

Decision Trees and Boosting

Lecture 22 - CS 689, Spring 2023

- This week

- Decision Trees

- Ensemble (additive) decision trees

- Today: optimization with boosting

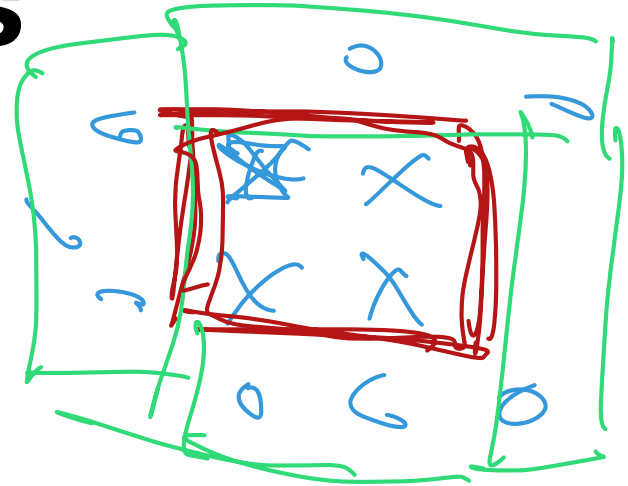
- Thursday: Bayesian learning with MCMC

- Discrete parameters; no latent representations

Axis-parallel regions

- Feature thresholds

$$1\{x_j < t\}$$



$$T_j = \text{unique values of } x_j \quad \forall x_i \in D^T$$

- Conjunctive form

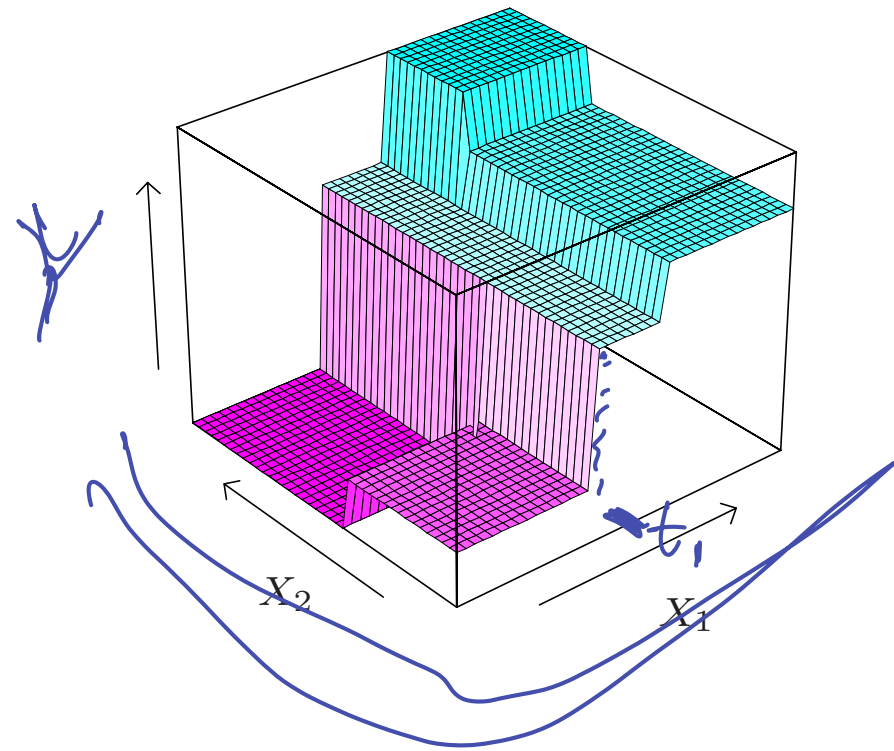
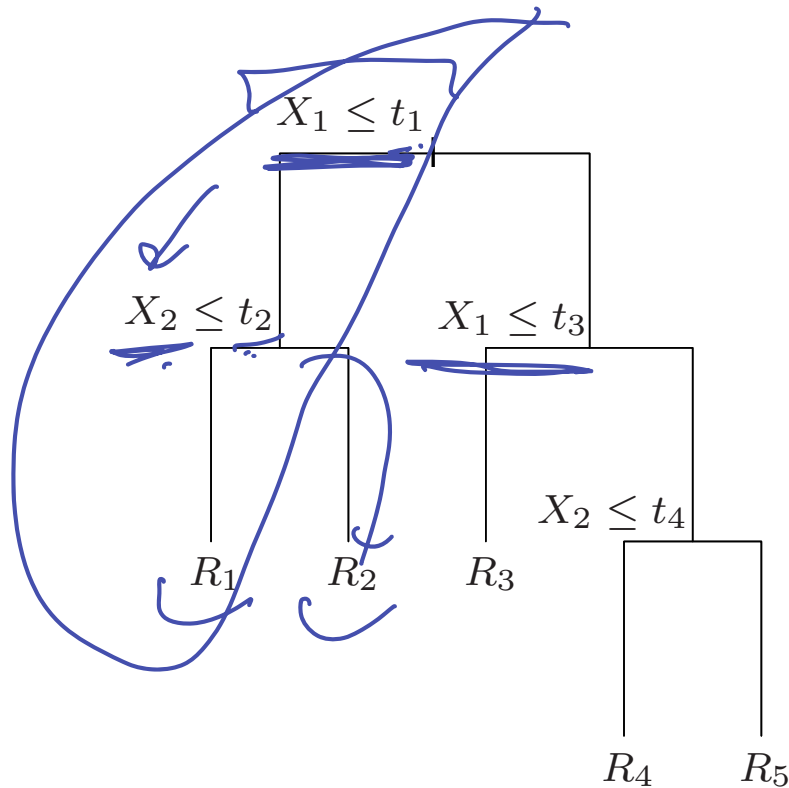
$$1\{x_1 < \underline{t}\} 1\{x_1 > \underline{s}\} 1\{x_2 < \underline{u}\} 1\{x_2 > \underline{v}\}$$

New feature

Learning? Intractable ☹️

Recursive splitting

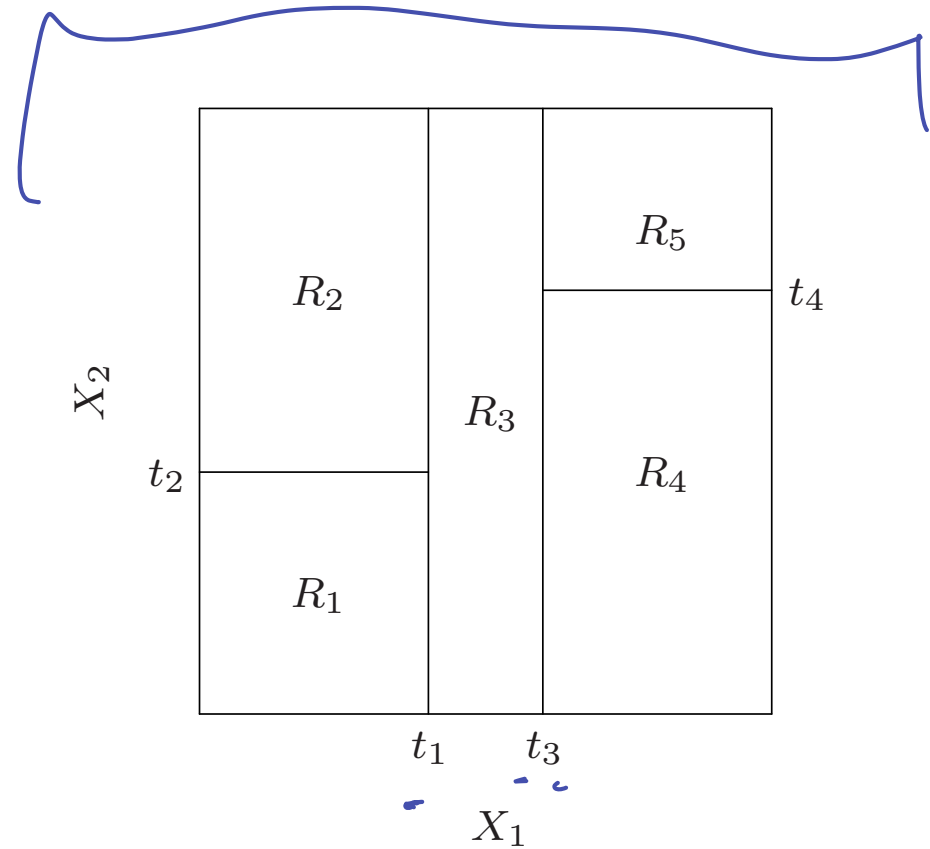
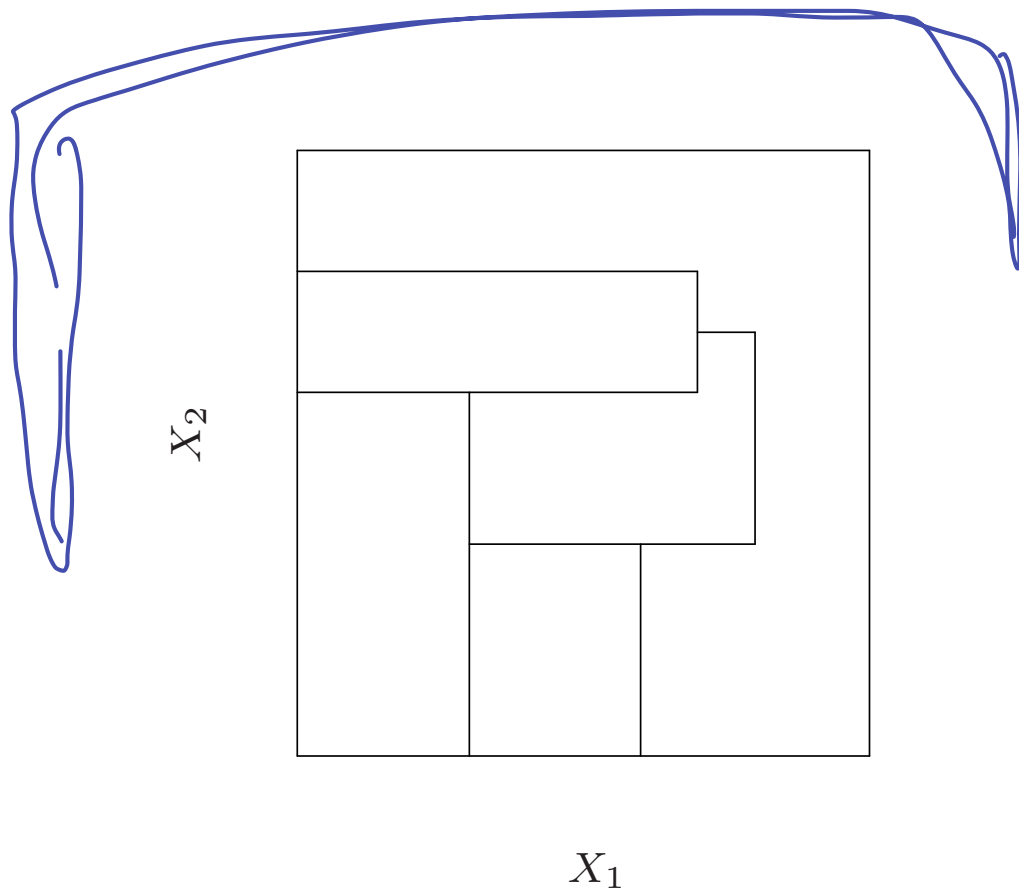
$$R_1 \Leftrightarrow (x_1 \leq t_1) \wedge (x_2 \leq t_2)$$



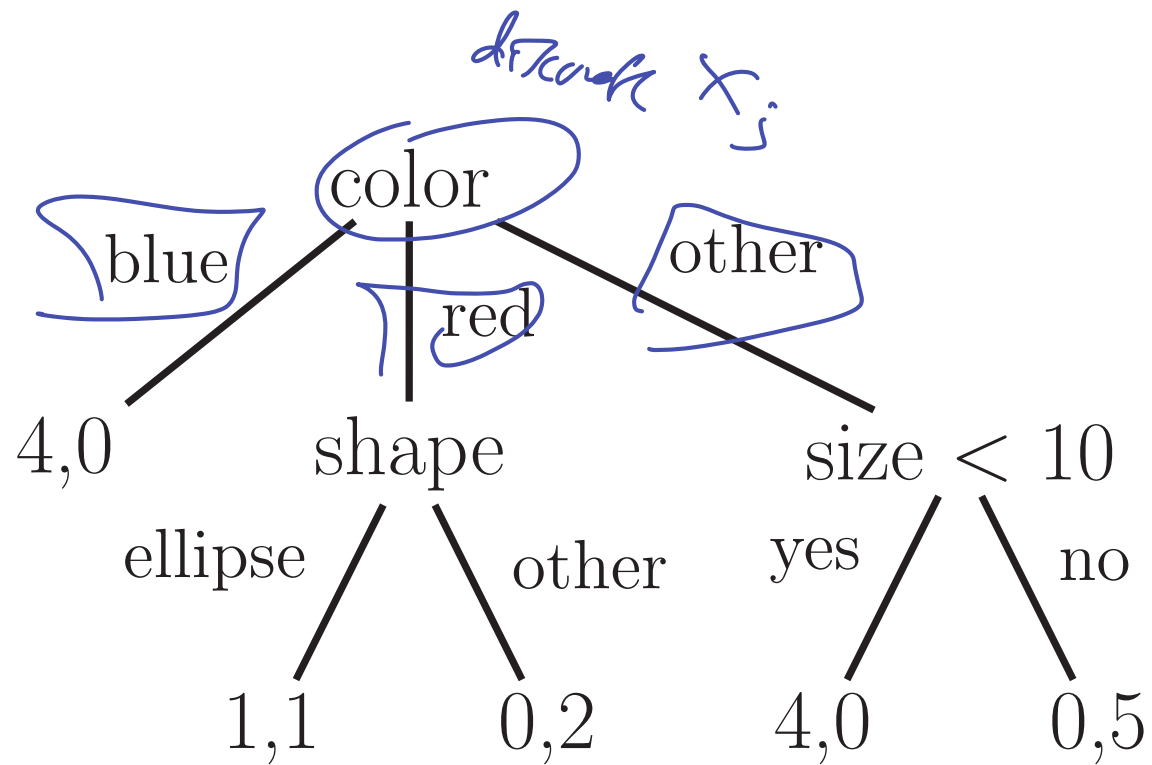
$\hookrightarrow M_4 \hookrightarrow M_5$

- From Hastie et al. 2009: *Elements of Statistical Learning*

list representable
by DTree



- From Hastie et al. 2009: *Elements of Statistical Learning*



Greedy decision tree learning

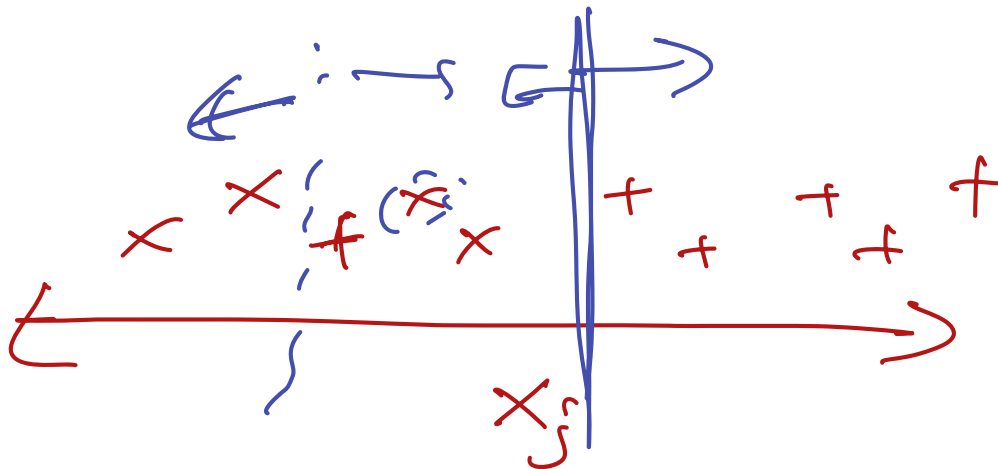
split(D): feature j , thresholds T_j

$$(j^*, t^*) = \arg \min_{j \in \{1, \dots, D\}} \min_{t \in T_j} \text{cost}(\{\mathbf{x}_i, y_i : x_{ij} \leq t\}) + \text{cost}(\{\mathbf{x}_i, y_i : x_{ij} > t\})$$

Handwritten notes:
 - Blue arrows point to j^* and t^* .
 - Red arrow points from "R-squared MSE" to the cost function.
 - Red bracket under $\{ \mathbf{x}_i, y_i : x_{ij} \leq t \}$ is labeled "dim."
 - Red bracket under $\{ \mathbf{x}_i, y_i : x_{ij} > t \}$ is labeled "dim."

MSE :

$$\sum_{i \in D_L} (y_i - \bar{y}_{D_L})^2 + \sum_{i \in D_R} (y_i - \bar{y}_{D_R})^2$$

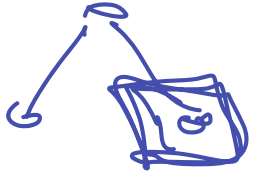


Greedy decision tree learning

split(D): feature j , thresholds T_j

$$(j^*, t^*) = \arg \min_{j \in \{1, \dots, D\}} \min_{t \in T_j} \text{cost}(\{\mathbf{x}_i, y_i : x_{ij} \leq t\}) + \text{cost}(\{\mathbf{x}_i, y_i : x_{ij} > t\})$$

Algorithm 16.1: Recursive procedure to grow a classification/ regression tree



```
1 function fitTree(node,  $\mathcal{D}$ , depth) ;
2   node.prediction = mean( $y_i : i \in \mathcal{D}$ ) // or class label distribution ;
3    $(j^*, t^*, \mathcal{D}_L, \mathcal{D}_R) = \text{split}(\mathcal{D})$ ;
4   if not worthSplitting(depth, cost,  $\mathcal{D}_L, \mathcal{D}_R$ ) then
5     return node
6   else
7     node.test =  $\lambda \mathbf{x}. x_{j^*} < t^*$  // anonymous function;
8     node.left = fitTree(node,  $\mathcal{D}_L$ , depth+1);
9     node.right = fitTree(node,  $\mathcal{D}_R$ , depth+1);
10    return node;
```

max depth

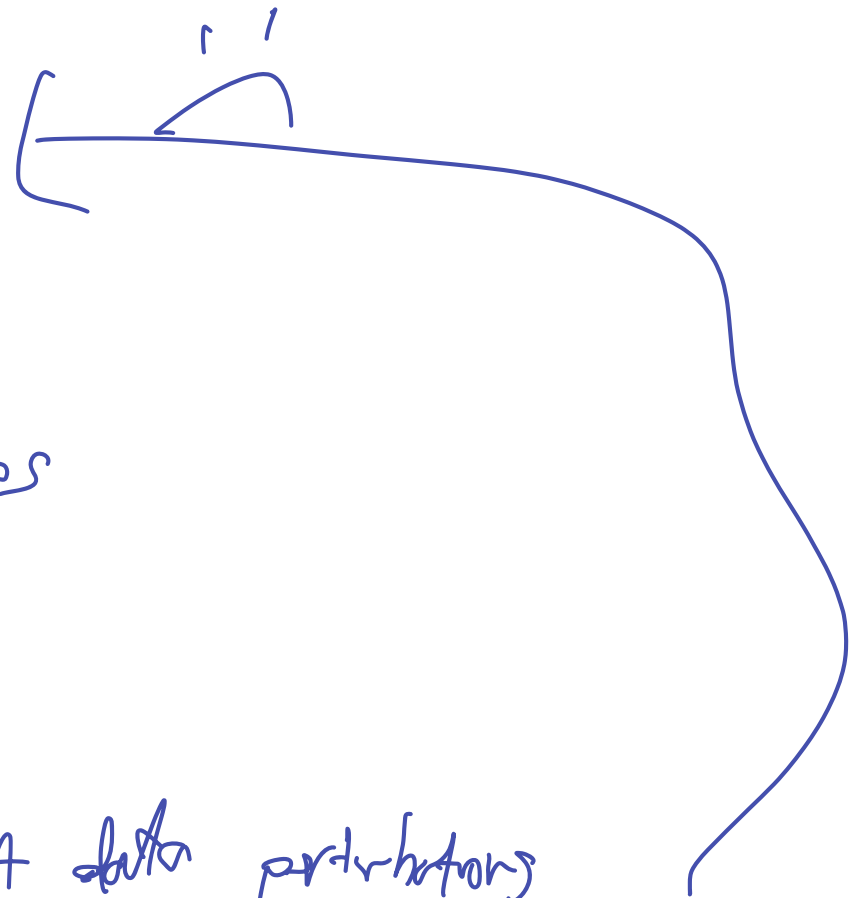
Other details

- "Not worth splitting": overfitting control
 - Min # examples in node
 - Δ improvement to loss
- Max depth
- Pruning
- CART, C4.5, ...

Issues for single dtree learning

Pro

- Nonlinear
- Simple
- Interpretable (??)



Con

- Linear & gradual relationships
- Not axis-aligned
- Combustion?? $\gamma = 5/N_0$
- High variance / unstable wrt data perturbations

Generalized additive models

$$f(x) = \alpha + \sum_{m=1}^M f_m(x)$$

↑
Basic nonlinear model

Ensemble models

Random forests

ensemble of trees

M = num trees you want

for $m=1 \dots M$:

- Draw bootstrap sample $D^{(m)}$ of D :

$i=1 \dots N$: $D_i^{(m)} \sim D^{\text{Train}}$

- Choose random subset of features $J^{(m)}$

- Fit tree f_m from $D^{(m)}, J^{(m)}$

Predict:
$$f(x) = \frac{1}{M} \sum_{m=1}^M f_m(x)$$

Idea:

Average over
Tree instability

Works well!

Boosting

→ "weak learner"

- Use optimization ideas to formulate/explain iterative learning of successive decision trees
- Basic idea: learn to predict residuals from sum of previous trees

$$y \Rightarrow (y - \hat{y}) = r$$

Forward Stagewise Learning

Greedy learning for multiple dtrees

DTree

Split & preds



Init: $f_0(x) = \arg\min_{\gamma} \sum_i L(y_i, \phi(x_i; \gamma))$

↑
CART or greedy rect. splitter

for $m=1 \dots M$

$(\beta_m, \gamma_m) = \arg\min_{\beta, \gamma} \sum_i L(y_i, f_{m-1}(x_i) + \beta \phi(x_i; \gamma))$

DTree



$f_m(x) = f_{m-1}(x) + \beta_m \phi(x; \gamma_m)$

Better: $f_m(x) = f_{m-1}(x) + \bigvee_{f(0,1)} \beta_m \phi(x; \gamma_m)$

$f(0,1)$ step size!

L2Boosting

- Apply forward stagewise learning to squared error loss.

Gradient Boosting

Algorithm 16.4: Gradient boosting

- 1 Initialize $f_0(\mathbf{x}) = \operatorname{argmin}_{\gamma} \sum_{i=1}^N L(y_i, \phi(\mathbf{x}_i; \gamma))$;
 - 2 **for** $m = 1 : M$ **do**
 - 3 Compute the gradient residual using $r_{im} = - \left[\frac{\partial L(y_i, f(\mathbf{x}_i))}{\partial f(\mathbf{x}_i)} \right]_{f(\mathbf{x}_i)=f_{m-1}(\mathbf{x}_i)}$;
 - 4 Use the weak learner to compute γ_m which minimizes $\sum_{i=1}^N (r_{im} - \phi(\mathbf{x}_i; \gamma_m))^2$;
 - 5 Update $f_m(\mathbf{x}) = f_{m-1}(\mathbf{x}) + \nu \phi(\mathbf{x}; \gamma_m)$;
 - 6 Return $f(\mathbf{x}) = f_M(\mathbf{x})$
-