



UNIVERSITÄT **BONN**

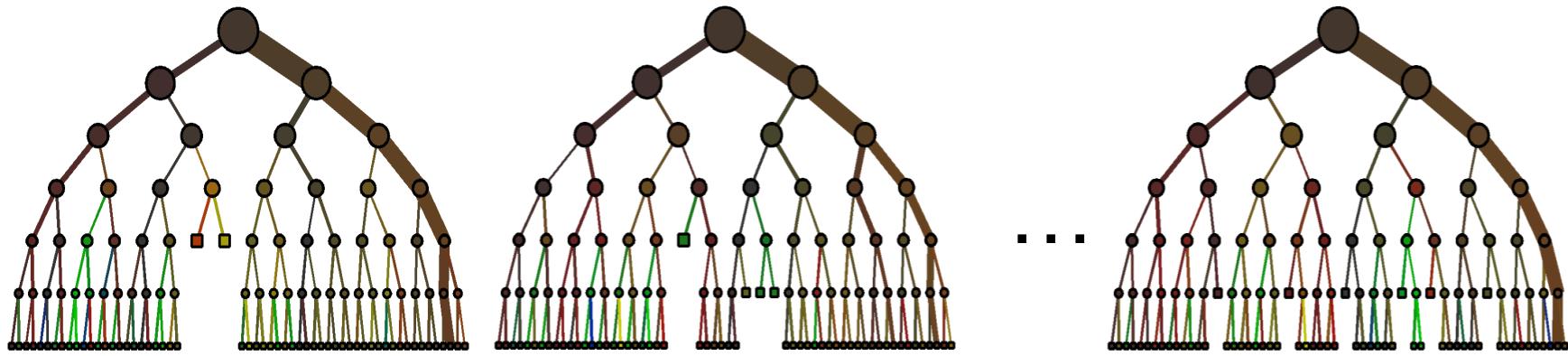
Prof. Dr. Juergen Gall

Random Forests
MA-INF 2213 - Advanced Computer Vision
SS25

Organization

- Exercise:
 - First exercise on **2.5.**
 - Tentative: **5 Sheets**
 - Tutorial: PyTorch and GPU cluster Bender on **16.5.** (tentative)
<https://www.hpc.uni-bonn.de/en/systems/bender>

Random Forests



[L. Breiman, J. Friedman, C.J. Stone, and R.A. Olshen. **Classification and Regression Trees (CART)**. 1984]

[T. Ho. **Random Decision Forests**. 1995]

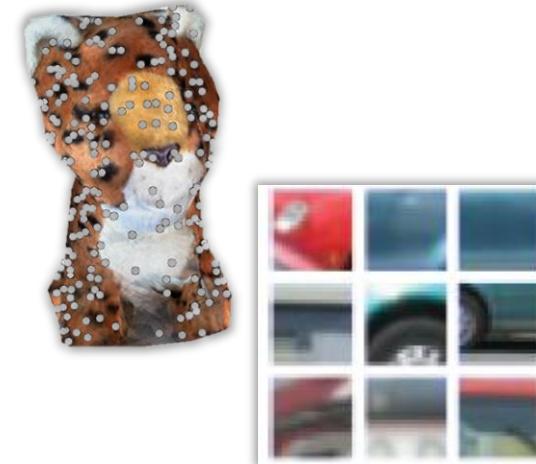
[Y. Amit and D. Geman. **Shape quantization and recognition with randomized trees**. 1997]

[L. Breiman. **Random forests**. 2001]

[A. Criminisi and J. Shotton. **Decision Forests in Computer Vision and Medical Image Analysis**. 2013]

Computer vision applications

[V. Lepetit and P. Fua. **Keypoint recognition using randomized trees.** 2005]



[F. Moosman, B. Triggs and F. Jurie. **Fast discriminative visual codebooks using randomized clustering forests.** 2006]

[J. Shotton, M. Johnson and R. Cipolla. **Semantic Texton Forests for Image Categorization and Segmentation.** 2008]



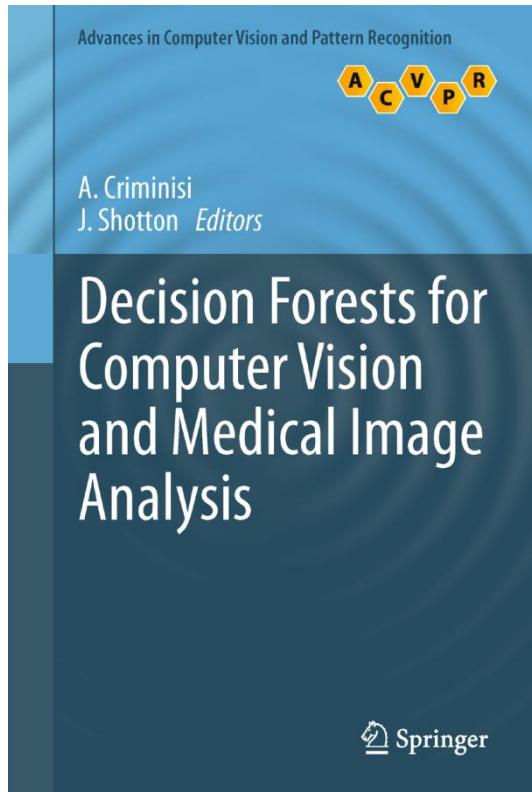
[J. Gall and V. Lempitsky. **Class-Specific Hough Forests for Object Detection.** 2009]

[J. Shotton et al. **Real-Time Human Pose Recognition in Parts from a Single Depth Image.** 2011]



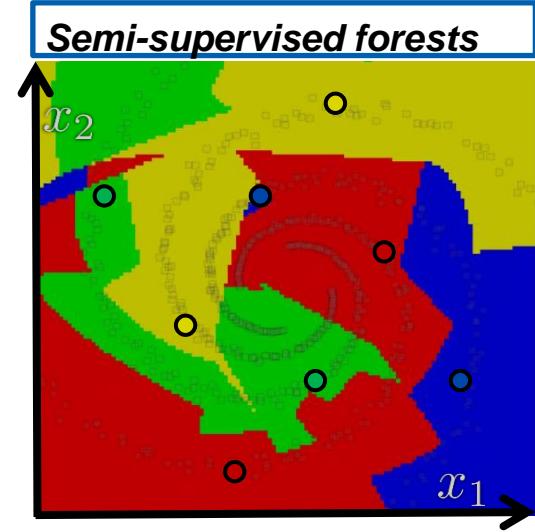
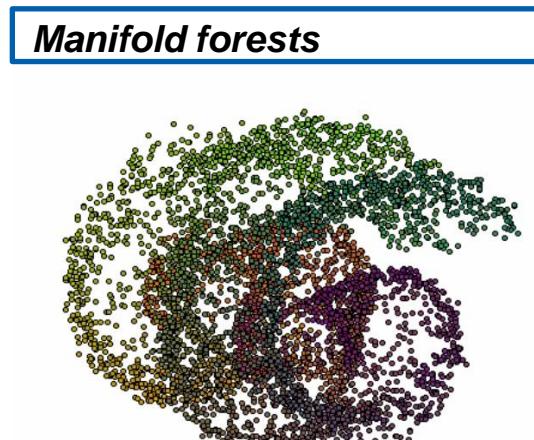
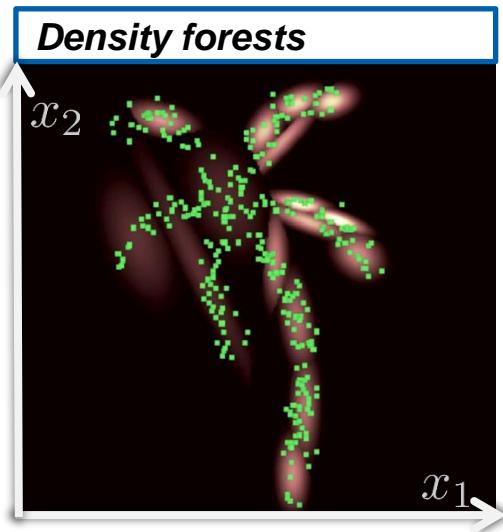
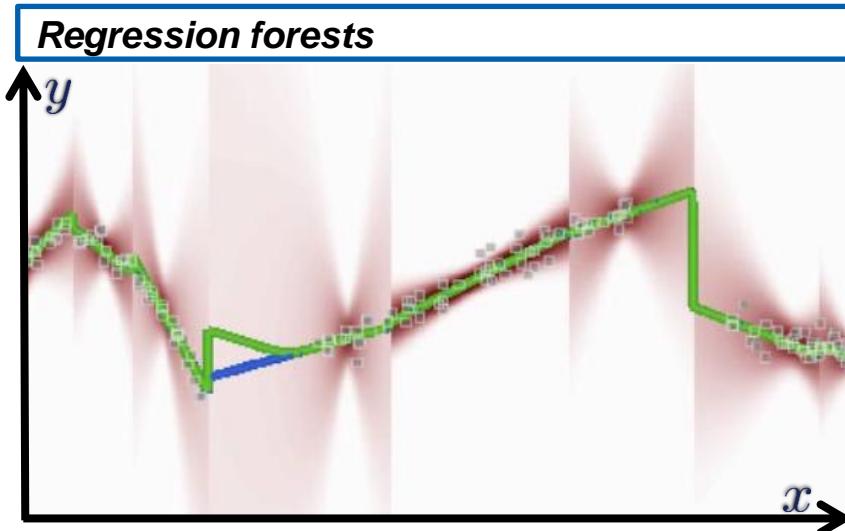
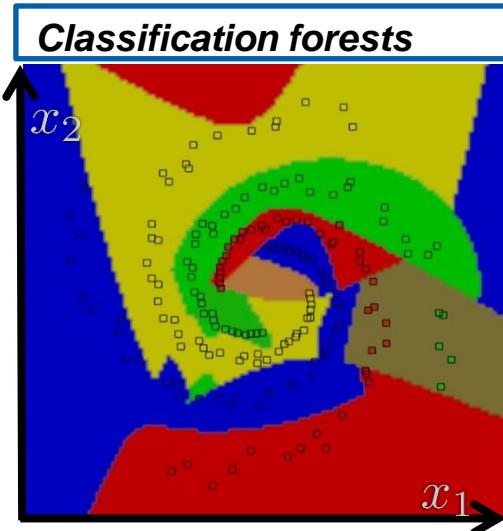
Literature

<http://research.microsoft.com/en-us/projects/decisionforests/>



A. Criminisi el al. **Decision
Forests for Computer Vision
and Medical Image Analysis**,
Springer 2013

What can decision forests do?

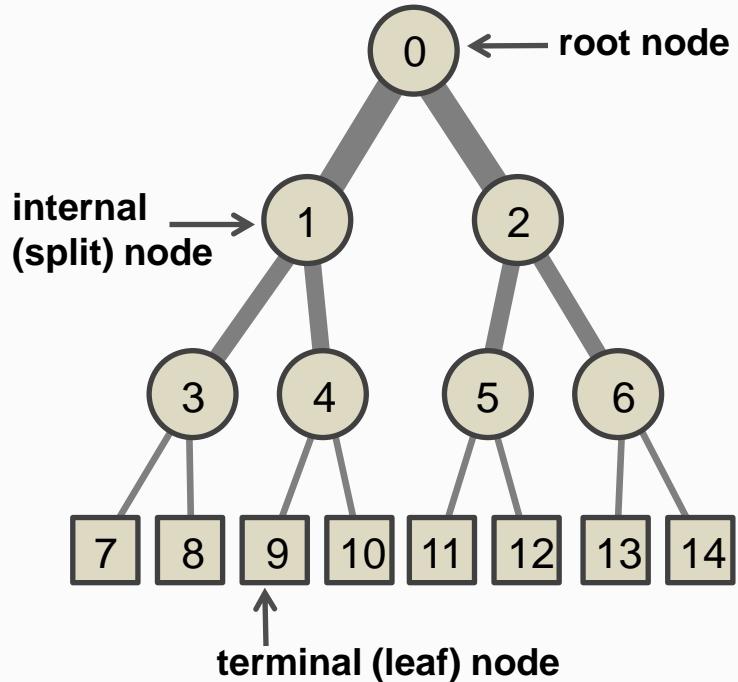


Overview

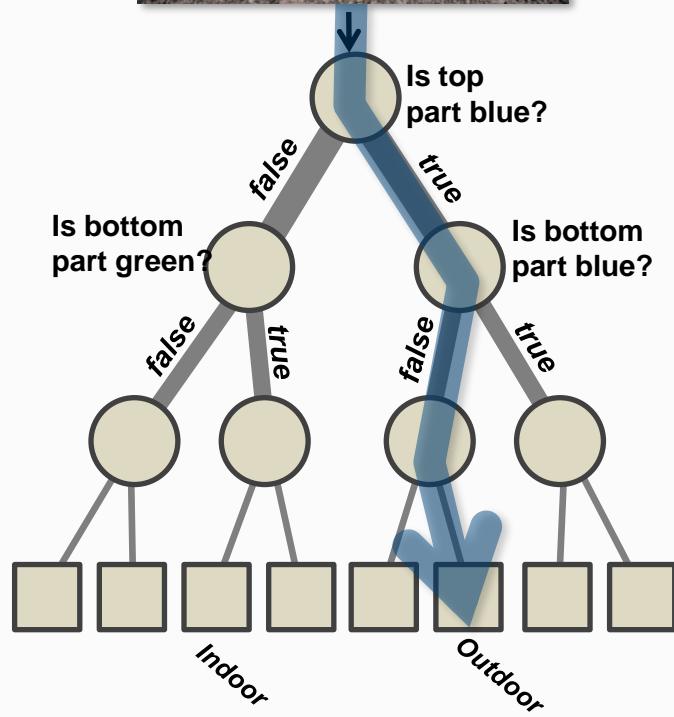
- Classification forests
- Regression forests
- Other variants

Decision trees

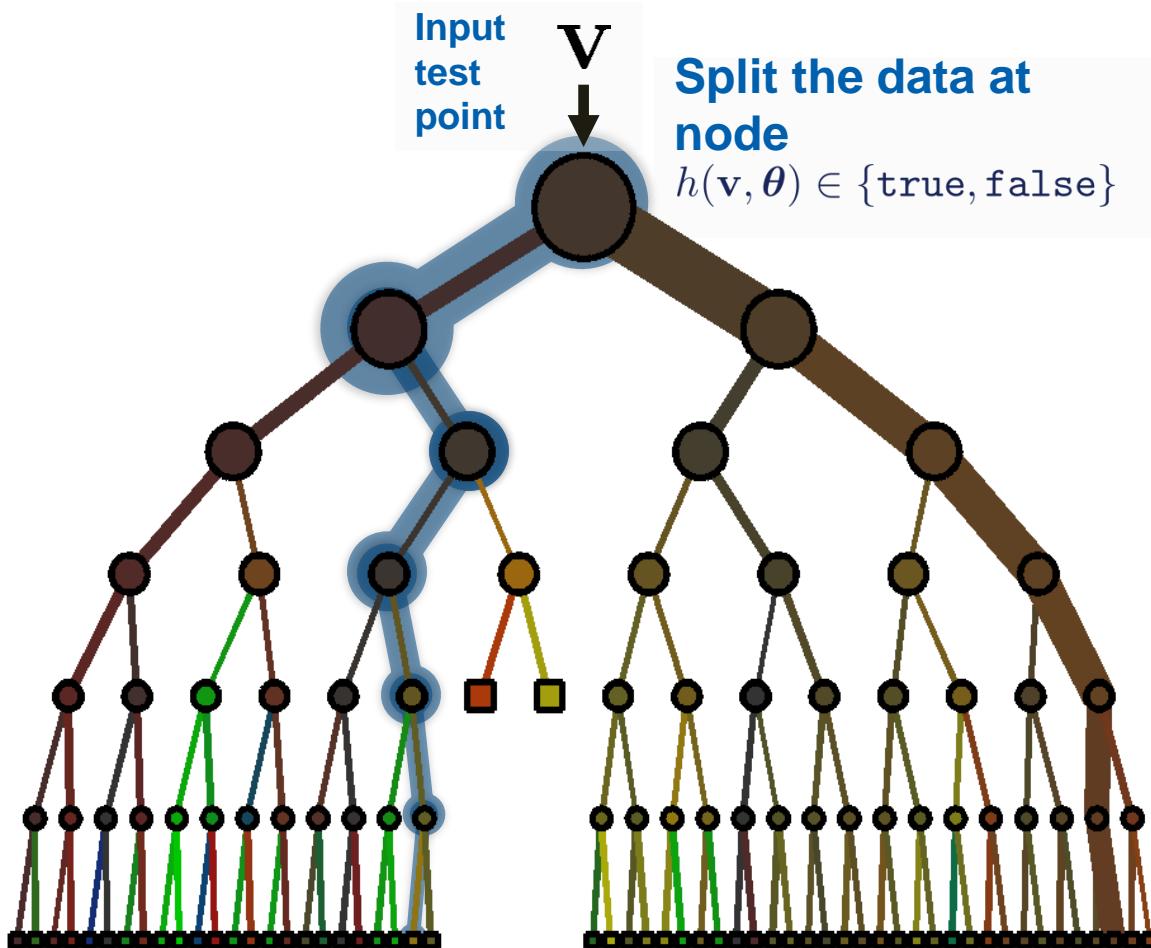
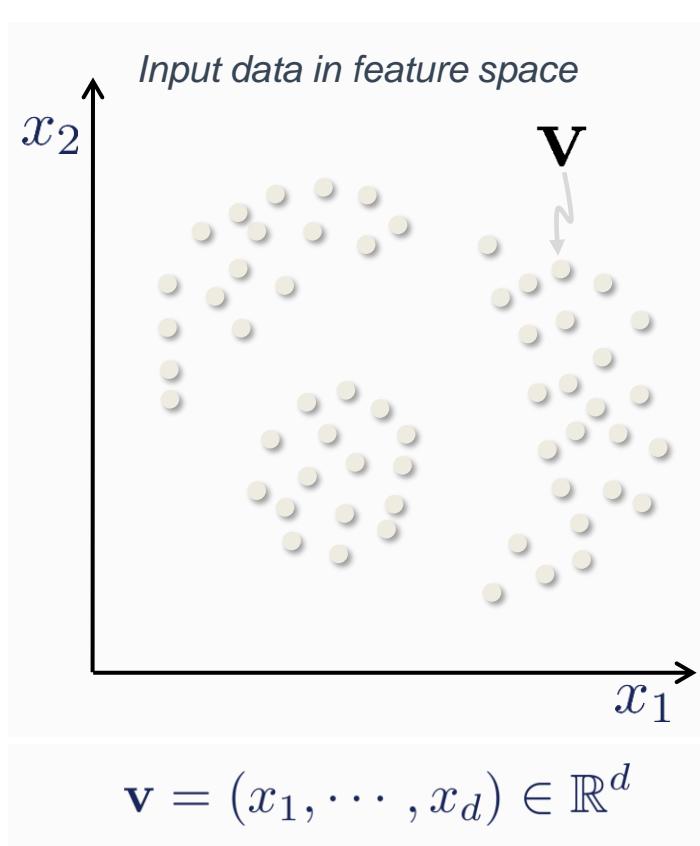
A general tree structure



A decision tree

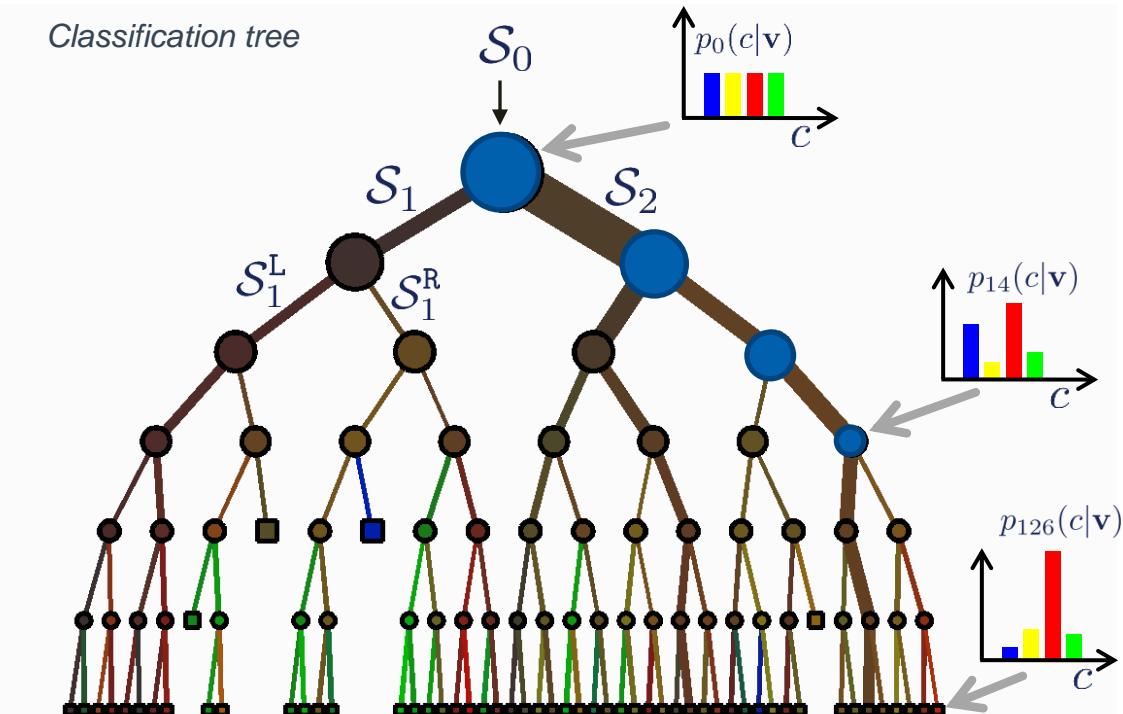
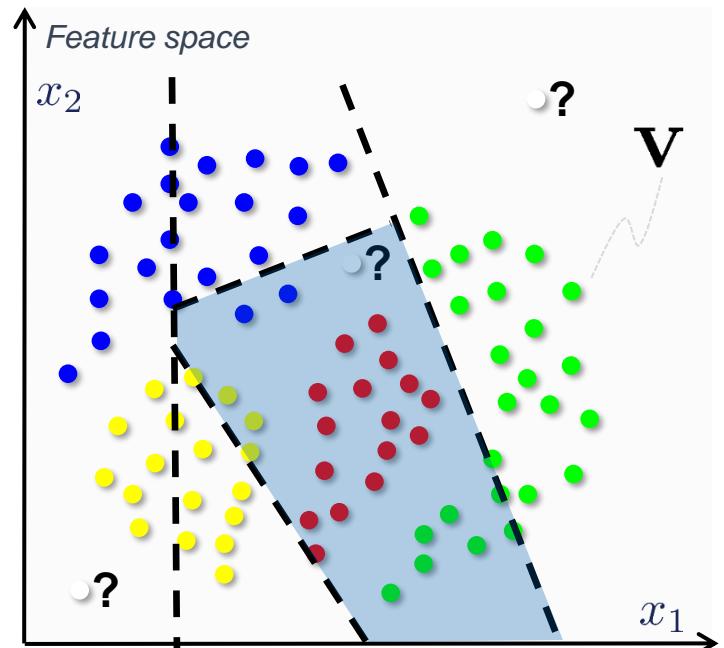


Decision tree



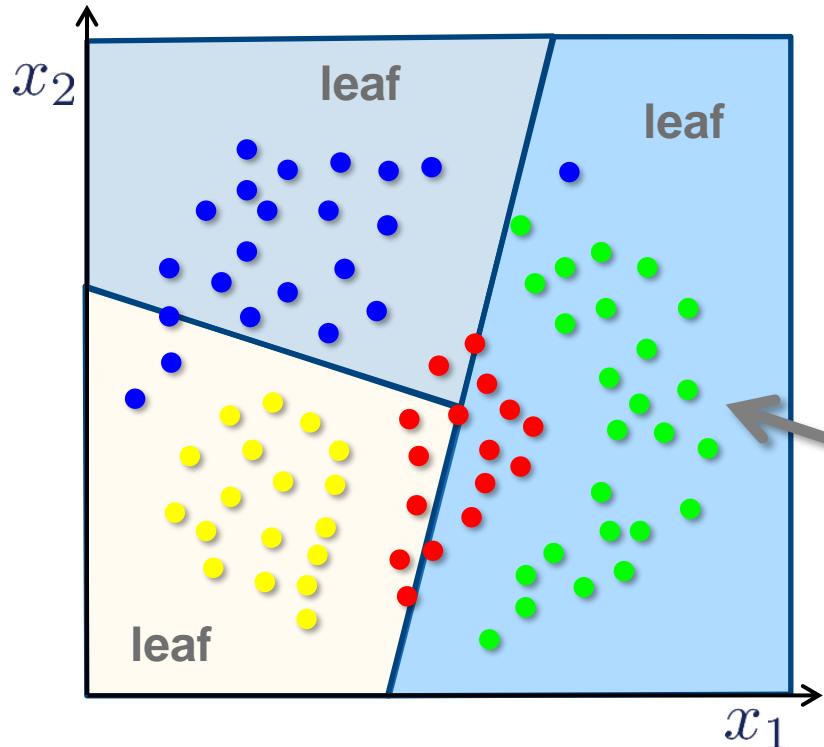
Prediction at leaf
 $p(c|\mathbf{v})$

Classification tree

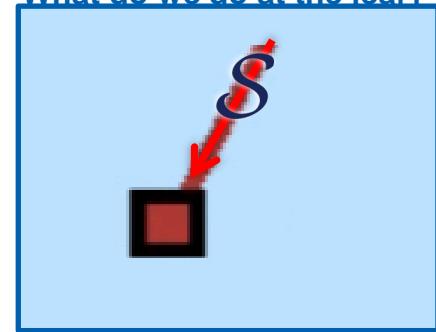


Figures adapted from A. Criminisi

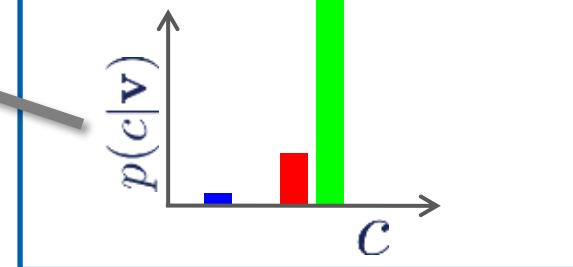
Classification forest: the prediction model



What do we do at the leaf?

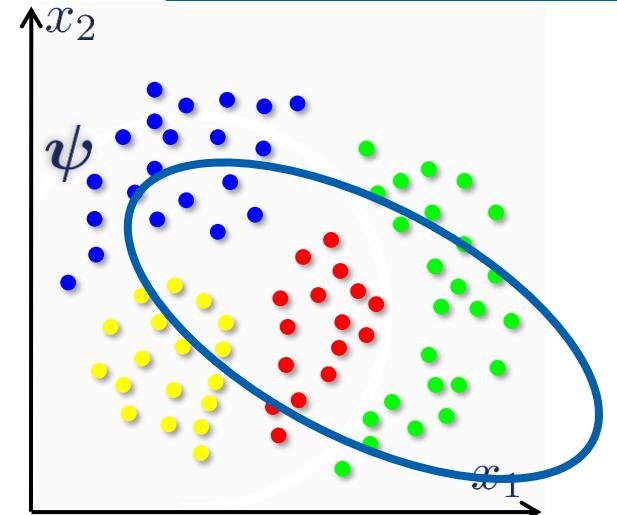
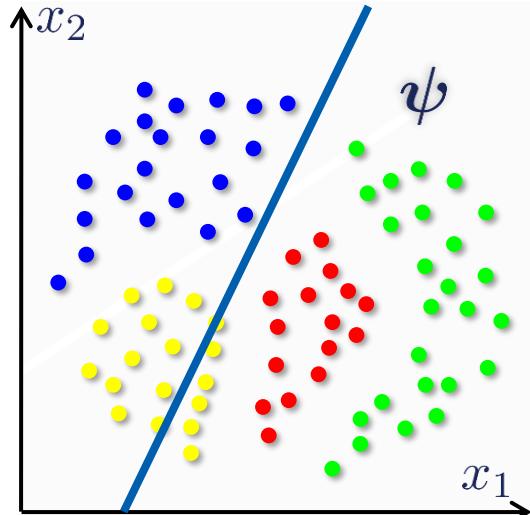
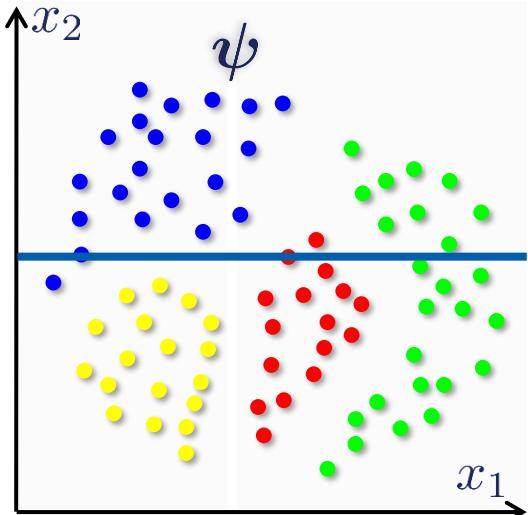


Prediction model: probabilistic



Classification forest: the weak learner model

Examples of weak learners



Weak learner: axis aligned

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

Feature response
for 2D example.
 $\phi(\mathbf{v}) = (x_1 \ x_2 \ 1)^\top$
With $\psi = (1 \ 0 \ \psi_3)$ or $\psi = (0 \ 1 \ \psi_3)$

In general ϕ may select only a very small subset of features

Weak learner: oriented line

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

Feature response
for 2D example.
 $\phi(\mathbf{v}) = (x_1 \ x_2 \ 1)^\top$
With $\psi \in \mathbb{R}^3$ a generic line in homog. coordinates.

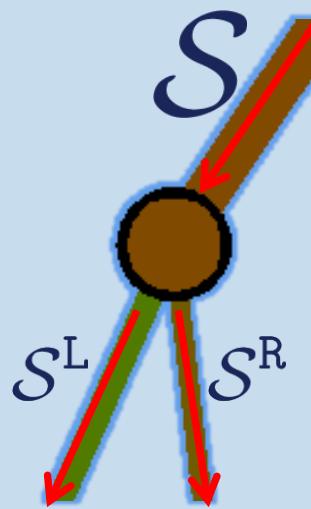
$\phi(\mathbf{v}) : \mathbb{R}^d \rightarrow \mathbb{R}^{d'+1}, \ d' \ll d$

Weak learner: conic section

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

Feature response
for 2D example.
 $\phi(\mathbf{v}) = (x_1 \ x_2 \ 1)^\top$
With $\psi \in \mathbb{R}^{3 \times 3}$ a matrix representing a conic.

Decision forest model: training and information gain



Information gain

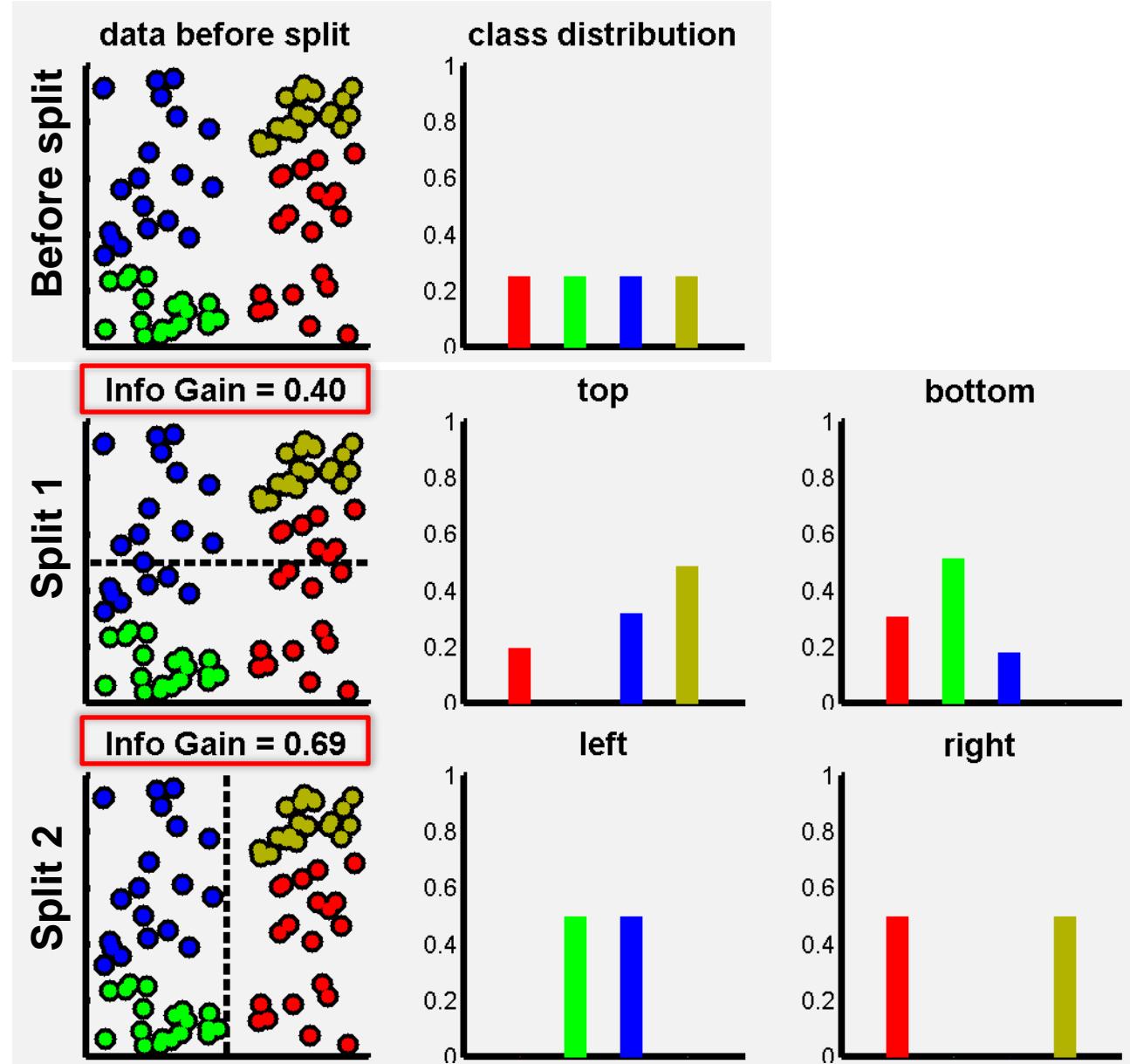
$$I(\mathcal{S}, \theta) = H(\mathcal{S}) - \sum_{i \in \{\text{L,R}\}} \frac{|\mathcal{S}^i|}{|\mathcal{S}|} H(\mathcal{S}^i)$$

Shannon's entropy

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

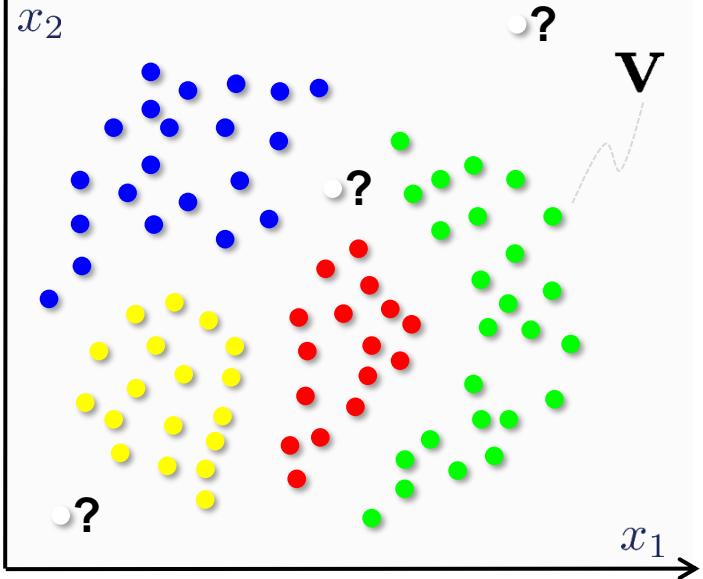
Node training

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

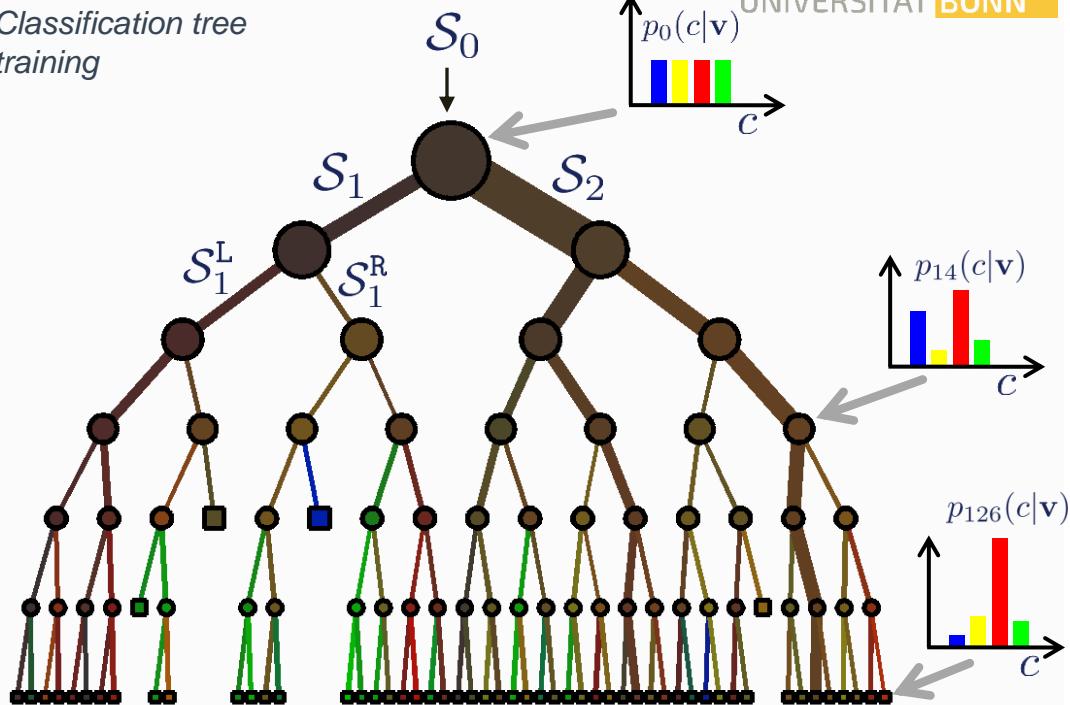


Classification forest

Training data in feature space



Classification tree
training



Input:

Output:

Information gain:

Entropy of discrete distribution:

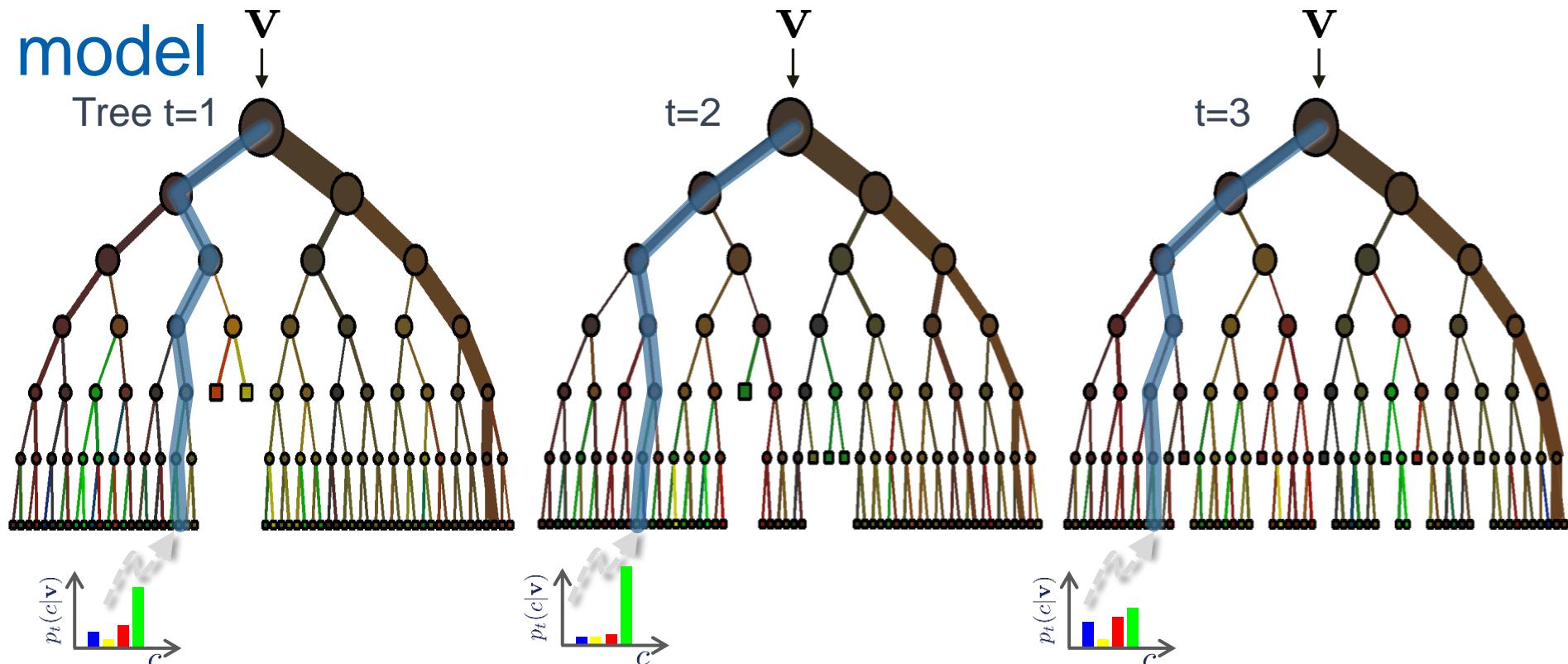
$$\mathbf{v} = (x_1, \dots, x_d) \in \mathbb{R}^d$$

$$p(c|\mathbf{v}) \quad c \in \mathcal{C} \quad \text{with} \quad \mathcal{C} = \{c_k\}$$

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

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

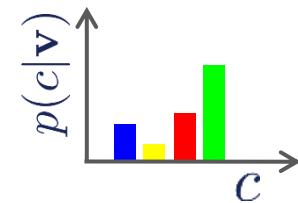
Classification forest: the ensemble model



The ensemble model

Forest output probability

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



Decision forest model: the randomness model

1) Bagging (randomizing the training set)

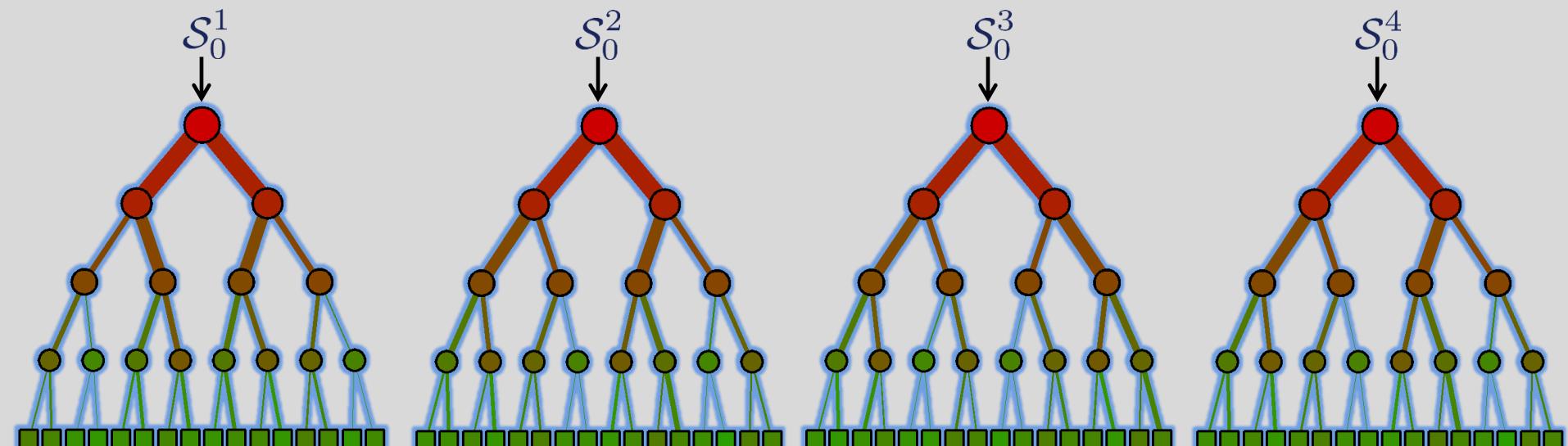
 \mathcal{S}_0

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 

Efficient training

Decision forest model: the randomness model

2) Randomized node optimization (RNO)

 \mathcal{T}

The full set of all possible node test parameters

 $\mathcal{T}_j \subset \mathcal{T}$

For each node the set of randomly sampled features

 $\rho = |\mathcal{T}_j|$

Randomness control parameter.

For $\rho = |\mathcal{T}|$ no randomness and maximum tree correlation.

For $\rho = 1$ max randomness and minimum tree correlation.

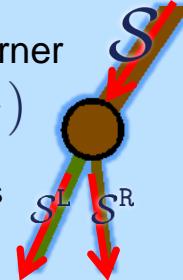
Node training

Node weak learner

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

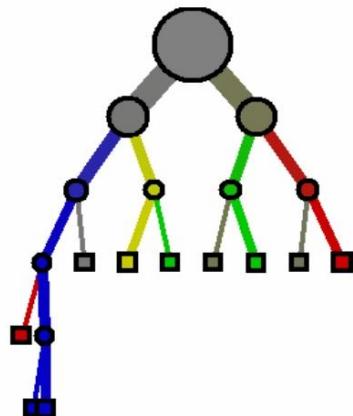
Node test params

$$\boldsymbol{\theta} \in \mathcal{T}_j$$

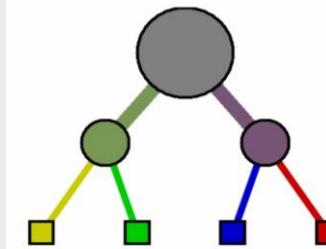


The effect of ρ

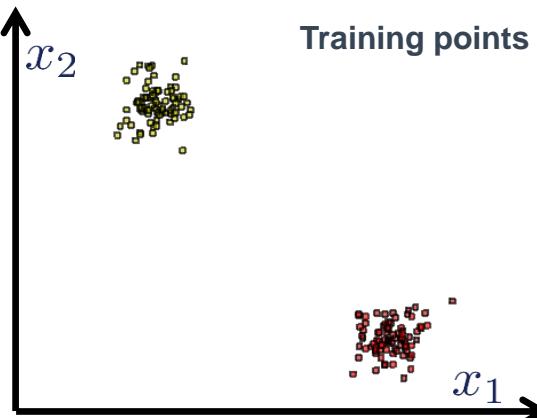
Small value of ρ ; little tree correlation.



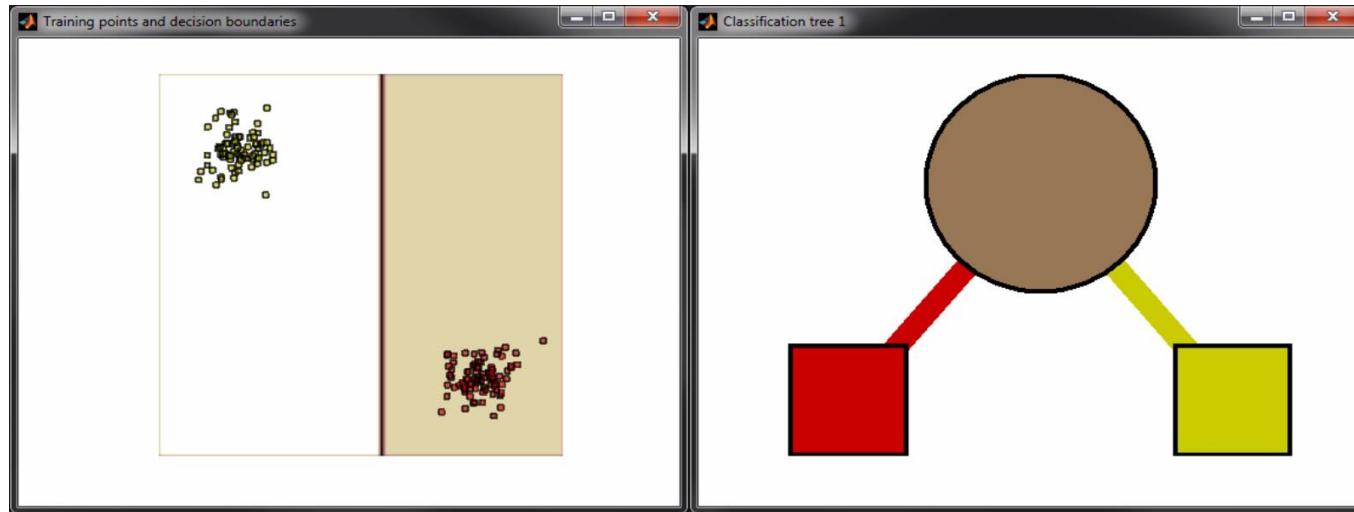
Large value of ρ ; large tree correlation.



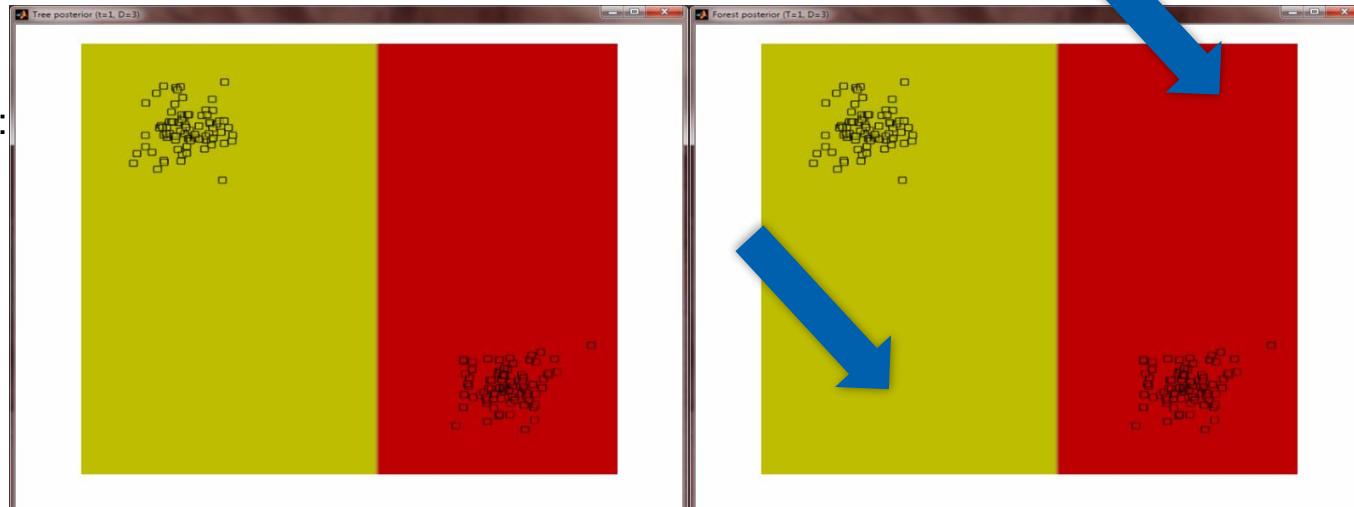
Classification forest: effect of the weak learner model



Training different trees in the forest



Testing different trees in the forest

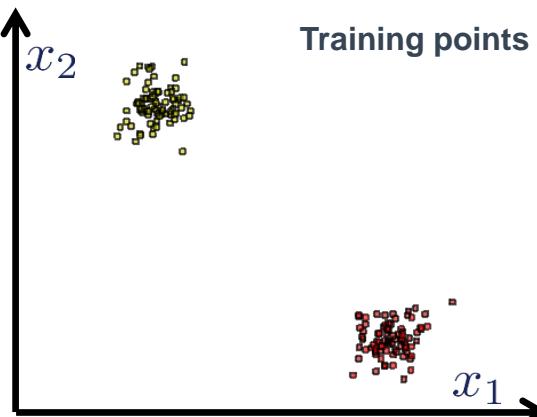


Three concepts to keep in mind:

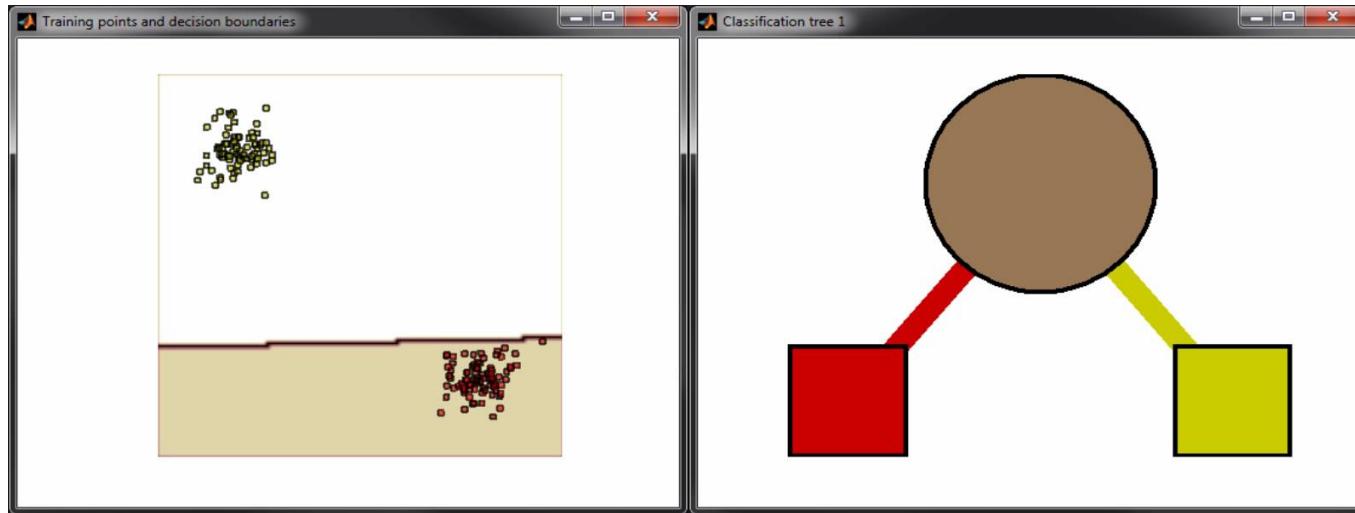
- “Accuracy of prediction”
- “Quality of confidence”
- “Generalization”

Parameters: T=200, D=2, weak learner = aligned, leaf model = probabilistic

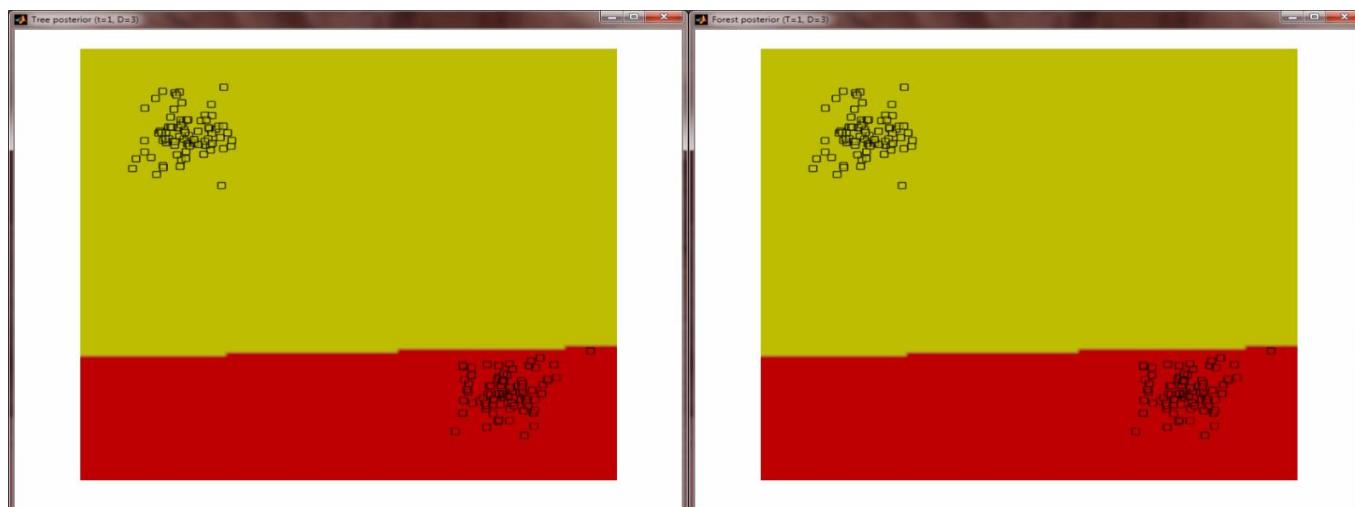
Classification forest: effect of the weak learner model



Training different trees in the forest

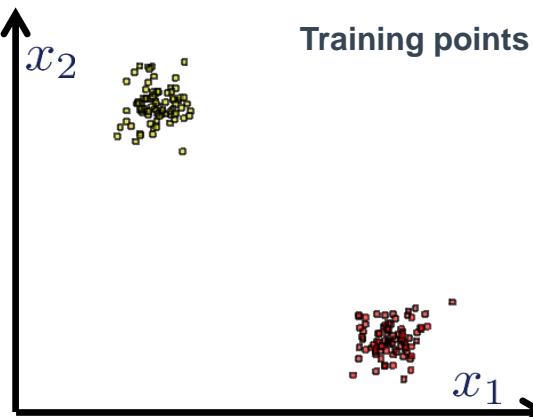


Testing different trees in the forest

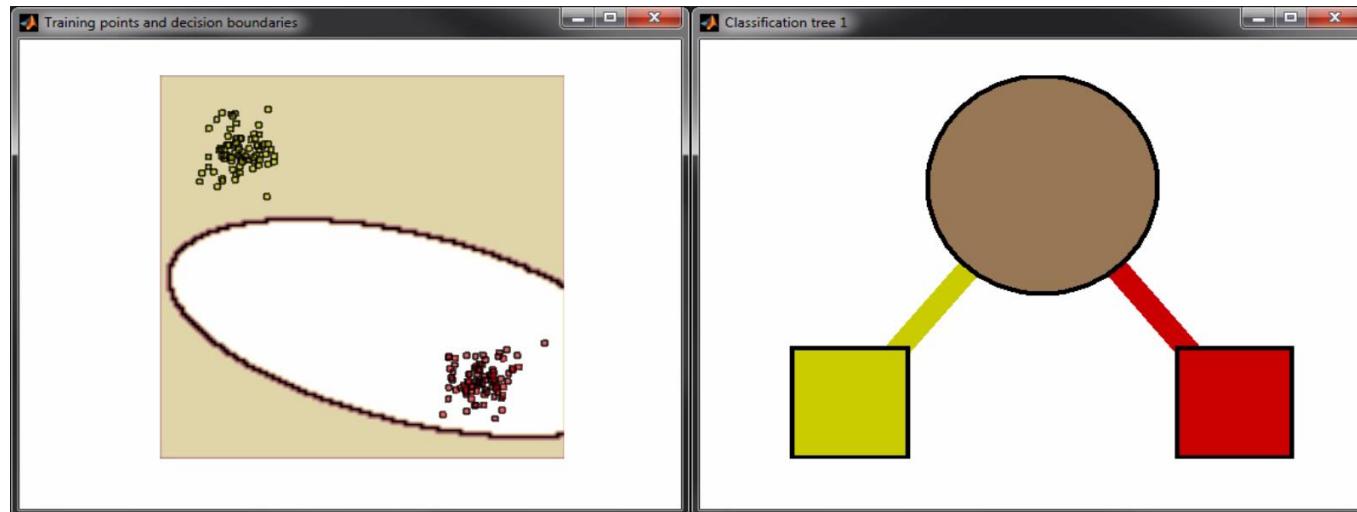


Parameters: T=200, D=2, weak learner = linear, leaf model = probabilistic

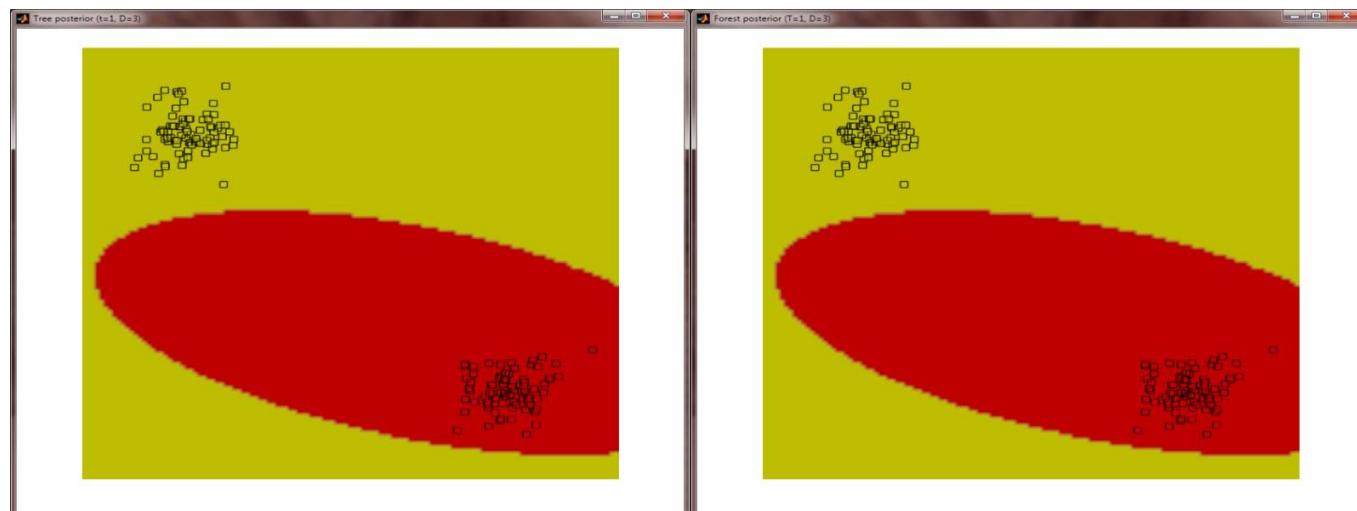
Classification forest: effect of the weak learner model



Training different trees in the forest

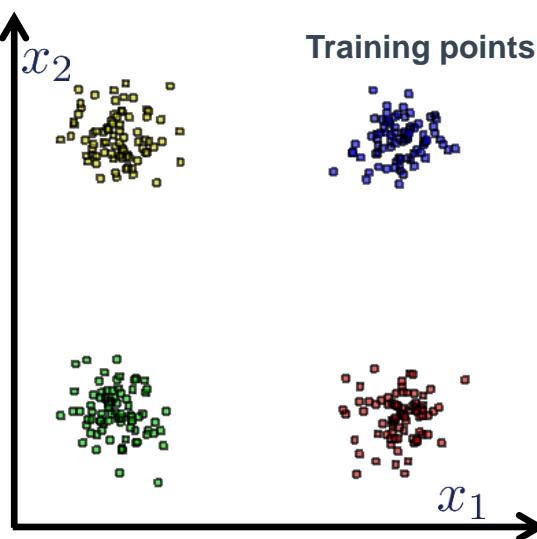


Testing different trees in the forest

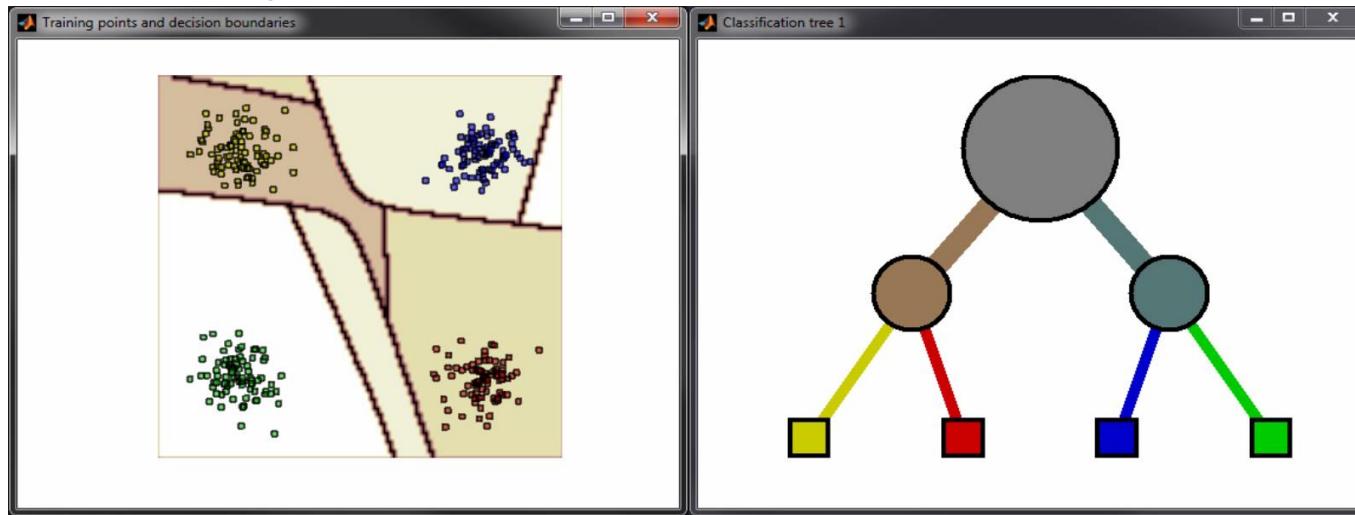


Parameters: T=200, D=2, weak learner = conic, leaf model = probabilistic

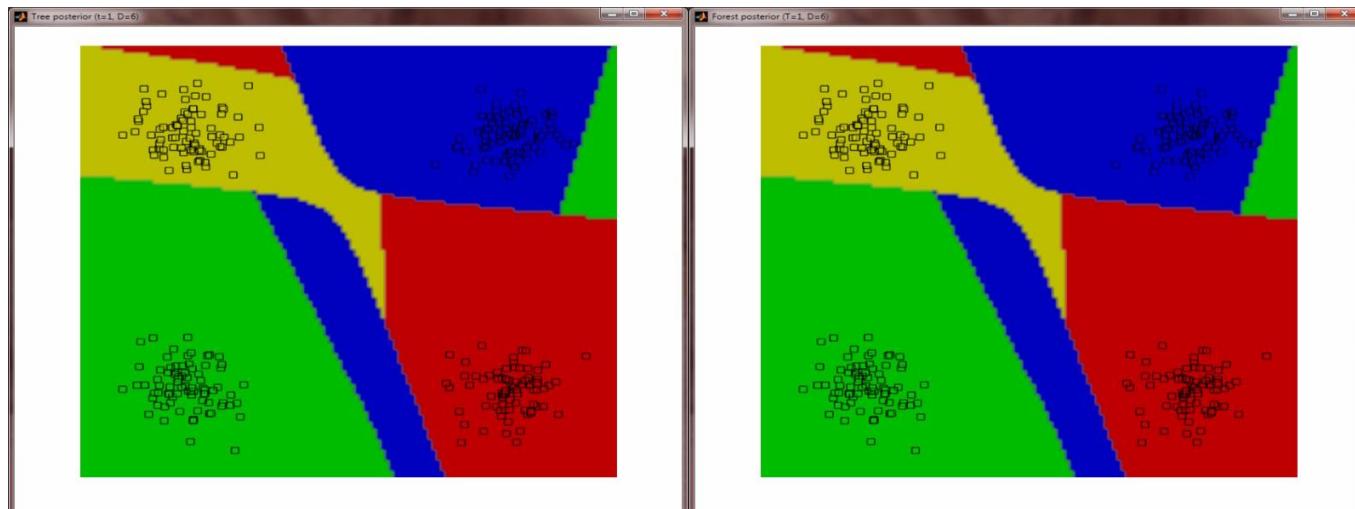
Classification forest: with >2 classes



Training different trees in the forest

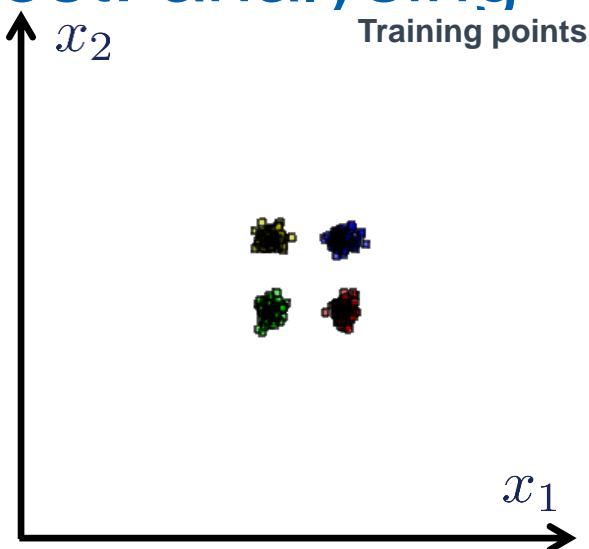


Testing different trees in the forest

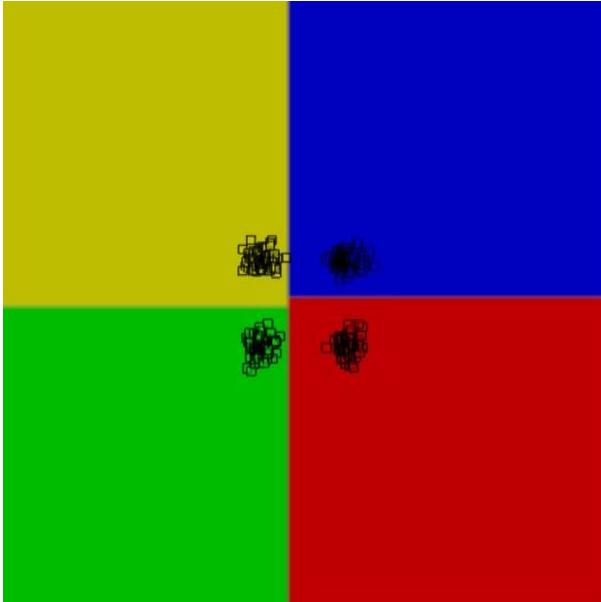


Parameters: T=200, D=3, weak learner = conic, leaf model = probabilistic

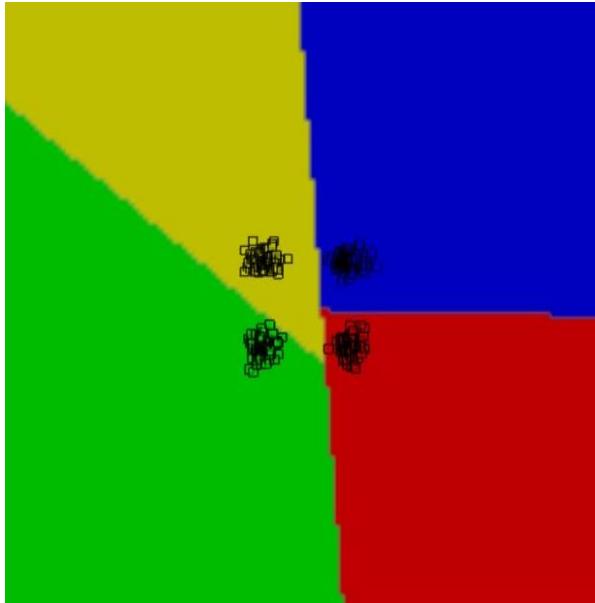
Classification forest: analysing generalization



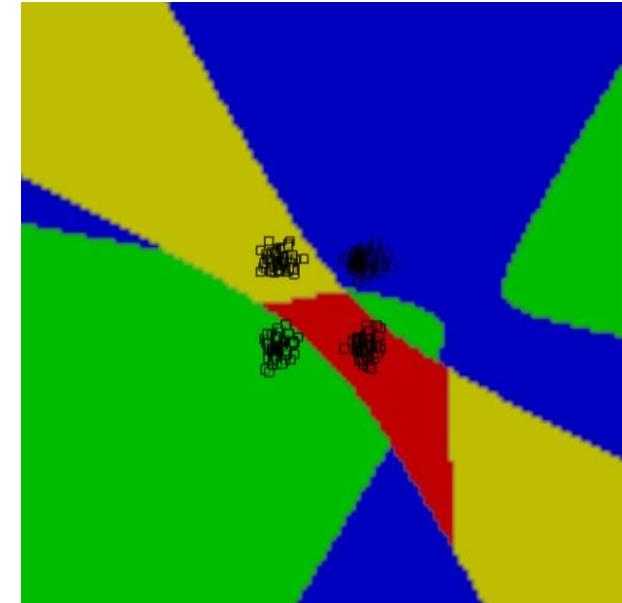
Weak learner: axis aligned



Weak learner: oriented line

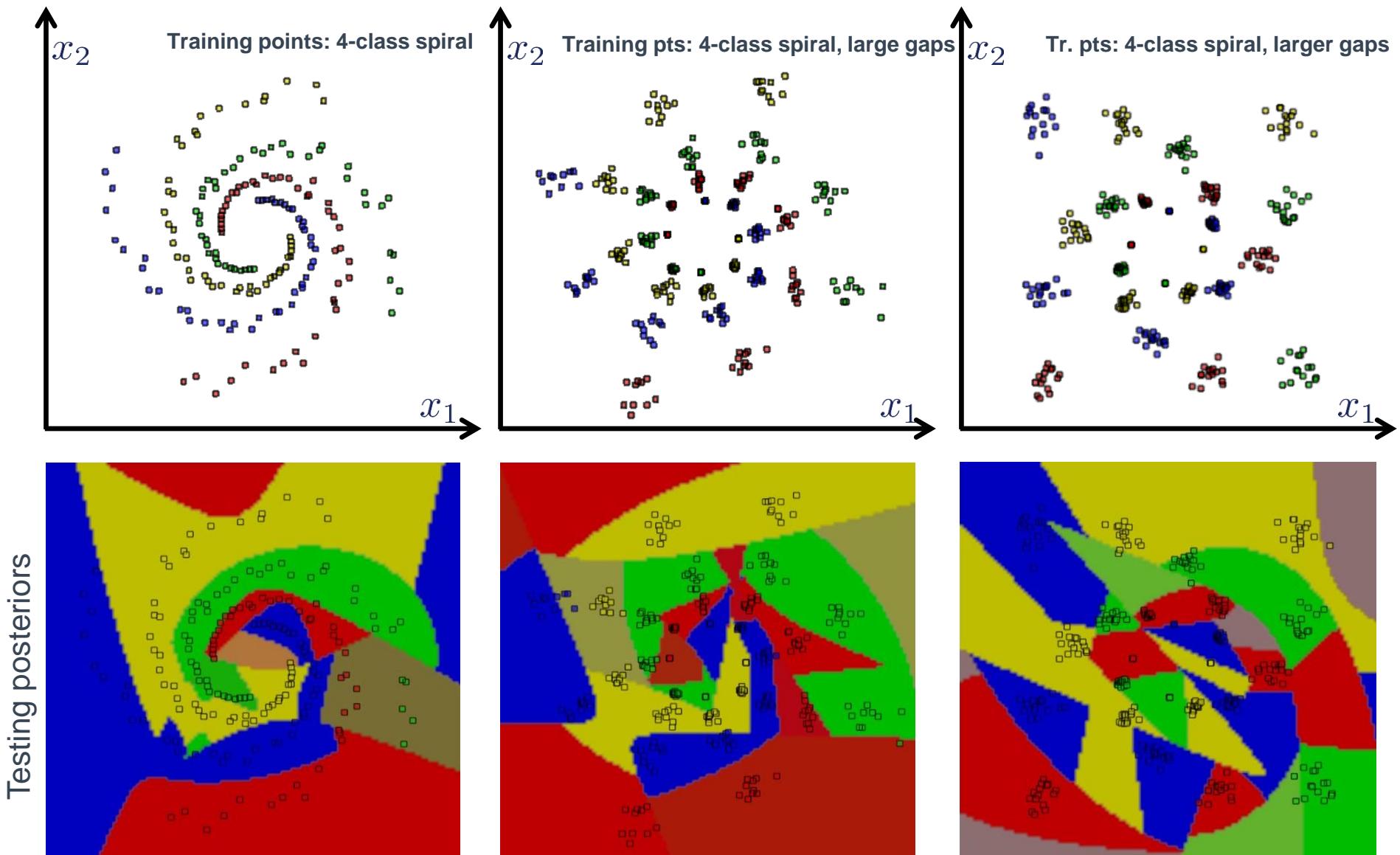


Weak learner: conic section



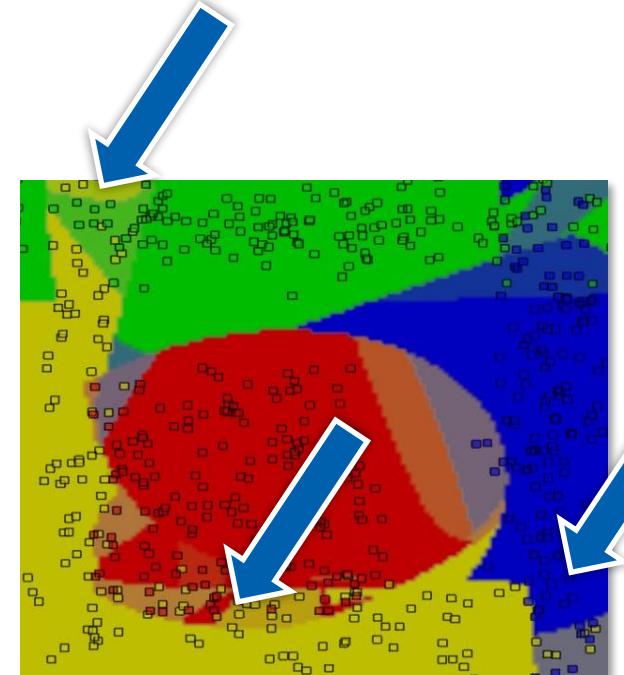
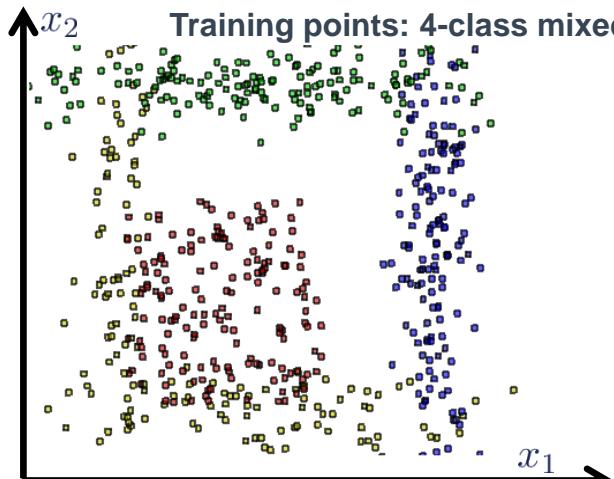
Parameters: T=200, D=3, leaf model = probabilistic

Classification forest: analysing generalization



Parameters: T=200, D=13, w. l. = conic, predictor = prob.

Classification forest: effect of tree depth



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

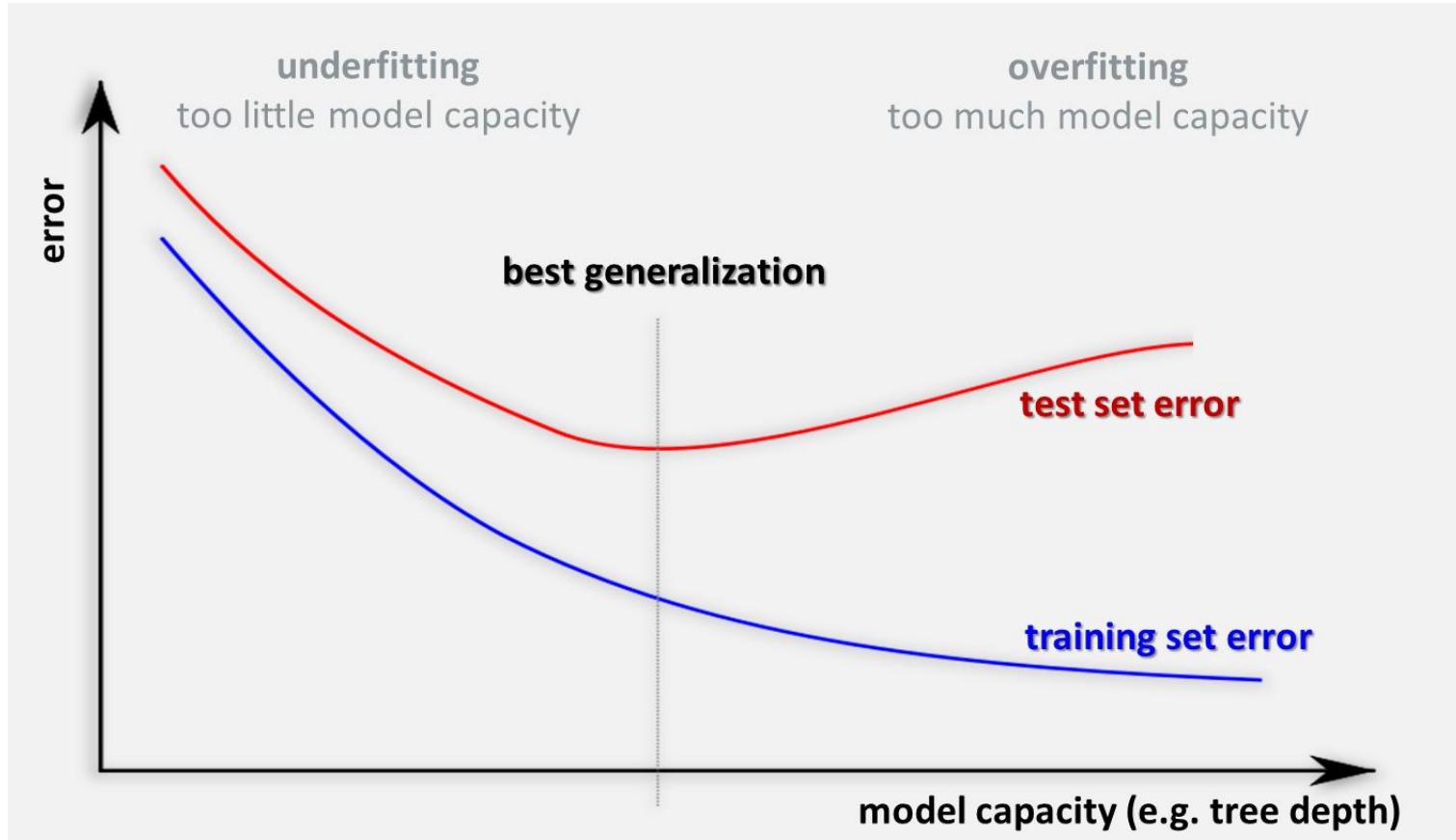
T=200, D=6, w. l. = conic

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

underfitting

max tree depth, D
overfitting
Predictor model = prob.

Background: overfitting and underfitting



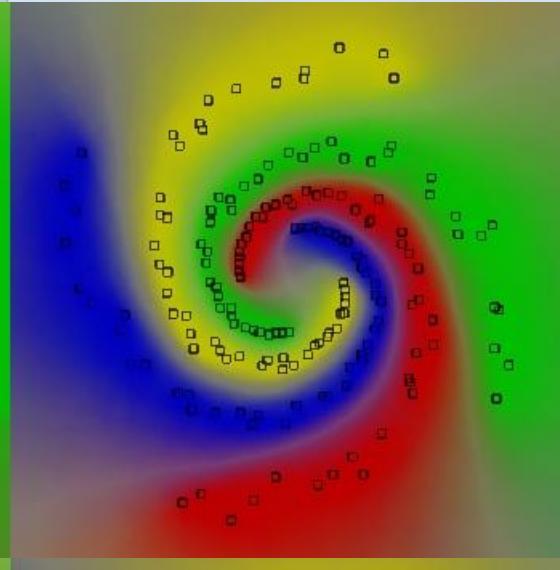
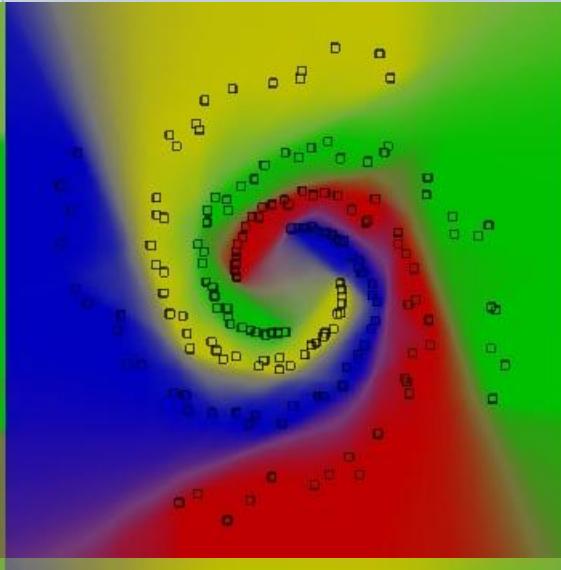
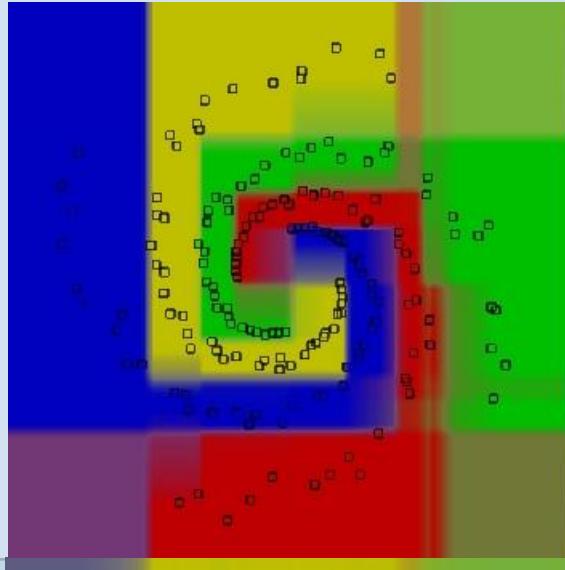
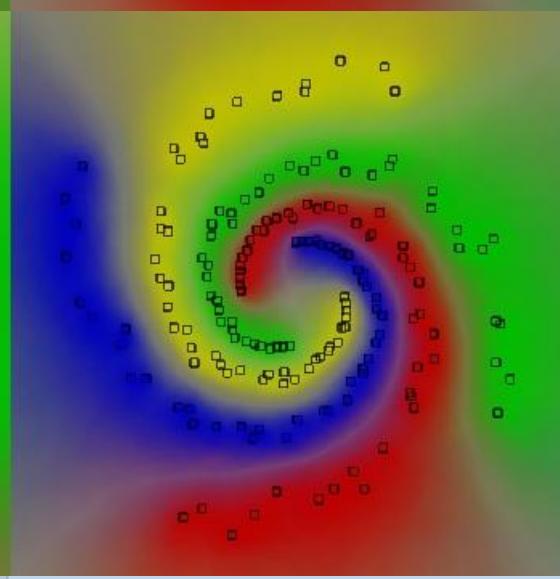
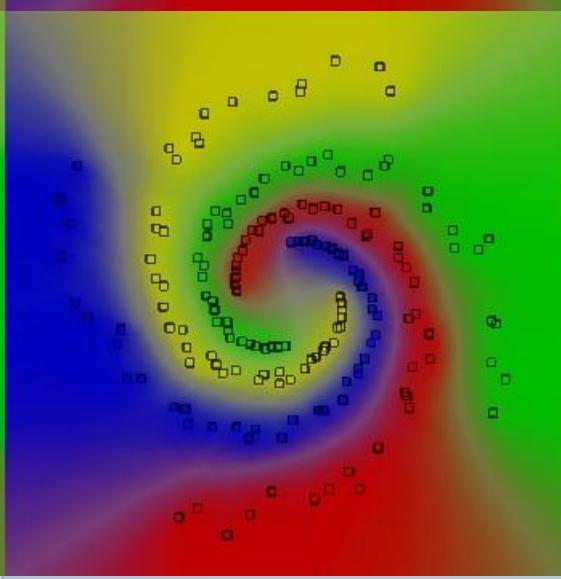
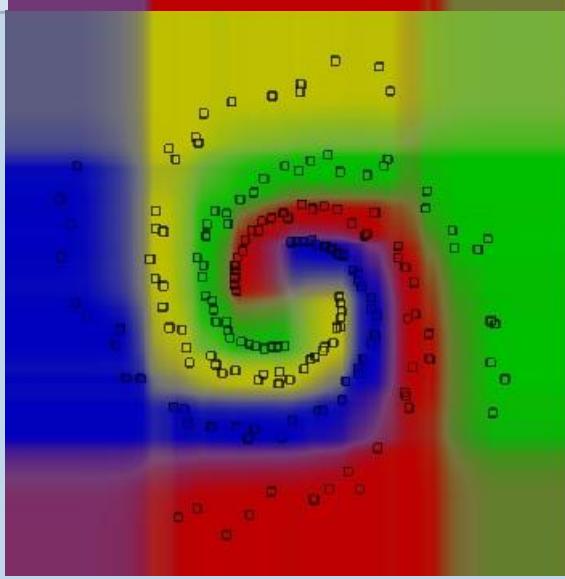
Classification forest: effect of weak learner model and randomness

Testing posteriors

Weak learner: axis aligned

Weak learner: oriented line

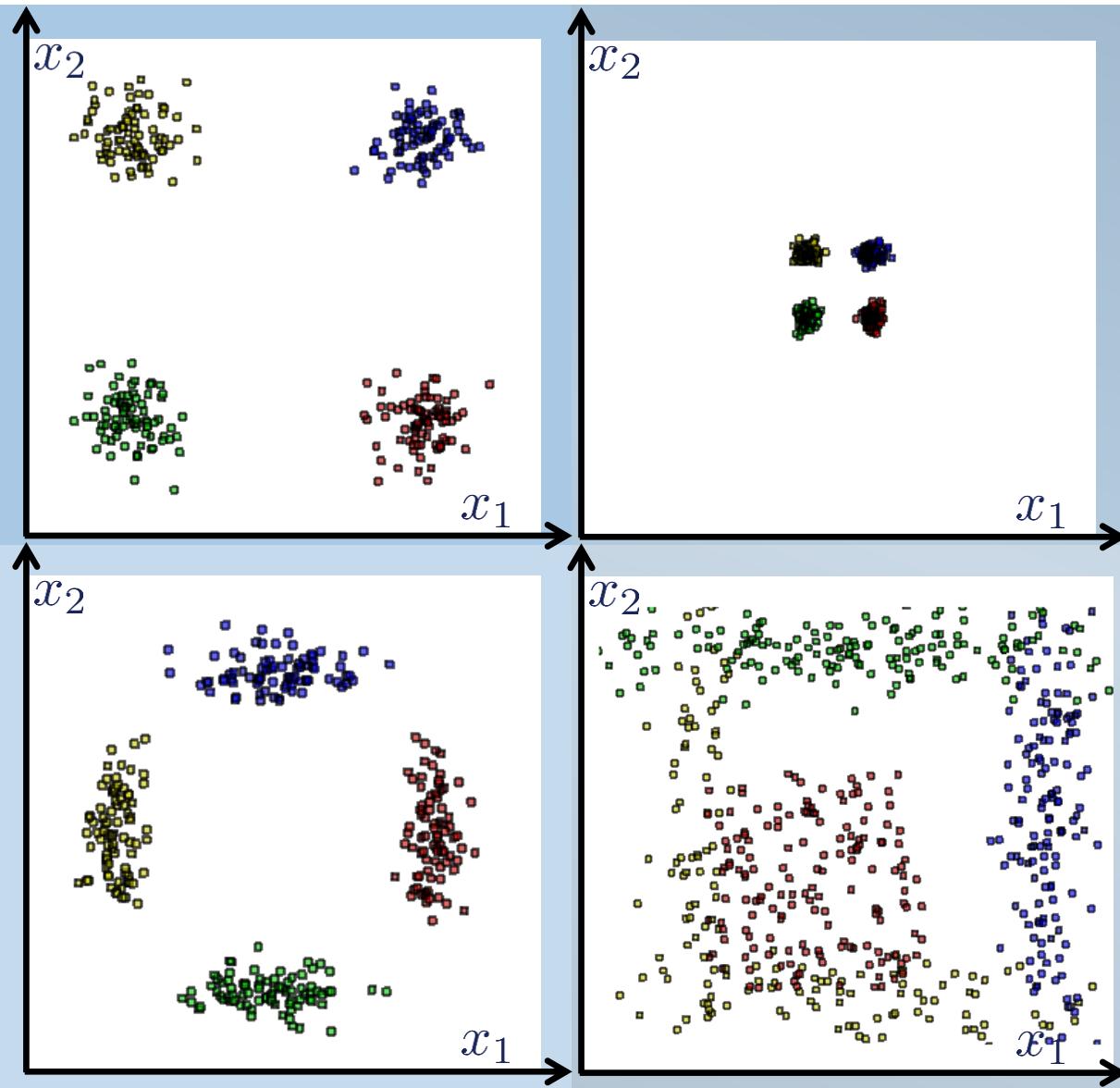
Weak learner: conic section

Randomness: $\rho = 500$ Randomness: $\rho = 5$ 

D=13

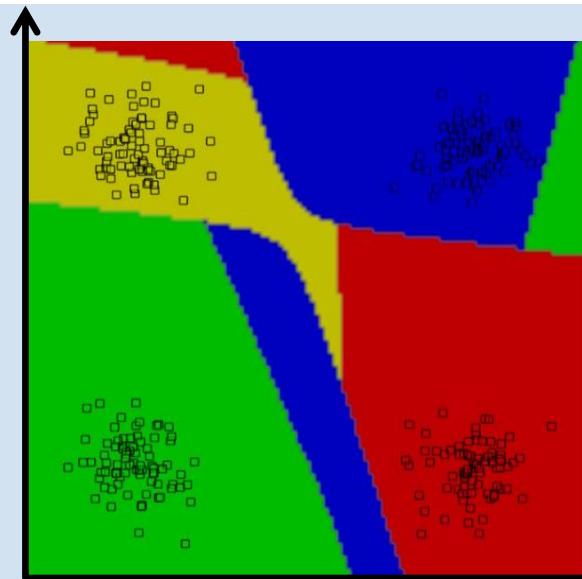
Parameters: T=400 predictor model = prob.

Classification forest: comparison with SVM

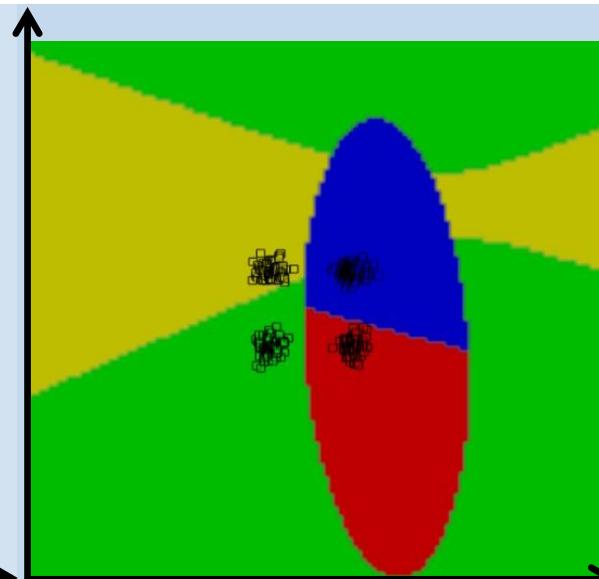


Classification forest: comparison with SVM

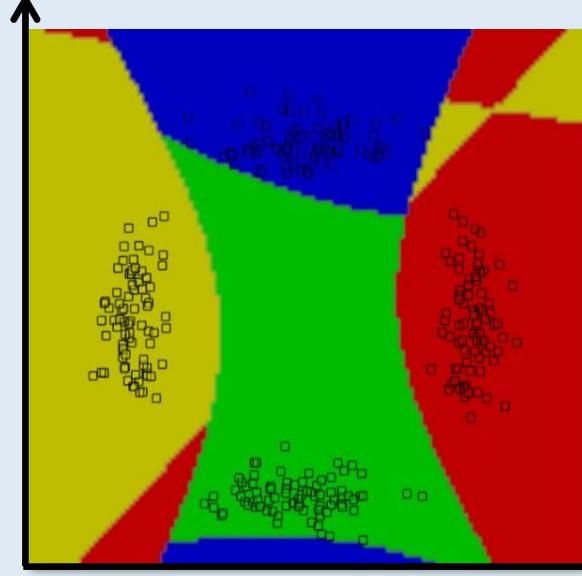
Max-margin like behaviour for multi-class problem



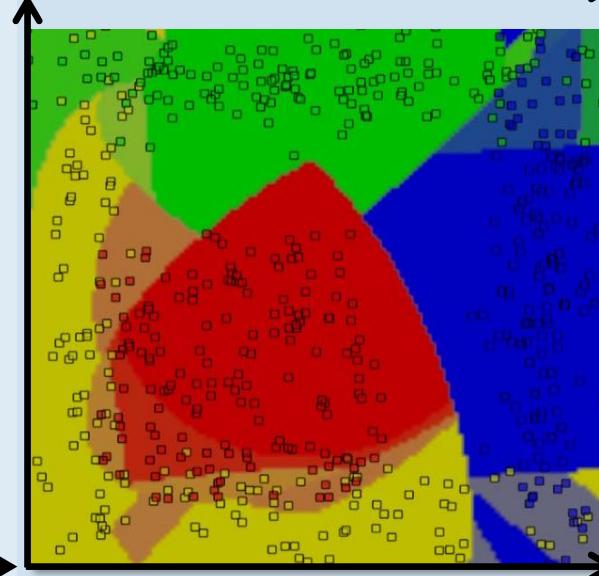
Increased uncertainty away from training points.



Max-margin like behaviour for multi-class problem

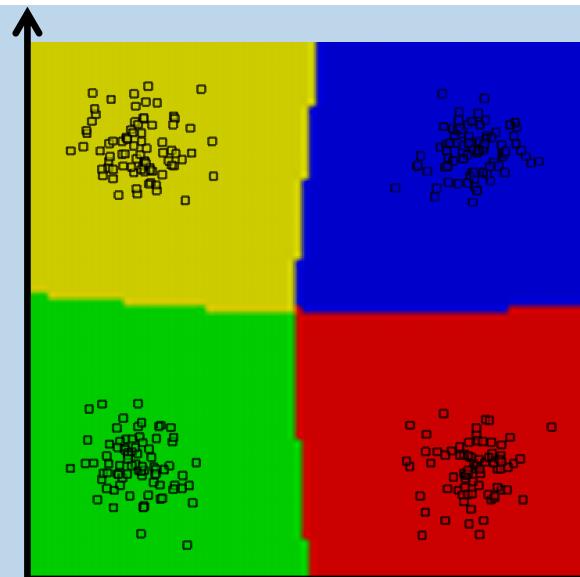


Increased uncertainty in mixed regions

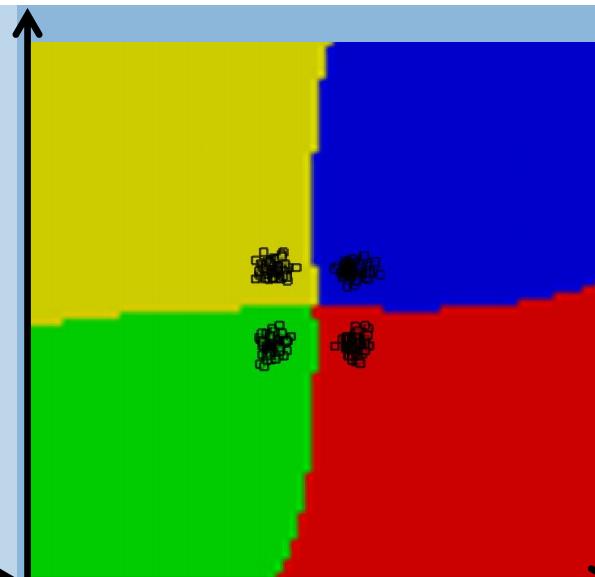


Classification forest: comparison with SVM

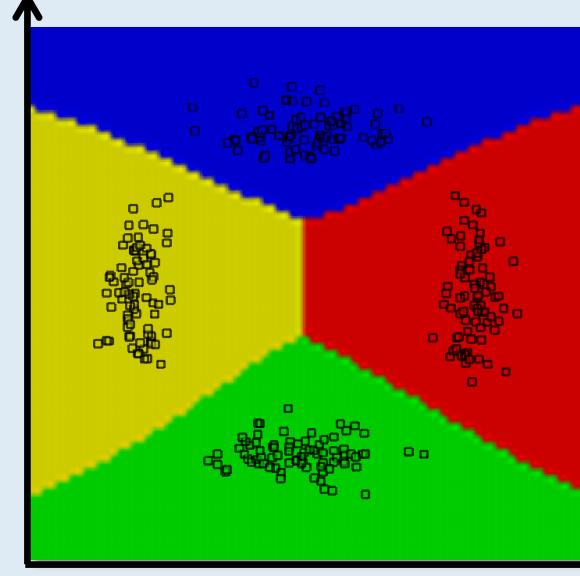
SVM produces nice separation but no confidence information



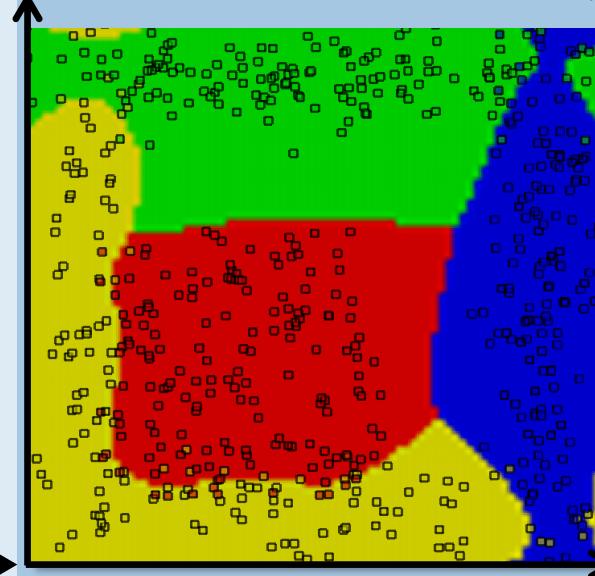
Same high confidence away from training data



Lack of symmetry

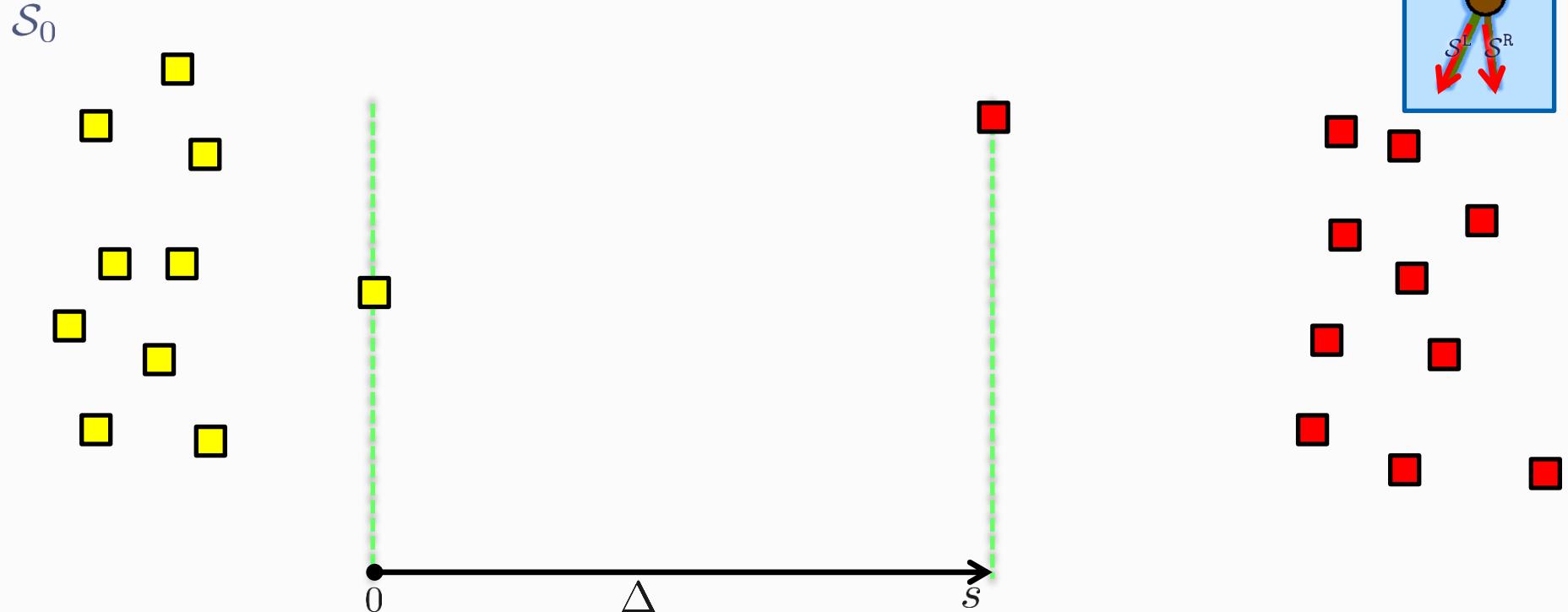


Note overfitting
+
overly confident



Classification forest: max-margin properties

Randomness type = Randomized node optimization



Node parameters

$$\boldsymbol{\theta} = (\phi, \psi, \tau) \in \mathcal{T}$$

ϕ = which features,
 ψ = what geometric model,
 τ = thresholds

Vertical weak learner

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

$$\psi = (1 \ 0 \ 0)$$

$$\phi(\mathbf{v}) = (x_1 \ x_2 \ 1)^\top$$

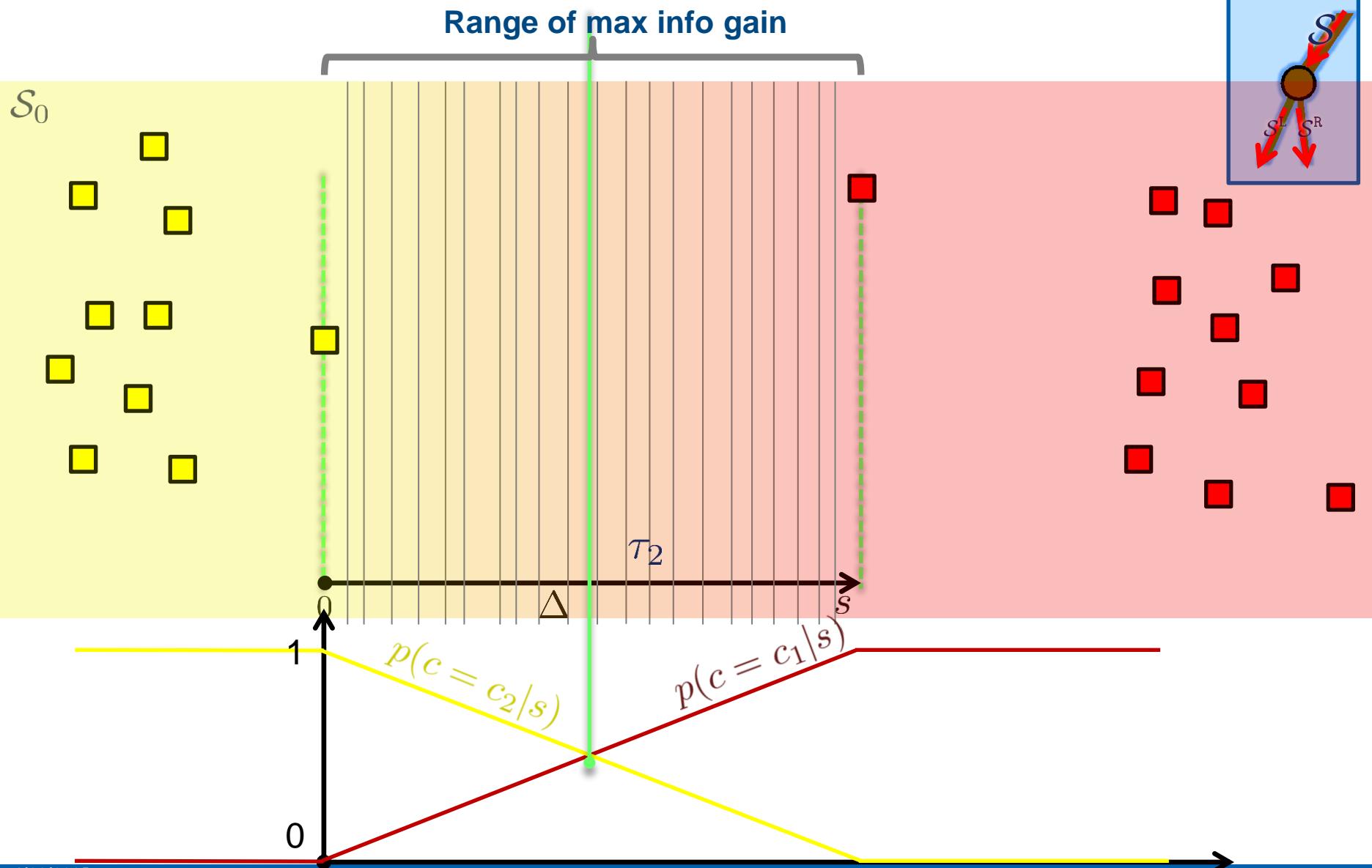
Node training

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

$$\tau_2^* = \arg \max_{\tau_2 \in \mathcal{T}_j} I$$

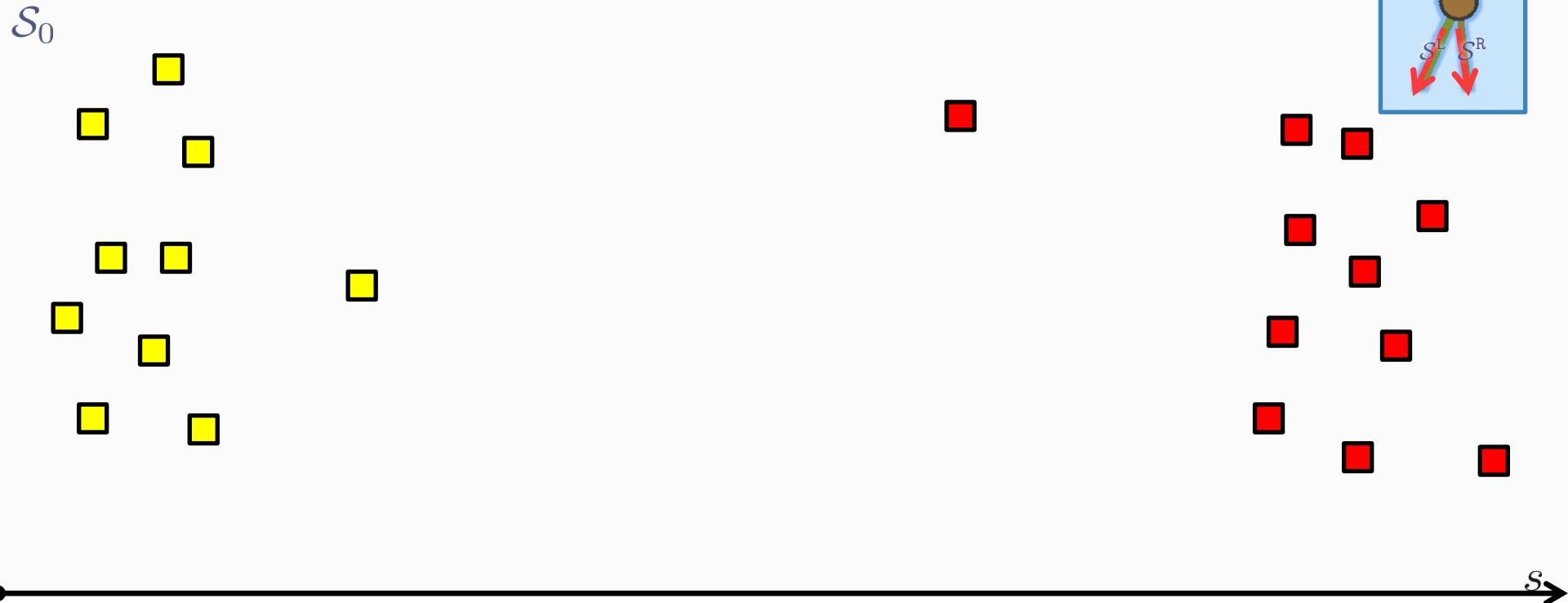
Classification forest: max-margin properties

Randomness type = Randomized node optimization



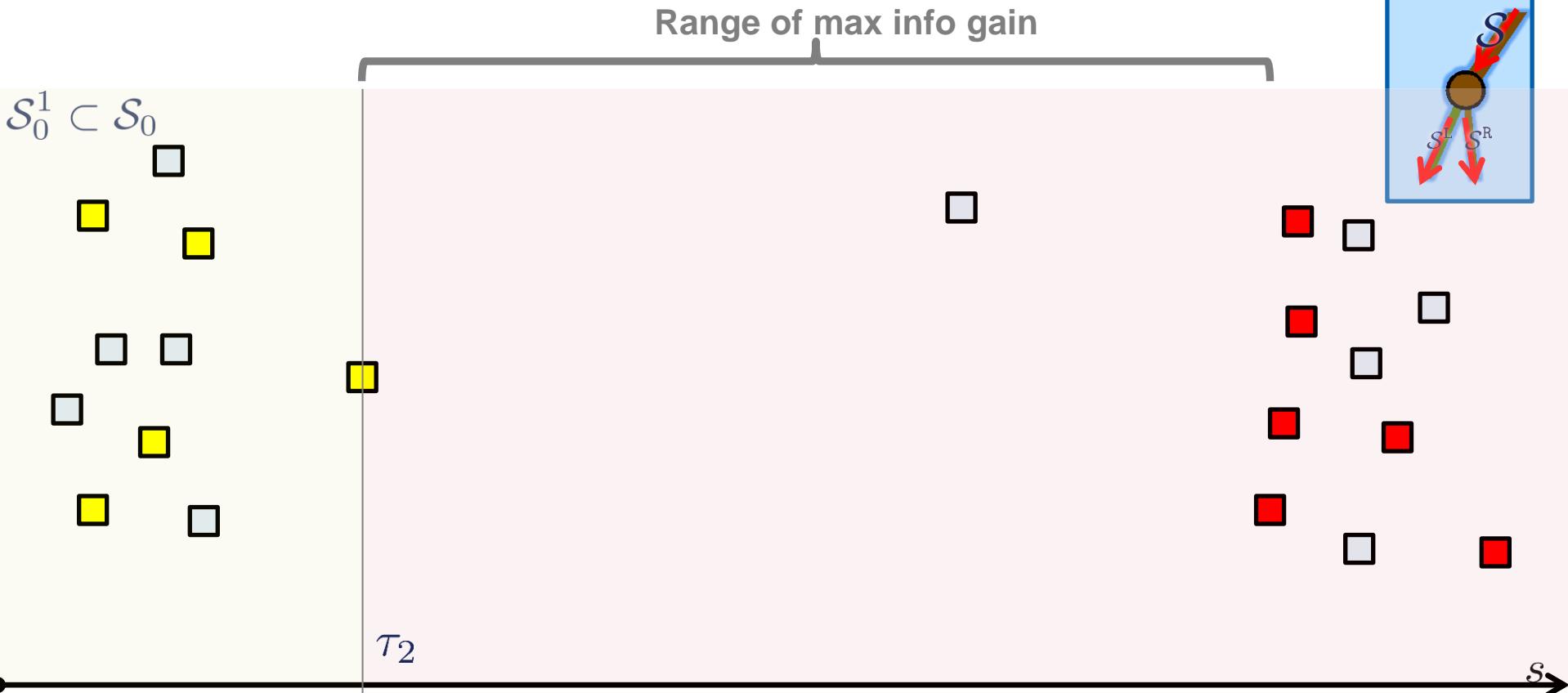
Classification forest: max-margin properties

Randomness type = Bagging (only)



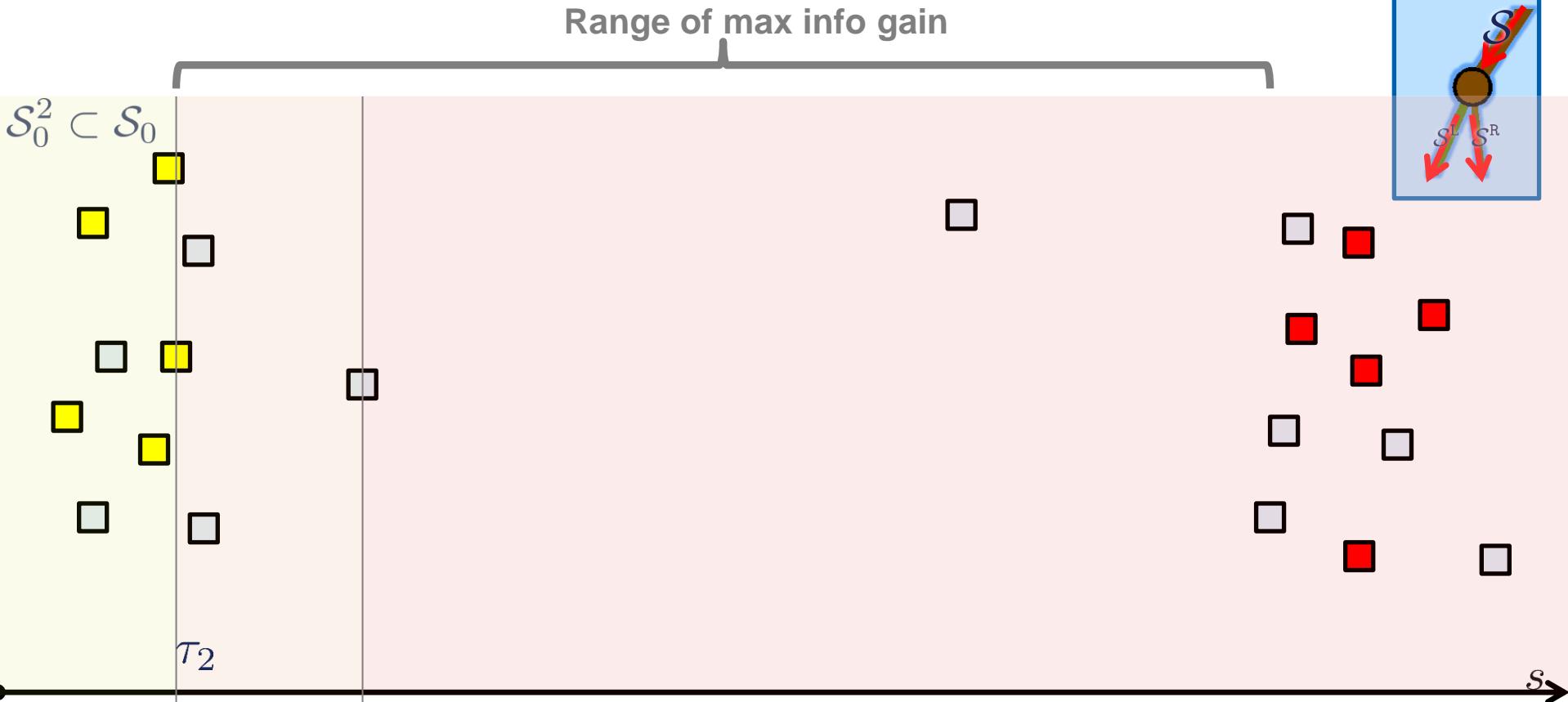
Classification forest: max-margin properties

Randomness type = Bagging (only)



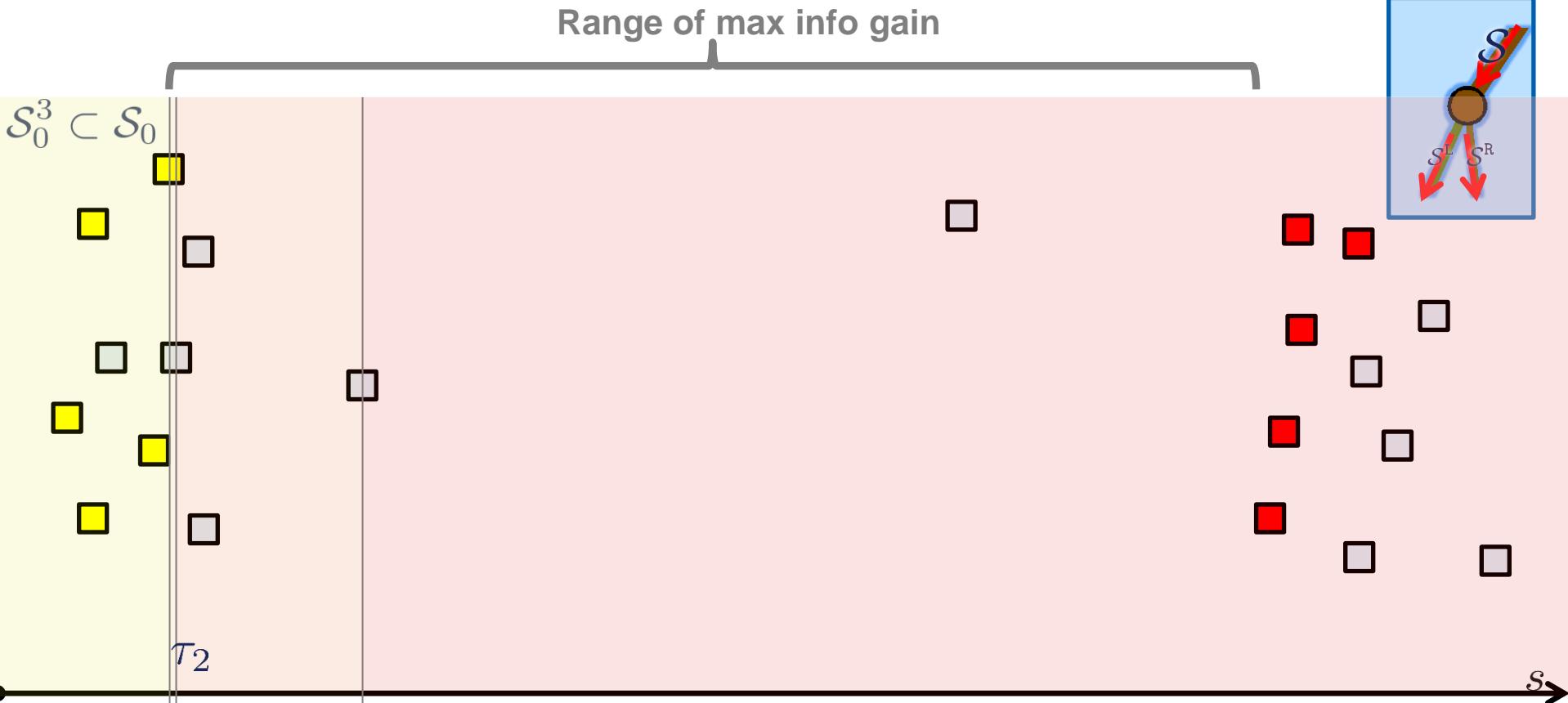
Classification forest: max-margin properties

Randomness type = Bagging (only)



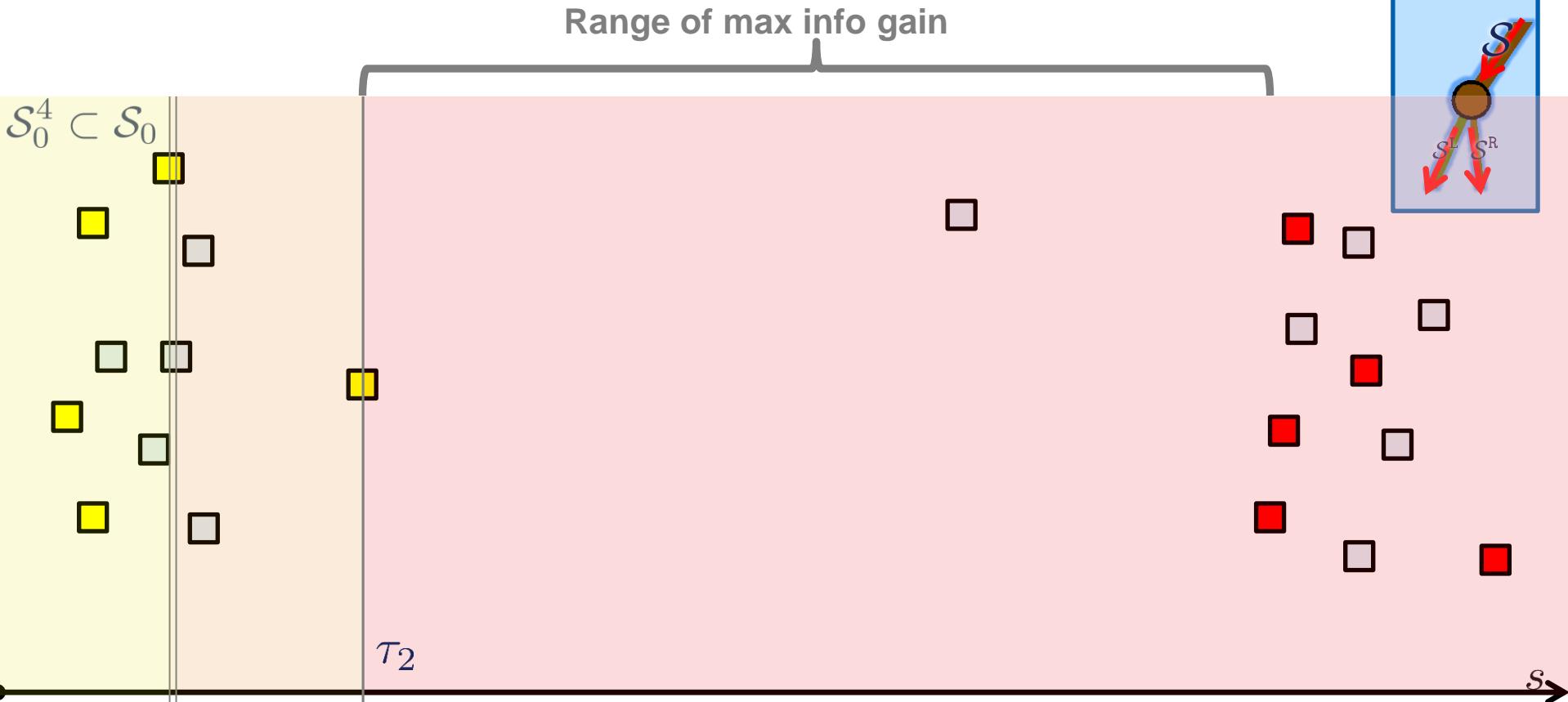
Classification forest: max-margin properties

Randomness type = Bagging (only)



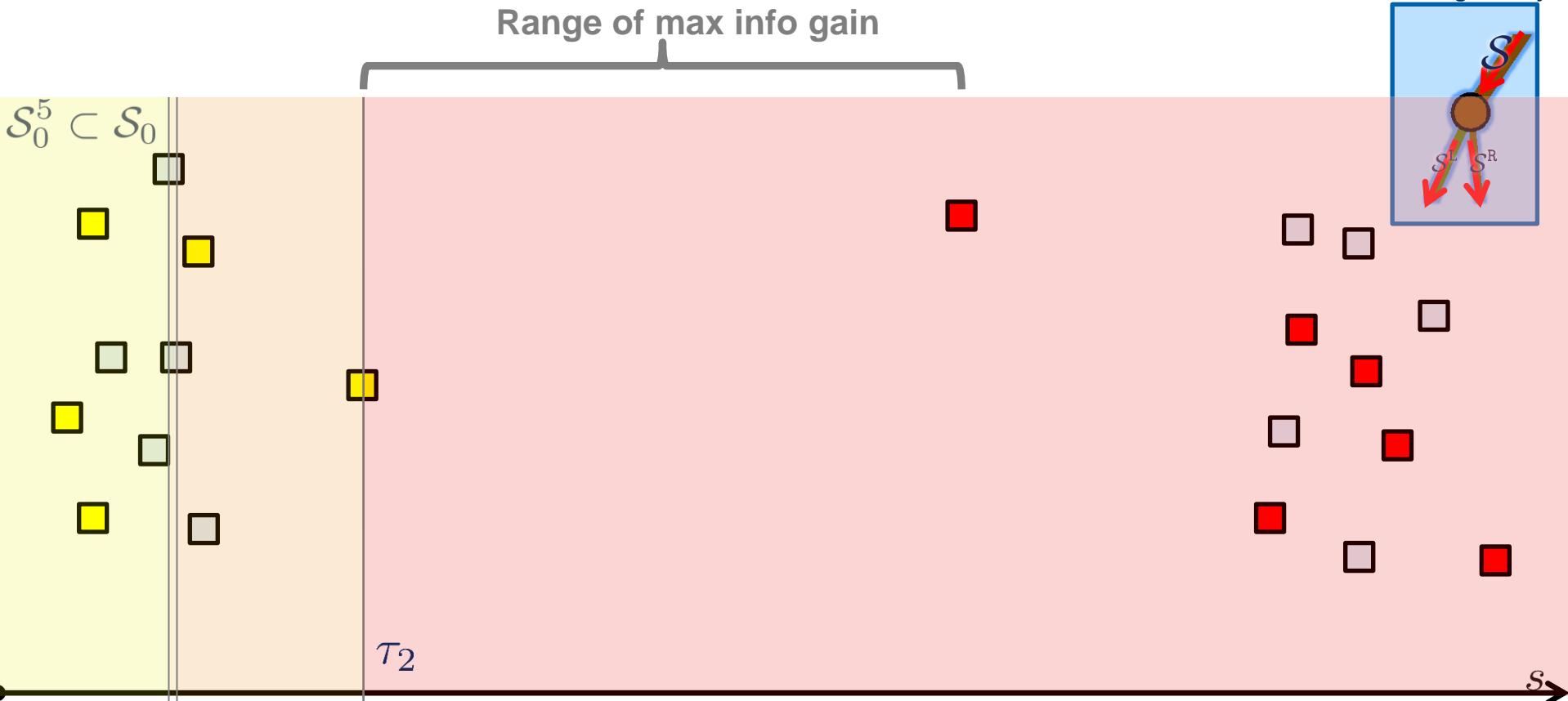
Classification forest: max-margin properties

Randomness type = Bagging (only)



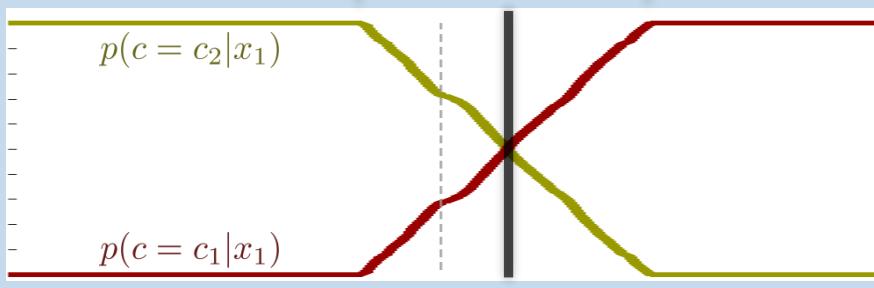
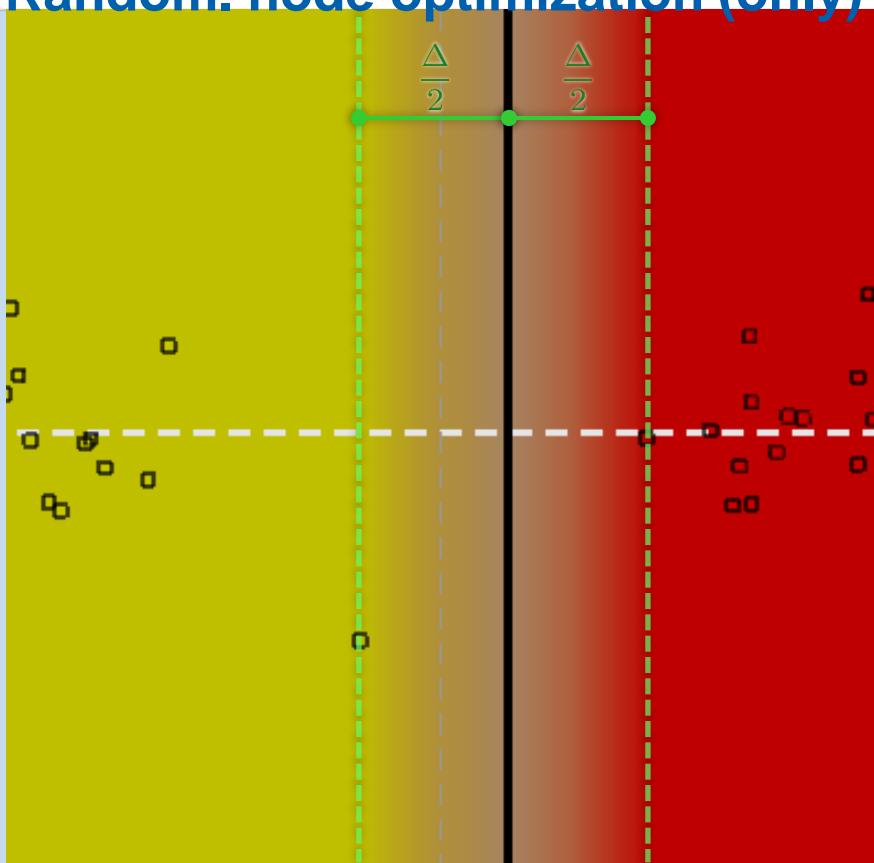
Classification forest: max-margin properties

Randomness type = Bagging (only)

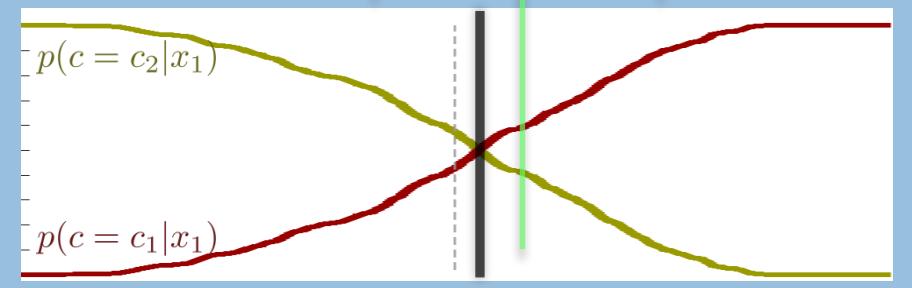
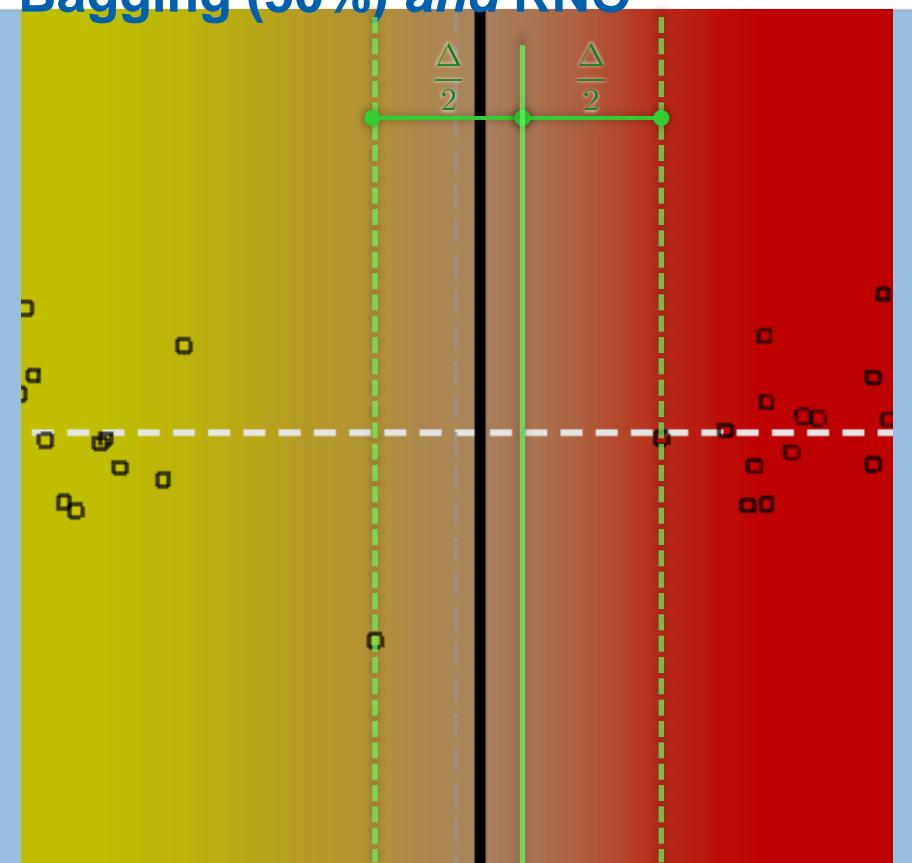


Classification forest: max-margin – comparing with bagging

Random. node optimization (only)

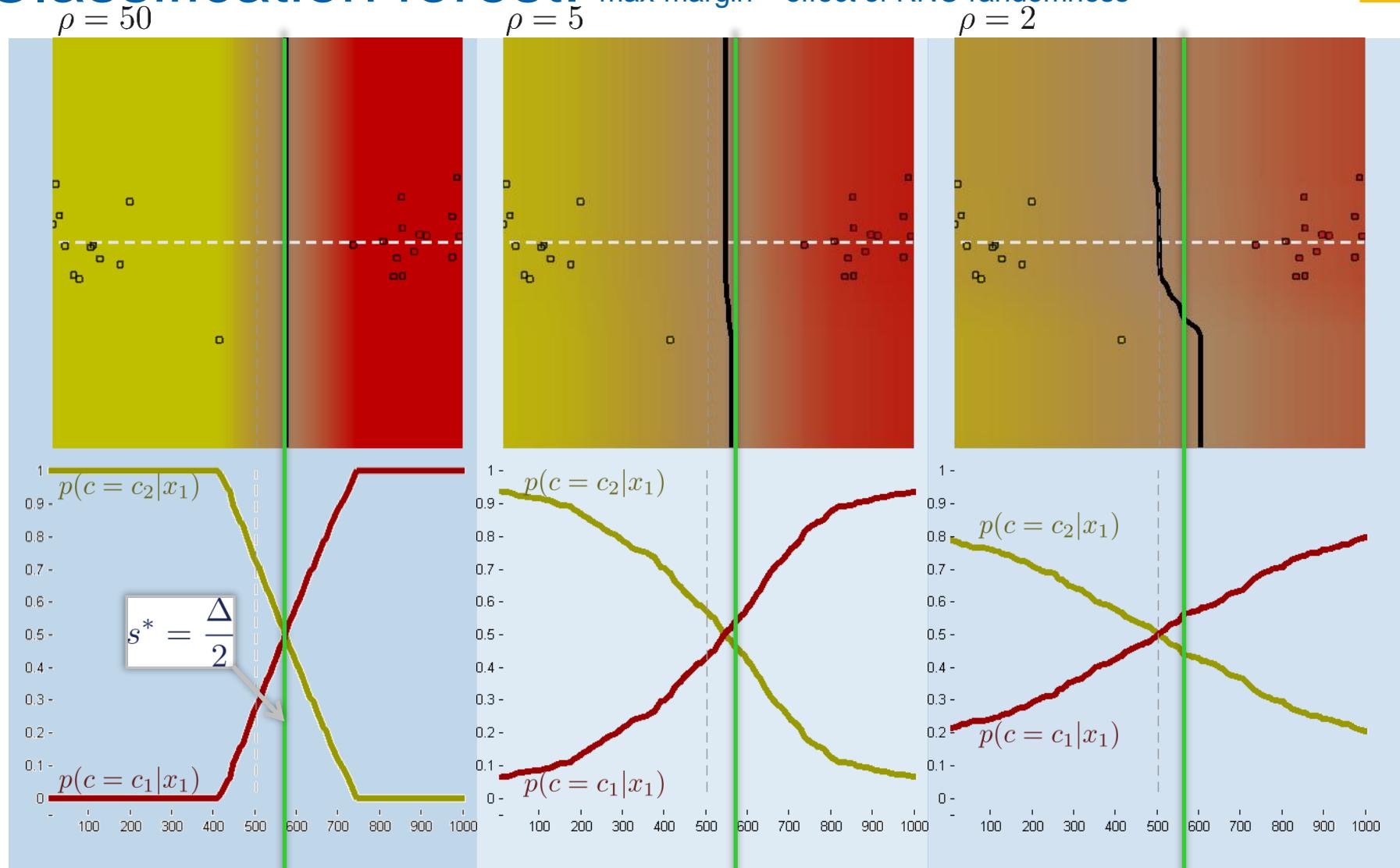


Bagging (50%) and RNO



Parameters: $\rho = 500, D = 2, T = 500$, w.l. axis aligned

Classification forest: max-margin – effect of RNO randomness



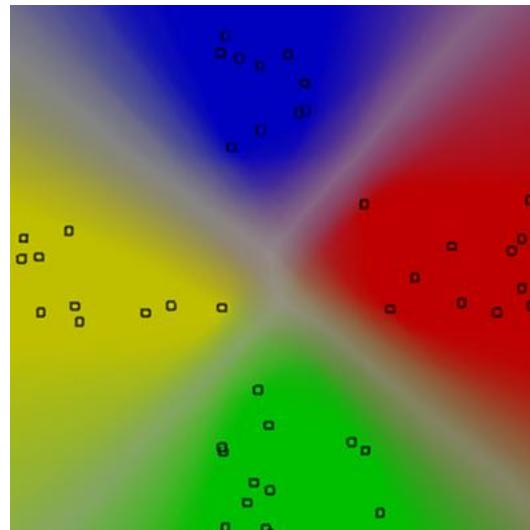
Increased randomness decreases max-margin behaviour.
The separating surface is less affected by individual points.

Parameters: $D = 2, T = 500$, w.l. axis aligned

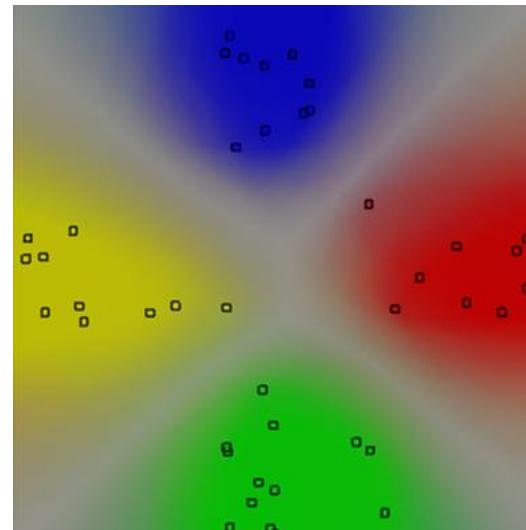
Classification forest: max-margin – for multiple classes



weak learner: oriented line



weak learner: conic section

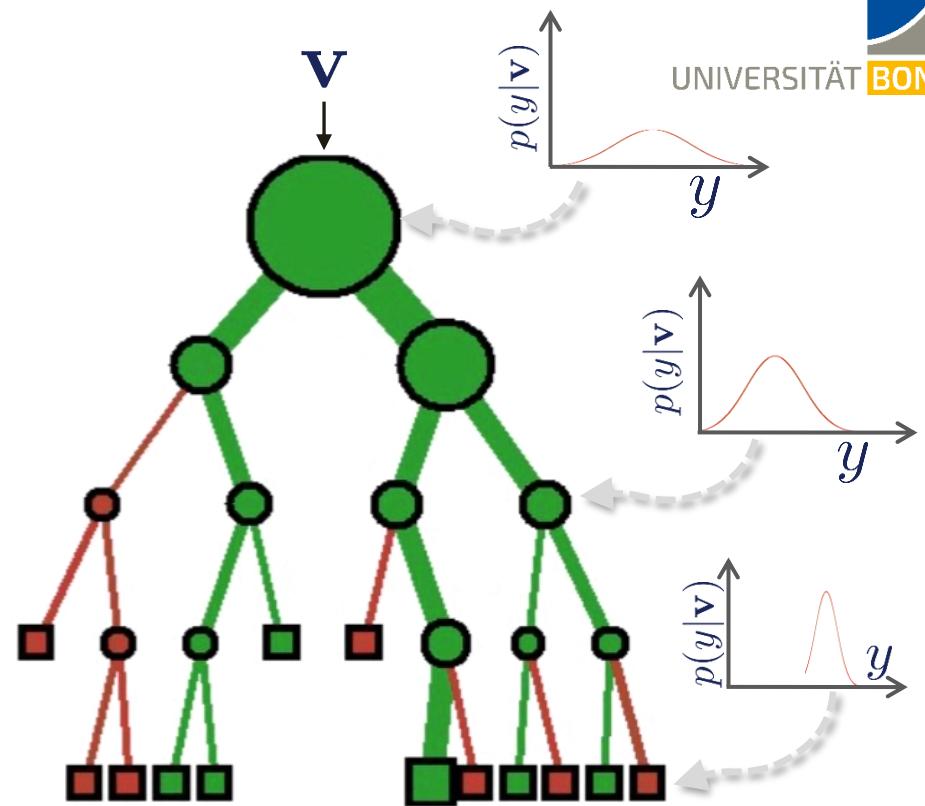
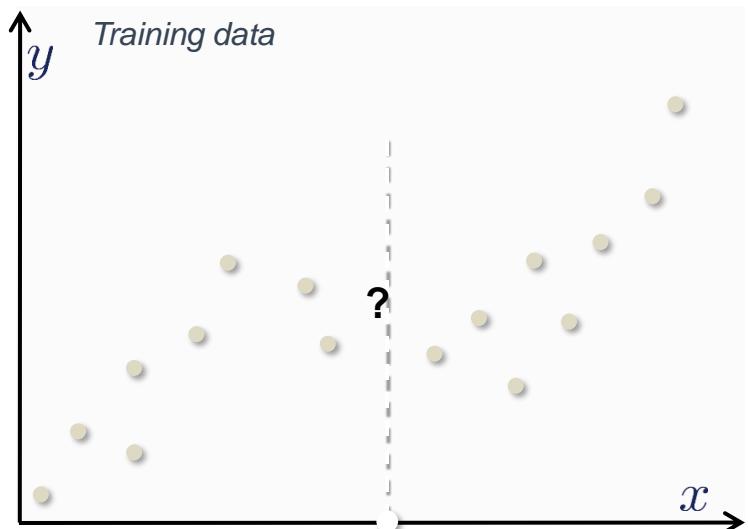


Parameters: $\rho = 50, D = 3, T = 500$

Overview

- Classification forests
- Regression forests
- Other variants

Regression forest



Input:

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

Output:

$$p(\mathbf{y}|\mathbf{v}) \quad \mathbf{y} \in \mathcal{Y} \subset \mathbb{R}^n$$

Information gain:

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

Differential entropy:

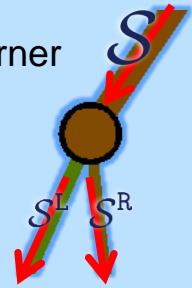
$$H(\mathcal{S}) = \frac{1}{|\mathcal{S}|} \sum_{x \in \mathcal{S}} \int_{\mathcal{Y}} p(y|x) \log p(y|x) dy$$

Regression forest: the node weak learner model

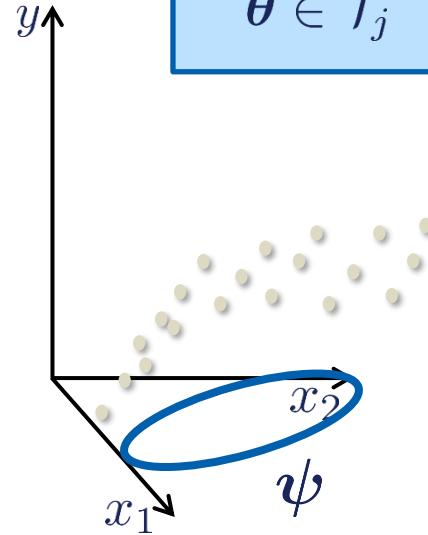
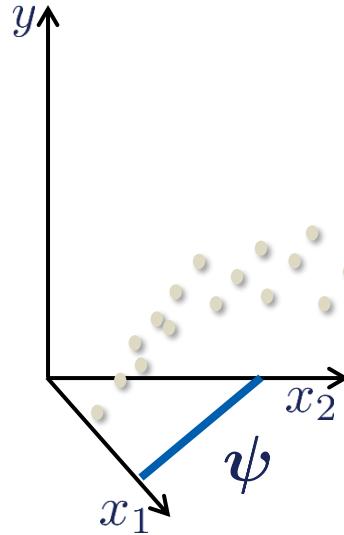
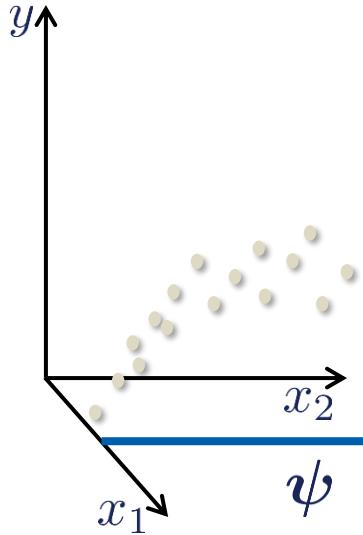
Splitting data at node j

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

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



Examples of node weak learners



Weak learner: axis aligned

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

Feature response
for 2D example.
 $\phi(\mathbf{v}) = (x_1 \ x_2 \ 1)^\top$
With $\psi = (1 \ 0 \ \psi_3)$ or $\psi = (0 \ 1 \ \psi_3)$

Weak learner: oriented line

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

Feature response
for 2D example.
 $\phi(\mathbf{v}) = (x_1 \ x_2 \ 1)^\top$
With $\psi \in \mathbb{R}^3$ a generic line in homog. coordinates

Weak learner: conic section

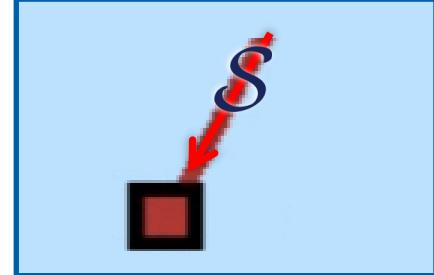
$$h(\mathbf{v}, \theta) = [\tau_1 > \phi^\top(\mathbf{v}) \psi \phi(\mathbf{v}) > \tau_2]$$

Feature response
for 2D example.
 $\phi(\mathbf{v}) = (x_1 \ x_2 \ 1)^\top$
With $\psi \in \mathbb{R}^{3 \times 3}$ a matrix representing a conic.

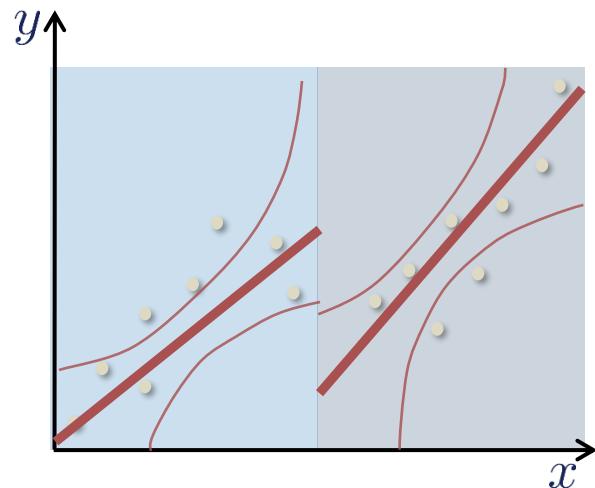
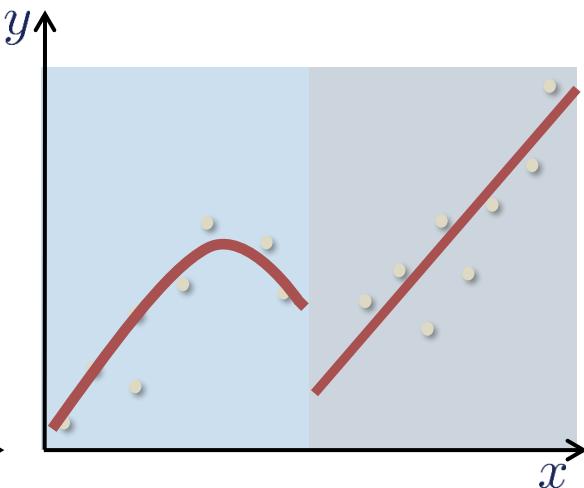
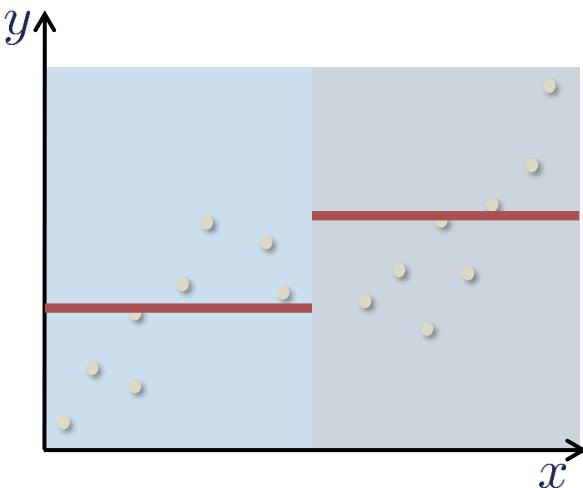
In general ϕ may select only a very small subset of features $\phi(\mathbf{v}) : \mathbb{R}^d \rightarrow \mathbb{R}^{d'+1}$, $d' \ll d$

Regression forest: the predictor model

What do we do at the leaf?



Examples of leaf (predictor) models



Predictor model: constant

$$y = \text{const}$$

Predictor model: polynomial

$$y = \sum_{i=0}^n w_i x^i$$

(note: linear for n=1, constant for n=0)

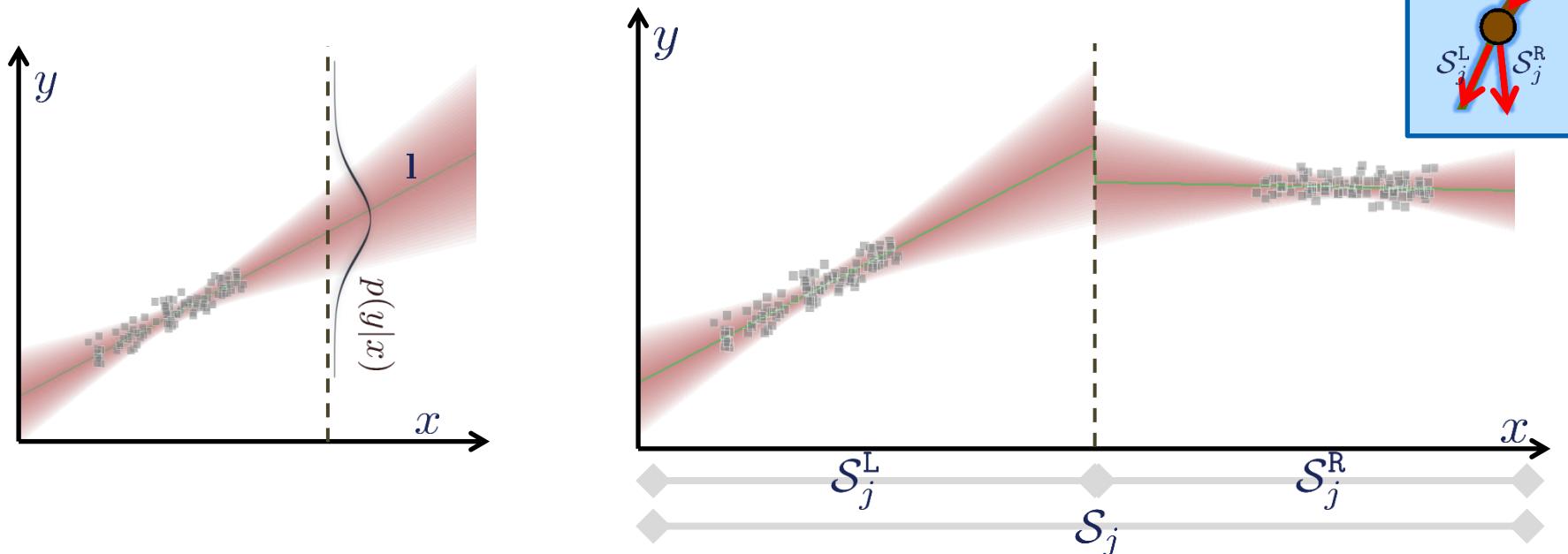
Predictor model: probabilistic-linear

$$y = l_1 x + l_2 \quad p(y|x)$$



Regression forest: objective function

Computing the regression information gain at node j



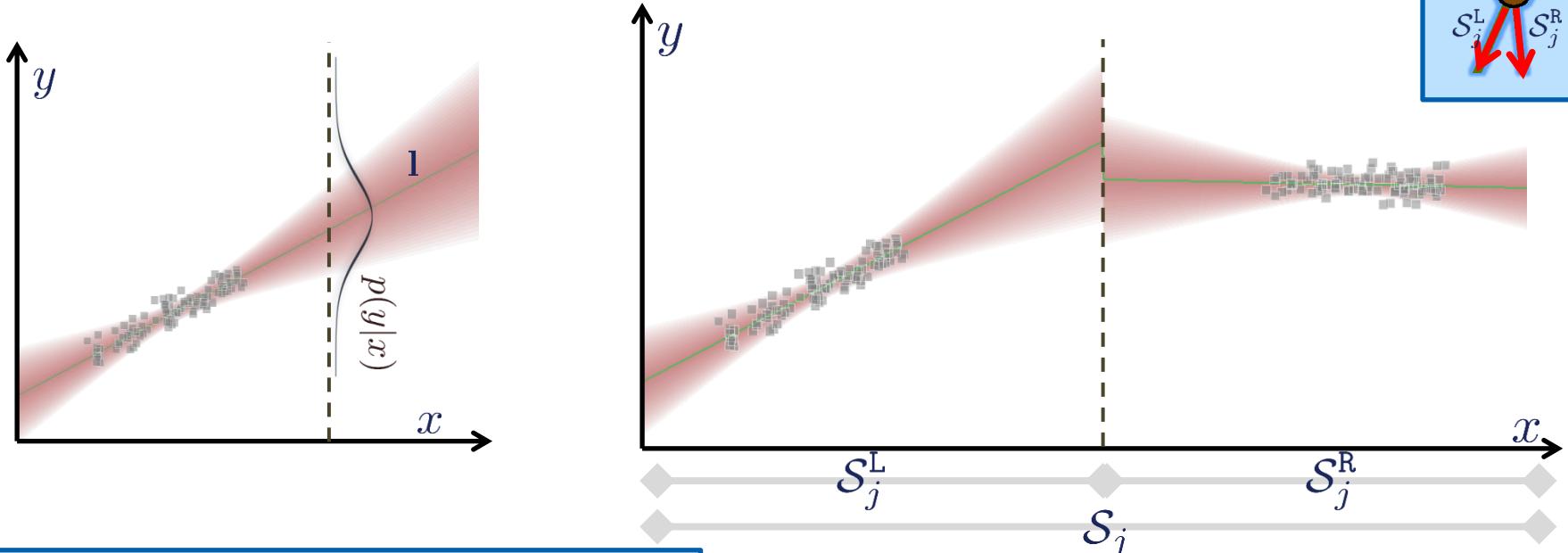
$$I(\mathcal{S}_j, \theta) = H(\mathcal{S}_j) - \sum_{i \in \{\text{L}, \text{R}\}} \frac{|\mathcal{S}_j^i|}{|\mathcal{S}_j|} H(\mathcal{S}_j^i)$$

$$H(\mathcal{S}) = \frac{1}{|\mathcal{S}|} \sum_{x \in \mathcal{S}} \int_y p(y|x) \log p(y|x) dy \quad p(y|x) \sim N(y; \bar{y}, \sigma_y^2(x))$$

$$H(\mathcal{S}) = \frac{1}{|\mathcal{S}|} \sum_{x \in \mathcal{S}} \frac{1}{2} \log ((2\pi e)^2 \sigma_y^2(x)) \quad \text{Differential entropy of Gaussian}$$

$$|\mathcal{S}_j| I = \sum_{(x,y) \in \mathcal{S}_j} \log (\sigma_y(x)) - \sum_{i \in \{\text{L}, \text{R}\}} \left(\sum_{(x,y) \in \mathcal{S}_j^i} \log (\sigma_y(x)) \right) \quad \text{Regression information gain}$$

Regression forest: objective function



Comparison with Breiman's error of fit (CART)

$$I = \sum_{(x,y) \in \mathcal{S}_j} \log(\sigma_y(x)) - \sum_{i \in \{\text{L,R}\}} \left(\sum_{(x,y) \in \mathcal{S}_j^i} \log(\sigma_y(x)) \right)$$

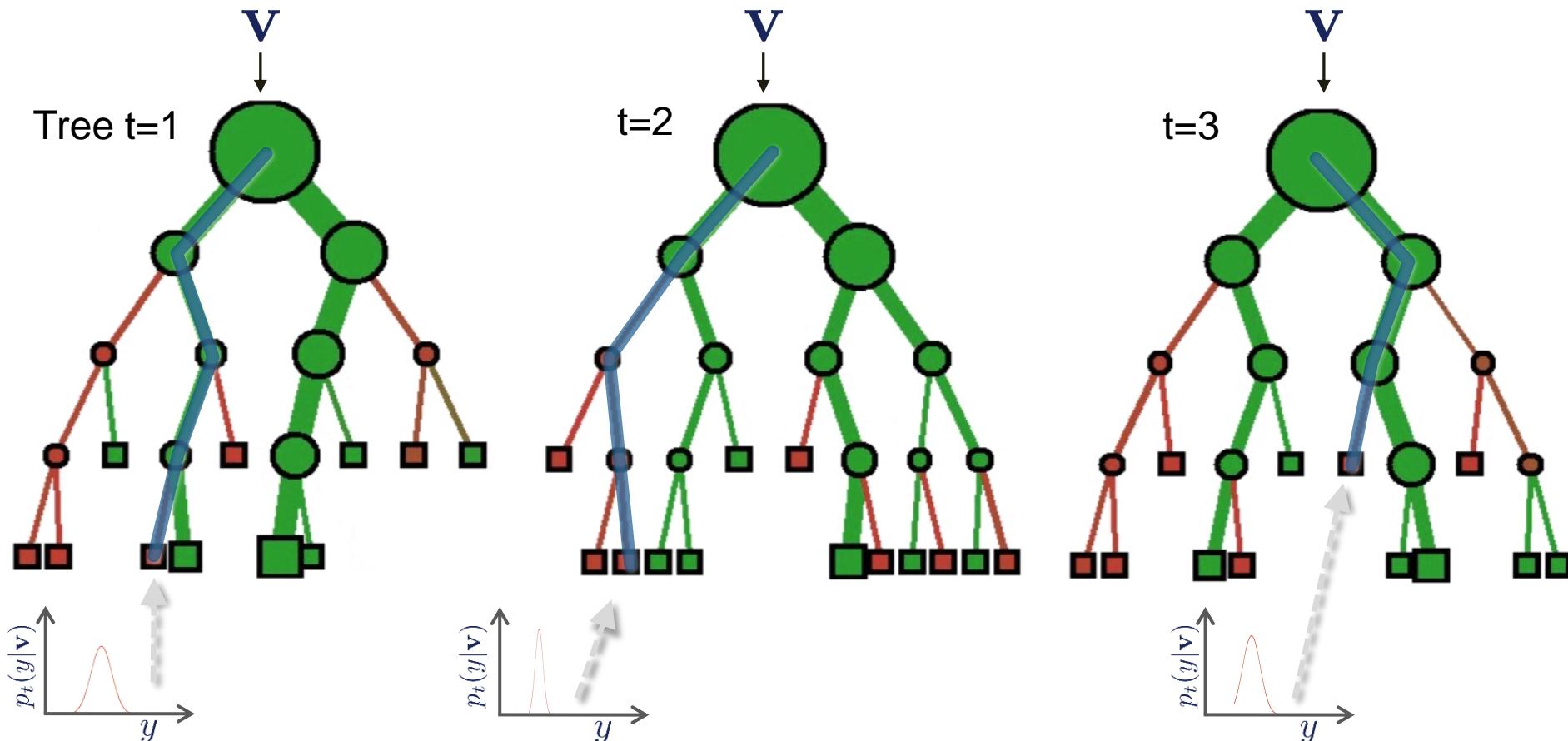
Differential entropy of Gaussian

$$E = \sum_{(x,y) \in \mathcal{S}_j} (y - \bar{y}_j)^2 - \sum_{i \in \{\text{L,R}\}} \left(\sum_{(x,y) \in \mathcal{S}_j^i} (y - \bar{y}_j)^2 \right)$$

Error of fit

The error of fit objective is a special instance of our more general information-theoretical objective function.

Regression forest: the ensemble model



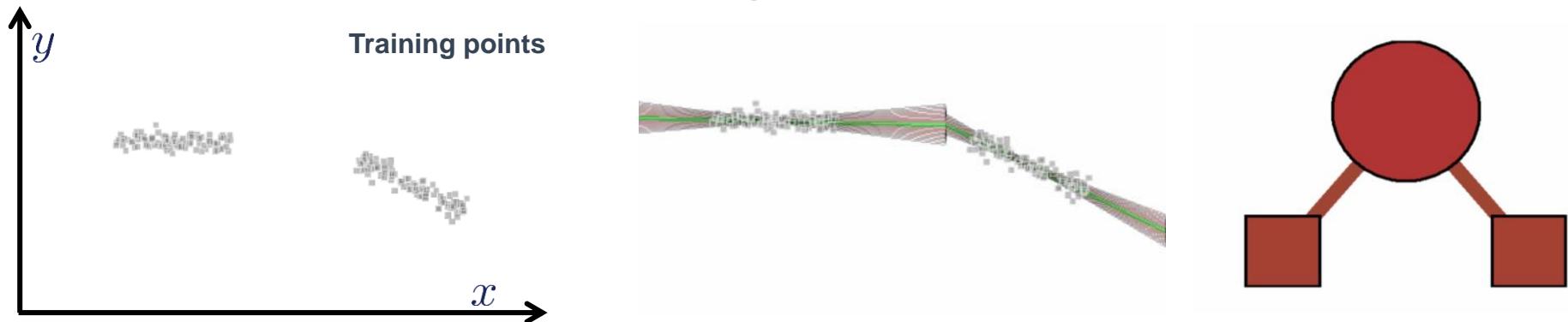
The ensemble model

Forest output probability

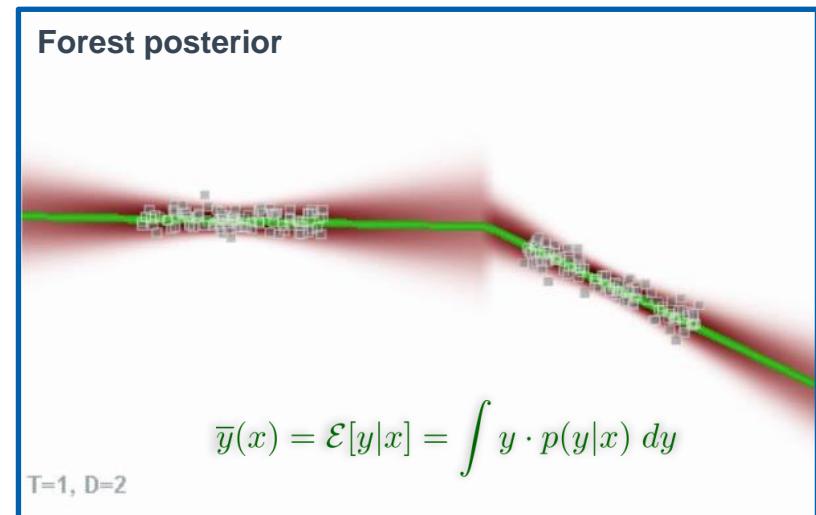
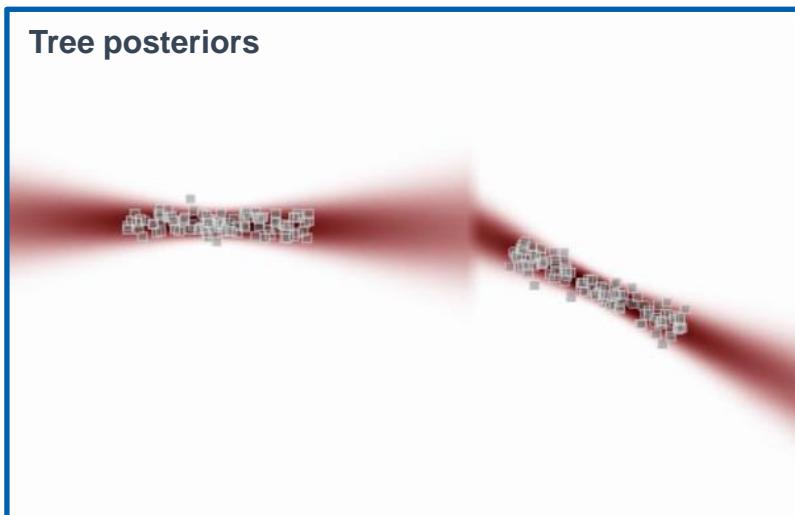
$$p(\mathbf{y}|\mathbf{v}) = \frac{1}{T} \sum_t^T p_t(\mathbf{y}|\mathbf{v})$$

Regression forest: probabilistic, non-linear regression

Training different trees in the forest (max tree depth D=2)



Testing different trees in the forest

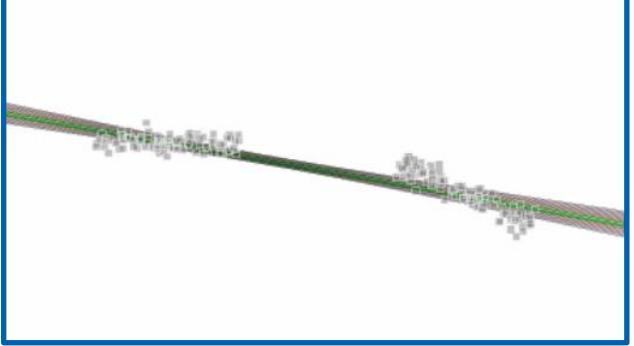
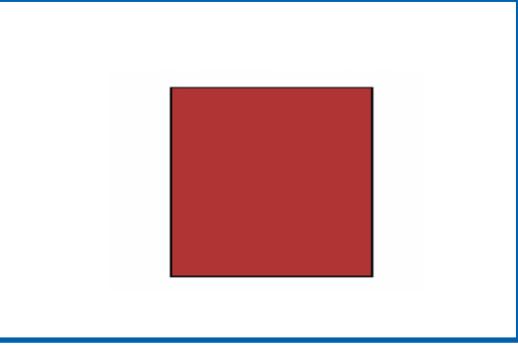
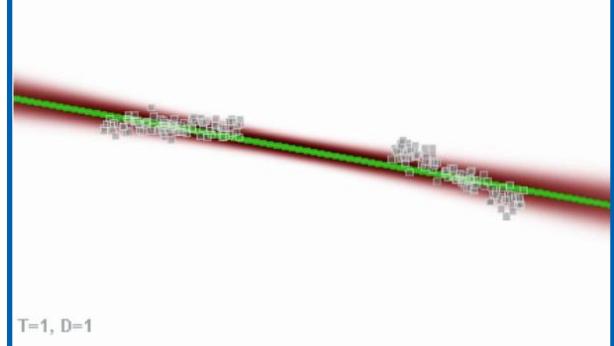
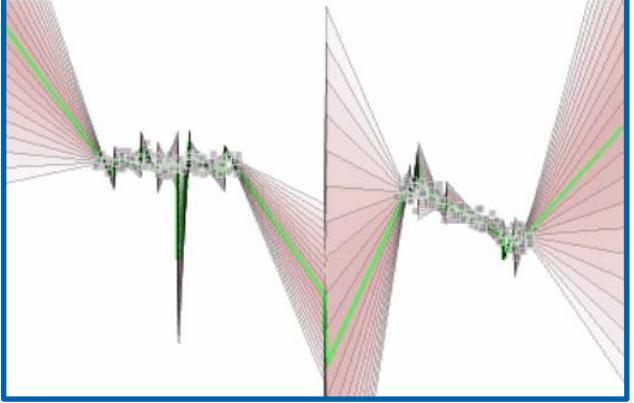
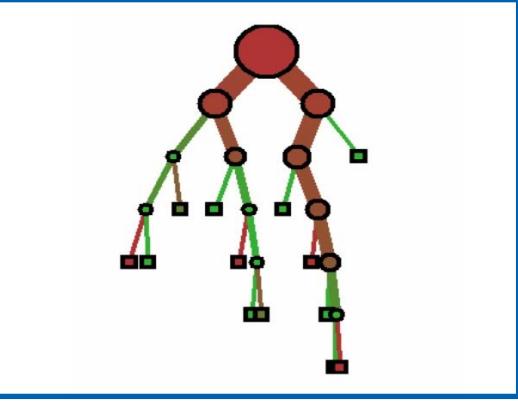
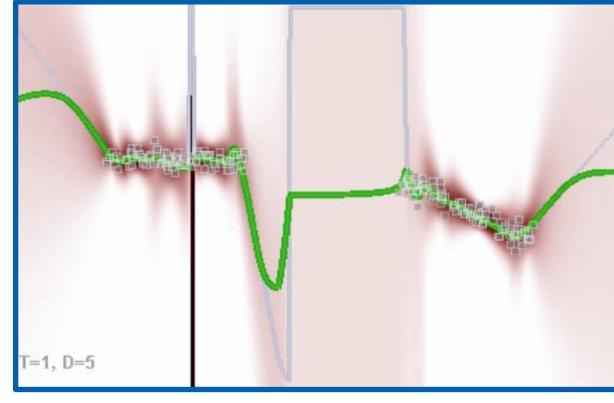


Generalization properties:

- Smooth interpolating behaviour in gaps between training data.
- Uncertainty increases with distance from training data.

Parameters: T=400, D=2, weak learner = aligned, leaf model = prob. line

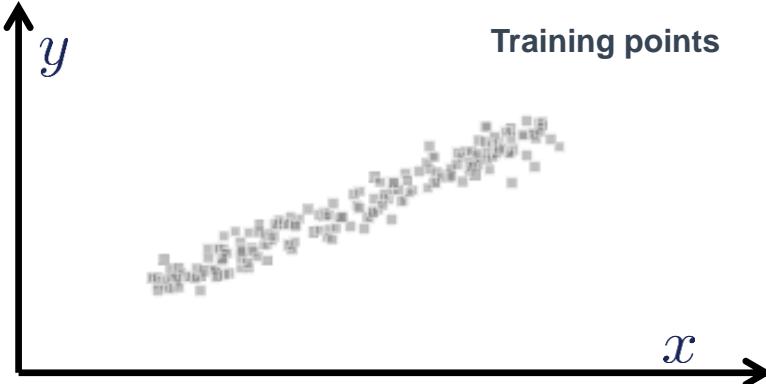
Regression forest: probabilistic, non-linear regression

	Training	Testing
underfitting	<p>max tree depth $D = 1$</p>  	<p>max tree depth $D = 1$</p>  <p>T=1, D=1</p>
overfitting	<p>max tree depth $D = 5$</p>  	<p>max tree depth $D = 5$</p>  <p>T=1, D=5</p>

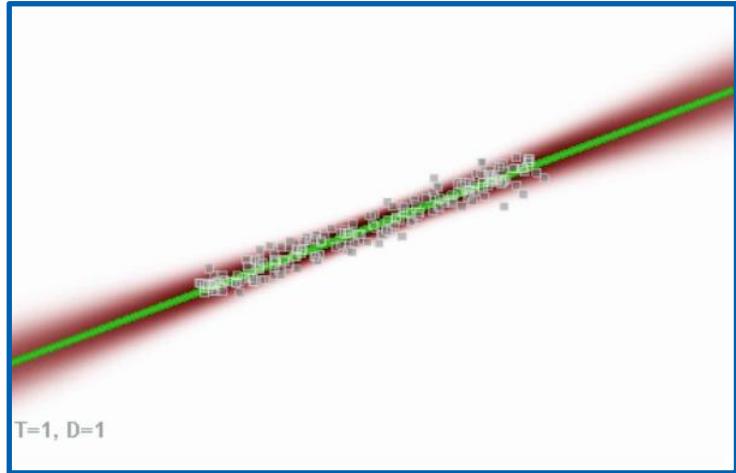
Parameters: T=400, w.l. = aligned, l.m. = prob. line

Regression forest: probabilistic, non-linear regression

Generalizing conventional linear regression



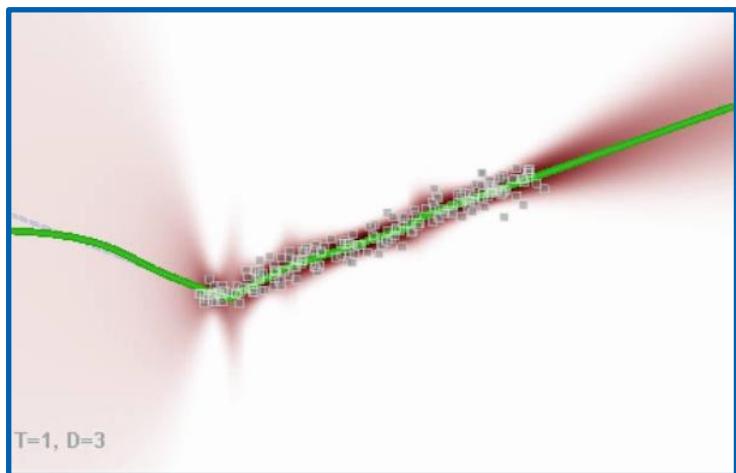
max tree
depth,
 $D = 1$



mean
 $\bar{y}(x) = \mathcal{E}[y|x] = \int y \cdot p(y|x) dy$

mode
 $\hat{y}(x) = \arg \max_y p(y|x)$

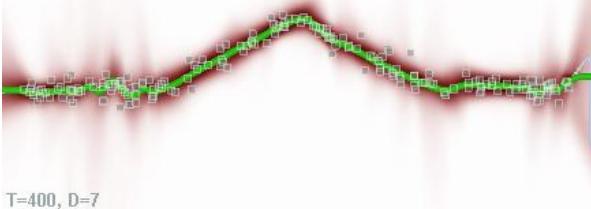
max tree
depth,
 $D = 3$



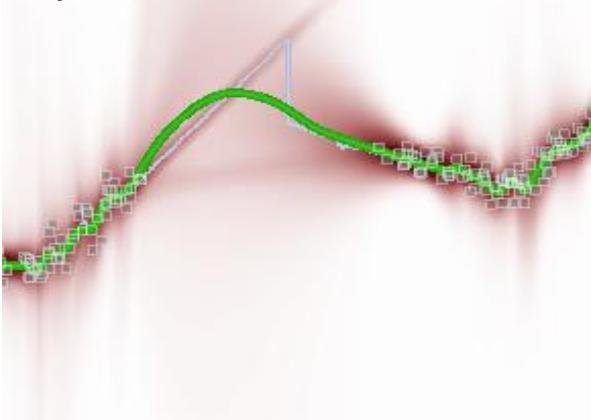
Parameters: $T=400$, w.l. = aligned, l.m. = probabilistic line

Regression forest: varying tree depth D

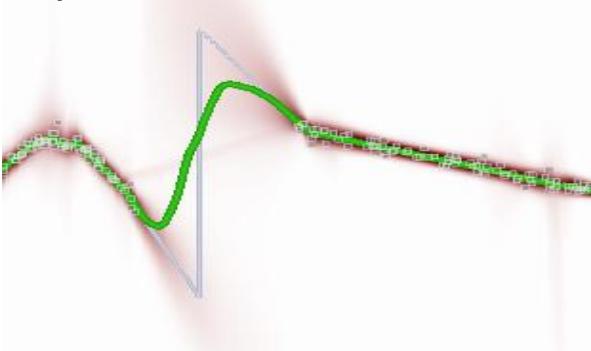
Experiment 5



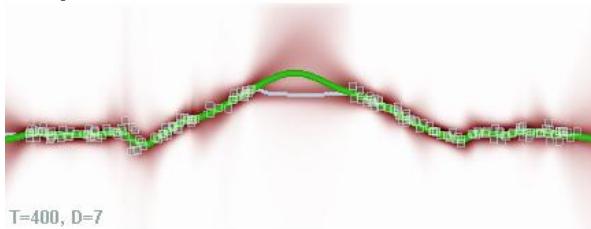
Exp. 7



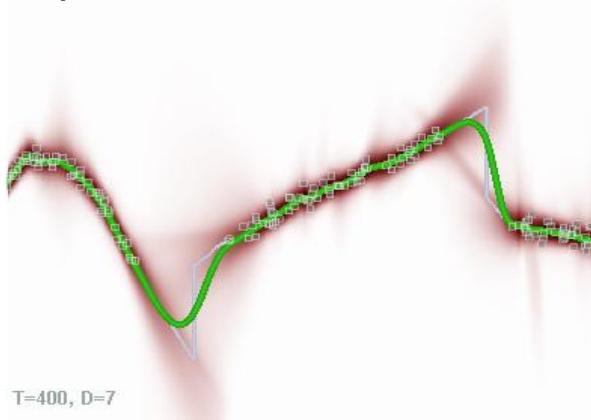
Exp. 9



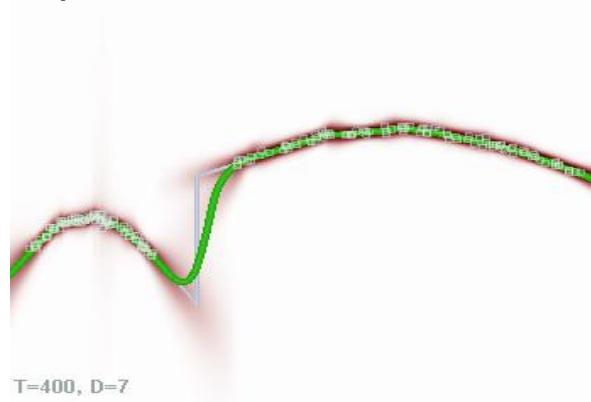
Exp. 6



Exp. 8



Exp. 10



Forest posteriors for six different experiments and varying depth D.

mean

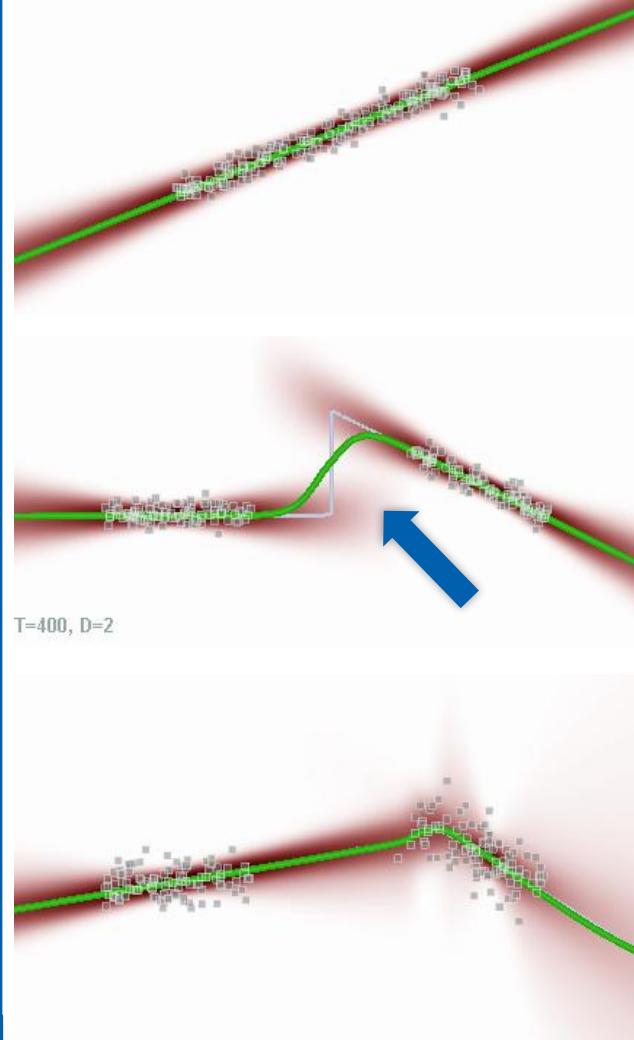
$$\bar{y}(x) = \mathcal{E}[y|x] = \int y \cdot p(y|x) dy$$

mode

$$\hat{y}(x) = \arg \max_y p(y|x)$$

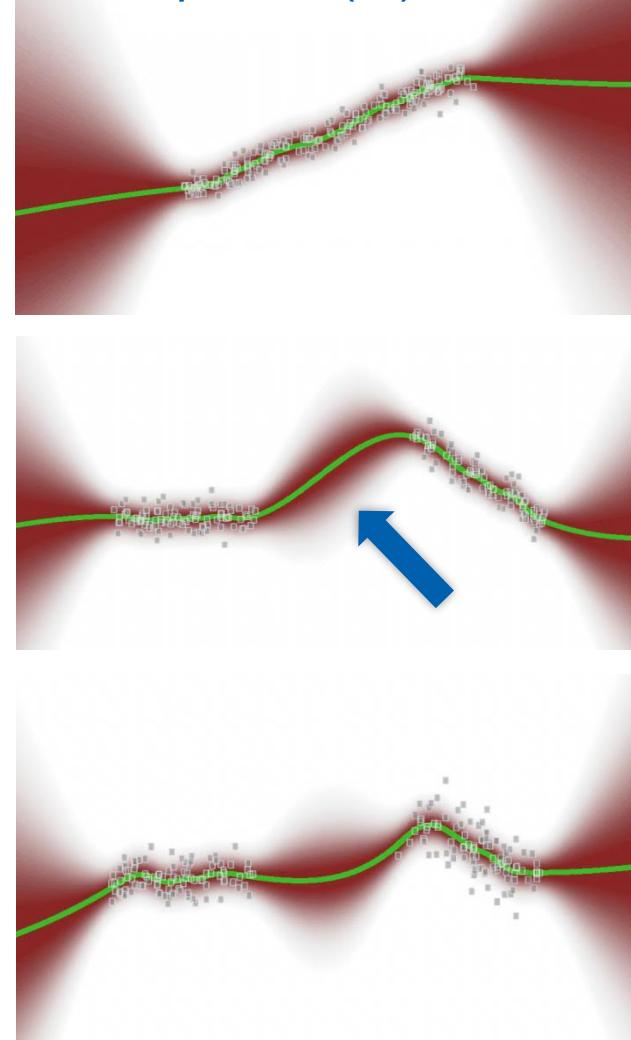
Regression forest: comparison with Gaussian processes

Regression forests



In regr.
forests
the posterior
can be
multimodal

Gaussian processes (GP)



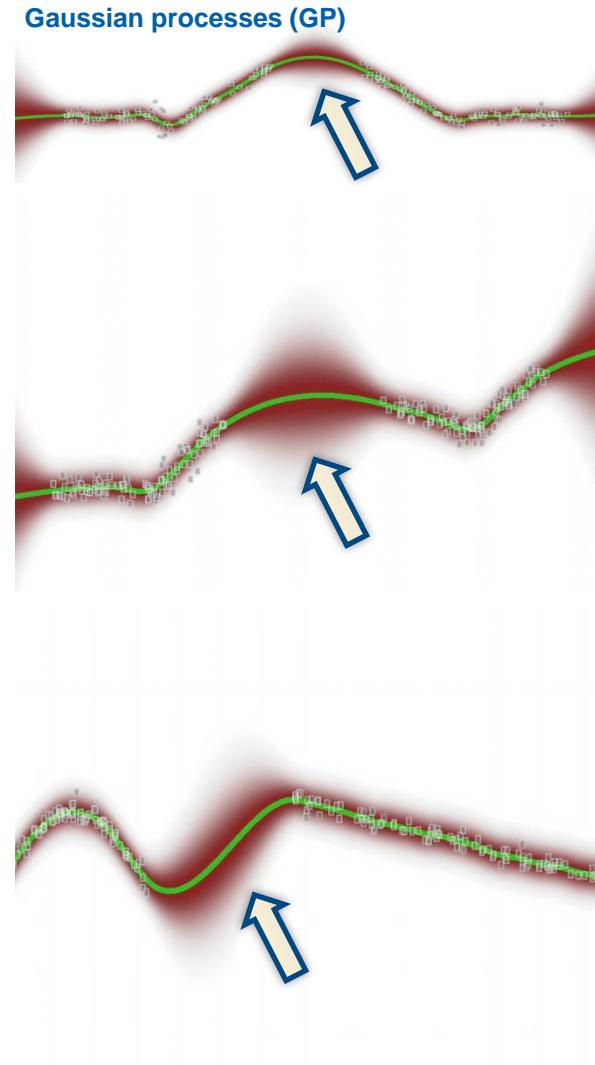
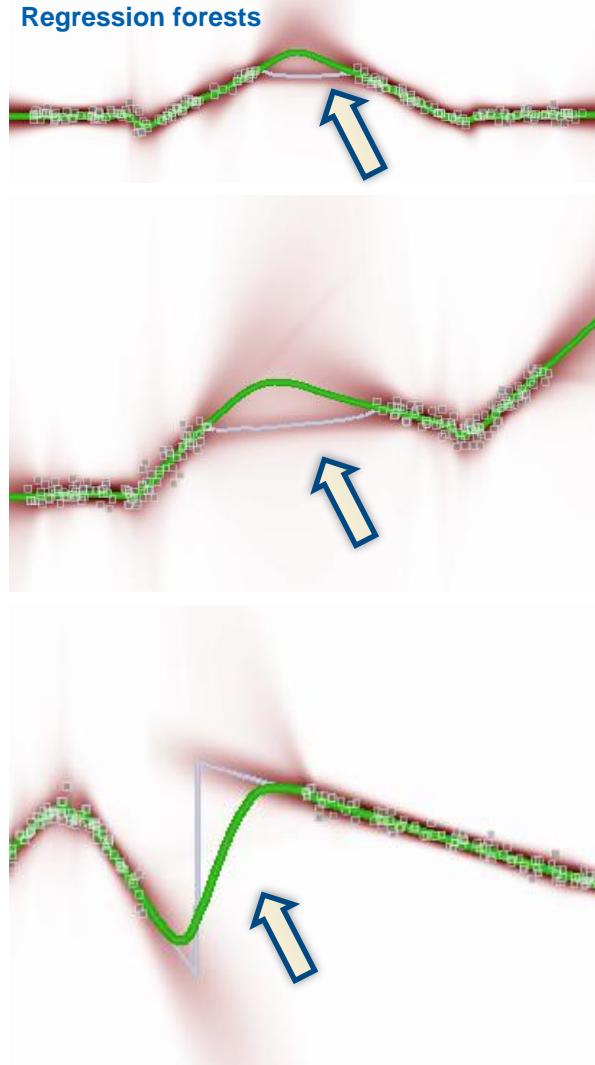
In GP the
posterior
is Gaussian
(unimodal)

Observations

- Similar behaviour of the mean. – Forests can capture multimodality.
- Forests are a parallel algorithm that does not need large matrix inversion.
- In forests the quality and shape of the uncertainty depends on the predictor model.

[GP code from <http://gaussianprocess.org/gpml/>]

Regression forest: comparison with Gaussian processes

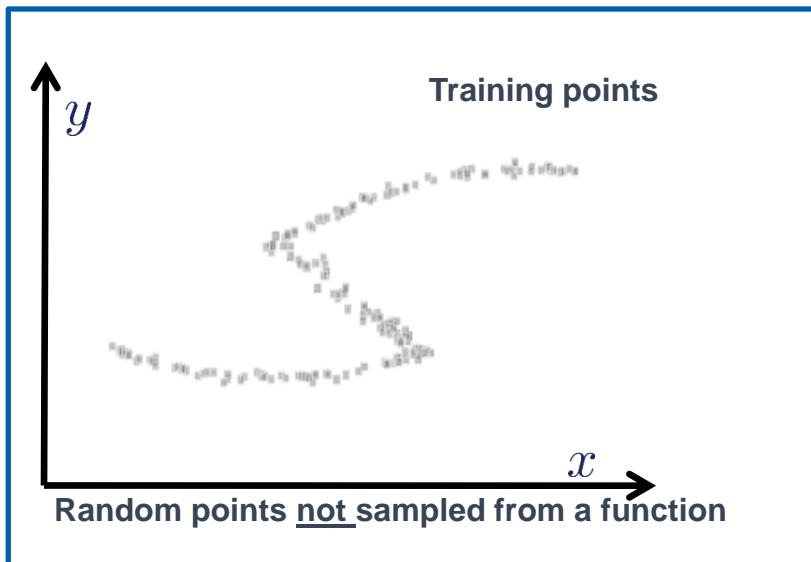


Observations

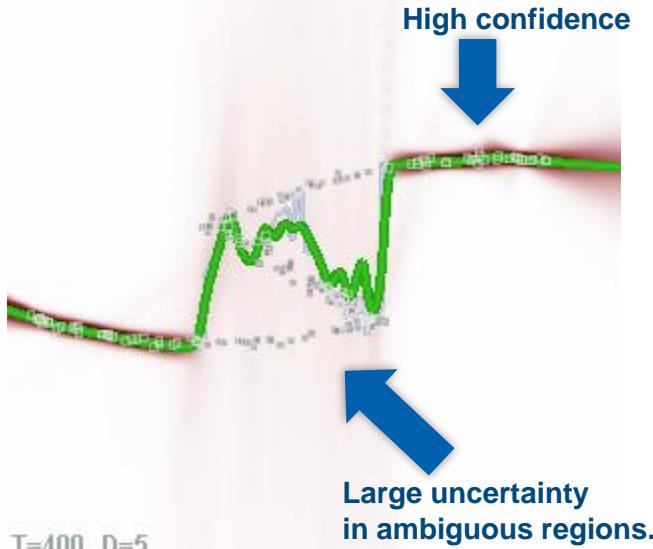
- Similar behaviour of the mean. – Forests can capture multimodality.
- Forests are a parallel algorithm that does not need large matrix inversion.
- In forests the quality and shape of the uncertainty depends on the leaf model.

[GP code from <http://gaussianprocess.org/gpml/>]

Regression forest: dealing with non-functions



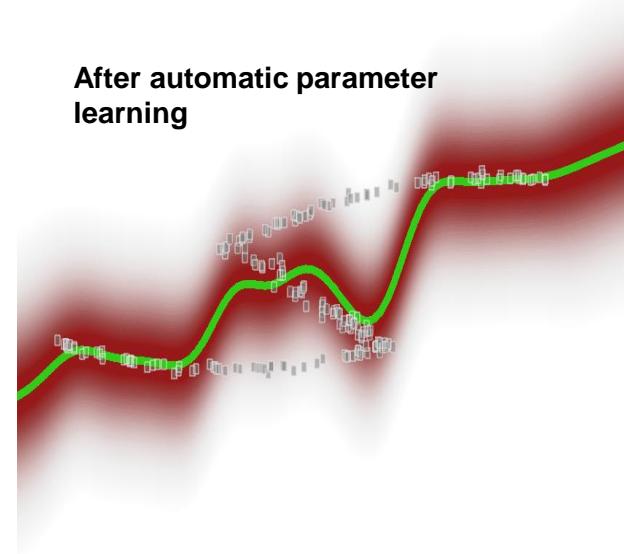
Regression forest (for varying D)



T=400, D=5

Gaussian process (for varying hyper parameters)

After automatic parameter learning



GPs do not seem able to capture the large ambiguity in the middle.

Theory

- Assuming that
 - training samples are i.i.d.
 - trees are trained independently of estimating leaf statistics
- Number of splitting functions per node are sampled from Poisson distribution
- Stopping criteria: Minimum samples per leaf $k(n)$
- Regression forests converge as the amount of training data $n \rightarrow \infty$

[M. Denil et al. **Narrowing the Gap: Random Forests In Theory and In Practice**. JMLR 2014]

Convergence

Theorem 1 Suppose that feature X is supported on \mathbb{R}^d and has a density which is non-zero almost everywhere. Moreover, suppose that $\mathbb{E}[\theta|X = x]$, which is the (unknown) mean prediction given x , is bounded and that $\mathbb{E}[\theta^2] < \infty$. Let $A(n)$ denote the training data and Z_t the randomness of a tree. Then

$$\mathbb{E}_{X, A(n), Z_t} \left[\left| \frac{1}{T} \sum_t^T \mathbb{E} [\theta | X, A(n), Z_t] - \mathbb{E} [\theta | X] \right|^2 \right] \rightarrow 0$$

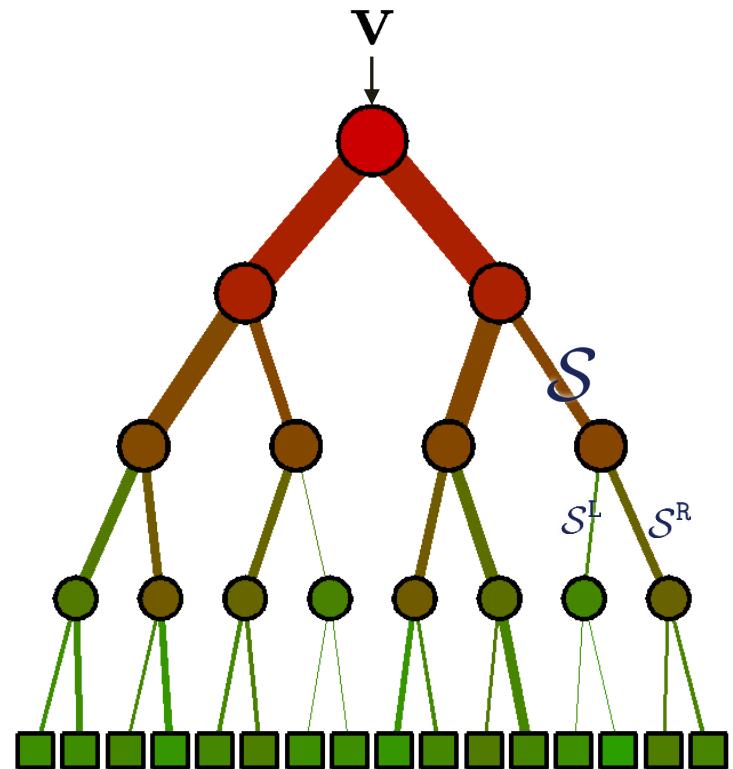
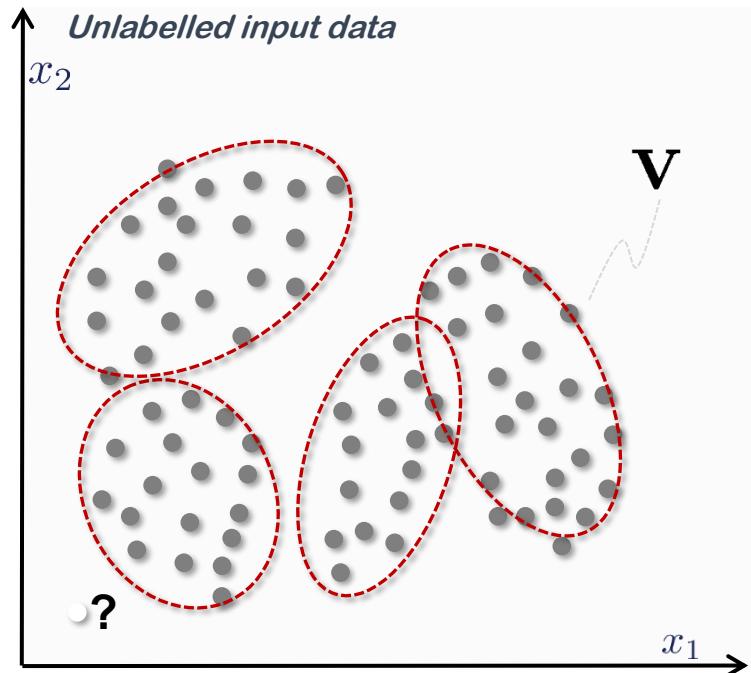
provided that $k_n \rightarrow \infty$ and $k_n/n \rightarrow 0$ as $n \rightarrow \infty$.

[M. Denil et al. **Narrowing the Gap: Random Forests In Theory and In Practice.** JMLR 2014]

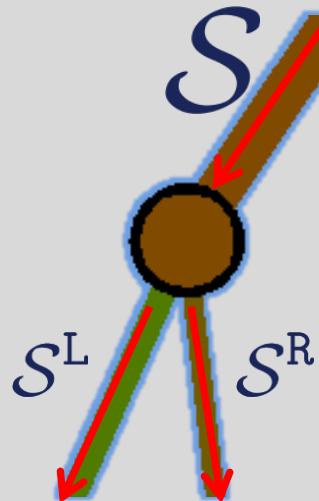
Overview

- Classification forests
- Regression forests
- Other variants

Density forest model



Training and information gain



Information gain

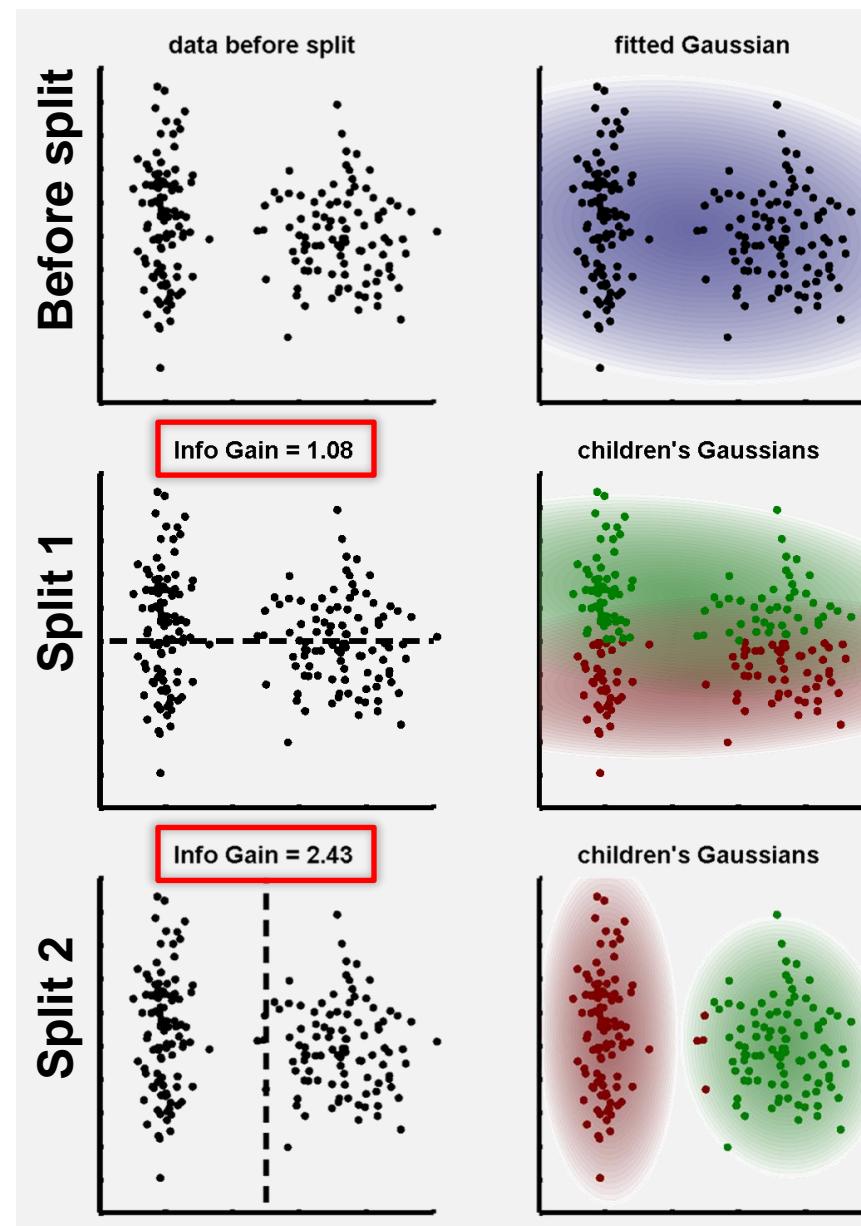
$$I = H(\mathcal{S}_j) - \sum_{i \in \{L,R\}} \frac{|\mathcal{S}_j^i|}{|\mathcal{S}|} H(\mathcal{S}_j^i)$$

Differential entropy of Gaussian

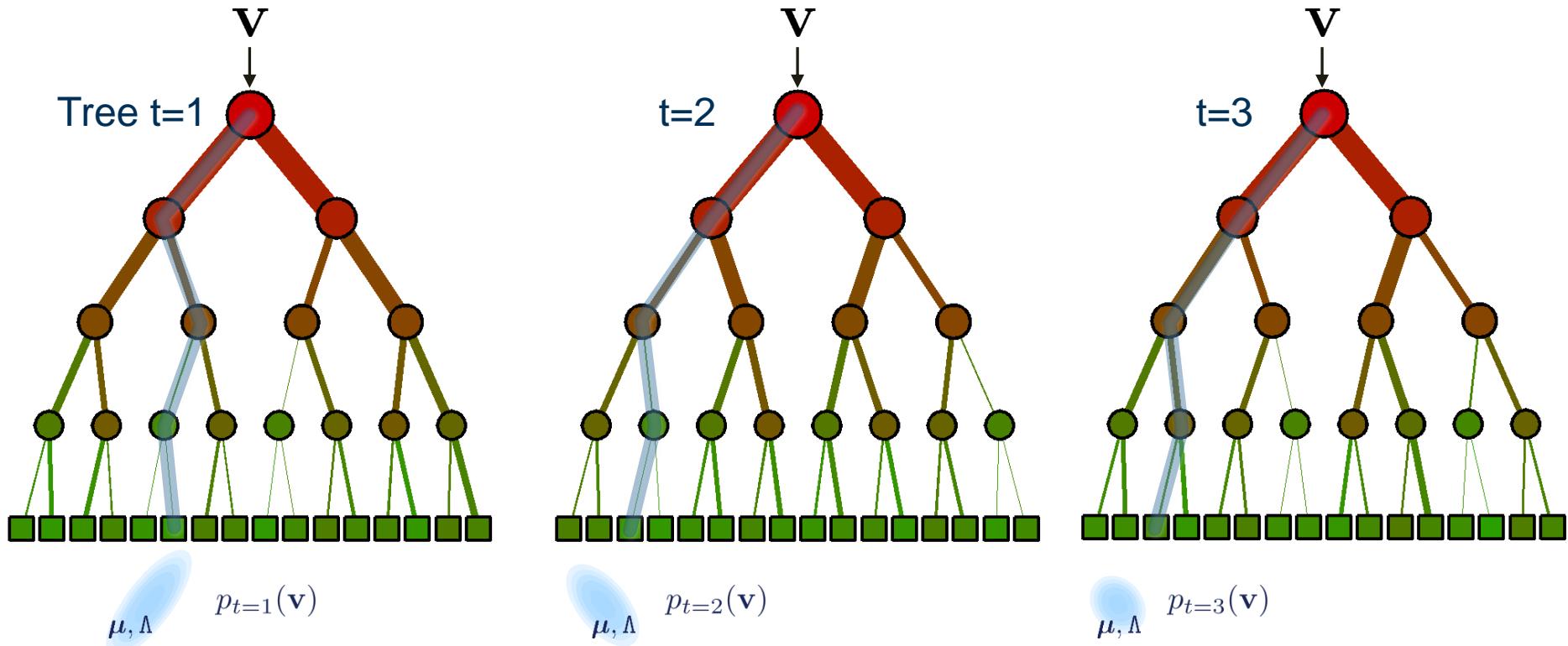
$$H(\mathcal{S}) = \frac{1}{2} \log ((2\pi e)^d |\Lambda(\mathcal{S})|)$$

Node training

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



Density forest model



The prediction and ensemble models

$$\text{Output of tree } p_t(\mathbf{v}) = \frac{\pi_{l(\mathbf{v})}}{Z_t} \mathcal{N}(\mathbf{v}; \boldsymbol{\mu}_{l(\mathbf{v})}, \boldsymbol{\Lambda}_{l(\mathbf{v})})$$

$$\text{with } \pi_{l(\mathbf{v})} = |\mathcal{S}_{l(\mathbf{v})}| / |\mathcal{S}_0|$$

Partition function

$$Z_t = \int_{\mathbf{v}} \left(\sum_l \pi_l \mathcal{N}(\mathbf{v}; \boldsymbol{\mu}_l, \boldsymbol{\Lambda}_l) p(l|\mathbf{v}) \right) d\mathbf{v}$$

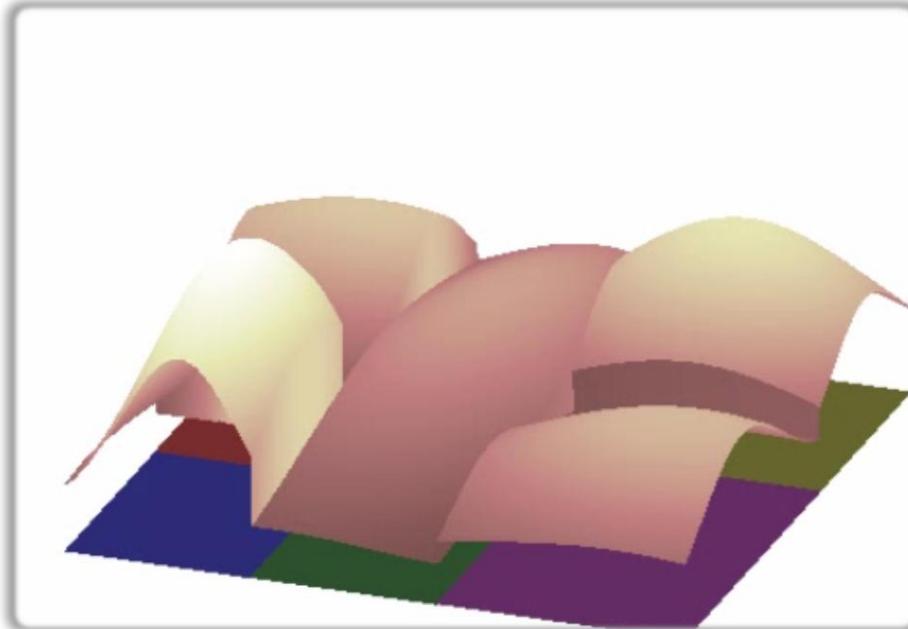
$$\text{The ensemble model } p(\mathbf{v}) = \frac{1}{T} \sum_{t=1}^T p_t(\mathbf{v})$$

$$\text{with } p(l|\mathbf{v}) = [\mathbf{v} \in l(\mathbf{v})]$$

Partition function

$$Z_t = \int_{\mathbf{v}} \pi_{l(\mathbf{v})} \mathcal{N} \left(\mathbf{v}; \boldsymbol{\mu}_{l(\mathbf{v})}, \Lambda_{l(\mathbf{v})} \right) d\mathbf{v}$$

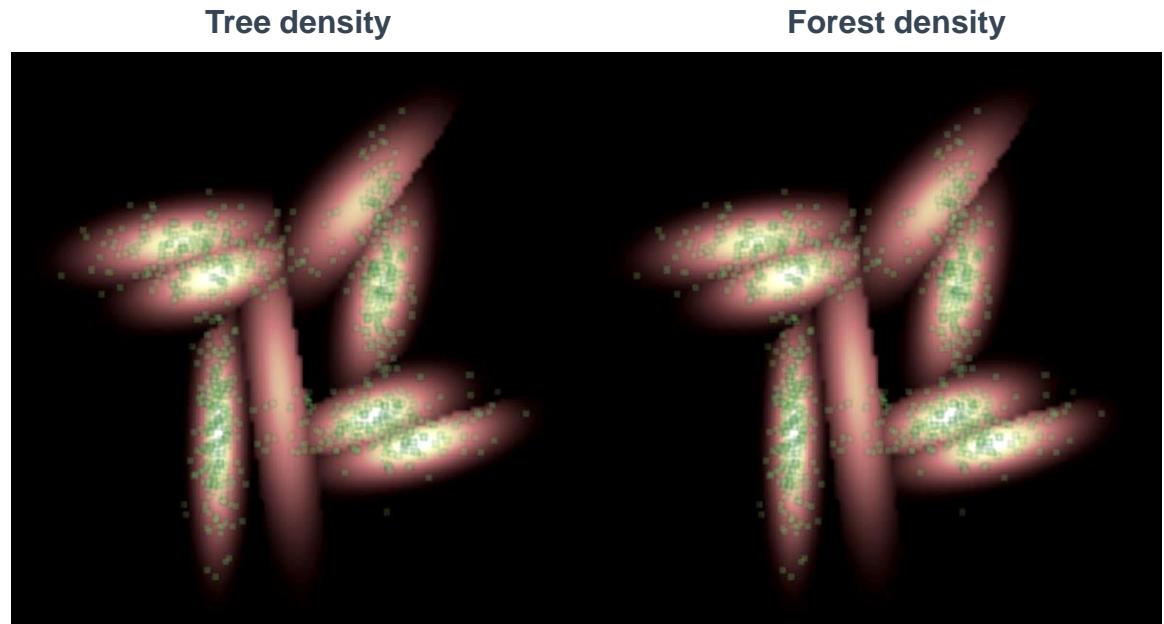
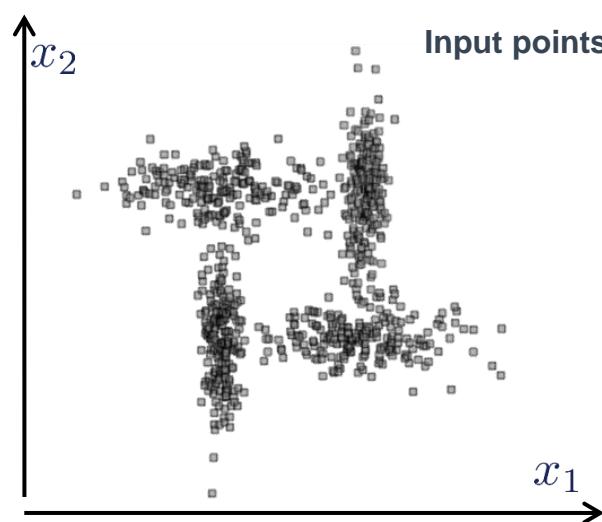
1. Compute via cumulative multivariate normal for ax.-align. w. learn.



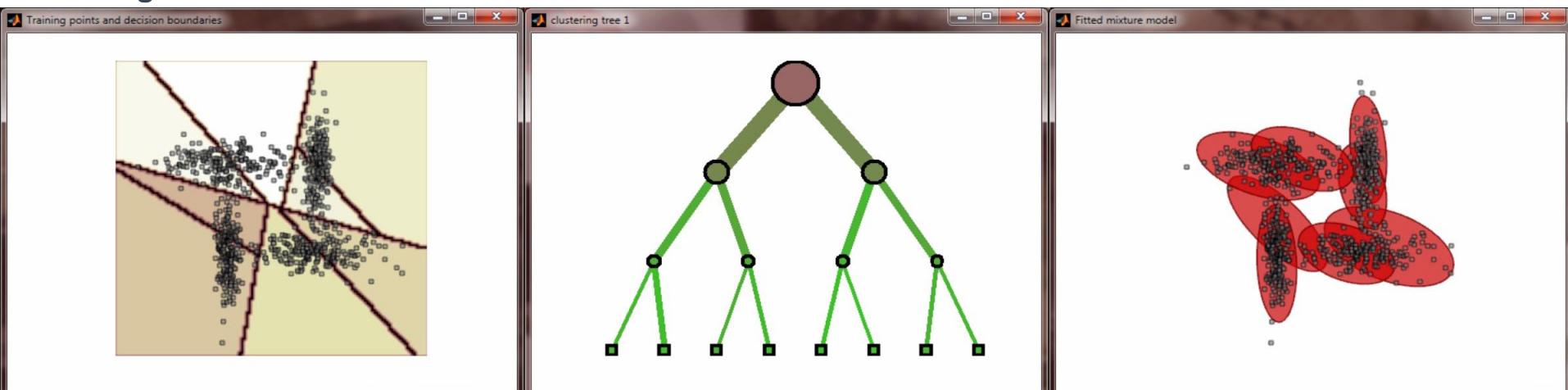
2. Grid-based numerical integration

$$Z_t \approx \Delta \sum_i \pi_{l(\mathbf{v}_i)} \mathcal{N}(\mathbf{v}_i; \boldsymbol{\mu}_{l(\mathbf{v}_i)}, \Lambda_{l(\mathbf{v}_i)})$$

Density forest: evaluation

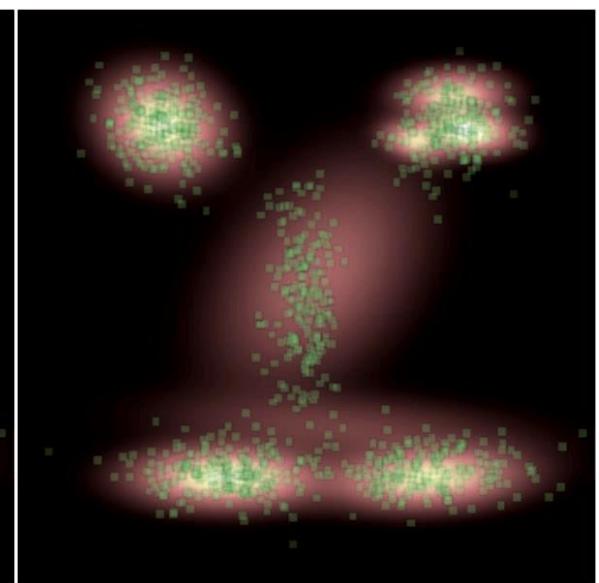
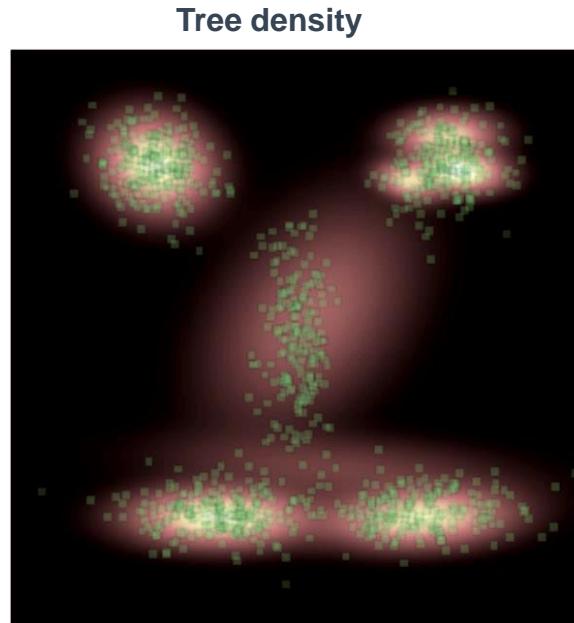
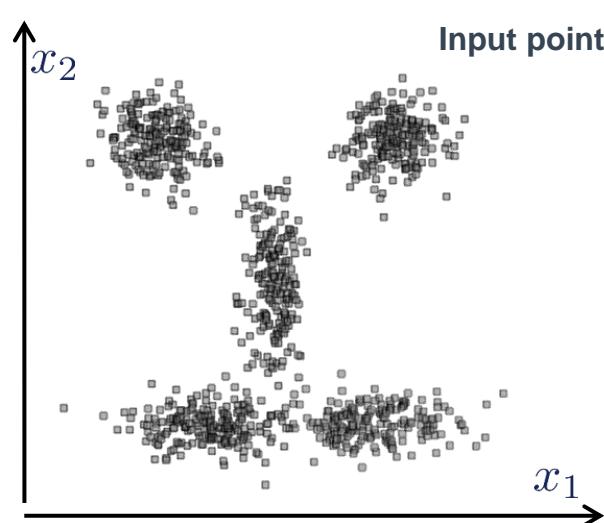


Training different trees in the forest

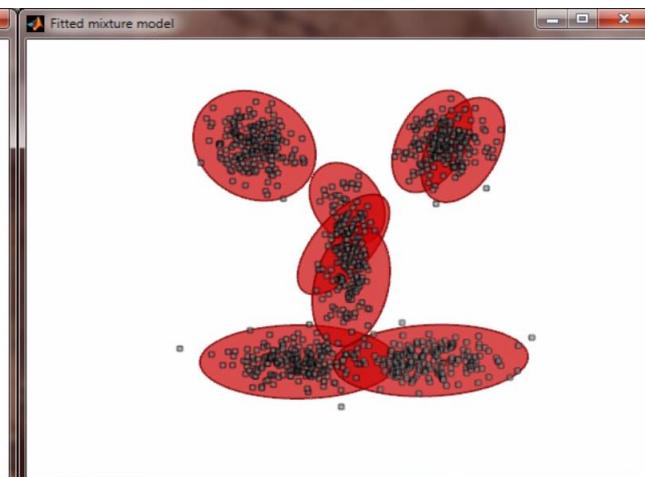
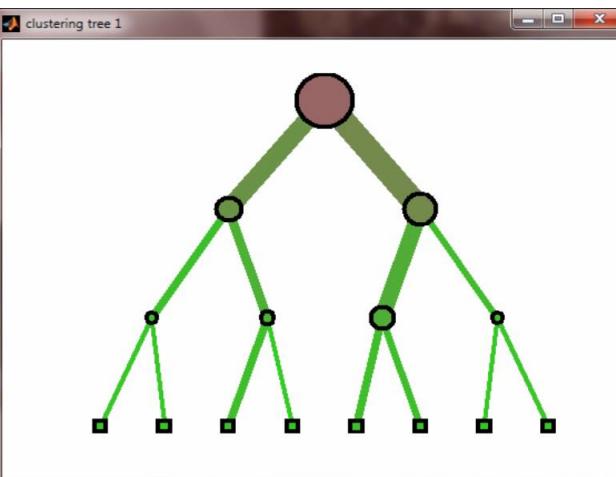
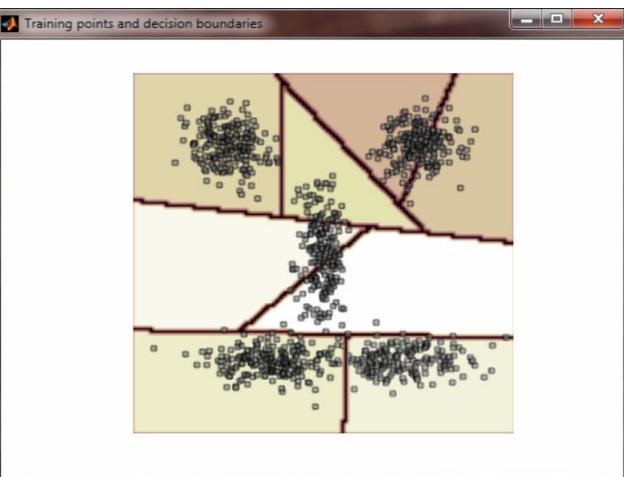


Parameters: T=400, D=4, weak learner = linear, predictor = Gaussian

Density forest: evaluation

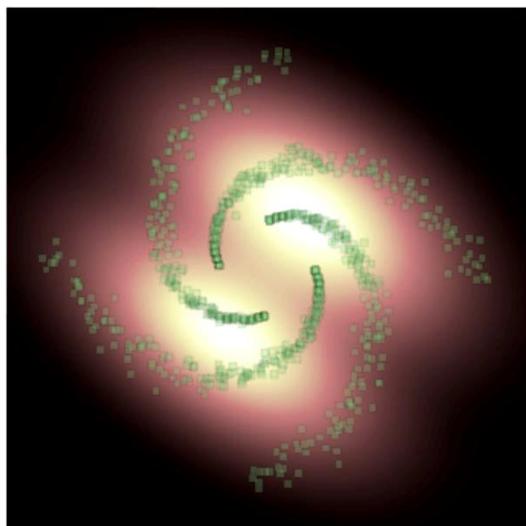
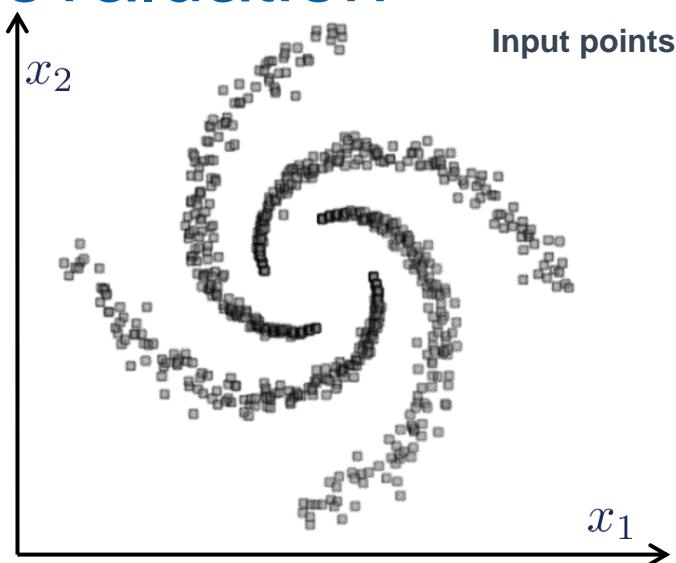


Training different trees in the forest



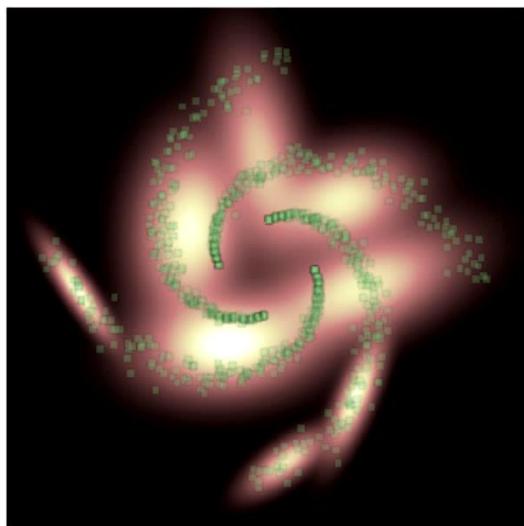
Parameters: T=400, D=4, weak learner = linear, predictor = Gaussian

Density forest: evaluation

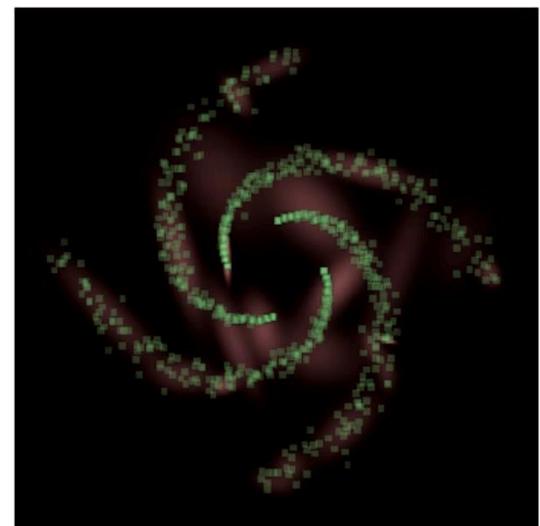


$D=2$

← underfitting



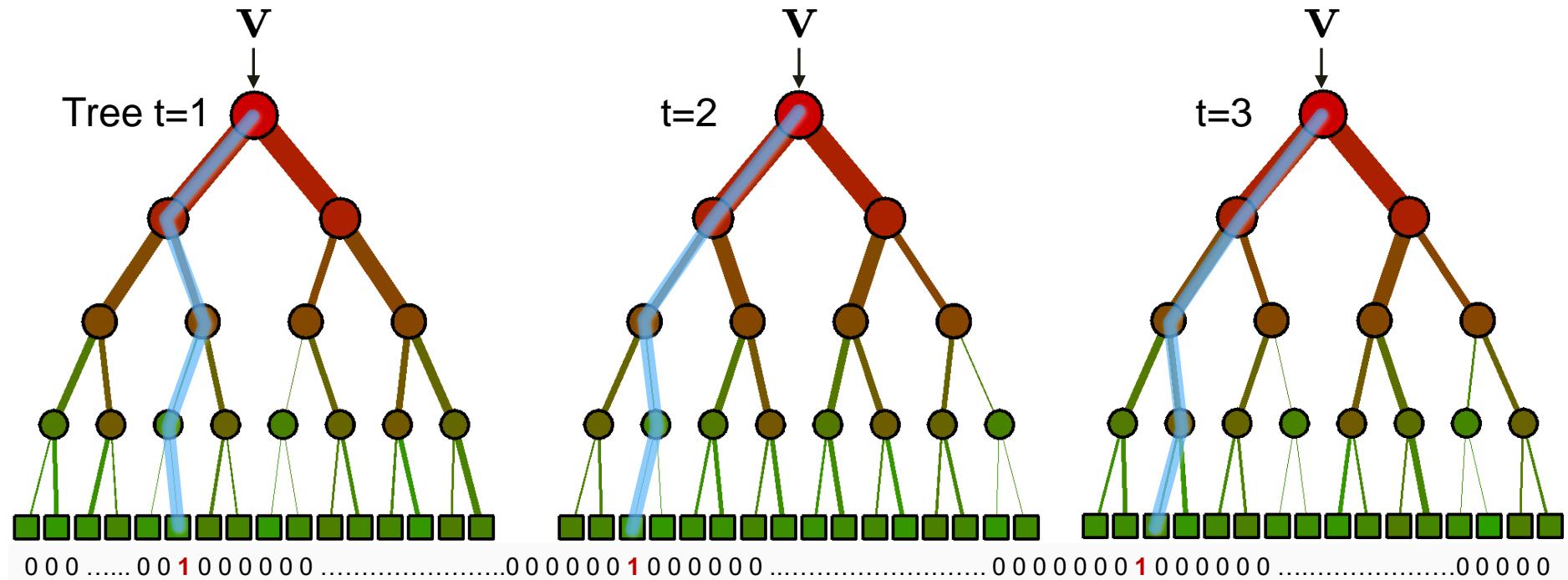
$D=4$



$D=6$

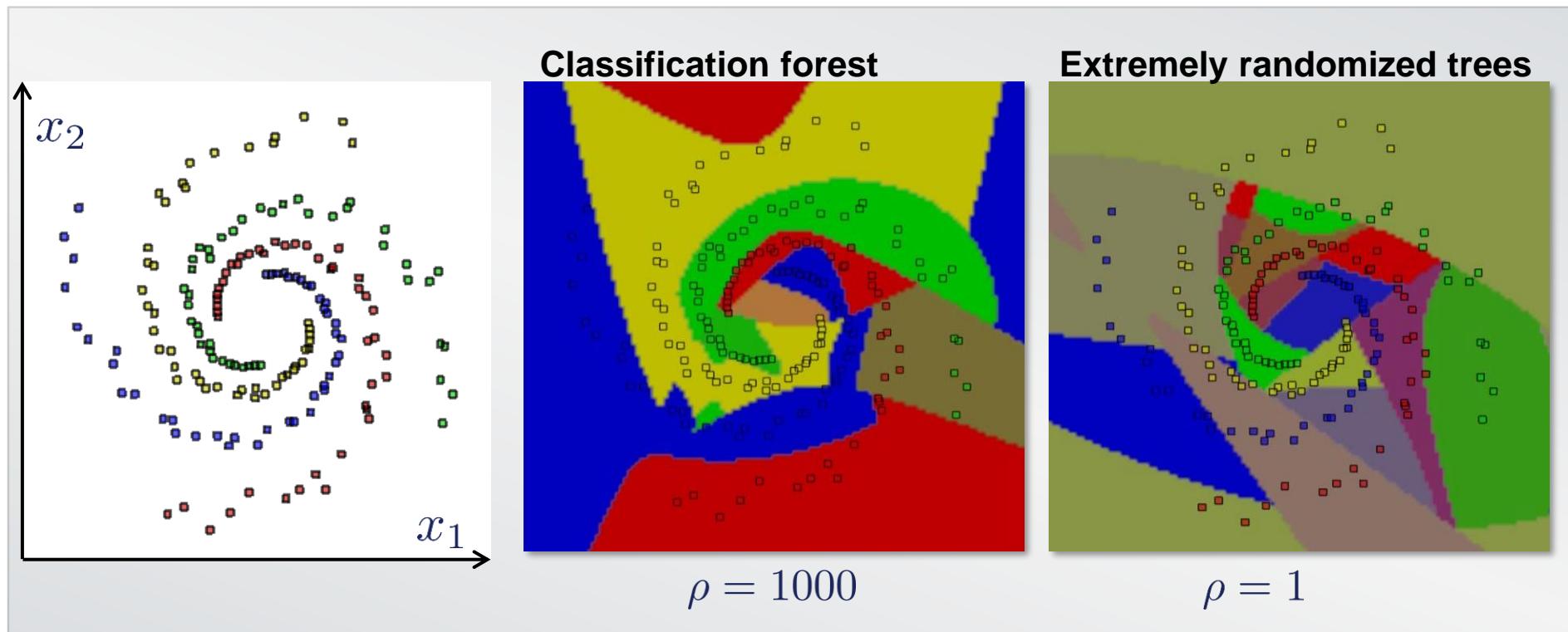
Parameters: T=400, D=4, weak learner = linear, predictor = Gaussian

Clustering trees for coding



[F. Moosmann, E. Nowak, F. Jurie. **Randomized Clustering Forests for Image Classification**. IEEE Trans. Pattern Anal. Mach. Intell. (PAMI) 30(9): 1632-1646 (2008)]

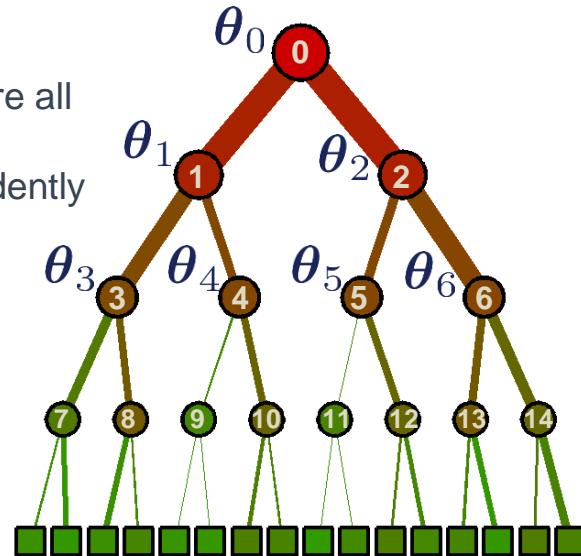
Extremely randomized trees (ERT)



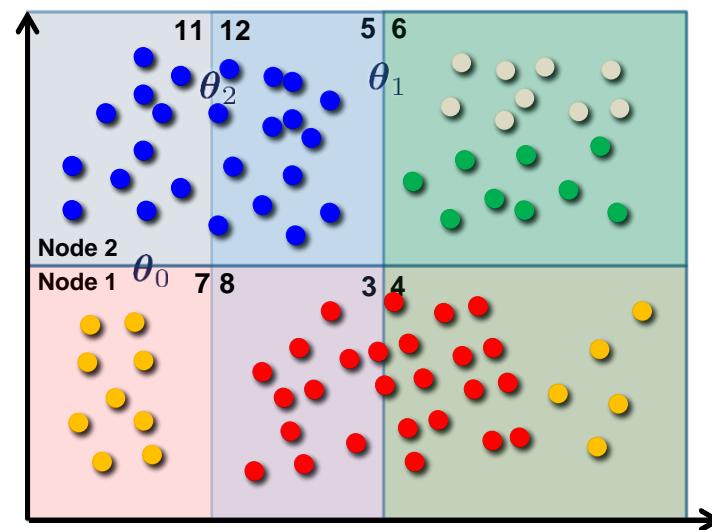
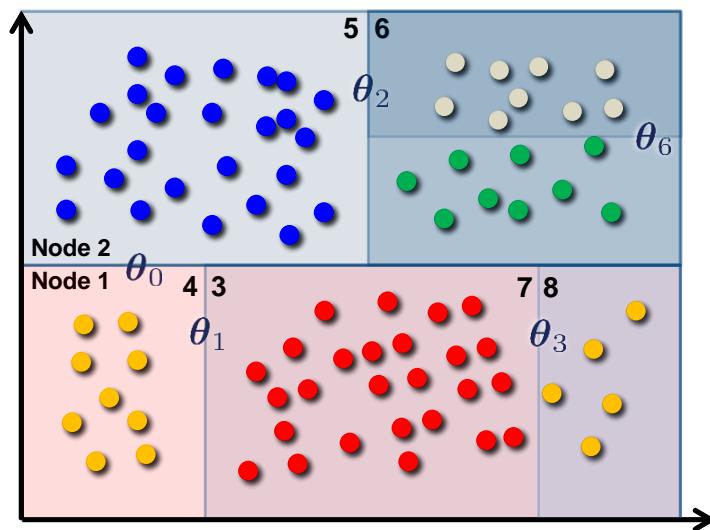
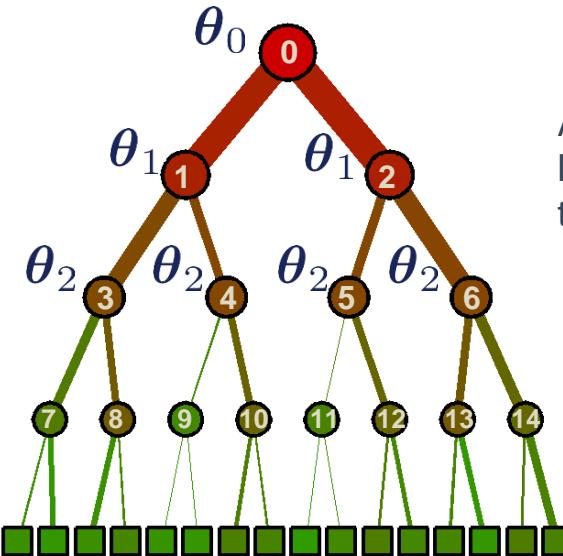
- ERTs take longer to converge (more trees) and under-confident.

Forests v. ferns

Nodes are all trained independently

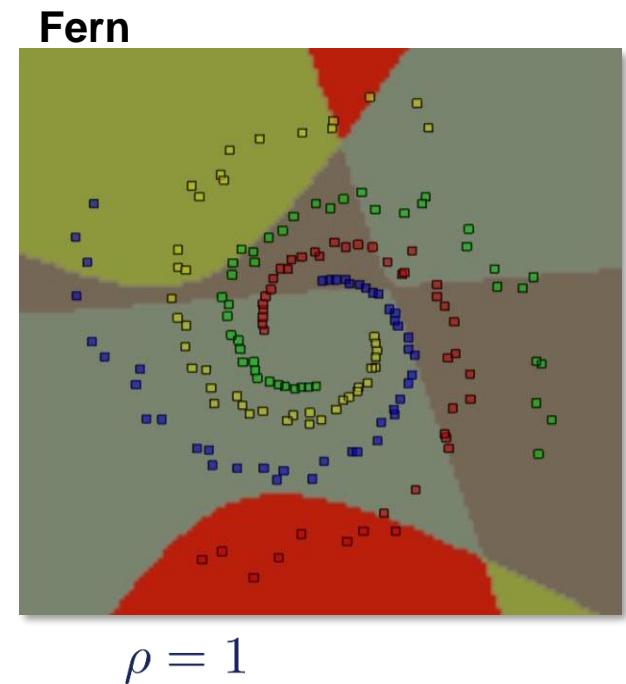
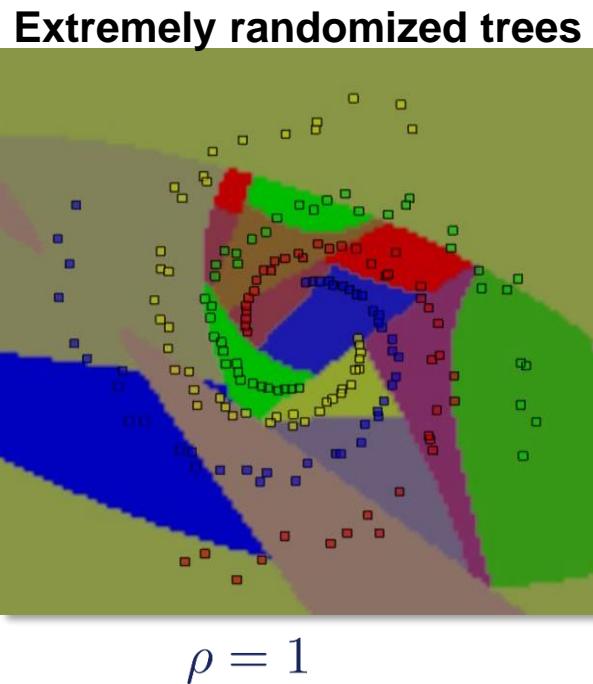
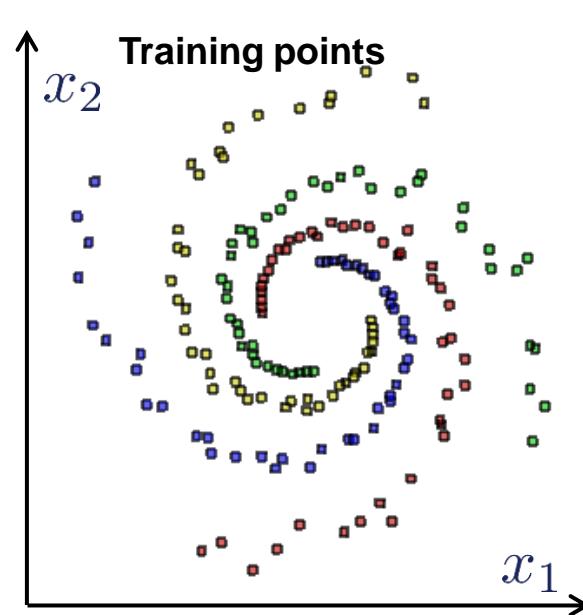


All nodes in the same level l are assigned the same test θ_l



Ferns are simpler but less rich and less flexible.
By level 2 the tree has done more correct data splitting than the fern.

Forests v. ferns



- Ferns take longer to converge (more trees) and under-confident.
- Ferns and ERTs have less parameters than forests and thus may over-fit less.

[M. Ozuysal, M. Calonder, V. Lepetit, and P. Fua.
Fast keypoint recognition using random ferns.
IEEE Trans. PAMI 32(3), 2001]

UNIVERSITÄT BONN

