

Motivating example

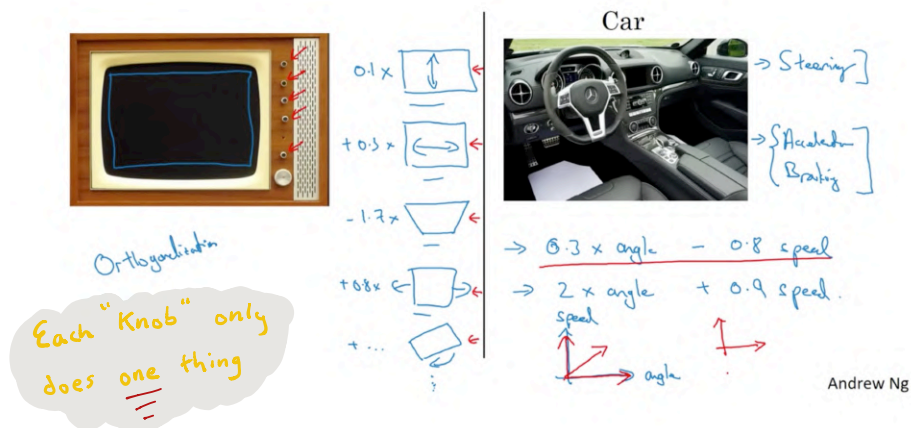


Ideas:

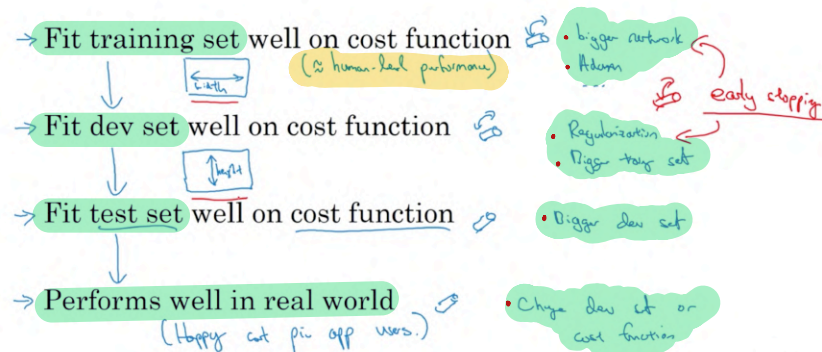
- Collect more data
- Collect more diverse training set
- Train algorithm longer with gradient descent
- Try Adam instead of gradient descent
- Try bigger network
- Try smaller network
- Try dropout
- Add L_2 regularization
- Network architecture
 - Activation functions
 - # hidden units
 - ...

Andrew Ng

TV tuning example

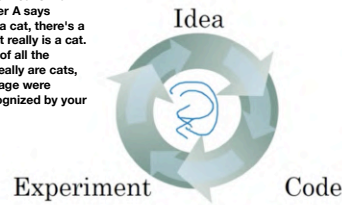


Chain of assumptions in ML



Using a single number evaluation metric

So if classifier A has 95% precision, this means that when classifier A says something is a cat, there's a 95% chance it really is a cat. And recall is, of all the images that really are cats, what percentage were correctly recognized by your classifier?



Of samples recognized as cat, what % actually are cats?
 What % of actual cats are correctly recognized?

Classifier	Precision	Recall	F1 Score
A	95%	90%	92.4%
B	98%	85%	91.0%

F1 score = "Average" of P and R.

$$\left(\frac{2}{\frac{1}{P} + \frac{1}{R}} \right)$$
 "Harmonic mean"

Real set + Single number evaluation metric
 real speed up iterating

Andrew Ng

Another example

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%

Error

Andrew Ng

Another cat classification example

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

Somewhat artificial

"Better..."

$$\text{Cost} = \text{accuracy} - 0.5 \times \text{Running Time}$$

Maximize Accuracy
 Subject to Running Time ≤ 100 ms.

N metrics: 1 optimizing
 N-1 satisfying

Wakewords / Trigger words

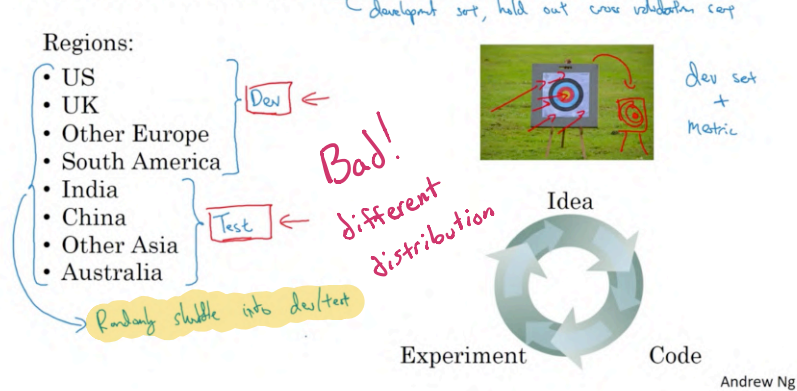
Alexa, OK Google,
 Hey Siri, 嘿 Siri, 你好 Siri
 你好 Siri

Accuracy.
 #False positive

Minimize accuracy ← Optimizing
 s.t. ≤ 1 false positive
 every 24 hours. ← Satisficing

Andrew Ng

Cat classification dev/test sets



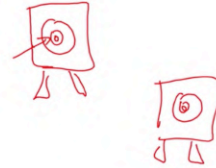
True story (details changed)

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

$x \rightarrow y$ (repay loan?)

Tested on low income zip codes

~3 months



Andrew Ng

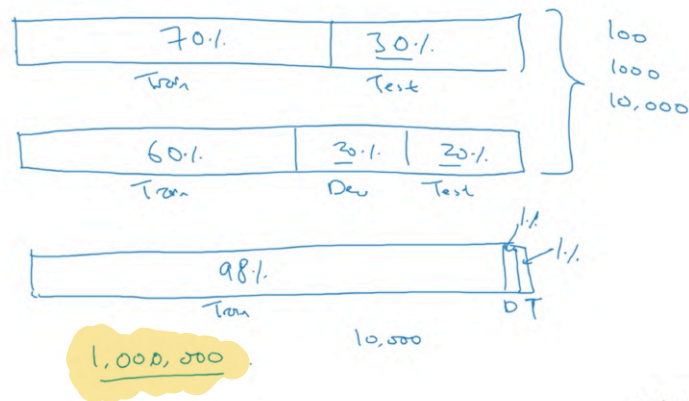
Guideline

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

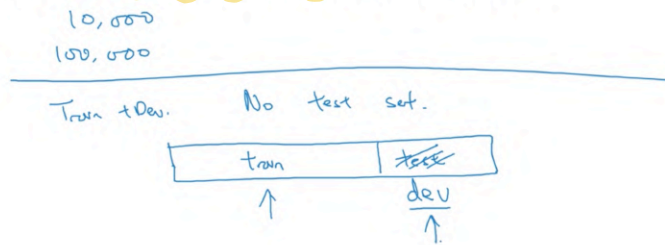
Old way of splitting data



Andrew Ng

Size of test set

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



Andrew Ng

Cat dataset examples

Metric + Dev : Prefer A
You/users : Prefer B.

- Metric: classification error

Algorithm A: 3% error → pornographic

✓ Algorithm B: 5% error

$$\text{Error} = \frac{1}{\sum w^{(i)}} \sum_{i=1}^{m_{\text{dev}}} w^{(i)} \mathbb{I}\{y_{\text{pred}}^{(i)} \neq y^{(i)}\}$$

→ $w^{(i)} = \begin{cases} 1 & \text{if } x^{(i)} \text{ is non-porn} \\ 10 & \text{if } x^{(i)} \text{ is porn} \end{cases}$

↑ predicted value (0/1)

Sign that you should change the evaluation metric or dev/test set

high level of take away is, if you find that evaluation metric is not giving the correct rank order preference for what is actually better algorithm, then there's a time to think about defining a new evaluation metric.

Andrew Ng

Orthogonalization for cat pictures: anti-porn

- 1. So far we've only discussed how to define a metric to evaluate classifiers. ← Place target
- 2. Worry separately about how to do well on this metric. ← Aim (shoot at target)

$$\rightarrow J = \frac{1}{\sum w_i} \sum_{i=1}^n w_i \ell(y^{(i)}, \hat{y}^{(i)})$$



Andrew Ng

Another example

Algorithm A: 3% error

✓ Algorithm B: 5% error ←

→ Dev/test



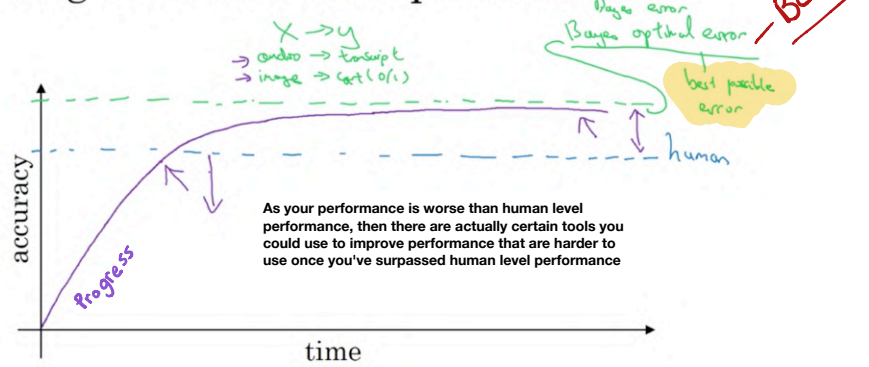
→ User images



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.

Andrew Ng

Comparing to human-level performance



Andrew Ng

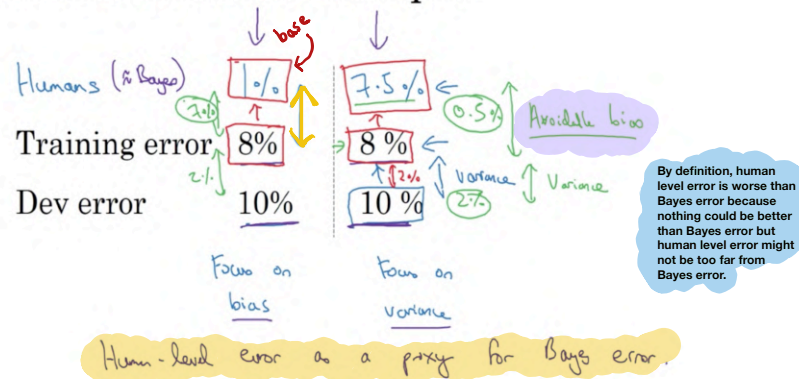
Why compare to human-level performance

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?
- Better analysis of bias/variance.

Andrew Ng

Cat classification example



Andrew Ng

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 ←

Be clear when defining human level error

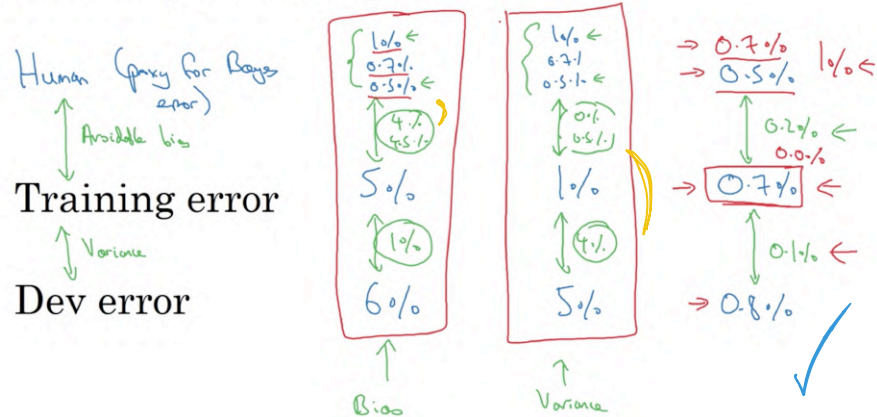


Bayes error $\leq 0.5\%$

What is "human-level" error?

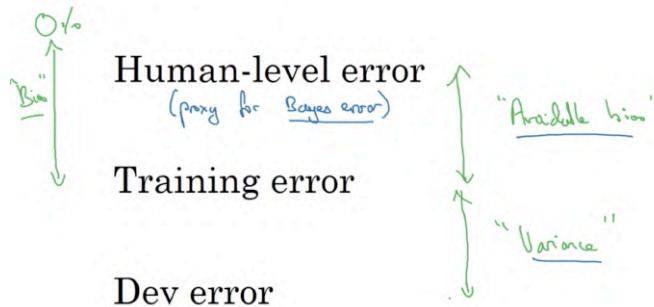
Andrew Ng

Error analysis example



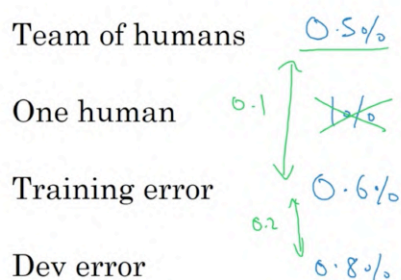
Andrew Ng

Summary of bias/variance with human-level performance



Andrew Ng

Surpassing human-level performance



What is avoidable bias?

Not enough info

Andrew Ng

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
 - ECG, Skin cancer, ...

Andrew Ng

The two fundamental assumptions of supervised learning

1. You can fit the training set pretty well.

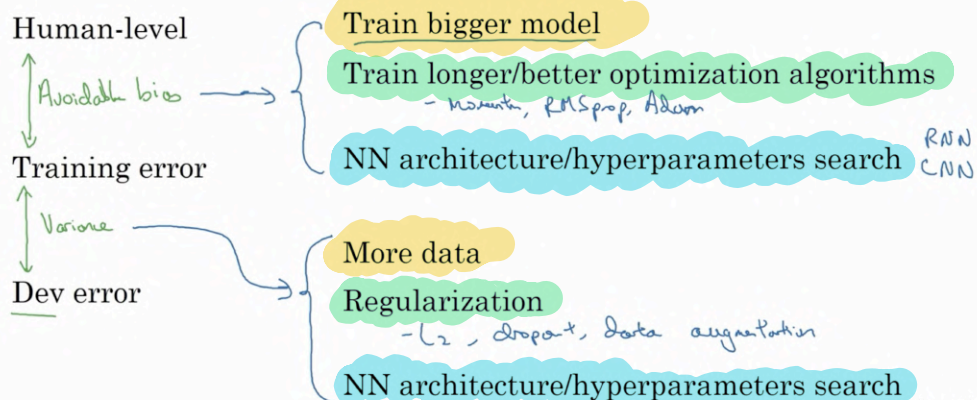
^{low}
~ Avoidable bias

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

~ Variance

Andrew Ng

Reducing (avoidable) bias and variance



Andrew Ng