# Chain of assumptions in ML

→ Fit training set well on cost function ✍ bigger network, Adam, ... *(≈ human-level performance)*
  If not, try bigger network, Adam ...

  ↓ ☐

→ Fit dev set well on cost function ✍ Regularization, Bigger traing set
  If not, try regularization, bigger training set ...

  ↓ ☐

→ Fit test set well on cost function ✍ Bigger dev set
  If not, try bigger dev set

  ↓

→ Performs well in real world ✍ Change dev set or cost function
  *(Happy cat pic app users.)*
  If not, change dev set or change cost function

Single real number evaluation metric:

# Using a single number evaluation metric

→ Of examples recognized as cats, what % actually are cats?.

→ what % of actual cats are correctly recognized

P and R is a pair of trade-off

Idea

Experiment     Code

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

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

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

calculated based on dev set

Dev set + Single number evaluation metric → real Speed up iterating

## 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% |

Optimizing and satisfying metric:

## Another cat classification example

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

optimizing    Sastisficing

$Cost = accuracy - 0.5 \times running\ Time$

try to **Maximize** accuracy

with **Subject to** running Time $\leq 100$ ms.

N metrics: 1 optimizing
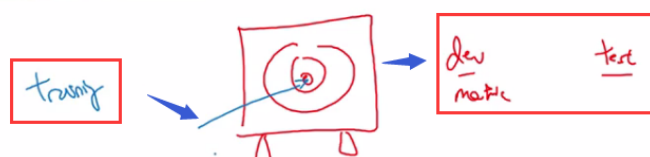N-1 Sastisficing

Wakewords / Trigger words

Alexa, OK Google,
Hey Siri, ni hao baidu
你好百度

accuracy.
#false positive

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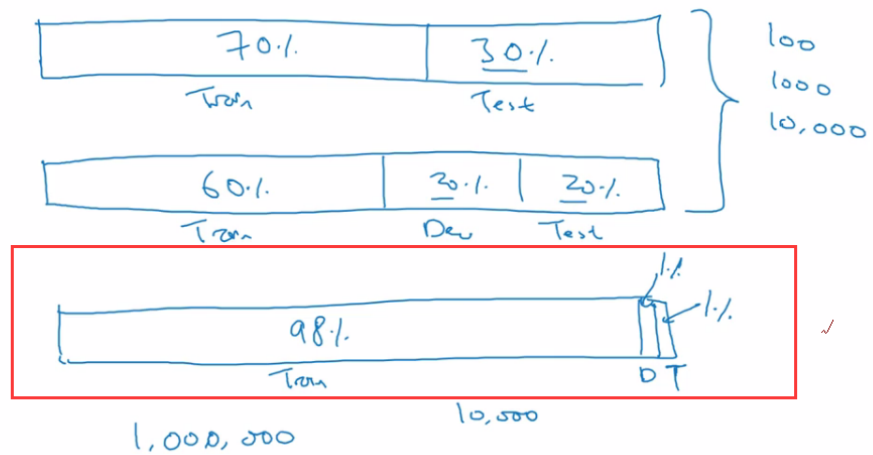s.t. $\leq 1$ false positive
every 24 hours.

Train, Dev and Test set splitting:

## 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.
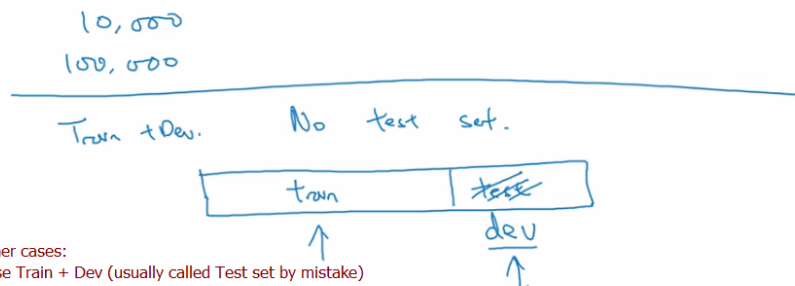
train → dev/metric → test

# Old way of splitting data

| 70% | 30% |
| :---: | :---: |
| Train | Test |

100
1000
10,000

| 60% | 30% | 20% |
| :---: | :---: | :---: |
| Train | Dev | Test |

| 98% | | 1% / 1% |
| :---: | :---: | :---: |
| Train | | D T |

10,500

1,000,000

# 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 |
| :---: | :---: |
| ↑ | ↑ |

And some other cases:
people just use Train + Dev (usually called Test set by mistake)

When to change Dev/Test set and metric:

# Cat dataset examples

Metric + Dev : Prefer A

But You/users : Prefer B.

### Metric: classification error

Algorithm A: 3% error $\longrightarrow$ pornographic

Because model A treats porn and non-porn equally

✓ Algorithm B: 5% error

function counts the number of non-match cases

Error: $\dfrac{1}{m_{dev}} \displaystyle\sum_{i=1}^{m_{dev}} \mathbb{I}\{ y_{pred}^{(i)} \neq y^{(i)} \}$

↳ predicted value (0/1)

that's the sign that we need to change evaluation metric

# Cat dataset examples

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

→ Metric: classification error

Algorithm A: 3% error $\longrightarrow$ pornographic

✓ Algorithm B: 5% error

Error: $\dfrac{1}{\sum \omega^{(i)}} \; \xcancel{\dfrac{1}{m_{dev}}} \displaystyle\sum_{i=1}^{m_{dev}} \omega^{(i)} \, \mathbb{I}\{ y_{pred}^{(i)} \neq y^{(i)} \}$

↳ predicted value (0/1)

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

Add more weight to the Error term

# Orthogonalization for cat pictures: anti-porn

→ 1. So far we've only discussed how to define a <u>metric</u> to evaluate classifiers. ← Place target 🎯

1st step, place the target

→ 2. Worry separately about how to do well on this metric. 🎯

↖ Aim (shot at target)

→ $J = \dfrac{1}{\sum \omega^{(i)}} \displaystyle\sum_{i=1}^{m} \omega^{(i)} \, \mathcal{L}(\hat{y}^{(i)}, y^{(i)})$

2nd step, define a cost function to shot the target

## Another example

Algorithm A: 3% error
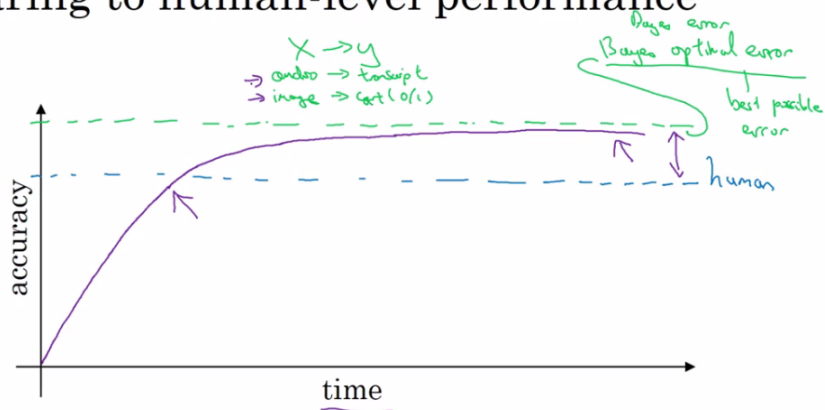✓ Algorithm B: 5% error ←

→ 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.
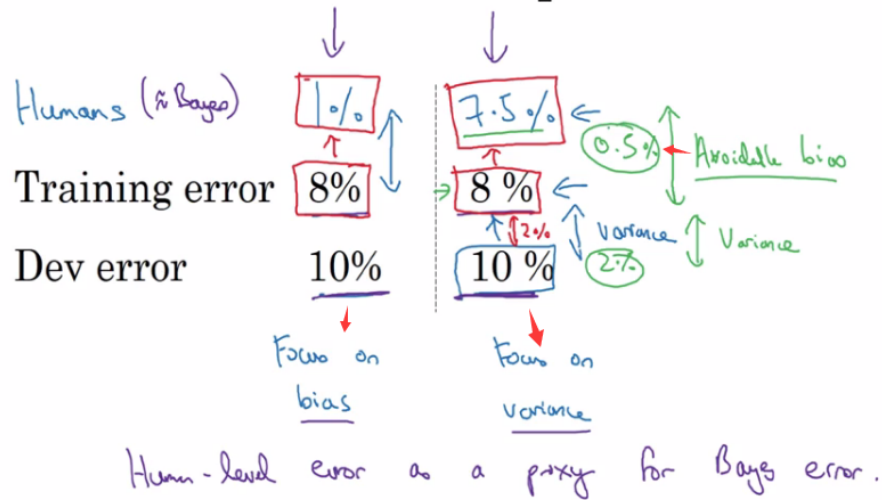
Human Level Performance:

## Comparing to human-level performance



Avoidable Bias:

# Cat classification example

Humans (≈ Bayes)   1%      7.5% ← 0.5% ← Avoidable bias

Training error   8%    → 8% ←    ↑2% Variance ↕ Variance

Dev error   10%      10%    2%

Focus on        Focus on
bias            Variance

Humm-level error as a proxy for Bayes error.

---

# Error analysis example

It depends which Bayes Error you choose

Human (proxy for Bayes error)
Avoidable bias

{ 1% ←      { 1% ←       0.7%
0.7%         6.7%        0.5%
0.5% ←       0.5% ←

4% / 4.5%    0%/ 6.5%     0.2% ←
                          0.0%

Training error   5%      1%      0.7%

Variance   1%      4%      0.1% ←

Dev error   6%      5%      0.8%

Focus on Bias        Focus on Variance

---

Improving Model Performance:

# The two fundamental assumptions of supervised learning

$\rightarrow$ 1. You can fit the training set pretty well.

$\sim$ Avoidable bias

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

$\sim$ Variance

# Reducing (avoidable) bias and variance

| Human-level | Train bigger model |
|---|---|
| $\uparrow$ Avoidable bias $\longrightarrow$ $\downarrow$ | Train longer/better optimization algorithms<br>— momentum, RMSprop, Adam |
| Training error | NN architecture/hyperparameters search RNN CNN<br><span style="color:red">Change the activation functions, cost functions</span> |
| $\uparrow$ Variance $\downarrow$ | More data |
| Dev error | Regularization<br>— $L_2$, dropout, data augmentation |
|  | NN architecture/hyperparameters search |