



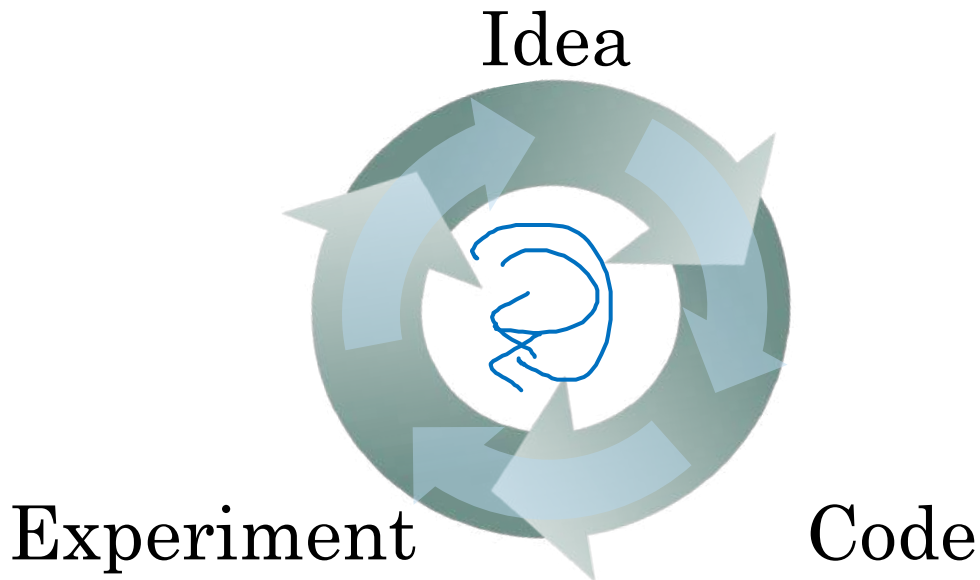
deeplearning.ai

Setting up  
your goal

---

Single number  
evaluation metric

# Using a single number evaluation metric



→ Of examples recognized as cost, what % actually are costs?

→ what % of actual costs are correctly recognized

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

F<sub>1</sub> score = "Average" of P and R.

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

Dev set + Single number evaluation metric  
real 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%





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

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

maximize accuracy

subject to Running Time  $\leq$  100 ms.

N metrics : 1 optimizing  
N-1 satisfying

Wakewords / Trigger words

Alexa, OK Google,

Hey Siri, nihao baidu  
你好 百度

accuracy.

#false positive

maximize accuracy.

s.t.  $\leq$  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

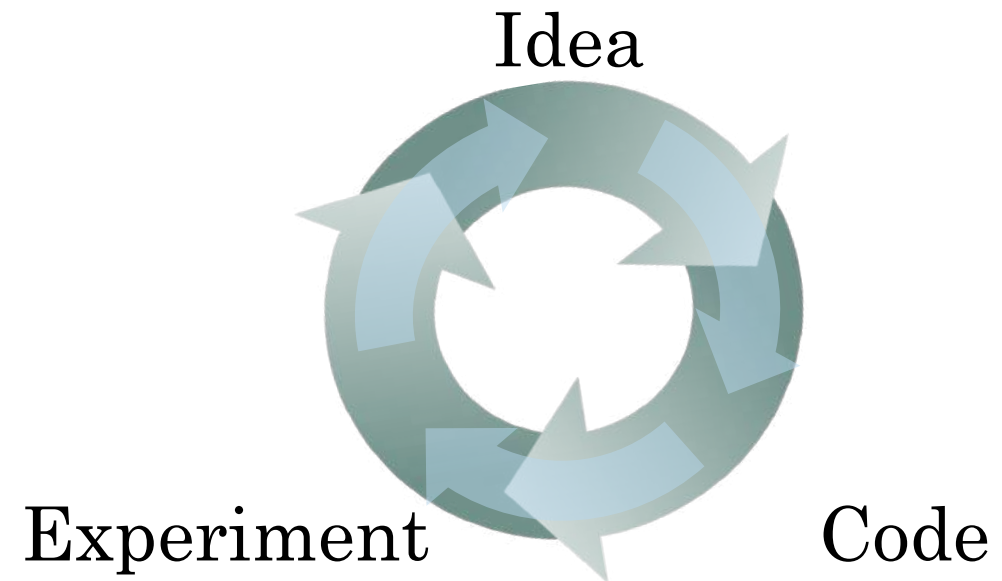
Dev

Test

→ Randomly shuffle into dev/test



dev set  
+  
metric



# True story (details changed)

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

↑

$x \rightarrow y$  (repay loan?)



[ Tested on low income zip codes

~ 3 month

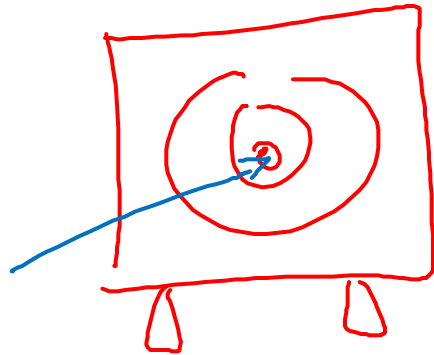




# Guideline

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



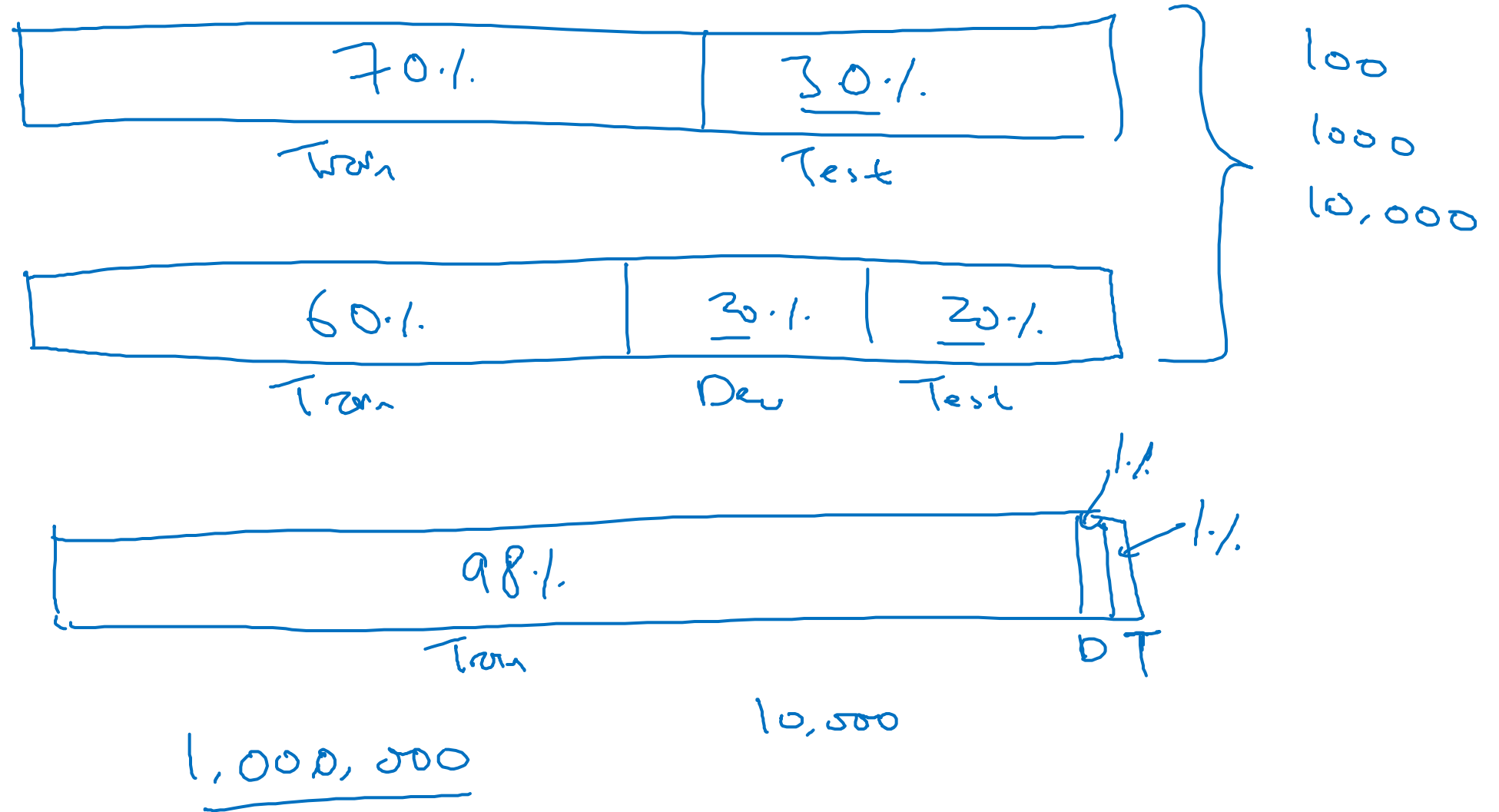
deeplearning.ai

Setting up  
your goal

---

Size of dev  
and test sets

# Old way of splitting data



# 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

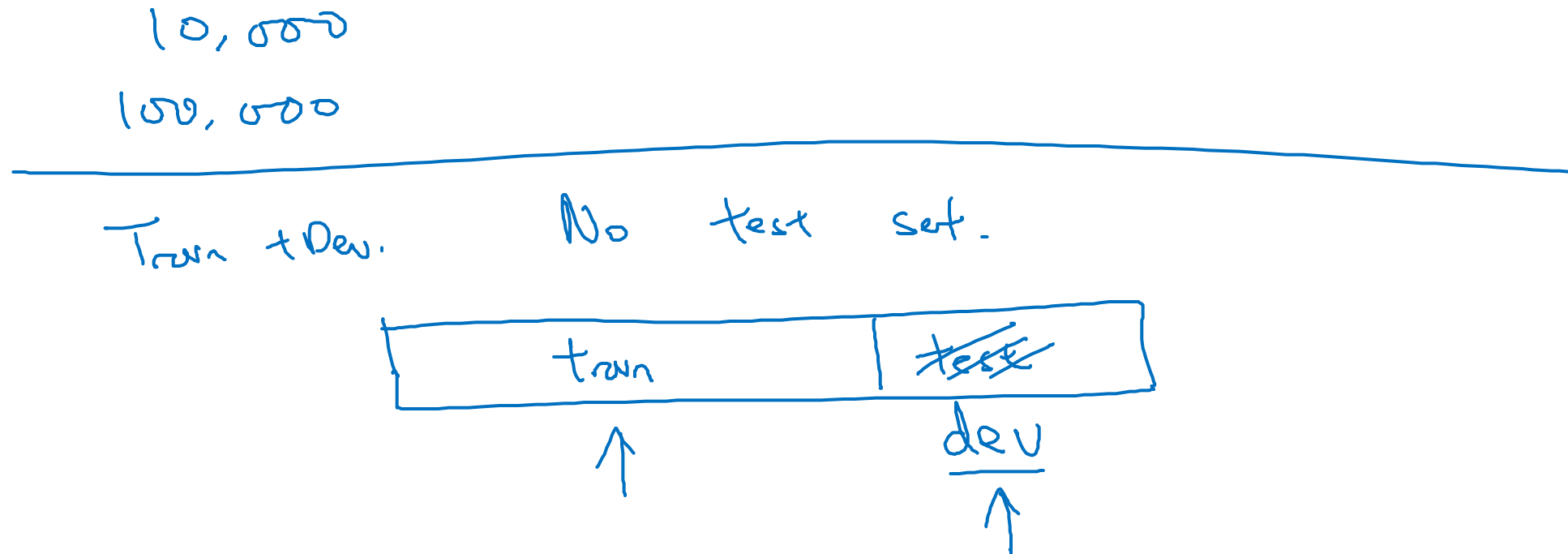
<sup>A</sup> 97% → <sup>B</sup> 97.1%  
0.1%  
└

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.





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.

→ Metric: classification error

Algorithm A: 3% error

→ pornographic

✓ Algorithm B: 5% error

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

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

$\mathbb{I}\{y_{pred}^{(i)} \neq y^{(i)}\}$   
predicted value (0/1)

# 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)

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





# 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.