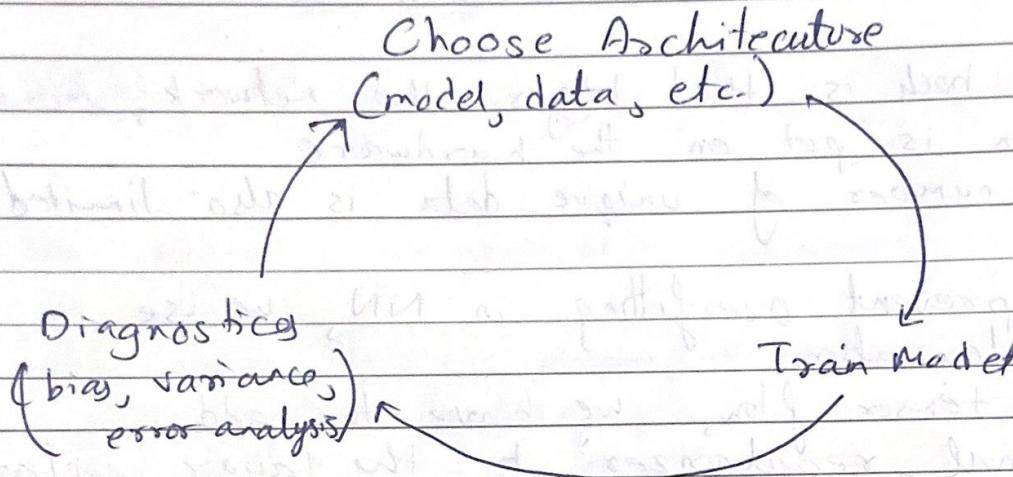


Iterative Loop of ML Development



Error Analysis

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

Algorithm misclassifies 100 of them

Mmanually examining 100 examples and categorize them based on common traits

Ex. if we have a spam mail catcher and get following errors.

Note: Some errors might be overlapping

Pharma: 21

Deliberate misspelling: 3

Unusual routing: 7

Steal Passwords (phishing): 18

Spam message embedded in image: 5

here 2 type of mails are giving most issues. So it is better to work on them than the rest

After detection of errors, we can fix them by adding more related features, like time and location and removing the old ones or adding new ones.

Once we find all common errors, we can make changes like → get more data, adding more features, etc. for those specific errors.

Note: In real projects, the data size are huge, so we can take random 100 ratio errors to do analysis.

Error Analysis is better for thing humans are good at.

Adding Data

When working, it is tempting to get more data on everything. However, it is important to get data of types where error analysis has indicated it might be of help.

Eg. going to unlabeled data & finding more examples of Pharma related spam.

Other than getting brand new training examples, we can use 'Data Augmentation'

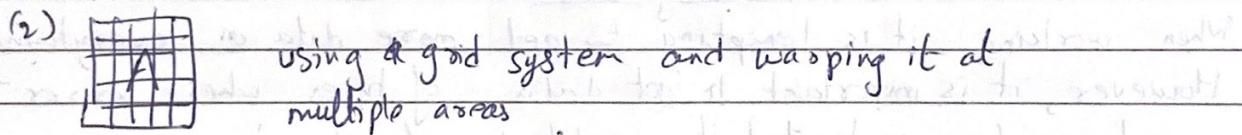
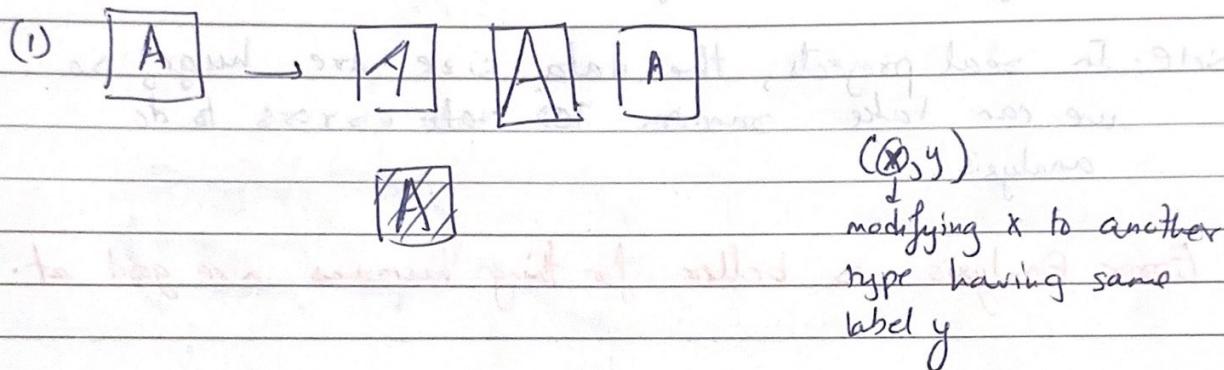
Note: Only augment data in ways which might be present in the test set.

It is not helpful to add purely random noise to our data.

Data.

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

Eg. image of 'A' detection



Eg. Speed recognition example

- | | |
|--------------------------------|------------------|
| [1] → original audio | original audio |
| [2] + Crowd noise | background noise |
| [3] + Car noise | noise |
| [4] + bad cellphone connection | |

we can combine the original audio with each of the background noise to create more training data

Data

Synthesis : using artificial data inputs to create a new training example.

Eg. OCR - optical character recognition

When seeing an image, algorithm will detect characters.

For training, we have images of billboards that the computer uses to detect/learn how to detect characters.

Instead of getting more data, we can use a text editor and spam letters in different fonts to get more data.

Note: Algorithms to produce realistic synthetic data can take a large amount of time to develop but help a lot when done.

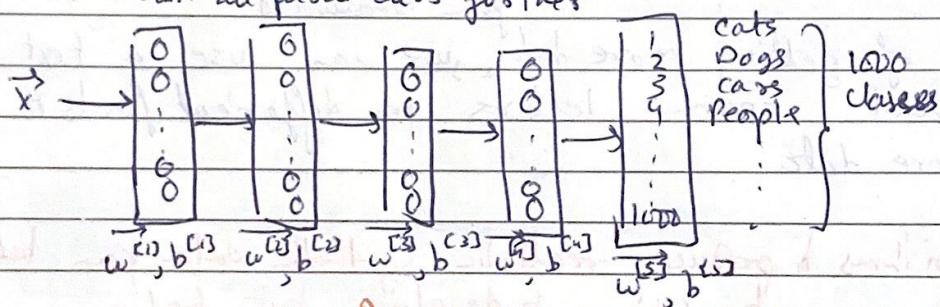
Note: Most notably seen in computer vision tasks. development.

Transfer Learning

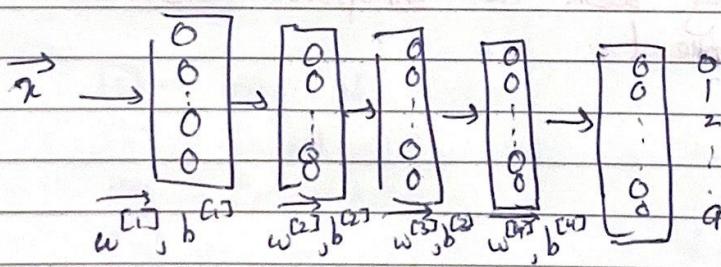
1. Download NN parameters pretrained on a large dataset with same input type (eg. images, audio, text) as your application (or train your own)

2. Further train (fine tune) the network on your own data.

Ex ↗ to train only output layer (for small data set)
 ↗ to train all parameters further



↑
Supervised pretraining



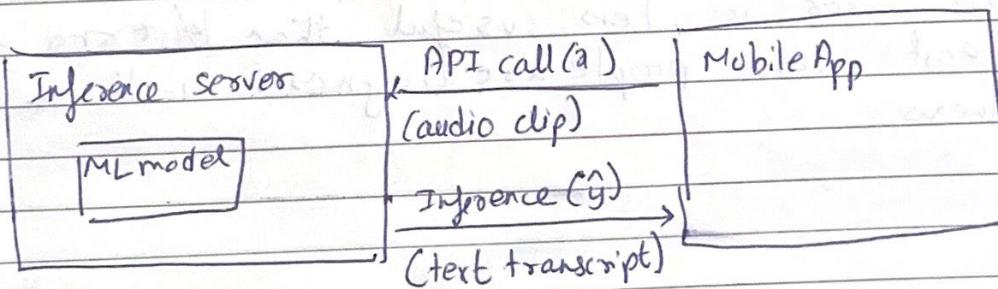
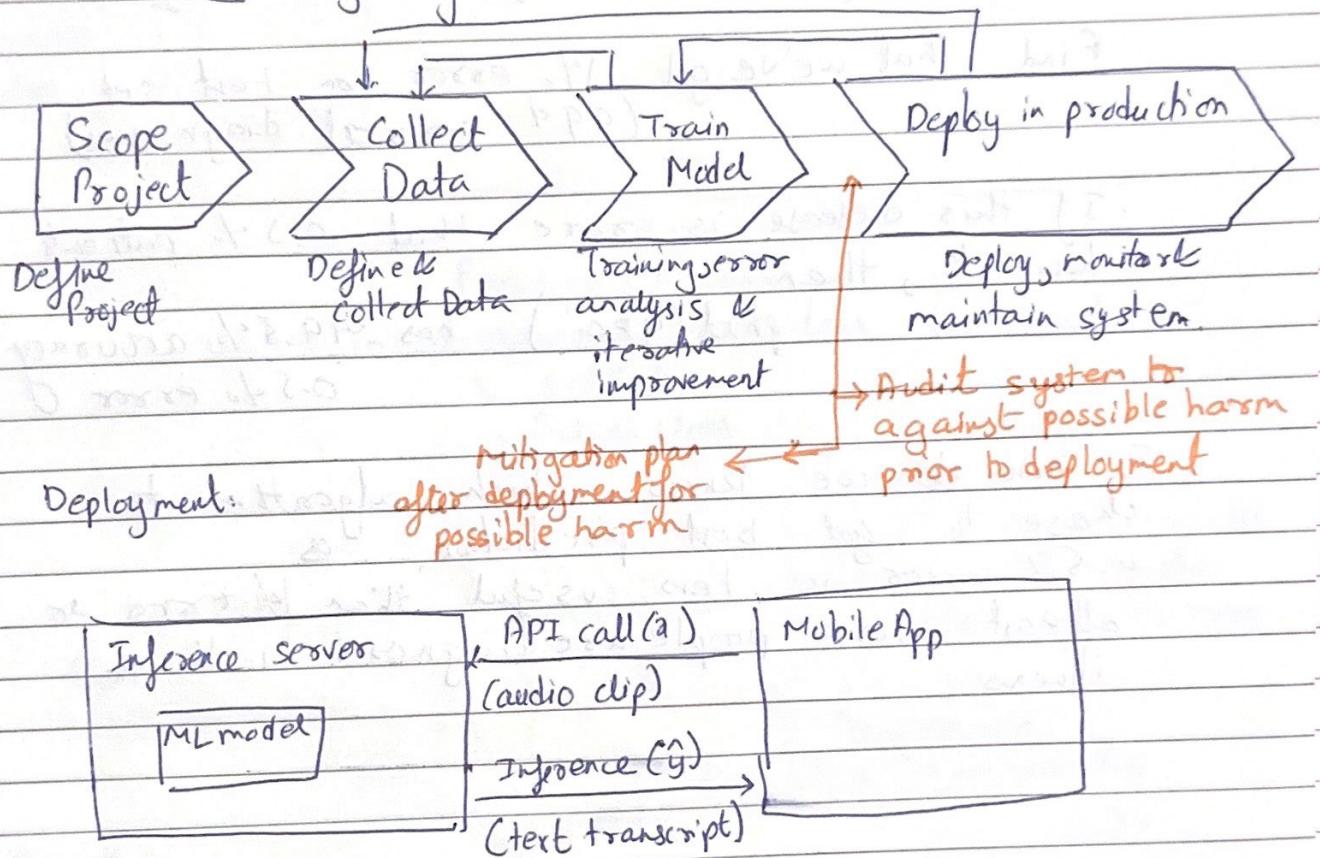
↑
fine tuning

Using a pretrained
NN helps as
several things like
edge detection / corners
and other features
are common amongst
all NN

transferring parameters of each layer except final layer
from trained NN to our NN.

After transfer, training our NN helps smoothen out the NN for its tasks.

Full Cycle of a Machine Learning Project



Software Engineering may be needed for:

- Ensure reliable & efficient predictions
- Scaling
- Logging
- System monitoring
- Model updates

Skewed Data Set

Skewed data sets are data sets that do not have a 50-50 classification training set.

Eg. Train classifier $f_{\text{ab}}(\vec{x})$ ($y=1$ if disease is present,
 $y=0$ otherwise)

Find that we've got 1% error on test set
(99% correct diagnoses)

If this disease is rare that 0.5% patients have it, then
point ($y=0$) has 99.5% accuracy
0.5% error

So how do we know which algorithm to choose to get best prediction, as 0.5% errors is less useful than 1% error so atleast some people are diagnosed with the illness.

Precision / Recall

$y = 1$ in presence of rare class we want to detect

Actual Class

Predicted class

		1	0	Precision:
		True Positive	False Positive	of all patients predicted $y = 1$, what fraction actually have the rare disease?
1	0	15	5	
	1	False Negative	True Negative	$\frac{\text{True Positives}}{\text{True Pos} + \text{False Pos}}$
		10	70	$= \frac{15}{15+5} = 0.75$

Recall:

of all patients that actually have the rare disease, what fraction did we correctly detect

$$\frac{\text{True Positives}}{\text{True Pos} + \text{False Neg}} = \frac{15}{15+10} = 0.6$$

Precision and Recall

Tradeoff

Suppose logistic regression: $0 < f_{\vec{w}, b}(\vec{x}) < 1$

- ↳ Predict 1 if $f_{\vec{w}, b}(\vec{x}) \geq 0.5$
- ↳ Predict 0 if $f_{\vec{w}, b}(\vec{x}) < 0.5$

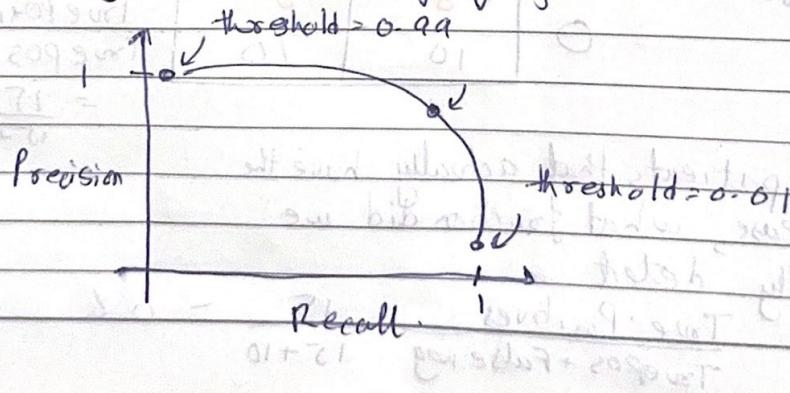
Suppose we want to predict $y=1$ (rare disease)
only if very confident

higher precision, lower recall

Suppose we want to avoid missing too many
cases of rare disease (when in doubt predict $y=1$)

lower precision, higher recall

More generally, predict 1 if: $f_{\vec{w}, b}(\vec{x}) \geq \text{threshold}$



We can choose the threshold manually or using 'F1 score'

Eg.

	Precision(P)	Recall(R)	F ₁ Score
Algorithm 1	0.5	0.4	0.444 ← better
Algorithm 2	0.7	0.1	0.175
Algorithm 3	0.02	1.0	0.0392

$$F_1 \text{ Score} = \frac{1}{\frac{1}{2} \left(\frac{1}{P} + \frac{1}{R} \right)}$$

$$= \boxed{\frac{2PR}{P+R}}$$

Harmonic
Mean Formula in
Maths.