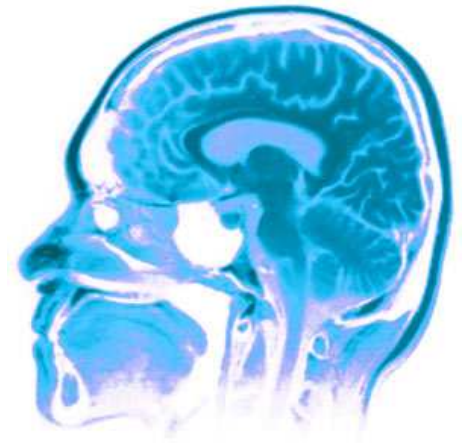




CPS C540



Random Forests



Nando de Freitas

February, 2013

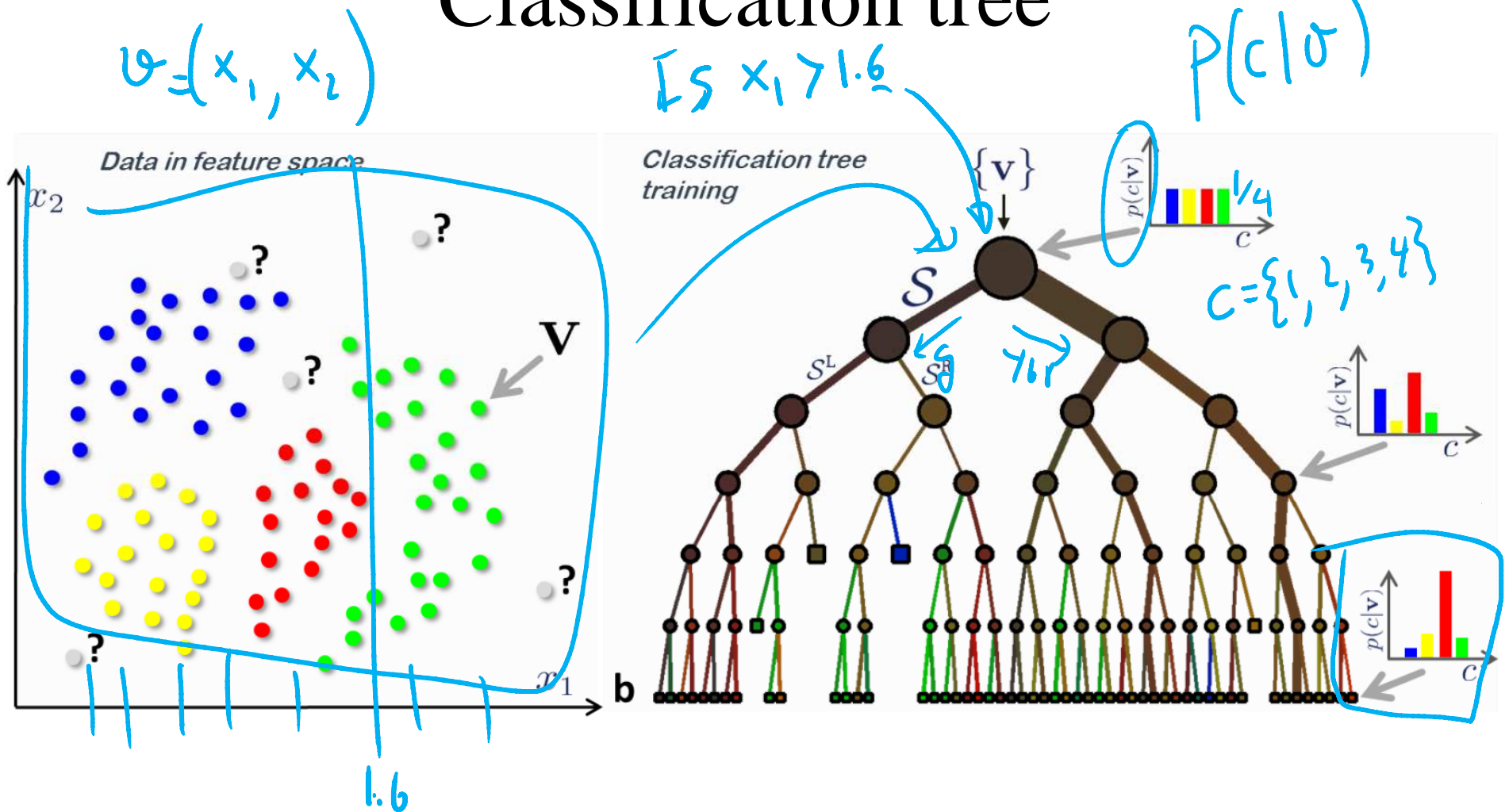
University of British Columbia

Outline of the lecture

This lecture discusses classification trees and how to incorporate them into an ensemble (random forest). It discusses:

- Random trees
- Random forests
- Object detection
- Kinect

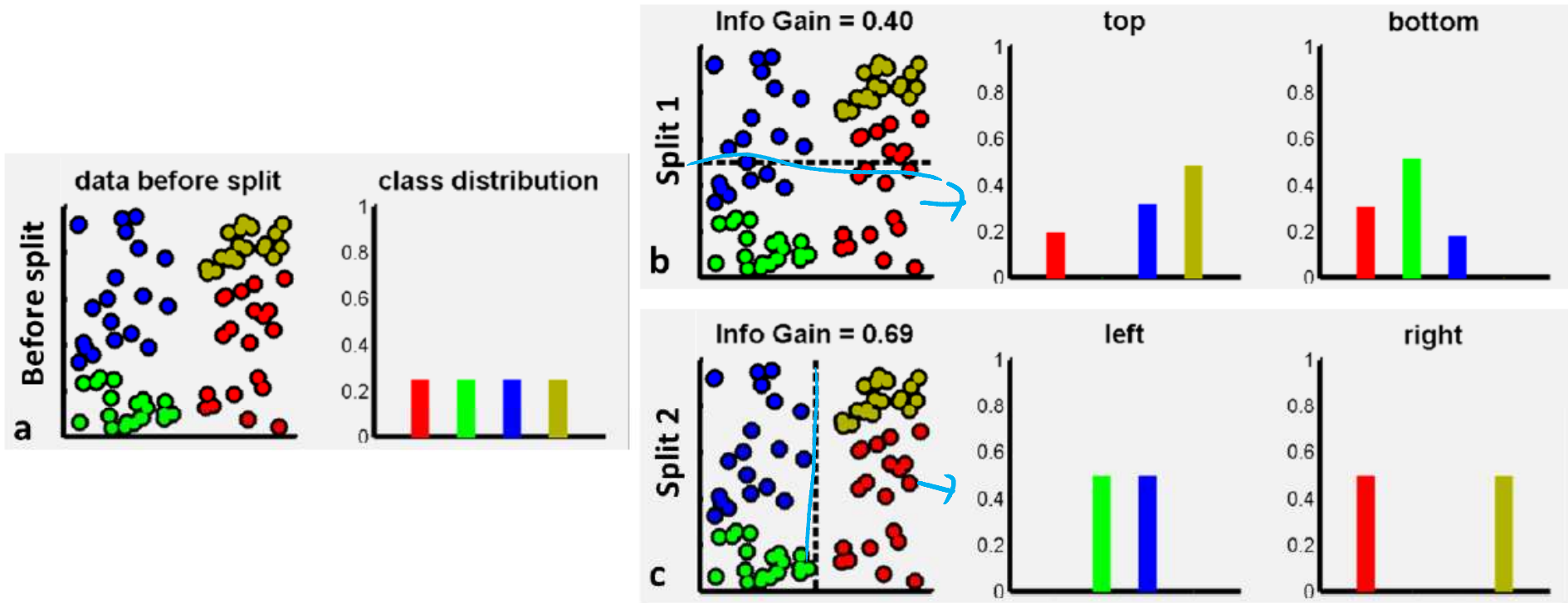
Classification tree



A generic data point is denoted by a vector $\mathbf{v} = (x_1, x_2, \dots, x_d)$

$$S_j = S_j^L \cup S_j^R$$

Use information gain to decide splits



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

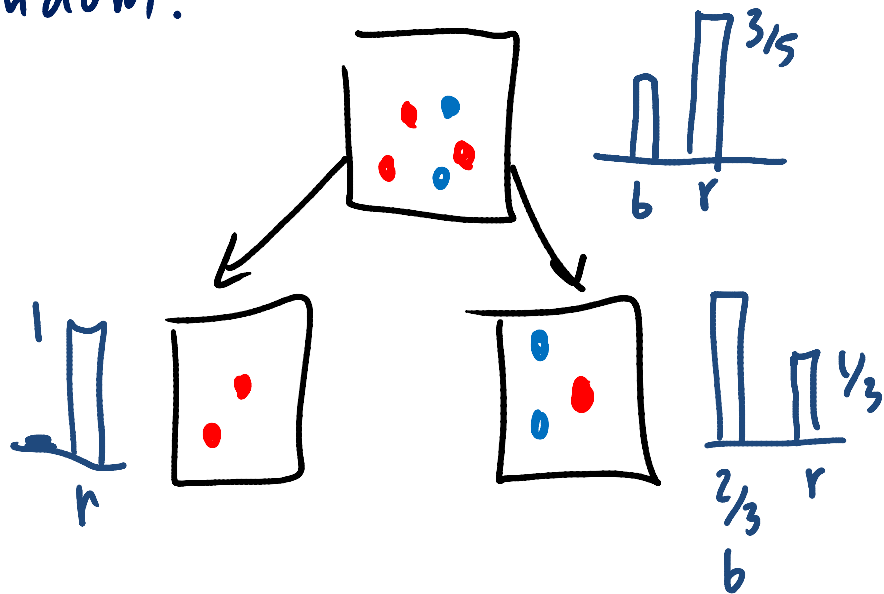
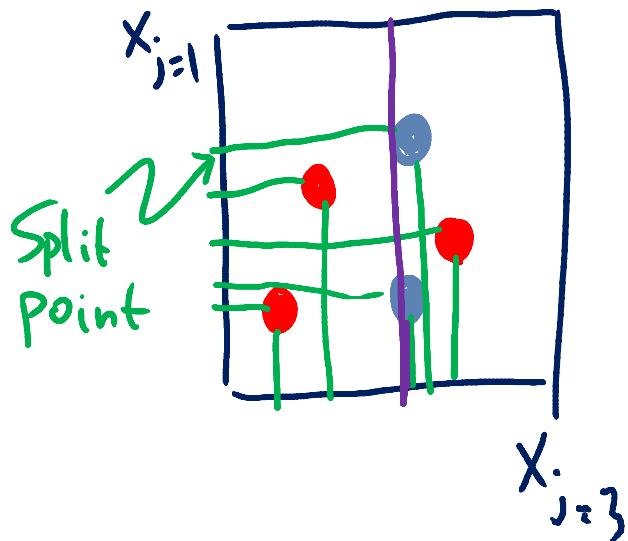
$d=3$ features
 $n=5$ data

Building a random tree

$$X = \begin{matrix} & i=1 & j=2 & & i=n=5 \\ \begin{matrix} j=1 \\ j=2 \\ j=3 \end{matrix} & \begin{pmatrix} 1 \\ 0 \\ 2 \end{pmatrix} & \begin{pmatrix} 3 \\ 6 \\ 1 \end{pmatrix} & \begin{pmatrix} 0 \\ 2 \\ 4 \end{pmatrix} & \begin{pmatrix} 8 \\ 9 \\ 0 \end{pmatrix} & \begin{pmatrix} 5 \\ 5 \\ 1 \end{pmatrix} \end{matrix}$$

$$Y = \begin{matrix} & i=1 & i=2 & & i=5 \\ & \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} & \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix} & & \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix} \end{matrix}$$

Pick 2 features at random.



Random Forests algorithm

1. For $b = 1$ to B : ^{trees}
- $data = \{ (x_1, y_1), \dots, (x_n, y_n) \}$
- (a) Draw a bootstrap sample Z^* of size N from the training data. ^{↪ bagging}
- (b) Grow a random-forest tree T_b to the bootstrapped data, by recursively repeating the following steps for each terminal node of the tree, until the minimum node size n_{min} is reached.
- i. Select m variables at random from the p variables. ←
 - ii. Pick the best variable/split-point among the m . ←
 - iii. Split the node into two daughter nodes. ←
2. Output the ensemble of trees $\{T_b\}_1^B$.

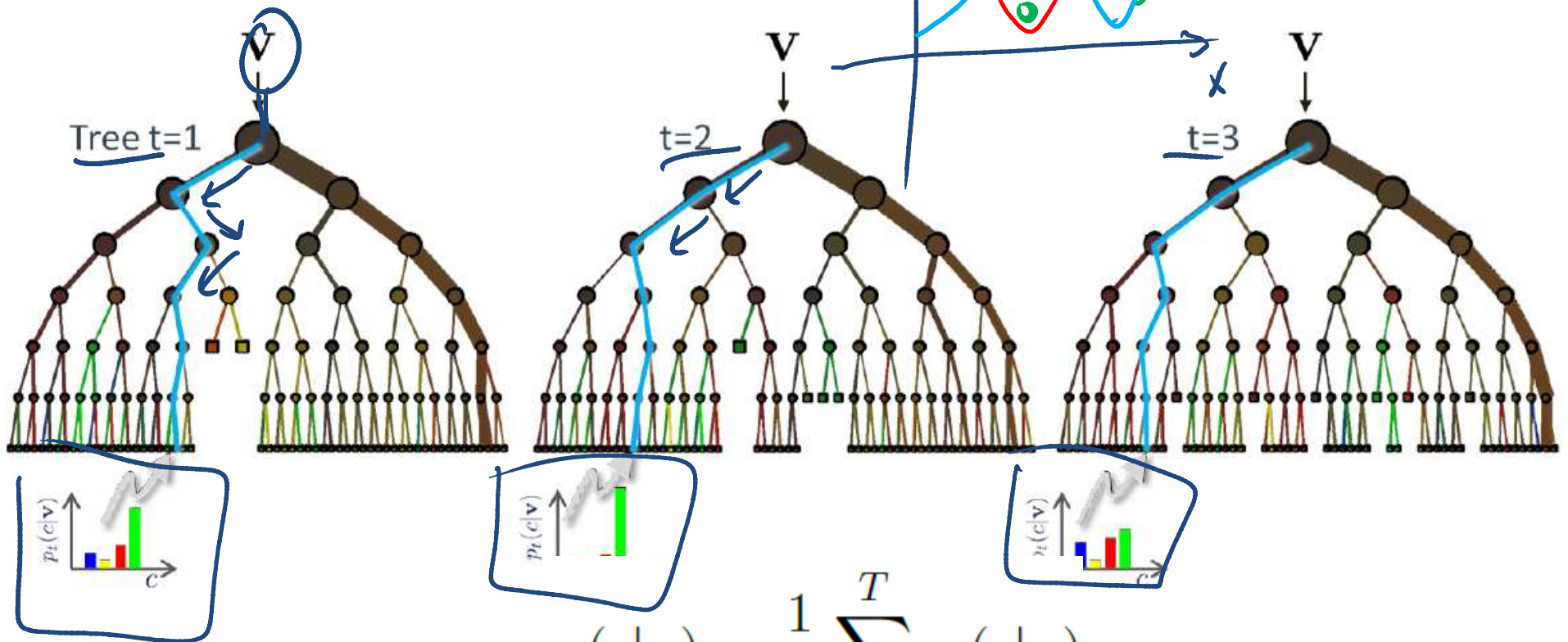
Randomization

Randomized node optimization. If \mathcal{T} is the entire set of all possible parameters θ then when training the j^{th} node we only make available a small subset $\mathcal{T}_j \subset \mathcal{T}$ of such values.

$$\theta_j^* = \arg \max_{\theta_j \in \mathcal{T}_j} I_j.$$

Building a forest (ensemble)

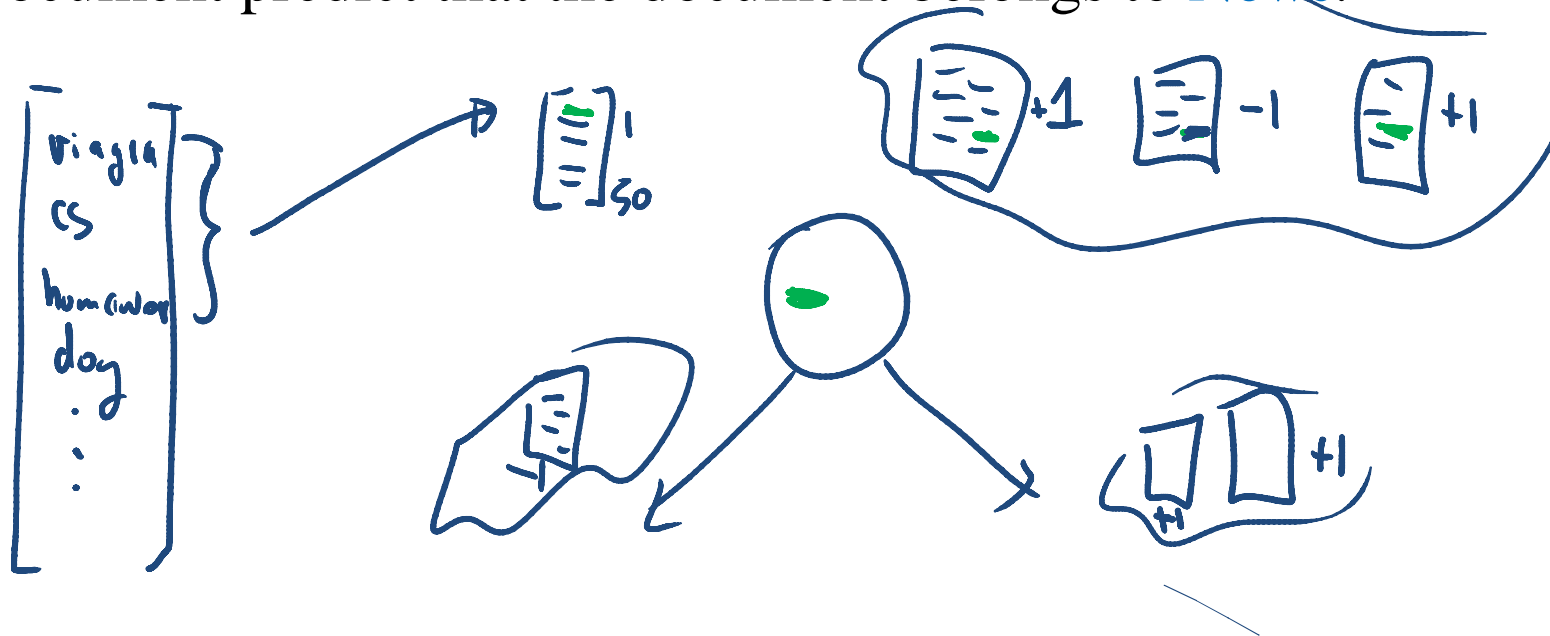
In a forest with T trees we have $t \in \{1, \dots, T\}$. All trees are trained independently (and possibly in parallel). During testing, each test point \mathbf{v} is simultaneously pushed through all trees (starting at the root) until it reaches the corresponding leaves.



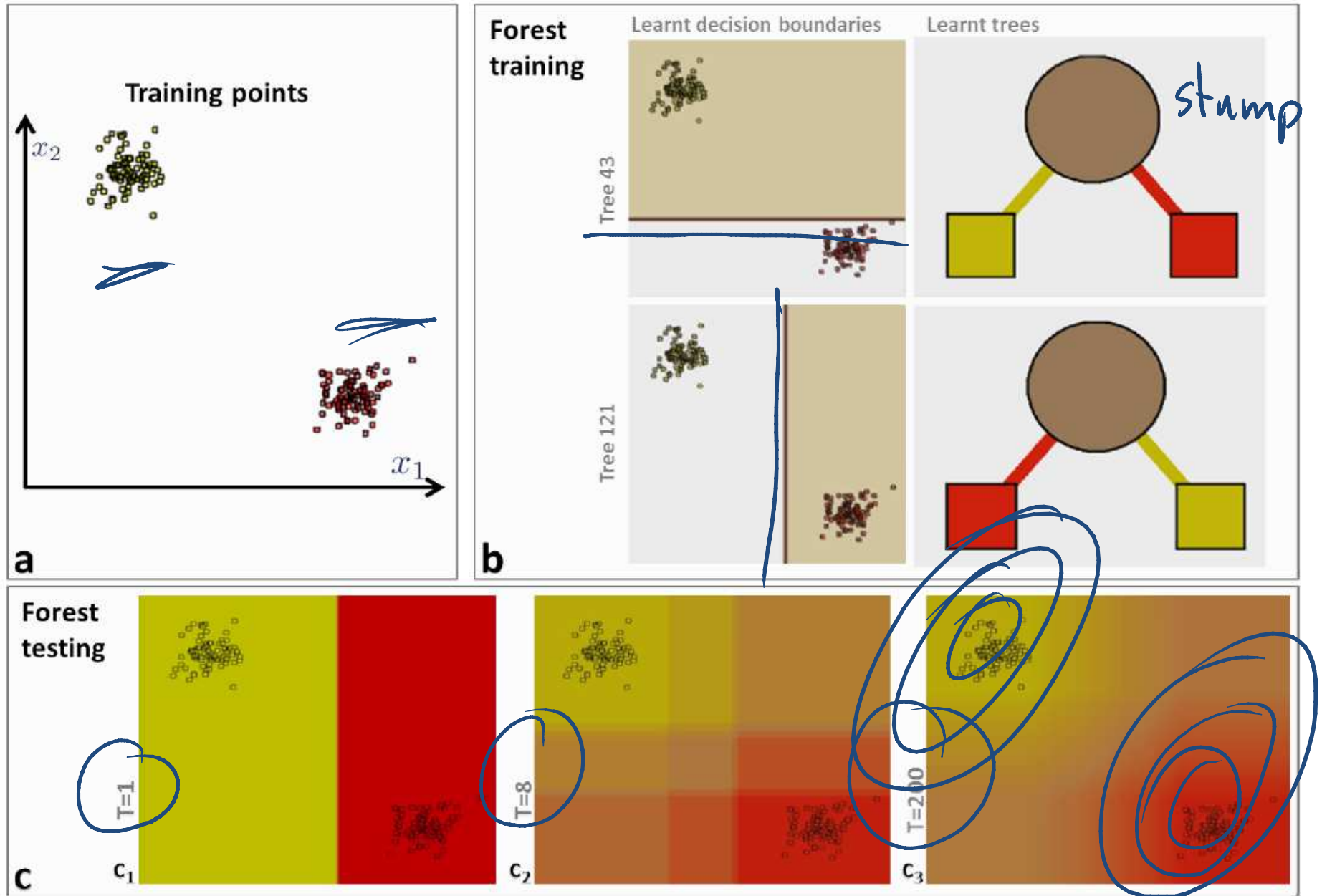
$$p(c|\mathbf{v}) = \frac{1}{T} \sum_{t=1}^T p_t(c|\mathbf{v})$$

Text classification example

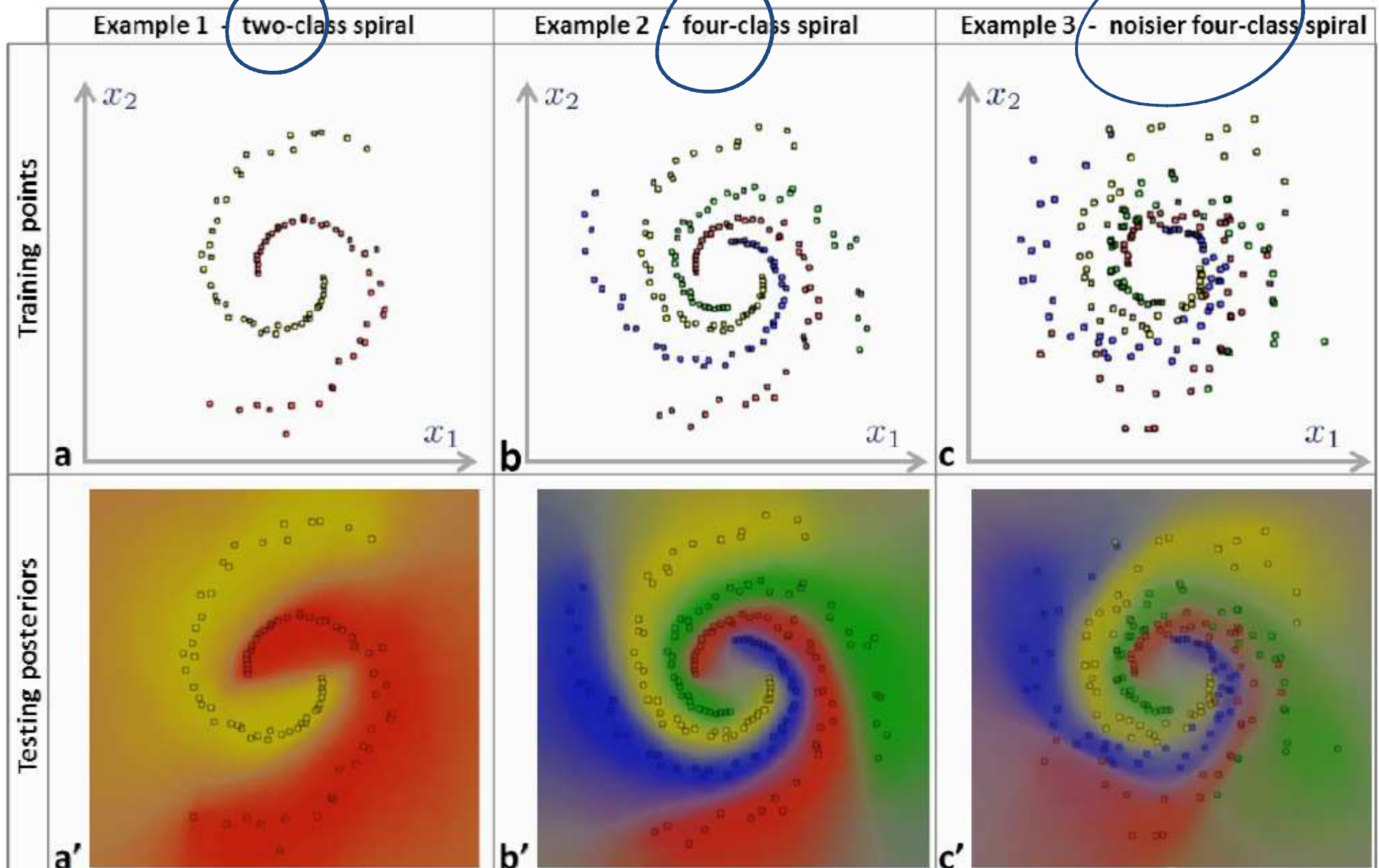
In news categorization, a possible term is *Bill Clinton*. A corresponding **weak learner (node)** is: If the term *Bill Clinton* appears in the document predict that the document belongs to *News*.



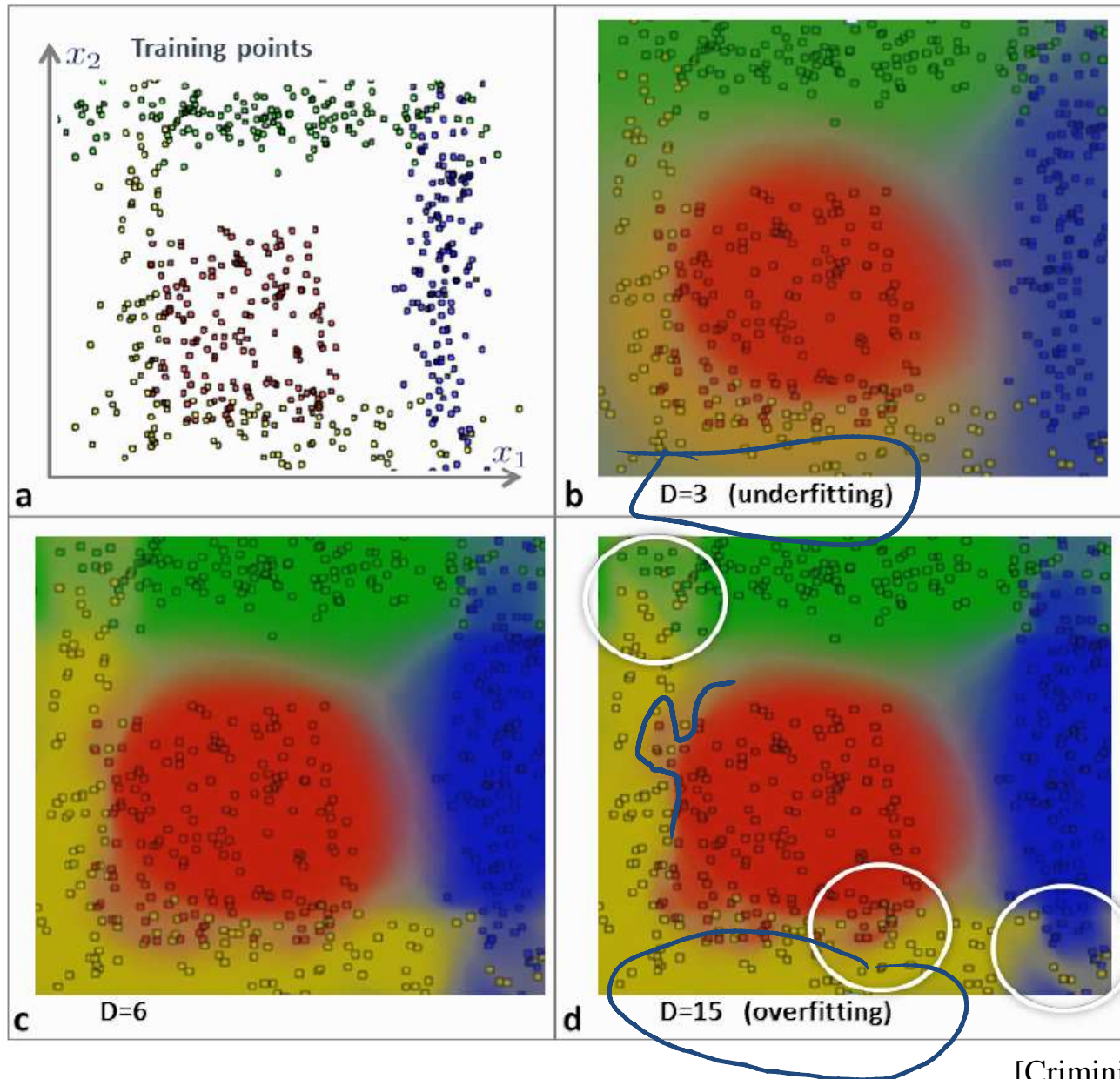
Effect of forest size



Effect of more classes and noise

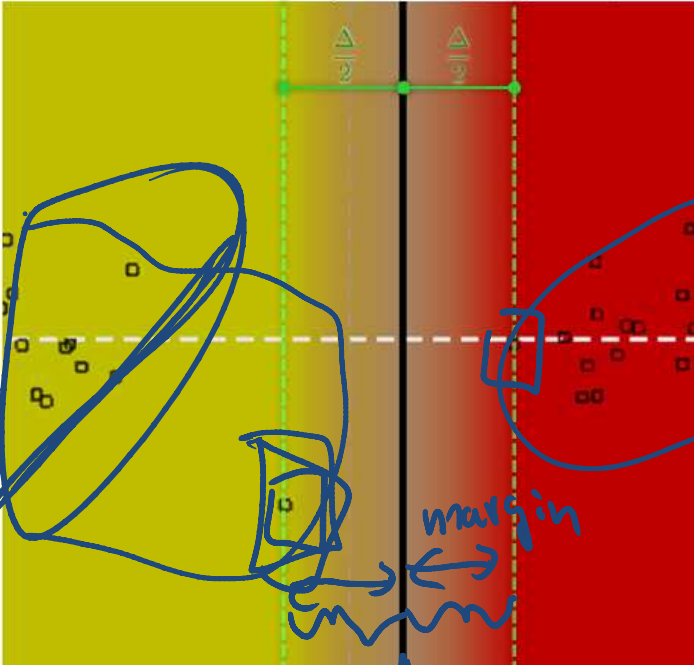


Effect of tree depth (D)



Effect of bagging

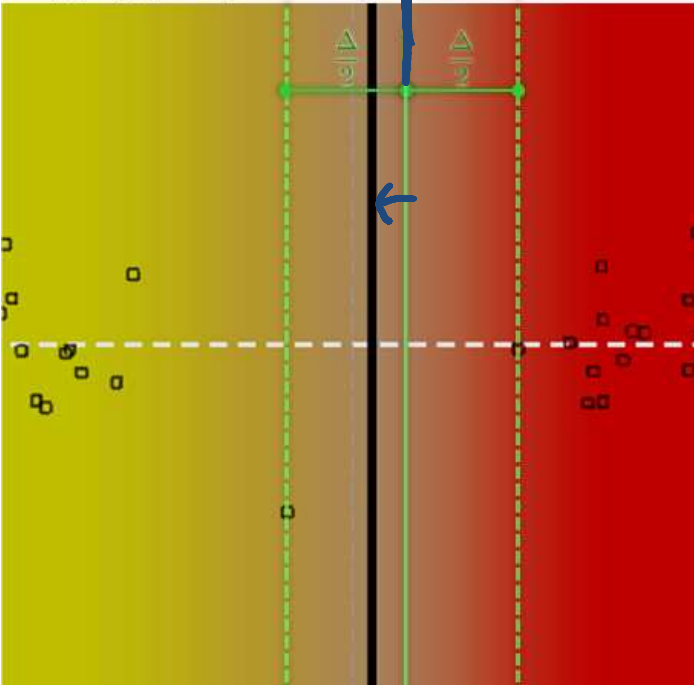
Randomized node optimization (RNO)



No bagging

max margin

Bagging (50%) and RNO



Application to face detection

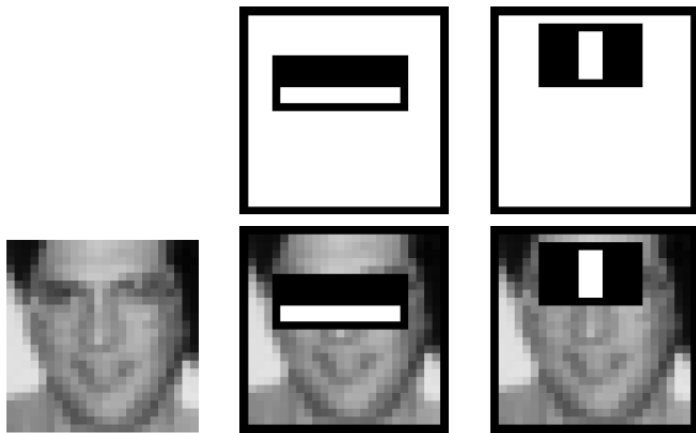
Training Data

- 5000 faces
 - All frontal
- 300 million non faces
 - 9400 non-face images



Object detection

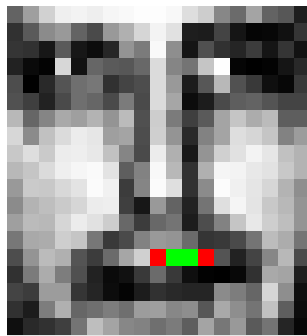
Idea: Extract simple features from all 24 by 24 pixel patches x_i . E.g., the value of a *two-rectangle feature* is the difference between the sum of the pixels within two rectangular regions. Then compare the level of activation (value of the feature f) with respect to a threshold (θ).



$$h_t(x_i) = \begin{cases} 1 & \text{if } f_t(x_i) > \theta_t \\ 0 & \text{otherwise} \end{cases}$$

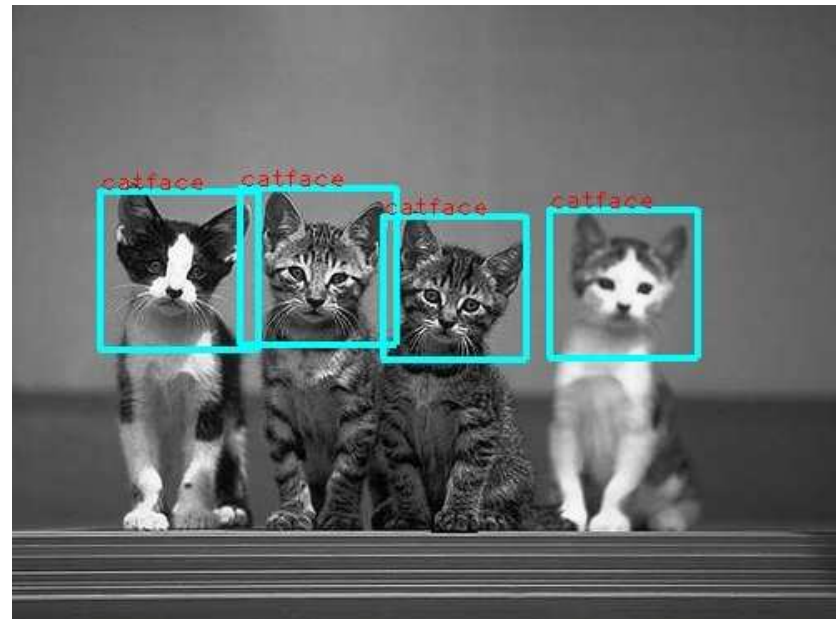
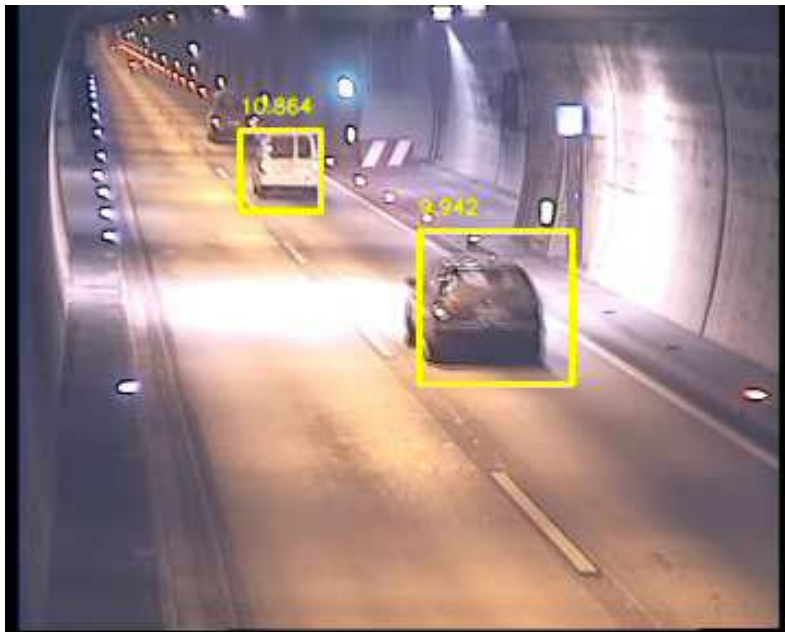
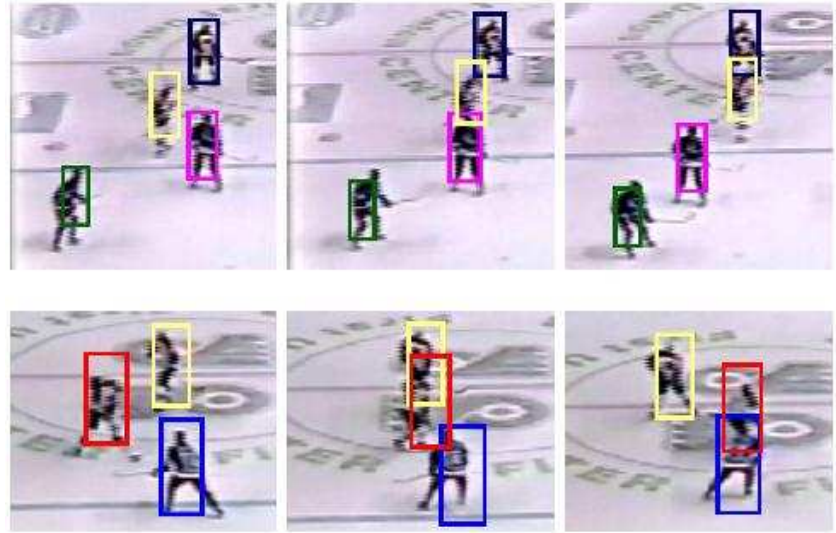
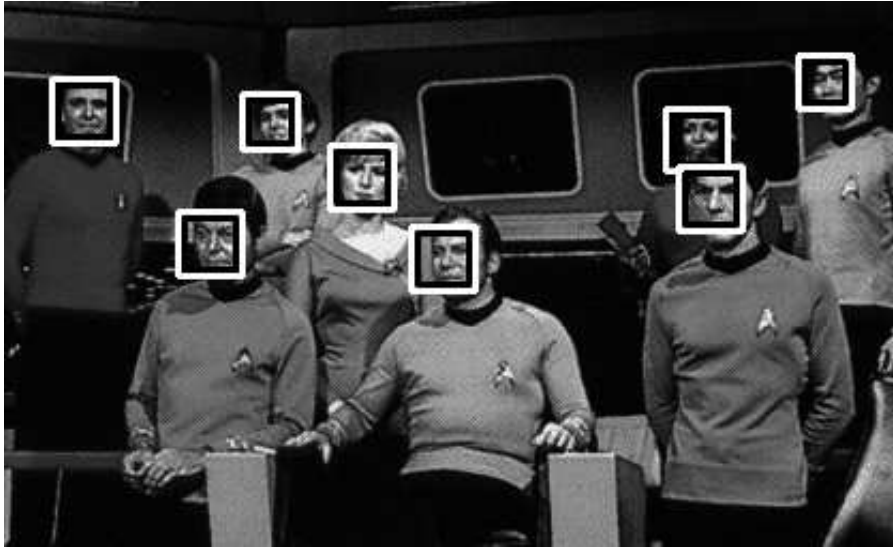


Relevant feature

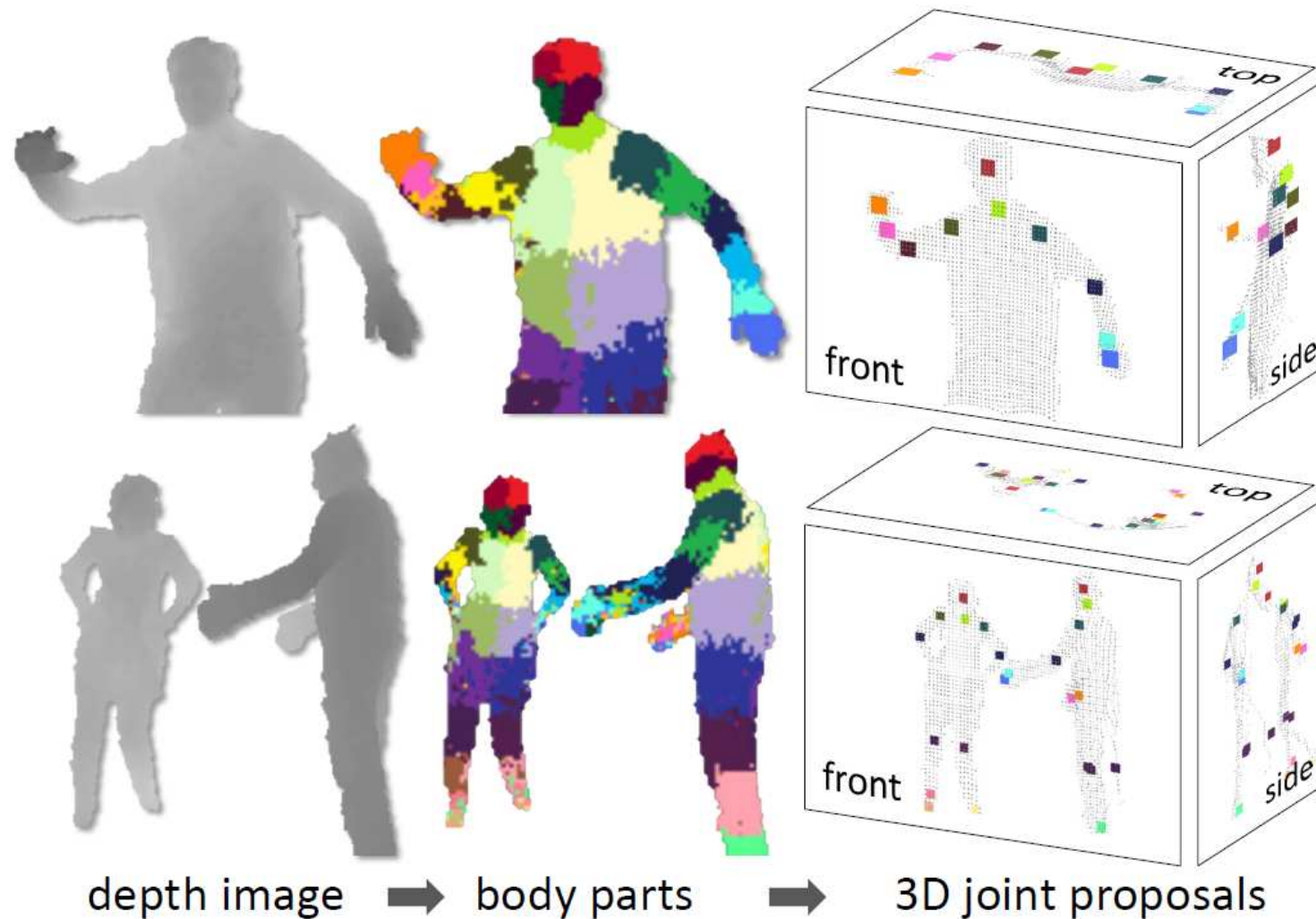


Irrelevant feature

Object detection



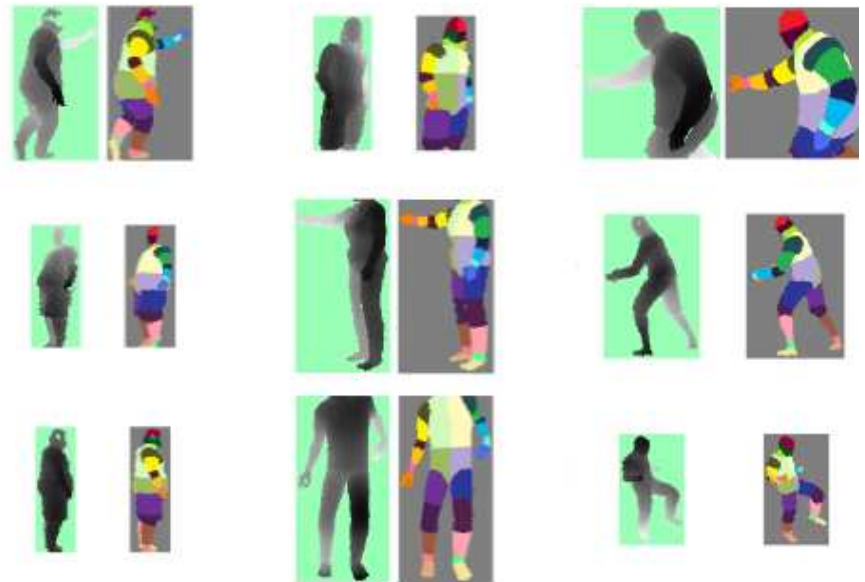
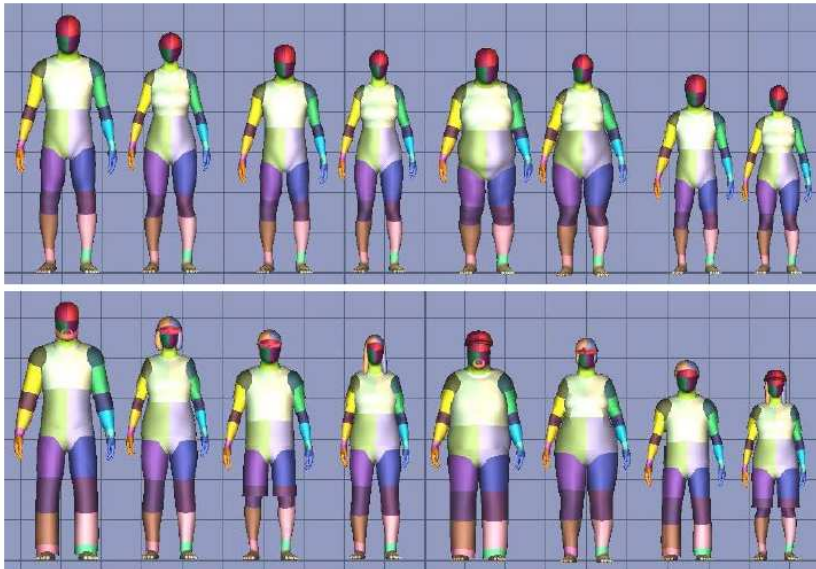
Random Forests and the Kinect



[Jamie Shotton et al 2011]

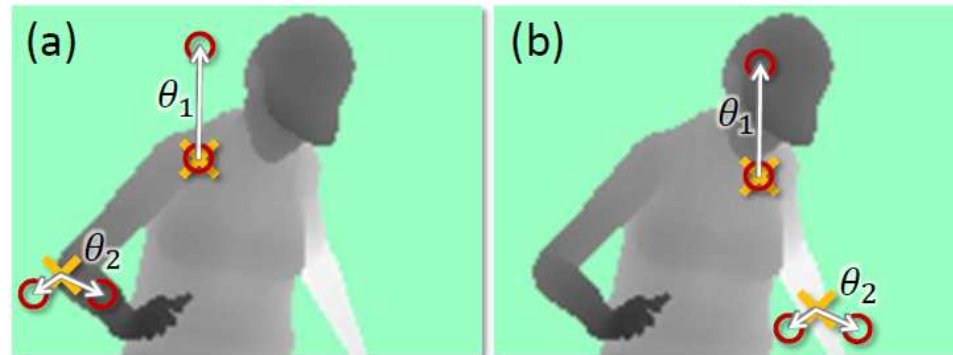
Random Forests and the Kinect

Lesson 1: Use computer graphics to generate plenty of data.



Random Forests and the Kinect

Lesson 2: Use simple depth features within random forests algorithm.



□ For each pixel x , compute the feature:

$$f_{\theta}(I, \mathbf{x}) = d_I \left(\mathbf{x} + \frac{\mathbf{u}}{d_I(\mathbf{x})} \right) - d_I \left(\mathbf{x} + \frac{\mathbf{v}}{d_I(\mathbf{x})} \right)$$

$d_I(x)$ is the depth at pixel x in image I

Parameters $\theta = (u; v)$ describe offsets u and v .

□ The normalization of the offsets ensures the features are depth invariant: At a given point on the body, a fixed world space offset will result whether the pixel is close or far from the camera.

Tree algorithm

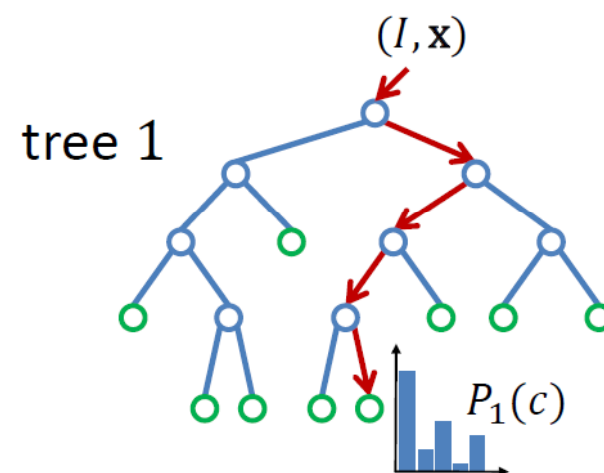
1. Randomly propose a set of splitting candidates $\phi = (\theta, \tau)$ (feature parameters θ and thresholds τ).
2. Partition the set of examples $Q = \{(I, \mathbf{x})\}$ into left and right subsets by each ϕ :

$$Q_l(\phi) = \{ (I, \mathbf{x}) \mid f_\theta(I, \mathbf{x}) < \tau \}$$
$$Q_r(\phi) = Q \setminus Q_l(\phi)$$

3. Compute the ϕ giving the largest gain in information:

$$\phi^* = \underset{\phi}{\operatorname{argmax}} G(\phi)$$
$$G(\phi) = H(Q) - \sum_{s \in \{l, r\}} \frac{|Q_s(\phi)|}{|Q|} H(Q_s(\phi))$$

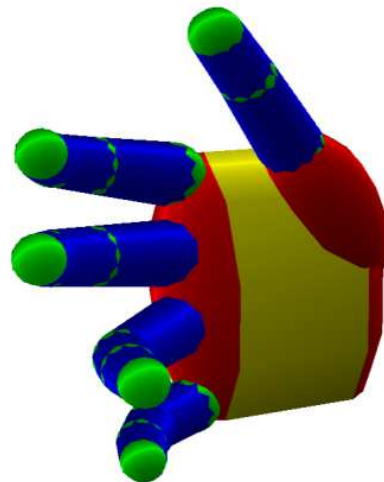
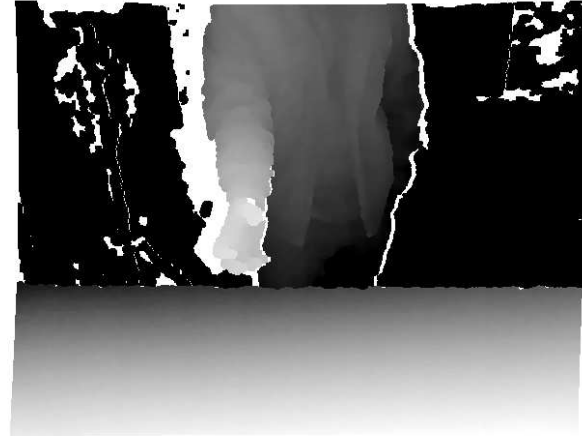
4. If the largest gain $G(\phi^*)$ is sufficient, and the depth in the tree is below a maximum, then recurse for left and right subsets $Q_l(\phi^*)$ and $Q_r(\phi^*)$.



Performance on train and test data

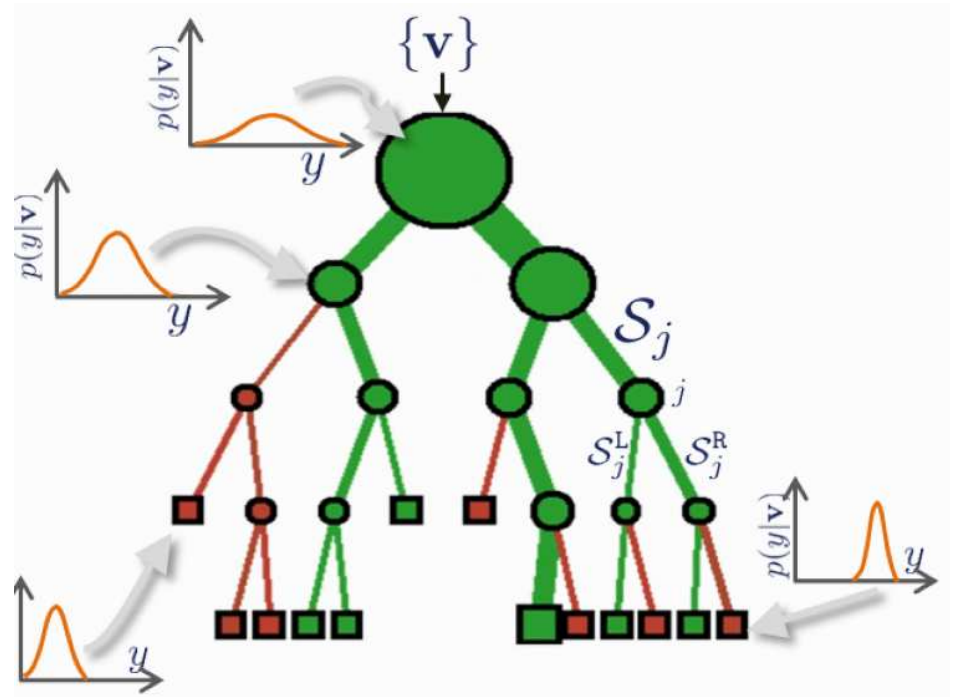
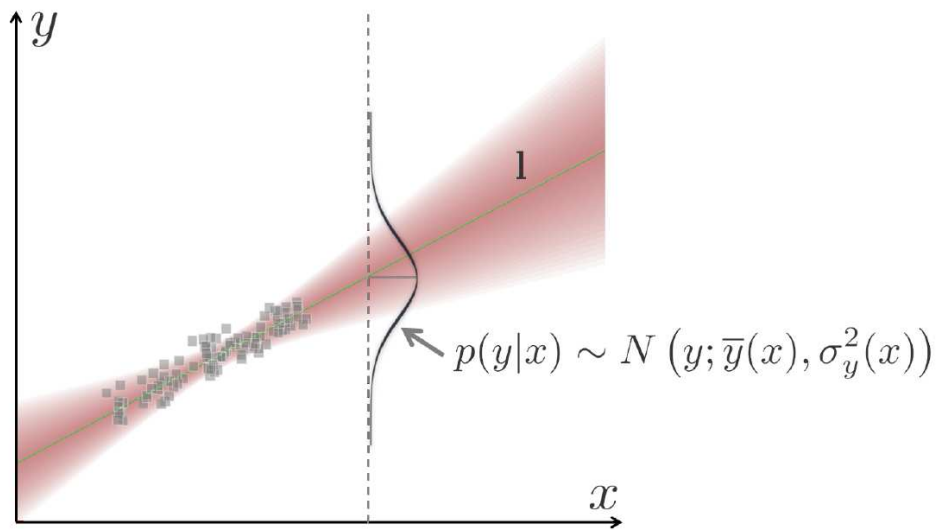
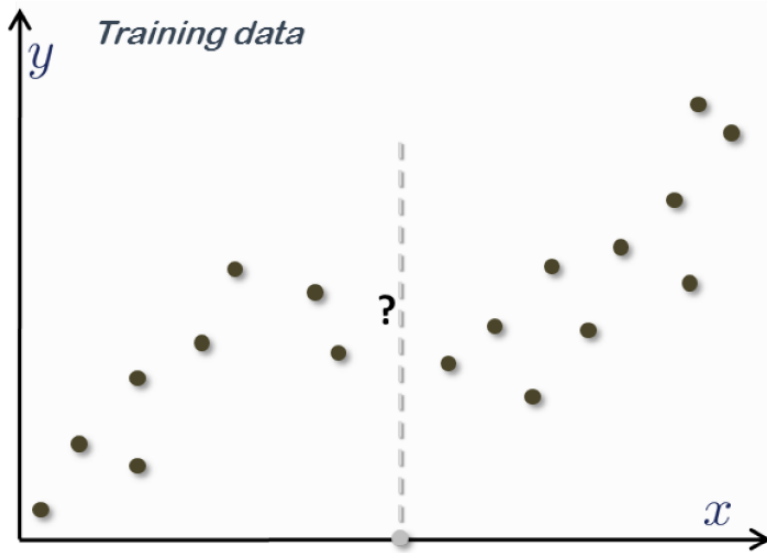


Applications: Interfaces

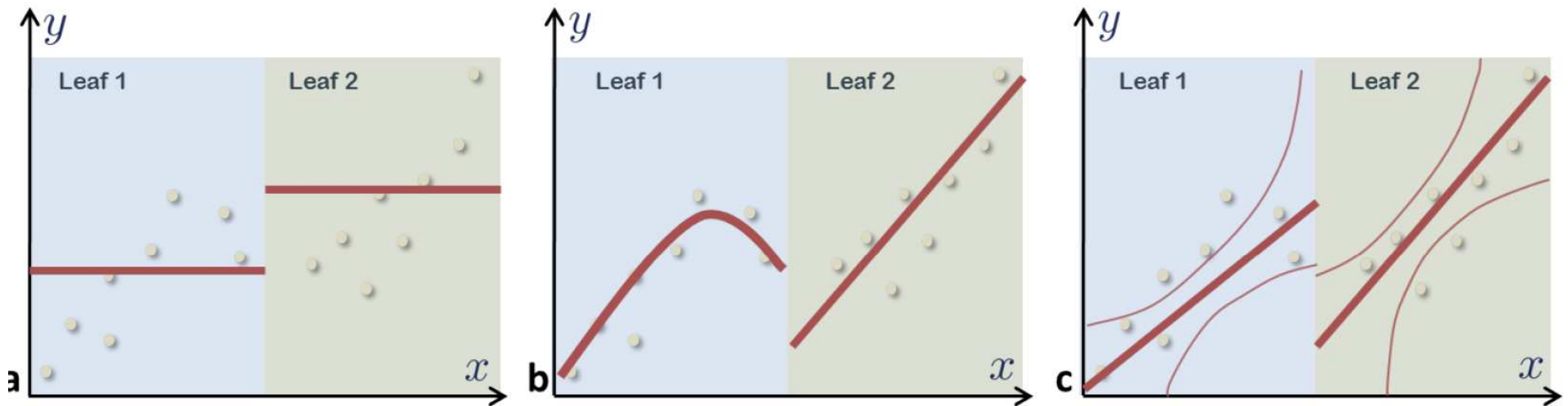


[Jason Oikonomidis et al 2011]

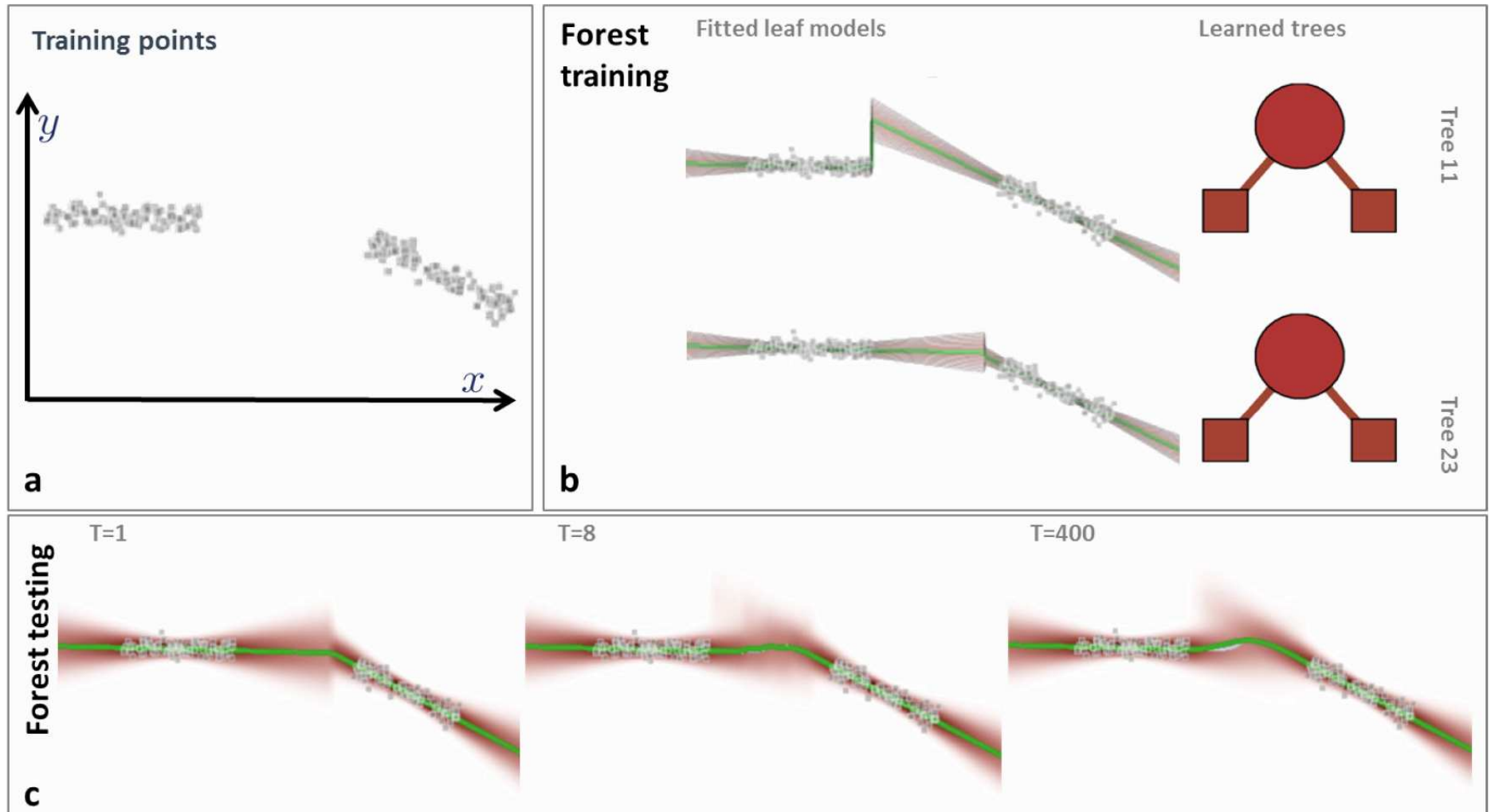
Trees for regression



Regression trees



Regression forests



Next lecture

On Thursday, we will embark on unconstrained optimization. This technique will enable us to train neural networks.