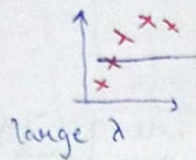
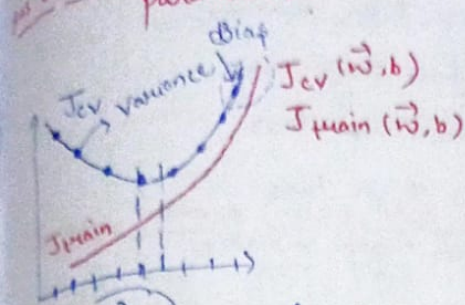




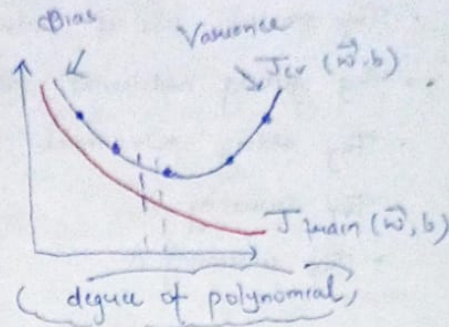
Week_2 Machine learning

Choosing Model

Plot and variance as a function of regularization parameters?



Comparison



Speech recognition example

Human level performance: 10.6%
 Training error J_{train} : 10.8%
 Cross validation error J_{cv} : 14.8%
 0.2% that's not bad
 4% it's high

Establishing a baseline level of performance

What is the level of error you can reasonably hope to get to?

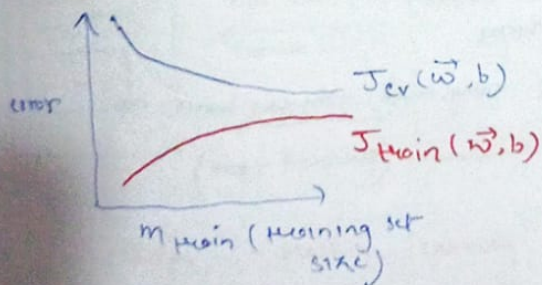
- Human level performance
- Competing algorithms performance
- Guess based on experience.

Learning Curves

$$f_{\vec{w}, b}(x) = w_1 x + w_2 x^2 + b$$

J_{train} = training error

J_{cv} = cross validation error



Conclusion

If your problem is high Bias
 Then giving more training data does not show much effect.

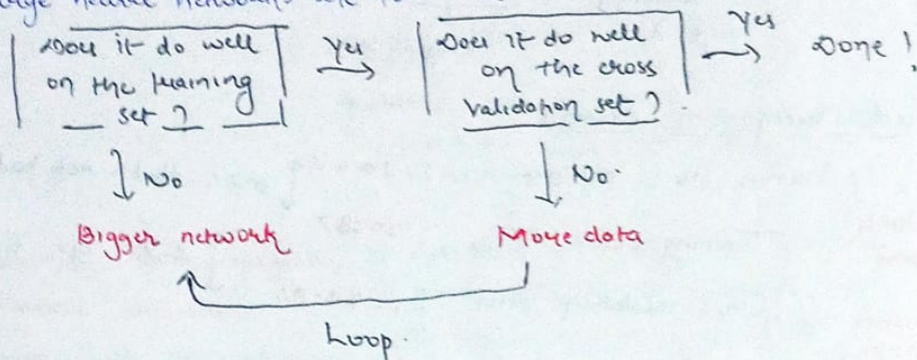
In case of high variance then increasing dataset size helps a lot.

Q → what do you try next?

- Get more training examples → fixes high variance
- Try smaller sets of features → fixes high variance
- Try getting additional features → fixes high bias
- Try adding polynomial features → fixes high bias
- Try decreasing λ → fixes high bias
- Try increasing λ → fixes high variance

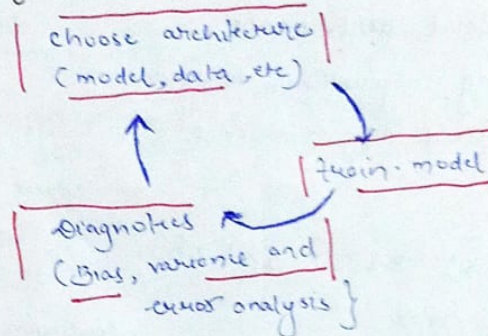
• Neural networks and bias variance

Large neural networks are low bias machines



• ML Development process

⇒ Iterative loop of ML development



• Error analysis:

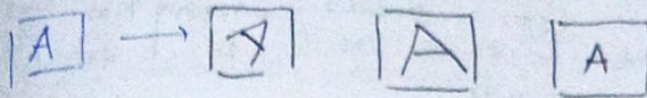
$m_{cv} = 500$ examples in cross validation set

Algorithm misclassifies 100 of them

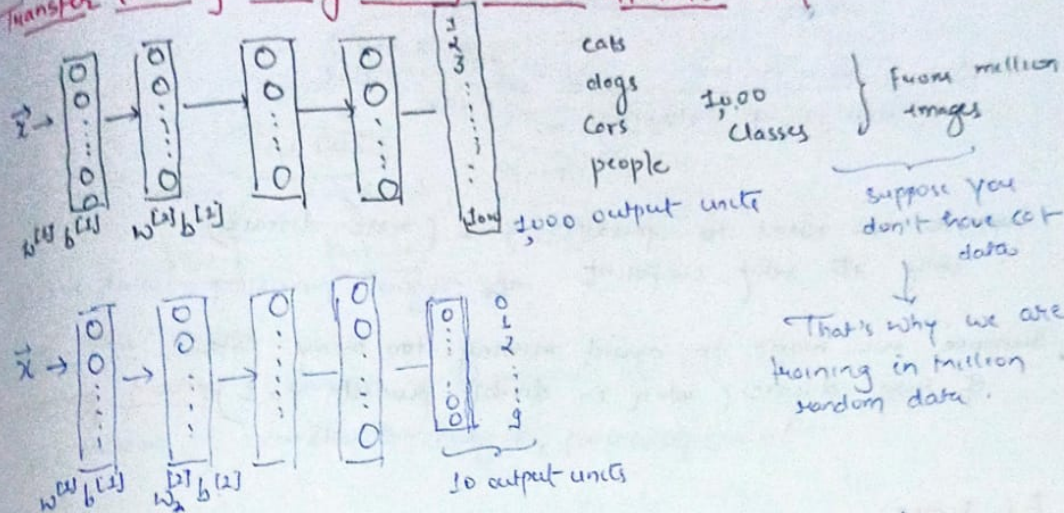
Manually examine 100 ex and categorize them based on common traits.

data augmentation:

Augmentation: modifying an existing training example to create a new training examples.

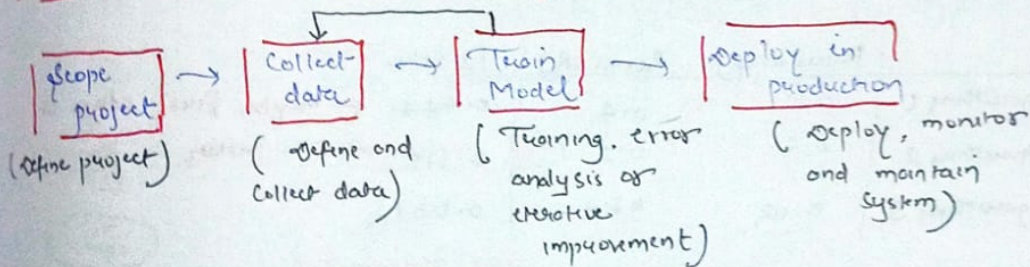


Transfer learning: using data from a different task

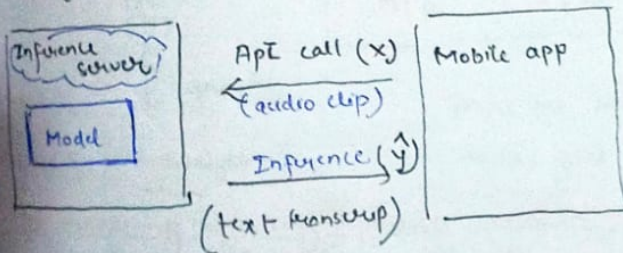


option 1: only train output layer parameters
option 2: train all parameters.

full cycle of a machine learning project



Deployment:



MLOps :- Machine learning operations.

• Precision / Recall

Precision: $\frac{\text{True positive}}{\text{predicted positive}}$

$$\Rightarrow \frac{\text{True positive}}{\text{True pos + false pos}} = \frac{15}{15+5} = 0.75$$

ex. $Y=1$ in presence of rare class we want to detect

	Actual	
	1	0
predicted class	T.P 15	F.P 5
	P. negative 10	T.N 70

Recall: $\frac{\text{True positive}}{\text{actual positive}}$

$$\Rightarrow \frac{\text{True positive}}{\text{True pos + false neg}} = \frac{15}{15+10} = 0.6$$

Q. Suppose we want to predict $Y=1$ (rare disease) only if very confident \rightarrow higher precision, lower recall

Suppose we want to avoid missing too many case of rare disease (when in doubt predict $Y=1$)
lower precision, higher recall

F1 score:

How to compare precision / recall numbers.

$$F_1 \text{ score} = \frac{1}{\frac{1}{2} \left(\frac{1}{P} + \frac{1}{R} \right)} \Rightarrow \frac{2PR}{P+R} \rightarrow \text{Harmonic Mean}$$

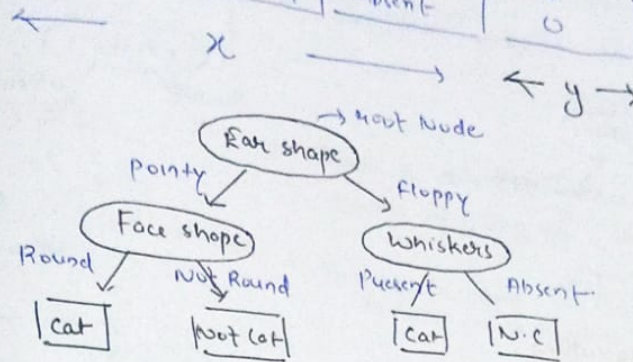
	Precision (P)	Recall (R)	F1 score
Algorithm 1	0.5	0.4	0.444
Algorithm 2	0.7	0.1	0.175
Algorithm 3	0.02	0.0	0.0392

\leftarrow maybe first algo is better.

Decision Tree Model

cat classification example

Ear shape	Face shape	Whiskers	Cat
Pointy	Round	Present	1
Pointy	Not Round	Absent	0
Floppy			



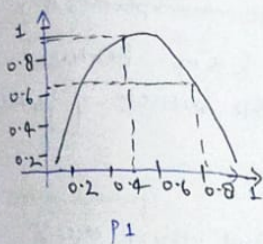
Decision Tree Learning

Measuring Purity: Entropy as a measure of impurity.

P_1 = fraction of examples that are cat $P_1 = 3/6$ $H(P_1) = 1$

$$P_1 = 2/6 \quad H(P_1) = 0.92$$

$$P_1 = 5/6 \quad H(P_1) = 0.65$$



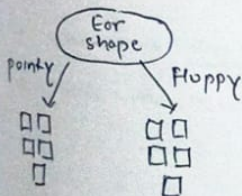
$$P_0 = 1 - P_1$$

$$H(P_1) = -P_1 \log_2(P_1) - P_0 \log_2(P_0)$$

$$= -P_1 \log_2(P_1) - (1 - P_1) \log_2(1 - P_1)$$

Note " $0 \log(0)$ " = 0

choosing a split: Information gain



$$P_{\text{left}} = 4/5 \quad P_{\text{right}} = 1/5$$

$$w_{\text{left}} = 5/10 \quad w_{\text{right}} = 1/10$$

Information gain

$$= H(P_1^{\text{root}}) - (w_{\text{left}} H(P_1^{\text{left}}) + w_{\text{right}} H(P_1^{\text{right}}))$$

Measure the reduction in entropy
more good it is to split feature.

Keep repeating splitting process until stopping criteria is met.

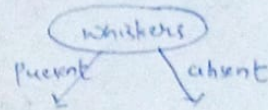
- when a node is 100% one class
- when splitting a node will result in the tree exceeding a maximum depth.
- Information gain from additional split is less than threshold
- when number of examples in a node is below a threshold.

- Using one-hot encoding of categorical features.

If a categorical features can take on k values, create k binary features (0 to 1 valued)

- Tree ensembles.

At Trees are highly sensitive to small changes of the data



- Collection of Trees

- Random Forest Algorithm

Given training set of size m

for $b = 1$ to B :

use sampling with replacement to create a new training set of size m Train a decision tree on the new dataset

Randomizing the feature choice

At each node, when choosing a feature to use to split. if n features are available, pick a random subset of $k \leq n$ features and allow the algorithm to only choose from the subset of features

$$k = \sqrt{n}$$

XGBoost:

we will focus on subset not doing properly. And then create a new decision tree boosting that particular part to perform well.

{ extreme gradient Boosting }