

Motivating example



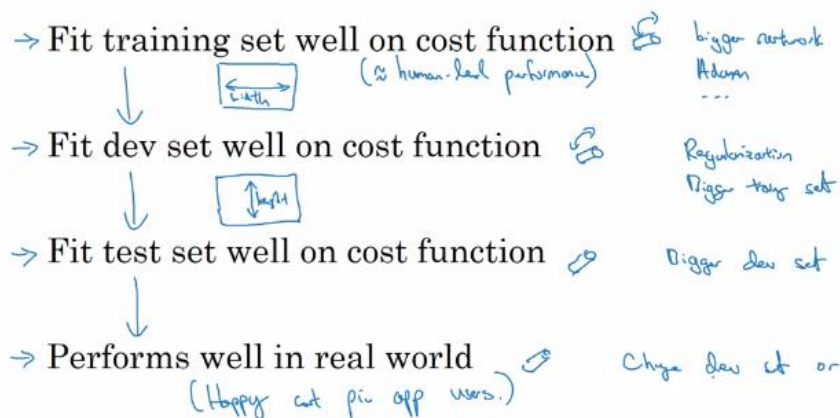
90%.

Ideas:

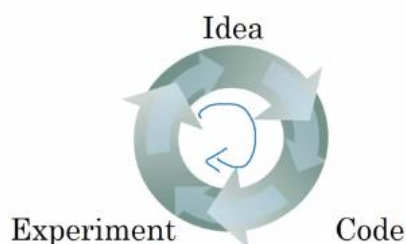
- Collect more data
- Collect more diverse training set
- Train algorithm longer with gradient descent
- Try Adam instead of gradient descent
- Try bigger network
- Try smaller network
- Try dropout
- Add L_2 regularization
- Network architecture
 - Activation functions
 - # hidden units
 - ...

Andrew Ng

Chain of assumptions in ML



Using a single number evaluation metric



Of examples recognized as cats, 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"

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%

How could we combine an optimizing and satisficing metrics?

Another cat classification example

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

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

Maximize Accuracy
 Subject to Running Time \leq 100 ms.

Similar examples of wake words/trigger words

Wakewords / Trigger words

Okra, OK Google,

Hey Siri, nihao baidu
 你好百度

accuracy.
 #false positive

Maximize accuracy.
s.t. \leq 1 false positive
every 24 hours.

Train/dev/test distributions

Cat classification dev/test sets

development set, hold out cross validation set

Regions:

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

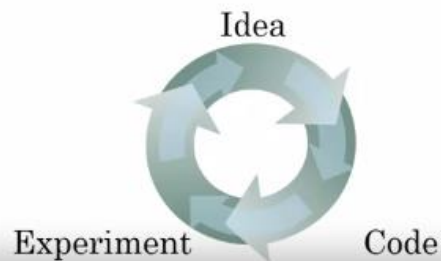
Dev

Test

Randomly shuffle into dev/test



dev set
+
metric



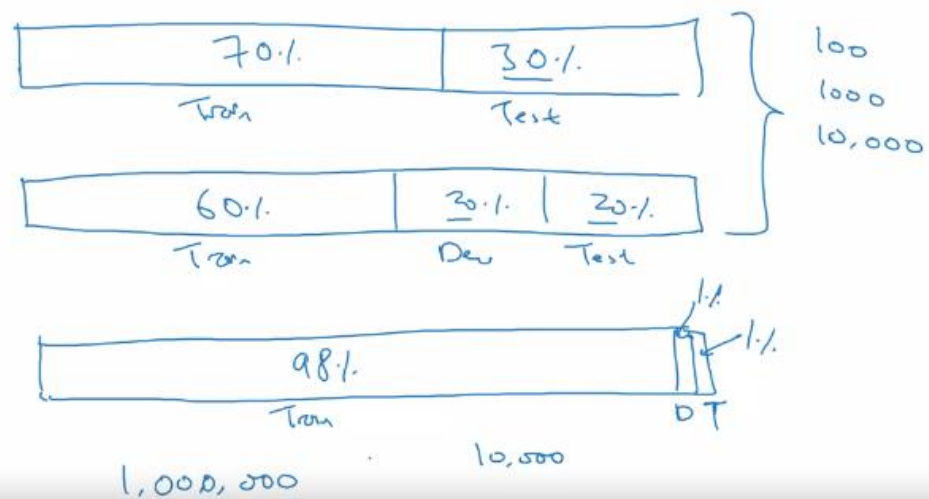
Andrew Ng

Guideline

Same distribution

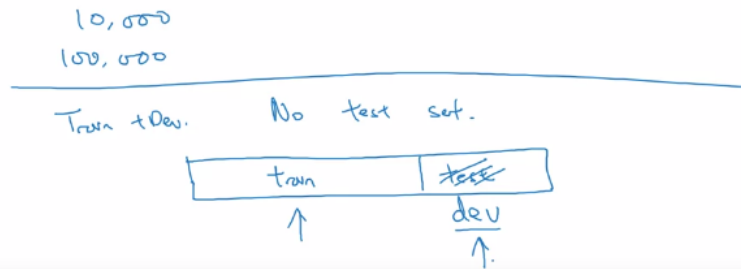
Choose a dev set and test set to reflect data you expect to get in the future and consider important to do well on.

Old way of splitting data



Size of test set

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



An unconventional option might be to only have train & dev sets, with dev set large enough that you won't overfit the train set.

In the blow, algo A is doing better with lower error, but we may also want to cost it for displaying the occasional pornographic image:

When to change dev/test sets and metrics

Cat dataset examples

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

Metric: classification error

Algorithm A: 3% error → pornographic

✓ Algorithm B: 5% error

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

↙ predicted value (0/1)

Weigh it with a w_i which will magnify the error by 10 if the image is pornographic

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

↙ predicted value (0/1)

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

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. ⌋
- ↑ Am (how) at target

$$\rightarrow J = \frac{1}{\sum w^{(i)}} \sum_{i=1}^M w^{(i)} \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$$



Another example

Algorithm A: 3% error

✓ Algorithm B: 5% error ←

→ Dev/test

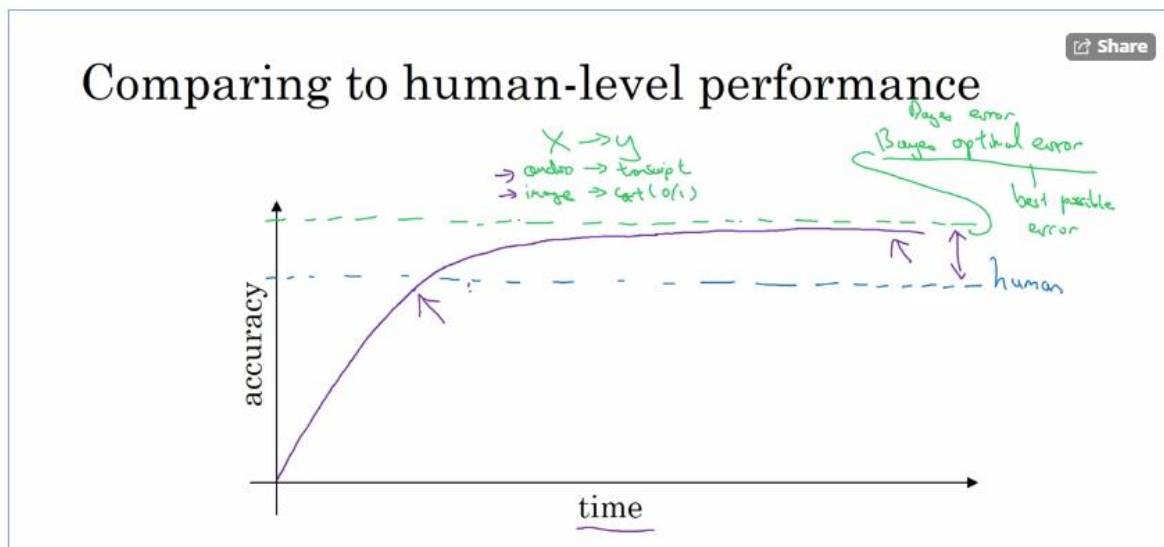


→ User images



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.

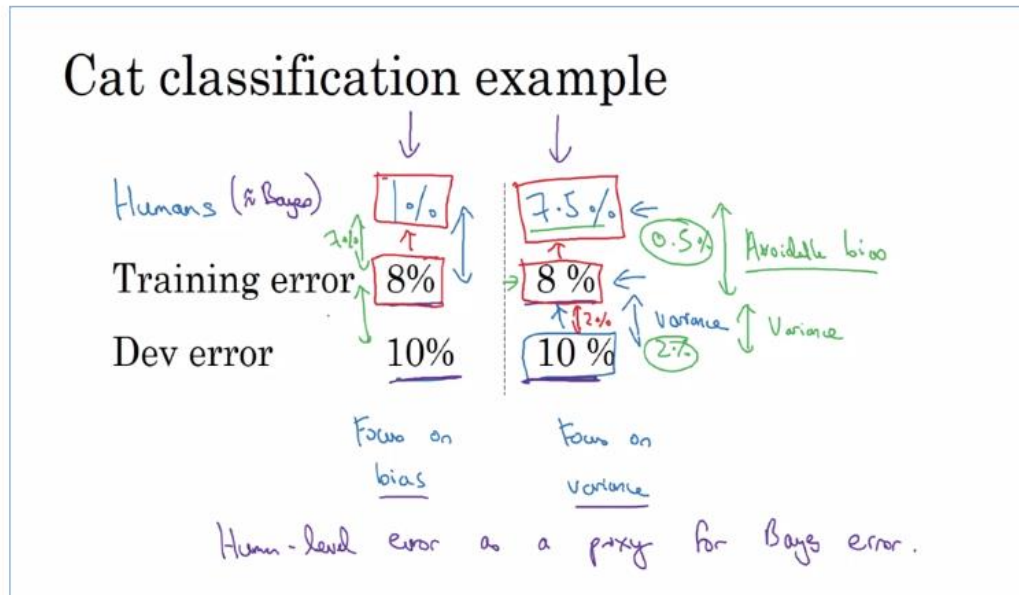
Why human-level performance?



Why the slowdown after surpassing human level performance?

1. For many tasks, human level performance is not that far from Bayes' optimal error. So, by the time you surpass human level performance maybe there's not that much head room to still improve.
2. As long as your performance is worse than human level performance, you can:
 - a. Get labeled data from humans
 - b. Gain insight from a manual error analysis
 - c. Better analysis of bias/variance

Avoidable bias



Avoidable bias: Diff between training error and the achievable min error level (Bayes error)

In scenario 2 on the right, much more scope to reduce the 2% variance than reducing the 0.5% avoidable bias.

Understanding human-level performance

Human-level error as a proxy for Bayes error

Medical image classification example:

Suppose:

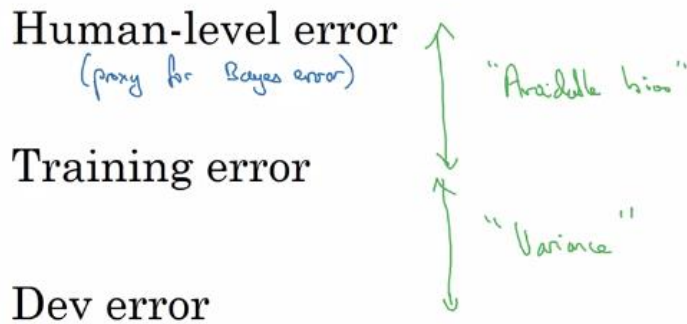
- (a) Typical human 3 % error
- (b) Typical doctor 1 % error
- (c) Experienced doctor 0.7 % error
- (d) Team of experienced doctors .. 0.5 % error



What is “human-level” error?

Bayes error \leq 0.5 %

Summary of bias/variance with human-level performance



Surpassing human-level performance

Share

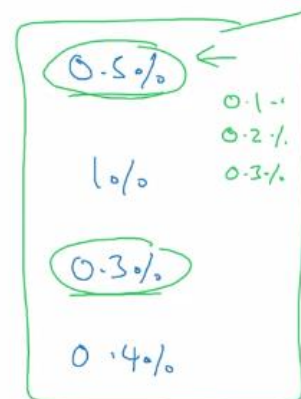
Surpassing human-level performance

Team of humans 0.5%

One human 0.1% ~~1.0%~~

Training error 0.6%

Dev error 0.8%



What is avoidable bias?

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, ...

Reducing (avoidable) bias and variance

