

INF-477. REDES NEURONALES ARTIFICIALES. I-2017.

TAREA 1 - PERCEPTRONES MULTICAPA Ó REDES FEED-FORWARD

Temas

- Diseño de Redes Feed-Forward para problemas de clasificación y regresión.
- Entrenamiento de Redes Feed-Forward vía GD y variantes (SGD, mini-batches).
- Momentum y estrategias de selección de paso en SGD.
- Evaluación de Redes Feed-Forward vía validación cruzada (cross-validation).
- Regularización de Redes Feed-Forward.

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.
- Se debe preparar una presentación de 20 minutos. Presentador será elegido aleatoriamente.
- Se debe mantener un respaldo de cualquier tipo de código utilizado, informe y presentación en Github.
- Fecha de entrega y discusión: definida en calendario publicado en moodle.
- Formato de entrega: envío de link Github al correo electrónico del ayudante[†], incluyendo al profesor en copia[‡]. Por favor especificar el siguiente asunto: [Tarea1-INF477-I-2017].

*Cada uno debe estar en condiciones de realizar una presentación y discutir sobre cada punto del trabajo realizado.

[†]joaquin.velasquez@alumnos.inf.utfsm.cl

[‡]jnancu@inf.utfsm.cl

1 Rol de las Capas Ocultas de una ANN

Como hemos discutido en clases, el modelo de neurona utilizando en la mayoría de las redes neuronales artificiales corresponde esencialmente a un discriminante lineal sobre el espacio original de atributos. Por lo tanto, si se entrena una neurona para resolver un problema cuya solución no es una función lineal de esos atributos debiésemos esperar un pobre rendimiento del modelo en la tarea. Un perceptrón multicapas en cambio agrega un nivel adicional de procesamiento entrenado para aprender automáticamente una representación adecuada al problema. En esta sección experimentaremos esta idea revisitando un problema “de juguete” denominado *moons*.

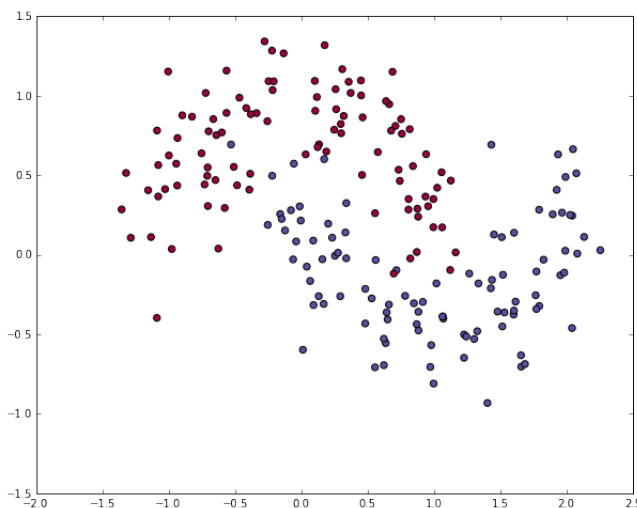


Fig. 1: Distribución deseada para la actividad 1. Los 2 colores representan 2 clases distintas.

- (a) Implemente una función que genere 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 a la izquierda de la Fig.1 (se trata esencialmente de dos semi-círculos de radio $\sqrt{2}$, uno centrado en $(0.0, 0.5)$ y el otro en $(1.0, 0.0)$). Utilice la función para generar dos subconjuntos distintos de 500 puntos a ser usados como conjuntos de entrenamiento y pruebas (por el objetivo de la actividad, en este caso no utilizaremos un conjunto de validación). Haga un gráfico que muestre los datos obtenidos.

```
1 from sklearn.datasets import make_moons
2 X, y = make_moons(200, noise=0.20)
3 plt.scatter(X[:,0], X[:,1], s=40, c=y, cmap=plt.cm.Spectral)
```

- (b) Demuestre experimentalmente que una neurona artificial individual no puede resolver satisfactoriamente el problema anterior. Puede utilizar la función de activación (no-lineal) y el método de entrenamiento que prefiera. Sea convincente. Describa y explique lo que observa.
- (c) Demuestre experimentalmente que un perceptrón multicapas puede resolver satisfactoriamente el problema obtenido en (a). Puede utilizar la arquitectura y el método de entrenamiento que prefiera. Sea convincente. Describa y explique lo que observa.

2 Back-propagation (BP) from Scratch

BP (Back-propagation) es sin duda el paradigma dominante para entrenar redes neuronales feed-forward. En redes grandes, diseñadas para problemas reales, implementar BP eficientemente puede ser una tarea delicada que puede ser razonable delegar a una librería especializada. Sin embargo, construir BP *from scratch* es muy útil con fines pedagógicos.

- (a) Escriba un programa que permita entrenar una red FF con 1 capa escondida (H neuronas) y O neuronas de salida, sin usar librerías, excepto eventualmente *numpy* para implementar operaciones básicas de álgebra lineal. Por simplicidad, asuma que todas las neuronas implementan una función de activación diferenciable y que la función de error (loss function) también lo es. Especifique explícitamente las funciones anteriores, así como sus gradientes. Escriba funciones para: (i) dar valores iniciales a los pesos de la red, (ii) implementar el *forward pass*, (iii) implementar el *backward pass* y (iv) implementar la rutina principal de entrenamiento, adoptando, por simplicidad, la variante cíclica de SGD (un ejemplo a la vez, pero iterando cíclicamente sobre el conjunto de entrenamiento) con una tasa de aprendizaje y número de ciclos fijos (epochs).
- (b) Escriba una función que permita hacer predicciones mediante una red FF con 1 capa escondida (H neuronas) y O neuronas de salida, sin usar librerías, excepto eventualmente *numpy*. Escriba una función vectorizada que implemente el *forward pass* sobre un conjunto de n_{test} ejemplos.
- (c) Demuestre que sus programas funcionan en un problema de clasificación, eligiendo funciones de error y de activación apropiadas. Para esto utilice el dataset *seeds*, disponible en UCI [1] y correspondiente a la clasificación de distintos tipos de semillas (3 clases) (recuerde que usualmente es conveniente normalizar los datos antes de trabajar con el modelo).

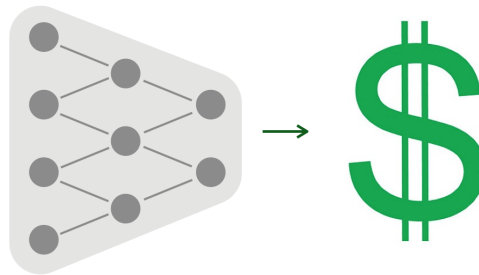
```
1 from sklearn.preprocessing import StandardScaler
2 import pandas as pd
3 url = 'http://archive.ics.uci.edu/ml/machine-learning-databases/00236/seeds_dataset.txt'
4 df = pd.read_csv(url, sep=r'\s+', header=None)
5 X_train = df.ix[:,0:6]
6 y_train = df.ix[:,7]
7 scaler = StandardScaler().fit(X_train)
8 X_train = scaler.transform(X_train)
```

Para evaluar los resultados, construya un gráfico correspondiente al error de clasificación versus número de epochs, utilizando sólo el conjunto de entrenamiento (el objetivo de esta sección es familiarizarse con el algoritmo BP, no encontrar la mejor red). Grafique también la evolución de la función objetivo utilizada para el entrenamiento.

- (d) Construya, sin usar librerías, excepto eventualmente *numpy* para implementar operaciones básicas de álgebra lineal, una variante de su programa anterior que entrene la red utilizando *weight-decay*.

3 Predicción del Precio de una Casa

En esta sección trabajaremos con un pequeño dataset conocido como *Boston Housing* que nos permitirá experimentar de modo más completo y exhaustivo con las técnicas bajo estudio. El problema consiste en predecir el precio de una casa en una zona/barrio de Boston (USA) a partir de una serie de atributos que describen el lugar que éste se ubica: tasa de criminalidad, proporción de zona residencial, proporción de zona industrial, si se encuentra junto al río ó no, contaminación atmosférica medida como la concentración de óxidos nítricos en el aire, etc. Para ver en detalle la descripción de la semántica asociada a los atributos de este problema, puede consultar <https://archive.ics.uci.edu/ml/datasets/Housing>.



- (a) Construya un dataframe con los datos a analizar descargando los datos desde la URL mantenida por los autores de [2]. Describa brevemente el dataset utilizar.

```
1 import pandas as pd
2 url = 'http://mldata.org/repository/data/download/csv/regression-datasets-housing/'
3 df = pd.read_csv(url, sep=',', header=None, names=['CRIM', 'ZN', 'INDUS', 'CHAS', 'NOX',
4           'RM', 'AGE', 'DIS', 'RAD', 'TAX', 'PTRATIO', 'B', 'LSTAT', 'MEDV'])
5 from sklearn.cross_validation import train_test_split
6 df_train, df_test = train_test_split(df, test_size=0.25, random_state=0)
7 df.shape
8 df.info()
9 df.describe()
```

- (b) Normalice los datos antes de trabajar. Determine la conveniencia de realizar esta operación.

```
1 from sklearn.preprocessing import StandardScaler
2 scaler = StandardScaler().fit(df_train)
3 X_train_scaled = pd.DataFrame(scaler.transform(df_train), columns=df_train.columns)
4 y_train = df_train.pop('MEDV')
```

- (c) Muestre en un gráfico el error cuadrático (MSE) vs número de *epochs* de entrenamiento, para una red *feedforward* de 3 capas, con 200 unidades ocultas y función de activación *sigmoidal*. Entrene la red usando gradiente descendente estocástico con tasa de aprendizaje (*learning rate*) 0.01 y 300 epochs de entrenamiento, en el conjunto de entrenamiento y de test. Comente. Si observara divergencia durante el entrenamiento, determine si esto ocurre para cada repetición del experimento.

```
1 from keras.models import Sequential
2 from keras.layers.core import Dense, Activation
3 from keras.optimizers import SGD
4
5 model = Sequential()
6 model.add(Dense(200, input_dim=X_train_scaled.shape[1], init='uniform'))
7 model.add(Activation('sigmoid'))
```

```

8  model.add(Dense(1, init='uniform'))
9  model.add(Activation('linear'))
10
11  sgd = SGD(lr=0.01)
12  model.compile(optimizer=sgd,loss='mean_squared_error')
13
14  hist = model.fit(X_train_scaled.as_matrix(), y_train_scaled.as_matrix(), nb_epoch=300,
15                  verbose=1, validation_data=(X_test_scaled.as_matrix(), y_test_scaled.as_matrix()))

```

(d) Repita el paso anterior, utilizando 'Relu' como función de activación y compare con lo obtenido en d).

(e) Repita c) y d) variando la tasa de aprendizaje (learning rate) en un rango sensible. Comente. Si observara divergencia durante el entrenamiento, determine si esto ocurre para cada repetición del experimento.

```

1  import numpy as np
2  n_lr = 20
3  lear_rate = np.linspace(0,1,n_lr)

```

(f) Estime el error de predicción de los modelos c) y d) usando validación cruzada con un número de folds igual a $K = 5$ y $K = 10$. Recuerde que para que la estimación sea razonable debe ajustar los pesos del modelo de nuevo, cada vez que trabaja sobre un determinado fold. Mida el error real del modelo sobre el conjunto de pruebas, compare y concluya.

```

1  from sklearn import cross_validation
2  Xm = X_train_scaled.as_matrix()
3  ym = y_train_scaled.as_matrix()
4  kfold = cross_validation.KFold(len(Xm), 10)
5  cvscores = []
6  for i, (train, val) in enumerate(kfold):
7      # create model
8      model = Sequential()
9      model.add(Dense(200, input_dim=Xm.shape[1], init='uniform'))
10     model.add(Activation('relu'))
11     model.add(Dense(1, init='uniform'))
12     model.add(Activation('linear'))
13     # Compile model
14     sgd = SGD(lr=0.2)
15     model.compile(optimizer=sgd,loss='mean_squared_error')
16     # Fit the model
17     model.fit(Xm[train], ym[train], nb_epoch=300)
18     # evaluate the model
19     scores = model.evaluate(Xm[val], ym[val])
20     cvscores.append(scores)
21  mse_cv = np.mean(cvscores)

```

(g) Entrene los modelos considerados en c) y d) usando *progressive decay* Compare y comente.

```

1  n_decay = 10
2  lear_decay = np.logspace(-6,0,n_decay)
3  sgd = SGD(lr=0.2, decay=1e-6)

```

(h) Entrene los modelos considerados en c) y d) usando momentum. Experimente usando momentum clásico y momentum de Nesterov. ¿Observa un mejor resultado final? ¿Observa una mayor velocidad de convergencia sobre el dataset de entrenamiento? ¿Sobre el dataset de pruebas?

```

1  n_decay = 21
2  momentum = np.linspace(0,1,n_decay)
3  sgd = SGD(lr=0.2,momentum=0.9)

```

- (i) Entrene los modelos considerados en c) y d) utilizando SGD en mini-batches. Experimente con diferentes tamaños del batch. Comente.

```

1  n_batches = 21
2  batch_sizes = np.round(np.linspace(1,X_train_scaled.shape[0],n_batches))
3  model.fit(X_train_scaled.as_matrix(),y_train_scaled.as_matrix(),batch_size=50,nb_epoch=300)

```

- (j) Entrene los modelos obtenidos en c) y d) utilizando estrategias modernas para adaptar la tasa de aprendizaje. Compare los desempeños de *adagrad*, *adadelata*, *RMS prop* y *adam*. ¿Se observa en algún caso un mejor resultado final? ¿Se observa en algún caso una mayor velocidad de convergencia sobre el dataset de entrenamiento? ¿Sobre el dataset de pruebas?

```

1  from keras.optimizers import SGD, Adam, RMSprop, Adagrad, Adadelata
2  moptimizer = Adagrad(lr=0.01)
3  model.compile(optimizer=moptimizer)
4  model.fit(X_train_scaled.as_matrix(),y_train_scaled.as_matrix())

```

- (k) Entrene los modelos obtenidos en c) y d) utilizando regularizadores ℓ_1 y ℓ_2 (weight decay). Compare los desempeños de prueba obtenidos antes y después de regularizar. Experimente con distintos valores del parámetro de regularización y comente. Luego, agregue una capa al modelo y aumente significativamente el número de neuronas escondidas en cada capa. Evalúe nuevamente el efecto de los regularizadores.

```

1  model = Sequential()
2  #la regularization se debe incorporar a cada capa separadamente
3  idim=X_train_scaled.shape[1]
4  model.add(Dense(200,input_dim=idim,init='uniform',W_regularizer=l2(0.01)))
5  model.add(Activation('sigmoid'))
6  model.add(Dense(1, init='uniform',W_regularizer=l2(0.01)))
7  model.add(Activation('linear'))

```

- (l) Entrene los modelos obtenidos en c) y d) utilizando *Dropout*. Compare los desempeños de prueba obtenidos antes y después de regularizar. Experimente con distintos valores del parámetro de regularización y comente. Luego, agregue una capa al modelo y aumente significativamente el número de neuronas escondidas en cada capa. Evalúe nuevamente el efecto de la regularización, evaluando también el efecto de regularizar sólo la capa inicial versus la alternativa de regularizar todas las capas.

```

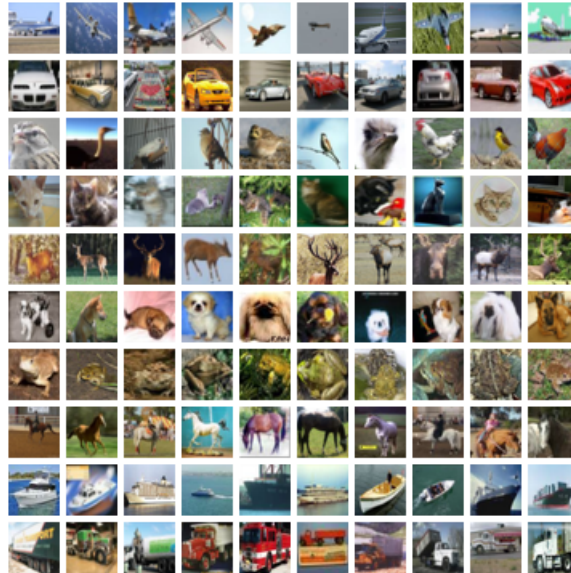
1  from keras.layers import Dropout
2  model = Sequential()
3  model.add(Dropout(0.2))
4  idim=X_train_scaled.shape[1]
5  model.add(Dense(200,input_dim=idim,init='uniform'))
6  model.add(Activation('sigmoid'))
7  model.add(Dense(1, init='uniform'))
8  model.add(Activation('linear'))

```

- (m) Fijando todos los demás hiper-parámetros del modelo definido en c) y d), utilice validación cruzada para determinar el mejor valor correspondiente a un parámetro que usted elija (tasa de aprendizaje, número de neuronas, parámetro de regularización, etc).

4 Reconocimiento de Imágenes en CIFAR10

En esta sección trabajaremos con un dataset bastante conocido y utilizado por la comunidad para experimentar con reconocimiento de objetos en imágenes: CIFAR10. Se trata de un conjunto de 60.000 imágenes RGB de 32×32 píxeles que contiene 10 clases de objetos y 6000 ejemplos por clase. La versión utilizada se atribuye a A. Krizhevsky, V. Nair y G. Hinton [4] 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: gatos, perros, ranas, caballos, pájaros, ciervos, aviones, automóviles, camiones y barcos.



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 un 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 5000 imágenes.

```

1 import cPickle as pickle
2 import numpy as np
3 import os
4 from scipy.misc import imread
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        Y = np.array(Y)
12        return X, Y
13
14 def load_CIFAR10(PATH):
15     xs = []
16     ys = []
17     for b in range(1,6):
18         f = os.path.join(PATH, 'data_batch_%d' % (b, ))
19         X, Y = load_CIFAR_one(f)
20         xs.append(X)
21         ys.append(Y)
22     Xtr = np.concatenate(xs)
23     Ytr = np.concatenate(ys)
24     del X, Y
25     Xte, Yte = load_CIFAR_batch(os.path.join(PATH, 'test_batch'))
26     return Xtr, Ytr, Xte, Yte

```

- (b) Construya una función que escale apropiadamente las imágenes antes de trabajar. Experimente sólo centrando los datos y luego centrando y escalándolos como en actividades anteriores.
- (c) Diseñe, entrene y evalúe una red neuronal con salida softmax 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.
- (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 [7], (ii) descriptores HOG [6]. Reporte y discuta los resultados obtenidos utilizando las distintas representaciones por separado o todas simultáneamente. La función *extract.features.py* estará definida

[§]Para empezar, puede considerar una pequeña arquitectura de 1 capa escondida con 50 neuronas y con funciones de activación ReLu. Como método de entrenamiento puede usar BP estocástico o en mini-batches, con *weight decay* y tasa de aprendizaje decreciente.

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
```

References

- [1] M. Lichman. UCI Machine Learning Repository. <http://archive.ics.uci.edu/ml>
- [2] Hastie, T.; Tibshirani, R., Friedman, J. (2009), The Elements of Statistical Learning, Second Edition. Springer New York Inc.
- [3] Bishop, Christopher M. (1995). Neural Networks for Pattern Recognition, Clarendon Press.
- [4] Krizhevsky, A., Hinton, G. (2009). Learning multiple layers of features from tiny images.
- [5] Harrison, D. and Rubinfeld, D. (1978). Hedonic prices and the demand for clean air, Journal of Environmental Economics and Management, 5, 81-102
- [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.