

Orthogonalization

Focus is on parameters that are independent to expedite results

Fitting training set: Larger networks, better algorithm

Fitting dev set: Regularization, larger training set

Fitting test set: Bigger dev set

Real world performance: Change dev set or cost function

Precision: $\frac{\text{correct classifications}}{\text{positive classifications}}$, Recall: $\frac{\text{correct classifications}}{\text{positive examples in dataset}}$

Multi-class: Take avg. of error over each class

Maximize accuracy under constraint $t_{run} \leq 100$ ms

optimizing

Satisficing method

Ex. Wake words / Trigger words

"Alexa", "ok Google", "Ni hao, Baibu", etc.

Maximize accuracy, ≤ 1 false positive per 24 hrs.

Training / Dev / test Set Distributions

Regions

Bad

Good

US
UK
Europe
S. America
India
China
Other Asia
Australia

Dev

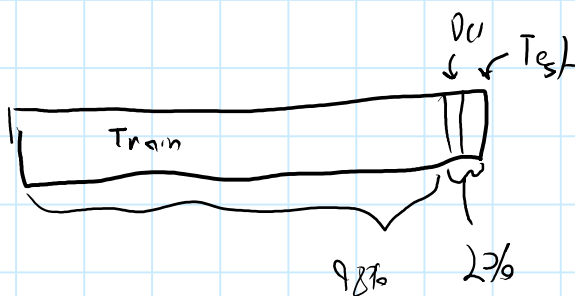
Test

Shuffle
&
divide

Dev & test sets should
come from same distribution

Size of Dev & Test sets

10^7 examples



When to Change Dev / Test set metrics

e.x, Cat image app porn delivery

metric: classification error

Algorithm A: 3% error, errors contain porn

Algorithm B: 5% error, no porn
MLU

Algorithm B: 5% error, no porn

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

problem: mislabels of porn should be **weighted** more harshly

$$w^{(i)} = \begin{cases} 1 & \text{if } x^{(i)} \neq \text{porn} \\ 10 & \text{if } x^{(i)} = \text{porn} \end{cases} \quad (\text{must label images as porn or not})$$

ex. user images

Al, A: 3% error, only ID's high quality images

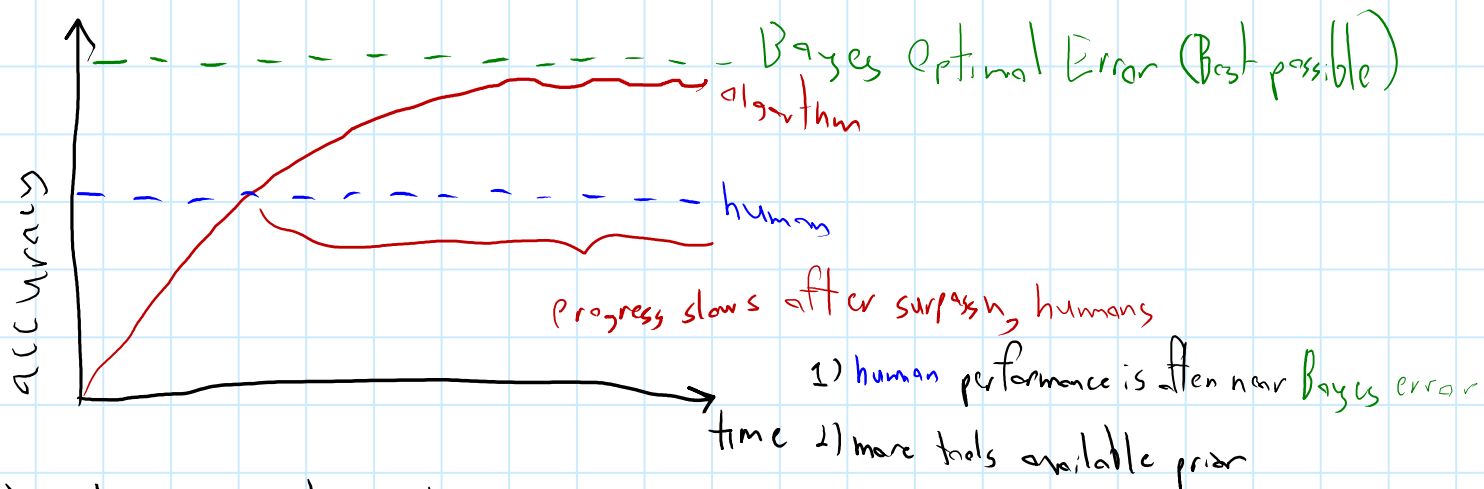
Al, B: 5% error, ID's user uploaded low quality images

Comparing to Human Performance

Sunday, November 19, 2017

11:23 AM

Why human level performance?



While worse than humans!

- Get human labeled data
- Manual error analysis
- Better analysis of bias & variance

Avoidable Bias

ex. cat classification

Human	1	7.5	human-level error is a proxy for Bayes error
Training error	8	8	
Dev error	10	10	
data	focus on bias	focus on variance	

"Avoidable bias"

Variance

Understanding Human-level Performance

ex. Medical image classification

	error(%)	
a) typical human	3	
b) typical doctor	1	← useful result, still deployable
c) experienced doctor	0.7	
d) team of c	0.5	← standard for human level (Borges proxy)

Human	1 / 0.5	0.5
Training	5	0.7
Dev	6	0.8
	reduce bias	reduce bias

Borges error is NOT always zero! (noisy data, etc.)

Surpassing human performance

hard to quantify avoidable bias when better than human

ML already does better with structured data such as ads & logistics where there are massive amounts of data

Improving model performance

2 assumptions

1) Fit training set well (Avoidable bias)

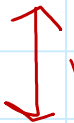
2) Training set generalizes well to dev/test (variance)

Human level



Avoidable bias

Training error



variance

Dev error

Train bigger model, train longer, better algorithms
NN architecture, hyperparameter tuning

More data, regularization (L2, dropout, data augmentation)
NN architecture, hyperparameter tuning

Week 2

Wednesday, November 22, 2017

4:30 PM

Carrying Out Error Analysis

Look at dev examples to eval ideas

(e.g. classification, learn dogs)

Get ~100 mislabeled dev set examples

Count how many are dogs

5/100, potential improvement is very small

50/100, half error, worthwhile effort

Evaluate multiple ideas in ||

Image	Idea 1	Idea 2	Idea 3	Comments
1				
2				
3				
⋮				
%	8	43	(61)	
			prioritize	

Cleaning up Mislabeled data (human labeled)

DL algorithms are fairly robust to random error if they're a small amount in training set

Add column to error analysis for mislabeled examples

Is it worthwhile to correct?

			A	B
Overall dev error	1.0	2	2.1	1.9
Incorrect labels	0.6	0.6	0.6	0.1
Other	0.4	1.4	1.5	1.8
	X	✓	higher ceiling, A can be better	

Build Fast, then Iterate

Use Bias/Variance + error to prioritize next steps

Exceptions: large amounts of academic lit on topic

Mismatched Training and Dev/test Set

Thursday, November 23, 2017

1:34 PM

Training & testing on different distributions

Cat app example

Data from web (200k)
professionally shot

Data from mobile app (10k)
more common

combine & shuffle

X Option 1: train 205k

dev/test 2.5k ea.

dev set still optimizes for web

✓ Option 2: train all 200k web + 5k user

dev/test rest of mobile app

testing app images

Speech Recognition ex. Car

Training (500k)
Purchased Data
Smart speaker
Voice keyboard
...

Dev/test (20k)
Speech activated carview mirror

train 500k
500k + 10k

dev/test 10k ea.
5k ea.

both acceptable

Bias & Variance w/ mismatched data distributions

human error ~ 0
 Train $\frac{1}{10}$
 Dev $\frac{1}{10}$

Are they the same distribution? If not, may be nothing wrong other than poor generalization due to data mismatch

Training - dev set: same dist. as training, but not used for training

Train	1	1	10	10
Train-dev	1	1.5	11	11
Dev	10	10	12	20

Annotations:
 - Between Train and Train-dev: variance $\swarrow \searrow$
 - Between Train-dev and Dev: data mismatch $\swarrow \searrow$
 - Above Dev: available bins \uparrow

General Principles

Human level	
Training Set error	\uparrow available bins
Training-dev set error	\downarrow variance
Dev error	\uparrow data mismatch
Test error	\uparrow degree of tuning to dev set (need more dev set data)

Addressing Data Mismatch

Manual error analysis to understand differences between training & dev / test sets

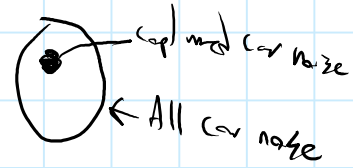
Make training data more similar (e.g. add car noise), or collect more data similar to dev/test sets

Data synthesis:
voice

looking

Car noise

I hear



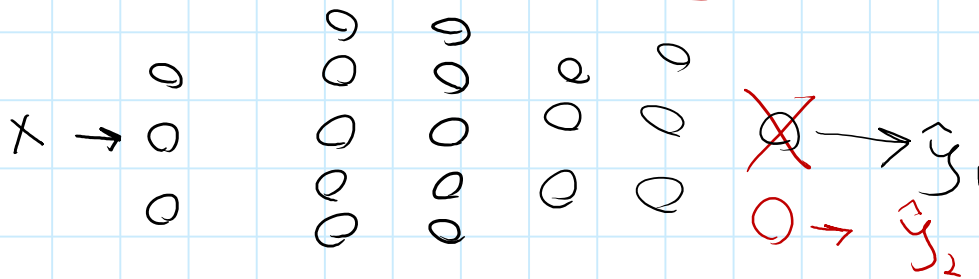
Same location, but not machines

Car recognition

Computer graphics? Game w/ 2D cars looks good, but not realistic representation

Transfer Learning

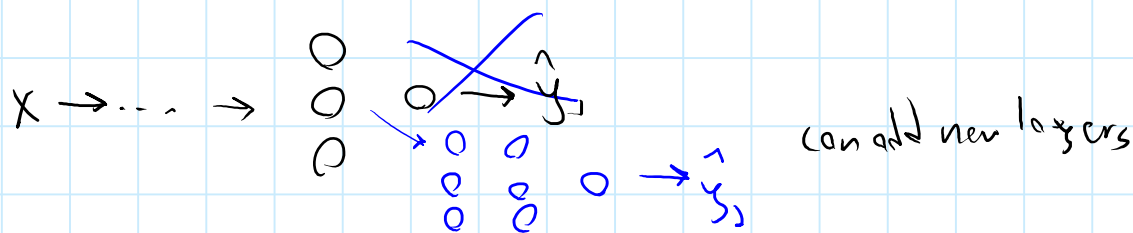
image recognition \rightarrow radiology change final weights



Retrains new (x, y) pairs

Little **data**: only retrain last layer

Another ex: speech recognition \rightarrow wake word



This tactic is useful when target area (like radiology) has fewer examples, but being a subset of images means transfer learning useful

abstract fine
10k 100 or 1
100 1000 or 100

Use transfer learning
May not hurt, but likely not very helpful

Summarized rules for when to use

Task A & B have same input

$m_A \gg m_B$

Low level features of A are helpful for B

Multi-Task Learning

e.g. self driving car images w/ multiple components present,
(car & stop sign)

$x \rightarrow \dots \rightarrow \begin{matrix} 0 \\ 0 \\ 0 \\ 0 \end{matrix} \rightarrow \hat{y} \begin{matrix} \text{person} \\ \text{car} \\ \text{sign} \\ \text{light} \end{matrix} \quad y \in \mathbb{R}^{4,1}$

loss: $\frac{1}{n} \sum_{i=1}^n \sum_{j=1}^4 \ell(\hat{y}_{ij}^{(i)}, y_{ij}^{(i)})$

\uparrow logistic loss

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

Unlike softmax, multiple labels are ok

Training 1 NN to ID multiple aspects can perform better than separate NNs provided underlying structures are the same (e.g. images)

Generalizes ok to missing data in y , (e.g. road images only marked w/ presence of a car) & just skips NaN values

1 1 1 >

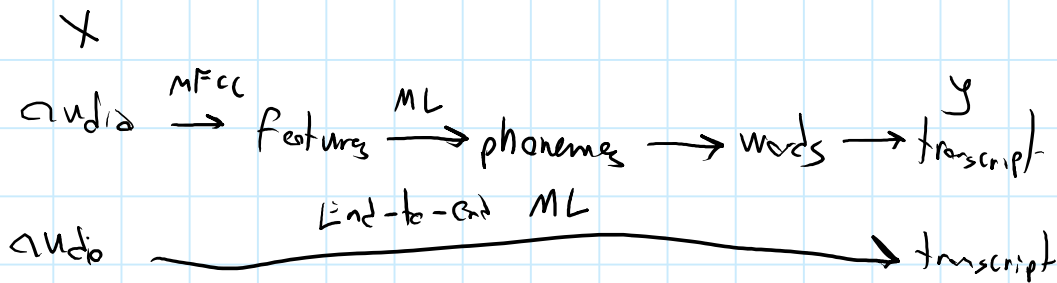
When to use?

Share low-level features
Amount of data for each task is similar
Can train large enough NN to do well on all tasks

End-to-End Deep Learning

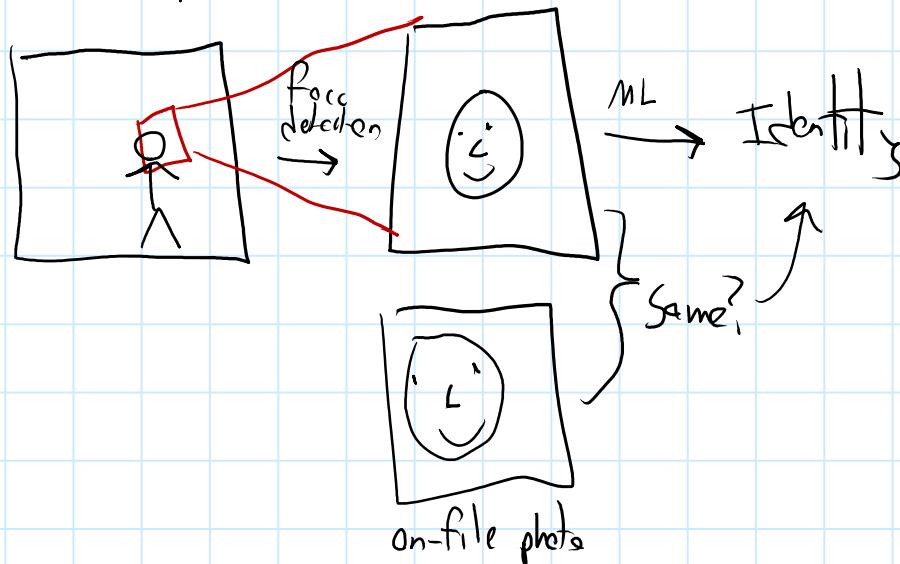
Saturday, November 25, 2017 6:24 PM

What is it?
Simplifies process



Need large amounts of data to be useful

facial rec, end-to-end not well suited



Machine translation

English → text analysis → → French

ML on
(X, y) values

Child's age based on hand X-ray

Image \rightarrow bones \rightarrow bone lengths \rightarrow table of
avg lengths for ages \rightarrow age

not enough pictures to train on

When to use end-to-end?

Pros:

- Lets data speak, no human preconceptions getting in the way of optimization
- Less hand-designing of components needed

Cons:

- Large amounts of data needed
- Excludes potentially useful hand designed components

Key question: Is there enough data for how complex your problem is?

end-to-end is not at all suited to self-driving cars