# Rensselaer

# Lecture 13: Bagging, Random Forests; Boosting (1)

Dr. Chengjiang Long

Computer Vision Researcher at Kitware Inc.

Adjunct Professor at RPI.

Email: **longc3@rpi.edu**

# Progress

Attendence

2.5% 2.5%

**62.5%**

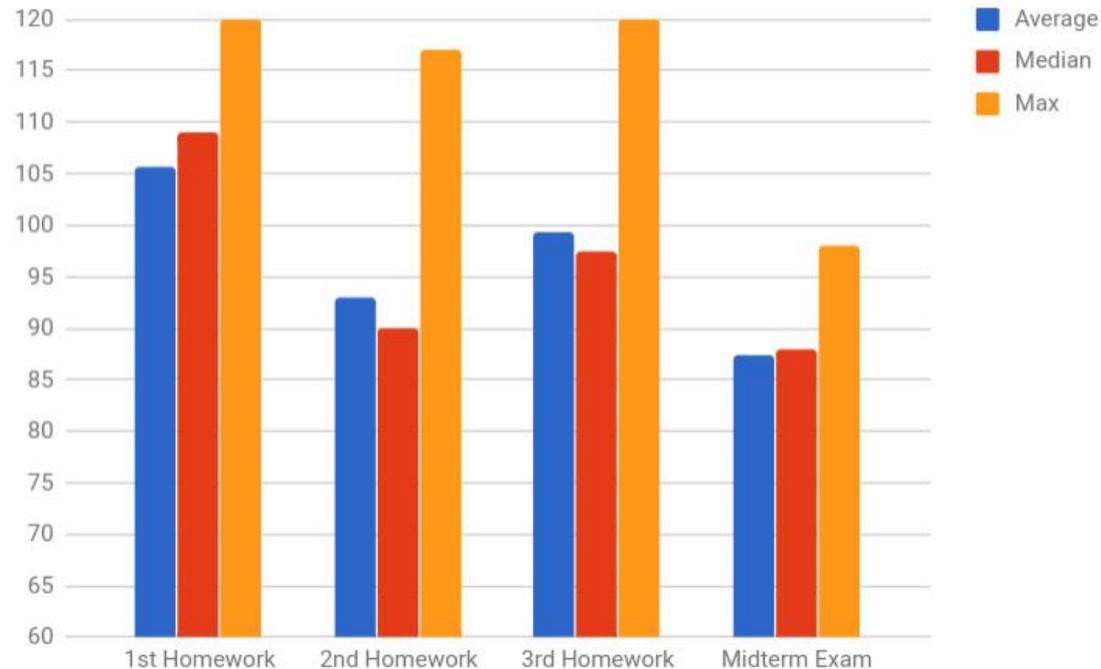Homework

15% 5%

Midterm & Final Exams

20% 20%

Final Project

35%

# Progress

| | Average | Meidan | Max |
|---|---|---|---|
| 1st Homework | 105.7 | 109 | 120 |
| 2nd Homework | 93 | 90 | 117 |
| 3rd Homework | 99.4 | 97.5 | 120 |
| Midterm Exam | 87.4 | 88 | 98 |

# Final Project

- The final project is an open-ended assignment, with the goal of gaining experience applying the techniques presented in class to real-world datasets.



MINST Dataset



LFW dataset

# Final Project

- Students should work either individually or in groups of 2-3. By 11:59 p.m. on March 23, you will submit a 1–3 paragraph proposal on your project.
- You must do this in order to get full credit for your project, and you must get my approval on it before presentating project proposals to the class on Mar 30, and starting work on your project.
- The proposal presentation should describe the problem you are solving, what data is being used, the proposed technique you are applying in addition to what baseline is used to compare against.
- The final project report should be at least 4 pages and is due on April 30, as well as the data and code.

# Final Project - Timeline

- 03/23/2018: 1-3 paragraph proposal for approval.
- 03/30/2018: proposal presentation to class.
- 04/30/2018: submit the data, code and project report.
- 05/01/2018: oral presentation to class.

# Outline
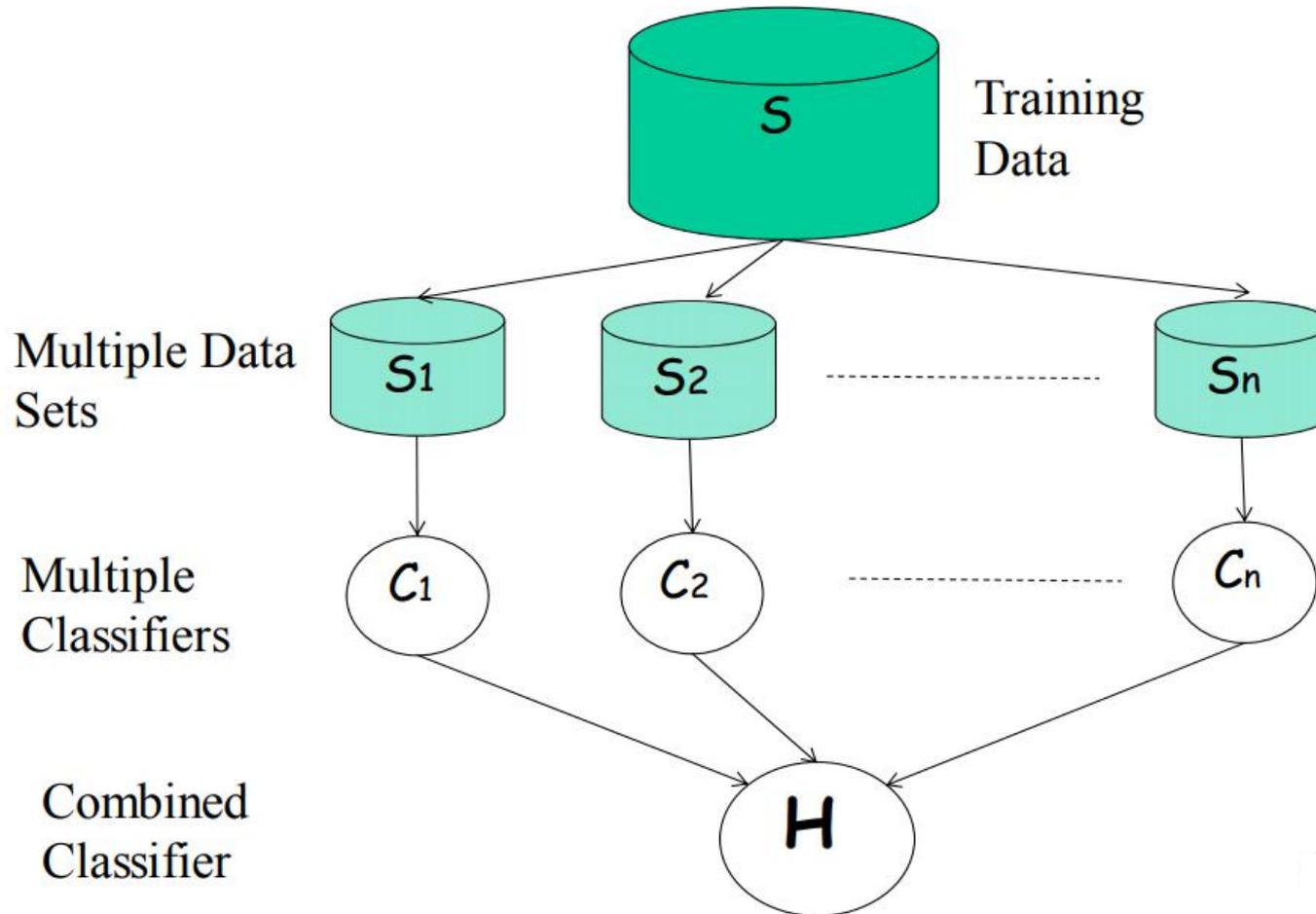
- Bagging

- Random Forest

- Boosting

# Outline

- **Bagging**

- Random Forest

- Boosting

# Ensemble Methods

- Bagging (Breiman $1994$,…)
- Random forests (Breiman $2001$,…)
- Boosting (Freund and Schapire $1995$, Friedman et al. $1998$,…)

- Predict class label for unseen data by aggregating a set of predictions (classifiers learned from the training data)

# General Idea



Training Data — $S$

Multiple Data Sets — $S_1$, $S_2$, ---------------------, $S_n$

Multiple Classifiers — $C_1$, $C_2$, --------------------- $C_n$

Combined Classifier — $H$

# Ensemble Classifiers

- Basic idea: Build different "experts" and let them vote
- Advantages:
  - Improve predictive performance
  - Different types of classifiers can be directly included
  - Easy to implement
  - Not too much parameter tuning
- Disadvantages:
  - The combined classifier is not transparent (black box)
  - Not a compact representation

# Why do they work?

- Suppose there are $25$ base classifiers
- Each classifier has error rate $\varepsilon = 0.35$
- Assume independence among classifiers
- Probability that the ensemble classifier makes a wrong prediction:

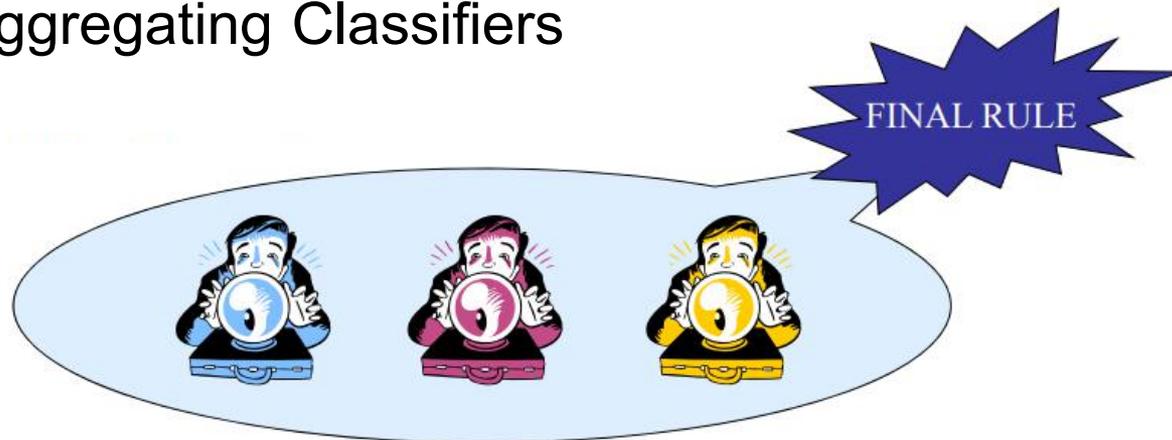$$\sum_{i=13}^{25} \binom{25}{i} \varepsilon^i (1-\varepsilon)^{25-i} = 0.06$$

# Bagging

- Training
  - □ Given a dataset S, at each iteration i, a training set $S_i$ is sampled with replacement from S (i.e. bootstraping)
  - □ A classifier $C_i$ is learned for each $S_i$
- Classification: given an unseen sample X
  - ❑ Each classifier $C_i$ returns its class prediction
  - ❑ The bagged classifier H counts the votes and assigns the class with the most votes to X

# Bagging

- Bagging works because it reduces variance by voting/averaging

  ❑ In some pathological hypothetical situations the overall error might increase

  ❑ Usually, the more classifiers the better

- Problem: we only have one dataset

- Solution: generate new ones of size n by bootstrapping, i.e. sampling with replacement

- Can help a lot if data is noisy

# Aggregating Classifiers

- Bagging and Boosting

  ❑ Aggregating Classifiers

  FINAL RULE

- Breiman $(1996)$ found gains in accuracy by aggregating predictors built from reweighed versions of the learning set

# Bagging

- Bagging = Bootstrap Aggregating
- Reweighing of the learning sets is done by drawing at random with replacement from the learning sets
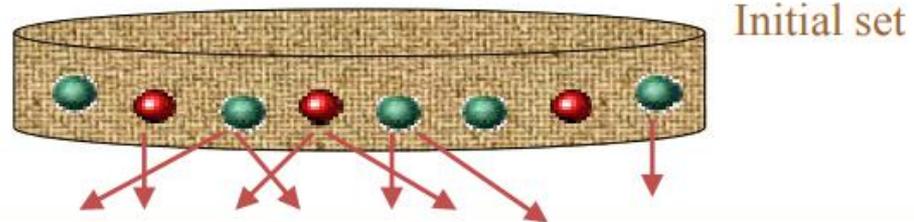- Predictors are aggregated by plurality voting

# The Bagging Algorithm

- B bootstrap samples
- From which we derive:

  - $B$ Classifiers $\in \{-1, 1\}: c^1, c^2, c^3, ..., c^B$

  - $B$ Estimated probabilities $\in [0,1]: p^1, p^2, p^3, ..., p^B$

- The aggregate classifier becomes:

$$c_{bag}(x) = sign\left(\frac{1}{B}\sum_{b=1}^{B} c^b(x)\right) \quad \text{or} \quad p_{bag}(x) = \frac{1}{B}\sum_{b=1}^{B} p^b(x)$$

# Example (1)

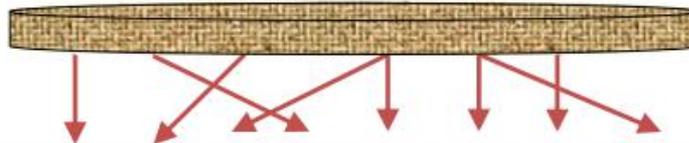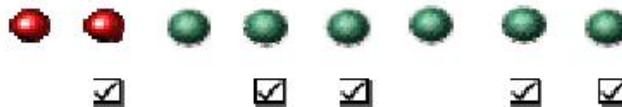# Example (2)



Testing set

Aggregation

Sign

Classifier 1

+

Classifier 2

+

Classifier 3

+

...

+

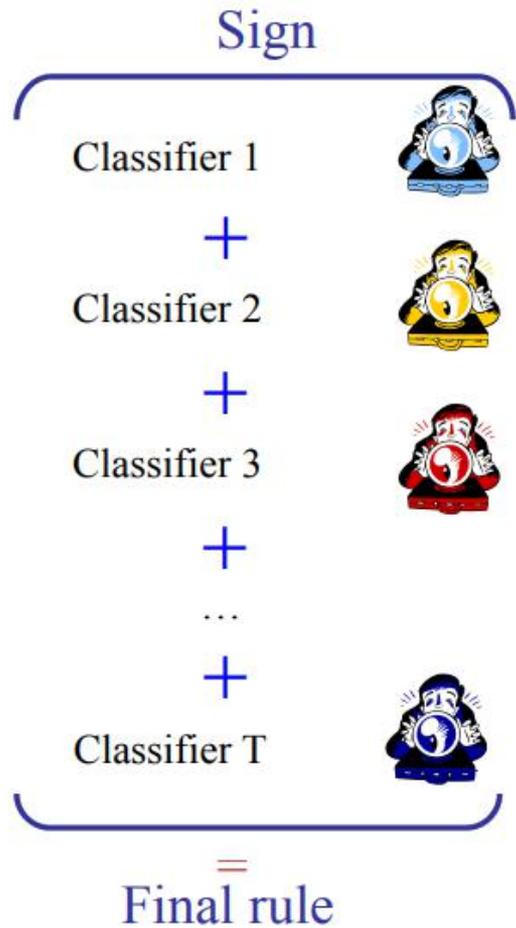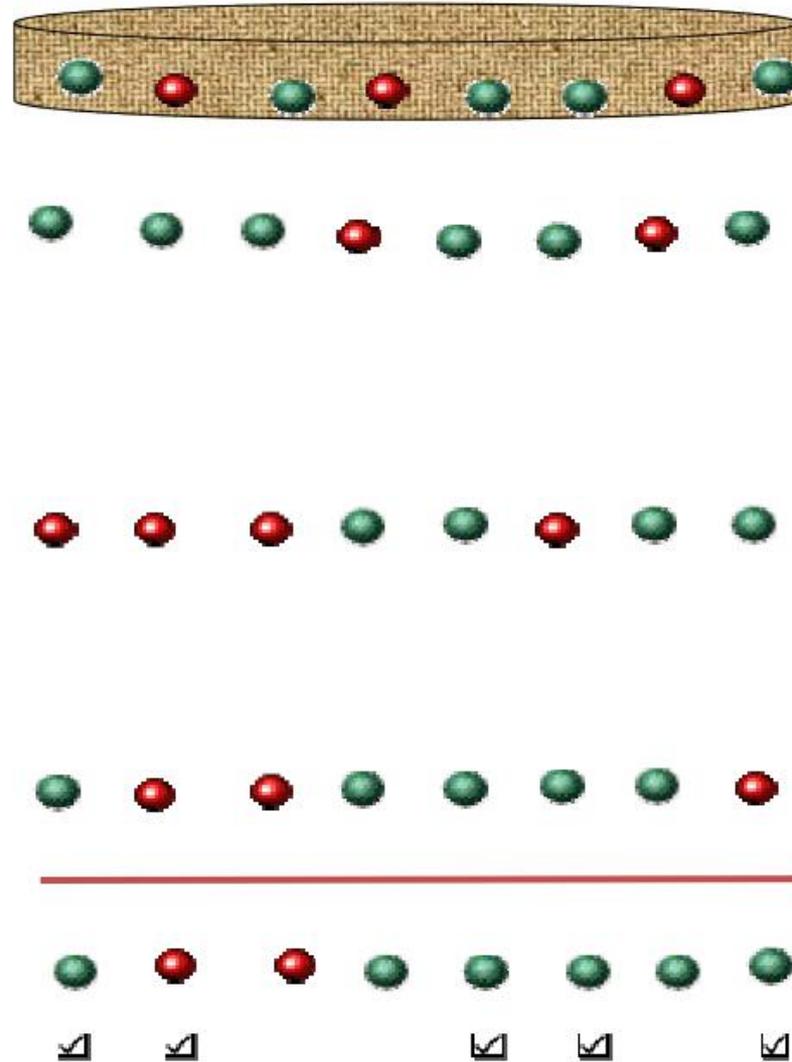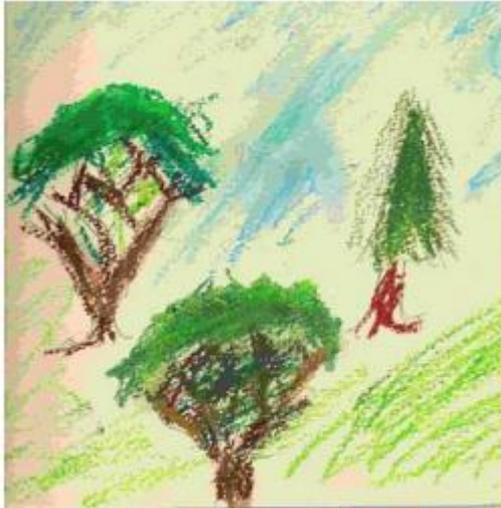Classifier T

=

Final rule

# Outline

- Bagging

- **Random Forest**

- Boosting

# Random Forests

- Breiman L. Random forests. Machine Learning 2001;45(1):5–32

- http://statwww.berkeley.edu/users/breiman/RandomForests/

- For computer vision and medical image analysis



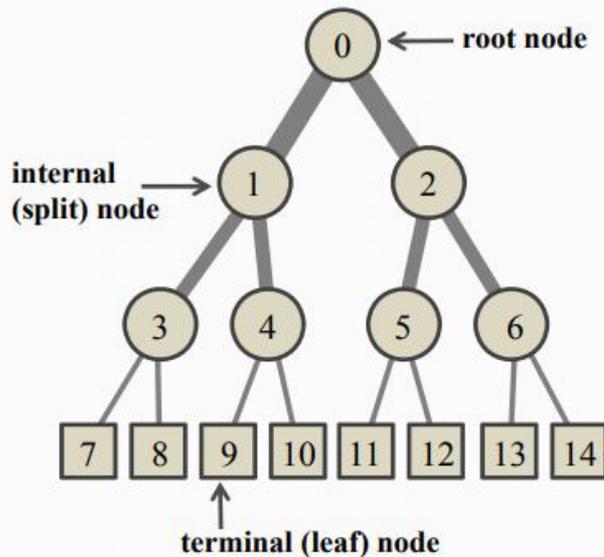A. Criminisi, J. Shotton and E. Konukoglu

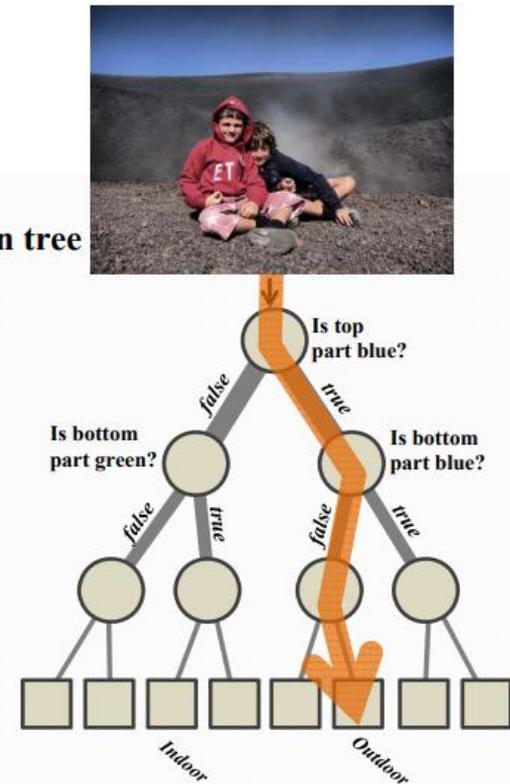http://research.microsoft.com/projects/decisionforests

# Decision Trees and Decision Forests

- A forest is an ensemble of trees. The trees are all slightly different from one another



A general tree structure
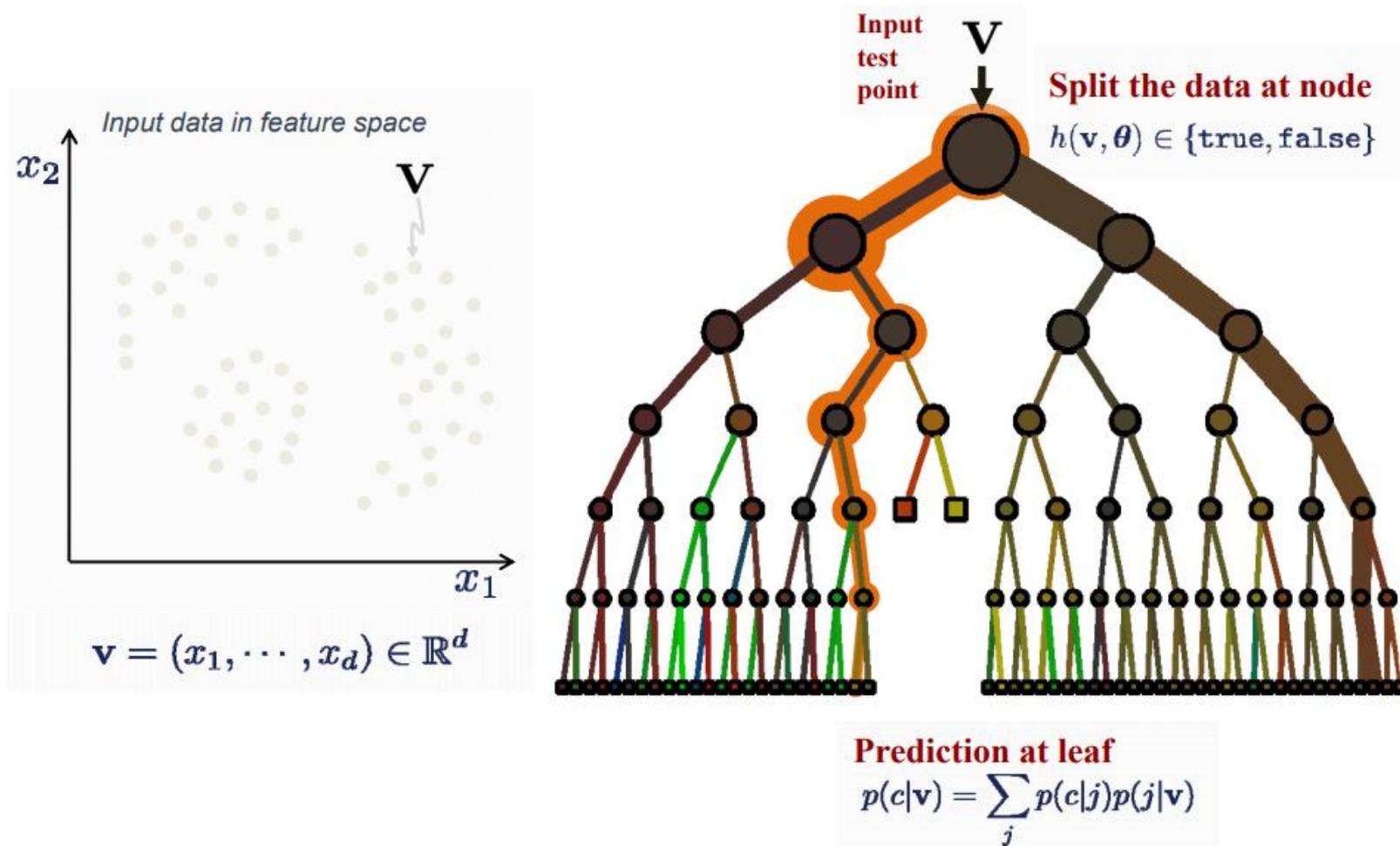
root node

internal (split) node

terminal (leaf) node

A decision tree

Is top part blue?

false    true

Is bottom part green?    Is bottom part blue?

false    true    false    true

Indoor    Outdoor

# Decision Tree Testing (Runtime)



Input data in feature space

$$\mathbf{v} = (x_1, \cdots, x_d) \in \mathbb{R}^d$$

Input test point **V**

Split the data at node

$$h(\mathbf{v}, \boldsymbol{\theta}) \in \{\texttt{true}, \texttt{false}\}$$

**Prediction at leaf**

$$p(c|\mathbf{v}) = \sum_j p(c|j)p(j|\mathbf{v})$$

# Decision Tree Training (Offline)



Input data in feature space

$x_2$

**V**

$\mathbf{v} = (x_1, \cdots, x_d) \in \mathbb{R}^d$

$x_1$

**How to split the data?**

$h(\mathbf{v}, \boldsymbol{\theta}) \in \{\text{true}, \text{false}\}$

$$\boldsymbol{\theta}_j = \arg \max_{\boldsymbol{\theta} \in \mathcal{T}_j} I$$

$I = I(\mathcal{S}_j, \boldsymbol{\theta})$

**Input training data**

$\mathcal{S}_0 = \{\mathbf{v}, c\}$

$\mathcal{S}_1$ $\mathcal{S}_2$

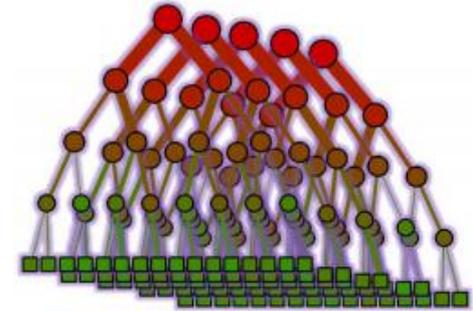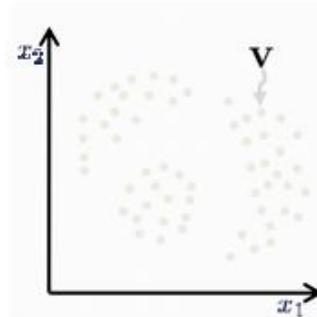$\mathcal{S}_1^{\mathbf{L}}$ $\mathcal{S}_1^{\mathbf{R}}$

**Binary tree? Ternary?**
**How big a tree?**
**What tree structure?**

# Decision Tree Training (Offline)



- How many trees？
- How different？
- How to fuse their outputs？

# Decision Forest Model



## Basic notation

| | | | |
|---|---|---|---|
| Input data point | e.g. | $\mathbf{v} = (x_1, \cdots, x_d) \in \mathbf{R}^d$ | Collection of feature responses $x_d = ?$ |
| Output/label space | e.g. | $\in \{c_k\}$ ? $\mathbb{R}$ ? | Categorical, continuous? |
| Feature response selector | | $\phi$ | Features can be e.g. wavelets? Pixel intensities? Context? |

## Forest model

| | | | |
|---|---|---|---|
| Node test parameters | | $\boldsymbol{\theta} \in \mathcal{T}$ | Parameters related to each split node: i) which features, ii) what geometric primitive, iii) thresholds. |
| Node objective function (train.) | e.g. | $I = I(\mathcal{S}_j, \boldsymbol{\theta})$ | The "energy" to be minimized when training the j-th split node |
| Node weak learner | e.g. | $h(\mathbf{v}, \boldsymbol{\theta}_j) \in \{\text{true}, \text{false}\}$ | The test function for splitting data at a node j. |
| Leaf predictor model | e.g. | $p(c|\mathbf{v})$ | Point estimate? Full distribution? |
| Randomness model (train.) | e.g. 1. Bagging, 2. Randomized node optimization | | How is randomness injected during training? How much? |
| Stopping criteria (train.) | e.g. max tree depth = $D$ | | When to stop growing a tree during training |
| Forest size | | $T$ | Total number of trees in the forest |
| The ensemble model | e.g. | $p(c|\mathbf{v}) = \frac{1}{T}\sum_t^T p_t(c|\mathbf{v})$ | How to compute the forest output from that of individual trees? |

# Random Forest

**Algorithm 1** Random Forest

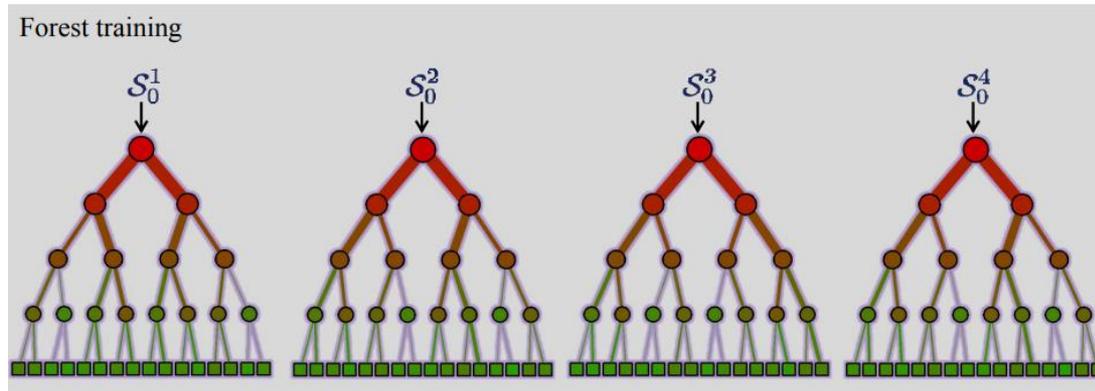**Precondition:** A training set $S := (x_1, y_1), \ldots, (x_n, y_n)$, features $F$, and number of trees in forest $B$.

1 **function** RANDOMFOREST($S$, $F$)
2      $H \leftarrow \emptyset$
3      **for** $i \in 1, \ldots, B$ **do**
4          $S^{(i)} \leftarrow$ A bootstrap sample from $S$
5          $h_i \leftarrow$ RANDOMIZEDTREELEARN($S^{(i)}$, $F$)
6          $H \leftarrow H \cup \{h_i\}$
7      **end for**
8      **return** $H$
9 **end function**
10 **function** RANDOMIZEDTREELEARN($S$, $F$)
11      At each node:
12          $f \leftarrow$ very small subset of $F$
13          Split on best feature in $f$
14      **return** The learned tree
15 **end function**

# Randomness Model

- 1) Bagging (randomizing the training set)

Efficient training



| | |
|---|---|
| $\mathcal{S}_0$ | The full training set |
| $\mathcal{S}_0^t \subset \mathcal{S}_0$ | The randomly sampled subset of training data made available for the tree $t$ |

Forest training

$\mathcal{S}_0^1$      $\mathcal{S}_0^2$      $\mathcal{S}_0^3$      $\mathcal{S}_0^4$

**Precondition:** A training set $S := (x_1, y_1), \ldots, (x_n, y_n)$, features $F$, and number of trees in forest $B$.

1 **function** RANDOMFOREST($S$, $F$)
2     $H \leftarrow \emptyset$
3     **for** $i \in 1, \ldots, B$ **do**
4        $S^{(i)} \leftarrow$ A bootstrap sample from $S$
5        $h_i \leftarrow$ RANDOMIZEDTREELEARN($S^{(i)}$, $F$)
6        $H \leftarrow H \cup \{h_i\}$
7     **end for**
8     **return** $H$
9 **end function**

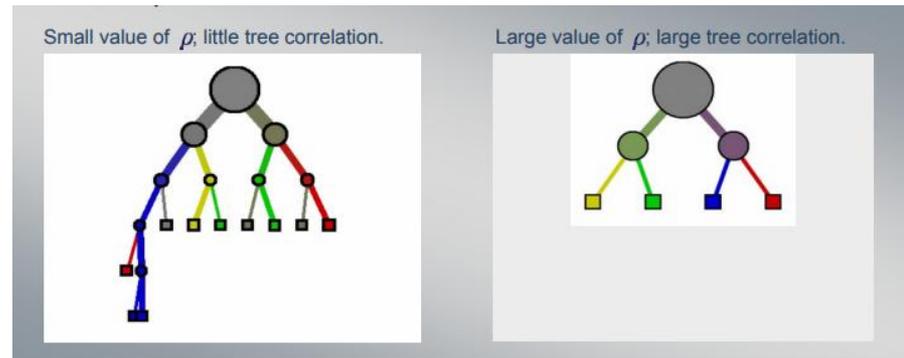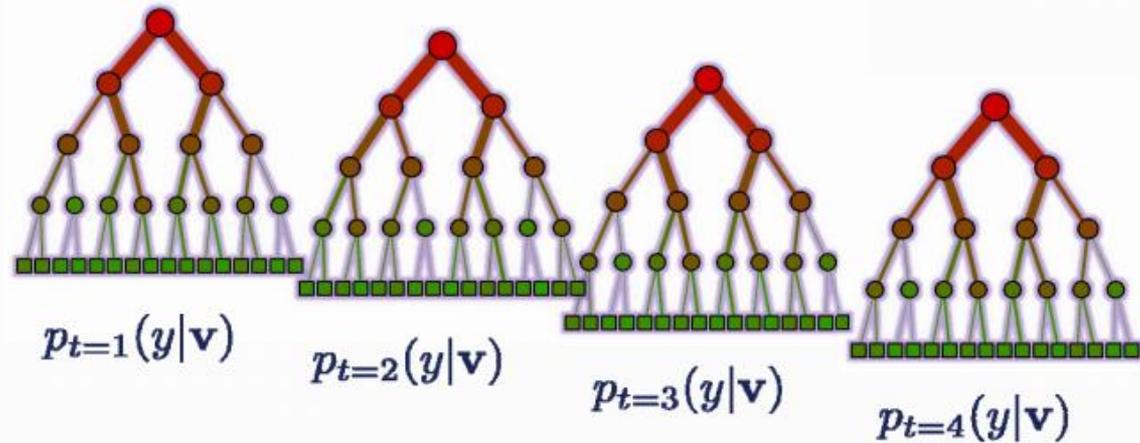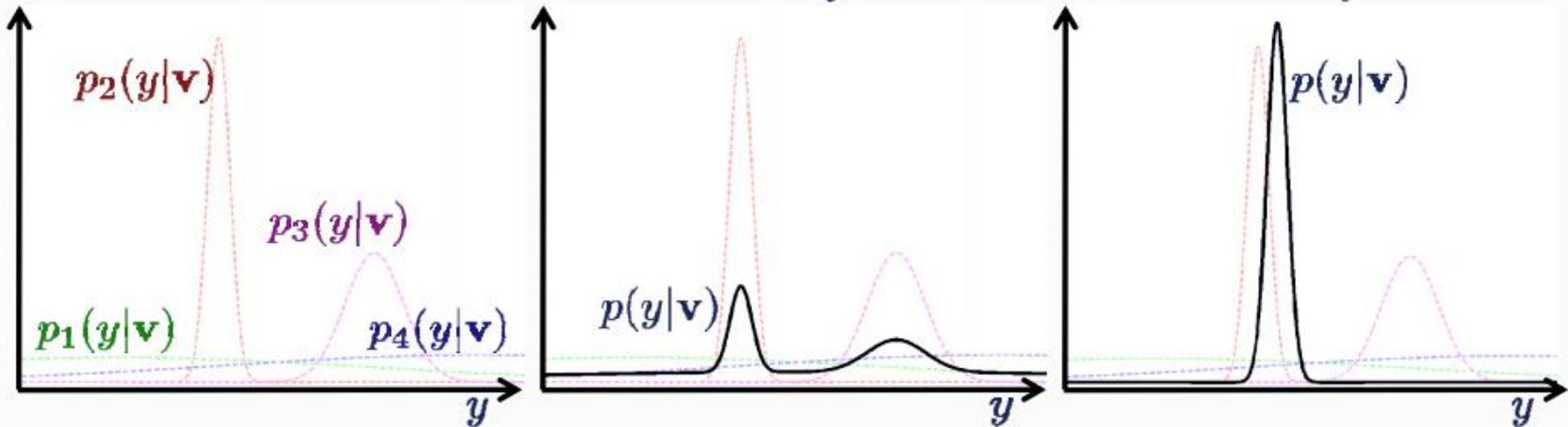# Randomness Model

- 2) Randomized node optimization (RNO)

| | | |
|---|---|---|
| $\mathcal{T}$ | The full set of all possible nodes | **Node training** |
| $\mathcal{T}_j \subset \mathcal{T}$ | For each node the set of randomly sampled features | Node weak learner $h(\mathbf{v}, \boldsymbol{\theta}_j)$ |
| $\rho = |\mathcal{T}_j|$ | Randomness control parameter. For $\rho = |\mathcal{T}|$ no randomness and maximum tree correlation. For $\rho = 1$ minimum tree correlation. | Node test params $\boldsymbol{\theta} \in \mathcal{T}_j$ |

**Precondition:** A training set $S := (x_1, y_1), \ldots, (x_n, y_n)$, features $F$, and number of trees in forest $B$.

```
1  function RANDOMFOREST(S, F)
2      H ← ∅
3      for i ∈ 1, ..., B do
4          S^(i) ← A bootstrap sample from S
5          h_i ← RANDOMIZEDTREELEARN(S^(i), F)
6          H ← H ∪ {h_i}
7      end for
8      return H
9  end function
```

Here B = ρ.

Small value of $\rho$; little tree correlation.

Large value of $\rho$; large tree correlation.

# The Ensemble Model



An example forest to predict continuous variables

$$p_{t=1}(y|\mathbf{v}) \qquad p_{t=2}(y|\mathbf{v}) \qquad p_{t=3}(y|\mathbf{v}) \qquad p_{t=4}(y|\mathbf{v})$$

$$p(y|\mathbf{v}) = \frac{1}{T}\sum_{t}^{T} p_t(y|\mathbf{v}) \qquad p(y|\mathbf{v}) = \frac{1}{Z}\prod_{t}^{T} p_t(y|\mathbf{v})$$

$p_2(y|\mathbf{v})$

$p_3(y|\mathbf{v})$

$p_1(y|\mathbf{v})$ $\qquad p_4(y|\mathbf{v})$

$p(y|\mathbf{v})$

$p(y|\mathbf{v})$

# Training and Information Gain



Information gain

$$I(\mathcal{S}, \boldsymbol{\theta}) = H(\mathcal{S}) - \sum_{i \in \{L,R\}} \frac{|\mathcal{S}^i|}{|\mathcal{S}|} H(\mathcal{S}^i)$$

Shannon's entropy

$$H(\mathcal{S}) = - \sum_{c \in \mathcal{C}} p(c) \log(p(c))$$

Node training

$$\boldsymbol{\theta} = \arg\max_{\boldsymbol{\theta} \in \mathcal{T}_j} I(\mathcal{S}_j, \boldsymbol{\theta})$$

# Calculating Information Gain

- Information Gain = entropy(parent) – [average entropy(children)]

$$\text{child entropy} \quad -\left(\frac{13}{17}\cdot\log_2\frac{13}{17}\right)-\left(\frac{4}{17}\cdot\log_2\frac{4}{17}\right)=0.787$$

Entire population (30 instances)

17 instances

$$\text{child entropy} \quad -\left(\frac{1}{13}\cdot\log_2\frac{1}{13}\right)-\left(\frac{12}{13}\cdot\log_2\frac{12}{13}\right)=0.391$$

$$\text{parent entropy} \quad -\left(\frac{14}{30}\cdot\log_2\frac{14}{30}\right)-\left(\frac{16}{30}\cdot\log_2\frac{16}{30}\right)=0.996$$

13 instances

$$(\text{Weighted}) \text{ Average Entropy of Children} = \left(\frac{17}{30}\cdot0.787\right)+\left(\frac{13}{30}\cdot0.391\right)=0.615$$

**Information Gain= 0.996 - 0.615 = 0.38  for this split.**

# Classification Forest



Training data in feature space

$x_2$

?

**V**

?

?

$x_1$

Classification tree training

$\mathcal{S}_0$

$p_0(c|\mathbf{v})$

$c$

$\mathcal{S}_1$    $\mathcal{S}_2$

$\mathcal{S}_1^{\mathbf{L}}$    $\mathcal{S}_1^{\mathbf{R}}$

$p_{14}(c|\mathbf{v})$

$c$

$p_{126}(c|\mathbf{v})$

$c$

Model specialization for classification

| | | |
|---|---|---|
| Input data point | $\mathbf{v} = (x_1, \cdots, x_d) \in \mathbb{R}^d$ | ( $x_i$ is feature response) |
| Output is categorical | $c \in \mathcal{C}$   with   $\mathcal{C} = \{c_k\}$ | (discrete set) |
| Node weak learner | $h(\mathbf{v}, \boldsymbol{\theta}) \in \{\mathbf{true}, \mathbf{false}\}$ | |
| Obj. funct. for node j | $I = H(\mathcal{S}_j) - \sum_{i=L,R} \frac{|\mathcal{S}_j^i|}{|\mathcal{S}_j|} H(\mathcal{S}_j^i)$ | (information gain) |
| Training node j | $\boldsymbol{\theta}_j = \arg\max_{\boldsymbol{\theta} \in \mathcal{T}_j} I(\mathcal{S}_j, \boldsymbol{\theta})$ | |
| Predictor model | $p(c|\mathbf{v}) = \sum_j p(c|j)p(j|\mathbf{v})$ | (class posterior) |

Entropy of a discrete  distribution

$$H(\mathcal{S}) = - \sum_{c \in \mathcal{C}} p(c) \log(p(c))$$

with   $c(\mathbf{v}) : \mathbb{R}^d \to \mathcal{C}$

# Weak Learners

Splitting data at node j

Node weak learner $h(\mathbf{v}, \boldsymbol{\theta}_j)$

Node test params $\boldsymbol{\theta} \in \mathcal{T}_j$

$\mathcal{S}$
$\mathcal{S}^L$ $\mathcal{S}^R$

Examples of weak learners



Weak learner: axis aligned

$$h(\mathbf{v}, \boldsymbol{\theta}) = [\tau_1 > \phi(\mathbf{v}) \cdot \boldsymbol{\psi} > \tau_2]$$

Feature response for 2D example. $\phi(\mathbf{v}) = (x_1 \; x_2 \; 1)^\top$

With $\boldsymbol{\psi} = (1 \; 0 \; \psi_3)$ or $\boldsymbol{\psi} = (0 \; 1 \; \psi_3)$

Weak learner: oriented line

$$h(\mathbf{v}, \boldsymbol{\theta}) = [\tau_1 > \phi(\mathbf{v}) \cdot \boldsymbol{\psi} > \tau_2]$$

Feature response for 2D example. $\phi(\mathbf{v}) = (x_1 \; x_2 \; 1)^\top$

With $\boldsymbol{\psi} \in \mathbb{R}^3$ a generic line in homog. coordinates.

Weak learner: conic section

$$h(\mathbf{v}, \boldsymbol{\theta}) = \left[\tau_1 > \phi^\top(\mathbf{v}) \, \boldsymbol{\psi} \, \phi(\mathbf{v}) > \tau_2\right]$$

Feature response for 2D example. $\phi(\mathbf{v}) = (x_1 \; x_2 \; 1)^\top$

With $\boldsymbol{\psi} \in \mathbb{R}^{3 \times 3}$ a matrix representing a conic.

In general $\boldsymbol{\phi}$ may select only a very small subset of features $\phi(\mathbf{v}) : \mathbb{R}^d \to \mathbb{R}^{d'+1}, \; d' << d$

# Prediction Model



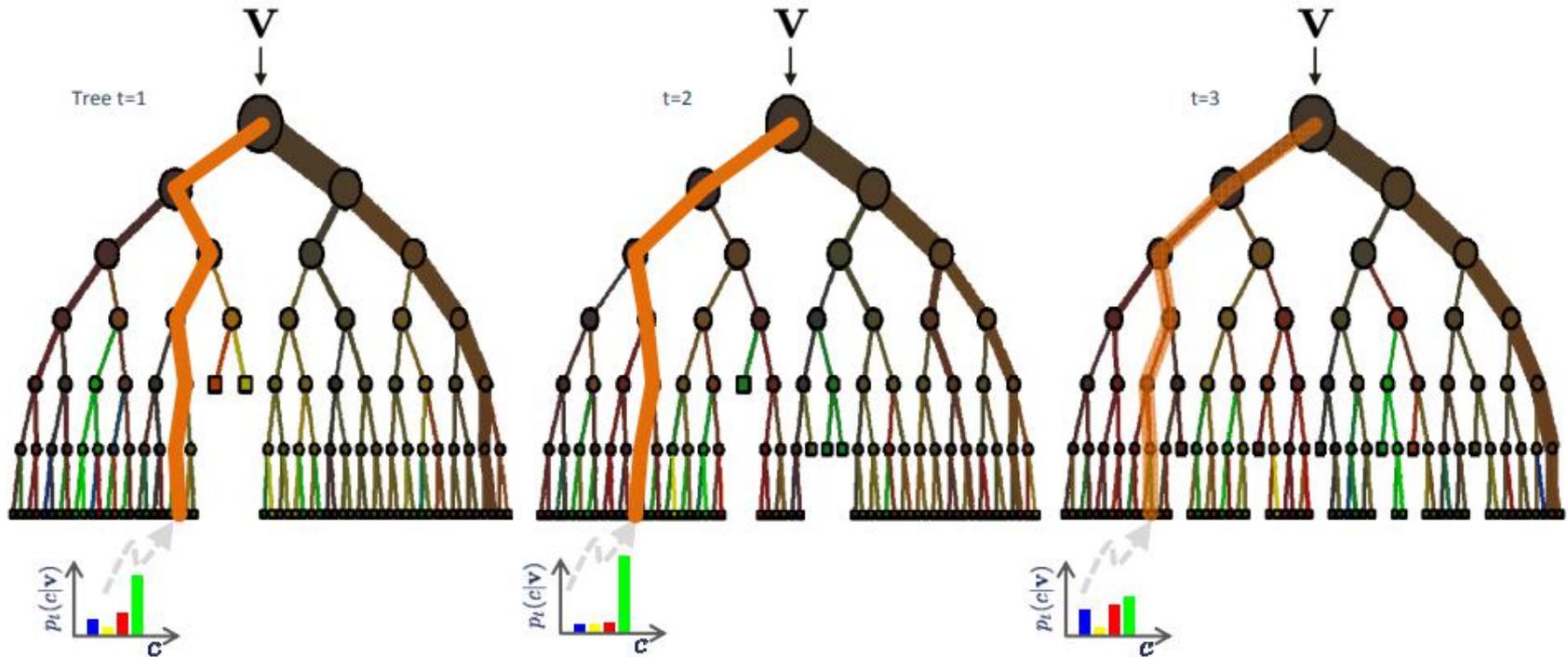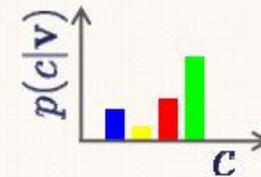What do we do at the leaf?

Prediction model: probabilistic

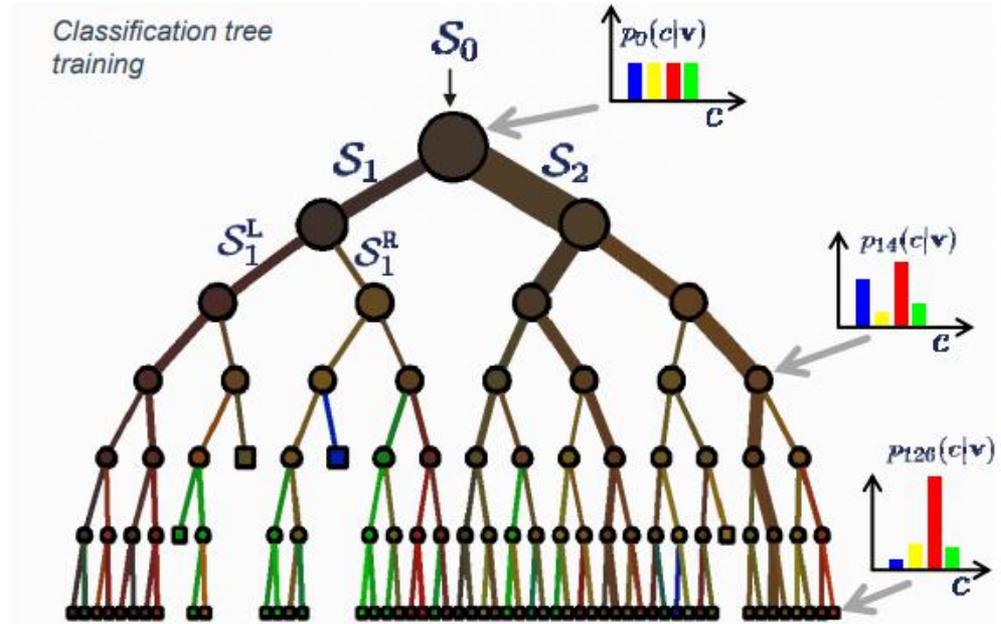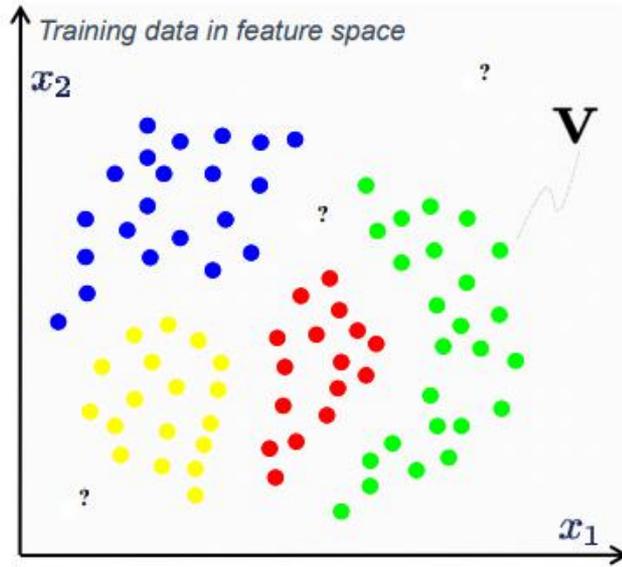$p(c|\mathbf{v})$

$c$

# Classification Forest: Ensemble Model



The ensemble model

Forest output probability $\quad p(c|\mathbf{v}) = \dfrac{1}{T}\sum_{t}^{T} p_t(c|\mathbf{v})$

# Classification Forest



Training data in feature space

$x_2$

$x_1$

V

?

?

?

Classification tree training

$S_0$

$S_1$    $S_2$

$S_1^L$    $S_1^R$

$p_0(c|\mathbf{v})$

$p_{14}(c|\mathbf{v})$

$p_{126}(c|\mathbf{v})$

Model specialization for classification

| | | |
|---|---|---|
| Input data point | $\mathbf{v} = (x_1, \cdots, x_d) \in \mathbb{R}^d$ | ( $x_i$ is feature response) |
| Output is categorical | $c \in \mathcal{C}$ with $\mathcal{C} = \{c_k\}$ | (discrete set) |
| Node weak learner | $h(\mathbf{v}, \boldsymbol{\theta}) \in \{\mathbf{true}, \mathbf{false}\}$ | |
| Obj. funct. for node j | $I = H(\mathcal{S}_j) - \sum_{i=\mathrm{L,R}} \dfrac{|\mathcal{S}_j^i|}{|\mathcal{S}_j|} H(\mathcal{S}_j^i)$ | (information gain) |
| Training node j | $\boldsymbol{\theta}_j = \underset{\boldsymbol{\theta} \in \mathcal{T}_j}{\arg\max}\, I(\mathcal{S}_j, \boldsymbol{\theta})$ | |
| Predictor model | $p(c|\mathbf{v}) = \sum_j p(c|j) p(j|\mathbf{v})$ | (class posterior) |

Entropy of a discrete distribution

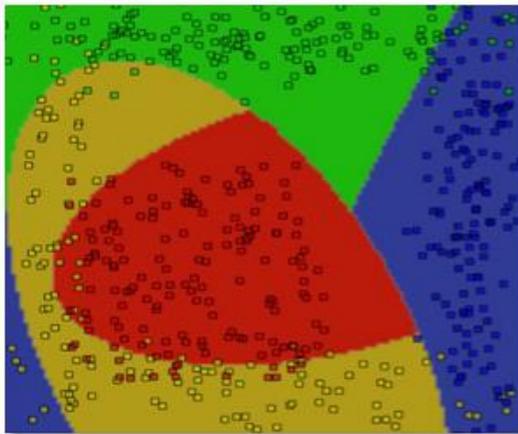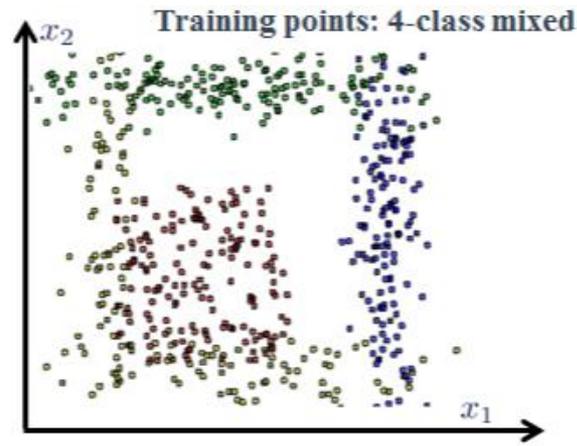$$H(\mathcal{S}) = - \sum_{c \in \mathcal{C}} p(c) \log(p(c))$$

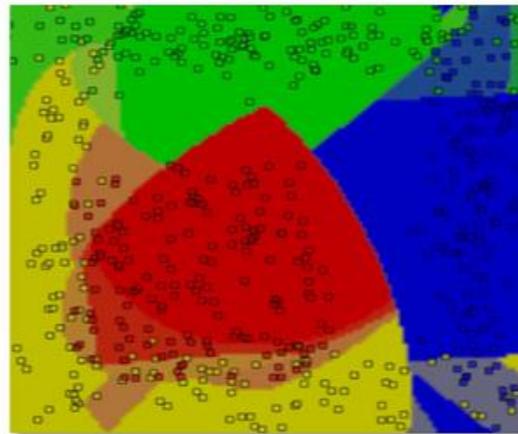with $c(\mathbf{v}) : \mathbb{R}^d \to \mathcal{C}$

# Overfitting and Underfitting

# Effect of Tree Depth
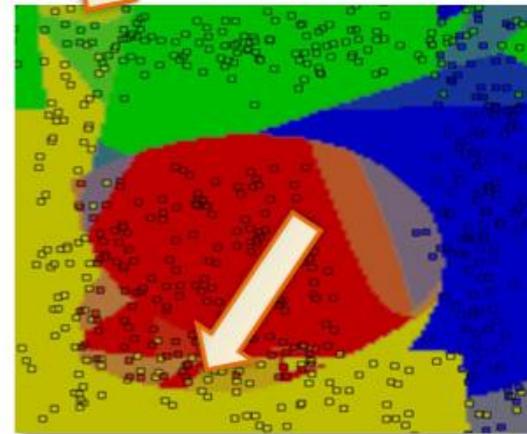


Training points: 4-class mixed

T=200, D=3, w. l. = conic

T=200, D=6, w. l. = conic

T=200, D=15, w. l. = conic

max tree depth, D

underfitting

overfitting

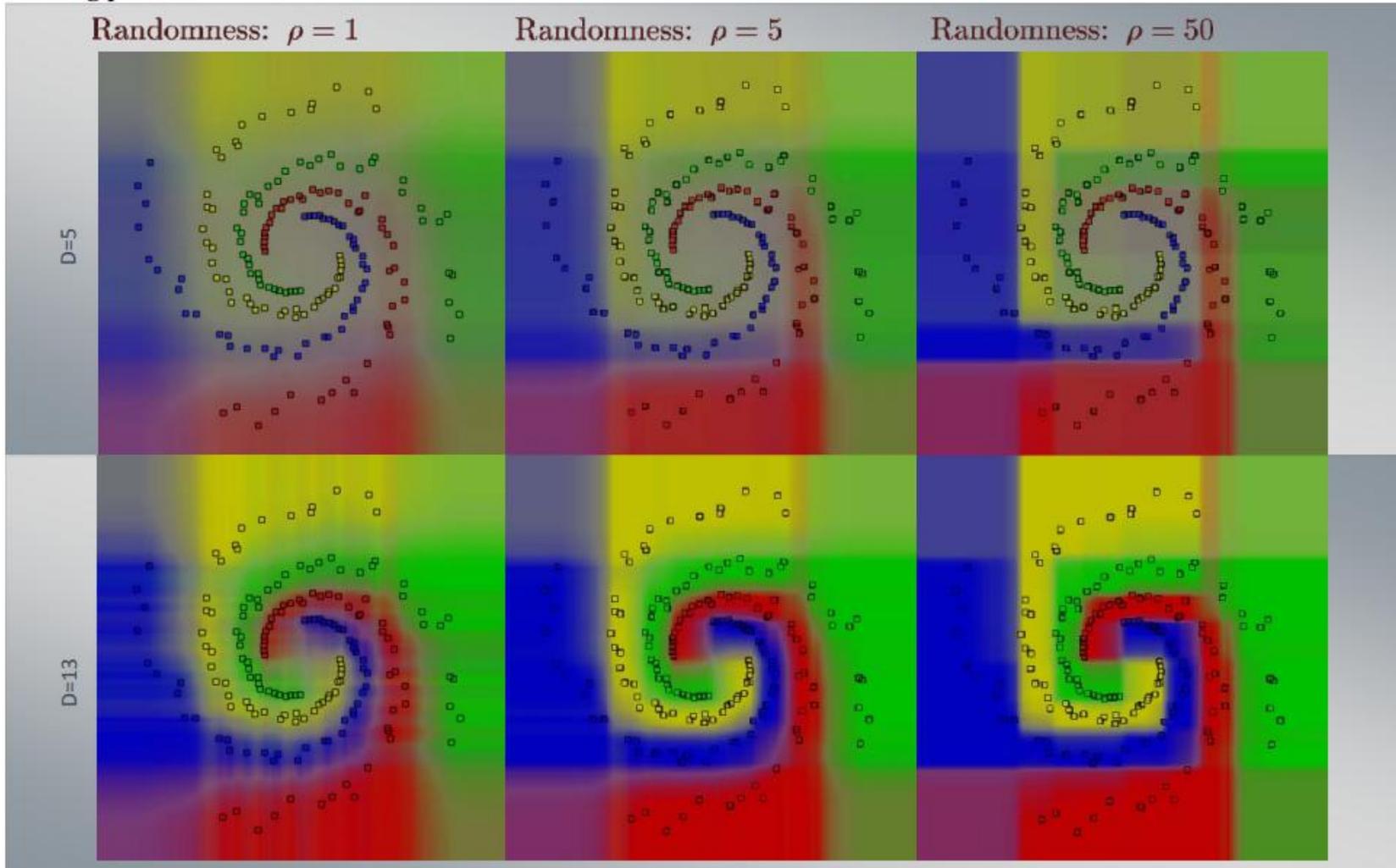# Effect of Weak Learner Model and Randomness

Testing posteriors



Weak learner: axis aligned     Weak learner: oriented line     Weak learner: conic section

D=5

D=13

Randomness: $\rho = 500$

Parameters: T=400 predictor model = prob.
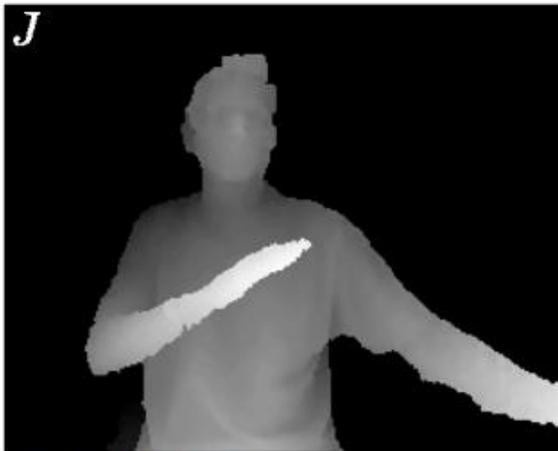
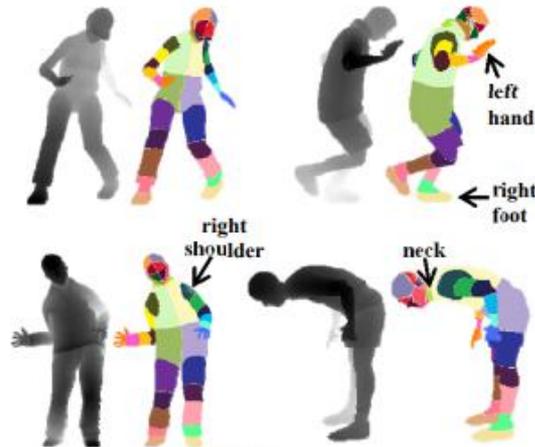# Effect of Weak Learner Model and Randomness

Testing posteriors



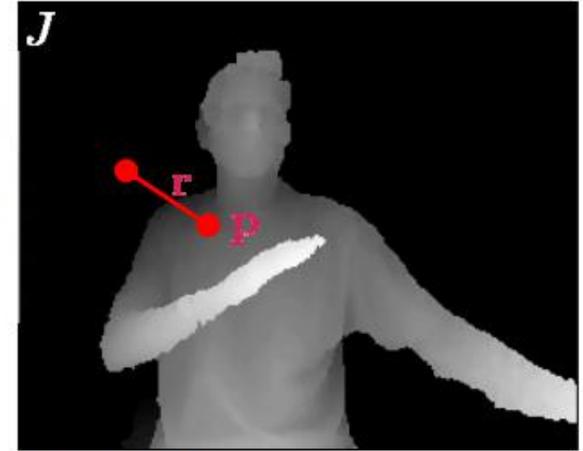Weak learner: axis aligned

Parameters: T=400 predictor model = prob.

# Body tracking in Microsoft Kinect for XBox 360



Input depth image

Training labelled data

Visual features

**Classification forest**

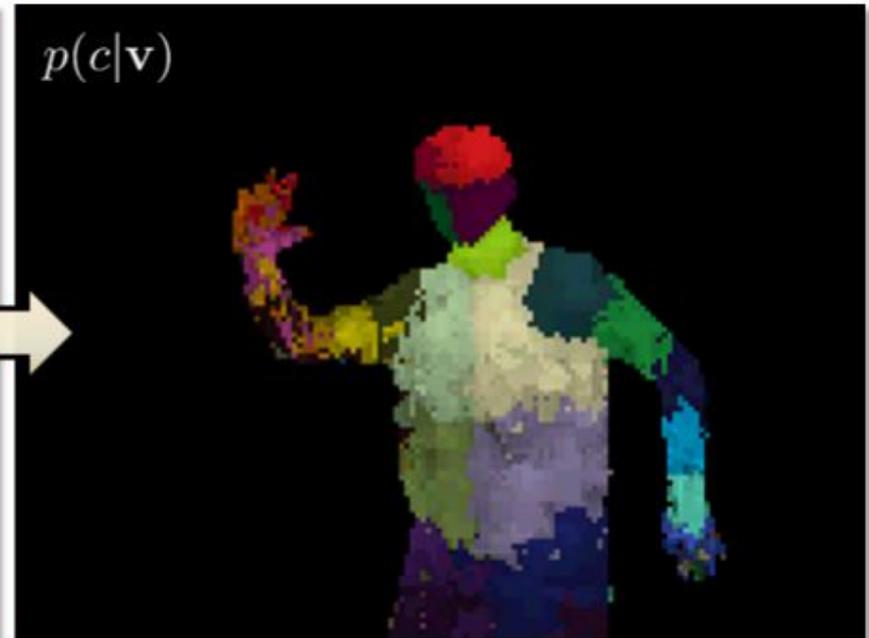| | | | |
|---|---|---|---|
| Labels are categorical | $c \in \{\mathbf{l.hand}, \mathbf{r.hand}, \mathbf{head}, \ldots\}$ | Objective function | $I = H(\mathcal{S}_j) - \sum_{i=L,R} \frac{|\mathcal{S}_j^i|}{|\mathcal{S}_j|} H(\mathcal{S}_j^i)$ |
| Input data point | $\mathbf{p} \in \mathbb{R}^2$ | Node parameters | $\boldsymbol{\theta} = (\mathbf{r}, \tau)$ |
| Visual features | $\mathbf{v}(\mathbf{p}) = (x_1, \ldots, x_i, \ldots, x_d) \in \mathbb{R}^d$ | Node training | $\boldsymbol{\theta}_j = \arg \max_{\boldsymbol{\theta} \in \mathcal{T}_j} I(\mathcal{S}_j, \boldsymbol{\theta})$ |
| Feature response | $x_i = J(\mathbf{p}) - J\left(\mathbf{p} + \frac{\mathbf{r}_i}{J(\mathbf{p})}\right)$ | Weak learner | $h(\mathbf{v}, \boldsymbol{\theta}) = [\phi(\mathbf{v}, \mathbf{r}) > \tau]$ |
| Predictor model | $p(c|\mathbf{v})$ | | |

# Body tracking in Microsoft Kinect for XBox 360



$J$

Input depth image (bg removed)

$p(c|\mathbf{v})$

Inferred body parts posterior

# Advantages of Random Forests

- Very high accuracy – not easily surpassed by other algorithms
- Efficient on large datasets
- Can handle thousands of input variables without variable deletion
- Effective method for estimating missing data, also maintains accuracy when a large proportion of the data are missing
- Can handle categorical variables
- Robust to label noise
- Can be used in clustering, locating outliers and semisupervised learning

# Outline

- Bagging

- Random Forest

- **Boosting**

# Boosting

- Given a set of weak learners, run them multiple times on (reweighted) training data, then let learned classifiers vote

- At each iteration t:
  - Weight each training example by how incorrectly it was classified
  - Learn a hypothesis – $h_t$

The one with the smallest error

  - Choose a strength for this hypothesis – $\alpha_t$

- Final classifier: weighted combination of weak learners

# Learning from Weighted Data

- Sometimes not all data points are equal

  ❑ Some data points are more equal than others

- Consider a weighted dataset

  ❑ $D(i)$ – weight of i−th training example $(x_i, y_i)$

  ❑ Interpretations:

    - i−th training example counts as $D(i)$ examples
    - If I were to "resample" data, I would get more samples of "heavier" data points

- Now, in all calculations the i−th training example counts as $D(i)$ "examples"

# Definition of Boosting

- Given training set $(x_1, y_1), \ldots, (x_m, y_m)$
- $y_i \in \{-1, +1\}$ correct label of instance $x_i \in X$
- For $t = 1, \ldots, T$
- – construct distribution $D_t$ on $\{1, \ldots, m\}$
- – find weak hypothesis
- – $h_t: X \rightarrow \{-1, +1\}$ with smallest error $\varepsilon_t$ on $D_t$

$$h_t = \arg \min_{h_j \in \mathcal{H}} \epsilon_j = \sum_{i=1}^{m} D_t(i)[y_i \neq h_j(x_i)]$$

- Output final hypothesis $H_{final}$

# AdaBoost

- Constructing Dt

$$D_1(i) = \frac{1}{m}$$

$$\Longrightarrow$$

$$D_{t+1} = \frac{D_t(i)}{Z_t} c(x)$$

$$c(x) = \begin{cases} e^{-\alpha_t} & : y_i = h_t(x_i) \\ e^{\alpha_t} & : y_i \neq h_t(x_i) \end{cases}$$

$$D_{t+1} = \frac{D_t(i)}{Z_t} e^{-\alpha_t y_i h_t(x_i)}$$

where Zt is a normalization constant

- Final hypothesis:

$$H(x) = sign\left(\sum_{t=1}^{T} \alpha_t h_t(x)\right)$$

$$\alpha_t = \frac{1}{2} \ln \frac{1 - \epsilon_t}{\epsilon_t} > 0$$

# The AdaBoost Algorithm

Given: $(x_1, y_1), \ldots, (x_m, y_m)$ where $x_i \in X$, $y_i \in Y = \{-1, +1\}$
Initialize $D_1(i) = 1/m$.
For $t = 1, \ldots, T$:

- Train base learner using distribution $D_t$.
- Get base classifier $h_t : X \to \mathbb{R}$.
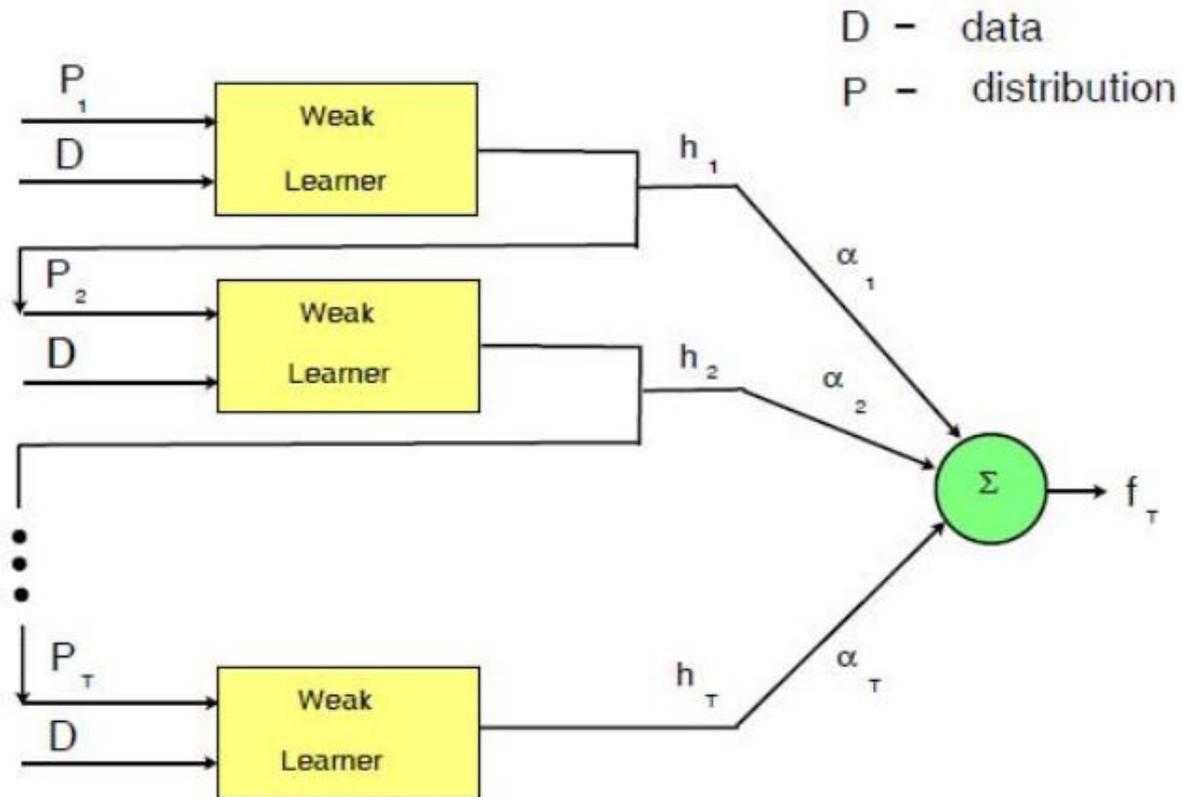- Choose $\alpha_t \in \mathbb{R}$.
- Update:

$$D_{t+1}(i) = \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

where $Z_t$ is a normalization factor (chosen so that $D_{t+1}$ will be a distribution).

Output the final classifier:

$$H(x) = \operatorname{sign}\left(\sum_{t=1}^{T} \alpha_t h_t(x)\right).$$

# The AdaBoost Algorithm

Given: $(x_1, y_1), \ldots, (x_m, y_m)$ where $x_i \in X$, $y_i \in Y = \{-1, +1\}$
Initialize $D_1(i) = 1/m$.
For $t = 1, \ldots, T$:

- Train base learner using distribution $D_t$.
- Get base classifier $h_t : X \to \mathbb{R}$.     with minimum $\epsilon_t$
- Choose $\alpha_t \in \mathbb{R}$.
- Update:

$$D_{t+1}(i) = \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

where $Z_t$ is a normalization factor (chosen so that $D_{t+1}$ will be a distribution).

Output the final classifier:

$$H(x) = \text{sign}\left(\sum_{t=1}^{T} \alpha_t h_t(x)\right).$$

# The AdaBoost Algorithm

Given: $(x_1, y_1), \ldots, (x_m, y_m)$ where $x_i \in X$, $y_i \in Y = \{-1, +1\}$
Initialize $D_1(i) = 1/m$.
For $t = 1, \ldots, T$:

- Train base learner using distribution $D_t$.
- Get base classifier $h_t : X \to \mathbb{R}$.
- Choose $\alpha_t \in \mathbb{R}$. $\longleftarrow$ $\boxed{\alpha_t = \frac{1}{2} \ln \left( \frac{1 - \epsilon_t}{\epsilon_t} \right)}$
- Update:

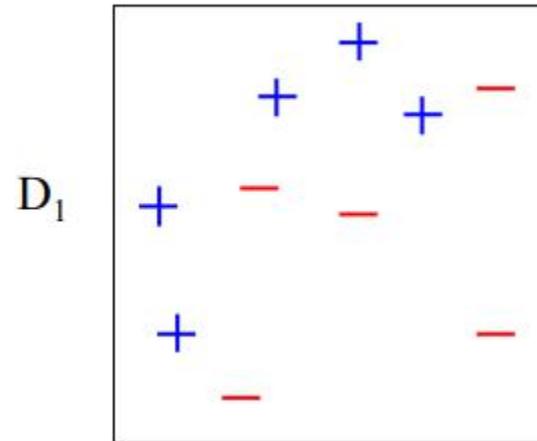$$D_{t+1}(i) = \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$$

$$\epsilon_t = \Pr_{i \sim D_t} \left[ h_t(x_i) \neq y_i \right]$$

$$\boxed{\epsilon_t = \frac{1}{\sum_{i=1}^n D_t(i)} \sum_{i=1}^m D_t(i) \delta(h_t(x_i) \neq y_i)}$$
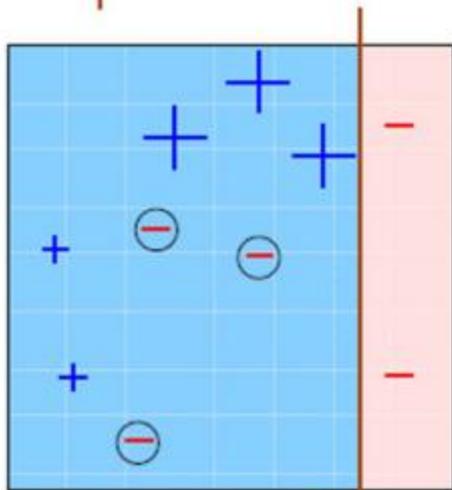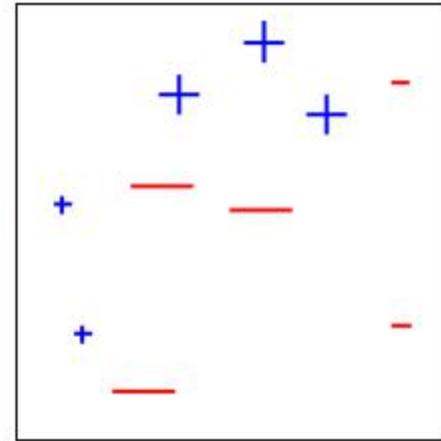
# The AdaBoost Algorithm

# Toy Example



$D_1$

# Toy Example: Round 1



$\varepsilon_1 = 0.3$
$\alpha_1 = 0.42$

$h_1$

$D_2$

# Toy Example: Round 2



$\varepsilon_2 = 0.21$
$\alpha_2 = 0.65$

$D_3$

$h_2$

# Toy Example: Round 3



$$\varepsilon_3 = 0.14$$
$$\alpha_3 = 0.92$$

$h_3$

# Toy Example: Final Hypothesis

$$H_{final} = \text{sign}\left( 0.42 \quad + 0.65 \quad + 0.92 \quad \right)$$

*Q & A*