

Ensemble Learning: Gradient Boosting Machine (GBM)

Pilsung Kang

School of Industrial Management Engineering

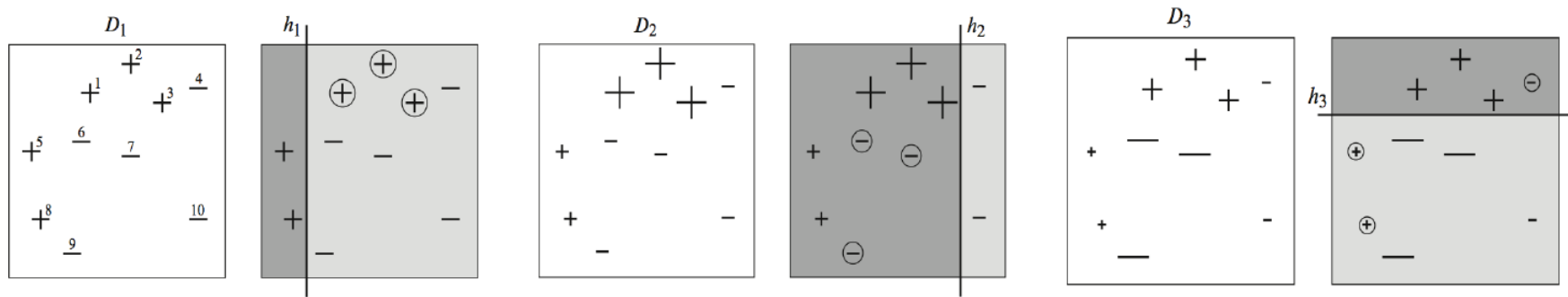
Korea University

Gradient Boosting Machine: GBM

Friedman (2001), Natekin and Knoll (2013)

Gradient Boosting = Gradient Descent + Boosting

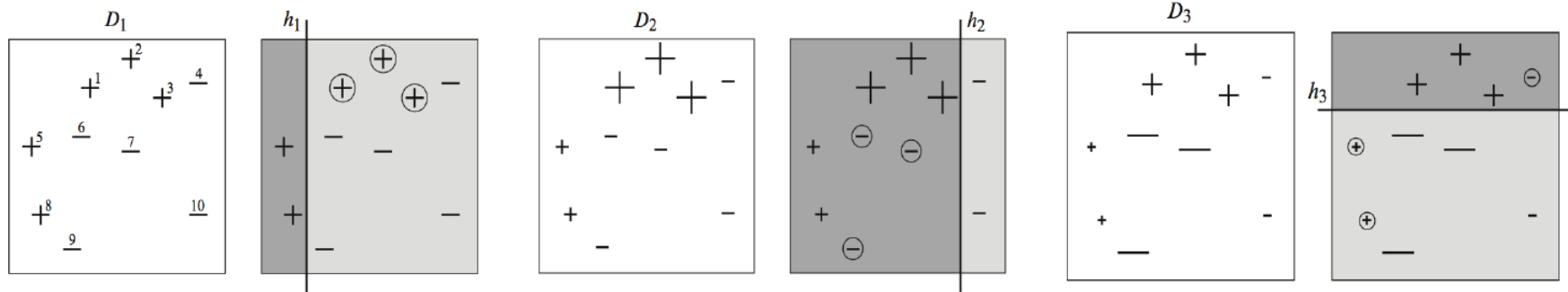
- Adaboost



- ✓ Fit an additive model (ensemble) $\sum_t \rho_t h_t(x)$ in a forward stage-wise manner.
- ✓ In each stage, introduce a weak learner to compensate the shortcomings of existing weak learners.
- ✓ In Adaboost, “shortcomings” are identified by high-weight data points.

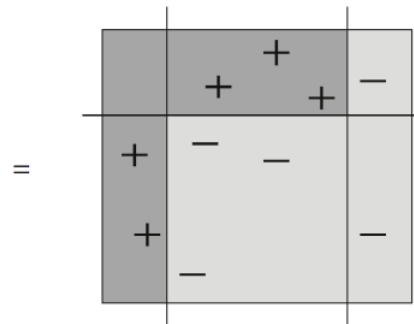
Gradient Boosting Machine: GBM

- Adaboost



$$H = \text{sign} \left(0.42 \begin{array}{|c|} \hline \text{[Diagram of } h_1 \text{]} \\ \hline \end{array} + 0.65 \begin{array}{|c|} \hline \text{[Diagram of } h_2 \text{]} \\ \hline \end{array} + 0.92 \begin{array}{|c|} \hline \text{[Diagram of } h_3 \text{]} \\ \hline \end{array} \right)$$

$$H(x) = \sum_t \rho_t h_t(x)$$

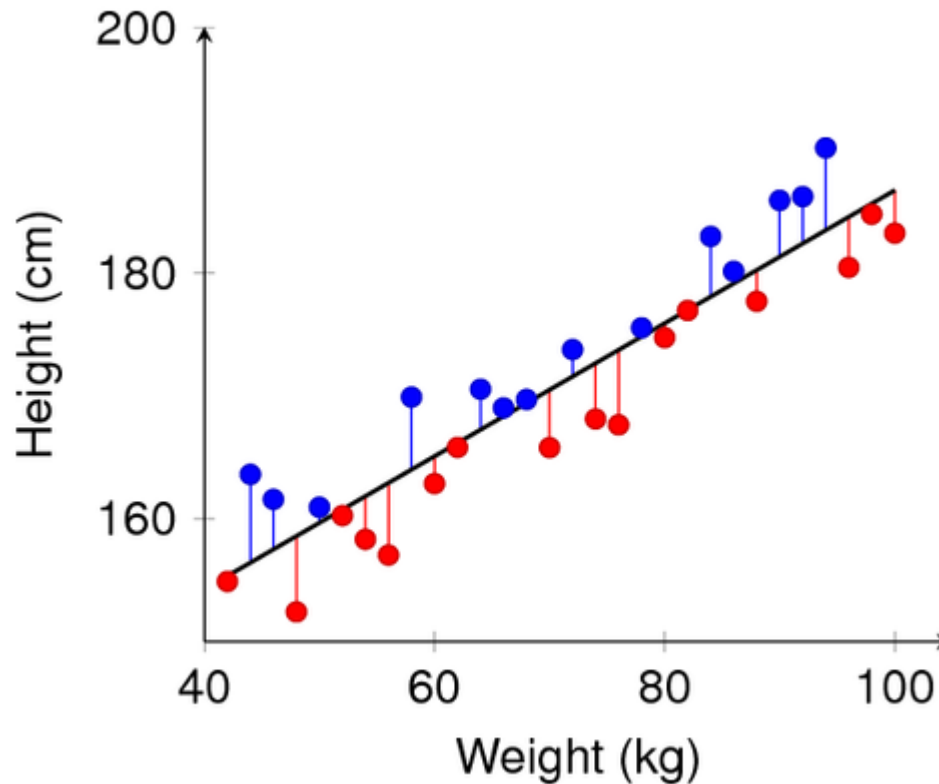


Gradient Boosting Machine: GBM

- Gradient Boosting
 - ✓ Fit an additive model (ensemble) $\sum_t \rho_t h_t(x)$ in a forward stage-wise manner.
 - ✓ In each stage, introduce a weak learner to compensate the shortcomings of existing weak learners.
 - ✓ In Gradient Boosting, “shortcomings” are identified by gradients.
 - ✓ Both high-weight data points and gradients tell us how to improve our model.
- Gradient Boosting for Different Problems
 - ✓ Difficulty: Regression < Classification < Ranking
 - Associated with the complexity of the derivative of a loss function

Gradient Boosting Machine: GBM

- Motivation (for regression problem)
 - ✓ What if we attempt to predict the residuals with the additional regression model?

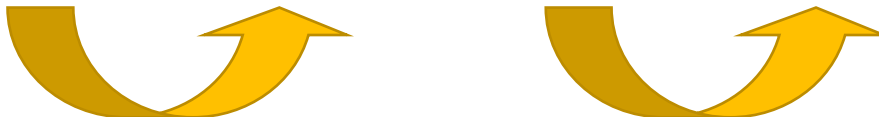


Gradient Boosting Machine: GBM

- Main idea

Original Dataset		Modified Dataset 1		Modified Dataset 2	
x^1	y^1	x^1	$y^1 - f_1(x^1)$	x^1	$y^1 - f_1(x^1) - f_2(x^1)$
x^2	y^2	x^2	$y^2 - f_1(x^2)$	x^2	$y^2 - f_1(x^2) - f_2(x^2)$
x^3	y^3	x^3	$y^3 - f_1(x^3)$	x^3	$y^3 - f_1(x^3) - f_2(x^3)$
x^4	y^4	x^4	$y^4 - f_1(x^4)$	x^4	$y^4 - f_1(x^4) - f_2(x^4)$
x^5	y^5	x^5	$y^5 - f_1(x^5)$	x^5	$y^5 - f_1(x^5) - f_2(x^5)$
x^6	y^6	x^6	$y^6 - f_1(x^6)$	x^6	$y^6 - f_1(x^6) - f_2(x^6)$
x^7	y^7	x^7	$y^7 - f_1(x^7)$	x^7	$y^7 - f_1(x^7) - f_2(x^7)$
x^8	y^8	x^8	$y^8 - f_1(x^8)$	x^8	$y^8 - f_1(x^8) - f_2(x^8)$
x^9	y^9	x^9	$y^9 - f_1(x^9)$	x^9	$y^9 - f_1(x^9) - f_2(x^9)$
x^{10}	y^{10}	x^{10}	$y^{10} - f_1(x^{10})$	x^{10}	$y^{10} - f_1(x^{10}) - f_2(x^{10})$

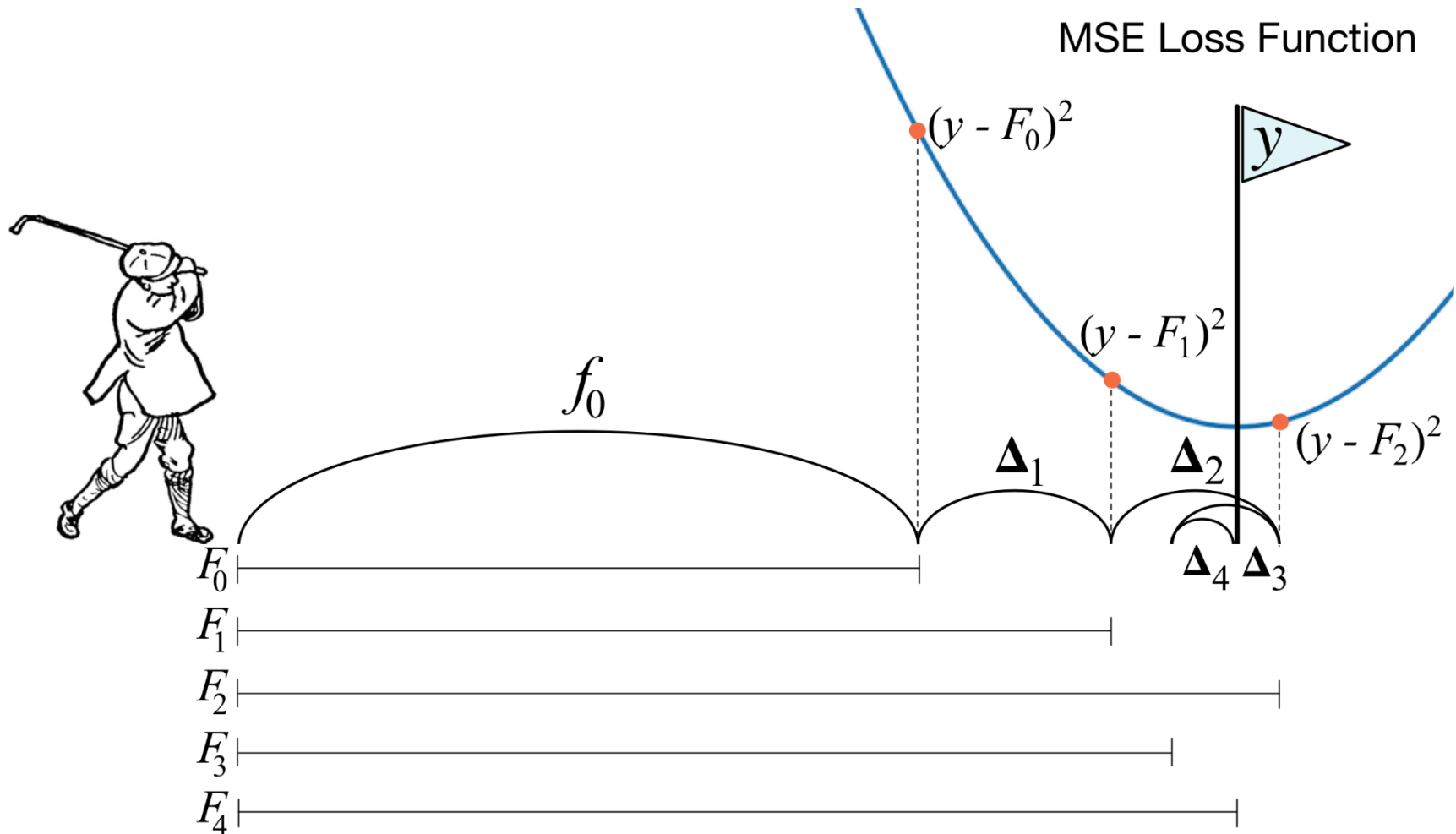
...



$$y = f_1(\mathbf{x}) \quad y - f_1(\mathbf{x}) = f_2(\mathbf{x}) \quad y - f_1(\mathbf{x}) - f_2(\mathbf{x}) = f_3(\mathbf{x})$$

Gradient Boosting Machine: GBM

- Illustrative Example



<https://github.com/topics/gradient-boosting-machine>

Gradient Boosting Machine: GBM

- How is this idea related to the gradient?
 - ✓ Loss function of the ordinary least square (OLS)

$$\min L = \frac{1}{2} \sum_{i=1}^n (y_i - f(\mathbf{x}_i))^2$$

- ✓ Gradient of the Loss function

$$\frac{\partial L}{\partial f(\mathbf{x}_i)} = f(\mathbf{x}_i) - y_i$$

- ✓ Residuals are the negative gradient of the loss function

$$y_i - f(\mathbf{x}_i) = -\frac{\partial L}{\partial f(\mathbf{x}_i)}$$

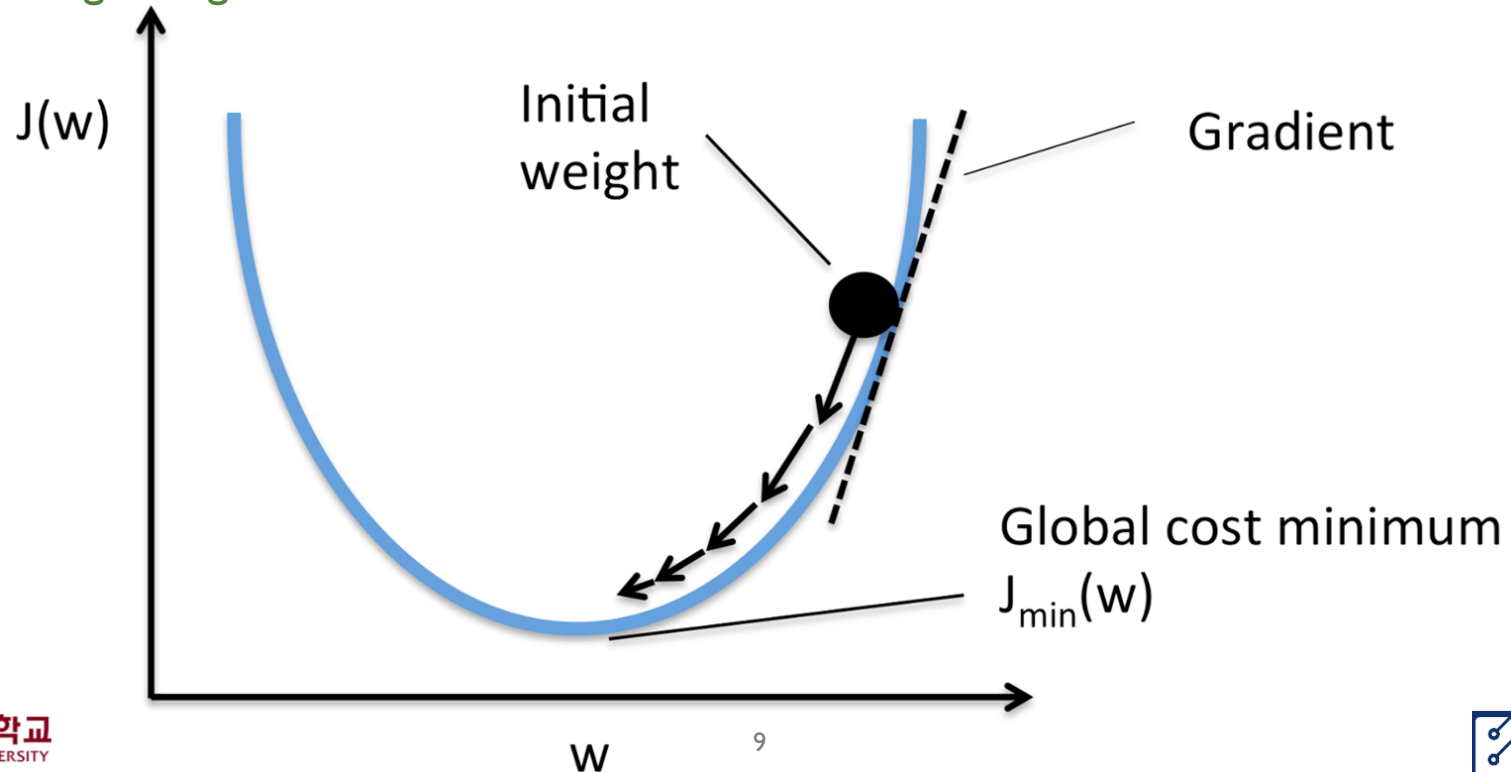
Gradient Boosting Machine: GBM

- Gradient Descent Algorithm

- ✓ Blue line: value of loss function with a given parameter

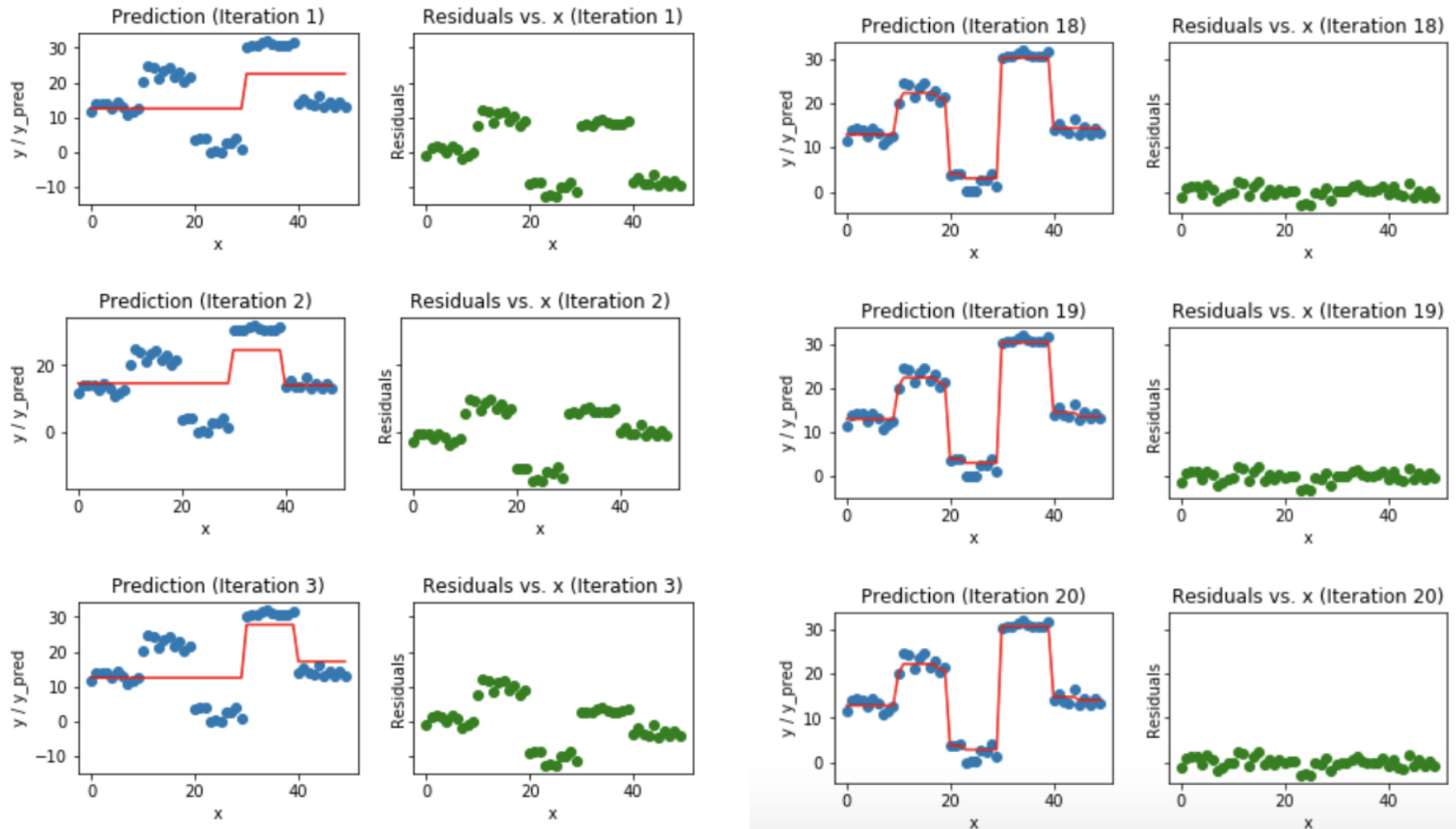
- ✓ Black point: current state

- ✓ Arrows: the direction that the parameter should follow to minimize the loss function
= negative gradient of the loss function



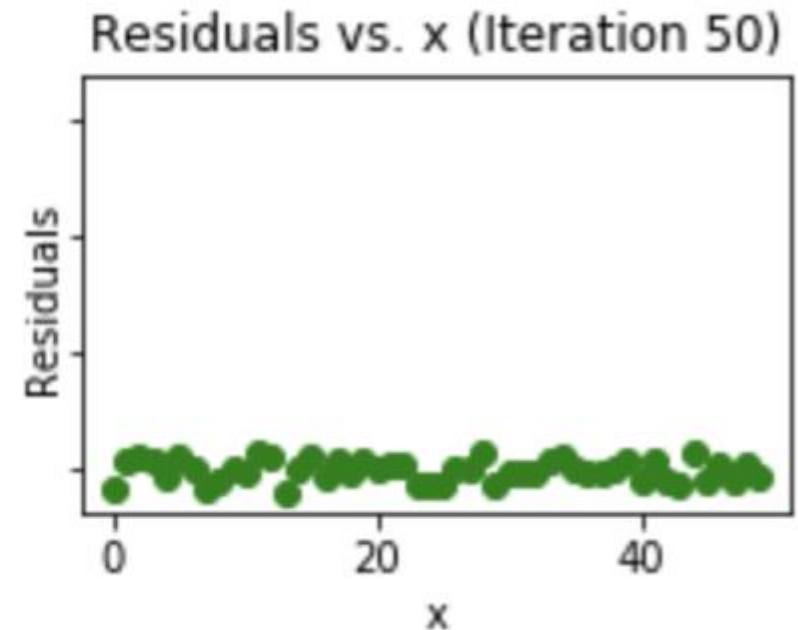
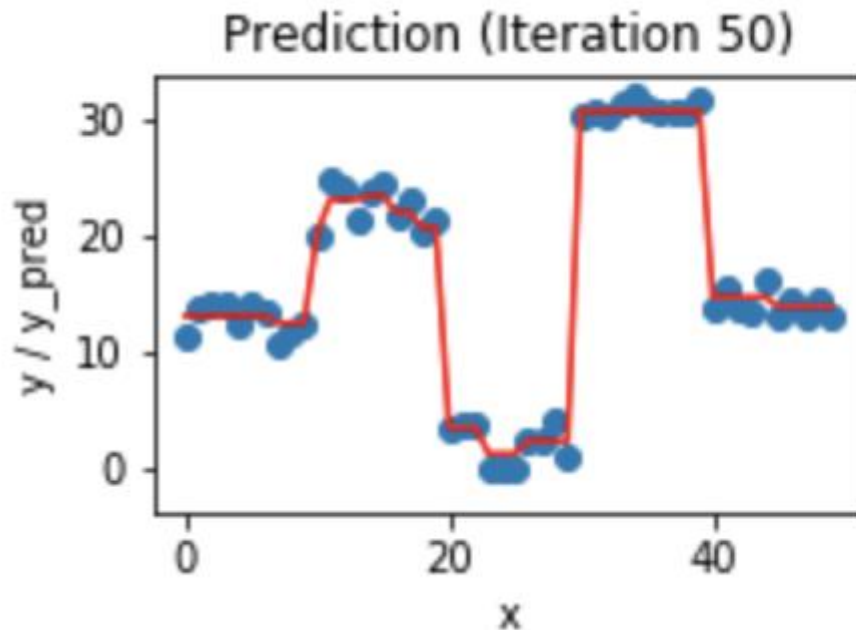
Gradient Boosting Machine: GBM

- GBM Regression Example I



Gradient Boosting Machine: GBM

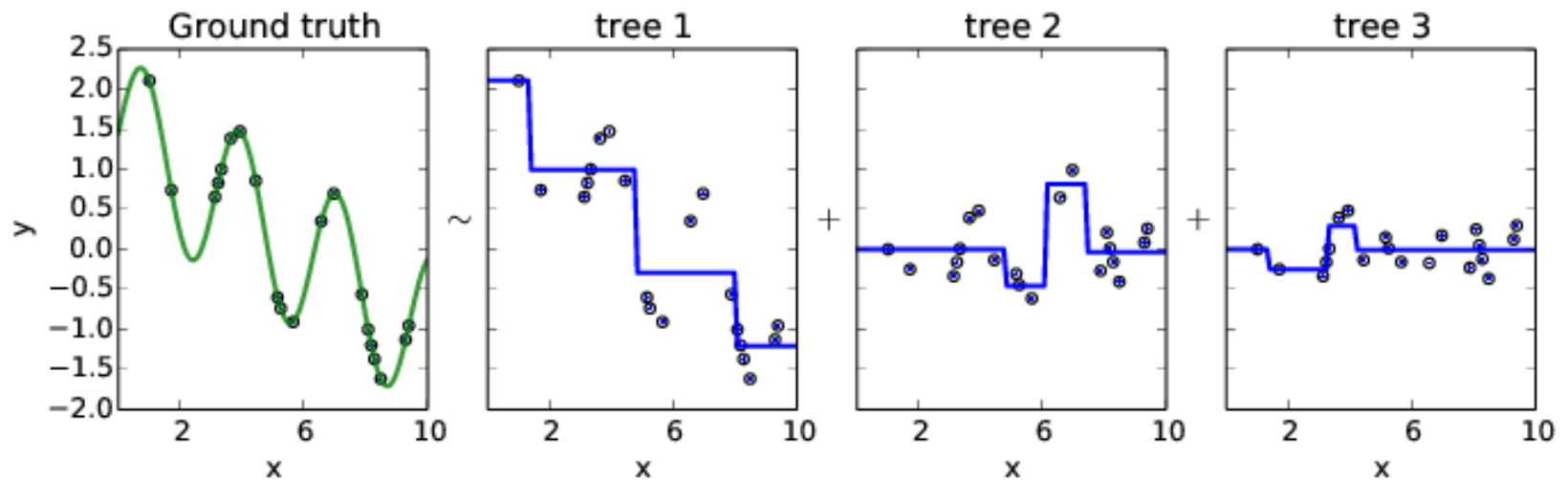
- GBM Regression Example I



<https://medium.com/mlreview/gradient-boosting-from-scratch-1e317ae4587d>

Gradient Boosting Machine: GBM

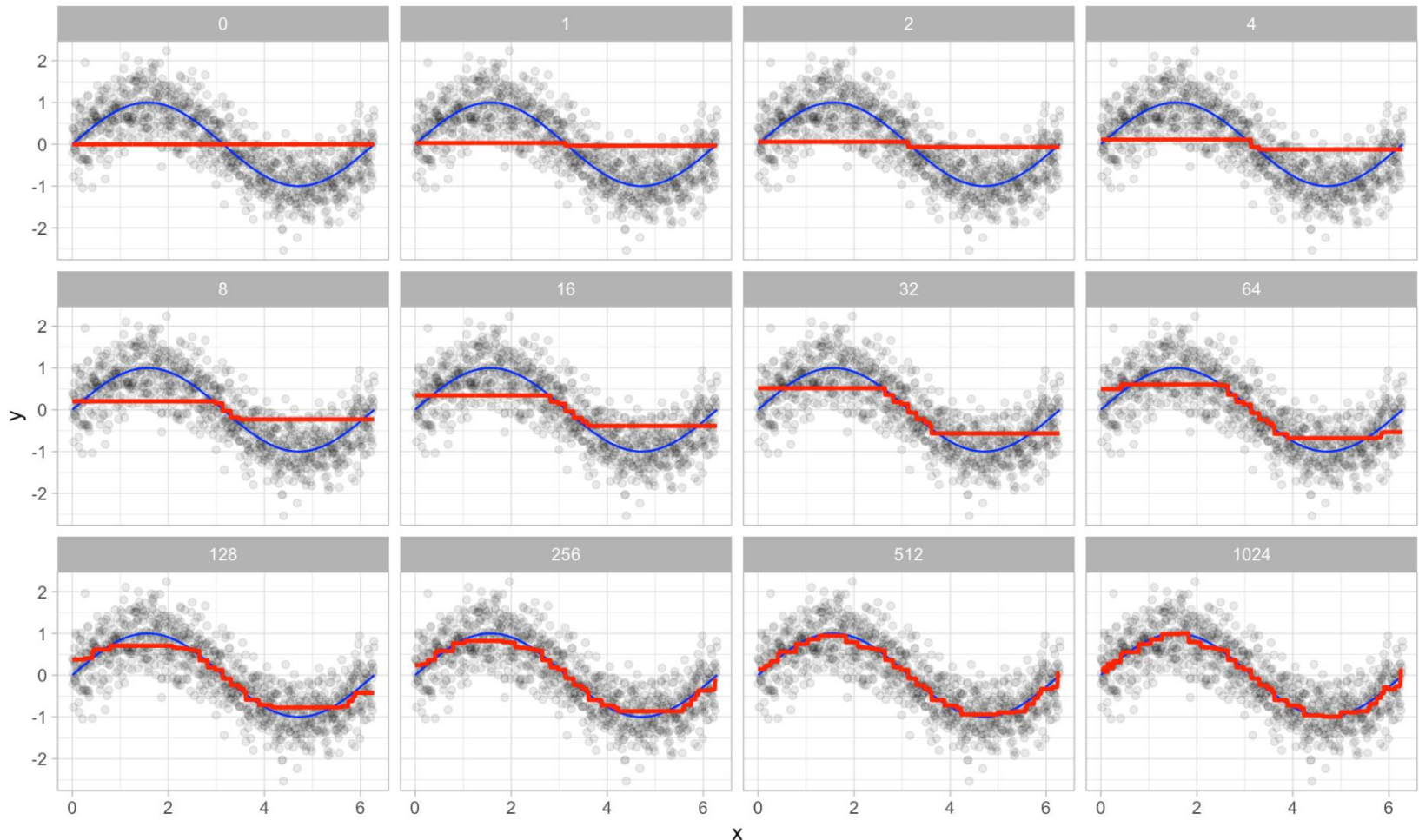
- GBM Regression Example 2



<https://www.quora.com/How-would-you-explain-gradient-boosting-machine-learning-technique-in-no-more-than-300-words-to-non-science-major-college-students>

Gradient Boosting Machine: GBM

- GBM Regression Example 3



<https://docs.paperspace.com/machine-learning/wiki/gradient-boosting>

Gradient Boosting Machine: GBM

- Gradient Boosting: Algorithm

1. Initialize $f_0(x) = \arg \min_{\gamma} \sum_{i=1}^N L(y_i, \gamma)$.

2. For $m = 1$ to M :

- 2.1 For $i = 1, \dots, N$ compute

$$g_{im} = \left[\frac{\partial L(y_i, f(x_i))}{\partial f(x_i)} \right]_{f(x_i) = f_{m-1}(x_i)}$$

- 2.2 Fit a regression tree to the targets g_{im} giving terminal regions $R_{jm}, j = 1, \dots, J_m$.

- 2.3 For $j = 1, \dots, J_m$ compute

$$\gamma_{jm} = \arg \min_{\gamma} \sum_{x_i \in R_{jm}} L(y_i, f_{m-1}(x_i) + \gamma)$$

- 2.4 Update $f_m(x) = f_{m-1}(x) + \sum_{j=1}^{J_m} \gamma_{jm} I(x \in R_{jm})$

3. Output $\hat{f}(x) = f_M(x)$.

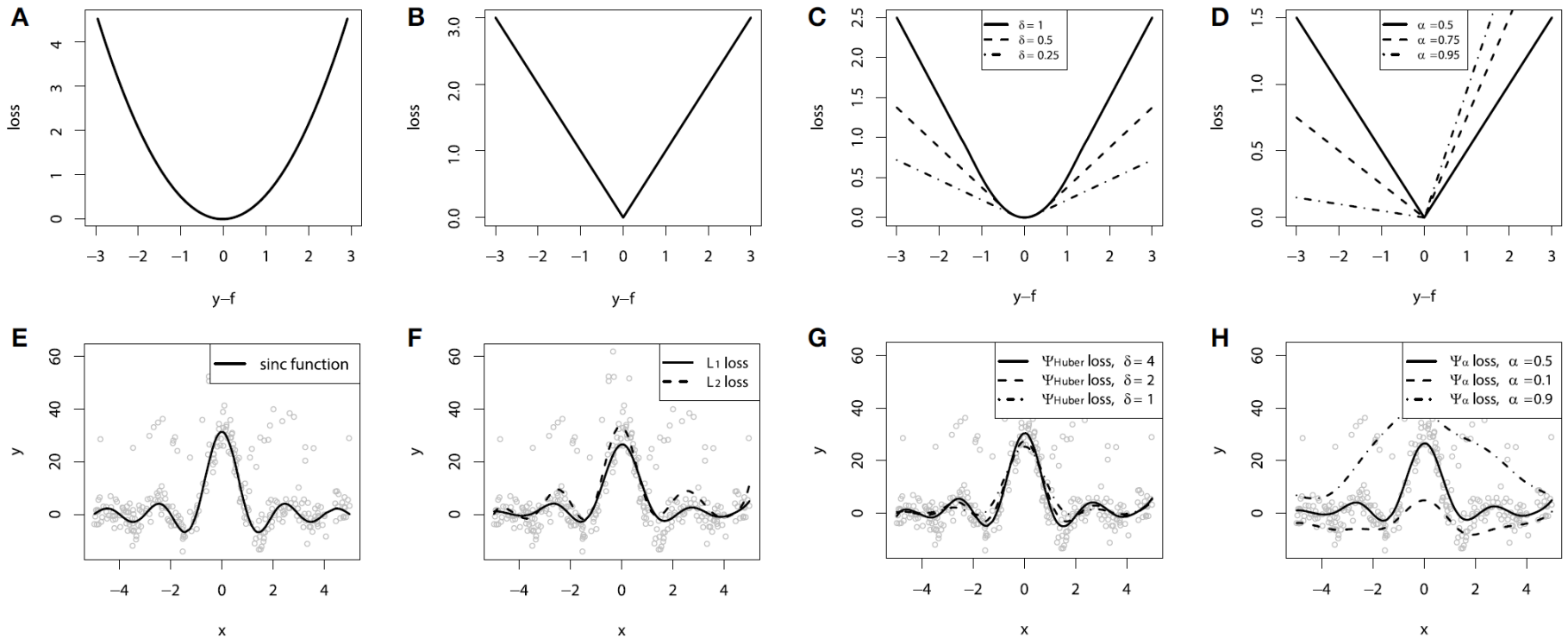
Gradient Boosting Machine: GBM

- Loss Functions for Regression

Loss Function	Formula
Squared loss (L_2)	$\Psi(y, f)_{L_2} = \frac{1}{2}(y - f)^2$
Absolute loss (L_1)	$\Psi(y, f)_{L_1} = y - f $
Huber loss	$\Psi(y, f)_{\text{Huber}, \delta} = \begin{cases} \frac{1}{2}(y - f)^2 & y - f \leq \delta \\ \delta(y - f - \delta/2) & y - f > \delta \end{cases}$
Quantile loss	$\Psi(y, f)_\alpha = \begin{cases} (1 - \alpha) y - f & y - f \leq 0 \\ \alpha y - f & y - f > 0 \end{cases}$

Gradient Boosting Machine: GBM

- Loss Functions for Regression



Gradient Boosting Machine: GBM

- Loss Functions for Classification

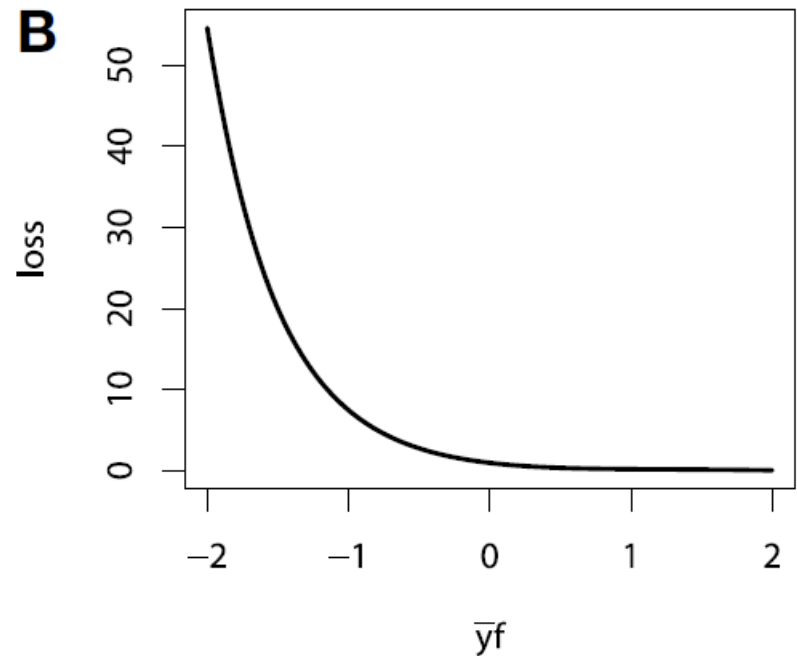
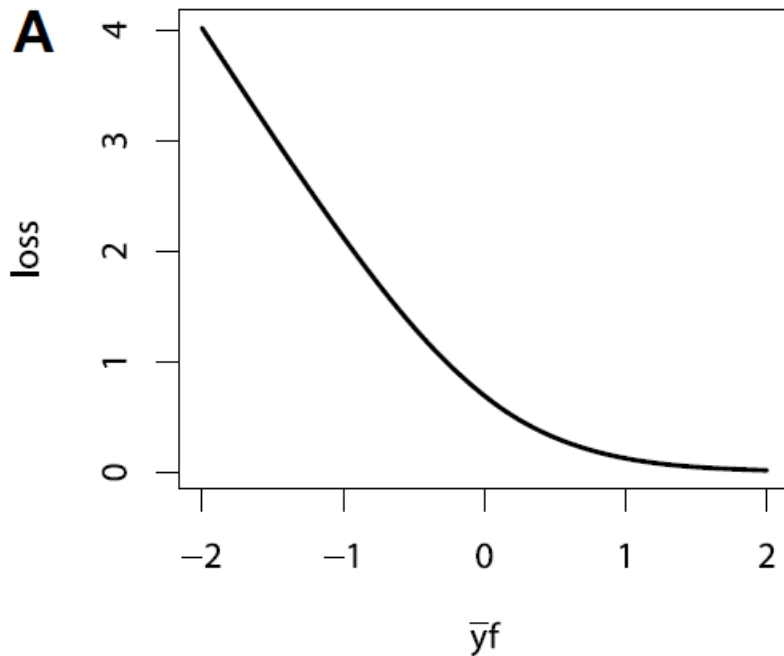
Loss Function	Formula
Bernoulli loss	$\Psi(y, f)_{\text{Bern}} = \log(1 + \exp(-2\bar{y}f))$
Adaboost loss	$\Psi(y, f)_{\text{Ada}} = \exp(-\bar{y}f)$

(Note)

In binary classification, the target is usually defined by $y \in \{0,1\}$, but here we define $\bar{y} = 2y - 1$ so that $\bar{y} \in \{-1,1\}$

Gradient Boosting Machine: GBM

- Loss Functions for Classification



(A) Bernoulli loss function. **(B)** Adaboost loss function.

Gradient Boosting Machine: GBM

- Overfitting problem in GBM

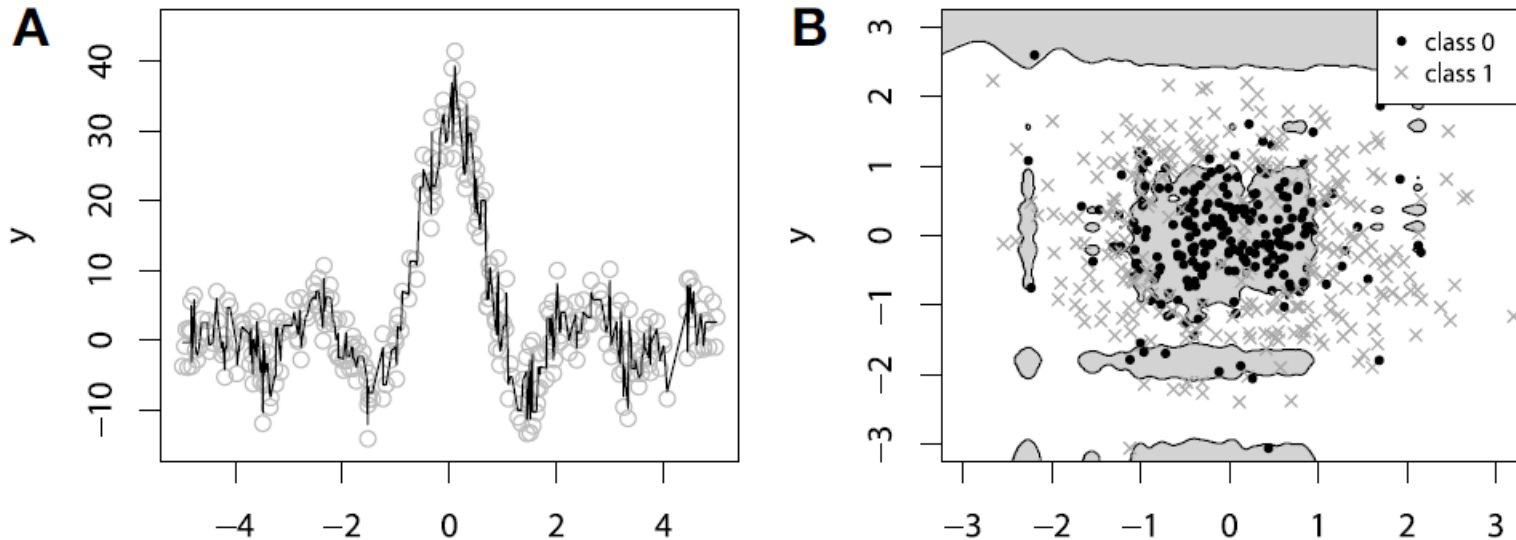


FIGURE 4 | Examples of overfitting in GBMs on: (A) regression task; (B) classification task. Demonstration of fitting a decision-tree GBM to a noisy $\text{sinc}(x)$ data: **(C)** $M = 100$, $\lambda = 1$; **(D)** $M = 1000$, $\lambda = 1$; **(E)** $M = 100$, $\lambda = 0.1$; **(F)** $M = 1000$, $\lambda = 0.1$.

Gradient Boosting Machine: GBM

- Regularization

- ✓ Subsampling

- At each learning iteration, only a random part of the training data is used to fit a consecutive base-learner.
 - The training data is typically sampled without replacement, but bagging can be also acceptable.

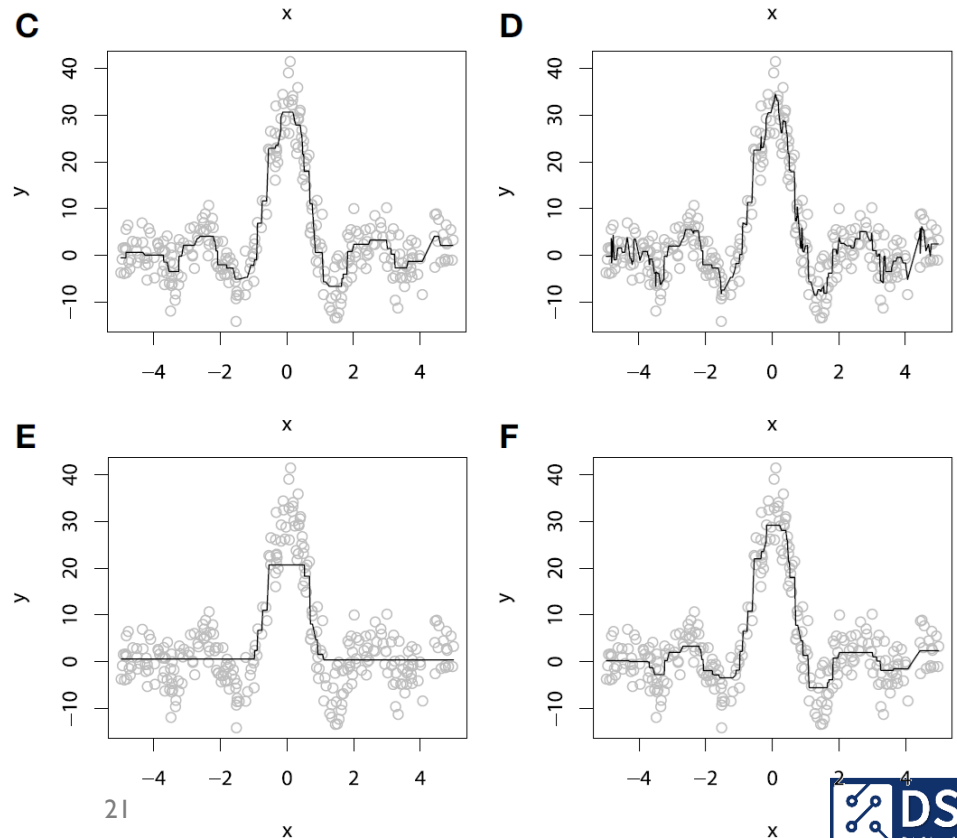
Gradient Boosting Machine: GBM

- Regularization

- ✓ Shrinkage

- Used for reducing/shrinking the impact of each additional fitted base-learners.
 - Better to improve a model by taking many small steps than by taking fewer large steps.

$$\hat{f}_t \leftarrow \hat{f}_{t-1} + \lambda \rho_t h(x, \theta_t)$$



Gradient Boosting Machine: GBM

- Regularization

- ✓ Early Stopping

- Use the validation error

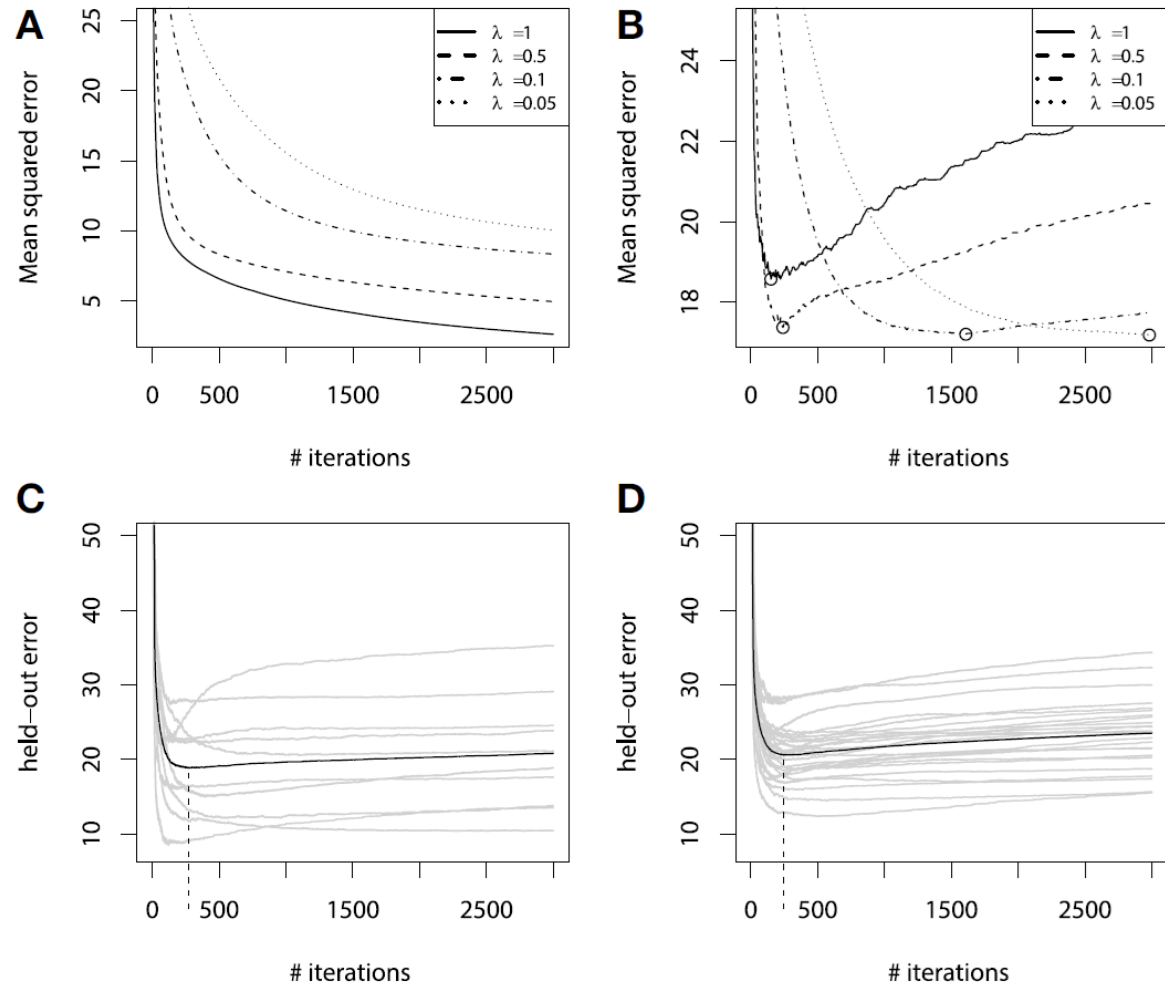


FIGURE 5 | Error curves for GBM fitting on $\text{sinc}(x)$ data: (A) training set error; (B) validation set error. Error curves for learning simulations and number of base-learners M estimation: (C) error curves for cross-validation; (D) error curves for bootstrap estimates.

Gradient Boosting Machine: GBM

- Variable Importance in Tree-based Gradient Boosting

- ✓ $Influence_j(T)$: importance of the variable j in a single tree T .

- ✓ Assume that there are L terminal nodes $\rightarrow L - 1$ splits.

$$Influence_j(T) = \sum_{i=1}^{L-1} (IG_i \times \mathbf{1}(S_i = j))$$

- ✓ Variable importance of Gradient boosting

$$Influence_j = \frac{1}{M} \sum_{k=1}^M Influence_j(T_k)$$

