

Couirse 3: Structuring Machine Learning Projects.

start date

15-01-19

Week 1: ML Strategy (1).

P1: Introduction to ML Strategy.

L1: why ML strategy. [10.30 am 15-01-19]

Some ideas of machine learning. (deep learning)

- 1. Collect more data.
- 2. Collect more diverse training set.
- 3. Train algorithm longer.
- 4. Try adam instead of Grd.
- 5. Try bigger network
- 6. Try smaller network.
- 7. Try dropout.
- 8. Add L2 regularization.
- 9. Change network architecture.
 - # activation function.
 - # hidden units.

[we will learn how to effectively use these ideas to solve our problem.]

L2: Orthogonalization

→ Orthogonalization is a system design property that ensures that modifying an instruction or a component of an algorithm will not create side effects to other components of the system.

when a supervised learning system is designed these 4 assumptions are needed to be true and orthogonal.

1. Fit training set well on cost function.

[if don't fit well, (i) use bigger neural network
(ii) use better optimization algorithm.]

2. Fit development set well on cost function.
→ if doesn't fit well, regularization or using bigger network might help.
3. Fit test set well on cost function.
→ [the use of bigger development set might help]
4. Performs well in real world.

P2: Setting up your goal.

L1: Single number evaluation metric. [8.10 pm 15-01-18]
[Precision, Recall, $\overbrace{F_1\text{-Score}}$. average of precision and recall.]

→ we have to chose a single number evaluation metric by which we can correctly separate results of our classifiers.

[such as using $F_1\text{-Score}$ as evaluation metric which averages precision and recall].

** To chose a classifier from a set of classifiers,

i) a well defined development set

ii) a single number evaluation metric

are useful to speed up iteration process.

L2 :- Satisfying and optimizing metric. [8.20 pm 15-08-19]

optimizing metrics

[The evaluation metric that we have to optimize]

[Accuracy, Cost minimization etc]

Evaluation Metrics. [used to evaluate performance of a classifier].

satisfying metrics.

[ok with some evaluation metric has satisfactory results]

[Running time, etc.]

Classifier	Accuracy	running time.
A	90%	80 ms
B	92%	90 ms
C	95%	1500 ms

optimizing metric.

satisfying metric.

$$N_{\text{metric}} = \begin{cases} 1 & \text{optimizing metric} \\ N_{\text{metric}} - 1 & \text{satisfying metric} \end{cases}$$

L3 :- Train/dev/Test distributions.

** choose the development and test sets from the same distribution. It must be taken randomly from all data.

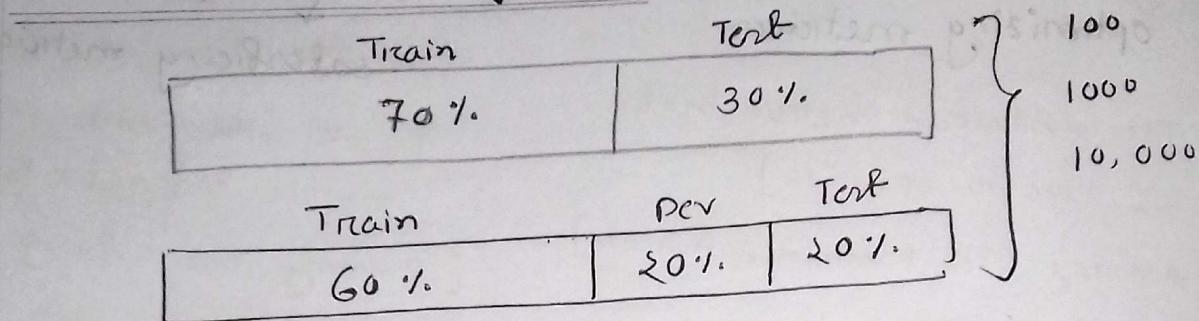
→ [choose some data once to predict one at a time for loan pay or not].

Guideline:-

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

L4 : size of the dev and test sets. [8.50 pm 15.01.19]

[From large amount of data we can use small dev and test set]
old way of splitting data.



- ** Set your test set to be enough to give high confidence in the overall performance of your system.
- ** Test set helps evaluate the performance of the final classifier which could be less 30% of the dataset.
- ** The development set has to be big enough to evaluate different ideas.

** modern era - Big data.
Because a large amount of data is available, we don't have to compromise as much and can use a greater portion to train the model.

LS :- when to change dev/test sets
and metrics.
(Ans)

[9:05 pm 15.01.19]

The misclassification error metric.

$$\text{Error} : \frac{1}{m_{\text{dev}}} \sum_{i=1}^{m_{\text{dev}}} \{ h \neq \hat{y}^{(i)} \neq y^{(i)} \}$$

This function counts up the number of misclassified examples.

→ If evaluation metric fails to correctly rank preferences between algorithms. The evaluation metric can be changed by adding up the weight term.

$$w^{(i)} = \begin{cases} 1 \\ 10 \end{cases} \quad \text{Adding different weight to rank preferences correctly.}$$

→ Final function :-

$$\text{or Error} : \frac{1}{\sum w^{(i)}} \sum_{i=1}^{m_{\text{dev}}} w^{(i)} \{ h \neq \hat{y}^{(i)} \neq y^{(i)} \}$$

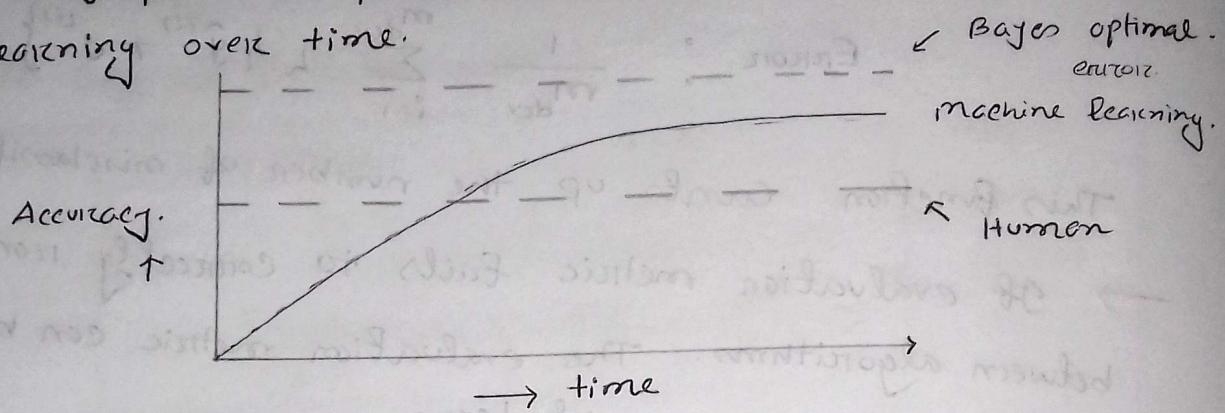
Guideline :-

1. Define correctly an evaluation metric that helps better rank order classifications.
2. Optimize evaluation metric.

P3: Comparing to human-level performance.

L1: why human level performance. [11:25 pm 15.01.19]

This graph shows performance of humans and machine learning over time.



Bayes optimal error is defined as the best possible error.

Any function mapping from X to J can't surpass this level of accuracy. (Algorithm that always correctly predicts)

** Machine learning programs surpass human -level performance.

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

- Get labeled data from humans.
- Gain insight from manual error analysis: why did a person get this right?
- Better analysis of Bias and variance.

L2 :- Avoidable Bias. [11.40 PM 15-01-19].

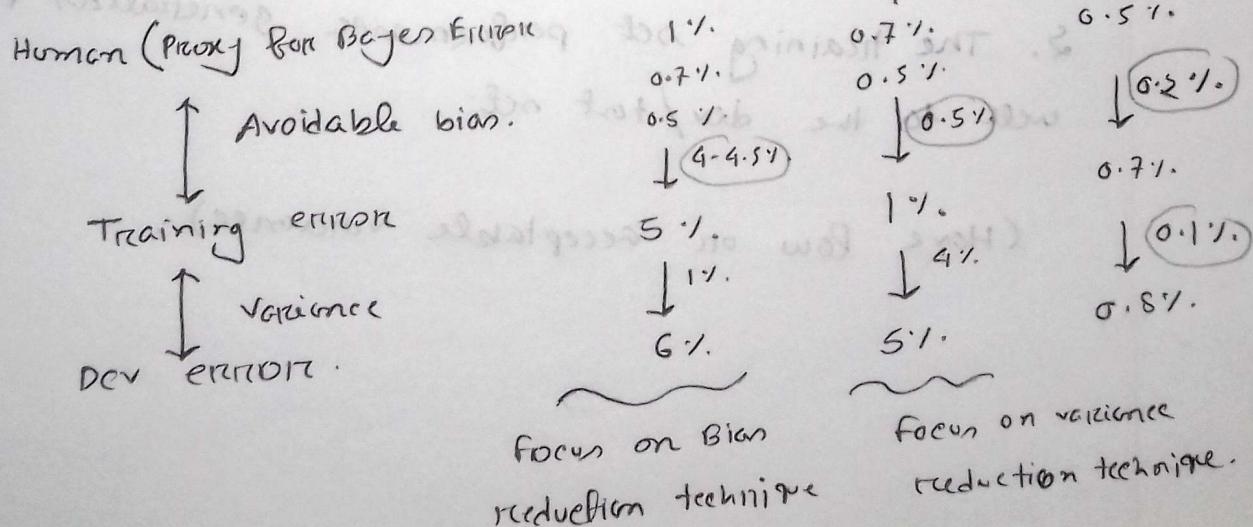
	A	B	→ ** The difference between training set and human level error is called <u>avoidable bias</u> . Try to reduce it.
Human	1%	7.5%	
Training error	8%	8%	
Dev error	10%	10%	

- (A) A gap between human error and training error is high use bias reduction technique. [Train a bigger neural network, running the training set longer]
- (B) A difference between human error and training error is low use variance reduction technique. [Regularization, Have a big training set]

L3 :- Understanding Human-level performance. [10.15 AM 16-01-19].

Human level error gives an estimate of Bayes error. The definition of human-level error depends on the purpose of the analysis. [If my machine can work very good as human of the analysis. Then it is really good.]

Error Analysis Example :-



L-4 Surpassing Human-level performance. [10:30 AM 16-01-19]

→ when training error is less than human level - errors or avoidable bias is (≤ 0) then we can call our machine has surpassed human level performance.

→ There are many problems where machine learning significantly surpass human level performance, especially with structured data.

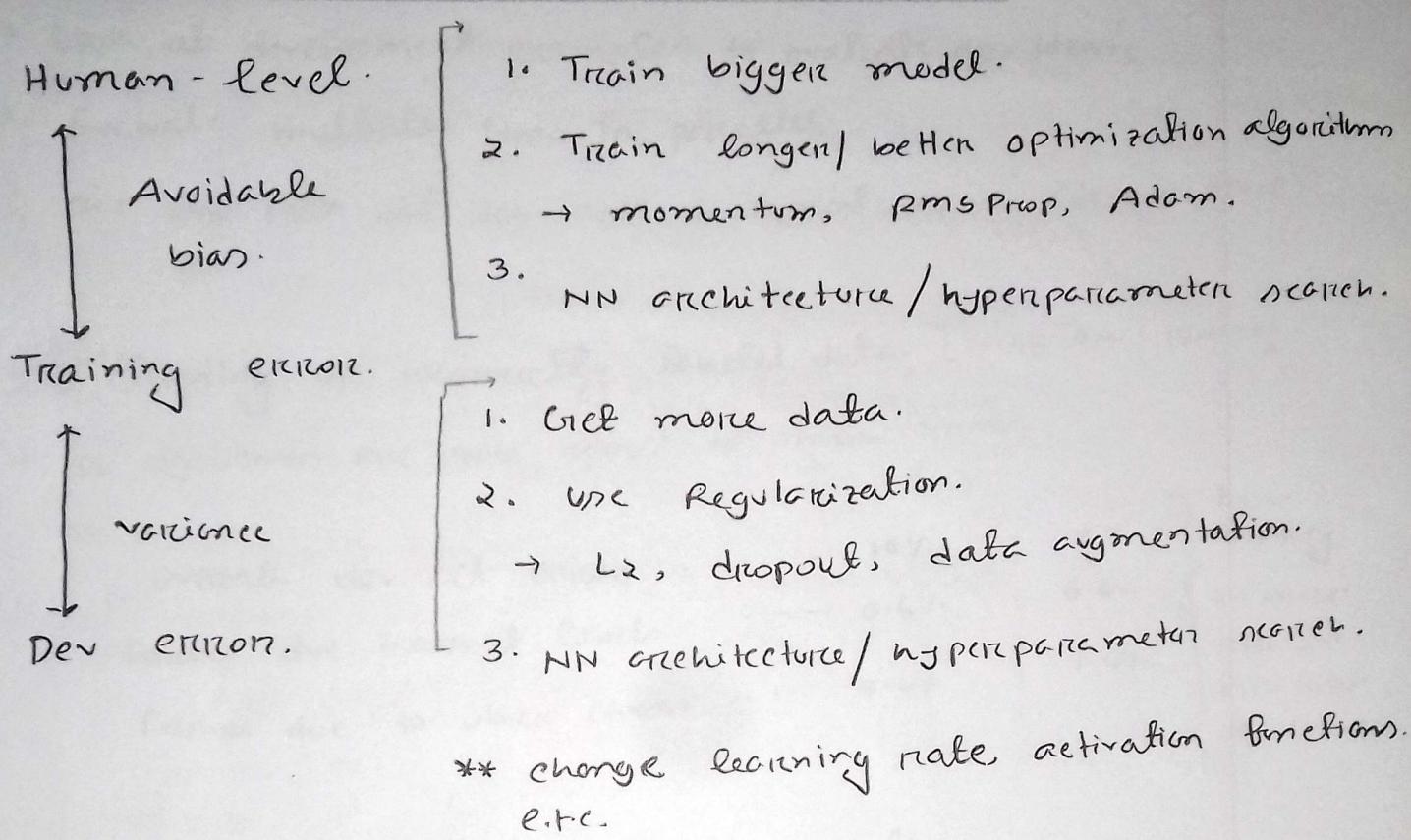
- online advertising. ** Structure data.
 - Product recommendations. ** Un-structure data.
 - Logistics (predict ticomit time).
 - Loan approvals. etc.

[10.40 AM 16.01.19] E5 Improving Your model performance

Two fundamental assumptions of supervised learning.

1. You can fit training set pretty well.
(Or have low avoidable bias)
 2. The training set performance generalizes pretty well to the dev/test set.
(Have low or acceptable variance)

Reducing (available) bias and variance.



quiz :-

Week 2: ML Strategy (2)

P1: Error Analysis

(Implementation in RnD Group)

L-1: carrying out error analysis.

(Theoretical part)

[10.30 AM 19.01.19]

- Look at development examples to evaluate new ideas.
- Evaluate multiple ideas in parallel.

(Unit idea costs one idea requires implementation with 5% cost & 20% E)

L-2: cleaning up incorrectly labeled data.

[10.45 AM 19.01.19]

- DL algorithms are quite robust to random errors.

	Overall dev set errors	Test set errors	Errors normally
Errors due to incorrect labels	10%	2%	20% more
Errors due to other causes	0.6%	0.6%	rather 20%
	9.4%	1.4%	further 10% into 20%

Correcting incorrect dev/test set examples:

1. Apply some process to your dev and test sets to make sure they come from same distribution.
2. Consider examining examples your algorithm got right as well as ones it got wrong.
3. Train and dev/test data may now come from different distributions.

L3: Build your first system quickly, then iterate. [12.00 AM 19.01.19]

Build your first system quickly,

1. Set up development/test set and metrics.

- set up a target.

2. Build an initial system quickly.

- Train training set quickly : fit the parameters.

- Development set : Tune the parameters.

- Test set : Assess the performance.

3. Use Bias/variance analysis and error analysis to prioritize next steps.

P2: mismatched training and dev/test set.

L1: Training and Testing on different distributions. [1.15 PM 19.01.19]

→ we have to choose a development set and test set to reflect data you expect to get in the future and consider important to do well.

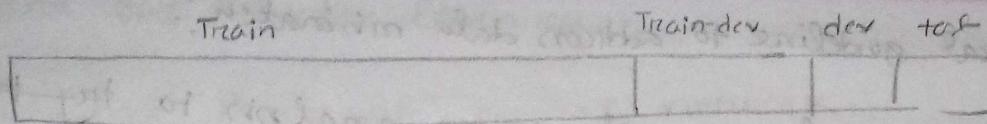
Training set test/dev set

similar

different

L2: Bias and variance with mismatched data distribution.

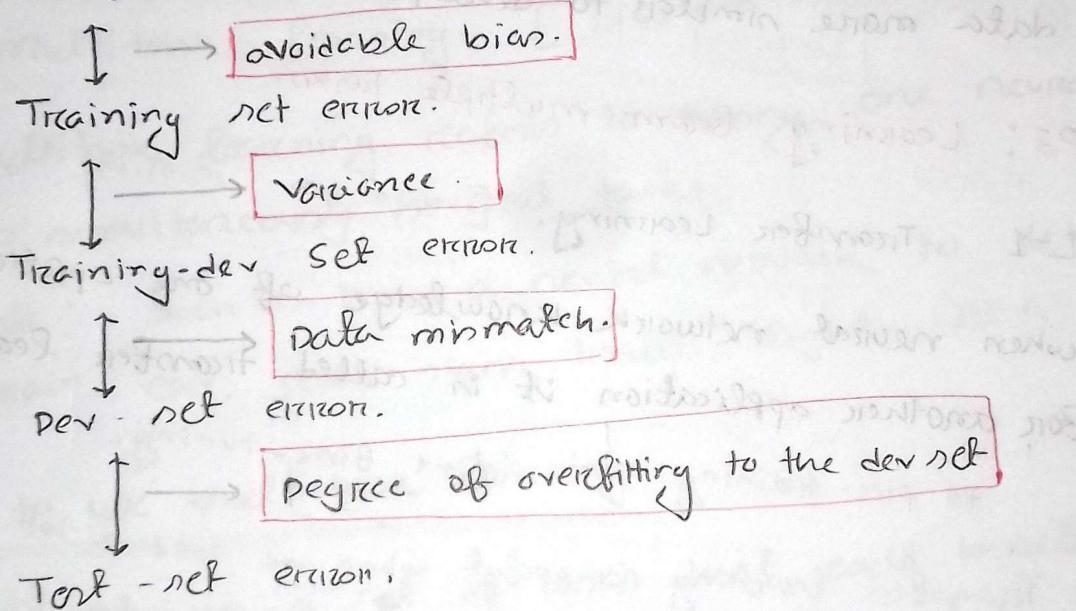
[4:20 PM 19.01.10]



- ** Train-dev net and train set comes from same distribution.
- Train-dev set is used find whether there is Bias, variance or data mismatch.

General formulation of error analysis.

Human level (Bayes errors)



changes

we learn how to analyze bias and variance efficiently when training set is from a different distribution than the development and test sets.

[we have to take decisions by keeping in mind different error analysis techniques.

L-3 Addressing Data mismatch. [7.45 pm 19.01.19]

General guideline to address data mismatch:-

- carry out manual error analysis to ~~try to~~ understand the error differences between training, dev/test sets.
- make training data or collect data similar to development and test sets.
- * Artificial data synthesis can be used to make training data more similar to development set.

P3: Learning from multiple tasks. [11.20 pm 19.01.19]

L-1 Transfer Learning. [11.20 pm 19.01.19]

when neural network knowledge of one application used for another application it is called transfer learning.

* pre-training
* fine-tuning

- Transfer learning is useful when we have a lot of data for the problem we transfer from and less data for the problem we transfer to:

when to use transfer learning. [Transfer $A \rightarrow B$].

1. Task A and B have the same input x .

2. You have a lot more data for Task A than Task B.

3. Low level features from task A could be helpful for task B.

Guideline for transfer learning.

- delete last layer of neural network.
- delete weights feeding into the last output layer of the neural network.
- Create a new set of randomly initialized weights for the last layer only.
- New data (x, y).

L-2 multitask learning.

multitask learning refers to having one neural network to do simultaneously several tasks.

Example :- Such as train a neural network to detect pedestrian, car, road sign, traffic lights etc.

when to use multi-task learning.

1. Training on a set of tasks that could benefit from having shared lower-level features.
2. usually amount of data for each task is quite similar.
3. Can train a big enough neural network to do well on all the tasks.

[Transfer learning use more often than transfer learning]

P4 : End to end deep learning.

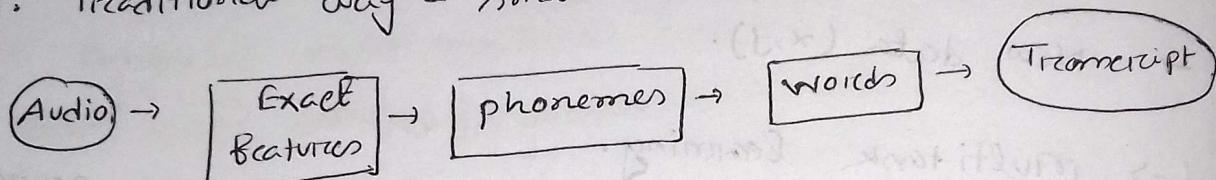
L-1: what is end to end deep learning? [1.20 AM 20.01.19]

End to end deep learning is the simplification of a ~~process~~ learning systems into one neural network.

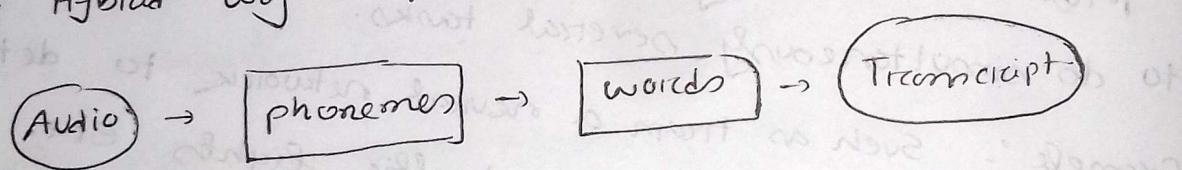
** End to end deep learning can not be used for every problem because it needs a lot of labeled data.

Example:- speech recognition model.

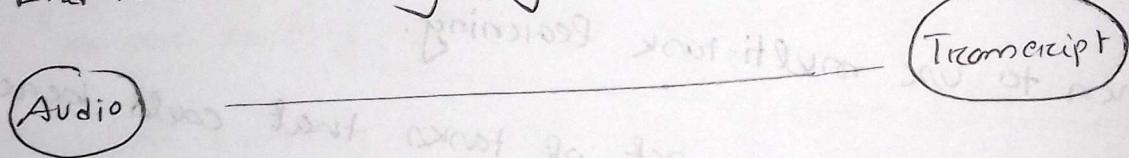
1. Traditional way - small data net



2. Hybrid way - medium data net



3. End to end learning way - large data net



Applications:-

1. Audio transcript / speech recognition.

2. image capture.

3. image synthesis.

4. machine translation.

5. face recognition.

6. self-driving car etc.

L2:- whether to use end-to-end
deep learning.

[11:15 AM 20-01-19]

If one has enough data to learn a function of then he can use end-to-end deep learning.

Pros:-

1. Let data speak :- machine will able to find which statistics are in the data, rather than being forced to reflect human preconception.
2. Less hand-designing of components needed :- gt simplification design work flow.

Cons:-

1. Require large amount of labeled data.
→ gt can not be used for every problem as it needs a lot of labeled data.
2. Excludes potentially useful hand-designed components.