CS 4100: Introduction to AI

Wayne Snyder Northeastern University

Lecture 21: Deep Learning – Generalization; Convolutional Networks



Plan for this Lecture:

- Image Classification: MNIST Digit Classification
- Generalization
- Convolutional Networks for Image Processing

Monday's Lecture:

- Sequence Models
- Recurrent Networks

Binary Classification

Recall: In Binary Classification, our goal is to label the data (e.g., text, images, sounds) with one of 2 labels.



Multiclass and Multilabel Classification

In Multiclass Classification, we have more than 2 labels, and our task is to assign a single label to each sample:

In Multilabel Classification, we have more than 2 labels, and our task is to assign any appropriate labels (not just one):



Let's look at an example of multiclass classification, using the MNIST handwritten digits database....

We will consider the MNIST database of handwritten digits:



The MNIST digit database consists of 70,000 28x28 BW pixel images, stored as 28*28=784 floating-point numbers in the range [0..1]; labels are integers 0..9:

In [36]:





3 0 2 1 1 7 9 0 2 6 7 8 3 9 0 4 6 7 4 6 8 0 7 8 3 1

Note that the array is sparse, it is mostly 0's.

As usual, the database is stored in four parts, split into training and testing sets already:

In [1]: 1 **import** tensorflow as tf 2 from tensorflow import keras 3 from keras import layers 4 from keras.datasets import mnist 5 6 import numpy as np 7 import matplotlib.pyplot as plt 1 (train images, train labels), (test images, test labels) = mnist.load data() In [2]: In [3]: 1 train_images.shape Out[3]: (60000, 28, 28) 1 train_labels.shape In [4]: Out[4]: (60000,) 1 test_images.shape In [5]: Out[5]: (10000, 28, 28) 1 test_labels.shape In [6]: Out[6]: (10000,)

However, just to show the complete process, let's combine training and test sets, and go through the process of creating training, validation, and test sets for data and labels:

```
1 # Just to show, let's put all the data together, shuffle, and separate
In [25]: V
            2 # into training, validation, and test sets.
            3
            4 images = np.concatenate([train_images,test_images])
            5 labels = np.concatenate([train_labels,test_labels])
            6
            7 num_data_instances = len(images)
            8
              shuffled indices = np.random.permutation(num data instances)
            9
                                                                                         Important: randomize the
           10
           11 images = images[shuffled indices]
                                                                                         ordering of your data!
           12 labels = labels[shuffled_indices]
In [11]:
            1 validation_percent = 0.1
            2 testing percent = 0.2
            3
            4 training_end = int( len(images) * (1 - validation_percent - testing_percent ) )
            5 validation end = int( len(images) * (1 - testing percent ) )
            6
            7 training images = images[:training end]
            8 training labels = labels[:training end]
            9
           10 validation images = images[training end:validation end]
           11 validation labels = labels[training end:validation end]
           12
           13 testing_images = images[validation_end:]
           14 testing labels = labels[validation end:]
In [34]:
            1 print(training images.shape)
            2 print(training labels.shape)
            3 print(validation images.shape)
            4 print(validation labels.shape)
            5 print(testing images.shape)
            6 print(testing labels.shape)
         (49000, 784)
         (49000,)
         (7000, 784)
         (7000,)
         (14000, 784)
         (14000,)
```

In [17]: 1 1	<pre>model = keras.Sequential([layers.Dense(512,activation="relu"),</pre>							
3	<pre>layers.Dense(10,activation='softmax')</pre>							
4	1)							
5								
▼ 6	<pre>model.compile(optimizer = 'rmsprop',</pre>							
7	<pre>loss = 'sparse_categorical_cr</pre>							
8	metrics = ['accuracy'])							
9								
▼ 10 11	<pre>nistory = model.fit(training_images, training_labels</pre>							
12	verbese 0							
12	epochs=100							
14	batch size=64							
15	validation data=(validat	ion images.validati	on labels))					
16	parrante (varrante		0102010//					
17	<pre>model.summary()</pre>							
Epoc	h 1/100							
Epoc	n 2/100							
Epoc	h = 3/100							
Epoc	h $\frac{1}{100}$ In [27]:	<pre>1 model.summary()</pre>						
Epoc	h 6/100	Madal, "gamential 1"						
Epoc	h 7/100	Model: sequential_1						
Epoc	h 8/100	Layer (type)	Output Shape	Param #				
Epoc	h 9/100							
Epoc	h 10/100	dense_2 (Dense)	(None, 512)	401920				
Epoc	h 11/100	dense 3 (Dense)	(None, 10)	5130				
Epoc	h 12/100							
Epoc	h 13/100	Total params: 407,050						
Epoc	h 14/100	Trainable params: 407	,050					
Epoc	n 15/100	Non-crainable params:	0					
Epoc.	II 10/100							
Epoc.	n 17/100 h 19/100							
±poc.	h 19/100							
Epoc.	L 20/100							
In [18]: 1	<pre>model.evaluate(testing_images, testing_label</pre>	s,batch_size = 128)						

110/110 [========================] - 0s 2ms/step - loss: 0.2465 - accuracy: 0.9795

Out[18]: [0.24646691977977753, 0.9794999957084656]



Best Validation Accuracy: 0.9813 at epoch 40

Remember: Always consider your baseline!

In [32]:

1 for d in range(10): print("Count of",d,":",list(train labels).count(d))

Count	of	0	:	5923
Count	of	1	:	6742
Count	of	2	:	5958
Count	of	3	:	6131
Count	of	4	:	5842
Count	of	5	:	5421
Count	of	6	:	5918
Count	of	7	:	6265
Count	of	8	:	5851
Count	of	9	:	5949

2

In [37]:

1 model.evaluate(testing_images, randint(10,size=len(testing_labels)), batch_size = 128)

110/110 [==================] - 0s 2ms/step - loss: 4.1812 - accuracy: 0.1021

Out[37]: [4.181239128112793, 0.10207142680883408]

So 98% looks really great compared with random labels!

Generalization – the ability of a NN to learn the patterns in a data set so as to perform well on data it has never seen – is the most important issue in deep learning.

The problem is overfitting – the NN is starting to "memorize" the training set without learning the most important patterns which characterize the essential information present in the data.

Overfitting can be seen when the training loss goes down, but the validation loss goes up. In general, you will see the validation accuracy peak at some epoch and then goes down (generally not as noticably as the rise in the validation loss):









Best Validation Accuracy: 0.9813 at epoch 40

The problem is that a NN can learn ANY data set you give it, essentially by memorizing the exact training set. Here is a dramatic example: we randomly permute the labels, so that there is no correspondence between data and labels. The model continues to "learn" the training set, but the validation accuracy remains around the baseline of 10%.



Best Validation Accuracy: 0.1139 at epoch 0

Digression: This is an issue throughout machine learning. In the IPD experiment in HW 05, agents can learn how to do pretty well in an environment of random agents:



Overfitting is often due to data which is

- Noisy (non-data, ambigious, or outliers)
- o Mislabeled
- Or has rare features or spurious correlations.

Noisy Data in MNIST:

Mislabeled data in MNIST:

Label: 4





Label: 7

Label: 9



Label: 3

Label: 5

Figure 5.5 Robust fit vs. overfitting giving an ambiguous area of the feature space

Overfitting is often due to data which is

- Noisy (non-data, ambigious, or outliers)
- \circ Mislabeled
- Or has rare features or spurious correlations.

Rare features

If your data contains "one-off" features (e.g., a "Getty Images" logo in one image, or a unique or misspelled word in an email), the NN will learn to associate that feature with its label – it is overfitting!

Spurious Correlations

This is actually worse—and more common—than rare features. A word may occur 100s of time in movie reviews, but by a statistical fluke, it occurs in 58% of the positive reviews, and 42% of the negative reviews. The NN will give this word undue weight in learning the data set.

Overfitting is often due to data which is

- Noisy (non-data, ambigious, or outliers)
- \circ Mislabeled
- Or has rare features or spurious correlations.

Rare features

If your data contains "one-off" features (e.g., a "Getty Images" logo in one image, or a unique or misspelled word in an email), the NN will learn to associate that feature with its label – it is overfitting!

Spurious Correlations

This is actually worse—and more common—than rare features. A word may occur 100s of time in movie reviews, but by a statistical fluke, it occurs in 58% of the positive reviews, and 42% of the negative reviews. The NN will give this word undue weight in learning the data set.

Generalization: A Deep Dive into the Math

The Manifold Hypothesis

A manifold in an N dimensional space is a set of points which is isomophic to a lowerdimensional space that is Euclidean, i.e., is continuous and has a notion of "distance."

Ex 1: A curved line is literally in 2 D, but can be mapped 1-to-1 (isomophic) to a 1 D line:



Ex 2: A crumpled piece of paper is 3 D, but is isomophic to a 2D (flat) piece of paper:



Möbius strip



Figure 5.9 Uncrumpling a complicated manifold of data



Generalization: A Deep Dive into the Math

The Manifold Hypothesis is "that many high-dimensional data sets that occur in the real world actually lie along low-dimensional latent manifolds inside that high-dimensional space. As a consequence of the manifold hypothesis, many data sets that appear to initially require many variables to describe, can actually be described by a comparatively small number of variables, likened to the local coordinate system of the underlying manifold. It is suggested that this principle underpins the effectiveness of machine learning algorithms in describing high-dimensional data sets by considering a few common features." -Wikipedia

Your model is searching in a high-dimensional space (= number of parameters attached as weights to neurons) for a representation of the data (lower dimensional manifold). The spaces are continuous and have a notion of distance, which are intrinsic to the gradient descent algorithm:





Figure 5.11 A dense sampling of the input space is necessary in order to learn a model capable of accurate generalization.

Generalization: A Deep Dive into the Math

Before training: the model starts with a random initial state.

Beginning of training: the model gradually moves toward a better fit.

Further training: a robust fit is achieved, transitively, in the process of morphing the model from its initial state to its final state.

Final state: the model overfits the training data, reaching perfect training loss.

Test time: performance of robustly fit model on new data points

Test time: performance of overfit model on new data points

Figure 5.10 Going from a random model to an overfit model, and achieving a robust fit as an intermediate state

Generalization: Underfitting and Overfitting

Overfitting is not a sign that something is wrong with your model, in fact, it shows that your model has sufficient power to represent the patterns that characterize the true "meaning" of the data. You just have to find ways to control this power.

Chollet, p.138: "The first big milestone of a machine learning project: getting a model that has some generalization power (it can beat a trivial baseline) and that is able to overfit." p.141: "Remember that it should always be possible to overfit."

Figure 5.1 Canonical overfitting behavior

Improving generalization can be accomplished by various techniques.

Getting more data, improving your data: more data is almost always better; make sure there are minimal labeling errors, reconsider your data normalization.

If you can not get more data, consider data augmentation: manipulating your existing data in ways that produce different samples with the same essential information.

Improving generalization can be accomplished by various techniques.

Tuning hyperparameters: Play with the hyperparameters, including learning rate and batch size.

Better feature engineering: Use domain knowledge about the data, and experience with the model you are using to better represent the data. Tools can help with feature selection (find out which features are making the most difference).

Example: How to represent a clock face: pixels, clock hands' coordinates, angles:

Improving generalization can be accomplished by various techniques.

Early Stopping: Stop training when a robust fit is achieved. This can often be done automatically by setting a parameter in your model.

Here is a naive example of early stopping, which does not do so well:

[0.11679539084434509, 0.9688571691513062]

Improving generalization can be accomplished by various techniques.

Early Stopping: Stop training when a robust fit is achieved. This can often be done automatically by setting a parameter in your model.

Tuning the early stopping callback results in better results:

Regularization: Various techniques which "actively impede the model's ability to fit perfectly to the training data, with the goal of making the model perform better during validation." The model is simpler, more "regular."

Reduce model size (but not too small):

Figure 5.17 Original model vs. smaller model on IMDB review classification

Figure 5.18 Original model vs. much larger model on IMDB review classification

Regularization: Various techniques which "actively impede the model's ability to fit perfectly to the training data, with the goal of making the model perform better during validation." The model is simpler, more "regular."

Weight Regularization: Place limits on how large the weights in the model can become, so that the model is forced to be simpler (having fewer possibilities of weights). There are two flavors:

- L1 regularization—The cost added is proportional to the *absolute value of the* weight coefficients (the L1 norm of the weights).
- L2 regularization—The cost added is proportional to the square of the value of the weight coefficients (the L2 norm of the weights). L2 regularization is also called weight decay in the context of neural networks. Don't let the different name confuse you: weight decay is mathematically the same as L2 regularization.

Figure 5.19 Effect of L2 weight regularization on validation loss

Regularization: Various techniques which "actively impede the model's ability to fit perfectly to the training data, with the goal of making the model perform better during validation." The model is simpler, more "regular."

Adding Dropout: Dropout is applied to a layer, and is very simple: with some probability p, drop out (as in, setting to 0.0) the outputs from the layer.

0.3	0.2	1.5	0.0	50% dropout	0.0	0.2	1.5	0.0	
0.6	0.1	0.0	0.3		0.6	0.1	0.0	0.3	*0
0.2	1.9	0.3	1.2		0.0	1.9	0.3	0.0	2
0.7	0.5	1.0	0.0		0.7	0.0	0.0	0.0	

Figure 5.20 Dropout applied to an activation matrix at training time, with rescaling happening during training. At test time the activation matrix is unchanged.

This is one of the weirdest great ideas in Deep Learning: it seems like it can't possibly help, but it is one of the most effective and most common ways to regularize your model.

Figure 5.21 Effect of dropout on validation loss

Regularization: Various techniques which "actively impede the model's ability to fit perfectly to the training data, with the goal of making the model perform better during validation." The model is simpler, more "regular."

Reconsidering your choice of architecture: Add more layers, or fewer, or of different widths. Consider starting with wider layers, and getting narrower as you go deeper. Consider different kinds of layers better suited to your data. Google around to see what others have done successfully with similar data.

For example, complex image data is almost always processed using a completely different kind of layer.....

CNNs add convolution and pooling layers to focus on small regions of the data (here, images). Each input to the next layer is calculated as the dot product of the convolution kernel with a small region of the image.

The convolution kernels are moved around the image, perhaps by some skip....

Since data is shared in the region covered by the kernel, various "features" of the image can be recognized in multiple places around the image:

And layers can be stacked.....

And layers can be stacked.....

A pooling layer reduces the dimensionality by averaging (or taking max) of small regions in the previous layer:

It is typical to alternate convolution and pooling layers and end with fully connected layers before output:

Input

Table 14-1. LeNet-5 architecture

Layer	Туре	Maps	Size	Kernel size	Stride	Activation
Out	Fully Connected	-	10	-	-	RBF
F6	Fully Connected	-	84	-	-	tanh
(5	Convolution	120	1×1	5×5	1	tanh
S4	Avg Pooling	16	5×5	2×2	2	tanh
ß	Convolution	16	10 imes10	5 × 5	1	tanh
S2	Avg Pooling	6	14 × 14	2×2	2	tanh
C1	Convolution	6	28 × 28	5×5	1	tanh
In	Input	1	32 × 32	-	-	-

Convolutional networks are built similarly to FF networks, but with different components; you can specify all relevant hyperparameters:

Listing 8.1 Instantiating a small convnet

```
from tensorflow import keras
from tensorflow.keras import layers
inputs = keras.Input(shape=(28, 28, 1))
x = layers.Conv2D(filters=32, kernel_size=3, activation="relu")(inputs)
x = layers.MaxPooling2D(pool_size=2)(x)
x = layers.Conv2D(filters=64, kernel_size=3, activation="relu")(x)
x = layers.MaxPooling2D(pool_size=2)(x)
x = layers.Conv2D(filters=128, kernel_size=3, activation="relu")(x)
x = layers.Flatten()(x)
outputs = layers.Dense(10, activation="softmax")(x)
model = keras.Model(inputs=inputs, outputs=outputs)
```