

STRUCTURE<sup>o</sup>

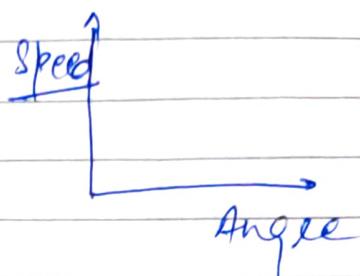
M

PROJECTS<sup>o</sup> →

# Week - 01

Camlin	Page
Date	/ /

Orthogonalization :-



Chain of Assumption of MC

1. Fit training set well on cost function [BISSECT AND DIVIDE]
2. Fit dev set well on cost-function [Regular Dropout? BISSECT AND DIVIDE]
3. Fit test set well on cost-function [BISSECT AND DIVIDE]
4. Perform well in real-world [change dev set cost function]

$$\underline{x_1 + \underline{x_2} - \underline{x_3}}$$

2

$$\underline{x_1 + \underline{x_2} - \underline{x_3}}$$

2

$$\underline{\underline{x_1 + x_2 - x_3}}$$

# Single Number Evaluation Metric

Class	Precision	Recall
A	95%	90%
B	98%	95%

combine them by f1 Score =  $\frac{2PR}{P+R}$

## Dev Set + Single Real Number Evaluate Metric

Algo	US	Chu	Endu	other	Avg
A	3%	7%	5%	6%	6%
B	51%	6%	5	10	64%

Take Avg to achieve Single Number Evaluation Metric

Classified	Acc.	RunTime
A	90%	80 ms
B	92%	95 ms
C	95%	1500 ms

we also care for Runtime :-

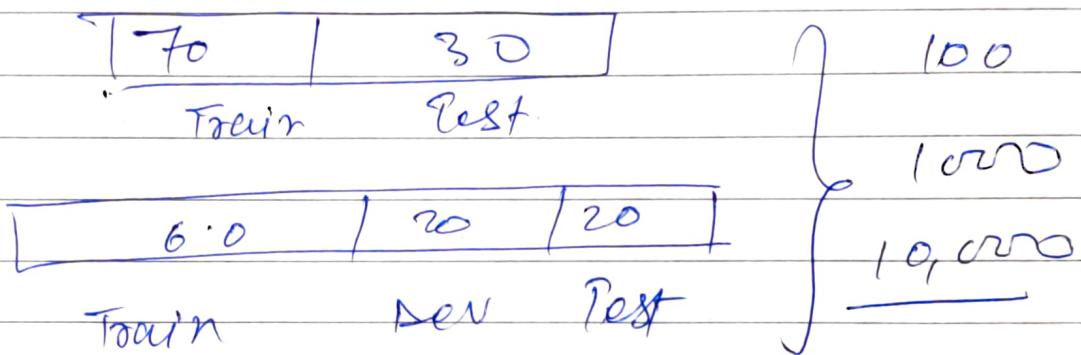
choose classified MAX ACCURCY

subject to  $\text{Runtime} \leq 100 \text{ ms}$

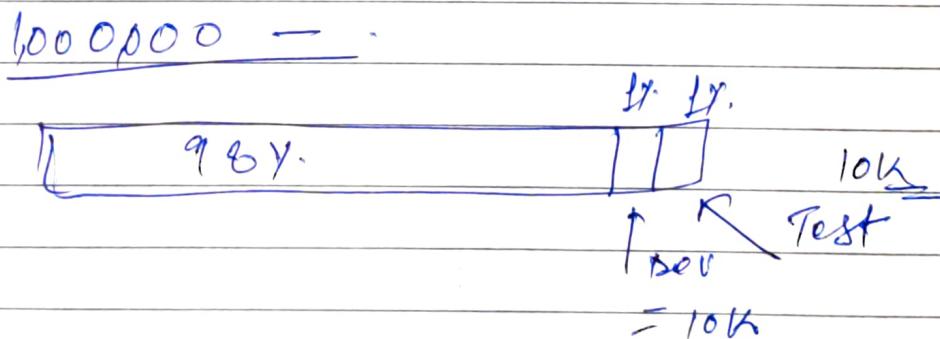
Ans :- Optimizing / Running : Satisfying Metrics

In general if N metric I optimizing  
of N-1 will be satisfying :-

Same distribution for Dev / Test test



If we have 1 million train example:-



If we have too large data-set then use  
much large data-set subset can be  
used for TRAIN!

## 2 Improving Model Performance

Two fundamental assumptions of supervised learning

1. You can fit the training set pretty well  
that means low avoidable bias
2. The training set performance should be generalize pretty well to the dev/test set. that means low variance

Avoidable Bias

- 15 Train Bigger Model
- Train longer Better Optimization ALGO
- Hyperparam Search

Variance

- 25 More Data
- Regularization - L2 / Dropout
- Data - Augmentation

10,000,000 (10 million, 420  $\Omega$ )  
(NoBird | Bird?)

Metric Success :- 1. High Accuracy 2. Runtime less  
3. less memory.

5

10

15

20

25

30

## Carrying Out Error Analysis

Ex. A classifier to ~~classify~~ <sup>CAT</sup> ~~dog~~, having 90% accuracy 10% error. How to do error analysis:

1. Take 100 mislabeled examples from DBV SBT

2. Count how many dogs has been missclassified. Say it turns out 5% are dog

3. So even if do very good on dog problem so my accuracy will increase error will improve by 0.5%. i.e 9.5% now.

④ Evaluate multiple Ideas in parallel

1. Ideas for cat detection:-

1. Fix pictures of dogs being recognized as cats.

2. Fix great cats (lion, panther etc) being misrecognized.

3. Improve performance on blurry images.

Image	Dog	Cheetah	Blur	Wrong Label	Comments
1				1	Pitbull
2				1	
3				1	
4				1	
5				1	Rain in zoo.
100	84	43%	61%	67%	
% of Total	84%				

Cleaning Up Incorrect Label: Can be deal as

on more column here. If OK leave it  
if % of incorrect label is not more  
than RANDOM,

Overall Dev Set Error : 10%.

Error Due to Incorrect Lbl: 0.6%.

Error Due to other cause : 9.4%.

L Dev / Test > must from same distribution

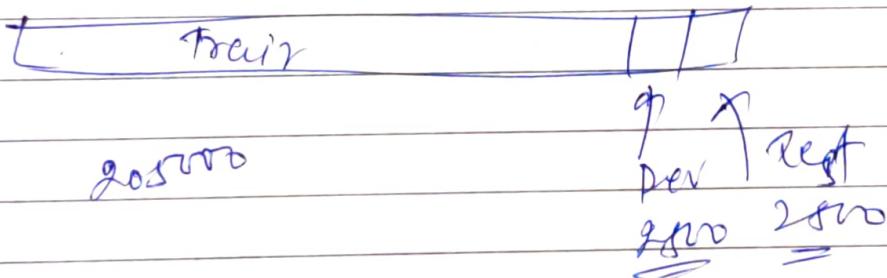
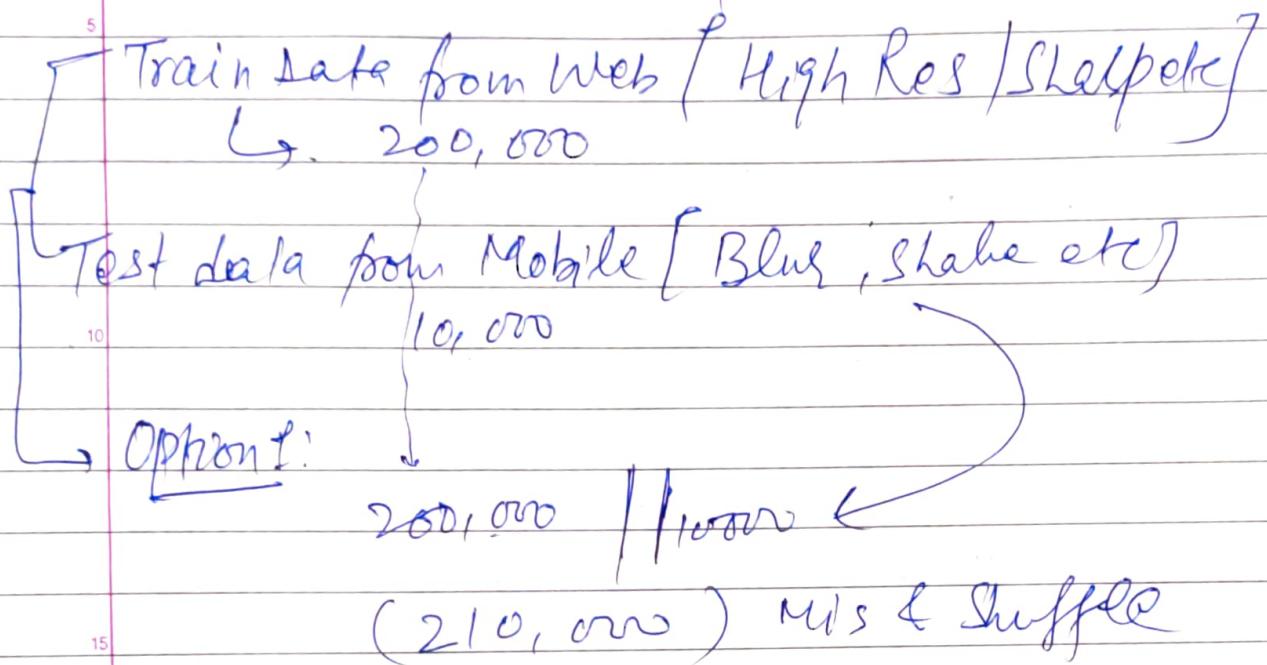
L Train // L Dev/Test may now come  
from slightly different distribution.

② Build Your First System Quickly  
then Iterate.

# Mismatches in Train / Dev / Test

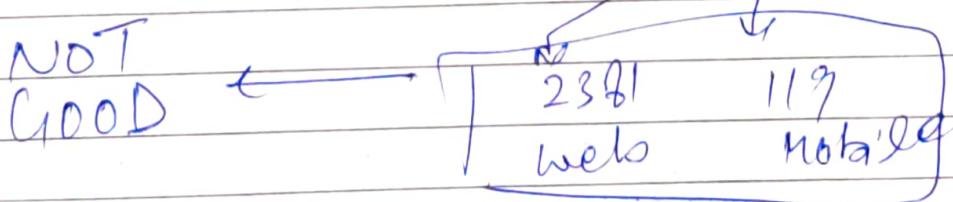
## Training & Testing on different DTS

Cat app example



Advantage: Train / Dev / Test come from same dist.

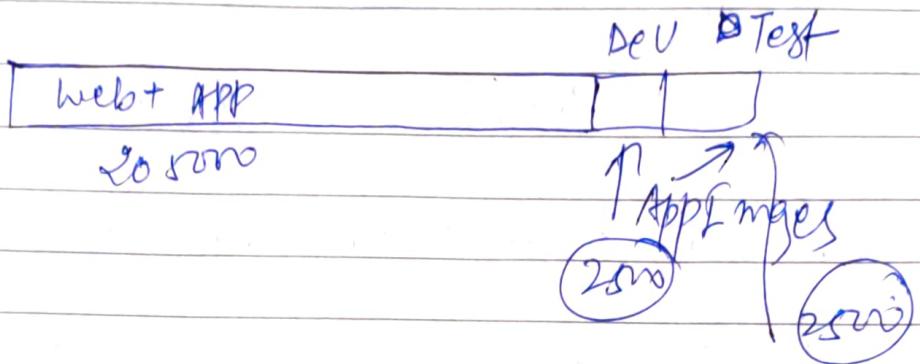
DisAdvantage: Dev Set (2500 approx)



We are saying to do GOOD web images. That is not MOTO.

# OPTION 2 NOT OK

option 2:



You are forgetting what we really care about.

& Bias & Variance with mismatched data  
distributions ↴

Cat classifier: Assume human get  
 $\approx 0\%$  error.

Train Error :- 1% } if train &  
Dev Error :- 10% } Dev have

same distribution I will say.  
there is a variance problem.

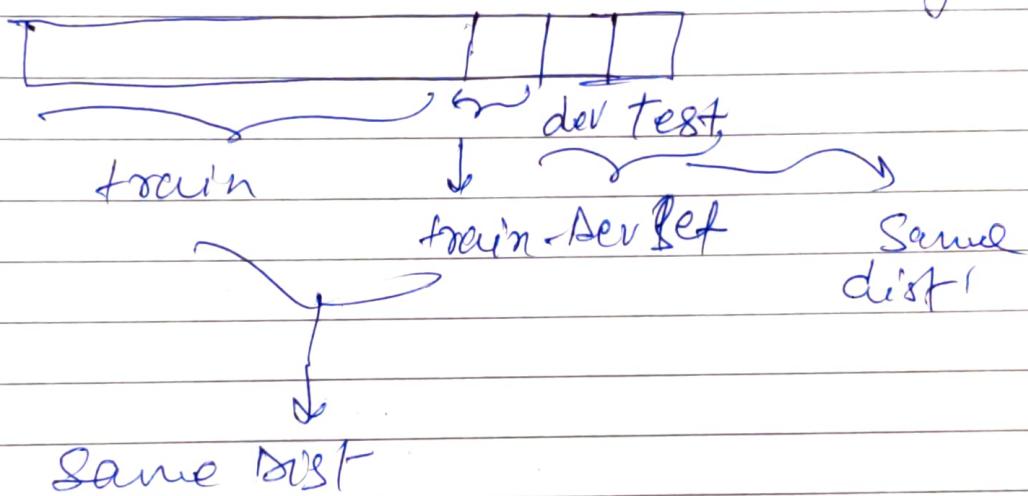
It is not generalizing well.

But if Dev & Train is from "DIFFERENT DISTRIBUTION"

So we coin a new term

Training-Dev Set :- Same dist as

train but not used for training.



Train 1Y.  
Train-Dev 9Y.  
Dev Err

so conclusion is  
we have VARIA  
ncs problem

Train 1Y.  
Train-Dev 1.5Y.  
Dev Err 10Y.

↑ Data miss match  
problem.

Human Err  $\approx 0\%$ .

Train Err  $\approx 10\%$ .

Train - Dev Err  $\approx 11\%$ .

Dev Err  $\approx 12\%$ .

Avoidable

Bias problem

Human Err =  $0\%$ .  $\uparrow \rightarrow$  Avoid Bias

Train Err =  $10\%$ .  $\downarrow$  problem as well as data

Train Dev Err =  $11\%$ . MESS MATCH

Dev Err =  $12\%$ .

Human Dev Err  $12\%$ .

Avoidable Bias

Train Set Err  $\downarrow 7\%$

Train - Dev Err  $\downarrow 10\%$ .

Dev Err  $\downarrow 12\%$ .  $\rightarrow$  data miss match

Test Err  $\downarrow 12\%$ .  $\rightarrow$  degree of overfitting.

If test Err higher than maybe we have to find bigger "Dev. Set".

## Addressing Data MissMatch

- Chorng out manual error analysis fo  
try to understand difference  
b/w training Dev / Test  
sets.
- Make training data more similar

=)) + Car Noise = Synthe  
sized  
The quick brown  
fox jump  
over the lazy dog

↓  
1 hrs

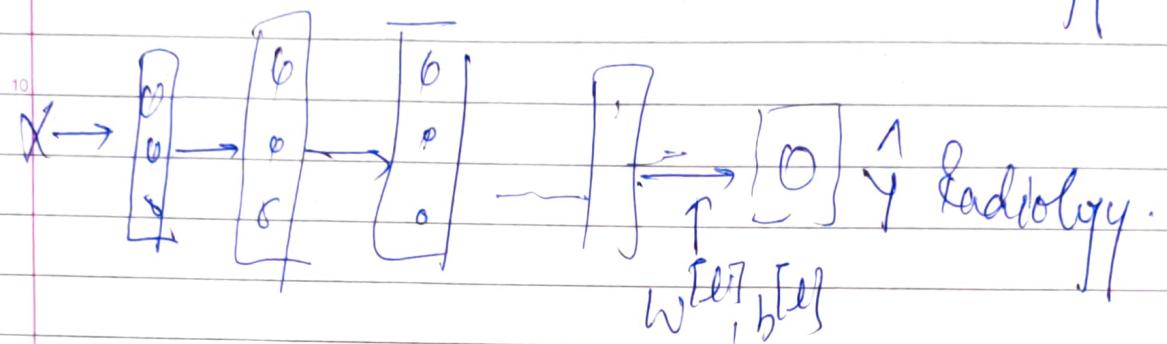
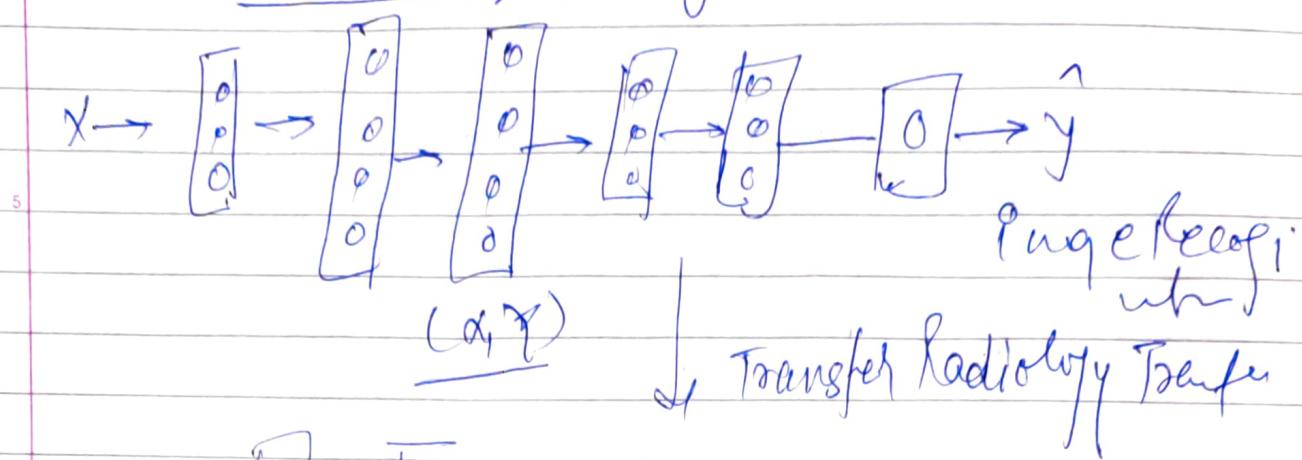
10,000 hrs

So we can  
overshoot to 1 hr of  
car noise

1 sec / hr  
Car Noise

Set of all audio of a car.

## Transfer Learning



15 1,000,000 for cat  
 1/1000/100 Radiology Image  $\rightarrow$  good transfer case

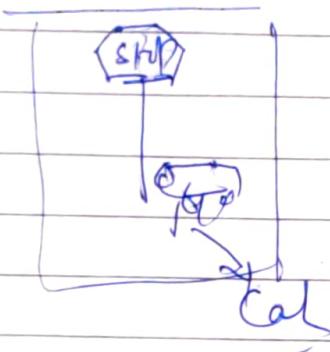
Task A  $\rightarrow$  Task B

- 20 \* Task A & B have same input  $x$ .  
 \* You have lot more data for A ~~and~~ B.  
 \* Low level feature from task A useful for task B.

# MultiTask Learning

Target : Pedestrian, Car, Sign, Traffic light

$x^{(i)}$

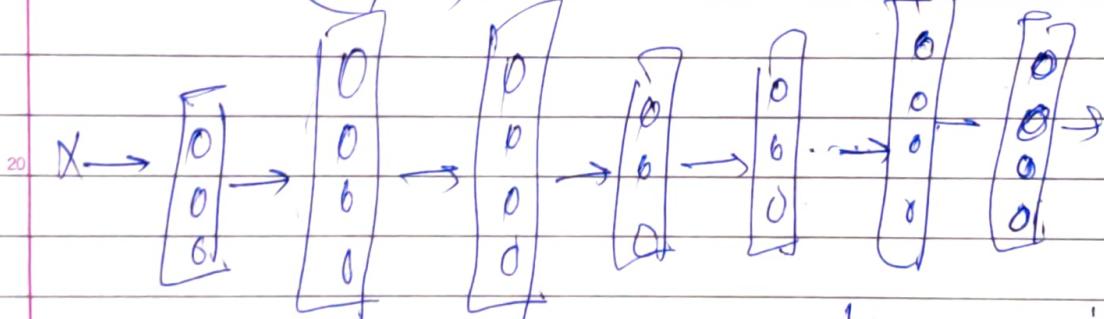


$x^{(i)} \rightarrow y^{(i)}$

$x^{(i)} \rightarrow \begin{cases} 0 \\ 1 \\ 1 \\ 0 \end{cases} \quad (\text{y}_{\text{all}})$

$$y = [y^{(1)}, y^{(2)} - y^{(\text{out})}]$$

$(y_{\text{all}})$



Loss  $y^{(i)}$   
 $(y_{\text{all}})$

$$\frac{1}{m} \sum_{i=1}^m \sum_{j=1}^4 L(y_j, \hat{y}_j)$$

$$- y_j^{(i)} \log \hat{y}_j^{(i)}$$

logistic loss

$$- (1 - y_j^{(i)}) \log(1 - \hat{y}_j^{(i)})$$

Unlike softmax Regression one image can have multiple labels.

When multi-task make sense?

- ① Training on a set of tasks that could benefit from having shared lower-level feature.
- ② Usually amount of data we use should be quite similar.
  - $A_1, 1000$
  - $A_2, 1000$
  - $A_3, 1000$
  - $A_{100}, 1000$
- ③ Can train a big enough neural-network to do well on all the tasks.

End-to-end Deep Learning

## Where to use deep learning?

Pros:-  
1. Let the data speak  $x \rightarrow y$

2. Less hand design hand-components.

3.

Cons:-  
1. Need large amount of data

2. Excludes potentially useful hand-design system.