

INF-477. REDES NEURONALES ARTIFICIALES.

TAREA 2 - AUTOENCODERS, RBMs Y CONVNETS

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

Temas

- Construcción y entrenamiento de autoencoders (AEs) - usando *keras*.
- Construcción y entrenamiento de RBMs - usando *keras*.
- Apilamiento de AEs y RBMs - usando *keras*.
- Denoising y DR usando AEs y RBMs *keras*.
- Pre-entrenamiento - usando AEs y RBMs *keras*.
- Construcción y entrenamiento de Redes Convolucionales (ConvNets) - usando *keras*.

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.
- Los ítemes identificados con los símbolos ★ (importante) y ★★ (muy importante) tendrán un peso mayor y mucho mayor (respectivamente) en la evaluación.
- 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: Lunes 12 de Octubre.
- Formato de entrega: envío de link Github al correo electrónico del ayudante (joaquin.velasquez@alumnos.inf.utfsm.cl), incluyendo a todos los profesores en copia y especificando asunto: [Taller02-INF477-02-2016].

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

1 Entrenamiento de Autoencoders (AEs) y RBMs en MNIST

Como hemos discutido en clases, las RBM's y posteriormente los AE's, fueron un componente crucial en el desarrollo de los modelos que entre 2006 y 2010 vigorizaron el área de las redes neuronales artificiales con logros notables de desempeño en diferentes tareas de aprendizaje automático.

En esta sección aprenderemos a utilizar estos modelos en tres escenarios clásicos: reducción de dimensionalidad, denoising y pre-entrenamiento. Con este objetivo en mente, utilizaremos un dataset denominado *MNIST*, bastante conocido en el área e introducido por Yann LeCunn hacia 1998 [7] en un trabajo que, junto al Neocognitron de Fukushima [8], se considera uno de los principales antecedentes de las redes convolucionales modernas. Se trata de una colección de 70.000 imágenes de 28×28 pixeles correspondientes a dígitos manuscritos (números entre 0 y 9). En su versión tradicional, la colección se encuentra separada en dos subconjuntos: uno de entrenamiento de 60.000 imágenes y otro de test de 10.000 imágenes. La tarea consiste en entrenar un programa para que aprenda a identificar correctamente el dígito representado en la imagen.

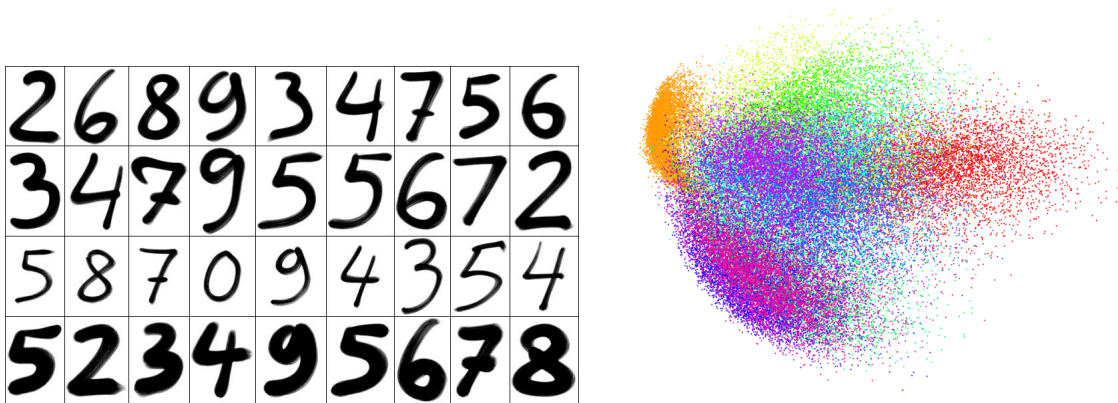


Fig. 1: Dataset MNIST y visualización obtenida usando las primeras dos componentes principales.

- (a) Escriba una función que cargue los datos desde el repositorio de *keras*, normalice las imágenes de modo que los pixeles queden en $[0, 1]$, transforme las imágenes en vectores ($\in \mathbb{R}^{784}$) y devuelva tres subconjuntos disjuntos: uno de entrenamiento, uno de validación y uno de pruebas. La normalización permite interpretar cada valor como una probabilidad de “activación” del un pixel. El conjunto de pruebas será aquel por defecto. Para la construcción del subconjunto de validación su función recibirá un parámetro `NVAL`, cuyo valor por defecto será 1000. El conjunto de validación se construirá utilizando los últimos `NVAL` casos del conjunto del entrenamiento por defecto. El conjunto de entrenamiento consistirá en las primeras $60000 - \text{NVAL}$ imágenes.

```
1 from keras.datasets import mnist
2 import numpy as np
3 (x_train, y_train), (x_test, y_test) = mnist.load_data()
4 x_train = x_train.astype('float32') / 255.
5 x_test = x_test.astype('float32') / 255.
6 x_train = x_train.reshape((len(x_train), np.prod(x_train.shape[1:])))
7 x_test = x_test.reshape((len(x_test), np.prod(x_test.shape[1:])))
8 x_val = x_train[-nval:]
9 y_val = y_train[-nval:]
10 x_train = x_train[:-nval]
11 y_train = y_train[:-nval]
12 Y_train = np_utils.to_categorical(y_train, 10)
```

```

13 Y_val = np_utils.to_categorical(y_val, 10)
14 Y_test = np_utils.to_categorical(y_test, 10)

```

1.1 Reducción de Dimensionalidad

Construir una representación de menor dimensionalidad de un objeto en \mathbb{R}^d , consiste en construir una función $\phi: \mathbb{R}^d \rightarrow \mathbb{R}^{d'}$, con $d' \ll d$ que preserve lo mejor posible la “información” original. Obtener tal representación es útil desde un punto de vista computacional (compresión) y estadístico (permite construir modelos con un menor número de parámetros libres). Una técnica de reducción de dimensionalidad se denomina *no supervisada* cuando no hace uso de información acerca de las clases a las que pertenecen los datos de entrenamiento, marco de trabajo útil cuando dicha información no está disponible. Un AE y una RBM se pueden considerar métodos no-supervisados de reducción de dimensionalidad.

- (a) Entrene un AE básico (1 capa escondida) para generar una representación de MNIST en $d' = 2, 8, 32, 64$ dimensiones. Determine el porcentaje de compresión obtenido y el error de reconstrucción en cada caso. ¿Mejora el resultado si elegimos una función de activación ReLU para el Encoder? ¿Podría utilizarse una ReLU en el decoder?

```

1 from keras.layers import Input, Dense
2 from keras.models import Model
3 from keras.optimizers import SGD
4 input_img = Input(shape=(784,))
5 encoded = Dense(32, activation='sigmoid')(input_img)
6 decoded = Dense(784, activation='sigmoid')(encoded)
7 autoencoder = Model(input=input_img, output=decoded)
8 encoder = Model(input=input_img, output=encoded)
9 encoded_input = Input(shape=(32,))
10 decoder_layer = autoencoder.layers[-1]
11 decoder = Model(input=encoded_input, output=decoder_layer(encoded_input))
12 autoencoder.compile(optimizer=SGD(lr=1.0), loss='binary_crossentropy')
13 autoencoder.fit(x_train,x_train,nb_epoch=50,batch_size=25,shuffle=True,
14                 validation_data=(x_val, x_val))
15 autoencoder.save('basic_autoencoder_768x32.h5')
16 #save other stuff ...

```

- (b) Para verificar la calidad del modelo obtenido, compare visualmente la reconstrucción que logra hacer el autoencoder desde la representación en $\mathbb{R}^{d'}$ para algunas imágenes del conjunto de pruebas. Determine si la percepción visual se corresponde con el error de reconstrucción observado. Comente.

```

1 from keras.models import load_model
2 autoencoder = load_model('basic_autoencoder_768x32.h5')
3 #load other stuff ...
4 encoded_test = encoder.predict(x_test)
5 decoded_test = decoder.predict(encoded_test)
6 import matplotlib
7 n = 10
8 plt.figure(figsize=(20, 4))
9 for i in range(n):
10     ax = plt.subplot(2, n, i + 1)
11     plt.imshow(x_test[i].reshape(28, 28))
12     plt.gray()
13     ax.get_xaxis().set_visible(False)
14     ax.get_yaxis().set_visible(False)
15     ax = plt.subplot(2, n, i + 1 + n)
16     plt.imshow(decoded_test[i].reshape(28, 28))
17     plt.gray()

```

```

18     ax.get_xaxis().set_visible(False)
19     ax.get_yaxis().set_visible(False)
20 plt.show()

```

- (c) Para verificar la calidad de la representación obtenida, implemente el siguiente clasificador, denominado *kNN* (*k*-nearest neighbor): dada una imagen \mathbf{x} , el clasificador busca las $k = 10$ imágenes de entrenamiento más similares $N_{\mathbf{x}} = \{\mathbf{x}^{(k_i)}\}_{i=1}^{10}$ (de acuerdo a una distancia, e.g. euclidiana) y predice como clase, la etiqueta más popular entre las imágenes $N_{\mathbf{x}}$. Mida el error de pruebas obtenido construyendo este clasificador sobre la data original y luego sobre la data reducida. Compare además los tiempos medios de predicción en ambos escenarios.

```

1 encoded_train = encoder.predict(x_train)
2 encoded_test = encoder.predict(x_test)
3 from sklearn.neighbors import KNeighborsClassifier
4 clf = KNeighborsClassifier(10)
5 clf.fit(encoded_train, y_train)
6 clf.fit(encoded_train, y_train)
7 score = clf.score(encoded_test, y_test)
8 print 'Classification Accuracy %.2f' % score

```

- (d) Para verificar la calidad de la representación obtenida, implemente *k-means* (un método básico de agrupamiento) sobre la representación obtenida por el autoencoder. Mida la calidad del agrupamiento obtenido sobre los datos reducidos utilizando la métrica denominada ARI (*Adjusted Rand Index*) y la función de desempeño (que llamaremos *clustering accuracy*) definida el código de ejemplo que se proporciona más abajo. Compare el resultado con el agrupamiento obtenido sobre los datos originales.

```

1 def clustering_accuracy(pred_labels, y, nclusters=10):
2     true_pred = 0.0
3     for i in range(0, nclusters):
4         mvlabel = np.argmax(np.bincount(y[pred_labels==i]))
5         true_pred += sum(y[pred_labels==i] == mvlabel)
6     return true_pred/len(y)
7
8 from sklearn.cluster import KMeans
9 from sklearn import metrics
10 model = KMeans(n_clusters=10)
11 labels_pred = model.fit_predict(encoded_train)
12 score = metrics.adjusted_rand_score(y_train, labels_pred)
13 print 'Clustering ARI %.2f' % score
14 print 'Clustering ACC %.2f' % clustering_accuracy(labels_pred, y_train)

```

- (e) ★ Compare la calidad de la representación reducida obtenida por el autoencoder básico con aquella obtenida vía PCA utilizando el mismo número de dimensiones d' . Considere los 4 criterios que hemos utilizado hasta el momento, i.e., error de reconstrucción, visualización de la reconstrucción, desempeño en clasificación (vía kNN) y desempeño en agrupamiento (vía kMeans). Comente.

```

1 from sklearn.decomposition import PCA
2 from sklearn.neighbors import KNeighborsClassifier
3 pca = PCA(n_components=32)
4 pca.fit(x_train)
5 pca_train = pca.transform(x_train)
6 pca_test = pca.transform(x_test)
7 clf = KNeighborsClassifier(10)
8 clf.fit(pca_train, y_train)
9 score = clf.score(pca_test, y_test)
10 print 'PCA SCORE %.2f' % score

```

- (f) Entrene una RBM binaria básica para generar una representación de MNIST en $d' = 2, 8, 32, 64$ dimensiones. Determine el porcentaje de compresión obtenido y el error de reconstrucción en cada caso. Compare con los resultados obtenidos por el autoencoder utilizando los 3 criterios que hemos utilizado hasta el momento, i.e., error de reconstrucción, desempeño en clasificación (vía kNN) y desempeño en agrupamiento (vía kMeans)[†].

```

1 from sklearn.neural_network import BernoulliRBM
2 import numpy as np
3 import pickle ##to save trained models
4 model = BernoulliRBM(n_components=32, batch_size=25,
5                       learning_rate=0.05, verbose=1, n_iter=50) ##n_components is d'
6 model.fit(x_train) ##Train using persistent Gibbs chains
7 fileo = open('basicRBM.pickle', 'wb')
8 pickle.dump(model, fileo)
9 fileo.close()

```

- (g) ★★ Modifique el autoencoder básico construido en (a) para implementar un *deep autoencoder* (deep AE), es decir, un autoencoder con al menos dos capas ocultas. Demuestre experimentalmente que este autoencoder puede mejorar significativamente la compresión obtenida por PCA utilizando el mismo número de dimensiones d' . Experimente con $d' = 2, 4, 8, 16, 32$ y distintas profundidades ($L = 2, 3, 4$). Considere en esta comparación los 3 criterios que hemos utilizado hasta el momento, i.e., error de reconstrucción, desempeño en clasificación (vía kNN) y desempeño en agrupamiento (vía kMeans). Comente.

```

1 target_dim = 2 #try other and do a nice plot
2 input_img = Input(shape=(784,))
3 encoded1 = Dense(1000, activation='relu')(input_img)
4 encoded2 = Dense(500, activation='relu')(encoded1)
5 encoded3 = Dense(250, activation='relu')(encoded2)
6 encoded4 = Dense(target_dim, activation='relu')(encoded3)
7 decoded4 = Dense(250, activation='relu')(encoded4)
8 decoded3 = Dense(500, activation='relu')(decoded4)
9 decoded2 = Dense(1000, activation='relu')(decoded3)
10 decoded1 = Dense(784, activation='sigmoid')(decoded2)
11 autoencoder = Model(input=input_img, output=decoded1)
12 encoder = Model(input=input_img, output=encoded3)
13 autoencoder.compile(optimizer=SGD(lr=1.0), loss='binary_crossentropy')
14 autoencoder.fit(x_train, x_train, nb_epoch=50, batch_size=25, shuffle=True,
15               validation_data=(x_val, x_val))
16 autoencoder.save('my_autoencoder_768x1000x500x250x2.h5')
17 from sklearn.decomposition import PCA
18 from sklearn.neighbors import KNeighborsClassifier
19 pca = PCA(n_components=target_dim)
20 pca.fit(x_train)

```

- (h) Para el caso $d' = 2$ de los experimentos anteriores, genere un gráfico que muestre la representación aprendida. Con este fin, utilice por ejemplo la herramienta de visualización TSNE disponible en *sklearn*. Compare cualitativamente el resultado con aquel obtenido usando PCA.

```

1 nplot=5000 #warning: mind your memory!
2 encoded_train = encoder.predict(x_train[:nplot])
3 from sklearn.manifold import TSNE
4 model = TSNE(n_components=2, random_state=0)
5 encoded_train = model.fit_transform(encoded_train)

```

[†]El código de ejemplo que se proporciona a continuación requiere la instalación de la librería *sklearn*. En [12] encontrará instrucciones precisas para diferentes sistemas.

```

6  colors={0:'b',1:'g',2:'r',3:'c',4:'m',5:'y',6:'k',7:'orange',8:'darkgreen',9:'maroon'}
7  markers={0:'o',1:'+',2:'v',3:'<',4:'>',5:'^',6:'s',7:'p',8:'*',9:'x'}
8  plt.figure(figsize=(10, 10))
9  for idx in xrange(0,nplot):
10     label = y_train[idx]
11     line = plt.plot(encoded_train[idx][0], encoded_train[idx][1],
12                     color=colors[label], marker=markers[label], markersize=6)
13  pca_train = pca.transform(x_train)
14  encoded_train = pca_train[:nplot]
15  ... #plot PCA

```

- (i) Construya una función que permita visualizar algunos de los pesos aprendidos por las neuronas de la primera capa del autoencoder. Muestre el resultado para las mejores redes conseguidas en los ítems anteriores.
- ji) Estudie como cambian los resultados del modelo construido en (a) si se impone simetría, es decir, si se trabaja con *tied weights*.

1.2 Denoising

Como hemos discutido en clases, un *denoising autoencoder* (dAE) es esencialmente un *autoencoder* entrenado para reconstruir ejemplos parcialmente corruptos. Varios autores han demostrado que mediante esta modificación simple es posible obtener representaciones latentes más robustas y significativas que aquellas obtenidas por un AE básico. En esta sección exploraremos la aplicación más “natural” o “directa” del método.

- (a) Genere artificialmente una versión corrupta de las imágenes en MNIST utilizando el siguiente modelo de ruido (masking noise): si $\mathbf{x} \in \mathbb{R}^d$ es una de las imágenes originales, la versión ruidosa $\tilde{\mathbf{x}}$ se obtiene como $\tilde{\mathbf{x}} = \mathbf{x} \odot \boldsymbol{\xi}$ donde \odot denota el producto de Hadamard (componente a componente) y $\boldsymbol{\xi} \in \mathbb{R}^d$ es un vector aleatorio binario con componentes $\text{Ber}(p)$ independientes.

```

1  from numpy.random import binomial
2  noise_level = 0.1
3  noise_mask = binomial(n=1,p=noise_level,size=x_train.shape)
4  noisy_x_train = x_train*noise_mask
5  noise_mask = binomial(n=1,p=noise_level,size=x_val.shape)
6  noisy_x_val = x_val*noise_mask
7  noise_mask = binomial(n=1,p=noise_level,size=x_test.shape)
8  noisy_x_test = x_test*noise_mask

```

- (b) Entrene un autoencoder para reconstruir las imágenes corruptas generadas en el ítem anterior. Mida el error de reconstrucción y evalúe cualitativamente (visualización de la imagen corrupta y reconstruida) el resultado para un subconjunto representativo de imágenes. Experimente diferentes valores de p en el rango $(0, 1)$.

```

1  # DEFINE YOUR AUTOENCODER AS BEFORE
2  autoencoder.fit(noisy_x_train, x_train, nb_epoch=50, batch_size=25,
3                  shuffle=True, validation_data=(noisy_x_val, x_val))

```

- (c) Genere artificialmente una versión corrupta de las imágenes en MNIST utilizando el siguiente modelo de ruido (Gaussian noise): si $\mathbf{x} \in \mathbb{R}^d$ es una de las imágenes originales, la versión ruidosa $\tilde{\mathbf{x}}$ se obtiene como $\tilde{\mathbf{x}} = \mathbf{x} + \boldsymbol{\xi}$ donde $\boldsymbol{\xi} \in \mathbb{R}^d$ es un vector aleatorio binario con componentes $\mathcal{N}(0, \sigma^2)$ independientes.

```

1  from numpy.random import standard_normal
2  devst = 0.5
3  noise_mask = devst*standard_normal(size=x_train.shape)
4  noisy_x_train = x_train*noise_mask
5  noise_mask = devst*standard_normal(size=x_val.shape)

```

```

6  noisy_x_val = x_val*noise_mask
7  noise_mask = devst*standard_normal(size=x_test.shape)
8  noisy_x_test = x_test*noise_mask

```

- (d) Entrene un autoencoder para reconstruir las imágenes corruptas generadas en el ítem anterior. Mida el error de reconstrucción y evalúe cualitativamente (visualización de la imagen corrupta y reconstruida) el resultado para un subconjunto representativo de imágenes. Experimente diferentes valores de σ .
- (e) Escriba una función que permita visualizar los pesos aprendidos por el dAE y compárelos con aquellos aprendidos por un AE ordinario. ¿Observa diferencias?
- (f) Suponga que su objetivo es aprender una representación de menor dimensionalidad del conjunto de ejemplos. ¿Es posible mejorar los resultados de reconstrucción obtenidos con un AE ordinario entrenándolo con datos artificialmente corruptos? Proyecte un conjunto de experimentos que permita evaluar esta hipótesis. Note que en este caso debe evaluar el AE sobre los datos de prueba no corruptos.

1.3 Pre-entrenamiento

En esta sección utilizaremos los modelos de las secciones anteriores (autoencoders y RBMs) para pre-entrenar redes profundas. Como hemos discutido en clases, el efecto esperado es regularizar el modelo, posicionando el modelo de partida en una buena zona del espacio de parámetros.

- (a) Construya y entrene una red FF para clasificar las imágenes de MNIST. Utilice BP sin ningún tipo de pre-entrenamiento. Para empezar, utilice una arquitectura $768 \times 1000 \times 1000 \times 10$ y funciones de activación sigmoideas. Determine la *accuracy* (fracción de clasificaciones correctas) alcanzada por el modelo en el conjunto de test.

```

1  ## PARAMETERS ...
2  n_hidden_layer1 = 1000
3  activation_layer1 = 'sigmoid'
4  decoder_activation_1 = 'sigmoid'
5  n_hidden_layer2 = 1000
6  activation_layer1 = 'sigmoid'
7  decoder_activation_2 = 'sigmoid'
8  loss_ = 'binary_crossentropy'
9  optimizer_ = SGD(lr=1.0)
10 epochs_ = 50, batch_size_ = 25
11
12 ## Load and preprocess MNIST as usual
13 from keras.datasets import mnist
14 ## HERE YOU NEED: Y_train, Y_val, Y_test produced in 1(a)
15
16 from keras.models import Sequential
17 model = Sequential()
18 model.add(Dense(n_hidden_layer1, activation=activation_layer1, input_shape=(784,)))
19 model.add(Dense(n_hidden_layer2, activation=activation_layer2))
20 model.add(Dense(10, activation='softmax'))
21 model.summary()
22
23 model.compile(optimizer=optimizer_, loss='binary_crossentropy', metrics=['accuracy'])
24 model.fit(x_train, Y_train, nb_epoch=50, batch_size=25,
25         shuffle=True, validation_data=(x_val, Y_val))
26 model.save('ReluNet-768x1000x1000x10-NFT-50epochs.h5') #USEFUL WHEN TRAINING IS SLOW
27 #TRAINING CAN THEN BE RESUMED FROM THIS POINT :-)

```

- (b) Construya y entrene una red neuronal profunda para clasificar las imágenes de MNIST utilizando la arquitectura propuesta en (a) y pre-entrenando los pesos de cada capa mediante un autoencoder básico.

Proceda en modo clásico, es decir, entrenando en modo no-supervisado una capa a la vez y tomando como input de cada nivel la representación (entrenada) obtenida en el nivel anterior. Después del entrenamiento efectúe un entrenamiento supervisado convencional (*finetunning*). Compare los resultados de clasificación sobre el conjunto de pruebas con aquellos obtenidos en (a), sin pre-entrenamiento. Evalúe también los resultados antes del *finetunning*. Comente.

```

1  ## Load and preprocess MNIST as usual
2  from keras.datasets import mnist
3
4  ###AUTOENCODER 1
5  input_img1 = Input(shape=(784,))
6  encoded1 = Dense(n_hidden_layer1,activation=activation_layer1)(input_img1)
7  decoded1 = Dense(784, activation=decoder_activation_1)(encoded1)
8  autoencoder1 = Model(input=input_img1, output=decoded1)
9  encoder1 = Model(input=input_img1, output=encoded1)
10 autoencoder1.compile(optimizer=optimizer_, loss=loss_)
11 autoencoder1.fit(x_train, x_train, nb_epoch=epochs_, batch_size=batch_size_,
12                 shuffle=True, validation_data=(x_val, x_val))
13 encoded_input1 = Input(shape=(n_hidden_layer1,))
14 autoencoder1.save('autoencoder_layer1.h5')
15 encoder1.save('encoder_layer1.h5')
16
17 ###AUTOENCODER 2
18 x_train_encoded1 = encoder1.predict(x_train) #FORWARD PASS DATA THROUGH FIRST ENCODER
19 x_val_encoded1 = encoder1.predict(x_val)
20 x_test_encoded1 = encoder1.predict(x_test)
21
22 input_img2 = Input(shape=(n_hidden_layer1,))
23 encoded2 = Dense(n_hidden_layer2, activation=activation_layer2)(input_img2)
24 decoded2 = Dense(n_hidden_layer2, activation=decoder_activation_2)(encoded2)
25 autoencoder2 = Model(input=input_img2, output=decoded2)
26 encoder2 = Model(input=input_img2, output=encoded2)
27 autoencoder2.compile(optimizer=optimizer_, loss=loss_)
28 autoencoder2.fit(x_train_encoded1,x_train_encoded1,nb_epoch=epochs_,batch_size=batch_size_,
29                 shuffle=True, validation_data=(x_val_encoded1, x_val_encoded1))
30 encoded_input2 = Input(shape=(n_hidden_layer2,))
31 autoencoder2.save('autoencoder_layer2.h5')
32 encoder2.save('encoder_layer2.h5')
33
34 #FINE TUNNING
35 from keras.models import Sequential
36 model = Sequential()
37 model.add(Dense(n_hidden_layer1, activation=activation_layer1, input_shape=(784,)))
38 model.layers[-1].set_weights(autoencoder1.layers[1].get_weights())
39 model.add(Dense(n_hidden_layer2, activation=activation_layer2))
40 model.layers[-1].set_weights(autoencoder2.layers[1].get_weights())
41 model.add(Dense(10, activation='softmax'))
42 model.summary()
43 model.compile(optimizer=optimizer_,loss='binary_crossentropy', metrics=['accuracy'])
44 model.fit(x_train, Y_train,nb_epoch=20, batch_size=25,
45           shuffle=True, validation_data=(x_val, Y_val))
46 model.save('Net-768x1000x1000x10-finetunned.h5')

```

- (c) Construya y entrene una red neuronal profunda para clasificar las imágenes de MNIST utilizando la arquitectura propuesta en (a) y pre-entrenando los pesos de cada capa mediante una RBM binaria

básica. Compare los resultados con aquellos obtenidos en (a) y (b). Comente.

```

1 from sklearn.neural_network import BernoulliRBM
2 rbm1 = BernoulliRBM(n_components=n_hidden_layer1, batch_size=25,
3                     learning_rate=0.05, verbose=1, n_iter=50)
4 rbm1.fit(x_train) ## Train using persistent Gibbs chains
5 encoded_train1 = rbm1.transform(x_train)
6 encoded_val1 = rbm1.transform(x_val)
7 encoded_test1 = rbm1.transform(x_test)
8
9 rbm2 = BernoulliRBM(n_components=n_hidden_layer2, batch_size=25,
10                    learning_rate=0.05, verbose=1, n_iter=50)
11 rbm2.fit(encoded_train1)
12 encoded_train2 = rbm2.transform(encoded_train1)
13 encoded_val2 = rbm2.transform(encoded_val1)
14 encoded_test2 = rbm2.transform(encoded_test1)

```

- (d) Construya y entrene una red neuronal profunda para clasificar las imágenes de MNIST utilizando la arquitectura propuesta en (a) y pre-entrenando los pesos de cada capa mediante un *denoising autoencoder*. Compare los resultados con aquellos obtenidos en (a), (b) y (c). Comente.
- (e) Evalúe el efecto de incorporar un regularizador ℓ_2 y/o ℓ_1 (elija usted) al entrenamiento de la red neuronal final. Comente.
- (f) ★ Repita los experimentos (a)-(d) utilizando funciones de activación *Tanh* y *ReLU*. Comente.

```

1 ## PARAMETERS ...
2 n_hidden_layer1 = 1000
3 activation_layer1 = 'relu'
4 decoder_activation_1 = 'sigmoid'
5 n_hidden_layer2 = 1000
6 activation_layer1 = 'relu'
7 decoder_activation_2 = 'sigmoid'

```

- (g) ★★ Evalúe el efecto de cambiar el número de neuronas ocultas en cada capa del modelo. Por simplicidad y aún si no es la arquitectura óptima para este problema puede fijar el número de capas ocultas a $L = 3$ y experimente, al menos, con un número de neuronas igual a 500, 1000, 2000, 4000.
- (h) ★★ Evalúe el efecto de aumentar la profundidad de 1, 2, 3 niveles en el modelo. Por simplicidad y aún si no es la arquitectura óptima para este problema mantenga fijo el número de neuronas ocultas.

2 Aprendizaje Semi-Supervisado en NORB

Como hemos discutido en clases, uno de los problemas más relevantes a la hora de aplicar técnicas de aprendizaje automático a problemas reales es el requisito de disponer de un gran número de datos etiquetados, es decir, ejemplos para los que se conoce la respuesta deseada del sistema. Un problema de aprendizaje para el que existen pocos datos etiquetados y muchos datos no etiquetados se denomina *semi-supervisado*. En esta sección, utilizaremos la idea de pre-entrenar una red en modo no supervisado para atacar problemas de aprendizaje semi-supervisado. Con este objetivo en mente, trabajaremos con un dataset denominado NORB, introducido en [9] y utilizado en [10], que corresponde a imágenes estéreo de juguetes clasificados en 6 categorías. Se tienen 291.600 ejemplos de entrenamiento y 58.320 ejemplos de pruebas.

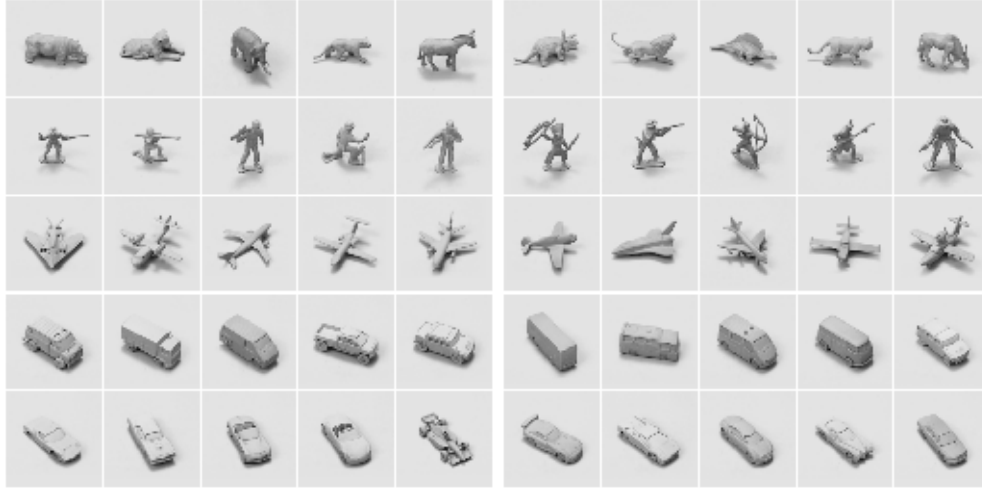


Fig. 2: Dataset NORB.

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.utfsm.cl/~ricky/data_batch_1
2 wget http://octopus.inf.utfsm.cl/~ricky/data_batch_2
3 wget http://octopus.inf.utfsm.cl/~ricky/data_batch_3
4 wget http://octopus.inf.utfsm.cl/~ricky/data_batch_4
5 wget http://octopus.inf.utfsm.cl/~ricky/data_batch_5
6 wget http://octopus.inf.utfsm.cl/~ricky/data_batch_6
7 wget http://octopus.inf.utfsm.cl/~ricky/data_batch_7
8 wget http://octopus.inf.utfsm.cl/~ricky/data_batch_8
9 wget http://octopus.inf.utfsm.cl/~ricky/data_batch_9
10 wget http://octopus.inf.utfsm.cl/~ricky/data_batch_10
11 wget http://octopus.inf.utfsm.cl/~ricky/data_batch_11
12 wget http://octopus.inf.utfsm.cl/~ricky/data_batch_12
13 wget http://octopus.inf.utfsm.cl/~ricky/data_batch_13

```

Los primeros 10 batches corresponden a los datos de entrenamiento y los últimos 2 a los datos de pruebas. 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 $2048 \times n$ (numpy array). Cada columna de esa matriz corresponde a una imagen estéreo de un juguete: los primeros 1024 valores vienen de una de las cámaras/vistas y los siguientes 1024 de la otra. Por otro lado, el elemento (*labels*) del diccionario contiene una lista de n valores enteros entre 0 y 5 que identifican las clases antes a las que pertenecen los juguetes.

- (a) Construya una función que cargue todos los bloques de entrenamiento y pruebas del problema NORB 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 seleccionando 5832 casos de cada batch.

```

1 def load_NORB_train_val(PATH):
2     xtr = []
3     ytr = []
4     xval = []
5     yval = []
6
7     for b in range(1,11):
8         f = os.path.join(PATH, 'data_batch_%d' % (b, ))
9         datadict = unpickle(f)
10        X = datadict['data'].T
11        Y = np.array(datadict['labels'])
12        Z = np.concatenate((X,Y),axis=1)
13        Z = np.random.shuffle(Z)
14        xtr.append(Z[5832:,0:-1])
15        ytr.append(Z[5832:,-1])
16        xval.append(Z[:5832,0:-1])
17        yval.append(Z[:5832,-1])
18
19    Xtr = np.concatenate(xtr)
20    Ytr = np.concatenate(ytr)
21    Xval = np.concatenate(xval)
22    Yval = np.concatenate(yval)
23    del xtr,ytr,xval,yval
24    return Xtr, Ytr, Xval, Yval

```

- (b) Construya una función que escale apropiadamente las imágenes antes de trabajar. Experimente escalando linealmente los datos de tal forma que cada pixel quede en el intervalo $[-1,1]$ con el máximo y mínimo valor observado en los extremos del intervalo. Evalúe más tarde la ventaja de centrar y escalar los datos para que cada atributo (pixel) tenga desviación estándar 1 y media nula.
- (c) Su objetivo será ahora evaluar el desempeño de una red FF en un escenario semi-supervisado. Para ello simulará una situación en la que se tienen n_s ejemplos de entrenamiento para los cuales se conoce la etiqueta correcta y $n_{ns} = n_{tr} - n_s$ ejemplos para los cuales no se tiene esta información (n_{tr} es el número total de ejemplos de entrenamiento). Para empezar, deberá entrenar una red FF con salida softmax para el problema NORB. Considere la inclusión de dos capas escondidas (de 4000 y 2000 unidades) y funciones de activación *relu*. Como parámetros de referencia considere: BP con tasa de aprendizaje constante, función de pérdida *cross-entropy binaria*, y mini batches de tamaño 10. Puede utilizar el conjunto de validación para mejorar el entrenamiento. Construya un gráfico que muestre cómo evoluciona el error de pruebas como función de $\theta_s = n_s/n_{tr}$. Experimente con $\theta_s = 0.1, 0.2, \dots, 1$.
- (d) Su objetivo será ahora construir un gráfico similar al anterior que muestre cómo evoluciona el error de pruebas como función de $\theta_s = n_s/n_{tr}$ cuando la red se pre-entrena utilizando los datos no supervisados. ¿Mejora el resultado con respecto a la red entrenada utilizando sólo los casos para los que se conoce la etiqueta? Experimente pre-entrenando con distintas estrategias (por ejemplo AE's versus dAE's ó AE's versus RBM's).
- (e) Repita el experimento anterior cambiando las funciones de activación a *sigmoidales* y *tanh*.

3 Redes Convolucionales en Google Street View

En esta sección, experimentaremos con redes convolucionales, popularmente conocidas como *CNNs* ó *ConvNets*. Para ello trabajaremos con un dataset denominado *SVHN* (Street View House Numbers) correspondiente a imágenes naturales de dígitos de direcciones de casas obtenidos desde Google Street View. El dataset contiene más de 600.000 imágenes de entrenamiento y 26.032 imágenes de test. Para facilitar la realización de pruebas arquitecturales, el dataset de entrenamiento se divide usualmente en un conjunto pequeño de 73.257 imágenes y un conjunto “extra” de 531.131 imágenes.



Los datos pueden ser obtenidos (en formato Matlab) ejecutando los siguientes comandos

```
1 wget http://ufldl.stanford.edu/housenumbers/train_32x32.mat
2 wget http://ufldl.stanford.edu/housenumbers/extra_32x32.mat
3 wget http://ufldl.stanford.edu/housenumbers/test_32x32.mat
```

- (a) Cargue los datos de entrenamiento y pruebas. Empiece trabajando con el conjunto de entrenamiento más pequeño (“train_32x32.mat”). Determine el tamaño de las imágenes, el número de clases posibles y el número de ejemplos en cada categoría. Finalmente, visualice 5 imágenes de entrenamiento y 5 de test (elegidas aleatoriamente). Comente.

```
1 import scipy.io as sio
2 import numpy as np
3 train_data = sio.loadmat('train_32x32.mat')
4 test_data = sio.loadmat('test_32x32.mat')
5 X_train = train_data['X'].T
6 y_train = train_data['y'] - 1
7 X_test = test_data['X'].T
8 y_test = test_data['y'] - 1
9 X_train = X_train.astype('float32')
10 X_test = X_test.astype('float32')
11 n_classes = len(np.unique(y_train))
12 print np.unique(y_train)
```

- (b) Normalice las imágenes de entrenamiento y pruebas, dividiendo las intensidades originales de pixel en cada canal por 255. Represente adecuadamente la salida deseada de la red de modo de tener un vector de tamaño igual al número de clases. Antes de seguir, se le recomienda también revisar su archivo de configuración *keras.json*, para verificar que el índice correspondiente a los canales de una imagen de entrada sea por defecto el 1 y no el 3.

```

1 from keras.utils import np_utils
2 X_train /= 255
3 X_test /= 255
4 Y_train = np_utils.to_categorical(y_train, n_classes)
5 Y_test = np_utils.to_categorical(y_test, n_classes)

1 vi ~/.keras/keras.json

```

- (c) Defina una CNN con 2 niveles convolucionales (cada una de ellas seguida de una capa de pooling) y 2 capas MLP clásicas completamente conectadas (una oculta y una de salida). Para la primera capa convolucional utilice 16 filtros de 5×5 y para la segunda 512 filtros de 7×7 . Para la capa MLP escondida use 20 neuronas (esta arquitectura, con algunas diferencias, fue una de las primera CNNs entrenadas sobre SVHN y consiguió una accuracy de 94.28% [11]). Genere un esquema lo más compacto posible que muestre los cambios de forma (dimensionalidad) que experimenta un patrón de entrada a medida que se ejecuta un forward-pass.

```

1 from keras.models import Sequential
2 from keras.layers.core import Dense, Dropout, Activation, Flatten
3 from keras.layers.convolutional import Convolution2D, MaxPooling2D, AveragePooling2D
4
5 model = Sequential()
6 model.add(Convolution2D(16, 5, 5, border_mode='same', activation='relu',
7                        input_shape=(n_channels, n_rows, n_cols)))
8 model.add(MaxPooling2D(pool_size=(2, 2)))
9 model.add(Convolution2D(512, 7, 7, border_mode='same', activation='relu'))
10 model.add(MaxPooling2D(pool_size=(2, 2)))
11 model.add(Flatten())
12 model.add(Dense(20, activation='relu'))
13 model.add(Dense(n_classes, activation='softmax'))
14 model.summary()

```

- (d) Entrene la red anterior un máximo de 12 epochs. ¿Logra mejorar o al menos igualar el resultado reportado en la literatura? Si no alcanzó el desempeño mencionado, o solo por diversión, proponga modificaciones a los criterios de entrenamiento y/o cambios arquitectónicos menores (por ejemplo, aumento del número de neuronas de la última capa escondida). *Hint: Es posible superar el 97% de accuracy sin modificar las capas convolucionales.*

```

1 from keras.optimizers import SGD, Adadelta, Adagrad
2 model.compile(loss='binary_crossentropy', optimizer=adagrad, metrics=['accuracy'])
3 model.fit(X_train, Y_train, batch_size=1280, nb_epoch=12, verbose=1, \
4         validation_data=(X_test, Y_test))

```

- (e) Evalúe el efecto de modificar el tamaño de los filtros de las capas convolucionales (y pooling) sobre los tiempos de entrenamiento y el desempeño de la red, transcurridas 10 epochs. Para simplificar, considere el mismo tamaño para todas las capas convolucionales (y el mismo tamaño para todas las capas de pooling). Por ejemplo puede considerar los casos 3×3 , 5×5 , 7×7 , 9×9 para las capas convolucionales y 2×2 , 4×4 para pooling. Estudie además los cambios de forma que experimenta un patrón de entrada a medida que se ejecuta un forward-pass (*Hint: use la función model.summary()*).

```

1 # import staff ...
2 cf_size = 3
3 cp_size = 2
4 model = Sequential()
5 model.add(Convolution2D(16, cf_size, cf_size, border_mode='same', activation='relu',
6                        input_shape=(n_channels, n_rows, n_cols)))
7 model.add(MaxPooling2D(pool_size=(cp_size, cp_size)))

```

```

8 model.add(Convolution2D(512,cf_size,cf_size, border_mode='same',activation='relu'))
9 model.add(MaxPooling2D(pool_size=(cp_size,cp_size)))
10 # compile + train!

```

- (f) Evalúe el efecto de modificar el número de filtros para las capas convolucionales sobre los tiempos de entrenamiento y el desempeño de la red, transcurridas 10 epochs. Para simplificar, puede considerar el mismo tamaño para todas las capas convolucionales.

```

1 # same as before ....
2 n_filters_l1 = 32
3 n_filters_l2 = 64
4 model = Sequential()
5 model.add(Convolution2D(n_filters_l1,cf_size,cf_size,border_mode='same',activation='relu',
6                         input_shape=(n_channels, n_rows, n_cols)))
7 model.add(MaxPooling2D(pool_size=(cp_size,cp_size)))
8 model.add(Convolution2D(n_filters_l2,cf_size,cf_size,border_mode='same',activation='relu'))
9 model.add(MaxPooling2D(pool_size=(cp_size,cp_size)))
10 # compile + train!

```

- (g) Se ha sugerido que la práctica bastante común de continuar una capa convolucional con una capa de pooling puede generar una reducción prematura de las dimensiones del patrón de entrada. Evalúe el efecto de modificar la arquitectura que hemos venido utilizando, re-definiendo un “nivel convolucional” como 2 capas de filtros convolucionales seguidas de una capa de pooling. ¿Se mejora significativamente el resultado? ¿Se observa un aprendizaje más lento en este caso? ¿Observa overfitting?

```

1 model = Sequential()
2 model.add(Convolution2D(32, 3, 3,border_mode='same',activation='relu',input_shape=(3,32,32)))
3 model.add(Convolution2D(32, 3, 3,activation='relu'))
4 model.add(MaxPooling2D(pool_size=(2, 2)))
5 model.add(Convolution2D(256, 3, 3, border_mode='same', activation='relu'))
6 model.add(Convolution2D(256, 3, 3, activation='relu'))
7 model.add(MaxPooling2D(pool_size=(2, 2)))
8 model.add(Flatten())
9 model.add(Dense(20, activation='relu'))
10 model.add(Dense(n_classes, activation='softmax'))
11 # compile + train!

```

- (h) Elija una de las redes entrenadas (preferentemente una con buen desempeño) y visualice los pesos correspondientes a los filtros de la primera capa convolucional. Visualice además el efecto del filtro sobre algunas imágenes de entrenamiento. Comente.

- (i) (Opcional, Bonus +5) Evalúe la conveniencia de regularizar el entrenamiento de la red anterior mediante *Dropout*.

```

1 model = Sequential()
2 model.add(Convolution2D(32, 3, 3,border_mode='same',activation='relu',input_shape=(3,32,32)))
3 model.add(Convolution2D(32, 3, 3,activation='relu'))
4 model.add(MaxPooling2D(pool_size=(2, 2)))
5 model.add(Convolution2D(256, 3, 3, border_mode='same', activation='relu'))
6 model.add(Convolution2D(256, 3, 3, activation='relu'))
7 model.add(MaxPooling2D(pool_size=(2, 2)))
8 model.add(Flatten())
9 model.add(Dense(20, activation='relu'))
10 model.add(Dense(n_classes, activation='softmax'))
11 sgd = SGD(lr=0.1, decay=1e-6, momentum=0.9, nesterov=True)
12 model.compile(loss='binary_crossentropy', optimizer=sgd, metrics=['accuracy'])
13 batch_size = 1280

```

```

14  n_epoch = 20
15  model.fit(X_train, Y_train, batch_size=batch_size, nb_epoch=n_epoch, verbose=1,
16           validation_data=(X_test, Y_test))

```

- (j) (Opcional, Bonus +10) Evalúe la conveniencia de utilizar todo el dataset (“extra_32x32.mat”) en el entrenamiento de la red en (f) ó (d).

References

- [1] Hastie, T.; Tibshirani, R., Friedman, J. (2009), The Elements of Statistical Learning, Second Edition. Springer New York Inc.
- [2] Bishop, Christopher M. (1995). Neural Networks for Pattern Recognition, Clarendon Press.
- [3] Krizhevsky, A., Hinton, G. (2009). Learning multiple layers of features from tiny images.
- [4] Harrison, D. and Rubinfeld, D. (1978). Hedonic prices and the demand for clean air, Journal of Environmental Economics and Management, 5, 81-102
- [5] 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.
- [6] Forsyth, D. A., Ponce, J. (2002). Computer vision: a modern approach. Prentice Hall Professional Technical Reference.
- [7] Yann LeCun, Leon Bottou, Yoshua Bengio, Patrick Haffner. (1998). *Gradient-based Learning Applied to Document Recognition*. Proceedings of the IEEE, 86(11), 2278-2324.
- [8] Kuniyiko Fukushima, Sei Miyake, Takayuki Ito. *Neocognitron: A neural network model for a mechanism of visual pattern recognition*. IEEE Transactions on Systems, Man, and Cybernetics 5 (1983): 826-834.
- [9] Yann LeCun, Fu Jie Huang, and Leon Bottou. *Learning methods for generic object recognition with invariance to pose and lighting*. Proceedings of the 2004 Computer Vision and Pattern Recognition Conference. CVPR 2004. IEEE Computer Society, 2004.
- [10] Xavier Glorot, Antoine Bordes, and Yoshua Bengio. *Deep Sparse Rectifier Neural Networks*. International Conference on Artificial Intelligence and Statistics. 2011.
- [11] Pierre Sermanet, Soumith Chintala, and Yann LeCun. *Convolutional neural networks applied to house numbers digit classification*. International Conference on Pattern Recognition (ICPR), 2012. IEEE, 2012.
- [12] *Scikit-learn: Machine Learning in Python*. <http://scikit-learn.org/stable/>