

INF-493. MACHINE LEARNING. TAREA 3.

MÉTODOS NO-LINEALES.

Prof. Ricardo Ñanculef & Carlos Valle
jnancu@inf.utfsm.cl

Temas

- Manipulación de dataframes en *pandas*.
- Pre-procesamiento de datos y extracción de características.
- SVMs no-lineales en *sklearn*.
- Árboles de clasificación en *sklearn*.
- Redes Neuronales feed-forward en *keras*.
- Ensamblados de máquinas de aprendizaje en *sklearn*.
- Selección de hiper-parámetros estructurales vía validación cruzada.

Formalidades

- Equipos de trabajo de: 2 personas.¹.
- Se debe preparar un (breve) Jupyter/IPython notebook que explique la actividad realizada y las conclusiones del trabajo.
- Se debe preparar una presentación de 20 minutos.
- Se debe mantener un respaldo de cualquier tipo de código utilizado, informe y presentación en Github.
- Deadline: 2 semanas desde publicación.
- 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].

Paquetes/Librerías

Como es usual utilizaremos *numpy*, *scipy*, *matplotlib* y *sklearn*.

Para esta tarea necesitará además instalar *keras*, una librería para prototipado rápido de modelos basados en redes neuronales usando python, muy similar en espíritu a *sklearn*. La librería puede usar *Tensor Flow* o *Theano* como backend, siendo éstas las librerías más populares para desarrollar nuevos modelos de redes neuronales o implementar eficientemente modelos conocidos con fines aplicativos. Para detalles sobre la instalación puede revisar [1] o escribir un email a su ayudante.

¹Cada miembro del equipo debe estar en condiciones de realizar una presentación y discutir sobre cada punto del trabajo realizado.

1 El Viejo XOR: Métodos No-lineales para Problemas No-lineales

El objetivo de esta sección es experimentar con algunos modelos no-lineales sobre un problema de juguete, pero bastante famoso en la historia del aprendizaje automático, que denominaremos *XOR*. 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 $y \in \{0,1\}$ 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*.

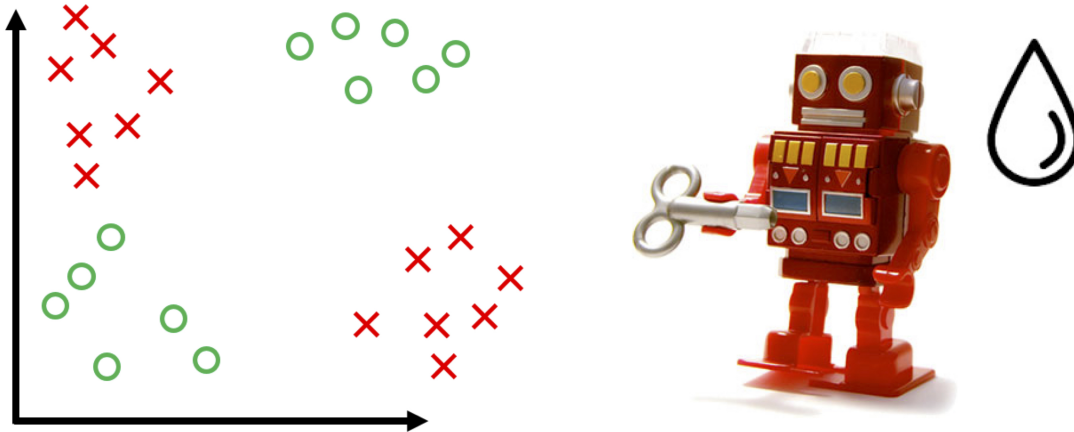


Figure 1: Distribución deseada para la actividad 1. Los 2 colores representan 2 clases distintas.

- (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.1². 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 10% de ruido al dataset de entrenamiento, generando \mathbf{x} 's cercanos a la frontera e invirtiendo la etiqueta correcta de este subconjunto. Genere un gráfico que muestre sus datos de entrenamiento y pruebas, identificando cada clase con un color diferente. ¿Porqué tiene sentido denominar este problema *XOR*?

```
1 import numpy as np
2
3 def do_XOR(n=1000, noisy_n=100, svm=True):
4     rng = np.random.RandomState(0)
5     X_train = rng.uniform(low=-1.0, high=1.0, size=(n, 2))
6     Y_train = np.logical_xor(X_train[:, 0] > 0, X_train[:, 1] > 0)
7     Y_train = 2*Y_train-1 if svm else Y_train
8     X_noisy = rng.uniform(low=-0.8, high=0.2, size=(noisy_n, 2))
9     Y_noisy = -1*np.logical_xor(X_noisy[:, 0] > 0, X_noisy[:, 1] > 0) + 1
10    Y_noisy = 2*Y_noisy-1 if svm else Y_noisy
11    X_train = np.vstack((X_train, X_noisy))
12    Y_train = np.hstack((Y_train, Y_noisy))
13    X_test = rng.uniform(low=-1.0, high=1.0, size=(n, 2))
14    Y_test = np.logical_xor(X_test[:, 0] > 0, X_test[:, 1] > 0)
15    Y_test = 2*Y_test - 1 if svm else Y_test
16    return X_train, Y_train, X_test, Y_test
17
18 X_train, Y_train, X_test, Y_test = do_XOR()
```

²Puede generar datos aleatorios distribuidos uniformemente en $[-1, +1] \times [-1, +1]$ para luego etiquetar aquellos ubicados en el primer y tercer cuadrante como 0 y aquellos ubicados en el segundo y cuarto cuadrante como 1.

- (b) Demuestre experimentalmente que una SVM lineal no puede resolver satisfactoriamente el problema anterior ³. 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. Escriba también una función, denominada *plot_classifier* que represente gráficamente la solución encontrada por la máquina de aprendizaje, representando los datos de entrenamiento y pruebas (asociando un color a la clase que les corresponde), además de la frontera de clasificación. Describa y explique lo que observa, reportando gráficos de la solución sólo para algunos casos representativos.

```

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

```

- (c) Demuestre experimentalmente que una SVM con kernel no-lineal puede resolver satisfactoriamente el problema. Para esta actividad, use los hiper-parámetros que se entregan como referencia en el código de ejemplo (muestre resultados tanto para el kernel *Polynomial* como para el kernel *RBF*). Cambie el valor del parámetro C eligiendo κ en el intervalo $[-2, 4]$ (con saltos de 1 unidad) y asignando al parámetro el valor $C = 2^\kappa$. Describa y explique lo que observa, graficando el error de entrenamiento y pruebas como función de C . Utilice la función *plot_classifier*, diseñada en el ítem anterior, para construir gráficos de la solución con distintos valores de C .

```

1  from sklearn.svm import SVC
2  clf = SVC(C=100, kernel='rbf')
3  clf = SVC(C=10, kernel='poly',degree=2, coef0=1)
4  clf.fit(X_train, Y_train)
5  print "Test Accuracy = %f"%clf.score(X_test,Y_test)
6  plot_classifier(clf,X_train,Y_train,X_test,Y_test,'svm')

```

- (d) 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()

```

³Defina satisfactorio como al menos 50% mejor que una predicción aleatoria de 1 de las clases.

```

6 model.add(Dense(1, input_dim=X_train.shape[1], init='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, nb_epoch=50, batch_size=100, verbose=1)
10 scores = model.evaluate(X_test, Y_test)
11 test_acc = scores[1]

```

- (e) 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=8
3 model = Sequential()
4 model.add(Dense(1, input_dim=X_train.shape[1], init='uniform', activation='relu'))
5 model.add(Dense(n_h, init='uniform', activation='sigmoid'))

```

- (f) 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')

```

- (g) 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)

```

2 Bike Sharing: Predicción de Demanda Horaria

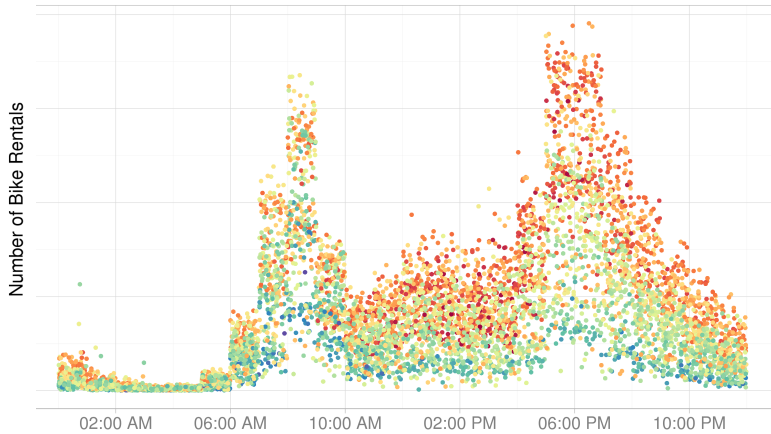
En esta sección simularemos nuestra participación en el desafío *Bike Sharing Demand* de *Kaggle* [2]. El objetivo es predecir la demanda de bicicletas sobre la red *Capital Bikeshare* de la ciudad de Washington, D.C., en función de la hora del día y otras variables descritas en la tabla 1. En principio, 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 2: Mist + Cloudy, Mist + Broken clouds, Mist + Few clouds, Mist 3: Light Snow, Light Rain + Thunderstorm + Scattered clouds, Light Rain + Scattered clouds 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

Table 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 dejaría 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 (porque 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 Y_pred_val = model.predict(X_val)
3 kagg_val = eval_bikemodel(Y_pred_val,Y_val)
4 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 , ϵ 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 import cross_validation
2 k_fold = cross_validation.KFold(len(X_train),10)
3 score_cv = 0
4 for k, (train, val) in enumerate(k_fold):
5     model = #define your model, e.g. model = SVR(C=candidate)
6     model.fit(X_train[train], Y_train[train])
7     Ypred_val = model.predict(X_train[val])
8     Ytrue_val = Y_train[val]
9     score_fold = eval_bikemodel(Ypred_val,Ytrue_val)
10    score_cv += score_fold
11 score_cv = score_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 especializaría 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*), pero por favor estudie y discuta la evolución del rendimiento (de entrenamiento y pruebas) a medida que aumenta el número de máquinas.

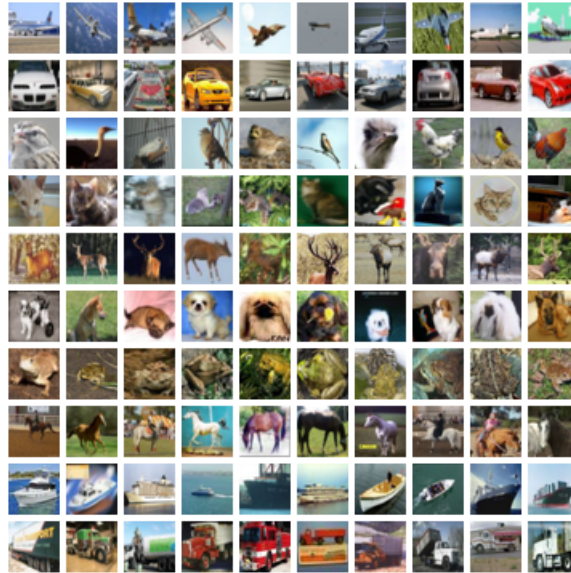
```

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

```


3 Reconocimiento de Imágenes en CIFAR10

En esta sección trabajaremos con un dataset muy utilizado para experimentar con reconocimiento de imágenes: CIFAR10. Se trata de un conjunto de 60.000 imágenes RGB de 32×32 píxeles que contiene 10 clases de objetos (6000 ejemplos por clase). La versión utilizada se atribuye a [5] y viene separada en 50000 ejemplos de entrenamiento y 10000 casos de prueba. El conjunto de pruebas fue obtenido seleccionando 1000 imágenes aleatorias de cada clase. Los datos restantes han sido ordenados aleatoriamente y están organizados en 5 bloques de entrenamiento (batches). Las clases son mutuamente excluyentes y corresponden a las siguientes categorías: gato, perro, rana, caballo, pájaro, ciervo, avión, automóvil, camión y barco.



Los datos asociados a esta actividad podrán ser obtenidos utilizando los siguientes comandos en la línea de comandos (sistemas UNIX)

```
1 wget http://octopus.inf.utfsml.cl/~ricky/data.tar.gz
2 tar -xzf data.tar.gz
3 rm data.tar.gz
```

En la carpeta generada encontrarán 5 archivos denominados 'data_batch_1', 'data_batch_2', 'data_batch_3', 'data_batch_4', 'data_batch_5' y 'test_batch' correspondientes a los 5 bloques de entrenamiento y al conjunto de pruebas respectivamente. Los archivos corresponden a diccionarios serializados de python y pueden ser “extraídos” utilizando la siguiente función:

```
1 def unpickle(file):
2     import cPickle
3     fo = open(file, 'rb')
4     dict = cPickle.load(fo)
5     fo.close()
6     return dict
```

Una vez extraído, cada diccionario contendrá 2 elementos importantes: *data* y *labels*. El primer elemento (*data*) es una matriz de 10000×3072 (numpy array). Cada fila de esa matriz corresponde a una imagen *RGB*: los primeros 1024 valores vienen del canal *R*, los siguientes 1024 del canal *G*, y los últimos 1024 del canal *B*. Para cada canal, las imágenes han sido vectorizadas por filas, de modo que los primeros 32 valores del canal *R* corresponden a la primera fila de la imagen. Por otro lado, el elemento (*labels*) del diccionario contiene una lista de 1000 valores enteros entre 0 y 9 que identifican las clases antes enumeradas.

```
1 label_names = ['airplane', 'automobile', 'bird', 'cat', 'deer', 'dog', \
2                'frog', 'horse', 'ship', 'truck']
```

- (a) Construya una función que cargue todos los bloques de entrenamiento y pruebas del problema CIFAR 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 10000 imágenes.

```

1  from scipy.misc import imread
2  import cPickle as pickle
3  import numpy as np
4  import os
5
6  def load_CIFAR_one(filename):
7      with open(filename, 'rb') as f:
8          datadict = pickle.load(f)
9          X = datadict['data']
10         Y = datadict['labels']
11         return X, np.array(Y)
12
13  def load_CIFAR10(PATH):
14      xs = [], ys = []
15      for b in range(1,6):
16          f = os.path.join(PATH, 'data_batch_%d' % (b, ))
17          X, Y = load_CIFAR_one(f)
18          xs.append(X)
19          ys.append(Y)
20      Xtr = np.concatenate(xs)
21      Ytr = np.concatenate(ys)
22      del X, Y
23      Xte, Yte = load_CIFAR_batch(os.path.join(PATH, 'test_batch'))
24      return Xtr, Ytr, Xte, Yte
25      #you need to add Xval
26  load_CIFAR10('..')

```

- (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 CIFAR a partir de la representación original de las imágenes (píxeles RGB). Experimente con distintas arquitecturas y métodos de entrenamiento, midiendo 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, evaluará 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 0.5.

```

1  from keras.utils.np_utils import to_categorical
2  from sklearn.preprocessing import StandardScaler
3
4  def scaler_function(Xtr,Xt,scale=True):
5      scaler = StandardScaler(with_std=scale).fit(Xtr)
6      Xtr_scaled = scaler.transform(Xtr)
7      Xt_scaled = scaler.transform(Xt)
8      return Xtr_scaled, Xt_scaled
9
10  Xtr,Xte = scaler_function(Xtr,Xte)
11  Ytr = to_categorical(Ytr)
12  Yte = to_categorical(Yte)
13

```

```

14 from keras.models import Sequential
15 from keras.layers import Dense, Activation
16 from keras.optimizers import SGD
17
18 model = Sequential()
19 model.add(Dense(100, input_dim=Xtr.shape[1], init='uniform', activation='relu'))
20 model.add(Dropout(0.1))
21 model.add(Dense(10, init='uniform', activation='softmax'))
22 model.compile(optimizer=SGD(lr=0.05), loss='binary_crossentropy', metrics=['accuracy'])
23 model.fit(Xtr, Ytr, nb_epoch=50, batch_size=32, verbose=1, validation_data=(Xte,Yte))
24 scores = model.evaluate(Xte, Yte)
25 test_acc = scores[1]

```

- (d) Repita la actividad anterior, pero mejorando los atributos utilizados para representar las imágenes. Para esta parte, se distribuirá junto a esta tarea una función denominada *extract_features.py* que extraerá 2 tipos de representaciones sobre una imagen y conjunto de imágenes: (i) histogramas de tono [4], (ii) descriptores HOG [3]. Reporte y discuta los resultados obtenidos utilizando las distintas representaciones por separado o todas simultáneamente. La función *extract_features.py* estará definida en un script denominado *top_level_features.py* y puede ser importada y utilizada como se muestra a continuación.

```

1 from top_level_features import hog_features
2 from top_level_features import color_histogram_hsv
3 from top_level_features import extract_features
4 Xtr, Ytr, Xte, Yte = load_CIFAR10("datasets/")
5 features = extract_features(Xtr,[hog_features]) #extrae hog features
6 features = extract_features(Xtr,[color_histogram_hsv]) #extrae histogramas de color
7 features = extract_features(Xtr,[hog_features, color_histogram_hsv]) #extrae todo
8 print Xtr.shape
9 print features.shape

```

- (e) Entrene una SVM no lineal sobre los pixeles originales y sobre los atributos de alto nivel utilizados para representar las imágenes en el ítem anterior. 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 un árbol de clasificación sobre los pixeles originales y sobre los atributos de alto nivel utilizados para representar las imágenes en el ítem anterior. Puede utilizar el conjunto de validación para seleccionar hiper-parámetros, como la profundidad máxima del árbol.

References

- [1] Keras: Deep Learning library for Theano and TensorFlow. <https://keras.io/>
- [2] <https://www.kaggle.com/c/bike-sharing-demand>
- [3] 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.
- [4] Forsyth, D. A., Ponce, J. (2002). Computer vision: a modern approach. Prentice Hall Professional Technical Reference.
- [5] Krizhevsky, A., Hinton, G. (2009). Learning multiple layers of features from tiny images.