Introduction to ML strategy

Why ML Strategy?

deeplearning.ai

# 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
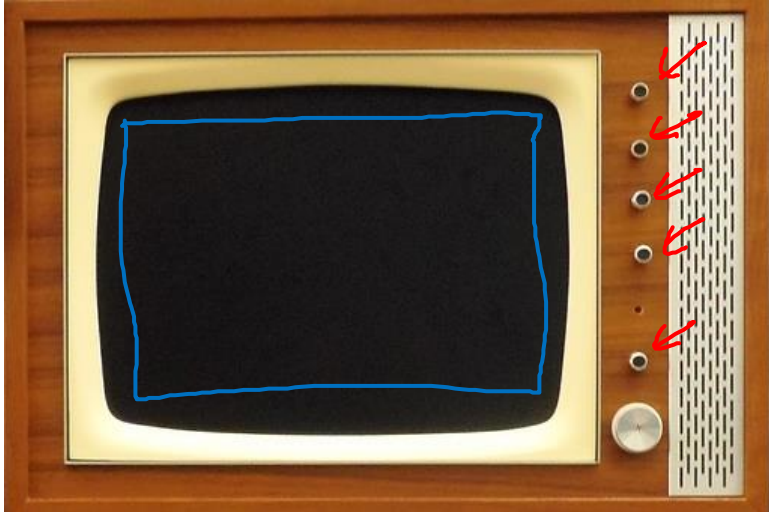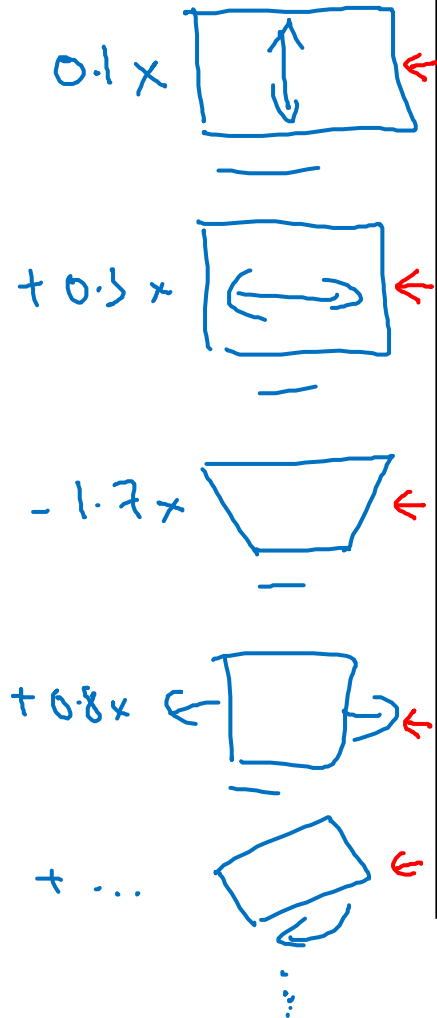
  - Activation functions

  - # hidden units

  - ...

Introduction to ML strategy

Orthogonalization

deeplearning.ai

# TV tuning example



Car

$0.1 \times$ → Steering ]

$+ 0.3 \times$ → { Accelerator
                  { Braking ]

$- 1.7 \times$

Orthogonalization

$+ 0.8 \times$

→ $0.3 \times$ angle   −   $0.8$ speed

→ $2 \times$ angle   +   $0.9$ speed.

$+ \ldots$

speed

angle

# Chain of assumptions in ML

→ Fit training set well on cost function (≈ 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
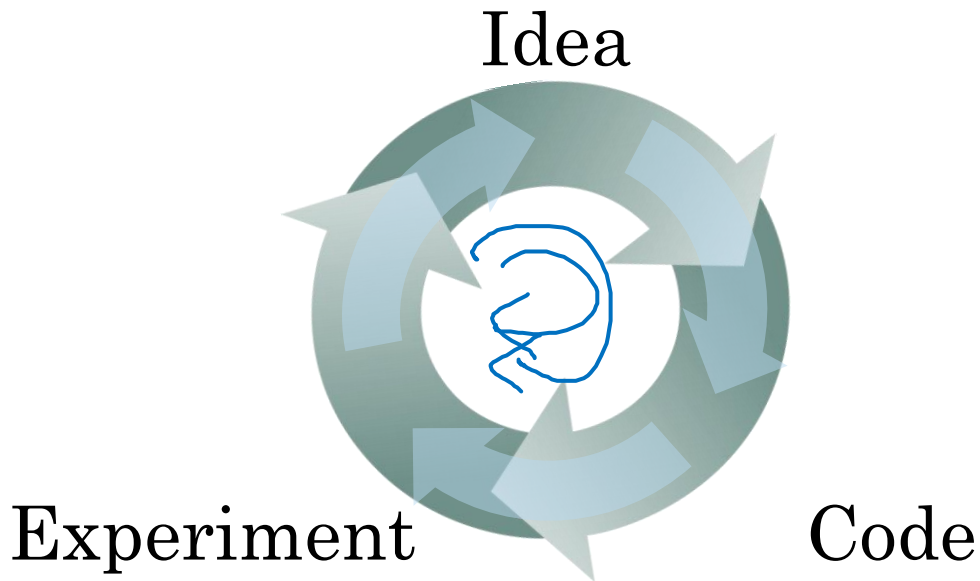Bigger traing set

Bigger dev set

Change dev set or cost function

deeplearning.ai

Setting up
your goal

Single number
evaluation metric

# Using a single number evaluation metric

Idea

Experiment

Code

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

what % of actual cats are correctly recognized

| 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 |
|---|---|---|---|---|
| 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% |

deeplearning.ai

Setting up
your goal

Satisficing and
optimizing metrics

# 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.

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.

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
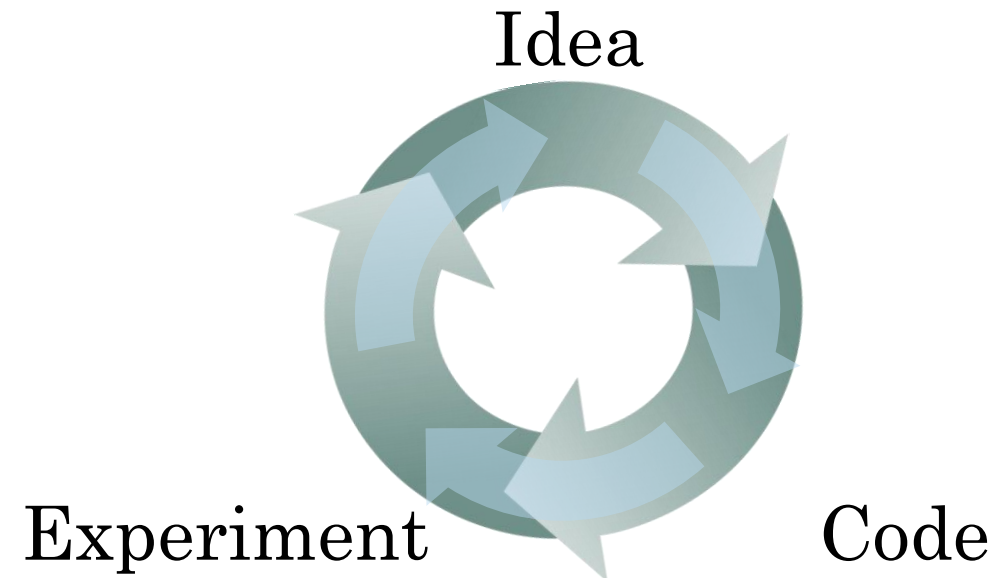- India
- China
- Other Asia
- Australia

Dev ←

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

$$x \longrightarrow y \ (\text{repay loan?})$$
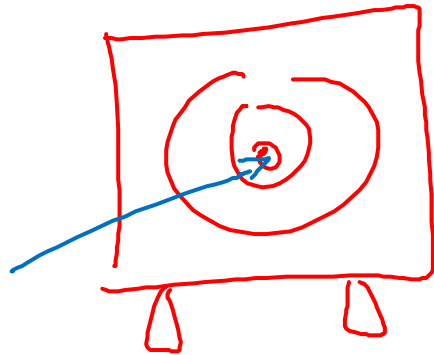
Tested on low income zip codes

~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.
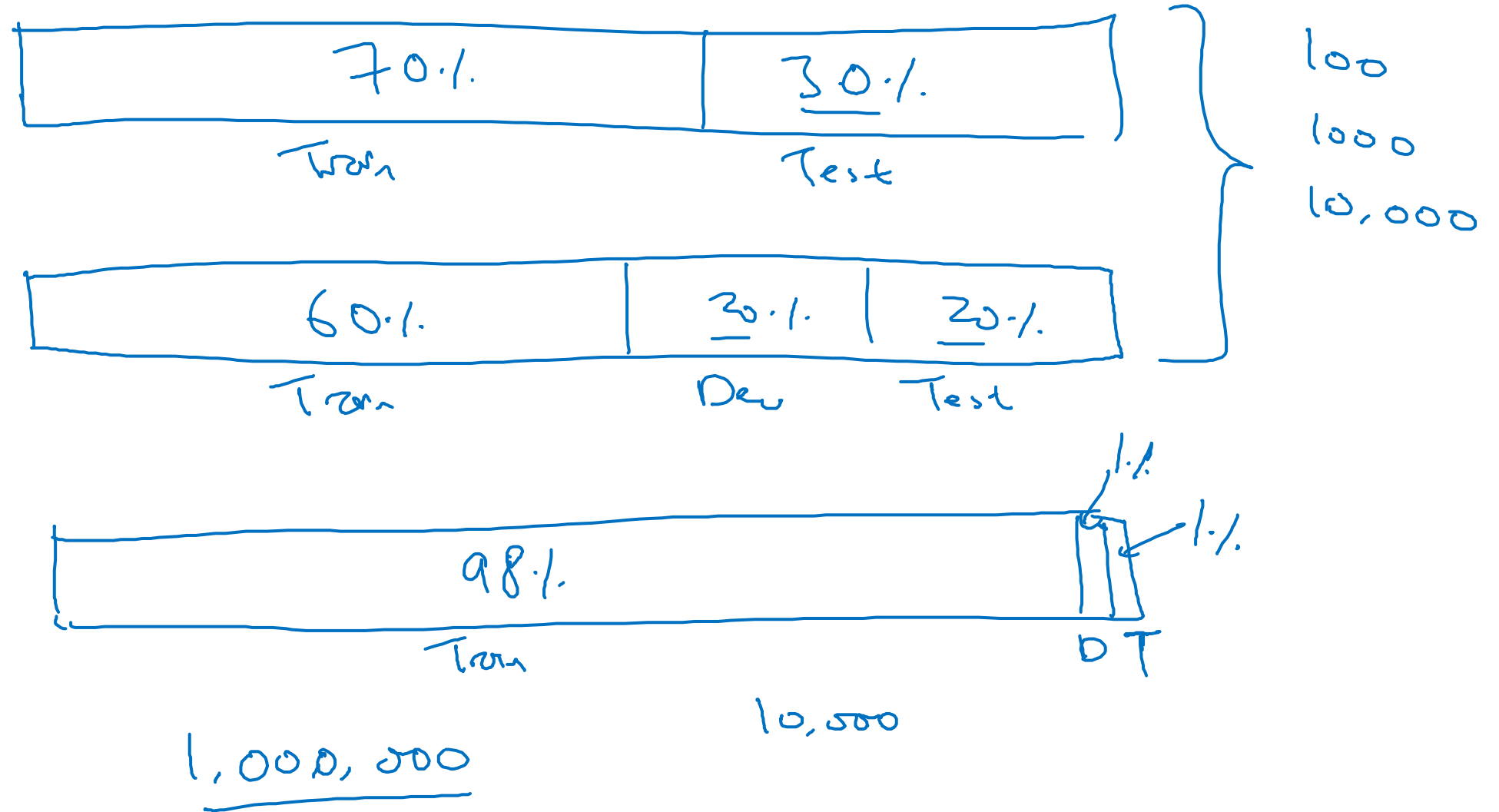
training

dev
matric

test

deeplearning.ai

Setting up
your goal

Size of dev
and test sets

# 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.
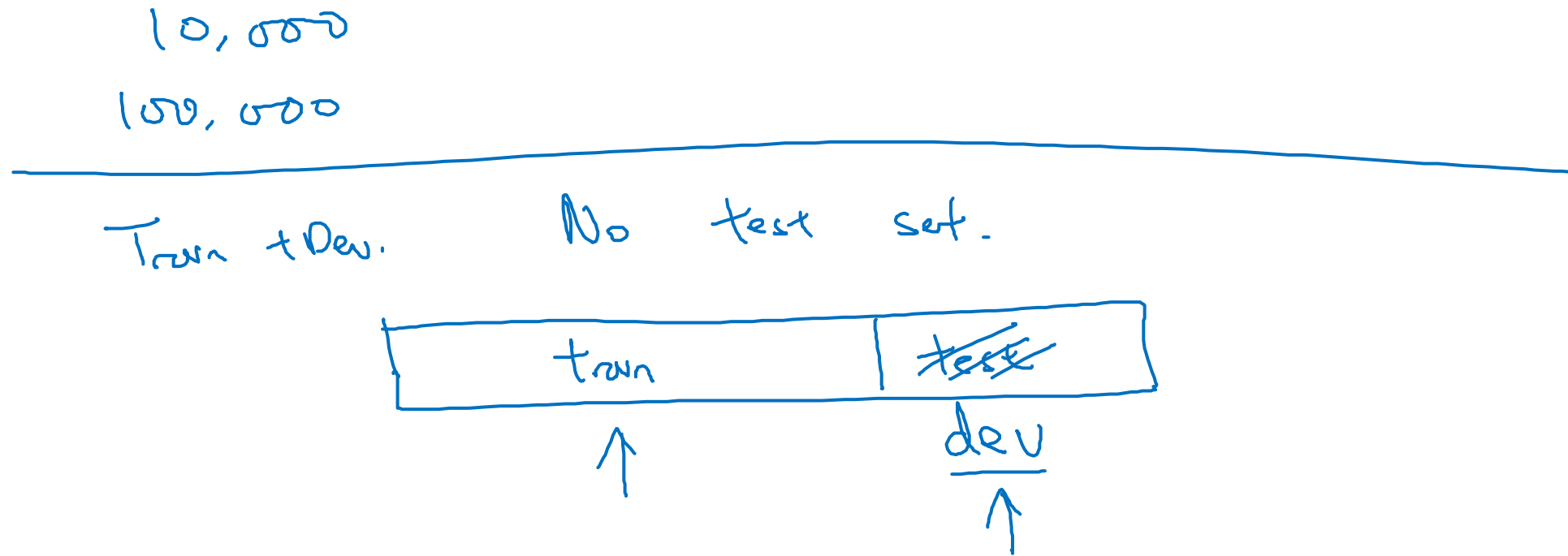
100 : Small

    $\leftarrow$ 1%

1,000

10,000

100,000

$$\overset{A}{97\%} \longrightarrow \overset{B}{97.1\%}$$

$$\frac{0.1\%}{\nwarrow}$$

$$\frac{0.01\%}{0.001\%}$$

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 | test̶ (crossed out)

↑ (under train)

dev
↑ (under dev)

Setting up
your goal

When to change
dev/test sets and
metrics

deeplearning.ai

# Cat dataset examples

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

→ Metric: classification error

Algorithm A: 3% error $\longrightarrow$ pornographic

✓ Algorithm B: 5% error

$$\text{Error}: \frac{1}{\sum w^{(i)}} \;\; \cancel{\frac{1}{M_{dev}}} \;\; \sum_{i=1}^{M_{dev}} w^{(i)} \, \mathbb{1}\left\{ \underbrace{y_{pred}^{(i)}}_{\text{predicted value } (0/1)} \neq y^{(i)} \right\}$$

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

# 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.

Aim (shoot at target)

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

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

# 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.
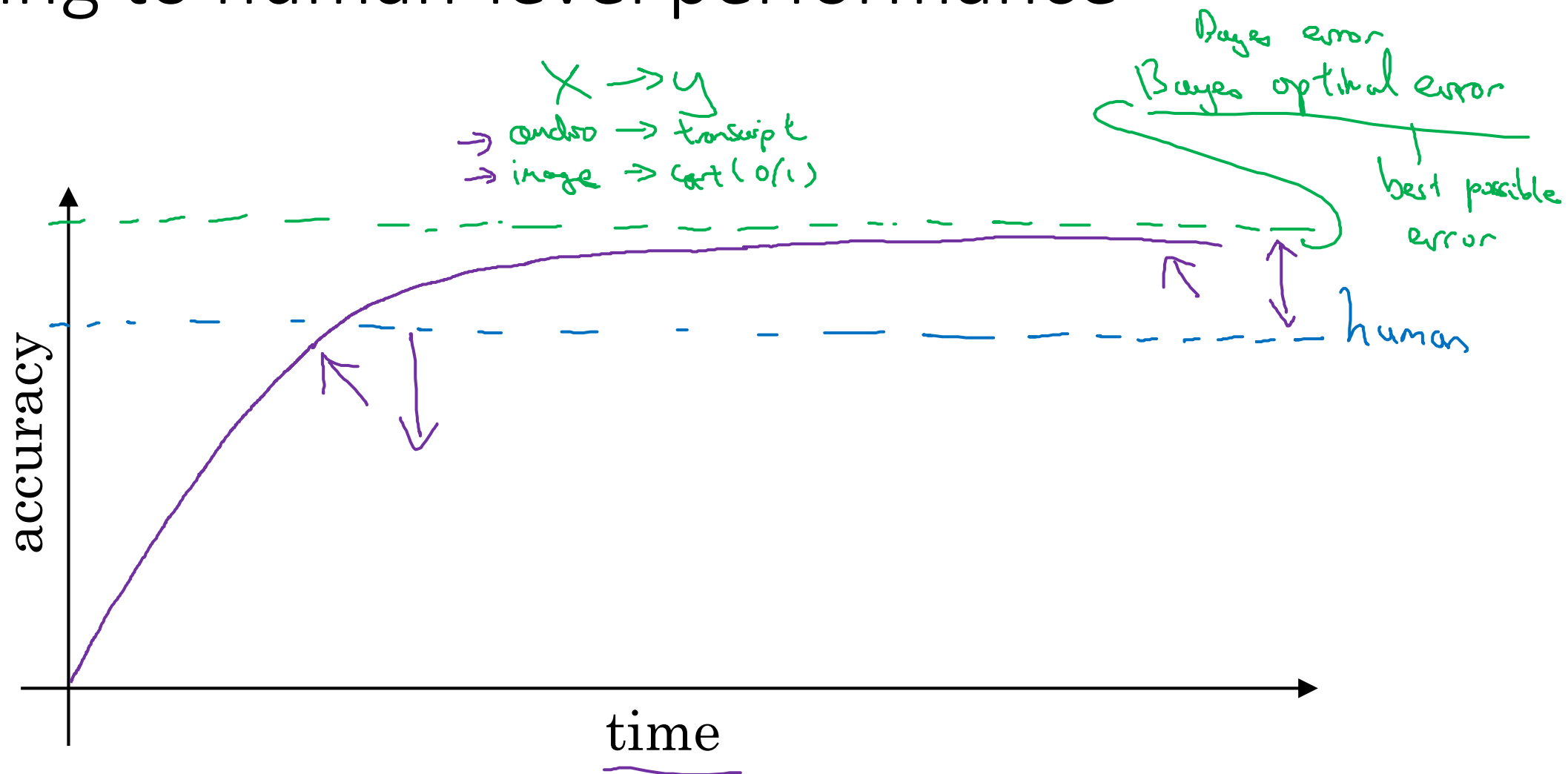
deeplearning.ai

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

# 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?
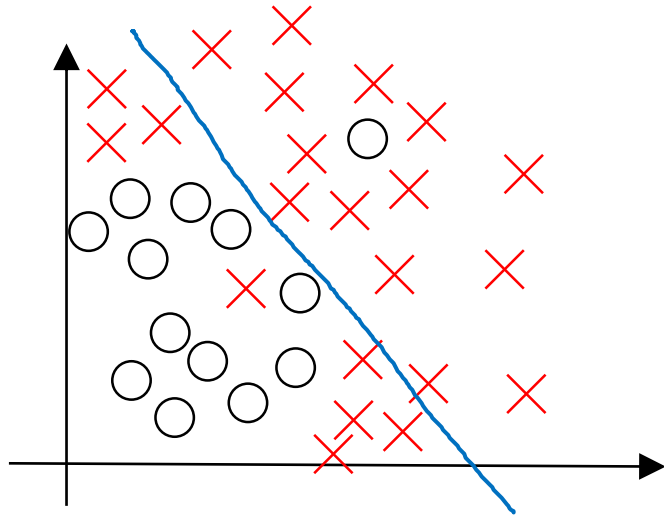
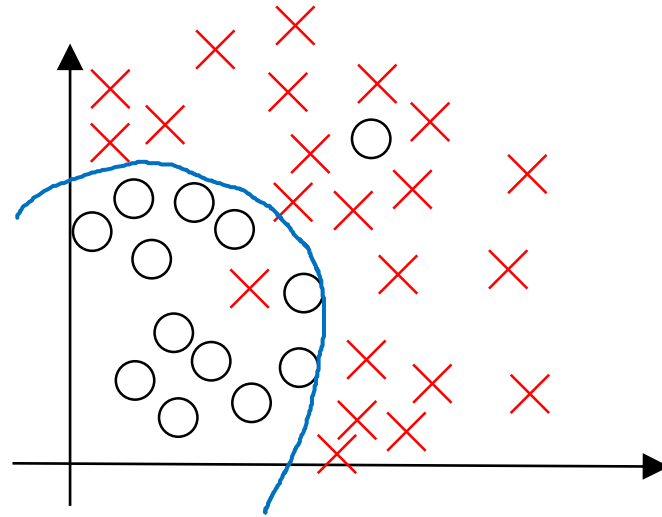→ - Better analysis of bias/variance.

Comparing to human-level performance

Avoidable bias

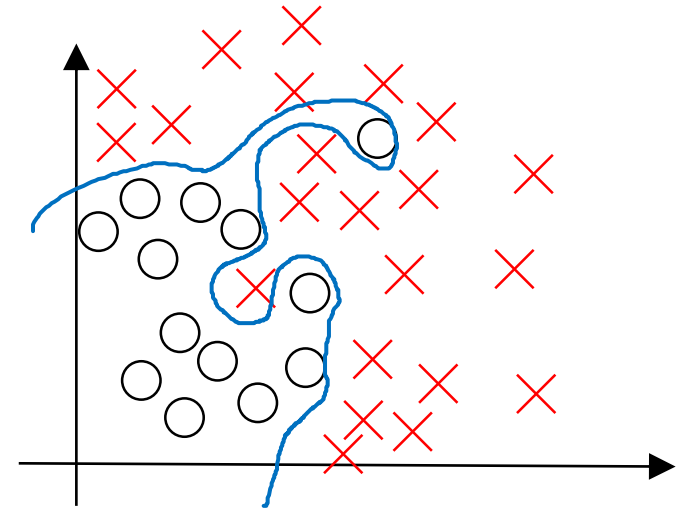deeplearning.ai

# Bias and Variance



high bias

*underfitting*

"just right"

high variance
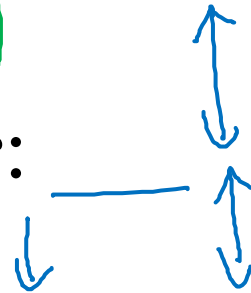
*overfitting.*

# Bias and Variance

## Cat classification



Human-level ≈ 0% ....

Training set error: ___
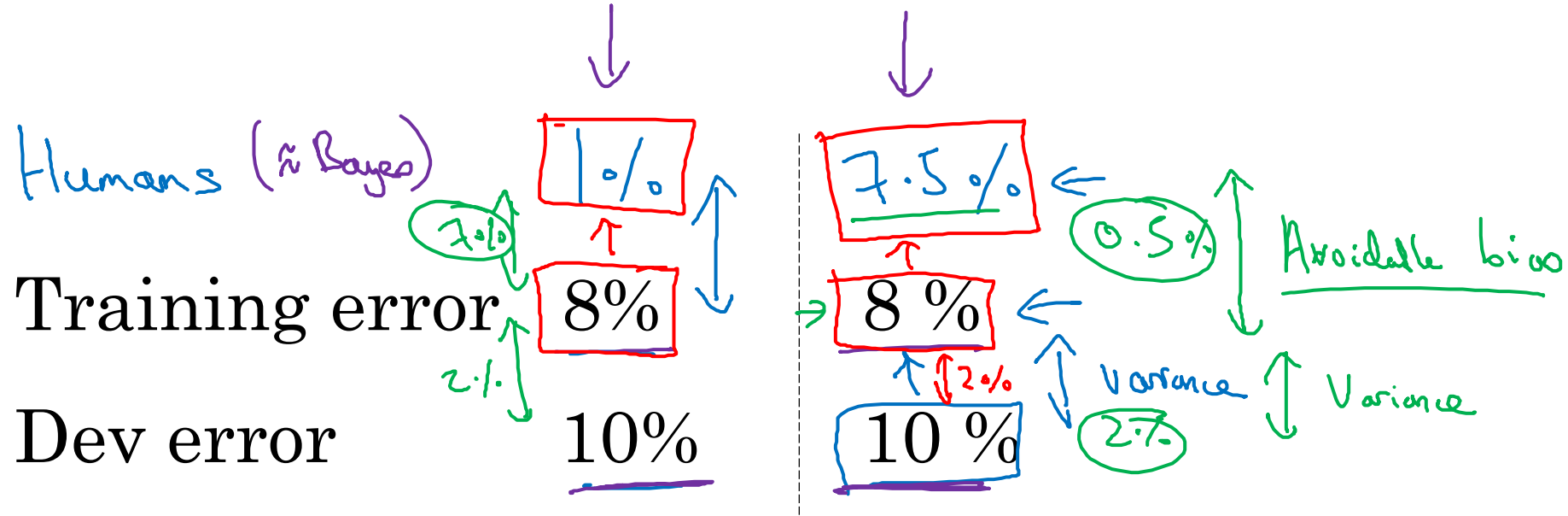
Dev set error:

high variance     high bias     high bias     low bias

high variance     low variance

# Cat classification example

Humans (≈ Bayes)

Training error

Dev error

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

7 %

2 %

0.5 %    Avoidable bias

2 %

Variance    Variance

2 %

Focus on bias

Focus on variance

Human-level error as a proxy for Bayes error.

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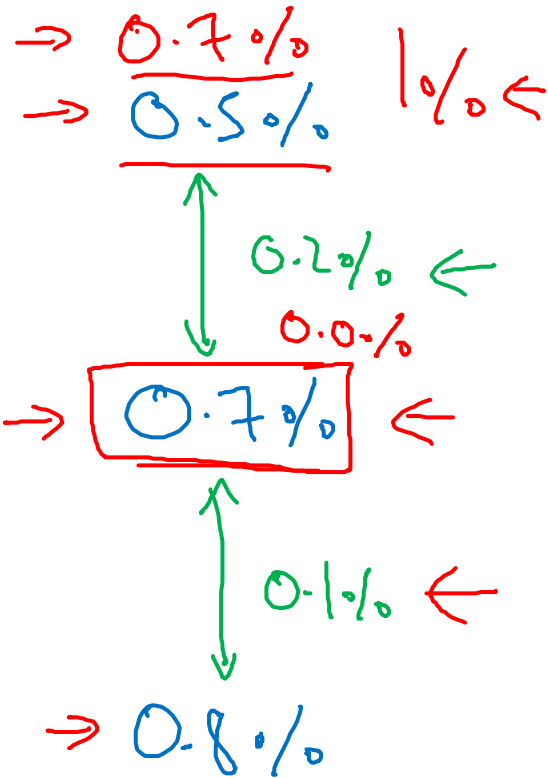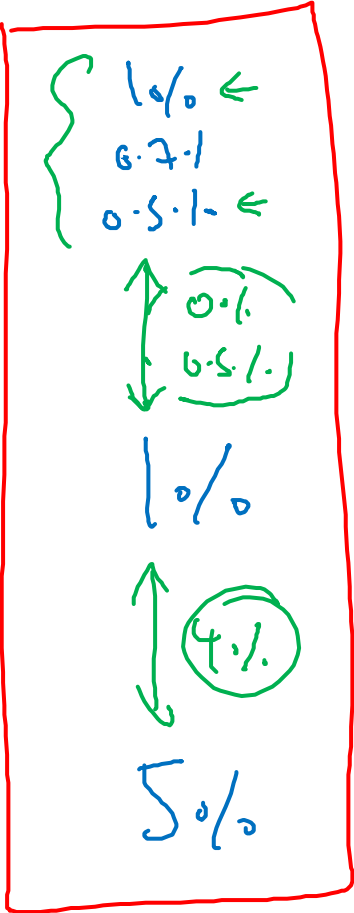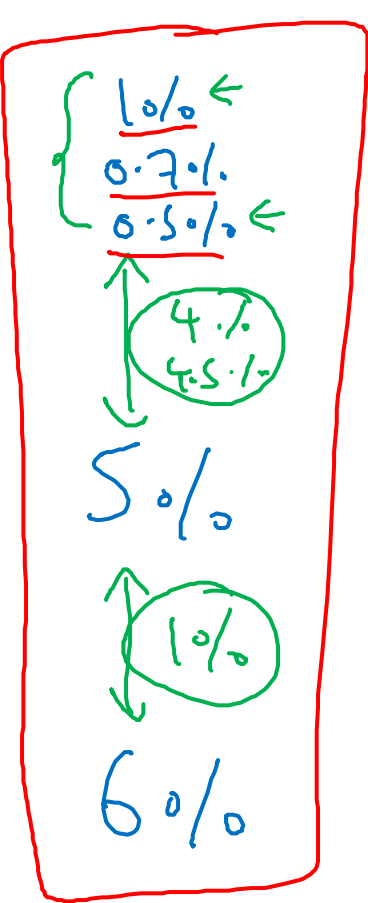
What is "human-level" error?

# Error analysis example

# Summary of bias/variance with human-level performance

0 %

"Bias"

Human-level error
(proxy for Bayes error)

Training error

Dev error

"Avoidable bias"

"Variance"

deeplearning.ai

Comparing to human-level performance

Surpassing human-level performance

# Surpassing human-level performance
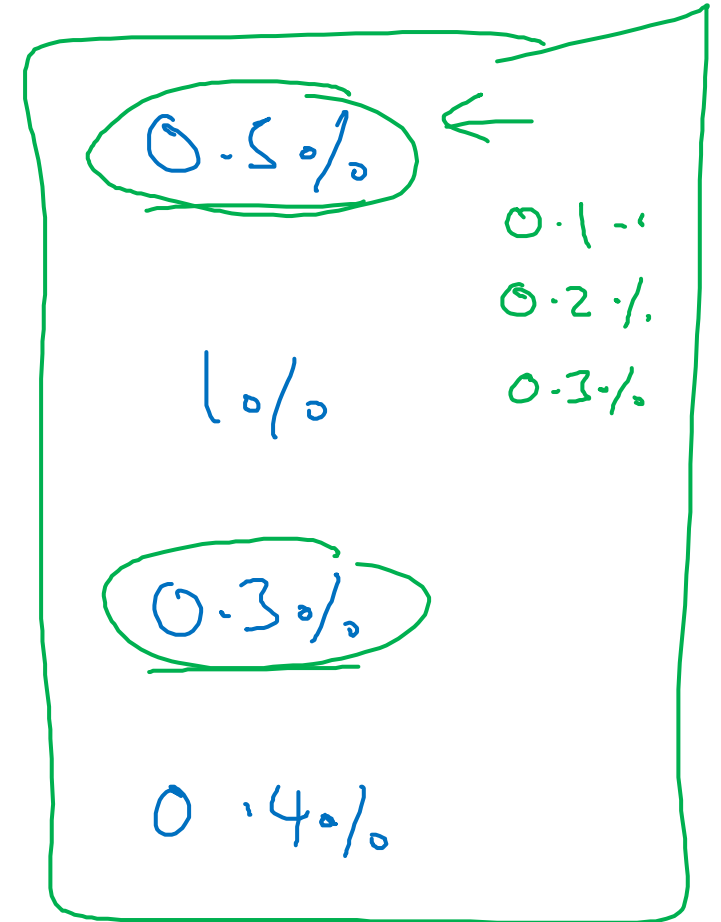
Team of humans    0.5%

One human    ~~1%~~

    0.1

Training error    0.6%

    0.2

Dev error    0.8%

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

0.5%

0.1%
0.2%
0.3%

1%

0.3%

0.4%

# 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, . . .

Comparing to human-level performance

Improving your model performance

deeplearning.ai

# The two fundamental assumptions of supervised learning
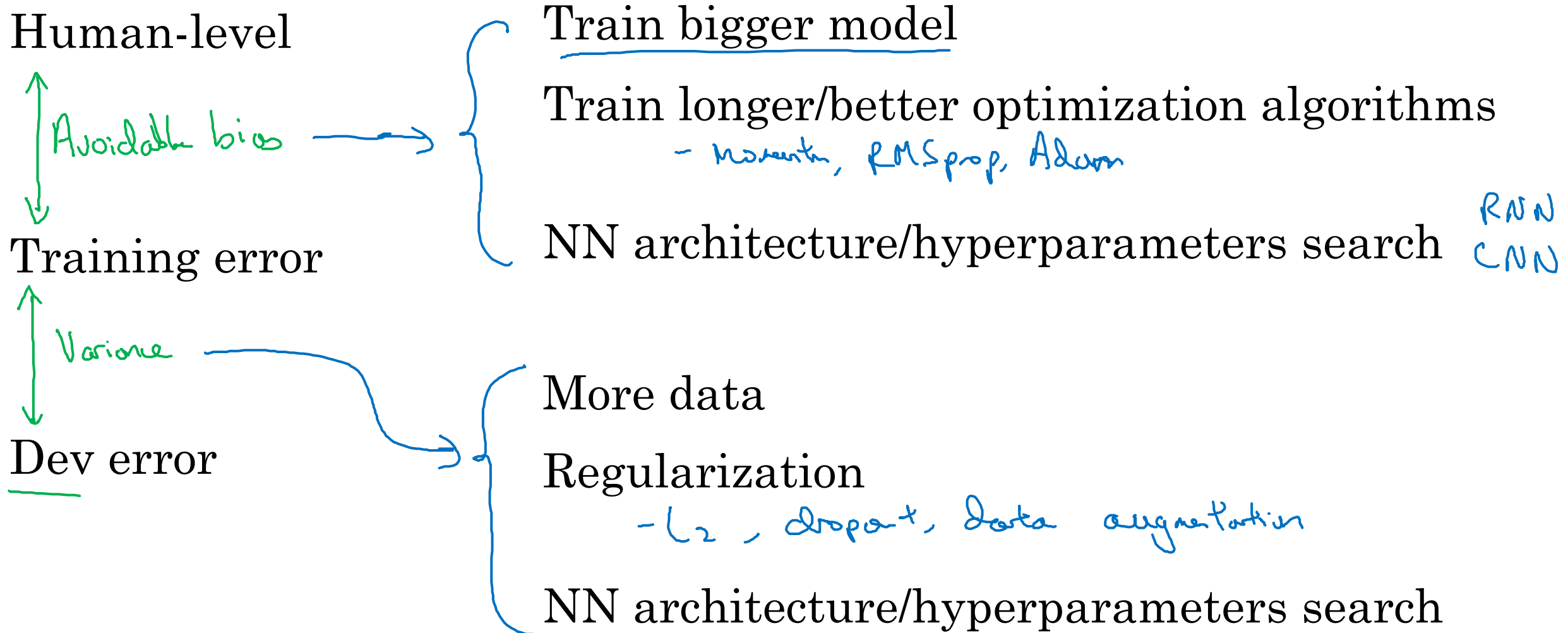
1. You can fit the training set pretty well.

   ~ Avoidable bias

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

   ~ Variance

# Reducing (avoidable) bias and variance

Human-level

Avoidable bias

Training error

Variance

Dev error

Train bigger model

Train longer/better optimization algorithms

   – Momentum, RMSprop, Adam

NN architecture/hyperparameters search

RNN
CNN

More data

Regularization

   – $L_2$, dropout, data augmentation

NN architecture/hyperparameters search

deeplearning.ai

Error Analysis

Carrying out error analysis

# Look at dev examples to evaluate ideas



90% accuracy
→ 10% error

Should you try to make your cat classifier do better on dogs? ←

Error analysis: → 5-10 min
- Get ~100 mislabeled dev set examples.
- Count up how many are dogs.

→ 5% 10%
5/100 9.5%

"ceiling"
→ 50% 10%
50/100 ↓ 5%

# Evaluate multiple ideas in parallel

Ideas for cat detection:

- Fix pictures of dogs being recognized as cats ←

- Fix great cats (lions, panthers, etc..) being misrecognized ←

- Improve performance on blurry images ←

| Image | Dog | Great Cats | Blurry | Instagram | Comments |
|-------|-----|------------|--------|-----------|----------|
| 1 | ✓ | | | ✓ | Pitbull |
| 2 | | | ✓ | ✓ | |
| 3 | | ✓ | ✓ | | Rainy day at zoo |
| ⋮ | ⋮ | ⋮ | ⋮ | | |
| % of total | 8% | 43% | 61% | 12% | |

deeplearning.ai

Error Analysis

---

Cleaning up Incorrectly labeled data

# Incorrectly labeled examples

x

y    1      0      1      1      0      1      1

Training Set.

DL algorithms are quite robust to <u>random errors</u> in the training set.

Systematic errors

# Error analysis

| Image | Dog | Great Cat | Blurry | Incorrectly labeled | Comments |
|---|---|---|---|---|---|
| ... | | | | | |
| 98 | | | | ✓ (circled) | Labeler missed cat in background ← |
| 99 | | ✓ | | | |
| 100 | | | | ✓ (circled) | Drawing of a cat; Not a real cat. ← |
| % of total | 8% | 43% | 61% | 6% (circled) | |

Overall dev set error ‑‑‑‑‑‑‑‑‑‑‑‑‑ 10%        2%

Errors due incorrect labels ‑‑‑‑‑‑‑ 0.6% ←        0.6%

Errors due to other causes ‑‑‑‑‑‑‑ 9.4% ←        1.4%

2.1%        1.9%

Goal of dev set is to help you select between two classifiers A & B.

# Correcting incorrect dev/test set examples

- Apply same process to your dev and test sets to make sure they continue to come from the same distribution

- Consider examining examples your algorithm got right as well as ones it got wrong.

- Train and dev/test data may now come from slightly different distributions.

# Speech recognition example

- Noisy background
  - Café noise
  - Car noise
- Accent
- Far from
- Young
- Stutter
- ...

Guideline:

**Build your first system quickly, then iterate**

- Set up dev/test set and metric
- Build initial system quickly
- Use Bias/Variance analysis & Error analysis to prioritize next steps.

Mismatched training and dev/test data

Training and testing on different distributions

deeplearning.ai

# Cat app example

## Data from webpages

## Data from mobile app

Care about this



≈ 200,000 → 210,000 ← ≈ 10,000

(shuffle

**X Option 1:**

train | dev | test
205,000 → 2,500
→ 2,500

200K / 210K

2381 – web
119 – mobile app

**Option 2:**

train
web

train: 205,000

app
app 2500 | app 2500

# Speech recognition example

Speech activated rearview mirror

## Training

Purchased data $x, y$

Smart speaker control

Voice keyboard

...

500,000 utterances

## Dev/test

Speech activated rearview mirror

→ 20,000

10K    5K 5K D T

train
500 K

510K    D T

10K mirror    5k 5k

Mismatched training and dev/test data

Bias and Variance with mismatched data distributions

deeplearning.ai

# Cat classifier example

Assume humans get ≈ 0% error.

Training error ···· $1\%$ ↓ $9\%$
Dev error ······· ··· $10\%$

Training-dev set: Same distribution as training set, but not used for training

→ train
→ train-dev
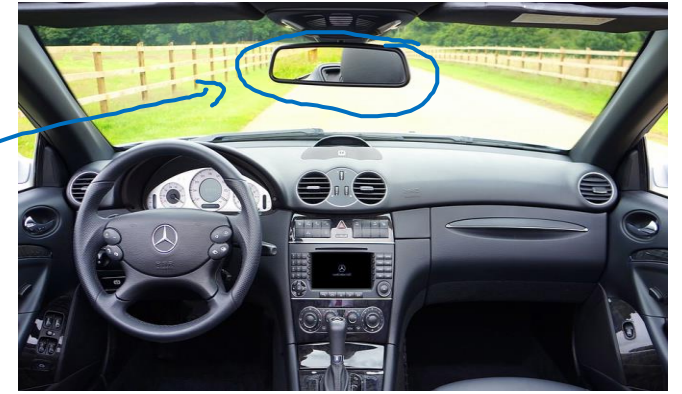→ dev test

NN saw this

| | | |
|---|---|---|
| Traing error | $1\%$ | $1\%$ |
| → Traing-dev error | $9\%$ Variance | $1.5\%$ data |
| → Dev error | $10\%$ | $10\%$ mismatch |

Variance

| | | | |
|---|---|---|---|
| Human error ···· | $0\%$ Avoidable bias | | $10\%$ Avoidable bias |
| Traing error | $10\%$ | | |
| Traing-dev error | $11\%$ | | $11\%$ Variance |
| Dev error | $12\%$ | | $20\%$ Data mismatch |
| | Bias | | Bias + Data mismatch |

# Bias/variance on mismatched training and dev/test sets

Human level      4%    ⟩ avoidable bias

Traing set error      7%    ↑ variance

Traing - dev set error    10%    ↑ data mismatch

→ Dev error      12%    ⟩ degree of overfity to dev set.

→ Test error      12%

4%

7% }

10% }

6% }

6% }

# More general formulation



Rearview Mirror

|  | General speech recognition | Rearview mirror speech data. |
|---|---|---|
| Human level | "Human level" 4% | 6% |
| Error on examples trained on | "Training error" 7% | 6% |
| Error on examples **not** trained on | "Training-dev error" 10% | Dev/Test error" 6% |

Avoidable bias

Variance

data mismatch

deeplearning.ai

Mismatched training and dev/test data

Addressing data mismatch

# Addressing data mismatch

→ • Carry out manual error analysis to try to understand difference between training and dev/test sets

E.g. noisy — car noise          street numbers

→ • Make training data more similar; or collect more data similar to dev/test sets

E.g. Simulate noisy in-car data

# Artificial data synthesis



"The quick brown fox jumps over the lazy dog."     +     Car noise     =     Synthesized in-car audio

10,000 hours

1 hour
of car noise

Overfit to 1 hour of car noise
10,000 hours

Synthesize
Set of all audio in car

# Artificial data synthesis

Car recognition:



$\approx 20$ cars

Synthesized — All cars

Learning from
multiple tasks

Transfer learning

deeplearning.ai

# Transfer learning



image recognition $(x, y)$ — pre-training

$\rightarrow (x, y)$ — fine-tuning

radiology image    diagnoses

$\hat{y}$

image recognition
$\rightarrow \boxed{1,000,000}$   $\boxed{100} \leftarrow$

$\rightarrow$ radiology diagnosis
$\rightarrow \boxed{100}$   $\boxed{1000}$

$W^{[L]}, b^{[L]}$

$\hat{y}$

x

audio

$\hat{y}$

Speech recognition
10h   10,000h

wakeword/triggerword
detection   1h

50h

# When transfer learning makes sense

Transfer from A → B

- Task A and B have the same input x.

- You have a lot more data for Task A than Task B.

- Low level features from A could be helpful for learning B.

deeplearning.ai

Learning from
multiple tasks

---

Multi-task
learning

# Simplified autonomous driving example



$x^{(i)}$

$$y^{(i)} \quad (4,1)$$

Pedestrians

Cars

Stop signs

Traffic lights

$$y^{(i)} = \begin{bmatrix} 0 \\ 1 \\ 1 \\ 0 \end{bmatrix}$$

$$Y = \begin{bmatrix} y^{(1)} & y^{(2)} & y^{(3)} & \cdots, & y^{(m)} \end{bmatrix}$$

$(4, m)$

# Neural network architecture



$x \longrightarrow$ [layers] $\longrightarrow \hat{y} \in \mathbb{R}^4$

pedestrian $\leftarrow 0$

Car $\leftarrow 1$

stop sign $\leftarrow 1$

traffic light $\leftarrow 0$

Loss: $\underset{(4,1)}{\overset{\hat{y}^{(i)}}{y}}$ $\longrightarrow \frac{1}{m} \sum_{i=1}^{m} \boxed{\sum_{j=1}^{4}} \mathcal{L}(\hat{y}_j^{(i)}, y_j^{(i)})$

Sum only over value of $j$ with 0/1 label.

$\longrightarrow$ Usual logistic loss

$-y_j^{(i)} \log \hat{y}_j^{(i)} - (1 - y_j^{(i)}) \log (1 - \hat{y}_j^{(i)})$

Unlike softmax regression:

One image can have multiple labels

Multi-task learning $\leftarrow$

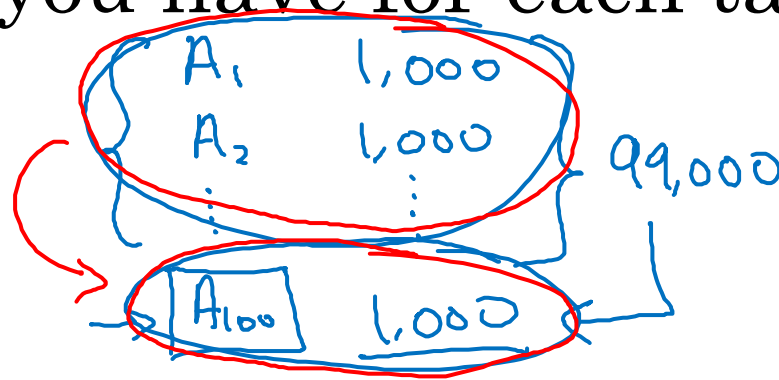$$Y = \begin{bmatrix} 1 & 1 & 0 & ? \\ 0 & 1 & 1 & 1 \\ ? & ? & 1 & ? \\ ? & ? & 0 & ? \end{bmatrix} \leftarrow$$

# When multi-task learning makes sense

- Training on a set of tasks that could benefit from having shared lower-level features.
- Usually: Amount of data you have for each task is quite similar.

$A \quad 1,000,000$

$\downarrow \qquad \downarrow$

$B \quad 1,000$

$A_1 \quad 1,000$

$A_2 \quad 1,000$

$\vdots \qquad \vdots$

$A_{100} \quad 1,000$

$99,000$

- Can train a big enough neural network to do well on all the tasks.

End-to-end deep learning

What is end-to-end deep learning

deeplearning.ai
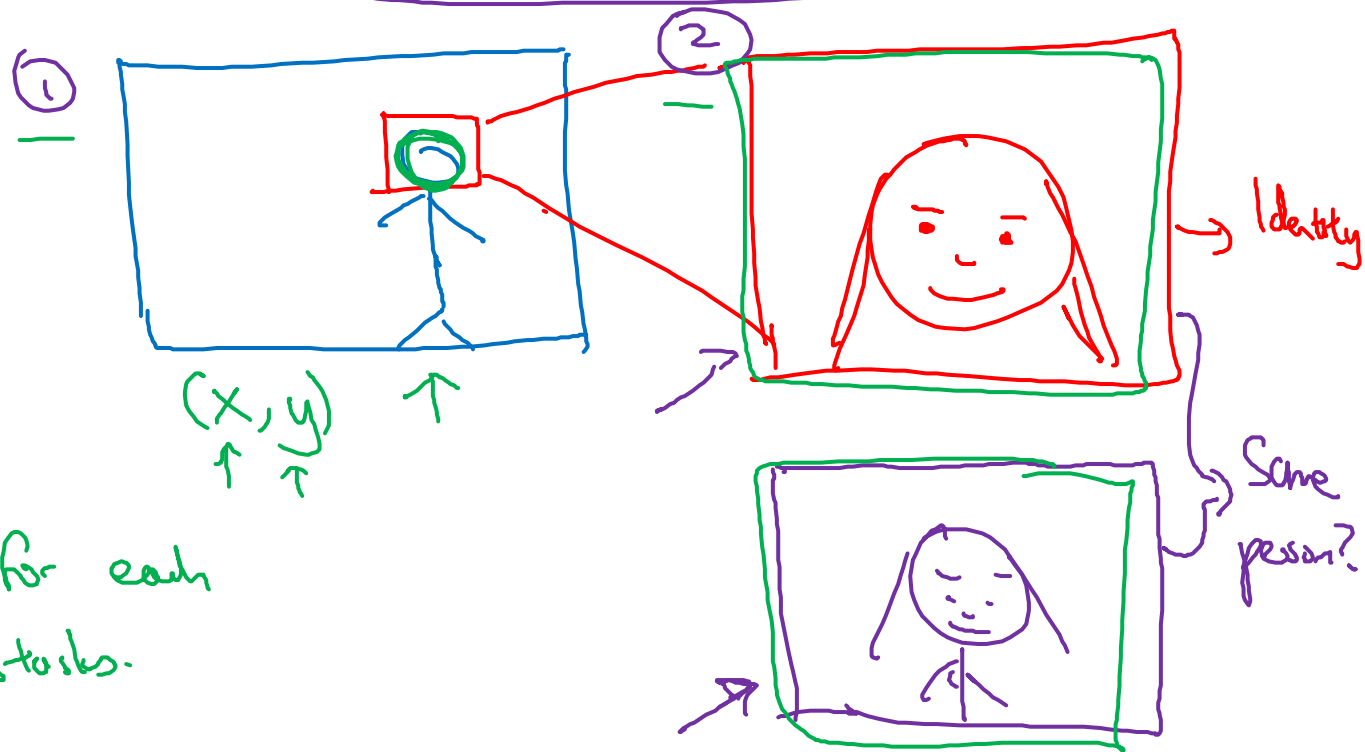
# What is end-to-end learning?

Speech recognition example

"cat"

$x$

$y$

$\rightarrow$ audio $\xrightarrow{\text{MFCC}}$ features $\xrightarrow{\text{ML}}$ Phonemes $\rightarrow$ Words $\rightarrow$ transcript

audio

$\rightarrow$ audio $\longrightarrow$ phonemes $\rightarrow$ $\cdots$ $\rightarrow$ transcript

3,000h

10,000h
$\vdots$
100,000h

# Face recognition



[Image courtesy of Baidu]

Image (x) $\longrightarrow$ y

Identity

(x, y)

① (x, y)

② Identity

Some person?

Have data for each of 2 subtasks.

# More examples

Machine translation

$(x, y)$

English ↗   ↖ French

Eglsh

English → text analysis → ... → French

English ————————————→ French

Estimating child's age:



Image ①→ bones ②→ age

Image ————————————→ age ←

deeplearning.ai

End-to-end deep learning
_____
Whether to use
end-to-end learning

# Pros and cons of end-to-end deep learning

Pros:

- Let the data speak $\qquad$ $X \longrightarrow y$ $\qquad$ $\rightarrow$ "phonemes" $\underline{c\ a\ t}$
- Less hand-designing of components needed

$X\ -\ -\ -\ -\ -\ \rightarrow y$ $\qquad$ input end $\qquad$ output end

Cons:

- May need large amount of data $\qquad$ $\underline{X \longrightarrow y}$ $\qquad$ $\underline{(x,y)}$
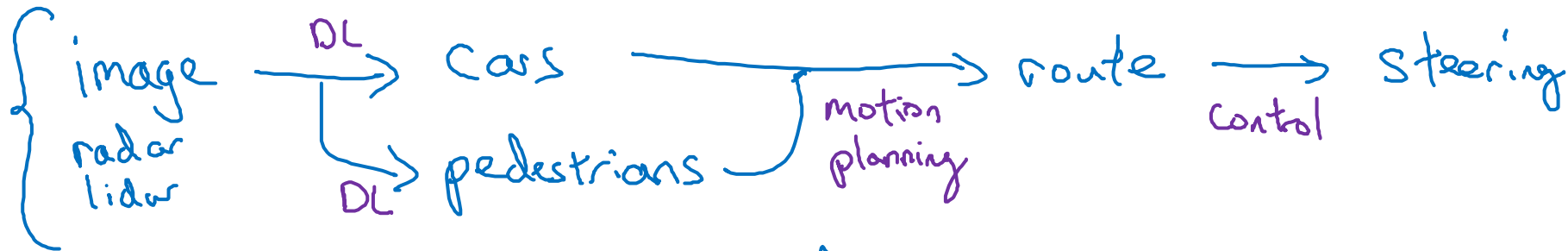- Excludes potentially useful <u>hand-designed components</u> $\qquad$ <u>Data.</u> $\qquad$ <u>Hard-design.</u>
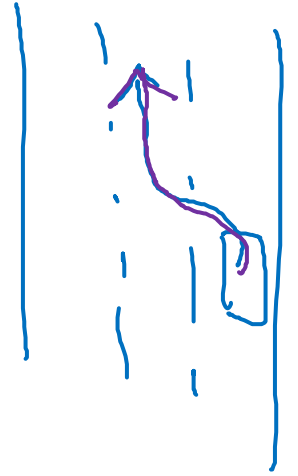
# Applying end-to-end deep learning

Key question: Do you have sufficient data to learn a function of the complexity needed to map x to y?

$x \rightarrow y$

$\rightarrow$ age

image
radar
lidar

$\xrightarrow{DL}$ cars

$\xrightarrow{DL}$ pedestrians

motion planning

$\rightarrow$ route

control

$\rightarrow$ steering

- Use DL to learn individual components
- Carefully choose $x \rightarrow y$ depending what tasks you can get data for.

$\rightarrow$ image $\longrightarrow$ steering