

Course 3, week 2: ML strategy (2)

Error Analysis:

Cleaning up Incorrectly Labeled Data

- what if we have data that is incorrectly classified/labeled?
Is it better while to go in and fix it?

- Deep learning algorithms are quite robust to random errors in the training set.

↳ so long as the incorrectly labeled examples are not too far from random, then it's probably OK to leave the errors as they are and not spending too much time fixing them. The algorithm will be OK so long as the training set is big enough and the percentage of errors is not too high.

- Caution: deep learning algorithms are less robust to systematic errors

↳ if a toaster consistently classifies a white dog as a cat, that will cause some problems for the learning algorithm as it will learn that relation.

- what about incorrectly labeled examples in the **Dev or Test sets?**

↳ do some error analysis where we also count of the # of incorrectly labeled examples in the set.

↳ analyze the misclassified examples by the classifier to understand which ones are truly misclassified and which ones are incorrectly labeled

(By a human, therefore the label is not the correct one).

* if the incorrectly labeled data on the dev/test set impacts our ability to faithfully assess the performance of the model then maybe we should spend a bit of time fixing the labels.

↳ to decide, we can look at the,

i. Overall dev set error: 10% incorrectly labeled percentage

ii. Error due to incorrect labels: $6\% \cdot 10\% = 0.6\%$

iii. Errors due to other causes: 9.4%.

* as in ii, it might make sense to look into the other sources of error than incorrect labels.

• when concerning incorrect dev/test set examples,

↳ after the same process/correction to both the dev/test sets to make sure they continue coming from the same distribution.

↳ consider examining the examples that the algorithm got right as well.

↳ here some examples that we might have gotten right that might also be worth fixing.

↳ if we only fix the ones that we got wrong (By fixing the labels) that will lead to more biased estimate of the error.

↳ the train and dev/test data may come from slightly different distributions.

↳ this is ok as the algorithm is robust to the

Training set.

Build your first model quickly, then iterate

- for almost any machine learning application, there could be 50 different directions we could go in that are reasonable.

The challenge is how do we pick one to focus on.

The best methodology is to pick a first system quickly and then iterate.

↳ common is to use imbalanced dev/test set and metrics selected.

↳ then use loss/accuracy and error analysis to decide where to focus next.

Error Analysis

Training and Testing on Different Distributions

- teams often combine data from different sources to enrich the train set. These additional data sets might not come from the same distribution as the dev/test sets.

Data from webpages



- crawled from the web

- 200k examples

Data from mobile app



- 600k from Victoria

↳ less noisy

↳ less professional

↳ less well framed

- 10k examples

- we want the final system to do well on the

mobile app distribution

↳ but the dataset :)

relatively small

- the dilemma is that we don't only want to use the 10k mobile pictures as it's too small, but the 200k crawled images don't come from the exact distribution that we're trying to predict.

- **Option One:** combine the datasets and randomly shuffle them into a train, dev, test set.

↳ advantages: all our sets will now come from the same test distribution so the training process will be easier to manage

↳ disadvantage (huge): the majority of the examples in the dev set will come from the web page distribution of images (rather than what we care about — the mobile app distribution).

↳ only $\frac{10k}{210k} \approx 5\%$ of dev set that will come

from the mobile app (for the dev set).

↳ we would be spending too much time optimizing on the web set distribution of images in this case.

↳ hence, option one is not recommended

- **Option Two:** have most of the training data be from the web crawler (200k web + 5k mobile) and then the dev and test sets should come from mobile app

exclusively (2.5K mobile each).

↳ advantage: how we are going to be tested accurately.

The dev set data comes from the distribution that we care about and now we are optimizing our algorithm to that set of data

↳ disadvantage: train vs. dev/test distribution how differ

↳ however, this way of splitting will provide better performance over the long term.

Bias and variance with mismatched data distributions

- the way we analyze bias and variance changes when the test set distribution is different from the dev/test set.

- car classifier example

↳ humans $\approx 0\%$ error (Charles' optimal error is nearly 0%)

↳ training error: 1%.

↳ dev error: 10%.

If our sets came from the same distribution, we could say here that we have a large variance problem (our algorithm is not generalizing well from the training set).

In the setting where our distributions come from different places, we can no longer draw this conclusion.

↳ this disparity might just be reflecting that the images in the dev set are harder to classify (e.g. they're of a lower resolution for any other reason).

It's difficult to quantify how much of the additional error is from overfitting the training data and how much is caused by the difference in distribution.

- A way to tease out this difference is to create a new set (called the training-dev set) where the data has the same distribution as the train set, but hasn't been used for training.

- Example (1)

↳ training error: 9%

↳ training - dev error: 8%.

↳ dev error: 10%.

we can conclude a various problem (as the model has overfitted to the training data)

- Example (2)

↳ training error: 1%

↳ training - dev error: 1.5%.

↳ dev error: 10%.

data mismatch problem.

* In data mismatch, we are learning to predict well on the training set distribution, while not doing so well in the dev set (which is composed of a slightly different distribution).

- Example (3)

↳ human level: ~ 0%.

↳ training errors: 10%.

available (i.e.) problem underfitting).

↳ train-dev error: 11%.

↳ dev error: 12%.

- Example (4) (two problems)

- └ human level: ~ 0%] avoidable bias
- └ training error: 10%
- └ train-dev error: 11%] data mismatch
- └ dev error: 20%.

* General principles:

- └ human level] avoidable bias
- └ train error] Variance (overfitting)
- └ train-dev error] data mismatch
- └ dev error] Overfitting to the dev set
(find a bigger dev set)
- └ test error

addressing data mismatch

- If error analysis shows that we have a data mismatch problem, then what can we do?
- We can try to fix this by making the training data more similar to the dev set data or alternatively, we could collect more data that's similar to dev/test sets.

Latent Data Synthesis

Learning from multiple tasks

Transfer Learning

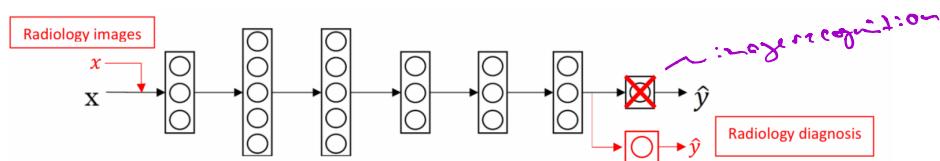
- one of the most powerful ideas from deep learning is that we can sometimes take knowledge that the network has learned

from one task and apply it to a separate task.

Such as training a neural net to recognize cats

and then using the knowledge for X-rays

- if we want to transfer a neural net from image recognition to radiology, we can do that by deleting the output layer of the network, and the weights and biases connecting to it, and then creating a new output layer with randomized weights and biases.



i. train the image recognition network as usual, training the usual parameters through forward pass and backprop.

ii. having trained the parameters, set the output layer (by initializing its weights)

iii. train on the desired data set to which we want to transfer the learned network (through radiology scans — for diagnosis).

↳ if we have a small dataset (for radiology)

we might want to train only the output layer

(the output layer) (this is known as fine tuning)

↳ if we have a lot of data we might want to re-train all the parameters of the network.

↳ if we want to do this, part i. is known

of pre-training because we are using image recognition data to initialize the weights of the network.

* we are taking knowledge learned through image recognition and transferring it to radiology diagnosis.

L this knowledge helps us to originally learn how to detect low level features, such as edges, lines, (conv) or positive objects (from a large image recognition dataset).

L helping the diagnosis learn faster and/or with less data later on.

- Sometimes, instead of only creating a new output node, we can add a few extra layer at the end of our original network to solve a different kind of problem.

- When does transfer learning make sense?

L when we have a lot of data for the data we are transferring from and relatively less data for the problem we are transferring to.

- Transfer learning does not make sense in the cases where you have a few examples in the task you're transferring from and a relatively large # of examples in the task you're transferring to.

L doing it will not work, but it will not gain any meaningful advantage if you do it in this way.

Multitask Learning

- We start off simultaneously trying to have one neural network do several things at the same time and note that the network can generalize to individual tasks (after being trained on different tasks)

- Imagine we are trying to train a neural network that has an output $(4, 1)$ for example (four labels in an image)

↳ we have to define a loss function that can take into account these four outputs

$$\text{Loss} : \hat{y}_j^{(i)} = \frac{1}{m} \cdot \sum_{i=1}^m \sum_{j=1}^4 \mathcal{L}(y_j^{(i)}, \hat{y}_j^{(i)})$$

$$\text{where } \mathcal{L}(y_j^{(i)}, \hat{y}_j^{(i)}) = -y_j^{(i)} \log(\hat{y}_j^{(i)}) - (1-y_j^{(i)}) \log(1-\hat{y}_j^{(i)})$$

* the main difference between this and the original loss function is that we are now summing from $j=1 \dots 4$ for the different outputs. Unlike softmax regression, which assigned a single label to a single example, this one example will have multiple labels (0-4 # classification vs. detection multiple objects in a image (car, person, vehicle)).

↳ multiple objects can appear in the same image.

- By using this cost function, we are looking at a image and trying to solve four problems simultaneously
↳ this could be replaced by four different neural

networks that can recognize one object in the image. The advantage of doing it in a multi-task learning manner is that the low-level features, in the early layers, are often shared by the multiple tasks we're trying to solve).

↳ by having more objects to recognize, we can end up having a better performance than if we were to train the networks separately.

• In multi-task learning, we can train the network even when some images have only a subset of the labels (the presence of a specific object in some examples of the dataset are not labeled at all). We can omit the object that aren't present from the sum of the loss function (removing the summation of outputs term from example 10 to example when the presence of certain objects in an image is a question mark).

• When does multi-task learning make sense?

↳ when we train on a set of tasks that could benefit from having shared low-level features).

Usually: the amount of data we get for each task is similar.

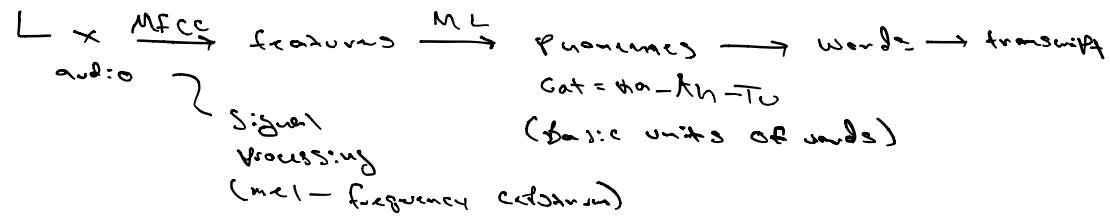
↳ can train a big enough neural network to do well on all the tasks.

End-to-End Deep Learning

What is end-to-end deep learning?

- End-to-end deep learning encapsulates all the intermediate processing steps into a single neural network to solve a specific problem.

↳ take for example traditional speech recognition



↳ in end-to-end deep learning, we skip all these intermediate stages and go directly from audio to the transcript.

End-to-end deep learning works in the cases where we have an abundance of data. For smaller datasets, the traditional methodology works better.

When to use end-to-end deep learning

- pros and cons of end-to-end deep learning

↳ pros

↳ let the data speak

↳ the mapping from $x \rightarrow y$ will be figured out by the neural network. The model will reflect less human pre-conceptions.

↳ less hand-designing of components needed

↳ simpler design workflow

↳ Cons

- ↳ may need a large amount of data to learn the appropriate mapping.
- ↳ excludes potentially useful hand designed components
 - ↳ there's an inability to inject manual knowledge into the model in end-to-end deep learning.