dalpozz/creditcardfraud>
- [4] https://en.wikipedia.org/wiki/Stop_words
- [5] <https://en.wikipedia.org/wiki/Stemming>
- [6] <https://en.wikipedia.org/wiki/Lemmatisation>
- [7] https://en.wikipedia.org/wiki/Oversampling_and_undersampling_in_data_analysis
- [8] <https://www.irjet.net/archives/V4/i8/IRJET-V4I857.pdf>