

# 3 ML strategies

Andrew Ng

Note by:  
Hossein Mamaghanian  
amirhm@gmail.com



## "ML Strategies"

better to define an evaluation metrics

Precision: recognized as cat what% are really cat

Recall: what % of all cats are correctly recognized.

→ F1 score: (standard way of combining Precision and Recall)

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

having a single evaluation metrics speed up iterating

or in other situation may be an average could be solution for a simple metrics.

---

example :

	Accuracy	Running time
A	92	80ms
B	93	95ms
C	95	150ms

cost = Acc - 0.5 × running Time

Maximize Accuracy subject to running time < 100ms

---

optimizing metrics.

---

satisfying metrics

if we have  $N$  metrics  
better to define  $\rightarrow$

: 1 optimizing	: $(N-1)$ satisfying
----------------	----------------------

→ example of google assistant or Amazon Alexa:

→ Wake words / Trigger words.

\* Alexa, ok google, hey siri

(optim) accuracy.

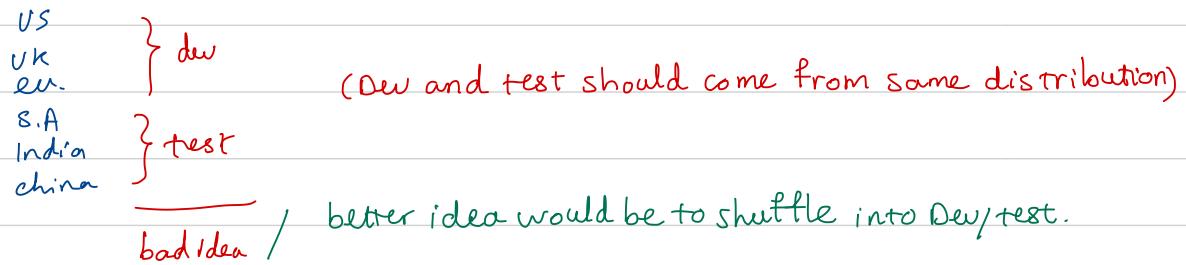
(satisfy) \* false positive:

Maximize accuracy s.t.  $\leq 1$  false positive in 24 h

## Train / Dev / test distributions

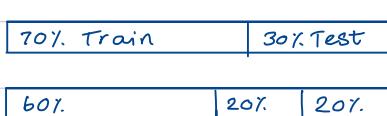
↳ development set / hold out cross validation set.

for example if we have data set from different geographic regions:



Dev/Test should reflect data you expect to get in the future.  
(should be from same distribution)

## old way of splitting data



} if we have less data (old way)  
100 ~ 10K



In modern era of ML, where we need  
huge data for training

size of test set.

big enough to give high confidence in the overall performance of system.

→ it is ok if only Train + Dev is used. (No test).

(Not usually recommended, better to keep test as an separate set).

→ when to change dev/test sets and metrics.

sometimes, since the metric does not reflect the what is preferred to rank different algorithms.

→ Alg A : 0.3% error(cat) → but higher F1P (Pornographic images).

Alg B : 0.5% error(cat)

→ then: metric + Dev → Alg A

You / user → prefer B

→ if metrics defined as:

$$\sum_{\text{index}}^{\text{num}} L(Y_{\text{pred}}^{(i)}, Y^{(i)})$$

Take away: if metrics you are using actually is not ranking on the preference order, define new metrics

→ adding a weight to penalize the F1P for Pornographic

$$\text{error} = \sum_{\text{index}}^{\text{num}} w_i L(Y_{\text{pred}}^{(i)}, Y^{(i)})$$

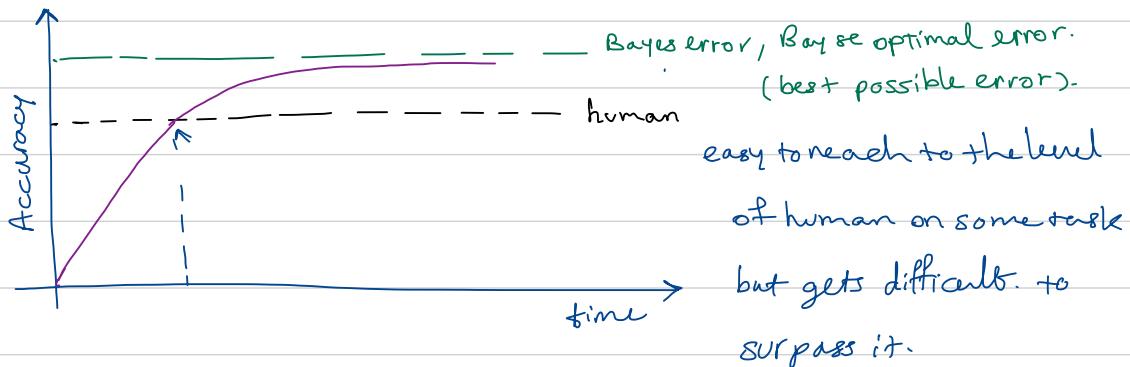
"Orthogonalization"

Wise  $\begin{cases} 1 \rightarrow X^{(i)} \text{ non porn} \\ 10 \rightarrow X^{(i)} \text{ porn} \end{cases}$

① Define a better metrics.

② Worry separately, how to do well on this metric.

## Human level / Bayesian error.

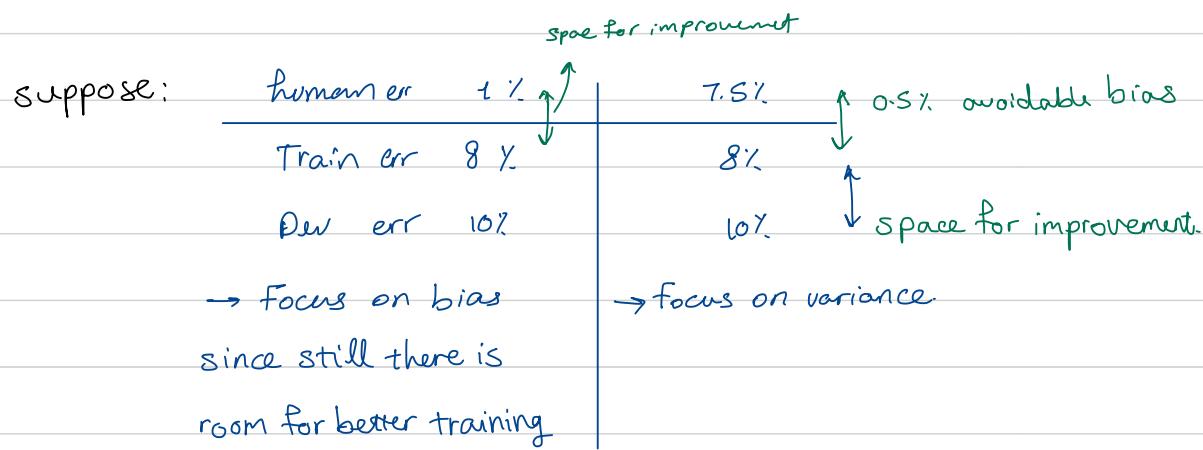


if performance is below the human level :

- ①. get more labeled data from humans.
- ②. get insight from manual error Analysis (why human get this right?)
- ③. Bias/ Variance Analysis (considering human level).

---

Bias/ Variance considering human level



→ human level error as a proxy for Bayes error.

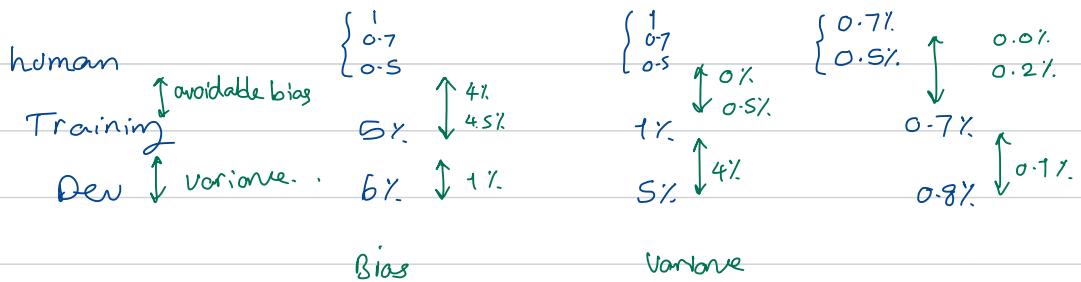
Imagine a radiology task:

Typical human	3%
Doctor	1%
Experienced doctor	0.7%
team of ..	0.5%

what is the human-level?

→ Bayes error  $\leq 0.5\%$ .

→ for the purpose of publication, maybe other options could be used.



Regardless of which human-level is used.

here, as it gets better it's relevant which one to choose.

Surpassing human-level.

→ team of human	0.5%	0.5%	when it's surpassed,
one human	1%	0.1%	it's difficult to rely
Train error	0.6	0.3	on human intuition.
Dev error	0.8	0.4	

Problems where ML significantly surpasses human-level

→ online ads

→ product recommendation

→ logistics.

→ loan approvals.

①. all learned from structured data

②. No natural perception (if there is, human does better)

③. lots of data are available.

---

Some other perception task

which ML is close or surpassing

human-level.

• speech recognition

• some image recognition

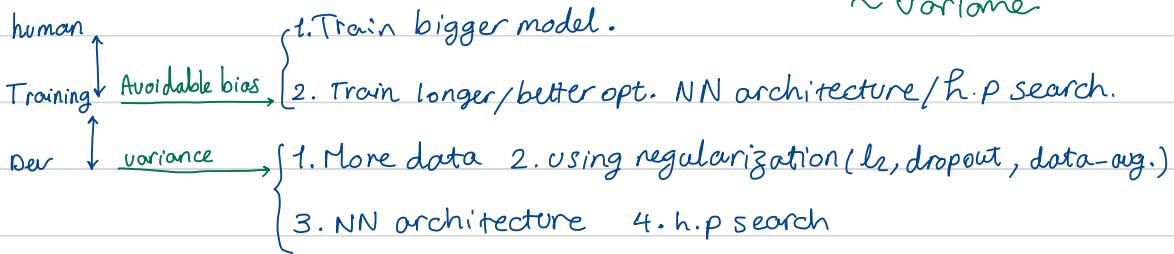
• medical (ecg, skin images,...)

---

Fundamental assumption of unsupervised learning

1. You can fit the training set, pretty well. ~ Avoidable bias

2. Training set performance generalize pretty well to the dev/test sets.



## "error Analysis"

For example using error Analysis to check whether an idea worth to work on

→ Imagine → 90% Accur. in cat/dog classifier  
10% error

→ ex. idea: You want to improve the cases that cats are wrongly classified  
→ Take 100 ex from Dev set, check how many of mislabels were dog  
3% → very little effect  
50% → worth to check.

→ evaluate several ideas in parallel.

①. Fix pic of dog recognized as cat

②. Fix for big cats.

③. Improvement of performance on blurry images.

Image	idea 1	idea 2	idea 3	Comments
1.				
2.	✓			Pitbulls, ...
3.		✓		
4.			✓	
:		✓		:
% of total	8%	43%	61%	

helps to see which idea worth to work and has the greatest impact.

cleaning up incorrectly labeled data.

→ if the error exist in labeled data (Train/ Dev/ Test) data sets.

Normally DL. algorithms are robust to random errors in the training set

(Most of the time it's ok, to leave it and not spend time to fix as long as it's random, small and size of dataset is big).

"But Not systematic errors"

↳ example if in cat/Dog problem all the white dogs are mislabeled as cat (This is called systematic error).

How about Dev/Test set

→ "Do error analysis"

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



does it make sense to go and Fix?

case1

case2

all errors

10%

- - 2%

{ other cause

9.4%

1.4%

mislabeled

0.6%

0.6%

Not a big

Big impact

impact

→ Make sense to go  
and Fix mislabel

If this was in dev set and goal was to  
select between two classifier A & B  
and range of errors between two were

in affected range then make sense

to fix mislabeled info.

Strategies for fixing incorrect Dev/test examples.

- Apply same process (To make sure they continue to come from same dist)
- consider examining both examples algo got right and as well wrong.
- Now after fixing Dev/Test, Training set is slightly from different distribution  
(This is less problematic, But O.K. Since Training size is huge and maybe not possible or less priority).

→ Build your system quickly, then iterate.

\* Set up Dev/Test set and metrics.

\* Build system quickly,

\* Use Bias/Variance analysis & error analysis

To prioritize next steps.

↳ For example in speech problem we may find that more error is due to "Far from Mic" then we can work on that.

→ Build Fast quick/ Dirty, Not over think just to do analysis and prioritize.

→ less strongly if there are significant prior experience / or literature!

example speech project

1. Noisy background  
→ cafe/car

2. Accented speech

3. Far from Mic

4. Young children's speech

5. Stuttering

Mis matched Training and dev/test data,

example : We want to design app for cat/dog problem

200K image from web (clearly framed)

10K from cellphone camera

→ There are two way to use this data.

Option 1: shuffle all 210K image and draw Train/dev/Test from that

Train 205K Pros: Train/Dev/Test are all from same distribution

Dev 2.5K Cons: only  $\frac{10}{210}$  % of Dev/Test set are from our target

Test 2.5K distribution (cellphone). → Aiming at wrong direction

Option 2:

Dev/Test all are from target application.

Train 205K (200K + 5K cellphone). Cons: Train still is from different distribution

Dev 2.5K (all cellphone) Pros: Dev/Test are from correct target distribution

Test ✓ ✓ ✓ Aiming at right direction.

→ Option 2 is better in longer term, since Training is done for correct target distribution.

Bias / Variance analysis with mismatched data distribution.

→ This is very useful analysis to prioritize next steps but how if datasets are not from same distributions?

→ Imagine Cat/Dog problem

Human err 0%. if dev is not from same distributions now is not  
Training err 1%.  
Dev err 10%. ↓ 9% } Possible to say that is it really variance problem  
or data mismatch (simply Dev set is much difficult).

→ Now, we could introduce new Train-dev set. (drawn from Train set)  
→ Same as Training distribution, but not trained on them, (not seen by algorithm)

Now human error 0%

Train error	1%	1%	0%
Train-dev n	9%	10%	10%
Dev v	10%	10%	11%
Main is variance Problem	Not data mismatch	main Problem is data mismatch,	available bias Problem.
			both bias and data mismatch

Human Level	4%	available bias
Train error	7%	Variance
Train Dev	10%	data mis-match
Dev_err	12%	degree of overfitting
Test_err	12%	+ to dev set

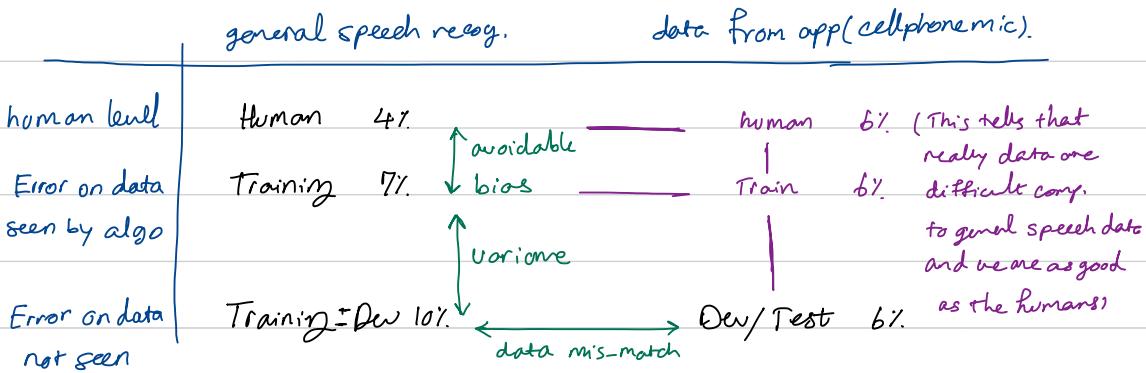
generally this trends of errors getting worse as we go down but in some cases maybe Dev/Test is much easier than Train set and error for them might be lower.

if this problem exist means that you should go back and get bigger dev set.

## More general Formulation

Imagine an application that we build a speech recognition system.

where Training is done on general data but target is from app.



Main Path for analysis is the green line but sometimes it is

worth to fill the violet parts as well to get some more insight.

## Addressing data mis-match.

→ carry out manual error analysis between Training, dev/test sets

For example car noise

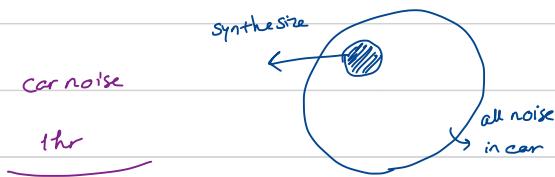
→ Make Training data more similar to dev/test sets

add artificial generated car noise ..

Artificial data synthesize work but there's a catch that maybe we overfit

to the artificial data → example

real speech + car noise  
100,000 hr      thr



chance that learning algo overfits to  
this subset

→ Most of the time for human ear this look fine but not in reality.

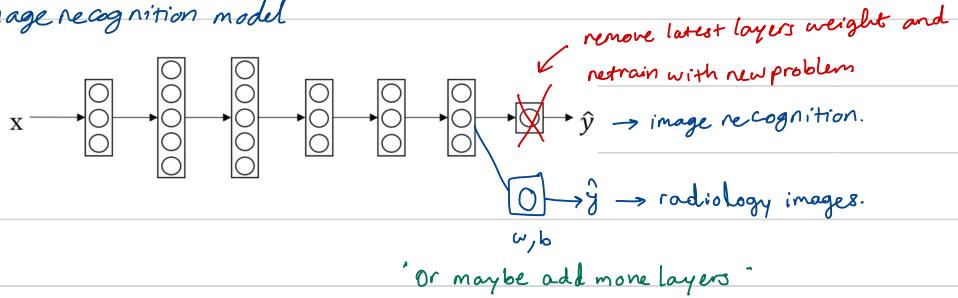
same for example of taking car images from games. (very small subset  
of cars exist in games, to human eyes is fine but not in reality).

## "Transfer learning"

Taking the knowledge neural network learned from one task and apply to separate task.

Imagine Problem of radiology images.

Take an image recognition model



→ depending on the amount of data we have for the problem at hand we may use the knowledge learned in image recognition task as initial value and Train full network or just the latest layers which are added

→ if full net is trained the training on image recognition is called "pre-training" and on radiology called "fine tuning"

→ The logic behind is that most of knowledge learned in early layers are general features in image (like edge / curvature,...) which is useful for radiology task

Transfer learning from Task A → B make sense if

1. both have same input (image in radiology example).
2. we have a lot more data for task A than B. (if reverse is true, it doesn't hurt but just wouldn't make sense)
3. low-level features from A could be helpful for learning B

## Learning from multiple task.

→ In multi-task learning idea is to start off simultaneously, trying to learn several task at the same time. each task helps other tasks.

example of self driving car

$$\text{goal: } \begin{cases} \text{Pedstrains.} \\ \text{cars} \\ \text{stop signs.} \\ \text{Traffic lights.} \end{cases} \quad \begin{matrix} y \\ \left[ \begin{matrix} x \\ x \\ x \\ x \end{matrix} \right] \end{matrix}$$

Output node is  $4 \times 1$  vector or in context of  
Training set  $Y \in \mathbb{R}^{4 \times M}$

The intuition is that all these tasks are related and knowledge learned could be useful for others as well.

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

↳ usual logistic loss

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

→ unlike softmax regression

one image could have multiple labels.

→ In multi-task learning, it is possible to train even if not all the labels are given for each image.

$$Y \in \begin{bmatrix} 1 & 1 & \dots \\ 1 & 0 & \dots \\ 0 & ? & \dots \\ ? & ? & \dots \end{bmatrix}$$

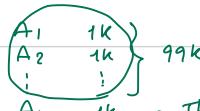
→ in such case in loss calculation sum over  $j$  is done just for available labels.

"When it make sense?"

①. a set of task that could benefit from having shared low-level features

②. Amount of data for each task is quite similar.

in T.L. A 1, M  
↓  
B 1k



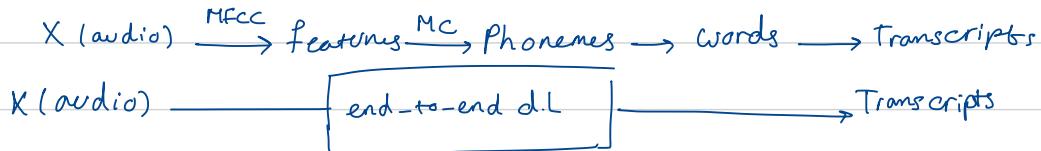
→ There are a lot more data to learn and augment on task A<sub>100</sub>

③. Can Train a big enough N.N. to do well on all.

→ Compared to alternative to multi-task learning which is independently Training different N.N for each task, — multi-task learning only hurts when the NN. is not big enough.

# end-to-end deep learning

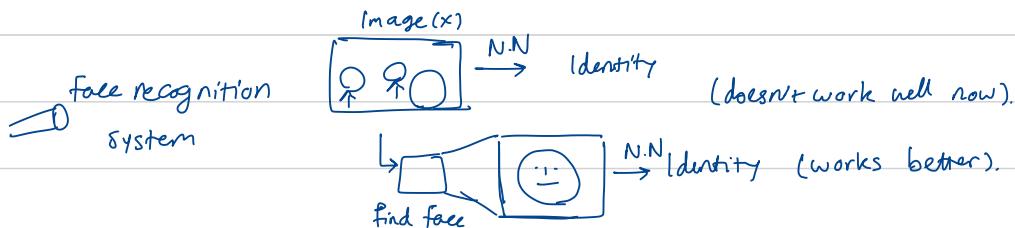
## Traditional speech recognition:



depending on the dataset size small : Traditional may work better

bigger : end-to-end system shines!

## example on intermediate solution



→ Face recognition N.N. takes two images and tells are the same person or not. (This Task needs 100,000s of images).

Machine Translation    english → text analysis --- → french

english    Trained on Pair of (F, E) sentence → french

## Estimating child age.



Image → extract bones → estimate age based on statistic  
(does not require a lot of data).

Image → age (does not work now since there are not enough data)

## Pros and Cons

Pros: ① let the data speak.

if big enough N.N is trained, more probable to learn all representation necessary without human pre-judgment.

② less hand-designing of components needed.

Cons:

① May need large amount of data.

② excludes potentially useful hand designed components.

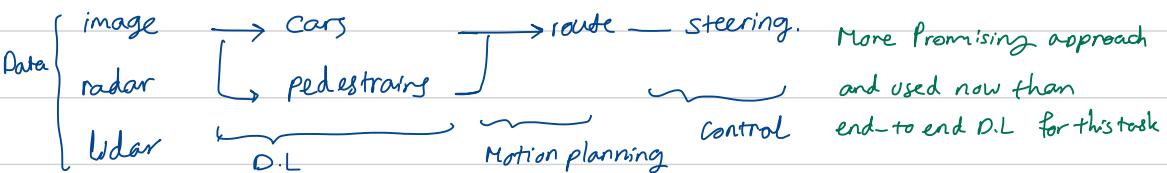
Data  
hand-designed features (Injecting human knowledge)

if less data exist, hand-designed features are useful

→ Applying end-to-end deep learning.

→ Do you have enough data to learn a function of the complexity needed to map X to Y?

Example of self-driving car:



→ Use D.L. to learn individual components.

→ Carefully choose  $X \rightarrow Y$  depending what task you can get data for.