

▼ Measure Feature Map Similarity

This notebook is an enhanced version of a notebook in the Keras examples: [Simple MNIST convnet](#)

Convolutional Neural Network (CNN) architectures use *feature maps* to capture aspects of an image.

Since the set of feature maps is the complete inventory of features of an image found by a CNN, a well-trained model should not have redundant feature maps- the feature maps should all be different. This notebook introduces a measurement of similarity across feature maps with the aim of avoiding redundant feature maps.

We will train a simple CNN against the standard MNIST stroke-digit dataset and will demonstrate how the mean similarity of feature maps slowly drops during training. We will also display feature map activations against an original MNIST image to illuminate how feature map similarity is a good measurement of the quality of a CNN model.

```
1 import random
2 import numpy as np
3 from tensorflow import keras
4 from tensorflow.keras import layers

1 # Model / data parameters
2 num_classes = 10
3 input_shape = (28, 28, 1)
4
5 # the data, split between train and test sets
6 (x_train, y_train), (x_test, y_test) = keras.datasets.mnist.load_data()
7
8 # Scale images to the [0, 1] range
9 x_train = x_train.astype("float32") / 255
10 x_test = x_test.astype("float32") / 255
11 # Make sure images have shape (28, 28, 1)
12 x_train = np.expand_dims(x_train, -1)
13 x_test = np.expand_dims(x_test, -1)
14 # convert class vectors to binary class matrices
15 y_train = keras.utils.to_categorical(y_train, num_classes)
16 y_test = keras.utils.to_categorical(y_test, num_classes)
```

▼ Build the model

The model in the original notebook is broken out into two models:

1. a sub-model which emits the output of the CNN

2. a parent model for training purposes

This allows us to extract feature maps and measure similarity in a callback function.

Remember, a Model is also a Layer.

```
1 cnn_model = keras.Sequential(  
2     [  
3         layers.InputLayer(input_shape=input_shape),  
4         layers.Conv2D(32, kernel_size=(3, 3), activation="relu"),  
5         layers.MaxPooling2D(pool_size=(2, 2)),  
6         layers.Conv2D(64, kernel_size=(3, 3), activation="relu"),  
7     ]  
8     , name='CNN_sub_model'  
9 )  
10 cnn_model.summary()  
11  
12 model = keras.Sequential(  
13     [  
14         layers.InputLayer(input_shape=input_shape),  
15         cnn_model,  
16         layers.MaxPooling2D(pool_size=(2, 2)),  
17         layers.Flatten(),  
18         layers.Dropout(0.5),  
19         layers.Dense(num_classes, activation="softmax"),  
20     ],  
21     name='Parent_model'  
22 )  
23  
24 model.summary()
```

Model: "CNN_sub_model"

Layer (type)	Output Shape	Param #
=====		
conv2d (Conv2D)	(None, 26, 26, 32)	320
=====		
max_pooling2d (MaxPooling2D)	(None, 13, 13, 32)	0
=====		
conv2d_1 (Conv2D)	(None, 11, 11, 64)	18496
=====		
Total params: 18,816		
Trainable params: 18,816		
Non-trainable params: 0		

Model: "Parent_model"

Layer (type)	Output Shape	Param #
=====		
CNN_sub_model (Sequential)	(None, 11, 11, 64)	18816
=====		
max_pooling2d_1 (MaxPooling2D)	(None, 5, 5, 64)	0
=====		
flatten (Flatten)	(None, 1600)	0

dropout (Dropout)	(None, 1600)	0
dense (Dense)	(None, 10)	16010
=====		
Total params: 34,826		
Trainable params: 34,826		
Non-trainable params: 0		

▼ Log image similarities during training

Add a Callback that fetches the final set of feature maps generated by the model, and calculates the average similarity of a random subset of pairs of feature maps.

There are various ways to calculate similarity. This multiplies the two feature maps together and counts the resulting "high" valued cells.

```

1 img_array = x_test[0:1]
2
3 def similarity_multiply(img1, img2):
4     # norm both to 0->1, multiply to produce 0->1
5     min1 = np.min(img1)
6     min2 = np.min(img2)
7     base1 = np.max(img1) - min1
8     base2 = np.max(img2) - min2
9     if base1 == 0:
10         base1 = 0.0001
11     if base2 == 0:
12         base2 = 0.0001
13     norm1 = (img1 - min1) / base1
14     norm2 = (img2 - min2) / base2
15     mult = norm1 * norm2
16     correlated = mult > np.mean(mult)
17
18     percentage = sum(correlated.flatten()) / len(img1.flatten())
19     return percentage
20
21 # While training, capture and log the mean similarity of the feature map pairs.
22 # This network only has 64 fmaps, so it's ok to just check every pair.
23 # This is using when training the complete network, but calls predict()
24 # on the sub-network to fetch the feature maps.
25
26 class LogSimilarities(keras.callbacks.Callback):
27     def __init__(self, cnn_model, img_array, simfunc):
28         super(LogSimilarities, self).__init__()
29         self._cnn_model = cnn_model
30         self._img_array = img_array
31         self._simfunc = simfunc

```

```

31         self._simfunc = simfunc
32
33     def on_epoch_end(self, epoch, logs=None):
34         maps = self._cnn_model.predict(self._img_array)[: , : , : , :]
35         preds = []
36         for i in range(maps.shape[3]):
37             preds.append(maps[0, :, :, i])
38         sims = []
39         for i in range(maps.shape[3]):
40             for j in range(i + 1, maps.shape[3]):
41                 measure = self._simfunc(preds[i], preds[j])
42                 sims.append(measure)
43         avg = sum(sims)/len(sims)
44         if logs:
45             if 'similarity' not in logs:
46                 logs['similarity'] = []
47                 logs['similarity'].append(avg)
48         else:
49             print('Epoch: ' + epoch + ', mean similarity: ' + avg)
50

```

▼ Train the model

```

1 batch_size = 512
2 epochs = 30
3
4 simfunc = similarity_multiply
5 logsim = LogSimilarities(cnn_model, img_array, simfunc=simfunc)
6
7 model.compile(loss="categorical_crossentropy", optimizer="adam", metrics=["accuracy"])
8
9 history = model.fit(x_train, y_train, batch_size=batch_size, epochs=epochs, validation_data=(x_val, y_val),
10                  callbacks=[logsim])

```

```

106/106 [=====] - 3s 8ms/step - loss: 0.7420 - accuracy: 0.0000
Epoch 2/30
106/106 [=====] - 1s 5ms/step - loss: 0.1904 - accuracy: 0.0000
Epoch 3/30
106/106 [=====] - 1s 5ms/step - loss: 0.1375 - accuracy: 0.0000
Epoch 4/30
106/106 [=====] - 1s 6ms/step - loss: 0.1104 - accuracy: 0.0000
Epoch 5/30
106/106 [=====] - 1s 6ms/step - loss: 0.0943 - accuracy: 0.0000
Epoch 6/30
106/106 [=====] - 1s 6ms/step - loss: 0.0852 - accuracy: 0.0000
Epoch 7/30
106/106 [=====] - 1s 6ms/step - loss: 0.0741 - accuracy: 0.0000
Epoch 8/30
106/106 [=====] - 1s 6ms/step - loss: 0.0701 - accuracy: 0.0000
Epoch 9/30
106/106 [=====] - 1s 6ms/step - loss: 0.0641 - accuracy: 0.0000

```

```

Epoch 10/30
106/106 [=====] - 1s 6ms/step - loss: 0.0610 - accura
Epoch 11/30
106/106 [=====] - 1s 5ms/step - loss: 0.0574 - accura
Epoch 12/30
106/106 [=====] - 1s 6ms/step - loss: 0.0551 - accura
Epoch 13/30
106/106 [=====] - 1s 6ms/step - loss: 0.0523 - accura
Epoch 14/30
106/106 [=====] - 1s 5ms/step - loss: 0.0496 - accura
Epoch 15/30
106/106 [=====] - 1s 6ms/step - loss: 0.0486 - accura
Epoch 16/30
106/106 [=====] - 1s 6ms/step - loss: 0.0470 - accura
Epoch 17/30
106/106 [=====] - 1s 6ms/step - loss: 0.0438 - accura
Epoch 18/30
106/106 [=====] - 1s 6ms/step - loss: 0.0432 - accura
Epoch 19/30
106/106 [=====] - 1s 6ms/step - loss: 0.0416 - accura
Epoch 20/30
106/106 [=====] - 1s 6ms/step - loss: 0.0398 - accura
Epoch 21/30
106/106 [=====] - 1s 6ms/step - loss: 0.0397 - accura
Epoch 22/30
106/106 [=====] - 1s 5ms/step - loss: 0.0388 - accura
Epoch 23/30
106/106 [=====] - 1s 6ms/step - loss: 0.0381 - accura
Epoch 24/30
106/106 [=====] - 1s 6ms/step - loss: 0.0368 - accura
Epoch 25/30
106/106 [=====] - 1s 6ms/step - loss: 0.0345 - accura
Epoch 26/30
106/106 [=====] - 1s 6ms/step - loss: 0.0333 - accura
Epoch 27/30
106/106 [=====] - 1s 6ms/step - loss: 0.0333 - accura
Epoch 28/30
106/106 [=====] - 1s 6ms/step - loss: 0.0322 - accura
Epoch 29/30
106/106 [=====] - 1s 6ms/step - loss: 0.0321 - accura
Epoch 30/30

```

▼ Analyze Similarity of Feature Maps

```

1 !pip install --force-reinstall -qq git+https://github.com/LanceNorskog/keract.git
2 import keract
3 from sklearn.preprocessing import MinMaxScaler
4

```

```

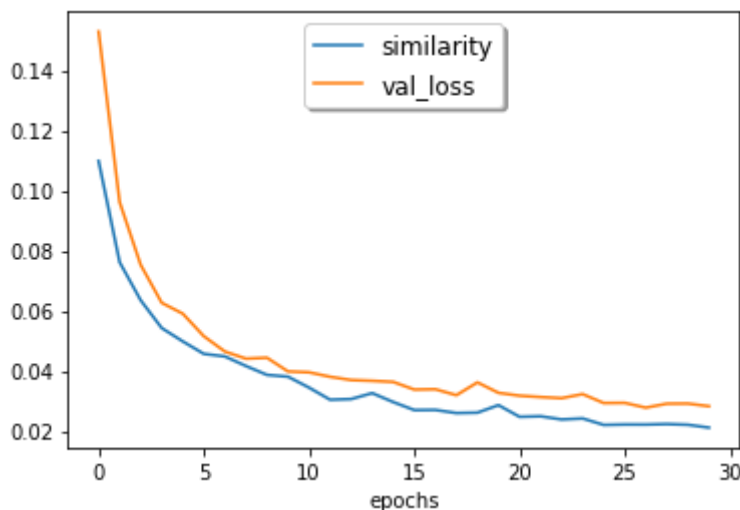
Building wheel for keract (setup.py) ... done
ERROR: tensorflow 2.5.0 has requirement numpy~=1.19.2, but you'll have numpy 1.2
ERROR: datascience 0.10.6 has requirement folium==0.2.1, but you'll have folium
ERROR: alumentations 0.1.12 has requirement imgaug<0.2.7,>=0.2.5, but you'll ha

```

▼ Plot mean similarity over epochs

Let's plot the **similarity** value gathered during training. This is the mean similarity between all pairs of the 64 feature maps generated by the CNN network.

```
1 import matplotlib.pyplot as plt
2
3 def plot_similarity_stats(history):
4     fig, ax = plt.subplots(nrows=1, ncols=1)
5     ax.plot(history['similarity'], label='similarity')
6     ax.plot(history['val_loss'], label='val_loss')
7     legend = ax.legend(loc='upper center', shadow=True, fontsize='large')
8     ax.set(xlabel='epochs', title='')
9
10 plot_similarity_stats(history.history)
```



This chart demonstrates how the drop in similarity tracks the improvement of the CNN (val_loss). CNN feature maps will slowly become decorrelated during a stable training cycle. Also notice how the similarity continues to drop as the network overtrains (val_loss starts increasing).

▼ Visualize the Feature Maps

```
1 def plot_heatmaps(img_array, fmap_i, fmap_j, similarity=None):
2     sim_label = ''
3     if similarity:
4         sim_label = "{:.2f}".format(similarity)
5     feature_maps = np.zeros((1, fmap_i.shape[0], fmap_i.shape[1], 3), dtype='float')
6     feature_maps[0, :, :, 0] = fmap_i
```

```

7     feature_maps[0, :, :, 1] = fmap_j
8     feature_maps[0, :, :, 2] = fmap_i[:, :] * fmap_j[:, :]
9     activations = {sim_label: feature_maps}
10    fig, axes = plt.subplots(1, 3, figsize=(12, 12))
11    keract.display_heatmaps_1(activations, img_array, in_fig=fig, in_axes=axes)

```

Calculate and display similarity over all pairs of feature maps. Keep the image pair with the maximum and minimum similarity.

Phillipe Remy's "Keract" library provides a very handy toolkit for fetching all of the feature maps generate for an image. It also will adorn the original image with data from a feature map to create a "heatmap", which superimposes the feature map onto the original image used to make the prediction.

```

1 maps = cnn_model.predict(img_array)[: , : , : ]
2 preds = []
3 for i in range(maps.shape[3]):
4     preds.append(maps[0, :, :, i])
5 preds = np.asarray(preds)
6 scaler = MinMaxScaler()
7 scaler.fit(preds.reshape(-1, 1))
8
9
10 sims = []
11 x = 0
12 min_i = -1
13 min_j = -1
14 max_i = -1
15 max_j = -1
16 min_sim = 100000
17 max_sim = -1
18
19 for i in range(len(preds)):
20     for j in range(i + 1, len(preds)):
21         measure = simfunc(preds[i], preds[j])
22         top_i = np.max(preds[i])
23         top_j = np.max(preds[j])
24         ratio = np.max([top_i, top_j])/np.min([top_i, top_j])
25         if measure < min_sim and ratio < 3:
26             min_sim = measure
27             min_i = i
28             min_j = j
29         if measure > max_sim:
30             max_sim = measure
31             max_i = i
32             max_j = j
33     sims.append(measure)

```

```

1 activations = {'': maps}
2 fig, axes = plt.subplots(8, 8, figsize=(12, 12))
3 axes[3][2].grid(color='r', linestyle='-', linewidth=2)
4 keract.display_heatmaps_1(activations, img_array, in_fig=fig, in_axes=axes)

```



These 64 heatmaps are "features" or "aspects" of what the CNN notices about the handwritten digit '7'. There are several different measurements of horizontal and diagonal strokes.

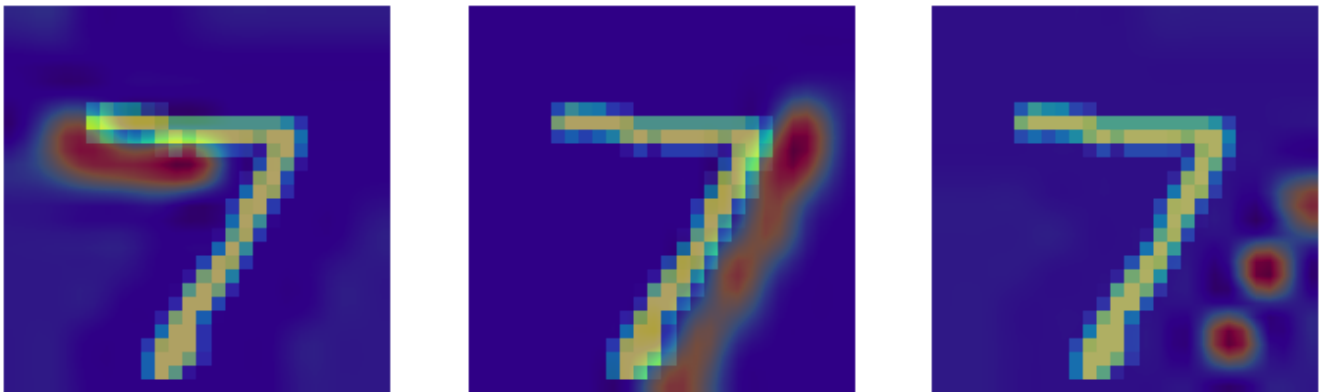
Note:

These images are a great demonstration of the "translation invariance" property of Convolutional Neural Networks. As the image is processed by a stack of Conv2D layers, activations "slide across" the image. Different input images with the same features in different places in the image can activate the same feature map. This is why a feature map might "light up" next to the handwritten stroke rather than on it: the handwritten digits are all roughly the same size, but they are placed differently inside the image. The feature maps pick an "average" placement for a horizontal or diagonal stroke.

▼ Similar Feature Maps

Next we will display the most similar pair of feature maps above, and then multiply the two feature maps together in Hadamard (cell-wise) mode to demonstrate their correlation. The left and middle images are the two feature maps, the rightmost image is the two feature maps multiplied together. This is the core idea of the similarity measure.

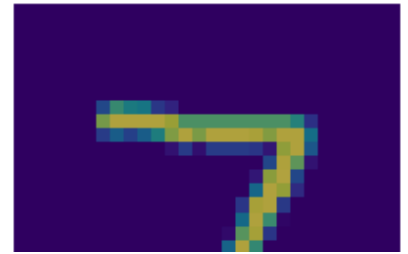
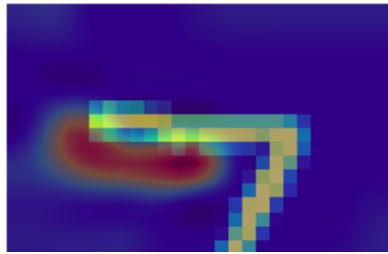
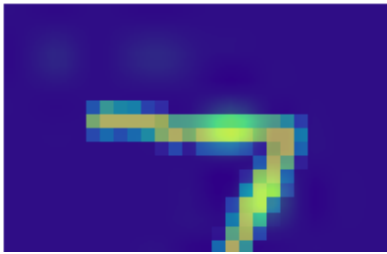
```
1 plot_heatmaps(img_array, preds[max_i], preds[max_j])
```



Here we do the same with the least similar feature maps. Since they have no common areas, there are no activated areas on the right.

The rightmost image has a darker background because high and low activations are exaggerated by multiplying.

```
1 plot_heatmaps(img_array, preds[min_i], preds[min_j])
```



▼ Conclusion



It is clear from this demonstration that feature map similarity is a useful measurement of the quality of a convolutional neural network: as the network improves, the mean similarity will drop.

The reason for this is simple: good feature maps are decorrelated. Feature maps are *independent captures* of features (parts of images) that happen over and over in the input images. If multiple feature maps describe the same feature, then processing power is being wasted. The **descriptive bandwidth** of the feature maps is optimized when no two feature maps describe the same feature.

Based on this insight, it should be possible to improve a CNN by measuring feature map similarity and providing feedback via the loss function. This is the core idea behind Wedge Dropout.