

# Introduction to ML strategy

---

Why ML  
Strategy?

deeplearning.ai



# 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



90%



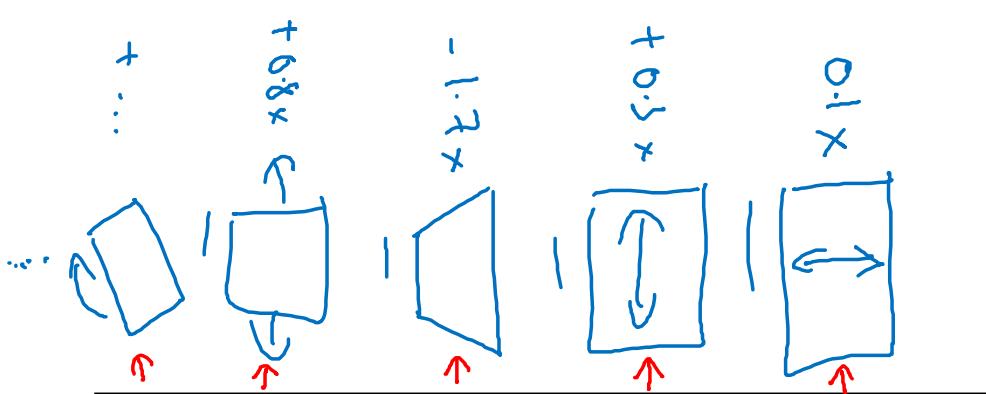
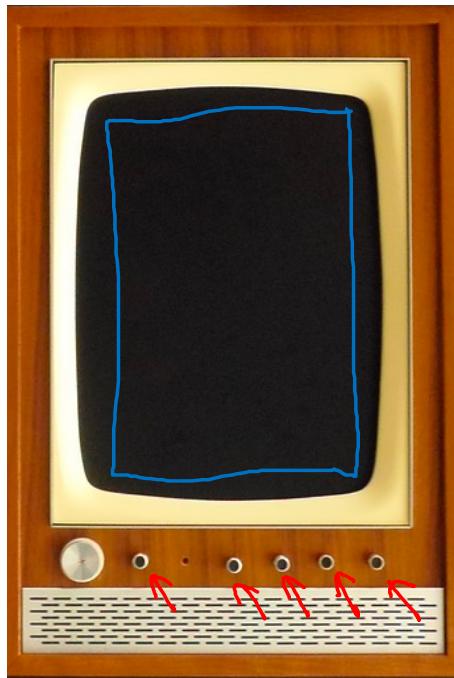
deeplearning.ai

# Orthogonalization

---

Introduction to  
ML strategy

# TV tuning example



Car

$\rightarrow 0.3 \times \underline{\text{angle}} - 0.8 \times \text{speed}$

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

Speed

angle

$\uparrow$

$\rightarrow \left\{ \begin{array}{l} \text{Accelerate} \\ \text{Braking} \end{array} \right\}$   $\rightarrow \text{Steering}$

# Chain of assumptions in ML

→ Fit training set well on cost function

(≈ human-level performance)

bigger network  
Adam  
...

early stopping

→ Fit dev set well on cost function

↪



→ Fit test set well on cost function

↪

Regularization  
Bigger dev set

Bigger dev set



→ Performs well in real world

(Happy cat pic app users.)

Change dev set or  
cost function



Setting up  
your goal

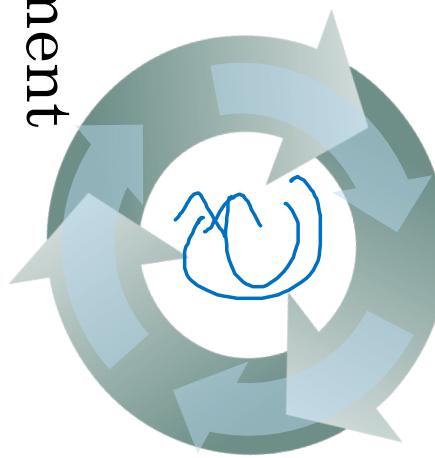
---

Single number  
evaluation metric

deeplearning.ai

# Using a single number evaluation metric

Idea



→ Of samples recognized as cat,  
what % actually are cats?

→ what % of actual cats  
are correctly recognized

Code

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

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

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

Dev set + Single number evaluation metric  
red speed up iterating

# Another example

Algorithm	US	China	India	Other
A	<u>3%</u>	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%





Setting up  
your goal

---

Satisficing and  
optimizing metrics

deeplearning.ai

# Another cat classification example

Optimizing ✓ Satisfying

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

$$\text{Cost} = \text{accuracy} - 0.5 \times \frac{\text{running time}}{\text{running time}}$$

Maximize Accuracy

Subject to running time  $\leq 100\text{ ms}$ .

N metric : 1 optimizing  
N-1 satisfying

Workwords / Trigger words

Alax, Ok Cough,  
Hey Siri, nihos baidin  
儿了百变

Maximize accuracy -  
# false positive

Maximize accuracy.

s.t.  $\leq 1$  false positive  
every 24 hours.



Setting up  
your goal

---

Train/dev/test  
distributions

deeplearning.ai

# Cat classification dev/test sets

development set, hold out cross validation set

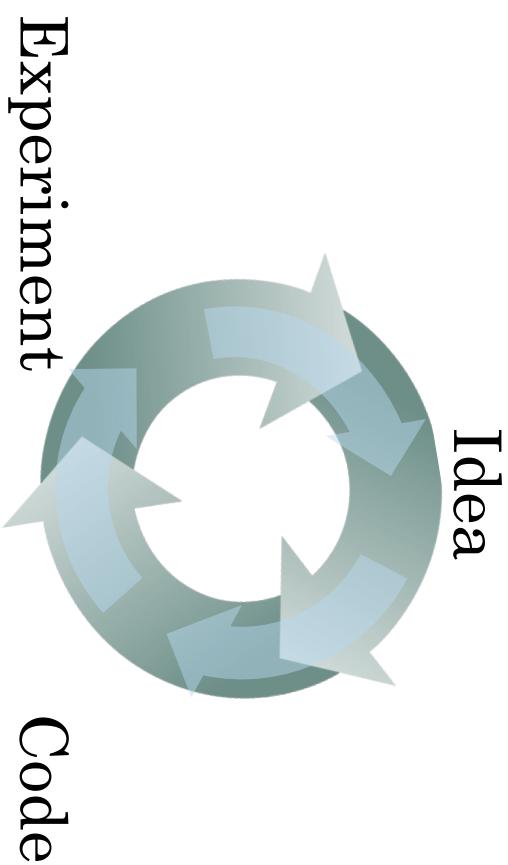
Regions:

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

Test

Dev

Randomly shuffle into dev/test



# True story (details changed)

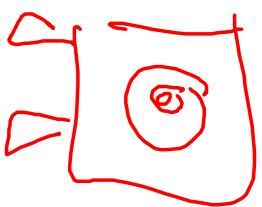
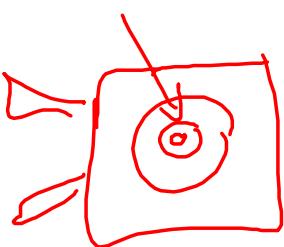
Optimizing on dev set on loan approvals for

medium income zip codes

$\uparrow$   
 $X \rightarrow y$  (repay loan?)

Tested on low income zip codes

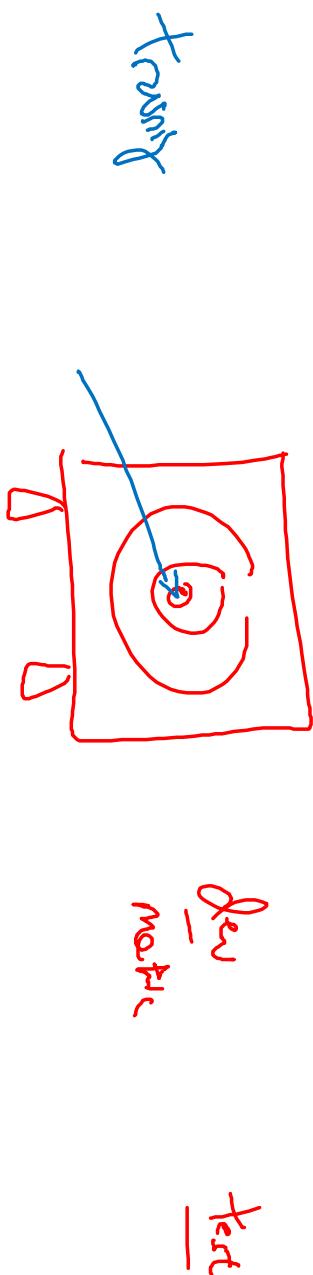
$\sim 3$  month



# 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.



Setting up  
your goal

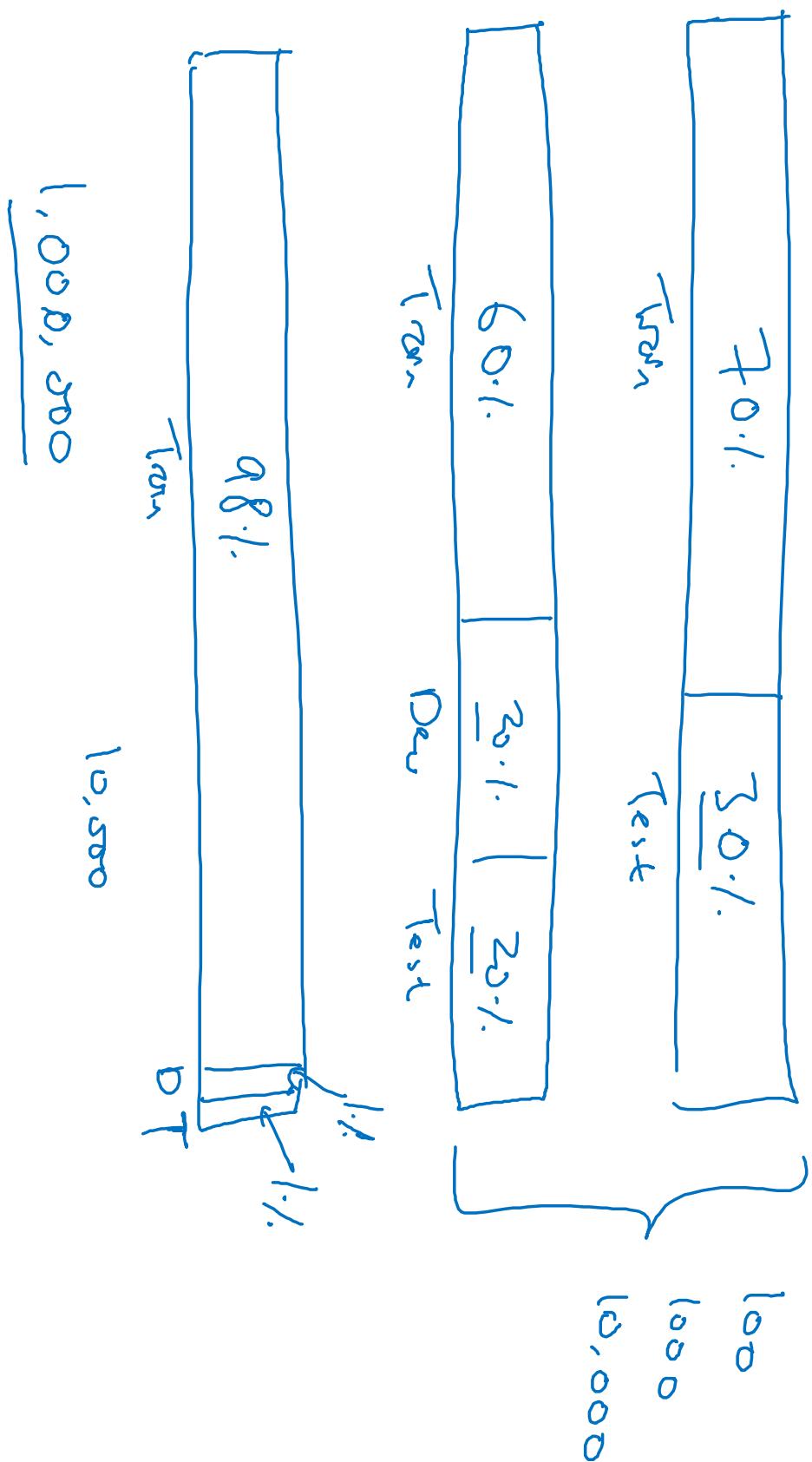
---

Size of dev  
and test sets



deeplearning.ai

# Old way of splitting data



# 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%

97%

97.1%

0.1%

1,000

10,000

100,000

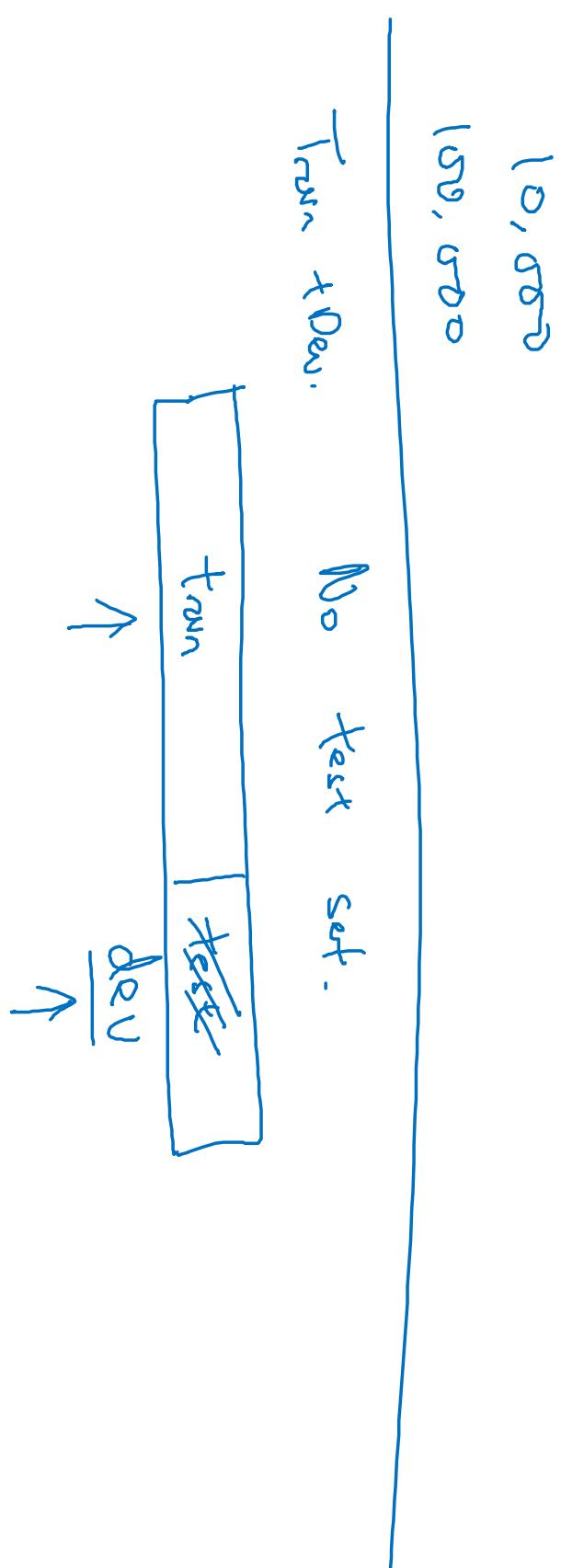
Only advertising

6.01%

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.



Setting up  
your goal



When to change  
dev/test sets and  
metrics

deeplearning.ai

# Cat dataset examples

→ Metric: classification error

Algorithm A: 3% error

Pornographic

✓ Algorithm B: 5% error

$$\text{Error} = \frac{1}{\sum_{i=1}^m w^{(i)}} \sum_{i=1}^m \left[ \begin{array}{l} \text{if } y^{(i)} \neq \hat{y}^{(i)} \\ \text{then } 1 \\ \text{else } 0 \end{array} \right]$$

Ypred = predicted value (0/1)

Ytrue = true value (0/1)

Metric + Dev : Pref A  
You (User) : Prefer B.

# 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. ↗

$$\rightarrow \mathcal{J} = \frac{1}{\sum_{i=1}^m w^{(i)} f(y^{(i)}, \hat{y}^{(i)})}$$



# Another example

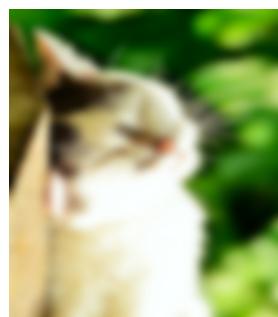
Algorithm A: 3% error

✓ 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.

---



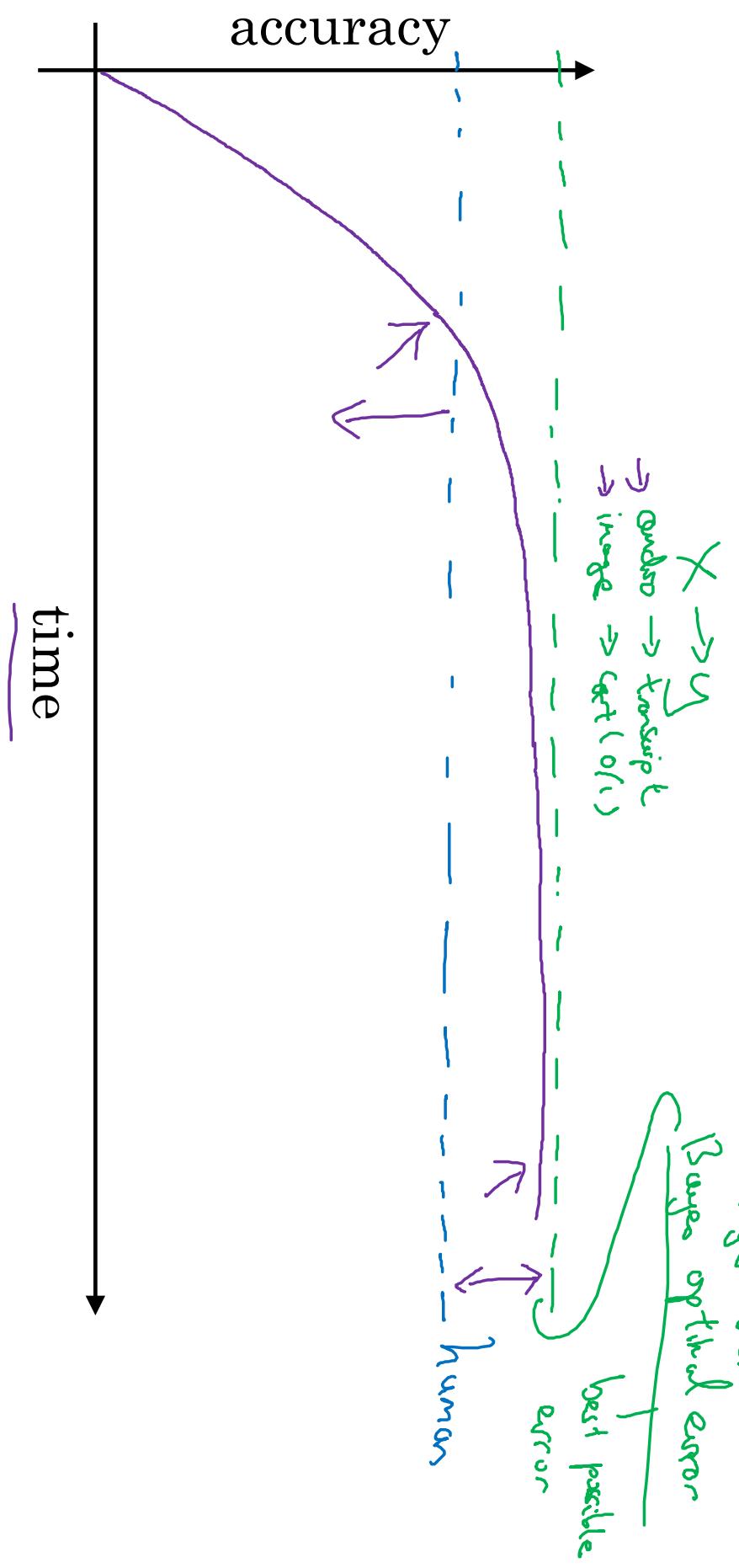
Comparing to human-level performance

---

Why human-level performance?

deeplearning.ai

# Comparing to human-level performance



# 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:

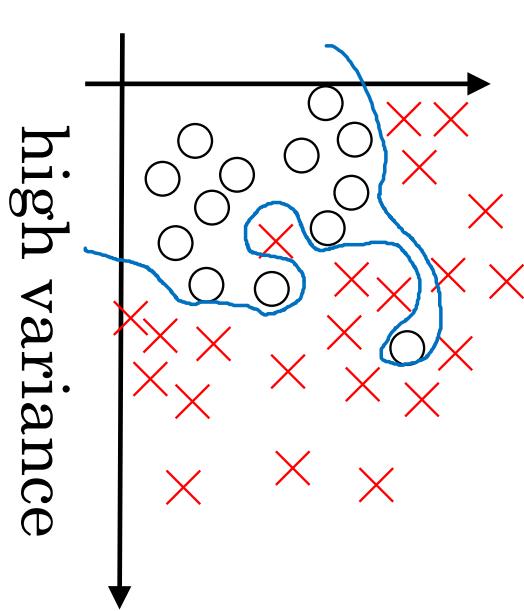
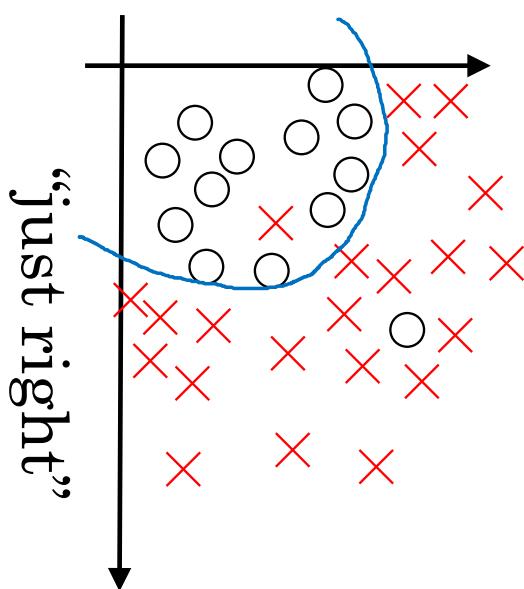
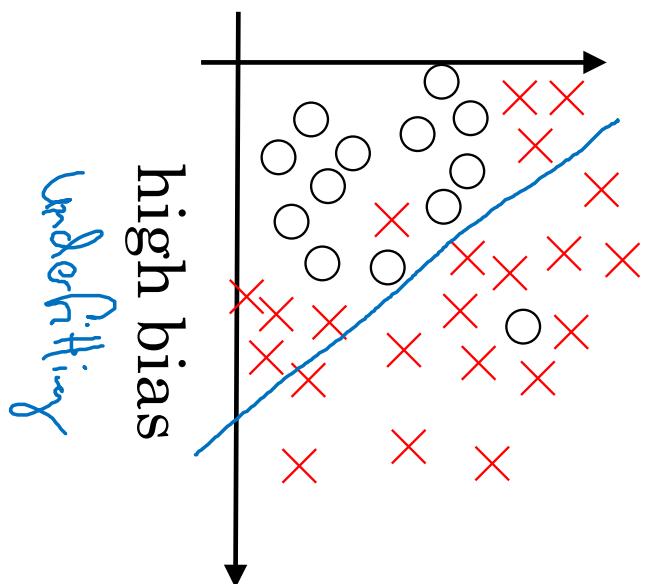
- - Get labeled data from humans.  $(x, y)$
- - Gain insight from manual error analysis:  
Why did a person get this right?
- - Better analysis of bias/variance.



Comparing to human-level performance

# Avoidable bias

# Bias and Variance

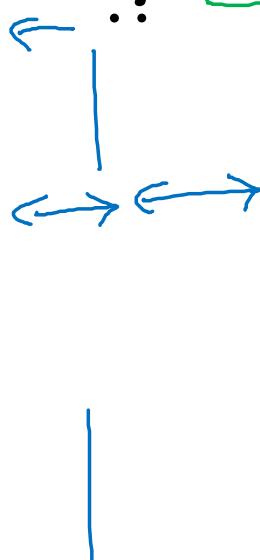


# Bias and Variance

## Cat classification

Human-level  $\approx 0\%$

Training set error:



high variance

high bias

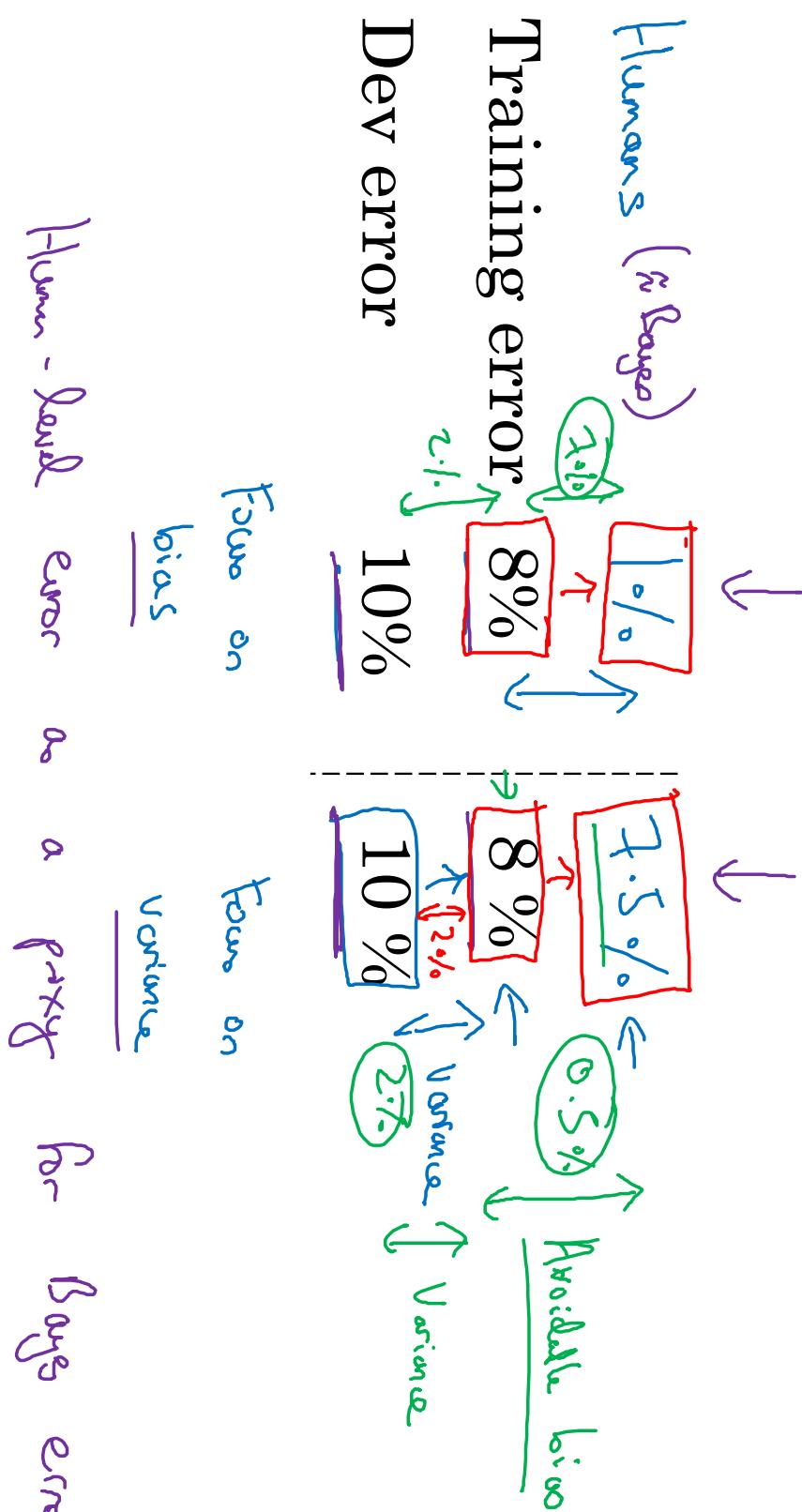
high bias  
high variance

low bias  
low variance



# Cat classification example

Human-level error as a proxy for Bayes error.





---

Comparing to human-level performance

Understanding  
human-level  
performance

deeplearning.ai

# Human-level error as a proxy for Bayes error

Medical image classification example:

Suppose:

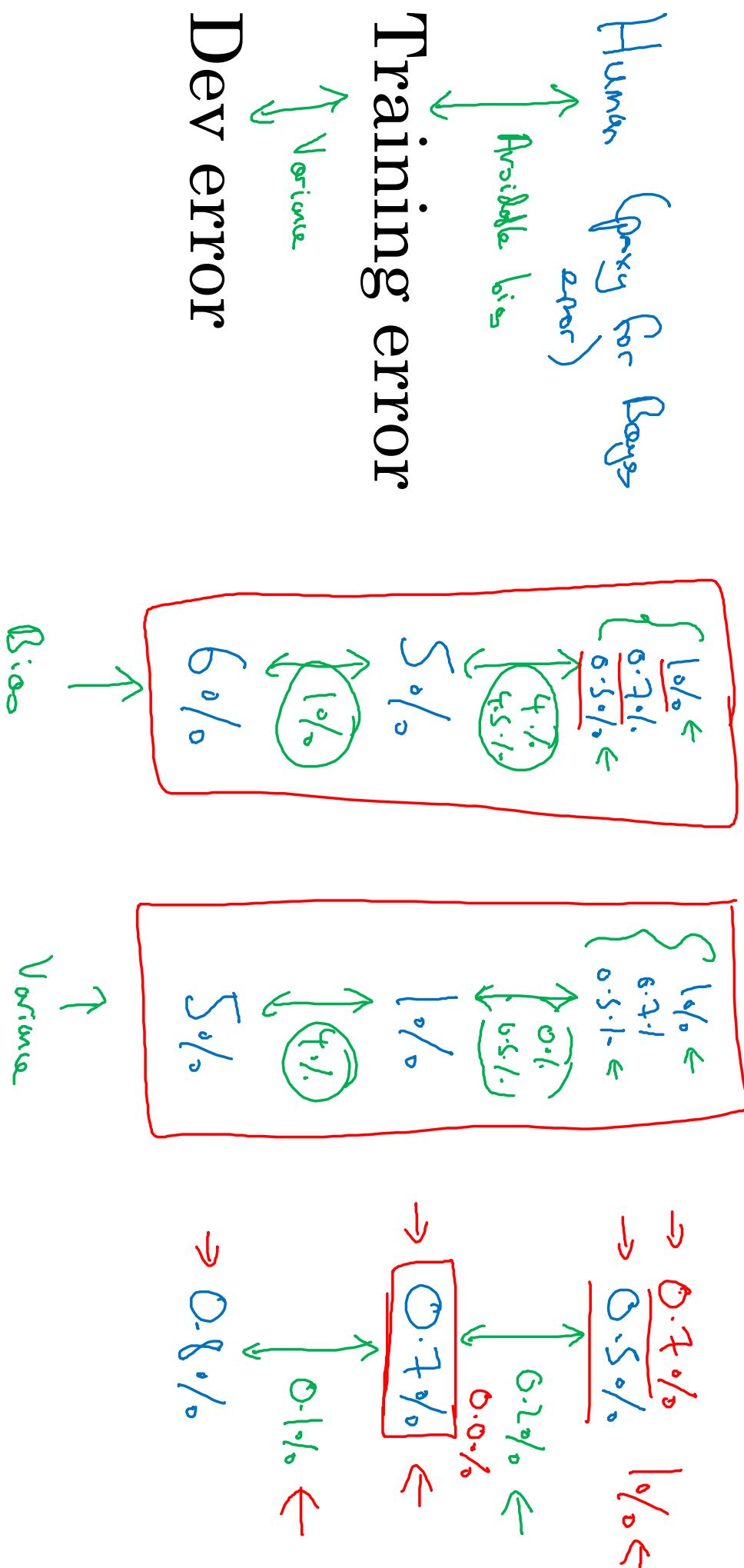
- (a) Typical human ..... 3 % error
- (b) Typical doctor ..... 1 % error
- (c) Experienced doctor ..... 0.7 % error
- (d) Team of experienced doctors .. 0.5 % error

$$\text{Base error} \leq 0.5\%$$

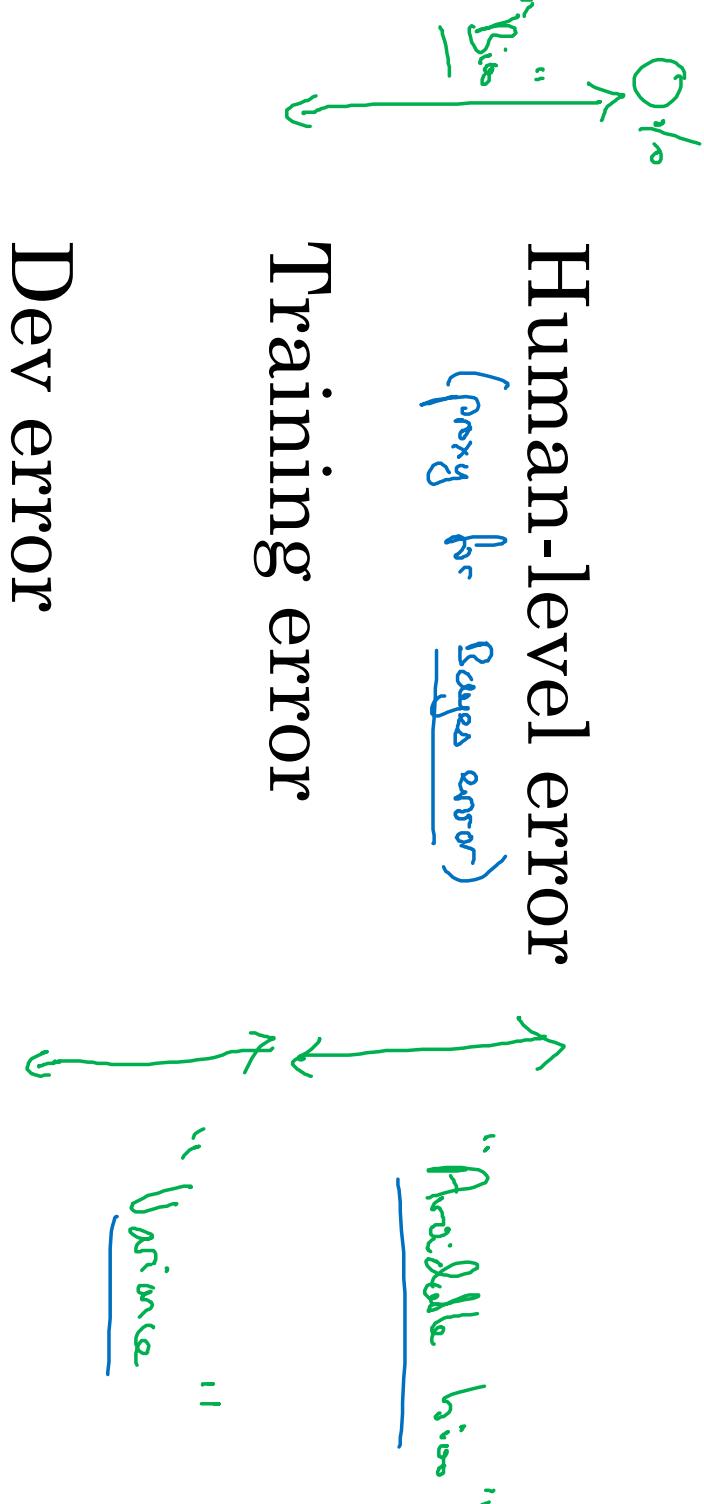
What is “human-level” error?



# Error analysis example



# Summary of bias/variance with human-level performance





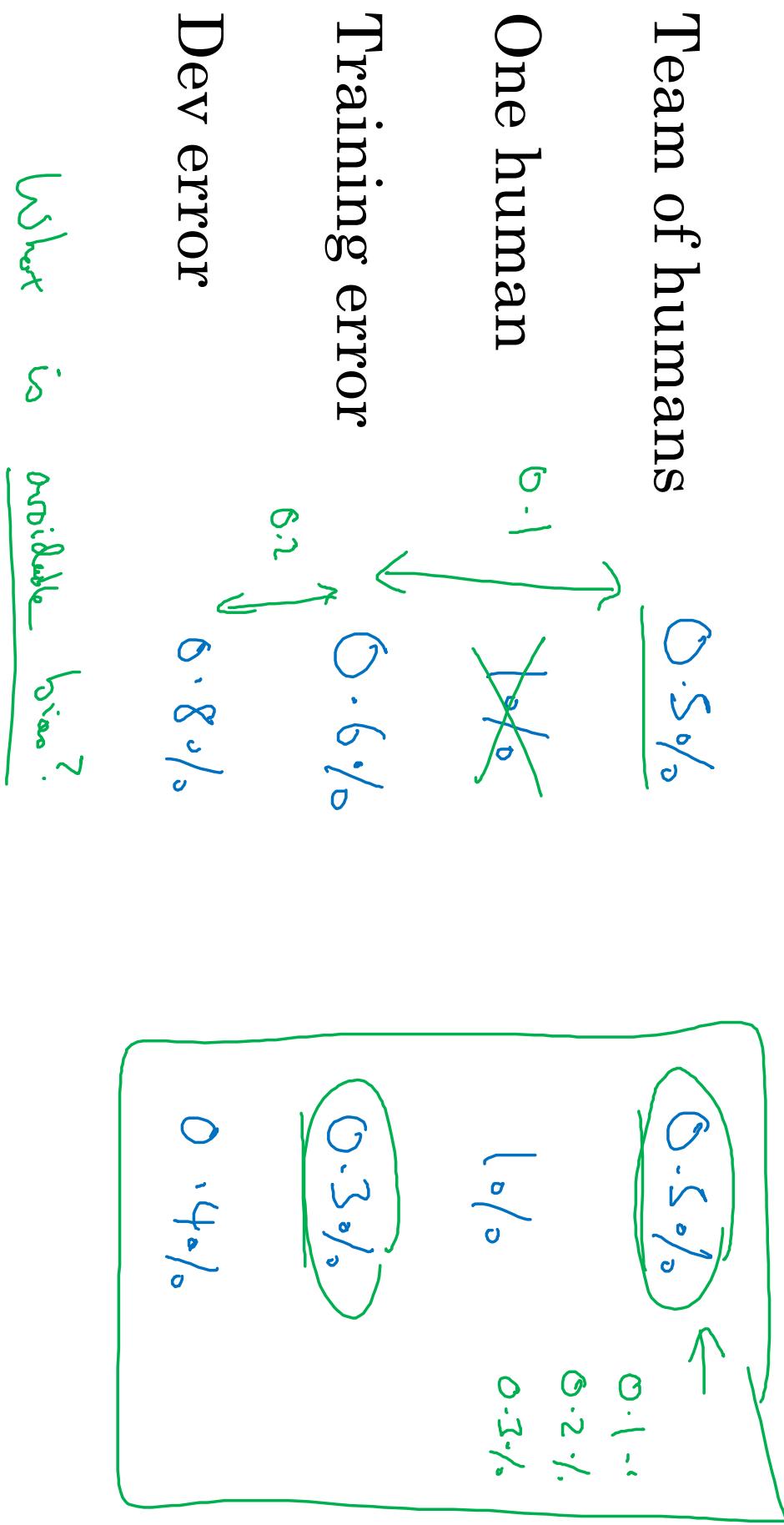
Comparing to human-level performance

---

Surpassing human-level performance

deeplearning.ai

# Surpassing human-level performance



# Problems where ML significantly surpasses human-level performance

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

- Speech recognition  
- Some image recognition  
- Medical

- E.g., skin cancer, ...

Structured Data

Not natural perception  
Lots of data



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.



~ Arbitrary bias

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



~ Variance

# Reducing (avoidable) bias and variance

