

¿CÓMO ESTRUCTURAR UN BUEN PROYECTO DE *MACHINE LEARNING*?

Anna Bosch Rue

VP Data Intelligence @ Launchmetrics

annaboschrue@gmail.com

Motivating example



90% Accuracy and you want to do better

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 regularization
- Network architecture:
 - Activation functions
 - # hidden units
 - ...

CONTENT

1. PART 1

- 1.1. Introduction to ML Strategy
- 1.2. Setting up your goal
- 1.3. Comparing to human-level performance

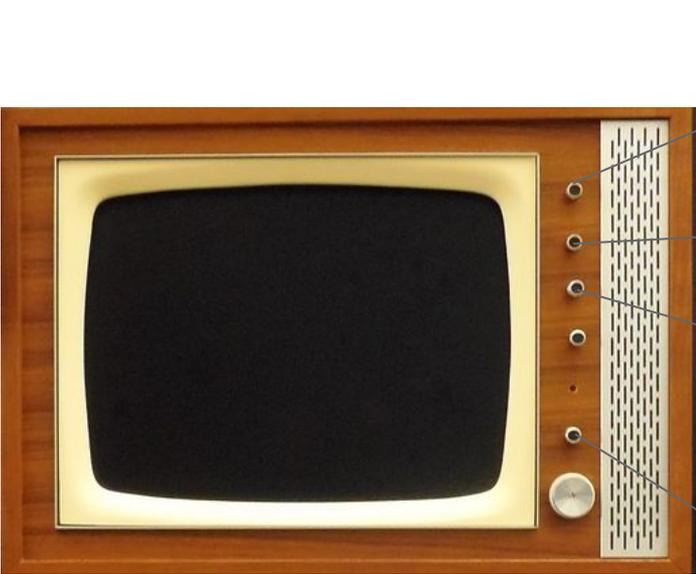
2. PART 2

- 2.1. Error analysis
- 2.2. Mismatched training and dev/test set
- 2.3. Learning from multiple tasks
- 2.4. End-to-end deep learning

PART I

INTRODUCTION TO ML STRATEGY

ORTHOGONALIZATION (I)



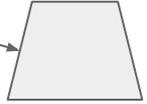
0.1 x width

+



0.3x height

-



1.7x trapez

+



0.8x rotation

ORTHOGONALIZATION (II)

Chain of assumptions in ML

Fit **training set** well on cost function

Bigger network
Adam

Fit **dev set** well on cost function

Regularization
Bigger training set

Fit **test set** well on cost function

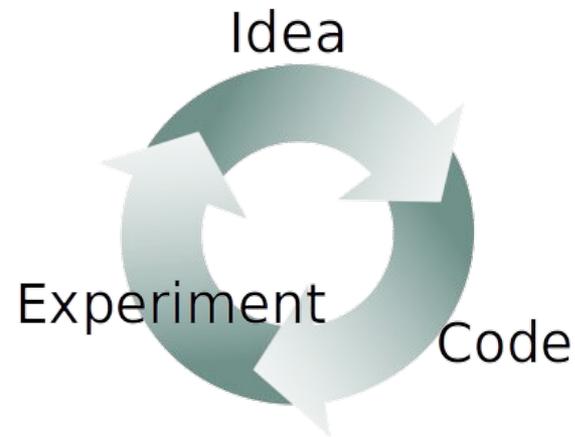
Bigger dev set

Performs well in **real world**

Change dev set
Change cost function

SETTING UP YOUR GOAL

SINGLE NUMBER EVALUATION METRIC



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

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%

SATISFICING AND OPTIMIZING METRIC

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

$$\text{Cost} = \text{Accuracy} - 0.5 \text{ Running time}$$

Maximize Accuracy

Subject to Running time \leq 100ms

N Metrics: 1 Optimizing
N-1 Satisficing

TRAIN/DEV/TEST DISTRIBUTIONS

REGIONS:

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

Dev

Test

Randomly shuffle
into dev/test



True story (By Andrew NG):

- Optimizing on dev set on loan approvals for medium income zip codes
- Tested on low income zip codes

Guideline:

Choose a dev set and test set (same distribution) to reflect data you expect to get in the future and consider important to do well on.

SIZE OF THE DEV & TEST SETS

OLD WAY OF SPLITTING DATA



100
1000
10000

Set your **dev set** to be big enough to detect differences in algorithm/models you're trying out

NEW WAY OF SPLITTING DATA



1000000

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

Train D/T

WHEN TO CHANGE DEV/TEST SETS AND METRICS

CAT DATASET EXAMPLES

Metric: classification error

Algorithm A: 3% \longrightarrow **Pornographic**

Algorithm B: 5%

Orthogonalization: (1) So far we've discussed how to define a metric to evaluate classifiers, (2) worry separately about how to do well on this metric

Error:

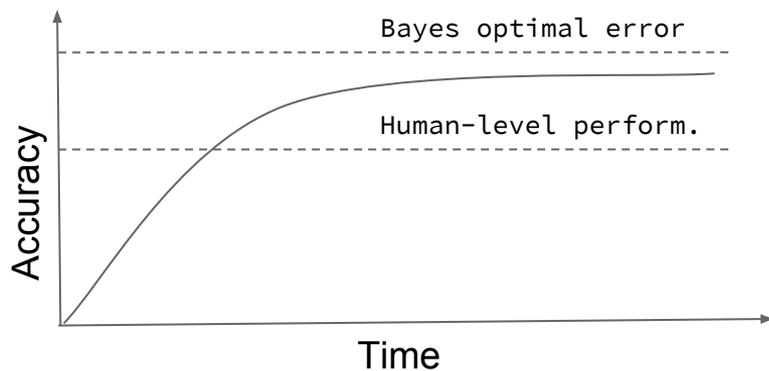
$$\frac{1}{\sum_i w(i)} \quad \cancel{\frac{1}{M_{dev}}} \quad M_{dev} \sum_{i=1} w(i) (\hat{y} - y)$$

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

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?



Humans are quite good at a lot of tasks. So long as ML is worse than humans, you can:

- Get labeled data from humans
- Gain insight from manual error
- Better analysis of bias/variance

AVOIDABLE BIAS

Use human level error as a proxy for Bayes error

BIAS & VARIANCE

Cat classification



Human level (aprox):

Training set error:

Dev set error:

0%	0%	0%	0%
15%	1%	15%	0.5%
16%	11%	30%	1%
High bias	high variance	high bias high variance	low bias low variance

Humans (Bayes) :

Training set error:

Dev set error:

1%	7.5%
8%	8%
10%	10%
Focus on bias	Focus on variance

Avoidable bias

Variance

UNDERSTANDING HUMAN -LEVEL PERFORMANCE (I)

HUMAN-LEVEL ERROR AS A PROXY FOR BAYES ERROR

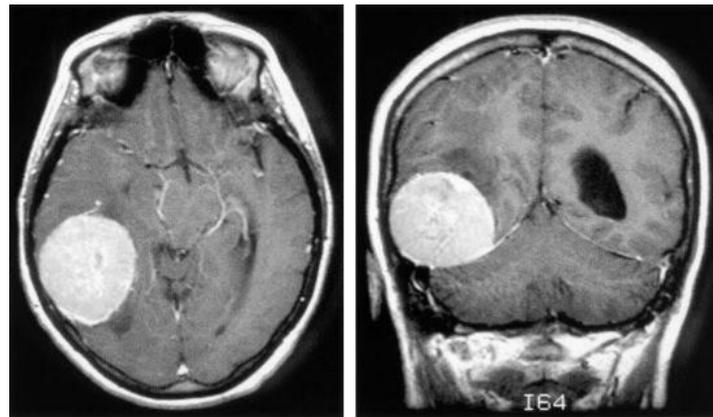
Medical image classification example:

Suppose:

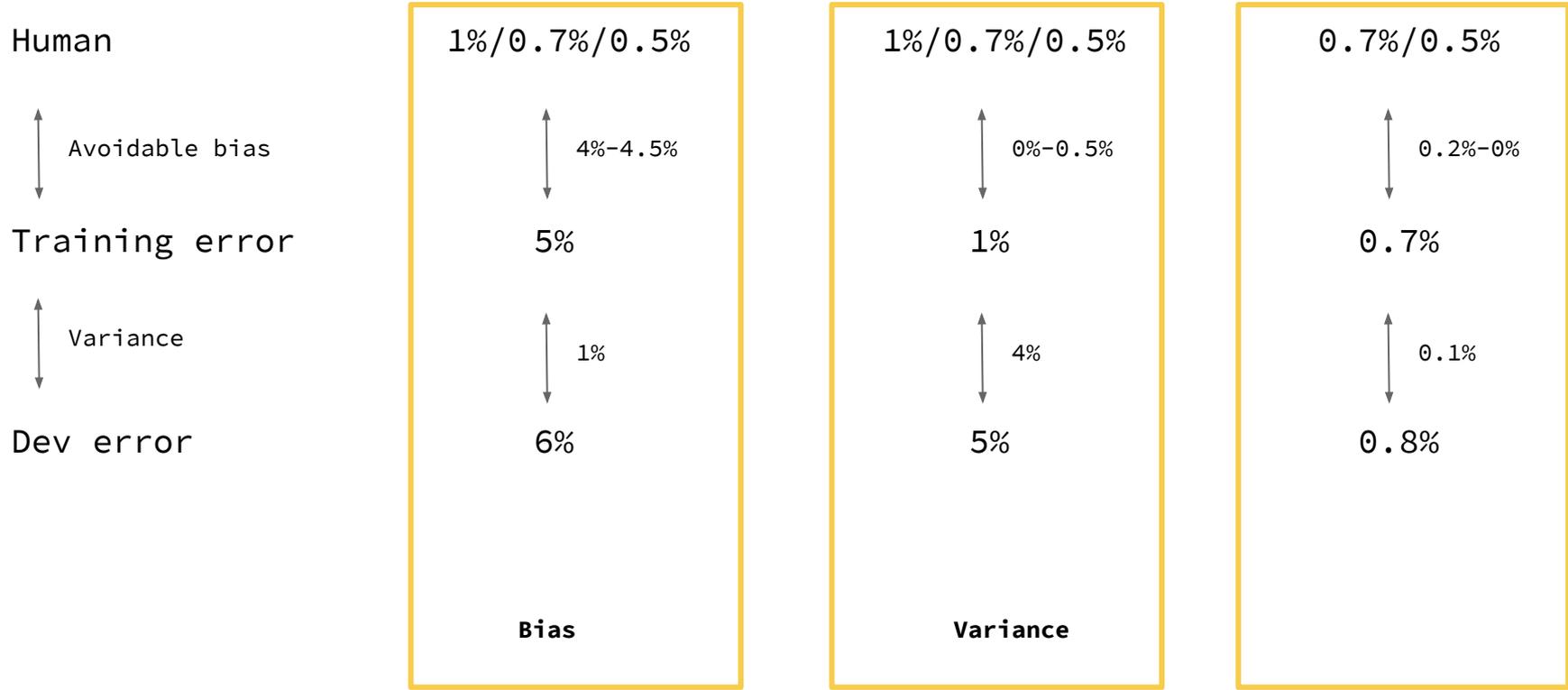
- | | |
|---------------------------------|------------|
| (a) Typical human | 3% error |
| (b) Typical doctor | 1% error |
| (c) Experienced doctor | 0.7% error |
| (d) Team of experienced doctors | 0.5 error |

What is the “human-level” error?

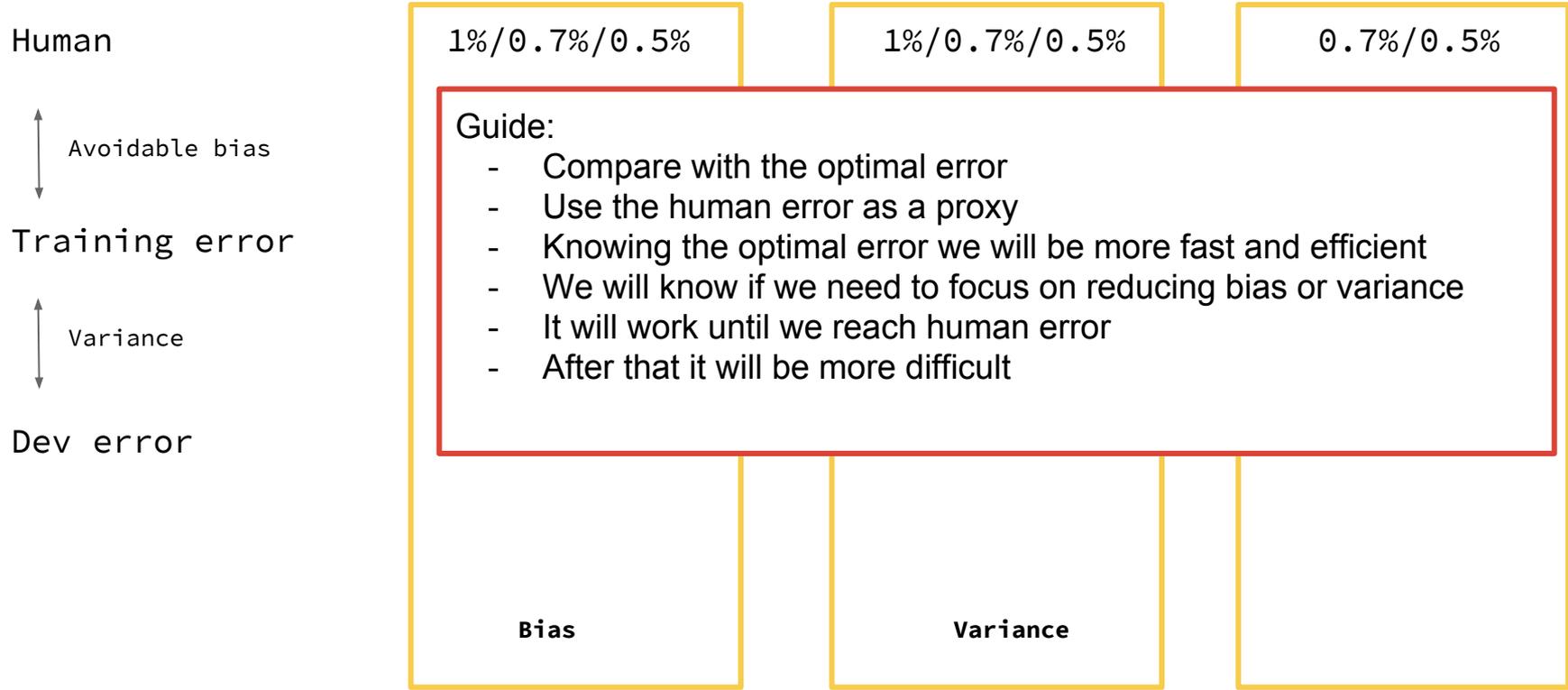
Bayes error $\leq 0.5\%$



UNDERSTANDING HUMAN - LEVEL PERFORMANCE (II)



UNDERSTANDING HUMAN - LEVEL PERFORMANCE (II)



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 (ECG, ...)

- Structured data
- Not natural perception
- Lots of data

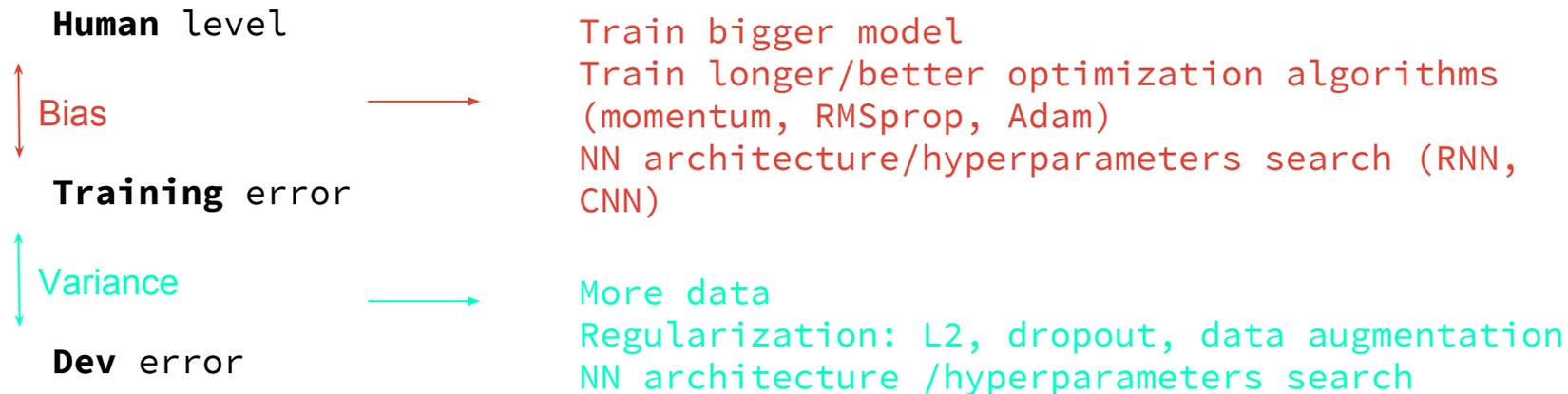
IMPROVING YOUR MODEL PERFORMANCE

THE TWO FUNDAMENTAL ASSUMPTIONS OF SUPERVISED LEARNING

- (1) You can fit the training set pretty well
- (2) The training set performance generalizes pretty well to the dev/test set



REDUCING (AVOIDABLE) BIAS & VARIANCE



Q&A

Anna Bosch Rue
VP Data Intelligence @ Launchmetrics
annaboschrue@gmail.com