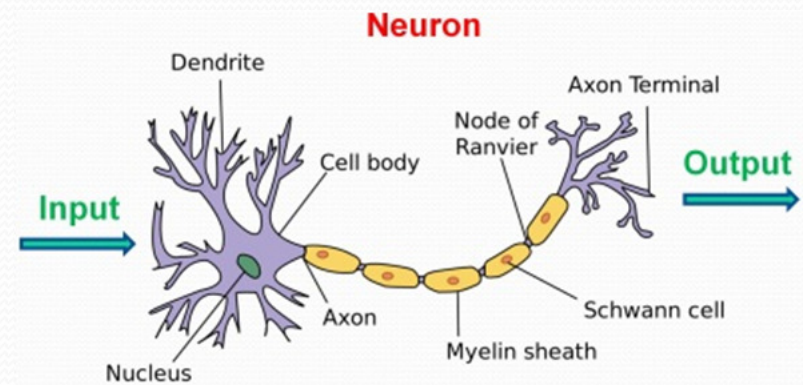# ESS2222

# Lecture 7 – Deep learning and Convolutional Networks

*Hosein Shahnas*
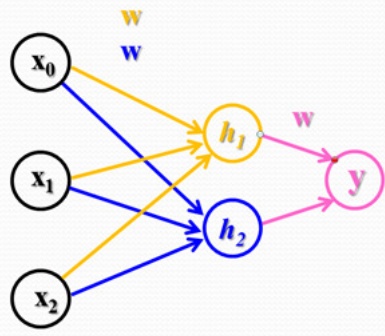
*University of Toronto, Department of Earth Sciences,*

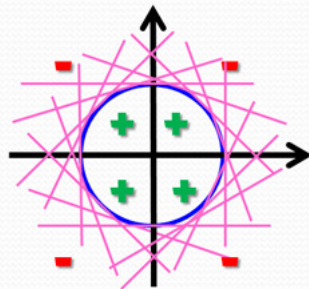# Outline

- ❑ **Deep Learning**
- ❑ **Convolutional Neural Networks**
- ❑ **Recurrent Neural Networks**

# Review of Lecture 6



w
w

w

$h_1$     $h_2$     y
*Or*     *Nand*     *And*

**16-Perceptrons**

## Dropout



**Standard neural network**     **Neural network with dropout**

## Neural Networks



input     output

## Back Propogation

$$\frac{\partial e(w)}{\partial w_{ij}^{(l)}} = x_i^{(l-1)} \, \delta_j^{(l)}$$

$$\delta_j^{(l-1)} = (1 - (x_i^{(l-1)})^2) \quad \sum_j^{d(l)} \delta_i^{(l)} \quad w_{ij}^{(l)}$$

3

# Deep Learning

**Single Neuron Binary Class**



**Shallow Learning Binary Class**



**Shallow Learning Multi Class**



**Multi layer Multi Class**



**Deep neural Network**

4

# Deep Learning
# Large Number of parameters

# Convolutional Networks (CNN)

**W**

$$(I * K)_{xy} = \sum_{i=1}^{h} \sum_{j=1}^{w} K_{ij} \cdot I_{x+i-1,y+j-1}$$

| | | | | | | |
|---|---|---|---|---|---|---|
| 0 | 1 | 1 | 1 | 0 | 0 | 0 |
| 0 | 0 | 1 | 1 | 1 | 0 | 0 |
| 0 | 0 | 0 | 1 | 1 | 1 | 0 |
| 0 | 0 | 0 | 1 | 1 | 0 | 0 |
| 0 | 0 | 1 | 1 | 0 | 0 | 0 |
| 0 | 1 | 1 | 0 | 0 | 0 | 0 |
| 1 | 1 | 0 | 0 | 0 | 0 | 0 |

**I**

**F**

| | | |
|---|---|---|
| 1 | 0 | 1 |
| 0 | 1 | 0 |
| 1 | 0 | 1 |

**K**

$W_c$

| | | | | |
|---|---|---|---|---|
| 1 | 4 | 3 | 4 | 1 |
| 1 | 2 | 4 | 3 | 3 |
| 1 | 2 | 3 | 4 | 1 |
| 1 | 3 | 3 | 1 | 1 |
| 3 | 3 | 1 | 1 | 0 |

**I \* K**

I: Image
K: 'Filter' or 'Kernel' or 'Feature Detector'
I*K: Convolved Feature (activation map)

3 Colour Channels

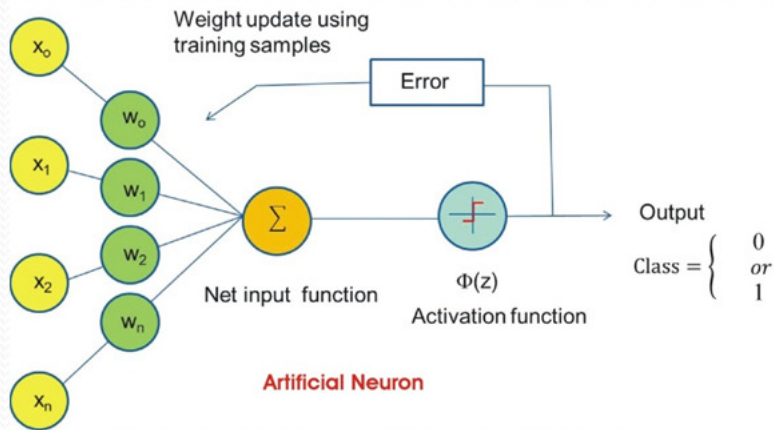| 4 | 6 | 1 | 3 |
|---|---|---|---|
| 0 | 9 | 7 | 3 | 2 |
| 2 | 26 | 35 | 19 | 25 | 6 |
| 1 | 15 | 13 | 22 | 16 | 53 |
| 8 | 4 | 3 | 7 | 10 |
| 0 | 8 | 1 | 3 |

Height: 4 Units (Pixels)

Width: 4 Units (Pixels)

**Volume image**

| | | | | |
|---|---|---|---|---|
| 1 | 1 | 1 | 0 | 0 |
| 0 | 1 | 1 | 1 | 0 |
| 0 | 0 | 1 | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

Image

| 4 | | |
|---|---|---|
| | | |
| | | |

Convolved Feature

$$W_c = \frac{W - F_W}{S_W} + 1$$   The size of the convolved field (I*K)

**S: stride**

**Ex.:** $W_c$ = (7-3)/1+1 = 5

7

# Convolutional Networks (CNN)
# Filters



Pixel representation of filter

Visualization of a curve detector filter

Test Image

Prediction from our network

No idea!?!

# Convolutional Networks (CNN)
## Filters



Identity $\begin{bmatrix} 0 & 0 & 0 \\ 0 & 1 & 0 \\ 0 & 0 & 0 \end{bmatrix}$



Box blur $\frac{1}{9}\begin{bmatrix} 1 & 1 & 1 \\ 1 & 1 & 1 \\ 1 & 1 & 1 \end{bmatrix}$



Edge detection $\begin{bmatrix} -1 & -1 & -1 \\ -1 & 8 & -1 \\ -1 & -1 & -1 \end{bmatrix}$



Sharpen $\begin{bmatrix} 0 & -1 & 0 \\ -1 & 5 & -1 \\ 0 & -1 & 0 \end{bmatrix}$



Gaussian blur $\frac{1}{16}\begin{bmatrix} 1 & 2 & 1 \\ 2 & 4 & 2 \\ 1 & 2 & 1 \end{bmatrix}$

In convolving an image:
1) The outputs shrink
2) The information on corners of the image is lost

This can be prevented by padding.

| 1 $_{\times 1}$ | 1 $_{\times 0}$ | 1 $_{\times 1}$ | 0 | 0 |
|---|---|---|---|---|
| 0 $_{\times 0}$ | 1 $_{\times 1}$ | 1 $_{\times 0}$ | 1 | 0 |
| 0 $_{\times 1}$ | 0 $_{\times 0}$ | 1 $_{\times 1}$ | 1 | 1 |
| 0 | 0 | 1 | 1 | 0 |
| 0 | 1 | 1 | 0 | 0 |

Image

| 4 | | |
|---|---|---|
| | | |
| | | |

Convolved
Feature

36

$32 \times 32 \times 3$

36

**Conv.**

32

32

| 0 | 0 | 0 | 0 | 0 | 0 | ... |
|---|---|---|---|---|---|---|
| 0 | 156 | 155 | 156 | 158 | 158 | ... |
| 0 | 153 | 154 | 157 | 159 | 159 | ... |
| 0 | 149 | 151 | 155 | 158 | 159 | ... |
| 0 | 146 | 146 | 149 | 153 | 158 | ... |
| 0 | 145 | 143 | 143 | 148 | 158 | ... |
| ... | ... | ... | ... | ... | ... | ... |

Input Channel #1 (Red)

| 0 | 0 | 0 | 0 | 0 | 0 | ... |
|---|---|---|---|---|---|---|
| 0 | 167 | 166 | 167 | 169 | 169 | ... |
| 0 | 164 | 165 | 168 | 170 | 170 | ... |
| 0 | 160 | 162 | 166 | 169 | 170 | ... |
| 0 | 156 | 156 | 159 | 163 | 168 | ... |
| 0 | 155 | 153 | 153 | 158 | 168 | ... |
| ... | ... | ... | ... | ... | ... | ... |

Input Channel #2 (Green)

| 0 | 0 | 0 | 0 | 0 | 0 | ... |
|---|---|---|---|---|---|---|
| 0 | 163 | 162 | 163 | 165 | 165 | ... |
| 0 | 160 | 161 | 164 | 166 | 166 | ... |
| 0 | 156 | 158 | 162 | 165 | 166 | ... |
| 0 | 155 | 155 | 158 | 162 | 167 | ... |
| 0 | 154 | 152 | 152 | 157 | 167 | ... |
| ... | ... | ... | ... | ... | ... | ... |

Input Channel #3 (Blue)

| -1 | -1 | 1 |
|---|---|---|
| 0 | 1 | -1 |
| 0 | 1 | 1 |

Kernel Channel #1

| 1 | 0 | 0 |
|---|---|---|
| 1 | -1 | -1 |
| 1 | 0 | -1 |

Kernel Channel #2

| 0 | 1 | 1 |
|---|---|---|
| 0 | 1 | 0 |
| 1 | -1 | 1 |

Kernel Channel #3

$$308 \quad + \quad -498 \quad + \quad 164 \quad + 1 = -25$$

Bias = 1

Output

| -25 | | | | ... |
|---|---|---|---|---|
| | | | | ... |
| | | | | ... |
| | | | | ... |
| ... | ... | ... | ... | ... |

$$H_c = \frac{H - F_H + P}{S_H} + 1$$

$$W_c = \frac{W - F_W + P}{S_W} + 1$$

http://machinelearninguru.com/computer_vision/basics/convolution/convolution_layer.html

# Stride and Padding

S    W        F

P

H



Stride = 1

Stride = 2

# Convolutional Networks
# Max Pooling/Downsampling with CNNs



$$\mathbf{Z} = \mathbf{X}.\mathbf{W} = \sum_1^n w_i\, x_i = w_1\, x_1 + w_2\, x_2 + \ldots + w_n\, x_n$$

$$\Phi_{\text{ReLU}}(z) = max(0, z)$$

Single depth slice

max pool with 2x2 filters and stride 2

Conv → ReLU → Pool → Conv → ReLU → Pool … … … …

Conv → ReLU → Pool → **Classfication**

$$\Phi_{\text{ReLU}}(z) = max\,(0, z)$$
Rectified Linear Unit (ReLU)
$$\Phi_{\text{SReLU}}(z) = log(1 + exp(z))$$
Analytic approx.

# Convolutional Networks (CNN)

$$W$$

$$W_c$$



**RGB image**

**Convolving**



We stack these up to get a "new image" of size 28x28x6!

**Six independent filters**

$$W = 32$$
$$W_c = 28$$



https://medium.com/technologymadeeasy/the-best-explanation-of-convolutional-neural-networks-on-the-internet-fbb8b1ad5df8

# Convolutional Networks (CNN)



**Sequence of CNN layers**



**Feature visualization**



**CNN Architecture**

# Convolutional Networks



**Softmax**

**ReLU**

**ReLU**

**Softmax function**

$$\Phi_{\text{Softmax}}(z_j) = \frac{e^{-zj}}{\sum_{m}^{k} e^{-z_m}}$$

**k: Number of classes**

$$\Phi_{\text{ReLU}}(z) = max(0, z)$$
Rectified Linear Unit (ReLU)
$$\Phi_{\text{SReLU}}(z) = log(1 + exp(z))$$
Analytic approx.

16

## Example 1: MNIST Database - Handwritten digits

**A simple CNN deep learning model for handwritten digits recognition using Keras,**



$3 \times 3$ *Conv*.

$2 \times 2$ *Pool*.

*flatten*

**28**

**10**

**Output**

**128 $\times 0.2$**

**128**

**Fully connected layers**

# A Simple CNN Model

```python
1 #https://towardsdatascience.com/image-classification-in-10-minutes-with-mnist-dataset-54c35b77a38d
2 import tensorflow as tf
3 import sys
4 #======================================================= import data
5 (x_train, y_train), (x_test, y_test) = tf.keras.datasets.mnist.load_data()  # Keras integrated into core Tensorflow
6 print('x_train.shape = ', x_train.shape)
7 print('y_train.shape = ', y_train.shape)
8 #======================================================= import data
9
10 #======================================================= sample image
11 import matplotlib.pyplot as plt
12 #%matplotlib inline # Only use this if using iPython
13 image_index = 7777 # You may select anything up to 60,000
14 print('y_train = ', y_train[image_index]) # The Label is 8
15 plt.imshow(x_train[image_index], cmap='Greys')
16 print('Sample image')
17 #======================================================= sample image
18
19 #======================================================= reshape the image for keras
20 # Reshaping the array to 4-dims so that it can work with the Keras API
21 x_train = x_train.reshape(x_train.shape[0], 28, 28, 1)     # grayscale image
22 x_test = x_test.reshape(x_test.shape[0], 28, 28, 1)       # grayscale image
23 print('x_train.shape[0] = ', x_train.shape[0])
24 print('x_train.shape2 = ', x_train.shape)
25 print('x_test.shape2 = ', x_test.shape)
26
27 input_shape = (28, 28, 1)
28
29 #======================================================= reshape the image for keras
30
```

```python
31 #================================================== change the pixel-values to floating poit
32 # Making sure that the values are float so that we can get decimal points after division
33 x_train = x_train.astype('float32')
34 x_test = x_test.astype('float32')
35 print('x_test[0] = ', x_test[0])
36 #================================================== change the pixel-values to floating poit
37
38 #================================================== normalize to (0.0-1.0)
39 # Normalizing the RGB codes by dividing it to the max RGB value.
40 x_train /= 255
41 x_test /= 255
42 print('x_train shape:', x_train.shape)
43 print('Number of images in x_train', x_train.shape[0])
44 print('Number of images in x_test', x_test.shape[0])
45 print('after normalization: x_test[0] = ', x_test[0])
46 #================================================== normalize to (0.0-1.0)
47
48 #================================================== model
49 # Importing the required Keras modules containing model and layers
50 from keras.models import Sequential
51 from keras.layers import Dense, Conv2D, Dropout, Flatten, MaxPooling2D
52 # Creating a Sequential Model and adding the layers
53 model = Sequential()
54 model.add(Conv2D(28, kernel_size=(3,3), input_shape=input_shape))   # set the number of filters, the size of kernel,
55 model.add(MaxPooling2D(pool_size=(2, 2)))                           # imput shape, pool (size and kind)
56 model.add(Flatten()) # Flattening the 2D arrays for fully connected layers
57 model.add(Dense(128, activation=tf.nn.relu))                       # set the number of the nodes just before the output
58 model.add(Dropout(0.2))                                            # set the amount of droupout
59 model.add(Dense(10,activation=tf.nn.softmax))                      #set the number of nodes for output
60
61 #================================================== model
62
63 #================================================== complie the model and train using training samples
64
65 model.compile(optimizer='adam',
66               loss='sparse_categorical_crossentropy',
67               metrics=['accuracy'])
68 model.fit(x=x_train,y=y_train, epochs=10)                # use 10 epochs for training
69 #================================================== complie the model and train using training samples
```

19

# A Simple CNN Model

```
70
71 #=============================================================== use the test samples for evaluation
72 print('--------------------------------------------------test')
73 print()
74 test_eval = model.evaluate(x_test, y_test)
75 print('test_eval', test_eval)
76 print()
77 print('--------------------------------------------------test')
78 #=============================================================== use the test samples for evaluation
79
80 #=============================================================== plot a sample and find the prediction
81 img_rows = 28
82 img_cols = 28
83 image_index = 4444
84 plt.imshow(x_test[image_index].reshape(28, 28),cmap='Greys')
85 pred = model.predict(x_test[image_index].reshape(1, img_rows, img_cols, 1))
86 print('--------------------------------------------- predit a sample')
87 print()
88 print('prediction = ', pred.argmax())
89 print()
90 print('--------------------------------------------- predit a sample')
91 #=============================================================== plot a sample and find the prediction
92 sys.exit()
```

**Example 2:** **Fashion-MNIST Database**

**CNN deep learning model using Keras,**



Input 28×28×1 → Conv. 32 - 3× 3 → Max Pool 2× 2 → Conv. 64 - 3× 3 → Max Pool 2× 2

Output layer 10 ← Dense layer 128 ← Flatten ← Max Pool 2× 2 ← Conv. 128 - 3× 3

# A Simple CNN Model

```python
1 # https://www.datacamp.com/community/tutorials/convolutional-neural-networks-python
2 import sys
3 #============================================================ import data (Fashion-MNIST)
4 #Load the Data
5 from keras.datasets import fashion_mnist
6 (train_X,train_Y), (test_X,test_Y) = fashion_mnist.load_data()
7 #============================================================ import data (Fashion-MNIST)
8
9 #============================================================ check data
10 print('=======================================1')
11 import numpy as np
12 from keras.utils import to_categorical
13 import matplotlib.pyplot as plt
14 #%matplotlib inline
15 print('Training data shape : ', train_X.shape, train_Y.shape)
16 print('Testing data shape : ', test_X.shape, test_Y.shape)
17
18 # Find the unique numbers from the train labels
19 classes = np.unique(train_Y)
20 nClasses = len(classes)
21 print('Total number of outputs : ', nClasses)
22 print('Output classes : ', classes)
23
24 # ploat
25 plt.figure(figsize=[15,15])
26
27 # Display the first image in training data
28 plt.subplot(121)
29 n = 0
30 plt.imshow(train_X[n,:,:], cmap='gray')
31 plt.title("Class : {}".format(train_Y[n]))
32
33 # Display the first image in testing data
34 plt.subplot(122)
35 n = 5
36 plt.imshow(test_X[n,:,:], cmap='gray')
37 plt.title("Class : {}".format(test_Y[n]))
38 print('=======================================1')
39 #============================================================ check data
40
```

# A Simple CNN Model

```python
40
41 #======================================================== preprocess data
42 print('=========================================2')
43 print('train_X.shape = ', train_X.shape)
44 print('test_X.shape = ', test_X.shape)
45 train_X = train_X.reshape(-1, 28,28, 1)          # reshape
46 test_X = test_X.reshape(-1, 28,28, 1)            # reshape
47
48 print('train_X.shape = ', train_X.shape)
49 print('test_X.shape = ', test_X.shape)
50
51
52 train_X = train_X.astype('float32')
53 test_X = test_X.astype('float32')
54 train_X = train_X / 255.
55 test_X = test_X / 255.
56
57 # Change the labels from categorical to one-hot encoding
58 train_Y_one_hot = to_categorical(train_Y)
59 test_Y_one_hot = to_categorical(test_Y)
60
61 # Display the change for category label using one-hot encoding
62 n = 0
63 print('Original label:', train_Y[n])
64 print('After conversion to one-hot:', train_Y_one_hot[n])
65 n = 7
66 print('Original label:', train_Y[n])
67 print('After conversion to one-hot:', train_Y_one_hot[n])
68 print('=========================================2')
69 #======================================================== preprocess data
70
71 #======================================================== test-train
72 print('=========================================3')
73 print('train_X.shape = ', train_X.shape)
74 print('train_Y_one_hot.shape = ', train_Y_one_hot.shape)
75 print()
76 from sklearn.model_selection import train_test_split
77 train_X,valid_X,train_label,valid_label = train_test_split(train_X, train_Y_one_hot, test_size=0.2, random_state=13)
78
```

23

# A Simple CNN Model

```python
78
79 print('train_X.shape = ', train_X.shape)
80 print('valid_X.shape = ', valid_X.shape)
81 print('train_label.shape = ', train_label.shape)
82 print('valid_label.shape = ', valid_label.shape)
83 print('=========================================3')
84 #======================================================= test-train
85
86 #======================================================= CNN model
87 import keras
88 from keras.models import Sequential,Input,Model
89 from keras.layers import Dense, Dropout, Flatten
90 from keras.layers import Conv2D, MaxPooling2D
91 from keras.layers.normalization import BatchNormalization
92 from keras.layers.advanced_activations import LeakyReLU
93
94
95 batch_size = 64
96 epochs = 20
97 num_classes = 10
98
99 #Neural Network Architecture
100 fashion_model = Sequential()
101 fashion_model.add(Conv2D(32, kernel_size=(3, 3),activation='linear',input_shape=(28,28,1),padding='same'))
102 fashion_model.add(LeakyReLU(alpha=0.1))
103 fashion_model.add(MaxPooling2D((2, 2),padding='same'))
104 #fashion_model.add(Dropout(0.25))
105 fashion_model.add(Conv2D(64, (3, 3), activation='linear',padding='same'))
106 fashion_model.add(LeakyReLU(alpha=0.1))
107 fashion_model.add(MaxPooling2D(pool_size=(2, 2),padding='same'))
108 #fashion_model.add(Dropout(0.25))
109 fashion_model.add(Conv2D(128, (3, 3), activation='linear',padding='same'))
110 fashion_model.add(LeakyReLU(alpha=0.1))
111 fashion_model.add(MaxPooling2D(pool_size=(2, 2),padding='same'))
112 #fashion_model.add(Dropout(0.4))
113 fashion_model.add(Flatten())
114 fashion_model.add(Dense(128, activation='linear'))
115 fashion_model.add(LeakyReLU(alpha=0.1))
116 #fashion_model.add(Dropout(0.3))
117 fashion_model.add(Dense(num_classes, activation='softmax'))
118 #======================================================= CNN model
```

# A Simple CNN Model

```
118#                                                          CNN model
119
120#============================================================= compile the model
121
122 fashion_model.compile(loss=keras.losses.categorical_crossentropy, optimizer=keras.optimizers.Adam(),metrics=['accuracy'])
123 fashion_model.summary()
124#============================================================= compile the model
125
126#============================================================= train the model
127 fashion_train = fashion_model.fit(train_X, train_label, batch_size=batch_size,epochs=epochs,verbose=1,
128                                   validation_data=(valid_X, valid_label))
129#============================================================= train the model
130
131#============================================================= evaluate the model
132 test_eval = fashion_model.evaluate(test_X, test_Y_one_hot, verbose=0)
133 print('Test loss:', test_eval[0])
134 print('Test accuracy:', test_eval[1])
135#============================================================= evaluate the model
136
137#============================================================= ploat the evaluation results
138 accuracy = fashion_train.history['acc']
139 val_accuracy = fashion_train.history['val_acc']
140 loss = fashion_train.history['loss']
141 val_loss = fashion_train.history['val_loss']
142 epochs = range(len(accuracy))
143 plt.plot(epochs, accuracy, 'bo', label='Training accuracy')
144 plt.plot(epochs, val_accuracy, 'b', label='Validation accuracy')
145 plt.title('Training and validation accuracy')
146 plt.legend()
147 plt.figure()
148 plt.plot(epochs, loss, 'bo', label='Training loss')
149 plt.plot(epochs, val_loss, 'b', label='Validation loss')
150 plt.title('Training and validation loss')
151 plt.legend()
152 plt.show()
153#============================================================= ploat the evaluation results
154
155
```

# A Simple CNN Model

**Example 3:** **Fashion-MNIST Database**
CNN deep learning model with **dropout** using Keras,



Input 28×28×1 → Conv. 32 - 3×3 → Max Pool 2×2 → Dropout 0.25 → Conv. 64 - 3×3

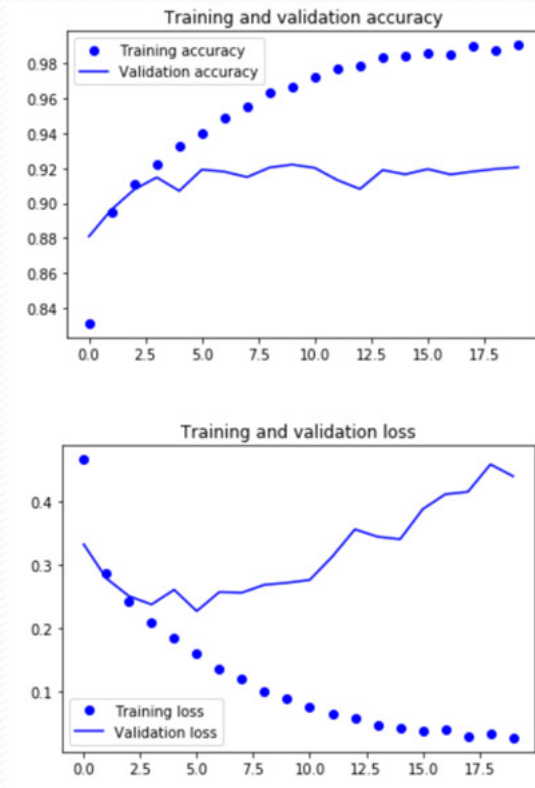Dropout 0.4 ← Max Pool 2×2 ← Conv. 128 - 3×3 ← Dropout 0.25 ← Max Pool 2×2

Flatten → Dense layer 128 → Dropout 0.3 → Output layer 10

# A Simple CNN Model

```
85
86 #========================================================== CNN model with dropout
87 import keras
88 from keras.models import Sequential,Input,Model
89 from keras.layers import Dense, Dropout, Flatten
90 from keras.layers import Conv2D, MaxPooling2D
91 from keras.layers.normalization import BatchNormalization
92 from keras.layers.advanced_activations import LeakyReLU
93
94 # Adding Dropout into the Network
95 batch_size = 64
96 epochs = 20
97 num_classes = 10
98
99 #Neural Network Architecture
100 fashion_model = Sequential()
101 fashion_model.add(Conv2D(32, kernel_size=(3, 3),activation='linear',padding='same',input_shape=(28,28,1)))
102 fashion_model.add(LeakyReLU(alpha=0.1))
103 fashion_model.add(MaxPooling2D((2, 2),padding='same'))
104 fashion_model.add(Dropout(0.25))
105 fashion_model.add(Conv2D(64, (3, 3), activation='linear',padding='same'))
106 fashion_model.add(LeakyReLU(alpha=0.1))
107 fashion_model.add(MaxPooling2D(pool_size=(2, 2),padding='same'))
108 fashion_model.add(Dropout(0.25))
109 fashion_model.add(Conv2D(128, (3, 3), activation='linear',padding='same'))
110 fashion_model.add(LeakyReLU(alpha=0.1))
111 fashion_model.add(MaxPooling2D(pool_size=(2, 2),padding='same'))
112 fashion_model.add(Dropout(0.4))
113 fashion_model.add(Flatten())
114 fashion_model.add(Dense(128, activation='linear'))
115 fashion_model.add(LeakyReLU(alpha=0.1))
116 fashion_model.add(Dropout(0.3))
117 fashion_model.add(Dense(num_classes, activation='softmax'))
118 #========================================================== CNN model with dropout
119
```

# Overcome Overfitting with Dropout

**Fashion MNIST Data**



**Without dropout**



**With dropout**

https://www.datacamp.com/community/tutorials/convolutional-neural-networks-python

**RNN:**

Unlike the regular neural networks in which the samples are assumed to be time independent, being inputted to the network as a whole, the recurrent neural networks, take their inputs from temporally distributed samples. They can be thought as multiple copies of the same network, each passing a message to a successor. These networks have loops in them, allowing time-dependent information to persist.

# Recurrent Neural Networks (RNN)

**Example 1:** Classification of events happening at every frame in a movie. RNN can use its reasoning about the previous events in the film to inform later ones.

**Example 2:** Earthquakes and their relevance to the previous earthquakes.

**Extension:** The application of recurrent neural networks is not limited the time dependent samples. The application can be extended to the other spaces (other than time) as can be seen from the following example.

**Example 3:** MNIST- digit recognition using RNN

In this example each $28 \times 28$-pixel image of handwritten digits, instead of being flattened to a 784-dimensional array as input for a regular neural network, is assumed as 28 one-dimensional images. Each row of pixels of the image is assumed to be a single 1D-image , and the next row another image, following the previous 1D-image. Each row of the image is then sent to one cell of the recurrent neural network in sequence.

The following code trains a model using RNN algorithm. Note that each row has some information to the next row which is not preserved in the regular neural networks.

# MNIST- Digit Recognition - RNN

```
1 #https://github.com/CSCfi/machine-learning-scripts/blob/master/notebooks/keras-mnist-rnn.ipynb
2 #===============================================================
3 #%matplotlib inline
4 from keras.models import Sequential
5 from keras.layers import Dense, Activation, Dropout
6 from keras.layers.recurrent import SimpleRNN, LSTM, GRU
7 from keras.utils import np_utils
8 from keras import backend as K
9
10 from distutils.version import LooseVersion as LV
11 from keras import __version__
12
13 from IPython.display import SVG
14 from keras.utils.vis_utils import model_to_dot
15
16 import numpy as np
17 import matplotlib.pyplot as plt
18 import seaborn as sns
19
20 print('Using Keras version:', __version__, 'backend:', K.backend())
21 assert(LV(__version__) >= LV("2.0.0"))
22 import sys
23 import os
24 import tensorflow as tf
25 #===============================================================
26
27 #=============================================================== import data
28 from keras.datasets import mnist
29 (X_train, y_train), (X_test, y_test) = mnist.load_data()
30 nb_classes = 10
31 img_rows, img_cols = 28, 28
32 my_batch_size = 128
33 epochs = 3
34
35 nb_units = 50 # Number of hidden units to use
36 nb_units = 128 # Number of hidden units to use
37
38 X_train = X_train.astype('float32')
39 X_test = X_test.astype('float32')
```

```
40 X_train /= 255
41 X_test /= 255
42
43 # one-hot encoding:
44 Y_train = np_utils.to_categorical(y_train, nb_classes)
45 Y_test = np_utils.to_categorical(y_test, nb_classes)
46
47 print()
48 print('MNIST data loaded: train:',len(X_train),'test:',len(X_test))
49 print('X_train:', X_train.shape)
50 print('y_train:', y_train.shape)
51 print('Y_train:', Y_train.shape)
52 #================================================================ import data
53
54 #================================================================ RNN model
55 model = Sequential()
56 # Recurrent layers supported: SimpleRNN, LSTM, GRU:
57 model.add(SimpleRNN(nb_units,
58                     input_shape=(img_rows, img_cols)))
59
60 # To stack multiple RNN layers, all RNN layers except the last one need
61 # to have "return_sequences=True".  An example of using two RNN layers:
62 #model.add(SimpleRNN(16,
63 #                    input_shape=(img_rows, img_cols),
64 #                    return_sequences=True))
65 #model.add(SimpleRNN(32))
66
67 model.add(Dense(units=nb_classes))
68 model.add(Activation('softmax'))
69
70 model.compile(loss='categorical_crossentropy',
71               optimizer='adam',
72               metrics=['accuracy'])
73
74 print(model.summary())
75 #================================================================ RNN model
76
```

```
76
77 #==========================================================================
78 SVG(model_to_dot(model, show_shapes=True).create(prog='dot', format='svg'))
79 #==========================================================================
80
81 #=================================================================== train
82 history = model.fit(X_train,
83                     Y_train,
84                     epochs=epochs,
85                     batch_size=my_batch_size,
86                     verbose=2)
87 #=================================================================== train
88
89 #======================================================= plot loss & accuracy
90 plt.figure(figsize=(5,3))
91 plt.plot(history.epoch,history.history['loss'])
92 plt.title('loss')
93
94 plt.figure(figsize=(5,3))
95 plt.plot(history.epoch,history.history['acc'])
96 plt.title('accuracy');
97 #======================================================= plot loss & accuracy
98
99 #=================================================================== scores
100 scores = model.evaluate(X_test, Y_test, verbose=2)
101 print("%s: %.2f%%" % (model.metrics_names[1], scores[1]*100))
102 #=================================================================== scores
103
104 #=============================================================== ploting function
105 def show_failures(predictions, trueclass=None, predictedclass=None, maxtoshow=10):
106     rounded = np.argmax(predictions, axis=1)
107     errors = rounded!=y_test
108     print('Showing max', maxtoshow, 'first failures. '
109           'The predicted class is shown first and the correct class in parenthesis.')
```

```python
110    ii = 0
111    plt.figure(figsize=(maxtoshow, 1))
112    for i in range(X_test.shape[0]):
113        if ii>=maxtoshow:
114            break
115        if errors[i]:
116            if trueclass is not None and y_test[i] != trueclass:
117                continue
118            if predictedclass is not None and predictions[i] != predictedclass:
119                continue
120            plt.subplot(1, maxtoshow, ii+1)
121            plt.axis('off')
122            plt.imshow(X_test[i,:,:], cmap="gray")
123            plt.title("%d (%d)" % (rounded[i], y_test[i]))
124            ii = ii + 1
125 #=================================================================== ploting function
126
127 #=================================================================== preditions for the test samples
128
129 predictions = model.predict(X_test)
130 print('predictions.shape = ', predictions.shape)
131 #=================================================================== preditions for the test samples
132
133 #=================================================================== accuracy
134 save_path = os.path.dirname(os.path.abspath(__file__))
135 print(save_path)
136 y_pred = model.predict(X_test)
137 name_of_file = 'Pred_Obs'
138 completeName = os.path.join(save_path, name_of_file+".dat")
139 file1 = open(completeName, "w")
140
141 acc = 0.0
142 print(y_test[5])
143 print(y_pred[5])
144 print('y_test.shape = ', y_test.shape)
145 print('y_Pred.shape = ', y_pred.shape)
146
```

```python
146
147 class_labels_pred = np.argmax(y_pred, axis=1)
148 print('=======class_labels_pred=========', class_labels_pred)
149 print('=======class_labels_pred.shape=========', class_labels_pred.shape)
150
151 for j in range (y_test.size):
152     Y_o = y_test[j]
153     Y_p = class_labels_pred[j]
154     file1.write("%6.3f  %6.3f  " % (Y_p, Y_o))
155     if (Y_p==Y_o):
156         acc = acc + 1.0
157     file1.write(" \n " )
158 file1.close();
159
160 print('acc = ', acc/y_test.size)
161 #===================================================== accuracy
162
163 #=========================================================
164 '''
165 The first 10 test digits the RNN classified to a wrong class.
166 '''
167 show_failures(predictions)
168 #=========================================================
169
170 #=========================================================
171 '''
172 Failures in which the true class is "6".
173 '''
174 show_failures(predictions, trueclass=6)
175 #=========================================================
176 sys.exit()
177 #=========================================================
```

## Optimizers are a crucial part of the neural networks

**Batch Gradient Descent (BGD)**

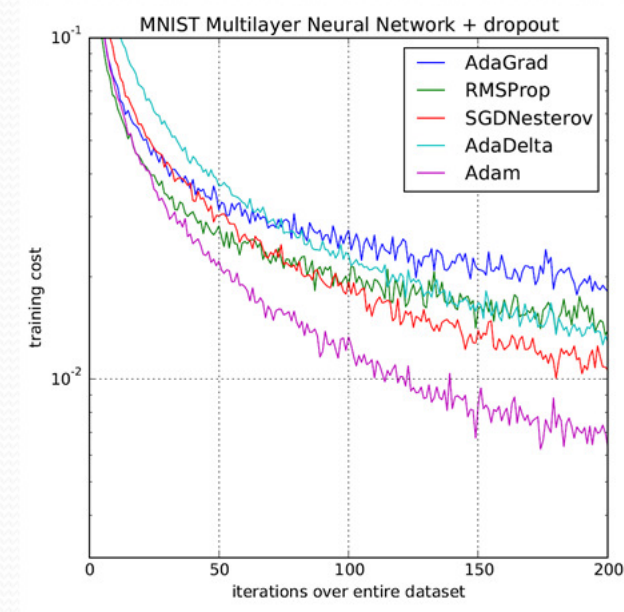$$\Delta w_j = -\eta \frac{\partial J}{\partial w_j} = \eta \sum_i (y^i - \phi(z^i)) x^i_j$$

**Mini-Batch Gradient Descent (BGD)**

$$\Delta w_j = -\eta \frac{\partial J}{\partial w_j} = \eta \sum_{i=1}^{k \ll n} (y^i - \phi(z^i)) x^i_j$$

**Stochastic Gradient Descent (SGD)**

$$\Delta w_j = -\eta \frac{\partial J}{\partial w_j} = \eta (y^i - \phi(z^i)) x^i_j$$

**i: sample, j: feature**



MNIST Multilayer Neural Network + dropout
- AdaGrad
- RMSProp
- SGDNesterov
- AdaDelta
- Adam

training cost / iterations over entire dataset

**AdaGrad**

**AdaGrad (adaptive gradient algorithm) is a modified stochastic gradient descent with per-parameter learning rate (adaptively tuned per parameter).**

$$w_j = w_j - \frac{\eta}{\sqrt{G_{jj}}} g_j, \qquad G_{jj} = \sum_{\tau=1}^{t} g^2_{\tau j}, \qquad g_\tau = \nabla j_i(w)$$

**RMSProp optimization algorithm**

**RMSProp (Root Mean Square Propagation) optimization algorithm (Kingma & Ba, 2015) is an update to the RMSProp optimizer.**

$$J(w) = \frac{1}{2} \sum_i (y^i - \phi(z^i))^2 \quad \text{sample objective function}$$

**w = w + Δw,**

$$s^{t+1}_{dw} = \beta \, s^{t}_{dw} + (1-\beta)\left(\frac{\partial j^t}{\partial w}\right)^2, \qquad\qquad \widehat{s_{dw}} = \frac{s^{t+1}_{dw}}{1 - \beta^{t+1}}$$

$$w^{t+1} = w^t - \eta \frac{\frac{\partial j^t}{\partial w}}{\sqrt{\widehat{s_{dw}}} + \epsilon}$$

**t: time step, i: sample, j: feature**

# Optimization

**Adam optimization algorithm**

Adam (Adaptive moment estimation) optimization algorithm (Kingma & Ba, 2015) is an update to the RMSProp optimizer.

1) Computationally efficient
2) Little memory requirements
3) Well suited for large data

$$v^{t+1}_{dw} = \beta_1 v^t_{dw} + (1-\beta_1)\frac{\partial j^t}{\partial w}, \qquad \widehat{v_{dw}} = \frac{v^{t+1}_{dw}}{1-\beta^{t+1}_1}$$

$$s^{t+1}_{dw} = \beta_2 s^t_{dw} + (1-\beta_2)\left(\frac{\partial j^t}{\partial w}\right)^2, \qquad \widehat{s_{dw}} = \frac{s^{t+1}_{dw}}{1-\beta^{t+1}_2}$$

$$w^{t+1} = w^t - \eta\frac{\widehat{v_{dw}}}{\sqrt{\widehat{s_{dw}}}+\epsilon}$$

# Optimization

**Momentum (Modified SGD)**

**Stochastic gradient descent with momentum remembers the update Δw at each iteration, and determines the next update as a linear combination of the gradient and the previous update.**

$$\Delta w^t = \alpha \Delta w^{t-1} - \eta \, \nabla J(w)^i$$
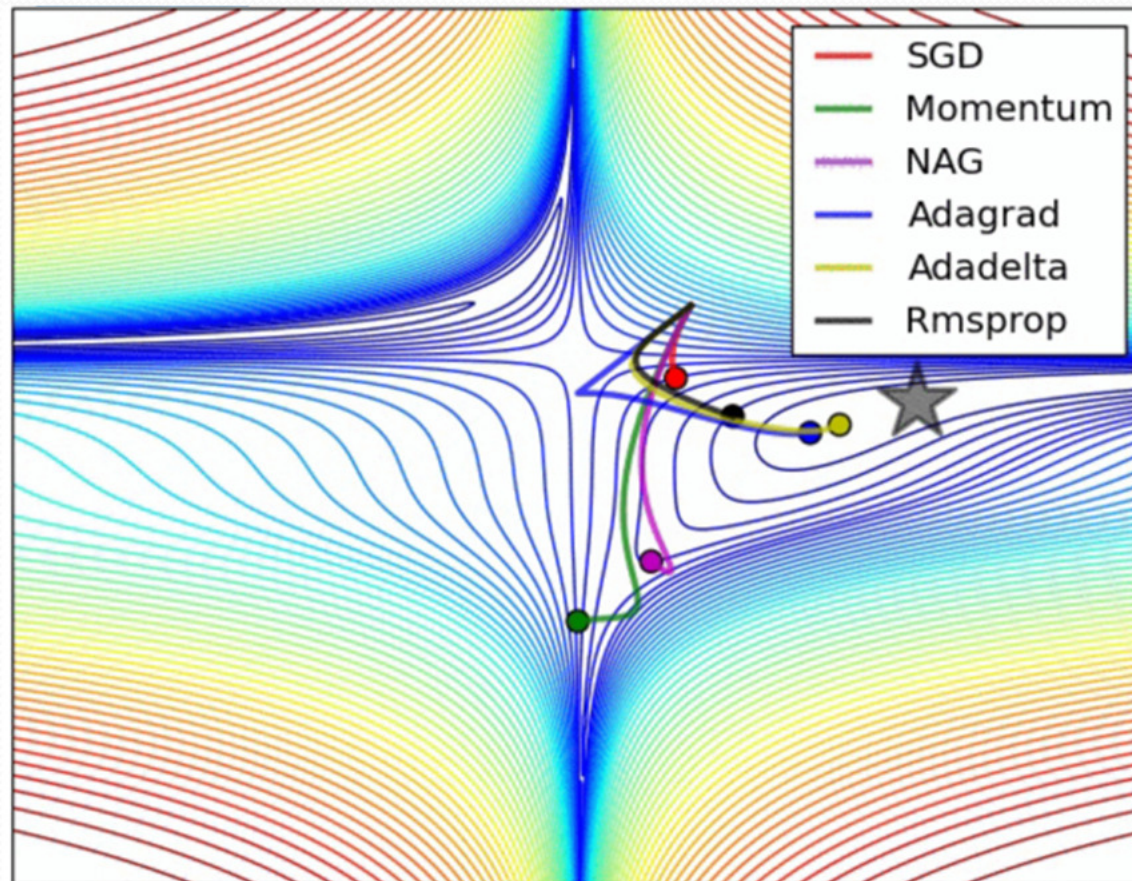$$w^t = w^{t-1} - (\alpha \Delta w - \eta \, \nabla J(w)^i)$$



**SGD**                              **Momentum**

**Implicit Stochastic Gradient Descent (ISGD)**

**SGD is generally sensitive to learning rate η. Fast convergence requires large learning rates but this may induce numerical instability. The problem can be largely solved by considering implicit updates whereby the stochastic gradient is evaluated at the next iterate rather than the current one.**
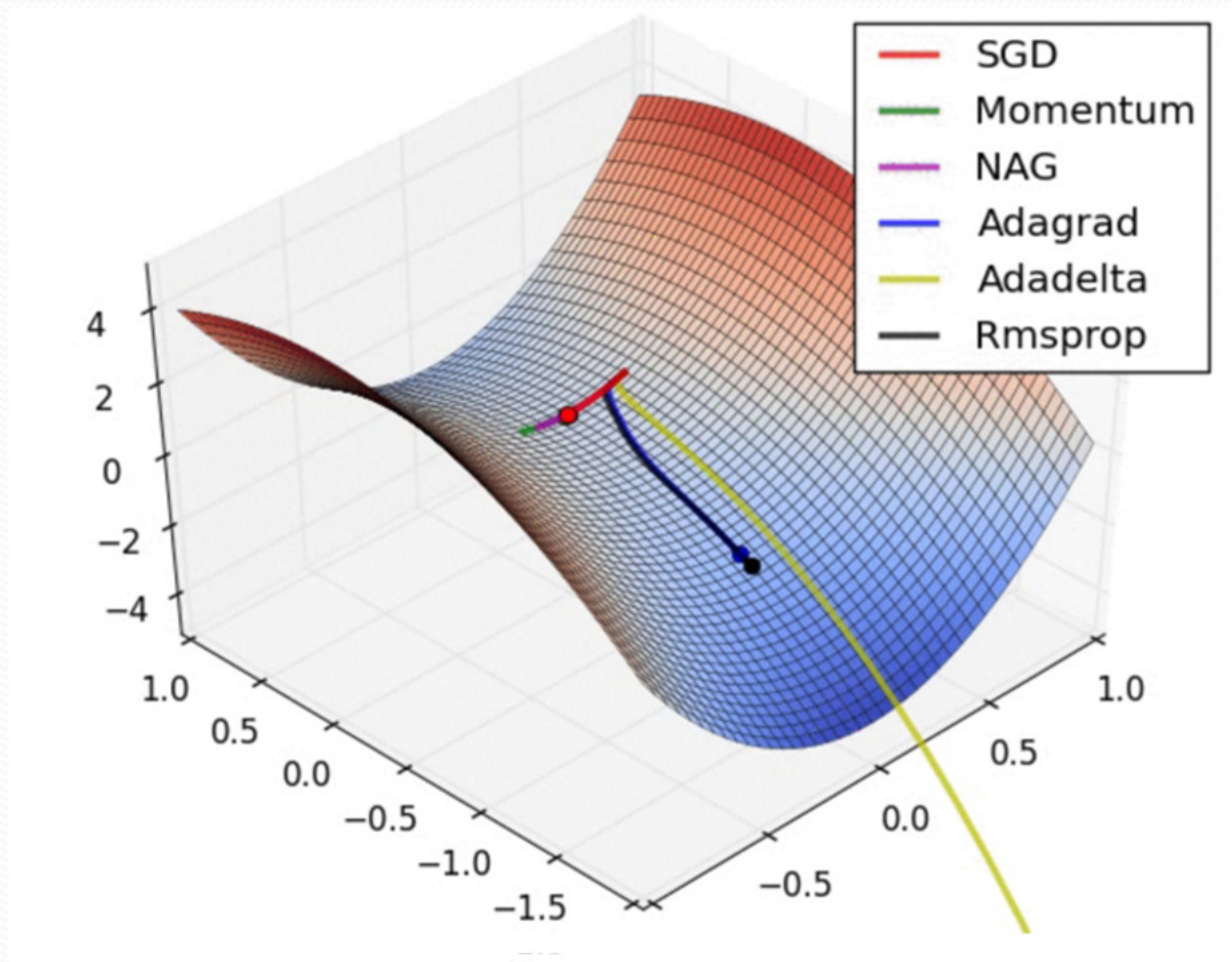
$$w^{new} = w^{old} + \Delta w^{new}$$

# Optimization

# Saddle Point Problem

$$\frac{\partial j}{\partial w} = 0$$