Objective of this course is to get the machine learnig system working

deeplearning.ai

# Introduction to ML strategy

---

# Why ML Strategy?

What is machine learning strategy. Thats start with a motivating example.

# Motivating example



say u get 90% accurency but this is not good enogh for you application
what u can do is:

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

There are a lot of ideas of what to change so u have to be carefull in what u choose to change,
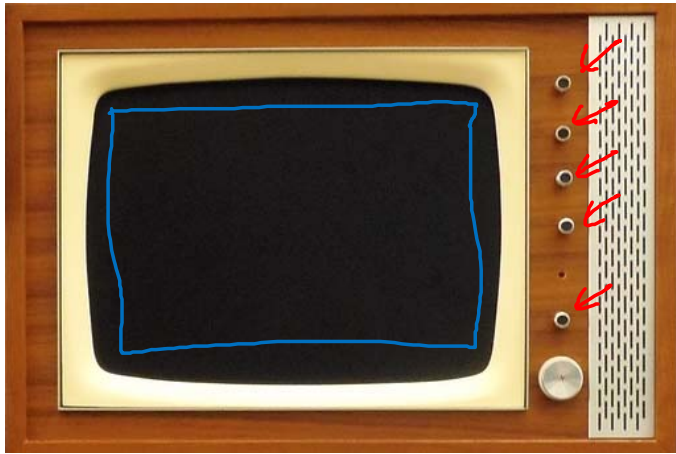it could take months to make these changes and then realize they are usles
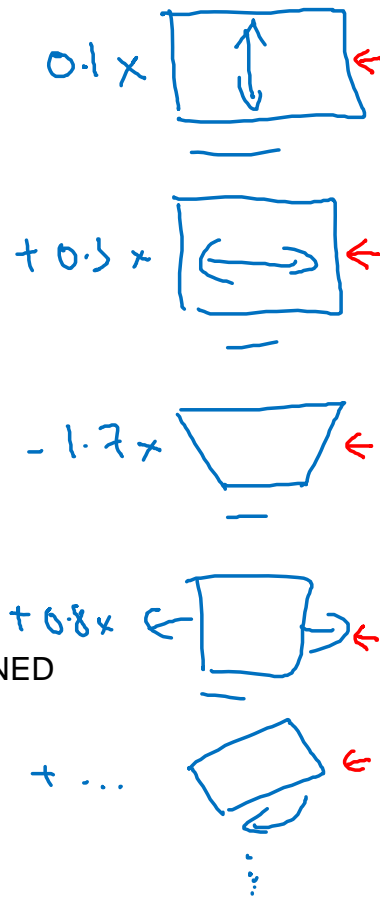
# Introduction to
# ML strategy

---

# Orthogonalization

what to tune to achive one effect, this is a process we call orthogonalization.

deeplearning.ai

# TV tuning example



Car

$0.1 \times$

$+ 0.3 \times$

$- 1.7 \times$

Orthogonalization

IT MEANS THAT EACH KNOB WAS DESIGNED
TO DO ONE THINK

$+ 0.8 \times$

$+ \ldots$

$\Rightarrow$ Steering ]

$\Rightarrow \begin{cases} \text{Accelerator} \\ \text{Braking} \end{cases}$

$\Rightarrow 0.3 \times \text{angle} \quad - \quad 0.8 \text{ speed}$

$\Rightarrow 2 \times \text{angle} \quad + \quad 0.9 \text{ speed}.$

speed

angle

Andrew Ng

# Chain of assumptions in ML

→ Fit training set well on cost function    ($\approx$ human-level performance)

width

→ Fit dev set well on cost function

height

→ Fit test set well on cost function

↓

→ Performs well in real world    (Happy cat pic app users.)

bigger network
Adam
...

early stopping

Regularization
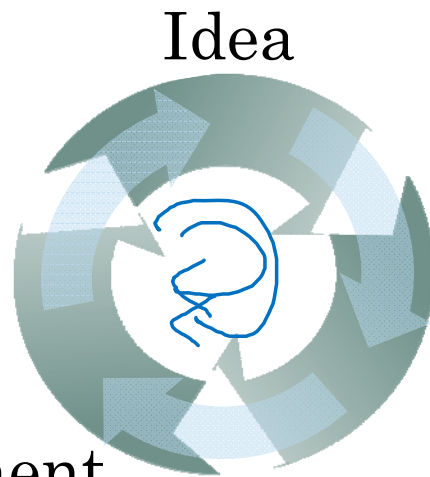Bigger trainy set

Bigger dev set

Change dev set or cost function

Andrew Ng

deeplearning.ai

Setting up
your goal

Single number
evaluation metric

# Using a single number evaluation metric

Of examples recognised as cats what % actually are cats?

Of examples recognised as cat, what % actually are cats?.

what % of actual cats are correctly recognized

what % of actual cats are correctly recognized

Idea

Experiment

Code

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

$F_1 \text{ Score} = \text{"Average" of } P \text{ and } R.$

its an avg of precision and recall

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

Dev 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% |

Andrew Ng

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 × running Time

Maximize accuracy

subject to Running Time ≤ 100 ms.

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.

Andrew Ng

deeplearning.ai

Setting up
your goal

Train/dev/test
distributions

# Cat classification dev/test sets

*development set, hold out cross validation set*

Regions:

- US
- UK
- Other Europe
- South America

→ Dev

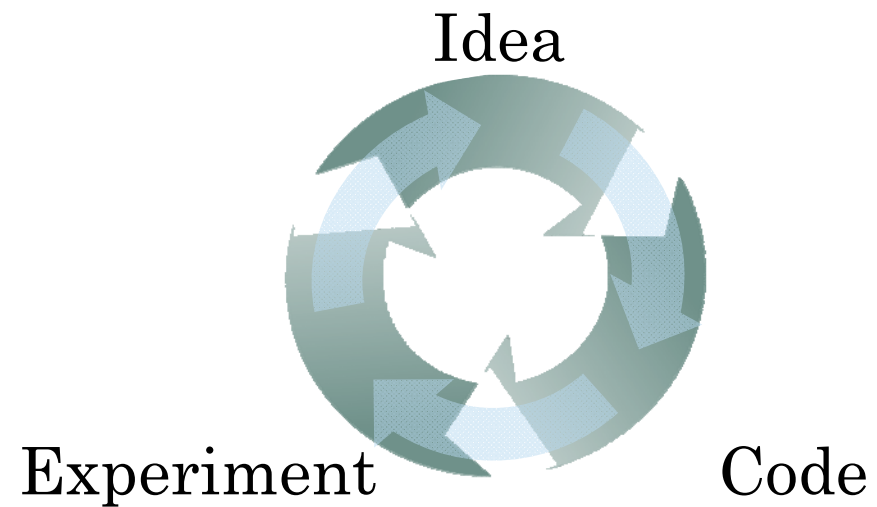- India
- China
- Other Asia
- Australia

→ Test

Randomly shuffle into dev/test



dev set
+
metric

Idea

Experiment

Code

# True story (details changed)

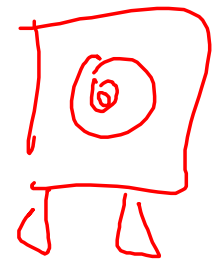Optimizing on dev set on loan approvals for medium income zip codes

$x \longrightarrow y$ (repay loan?)

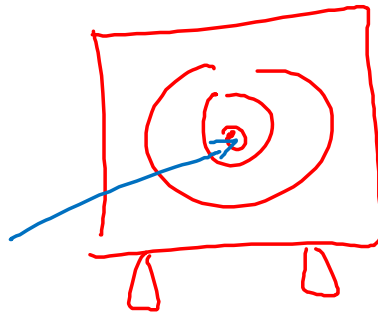Tested on low income zip codes

~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.
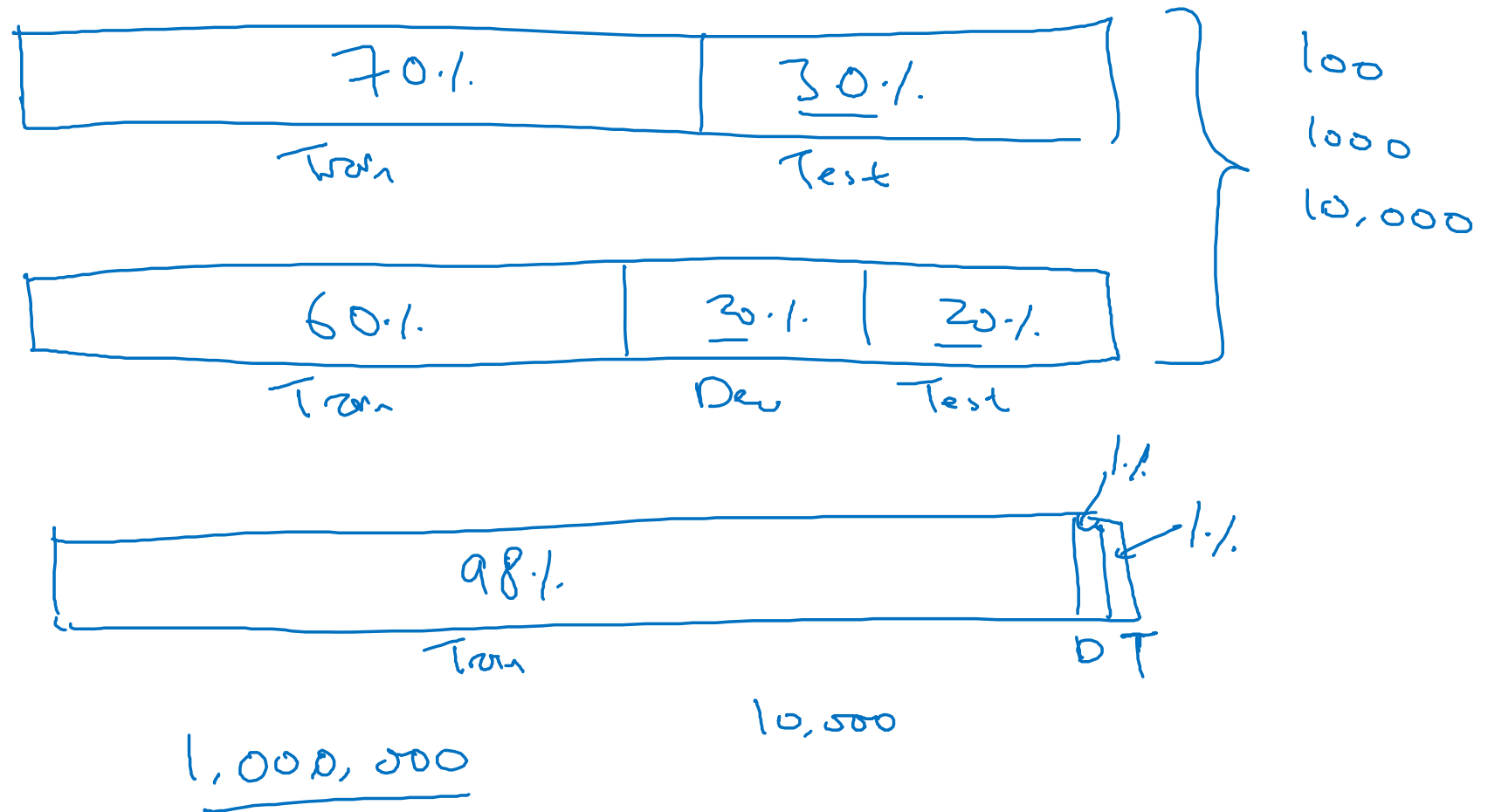
training

dev metric

test

Setting up
your goal

Size of dev
and test sets

deeplearning.ai

# Old way of splitting data



70%        30%
Train      Test

60%        30%    30%
Train      Dev    Test

100
1000
10,000

98%
Train          D  T
                1%  1%
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
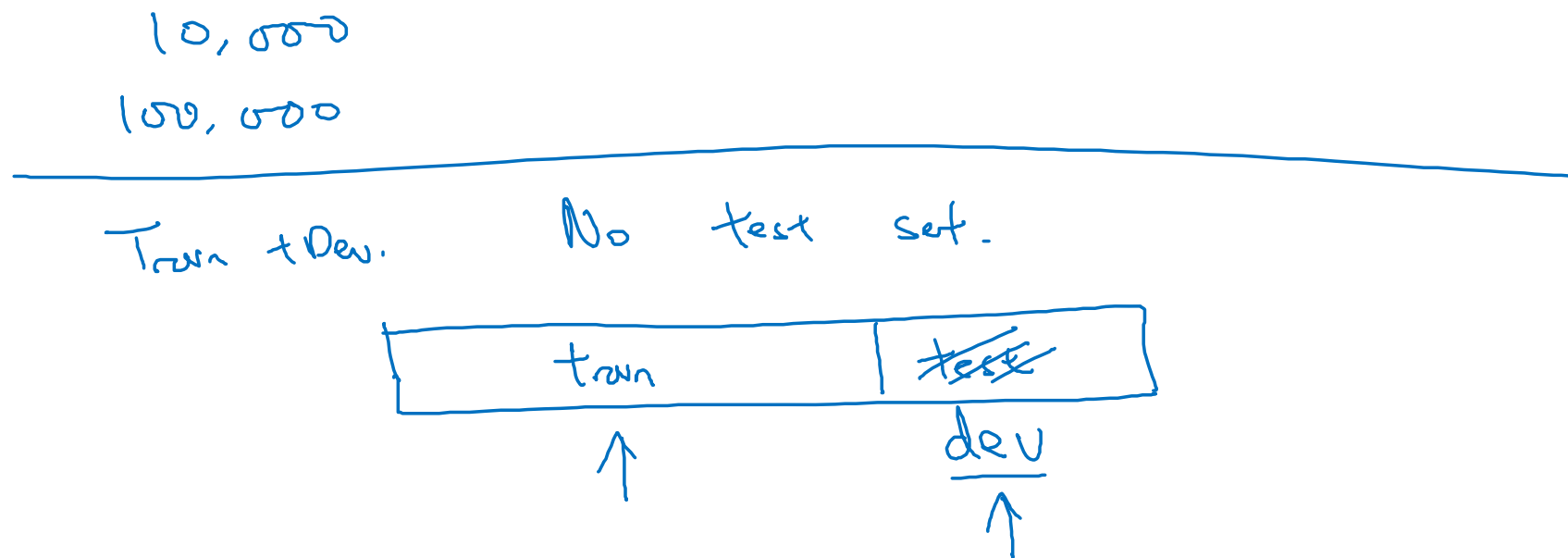  └ 1%
1,000

10, 000

100, 000

A
97% ⟶ 97.1 %

0.1%
↗

0.01%
↗

0.001%

Online advertising

# 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.

| train ↑ | test dev ↑ |

deeplearning.ai

Setting up
your goal

When to change
dev/test sets and
metrics

# Cat dataset examples

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

$\rightarrow$ Metric: classification error

Algorithm A: 3% error $\longrightarrow$ Pornographic

✓ Algorithm B: 5% error

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

predicted value (0/1)

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

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)

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

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.
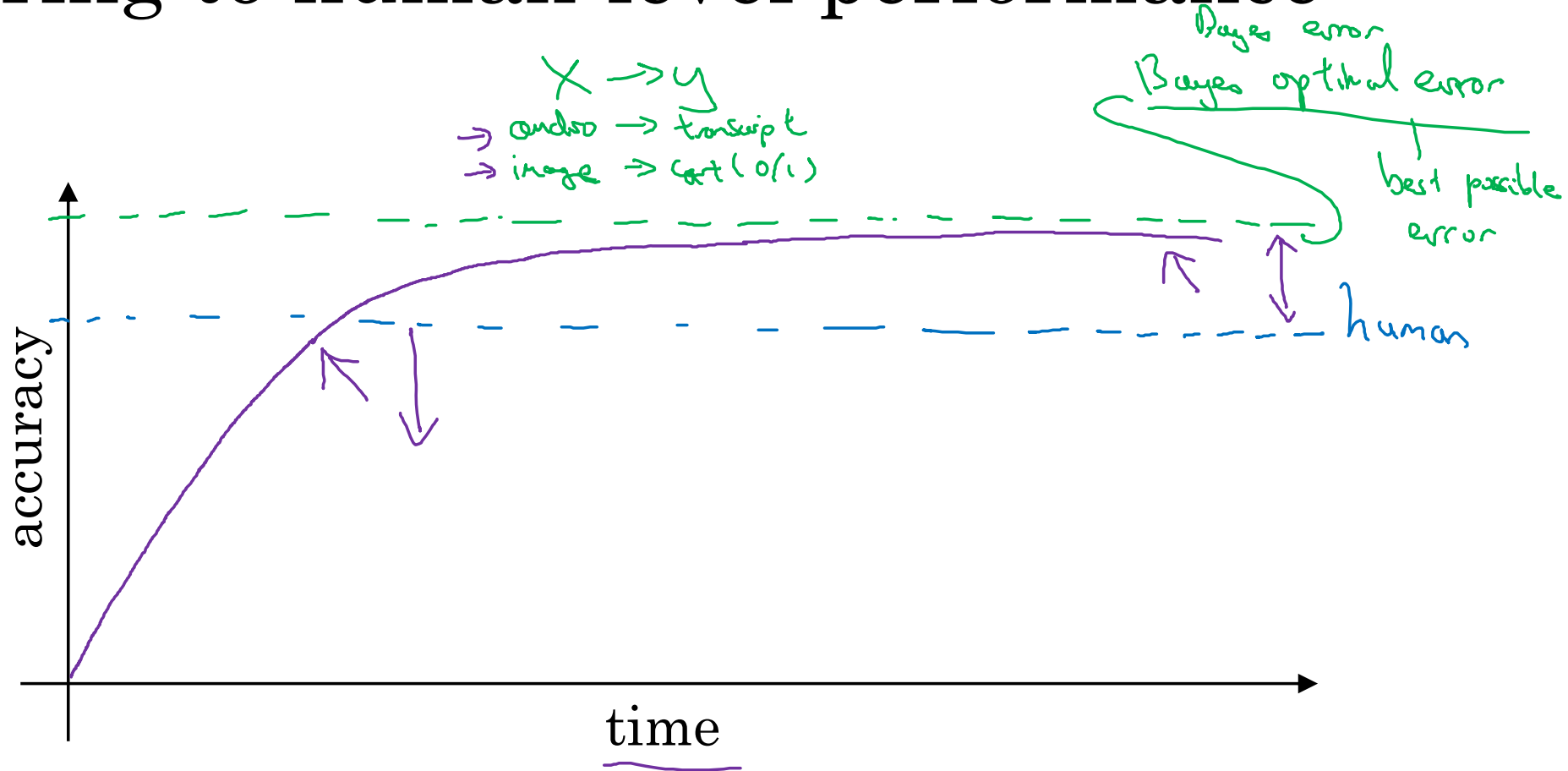
Andrew Ng

deeplearning.ai

Comparing to human-level performance

Why human-level performance?

# Comparing to human-level performance



X →> y

→ audio →> transcript
→ image ⇒ cat (0/1)

Bayes error
Bayes optimal error

best possible error

human

time

# 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:
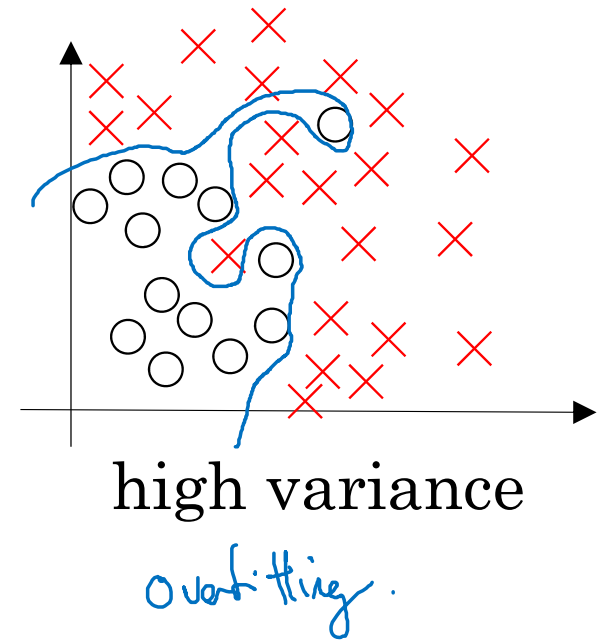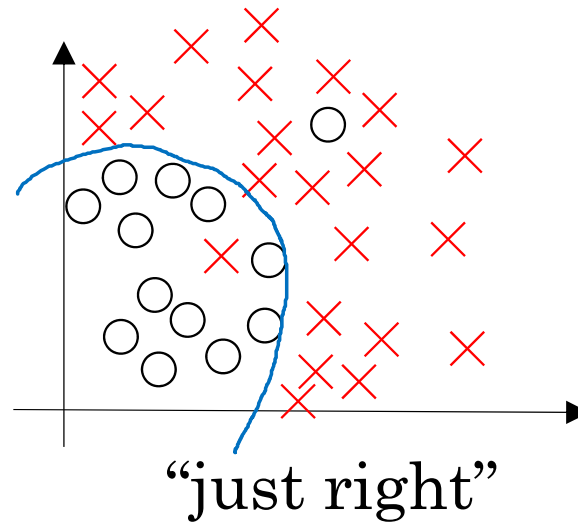  Why did a person get this right?

- Better analysis of bias/variance.

Andrew Ng

deeplearning.ai

Comparing to human-level performance

Avoidable bias

# 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 ($\approx$ Bayes)

Training error

Dev error

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

7% 2.1%

0.5% Avoidable bias

2% Variance Variance

Focus on bias

Focus on variance

Human-level error as a proxy for Bayes error.

Andrew Ng

deeplearning.ai

# Comparing to human-level performance

---

# Understanding human-level performance

# Human-level error as a proxy for Bayes error

Medical image classification example:



Suppose:

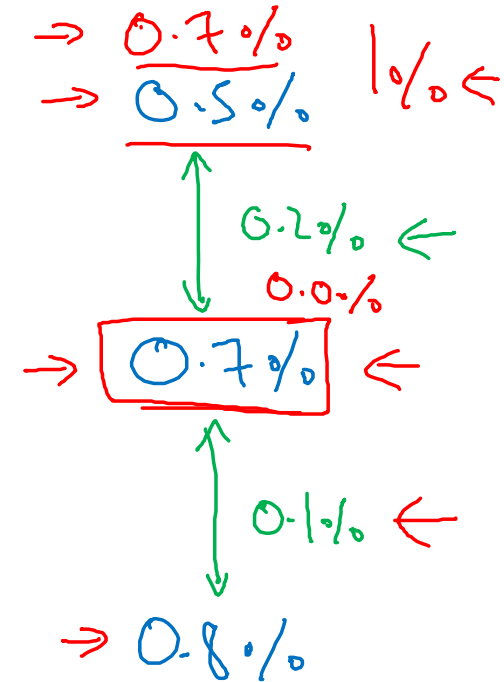    (a) Typical human ................... 3 % error

→  (b) Typical doctor ................... 1 % error

    (c) Experienced doctor ............... 0.7 % error
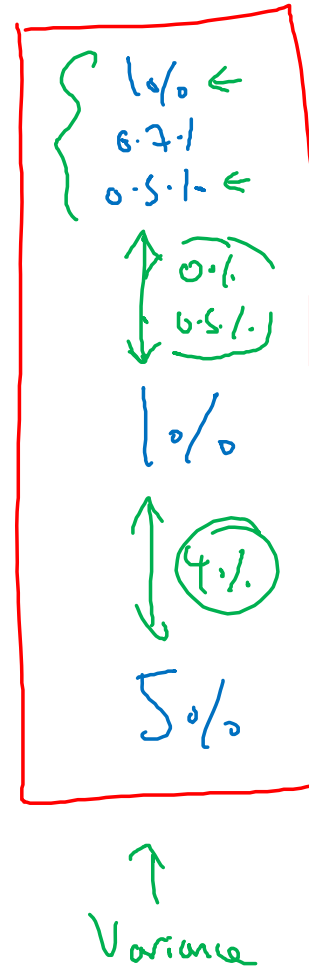
→  (d) Team of experienced doctors .. 0.5 % error ←

Bayes error ≤ 0.5%

What is "human-level" error?

# Error analysis example



Human (proxy for Bayes error)

Avoidable bias

Training error

Variance

Dev error

1%
0.7%
0.5%

4%
4.5%

5%

1%

6%

Bias

1%
0.7%
0.5%

0.1%
0.5%

1%

4%

5%

Variance

0.7%
0.5%    1%

0.2%
0.0%

0.7%

0.1%

0.8%

Andrew Ng

# Summary of bias/variance with human-level performance

$O \, \%$

"Bias"

Human-level error

(proxy for Bayes error)

"Avoidable bias"

Training error

"Variance"

Dev error

deeplearning.ai

Comparing to human-level performance

Surpassing human-level performance

# Surpassing human-level performance

Team of humans     $\underline{0.5\%}$           $\boxed{0.5\%}$ &larr;

                                     $0.1$

One human         ~~$1\%$~~           $0.1\%$

                                     $1\%$      $0.2\%$

                                                       $0.3\%$

Training error     $0.6\%$

                               $0.2$                                          $\boxed{0.3\%}$

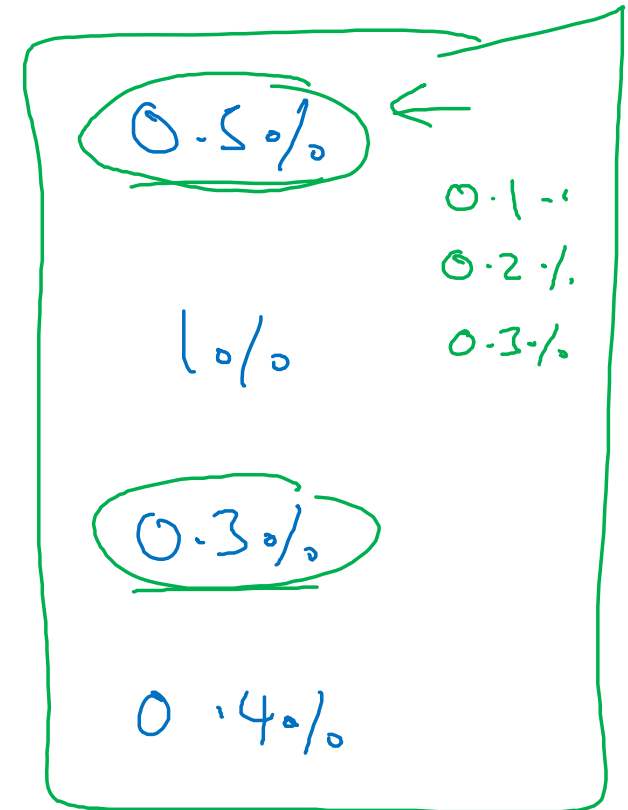Dev error        $0.8\%$                                $0.4\%$

What is $\underline{\text{avoidable bias}}$?

# Problems where ML significantly surpasses human-level performance

→ - Online advertising

→ - Product recommendations

→ - Logistics (predicting transit time)

→ - Loan approvals

Structural data

Not natural perception

Lots of data

- Speech recognition
- Some image recognition
- Medical
  - ECG, Skin cancer, ...

Andrew Ng

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.

   ~ Avoidable bias

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

   ~ Variance

Andrew Ng

# Reducing (avoidable) bias and variance

Human-level

Avoidable bias → Train bigger model

Train longer/better optimization algorithms
- Momentum, RMSprop, Adam

NN architecture/hyperparameters search    RNN
                                          CNN

Training error

Variance →

More data

Regularization
- $L_2$, dropout, data augmentation

Dev error

NN architecture/hyperparameters search

Andrew Ng