

- Support Vector Machines (VAPNIK, 1995)
 - Very good classifier
 - Can be adapted to ranking and multiclass problems
- Neural Nets
 - RankNet (BURGES et al., 2006)
- Tree Ensembles
 - Random Forests (BREIMAN and SCHAPIRE, 2001)
 - Boosted Decision Trees
 - Multiple Additive Regression Trees (FRIEDMAN, 1999)
 - LambdaMART (BURGES, 2010)
 - Used by AltaVista, Yahoo!, Bing, Yandex, ...

All top teams of the *Yahoo! Learning to Rank Challenge (2010)* used combinations of Tree Ensembles!

Yahoo! Learning to Rank Challenge

- Yahoo! Webscope dataset (CHAPELLE and CHANG, 2011):
36,251 queries, 883 k documents, 700 features, 5 ranking levels
 - set-1:
 - 473,134 feature vectors
 - 519 features
 - 19,944 queries
 - set-2:
 - 34,815 feature vectors
 - 596 features
 - 1,266 queries
- Winner used a combination of 12 models:
 - 8 Tree Ensembles (LambdaMART)
 - 2 Tree Ensembles (Additive Regression Trees)
 - 2 Neural Nets

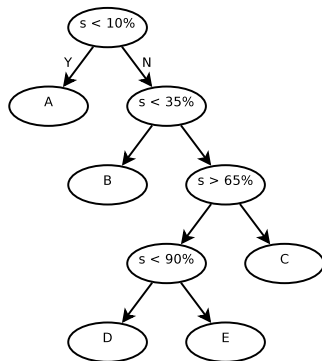
Decision Trees

Characteristics of a tree:

- Graph based model
- Consists of a root, nodes, and leaves

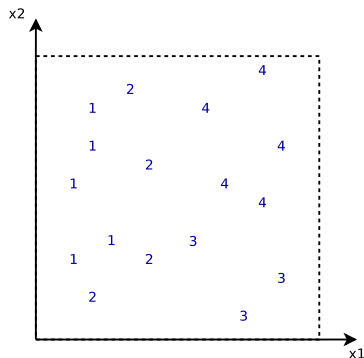
Advantages:

- Simple to understand and interpret
- *White box* model
- Can be combined with other techniques



Decision trees are basic learners for machine learning, e.g. *classification* or *regression trees*.

Learning a Regression Tree (I)

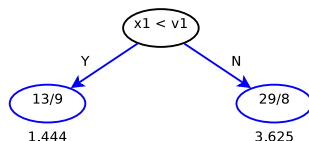
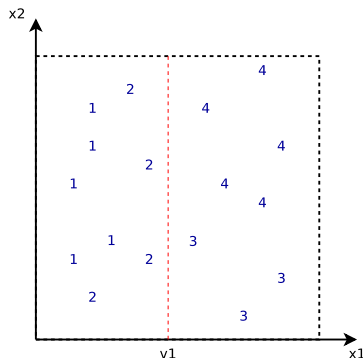


What is the prediction of this empty tree?
Mean value of all training instances

Consider a 2-dimensional space consisting of data points of the indicated values. We start with an empty root node (blue).

Learning a Regression Tree (II)

How to determine split threshold v_1 ?
How to find the prediction for each leaf?

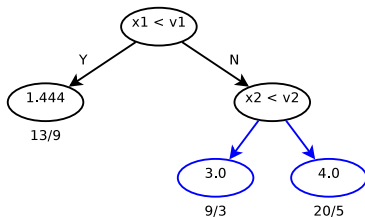
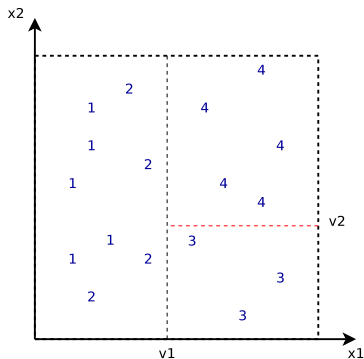


Use mean value of training instances
on left leaf as prediction:
 $(5*1+4*2)/9=1.444$
For right leaf, $(4*5+3*3)/8=3.625$.

Try all possible v_1 thresholds.

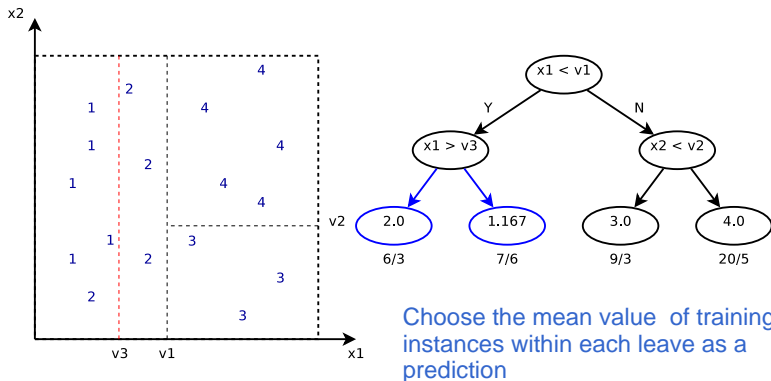
The algorithm searches for split variables and split points, x_1 and v_1 , that predict values minimizing the predicted error, e.g. $\sum (y_i - f(x_i))^2$.

Learning a Regression Tree (III)



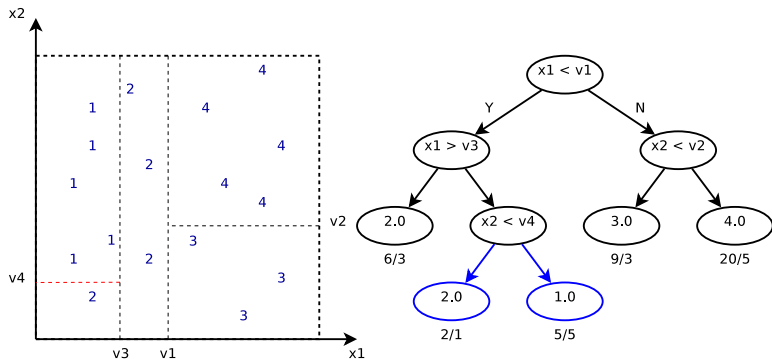
Here we examine the right side first: find a split variable and a split value that minimize the predicted error, i.e. x_2 and v_2 .

Learning a Regression Tree (IV)



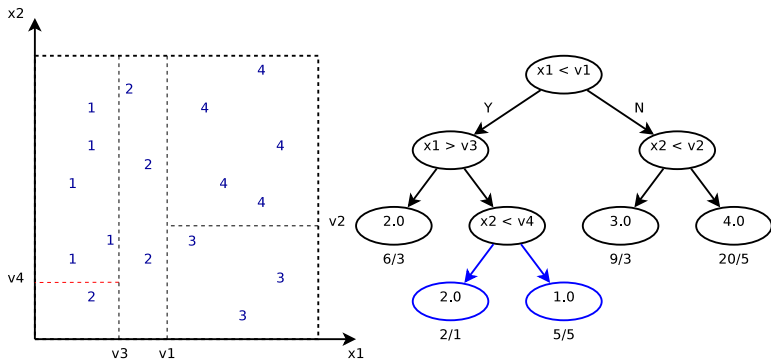
Now to the left side: Again, find a split variable and a split value that minimize the predicted error, i.e. x_1 and v_3 .

Learning a Regression Tree (V)



Once again, find a split variable and a split value that minimize the predicted error, here x_2 and v_4 .

Learning a Regression Tree (V)



Once again, find a split variable and a split value that minimize the predicted error, here x_2 and v_4 . The tree perfectly fits the data! Problem?

Formal Definition of a Decision Tree

A decision tree partitions the parameter space into disjoint regions R_k , $k \in \{1, \dots, K\}$, $K = \text{number of leaves}$. Formally, the regression model (1) predicts a value using a constant γ_k for each region R_k :

Prediction:

$$T(\mathbf{x}; \Theta) = \sum_{k=1}^K \gamma_k 1(\mathbf{x} \in R_k) \quad (1)$$

$\Theta = \{R_k, \gamma_k\}_1^K$ describes the model parameters, $1(\cdot)$ is the *characteristic function* (1 if argument is true, 0 otherwise), and $\hat{\gamma}_k = \text{mean}(y_i | \mathbf{x}_i \in R_k)$. Optimal parameters $\hat{\Theta}$ are found minimizing the empirical risk:

For all instances

Loss function

$$\hat{\Theta} = \arg \min_{\Theta} \sum_{k=1}^K \sum_{\mathbf{x}_i \in R_k} L(y_i, \gamma_k) \quad (2)$$

The combinatorial optimization problem (2) is usually split into two parts:
(i) *finding* R_k and (ii) *finding* γ_k given R_k .

Idea

Combine multiple weak learners to build a strong learner.

A weak learner is a learner with an error rate slightly better than random guessing. A strong learner is a learner with high accuracy.

Approach:

- Apply a weak learner to iteratively modified data
- Generate a sequence of learners
- For classification tasks: use majority vote
- For regression tasks: build weighted values

Loss function is sum of squared error for all instances

Function Estimation

Find a function $F^*(\mathbf{x})$ that maps \mathbf{x} to y , s.t. the expected value of some loss function $L(y, F(\mathbf{x}))$ is minimized:

$$F^*(\mathbf{x}) = \arg \min_{F(\mathbf{x})} \mathbb{E}_{y, \mathbf{x}} [L(y, F(\mathbf{x}))]$$

Boosting approximates $F^*(\mathbf{x})$ by an additive expansion

$$F(\mathbf{x}) = \sum_{m=1}^M \beta_m h(\mathbf{x}; \mathbf{a}_m)$$

What are parameters \mathbf{a} ?
In regression trees,
they represent how
each tree is built and
leave prediction is made.

where $h(\mathbf{x}; \mathbf{a})$ are simple functions of \mathbf{x} with parameters $\mathbf{a} = \{a_1, a_2, \dots, a_n\}$ defining the function h , and β are expansion coefficients.

For boosting trees, $F(\mathbf{x})$ is sum of tree classifiers.
 M is the number of trees.

Expansion coefficient typically includes a learning rate factor

Finding Parameters

Expansion coefficients $\{\beta_m\}_0^M$ and the function parameters $\{\mathbf{a}_m\}_0^M$ are iteratively fit to the training data:

- 1 Set $F_0(\mathbf{x})$ to initial guess
- 2 For each $m = 1, 2, \dots, M$

Minimize overall
loss for all N training instances.
 M is # of trees

$$(\beta_m, \mathbf{a}_m) = \arg \min_{\beta, \mathbf{a}} \sum_{i=1}^N L(y_i, F_{m-1}(\mathbf{x}_i) + \beta h(\mathbf{x}_i, \mathbf{a})) \quad (3)$$

and

$$F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \beta_m h(\mathbf{x}; \mathbf{a}_m) \quad (4)$$

For each iteration, we only determine one
tree classifier $h()$.
How to minimize loss for one tree?
Try to find a function close to residual $y_i - F$

Gradient Boosting

Gradient boosting approximately solves (3) for differentiable loss functions:

- 1 Fit the function $h(\mathbf{x}; \mathbf{a})$ by least squares

$$\mathbf{a}_m = \arg \min_{\mathbf{a}} \sum_{i=1}^N [\tilde{y}_{im} - h(\mathbf{x}_i, \mathbf{a})]^2 \quad (5)$$

to the “pseudo”-residuals

$$\tilde{y}_{im} = - \left[\frac{\partial L(y_i, F(\mathbf{x}_i))}{\partial F(\mathbf{x}_i)} \right]_{F(\mathbf{x})=F_{m-1}(\mathbf{x})} \quad (6)$$

- 2 Given $h(\mathbf{x}; \mathbf{a}_m)$, the β_m are Residual is $y_i - F$

$$\beta_m = \arg \min_{\beta} \sum_{i=1}^N L(y_i, F_{m-1}(\mathbf{x}_i) + \beta h(\mathbf{x}_i; \mathbf{a}_m)) \quad (7)$$

\Rightarrow Gradient boosting simplifies the problem to least squares (5).

Gradient Tree Boosting

Gradient tree boosting applies this approach on functions $h(\mathbf{x}; \mathbf{a})$ representing K -terminal node regression trees.

$$h(\mathbf{x}; \{R_{km}\}_1^K) = \sum_{k=1}^K \bar{y}_{km} 1(\mathbf{x} \in R_{km}) \quad (8)$$

With $\bar{y}_{km} = \text{mean}_{\mathbf{x}_i \in R_{km}}(\tilde{y}_{im})$ the tree (8) predicts a constant value \bar{y}_{km} in region R_{km} . Equation (7) becomes a prediction of a γ_{km} for each R_{km} :

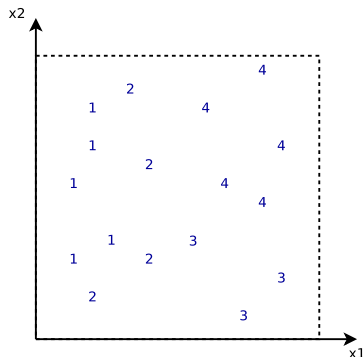
$$\gamma_{km} = \arg \min_{\gamma} \sum_{\mathbf{x}_i \in R_{km}} L(y_i, F_{m-1}(\mathbf{x}_i) + \gamma) \quad (9)$$

The approximation for F in stage m is then:

$$F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \eta \cdot \gamma_{km} 1(\mathbf{x}_i \in R_{km}) \quad (10)$$

The parameter η controls the *learning rate* of the procedure.

Learning Boosted Regression Trees (I)



$$F_0(x)$$
$$\textcircled{2.471}$$
$$\uparrow$$
$$5[1-x]^2 + 4[2-x]^2 + 3[3-x]^2 + 5[4-x]^2$$
$$\Rightarrow x=2.471$$

Mean value minimizes
the above expression

First, learn the most simple predictor that predicts a constant value minimizing the error for all training data.

Calculating Optimal Leaf Value for F_0

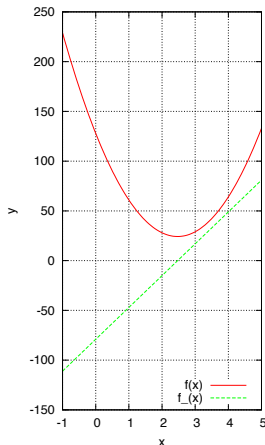
Recall the exp. coefficient: $\gamma_{km} = \arg \min_{\gamma} \sum_{\mathbf{x}_i \in R_{km}} L(y_i, F_{m-1}(\mathbf{x}_i) + \gamma)$

- Quadratic loss for the leaf (red):

$$f(x) = 5 \cdot (1 - x)^2 + 4 \cdot (2 - x)^2 \\ + 3 \cdot (3 - x)^2 + 5 \cdot (4 - x)^2$$

- $f(x)$ is quadratic, *convex*
 \Rightarrow Optimum at $f'(x) = 0$ (green)

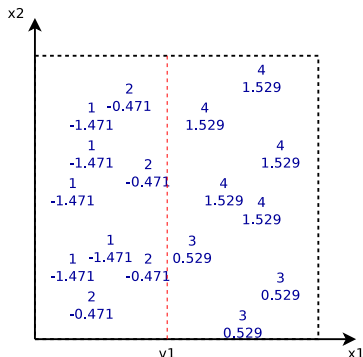
$$\frac{\partial f(x)}{\partial x} = 5 \cdot (-2 + 2x) + 4 \cdot (-4 + 2x)^2 \\ + 3 \cdot (-6 + 2x)^2 + 5 \cdot (-8 + 2x)^2 \\ = -84 + 34x = 32(x - 2.471)$$



Mean value minimizes the square error loss

Learning Boosted Regression Trees (II)

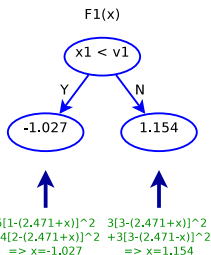
Try all possible v_1
as



$$F_0(x)$$

2.471

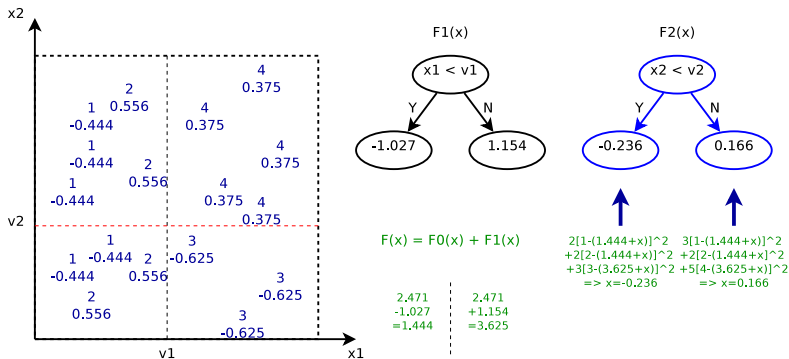
$$F(x) = F_0(x) = 2.471$$



Split root node based on least squares criterion to build a tree predicting the “pseudo”-residuals. Pseudo-residuals are: target value - previous prediction
 $= y_i - F_0(x)$.

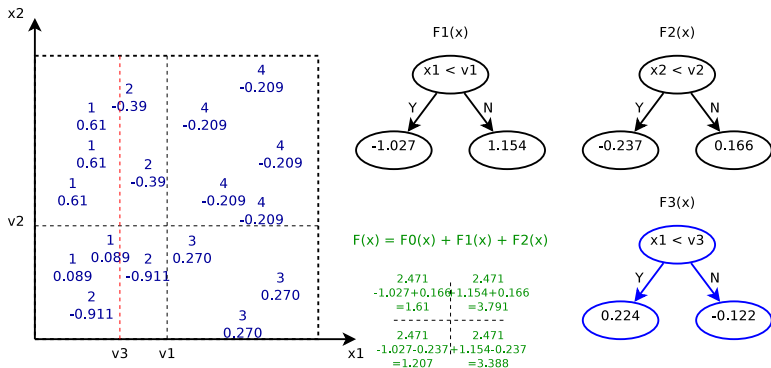
For example, target 1, residual = $1 - 2.471$

Learning Boosted Regression Trees (III)



In the next stage, another tree is created to fit the actual “pseudo”-residuals predicted by the first tree.

Learning Boosted Regression Trees (IV)



This is iteratively continued: in each stage, the algorithm builds a new tree based on the “pseudo”-residuals predicted by the previous tree ensemble.

Multiple Additive Regression Trees (MART)

Algorithm 1 Multiple Additive Regression Trees.

```
1: Initialize  $F_0(\mathbf{x}) = \arg \min_{\gamma} \sum_{i=1}^N L(y_i, \gamma)$ 
2: for  $m = 1, \dots, M$  do
3:   for  $i = 1, \dots, N$  do
4:      $\tilde{y}_{im} = - \left[ \frac{\partial L(y_i, F(\mathbf{x}_i))}{\partial F(\mathbf{x}_i)} \right]_{F(\mathbf{x})=F_{m-1}(\mathbf{x})}$ 
5:   end for
6:    $\{R_{km}\}_{k=1}^K$  // Fit a regression tree to targets  $\tilde{y}_{im}$ 
7:   for  $k = 1, \dots, K_m$  do
8:      $\gamma_{km} = \arg \min_{\gamma} \sum_{\mathbf{x}_i \in R_{jm}} L(y_i, F_{m-1}(\mathbf{x}_i) + \gamma)$ 
9:   end for
10:   $F_m(\mathbf{x}) = F_{m-1}(\mathbf{x}) + \eta \sum_{k=1}^{K_m} \gamma_{km} 1(\mathbf{x}_i \in R_{km})$ 
11: end for
12: Return  $F_M(\mathbf{x})$ 
```

Revised targets=
Pseudo-residual for each instance

Build a regression
tree for revised
targets

Compute the best
prediction for each
of K leaves

Add this new tree
to the MART multiplied
by a learning rate

- BREIMAN, LEO and E. SCHAPIRE (2001). *Random forests*. In *Machine Learning*, pp. 5–32.
- BURGES, CHRISTOPHER J. C. (2010). *From RankNet to LambdaRank to LambdaMART: An Overview*.
- BURGES, CHRISTOPHER J. C., R. RAGNO and Q. V. LE (2006). *Learning to Rank with Nonsmooth Cost Functions*. In SCHÖLKOPF, BERNHARD, J. PLATT and T. HOFFMAN, eds.: *NIPS*, pp. 193–200. MIT Press.
- CHAPELLE, OLIVIER and Y. CHANG (2011). *Yahoo! Learning to Rank Challenge Overview*. Journal of Machine Learning Research - Proceedings Track, 14:1–24.
- CROFT, W.B., D. METZLER and T. STROHMANN (2010). *Search Engines: Information Retrieval in Practice*. Pearson, London, England.

References II

- FRIEDMAN, JEROME H. (1999). *Greedy Function Approximation: A Gradient Boosting Machine*. *Annals of Statistics*, 29(5):1189–1232.
- GANJISAFFAR, YASSER (2011). *Tree Ensembles for Learning to Rank*. PhD thesis, University of California, Irvine.
- HASTIE, TREVOR, R. TIBSHIRANI and J. FRIEDMAN (2002). *The Elements of Statistical Learning*. Springer, New York.
- LIU, TIE-YAN (2010). *Learning to Rank for Information Retrieval..* Springer-Verlag New York Inc.
- MANNING, CHRISTOPHER D., P. RAGHAVAN and H. SCHÜTZE (2008). *Introduction to Information Retrieval*. Cambridge University Press.
- VAPNIK, VLADIMIR N. (1995). *The Nature of Statistical Learning Theory*. Springer New York Inc., New York, NY, USA.
- WU, QIANG, C. J. C. BURGESS, K. M. SVORE and J. GAO (2008). *Ranking, Boosting, and Model Adaptation*.