



ANÁLISIS INTELIGENTE DE DATOS - TAREA III

Prof. Ricardo Ñanculef
jnancu@inf.utfsm.cl

Instrucciones

- El objetivo de estas tareas es que pongan en práctica los métodos discutidos en clases. El énfasis debe ser puesto por lo tanto en la comprensión del problema abordado, la implementación computacional de los algoritmos correspondientes, las decisiones prácticas que haya tenido que tomar para completar el análisis y la interpretación de los resultados.
- Lo ideal es que hagan esta tarea en parejas. Cada miembro del equipo debe estar en condiciones de exponer cada punto del trabajo realizado.
- Se recomienda el uso de Python como lenguaje de programación, junto a todo el ecosistema de librerías que hemos mencionado en clases, incluyendo: pandas, numpy, sklearn, statsmodels y seaborn.
- La realización de la tarea consiste en preparar un breve informe y presentación (slides) del trabajo realizado. Por favor, incluya en estos productos sólo aquello esencial para dar a entender lo realizado. Los programas utilizados contendrán todos los detalles. Como de costumbre, estos deben ser entregados y mantenidos en un repositorio público como github.
- La fecha entrega de todo el material entregable corresponde al día Jueves 14 de Julio.
- El formato de entrega será electrónico vía email a *jnancu@inf.utfsm.cl*.

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, y no sólo porque 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 [3]. 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/~tibs/ElemStatLearn/datasets/vowel.train"
4 test_data_url = "http://statweb.stanford.edu/~tibs/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.ix[:, 'x.1': 'x.10'].values
3 y = train_df.ix[:, 'y'].values
4 X_std = StandardScaler().fit_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 1 y 2. Proponga un método para elegir una técnica de reducción de dimensionalidad.
- (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  Xtest = test_df.ix[:, 'x.1': 'x.10'].values
5  ytest = test_df.ix[:, 'y'].values
6  X_std_test = StandardScaler().fit_transform(Xtest)
7  lda_model = LDA()
8  lda_model.fit(X_std, y)
9  print lda_model.score(X_std, y)
10 print lda_model.score(X_std_test, ytest)
11 qda_model = QDA()
12 knn_model = KNeighborsClassifier(n_neighbors=10)

```

- (g) 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.
- (h) 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 [6]. 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.

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

Para aumentar la eficacia de las características extraídas es conveniente ejecutar algunas técnicas de pre-procesamiento 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 [7]). Otra técnica que suele ser útil para obtener buenas características (features) es la lematización [9], 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* [8].

- (a) Construya un dataframe con los datos a analizar descargando los datos desde la URL local [5]. 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.utfsm.cl/~jnancu/stanford-subset/polarity.train"
4  test_data_url = "http://www.inf.utfsm.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 train_df.shape
16 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      wordlemmatizer = WordNetLemmatizer()
7      commonwords = stopwords.words('english')
8      text = re.sub(r'([a-z])\1+', r'\1\1',text)#substitute multiple letter by two
9      words = ""
10     wordtokens = [ wordlemmatizer.lemmatize(word.lower()) \
11                    for word in word_tokenize(text.decode('utf-8', 'ignore')) ]
12     for word in wordtokens:
13         if word not in commonwords:
14             words+=" "+word
15     return words
16 word_extractor("I love to eat cake")
17 word_extractor("I love eating cake")
18 word_extractor("I loved eating the cake")
19 word_extractor("I do not love eating cake")
20 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     text = re.sub(r'([a-z])\1+', r'\1\1',text)#substitute multiple letter by two
5     words = ""
6     wordtokens = [ wordlemmatizer.lemmatize(word.lower()) \
7                     for word in word_tokenize(text.decode('utf-8','ignore')) ]
8     for word in wordtokens:
9         if word not in commonwords:
10             words+=" "+word
11     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[: -1])
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, random
2 def do_NAIVE_BAYES(x,y,xt,yt):

```

```

3         model = BernoulliNB()
4         model = model.fit(x, y)
5         score_the_model(model,x,y,xt,yt,"BernoulliNB")
6         return model
7     model=do_NAIVE_BAYES(features_train,labels_train,features_test,labels_test)
8     test_pred = model.predict_proba(features_test)
9     spl = random.sample(xrange(len(test_pred)), 15)
10    for text, sentiment in zip(test_df.Text[spl], test_pred[spl]):
11        print sentiment, text

```

- (g) Construya una función que entrene/ajuste un clasificador *Bayesiano Ingenuo Multinomial* (las características enteras no nulas no serán reducidas a 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 MultinomialNB,random
2    def do_MULTINOMIAL(x,y,xt,yt):
3        model = MultinomialNB()
4        model = model.fit(x, y)
5        score_the_model(model,x,y,xt,yt,"MULTINOMIAL")
6        return model
7    model=do_MULTINOMIAL(features_train,labels_train,features_test,labels_test)
8    test_pred = model.predict_proba(features_test)
9    spl = random.sample(xrange(len(test_pred)), 15)
10   for text, sentiment in zip(test_df.Text[spl], test_pred[spl]):
11       print sentiment, text

```

- (h) Construya una función que entrene/ajuste un modelo de *Regresión Logística Regularizado* (utilizando, como penalizador, la norma ℓ_2) 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. Explique 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.linear_model import LogisticRegression
2    def do_LOGIT(x,y,xt,yt):
3        start_t = time.time()
4        Cs = [0.01,0.1,10,100,1000]
5        for C in Cs:
6            print "Usando C= %f"%C
7            model = LogisticRegression(penalty='l2',C=C)
8            model = model.fit(x, y)
9            score_the_model(model,x,y,xt,yt,"LOGISTIC")
10   do_LOGIT(features_train,labels_train,features_test,labels_test)

```

- (h) 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 "El valor de C que se esta probando: %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)

```

- (i) Construya un gráfico que permita comparar los resultados obtenidos por los diferentes métodos de clasificación.

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] <http://statweb.stanford.edu/~tibs/ElemStatLearn/datasets/>
- [4] <https://www.dropbox.com/sh/8r1wrblyfokwuq0/AABUEvgcuMxyZht2-KYyBptUa?dl=0>
- [5] <http://www.inf.utfsm.cl/~jnancu/stanford-subset/>
- [6] <https://www.kaggle.com/c/sentiment-analysis-on-movie-reviews>
- [7] https://en.wikipedia.org/wiki/Stop_words
- [8] <https://en.wikipedia.org/wiki/Stemming>
- [9] <https://en.wikipedia.org/wiki/Lemmatisation>