

Course 3: Structuring Machine Learning Projects

WEEK 1 Machine Learning Strategy (1)

why ML strategy?

Motivating example e.g. Cat classifier



90% accuracy \hookrightarrow not good enough
 \hookrightarrow want to improve!

Ideas:

- Collect more data
- Collect more diverse training set
- Train algorithm longer with gradient descent
- Try Adam instead of gradient descent
- Try bigger network
- Try smaller network
- Try dropout
- Add L_2 regularization
- Network architecture
 - Activation functions
 - # hidden units
 - ...

Andrew Ng

which one to pursue? \rightarrow ML Strategy

Orthogonalisation

clear about what to tune, in order to try to achieve one effect

\hookrightarrow orthogonalisation

By having orthogonal controls (not coupled with other goal) that are ideally aligned with the things we want to control \rightarrow easier to tune.

chain of assumption in ML

For a supervised learning system to do well, we usually need to tune the knobs of our system to make sure 4 things hold true

① Fit training set well on cost function

need to at least fit training set well \rightarrow pass some acceptability assessment
 \approx human level performance

\hookrightarrow 1 knob to adjust so algorithm

can fit well on training set

\hookrightarrow separate set of knobs to tune
model

\hookrightarrow regularization

\hookrightarrow bigger data set

② Fit dev set well on cost function

③ Fit test set well on cost function

④ Performs well in real world

bigger dev set

\hookrightarrow change dev set or cost function

\neg if model performs well on test set but not real world / dev vary with dev set distribution / cost func not doing the right thing

when training NN, maybe not using Early Stopping

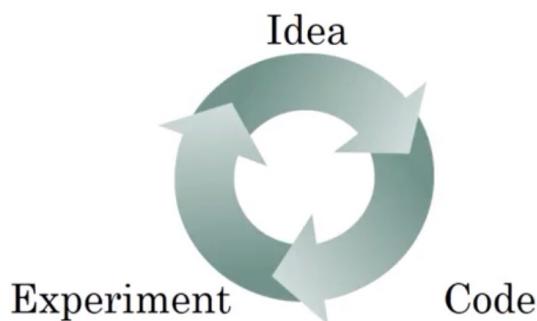
\hookrightarrow it affects fitting training set (fit less well)

and fitting dev set (fit better)

\hookrightarrow less orthogonalise

This course is about what to tune, to solve problems that limit the performance of ML in various places above.

Single number evaluation metric



95% chance it really is a car
from all cars pic in the set, classifier
finds and 90%
of them

Classifier	Precision	Recall	
A	95%	90%	72.4%
B	98%	85%	91.0%

F1 Score

Usually a trade-off

loosely:

Precision: of the examples that our classifier recognized as cars, what percentage are actually cars?

$$\text{Precision} = \frac{\text{True Positive}}{\text{Predicted Results}} \quad \text{or} \quad \frac{\text{True Positive}}{\text{True Positive} + \text{False Positive}}$$

Recall: of all the images that are really cars, what percentage are correctly recognized by our classifier

$$\text{Recall} = \frac{\text{True Positive}}{\text{Actual Results}} \quad \text{or} \quad \frac{\text{True Positive}}{\text{True Positive} + \text{False Negative}}$$

With 2 evaluation metrics, it may be hard to pick best model

↳ a metric that combines recall & precision → F₁ score

$$F_1 = \frac{2}{\frac{1}{P} + \frac{1}{R}} \quad \text{harmonic mean}$$

⇒ dev set + single real number evaluation metric

↳ speed up iterating

Percentage %.

Algorithm	US	China	India	Other	Average
A	3%	7%	5%	9%	6%
B	5%	6%	5%	10%	6.5%
C	2%	3%	4%	5%	3.5%
D	5%	8%	7%	2%	5.25%
E	4%	5%	2%	4%	3.75%
F	7%	11%	8%	12%	9.5%

hard to track 4 metrics
we average instead

Satisficing and Optimising metrics

It's not always easy to combine all the things we care into a single real number evaluation metric

↳ set up satisficing and optimising metrics

Classifier	Accuracy	Running time
A	90%	80ms
B	92%	95ms
C	95%	1,500ms

say we care about accuracy (can be F1 score) and running time
e.g. combine the 2 metrics

$$\text{cost} = \text{accuracy} - 0.5 \times \text{run_time}$$

maximises accuracy, but subject to running time ≤ 100 ms

↑
optimising
metric

T
satisficing metric
just have to be good enough
beyond which we
don't care

Generally, if we have N metrics that we care, can pick 1 of them to optimize, then $(N-1)$ to be satisfying.

as long as it reach some threshold

↳ we don't really care

Another example: Detection of hate words / trigger words

Care about accuracy

false positive (falsely triggered)

↳ maximize accuracy, subject to ≤ 1 false positive / day

Train/dev/test distribution

Cat classification dev/test sets

↳ development set, hold out cross validation set

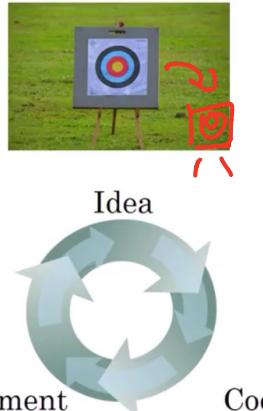
Regions:

- US
- UK
- Other Europe
- South America
- India
- China
- Other Asia
- Australia

{ Dev

{ Test
bad idea

↳ Dev and Test set from different distribution



dev set + Metric

Afford different distribution test set

Dev set ↳ hold out cross validation set
also called

↳ Randomly shuffle into dev/test

↳ after trained multiple model, use this dev set to pick one

↳ then improve on dev set performance

↳ use test set to do final evaluation

Guideline:

choose a dev set and test set to reflect data we expect to get in the future and consider important to do well on.

Dev set and test set should come from same distribution.

Setting Dev set + metric: Define what target we aim for

test set (same distribution): help to aim final target

training set: how well we hit the target

size of the Dev and Test Set

old way of splitting Data

70%	30%	
train		test
60%	20%	20%
train	dev	test

Reasonable when data are smaller

But in modern era, if we have 1 mil example (10,000 examples)

98%	1%	1%
train	dev	test

size of test set:

Set your test set to be big enough to give high confidence in the overall performance of our system

When to change dev/test sets and metrics

(c). Cat dataset examples

Metric: classification error

Algorithm A: 3% error \rightarrow but also with porn pic \rightarrow unacceptable

Algorithm B: 5% error \rightarrow no porn image

\hookrightarrow metric + Dev: Prefer A

\hookrightarrow User: Prefer B

$$\text{Error: } \frac{\sum_{i=1}^m w^{(i)} \left[l^{(i)} y_{\text{pred}}^{(i)} + y^{(i)} \right]}{\sum_{i=1}^m w^{(i)}} \quad \begin{array}{l} \text{just count misclassified example} \\ \text{predicted value (1/0)} \end{array}$$

\hookrightarrow problem: treat porn/non-porn images equally

\hookrightarrow can add weight term: $w^{(i)} = \begin{cases} 1 & \text{if } x^{(i)} \text{ is non-porn} \\ 10 & \text{if } x^{(i)} \text{ is porn} \end{cases}$

If we find the evaluation metric is not giving the correct rank order preference for what is actually better algorithm, then it's time to think of a new metric

* Goal of evaluation metric: accurately tells us, which model is better.

Orthogonalization for cat pictures: anti-porn

\rightarrow 1. So far we've only discussed how to define a metric to evaluate classifiers. \hookleftarrow place the target

Step 1

\rightarrow 2. Worry separately about how to do well on this metric.

'Criminally'

Step 2



Another example

Algorithm A: 3% error

✓ Algorithm B: 5% error ↵

→ Dev/test



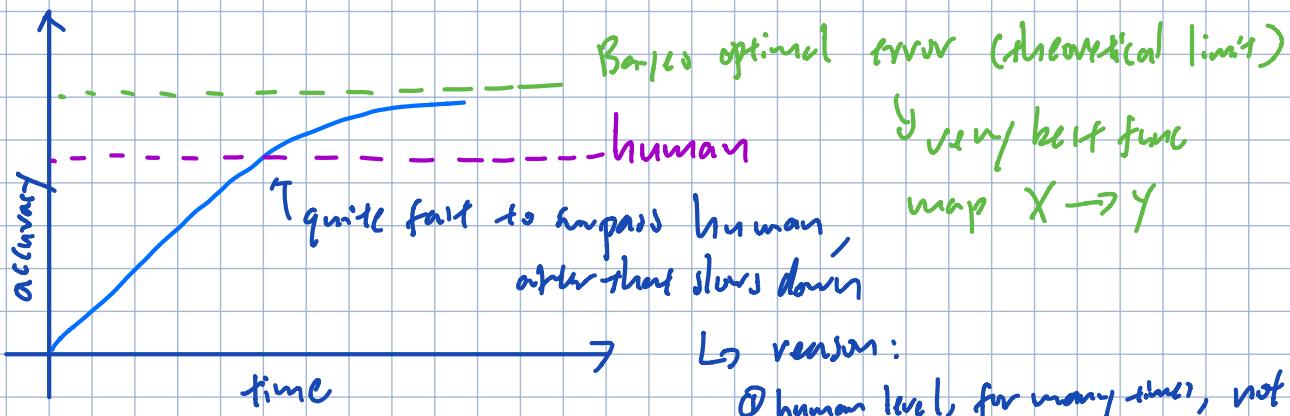
→ User images



If doing well on your metric + dev/test set does not correspond to doing well on your application, change your metric and/or dev/test set.

why human-level performance?

- ① Advancement in Deep Learning, ML algorithms are suddenly much better → more feasible in a lot of applications even to be competitive with human level performance
- ② workflow of design and building ML systems is much more efficient.



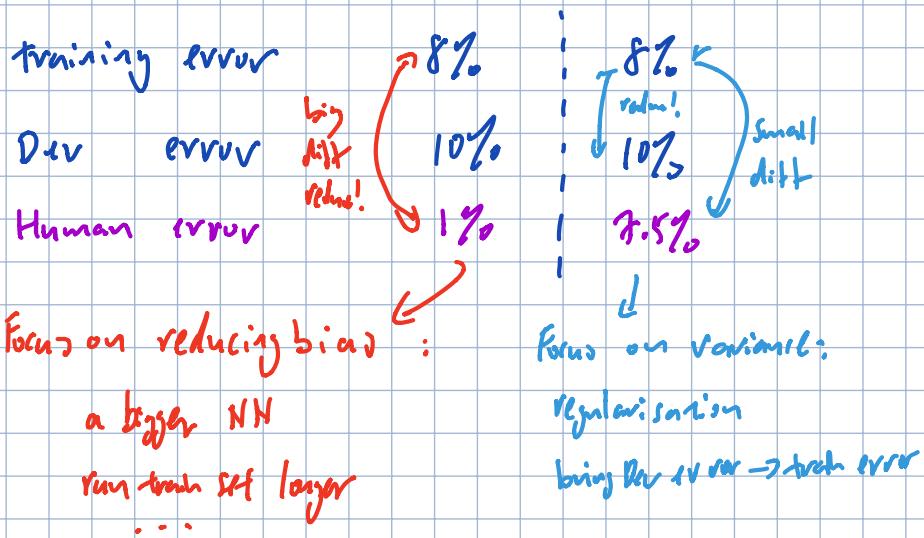
As long as ML is worse than humans, we can:

- Get labeled data from human
- Gain insight from manual error analysis:
Why did a person get this right?
- Better analysis of bias/variance

② so long your performance is lower than human level, there are two we can use to improve. However it surpasses human level

Avoidable bias

Cart classification example



Can think human level error as an estimate for Bayes error

Avoidable Bias: Difference between Bayes error and the training error

Variance : Difference between training error and dev error

understanding human level performance

Human-level error as a proxy for Bayes error

Medical image classification example:



Suppose:

- (a) Typical human 3 % error
- (b) Typical doctor 1 % error
- (c) Experienced doctor 0.7 % error
- (d) Team of experienced doctors .. 0.5 % error

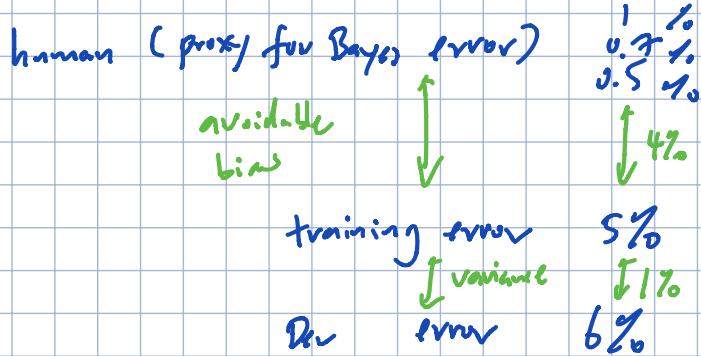
What is "human-level" error?

\Downarrow Bayes error $\leq 0.5\%$

$\Downarrow 0.5\% \because \sim$ Bayes error

need to ask what is the purpose of defining human-level error

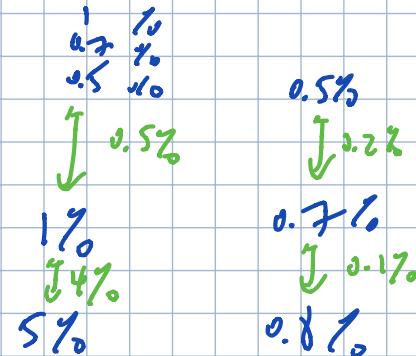
e.g. medical imaging diagnosis



focus on bias reduction:
e.g. bigger NH

→ human level we use

↓ doesn't matter which human level we use



does not matter

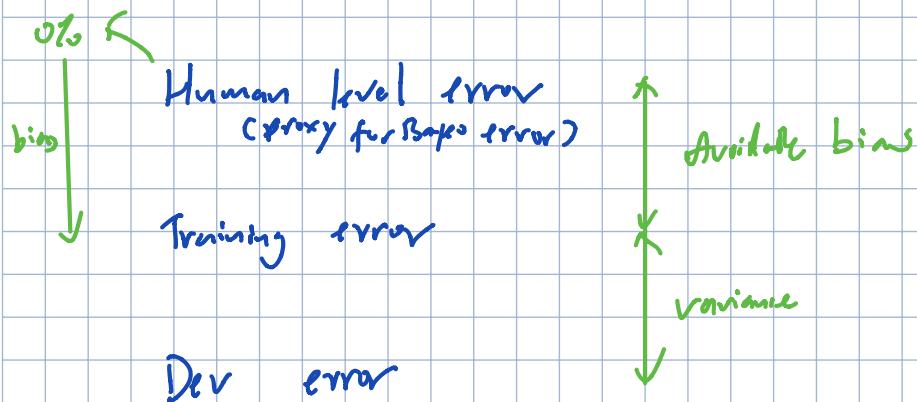
also

E focus on variance reduction

e.g. regularization

larger kernel size

Summary



Surpassing human-level performance

Team of humans

0.5%

One human

1%

0.4%

Training error

0.8%

Dev error

0.5% ↕

1%

0.3%

0.4%

not much into it

focus on bias or variance reduction

Towerfit! or Bayes error ~ 0.1-0.2%

what is available bias?

Problems where ML significantly surpasses human-level performance

- - Online advertising
- - Product recommendations
- - Logistics (predicting transit time)
- - Loan approvals

Structural data → need from DB
Not natural perception → not CV, NLP etc
Lots of data

- Speech recognition
- Some image recognition
- Medical
 - ECG, Skin cancer, ...

Andrew Ng

Improving your model performance

The two fundamental assumptions of supervised learning

1. You can fit the training set pretty well.

~ Avoidable bias (low)

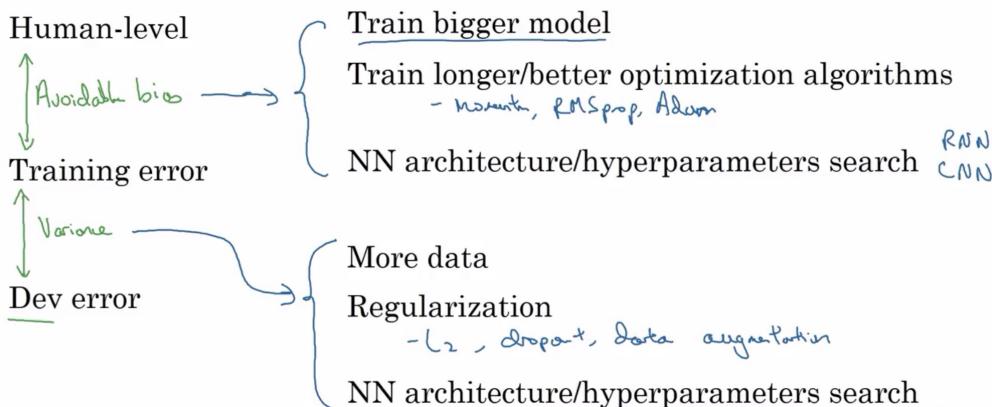
bigger NN
train layer

2. The training set performance generalizes pretty well to the dev/test set.

~ Variance (low)

orthogonal features
regularization, more training data

Reducing (avoidable) bias and variance



Week 2: Machine Learning Strategy (2)

Carrying out error analysis

If our learning algorithm is not yet at the performance of a human, we can manually examine mistakes that our algorithm is making
↳ error analysis

Look at dev examples to evaluate ideas



90% accuracy
10% error

← misclassify dogs as cats?

cat classifier

Should you try to make your cat classifier do better on dogs? ←

Error analysis:

- Get ~100 mislabeled dev set examples. ← train manually
- Count up how many are dogs.

↳ say 5% are dogs → 10% error → 9.5% error

if we solve the dog problem

↳ 5% relative decrease in error
from 10% → 9.5%

↳ "ceiling"



Say 50% are dogs → more confident to start on a dog problem

if we solve this. error 10% → 5%

Evaluate multiple ideas in parallel

Ideas for cat detection:

- Fix pictures of dogs being recognized as cats ↪
- Fix great cats (lions, panthers, etc..) being misrecognized
- Improve performance on blurry images ↪

Image	Dog	Great Cats	Blurry	Instagram	Comments
1	✓			✓	Pitbull
2			✓	✓	
3		✓	✓		Rainy day at 200
:	⋮	⋮	⋮	⋮	
% of total	8%	43%	61%	12%	

↑
↓

↳ ceiling ↗ for improvement

dataset
misclassified
examples

Cleaning up incorrectly labeled data

By right, labeled data map $X \rightarrow y$, what if y is mislabeled?

Incorrectly labeled examples

x						 1	 1
y	1	0	1	1	0		

wrong



incorrectly
labeled

if we found incorrectly labeled data, what we should do?

we first consider training set

DL algorithms are quite robust to random errors in the training set as long as the errors are not too far from random → probably ok to just leave the error. → actual % of error is not too high

DL algorithms are less robust to systematic errors.

e.g. consistently labels white dogs as cats \rightarrow classifier will learn to classify all white colored dogs as cats.

How about dev / test set?

during error analysis add in

Error analysis

Image	Dog	Great Cat	Blurry	Incorrectly labeled	Comments
...					
98				✓	Labeler missed cat in background
99		✓			
100				✓	Drawing of a cat; Not a real cat.
% of total	8%	43%	61%	6%	

Overall dev set error

10%

Errors due incorrect labels

0.6%

Errors due to other causes

9.4%

6% of 10%

relatively small compared to

maybe not worth the time

other example:

2%

0.6% \curvearrowleft know it is a larger fraction $0.6/2 \approx 20\%$

1.4%

\curvearrowleft it is much more worth while to fix label

Final goal of dev set: \rightarrow to help us select between (e.g.) 2 classifier A or B.

Correcting incorrect dev/test set examples (guideline)

- 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.

this has to be the same!

\rightarrow actually ok, Learning only with one is quite robust to this

Building your first system quickly, then iterate

Speech recognition example

Various directions

- Noisy background
 - Café noise
 - Car noise
- Accented speech
- Far from microphone
- Young children's speech
- Stuttering *uh, ah, um, ...*
- ...



- Set up dev/test set and metric
- Build initial system quickly
- Use Bias/Variance analysis & Error analysis to prioritize next steps.



Training and testing on different distributions

DL need a lot of data, sometimes not all training data come from the same distribution as dev/test set.

Cat app example

Data from webpages



$\approx 200,000$

I get from internet as more training data

Data from mobile app



$\approx 10,000$

we care about this mostly
↳ what user uploaded pic will be predicted in the app

Want to use both, how?

option 1: put both data sets together - so 210,000 pics, then shuffle

train	dev	test
205,000	2,500	2,500



advantage: all train/dev/test come from same distributions

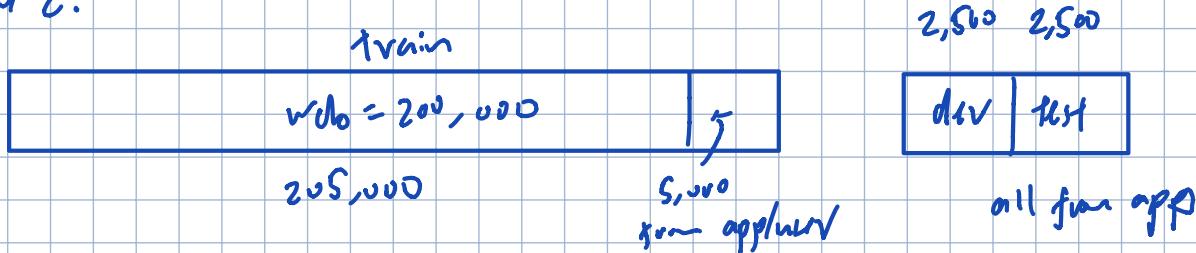
disadvantage (huge!): dev set (2,500) mostly came from internet, rather than app user (you actually care about).

Out of 2,500 - 2,381 from internet

↳ spent most time on wrong targets! (diff between dev and actual app's data!)

* This option is bad!

Option 2:



advantage: dev/test set align with final product

disadvantage: train has different distribution with dev/test set.

Option 2 will have better performance in the long term.

Speech recognition example

Speech activated rearview mirror



Training

audio | transcript



- Purchased data
- Smart speaker control
- Voice keyboard
- ...

500,000 utterances

Dev/test

Speech activated
rearview mirror

20,000

train

5,000

dev

test

{ contains
more street noises

5K each

virtual speech
from rearview mirror

Bias and Variance with mismatched data distribution

The way we analyse bias and variance changes when our training set comes from a different distribution than our dev/test set.

Cart classifier example

Assume human get $\approx 0\%$ error

To carry out error analysis, look at

training error	1%
Dev error	10%

If they come from same distribution \rightarrow large variance problem (overfit)

If they come from different distribution \rightarrow can't draw the same conclusion

Lets say training set is easy/dev set is hard

2 things change from train set \rightarrow dev set

- ① algorithm saw data in train set but not in dev set } hard to say where
- ② distribution in dev set is different } 9% error from

In order to tease out these 2 effects, define a new piece of data:

Training-dev set: Same distribution as training set, but not used for training

i.e.

some distribution $\xrightarrow{\text{shuffle again}}$ train-dev Dev test



from NN on training set proper Some distribution

e.g. training error	1%	variance problem (overfitting)
train-dev error	9%	
Dev error	10%	

e.g. training error	1%	low variance (not overfit, can generalize)
train-dev error	1.5%	
Dev error	10%	

e.g. training error	10%	\rightarrow available bias problem (model too simple)
train-dev error	11%	
Dev error	12%	

(recall human error $\sim 0\%$)

e.g. training error	10%	\rightarrow available bias problem
train-dev error	11%	
Dev error	20%	

} due to data mismatch

Bias/variance on mismatched training and dev/test sets

Key quantity:

Human-level error	4%	2 issues
training set error	7%	
train-dev error	10%	
dev set error	12%	
test set error	12%	

} available bias

} variance (generalization)

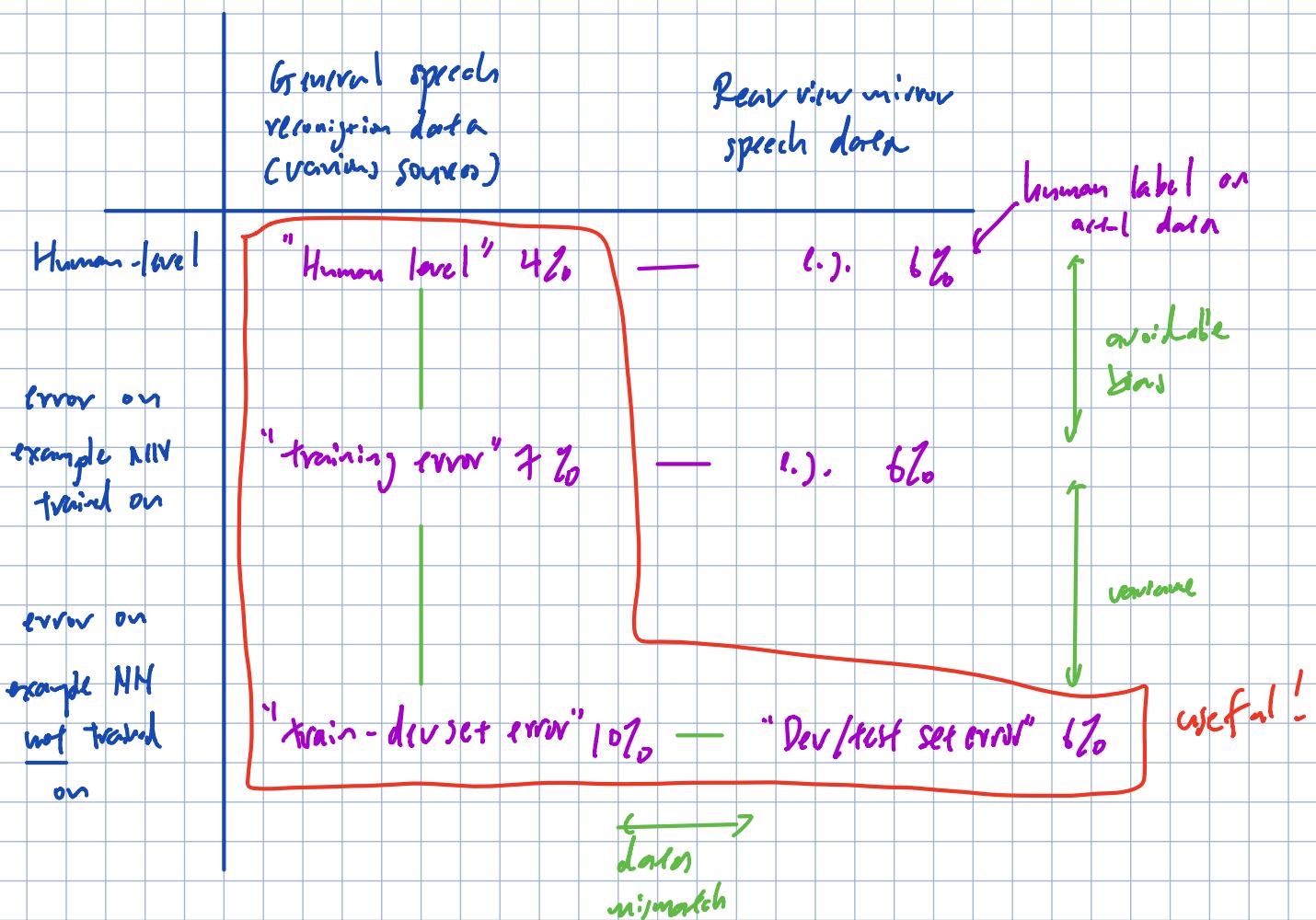
} data mismatch

} degree of overfitting to the dev set

4%	4%
7%	7% } train data hard
10%	10% } dev/test data
12%	12% } easy
12%	12% } easy

More general formulation

e.g. Rearview mirror speech recognition



Addressing Data Mismatch

Don't have systematical solutions, can try a few things

Addressing data mismatch

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

E.g. noisy - car noise street numbers

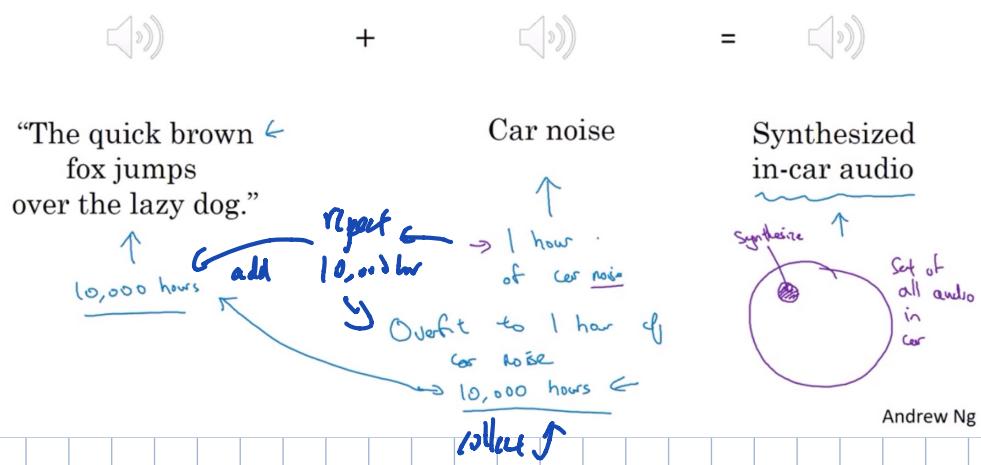
→ to dev/test set

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

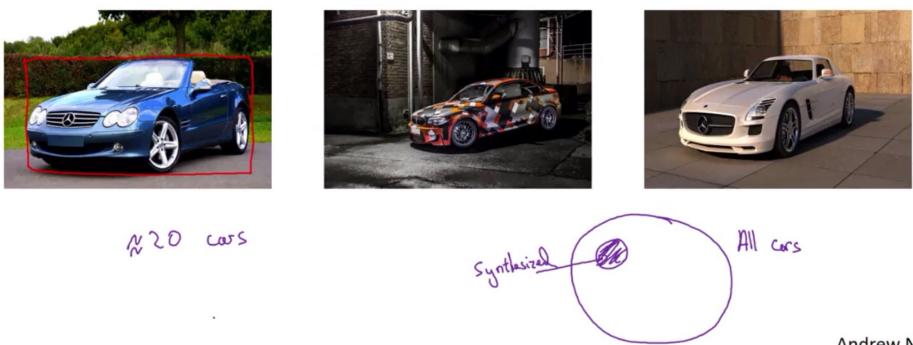
E.g. Simulate noisy in-car data

Can try artificial data synthesis

Artificial data synthesis



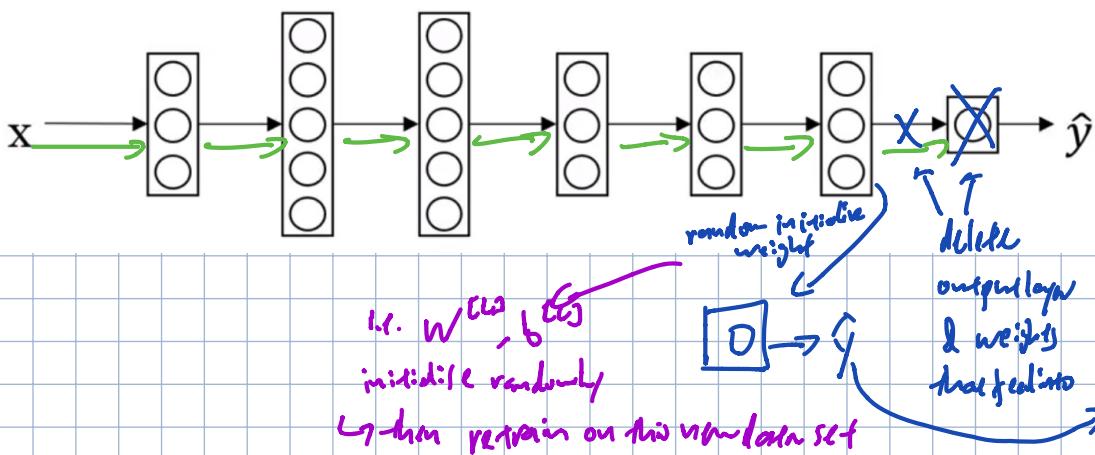
Car recognition:



Learning from multiple tasks: transfer learning

Sometimes we can take knowledge of a NN that learned from 1 task, apply that knowledge to a separate task.

radiology
new (X, y)
imaging
Diagnosis
 (X, y)
obj



e.g. image
recognition

↳ transfer to
radiology
diagnosis

if radiology data is small, might keep the previous layer weights and only train the last layer.

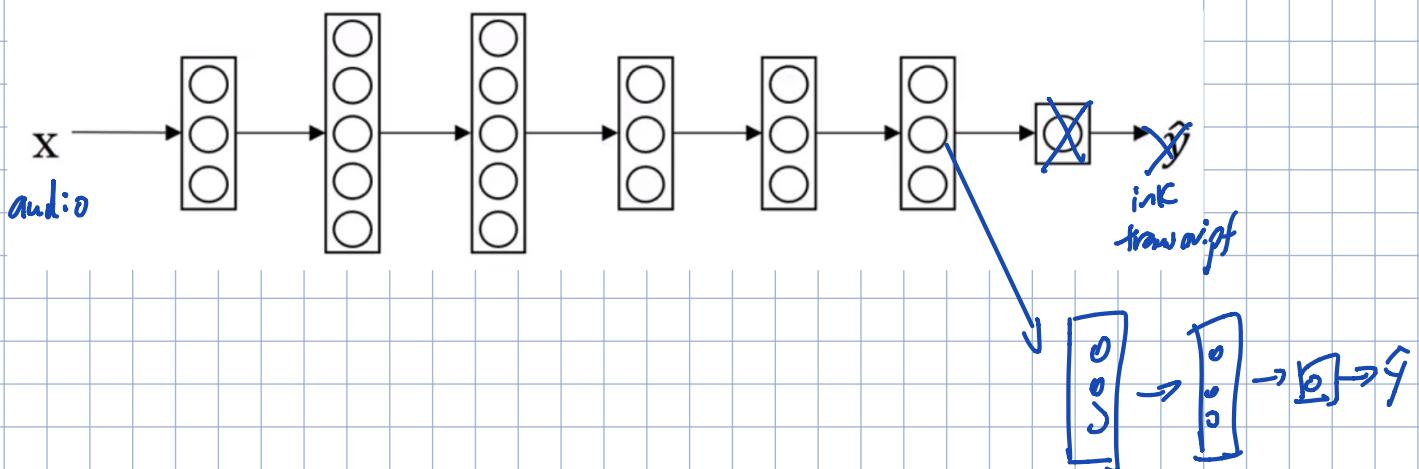
if we have enough data, can retrain all the layers.

if then the initial (X, y) is called pre-training

train on (X, y) is called fine-tuning

A lot of the low level features, such as detecting edges, detecting curves detection positive objects, learned from large set of image recognition, can help our learning algorithm do better in radiology diagnosis.

Another example



Transfer learning make sense when we have a lot of data for the problem you're transferring from and usually less data for the problem you're transferring to.

When transfer learning makes sense

Transfer from A \rightarrow B

- Task A and B have the same input x . *images / audio*

- You have a lot more data for Task A than Task B.

- Low level features from A could be helpful for learning B.

Multi-task learning

transfer learning: sequential process

multi-task learning: start off simultaneously

Simplified autonomous driving example



$x^{(i)}$

Pedestrians
Cars
Stop signs
Traffic lights
⋮

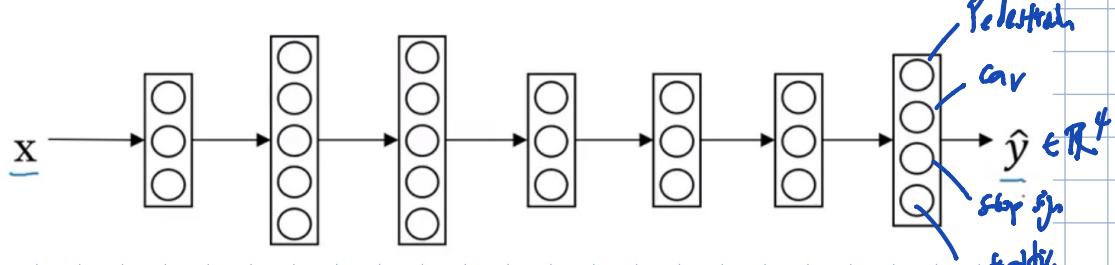
$y^{(i)}$
0
1
1
0
⋮

(4,1)

detect 4 things
at the same
time.

$$Y = \begin{bmatrix} y_1^{(1)} & y_2^{(1)} & y_3^{(1)} & \dots & y_n^{(1)} \end{bmatrix}$$

Neural network architecture



Loss: $\hat{y}^{(i)}$
(4,1)

$$\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^4 2(\hat{y}_j^{(i)}, y_j^{(i)})$$

sum only over values of j with 0/1 label

diff from softmax classification

$$-\hat{y}_j^{(i)} \log \hat{y}_j^{(i)} - (1 - \hat{y}_j^{(i)}) \log (1 - \hat{y}_j^{(i)})$$

Unlike softmax regression: (1 image | class/label)

this one image can have multiple class/labels.

↳ this multi-task learning

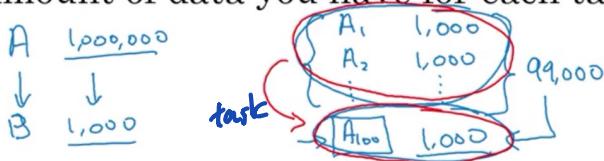
↳ 1 NN solving 4 problems

also works with unlabeled data

in Σ for y_j ^{unit}

$$T = \begin{bmatrix} 1 \\ 0 \\ \vdots \\ ? \\ \dots \end{bmatrix}$$

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.

 - A $\frac{1,000,000}{100}$
 - \downarrow \downarrow
 - B $1,000$
 - task
 - $A_1, 1,000$
 - $A_2, 1,000$
 - \vdots
 - $A_{100}, 1,000$
 - $99,000$
- Can train a big enough neural network to do well on all the tasks.

End-to-End deep learning

speech recognition example



End-to-End



Needed a lot of data to work well

3,000 h of data works well on traditional pipeline,

but need 10,000 \rightarrow 100,000 h of data for end-to-end pipeline

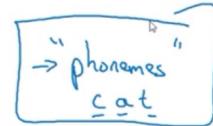
Also, sometimes we have more data for individual sub-task
but not many data that satisfy all steps together.

whether to use end-to-end learning

Pros and cons of end-to-end deep learning

Pros:

- Let the data speak $x \rightarrow y$
- Less hand-designing of components needed



Cons:

- May need large amount of data
- Excludes potentially useful hand-designed components

Data ↘ Hand-design.

2 source of knowledge of ML

Andrew Ng

Applying end-to-end deep learning

Key question: Do you have sufficient data to learn a function of the complexity needed to map x to y ?

