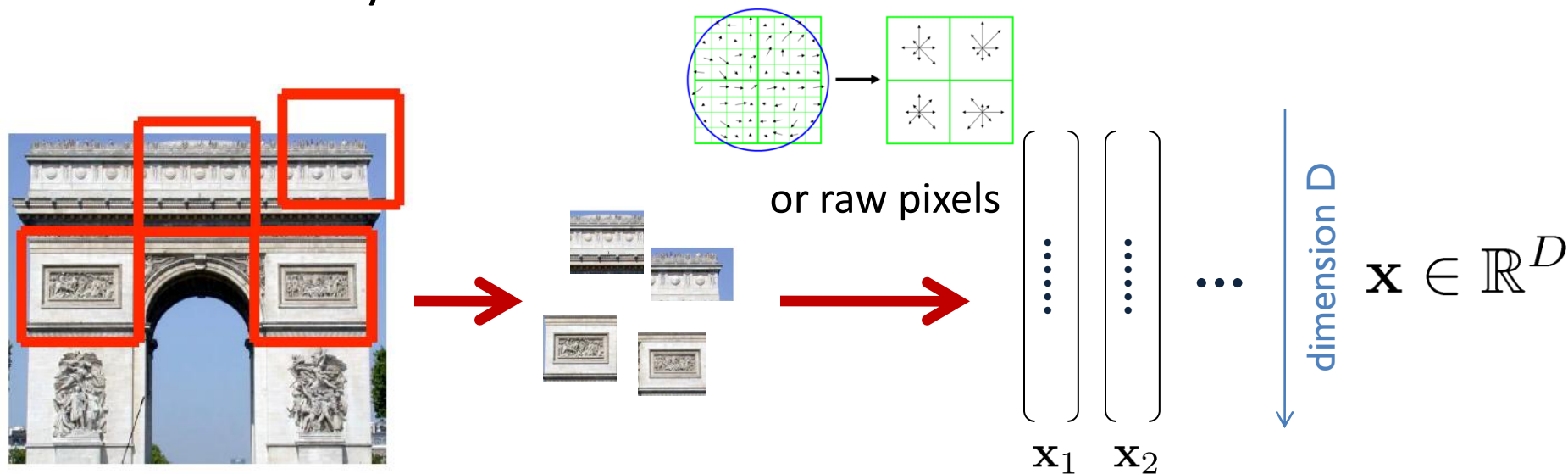# Visual Codebook

Tae-Kyun Kim

Sidney Sussex College

# Visual Words

- Visual words are base elements to describe an image.

- Interest points are detected from an image
  - Corners, Blob detector, SIFT detector

- Image patches are represented by descriptor
  - SIFT (Scale-Invariant Feature Transform) or Raw pixel intensity

or raw pixels

dimension D

$\mathbf{x} \in \mathbb{R}^D$

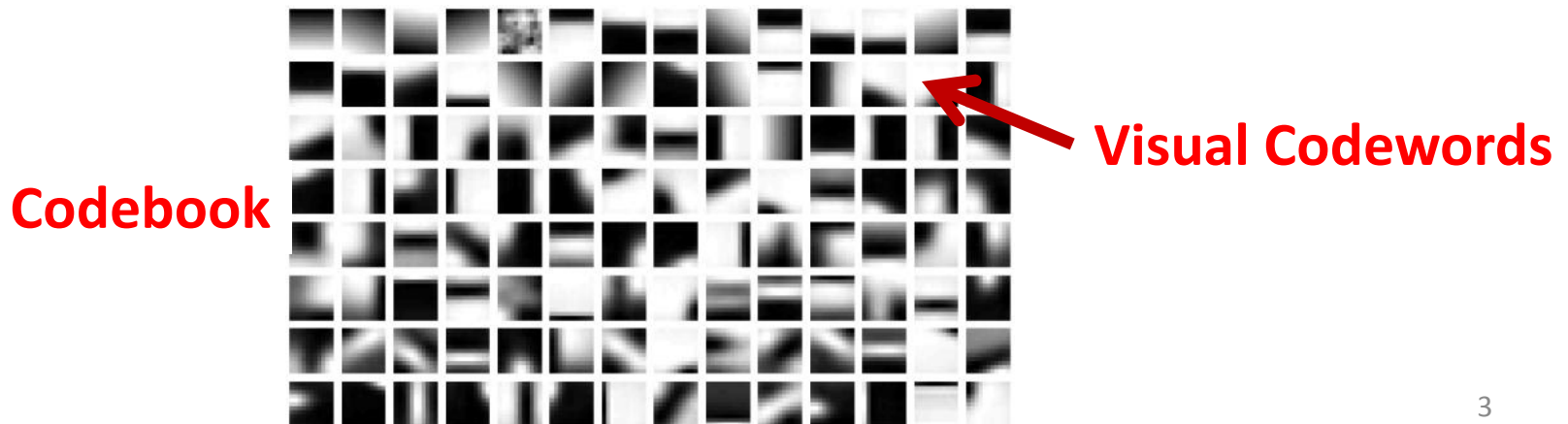$\mathbf{x}_1 \quad \mathbf{x}_2$

# Building a Visual Codebook (Dictionary)

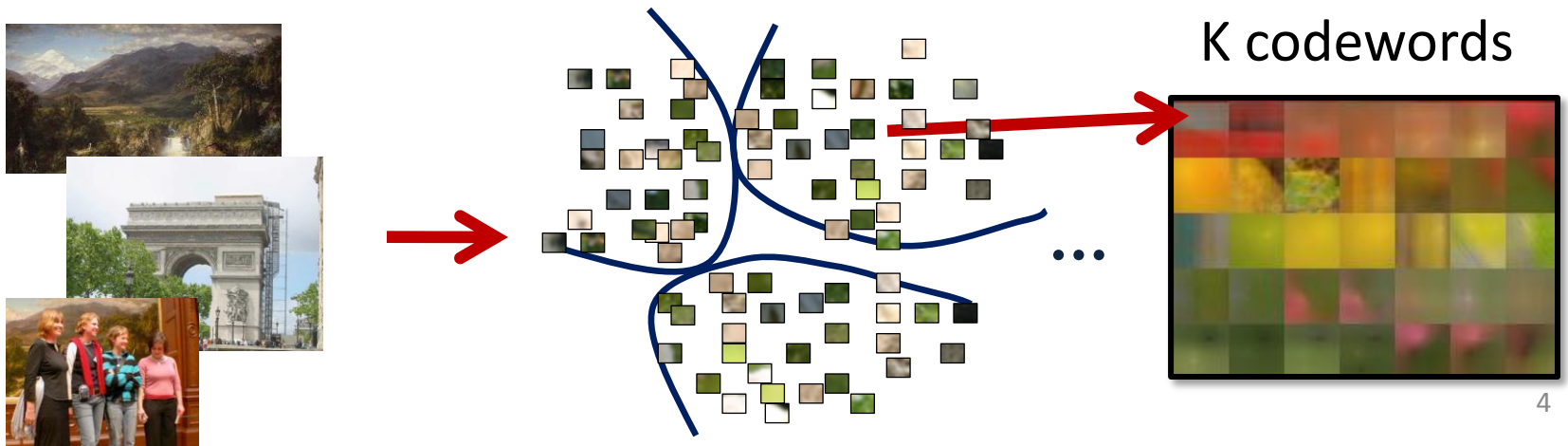- Visual words (real-valued vectors) can be compared using Euclidean distance:

$$E(\mathbf{x}_1, \mathbf{x}_2) = \sqrt{\sum_d (\mathbf{x}_1^d - \mathbf{x}_2^d)^2}$$

- These vectors are divided into groups which are similar, essentially clustered together, to form a codebook.
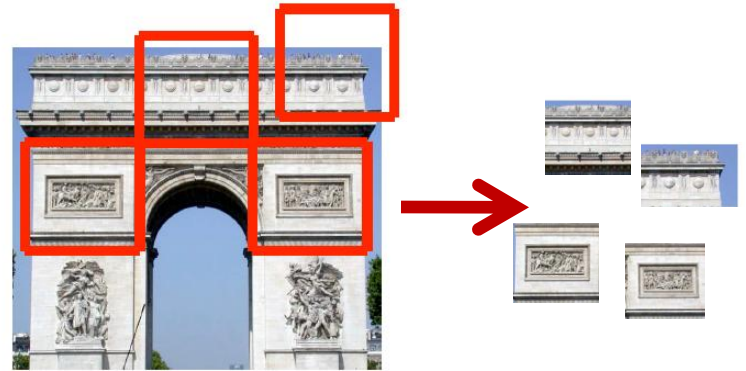
**Codebook**

**Visual Codewords**

3

# K-means Clustering

- The two steps repeated until no vector changes membership:

  1. Compute a cluster center for each cluster as the mean of the cluster members

  2. Reassign each data point to the cluster whose center is nearest

- The cluster centers (mean values) now form a visual dictionary.

K codewords

...
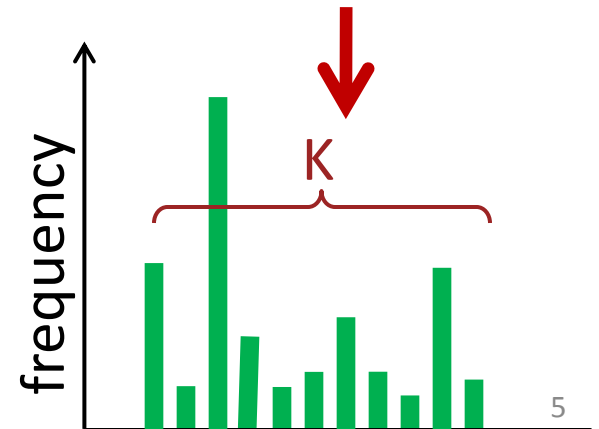
# Histogram of Visual Words

- Every visual word is compared with codewords and assigned to the nearest codeword.

- Histogram bins are codewords and each bin counts the number of words assigned to the codeword.
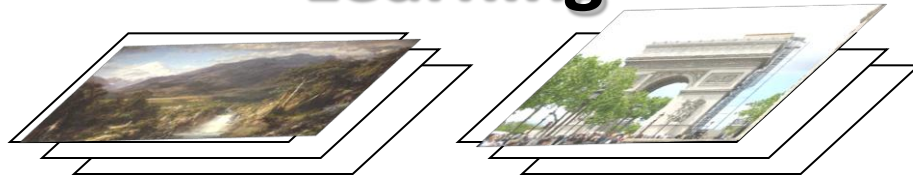
Nearest Neighbour Matching

**"bag of words"**

# Categorisation

**Learning**

**Recognition**

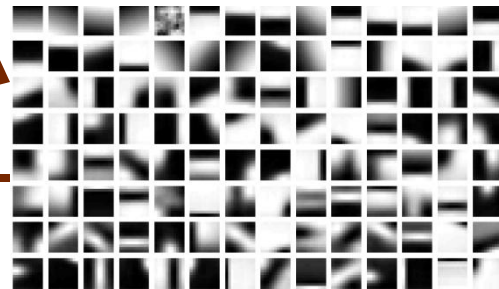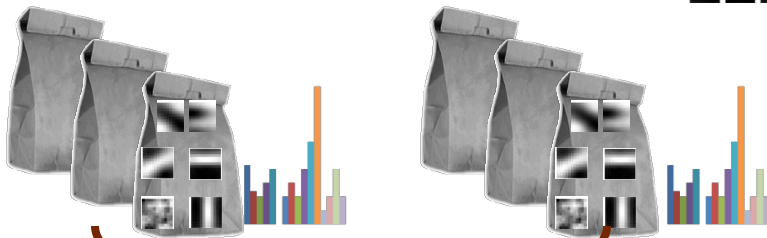feature detection & representation

**codewords dictionary**

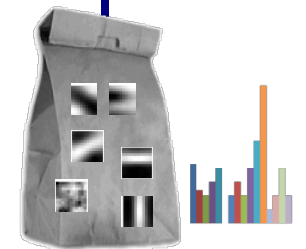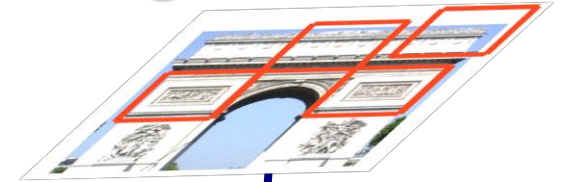image representation
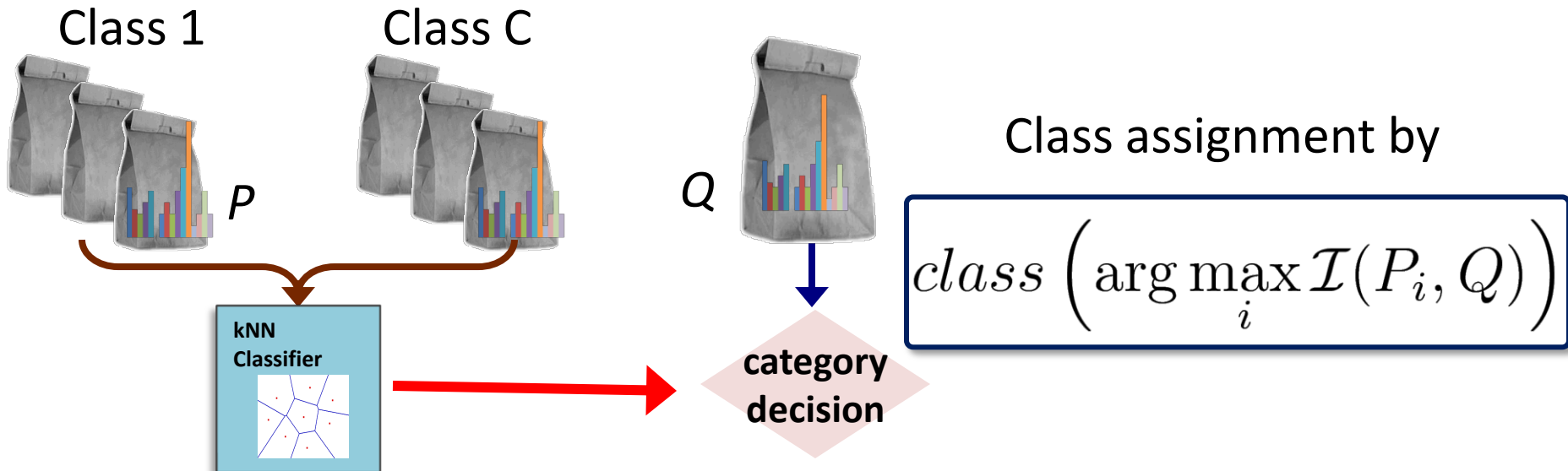
**category models (and/or) classifiers**

**category decision**

# Categorisation

- Histogram matching of a pair of images *P* and *Q* is

$$\mathcal{I}(P, Q) = \sum_{n=1}^{K} \min \left( H^{(n)}(P), H^{(n)}(Q) \right)$$

- Based on HMK, we can use various classifiers such as kNN, SVM, Naïve-Bayes classifiers or pLSA, LDA.

Class 1    Class C

*P*    *Q*

Class assignment by

$$class \left( \arg \max_{i} \mathcal{I}(P_i, Q) \right)$$

**kNN Classifier**

**category decision**

# Summary of K-means Codebook

- The histogram is 1-D, a highly compact and robust representation of images.

- The histogram representation greatly facilitates modelling images and their categories.

- Codeword assignment (quantisation process) by K-means is time-demanding.

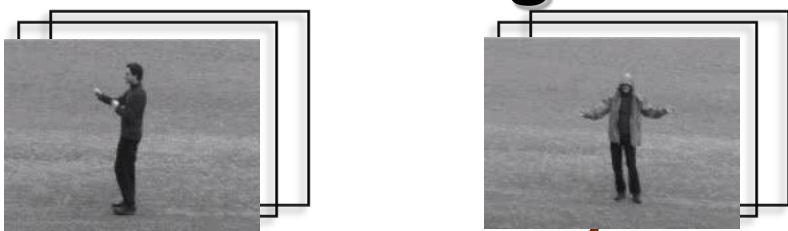- K-means is an unsupervised learning method.

# **Case Study:**
# Video (action) Categorisation by RF Codebook

In collaboration with Tsz Ho Yu
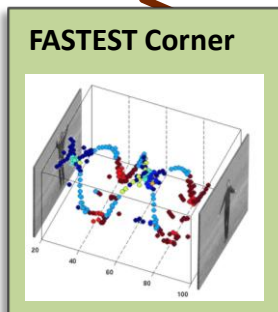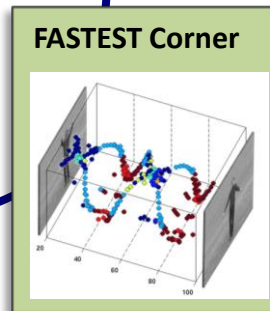
# Demo video: Action categorisation

# Learning

# Recognition
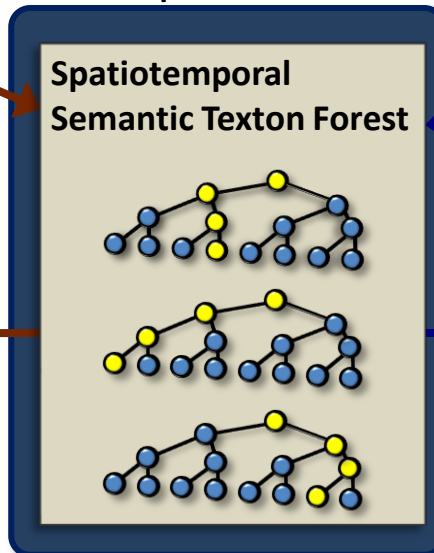
Interest point detection

**FASTEST Corner**

Descriptor/Codebook

**FASTEST Corner**

**Spatiotemporal Semantic Texton Forest**

Bag of semantic textons

✓ Powerful Discriminative Codebook
✓ Extremely Fast

Classification

**kNN Classifier**

**category decision**

# Relation to the previous

| | 2D → 3D | |
|---|---|---|
| **Input** | Images | Videos |
| **Interest Point** | Corners | 3D corners |
| **Descriptor** | Patches | Space-time volumes |
| | SIFT | Raw pixels |
| **Codebook** | K-means | *Randomised Decision Forests* |
| **Representation** | Histogram of visual words | |
| **Classifier** | kNN classifier | |

**Skip**

# Relation to the previous

## *Patches ↔ Space-Time Volumes*

**Image (2D)**

**Video (3D)**



↔

**Patch**

**Cuboid**

*or*

$$\mathbf{x} \in \mathbb{R}^D$$

**"visual words"**

# Data Set and Interest Point

# Action data set

- http://www.nada.kth.se/cvap/actions/
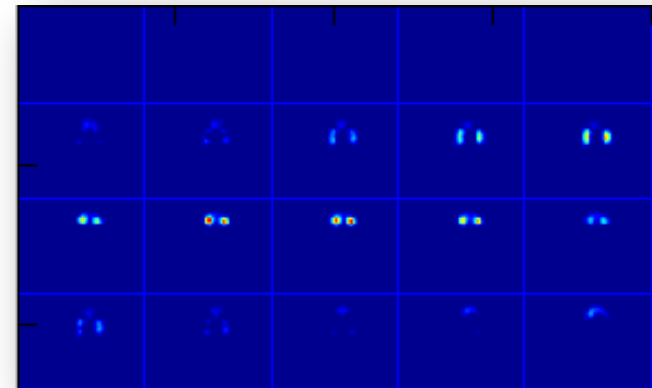- 25 subjects, 2391 sequences – the largest public action DB

# Space-Time Interest Point

- In Dollar et al. 2005, separable linear filters are applied.

- The response function is
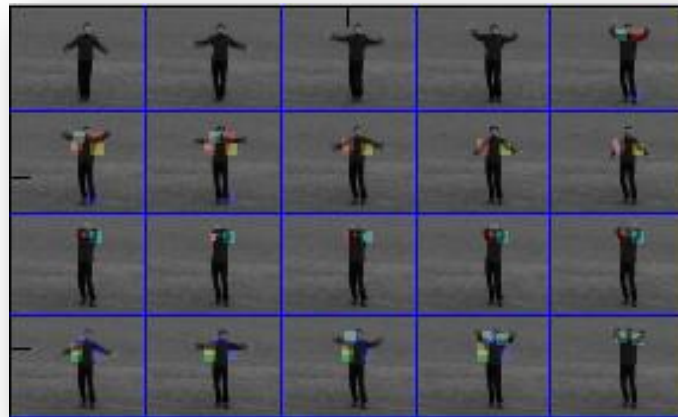
$$R = (I * G * h_{ev})^2 + (I * G * h_{od})^2$$

  - $G(x, y; \sigma)$ is **2D Gaussian smoothing** kernel
  - $h_{ev}/h_{od}$ are a pair of **1D Gabor filters temporally applied** as

$$h_{ev}(t; \tau, \omega) = -\cos(2\pi t\omega)e^{-t^2/\tau^2} \qquad h_{od}(t; \tau, \omega) = -\sin(2\pi t\omega)e^{-t^2/\tau^2}$$

# Space-Time Interest Point

- For multiple scales, it runs the detector over a set of spatial and temporal scales.

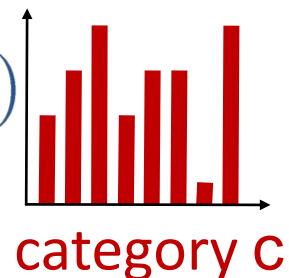- Interest points are detected as **local maxima** of the response-function.
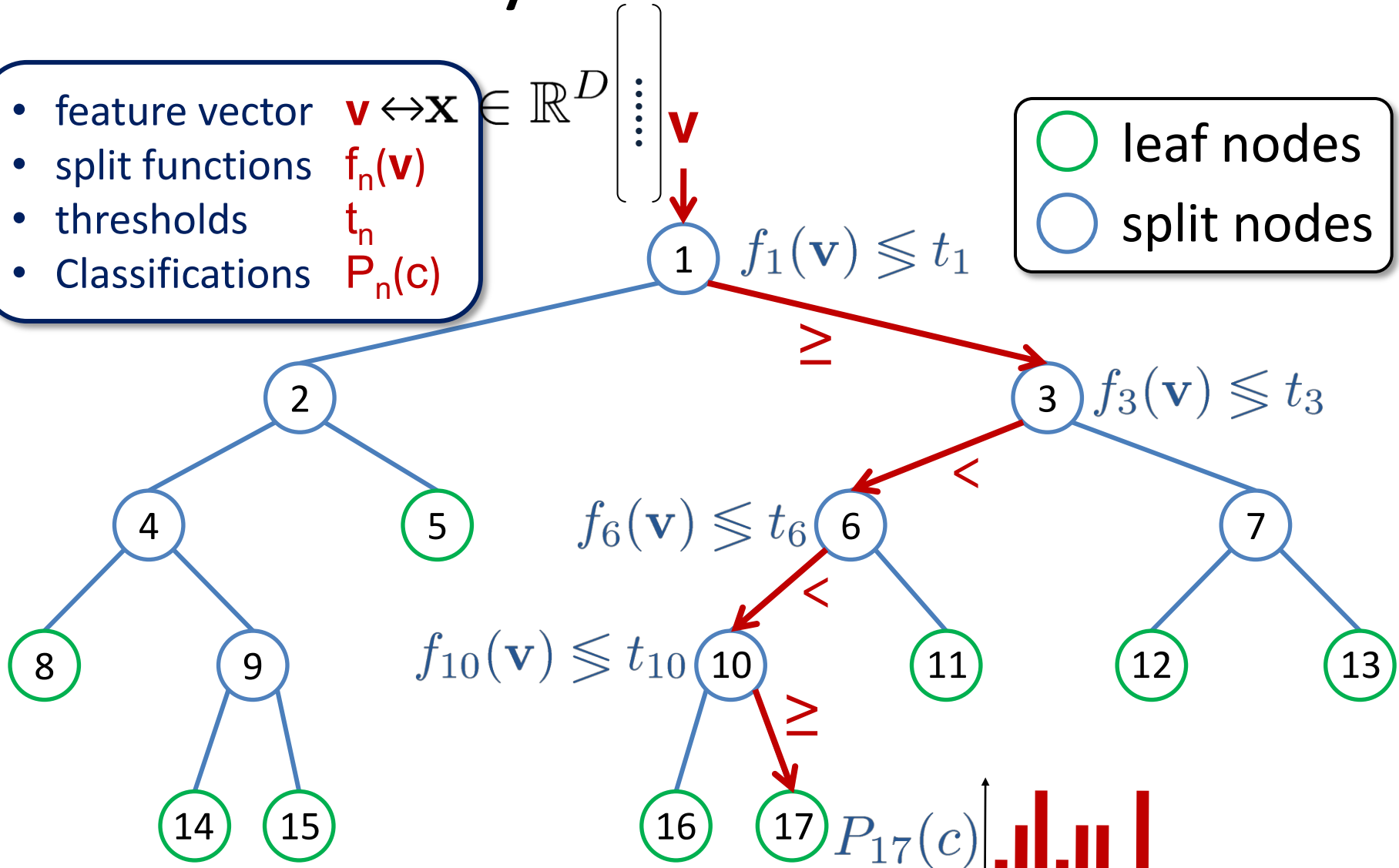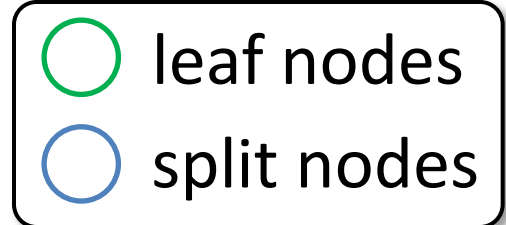
# Matlab Demo:

## *Space-Time Interest Points*
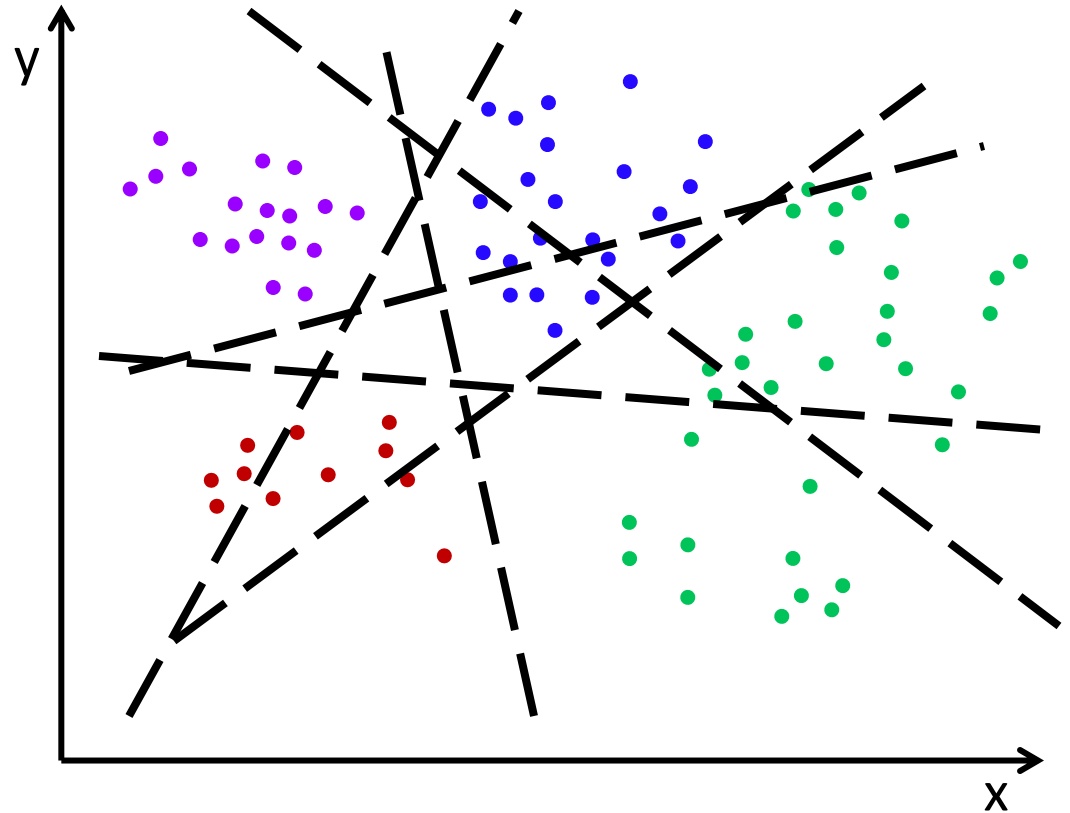
# Background:
# Randomised Decision Forest

# Binary Decision Trees

- feature vector $\mathbf{v} \leftrightarrow \mathbf{x} \in \mathbb{R}^D$ $\quad \mathbf{v}$
- split functions $f_n(\mathbf{v})$
- thresholds $t_n$
- Classifications $P_n(c)$

leaf nodes

split nodes

1 $\quad f_1(\mathbf{v}) \lessgtr t_1$

$\geq$

3 $\quad f_3(\mathbf{v}) \lessgtr t_3$

2

4

5

$f_6(\mathbf{v}) \lessgtr t_6$ 6 $\quad <$

7

8

9

$<$

$f_{10}(\mathbf{v}) \lessgtr t_{10}$ 10

11

12

13

$\geq$

14

15

16

17 $\quad P_{17}(c)$
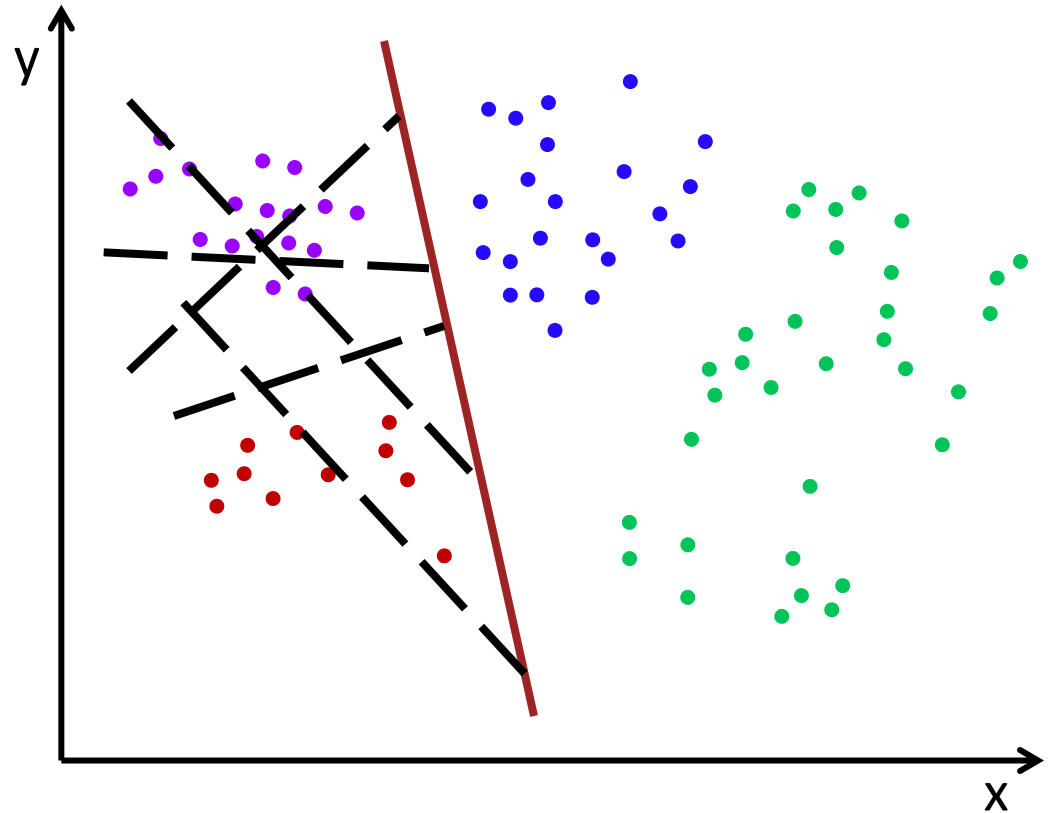
category c

# Toy Learning Example

- **Try several lines, chosen at random**

- **Keep line that best separates data**
  – information gain

- **Recurse**

- feature vectors are x, y coordinates: $\mathbf{v} = [x, y]^T$
- split functions are lines with parameters a, b: $f_n(\mathbf{v}) = ax + by$
- threshold determines intercepts: $t_n$
- four classes: purple, blue, red, green

Slide credits to J. Shotton

# Toy Learning Example

- **Try several lines, chosen at random**

- **Keep line that best separates data**
  - information gain

- **Recurse**



- feature vectors are x, y coordinates: $\mathbf{v} = [x, y]^T$
- split functions are lines with parameters a, b: $f_n(\mathbf{v}) = ax + by$
- threshold determines intercepts: $t_n$
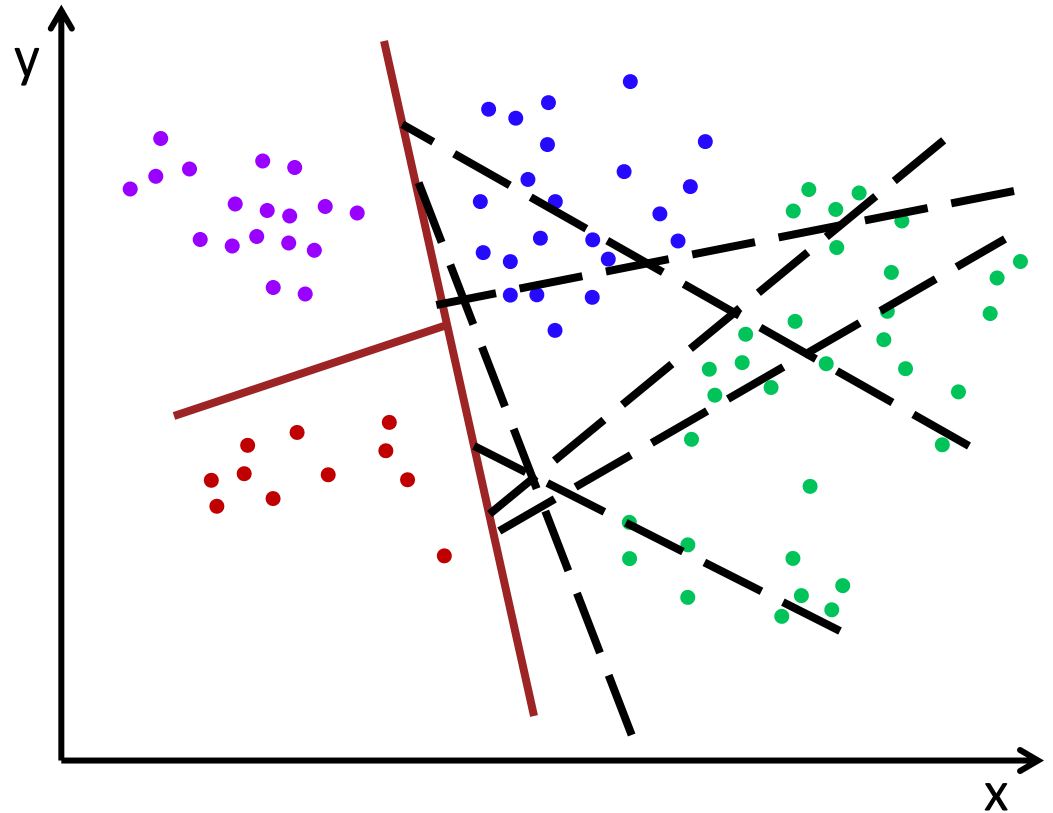- four classes: purple, blue, red, green
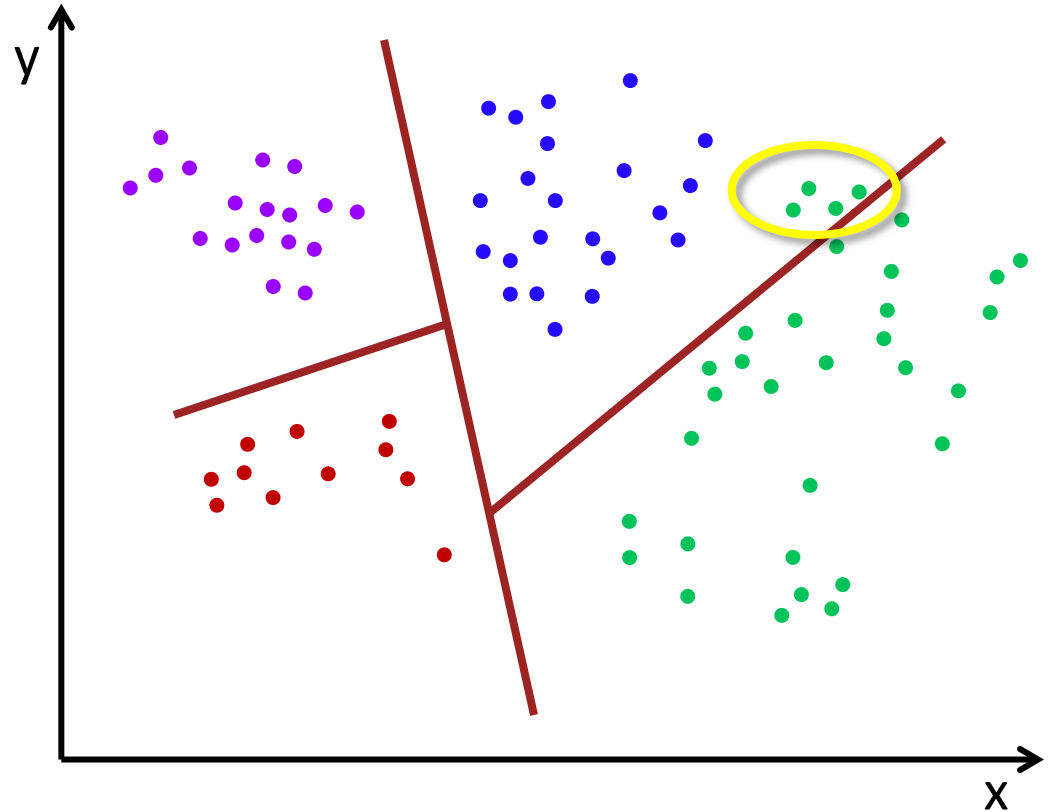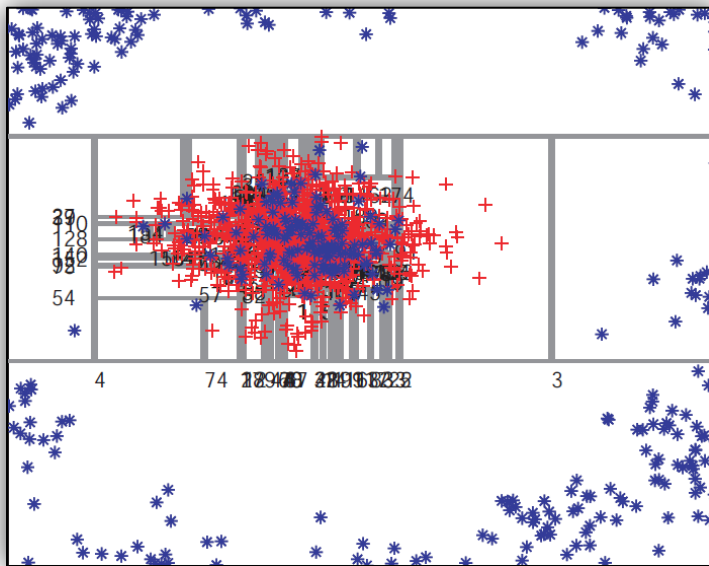
# Toy Learning Example

- **Try several lines, chosen at random**

- **Keep line that best separates data**
  - information gain

- **Recurse**

- feature vectors are x, y coordinates: $\mathbf{v} = [x, y]^T$
- split functions are lines with parameters a, b: $f_n(\mathbf{v}) = ax + by$
- threshold determines intercepts: $t_n$
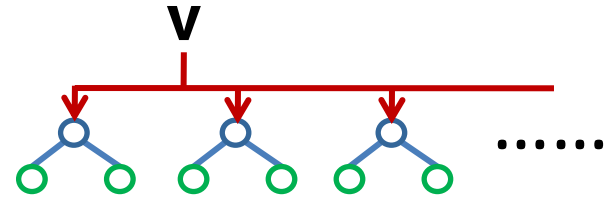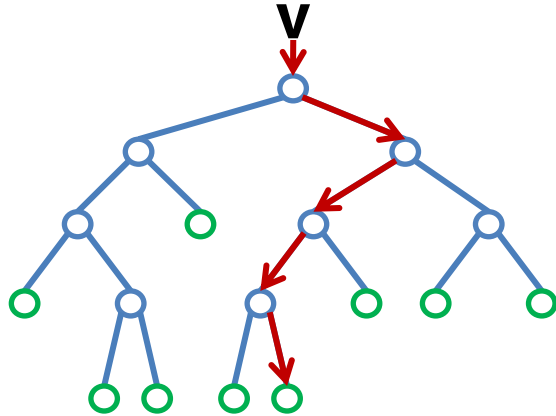- four classes: purple, blue, red, green

# Toy Learning Example

- Try several lines, chosen at random

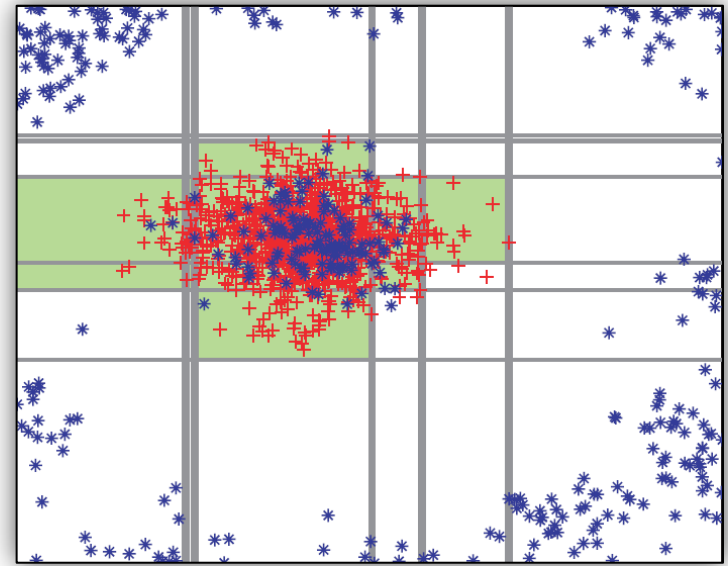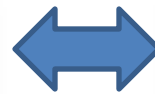- Keep line that best separates data
  - information gain

- Recurse



- feature vectors are x, y coordinates:  $\mathbf{v} = [x, y]^T$
- split functions are lines with parameters a, b:  $f_n(\mathbf{v}) = ax + by$
- threshold determines intercepts:  $t_n$
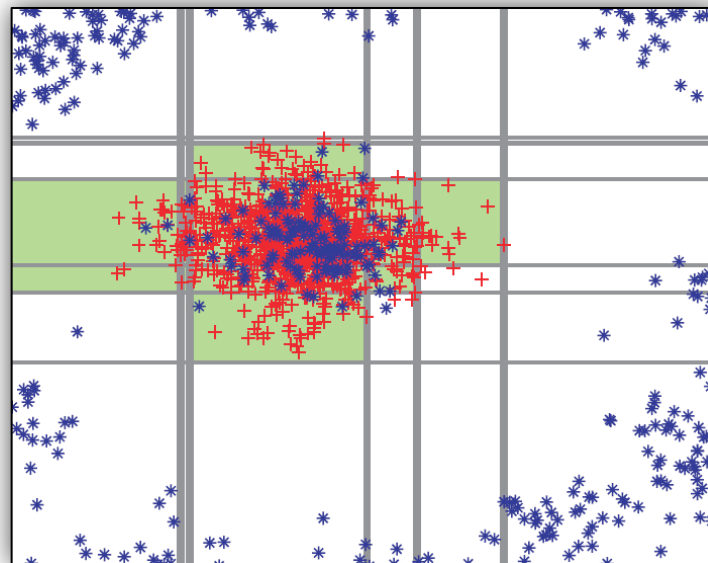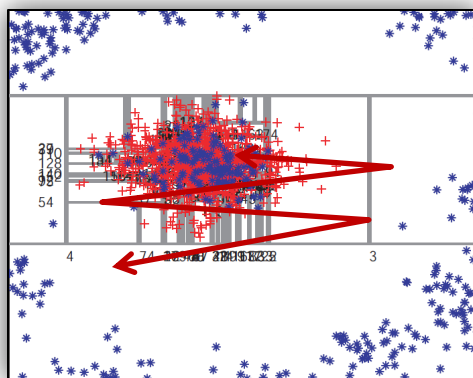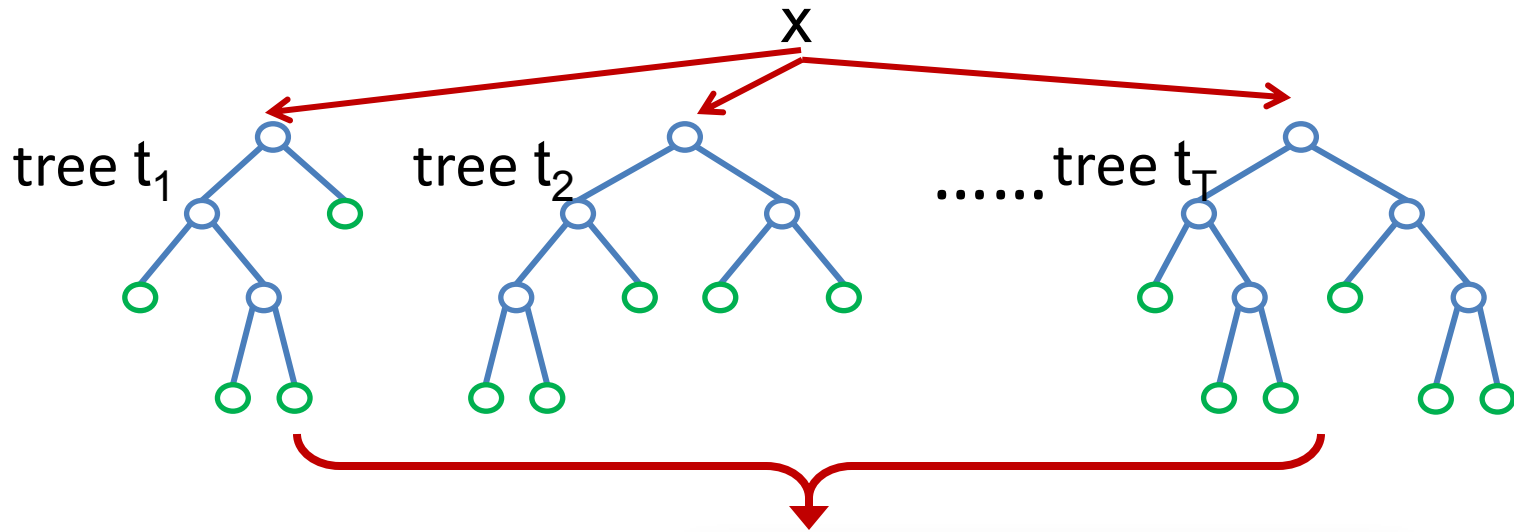- four classes: purple, blue, red, green

# Generalisation



**Overfit**

**Reasonably Smooth**

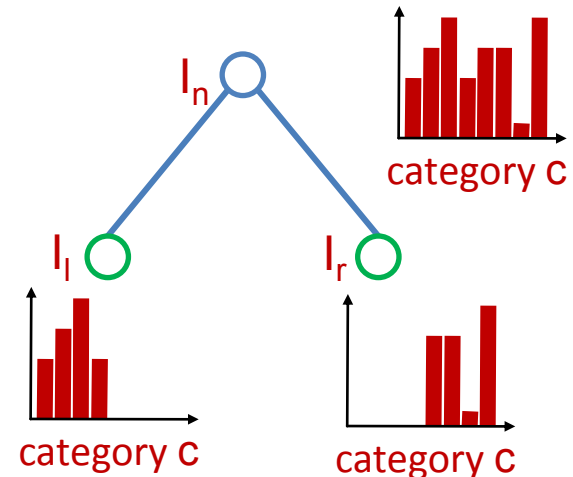# Generalisation



**Reasonably Smooth**

# Randomized Tree Learning

- Train data: $\{(\mathbf{x}_i, y_i)\}_{i=1}^N \quad y_i \in \{1, 2, ..., C\}$

- Recursive algorithm

    – set $I_n$ of training examples that reach node $n$ is split:

$$
\begin{aligned}
\text{left split} \quad I_l &= \{i \in I_n \mid f(\mathbf{v}_i) < t\} \\
\text{right split} \quad I_r &= I_n \setminus I_l
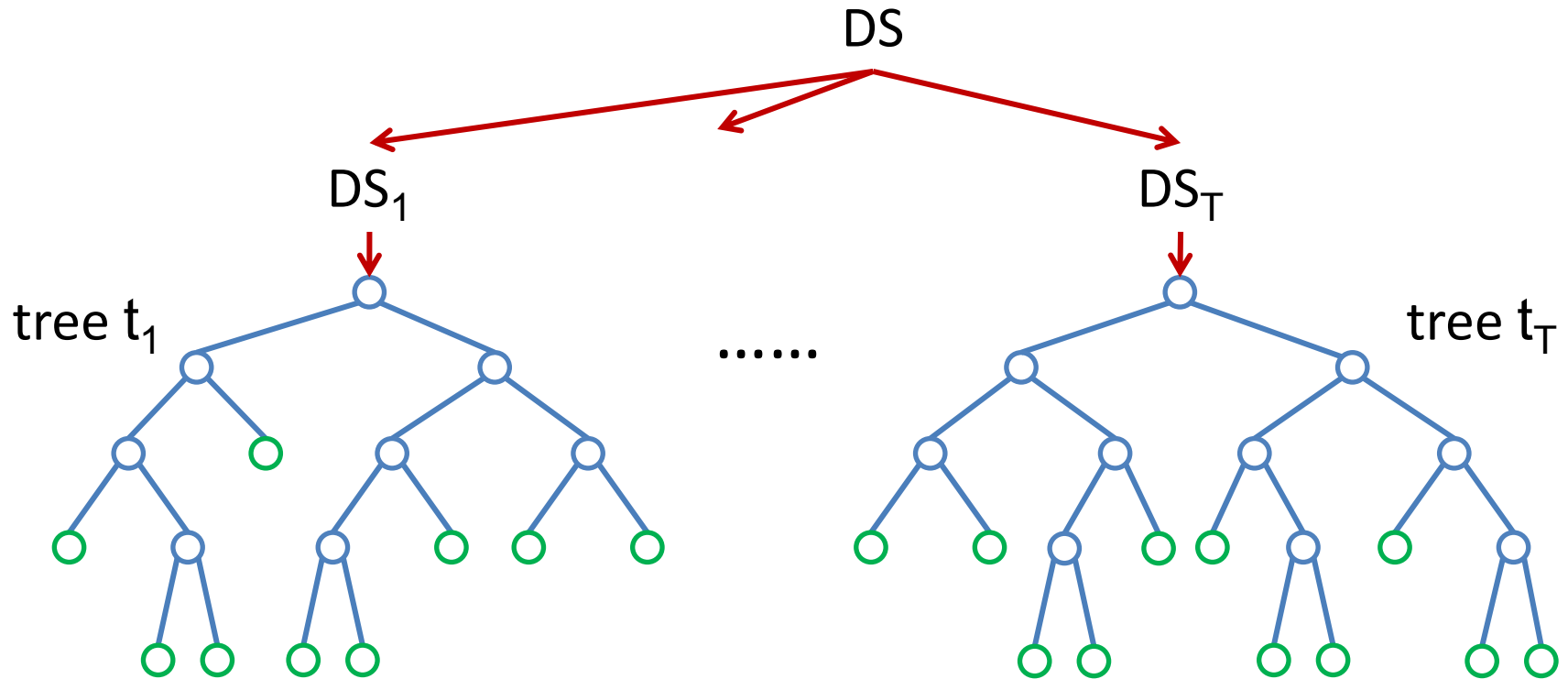\end{aligned}
$$

- Features $f$ and thresholds $t$ chosen at random

- Choose $f$ and $t$ to maximize gain in information

$$
\Delta E = -\frac{|I_l|}{|I_n|} E(I_l) - \frac{|I_r|}{|I_n|} E(I_r)
$$



category c

category c

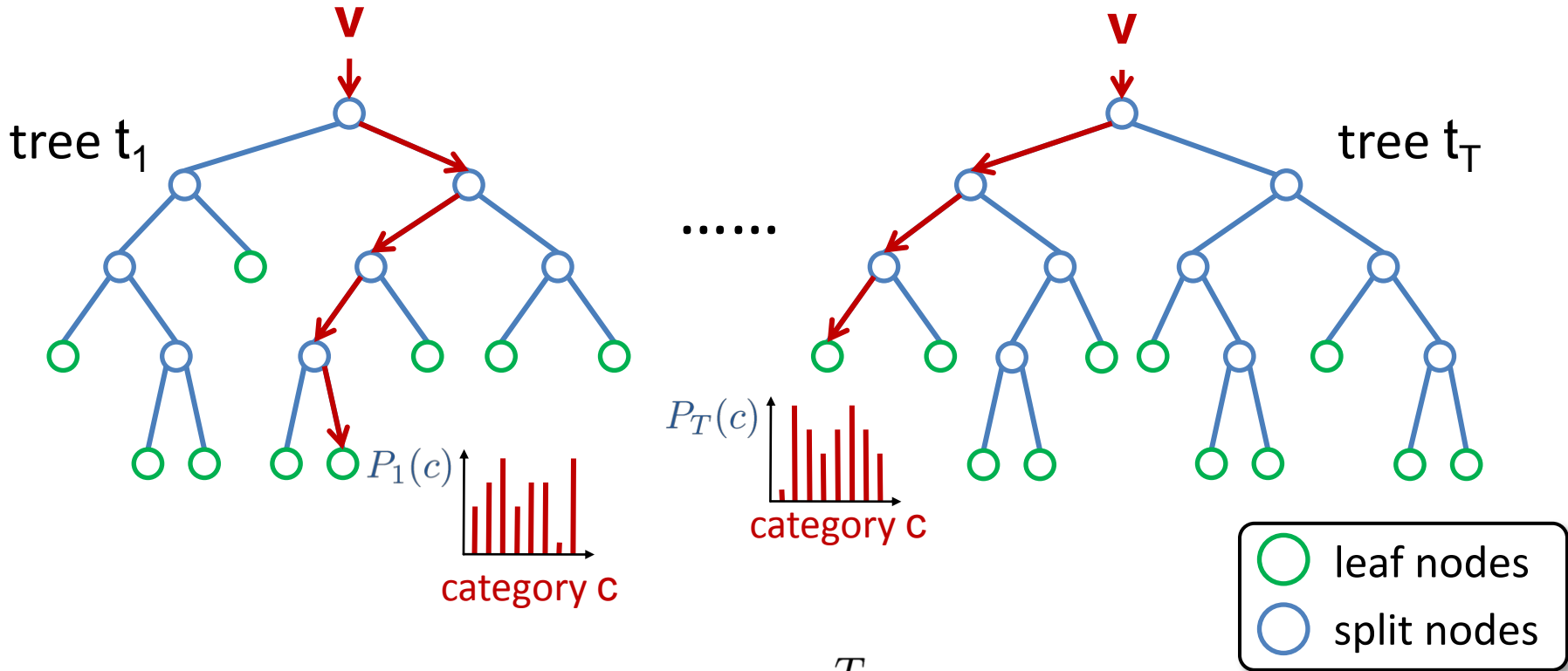category c

# A Forest of Trees : Learning

- Forest is an ensemble of decision trees



- Bagging (Bootstrap AGGregatING)
  - For each set, it randomly draws examples from the uniform dist. allowing duplication and missing.
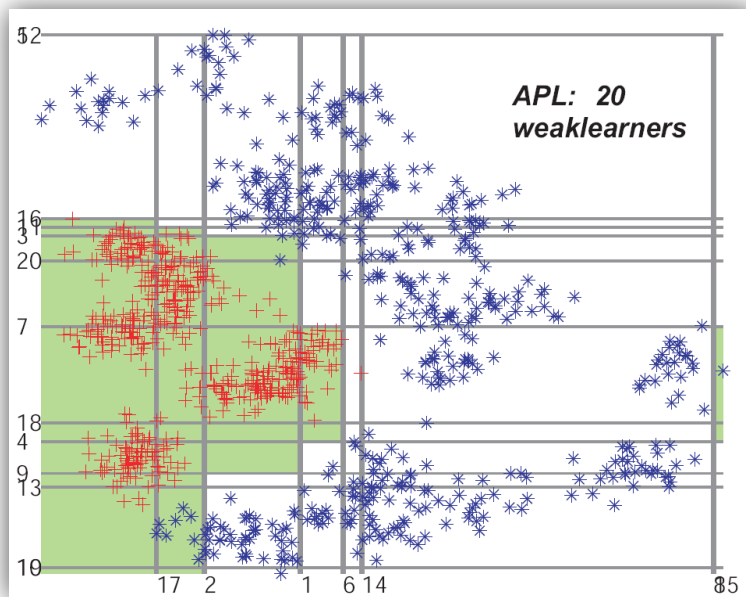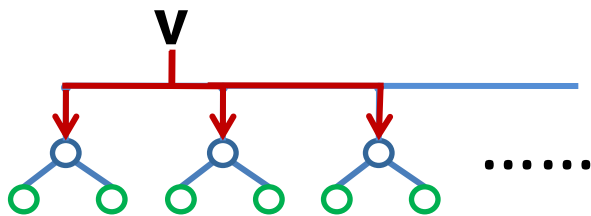
# A Forest of Trees : Recognition
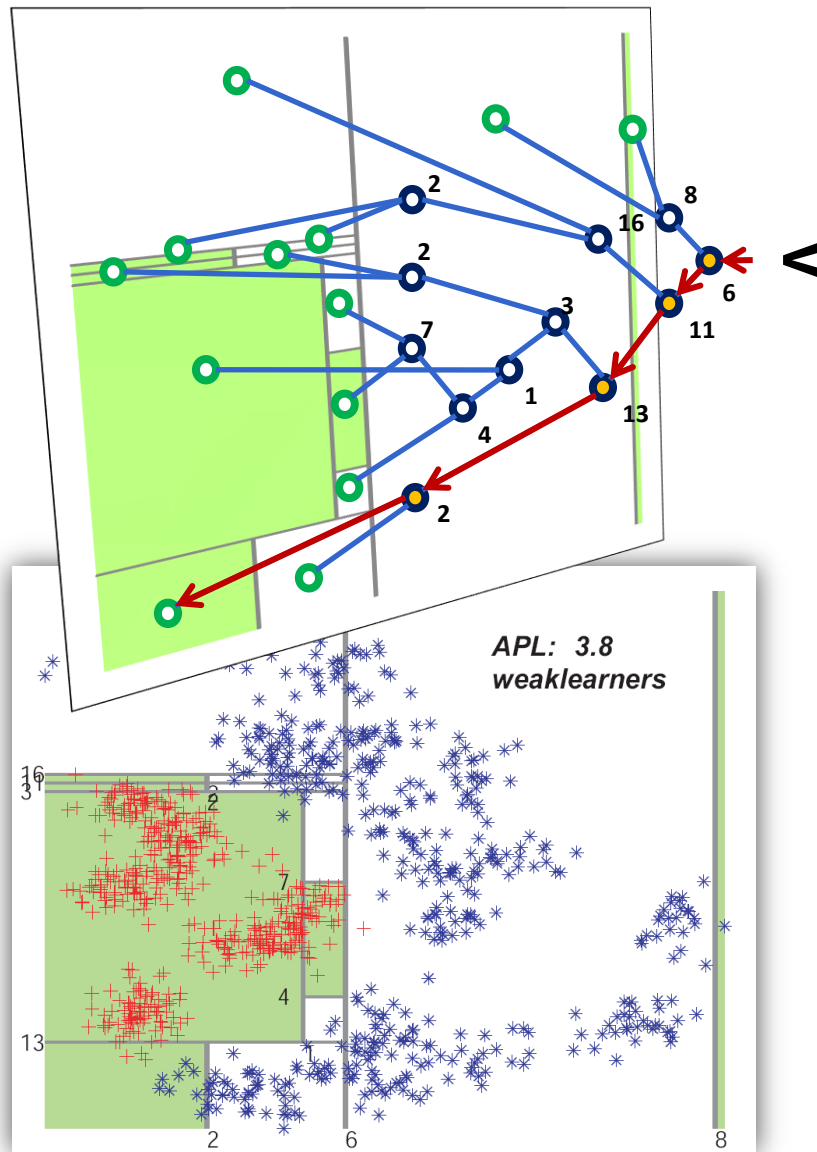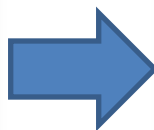
- Forest is an ensemble of decision trees



tree $t_1$     **V**       **V**     tree $t_T$

...... 

$P_1(c)$

$P_T(c)$

category **c**

category **c**

○ leaf nodes

○ split nodes

    – classification is     $P(c|\mathbf{v}) = \displaystyle\sum_{t=1}^{T} P_t(c|\mathbf{v})$

# Fast Evaluation



**V**

APL: 20 weaklearners
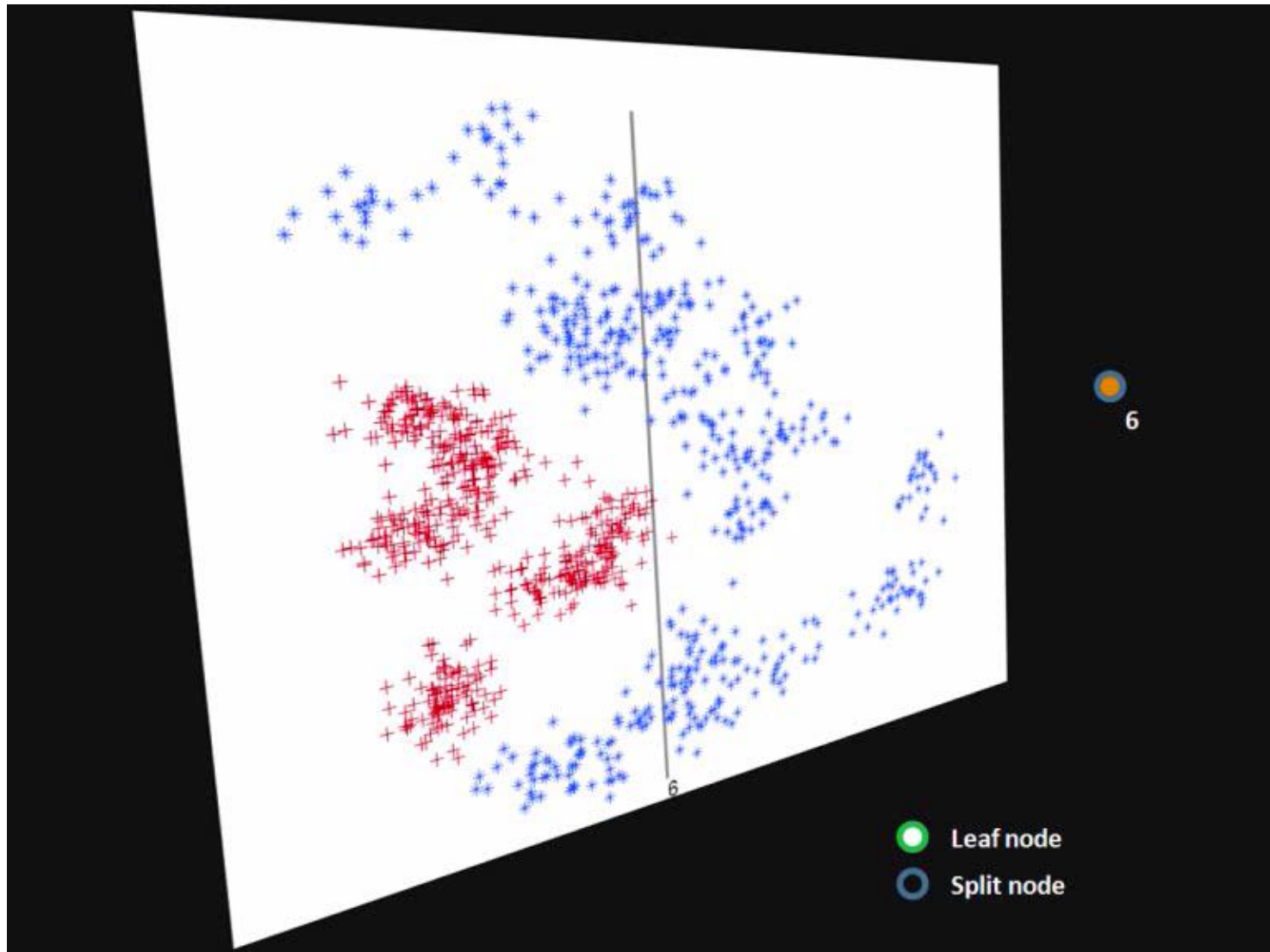
Flat structure

5 times speed up

**<**
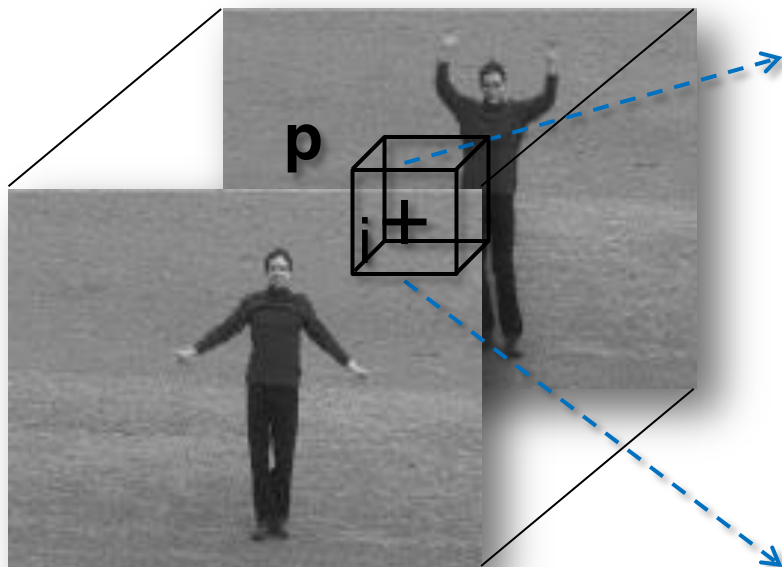
APL: 3.8 weaklearners

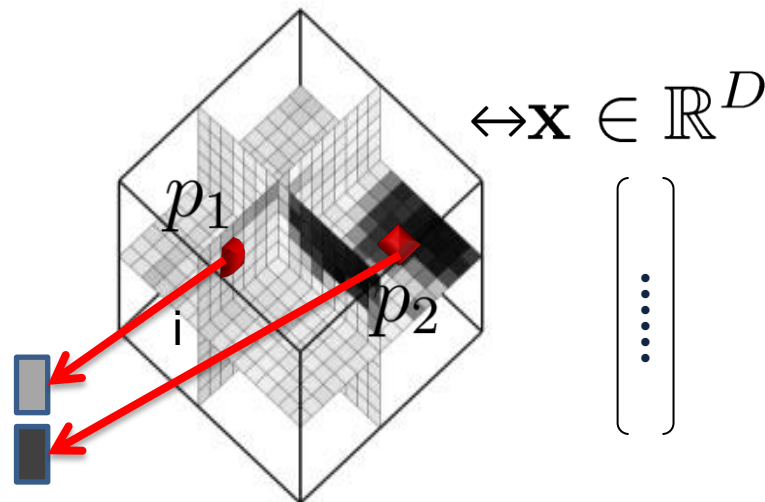Tree structure

# Demo video: Fast evaluation

# Random Forest Codebook

# Video Cuboid Features



**Video pixel i gives cuboid p**

(13x13x19 pixels in experiments)
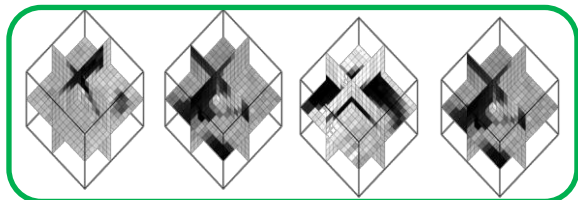
$\leftrightarrow \mathbf{x} \in \mathbb{R}^D$
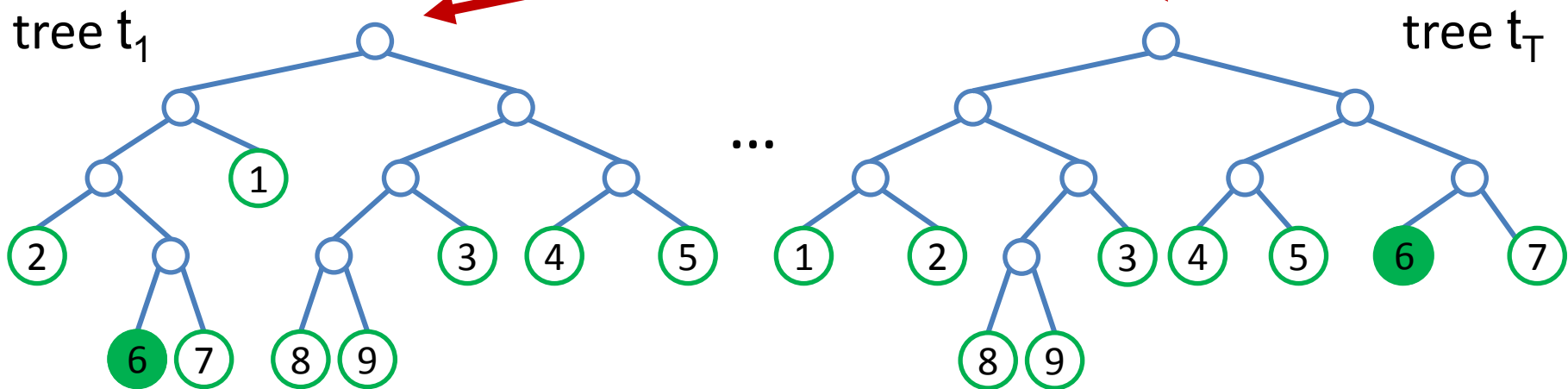
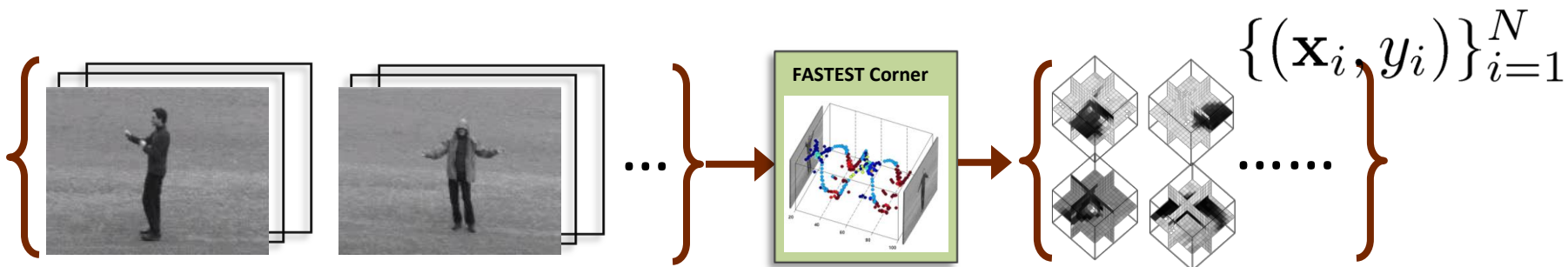$$f(\mathbf{p}) = p_{x_1,y_1,t_1}$$

$$f(\mathbf{p}) = p_{x_1,y_1,t_1} - p_{x_2,y_2,t_2}$$

$$f(\mathbf{p}) = p_{x_1,y_1,t_1} + p_{x_2,y_2,t_2}$$

tree split function

If f(**p**) > threshold

goto Left

**Else**

goto Right

# Randomized Forest Codebook



$\{(\mathbf{x}_i, y_i)\}_{i=1}^N$

FASTEST Corner

tree $t_1$

......

tree $t_T$

1
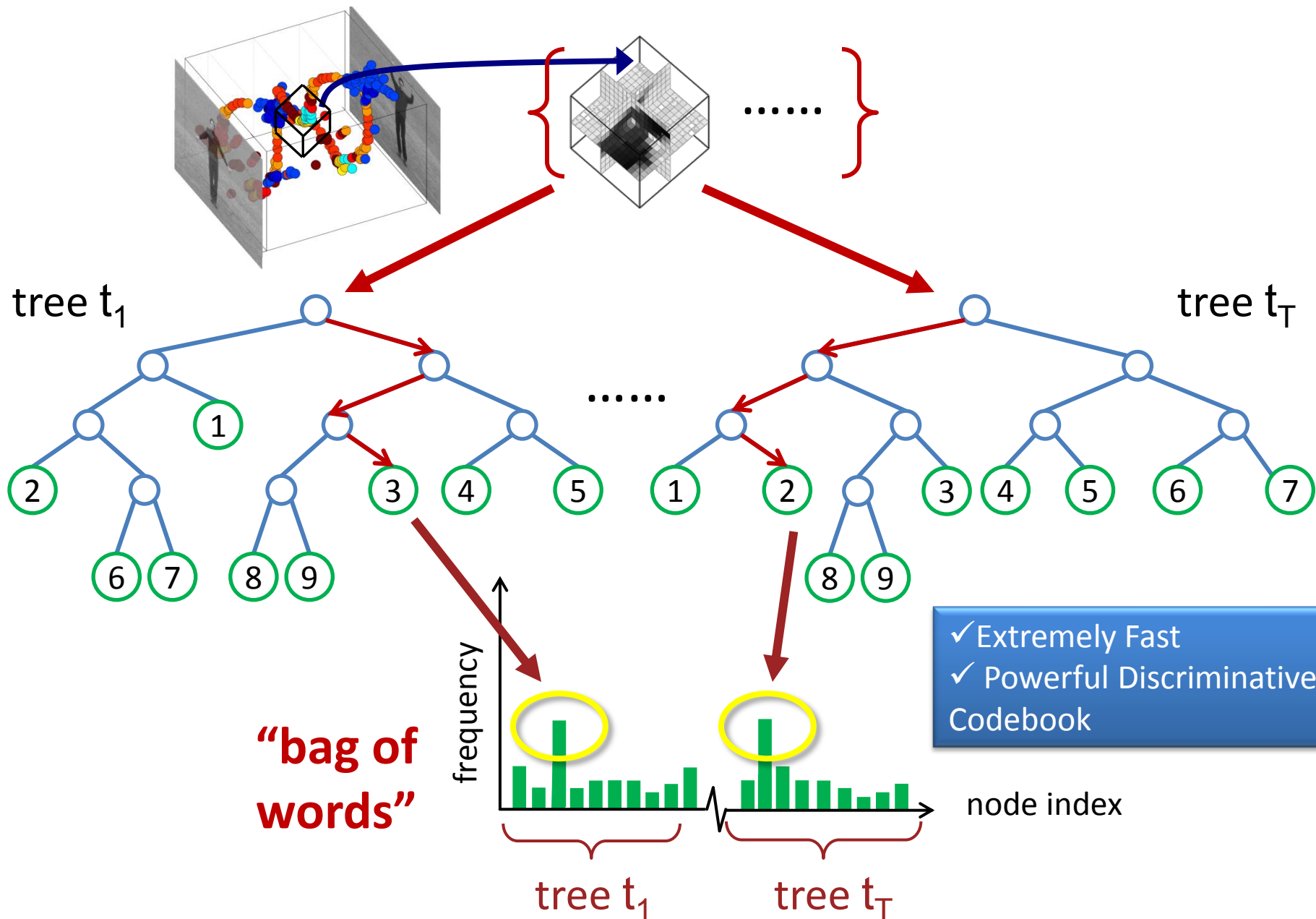2
3
4
5
6
7
8
9

1
2
3
4
5
6
7
8
9

Example Cuboids

Example Cuboids

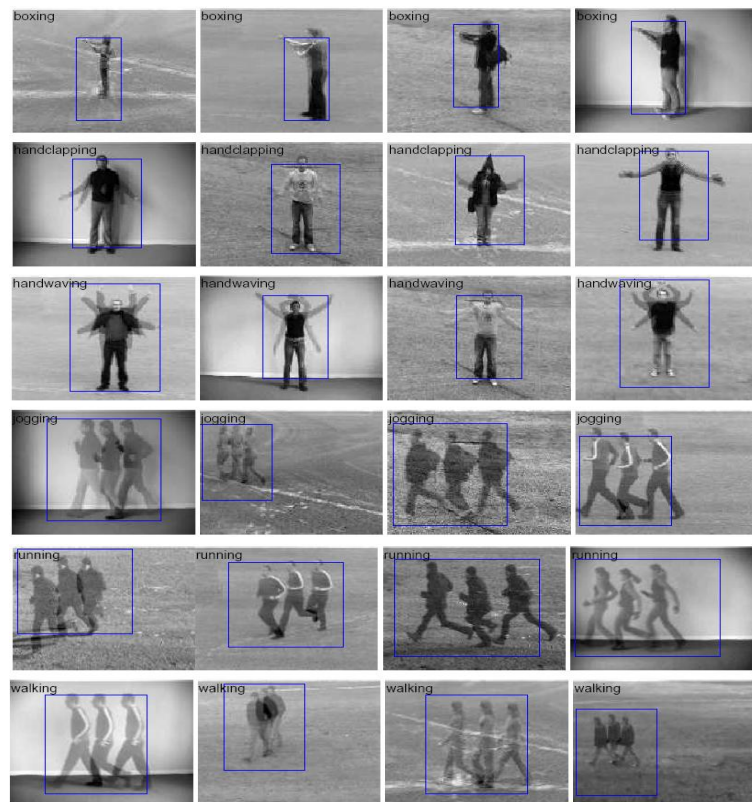# Histogram of Randomized Forest Codewords

# Matlab Demo:

## *Random Forest Codebook*

# Action Categorisation Result

- **Action categorization accuracy (%) on KTH dataset**
  - 6 types of actions, 25 subjects of 4 scenarios by leave-one-out protocol

# Matlab Demo:

## *Action Categorisation*

# Take-home Message

- Use of visual codebook greatly facilitates image and video categorisation tasks.
- The RF codebook is extremely fast as it only uses a small number of features in trees.
  - c.f. K-means codebook is in a flat structure. It compares an input with every codewords, which is time-demanding.
- The RF codebook is also a powerful discriminative codebook. It combines multiple decision trees showing good generalisation.
  - c.f. K-means is an unsupervised learning method.

# Questions?

# Action Recognition Result

- **Action categorization accuracy (%) on KTH dataset**
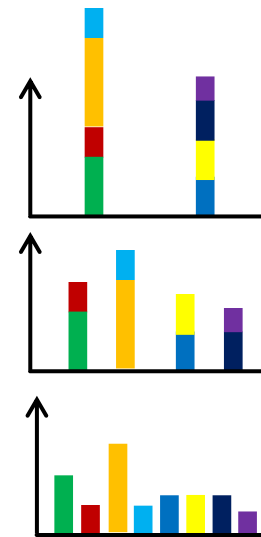  - 6 types of actions, 25 subjects of 4 scenarios by leave-one-out protocol

quantisation effect!

# Pyramid Match Kernel

- PMK acts on semantic texton histogram
  - matches P and Q in *learned* hierarchical histogram space
  - deeper node matches are more important

increased similarity at depth d

$$K(P,Q) = \sum_{d=1}^{D} \frac{1}{2^{D-d+1}} \left( \mathcal{I}_d - \mathcal{I}_{d+1} \right)$$

depth weight



leaf nodes
split nodes

d=1

d=2

d=3

# Categorisation

- Histogram matching of a pair of videos *P* and *Q* is

$$\mathcal{I}_d(P,Q) = \sum_{n=1}^{b_d} \min\left(H_d(P^{(n)}), H_d(Q^{(n)})\right)$$

- Based on HMK, we can use various classifiers such as kNN, SVM, Naïve-Bayes classifiers or pLSA, LDA.

Class 1    Class M

*P*

*Q*

**NN classification by**

$$\arg\max_c \mathcal{I}_d\left(P_c, Q\right)$$

**kNN Classifier**

**category decision**

# Action Recognition Result

- **Action categorization accuracy (%) on KTH dataset**
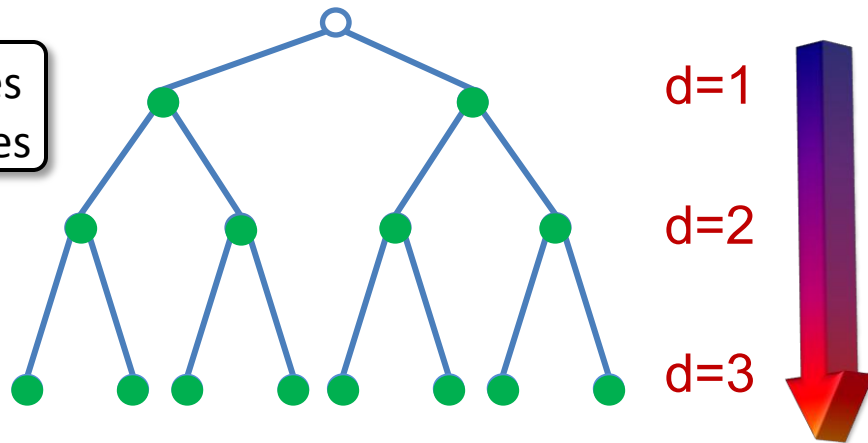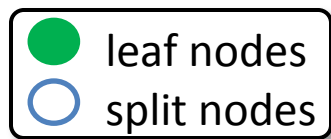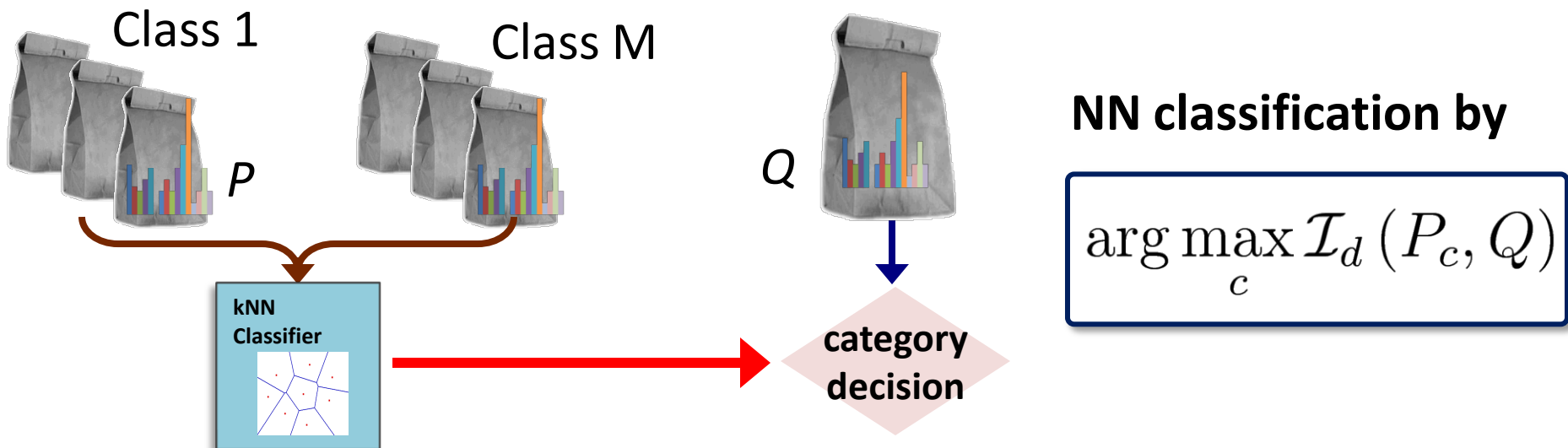  - 6 types of actions, 25 subjects of 4 scenarios by leave-one-out protocol



|        | box | hclp | hwav | jog | run | walk |
|--------|-----|------|------|-----|-----|------|
| box    | .98 | .02  | .00  | .00 | .00 | .00  |
| hclp   | .00 | 1.0  | .00  | .00 | .00 | .00  |
| hwav   | .01 | .02  | .97  | .00 | .00 | .00  |
| jog    | .00 | .00  | .00  | .90 | .10 | .00  |
| run    | .00 | .00  | .00  | .12 | .88 | .00  |
| walk   | .00 | .00  | .00  | .01 | .00 | .99  |