

## Regularization and Best Practices

Regularization:- Any modification we make to learning algorithm to reduce its generalization error rather than training error. Trade increased bias for variance by adding some prior knowledge.

### Strategies

- Extra constraints and penalties
- Restrict param values
- Data Augmentation
- Adversarial Training
- Ensemble methods | Dropout

→ Better to have an overcapacity model (overfit) and then regularize it to required level rather than searching for correct capacity.

### Parameter Norm Penalties

→ Limit capacity of model by adding parameter norm penalty  $\Omega(\theta)$  to objective function  $J$

$$\tilde{J} = J(\theta; X, y) + \alpha \Omega(\theta)$$

$\alpha=0 \rightarrow$  unregularized model  $\Rightarrow \alpha=\text{some large number} \rightarrow$  underfit

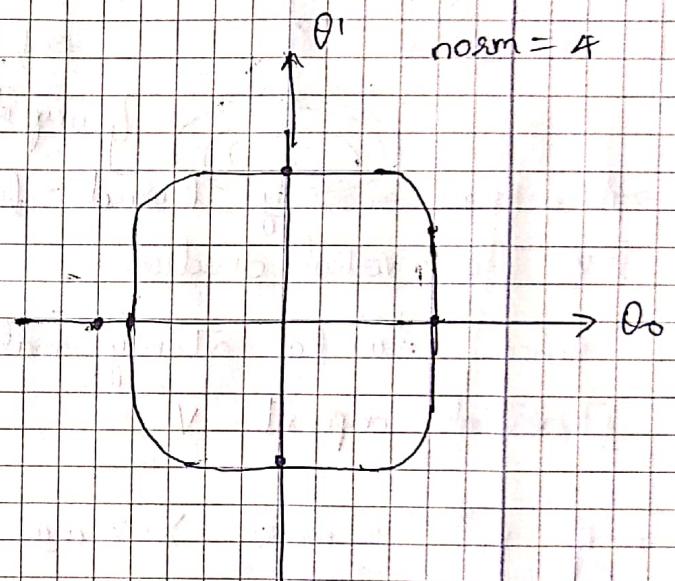
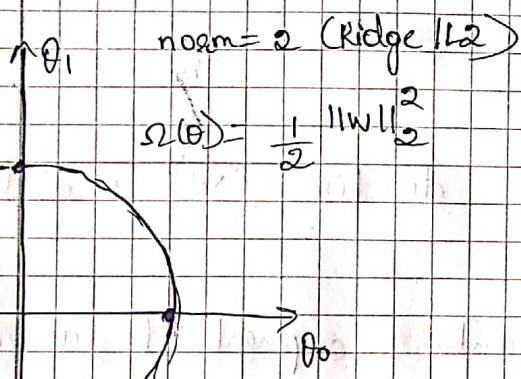
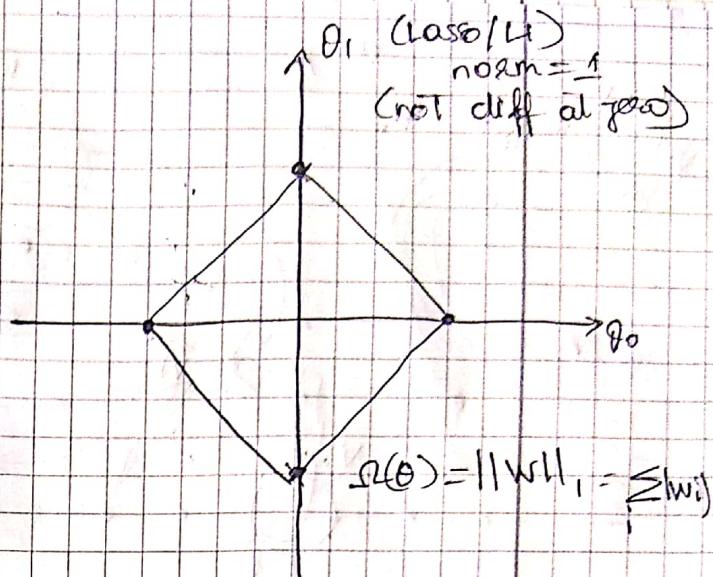
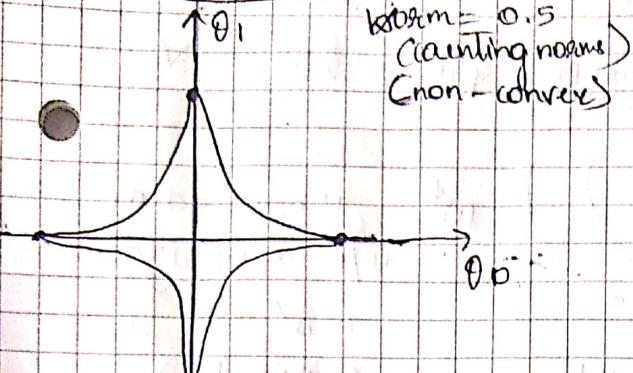
→ Typically we omit bias from regularization  $\Omega(\theta)$  because bias term is much more simpler & depends induced by one value rather than weight which is induced by multiple values

\*→ Usually Each layer must have a separate penalty ' $\alpha$ ' but we use a single penalty ' $\alpha$ ' To ease computation

Need for parameter sparsity :- if parameters are sparse (Mostly near one / zero and not between). It is easier to store and compute

→ Most useful params would be params near one & useless ones tend near zero. Norms tend to provide sparsity in various ways

Different Parameter norms



L<sub>2</sub> norm

$$\tilde{J}(\theta; x, y) = J(\theta; x, y) + \frac{\alpha \|w\|_2^2}{2}$$

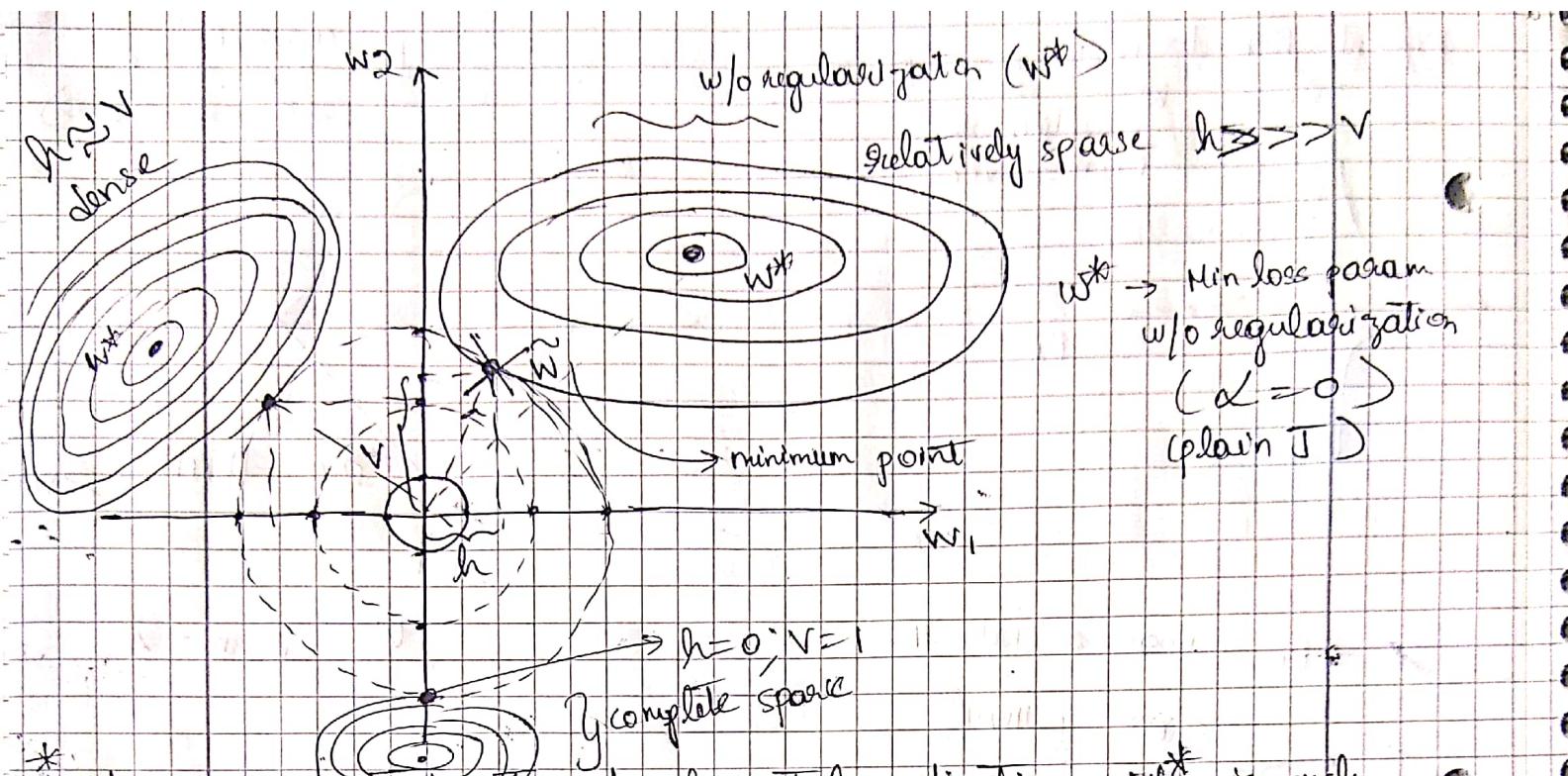
$$\nabla_w \tilde{J} = \alpha w + \nabla_w J(w; x, y)$$

$$w' = w - \epsilon \nabla_w \tilde{J}$$

$$w = w - \epsilon \alpha w - \epsilon \nabla_w J(w; x, y)$$

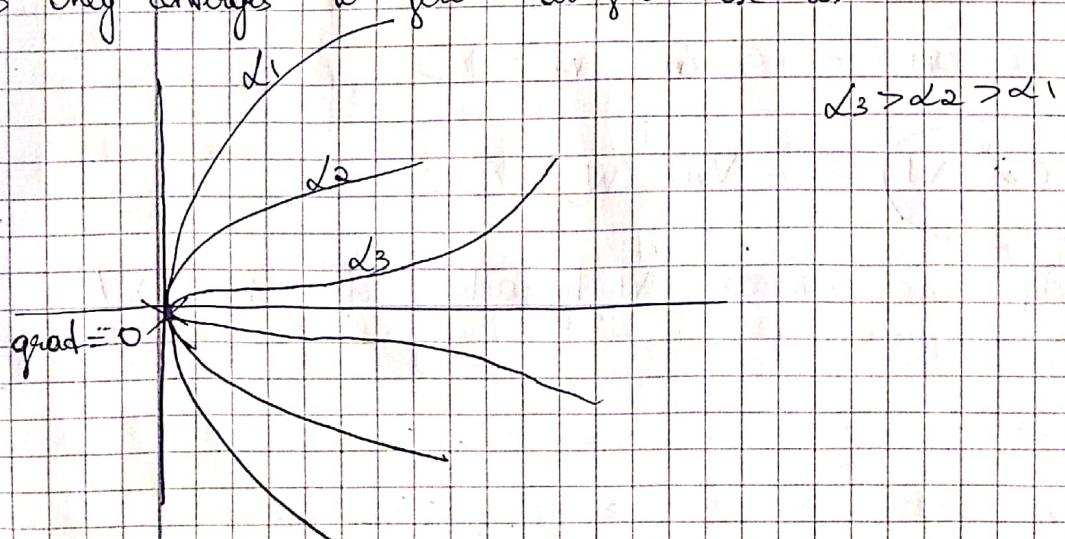
$$w = (1 - \epsilon \alpha) w - \epsilon \nabla_w J(w; x, y)$$

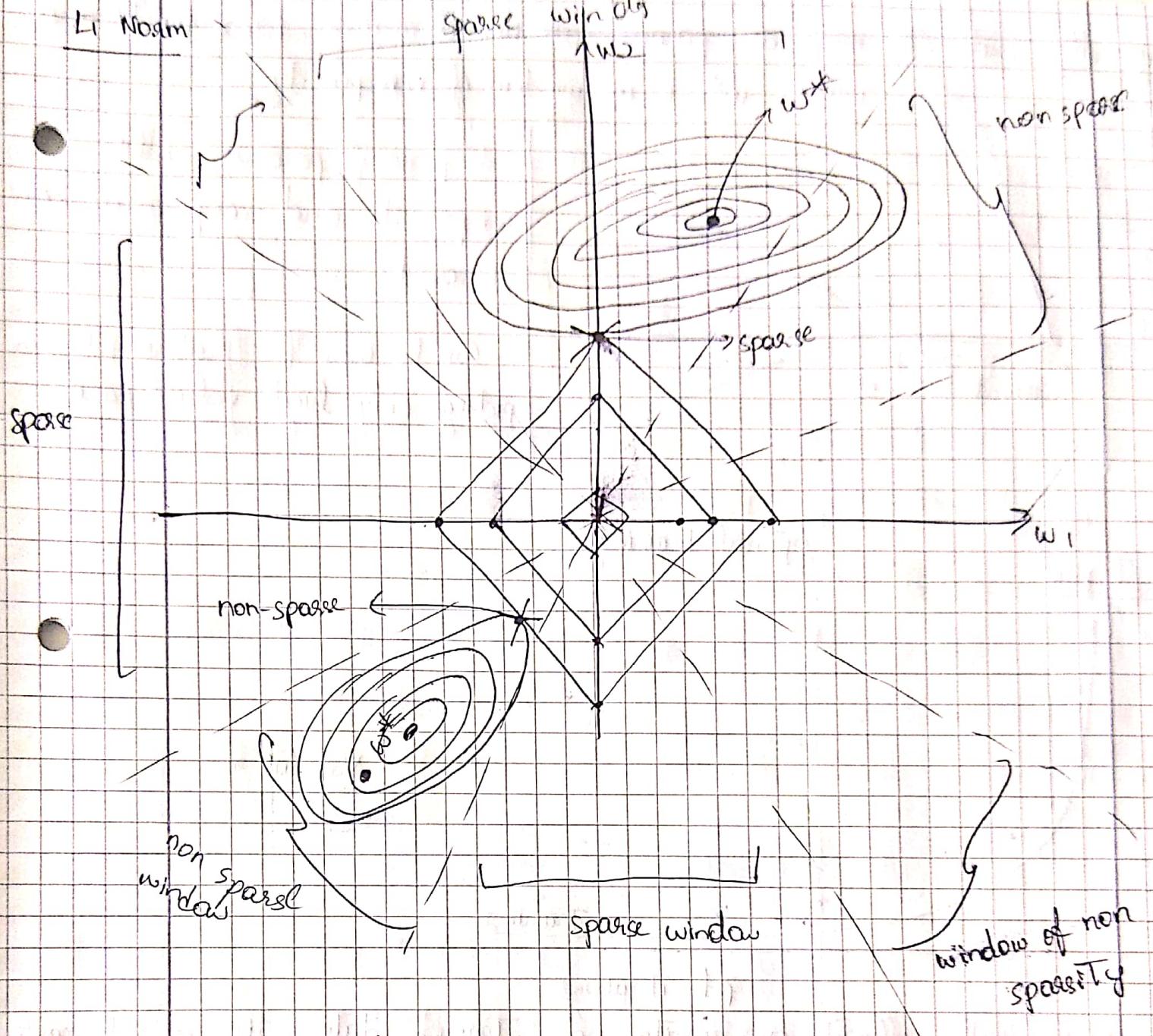
\* At each Training step we decrease w by factor of  $(1 - \epsilon \alpha)^{-1}$ .



- \* As we know sensitivity towards horizontal direction  $w^*$  is much lesser than vertical direction
- \* Hence  $\hat{w}$  has extremely small horizontal component 'h' as compared to vertical component 'v'

- if  $w^*$  is near X, Y axis we would get much sparser matrix
- if  $w^*$  is at  $45^\circ$  to center then we would get dense matrix
- L2 norm penalizes large weights due to " $(1 - \alpha \epsilon)w$ " term in gradient
- if  $w$  is large then  $-\alpha \epsilon w$  is large, hence it pulled down faster
- Only converges to zero at zero or when  $w^*$  is exactly on X or Y axis





$$L_{\text{Noam}} := \Omega(\mathcal{J}) = \|w\|_1 = \sum_i |w_i|$$

$$\tilde{\mathcal{J}}(w, x, y) = \mathcal{J}(w, x, y) + \alpha \|w\|_1$$

$$\nabla_w \tilde{\mathcal{J}}(w, x, y) = \nabla_w \mathcal{J}(w, x, y) + \alpha \text{sign}(\nabla \|w\|_1)$$

if sign is +ve

$\alpha$  → +ve sign grad

if sign is -ve

$-\alpha$  → -ve grad

undefined at zero

→ Undefined at exactly zero

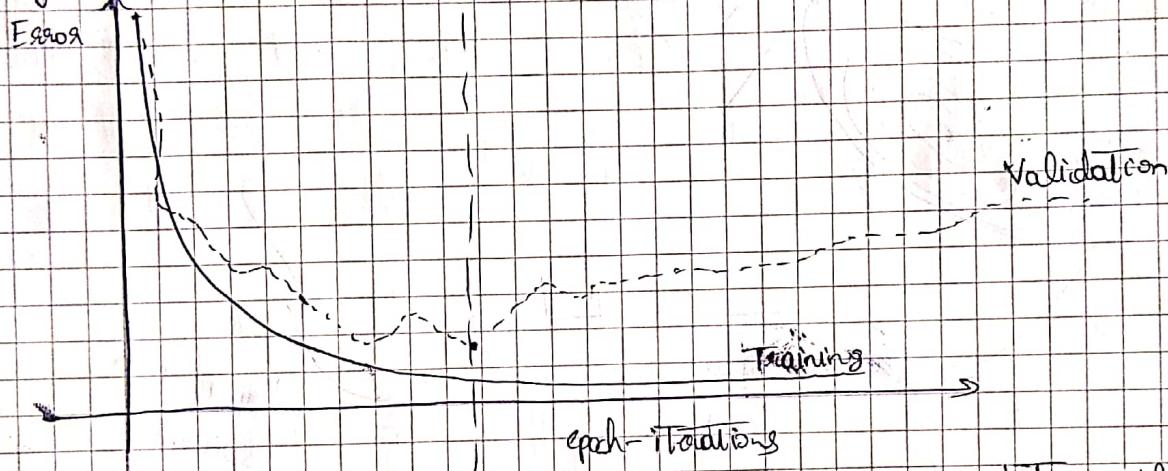
$L_1$  generally tends to non sparse params unless  $w^*$  is in window of non sparsity sign



→ Sparse of params could be compressed and computation can be saved

→ Can be used for feature selection picking only high valued params

## ii) Early Stopping



→ When training with sufficient overcapacity on training data the error decreases but validation error starts increasing after certain point/epoch

→ Why not monitor and stop at this point?

a) To monitor and make a stop decision we must have a good validation set with enough size else we cannot make decision on small validation set

b) Need to save copy of parameters until error decreases for some iterations

for (a) we cannot get away with high enough validation step and second stage retraining. (b) → save in eps/ inexpensive storage

## Retraining

→ if our data set was small and we had taken validation data from it to perform early stopping, Then we can reintegrate back in Two ways

- a) After early stopping, combine old validation + old training data to new Training data. Reinitializing & reset the n/w. Start training from scratch only till epoch where we stopped early
- b) Instead of reinitializing network since it's costly, simply treat validation as part of training set and keep on training till avg error decreases below certain limit. But this still doesn't guarantee minimal as (a).

How early stopping acts as regularizer : IT slows down network growth and keeps it to small neighbourhood of initial parameter  $\theta_0$

- \* Much better than L1/L2 as we are not changing weights and allowing natural fit of network. Also no hyperparameter tuning like ( $\alpha$ ) is required in early stopping

### Other regularizations

Ensemble, Dropout, Adversarial Training, parameter sharing etc

### Best Practices (Machine learning / Learning)

#### I) Setting up dev and test sets

Trainingset := On which we train our algorithm

Dev set := Tune parameters, select features, (hold out cross-validation)

Test set := Evaluate final performance, but don't use for making any decisions like tuning hyperparams

→ Dev & Test allows team to see how algo is doing.

→ 70/30 split may not be prevalent in modern deep learning era as Neural net are data hungry & perform well on large data ( $> 95\%$ )

\* Choose Dev and Test set to reflect data you expect to get in future and which you are planning to do well on. Collected better dev/test set

\* Choose dev and test from same distribution; if you don't have this data now, then proceed to deploy. Later once you have collected relevant data from user, retrain the model with new dev/test set and redeploy it.

Example :- Training set :- cat pictures from website + Japan + India + USA  
Dev set :- cat pictures from USA + India } Different distributions  
Test set :- cat pictures from China

if  $\text{Error}_{\text{dev}} \gg \text{Error}_{\text{test}}$  Then reasons could be :

- i) Overfit on dev set
- ii) Even though test set dev are from distro, they test set could be much harder to learn (cannot do anything about this)
- iii) Dev & Test from different distribution  $\Rightarrow$  (Wasted Optimization effort)

$$0.01 - 1 = \frac{0.01}{100} = 1 \times 10^{-4}$$

### ii) Size of Dev and Test set

→ Dev set should be large enough for us to detect improvements b/wn two classifiers. For example :- if Dev set has 1000 examples then classifier A = 90%. Classifier B = 90.01%.  $\Rightarrow$  (cannot be detected as  $0.01\% = 10^{-4}$ ) and we have only  $10^3$  examples we need atleast  $10^4$  examples to detect  $0.01\%$  change

→ Test set must be large enough to give us confidence about overall performance of your system

### III) Single number evaluation metric

→ Usually having a single number evaluation metric helps to iterate quicker and optimize better than multiple targets

Example :- Instead of Precision + Recall Optimize / Use F1-score

- In case we need or have some weird combination of metric like F1-Score, Running Time, Size, Mem required, False positives, False-neg etc
- keep only one metric as optimizing, on which you try to improve as much as possible like F1-Score
- Running Time, Memory, Model size could be satisfying metric and classifier just needs to be satisfying a cut-off criterion

Selected Class	Classifier	F1-Score	Running Time(sec)	Memory RAM (MB)	Model Size (GB)
✗	A	0.89	400	100	2
✗	B	0.92	450	200	3
✓	C	0.77	150	50	0.7
✓	D	0.85	170	70	1
✓	E	0.72	150	70	1.1
✗	F	0.71	100	90	0.7
✗	G	0.81	200	60	4

satisfying metric  
 $\text{Running Time} \leq 200 \text{ sec}$

$\text{Memory} \leq 75 \text{ MB}$

$\text{Model size} \leq 1.5 \text{ GB}$

Selected classifiers C, D, E whose F1-score would be optimized further

- \* \*
- Select dev/Test set and setup a basic metric asap unless you are changing any app already running in prod or really matured.

#### IV) Changing Dev/Test sets and metric

- Let's say you took initial cat images uploaded by users as your dev/Test set. Now users are uploading kitten images, hence our dev/Test set is no longer valid. Change them accordingly.

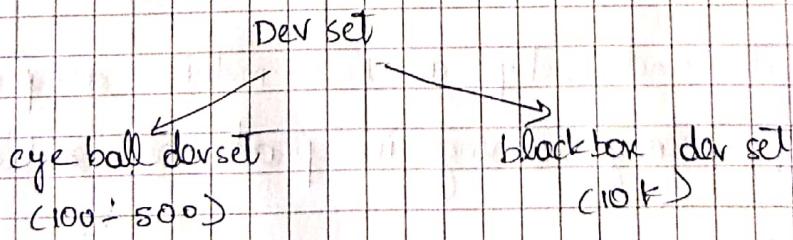
## Problems

- i) The actual distribution may changed from dev/test sets we used.
- ii) You have tuned so much on dev set such that there is significant error b/w dev & test sets indicating that you have overfit on dev set  
Throw away dev set and get new dev set
- iii) Business changed the metrics they told earlier.  
Ex:- Classifier A = 90% accuracy + lets through porno images  
Classifier B = 85% accuracy + doesn't let porno images  
→ We had used A as accuracy was our metric, but it's letting porno images through.  
→ Include penalty on porno in our metric and reiterate process  
$$\text{new metric} = \text{Accuracy} - \text{penalty} * \frac{\sqrt{\text{# of porno images}}}{\text{let in}}$$

## Basic Error Analysis

Error analysis :- Process of finding/examining dev set examples that our classifier got wrong To better understand its underlying cause.

- \* Only fix errors that may give you max boost in your metric
- Make an excel sheet of errors and their plausible reasons



- Instead of looking at each error manually, make a separate dev set called eyeball from existing one which is representative of all categories but still small enough to be manually examined.
- Analyse and make spread sheet on eyeball dev set

Image id	Mislabeled	Cat	Dog	Mouse	Error reason
1	✓	✓			Instagram filter
2			✓		Instagram filter
3	✓				
A					
:					
99					
100	✓			✓	Toosmall and behind elephant
Total misclassified	20/-	5/-	70/-	25/-	

→ Instead of fixing cat / mouse , if we fix dog related error we can get more gains and loose less

\* Use this as eyeball dev set and tune your hyper params on blackbox dev set . Don't tune it on eyeball set unless its absolutely necessary . If you see overfit . Then throw this away and get new dev set

→ For smaller devsets  $\leq 1000$  make eyeball & blackbox set the same

### Bias and Variance

~~Bias~~ ↗ avoidable → can be bridged by bigger model  
 ↗ unavoidable → even humans make mistake on this.

Variance : Difference of error b/wn dev v/s Train. Regularize well

Example : Audio recognition : Some audio clips are laden with noise that even us humans make mistakes . Hence human error = 8/-  
 $\Rightarrow$  accuracy = 92%. (Optimal bayes error  $\Rightarrow$  Unavoidable bias)

$$\text{Optimal error} = 8\%.$$

Training error = 20%

Dev error = 20%

Reason: High bias & low variance

Correction: Build bigger model

Training error = 10%

Dev error = 20%

Reason: High variance & low bias

Correction: Regularize model  
Add more training data

Training error = 20%

Dev error = 30%

Reason: High bias & low variance  
<sup>high</sup>

Correction: Bigger Model + Regularize

Training error = 10%

Dev error = 10%

Reason: Low bias & low variance

Correction: Not needed / good model

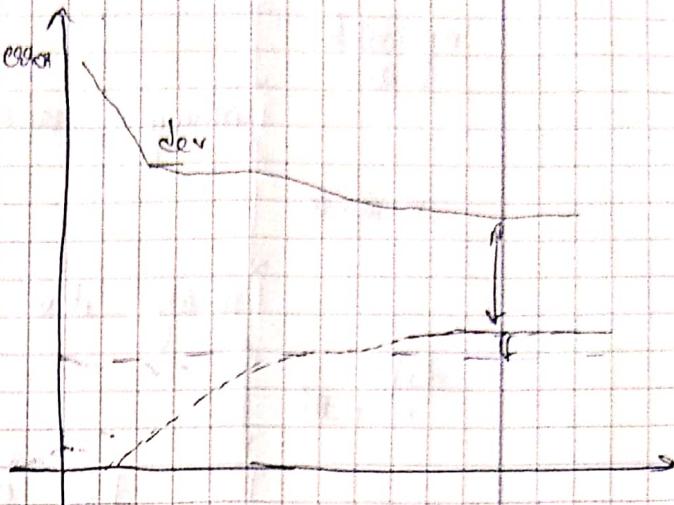
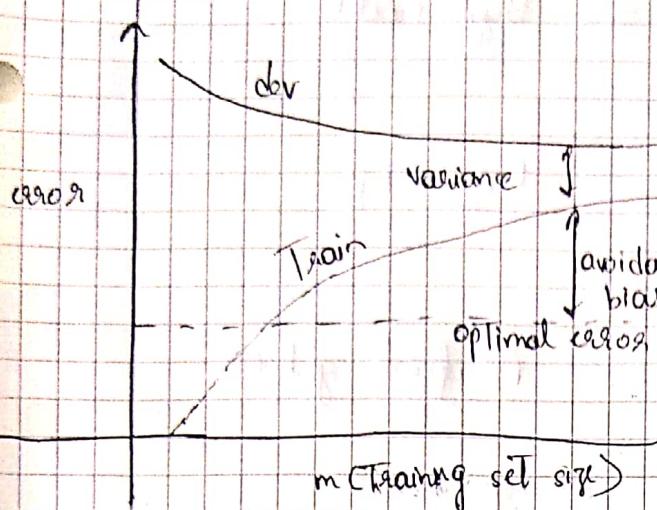
### Techniques To reduce avoidable bias

- Increase model size
- Modify i/p feature based on error analysis
- Turn down regularization
- Modify model architecture

### Techniques for reducing Variance

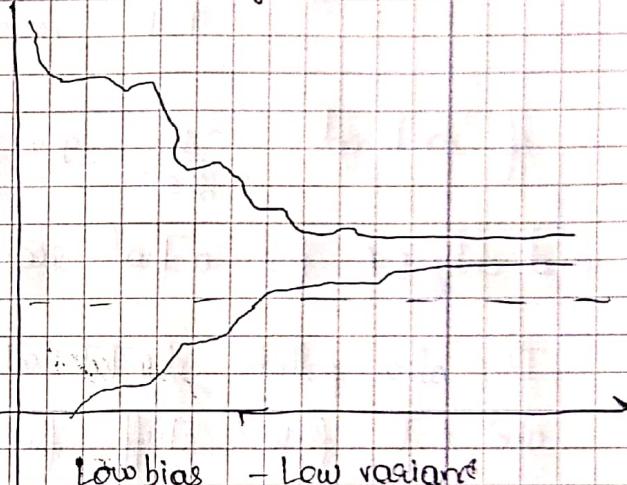
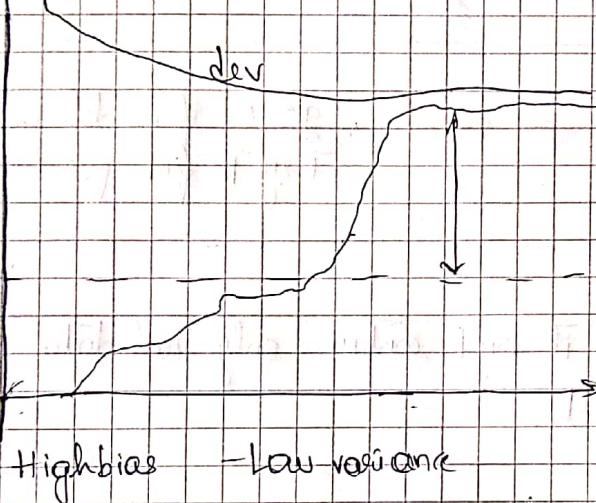
- Add more varied Training data
- Add regularizer
- Early stopping
- Feature selection (select lesser features)
- Decrease model size
- Modify model architecture
- Error analysis

## Diagnosing bias and variance



High bias - High variance

Low bias - High variance



## VII) Human level performance

Collect data from humans if they can solve tasks much better than machines (Image segmentation, Audio recog) and not in one they control (Recommendation, Ad search)

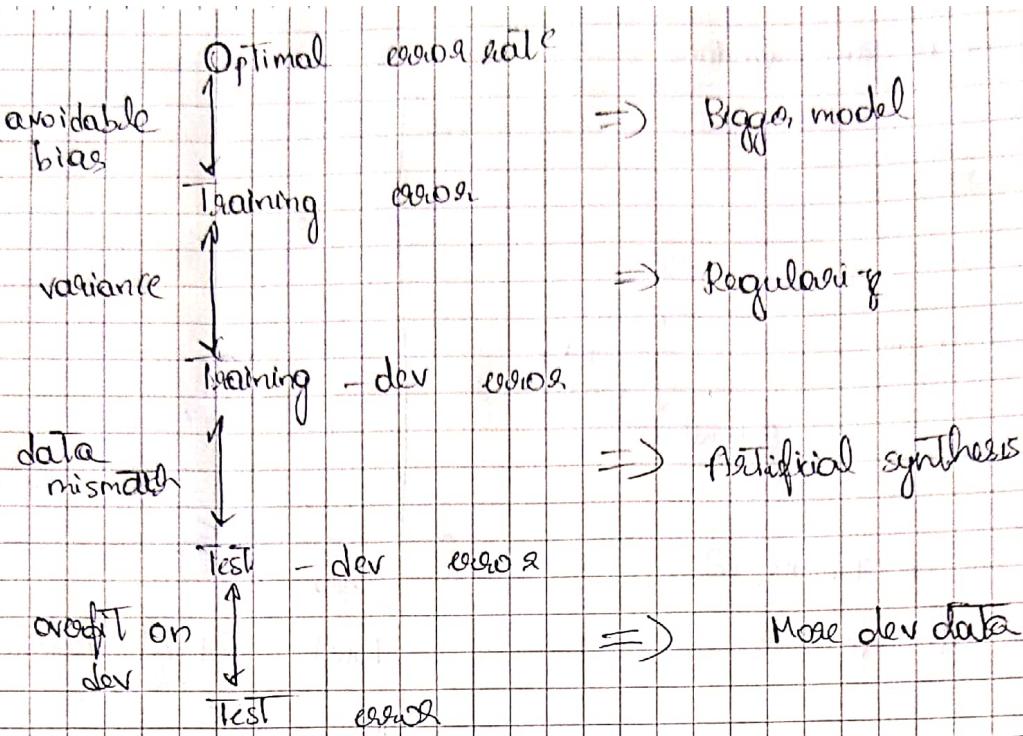
## III) Train and Test/dev from different distrs

\*\*

→ Carve out a subset of training data as train dev set to help see variance  
→ If it's better not to mix Train-dev-Test data and randomly select if  $\text{size}(\text{train data}) \gg \text{size}(\text{dev-test})$

Ex:- 200,000 infrared cat image

20k mobile cat image



→ if combined  $\frac{10200}{220} \approx 24$ ,  $\frac{80}{20} \approx 4$ .  $g2.1.$  of model would try to fit for internet images

and only 8-1% for actual required

→ This also pollutes dev-test set as it must contain only real data (mobile uploads) and not fake stuff from internet.

→ Even if Train is entirely different from dev-test, keep it this way and don't mix. Find any other method for fixing

if Penalty score (Weighing data)

→ if size (Training)  $\approx$  size of (dev-test) Then it's better to mix both and add a penalty weight more for dev-test rather than train

→ So if classifier misclassifies internet image it must get lesser loss compared to misclassifying mobile uploads.

Eg: 50k internet images 20k mobile images  
penalty  $> 1$

$$\text{Loss} = \sum_{\text{mobile image}} \text{Loss} * P + \sum_{\text{internet}} \text{Loss} * 1.0$$

## Artificial data synthesis

- Data mismatch can be addressed by artificially synthesizing dev like data  
(+) and use it to augment training

Ex:- Adding car noise to clear audio of training.

## Problems

- if we synthetically add same car noise (Audi / Volkswagen) to training  
Then if dev has different car (Ford / Ferrari), it still does give mismatch
- Careful synthesis required