# Combining Multiple Learners
## Part 1

*Ethem Chp. 15*
*Hastie Chp 8*
*Haykin Chp. 7, pp. 351-370*

# *Overview*

- **Introduction**
  - □ Motivation

- **Static structures:** the responses of several *experts* (individual networks) are combined in a way that **does not** involve the input signal.
  - □ *Ensemble averaging / Majority Voting*
  - □ *Boosting*
  - □ *Stacking*
  - □ *Error orrecting Output Codes*

- **Dynamic structures**: the input signal actuates the mechanism that combines the responses of the experts.
  - □ *Mixture of experts*
  - □ *Hierarchical mixture of experts*

# *Motivation*

- When designing a learning machine, we generally make some choices:
  - parameters of machine, training data, representation, etc…

- This implies some sort of variance in performance

- **Why not keep all machines and combine their predictions?**

- Intuition: Combining experts opinions, votes,...

- **Ensemble Learning / Classifier Combination:**
  - do not learn a single classifier but learn a set of base classifiers
  - combine the predictions of multiple classifiers

- **Hope to….:**
  - reduce variance: results are less dependent on peculiarities of a single training set
  - reduce bias: a combination of multiple classifiers may learn a more expressive concept class than a single classifier
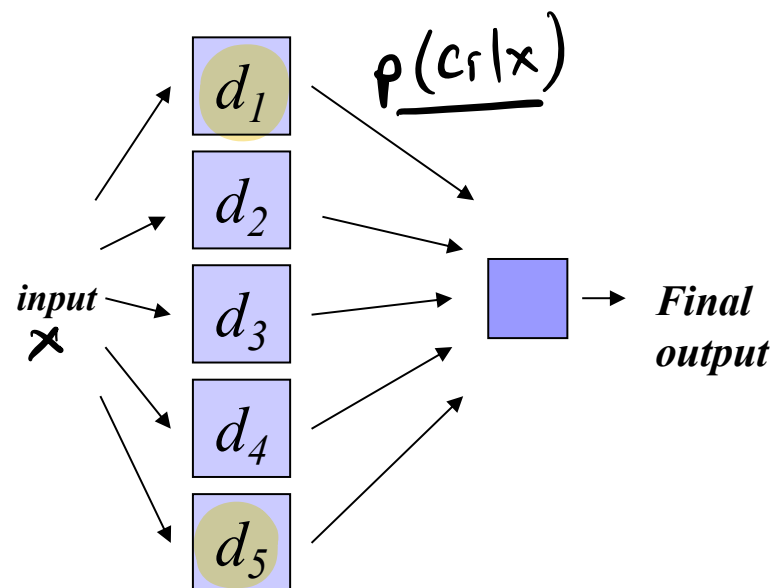
- **Important:**
  - formation of an ensemble of diverse classifiers from a single training set

# *Rationale*

- No Free Lunch thm: "There is no algorithm that induces the most accurate learner in any domain, all the time."
  - http://www.no-free-lunch.org/

- Generate a group of base-learners which when combined has higher accuracy

- Different learners use different
  - Algorithms: making different assumptions
  - Hyperparameters: e.g number of hidden nodes in NN, k in k-NN
  - Representations: diff. features, multiple sources of information
  - Training sets: small variations in the sets or diff. subproblems

# *Reasons to Combine Learning Machines*

**Terminology:**
Ensemble formed
from base classifiers

$$\frac{d_1}{} \quad p(c_r|x)$$

input $x$ → $d_1$, $d_2$, $d_3$, $d_4$, $d_5$ → *Final output*

Lots of different combination methods:

Most popular are *ensemble averaging* and *majority voting*.

$$\frac{1}{L} \sum_{J} p(c_{kJ}|x)$$

# *Ensemble Averaging*

- Regression

$$y = \sum_{j=1}^{L} w_j d_j$$

$$w_j \geq 0 \quad \text{and} \quad \sum_{j=1}^{L} w_j = 1$$

$d_j$ is the output of the jth base classifier in regression

$d_{ji}$ is the output of the jth base classifier for ith class in classification

- Classification

$$y_i = \sum_{j=1}^{L} w_j d_{ji}$$

*inversely*

$w_j$ often proportional to error rate of classifier:

Learned over a validation set

# *Ensemble Averaging*

**The output of the ensemble is the (weighted) average of the base classifiers` outputs.**
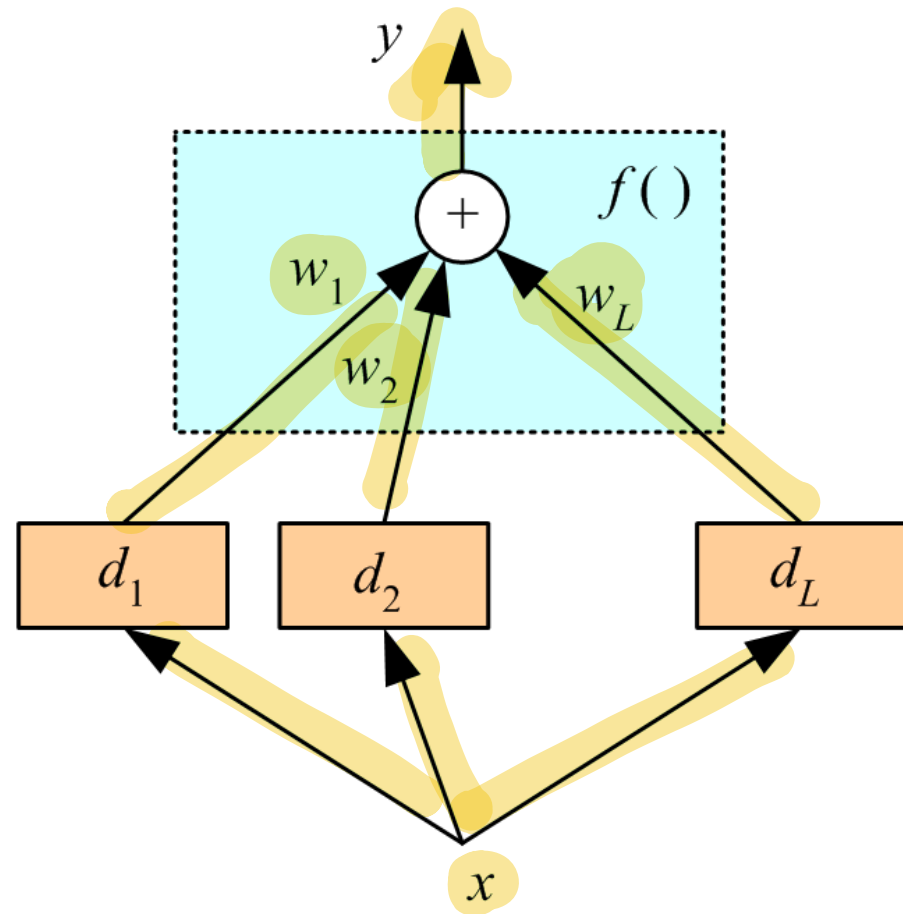
- Regression

$$y = \sum_{j=1}^{L} w_j d_j$$

$$w_j \geq 0 \text{ and } \sum_{j=1}^{L} w_j = 1$$

- Classification

$$y_i = \sum_{j=1}^{L} w_j d_{ji}$$

# Ensemble Averaging

If we use a committee machine f$_{com}$ whose output is the average:

$$f_{com} = \frac{1}{M} \sum_{i=1}^{i=M} d_i \quad \text{base learners}$$

*Error of combination is <u>guaranteed</u> to be lower than the <u>average</u> error:*

ensemble

$$(f_{com} - t)^2 = \frac{1}{M} \sum_i \underbrace{(d_i - t)^2} - \frac{1}{M} \sum_i \underbrace{(d_i - f_{com})^2}$$

average error

always positive

*(Krogh & Vedelsby 1995)*

$Avg(e_{d_1}, e_{d_2}, e_{d_3})$

$\geq e_{f_{comm}}$.

- Similarly, we can show that if $d_j$ are iid:

$$Bias^2 : \quad (E_D(d) - f)^2$$
$$Variance : E_D[(E_D(d) - d)^2]$$

$$E\left[f_{com}\right] = E\left[\sum_j \frac{1}{L} d_j\right] = \frac{1}{L} L \cdot E[d_j] = E[d_j]$$

$$Var\left(f_{com}\right) = Var\left(\sum_j \frac{1}{L} d_j\right) = \frac{1}{L^2} Var\left(\sum_j d_j\right) = \frac{1}{L^2} L \cdot Var(d_j) = \frac{1}{L} Var(d_j)$$
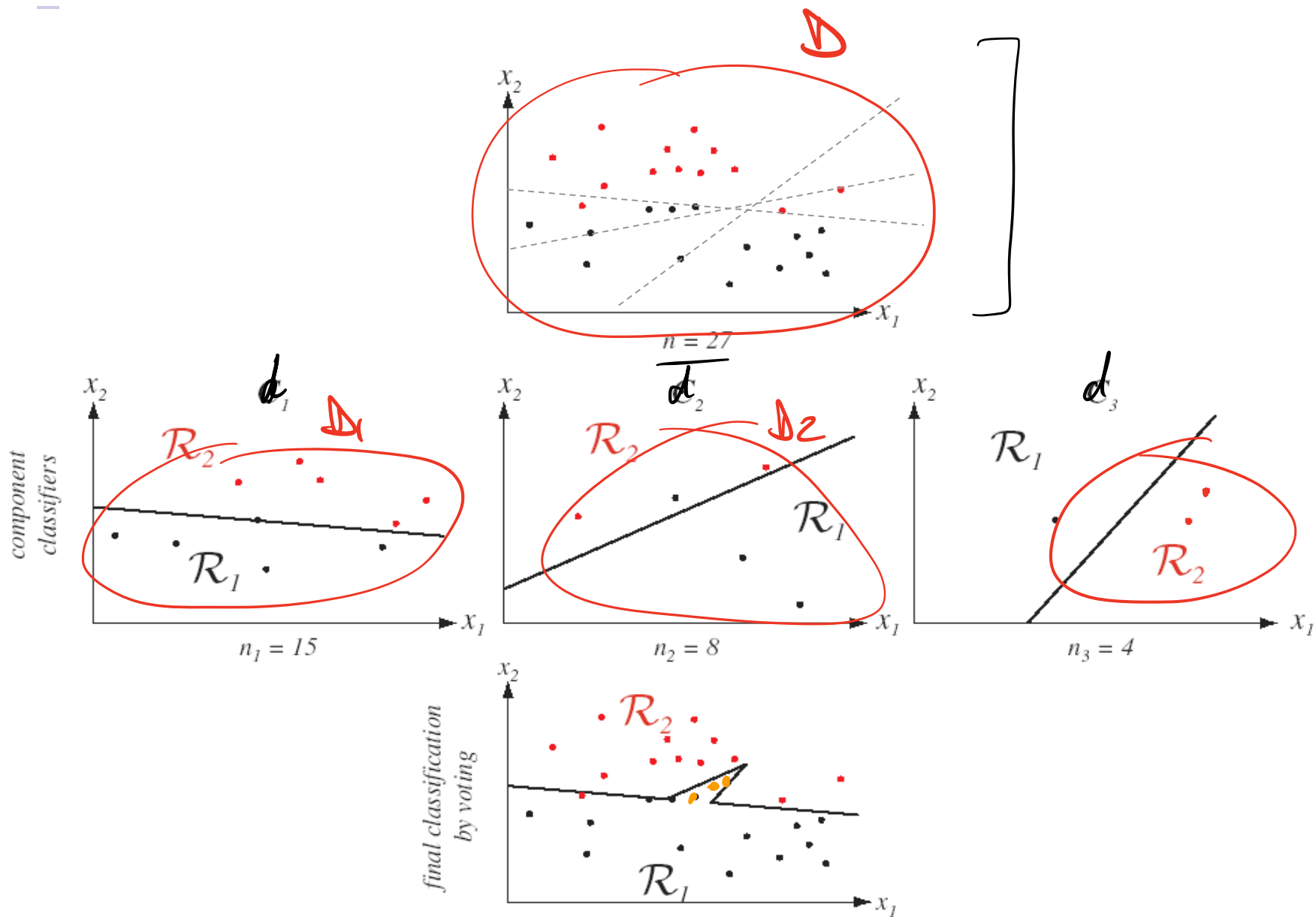
- Bias does not change, variance decreases by 1/$L$

# *Ensemble Averaging*

What we can exploit from this fact:

- Combine multiple experts with the same bias and variance, using ensemble-averaging
  - □ the bias of the ensemble-averaged system would the **same as the bias of one of the individual experts**
  - □ the variance of the ensemble-averaged system would be **less than the variance of one of the individual experts.**

➢ We can purposefully use complex/flexible models, the variance will be reduced due to averaging.

component classifiers

$x_2$    $D$

$n = 27$

$d_1$    $D1$

$\mathcal{R}_2$

$\mathcal{R}_1$

$n_1 = 15$

$d_2$    $D2$

$\mathcal{R}_2$

$\mathcal{R}_1$

$n_2 = 8$

$d_3$

$\mathcal{R}_1$

$\mathcal{R}_2$

$n_3 = 4$

final classification by voting

$\mathcal{R}_2$

$\mathcal{R}_1$

*Lecture Notes for E Alpaydın 2004 Introduction to Machine Learning © The MIT Press (V1.1)*

# *Majority Voting*

**In classification, majority voting chooses the class that gets the most votes from the base classifiers** (not averaging the output probabilities).

- Majority voting is also considered a form of ensemble averaging (**hard vs soft averaging**).
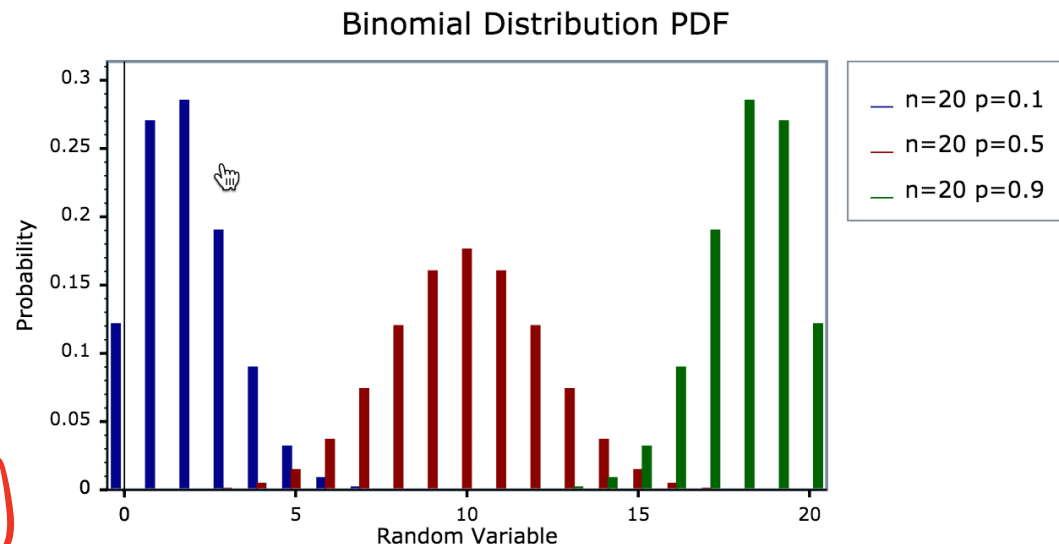
*averaging posterior probs* $\left\{ P(C_k \mid x) \right\}_J$ *scores*

*ensemble "*

Majority voting ensemble makes an error when more than half of the *base* classifiers make a mistake. What is this probability?

*confidences*

| | Class A | Class B |
|-----|---------|---------|
| $d_1$ | 1 | 0 |
| $d_2$ | 1 | 0 |
| $d_3$ | 0 | 1 |

Hint:

*ens. = 1 (chooses Class A w/ majority vote (hard averaging))*

*fcomm*



**Binomial Distribution PDF**

Legend:
- n=20 p=0.1
- n=20 p=0.5
- n=20 p=0.9

(x-axis: Random Variable, y-axis: Probability)
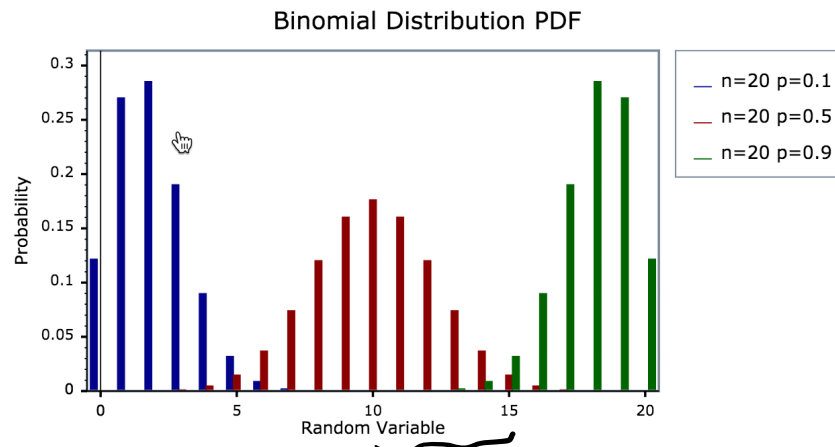
If the base classifiers are all independent, the probability of error of the **majority voting ensemble** is:

$$\text{ensemble} \quad P(error) = \sum_{k=\frac{N}{2}+1}^{N} \binom{N}{k} p^k (1-p)^{N-k}$$

$P$ : probability that $d_i$ makes an error.

In prev. example, sum goes from $k=2$ to $k=3$.

**Binomial Distribution PDF**
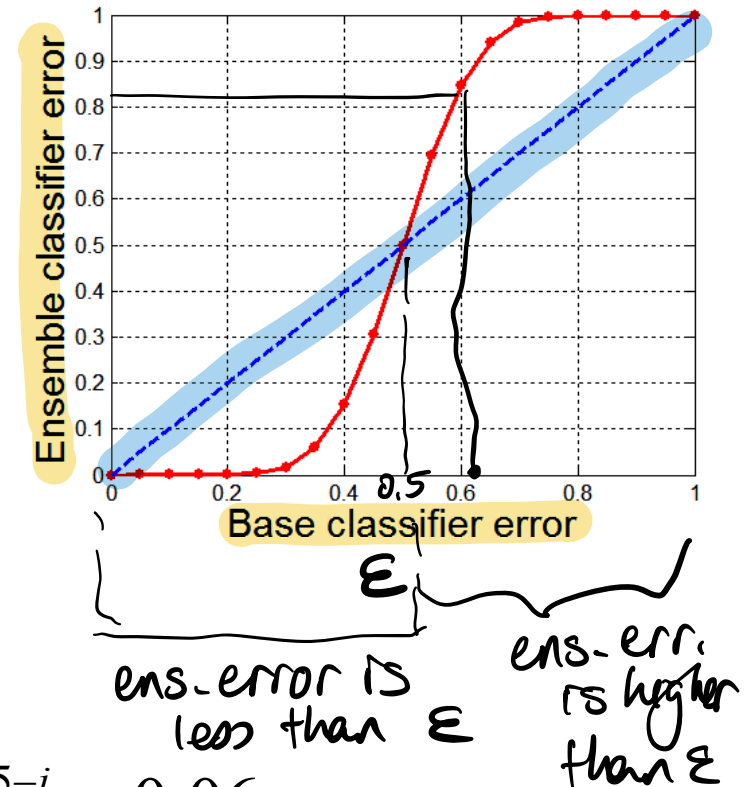


14

# *Example*

- Suppose there are **25 base classifiers**
  - Each classifier has error rate, $\varepsilon = 0.35$
  - Assume classifiers are independent
    - i.e., probability that a classifier makes a mistake does not depend on whether other classifiers made a mistake
    - **Note:** in practice they are not independent!
- Probability that the ensemble classifier makes a wrong prediction
  - The ensemble makes a wrong prediction if the majority of the classifiers makes a wrong prediction
  - The probability that 13 or more classifiers err is

$$\sum_{i=13}^{25} \binom{25}{i} \varepsilon^i (1-\varepsilon)^{25-i} \approx 0.06 \ll \varepsilon$$

# Why Majority Voting works?

- Suppose there are 25 base classifier, and
  - Each classifier has error rate, $\varepsilon = 0.35$
  - Errors made by classifiers are uncorrelated

- Probability that the ensemble classifier makes a wrong prediction:



Ensemble classifier error (y-axis), Base classifier error $\varepsilon$ (x-axis)

ens. error is less than $\varepsilon$

ens. err. is higher than $\varepsilon$

$$P(X \geq 13) = \sum_{i=13}^{25} \binom{25}{i} \varepsilon^i (1-\varepsilon)^{25-i} = 0.06$$

Note: Ensemble does not help when base classifier error is high!
(where red curve is higher than the blue line)

- **We want the base learners to be complementary**
  - What if they were all the same or very similar?
    - No use...

- **Reasonably accurate, but not necessarily very accurate**

# *Overview*

- **Introduction**
  - ☐ Rationale
- **Combination Methods**
  - ☐ Static Structures
    - Ensemble averaging / Majority Voting
    - **Bagging**
    - Boosting
    - Error Correcting Output Codes

  - ☐ Dynamic structures
    - Mixture of Experts
    - Hierarchical Mixture of Experts

# Ensemble Methods > *Bagging*

**Voting method where base-learners are made different by training over slightly different training sets.**

Bagging (Bootstrap Aggregating) - Breiman, 1996
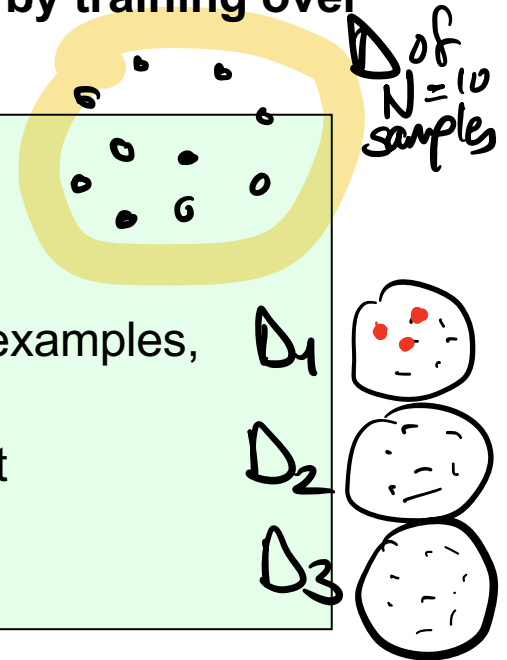
*take a training set D, of size N*

*for each network / tree / k-nn / etc...* base learner

    - build a new training set by sampling N examples,

        randomly with replacement, from D

    - train your machine with the new dataset

*end for*

*output is average/vote from all machines trained*

*D of N=10 samples*

$D_1$   $D_2$   $D_3$

- ☐ Resulting base-learners are similar because they are drawn from the same original sample

- ☐ Resulting base-learners are slightly different due to chance

- Generate new training sets using sampling with replacement

| Original Data | 1 | 2 | 3 | 4 | 5 | 6 | 7 | 8 | 9 | 10 |
|---|---|---|---|---|---|---|---|---|---|---|
| Bagging (Round 1) | 7 | 8 | 10 | 8 | 2 | 5 | 10 | 10 | 5 | 9 |
| Bagging (Round 2) | 1 | 4 | 9 | 1 | 2 | 3 | 2 | 7 | 3 | 2 |
| Bagging (Round 3) | 1 | 8 | 5 | 10 | 5 | 5 | 9 | 6 | 3 | 7 |

- some examples may appear in more than one set
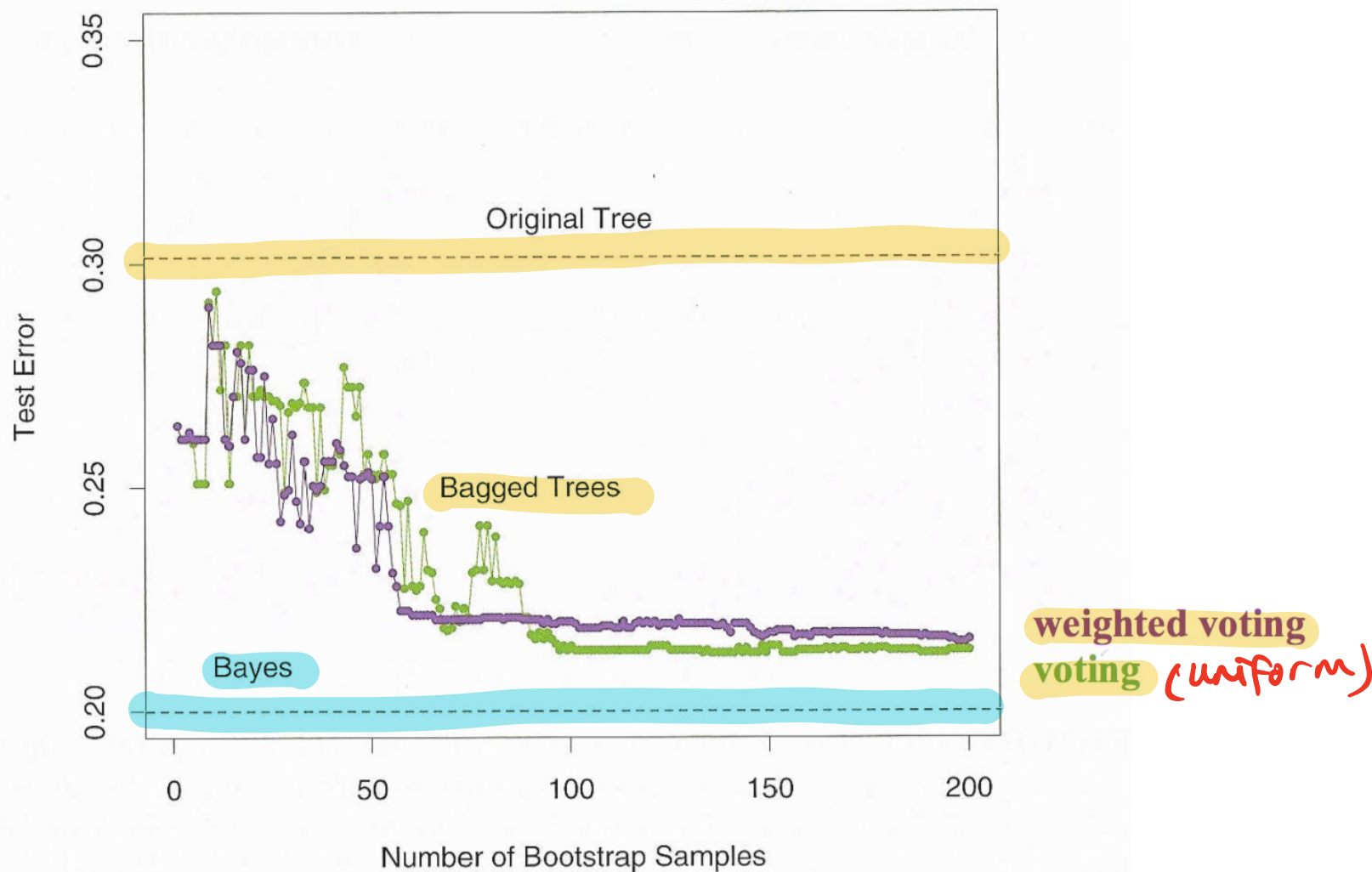- for each set, the probability that a given example doesn't appear in it is

$$\Pr(x \notin D_i) = \left(1 - \frac{1}{n}\right)^n \rightarrow 0.3678$$

- i.e., less than 2/3 of the examples appear in one bootstrap sample

# *Bagging*

- **Not all data points will be used for training**
  - ☐ Waste of training set
    - Each sample has a probability of 0.37 of not being selected in any one bootsrap training set.
    - I.e. only about 2/3rd of the samples are used in any one bootsrap sample.
    - They can be used for testing (see Out-of-Bag error in Random Forests)

- **Bagging is suitable for unstable learning algorithms**
  - ☐ Unstable algorithms change significantly due to small changes in the data.
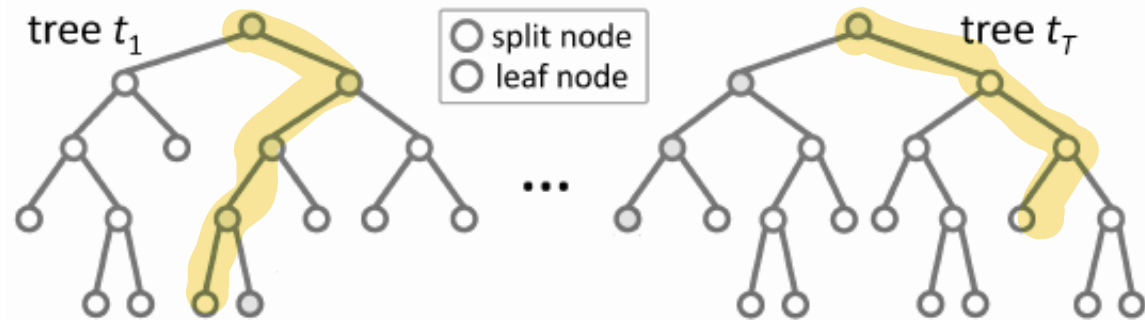  - ☐ Such as MLPs, decision trees

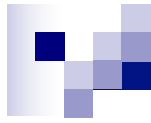from Hastie, Tibshirani, Friedman: The Elements of Statistical Learning, Springer Verlag 2001

# Random Forests

# Random Forests



- An improved method over simple bagged trees
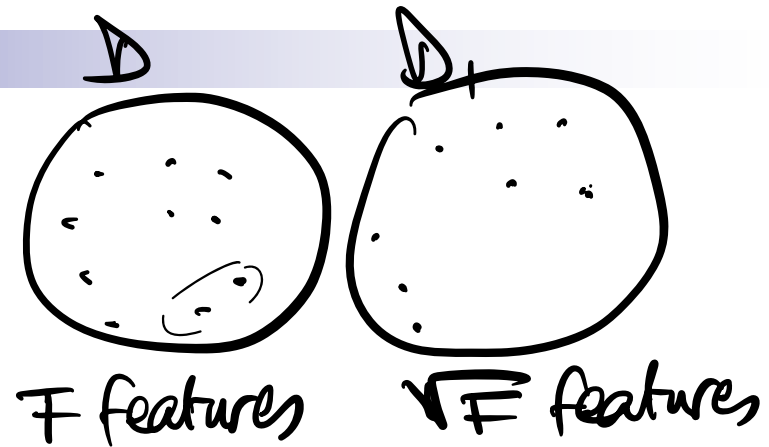- All trees vote to produce a final answer.

# *Motivation*

- **With decision trees, it is found that optimal cut points can strongly depend on the training set used.**

  - ☐ This suggested using multiple trees and use voting to combine the results.

- Averaging the outputs of the trees reduces overfitting.

- For the use of multiple trees to be most effective, the trees should  be independent as possible.

  - ☐ Splitting using a random subset of features hopefully achieves this.

  - ☐ If the trees really are independent, the performance should improve with more trees

## *The Random Forests Algorithm*

//Random forest with k trees

F features          √F features

Given a training set S

- For  i = 1 to k do:
- Build subset $S_i$ by sampling with replacement from S
  Learn tree $T_i$ from $S_i$ as:
  - **At each node:**
    - Choose best split from random subset of F features
  Each tree grows to the largest extend, no pruning

- Make predictions according to **majority vote** of the set of k trees.

- **Bagging**

Generate randomized training sets by sampling with replacement from the full training set (bootstrap sampling)

Full training set $\quad$ $D_1 \ D_2 \ D_3 \ D_4 \ D_5 \ D_6 \ D_7 \ D_8 \ D_9 \ D_{10} \ D_{11} \ D_{12}$

Random "bag" $\quad$ $D_4 \ D_9 \ D_3 \ D_4 \ D_{12} \ D_{10} \ D_{10} \ D_7 \ D_3 \ D_1 \ D_6 \ D_1$

- **Feature subset selection**

Choose different random subsets of the full feature vector to generate each tree

Full feature vector $\quad$ $f_1 \ f_2 \ f_3 \ f_4 \ f_5 \ f_6 \ f_7 \ f_8 \ f_9 \ f_{10}$ $\qquad F$

Feature subset $\quad$ $f_4 \ f_6 \ f_7 \ f_9 \ f_{10}$ $\qquad \sqrt{F}$

features

- Typically 5 – 100 trees are used.  Often only a few trees are needed.

- Results seem fairly insensitive to the number of random attributes that are tested for each split.  A common default is to use the square root of the number of attributes.

- Trees are fast to generate because fewer attributes have to be tested for each split and no pruning is needed.

# Out-of-Bag Error

- ## To grow one tree
  - Bootstrap sample set from learning set L
  - Remaining samples are called **out-of-bag samples**
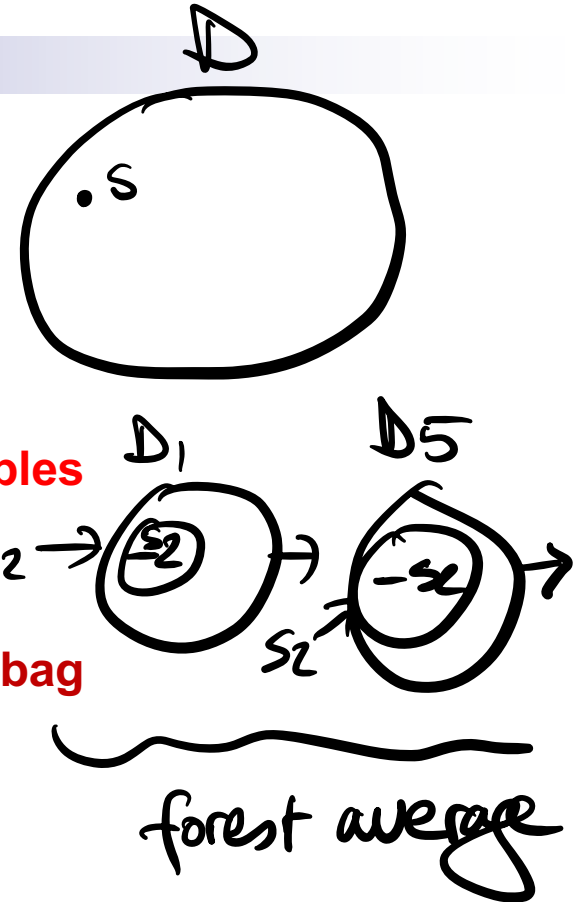
- ## For each sample S of the learning set
  - **Look for all the trees for which S was out-of-bag**
  - Build the corresponding sub-forest
  - Predict the class of S with it
  - Measure error on S

- ## **Out-of-bag error** = average over all samples of S
  - Predictions not made using the whole forest... but with some aggregation
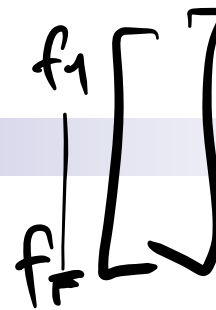
- ## Provides an estimation of the generalization error
  - Can be used to decide when to stop adding trees to the forest

## *Features of Random Forests*

- It gives some of the best accuracy among current algorithms.
- It is efficient.
- It generates an internal unbiased estimate of the generalization error as the forest is built.

Plus…

- It handles missing data effectively.
  - A tree doesn't use all the features, plus other methods when a feature value is unknown…

- It gives estimates of what variables are important in the classification.
  - Impurity beased importance (average info gain from that feature over all trees)
  - Permutation importance

- It has methods for balancing error in class population unbalanced data sets.
  - Bootstrap sampling introduces variation in class proportions, or one can use a sample with equal class priors