

STA 221: LECTURE 11

KRISHNA BALASUBRAMANIAN

(UNIVERSITY OF CALIFORNIA, DAVIS)

EXAMPLES OF CLASSIFIERS

Logistic Regression ✓

Support vector Machine

K-NN

Decision Tree

Random Forest

Gradient Boosted Decision Tree (GBDT)

Support Vector Machine

Support vector machines (SVM) were the best performing general purpose classifier before the recent deep learning revolution.

We assume binary classification with $y \in \{+1, -1\}$ and represent the classifier using the inner product notation $\hat{y} = \text{sign}\langle \mathbf{w}, \mathbf{x} \rangle$ where $\langle \mathbf{x}, \mathbf{z} \rangle = \mathbf{x}^\top \mathbf{z}$ in vector notation.

Note that a bias term may be included i.e.

$\hat{y} = \text{sign}(\mathbf{w}_0 + \langle \mathbf{w}, \mathbf{x} \rangle)$ using the notation $\langle \mathbf{w}, \mathbf{x} \rangle$ if \mathbf{x} is augmented with an always one component.

Similarly, note that \mathbf{x} may be a vector of non-linear features of the actual data \mathbf{x}' i.e. $\mathbf{x}_1 = \mathbf{x}'_1 \mathbf{x}'_2$ so the linear classifier in the space of \mathbf{x} is really non-linear in the original data space of \mathbf{x}' .

The linear classifier $f(\mathbf{x}) = \text{sign}\langle \mathbf{w}, \mathbf{x} \rangle$ is parameterized by a weight vector which is normal (i.e. perpendicular) to the decision boundary which is a subspace or a linear hyperplane passing through the origin (note that as described above this does not preclude having a bias term).

Any \mathbf{x} can be represented as a sum of its projection onto the subspace and its perpendicular component

$\mathbf{x} = \mathbf{x}_\perp + \mathbf{x}_\parallel = \mathbf{x}_\parallel + r \frac{\mathbf{w}}{\|\mathbf{w}\|}$. Since

$$\begin{aligned}\langle \mathbf{w}, \mathbf{x} \rangle &= \langle \mathbf{w}, \mathbf{x}_\parallel \rangle + \langle \mathbf{w}, r\mathbf{w}/\|\mathbf{w}\| \rangle = 0 + r\|\mathbf{w}\| \quad \Rightarrow \\ r &= \langle \mathbf{w}, \mathbf{x} \rangle / \|\mathbf{w}\|,\end{aligned}$$

we have that for correctly classified points \mathbf{x}_i, y_i the distance to the hyperplane is $|r_i| = y_i \langle \mathbf{w}, \mathbf{x}_i \rangle / \|\mathbf{w}\|$.

The idea of support vector machines in the context of linearly separable data is to choose \mathbf{w} that leads to the largest margin - defined as the distance of the closest data point to the hyperplane

$$\mathbf{w} = \arg \max_{\mathbf{w} \in \mathbb{R}^d} \left(\|\mathbf{w}\|^{-1} \min_{1 \leq i \leq n} y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \right). \quad (1)$$

The direct solution of (1) is difficult as the objective function is non-differentiable. We proceed by solving an equivalent optimization problem that is easier to solve.

We start by observing that rescaling the weight vector $\mathbf{w}' = c\mathbf{w}$, $c \in \mathbb{R}_+$ leaves the classifier $f(\mathbf{x})$ unchanged and does not change the distance r of points to the subspace. More importantly, it also leaves the objective function (1) unchanged. By rescaling \mathbf{w} so that the distance of the closest point to the hyperplane is 1 we get that $\min_{1 \leq i \leq n} y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1$ with the minimum achieved for one or more training points

$$\mathbf{w} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{w}\|^2 \quad \text{subject to} \quad y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 \quad i = 1, \dots, n. \quad (2)$$

A different way to see the equivalence between (1) and (2) is to note that as \mathbf{w} gets closer to the origin (as we try to do by minimizing $\|\mathbf{w}\|$) one or more of the constraints will be satisfied with equality rather than inequality in which case (1) and (2) are equivalent (i.e., at the solution one or more of the constraints must be active with equality).

Problem (2) is a quadratic program (minimization of a quadratic function subject to linear constraints) and is easier to solve than (1).

However, it involves a large number of linear inequality constraints. The dual problem which is yet another equivalent SVM formulation is the easiest to solve computationally.

It is obtained by optimizing the Lagrangian

$$\mathcal{L}(\mathbf{w}, \lambda) = \frac{1}{2} \|\mathbf{w}\|^2 - \sum_{i=1}^n \lambda_i (y_i \langle \mathbf{w}, \mathbf{x}_i \rangle - 1) \quad (3)$$

with respect to w first, substituting the solution into \mathcal{L} and then optimizing with respect to λ (recall for a point to be a solution of the constrained optimization problem both $\nabla_{\mathbf{w}} \mathcal{L}$ and $\nabla_{\lambda} \mathcal{L}$ need to be zero).

$$0 = \frac{\partial \mathcal{L}(\mathbf{w}, \lambda)}{\partial \mathbf{w}_i} \quad i = 1, \dots, n \quad \Rightarrow \quad \mathbf{w}^* = \sum_{j=1}^n \lambda_j^{(j)} \mathbf{x}_i \mathbf{x}_j \quad (4)$$

$$\begin{aligned} \mathcal{L}(\lambda, \mathbf{w}^*) &= \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle - \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle + \sum_{i=1}^n \lambda_i \\ &= \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle. \end{aligned} \quad (5)$$

We have thus shown, using convex duality that the equivalent dual formulation of SVM is

$$f(\mathbf{x}) = \text{sign}\langle \mathbf{w}, \mathbf{x} \rangle = \text{sign}\left\langle \sum_{j=1}^n \lambda_j y_j \mathbf{x}_j, \mathbf{x} \right\rangle \quad (6)$$

$$= \text{sign} \sum_{j=1}^n \lambda_j y_j \langle \mathbf{x}_j, \mathbf{x} \rangle \quad (7)$$

where

$$\lambda = \arg \max_{\lambda \in \mathbb{R}^d} \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \quad (8)$$

$$\text{subject to } \lambda_i \geq 0. \quad (9)$$

Notice that (6) is also a quadratic program, but in contrast to (2) the constraints are much simpler. Also the number of variables and constraints are n rather than d (data dimensionality) which favors (6) further in high dimensional cases. Furthermore, the solution of (6) will have non-zero λ_j only for constraints that hold with equality at the optimum (the corresponding training points are called support vectors). Thus, many of the λ_j will converge early on to zero effectively reducing the dimensionality of (6).

Non-Separable Case: Thus far, we have assumed that the training data is linearly separable (can be correctly classified using a linear decision surface).

We proceed as before in the separable case, only that this time some examples may be on the wrong side of the hyperplane. In these cases we introduce ξ_i that measure the amount of violation in the constraints $y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1$:

$$(\mathbf{w}, \xi) = \arg \min_{\mathbf{w} \in \mathbb{R}^d, \xi \in \mathbb{R}^n} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \xi_i \quad (10)$$

$$\text{subject to } y_i \langle \mathbf{w}, \mathbf{x}_i \rangle \geq 1 - \xi_i, \quad \xi_i \geq 0 \quad i = 1, \dots, n. \quad (11)$$

The parameter $C \geq 0$ is a regularization parameter controlling the relative importance of the two terms: maximizing margin for correctly classified examples and minimizing the misclassification penalty $\sum \xi_i$. Since we want small ξ_i (in order to minimize the objective function) we can easily determine $\xi_i = \max(0, 1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle)$ and remove ξ_i as optimization variables

$$\mathbf{w} = \arg \min_{\mathbf{w} \in \mathbb{R}^d} \frac{1}{2} \|\mathbf{w}\|^2 + C \sum_{i=1}^m \max(0, 1 - y_i \langle \mathbf{w}, \mathbf{x}_i \rangle). \quad (12)$$

Proceeding as before the dual formulation leads to

$$\lambda = \arg \max_{\lambda \in \mathbb{R}^d} \sum_{i=1}^n \lambda_i - \frac{1}{2} \sum_{i=1}^n \sum_{j=1}^n \lambda_i \lambda_j y_i y_j \langle \mathbf{x}_i, \mathbf{x}_j \rangle \quad (13)$$

$$\text{subject to } 0 \leq \lambda_i \leq C \quad i = 1, \dots, n \quad (14)$$

(with the classifier $f(\mathbf{x})$ given in terms of the dual variables λ as before).

An interesting observation is that (12) can be interpreted as L_2 regularized margin based minimizer

$$\mathbf{w} = \arg \min_{\mathbf{w}} \sum_{i=1}^n l_{\text{hinge}}(y_i \langle \mathbf{w}, \mathbf{x}_i \rangle) + c \|\mathbf{w}\|^2$$

where $l_{\text{hinge}}(z) = (1 - z)_+$. Viewed in this way, we see a striking similarity between SVM and other penalized likelihood methods such as regularized logistic regression.

The function $l_{\text{hinge}}(z)$, called the hinge loss, represents an empirical loss and replaces the logistic regression negative loglikelihood.

The term $\|\mathbf{w}\|^2$ represents regularization penalty.

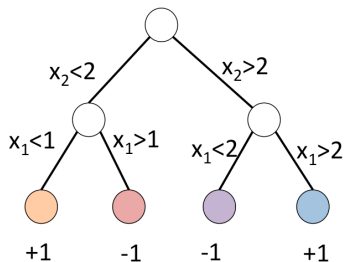
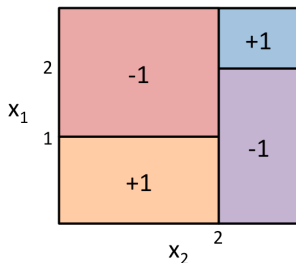
Decision Tree

DECISION TREE

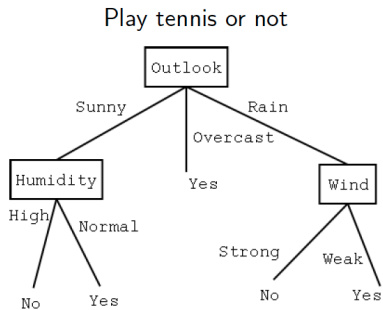
Each node checks one feature x_i :

Go left if $x_i < \text{threshold}$

Go right if $x_i \geq \text{threshold}$



A REAL EXAMPLE



DECISION TREE

Strength:

- It's a **nonlinear** classifier

- Better **interpretability**

- Can naturally handle **categorical** features

Computation:

- Training: **slow**

- Prediction: **fast**

 - h operations (h : depth of the tree, usually ≤ 15)

DECISION TREE

Strength:

- It's a **nonlinear** classifier

- Better **interpretability**

- Can naturally handle **categorical** features

Computation:

- Training: **slow**

- Prediction: **fast**

 - h operations (h : depth of the tree, usually ≤ 15)

SPLITTING THE NODE

Classification tree: Split the node to maximize entropy

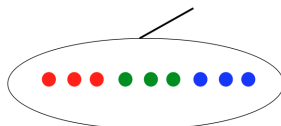
Let S be set of data points in a node, $c = 1, \dots, C$ are labels:

$$\text{Entropy} : H(S) = - \sum_{c=1}^C p(c) \log p(c),$$

where $p(c)$ is the proportion of the data belong to class c .

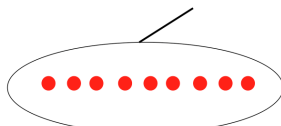
Entropy=0 if all samples are in the same class

Entropy is large if $p(1) = \dots = p(C)$



Entropy:
 $-(1/3) \cdot \log(1/3) - (1/3) \cdot \log(1/3) - (1/3) \cdot \log(1/3)$
 $= 1.58$

Bad split



Entropy:
 $-1 \cdot \log(1) = 0$

Good split

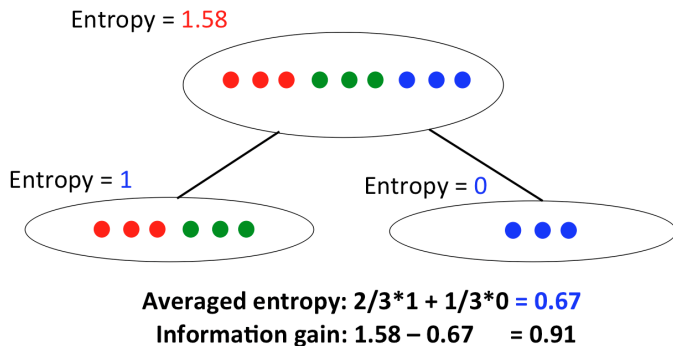
The averaged entropy of a split $S \rightarrow S_1, S_2$

$$\frac{|S_1|}{|S|}H(S_1) + \frac{|S_2|}{|S|}H(S_2)$$

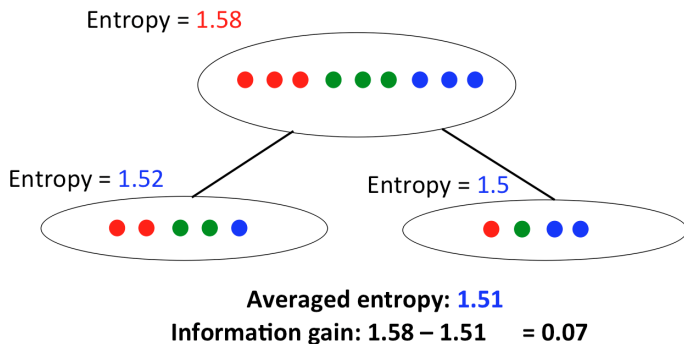
Information gain: measure how good is the split

$$H(S) - \left((|S_1|/|S|)H(S_1) + (|S_2|/|S|)H(S_2) \right)$$

INFORMATION GAIN



INFORMATION GAIN



SPLITTING THE NODE

Given the current node, how to find the **best split**?

For all the **features** and all the **threshold**

 Compute the information gain after the split

 Choose the best one (**maximal information gain**)

For n samples and d features: need $O(nd)$ time

SPLITTING THE NODE

Given the current node, how to find the **best split**?

For all the **features** and all the **threshold**

 Compute the information gain after the split

 Choose the best one (**maximal information gain**)

For n samples and d features: need $O(nd)$ time

SPLITTING THE NODE

Given the current node, how to find the best split?

For all the features and all the threshold

 Compute the information gain after the split

 Choose the best one (maximal information gain)

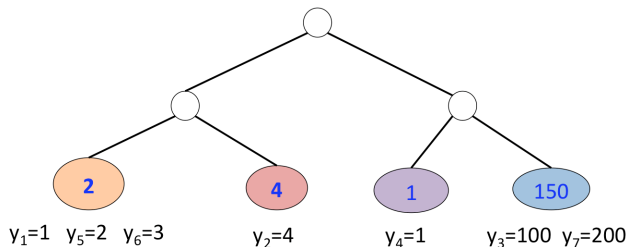
For n samples and d features: need $O(nd)$ time

REGRESSION TREE

Assign a real number for each leaf

Usually **averaged y values** for each leaf

(minimize square error)



Objective function:

$$\min_F \frac{1}{n} \sum_{i=1}^n (y_i - F(\mathbf{x}_i))^2 + (\text{Regularization})$$

The quality of partition $S = S_1 \cup S_2$ can be computed by the objective function:

$$\sum_{i \in S_1} (y_i - y^{(1)})^2 + \sum_{i \in S_2} (y_i - y^{(2)})^2,$$

where $y^{(1)} = \frac{1}{|S_1|} \sum_{i \in S_1} y_i$, $y^{(2)} = \frac{1}{|S_2|} \sum_{i \in S_2} y_i$

Find the best split:

Try all the features & thresholds and find the one with
minimal objective function

Objective function:

$$\min_F \frac{1}{n} \sum_{i=1}^n (y_i - F(\mathbf{x}_i))^2 + (\text{Regularization})$$

The quality of partition $S = S_1 \cup S_2$ can be computed by the objective function:

$$\sum_{i \in S_1} (y_i - y^{(1)})^2 + \sum_{i \in S_2} (y_i - y^{(2)})^2,$$

where $y^{(1)} = \frac{1}{|S_1|} \sum_{i \in S_1} y_i$, $y^{(2)} = \frac{1}{|S_2|} \sum_{i \in S_2} y_i$

Find the best split:

Try all the features & thresholds and find the one with
minimal objective function

PARAMETERS

Maximum depth: (usually ~ 10)

Minimum number of nodes in each node: (10, 50, 100)

Single decision tree is not very powerful...

Can we build multiple decision trees and **ensemble** them together?

PARAMETERS

Maximum depth: (usually ~ 10)

Minimum number of nodes in each node: (10, 50, 100)

Single decision tree is not very powerful...

Can we build multiple decision trees and **ensemble** them together?

Random Forest

RANDOM FOREST

Random Forest (Bootstrap ensemble for decision trees):

- Create T trees

- Learn each tree using a subsampled dataset S_i and subsampled feature set D_i

- Prediction: Average the results from all the T trees

Benefit:

- Avoid over-fitting

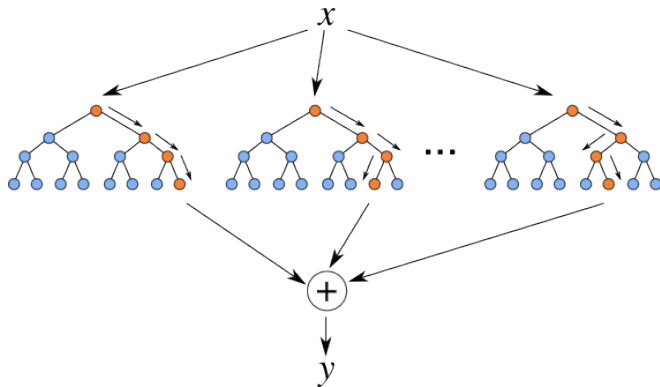
- Improve stability and accuracy

Good software available:

- R: “randomForest” package

- Python: sklearn

RANDOM FOREST



Gradient Boosted Decision Tree

BOOSTED DECISION TREE

Minimize loss $\ell(y, F(x))$ with $F(\cdot)$ being ensemble trees

$$F^* = \underset{F}{\operatorname{argmin}} \sum_{i=1}^n \ell(\mathbf{y}_i, F(\mathbf{x}_i)) \quad \text{with} \quad F(\mathbf{x}) = \sum_{m=1}^T f_m(\mathbf{x})$$

(each f_m is a decision tree)

Direct loss minimization: at each stage m , find the best function to minimize loss

$$\begin{aligned} \text{solve } f_m &= \underset{f_m}{\operatorname{argmin}} \sum_{i=1}^N \ell(y_i, F_{m-1}(\mathbf{x}_i) + f_m(\mathbf{x}_i)) \\ \text{update } F_m &\leftarrow F_{m-1} + f_m \end{aligned}$$

$F_m(\mathbf{x}) = \sum_{j=1}^m f_j(\mathbf{x})$ is the prediction of \mathbf{x} after m iterations.

Two problems:

Hard to implement for general loss

Tend to overfit training data

BOOSTED DECISION TREE

Minimize loss $\ell(y, F(x))$ with $F(\cdot)$ being ensemble trees

$$F^* = \underset{F}{\operatorname{argmin}} \sum_{i=1}^n \ell(\mathbf{y}_i, F(\mathbf{x}_i)) \quad \text{with} \quad F(\mathbf{x}) = \sum_{m=1}^T f_m(\mathbf{x})$$

(each f_m is a decision tree)

Direct loss minimization: at each stage m , find the best function to **minimize loss**

$$\begin{aligned} \text{solve } f_m &= \underset{f_m}{\operatorname{argmin}} \sum_{i=1}^N \ell(y_i, F_{m-1}(\mathbf{x}_i) + f_m(\mathbf{x}_i)) \\ \text{update } F_m &\leftarrow F_{m-1} + f_m \end{aligned}$$

$F_m(\mathbf{x}) = \sum_{j=1}^m f_j(\mathbf{x})$ is the prediction of \mathbf{x} after m iterations.

Two problems:

Hard to implement for general loss

Tend to overfit training data

BOOSTED DECISION TREE

Minimize loss $\ell(y, F(x))$ with $F(\cdot)$ being ensemble trees

$$F^* = \operatorname{argmin}_F \sum_{i=1}^n \ell(\mathbf{y}_i, F(\mathbf{x}_i)) \quad \text{with} \quad F(\mathbf{x}) = \sum_{m=1}^T f_m(\mathbf{x})$$

(each f_m is a decision tree)

Direct loss minimization: at each stage m , find the best function to **minimize loss**

$$\begin{aligned} \text{solve } f_m &= \operatorname{argmin}_{f_m} \sum_{i=1}^N \ell(y_i, F_{m-1}(\mathbf{x}_i) + f_m(\mathbf{x}_i)) \\ \text{update } F_m &\leftarrow F_{m-1} + f_m \end{aligned}$$

$F_m(\mathbf{x}) = \sum_{j=1}^m f_j(\mathbf{x})$ is the prediction of \mathbf{x} after m iterations.

Two problems:

Hard to implement for general loss

Tend to overfit training data

GRADIENT BOOSTED DECISION TREE (GBDT)

Approximate the current loss function by a quadratic approximation:

$$\begin{aligned}\sum_{i=1}^n \ell_i(\hat{y}_i + f_m(\mathbf{x}_i)) &\approx \sum_{i=1}^n (\ell_i(\hat{y}_i) + g_i f_m(\mathbf{x}_i) + \frac{1}{2} h_i f_m(\mathbf{x}_i)^2) \\ &= \sum_{i=1}^n \frac{h_i}{2} \|f_m(\mathbf{x}_i) - g_i/h_i\|^2 + \text{constant}\end{aligned}$$

where $g_i = \partial_{\hat{y}_i} \ell_i(\hat{y}_i)$ is gradient,
 $h_i = \partial_{\hat{y}_i}^2 \ell_i(\hat{y}_i)$ is second order derivative

Finding $f_m(\mathbf{x}, \theta_m)$ by minimizing the loss function:

$$\operatorname{argmin}_{f_m} \sum_{i=1}^N [f_m(\mathbf{x}_i, \theta) - g_i/h_i]^2 + R(f_m)$$

Reduce the training of any loss function to regression tree (just need to compute g_i for different functions)
 $h_i = \alpha$ (fixed step size) for original GBDT.

XGboost shows computing second order derivative yields better performance

Algorithm:

Computing the current gradient for each \hat{y}_i .

Building a base learner (decision tree) to fit the gradient.

Updating current prediction $\hat{y}_i = F_m(\mathbf{x}_i)$ for all i .

Finding $f_m(\mathbf{x}, \theta_m)$ by minimizing the loss function:

$$\operatorname{argmin}_{f_m} \sum_{i=1}^N [f_m(\mathbf{x}_i, \theta) - g_i/h_i]^2 + R(f_m)$$

Reduce the training of any loss function to regression tree (just need to compute g_i for different functions)

$h_i = \alpha$ (fixed step size) for original GBDT.

XGboost shows computing second order derivative yields better performance

Algorithm:

Computing the current gradient for each \hat{y}_i .

Building a base learner (decision tree) to fit the gradient.

Updating current prediction $\hat{y}_i = F_m(\mathbf{x}_i)$ for all i .

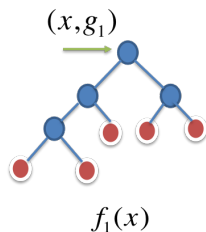
GRADIENT BOOSTED DECISION TREES (GBDT)

Key idea:

Each base learner is a decision tree

Each regression tree approximates the functional gradient

$$\frac{\partial \ell}{\partial F}$$



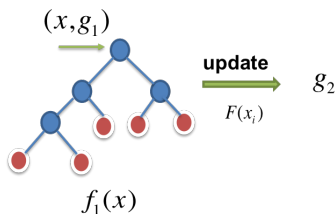
GRADIENT BOOSTED DECISION TREES (GBDT)

Key idea:

Each base learner is a decision tree

Each regression tree approximates the functional gradient

$$\frac{\partial \ell}{\partial F}$$



$$F_{m-1}(x_i) = \sum_{j=1}^{m-1} f_j(x_i) \quad g_m(x_i) = \frac{\partial \ell(y_i, F(x_i))}{\partial F(x_i)} \Big|_{F(x_i)=F_{m-1}(x_i)}$$

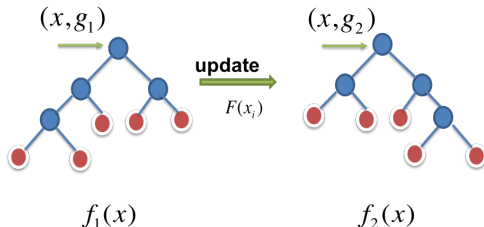
GRADIENT BOOSTED DECISION TREES (GBDT)

Key idea:

Each base learner is a decision tree

Each regression tree approximates the functional gradient

$$\frac{\partial \ell}{\partial F}$$



$$F_{m-1}(x_i) = \sum_{j=1}^{m-1} f_j(x_i) \quad g_m(x_i) = \left. \frac{\partial \ell(y_i, F(x_i))}{\partial F(x_i)} \right|_{F(x_i)=F_{m-1}(x_i)}$$

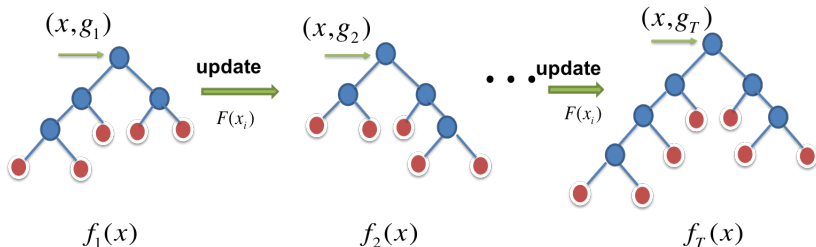
GRADIENT BOOSTED DECISION TREES (GBDT)

Key idea:

Each base learner is a decision tree

Each regression tree approximates the functional gradient

$$\frac{\partial \ell}{\partial F}$$



$$F_{m-1}(x_i) = \sum_{j=1}^{m-1} f_j(x_i) \quad g_m(x_i) = \left. \frac{\partial \ell(y_i, F(x_i))}{\partial F(x_i)} \right|_{F(x_i) = F_{m-1}(x_i)}$$

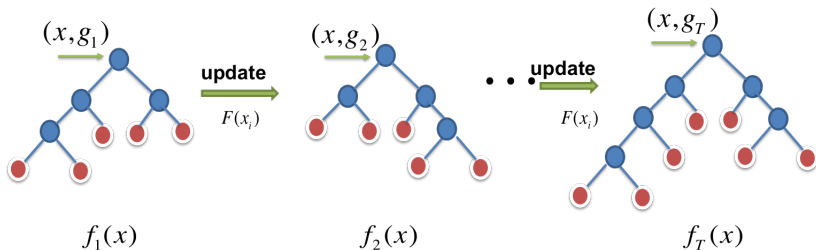
GRADIENT BOOSTED DECISION TREES (GBDT)

Key idea:

Each base learner is a decision tree

Each regression tree approximates the functional gradient

$$\frac{\partial \ell}{\partial f}$$



Final prediction
$$F(x_i) = \sum_{j=1}^T f_j(x_i)$$