# Error Analysis

Carrying out error analysis:

Error Analysis - manually examining the mislabelled data by a ML model. We can use this to analyze which idea is worth our time.

Consider an example where the cat classifier is classifying certain dogs as cats.

→ Get ~100 mislabeled dev set examples

→ Count up how many are dogs.

i) If 5/100 are dogs, Then error will only reduce from 10% → 9.5%.

∴ Not worth your time

ii) If 50/100 are dogs, Then error will reduce from 10% → 5%.

∴ Worth your time

Evaluating multiple ideas in parallel:
Ideas for cat detection,

→ Fix pics of dogs being recognised as cats

→ Fix great cats (lions, panthers, etc) being misrecognised

→ Improve performance on blurry images

→ Fix instagram filters that are ruining recognition

| Image | Dog | Great Cats | Blurry | Insta | Comments |
|-------|-----|-----------|--------|-------|----------|
| 1 | ✓ | | | ✓ | Pitbull |
| 2 | | | ✓ | ✓ | |
| 3 | | ✓ | ✓ | | Rainy day at zoo |
| ⋮ | ⋮ | ⋮ | ⋮ | ⋮ | |
| % of total | 8% | (43%) | (61%) | 12% | |

These are worth working on ↖

(leaning up incorrectly labeled examples, is it worth the time?

→ DL algorithms are quite robust to random errors (if the data is mislabelled randomly - like one white dog) in the training set but not to systematic errors (like all white dogs labelled as cats - the algo will learn that white dogs are cats).

→ Suppose we add a column 'incorrectly labeled' to this _____ and it has 6%. Then,

*after reducing other cause*

| | | |
|---|---|---|
| Overall dev set error .... | 10 % | 2% |
| Errors due to incorrect labels... | 0.6% | 0.6% |
| Errors due to other cause ... | 9.4% ↻ | 1.4% |

Its more worthwhile to solve this

Since it is 30% of the error, it is worthwhile to solve

correcting incorrect dev/test set examples.

→ Apply same process to your dev and test sets to make sure they continue to come from the same distribution

→ Consider examining examples your algorithm got right as well as ones it got wrong

→ Train and dev/test data may now come from slightly different distributions.

Build your first system quickly, then iterate

There are many ideas you can work on like in speech recognition - Noisy background, accent, far from microphone, etc.

A good strategy is to build something quick and then iterate:

  → Set up dev/test set and metric

  → Build initial system quickly

  → Use bias/variance analysis & error analysis to prioritize next steps.

# Mismatched training and dev/test set

Training and testing on different distributions,

Consider cat app example:

(good quality)
Data from webpages

$\approx 200k$

(blurry - bad quality)
Data from mobile phone

$\approx 10K$

↑

This is what user upload and we care about this

Then we split is as:

| Train | Dev | Test |
|---|---|---|
| 200K web + 5k app | 2.5 k app | 2.5K app |

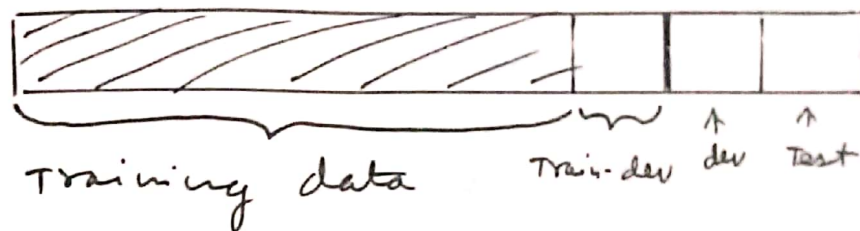This way we set the target to what our user will be uploading

How to evaluate if training data and dev/test data come from different distribution?

Assuming human error $\approx 0\%$ error, lets say we get

Training error ... 1% ⎱ 9%
Dev error ... 10% ⎰

This 9% error is a result of variance and data mismatch (training data may have been more easy to recognise than dev data - like website cats are easy training data but app cats are hard dev data)

So to find out if its variance or data mismatch we use a training-dev set - same distribution as training set, but not used for training.

Example:



Training data      Train-dev   dev   Test

Examples,

1. Variance Problem

Training error   1%.
train-dev error   9%. ⎱ Variance
Dev error   10%.

2. Data Mismatch Problem

Train error   1%.
Train-dev error   15% ⎱ Data Mismatch
Dev error   10%.

3. Bias Problem

Human error   0%. ⎱ Avoidable bias
Train error   10%.
train-dev error   11%.
Dev error   12%.

4. Bias & Data Mismatch error

Human error   0%. ⎱ Avoidable bias
Train error   10%.
Train-dev error   11%. ⎱ Data Mismatch
Dev error   20%.

Bias/variance on mismatched training & dev/test sets,

| | | |
|---|---|---|
| Human level | 4%. ↑ Avoidable bias | 4%. |
| Training set error | 7%. ↑ Variance | 7%. |
| Train-dev error | 10%. ↑ Data Mismatch | 10%. |
| Dev error | 12%. ↑ Degree of overfitting to dev set | 6%. |
| Test error | 12%. | 6%. |

It has decreased here since dev/test set is easier than train set

It may not always increase like this

More general formulation,

| | General speech recognition | New project ↓ Rearview mirror speech data |
|---|---|---|
| Data gathered from older projects → | | |
| Human level | "Human level" 4% | 6% |
| | ↑↓ Avoidable bias | |
| Error on examples trained on | "Training error" 7% | 6% |
| | ↑↓ Variance | |
| Error on examples **not** trained on | "Train-dev" error 10% ←→ "dev/test" error 6% Data Mismatch | |

Addressing data mismatch,

→ Carry out manual error analysis to try to understand difference between training and dev/test sets

   (g. noisy - car noise

→ Make training data more similar; or collect more data similar to dev/test sets

   (g. Simulate noisy 'in car' data

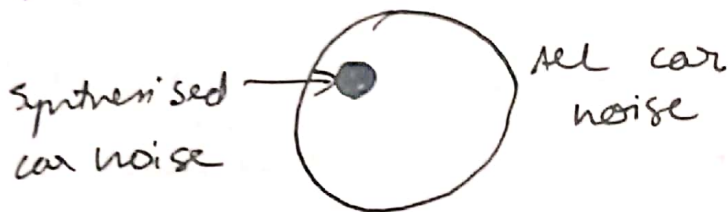One way to do this is through artificial data synthesis.

Clear Voice + Car Noise = Synthesised in-car audio

10,000 hours                    1 hour

The problem with this is, we may end up overfitting to the 1 hour of car noise. Although car noise seems the same to uh, 1 hour of the car noise we have may only be a small subset of all car noise



Synthesised car noise → ● ← all car noise

This is similar when people use computer generated images or video games to get more data of cars. (The video game may only have 20 models of cars
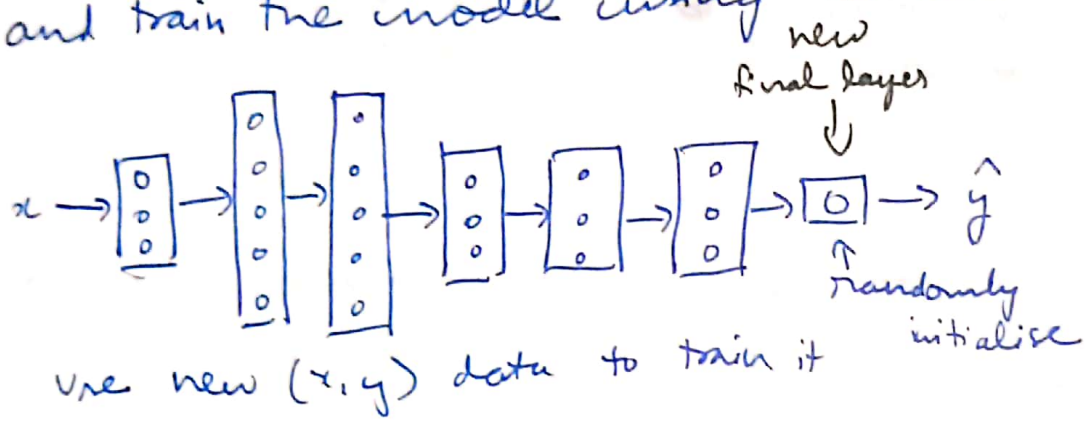
So the model will be better if we train on synthesised data made using 10,000 hours of car noise

# Learning from multiple tasks

## Transfer learning:

When we use a pre-trained model to train on new data (like using a cat classifier model to train a x-ray classifier).

We do this by using the old NN but removing the final layer. Instead we use a new randomised final layer and train the model using the new data

$$x \rightarrow \boxed{\vdots} \rightarrow \boxed{\vdots} \rightarrow \boxed{\vdots} \rightarrow \boxed{\vdots} \rightarrow \boxed{\vdots} \rightarrow \boxed{\vdots} \rightarrow \boxed{0} \rightarrow \hat{y}$$

new final layer ↓

randomly initialise ↑

Use new $(x, y)$ data to train it

When to use transfer learning,

→ Task A and B have same input x (like cat & x-ray are both images)

→ You have a lot more data for Task A (cat) than Task B (x-ray)

→ Low level features from A (like edges and curves) could be helpful for learning B.

# Multi-task Learning:

Where we design a neural network that performs multiple tasks.

Example, an autonomous driving car:

$$y^{(i)} \qquad - (4,1)$$

|  |  |
|---|---|
| Pedestrian | 0 |
| cars | 1 |
| Stop sign | 1 |
| Traffic lights | 0 |
| $\vdots$ | |

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \cdots & y^{(m)} \end{bmatrix} \qquad - (4,m)$$

This is different from softmax regression (the multiclass regression where an image has multiple labels) — in that it is either a pedestrian or a car or a stop sign.... but here a pedestrian and a car can be in an image (we are multitasking by recognising a car & pedestrian).

$$\text{Loss:} \quad \hat{y}^{(i)} \qquad \rightarrow \quad \frac{1}{m} \sum_{i=1}^{m} \sum_{j=1}^{4} \mathcal{L}\left(\hat{y}_j^{(i)}, y_j^{(i)}\right)$$
$$(4,1) \qquad\qquad \uparrow$$

we sum over only values of $j$ with 0/1 label

$$Y = \begin{bmatrix} 1 & 1 & & 0 & ? \\ 0 & 0 & & 1 & 0 \\ ? & ? & \cdots & 1 & ? \\ ? & ? & & 0 & \vdots \\ \vdots & \vdots & & \vdots & \end{bmatrix} \quad \left.\begin{array}{l} \text{it can have incomplete} \\ \text{data and it won't matter} \\ \text{since the Loss will only} \\ \text{take into account the} \\ \text{0/1 labels and not ?} \end{array}\right)$$

multitask learning labels

when multi-task learning makes sense:

→ Training on a set of tasks that could benefit from having shared lower-level features.

→ Usually: Amount of data you have for each task is quite similar

→ Can train a big enough neural network to do well on all the tasks

## End-to-End Deep Learning

You can sometimes replace a multiple set of tasks into one neural network.

Example:

(1) audio ⟶ feature → Phonemes → Words → Transcript

(2) audio ————————— END TO END DL —————————→ Transcript

But for this you need a lot of data.
However in some cases its better to have multiple tasks.

Exa: In face recognition, you can have one NN to identify where the face is in a picture, crop and centre it, and then have another NN to identify the person's face.

Pros and Cons of end-to-end DL,

pros:
→ Let the data speak
→ Less hand-designing of components needed

cons:
→ May need large amounts of data
→ Excludes potentially useful hand-designed components

The key question is,
Do you have sufficient data to learn a function of the complexity needed to map x to?
Example, in face recognition its hard to get a lot of data of every possible image about of a person's face position, so it makes more sense to first make it learn to detect where the face is and then detect the face's identity.