

dec1: Why ML strategy?

cat classifier - present system accuracy $\approx 90\%$.

Ideas to improve:

- collect more data
- collect more diverse data
- Training algorithm longer with GD
- Try Adam instead GD
- Try bigger network
- Try smaller network
- Try dropout
- add L_2 regularization
- change network architecture
 - activation function
 - hidden units.

Why one to try
for getting most
effective.

dec2: orthogonalization:

- each change of parameter should be controlled independently to output, which make tuning easier.

- Fit training set well on cost function (bigger network, adam) ^(not ovns) _(early stopping)
- Fit dev set well on cost function (more data, regularization) _(early stopping)
- Fit test set well on cost function (bigger dev set)
- performs well on real data (change dev set or cost function)

dec3: single Number evaluation metric:

classifier	Precision	Recall	F1 score	
A	95%	90%	92.4%	"select classifier A"
B	98%	85%	91.0%	

precision: of examples recognised as cats, what % actually cats?

recall: of all the images that are really cats, what % are correctly classified by classifier.

New evaluation metric

F1 score = average of precision and recall

$$F1 = \left[\frac{2}{\frac{1}{P} + \frac{1}{R}} \right] \quad \text{"Harmonic mean"}$$

Having well defined Devset + single real number evaluation metric
speed up iterate process in improving machine learning algorithm

example: 2:

Algorithm	US	china	India	other	Average
"select A"	3%	7%	5%	9%	6% (low error)
	5%	6%	5%	10%	6.5%

deci: satisficing and optimizing metric

classifier	Accuracy ^{optimizing}	Running time ^{satisficing}
A	90%	80ms
B	92%	95ms $\leq 100ms$
C	95%	150ms

"select B"

$$\text{cost} = \text{accuracy} - 0.5 \times \text{running time (not best)}$$

best:

maximize accuracy

subject to running time $\leq 100ms$

N metrics: 1 optimizing, N-1 satisficing

Example:

Wake word / Trigger word:-

Alexa, google, Hey Siri,

1. accuracy

2. no. of false positive

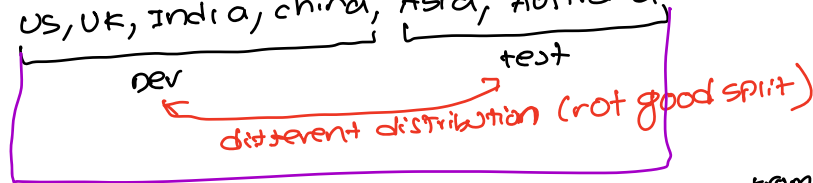
maximize accuracy | # false positive ≤ 1 every 24 hours.

lec 5: Train/Dev/Test Distributions:

cat classification dev / test sets
development set
validation set



Regions: US, UK, India, china, Asia, Australia



randomly shuffle (then dev, test are from same dist)

example:

optimizing on dev set on loan approvals for medium income zip code

$X \rightarrow y$ (renew loan?)

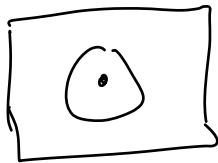
loan application

Tested on low income zip codes. (performance is really bad)

Guideline:

choose a dev set and test set to reflect data you expect to get in future and consider important to do well on.

dev set \approx test set distribution.



target = dev set + test set + metric

training set will help in hitting the target.

lec 6: size of Dev/ Test set:

old way:

train - 70%
test - 30%

train - 60%
dev - 20%
test - 20%

} data set
size are small
 $m \leq 10,000$

Bigdata era:

$m = 1,000,000$

train = 98%

dev = 1% 10,000

test = 1% 10,000

size of test set:

- big enough to give high confidence in overall predictor
 $m_{\text{test}} (10,000 \text{ to } 100,000)$

some application:

train + dev no test set.

dec7: When to change dev/test sets & Metric:

cat classifier

classifier

error

A

3%

letting pornographic images.

B

5%

doesn't have pornographic images.

evaluation metric + dev set preter classifier A

user preter classifier B.

⇒ change metric, dev set. Indicator

$$\text{previous error metric} = \frac{1}{m_{\text{dev}}} \sum_{i=1}^{m_{\text{dev}}} \mathbb{1} \{ y_{\text{pred}}^{(i)} \neq y^{(i)} \}$$

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

$$w^{(i)} = \begin{cases} 1 & \text{if } x^{(i)} \text{ is not porn} \\ 100 & \text{if } x^{(i)} \text{ is porn} \end{cases}$$

Example:

cat classifier

error dev/test

real world

A

3%

10%

B

5%

9% better.

dev/test images

- high res

- clear

real world

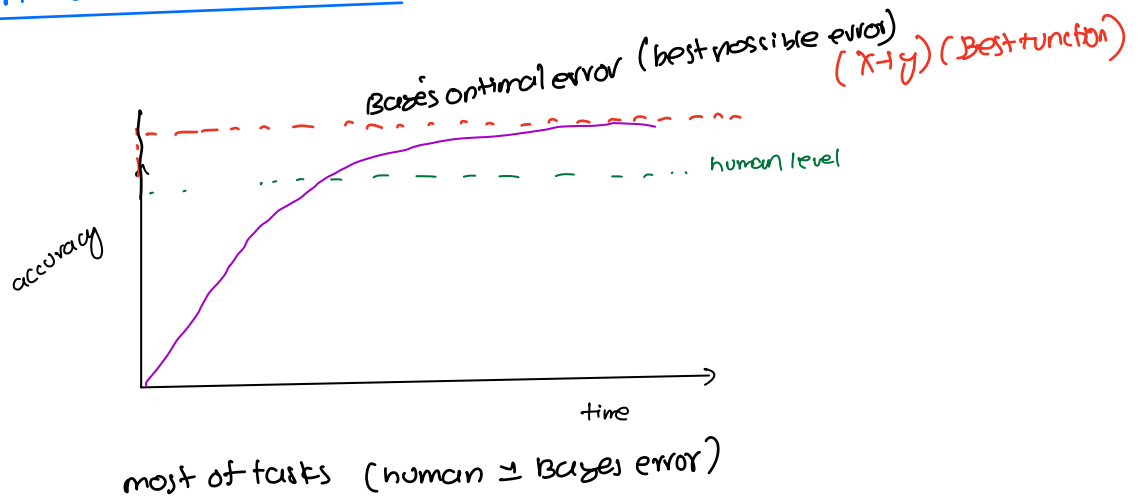
- low res

- blurry

it's doing well dev + metric doesn't correspond to real world data

⇒ change dev/test set + metric.

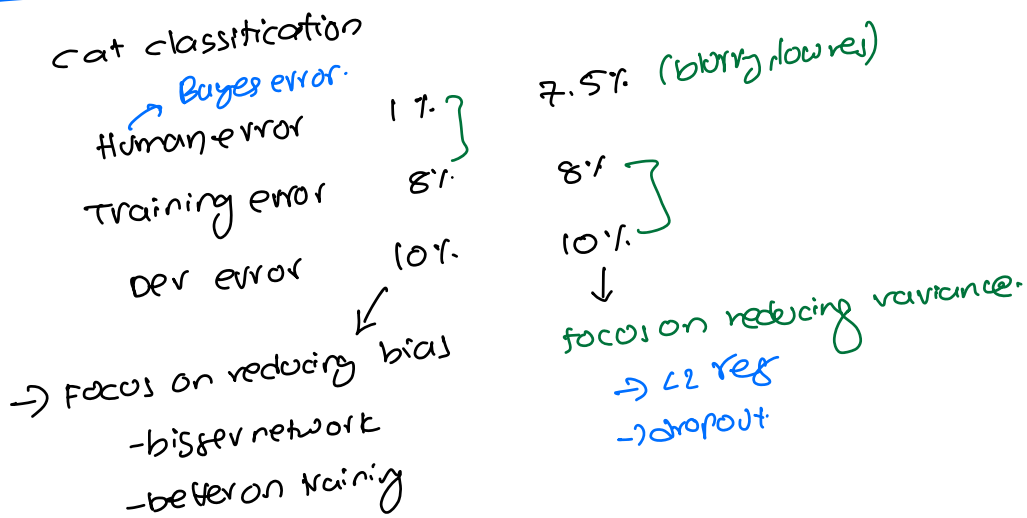
decs: Human level performance?



Human $>$ ML

- get label data from humans (x, y)
- gain insight from manual error analysis
- Better analysis of bias/variance.

deca: Avoidable Bias:-



for cv tasks Human \approx Bayes error

$$\text{Avoidable bias} = (\text{Bayes Error} - \text{Training Error})$$
$$\text{Variance} = (\text{Training} - \text{Dev Error})$$

lec 10: human level performance:

medical image classification example

(human level error proxy (Bayes Error))

Example:

(a) Typical human - 3% error

(b) Typical doctor - 1% error (human)

(c) Experienced doctor - 0.7% error

(d) team of experienced doctors - 0.5% error Bayes error.

Human level error = Bayes Error = 0.5% error.

Human / Bayes
↑ Error
↓ Avoidable bias
Training error
↓ variance
Dev error

same.
1%
0.7%
0.5% } 4 to 4.5%
5% } 1%
6%

problem is
to do better on
training data

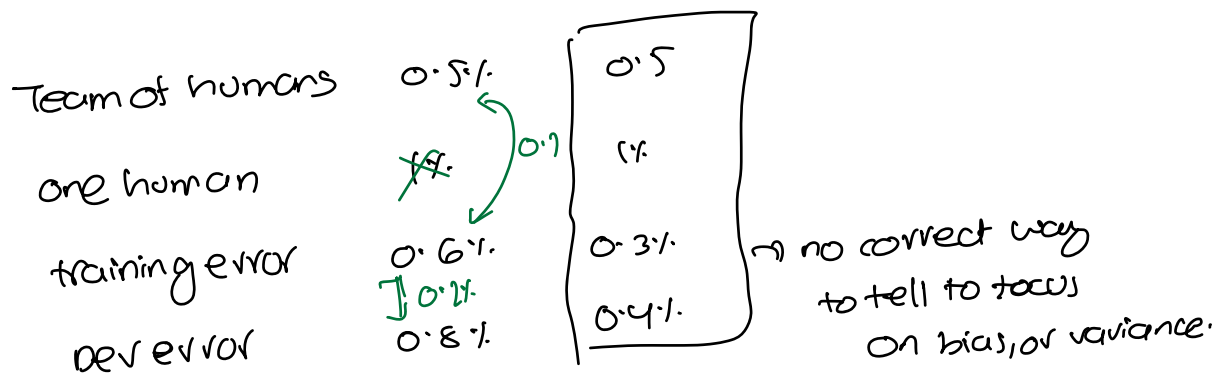
same
1%
0.7%
0.5% } 0.1%
0.5% } 0.5%
1% } 4%
5%

problem to
focus on
variance

Bayes error
0.5%
0.7% } 0.2%
0.8% } 0.1%

focus on
training data.

DECI: Surpassing human level performance:



avoidable bias = $0.6 - 0.5 = 0.1$

→ slows down efficiency.

ML > human

- online advertising
 - product recommendations
 - logistics
 - loan approvals
- } structural data.

single human {

- some audio
- some vision
- some medical (ECG, skin cancer)

lec12: Improve the performance

- ① Fit training set very well
- ② training set performance generalizes well to dev/test sets.

