

---

## INF-393/INF-578 MACHINE LEARNING. TAREA 3.

### MÉTODOS NO-LINEALES.

---

Prof. Ricardo Ñanculef - [jnancu@inf.ut fsm.cl](mailto:jnancu@inf.ut fsm.cl)  
Prof. Carlos Valle - [cvalle@inf.ut fsm.cl](mailto:cvalle@inf.ut fsm.cl)  
Pablo Ibarra S - [pablo.ibarras@alumnos.usm.cl](mailto:pablo.ibarras@alumnos.usm.cl)  
Francisco Mena - [francisco.mena.13@sansano.usm.cl](mailto:francisco.mena.13@sansano.usm.cl)

### Temas

- Pre-procesamiento de datos y extracción de características.
- Árboles de clasificación en *sklearn*.
- Redes Neuronales feed-forward en *keras*.
- Métodos de *Kernel*
- Ensamblados de máquinas de aprendizaje en *sklearn*.
- Selección de hiper-parámetros estructurales vía validación cruzada.

### Formalidades

- Equipos de trabajo de: 2 personas.<sup>1</sup>
- Se debe preparar un (breve) Jupyter/IPython notebook que explique la actividad realizada y las conclusiones del trabajo.
- En la clase Recuperativa del Martes 19 de Diciembre, 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: Martes 19 Diciembre.
- Formato de entrega: envío de link Github al correo electrónico del ayudante <sup>2</sup> incluyendo a todos los profesores en copia y especificando asunto: [TallerXXX-CursoXXX-Semestre-Año]

---

<sup>1</sup>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.

<sup>2</sup>[francisco.mena.13@sansano.usm.cl](mailto:francisco.mena.13@sansano.usm.cl)

## Paquetes/Librerías

Como es usual utilizaremos *numpy*, *scipy*, *matplotlib* y *sklearn*. Además de éstas se necesitará instalar *keras*, una librería en python para prototipado rápido de modelos basados en redes neuronales, muy similar en espíritu a *sklearn*. La librería puede usar *TensorFlow* o *Theano* como backend, siendo éstas las librerías más populares para desarrollar nuevos modelos de redes neuronales o implementar eficientemente modelos conocidos con fines prácticos. Para detalles sobre la instalación puede revisar [1] o escribir un email a su ayudante.

### 1. Small Circle inside Large Circle

El objetivo de esta sección es experimentar con algunos modelos no-lineales sobre un problema de juguete generado para visualizar algoritmos de *clustering*. Se trata de un problema de clasificación a todas luces *linealmente inseparable*, en el sentido que, si denotamos por  $\mathbf{x} \in \mathbb{R}^2$  un patrón de entrada y por  $y \in \{-1, 1\}$  su correspondiente etiqueta, no existen  $\mathbf{w} \in \mathbb{R}^2, b \in \mathbb{R}$  tal que  $y(\mathbf{w}^T \mathbf{x} + b) \geq \rho > 0$ . El problema nos permite hacer un recorrido rápido por las grandes ideas en la búsqueda de la *no-linealidad*.

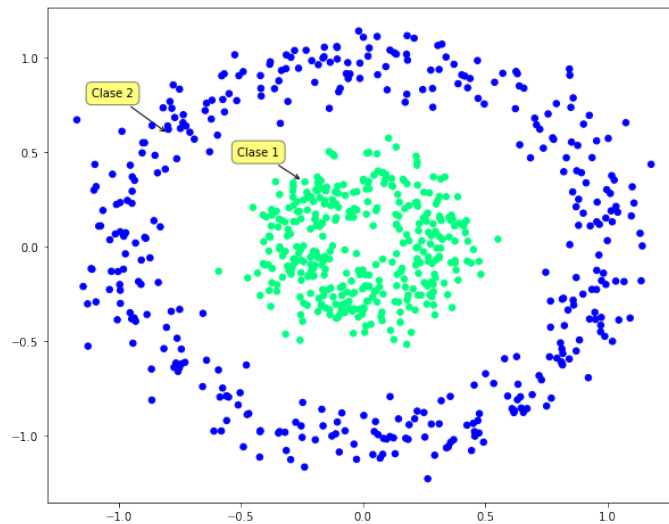


Figura 1: Distribución deseada para la actividad 1. Los 2 colores representan 2 clases distintas.

- (a) Escriba una función que genere (aleatoriamente)  $n$  datos etiquetados de la forma  $\{(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)\}$ ,  $\mathbf{x}_i \in \mathbb{R}^2, y_i \in \{0, 1\}$ , con una distribución de probabilidad que refleje la configuración linealmente inseparable que muestra la Fig. 1<sup>3</sup>. Utilice esta función para crear 1000 datos de entrenamiento y 1000 datos de pruebas. Para medir la tendencia de los modelos a *sobre-ajuste*, agregue un 5 % de ruido al dataset, generando  $\mathbf{x}$ 's cercanos a la frontera. Genere un gráfico que muestre datos de entrenamiento y pruebas, identificando cada clase con un color diferente (como lo muestra la Fig. 1).

```
1 import numpy as np
2 from sklearn.utils import check_random_state
3 from sklearn.model_selection import train_test_split
4 def do_circles(n=2000, noisy_n=0.05):
5     generator = check_random_state(10)
6     linspace = np.linspace(0, 2 * np.pi, n // 2 + 1)[:n]
7     outer_circ_x = np.cos(linspace)
8     outer_circ_y = np.sin(linspace)
```

<sup>3</sup>Puede generar datos aleatorios distribuidos de manera circular para luego etiquetar aquellos ubicados en el círculo interior como 1 y en el círculo exterior 0.

```

9     inner_circ_x = outer_circ_x * .3
10    inner_circ_y = outer_circ_y * .3
11    X = np.vstack((np.append(outer_circ_x, inner_circ_x),
12                    np.append(outer_circ_y, inner_circ_y))).T
13    y = np.hstack([np.zeros(n // 2, dtype=np.intp),
14                  np.ones(n // 2, dtype=np.intp)])
15    X += generator.normal(scale=noisy_n, size=X.shape)
16
17    X_train, X_test, y_train, y_test = train_test_split(X, y,
18                                                         test_size=0.5, random_state=42)
19    return X_train, y_train, X_test, y_test

```

Para lo que sigue de la actividad utilice la siguiente función para graficar las fronteras de clasificación en base a la probabilidad, definida por un algoritmo, de un ejemplo a pertenecer a una clase en particular.

```

1  import matplotlib.pyplot as plt
2  def plot_classifier(clf, X_train, Y_train, X_test, Y_test, model_type):
3      f, axis = plt.subplots(1, 1, sharex='col', sharey='row', figsize=(12, 8))
4      axis.scatter(X_train[:,0], X_train[:,1], s=30, c=Y_train, zorder=10, cmap='cool')
5      axis.scatter(X_test[:,0], X_test[:,1], s=20, c=Y_test, zorder=10, cmap='Greys')
6      XX, YY = np.mgrid[-2:2:200j, -2:2:200j]
7      if model_type == 'tree':
8          Z = clf.predict_proba(np.c_[XX.ravel(), YY.ravel()])[:,0]
9      elif model_type == 'ann':
10         Z = clf.predict(np.c_[XX.ravel(), YY.ravel()])
11     else: raise ValueError('model type not supported')
12     Z = Z.reshape(XX.shape)
13     Zplot = Z >= 0.5
14     axis.pcolormesh(XX, YY, Zplot, cmap='YlGn')
15     axis.contour(XX, YY, Z, alpha=1, colors=["k", "k", "k"], linestyles=["--", "-", "--"],
16                levels=[-2, 0, 2])
17     plt.show()

```

- (b) Demuestre experimentalmente que una red neuronal artificial correspondiente a 1 sola neurona (i.e. sin capas escondidas) no puede resolver satisfactoriamente el problema. Puede utilizar la función de activación y el método de entrenamiento que prefiera. Sea convincente: por ejemplo, intente modificar los parámetros de la máquina de aprendizaje, reportando métricas que permitan evaluar el desempeño del modelo en el problema con cada cambio efectuado. Adapte también la función *plot\_classifier* para que represente gráficamente la solución encontrada por la red neuronal. Describa y explique lo que observa, reportando gráficos de la solución sólo para algunos casos representativos.

```

1  from keras.models import Sequential
2  from keras.layers import Dense
3  from keras.optimizers import SGD
4  n_h=1
5  model = Sequential()
6  model.add(Dense(1, input_dim=X_train.shape[1], kernel_initializer='uniform', activation='relu'))
7  model.add(Dense(n_h, init='uniform', activation='sigmoid'))
8  model.compile(optimizer=SGD(lr=1), loss='binary_crossentropy', metrics=['accuracy'])
9  model.fit(X_train, Y_train, epochs=50, batch_size=100, verbose=1)
10 scores = model.evaluate(X_test, Y_test)
11 test_acc = scores[1]

```

- (c) Demuestre experimentalmente que una red neuronal artificial con 1 capa escondida puede resolver satisfactoriamente el problema obtenido en (a). Puede utilizar la arquitectura y el método de entrenamiento que prefiera, pero en esta actividad puede optar tranquilamente por usar los hiper-parámetros que se

entregan como referencia en el código de ejemplo. Cambie el número de neuronas  $N_h$  en la red entre 2 y 32 en potencias de 2, graficando el error de entrenamiento y pruebas como función de  $N_h$ . Describa y explique lo que observa. Utilice la función *plot\_classifier*, diseñada anteriormente, para construir gráficos de la solución en algunos casos representativos.

```

1  ...
2  n_h=32
3  model = Sequential()
4  model.add(Dense(n_h, input_dim=X_train.shape[1], kernel_initializer='uniform', activation='relu'))
5  model.add(Dense(1, kernel_initializer='uniform', activation='sigmoid'))

```

- (d) Demuestre experimentalmente que stump (árbol de clasificación de 1 nivel) no puede resolver satisfactoriamente el problema anterior. Puede utilizar el criterio y la función de partición que prefiera. Sea convincente: por ejemplo, intente modificar los parámetros de la máquina, reportando métricas que permitan evaluar el desempeño del modelo en el problema con cada cambio efectuado. Adapte también la función *plot\_classifier* para que represente gráficamente la solución encontrada por el árbol. Describa y explique lo que observa, reportando gráficos de la solución sólo para algunos casos representativos.

```

1  from sklearn.tree import DecisionTreeClassifier as Tree
2  clf=Tree(criterion='gini',splitter='best',random_state=0,max_depth=1)
3  clf.fit(X_train,Y_train)
4  acc_test = clf.score(X_test,Y_test)
5  print "Test Accuracy = %f"%acc_test
6  print clf.tree_.max_depth
7  plot_classifier(clf,X_train,Y_train,X_test,Y_test,'tree')

```

- (e) Demuestre experimentalmente que un árbol de clasificación de múltiples niveles puede resolver satisfactoriamente el problema estudiado. Puede utilizar el criterio y la función de partición que prefiera, pero puede optar tranquilamente por usar los hiper-parámetros que se entregan como referencia en el código de ejemplo. Cambie el número de niveles admitidos en el árbol  $N_t$  entre 2 y 20, graficando el error de entrenamiento y pruebas como función de  $N_t$ . Describa y explique lo que observa. Utilice la función *plot\_classifier*, diseñada anteriormente, para construir gráficos de la solución en algunos casos representativos.

```

1  ...
2  n_t=8
3  clf=Tree(criterion='gini',splitter='best',random_state=0,max_depth=n_t)
4  clf.fit(X_train,Y_train)

```

- (f) Como ya se demostró experimentalmente que este problema es linealmente inseparable, ahora se pide experimentar otra alternativa. Para ello deberá realizar una proyección de los datos a un nuevo espacio dimensional (*manifold*) en el cual se reconozcan sus patrones no lineales, para poder trabajarlos con fronteras lineales. Utilice la técnica de PCA con la ayuda de un *Kernel Gaussiano* ([2]) para extraer sus vectores con dimensión *infinita* de mayor varianza.

```

1  kpca = KernelPCA(n_components=2,kernel="rbf", gamma=5)
2  kpca = kpca.fit(X_train)
3  Xkpca_train = kpca.transform(X_train)
4  Xkpca_test = kpca.transform(X_test)

```

- (g) Ajuste un algoritmo de aprendizaje con fronteras lineal para los datos proyectados en este nuevo espacio que captura sus componentes no lineales, muestre gráficamente que el problema ahora puede ser resuelto con estos métodos. Reporte métricas para evaluar el desempeño, comente y concluya.

## 2. Bike Sharing: Predicción de Demanda Horaria

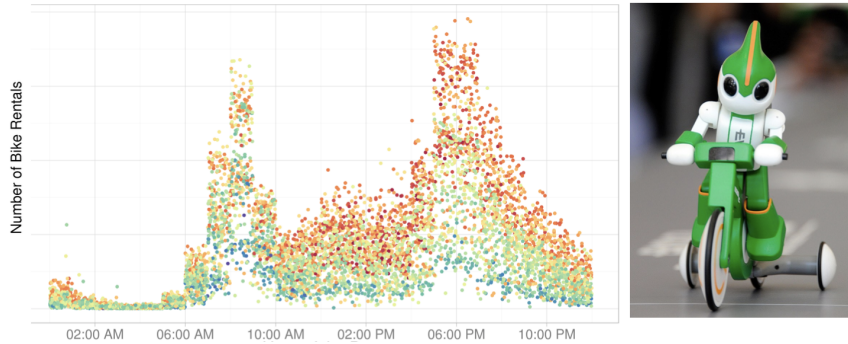
En esta sección simularemos nuestra participación en el desafío *Bike Sharing Demand* de *Kaggle* [3]. El objetivo es predecir la demanda de bicicletas sobre la red *Capital Bikeshare* de la ciudad de Washington, D.C., en función de la hora del día y otras variables descritas en la tabla 1. En principio, y como muestra la figura, la función es altamente no lineal y no determinista como función de la hora del día. Su objetivo será entrenar un modelo para obtener un puntaje correspondiente al top-100 del “leaderboard” final, es decir superior o igual a 0.37748. La función utilizada para evaluar este concurso Kaggle se proporciona en la siguiente ecuación:

$$E_{\text{bikes}}(\mathbf{y}, \hat{\mathbf{y}}) = \frac{1}{n} \sum_i (\ln(y_i + 1) - \ln(\hat{y}_i + 1))^2, \quad (1)$$

donde  $\mathbf{y}, \hat{\mathbf{y}} \in \mathbb{R}^n$  denotan los vectores de observaciones y predicciones respectivamente.

Como el dataset de pruebas original no está disponible se fabricará uno, correspondiente al 20 % de los datos de entrenamiento. Además, se pondrá a su disposición un subconjunto independiente de datos con propósitos de validación. Usted podrá descargar los archivos correspondientes al subconjunto de entrenamiento y pruebas a utilizar ejecutando los siguientes comandos:

```
1 wget http://octopus.inf.utfsm.cl/~ricky/bike_sharing_train.csv
2 wget http://octopus.inf.utfsm.cl/~ricky/bike_sharing_val.csv
3 wget http://octopus.inf.utfsm.cl/~ricky/bike_sharing_test.csv
```



| Atributo   | Descripción  |
|------------|--|
| datetime   | hourly date + timestamp  |
| season     | 1 = spring, 2 = summer, 3 = fall, 4 = winter   |
| holiday    | whether the day is considered a holiday  |
| workingday | whether the day is neither a weekend nor holiday   |
| weather    | 1: Clear, Few clouds, Partly cloudy, Partly cloudy<br>2: Mist + Cloudy, Mist + Broken clouds, Mist + Few clouds, Mist<br>3: Light Snow, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered clouds<br>4: Heavy Rain + Ice Pallets + Thunderstorm + Mist, Snow + Fog |
| temp       | temperature in Celsius   |
| atemp      | ?feels like? temperature in Celsius  |
| humidity   | relative humidity  |
| windspeed  | wind speed   |
| casual     | number of non-registered user rentals initiated  |
| registered | number of registered user rentals initiated  |
| count      | number of total rentals  |

Tabla 1: Atributos para el Problema 2 (*Bike Sharing*).

- (a) Cargue los datos de entrenamiento y pruebas como dataframes de *pandas*. Describa las variables involucradas en el problema, explorando el tipo de datos de que se trata, el número de valores distintos y, si corresponde, un gráfico (e.g. un histograma) que resuma su comportamiento. Su primera operación de pre-procesamiento de datos será obtener la hora del día desde el campo fecha (que en este momento es de tipo string), creando una nueva columna denominada *hour* y de tipo *int*. Para hacer esta operación se concatenarán los dataframes de entrenamiento y pruebas y luego se volverán a separar manteniendo la separación original.

```

1  import pandas as pd
2  import numpy as np
3  dftrain = pd.read_csv('bike_sharing_train.csv')
4  dfval = pd.read_csv('bike_sharing_val.csv')
5  dftest = pd.read_csv('bike_sharing_test.csv')
6  ntrain = len(dftrain)
7  nval = len(dftrain) + len(dfval)
8  df = pd.concat([dftrain, dfval, dftest])
9  print '\nSummary - dataframe completo:\n'
10 print df.describe()
11 df['hour'] = pd.to_datetime(df['datetime']).apply(lambda x: x.strftime('%H'))
12 df['hour'] = pd.to_numeric(df['hour'])

```

- (b) Entrene un árbol de regresión para resolver el problema usando parámetros por defecto. Con este fin, construya una matriz  $\mathbf{X}_{\text{train}}$  de forma  $n_{\text{train}} \times d_1$  que contenga los datos de entrenamiento en sus filas, seleccionando las columnas que desee/pueda utilizar para el entrenamiento. Implemente además, la función de evaluación que hemos definido anteriormente para este problema. Evalúe el árbol de regresión ajustado a los datos de entrenamiento sobre el conjunto de entrenamiento y pruebas. Construya un gráfico que compare las predicciones con los valores reales. En este punto usted debiese tener un modelo con puntaje del orden de 0.59, lo que lo dejará más o menos en la posición 2140 de la competencia.

```

1  from sklearn.tree import DecisionTreeRegressor as Tree
2  import matplotlib.pyplot as plt
3  def eval_bikemodel(y_predict, y_true):
4      diff = np.log(y_predict+1.0) - np.log(y_true+1.0)
5      return np.sqrt(np.sum(np.square(diff))/len(y_predict))
6  Xdf=df.ix[:, ['season', 'holiday', 'workingday', 'weather', 'temp', 'atemp',
7               'humidity', 'windspeed', 'hour']]
8  Ydf=df.ix[:, 'count']
9  X_train = Xdf[0:ntrain].values
10 X_val = Xdf[ntrain:nval].values
11 X_test = Xdf[nval:].values
12 Y_train = Ydf[0:ntrain].values
13 Y_val = Ydf[ntrain:nval].values
14 Y_test = Ydf[nval:].values
15
16 model = Tree(random_state=0)
17 model.fit(X_train, Y_train)
18 score_test = model.score(X_test, Y_test)
19 print "SCORE TEST=%f"%score_test
20
21 Y_pred_train = model.predict(X_train)
22 Y_pred_val = model.predict(X_val)
23 Y_pred_test = model.predict(X_test)
24 kagg_train = eval_bikemodel(Y_pred_train, Y_train)
25 kagg_val = eval_bikemodel(Y_pred_val, Y_val)
26 kagg_test = eval_bikemodel(Y_pred_test, Y_test)

```

```

27 print "KAGG EVAL TRAIN =%f"%kagg_train
28 print "KAGG EVAL TEST =%f"%kagg_test
29 plt.plot(Y_test,Y_pred_test, '.')
30 plt.show()

```

- (c) Mejore el árbol de regresión definido en el punto anterior haciendo modificaciones a los hiper-parámetros del modelo. Por ejemplo, como estos modelos tienden a sobre-ajustar, podría intentar limitar la profundidad del árbol (¿Por qué esto debiese ayudar?). Naturalmente, está absolutamente prohibido tomar este tipo de decisiones en función del resultado de pruebas. Debe realizar estas elecciones evaluando sobre el conjunto de validación. Si no desea utilizarlo, y prefiere implementar validación cruzada u otra técnica automática, tiene la ventaja de poder usar el conjunto de validación como parte del entrenamiento. Con estas modificaciones debiese poder mejorar su ranking en unas 300 posiciones.

```

1 model = Tree(random_state=0,max_depth=20)
2 model.fit(X_train,Y_train)
3 Y_pred_val = model.predict(X_val)
4 kagg_val = eval_bikemodel(Y_pred_val,Y_val)
5 print "KAGG EVAL VAL =%f"%kagg_val

```

- (d) Mejore el árbol de regresión definido en el punto anterior haciendo modificaciones sobre la representación utilizada para aprender desde los datos. Por ejemplo, los histogramas que construyó en el punto (a) así como la forma especial de la función de evaluación, sugieren una cierta transformación de la variable respuesta. Podría intentar también normalizando los datos o normalizando la respuesta. Otra opción es intentar rescatar algo más acerca de la fecha (anteriormente sólo se extrajo la hora), como por ejemplo el año o el día de la semana ('lunes','martes', etc) que corresponde. Sea creativo, este paso le debiese reportar un salto de calidad muy significativo. Una observación importante es que si hace una transformación a la variable respuesta (por ejemplo raíz cuadrada), debe invertir esta transformación antes de evaluar el desempeño con *eval\_bikemodel* (por ejemplo, elevar al cuadrado si tomó raíz cuadrada). Con modificaciones de este tipo, podría mejorar su ranking en unas 1000 posiciones, entrando ya al top-1000 con un score del orden de 0.45.

```

1 df['cday'] = pd.to_datetime(df['datetime']).dt.dayofweek#0:lunes,6:domingo
2 df['cday'] = pd.to_numeric(df['cday'])
3 Xdf=df.ix[:,['season','holiday','workingday','weather','temp','atemp',
4             'humidity','windspeed','hour','cday']]

```

- (e) Entrene una SVM no lineal para resolver el problema midiendo el efecto de las distintas representaciones que haya descubierto hasta este punto. Un detalle importante es que antes de entrenar la SVM sería aconsejable hacer dos tipos de pre-procesamiento adicional de los datos: (i) codificar las variables categóricas en un modo apropiado - por ejemplo como vector binario con un 1 en la posición del valor adoptado-, (ii) escalar los atributos de modo que queden centrados y con rangos comparables. Usando parámetros por defecto para la SVM debiese obtener un score del orden de 0.344, quedando definitivamente en el top-10 de la competencia.

```

1 #load dataframes as before ...
2 df = pd.concat([dftrain,dfval,dftest])
3 df['hour'] = pd.to_datetime(df['datetime']).apply(lambda x: x.strftime('%H'))
4 df['cday'] = pd.to_datetime(df['datetime']).dt.dayofweek
5 df['hour'] = pd.to_numeric(df['hour'])
6 df['cday'] = pd.to_numeric(df['cday'])
7 Xdf=df.ix[:,['season','holiday','workingday','weather','temp','atemp',
8             'humidity','windspeed','hour','cday']]
9 #PASO IMPORTANTE MAS ABAJO ...
10 Xdf = pd.get_dummies(Xdf,columns=['season','weather','hour','cday'])
11 Ydf=df.ix[:, 'count']
12

```

```

13 from sklearn.preprocessing import StandardScaler
14 scalerX = StandardScaler()
15 X_train = scalerX.fit_transform(X_train)
16 X_val = scalerX.fit_transform(X_val)
17 X_test = scalerX.transform(X_test)
18
19 from sklearn.svm import SVR
20 model = SVR()
21 model.fit(X_train,Y_train)
22 Y_pred_train = model.predict(X_train)
23 Y_pred_val = model.predict(X_val)
24 Y_pred_test = model.predict(X_test)

```

- (f) Mejore la SVM definida en el punto anterior haciendo modificaciones a los hiper-parámetros de la máquina ( $C$ ,  $\epsilon$  o la misma función de kernel). Naturalmente, está absolutamente prohibido tomar este tipo de decisiones de diseño mirando el resultado de pruebas. Debe realizar estas elecciones evaluando sobre el conjunto de validación. Si no desea utilizarlo, y prefiere implementar validación cruzada u otra técnica automática, tiene la ventaja de poder usar el conjunto de validación como parte del entrenamiento.

```

1 model = SVR(C=1,epsilon=0.01)
2 kagg_train = eval_bikemodel(Y_pred_train,Y_train)
3 kagg_val = eval_bikemodel(Y_pred_val,Y_val)
4 print "KAGG EVAL TRAIN =%f"%kagg_train
5 print "KAGG EVAL VAL =%f"%kagg_val

```

- (g) Evalúe el efecto de utilizar el dataset de validación para entrenamiento y seleccionar los parámetros estructurales del árbol de clasificación y la SVM usando validación cruzada. El código de ejemplo para esto ha sido proporcionado en las tareas 1 y 2, pero se adjunta de nuevo a continuación

```

1 from sklearn.model_selection import KFold
2 kf = KFold(n_splits=10)
3 mse_cv = 0
4 for train, val in kf.split(Xm):
5     model = #define your model
6     model.fit(Xm[train], ym[train])
7     yhat_val = model.predict(Xm[val])
8     ytrue_val = ym[val]
9     score_fold = eval_bikemodel(yhat_val,ytrue_val)
10    mse_cv += score_fold
11 mse_cv = mse_cv / 10

```

- (h) Evalúe el efecto de utilizar un ensamblado de 2 máquinas de aprendizaje para predecir la demanda total de bicicletas. Un modelo se especializará en la predicción de la demanda de bicicletas de parte de usuarios registrados y otra en la predicción de la demanda de usuarios casuales. Hay razones claras para pensar que los patrones son distintos.

```

1 Ydf=df.ix[:, 'count'] #demanda total
2 Ydf=df.ix[:, 'registered'] #demanda registrada
3 Ydf=df.ix[:, 'casual'] #demanda casual

```

- (i) Evalúe el efecto de utilizar un algoritmo genérico para ensamblar máquinas de aprendizaje para predecir la demanda total de bicicletas. Puede experimentar con una sola técnica (e.g. *Random Forest*), discuta la evolución a medida que aumenta el número de máquinas.

```

1 from sklearn.ensemble import RandomForestRegressor
2 model = RandomForestRegressor(n_estimators=10,max_depth=max_depth,random_state=0)

```



### 3. Calidad de un vino

Dentro de las muchas variedades de vino existentes, algunas gustan más que otras, esto es debido al gusto de una persona en particular o bien a la gran cantidad de químicos y procesos que se aplican a la producción de vino. Para el área de negocios, el estimar cuál es la calidad de un vino en base a la apreciación del público es una tarea bastante difícil.

Para esta actividad se trabajará con dos datasets asociados a las variantes tinto y blanco del vino portugués "Vinho Verde" [4]. Debido a temas privados solo se cuenta con las características fisicoquímicas asociadas a un vino en particular, los cuales corresponden a 11 atributos numéricos descritos en el siguiente link.

Este problema puede ser abordado como clasificación de 11 clases o de regresión, ya que el atributo a estimar, *quality*, es un valor entero entre 0 y 10.

- (a) Cargue los dos dataset en un único dataframe de pandas, además de agregar una columna indicando si es vino tinto o blanco. Describa el dataset a trabajar.

```
1 import pandas as pd
2 df_red = pd.read_csv("winequality-red.csv", sep=";")
3 df_white = pd.read_csv("winequality-white.csv", sep=";")
4 df = pd.concat([df_red, df_white], axis=0)
5 #genere atributo 'tipo'
```

- (b) Aborde este problema como si fuera de clasificación binaria para predecir si un vino es de buena calidad o no, es decir, utilice las distintas características fisicoquímicas presentes en los datos para estimar esta etiqueta. Para esto cree las matrices de entrenamiento y de pruebas, además de la etiqueta para ambos conjuntos, considerando como *quality* mayor a 5 un vino de buena calidad. El conjunto de pruebas (25 %) será utilizado únicamente para verificar la calidad de los algoritmos a entrenar.

```
1 df['good_quality'] = [1 if q>5 else 0 for q in df.quality]
2 #train and test split over df
```

- (c) Entrene un solo Árbol de Clasificación de múltiples niveles para resolver el problema. Puede variar los hiper-parámetros que prefiera, recuerde que las decisiones no pueden ser basadas mirando el conjunto de pruebas. Debido al desbalanceo que se produce en las dos clases mida la métrica *F1-score* [5] sobre el conjunto de entrenamiento y de pruebas.

- (d) Entrene un ensamblador de árboles de múltiples niveles, mediante la técnica de *Random Forest*. Varíe la cantidad de árboles de decisión utilizados en el ensamblado (*n\_estimators*), realice un gráfico resumen del *F1-score* de entrenamiento y de pruebas en función de este hiper-parámetro.

```
1 from sklearn.ensemble import RandomForestClassifier
2 model = RandomForestClassifier(n_estimators=, max_depth=, n_jobs=-1)
3 from sklearn.metrics import f1_score
4 f1_score(y_true, y_pred)
```

- (e) Entrene un ensamblador de árboles de múltiples niveles, mediante la técnica de *AdaBoost*. Varíe la cantidad de árboles de decisión utilizados en el ensamblado (*n\_estimators*), realice un gráfico resumen del *F1-score* de entrenamiento y de pruebas en función de este hiper-parámetro. Compare y analice con la técnica utilizada en d).

```
1 from sklearn.ensemble import AdaBoostClassifier
2 model = AdaBoostClassifier(base_estimator=Tree(max_depth=), n_estimators=)
```

- (f) Entrene alguna otra máquina de aprendizaje, elegida por usted, para resolver este problema. Elija los hiper-parámetros que estime convenientes intentando aumentar el *F1-score* obtenido por los algoritmos anteriores. Compare y analice estas 4 maneras de resolver el problema definido en b).

- (g) Defina un criterio para estimar la importancia de los distintos atributos en el ensamblado de *Random Forest*, implementelo sobre alguno de los ensambladores entrenados en d), haga un ranking de importancia de atributos ¿Es posible implementar este criterio sobre una técnica de *boost* como lo es *AdaBoost*?

## 4. Reconocimiento de Imágenes Sign Gestures

MNIST es un *dataset* muy popular de dígitos escrito a mano que a servido para probar distintos algoritmos de *Machine Learning* relacionados con *Computer Vision*. Buscando nuevos desafíos, investigadores generaron un dataset que podría usarse eventualmente en aplicaciones reales, Sign Gestures, consta de imágenes del lenguaje de señas, estas tienen una resolución de 28x28 pixeles representados en una escala de grises 0-255. La versión utilizada se atribuye a [8] y viene separada en 27455 ejemplos de entrenamiento y 7172 casos de pruebas. Las clases son mutuamente excluyentes y corresponden a las letras del alfabeto (ver imagen).



- (a) Construya una función que cargue todos los datos de entrenamiento y pruebas del problema generando como salida: (i) dos matrices  $X_{tr}, Y_{tr}$ , correspondientes a las imágenes y etiquetas de entrenamiento, (ii) dos matrices  $X_t, Y_t$ , correspondientes a las imágenes y etiquetas de pruebas, y finalmente (iii) dos matrices  $X_v, Y_v$ , correspondientes a imágenes y etiquetas que se usarán como conjunto de validación, es decir para tomar decisiones de diseño acerca del modelo. Este último conjunto debe ser extraído desde el conjunto de entrenamiento original y no debe superar las 7000 imágenes.

```
1 import pandas as pd
2 from sklearn.model_selection import train_test_split
3 def load_data():
4     train = pd.read_csv('sign_mnist_train.csv')
5     test = pd.read_csv('sign_mnist_test.csv')
6     y_tr = train['label']
7     x_tr = train.iloc[:,1:]
8     y_t = test['label']
9     x_t = test.iloc[:,1:]
10    #you need to add Xval: x_v, y_v
11    return(x_tr,x_v,x_t,y_tr,y_v,y_t)
12 x_tr, x_v, x_t, y_tr, y_v, y_t= load_data()
```

- (b) Construya una función que escale apropiadamente las imágenes antes de trabajar. Experimente sólo escalando los datos de acuerdo a la intensidad máxima de pixel (i.e., dividiendo por 255) y luego centrando y escalándolos como en actividades anteriores.
- (c) Diseñe, entrene y evalúe una red neuronal para el problema partir de la representación original de las imágenes. Experimente con distintas arquitecturas, pre-procesamientos y métodos de entrenamiento,

mediendo el error de clasificación sobre el conjunto de validación. En base a esta última medida de desempeño, decida qué modelo, de entre todos los evaluados, medirá finalmente en el conjunto de test. Reporte y discuta los resultados obtenidos. Se espera que logre obtener un error de pruebas menor o igual a 0.2.

```

1 from keras.models import Sequential
2 from keras.layers import Dense, Activation
3 from keras.optimizers import SGD
4 from keras.utils.np_utils import to_categorical
5
6 model = Sequential()
7 model.add(Dense(30, input_dim=x_tr.shape[1], init='uniform', activation='relu'))
8 model.add(Dense(30, init='uniform', activation='relu'))
9 model.add(Dense(25, init='uniform', activation='softmax'))
10 model.compile(optimizer=SGD(lr=0.05), loss='categorical_crossentropy', metrics=['accuracy'])
11 model.fit(x_tr.values, to_categorical(y_tr), nb_epoch=100, batch_size=128, verbose=1,
12         validation_data=(x_v.values, to_categorical(y_v)))

```

- (d) Para la mejor red entrenada anteriormente construya la matriz de confusión de las distintas clases, para así visualizar cuáles son las clases más difíciles de clasificar y con cuáles se confunden. Comente.
- (e) Entrene una SVM no lineal sobre los píxeles con y sin pre-procesamiento. Puede utilizar el conjunto de validación para seleccionar hiper-parámetros, como el nivel de regularización aplicado y/o la función de kernel a utilizar.
- (f) Entrene una árbol de clasificación sobre los píxeles con y sin pre-procesamiento. Puede utilizar el conjunto de validación para seleccionar hiper-parámetros, como la profundidad máxima del árbol.

## Referencias

- [1] Keras: Deep Learning library for Theano and TensorFlow. <https://keras.io/>
- [2] Bernhard Schoelkopf, Alexander J. Smola, and Klaus-Robert Mueller. 1999. Kernel principal component analysis. In *Advances in kernel methods*, MIT Press, Cambridge, MA, USA 327-352.
- [3] <https://www.kaggle.com/c/bike-sharing-demand>
- [4] <http://archive.ics.uci.edu/ml/datasets/Wine+Quality>
- [5] [http://scikit-learn.org/stable/modules/generated/sklearn.metrics.f1\\_score.html](http://scikit-learn.org/stable/modules/generated/sklearn.metrics.f1_score.html)
- [6] Dalal, N., Triggs, B. (2005, June). Histograms of oriented gradients for human detection. In *2005 IEEE Computer Society Conference on Computer Vision and Pattern Recognition (CVPR'05)* (Vol. 1, pp. 886-893). IEEE.
- [7] Forsyth, D. A., Ponce, J. (2002). *Computer vision: a modern approach*. Prentice Hall Professional Technical Reference.
- [8] <https://www.kaggle.com/datamunge/sign-language-mnist>