

# WEEK 1. ML Strategy

## 1. Introduction to ML Strategy

### • Why ML Strategy?

#### - Motivating example

Let's say you developed Cat classifier & achieved 90% of accuracy.

but let's say 90% is not enough for your application.

⇒ Then, you'll try following ideas.

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

⇒ 초기의 경우 오류시간이 점차로 개선이 되기 않을 수 있다.

여기 있는 아이디어 중 무엇이 충분한지 빠르게 표시하고 알아내는 법을 배워보자!

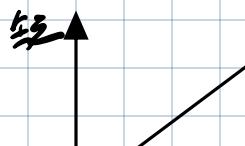
### • Orthogonalization

#### - TV tuning example

예를 들어 TV를 차를 작동하기 위한 여러개의 버튼이 있다고 해보자.

이 버튼들을 이용하면 [회전 가로길이가 0.1배 규칙으로, 동시에 회전도  $0.5^\circ$  한다  
오른쪽으로 0.5 방향이 틀리는 동시에 0.9의 속도로 간다] 고 하자.

⇒ 이 버튼들은 한 번의 조작에 여러 동작을 동시에 처리하는 것이기 때문에 TV/차를  
쉽게하게 조작하기 힘들다.



Orthogonalization은 TV/자동차  
같은 여러 가지 행동이 동시에 처리되

방향  
orthogonal

다양한 기능이 있는 것을  
기능만 기준을 만드는 것을 말한다.

## - Chain of assumptions in ML

→ 기계학습이 잘 되려면? 학습의 비율을 잘 조정해서 4가지를 만족해야 함.

① Fit training set well on cost function ( $\approx$  human-level performance)

↳ button [bigger network  
Adam  
...]



② Fit dev set well on cost function ← early Stopping은 less orthogonalized

↳ button [Regularization  
Bigger training set]

① all 다른 비율이 될 수 있지만 ②의 성능은 당연히 개선하기가 때문에 orthogonalization이 잘 안 될 수 있음.

③ Fit test set well on cost function.

↳ button [Bigger dev set]

④ Performs well in real world ( $\approx$  happy cat pic users)

↳ button [Change dev set or cost function]

## 2. Setting up your goal

### • Single number evaluation metric

→ 실수(real number)로 된 평가기준은 몇몇 기준을 고려해 계산된다.

### - Using a single number evaluation metric

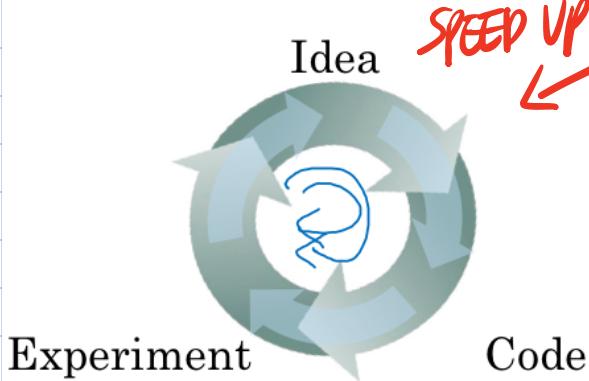
Classifier	Precision	Recall
A	95%	90%
B	98%	85%

) which one  
is better?

Redefine  
new evaluation metric

⇒ F1 Score  
= Average of P & R  
 $(\frac{1}{\frac{1}{P} + \frac{1}{R}})$ , "Harmonic mean")

dev set + single real # evaluation metric



# SPEED UP!

## - Another Example

Algorithm	US	China	India	Other
A	<del>errors</del> 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%

Make  
average of  
errors &  
pick one

- Satisficing (足りる) & optimizing metrics
    - ⇒ How to efficiently set up evaluation 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

⇒ 이러한 선형경향은 다소 인위적

⇒ Establish optimizing metric & satisficing metric

↳ metric you want to optimize      ↳ 어떤 조건만 만족하면 되는 것도 얼마나 좋은지 살피면 X

(e.g. Maximize accuracy)

(e.g. running time  $\leq 100$  ms, false negativity)

→ 모든 N개의 metric이 조율되면 1개의 optimizing metric & (N-1)개의 sofisticating metric으로 구성

- ## • Train / dev distribution

## - Old way of splitting data

70%	30%	
Train	Test	
60%	20%	20%
train	dev	Test
98%	1%	1%
Train	dev	Test

↳ allocate more on training set

for earlier era of ML  
 $(M=100, 1000, 10,000 \dots)$

$\Rightarrow$  modern ML  
 $(M=1,000,000, \dots)$

## - Size of dev set

$\Rightarrow$  Set your dev set to big enough to detect differences in algorithms models you're trying out. (A/B test 할 때 차이를 찾을 수 있을 만큼)

## - Size of test set

$\Rightarrow$  Set your test set to be big enough to give high confidence in the overall performance of your system

- 어떤 application에서는 train/dev가 70%, test set은 30% 정도로 되는 경우 (기본)

## • When to change dev/test sets and metrics

### - Cat dataset examples

Metric: classification error

Algorithm A: 3% error

Algorithm B: 5% error

but let's say...

also shows pornographic pictures to users

no porn

$\Rightarrow$  Metric + Dev : prefer A, You / Users : prefer B

→ 평가척도가 높아 더 나은 알고리즘인가 순위를 못 알려주는 것

→ appol 잘하는 평가척도 재정의 필요.

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

indicator function

$\hookrightarrow$  0 혹은 1이면 = 안 맞는 시험

$$\text{Error} = \frac{1}{\sum_i W^{(i)}} \sum_{i=1}^{m_{\text{dev}}} W^{(i)} I \{y_{\text{pred}}^{(i)} \neq y^{(i)}\}$$

indicator function  
 ↳ ok한 사진 = 양 아한 사진  
 ↳ pornographic picall 아한 기준  
 ↳ porn → error↑

- Orthogonalization for cat pictures: anti-porn

→ 모든 분류기의 손상을 예기할 수 있도록 평가准则를 정의

① So far we've only discussed how to define a metric to evaluate classifiers

⇒ 찾고자 하는 = Place target

② Worry separately about how to do well on this metric

⇒ How to aim the target accurately → 이 과정에서 찾고자 하는게 가능할 수 있음.  
 (각 찾고자 하는 어떻게 좋은 성능을 낼지에 대한 고민)

- Another example

[Algorithm A: 3% error

Algorithm B: 5% error

→ Dev/test



→ B가 오히려 실제 사용에서 더 잘 작동할 수 있음.

→ User images



⇒ Metric + dev/test set에서는 성능이 좋지만 application에서는 잘 작동하지 못하는 경우, Metric & dev/test를 바꿔야 함

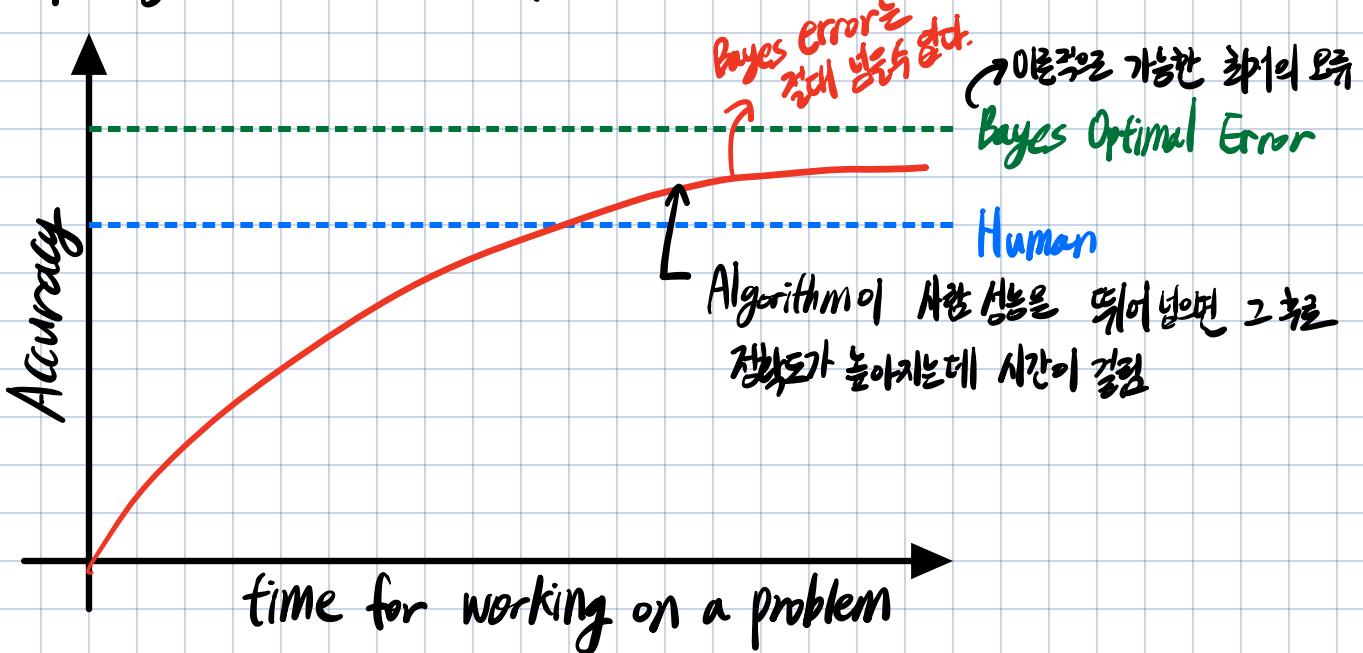
### 3. Comparing to human-level performance

• Why human-level performance?

① Deep learning → ML performance competitive to human

② 사람이 할 수 있는 일에 대해 ML을 디자인하는 게 효율적 → 비효율적으로 학습하기

## - Comparing to human-level performance



예를 들어,  $X \rightarrow Y$ 에서  $X$ 의 품질이 사람도 못 알아볼 정도로 떨어진다면 모델은 잘 알아보는 경향  
(Audio  $\rightarrow$  transcript)  
(image  $\rightarrow$  class (0/1))

## - Why compare to human-level performance

$\Rightarrow$  Humans are quite good at a lot of tasks.

if ML performance < humans, you can:

① Get labeled data from humans.

② Gain insight from manual error analysis:

Why did a person get this right?

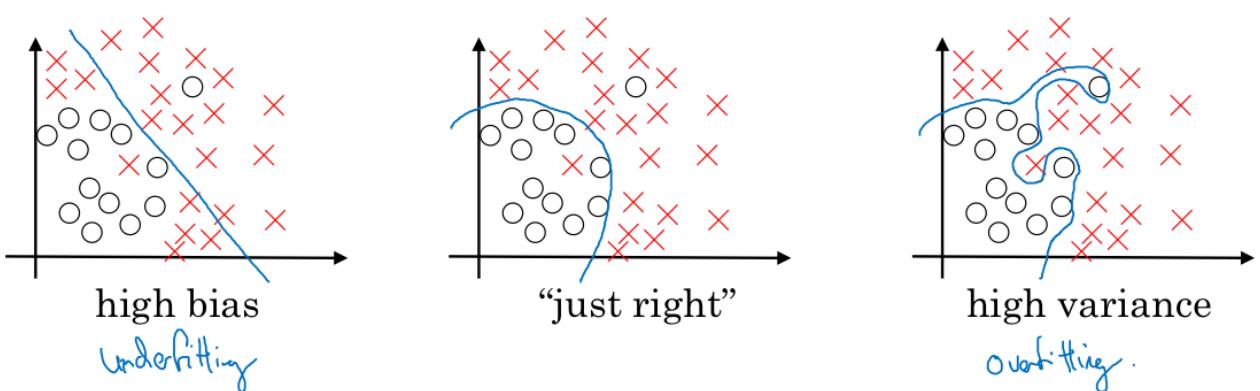
③ Better analysis of bias/Variance  $\leftarrow$  cannot work if model > humans

### Avoidable bias

- Sometimes we don't want ML to do too well

$\Rightarrow$  Human performance로 학습세트에 대한 알고리즘의 가능한 성능을 가능할 수 있음

- Bias & Variance



## - Cat Classification example

Humans ( $\approx$ Bayes)	1%	7.5%
Training Error	8%	8%
Dev Error	10%	10%
Bayes Error - Training error:	Focus on Avoidable Bias	Focus on Variance
	↓	↓
		. Training error - Test error

## • Understanding human-level performance

### - Human-level error as a proxy for Bayes error

⇒ What is "human-level" error? ⇒ Define purpose (proxy of bayes error)

e.g. Medical image classification example

① Typical human: 3% error

② Typical doctor: 1% error ⇒ system이 올바르면

③ Experienced doctor: 0.9% error

④ Team of experienced doctors: 0.5% error ⇒ Bayes Error all 대비  
추상화가 올바르면

### - Error analysis example

Humans ( $\approx$  Bayes)

1% / 0.9% / 0.5%

Training Error

5%

1 %

Dev Error

6 %

5 %

0.9% / 1% / 0.5%

0.2%

0.7% ) 0.1% X2

0.8% ) X1

↓  
Focus on  
Avoidable  
Bias

↓  
Focus on  
Variance

↓  
It's important to  
determine what's  
Bayesian Error

⇒ 지금 성능을 기沽을 했을 때 Variance / bias ≠ 0이면 광범위한 알고리즘 알기 어려움.

### - Summary of bias / variance with human-level performance

Human-level error (Proxy for Bayes Error)

사람이 알고 있는 데이터 사용률의 정도

Training error

Dev error

⇒ 0.5%까지 Bayes error  $\leq 0\% \leq 1\%$  but 0.7% Bayes error가 0.1% 광범위한 알고리즘 noise를 얻는 경우, use better estimate of Bayes error by human-performance

### • Surpassing human-level performance

#### - Why is it difficult to improve ML algorithm after surpassing human performance?

One human : 0.5%  
Team of humans : 1%  
Training error : 0.3%  
Dev error : 0.4%

→ overfitted? / Bayes error is 0.3%?

⇒ don't have enough info to reduce bias / variance in your algorithm

#### - Problems where ML significantly surpasses human-level performance.

① Structured data

② Not natural perception

e.g. Speech / Image recognition, medical (CG, skin cancer)

③ Lots of data

e.g. Online advertising, product recommendations, Logistics  
(predicting transit time), Loan approvals

## • Improving your model performance

- The two fundamental assumptions of supervised learning

- ① You can fit the **training set** pretty well  $\Rightarrow$  achieve low **avoidable bias**
  - ② The training set performance generalizes pretty well to the **dev/test set**  
 $\Rightarrow$  good **variance**
- orthogonalization**

- Reducing (avoidable) bias and Variance

$\Rightarrow$  look at the difference between performance (= proxy of bayes error)

