# Introduction to ML strategy

## Why ML Strategy?

- what's the best strategy to improve a certain DL system?

# 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

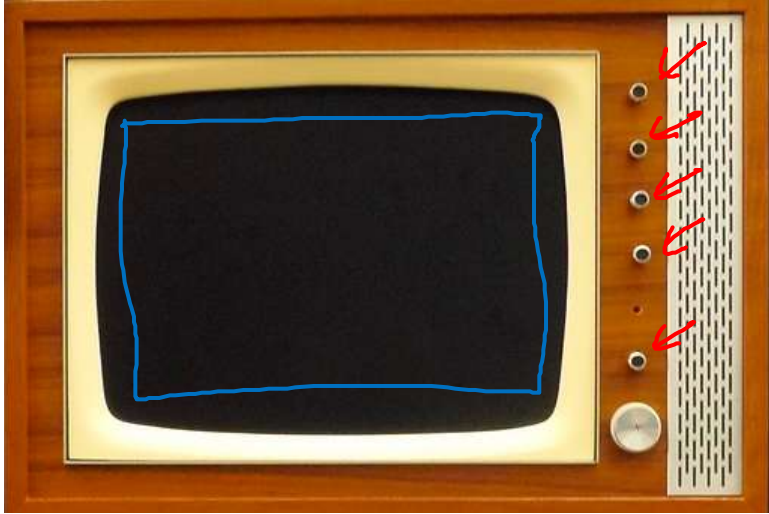  - Activation functions

  - # hidden units

  - ...

Andrew Ng

# Introduction to ML strategy

# Orthogonalization

- It's easier to change/tune the hyperparameters of the DL algorithm if their effects are independent from each other.

deeplearning.ai

# TV tuning example

Car

0.1 x

+ 0.3 x

- 1.7 x

Orthogonalization

+ 0.8 x

+ ...

→ Steering ]

→ { Accelerator
    Brarking ]

→ 0.3 x angle   -   0.8 speed

→ 2 x angle   +   0.9 speed.

speed

angle

Andrew Ng

# Chain of assumptions in ML

For ML applications we hope that

① basic fitting

→ Fit training set well on cost function

($\approx$ human-level performance)

② generalization

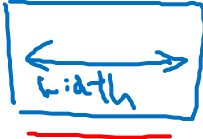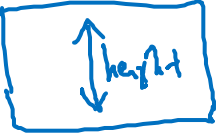→ Fit dev set well on cost function

③ performance

→ Fit test set well on cost function

④ real-world performance

→ Performs well in real world

(Happy cat pic app users.)

⑤ for ML we have orthogonal hyperparameters to try meeting each of the assumptions.

bigger network
Adam
- - -

DO NOT USE

early stopping

because it has effects on network size AND regularization and is therefore not orthogonal.

Regularization
Bigger traing set

Bigger dev set
(existing dev set is not representative enough of the training set).

Change dev set or cost function ?

Andrew Ng

- AT THE BEGINNING OF THE PROJECT ONE MUST IMMEDIATELY SET UP A DEV SET THAT REFLECTS THE REAL APPLICATION AND A PROPER EVALUATION METRICS.
- These should be reassessed periodically to compensate unforeseen side-effects.

deeplearning.ai

## Setting up your goal

---

## Single number evaluation metric

# Using a single number evaluation metric

Idea

Experiment

Code

Of examples recognized as ~~CATS~~ (cat) what % actually are cats? $\frac{TP}{TP+FP}$

What % of actual cats are correctly recognized $\frac{TP}{TP+FN}$

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

$F_1$ score = "Average" of P and R.

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

Dev set + Single number evaluation metric
real                                Speed up iterating

# Another example

| Algorithm | US | China | India | Other | AVERAGE |
|-----------|-----|-------|-------|-------|---------|
| 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

# Setting up your goal

---

# Satisficing and optimizing metrics

deeplearning.ai

Idea: one can have and combine different types of metrics. Some metrics need to be optimized (max/minimized), eg accuracy, others need only to be satisfied in a binary way (below/above a threshold), e.g. running time ≤ 100 ms.

# Another cat classification example

optimizing ↓ ↓ Satisficing

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

Cost = accuracy − 0.5 × running Time

Maximize accuracy

Subject to running Time ≤ 100 ms.
→ "Satisficing"

N metrics: 1 optimizing
N−1 satisficing

Wakewords / Trigger words

Alexa, OK Google,
Hey Siri, nihaobaidu
你好百度

accuracy.
# false positive

Maximize accuracy.
s.t. ≤ 1 false positive
every 24 hours.

e.g. for Alexa we might want to maximize the recall, while ensuring the #FP ≤ 1 every 24 hrs.

Andrew Ng

Setting up your goal

Train/dev/test distributions

deeplearning.ai

1- make sure that dev and test set come from the same distribution and represent well the real data the system will be used on.

# Cat classification dev/test sets

development set, hold out cross validation set

Regions:

- US
- UK
- Other Europe
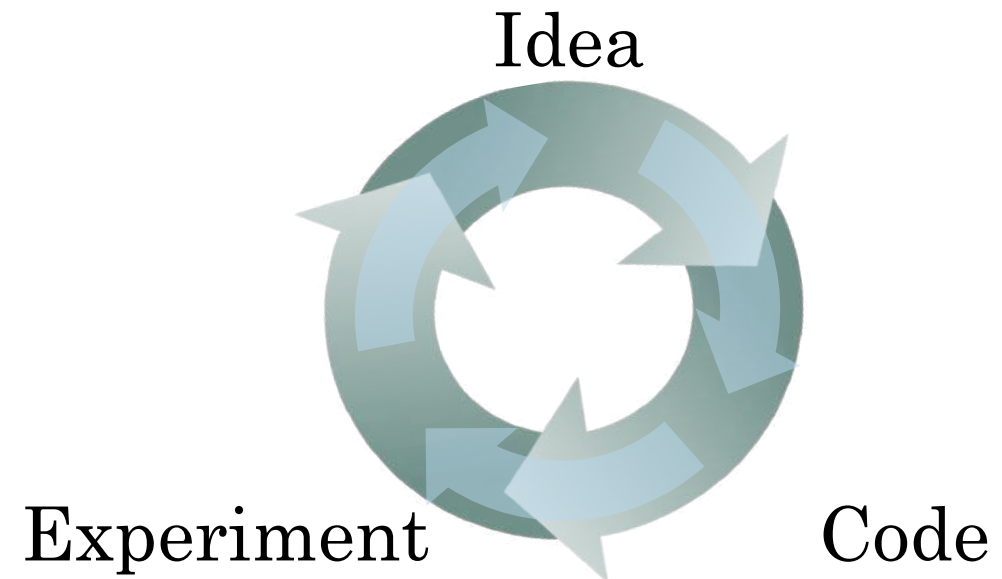- South America

→ Dev

- India
- China
- Other Asia
- Australia

→ Test

→ Randomly shuffle into dev/test

dev set + Metric

Idea

Experiment

Code

Andrew Ng

# True story (details changed)

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

$x \longrightarrow y$ (repay loan?)

Tested on low income zip codes

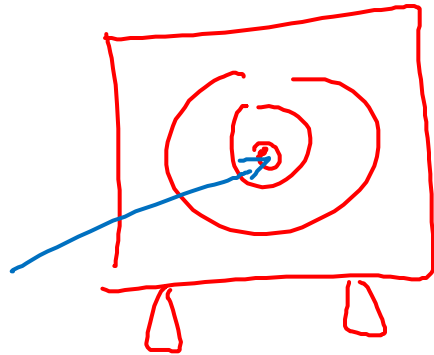~3 month

# 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.

training

dev
metric

test

deeplearning.ai

Setting up
your goal

Size of dev
and test sets

# Old way of splitting data



70%   30%
Train   Test

100
1000
10,000

60%   30%   20%
Train   Dev   Test

In huge datasets
(>1m samples) used
for DL, one can
split TR/DEV/TE
98%. 1%. 1%.

98%
Train
1%
1%
D T

1,000,000

10,500

Andrew Ng

# Size of dev set

A    B

Set your dev set to be big enough to detect differences in algorithm/models you're trying out.

100 : Small

    ⌐ 1%

1,000

10,000

100,000
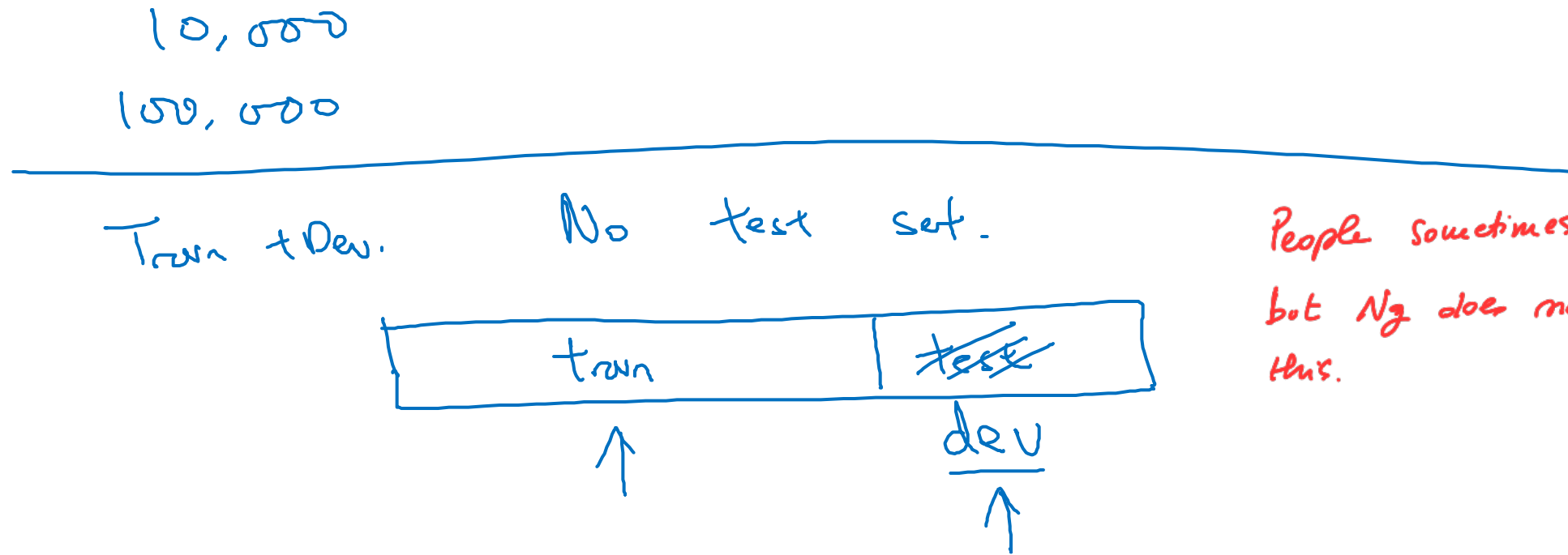
A
97% ⟶ 97.1%    B

0.1%
↑

0.01%    Online advertising
0.001%

# Size of test set

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

10,000

100,000

Train + Dev.     No test set.     People sometimes do this but Ng does not recommend this.

| train | test dev |
|-------|----------|

↑                ↑

Andrew Ng

deeplearning.ai

Setting up your goal

---

When to change dev/test sets and metrics

The performance metric and dev set must be corrected to catch unforseen side effects.

# Cat dataset examples

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

→ Metric: classification error

Algorithm A: 3% error ———→ pornographic

✓ Algorithm B: 5% error

$$\text{Error:} \frac{1}{\sum w^{(i)}} \underset{M_{dev}}{\cancel{\frac{1}{M_{dev}}}} \sum_{i=1}^{M_{dev}} w^{(i)} I\{ y^{(i)}_{pred} \neq y^{(i)} \}$$

predicted value (0/1)

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

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

$$J = \frac{1}{m} \sum_{i=1}^{m} \omega^{(i)} \mathcal{L}\left(\hat{y}^{(i)}, y^{(i)}\right)$$

$\sum \omega^{(i)}$

In case you wanna keep the metric (accuracy) you can decide to (orthogonally) modify the algorithm (loss in this case) so to counteract side effects.

⤷ (?) I don't fully get why one should change the loss but keep the old metric...

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.

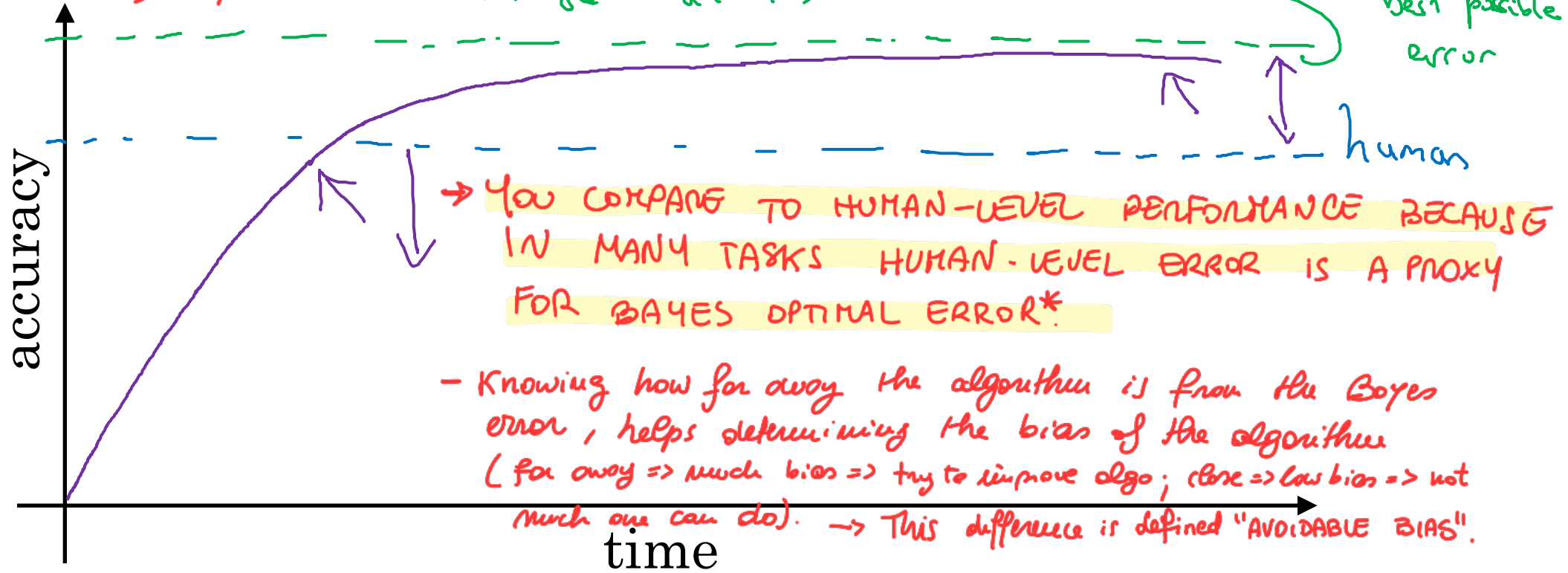Andrew Ng

deeplearning.ai

Comparing to human-level performance

---

Why human-level performance?

# Comparing to human-level performance

- It is useful to know the human level performance on tasks where humans are not great, such as speech recognition with very noisy audio.

$X \rightarrow y$

$\rightarrow$ audio $\rightarrow$ transcript

$\rightarrow$ image $\rightarrow$ cat (0/1)

Bayes error

Bayes optimal error

best possible error



human

$\rightarrow$ YOU COMPARE TO HUMAN-LEVEL PERFORMANCE BECAUSE IN MANY TASKS HUMAN-LEVEL ERROR IS A PROXY FOR BAYES OPTIMAL ERROR*

- Knowing how far away the algorithm is from the Bayes error, helps determining the bias of the algorithm (far away => much bias => try to improve algo; close => low bias => not much one can do). $\rightarrow$ This difference is defined "AVOIDABLE BIAS".

- Moreover, until the algorithm's error (training error) is below human's error, the human can try to help the algorithm further (better labeling, identifying sources of noise); when the algorithm performs better than the human, this is not possible anymore.

*BAYES OPTIMAL ERROR: prediction error of the best function theoretically possible to solve the problem.

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:

$\Rightarrow$ - Get labeled data from humans. $(x, y)$

$\Rightarrow$ - Gain insight from manual <u>error analysis:</u>
Why did a person get this right?

$\Rightarrow$ - Better analysis of bias/variance. $\rightarrow$ explained later

- Until the algorithm's performance are worse than humans, the human can do a lot to improve the algorithm's performance (labeling, identifying sources of imprecision in the algorithm, deciding if the bias of the algorithm is too high).
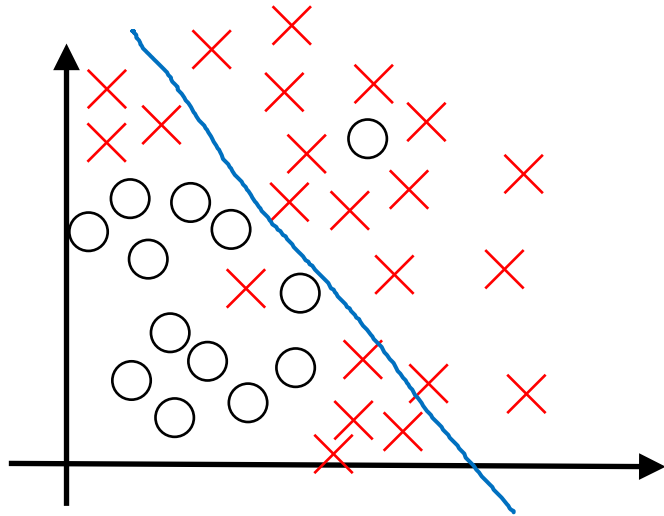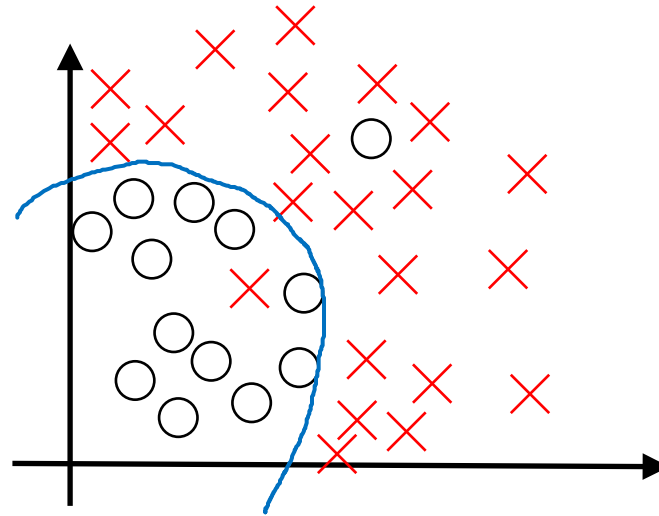
Andrew Ng

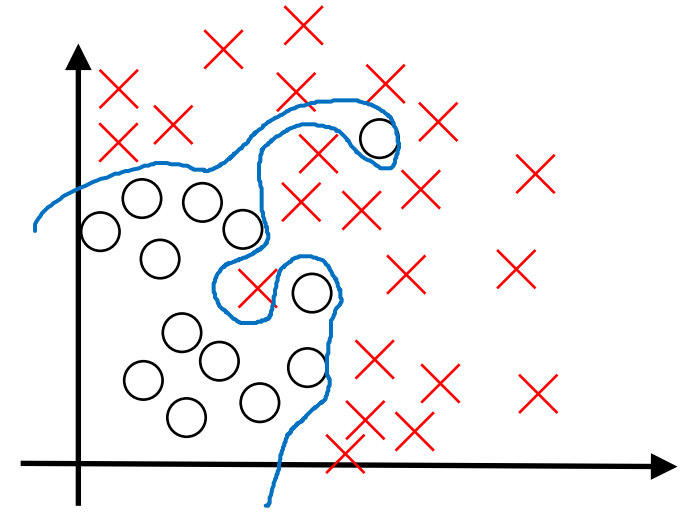Comparing to human-level performance

---

Avoidable bias

deeplearning.ai

# Bias and Variance



high bias

*underfitting*

"just right"

high variance

*overfitting.*

Andrew Ng

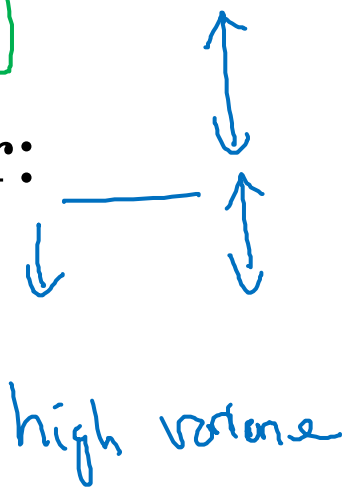# Bias and Variance



Cat classification

Human - level $\approx$ 0% ......

Training set error:

Dev set error:

high variance          high bias          high bias          low bias
                                           high variance      low variance

# Cat classification example

Humans ($\tilde{} Bayes$)

Training error

Dev error

| | | |
|---|---|---|
| | 1 % | 7.5 % |
| | 8% | 8 % |
| | 10% | 10 % |

7% · 0.5% Avoidable bias

2% · 2% Variance · Variance

Focus on bias

Focus on variance

AVOIDABLE BIAS: difference between training error and (approximately) Bayes error. This bias can likely be modified acting on the algorithm.

Human-level error as a proxy for Bayes error.

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

Bayes error $\leq 0.5\%$
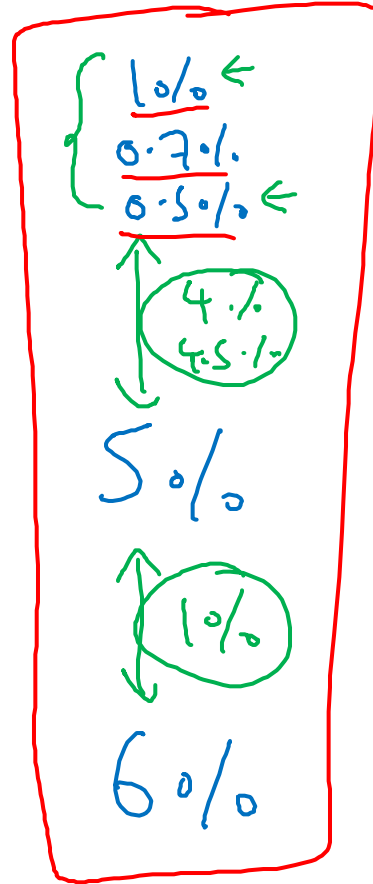
What is "human-level" error?

# Error analysis example
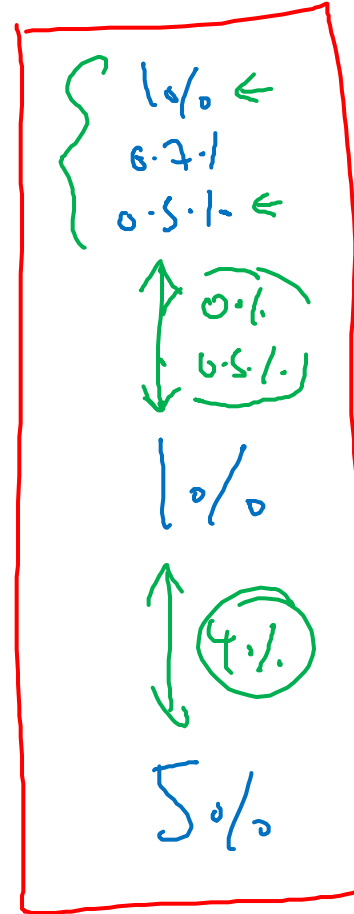


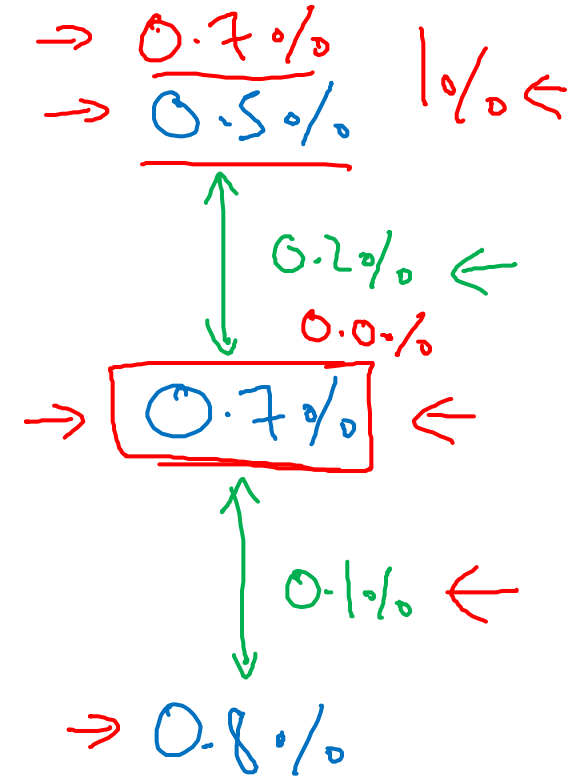Human (proxy for Bayes error)

Avoidable bias

Training error

Variance

Dev error

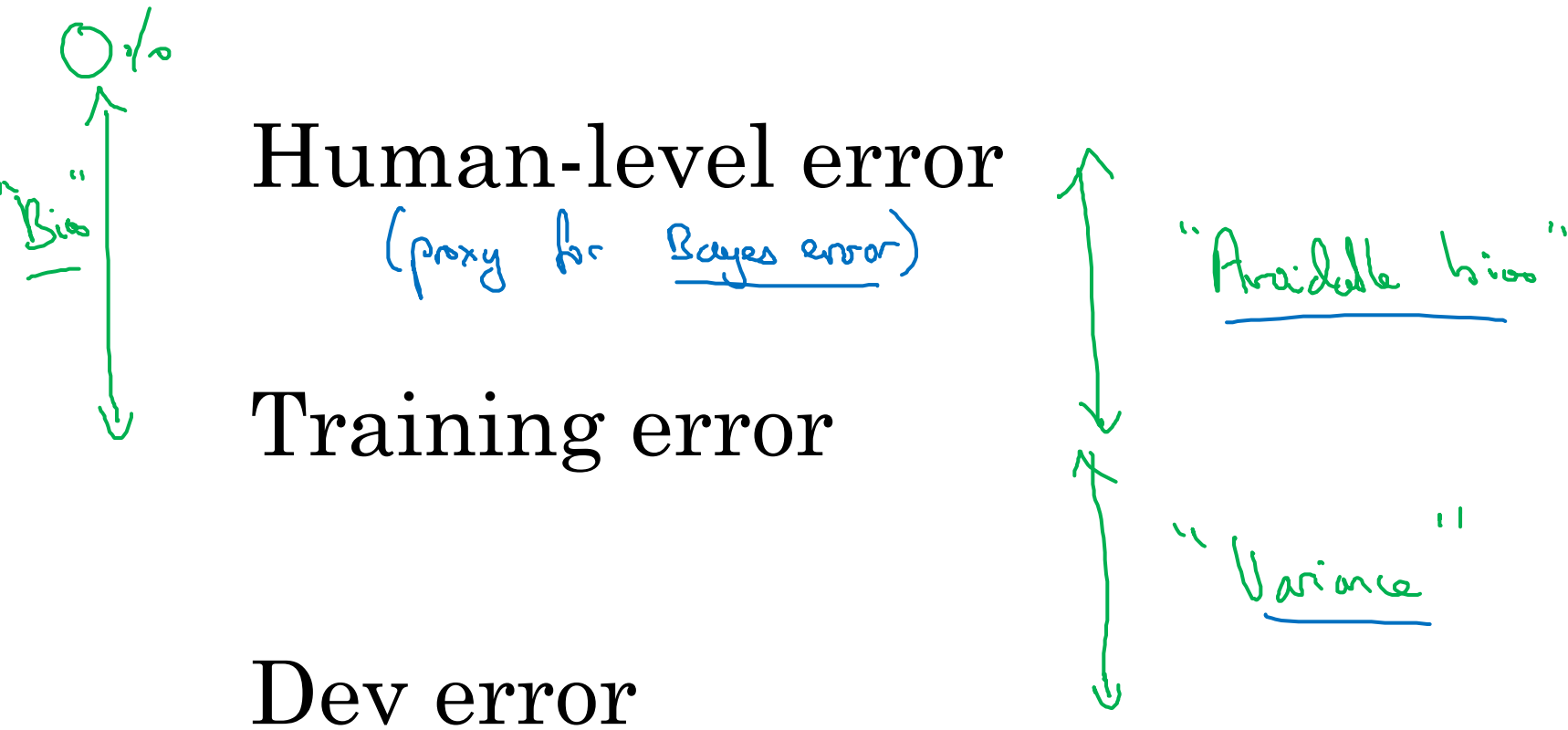| | | |
|---|---|---|
| 1% | 1% | → 0.7% 1% |
| 0.7% | 0.7% | → 0.5% |
| 0.5% | 0.5% | |
| 4% 4.5% | 0.1% 0.5% | 0.2% 0.0% |
| 5% | 1% | 0.7% |
| 1% | 4% | 0.1% |
| 6% | 5% | → 0.8% |
| Bias | Variance | |



Andrew Ng

# Summary of bias/variance with human-level performance

$0\%$

"Bias"

Human-level error

(proxy for Bayes error)

"Avoidable bias"

Training error

"Variance"

Dev error

Comparing to human-level performance

Surpassing human-level performance

deeplearning.ai

# Surpassing human-level performance
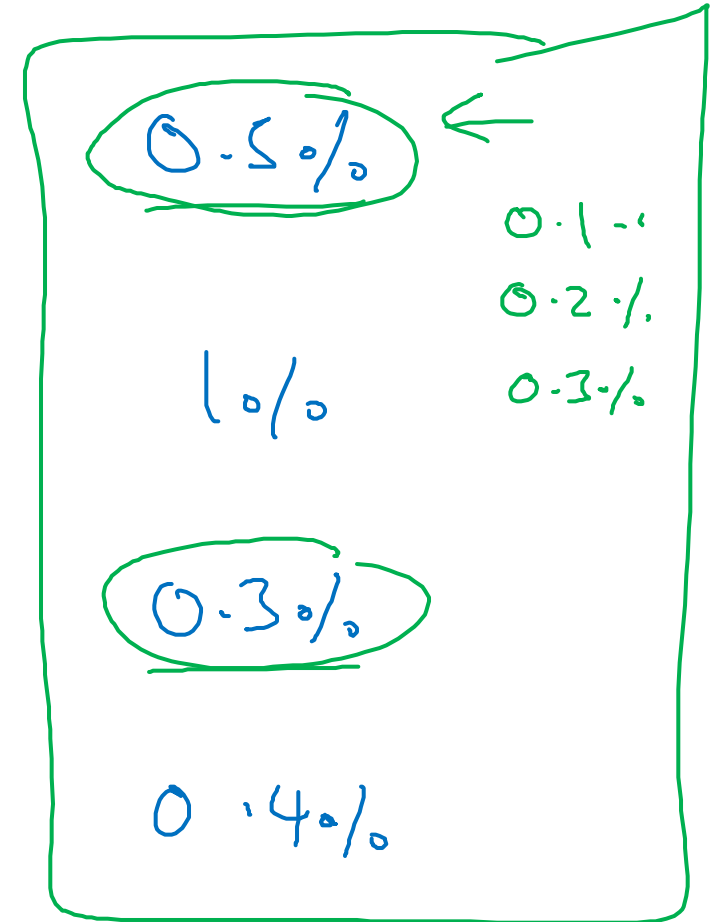
Team of humans    0.5%

One human    ~~1%~~

Training error    0.6%

Dev error    0.8%

0.1

0.2

What is <u>avoidable bias?</u>

0.5%    0.1%
        0.2%
1%      0.3%

0.3%

0.4%

# 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, ...

Comparing to human-level performance

Improving your model performance

deeplearning.ai

# The two fundamental assumptions of supervised learning

1. You can fit the training set pretty well.

   Reduce
   ~ Avoidable bias

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

   Reduce
   ~ Variance

# Reducing (avoidable) bias and variance

Human-level = Proxy for Bayes error

Avoidable bias

Training error

Variance

Dev error

Train bigger model

Train longer/better optimization algorithms
- Momentum, RMSprop, Adam
- step size, #epochs

NN architecture/hyperparameters search
- activation functions, #layers, # hidden units

RNN
CNN

More data

Regularization
- $L_2$, dropout, data augmentation

NN architecture/hyperparameters search

Andrew Ng