CS 505: Introduction to Natural Language Processing Wayne Snyder Boston University

Lecture 12: Deep Learning: Evaluating Classifiers; Generalization



Recall: Supervised Machine Learning Workflow



Training involves multiple phases of evaluation with a validation set to find optimal values for the hyperparameters.

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....

Binary Classification

Example: Spam or Not Spam?



	Total	$16,\!545$	$17,\!171$	33,716
	Enron 6	1,500	4,500	6,000
	Enron 5	1,500	$3,\!675$	$5,\!175$
7	Enron 4	1,500	4,500	6,000
-	Enron 3	4,012	1,500	5,512
	Enron 2	4,361	$1,\!496$	$5,\!857$
	Enron 1	$3,\!672$	1,500	$5,\!172$
•	Dataset	N ^o of Legitimate	$\rm N^o of \ Spam$	Total

Labels: 1 =Spam 0 =Not Spam

🗸 🚞 enron1		Today, 3:06 PM	↑ 1.8 MB	Folder
> 🚞 ham	\bigcirc	4/18/20, 1:18 PM	↑ 827 KB	Folder
🛄 README.txt		10/29/05, 4:06 PM	430 bytes	text
> 🚞 spam	\bigcirc	Today, 3:06 PM	↑ 985 KB	Folder
enron1.zip		4/19/20, 2:05 PM	5.4 MB	ZIP archive
enron2.zip		4/19/20, 2:04 PM	7.6 MB	ZIP archive
enron3.zip		4/19/20, 2:04 PM	8.9 MB	ZIP archive
enron4.zip		4/19/20, 2:06 PM	6.7 MB	ZIP archive
enron5.zip		4/19/20, 2:05 PM	6.4 MB	ZIP archive
enron6.zip		4/19/20, 2:05 PM	7.5 MB	ZIP archive

Notice that this data set is fairly well balanced:

51% Spam vs. 49% Not Spam

Multiclass Classification

The well-known MNIST dataset of handwritten digits is a classic multiclass problem.



Multiclass Classification

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]:





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

Multilabel Classification

Multilabel classification is a very important problem in medical diagnosis and object recognition in images:

Enhancing Medical Multi-Label Image Classification Using PyTorch & Lightning





Evaluation of Classifiers: "It's complicated!"

Precise evaluation of classifier performance is complex, even for the simplest case of binary classification:

We have four different outcomes for predictions:









Intuitively, these make good sense:

- Accuracy is the percentage of *all* predictions which were correct;
- Precision is the percentage of *positive* predictions which were correct;
- Sensitivity (or Recall or "true positive rate") is the percentage of the actual positives identified correctly; and
- **Specificity** (or "true negative rate") is the percentage of the actual negatives identified correctly.

The idea of accuracy is fairly straight-forward: how many did it get correct? This is the most common in NLP.

However, the others are useful in many circumstances where the cost of errors may be unacceptable:



- If <u>positive = the patient has cancer</u>, then a FN is disasterous, since you have missed diagnosing a deadly disease, so we want to increase *sensitivity*;
- If <u>positive = the patient does NOT have cancer</u>, then a FP is disasterous, since again you have missed diagnosing a deadly disease, so we want to increase *precision*;
- If you are querying a database of documents, and positive = document retrieved, and you
 want to retrieve all possible documents, then you need to increase *recall* (hence the name);

Precision and recall often in conflict: increasing one will decrease the other; therefore a composite measure is often used, the "harmonic mean" of prediction and recall, which attempts to equalize the error between these two:

 F_1 score = $\frac{2 \times \text{Precision} \times \text{Recall}}{\text{Precision} + \text{Recall}}$

For a very complete listing of metrics for confusion matrices, see: https://en.wikipedia.org/wiki/Precision_and_recall#Definition_(classification_context)

But even accuracy is not that simple!

Consider two classifiers with accuracy metrics as follows:

- A. Accuracy of 4-way classification of blobs: 90%
- B. Accuracy of 10-way classification of digits: 85%

Which is the more accurate classifier?

Consider two classifiers with accuracy metrics as follows:

A. Accuracy of 4-way classification of blobs: 90%

B. Accuracy of 10-way classification of digits: 85%

Which is the more accurate classifier?

Answer 1: 90% > 85% A is more accurate

But wait! A baseline classifier which chooses randomly gets

Alternately: Your baseline model is your model before training, assuming parameters are initialized randomly.

<u>25% accuracy for A</u> and <u>10% accuracy for B</u> !!

[Assuming equal distribution of classes – else it is the percentage of biggest class.] Punchline:

Always compare your model with a suitable baseline model, for example a random model which always chooses the largest class.

Example: Suppose you get the following results:

A. Accuracy of 4-way classification of blobs: 90%

B. Accuracy of 10-way classification of digits: 85%

Ok, how to precisely compare a baseline model?

Answer 2 (most common in ML): Measure improvement above baseline:

A improves on random (25%) by +65%; B improves on random (10%) by +75%, so B is a better model.

Answer 3: Measure improvement normalized to inaccuracy of baseline model:

Cohen's Kappa:
$$\kappa = \frac{\text{accuracy} - \text{baseline accuracy}}{1.0 - \text{baseline accuracy}}$$

Now A wins:
$$\kappa(A) = \frac{0.9 - 0.25}{1.0 - 0.25} = 0.87$$
 $\kappa(B) = \frac{0.85 - 0.1}{1.0 - 0.1} = 0.83$

Which is correct? Well, "it depends"! Statisticians have been arguing about the right way to do this for at least a century, here is a recent paper which compares the principal methods:



$$MCC = \frac{TP \cdot TN - FP \cdot FN}{\sqrt{(TP + FP) \cdot (TP + FN) \cdot (TN + FP) \cdot (TN + FN)}}$$
(worst value = -1; best value = +1)

ML people seem to like the simplicity of Answer 2 : percent improvement over baseline model which always chooses the largest class.

Quick Check:

Suppose you have a spam/not spam binary classification task, but your data set is heavily unbalanced: 30% is spam, and 70% is not spam.

You train two different models, A gets 88% accuracy and B gets 92%.

How much better is B than A?

Quick Check:

Suppose you have a spam/not spam binary classification task, but your data set is heavily unbalanced: 30% is spam, and 70% is not spam.

You train two different models, A gets 88% accuracy and B gets 92%.

How much better is B than A?

A baseline model would simply say that all samples are not spam, and would have 70% accuracy.

Standard answer: B is 22% above the baseline, and A is 18% above.

This doesn't sound that great, compared with 92% for your best model!

You can also see why balanced data sets are preferred!

Evaluation of Multiclass Classifiers (yup, complicated!)

Large Confusion Matrices

If you are doing multiclass classification, you can extend the confusion matrix to have as many rows and columns as the number of classes.

Here is a simple example of a large Confusion Matrix showing how two strings match up (as if each character were a label!):

>>> reference = 'This is the reference data. Testing 123. aoaeoeoe 'Thos iz the rifirenci data. Testeng 123. >>> test = aoaeoeoe >>> print(ConfusionMatrix(reference, test)) The same measures are used as .123T_acdefghinorstz --+-----+ in binary classification: Accuracy = % correct . . .<1>. 3<1>. т<2>. Sensitivity and specificity are<.>.<4>. the average for each row. Recall and NPV are the a h average for each column. 1 . . .<1>. 1 . . . r t _____ (row = reference; col = test)

Evaluation of Multiclass Classifiers (yup, complicated!)

Large Confusion Matrices

These displays can help you to understand your data and how your model is performing with respect to individual samples.

Here is a confusion matrix for the MNIST digit-recognition task:



Confusion Matrix

Classification	report rebu	ilt from o	confusion	matrix:
	precision	recall 1	f1-score	support
0	1.00	0.99	0.99	88
1	0.99	0.97	0.98	91
2	0.99	0.99	0.99	86
3	0.98	0.87	0.92	91
4	0.99	0.96	0.97	92
5	0.95	0.97	0.96	91
6	0.99	0.99	0.99	91
7	0.96	0.99	0.97	89
8	0.94	1.00	0.97	88
9	0.93	0.98	0.95	92
accuracy			0.97	899
macro avg	0.97	0.97	0.97	899
weighted avg	0.97	0.97	0.97	899
5 5				

Generalization and Overfitting

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 goal in developing deep learning models.

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):



Generalization: Overfitting

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

Generalization: Overfitting

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:



Figure 5.3 Mislabeled MNIST training samples



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



Generalization: Overfitting

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, and **it won't generalize well**.

Generalization: A Deep Dive into the Matrix..

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:

Ex 3:

Möbius strip





Figure 5.9 Uncrumpling a complicated manifold of data



Generalization: A Deep Dive into the Matrix..

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 Matrix..

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: performanceTest tof robustly fit modelofon new data pointson

Test time: performance of overfit model on new data points





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 **awesome 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*.









	1
-	



Unfortunately, data augmentation is a little tricky in NLP: how do you create new documents or text that has the same "essential information"?

Typical approaches: replace words by synonyms, translate to another language, then back, etc. (not very satisfying).

Reconsider 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 (this works better with images than in NLP). Google around to see what others have done successfully with similar data.

Tuning hyperparameters: Play with the hyperparameters, including type of optimizer, the learning rate, and the 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).

Most 3

Example: What kind of word vector?

TF? TF-IDF? DIY embeddings?

Glove embeddings? etc.

In some simple cases, we can observe which features were most important; here is a Naive Bayes classifier for identifying gender for names:

Informative	Features	5								
	suffix2	=	'na'	female	:	male	=	100.3	:	1.0
	suffix2	=	'la'	female	:	male	=	77.0	:	1.0
	suffix2	=	'ia'	female	:	male	=	41.3	:	1.0
	suffix3	=	'son'	male	:	female	=	36.5	:	1.0
	suffix2	=	'ld'	male	:	female	=	35.9	:	1.0
	suffix2	=	'sa'	female	:	male	=	34.8	:	1.0
	suffix2	=	'us'	male	:	female	=	29.6	:	1.0
	suffix3	=	'ana'	female	:	male	=	25.7	:	1.0
	suffix3	=	'tta'	female	:	male	=	25.7	:	1.0
	suffix2	=	'ta'	female	:	male	=	25.7	:	1.0
	suffix2	=	'rd'	male	:	female	=	25.2	:	1.0
	suffix2	=	'ra'	female	:	male	=	25.1	:	1.0
	suffix2	=	'do'	male	:	female	=	23.8	:	1.0
	suffix2	=	'rt'	male	:	female	=	22.5	:	1.0
	suffix3	=	'ard'	male	:	female	=	21.0	:	1.0
	suffix2	=	'os'	male	:	female	=	20.4	:	1.0
	suffix3	=	'nne'	female	:	male	=	20.1	:	1.0

Early Stopping: Stop training when a robust fit is achieved. This can often be done automatically by setting a parameter in your model. An elegant method is to save the best model after every K epoches, then refer back after you've gone too far....

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



[0.11679539084434509, 0.9688571691513062]

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, set each parameter in a layer to 0.0:

0.3	0.2	1.5	0.0		0.0	0.2	1.5	0.0	
0.6	0.1	0.0	0.3	50% dropout	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.

```
Listing 5.15 Adding dropout to the IMDB model
model = keras.Sequential([
    layers.Dense(16, activation="relu"),
    layers.Dropout(0.5),
    layers.Dense(16, activation="relu"),
    layers.Dropout(0.5),
    layers.Dense(1, activation="sigmoid")
])
model.compile(optimizer="rmsprop",
        loss="binary_crossentropy",
        metrics=["accuracy"])
history_dropout = model.fit(
        train_data, train_labels,
        epochs=20, batch_size=512, validation_split=0.4)
```



Figure 5.21 Effect of dropout on validation loss