

Boosting (Chapter 18 of Prob ML book)

- Model's form is weighted sum of component models, like in bagging, but the training method is different
- Form of model:
 - $f(x; \theta) = \sum_1^M \beta_m F_m(x, \theta_m)$
- Sequentially train (weak learners) where subsequent learners focus on residues of the previous stages → complicated decision boundaries from simple classifiers
- Simple classifier examples:
 - Naïve Bayes classifier
 - Decision stumps
 - Linear classifiers.

$$f(x) = \beta_1 F_1(x, \theta_1) + \dots + \beta_M F_M(x, \theta_M)$$

$\beta_1 \quad \beta_2 \quad \dots \quad \beta_M$
 $F_1(\theta_1) \quad F_2(\theta_2) \quad \dots \quad F_M(\theta_M)$
 $\uparrow x \quad \uparrow x \quad \dots \quad \uparrow x$

Forward stagewise additive boosting

- Sequentially optimize the overall training loss while keeping previous stages fixed.

- Given training data and loss: $f(x, \theta)$ $D \equiv \{(x^i, y_i) : i = 1 \dots N\}$

$$\min_{\theta} \sum_{i=1}^N l(y_i, f(x^i, \theta)) \equiv \min_{\theta_1, \theta_2, \dots, \theta_M} \sum_{i=1}^N l(y_i, \sum_{m=1}^M \beta_m F_m(x^i, \theta_m))$$

- Initially: $\underline{f_1} = \operatorname{argmin}_{\theta} \sum_{i=1}^N l(y_i, \underline{f}(x^i, \theta)) \equiv \operatorname{argmin}_{\theta} \sum_{i=1}^N l(y_i, \underline{F}(x^i, \theta))$
- For $\underline{m} = \underline{2}$ to \underline{M}

- $$\theta_m \equiv \operatorname{argmin}_{\theta} \sum_{i=1}^N l(y_i, \underbrace{f_{m-1}(x^i)}_{\text{previous stage}} + \underbrace{F_m(x^i, \theta)}_{\text{new stage}})$$

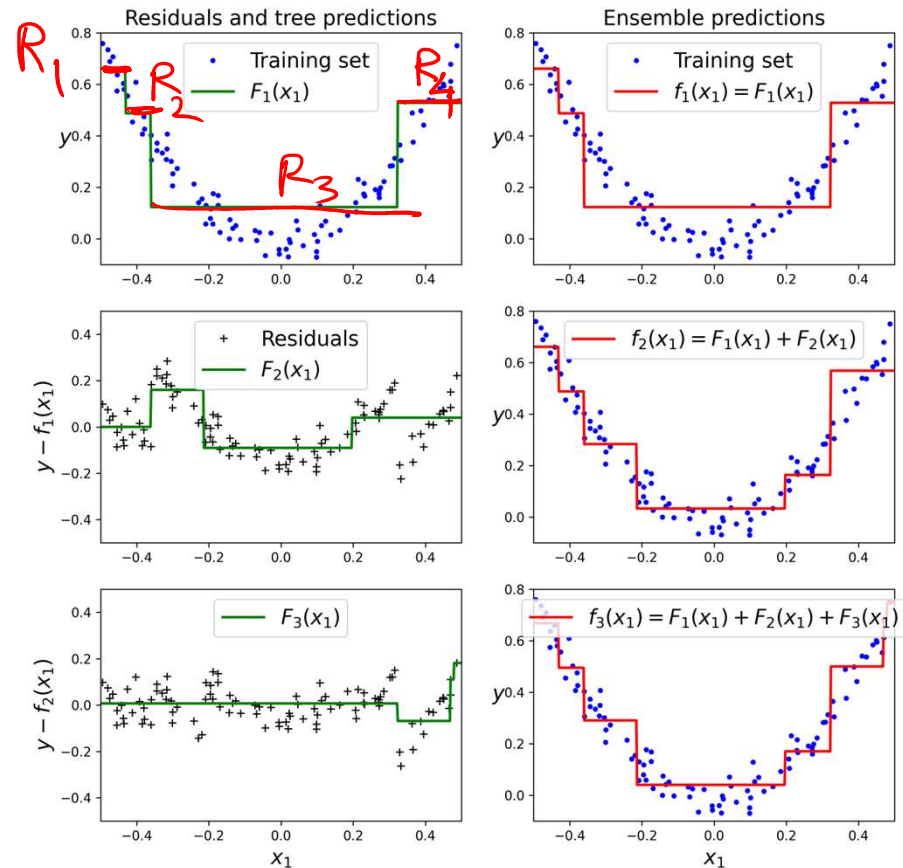
- $$f_m = \sum_{j=1}^m F_j(x; \theta_j)$$

$f = f_m$

$\beta_m \Leftarrow$ Algorithm specific to be covered.
default = 1.

Boosted regression trees with square loss

- When loss function is square loss



- <https://drive.google.com/file/d/1a7kfsrBtjwhP4x7XbwOtbWPGoxdP5WS/view?usp=sharing>

Gradient boosting algorithm.

- A general-purpose algorithm for any differentiable loss function.
- Given N training examples, view $f(x^i)$ as an independent parameter

$$w_1 = f(x^1) \quad w_2 = f(x^2) \quad \dots \quad w_N = f(x^N)$$

$$\min_{f(x^1), \dots, f(x^N)} \sum_{i=1}^N \ell(y_i - f(x^i))^2$$

- At each stage m , we have a current value of the function.

$$w_i^m = f_m(x^i) \quad f_m(x^i) = \sum_{j=1}^{m-1} \beta_j F_j(x^i)$$

- Find gradient of the loss with respect to each $f(x^i)$ independently.

$$\frac{\partial \ell(y_i - f(x^i))}{\partial f(x^i)} \bigg|_{f(x^i) = w_i^m = f_{m-1}(x^i)}$$

- Function at the m th stage F_m is trained to match the gradients.

Algorithm 18.3: Gradient boosting

```

1 Initialize  $f_0(\mathbf{x}) = \operatorname{argmin}_F \sum_{i=1}^N L(y_i, F(\mathbf{x}_i))$ 
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  $F_m = \operatorname{argmin}_F \sum_{i=1}^N (r_{im} - F(\mathbf{x}_i))^2$ 
5     Update  $f_m(\mathbf{x}) = f_{m-1}(\mathbf{x}) + \nu F_m(\mathbf{x})$ 
6 Return  $f(\mathbf{x}) = f_M(\mathbf{x})$ 

```

Gradient boosted regression trees.

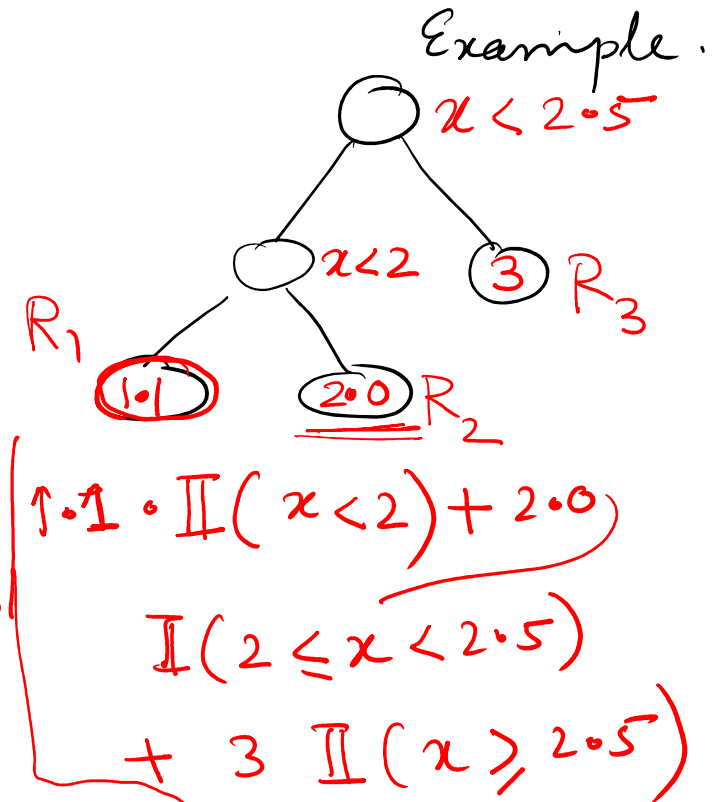
- Assume tree is of the form

$$\underline{F_m(x)} = \sum_{j=1}^{J_m} \underline{\omega_{jm}} \underline{\mathbb{I}(x \in R_{jm})}$$

- First find gradient of loss at m (residuals) and build the tree to find good regions

loss function = square loss

$$\frac{\partial}{\partial f(x_i)} \frac{1}{2} (y_i - f(x_i))^2 \Big|_{f(x_i) = f_{m-1}(x_i)} = (y_i - f_{m-1}(x_i))$$



- Then resolve for the weights of each leaf by residual

$$\underline{\hat{\omega}_{jm}} = \underset{\underline{\omega}}{\operatorname{argmin}} \sum_{x_i \in R_{jm}} \ell(y_i, f_{m-1}(x_i) + \omega)$$

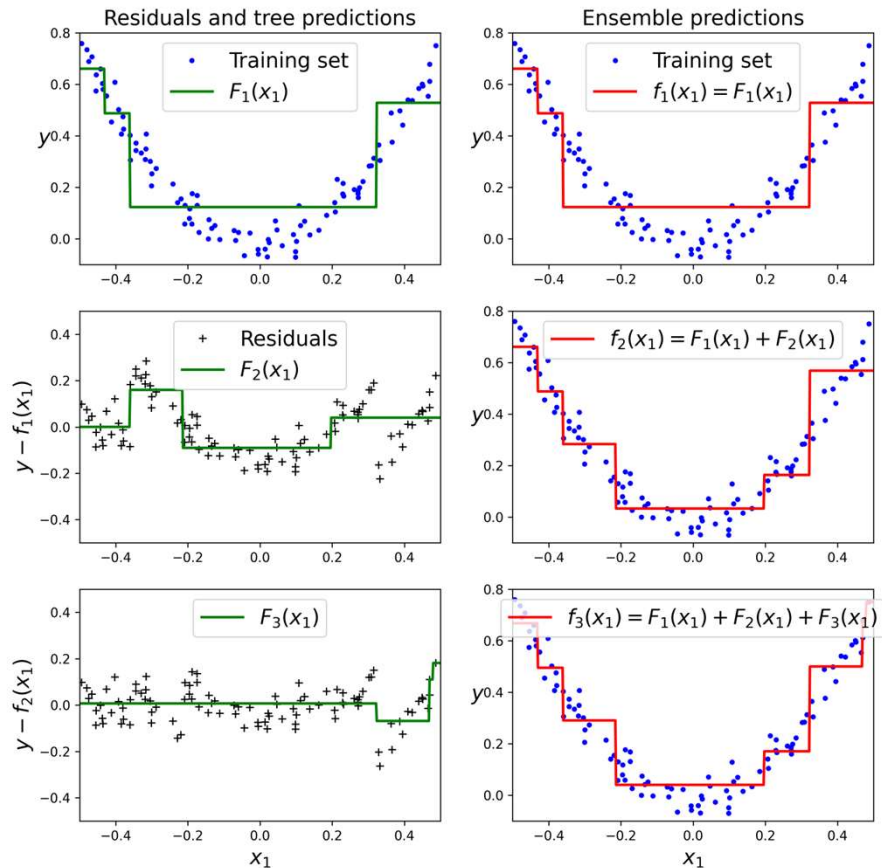
Boosted regression trees with square loss

- When loss function is square loss

- Residue is:

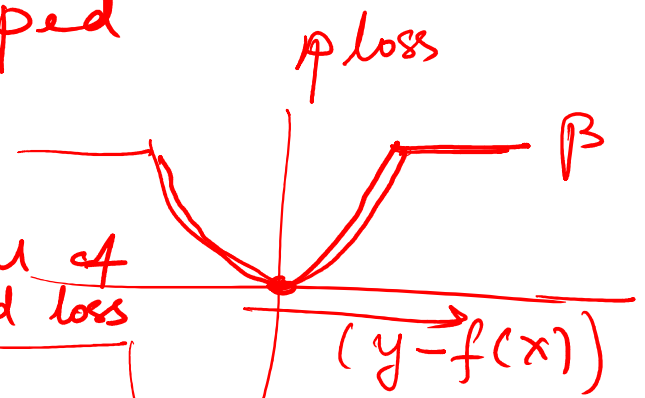
- Weight of each leaf

- <https://drive.google.com/file/d/1a7kfsrBtjwhP4x7XbwOtbWPG0xdP5WS/view?usp=sharing>

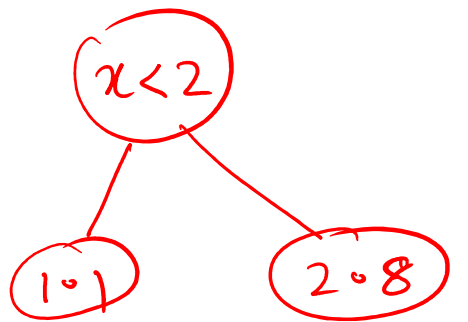


Fit regression tree with clipped

$$\min(\beta, (y - f(x))^2)$$



f_1



x	y	$y - f_1(x)$	gradient of clipped loss
1	1	-0.1	0.1
1.1	1.2	0.1	-0.1
2	2.1	-0.7	0
2.2	1.9	-0.9	0
4	3	0.2	-0.2
4.2	2.8	0	0

$$\beta = 0.1$$

$$\frac{\partial}{\partial f(x)} \min(\beta, \frac{(y - f(x))^2}{2})$$

$$= 0 \quad \text{if } \frac{(y - f(x))^2}{2} > \beta$$

$$= -(y - f(x))$$

Other variants

- Extreme gradient boosting (XGBoost).
- Significantly more evolved.
- But high accuracy with high speed training
- Extensively used