

## 1) Define the ML Recipe:

- i) Train the network and evaluate first on the training data
- if bias is high (underfitting) (bad Perf. in train)
    - ↳ use bigger Network [more hidden layers/nodes]
    - ↳ Train longer
  - Check the bias again and keep changing until you reach good bias

## ii) Check for Variance (Validation set Performance)

- if variance is high (overfitting) (bad Perf in valid.)
- ↳ use more data (if possible)
- ↳ use regularization
- Check for Bias, then variance, keep changing until good bias  
good var

## iii) search for better NN architecture, that better suits the problem!

## 4) Diff between WRAPPERS &amp; FILTER method for feature extraction?

Wrappers → Select features by taking the classifier into consideration

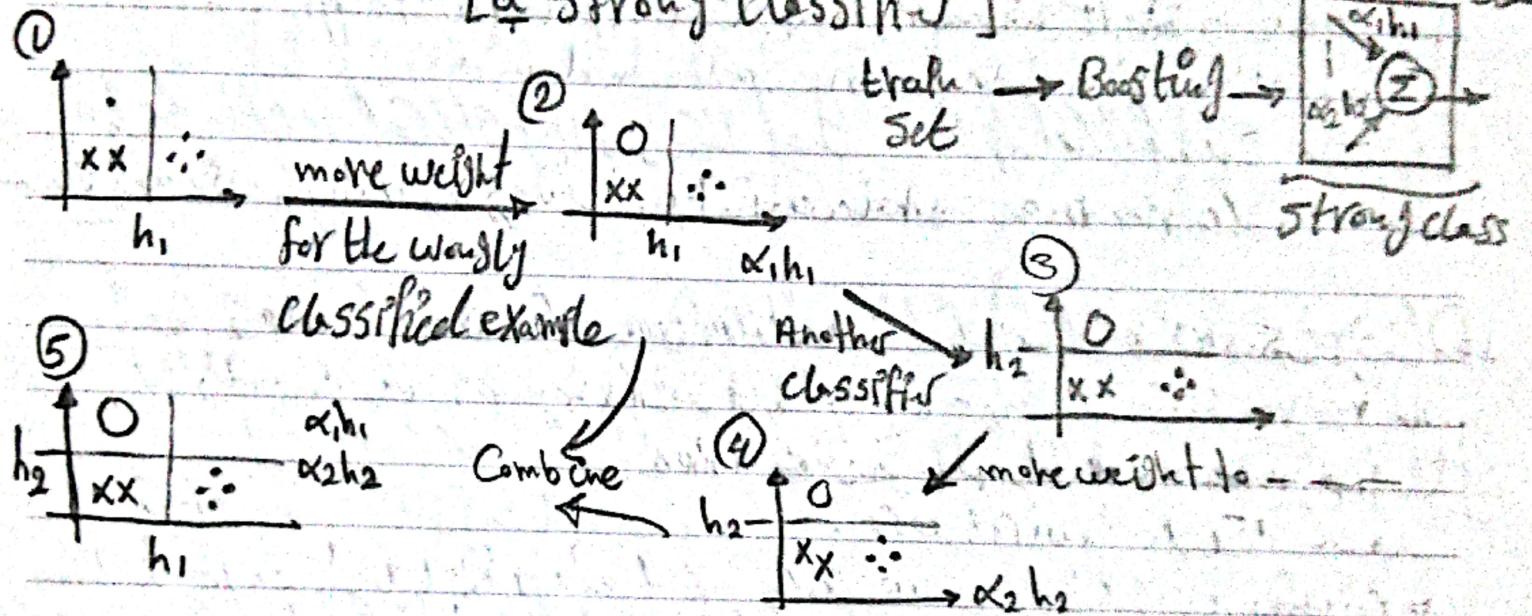
- e.g. SFS, SBS
- Adv: High Accuracy
- Dis: Slow Execution, lack of Generality

Filters → Select Features regardless of the classifier, based on

- e.g. PCA [Corr. betw. features, distance bet. class means]
- Adv: Fast exec, General
- Dis: Tendency to select large Subsets

## 2) Key Idea of AdaBoost, using diagrams?

- \* It tries to approx the Bayes classifier by Combining many weak classifiers.
- \* It selects the best weak classifiers to create  $h \rightarrow \text{weak class}$  [a strong classifier].



## 3) Why Can't AdaBoost work for more than binary classification?

- \* Belief assumes that the error of each weak classifier is less than (0.5). However, in case of multi-class, the error is  $\frac{K-1}{K} = 1 - \frac{1}{K}$ . It may be more than 0.5 in some cases, resulting in -ve weights ( $\alpha_t$ )

\* Model Fitting? - we added the  $[\log(K-1)]$  term, to make it +ve in case of  $K > 2$ .

$$-\alpha_t = \log\left(\frac{1-\text{err}_t}{\text{err}_t}\right) + \log(K-1)$$

$$-\alpha_t \text{ is } +\text{ve, iff: } (1-\text{err}_t) > \frac{1}{K}$$

3) Compare b/w AI, ML, DL!

AI  $\Rightarrow$  It's Categorized into 4 Categories

$\hookrightarrow$  thinking humanly, thinking rationally

$\hookrightarrow$  Acting humanly, Acting rationally

most scientists consider it as acting rationally

ML  $\Rightarrow$  Subfield of AI, that uses statistical methods to determine some outputs, corresponding to some law

DL  $\Rightarrow$  subfield of AI, where we use multi-layered

NN  
if U trained a NN, U got 54% train accuracy  
and 51% valid accuracy; what will U do next?

- \* TRY USING BIGGER Network and/or Train longer
- \* Check for the accuracy again and keep changing until reaching an acceptable accuracy.

However,

- \* if the validation accuracy, wasn't improved,  
try using more data and/or regularization
- \* check again for both accuracy and keep changing  
until reaching acceptable accuracies

However,

- \* if nothing improves, we may search for better  
NN architectures

— [ML Recipe]

7) Explain with examples the Linear Perceptron  
Update rule [Widrow-Hoff] ← Also the main Algo

8) Initialize the weights and bias randomly.

9) Present the feature vector of the  $m^{th}$  train example  $\underline{u}(m)$   
and its desired output  $d(m) = \begin{cases} 1, & u(m) \in C_1 \\ 0, & u(m) \in C_2 \end{cases}$

10) Calculate the actual output for  $m \rightarrow y(m) = f(\underline{w}^T \underline{u}(m))$

11) Weights are adapted to the [Widrow-Hoff] Rule.  
 $\underline{w}(\text{new}) = \underline{w}(\text{old}) + \alpha [d(m) - y(m)] \underline{u}(m)$  ↗ learning rate

12) Return to [10] until, all patterns are classified correctly

Example: let  $d(m) = 1$ , and  $y(m) = f(\underline{w}^T \underline{u}(m)) = 0$

so  $\underline{w}^T \underline{u}(m) < 0 \rightarrow$  but, we want  $y(m) = 1$

So, we need to make  $\underline{w}^T \underline{u}(m)$  more +ve.

Widrow-Hoff →  $\underline{w}(\text{new}) = \underline{w}(\text{old}) + \alpha [d(m) - y(m)] \underline{u}(m)$

$$\begin{aligned} y(\text{new}) &= f(\underline{w}^T(\text{new}) \cdot \underline{u}(m)) = f(\underline{w}^T \underline{u}(m) + \alpha [d(m) - y(\text{old})] \underline{u}(m)) \\ &= f(-\text{ve} + \alpha [1-0] \|\underline{u}(m)\|^2) = f(+\text{ve}) = 1 \end{aligned}$$

and vice versa!

On example | نتائج ملحوظة \*  
هو الذي يتحقق بالفعل Algo

## ⑧ Explain the Naïve Estimator, Compare it to histogram Analysis!

\* It's the same as histogram analysis, but instead of partitioning the feature space into pre-specified ranges, we take a range around each  $x$ . So the Naïve estimator is smoother than the histogram analysis.

\* Density estimation  $\hat{P}(x) = \frac{\text{no of Points in } [x-h/2, x+h/2]}{\text{total no of Points} \times Mh}$

\* DIS →  $\hat{P}(x)$  is discontinuous  
Adv → All data points are treated equally, regardless of their distance to  $x$  (estimation point).

⑨ Kernel Density → used to get the density of data by [Parzen Window] using Bump function for each point, then it gets the estimation density through the average of the sum.

$$q_n(x) = \frac{1}{M} g\left(\frac{x}{h}\right)$$

→ it's continuous and gives weights for all the points; Hence better than Naïve estimator.

$$\boxed{\text{bump}} \otimes \boxed{\text{PDF}} = \boxed{\text{sum}}$$

\* if Kernelwidth(h) → too small → too bumpy  
too large → lost details

$$h_{opt} = \frac{1}{N} \sum H_i$$

$$\rightarrow H_i = \tilde{G} \left[ \frac{4}{M(N+2)} \right]^{1/N+4}$$

1) Why ReLU is used instead of Sigmoid?

- \* Because, ReLU results in a faster learning.
- \* Sigmoid has fast saturation
- \* ReLU → used with Sigmoid → output layer in Emojis
- \* binary classification

2) Explain the least SE classifier [No Proof Needed]

- examples      weights
- \* In some cases, if  $M > N+1 \rightarrow$  we can't get exact equality  
Eq unknown to all the weights!
  - \* We will define an error func  $\rightarrow E = \sum_{m=1}^M (\underline{w}^T \underline{u}(m) - b_m)^2$   
actual      desired
  - \* We want to find the weights that will lead to the least error i.e.  $\rightarrow$  solve  $\frac{\partial E}{\partial w} = 0$ . Here, we can converge if the problem is not linearly separable.
  - \*  $E = [\underline{Y}\underline{w} - b]^T [\underline{Y}\underline{w} - b]$   
 $= \underline{w}^T \underline{Y}^T \underline{Y} \underline{w} - 2b^T \underline{Y} \underline{w} + b^T b \xrightarrow{\frac{\partial E}{\partial w}} 2\underline{Y}^T \underline{Y} \underline{w} - 2\underline{Y}^T b = 0$   
 $[\underline{w} = (\underline{Y}^T \underline{Y})^{-1} \underline{Y}^T b] \leftarrow$

3) Major Issues in GMM?

- \* Initialization
  - GMM is an iterative algo, that's very sensitive to initial conditions
  - we usually, use K-means to get a good initialization
- \* No of Gaussian Components
  - Try diff no. of comp., choose the best based on a validation set.

- 13) Where are the types of feature removed in Feature Selection method?
- Irrelevant Features, that don't help in discriminating bet. classes
  - Correlated Features, that are very close to each other, so, can be removed without much loss in information

14) Explain Regularization and why it's used?

- \* It's a technique used to prevent model overfitting
- It simplifies the network and makes it smaller.

$$\text{L1 Regularization} \Rightarrow J = \frac{1}{M} \sum_{m=1}^M E_m + \frac{\lambda}{2M} \|\underline{w}\|_1 \rightarrow \sum_j |w_j|$$

$$\text{L2 Regularization} \Rightarrow J = \frac{1}{M} \sum_{m=1}^M E_m + \frac{\lambda}{2M} \|\underline{w}\|_2^2 \rightarrow \underline{w}^T \underline{w}$$

- Dropout Reg. → Eliminate some nodes based on a probability
- 15) How Bayes Rule Can be used in Classification?

→ We calculate  $[P(x|C_i) P(C_i)]$  for each class and classify  $x$  for the class that gives the highest value

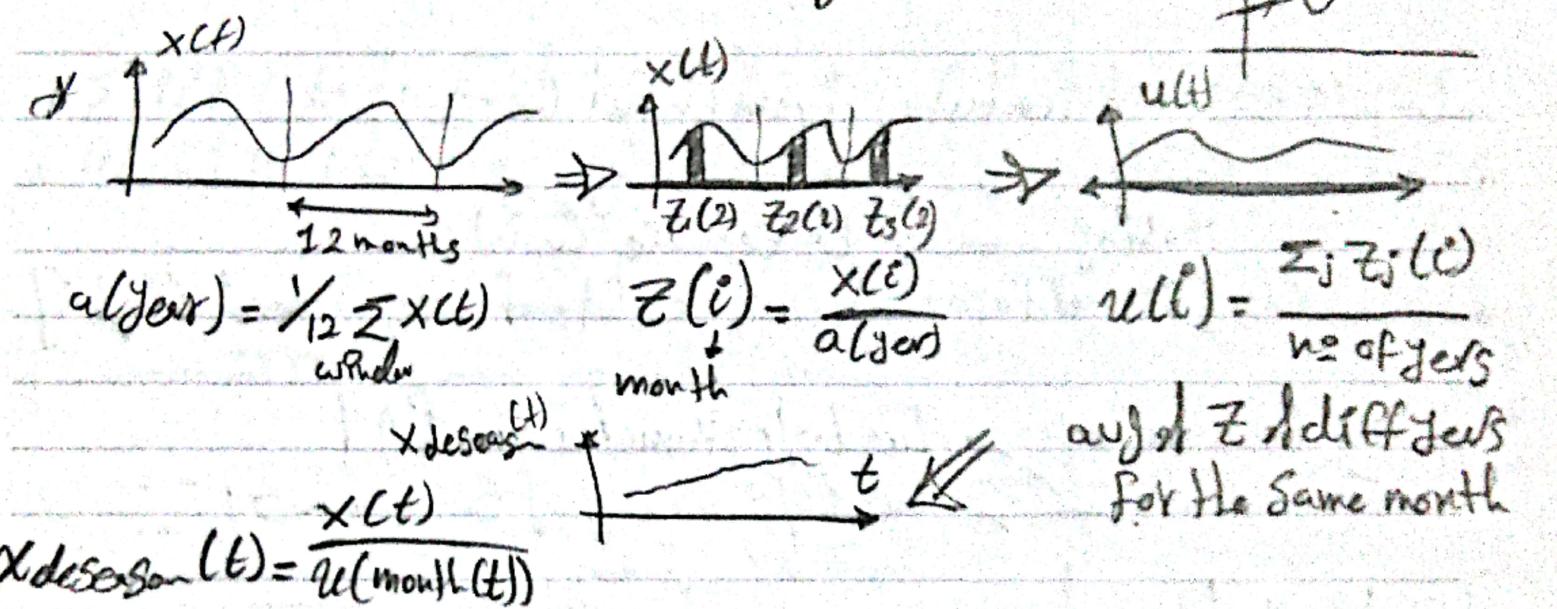
16) Drawbacks of Back Propagation?

→ Can be slow in reaching the min. [adjust  $\alpha$ ]

→ prone to get stuck in local minima  
 (local  $\rightarrow$  global)

## 17] What's the Importance of deSeasonalization?

\* Used to Remove Seasonal Periodicities. Here, we can predict correctly.



## 18] Three Methods to update weights?

\* Gradient Descent → Use Full Set of Training Examples  
 Adv: Consistent Optimization  
 Dis: Slow

\* Mini-Batch → SubSet of Training Examples  
 Adv: Fast

\* Stochastic  
 GD → One Training Example at a time  
 Adv: Faster than GD  
 Dis: Hard to Converge  
 I less Speedup from Vectorization

Final 2019

Same as 2018, but  
those were added

① Describe how CNNs work, why they have less footprints?

\* layers extract features from input images.

\* it uses Conv layers → for filtering

Pooling → to reduce inputs.

\* it leaves [filter parameters & weights of fully connected layers]

It has less memory → due to [Parameter Sharing]

Footprints → due to [Sparsity of Connections]

② 3-limitations of DL?

→ Needs large amount of labeled data

→ Datasets may be biased

→ Sensitive to std adversarial attacks

→ Over sensitive to changes in context.

→ Visual understanding is tricky

→ Unintended results from fitness functions

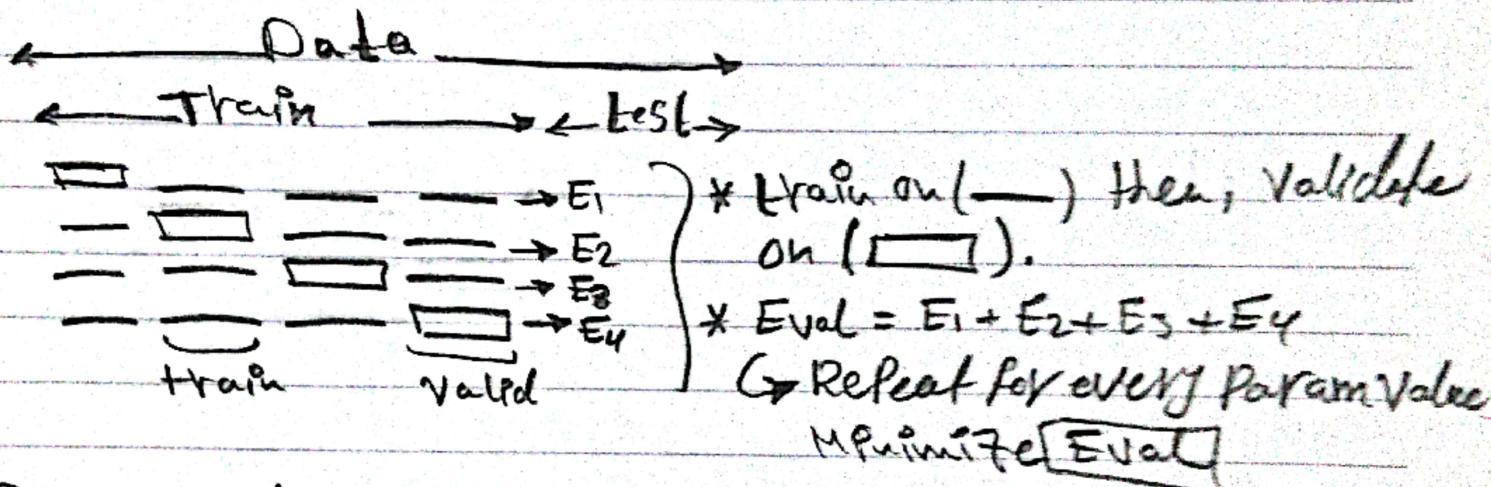
→ Not a magic tool [Not able to build general-intelligent]

→ cannot fit all real-world scenarios [machine]

3) Explain K-fold Validation, In what context is it used?

- \* Used for parameter tuning over the Train Data
- \* Better than convention [Train-Valid-Test Split]

- \* Split Into Train & test (No Need for validation)
- \* Apply K-fold to the Train set (usually  $K=5$ )



18) Steps of Back Propagation? [Alg]

i) Initialize all the weights  $\rightarrow$  small rand values [-1, 1]

ii) let  $u(m)$  &  $d(m)$  be the train input & output example

iii) for  $m = 1$  to  $M$

    → Present  $u(m)$  to the network  $\rightarrow$  Compute outputs

    → Using them  $\rightarrow$  calc  $\frac{\partial E}{\partial w}$  of each layer

    → Update the weights  $\rightarrow w_{new}^{(l)} = w_{old}^{(l)} - \alpha \frac{\partial E}{\partial w^{(l)}}$

iv) Compute the total error  $\rightarrow$  Stop in case of convergence