



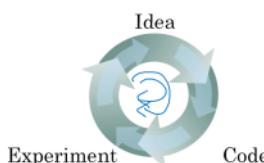
deeplearning.ai

## Setting up your goal

### Single number evaluation metric

set up a number that can evaluate performance

#### Using a single number evaluation metric



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

Dev set + Single number evaluation metric  
red speed up iterating

Andrew Ng

tell if Classifier A or B is better — speed up iterating.

ex: cat classifier  
→ of examples recognized as cat,  
what % actually are cats?  
→ what % of actual cats  
are correctly recognized

single # evaluation metric.

problem better precision, worse recall. or reverse.  
$$F_1 = \frac{2}{\frac{1}{P} + \frac{1}{R}}$$

#### Another example

Algorithm	US	China	India	Other
A	3%	7%	5%	9%
B	5%	6%	5%	10%
C	2%	3%	4%	5%
D	5%	8%	7%	2%
E	4%	5%	2%	4%
F	7%	11%	8%	12%

Andrew Ng

what if the algorithms have different performance (ranking)  
in different test sets (US, China, India, Other)?

Take average



deeplearning.ai

## Setting up your goal

### Satisficing and optimizing metrics

#### Another cat classification example

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

Cost = accuracy - 0.5 × runningTime

Maximize Accuracy  
Subject to runningTime ≤ 100 ms.

N metric:  
1 optimizing  
N-1 satisfying

optimizing → Accuracy  
Satisficing → Running time

Voice  
Wakewords / Trigger words  
Alexa, OK Google.  
Hey Siri, nihao.baidu.  
你多百度

accuracy.  
#false positive

Minimize accuracy.  
st. # false positive  
energy 24 hours.

What to do? Criteria: Accuracy, run time

what to do? Define cost function, then set constraint  
fn - run time  $\leq$  100ms, fix max accuracy.

Andrew Ng

traindevtest



deeplearning.ai

## Setting up your goal

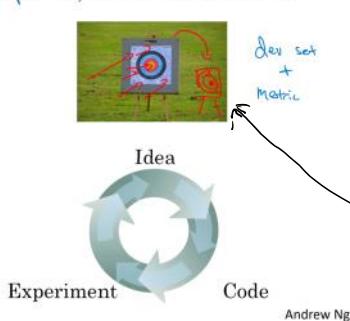
### Train/dev/test distributions

## Cat classification dev/test sets

Regions:

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

→ Randomly shuffle into dev/test



Dev and test sets should come from  
the same distribution

otherwise ML learn models that  
do not get good performance  
on the test set (WRONG TARGET)

True story (details changed)

A real world example

[ Optimizing on dev set on loan approvals for  
medium income zip codes

$$\uparrow \quad x \rightarrow y \text{ (repay loan?)}$$



[ Tested on low income zip codes

~3 month

Andrew Ng

## Guideline

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



Andrew Ng



deeplearning.ai

## Setting up your goal

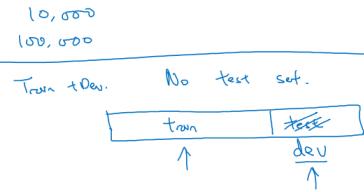
## Size of dev and test sets

rule of thumb:

reasonable for smaller size of data.

## Size of test set

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



can do train + dev. with No test set  
⇒ Not recommended

Andrew Ng

Take away: the 7/3 rule ~~is~~ no longer work  
trend: assignment ~~use~~ more data to training set.

changeset



deeplearning.ai

## Setting up your goal

### When to change dev/test sets and metrics

If your evaluation matrix don't give you results that make sense  
⇒ use a different matrix

### Cat dataset examples

Metric + Dev : Prate A  
You/users : Prater B.

Criteria  
在评价标准方面有什么需求？

→ Metric: classification error

Algorithm A: 3% error

pornographic

不注重分类把两个混淆起来？

✓ Algorithm B: 5% error

$$\left\{ \begin{array}{l} \text{Error: } \frac{1}{m} \sum_{i=1}^m \sum_{j=1}^{n_{\text{out}}} w^{(j)} \delta \left[ y_{\text{pred}}^{(i)} + y^{(i)} \right] \\ \rightarrow w^{(j)} = \begin{cases} 1 & \text{if } x^{(i)} \text{ is non-porn} \\ 0 & \text{if } x^{(i)} \text{ is porn} \end{cases} \end{array} \right.$$

penalize to porn

Andrew Ng

misclassify a cat as non-porn.

## Orthogonalization for cat pictures: anti-porn

- 1. So far we've only discussed how to define a metric to evaluate classifiers. ← Plane target ↗
- 2. Worry separately about how to do well on this metric. ↗  
↖ An (shot at target)

$$\rightarrow J = \frac{1}{\sum w^2} \sum_{i=1}^m w^T L(\hat{y}_i, y_i)$$



Andrew Ng

1. Define metrics
2. Optimize on this matrix

## Another example

Algorithm A: 3% error

✓ Algorithm B: 5% error ↗



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.

Andrew Ng

WRONG Target  
or. correct target, ill-defined in the evaluation metrics  
Even you can't decide: set up quickly some evaluation  
matrix. try ⇒ adjust the metrics, ⇒ train model  
again

compare



## Comparing to human-level performance

### Why human-level performance?

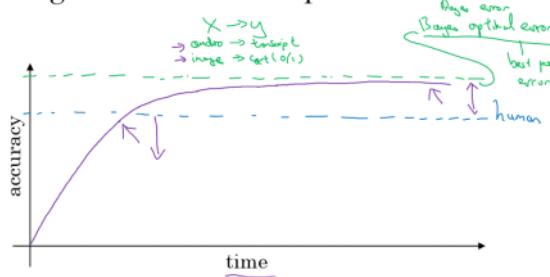
## Comparing to human-level performance

→  $x \rightarrow y$   
→ audio → transcript  
→ image → text (of it)

Bayes error  
Bayes optimal error  
loss function

Bayes Optimal error; the best possible error  
e.g. Some audio is just too noisy; some images are too blurry -  
... or reaches slow down after surpasses human performance

## Comparing to human-level performance



Andrew Ng

Bayes Optimal  
e.g. some  
prog reaches slow down after surpasses human performance  
if performance worse than human-level performance,  
can diagnose borrowing from human

## Why compare to human-level performance

Humans are quite good at a lot of tasks. So long as ML is worse than humans, you can:

- - Get labeled data from humans.  $(x, y)$
- - Gain insight from manual error analysis: Why did a person get this right?  
*manual error analysis*
- - Better analysis of bias/variance.

Andrew Ng

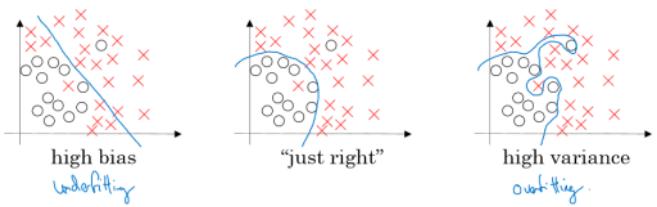
avoidablebias



Comparing to human-level performance

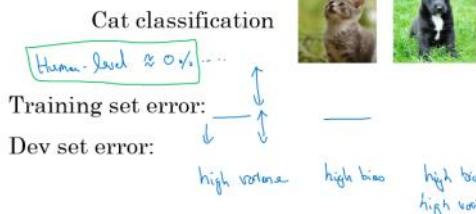
Avoidable bias

## Bias and Variance



Andrew Ng

## Bias and Variance



Andrew Ng

*Dev set - train err = Variance*

*High Variance  $\Rightarrow$  High Dev set error*

## Cat classification example

Humans (0% Bayes)	Training error	Dev error
1%	8%	10%
$\downarrow$	$\downarrow$	$\downarrow$
Focus on bias	Focus on variance	Focus on variance

Another example

Training error  $\downarrow$   $\Rightarrow$  Human-level error  $\downarrow$ .

① Focus on reducing bias.

②  $\downarrow$  Training error  $\Rightarrow$  Human-level error  $\downarrow$ .

③  $\downarrow$  Training error  $\Rightarrow$  Dev set error  $\downarrow$   $\Rightarrow$  Variance (dev set error)  $\downarrow$ .

④  $\downarrow$  Dev set error  $\Rightarrow$  focus on  $\downarrow$  Variance (regulation, ...).

*Vary and across tasks*

Andrew Ng



## Comparing to human-level performance

### 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?

How to define



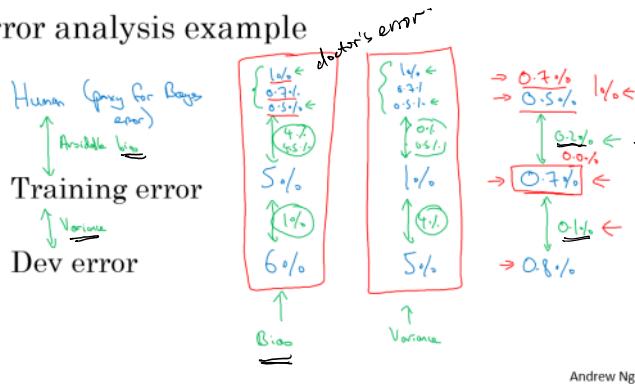
Radiologist

Andrew Ng

It is a matter of what you choose as the benchmark.

define human-level error as a proxy for Bayes error.

#### Error analysis example

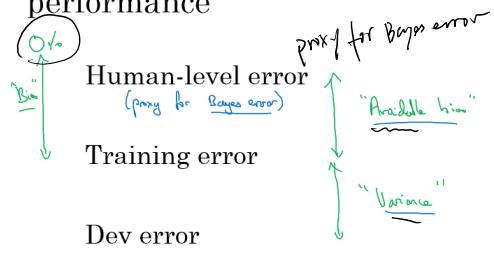


Andrew Ng

What if all bad?  
i.e. can't improve training error any more  
⇒ variance ⇒ avoidable bias is a bigger problem

(D<sup>bias</sup>)  
bad train set performance → human level performance  
test set performance  
→ work? Focus on bias & variance reduction.

## Summary of bias/variance with human-level performance



Instead of comparing [with 0% — bias] with human-performance — avoidable bias.

Human-level performance  $\Rightarrow$  Bayes Error  $\Rightarrow$   
Decide whether to focus on [bias variance] reduction

Andrew Ng  
↓  
↓ Human-level performance Th  
↓ W[Bay] 2.3.

~~Noisy data~~: Having a better estimate of Bayes error helps knowing what the "available error".



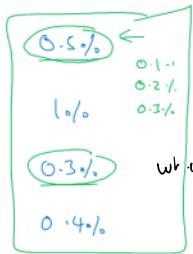
Comparing to human-level performance

Surpassing human-level performance

## Surpassing human-level performance

Team of humans	0.5%
One human	0.1%
Training error	0.6%
Dev error	0.8%

What is available bias?



Andrew Ng

what's it means? ① over-fitting  
② the Bayes error is actually lower.

↓ That slows down the speed to make progress

## Problems where ML significantly surpasses human-level performance

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

Structured data  
Not natural perception  
lots of data

{  
- Speech recognition  
- Some image recognition  
- Medical (Medical)  
- ECG, Skin cancer, ...  
↓  
in some cases, can  
surpass human-level  
performance  
Andrew Ng

ML & human {  
→ {  
}

what these share in common

structured data

e.g. what you bought.  
your loan record.

Not natural perception problem (?)

improve model



deeplearning.ai

## Comparing to human-level performance

## Improving your model performance

### The two fundamental assumptions of supervised learning

low avoidable bias

1. You can fit the training set pretty well.



→ Avoidable bias

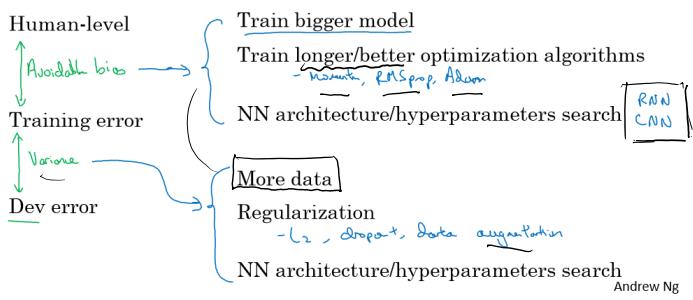
2. The training set performance generalizes pretty well to the dev/test set.



→ Variance not bad

Andrew Ng

## Reducing (avoidable) bias and variance

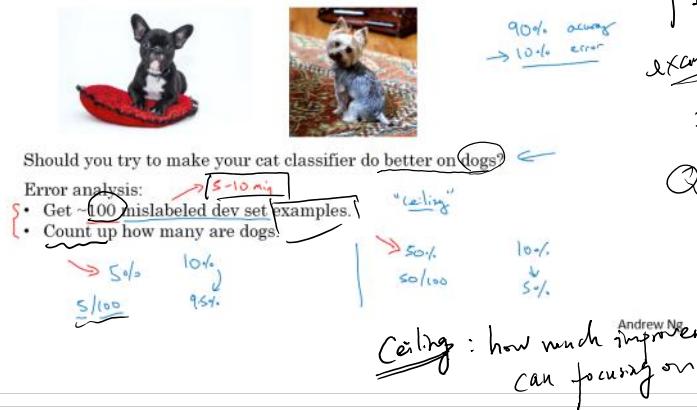


error analysis

## Error Analysis

### Carrying out error analysis

Look at dev examples to evaluate ideas



手写笔记:

Example: cat classifier to high accuracy.

Focus on "dog pictures".

Q. Is it worth doing it — focusing on dogs

What to do? Examine if dogs account for a large proportion of mislabeled pictures.

This won't take long

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

OK.. 找出并解决这些问题  
(是更重要的)

- 卷积神经网络
- Fix pictures of dogs being recognized as cats ←
  - Fix great cats (lions, panthers, etc..) being misrecognized ←
  - Improve performance on blurry images ← filter

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

Andrew Ng

work on them separately

Create a table for this  
将这些放入表格

清理 up



## Error Analysis

### Cleaning up Incorrectly labeled data

#### Incorrectly labeled examples

x							
y	1	0	1	1	0	1	1

Training set

mislabeled

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

Very robust  
Systematic errors

need not worry much

e.g. labeled all white dogs as cats

How about incorrectly  
labeled data in dev or test sets.

#### Error analysis

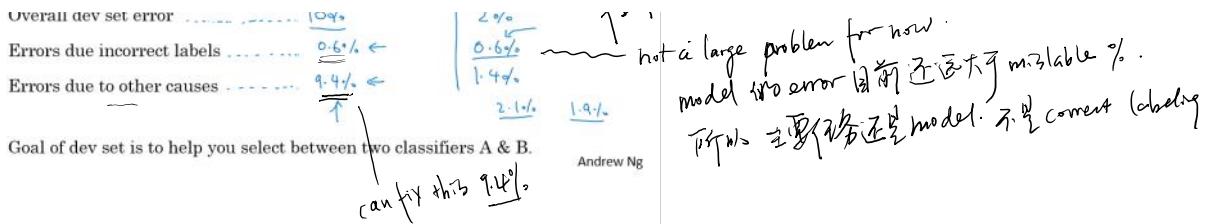
Image	Dog	Great Cat	Blurry	Incorrectly labeled	Comments
...					
98				✓	Labeled missed cat in background
99		✓			
100				✓	Drawing of a cat; Not a real cat.
% of total	8%	43%	61%	6%	

Overall dev set error 10%

Errors due incorrect labels 0.6%

Errors due to other causes 9.4%

not a large problem for now  
model has error 10% but only 6% mislabel %



## 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.

Training set is too large.  
 → sometimes you don't want  
 to manually correct advice  
 Andrew Ng

} if only fix those getting wrong - overconfidence or performance.

there are tricks to deal with this

Advice: Manual data analysis is important in ML practice  
 Sit down and look at the data.  
 good use of time → help to decide what direction to prioritize

build



## Error Analysis

Build your first system quickly, then iterate

lots of things to do...

### Speech recognition example



- • Noisy background
- • Café noise
- • Car noise
- • Accent
- • Far from build system

- • Set up dev/test set and metric
- Build initial system quickly
- Use Bias/Variance

... analysis helps to examine

- Car noise
  - Accent
  - Far from
  - Young
  - Stutter
  - ...
- Guideline:**  
**Build your first system quickly, then iterate**

- Build initial system quickly
- Use Bias/Variance analysis & Error analysis to prioritize next steps.

Andrew Ng

error analysis helps to examine where you get wrong

But if there's a literature telling you where to start, it helps to start with complex systems.

- Dont overthink when starting
- Build something that works (quick & dirty), then fix it using Bias/Variance analysis & error analysis

diff distribution



deeplearning.ai

## Mismatched training and dev/test data

### Training and testing on different distributions

#### Cat app example



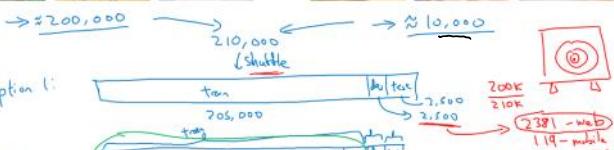
Data from webpages



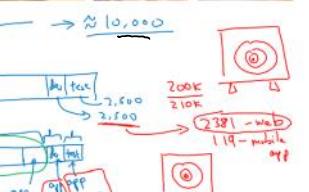
Data from mobile app

care about this  
train model and test it on data from a different distribution.

X Option 1:

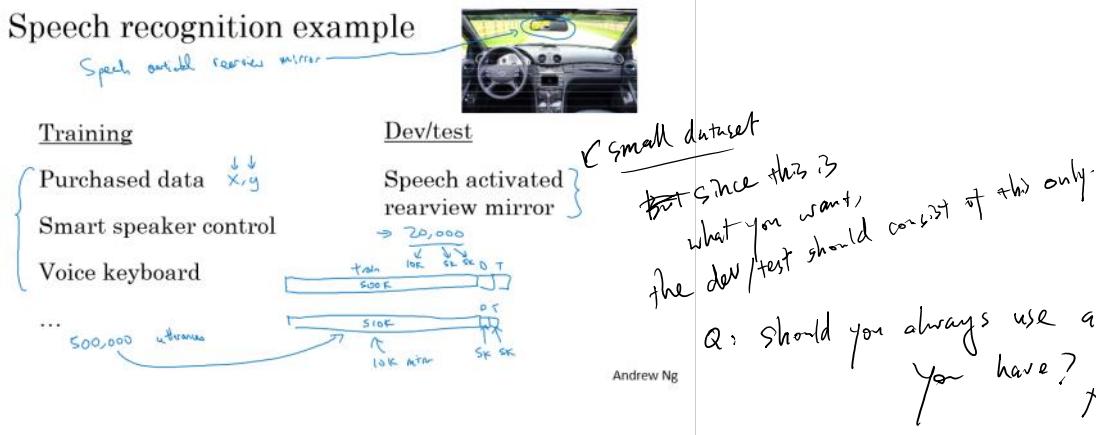


Option 2:



WRONG TARGET. optimize for a different dist of data from what you care about  
All test should be from the app (what you ultimately care about)

## Speech recognition example



Q: should you always use all the data you have?  
Not really

mismatch



## Mismatched training and dev/test data

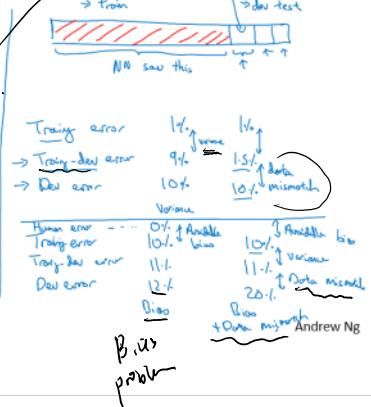
### Bias and Variance with mismatched data distributions

#### Cat classifier example

Assume humans get  $\approx 0\%$  error.

Training error ..... 1%  
Dev error ..... 10%

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



may not be variance problem  
maybe the dev set is just larger  
Data mismatch

这像似于数据分布不一.  
利用一些和分布相关的技巧  
拟合模型的技巧

## Bias/variance on mismatched training and dev/test sets

Human level	4%
Train set error	7% ↴ available bias

4% ↴

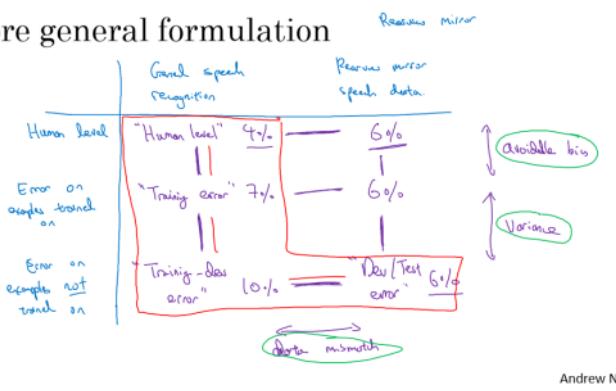
## Bias/variance on mismatched training and dev/test sets

Human level	4%	available bias
Training set error	7%	Variance
Training-dev set error	10%	Data mismatch
→ Dev error	12%	degree of similarity to dev set.
→ Test error	12%	

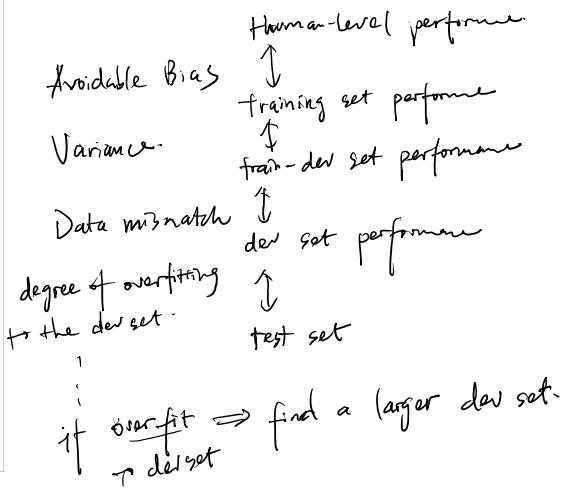
4%
7%
10%
6%
6%

Andrew Ng

## More general formulation

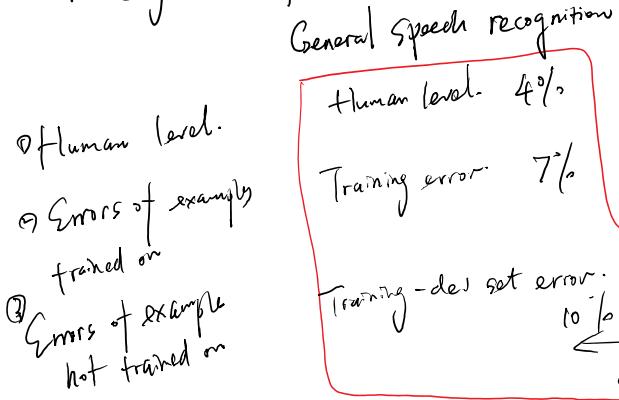


Andrew Ng

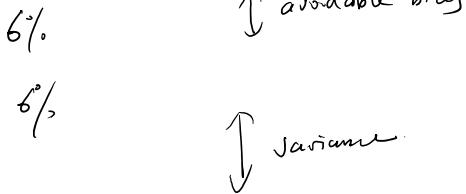


What if you do better in dev/test set than on training set?  
one possibility, dev/test set is easier tasks (e.g. clearer photo)

## More general formulation



Review mirror speech recognition.



How to deal with data mismatch? up next - not much -



## Mismatched training and dev/test data

### Addressing data mismatch

#### Addressing data mismatch

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

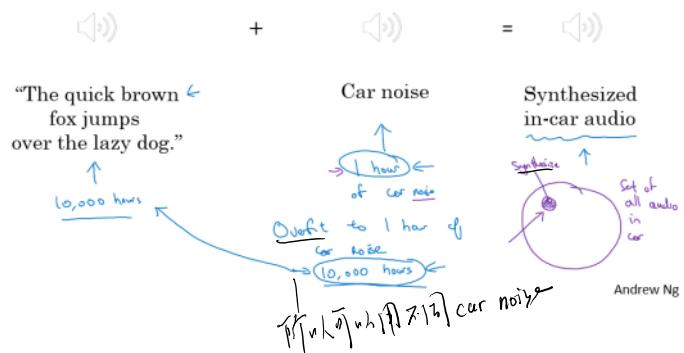
E.g. noisy - car noise      street numbers
- • Make training data more similar; or collect more data similar to dev/test sets
 

E.g. Simulate noisy in-car data

Andrew Ng

Key: manual analysis  
error  
 "How your dev/test sets are different from  
 your training set?"  
 e.g. ~~too~~ too noisy

#### Artificial data synthesis

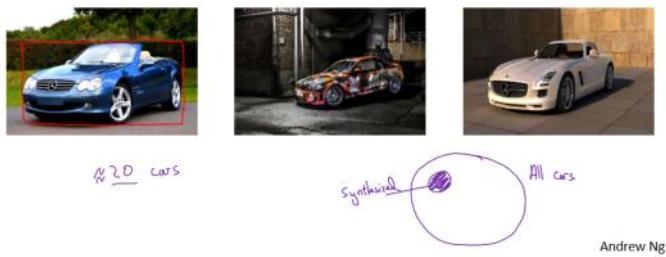


SYNTHESIS

trick:用一隻car noise  
 全部過fit到 certain  
 car noise

## Artificial data synthesis

Car recognition:



recognize car  
从 car video game 里找车  
从 game 里找车的样本

Take away: ① Error analysis  
② Artificial synthesis. — caution: overfit.

transfer learning

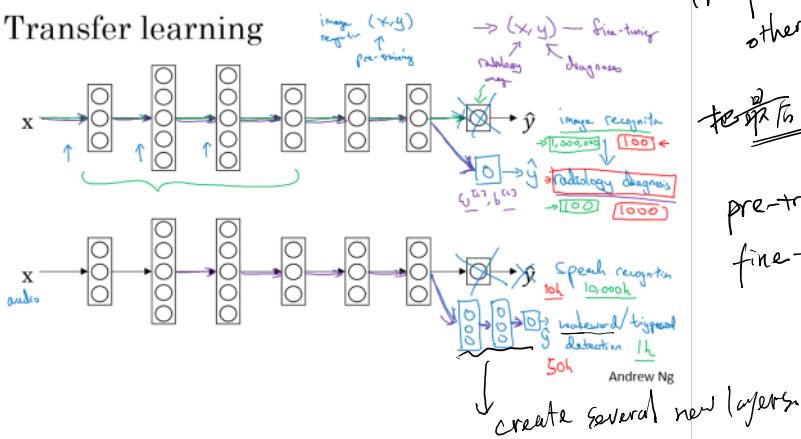
## Learning from multiple tasks

### Transfer learning



deeplearning.ai

### Transfer learning



transfer learned model to  
other task

pre-train -> re-train

pre-training: 用别的数据训练模型  
fine-tuning.

If you want to collect more data.  
1 . . . 1 - train ->

creatin over . . .

If you want to collect more data.  
Data of the task of interest  
more important.

## When transfer learning makes sense

Train from A  $\rightarrow$  B

- Task A and B have the same input  $x$ .
- You have a lot more data for Task A than Task B.
- Low level features from A could be helpful for learning B.

learn lower-level features

Andrew Ng

multitask



deeplearning.ai

## Learning from multiple tasks

### Multi-task learning

## Simplified autonomous driving example



outcome  $y = \begin{bmatrix} y_1^{(i)} & y_2^{(i)} & y_3^{(i)} & \dots & y_n^{(i)} \end{bmatrix}$

$y_1^{(i)}$   $y_2^{(i)}$   $y_3^{(i)}$   $\dots$   $y_n^{(i)}$

Pedestrians  
Cars  
Stop signs  
Traffic lights  
 $\vdots$

$(4, n)$

$x^{(i)}$

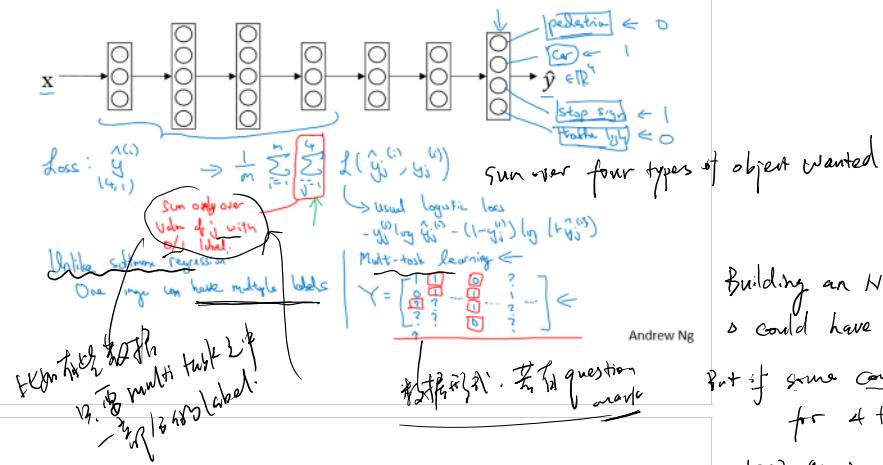
$y^{(i)}$

$(4, 1)$

$\underbrace{\quad}_{\text{Car}} \quad \underbrace{\quad}_{\text{Pedestrian}}$

Andrew Ng

## Neural network architecture



## 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.

Andrew Ng

When does multi-task learning make sense?  
If we have a lot of data for one task, we can train a large model.  
CV - multitask outperform single task  
⇒ NG (7.2.6)  
→ If we train a big model,  
we can't separate  
the different objects  
Hence, what's the objects?  
What's the task?  
Transfer learning 相互学习  
multi-task learning 效率高.

endtoend



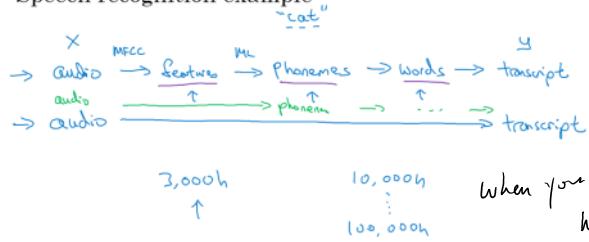
deeplearning.ai

End-to-end deep learning  
What is  
end-to-end  
deep learning

Multiple stages  $\Rightarrow$  one stage

## What is end-to-end learning?

Speech recognition example



Why break (1, 2, 3) stages for DL? ML - 3 stages

3,000h  
↑  
10,000h  
↓  
100,000h

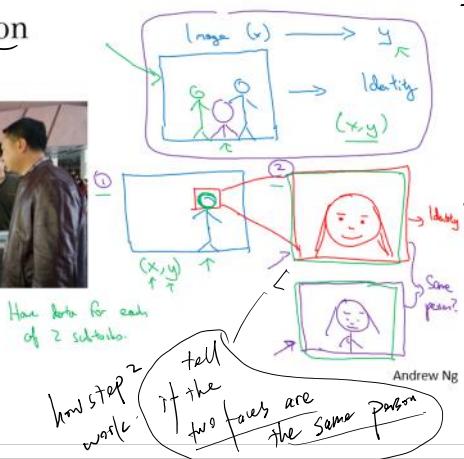
When you have a large dataset, the end-to-end approach has advantage.  
if small data, traditional approach is better.

Andrew Ng

## Face recognition



[Image courtesy of Baidu]



In this example, the two step

approach works better?

- ① Two tasks are both simpler (much simpler). (2/3)
- ② lots of data for each ~~step~~ of the steps respectively. (2/3)  
But much for the end-to-end task

multi step

- ① Detect where the person's face is
- ② Crop image, center the person's face
- ③ Recognize the face

Two tasks — this

## More examples

Machine translation



Estimating child's age: Image  $\xrightarrow{\text{① bones}} \xrightarrow{\text{② age}}$



Image  $\xrightarrow{\text{age}}$

Andrew Ng

End-to-end DL.  
works, but not a parallel.



## End-to-end deep learning

### Whether to use end-to-end learning

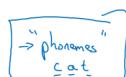
End-to-end learning

Just let the data speak.

#### Pros and cons of end-to-end deep learning

Pros:

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



no need to reflect human pre-conception.  
Some other pre conception may be ~~wrong~~ wrong

Cons:

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



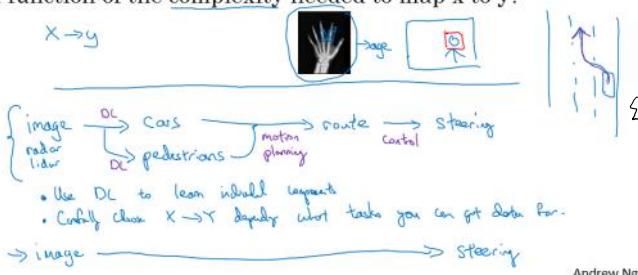
Need lots of data

An algorithm has two sources of knowledge

[data → when small human knowledge is more important.  
human design]

#### 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  $y$ ?



for simple task → easily using end-to-end DL with small datasets

Example: Autonomous driving

for this task, end-to-end dl does not work well.