

CPSC340



#### Decision trees & Random Forests



Nando de Freitas November, 2012 University of British Columbia

# Outline of the lecture

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

Classification trees
Random forests
Object detection
Kinect

### Classification tree



A generic data point is denoted by a vector  $\mathbf{v} = (x_1, x_2, \cdots, x_d)$  $S_j = S_j^{L} \cup S_j^{R}$ 



#### Advanced: Gaussian information gain to decide splits



$$H(\mathcal{S}) = \frac{1}{2} \log \left( (2\pi e)^d |\Lambda(\mathcal{S})| \right)$$

Each split node j is associated with a binary split function





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

#### Alternative node decisions



$$\mathbf{v} = (x_1 \ x_2) \in \mathbb{R}^2$$

#### Random Forests for classification or regression

- 1. For b = 1 to B:
  - (a) Draw a bootstrap sample  $\mathbf{Z}^*$  of size N from the training data.
  - (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  $\boldsymbol{\theta}$  then when training the  $j^{\text{th}}$  node we only make available a small subset  $\mathcal{T}_j \subset \mathcal{T}$  of such values.

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

### 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.



### Effect of forest size



### Effect of more classes and noise



### Effect of tree depth (D)





### Effect of bagging

# Application to face detection

#### Training Data

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



[Viola and Jones, 2001]

# **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 (theta).



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

 $\mathbf{2}$ 

$$h(x) = \begin{cases} 1 & \sum_{t=1}^{t} \alpha_t n \\ 0 & \text{otherwise} \end{cases}$$





**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.



 $\Box$ For each pixel *x*, compute the feature:

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

dI(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.

[Jamie Shotton et al 2011]

# Tree algorithm

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

 $Q_{l}(\phi) = \{ (I, \mathbf{x}) \mid f_{\theta}(I, \mathbf{x}) < \tau \}$  $Q_{r}(\phi) = Q \setminus Q_{l}(\phi)$ 

3. Compute the  $\phi$  giving the largest gain in information:

$$\phi^{\star} = \operatorname{argmax}_{\phi} G(\phi)$$
$$G(\phi) = H(Q) - \sum_{s \in \{1, r\}} \frac{|Q_s(\phi)|}{|Q|} H(Q_s(\phi))$$

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  $Q_1(\phi^*)$  and  $Q_r(\phi^*)$ .

[Jamie Shotton et al 2011]

 $(I,\mathbf{x})$ 

tree 1

0

### Performance on train and test data



[Jamie Shotton et al 2011]

# **Applications:** Interfaces





[Iason Oikonomidis et al 2011]

### Next lecture

There is no next lecture. Congrats!