

Course-3 Structuring ML Projects

Week-1 ML project might not work well in such a case,
we have following options

Collect more data

" diverse training set

Train algorithm longer w/ gradient descent

Try Adam instead of "

Try bigger n/w

" smaller "

" dropout

Add L₂ regularization

Network architecture

Activation fn

#hidden units

In this course, strategies will be taught on
how to analyze issues of a n/w

The process to tune, in order to achieve one effect,
is known as orthogonalization.

Using a single Number Evaluation metric

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

Precision = Of samples recognized as cats, what % actually are cats?

Recall = What % of actual cats are correctly recognized

There is certain trade-off going on between precision
and recall, thus, we take a single evaluation term
called F1 score. It is harmonic mean of Precision & Recall

$$F1\text{ score} = \frac{2}{\frac{1}{P} + \frac{1}{R}}$$

Better classifier \rightarrow ↑ F1 score

Similarly Cats APP for cat bavers in 4 geographies. Diff errors by classifiers	Algo	US	China	India	Other	Average
A		3.1.	7	5	9.1.	6.7
B		5.1.	6	5	10.1	6.5.1.
C		2.1.	3.8	4	5.1	(3.5.1)
D		5.1.	8	7	2.1.	5.25.1.
E		4.1.	5	2	4.1	3.75.1.
F		7.1.	11	8	12.1.	9.5.1.

We choose the one whose avg. error is least hence instead of evaluating the whole grid, we just need evaluate a single number matrix.

It is not always easy to make a simple single matrix each time.

Classifier	Accuracy	Running Time
A	90.1.	80ms
B	92.1.	95ms
C	95.1.	1.50ms

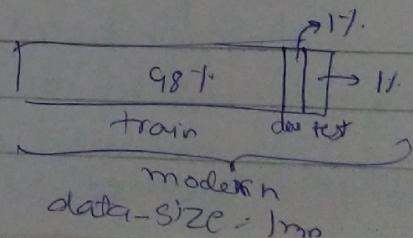
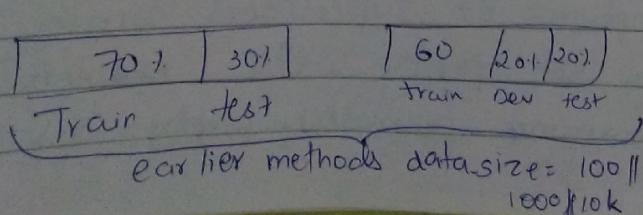
Instead of using some random formula to select the best classifier, what we do is

maximize accuracy

suffice to running Time $\leq 100\text{ms}$

In such a case, accuracy becomes optimizing matrix, and Running Time becomes satisfying matrix

\rightarrow Dev and test set should come from same distribution.
But, how large should they be.

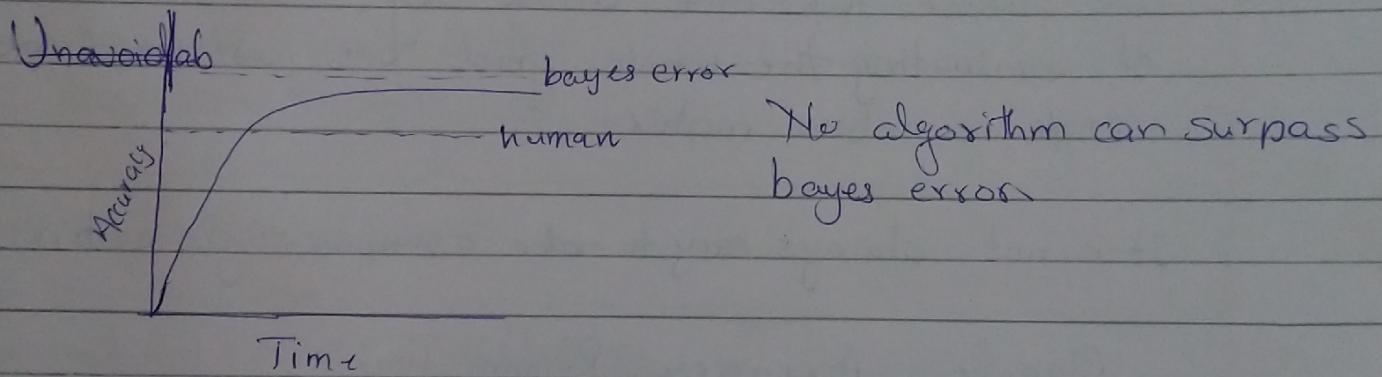


Many a times, it may happen that the classifier has less performance but does quite well in real life scenarios.

For e.g. cat images with high definitⁿ are used to train
In that case, ~~B~~ classifier A \rightarrow 3% error
classifier B \rightarrow 5%.

But it turns out that when low quality images are used, B works better.

Hence B is chosen



For ~~a lot~~ classifiers which identify cats:

	A	A	
Training Humans	1%	7.5%	Avoidable Bias
Training error	8%	8%	Variance
Dev error	10%	10%	

Focus on Bias Focus on variance

We consider human level error as a proxy for Bayes error

{My thoughts: Focus the one whose error is less.}

Two functional assumptions of supervised learning

1. You can fit the training set pretty well \rightarrow ~~Avoidable bias~~ ~ Avoidable bias
2. The training set performance generalizes pretty well to the dev/test set \rightarrow Variance

Reducing avoidable bias

Train bigger model

Train longer/better optimization algos (e.g. RMSprop, momentum, Adam)
NN archi./hyperparameter search

↳ RNN, CNN (model change)

#hidden units

activation change

Reducing Variance

More data

Regularization (e.g. L2, dropouts, data augmentation)

NN architecture / hyperparameters search

(some as above)