

Week 13: (9 - 9-5)

Combining Models

Learning objectives: • Advantages of combining models

- Ways to combine models
- Explain decision tree learning
- Explain adaptive boosting (adaboost)

Questions?:

1. Why to combine models?

Multiple available, none perfect

→ Many features extracted from given pattern

2 Issues:

- How many classifiers are needed?

Complementary properties among different models and features

- What kind of model should be used?

- Features used in classifier?

- How to combine results from different classifiers?

Stories of success:

Million dollar prize; Netflix top submissions combine models.

Data mining competitions: Classification problems,

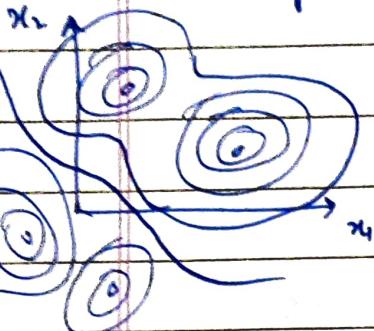
Winning team employ ensemble of classifiers

Ensemble or Averaging multiple models
↳ Gives global picture

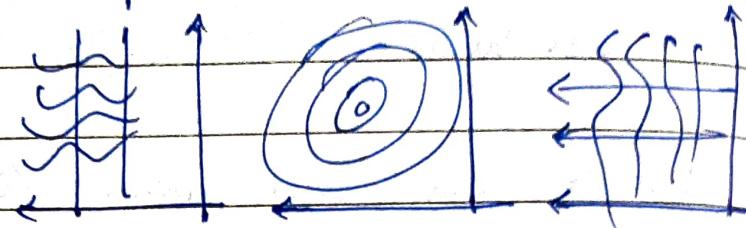
PAGE No.

DATE / / /

Ex.) Combining two 1D Gaussians



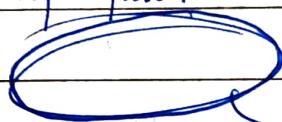
Combined decision boundary for ?
multiple models needed



Introduction to hypothesis class:

Space of functions : Set of functions modelled by 1 model

$h_0(x)$ richer set may be a superset of this

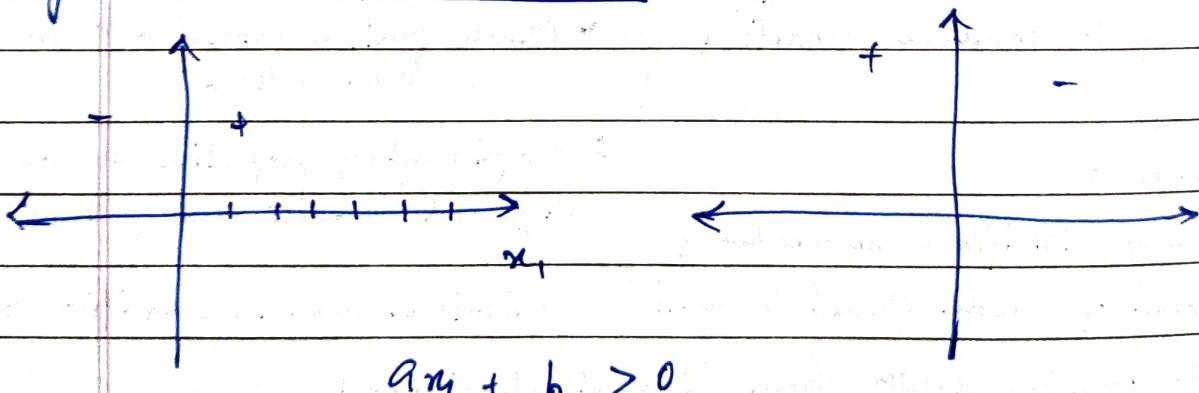


Comes close but not exactly models data

Eg. $h_0(x)$ all linear func. $\Rightarrow g_0(x)$ all polynomials

$$\Rightarrow g_0(x) \supset h_0(x)$$

Hypothesis class of a threshold:



$$ax_1 + b > 0$$

Only one variable can form threshold :

$$\text{Eg. } x_1 < \omega_1 \text{ or } x_1 > \omega_2$$

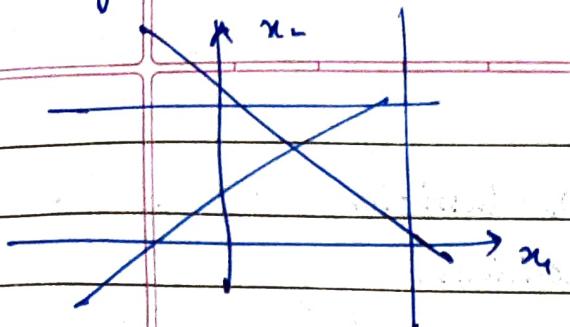
$$\text{or } x_1 < \omega_1 \text{ or } x_1 > \omega_2$$

Only line || to one of the axis

(No slant threshold boundary allowed)

Hypothesis class of a linear classifier.

PAGE NO.	
DATE	/ /



$$w^T x + b > 0$$

kind of decision thresholds

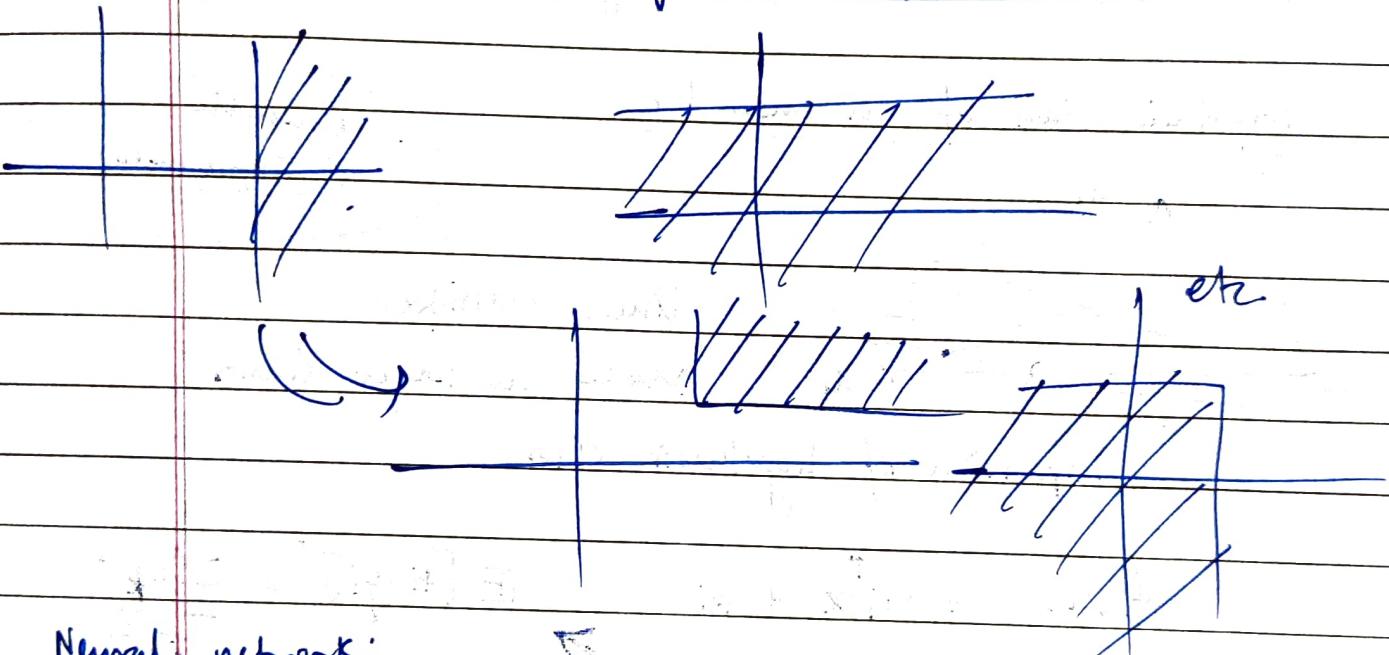
$$\text{loss} + \frac{\lambda}{2} \|w\|^2$$

linear classifiers

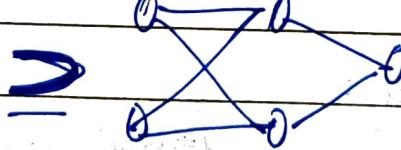
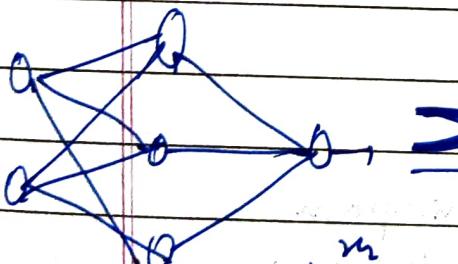
thresholds

reduces hypothesis class

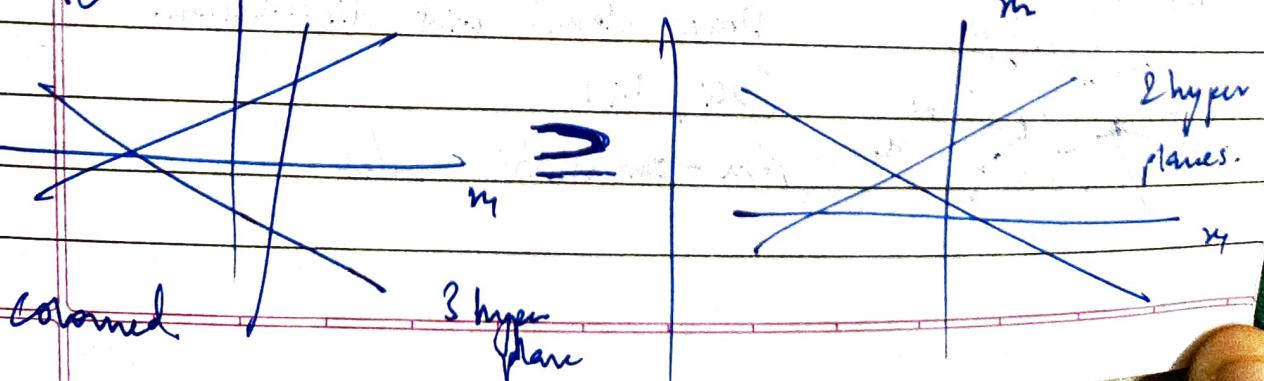
Nested hypothesis classes of single & multiple variables



Neural network:



Can get 2 hidden NN from 3 one

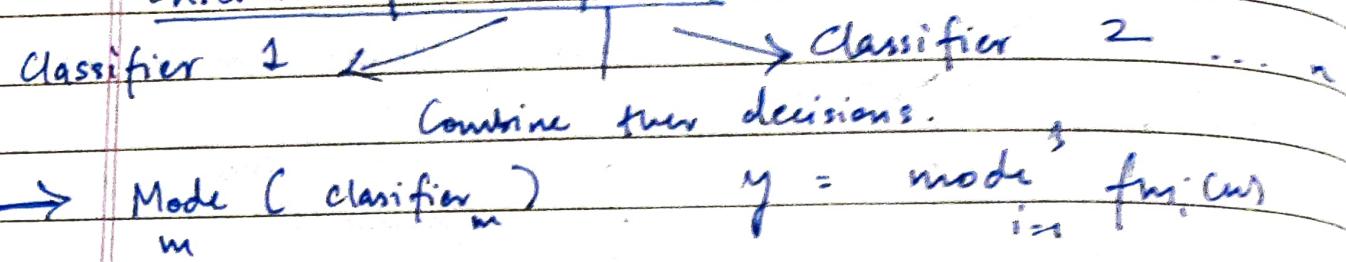


2 hyper planes.

(9-2) Ensembles :

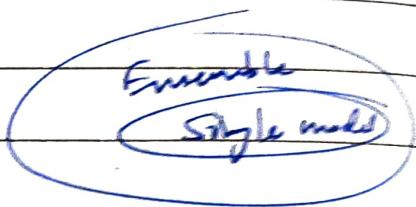
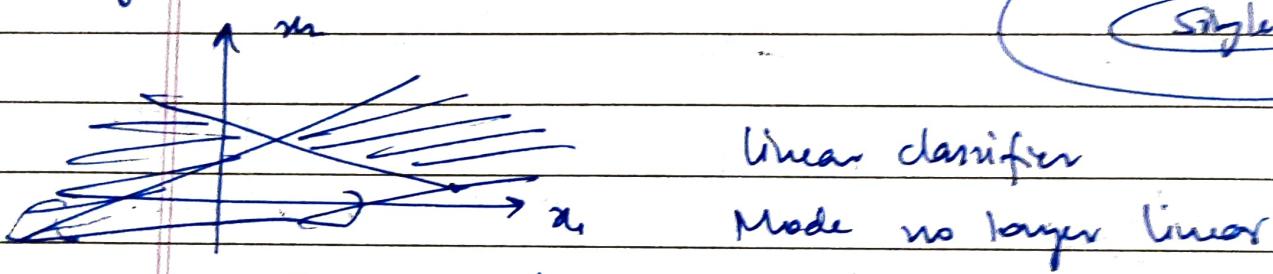
PAGE NO.	
DATE	/ / / /

Ensembles for classification



Work if they are independent of each other. (Assumption of independence)
Binomial distribution. (Obviously)

Hypothesis class of an ensemble:

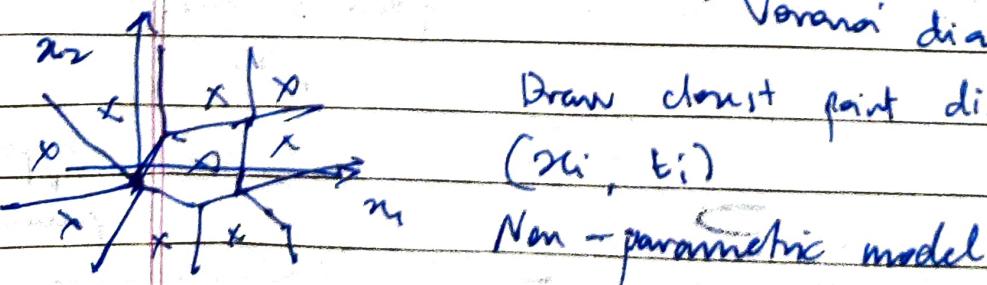


Expansion of hypothesis class

$$f_i(x) = w_i^T x + b_i \quad \mathbb{E} [t_i - y_i]^2 = \frac{\sum \sigma_k^2}{K^2}$$

$$\therefore y_i = \frac{\sum f_i(x)}{K}$$

Nearest neighbor model:



Combine K modes by:-

K-Nearest Neighbor model

If some value is wrong, it is corrected & averaged out.
Hypothesis class have a smoothing & self-correcting effect

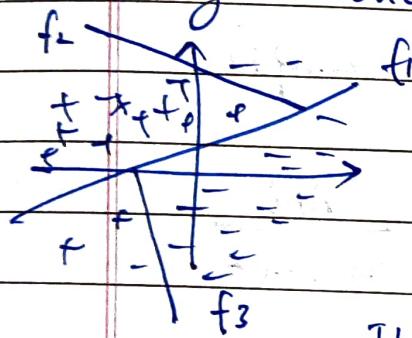
(9-3-1) Decision Trees 1

Voting

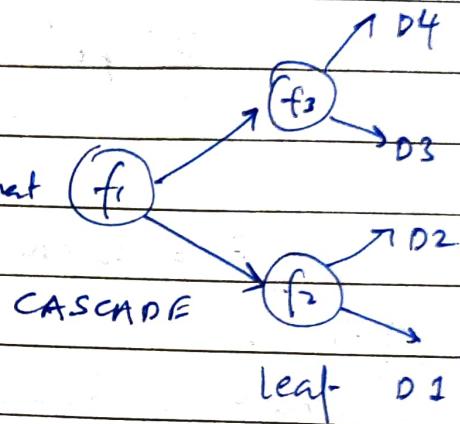
$$f_1 \quad f_2 \quad f_3 = \text{Ensemble}$$

In decision tree, we first find f_i such that

Each classifier may be weak but
combined may be better.



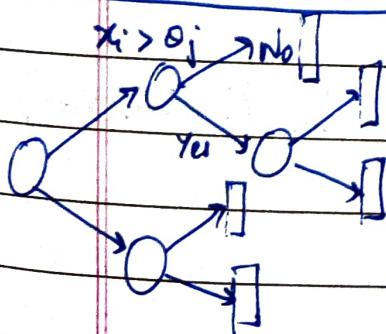
This is how cascade separates the classes.



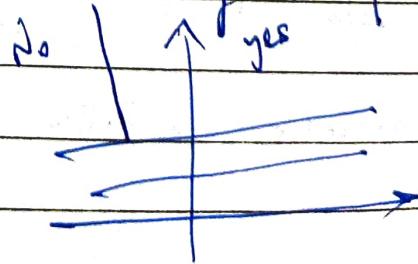
It is a tree as last decision nodes = leaf
& earlier levels are internal nodes.

- Need not be balanced, need not be binary either

Threshold-based decision tree:

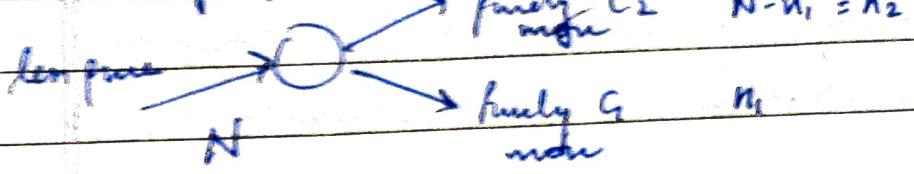


> Binary decision at every node of tree



Greedy algorithm for node splitting

- Purity of branches to be maximized.



Predominance of one class over the other classes.

$$\max_c (p_c) \quad \text{probability of class } c \quad c \in \{1, \dots, C\}$$

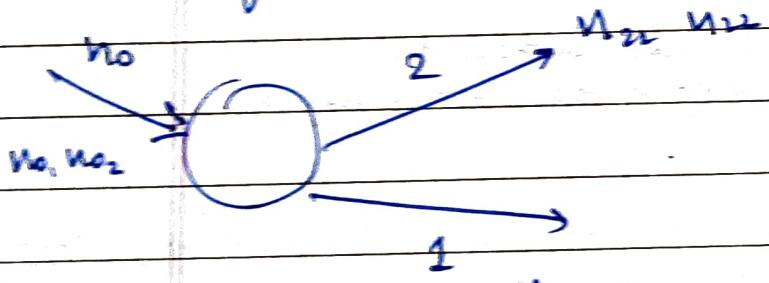
$$\sum_c p_c = 1 \quad p_c \geq 0 \quad (\text{imarity})$$

$$- \sum p_c \log(p_c) \quad \text{Entropy minimized?}$$

$$p_c = n_c / n \quad \text{Or sum } \sum p_c (1-p_c) \quad \text{Gini Index}$$

(g-3-2)

Greedy : Balanced purity (Need this concept)



$$p_{oi} = \frac{n_{oi}}{n_{o1} + n_{o2}}$$

$$p_{21} = \frac{n_{21}}{n_2 + n_{22}}$$

p_{11} & p_{21} need to encode MORE purity.

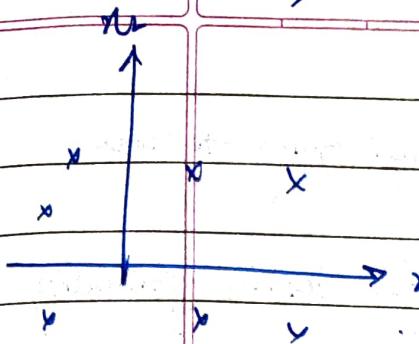
$$p_{11} = \frac{n_1}{n_1 + n_2}$$

$$\Rightarrow \underset{\text{node}}{H_1 H(p_{11}, p_{21})} + \underset{\text{node}}{H_2 H(p_{11}, p_{21})} < \underset{\text{node}}{H_0 H(n_1, n_2)}$$

Max decrease \Rightarrow Better decision update.

We have to look at all possible x, α relations to create thresholds

PAGE No.	
DATE	/ /



D features Gaps between N points = $N-1$

$\Rightarrow D \times (N-1)$ decisions

made at any node

$$D(N-1) \rightarrow D(N_1-1) + D(N_2-1)$$

Algorithm:

- Start at root
- Find ideal x, α relation to separate/max.
- If pure, label & terminate
- Recurse over childern node.

Issue: Tree gets deeper such that nodes $\propto N$ (samples)

Too many axis aligned boundaries. \Rightarrow Not a good hypothesis class

Depth of tree \propto Overfitting \Rightarrow Pruning needed.

Validation data can be used, etc.

(9-4) Combination of trees, Model randomization, Random Forests

Tree 1 + Tree 2 + Tree 3 / Average / mode

How to find threshold such that \Downarrow Ensemble of trees
to have min. mean square error on children \equiv Random Forest
= Ensemble of trees.

Randomizing Models

Non convex initialization

Variable subset selection

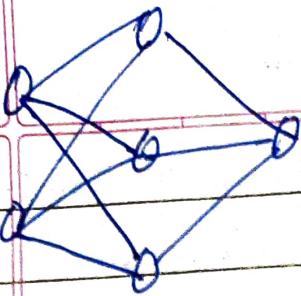
Training subset selection

\Rightarrow Bagging

Random initialization of Neural Network

PAGE No.

DATE / / /

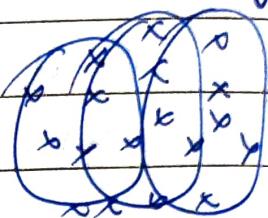


At every node we are not going to use all D variables

→ Randomly select $d < D$ variables for testing optimality

Random forest idea

Random training subsets

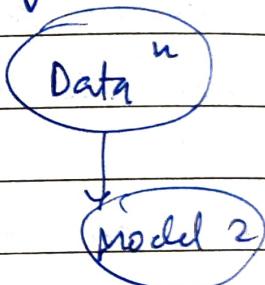


(9-5) Boosting:

Sequential refinement via boosting

Training data → Model 1

Weight ↑
↓ error
change of missclassified



Adaboost:

- uniform weight for all data points
- At each round
 - Bootstrap sample based on tree weights
 - Train classifier & apply to original train set
 - Wrongly classified \Rightarrow Wt increased
 - Correctly classified \Rightarrow Wt decreased

Error $> 50\%$ start over

- Final prediction is weighted avg of all classifiers with weights representing the training accuracy.

PAGE NO.	
DATE	/ /

Adaboost visualized:

S is training data

$$\text{Initialize } D^{(0)} = \left[\frac{1}{n}, \dots, \frac{1}{n} \right]$$

for t in $1 \dots T$ / m in $1 \dots M$

↳ Learn t^{th} WL on $(D^{(t)}, S)$

(Weighted loss)

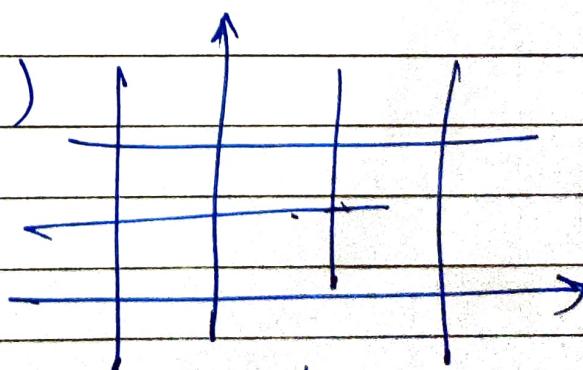
$$E_t = \sum_i D_i^{(t)} \mathbb{1}\{y_i \neq t_i\} \quad | \quad w_t = \frac{1}{2} \log [1/E_{t-1}]$$

All training points, recompute wts st we focus on better perf on uncorrected points

$$\Rightarrow D_i^{(t+1)} = \frac{D_i^{(t)}}{\sum_j D_j^{(t)} \exp(-w_t t_j y_j)}$$

Proceed for T steps.

$$\text{Sign} \left(\sum_{m=1}^M w_m y_m^{(t)} \right)$$



Adaboost :

Viola Jones face detector

Using Haar filters

Many learners combine & give ensemble