# DDA4210/AIR6002 Advanced Machine Learning
## Lecture 02 Advanced Ensemble Learning

Tongxin Li

School of Data Science, CUHK-Shenzhen

Spring 2024

# Ensemble Learning

- Ensemble Learning is the process where multiple machine learning models are combined to get better results
  - Example: Suppose classifier 1 predicts Class A, classifier 2 predicts Class B, and classifier 3 predicts Class B, then the final prediction result is Class B. $Aggregate(1, 2, 3) \rightarrow B$
- Ensemble learning methods achieved SOTA performance in many real cases (e.g. Kaggle competitions, Netflix competition, and many categorical datasets)

Outlook. (EL).

- Works well for
  1. ranking
  2. regression.
  3. clustering. etc.

- Image

# Ensemble Learning

- Ensemble Learning is the process where multiple machine learning models are combined to get better results
  - Example: Suppose classifier 1 predicts Class A, classifier 2 predicts Class B, and classifier 3 predicts Class B, then the final prediction result is Class B.

- Ensemble learning methods achieved SOTA performance in many real cases (e.g. Kaggle competitions, Netflix competition, and many categorical datasets)
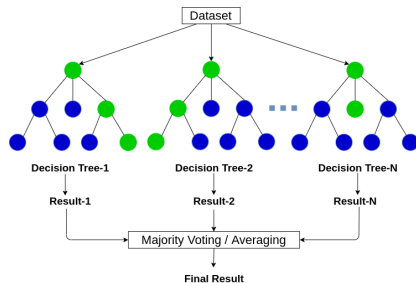
- What have we learned in basic machine learning courses?
  - Bagging (bootstrap aggregation)
  - Random forest (right figure)
  - * Connection and difference?
  - * Strength and weakness?

- Strength of RF
  - Do not require a lot of tuning.
  - Typically very accurate.
  - Handle heterogeneous features well.
  - Implicitly select the most relevant features. Why?
- Weakness of RF  $\nearrow^{vs}$ DT
  - Less interpretable, slower to train (but parallelizable)
  - Do not work well on high dimensional sparse data (e.g. text)

Q: What are the key hyper parameters in a RF?

1. # of trees.

2. depth of a tree.

3. p: = # of features related to each tree.

Example of heterogeneous data.

| Students ID | Gender | Height | Age |
|---|---|---|---|
| 1 | | | |
| 2 | | | |
| 3 | | | |
| ⋮ | | | |

# Strength and weakness of random forest (RF)

- Strength of RF
    - Do not require a lot of tuning.
    - Typically very accurate.
    - Handle heterogeneous features well.
    - Implicitly select the most relevant features. Why?
- Weakness of RF
    - Less interpretable, slower to train (but parallelizable)
    - Do not work well on high dimensional sparse data (e.g. text)
- The view of bias-variance trade-off for bagging and RF
    - Recall: Expected test error = Bias$^2$ + Variance + Noise$^2$
    - Bagging reduces variance by averaging
    - Bagging has little effect on bias
    - How can we reduce bias?

*Can <u>weak learners</u> be combined to generate a strong learner with low bias?*

—-Michael Kearns,1988

- "weak learner" (also called base learner): a learner (e.g. classifier, predictor, etc) that performs relatively poorly–its accuracy is above chance. $0 \cdot 55$
  - e.g., shallow decision trees, small neural networks
- "strong learner": a learner (e.g. classifier, predictor, etc) that achieves arbitrarily good performance, much better than random guessing. $0 \cdot 95$

*Can <u>weak learners</u> be combined to generate a strong learner with low bias?*
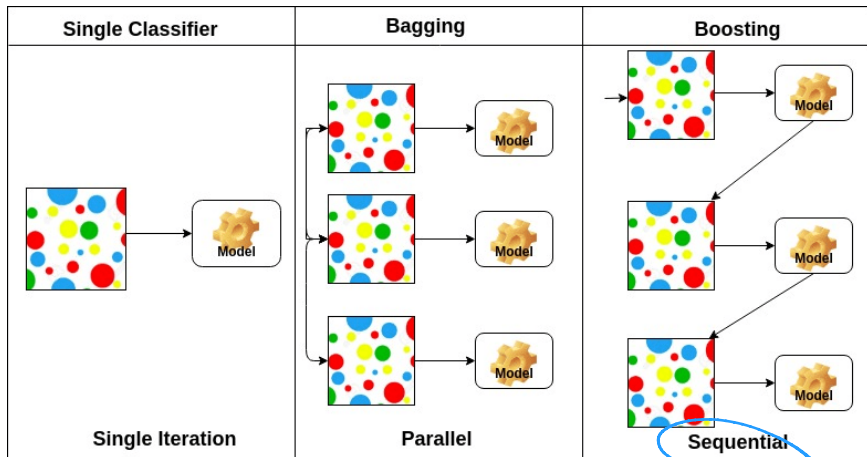
—-Michael Kearns,1988

- "weak learner" (also called base learner): a learner (e.g. classifier, predictor, etc) that performs relatively poorly–its accuracy is above chance.
    - e.g., shallow decision trees, small neural networks
- "strong learner": a learner (e.g. classifier, predictor, etc) that achieves arbitrarily good performance, much better than random guessing.
- Boosting can reduce bias!

    In machine learning, boosting is an ensemble meta-algorithm for primarily reducing bias, also variance in supervised learning, and a family of machine learning algorithms that convert weak learners to strong ones.—Wikipedia

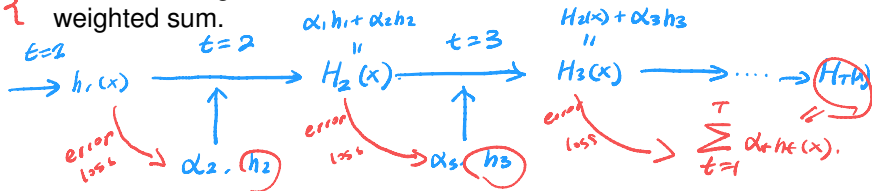Single classifier vs bagging vs boosting

# Boosting: main idea

- Like bagging, boosting is a general approach that can be applied to many machine learning methods for regression or classification
- Main idea of boosting $h_1(x), h_2(x), \dots, h_T(x)$ ← *weak learners*
  - Ensemble model (weighted sum of weak learners): *TBD*

*Idea: train $h_1(x) \dots - h_T(x)$ sequentially.*

*Training:*

$$strong\ learner \leftarrow H_T(\mathbf{x}) = \sum_{t=1}^{T} \alpha_t h_t(\mathbf{x}) \tag{1}$$

- This ensemble model is built in an iterative fashion.
- In iteration $t$ we add the weak learner $\alpha_t h_t(\mathbf{x})$ to the ensemble.
- At the test stage we evaluate all weak learners and return the weighted sum.

*Testing*

$\alpha_1 h_1 + \alpha_2 h_2$           $H_2(x) + \alpha_3 h_3$

$t=1$          $t=2$                          $t=3$

$\rightarrow h_1(x) \longrightarrow H_2(x) \longrightarrow H_3(x) \longrightarrow \cdots \rightarrow H_T(x)$

*error loss* $\rightarrow \alpha_2, h_2$   *error loss* $\rightarrow \alpha_3, h_3$   *error loss* $\rightarrow \sum_{t=1}^{T} \alpha_t h_t(x).$
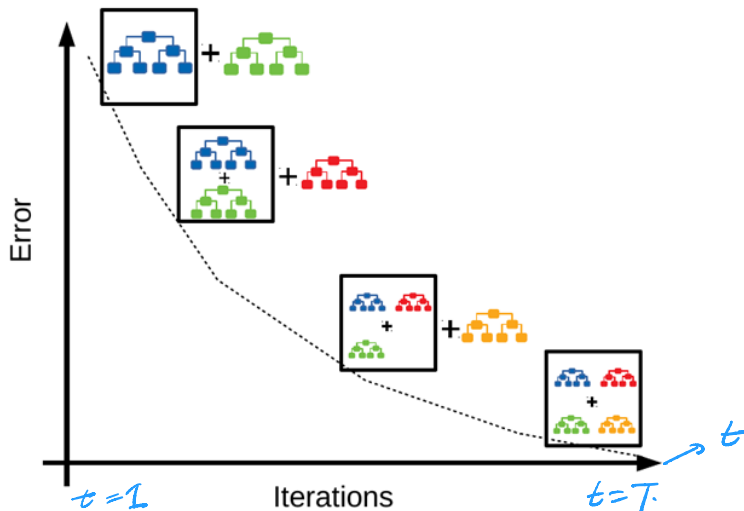
# Boosting: main idea

- Like bagging, boosting is a general approach that can be applied to many machine learning methods for regression or classification
- Main idea of boosting
  - Ensemble model (weighted sum of weak learners):

$$H_T(\mathbf{x}) = \sum_{t=1}^{T} \alpha_t h_t(\mathbf{x}) \tag{1}$$

  - This ensemble model is built in an iterative fashion.
  - In iteration $t$ we add the weak learner $\alpha_t h_t(\mathbf{x})$ to the ensemble.
  - At the test stage we evaluate all weak learners and return the weighted sum.
- Comparing boosting with bagging (suppose $h$ is a tree)
  - Bagging constructs all trees independently
  - Boosting constructs all trees sequentially

An intuitive example

- Ensemble model: $H_T(\mathbf{x}) = \sum_{t=1}^{T} \alpha_t h_t(\mathbf{x})$; data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^{n}$
- Details of constructing the ensemble model
  - Very similar to gradient descent.
  - However, instead of updating the model parameters in each iteration, we add functions to our ensemble.

Goal : Find / optimize $\{\alpha_t, h_t\}_{t=1}^{T}$ sequentially !

# Boosting: main idea

- Ensemble model: $H_T(\mathbf{x}) = \sum_{t=1}^{T} \alpha_t h_t(\mathbf{x})$; data $\mathcal{D} = \{(\mathbf{x}_i, y_i)\}_{i=1}^{n}$
- Details of constructing the ensemble model
  - Very similar to gradient descent.
  - However, instead of updating the model parameters in each iteration, we add functions to our ensemble.
  - Let $\ell$ denote a loss function and write    *e.g.* $\ell(x,y) = (x-y)^2$

    *square loss. Then*

$$\mathcal{L}(H) := \frac{1}{n} \sum_{i=1}^{n} \ell(H(\mathbf{x}_i), y_i) \qquad \mathcal{L}(H) = \mathcal{L}_{MSE}(H) \quad (2)$$

    $\downarrow$

    *mean square error*

    *Current time $t$. I have $\{h_\tau\}_{\tau=1}^{t}$.*

  - We compute $h_{t+1}$ (as well as $\alpha_{t+1}$) via

    *At time $t$, I know $H_t$.* $\longrightarrow$ *current model.*

$$h_{t+1}, \alpha_{t+1} = \underset{h \in \mathbb{H},\ \alpha}{\operatorname{argmin}} \mathcal{L}(H_t + \alpha h) \qquad (3)$$

    $\longrightarrow$ *I can customize*

    *$\alpha_{t+1} \cdot h_{t+1}$*

    - $\mathbb{H}$: hypothesis set
    - Once (3) is solved, we add $h_{t+1}$ to our ensemble, i.e.,

      $H_{t+1} := H_t + \alpha_{t+1} h_{t+1}$.
    - How to solve (3)?

- Taylor approximation on $\mathcal{L}(H + \alpha h)$
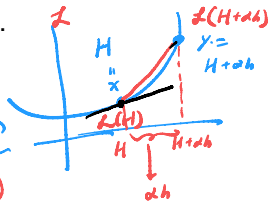
$$\mathcal{L}(H + \alpha h) \approx \mathcal{L}(H) + \alpha \langle \nabla \mathcal{L}(H), h \rangle \tag{4}$$

* $\alpha$ should be small enough; suppose $\alpha$ is fixed.

Remark:

1. We only focus on optimizing $h$.
   ($\alpha$ is more like a step size)

2. The learning rate (weight) $\alpha$ can also be optimized similarly (details omitted in gradient boosting).

# Gradient Boosting

- Taylor approximation on $\mathcal{L}(H + \alpha h)$

  *(doesn't depend on $h$)*

$$\mathcal{L}(H + \alpha h) \approx \boxed{\mathcal{L}(H)} + \alpha \langle \nabla \mathcal{L}(H), h \rangle \tag{4}$$

  * $\alpha$ should be small enough; suppose $\alpha$ is fixed.   *(fixed.)*

- Find $h$ via

$$
\begin{aligned}
h &= \operatorname{argmin}_{h \in \mathbb{H}} \mathcal{L}(H + \alpha h) \\
&\approx \operatorname{argmin}_{h \in \mathbb{H}} \langle \nabla \mathcal{L}(H), h \rangle \quad \text{\# of data samples.} \\
&= \operatorname{argmin}_{h \in \mathbb{H}} \sum_{i=1}^{n} \frac{\partial \mathcal{L}}{\partial [H(\mathbf{x}_i)]} h(\mathbf{x}_i) \quad \text{(Gradient Boosting)}
\end{aligned} \tag{5}
$$

$$
\nabla \mathcal{L}(H) = \begin{bmatrix} \dfrac{\partial \mathcal{L}}{\partial [H(x_1)]} \\ \vdots \\ \dfrac{\partial \mathcal{L}}{\partial [H(x_n)]} \end{bmatrix} \qquad h = \begin{bmatrix} h(x_1) \\ \vdots \\ h(x_n) \end{bmatrix}
$$

*column vector*

- Taylor approximation on $\mathcal{L}(H + \alpha h)$

  *(annotation: next loss)*    *(annotation: current loss)*    *(annotation: $< 0$)*

$$\mathcal{L}(H + \alpha h) \approx \mathcal{L}(H) + \alpha \langle \nabla \mathcal{L}(H), h \rangle \tag{4}$$

  * $\alpha$ should be small enough; suppose $\alpha$ is fixed.

- Find $h$ via

  *(annotation: $H_{t+i} = H + \alpha h$ new model.)*
  *(annotation: $H_t := H$ old model)*

$$
\begin{aligned}
h &= \arg\min_{h \in \mathbb{H}} \mathcal{L}(H + \alpha h) \\
&\approx \arg\min_{h \in \mathbb{H}} \langle \nabla \mathcal{L}(H), h \rangle \\
&= \arg\min_{h \in \mathbb{H}} \sum_{i=1}^{n} \frac{\partial \mathcal{L}}{\partial [H(\mathbf{x}_i)]} h(\mathbf{x}_i) \quad \text{(Gradient Boosting)}
\end{aligned} \tag{5}
$$

  *(annotation: $\mathcal{L}(\alpha h + H) \lesssim \mathcal{L}(H)$.)*

- Thus we can do boosting if we have an algorithm to compute
  $h_{t+1} = \arg\min_{h \in \mathbb{H}} \sum_{i=1}^{n} q_i h(\mathbf{x}_i)$, where $q_i = \frac{\partial \mathcal{L}}{\partial [H_t(\mathbf{x}_i)]}$, irrelevant to $h$.

  *(annotation: $\mathcal{L}(H + \alpha h) \lesssim \mathcal{L}(H)$.)*

  - $q_i$ is the weight for the training data $i$.
  - $h$ does not need to be perfect.
  - We make progress as long as $\sum_{i=1}^{n} q_i h(\mathbf{x}_i) < 0$. Why?

# Gradient Boosting

- Suppose $\mathcal{F}$ is the objective function we want to minimize to find the parameters $\Theta$ for a model (e.g., linear regression, neural network).
- We compare gradient boosting with gradient descent.

| gradient boosting | gradient descent |
|---|---|
| $h = \text{argmin}_{h \in \mathbb{H}} \mathcal{L}(H + \alpha h)$ $\approx \text{argmin}_{h \in \mathbb{H}} \langle \nabla \mathcal{L}(H), h \rangle$ | $\theta = \text{argmin}_\theta \mathcal{F}(\Theta + \eta \theta)$ $\approx \text{argmin}_\theta \langle \nabla \mathcal{F}(\Theta), \theta \rangle$ |
| $h = -\nabla \mathcal{L}(H)$ (ideal) | $\theta = -\nabla \mathcal{F}(\Theta)$ |
| learn $h$ on $-\frac{\partial \mathcal{L}}{\partial [H(\mathbf{x}_i)]}, i = 1, \ldots, n$ | $\theta$ is the negative gradient |

a parametrised function (weak learner) in $\mathbb{H}$.
$$\left[ h: X \to Y \quad \begin{array}{l} x_i \in X \quad i = 1 \ldots n \\ y_i \in Y \quad i = 1 \ldots n \end{array} \right]$$

# Gradient Boosting

- Suppose $\mathcal{F}$ is the objective function we want to minimize to find the parameters $\Theta$ for a model (e.g., linear regression, neural network).
- We compare gradient boosting with gradient descent.

| gradient boosting | gradient descent |
|---|---|
| $h = \text{argmin}_{h \in \mathbb{H}} \mathcal{L}(H + \alpha h)$ $\approx \text{argmin}_{h \in \mathbb{H}} \langle \nabla \mathcal{L}(H), h \rangle$ | $\theta = \text{argmin}_\theta \mathcal{F}(\Theta + \eta \theta)$ $\approx \text{argmin}_\theta \langle \nabla \mathcal{F}(\Theta), \theta \rangle$ |
| $h = -\nabla \mathcal{L}(H)$ (ideal) | $\theta = -\nabla \mathcal{F}(\Theta)$ |
| learn $h$ on $-\frac{\partial \mathcal{L}}{\partial [H(\mathbf{x}_i)]}$, $i = 1, \ldots, n$ | $\theta$ is the negative gradient |

- In gradient boosting, we require $h(\mathbf{x}_i) \approx -\frac{\partial \mathcal{L}}{\partial [H(\mathbf{x}_i)]}$ such that $\sum_{i=1}^{n} q_i h(\mathbf{x}_i) < 0$. For instance, when $h(\mathbf{x}_i) = -\frac{\partial \mathcal{L}}{\partial [H(\mathbf{x}_i)]}$, $\mathcal{L}(H + \alpha h) \approx \mathcal{L}(H) - \alpha \|\nabla \mathcal{L}(H)\|^2$, which ensures a decreasing $\mathcal{L}$.

# Generic algorithm of Gradient Boosting

*Big picture.*

Also called AnyBoost [Llew Mason et al. 2000]

---

**Input:** $\ell$, $\alpha$, $T$, $\{(\mathbf{x}_i, y_i)\}$, $\mathbb{H}$, $\mathbb{A}$(an algorithm)

1: $H_0 = 0$

2: **for** $t = 0 : T - 1$ **do**

3:      $q_i = \frac{\partial \mathcal{L}}{\partial H_t(\mathbf{x}_i)}$, $i = 1, 2, \ldots, n$    → Taylor apr···

                                   gradient boosting.

4:      $h_{t+1} = \operatorname{argmin}_{h \in \mathbb{H}} \sum_{i=1}^{n} q_i h(\mathbf{x}_i)$ via Algorithm $\mathbb{A}$

5:      **if** $\sum_{i=1}^{n} q_i h_{t+1}(\mathbf{x}_i) < 0$ **then**

6:          $H_{t+1} = H_t + \alpha h_{t+1}$    ($\alpha$ may also be optimized)

7:      **else**

8:          return ($H_t$)

9:      **end if**

10: **end for**

**Output:** $H_T$

---

# A brief history of Gradient Boosting

- [Schapire,1989] first provable boosting algorithm
- [Freund, 1990] "optimal" algorithm that "boosts by majority"
- [Drucker, Schapire & Simard, 1992] first experiments using boosting, limited by practical drawbacks
- [Freund et al., 1996, Freund and Schapire, 1997] invent AdaBoost (to be introduced later), the first successful boosting algorithm
- [Breiman et al., 1998, Breiman, 1999] formulate AdaBoost as gradient descent with a special loss function
- [Friedman et al., 2000, Friedman, 2001] generalize Adaboost to Gradient Boosting in order to handle a variety of loss functions

- Solve $h_{t+1} = \operatorname*{argmin}_{h \in \mathbb{H}} \sum_{i=1}^{n} q_i h(\mathbf{x}_i)$, where $q_i = \frac{\partial \mathcal{L}}{\partial [H(\mathbf{x}_i)]}$

# Gradient Boosting for Regression

- Solve $h_{t+1} = \text{argmin}_{h \in \mathbb{H}} \sum_{i=1}^{n} q_i h(\mathbf{x}_i)$, where $q_i = \frac{\partial \mathcal{L}}{\partial [H(\mathbf{x}_i)]}$
- Let $\sum_{i=1}^{n} h^2(\mathbf{x}_i) = \text{constant}$ (we can always normalize the predictions) and replace $q_i$ with $-2r_i$. We have

*Original Data*
$\{x_i, y_i\}_{i=1}^{n}$

$\downarrow$
$H_0$
$\downarrow$
*weighted Data.*
$\downarrow$
$h_1$
$\downarrow$
$H_1$

$$
\begin{aligned}
h_{t+1} =& \text{argmin}_{h \in \mathbb{H}} \sum_{i=1}^{n} q_i h(\mathbf{x}_i) \\
=& \text{argmin}_{h \in \mathbb{H}} -2 \sum_{i=1}^{n} r_i h(\mathbf{x}_i) \\
=& \text{argmin}_{h \in \mathbb{H}} \sum_{i=1}^{n} \left( r_i^2 - 2 r_i h(\mathbf{x}_i) + (h(\mathbf{x}_i))^2 \right) \\
=& \text{argmin}_{h \in \mathbb{H}} \sum_{i=1}^{n} (h(\mathbf{x}_i) - r_i)^2
\end{aligned}
\tag{6}
$$

$= \text{train } h \text{ on new weighted data } \{(x_i, r_i)\}_{i=1}^{n}.$

- We train $h_{t+1}$ to predict $r_i$, which are from the old model $H_t$.

*weighted Data.*

Special Case: MSE

$$\mathcal{L}(H) = \sum_{i=1}^{n} \left( H(x_i) - y_i \right)^2$$

what is $\dfrac{\partial \ell}{\partial [H(x_i)]} = 2(H(x_i) - y_i)$

- Note that $h_{t+1} = \operatorname{argmin}_{h \in \mathbb{H}} \sum_{i=1}^{n} (h(\mathbf{x}_i) - r_i)^2$, where $r_i = -q_i/2$.
- Let $\ell$ be the squared loss, i.e., $\mathcal{L}(H) = \sum_{i=1}^{n}(H(\mathbf{x}_i) - y_i)^2$, then

$$r_i = -\frac{1}{2} \times \frac{\partial \mathcal{L}}{\partial [H(\mathbf{x}_i)]} = y_i - H(\mathbf{x}_i)$$

called "pseudo-residuals"

- $r_i$ are just the residuals given by the old model.



Data
$x \cdot y$ → $H_0$ → $h_1$ → $H_1 = H_0 + \alpha \cdot h_1$ → $h_2$ → $H_2$ → ...

error 1
$r_1$

error 2
$r_2$

$r_n$

# Gradient Boosting for Regression

- Pseudo code [A special instance using MSE]

---

**Require:** $\ell$, $T$, $\alpha$, $\{(\mathbf{x}_i, y_i)\}_{i=1}^n$, $\mathbb{H}$

  1: $H = 0$

  *step size (can be optimized. later!)*

  2: **for** $t = 0 : T - 1$ **do**

  3:      $r_i = y_i - H(\mathbf{x}_i)$, $i = 1, 2, \ldots, n$

           *"Trade-off of selecting the step-size $\alpha$"*

  4:      $h = \operatorname{argmin}_{h \in \mathbb{H}} \sum_{i=1}^n (h(\mathbf{x}_i) - r_i)^2$

  5:      $H \leftarrow H + \alpha h$

  6: **end for**

**Ensure:** $H$

---

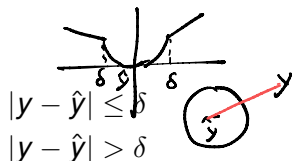- Characteristic: squared loss and fixed weight $\alpha$

- Square loss is easy to deal with mathematically but not robust to outliers.
- Absolute loss (more robust to outliers)

$$\ell(y, \hat{y}) = |y - \hat{y}|$$

$y_1 \quad h(x_1)$
$\downarrow$
$outlier$

- Square loss is easy to deal with mathematically but not robust to outliers.
- Absolute loss (more robust to outliers)

$$\ell(y, \hat{y}) = |y - \hat{y}|$$

"sub-gradient"

GBR with absolute loss:
- The gradient is $q_i = \frac{\partial \mathcal{L}}{\partial H(\mathbf{x}_i)} = -\text{sign}(y_i - H(\mathbf{x}_i))$.
- Then fit $h$ on $-q_i$, $i = 1, 2, \ldots, n$. (no longer the residuals, different from using the squared loss)

## Gradient Boosting for Regression

- Square loss is easy to deal with mathematically but not robust to outliers.
- Absolute loss (more robust to outliers)

$$\ell(y, \hat{y}) = |y - \hat{y}|$$

GBR with absolute loss:
- The gradient is $q_i = \frac{\partial \mathcal{L}}{\partial H(\mathbf{x}_i)} = -\text{sign}(y_i - H(\mathbf{x}_i))$.
- Then fit $h$ on $-q_i$, $i = 1, 2, \ldots, n$. (no longer the residuals, different from using the squared loss)
- Huber loss (more robust to outliers)

$$\ell(y, \hat{y}) = \begin{cases} \frac{1}{2}(y - \hat{y})^2 & |y - \hat{y}| \leq \delta \\ \delta(|y - \hat{y}| - \delta/2) & |y - \hat{y}| > \delta \end{cases}$$

The gradient is

$$\frac{\partial \mathcal{L}}{\partial H(\mathbf{x}_i)} = \begin{cases} -(y_i - H(\mathbf{x}_i)) & |y_i - H(\mathbf{x}_i)| \leq \delta \\ -\delta \text{sign}(y_i - H(\mathbf{x}_i)) & |y_i - H(\mathbf{x}_i)| > \delta \end{cases}$$

GBR with $h$ being a decision tree
$\longrightarrow$ Gradient Boosted Regression Tree (GBRT)

Remarks:
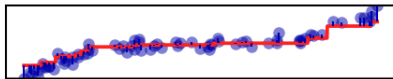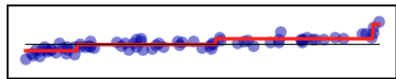
1. Regression w.r.t.
   $\ell_2$-norm.
   $$h_{t+1} = \min_{h \in \mathcal{H}} \sum_{i=1}^{n} (h(x_i) - r_i)^2$$



Residual prediction step 1
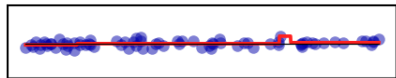
Total prediction step 1

$H(x_i)$
$\forall i = 1, \dots, 100$

Residual prediction step 4

Total prediction step 4

Residual prediction step 10

Total prediction step 10

- Predicting probability of each of K classes, namely

  soft max.

  $$p_k(\mathbf{x}) = \frac{\exp\left(h^{(k)}(\mathbf{x})\right)}{\sum_{c=1}^{K} \exp\left(h^{(c)}(\mathbf{x})\right)} \triangleq \hat{y}^{(k)}, \quad k = 1, 2, \ldots, K.$$

- Consider the loss $\mathcal{L}(H) = \sum_{i=1}^{n} \ell(\mathbf{y}_i, \hat{\mathbf{y}}_i)$ (e.g. multi-class cross-entropy or KL divergence)

  $$\text{e.g.} \quad y_1 = \begin{bmatrix} 0 \\ 1 \\ 0 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \Big\} K \qquad \hat{y}_1 = \begin{bmatrix} 0 \\ 0.9 \\ 0.05 \\ 0 \\ \vdots \\ 0.05 \end{bmatrix} \Big\} K$$

## Gradient Boosting for Classification

- Predicting probability of each of K classes, namely
  $$p_k(\mathbf{x}) = \frac{\exp\big(h^{(k)}(\mathbf{x})\big)}{\sum_{c=1}^{K} \exp\big(h^{(c)}(\mathbf{x})\big)} \triangleq \hat{y}^{(k)}, \quad k = 1, 2, \ldots, K.$$

- Consider the loss $\mathcal{L}(H) = \sum_{i=1}^{n} \ell(\mathbf{y}_i, \hat{\mathbf{y}}_i)$ (e.g. multi-class cross-entropy or KL divergence)

- Model initialization $H^{(1)}, H^{(2)}, \ldots, H^{(K)}$ and iterate until converge or reach maximum $T$

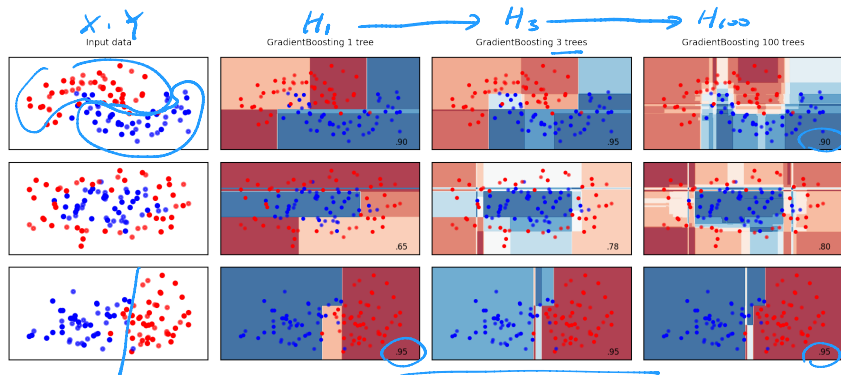  1. Calculate negative gradients for every class:

  $$-g_k(\mathbf{x}_i) = -\frac{\partial \mathcal{L}}{\partial [H^{(k)}(\mathbf{x}_i)]}, \ i = 1, 2, \ldots, n, \ k = 1, 2, \ldots, K$$

  2. Fit $h^{(k)}$ to the negative gradients $-g_k(\mathbf{x}_i)$, $k = 1, 2, \ldots, K$.
  3. Update the models via $H^{(k)} \leftarrow H^{(k)} + \alpha h^{(k)}$, $k = 1, 2, \ldots, K$.

GBC with *h* being a decision tree
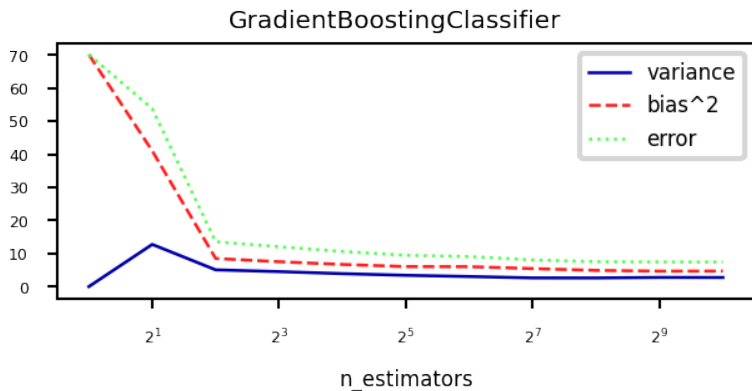$\longrightarrow$ Gradient Boosted Classification Tree (GBCT)



https://ml-course.github.io/master/notebooks/05%20-%20Ensemble%20Learning.html

# Numerical example of GBC

- Bias-Variance analysis
  - Gradient Boosting reduces bias (and a little variance)
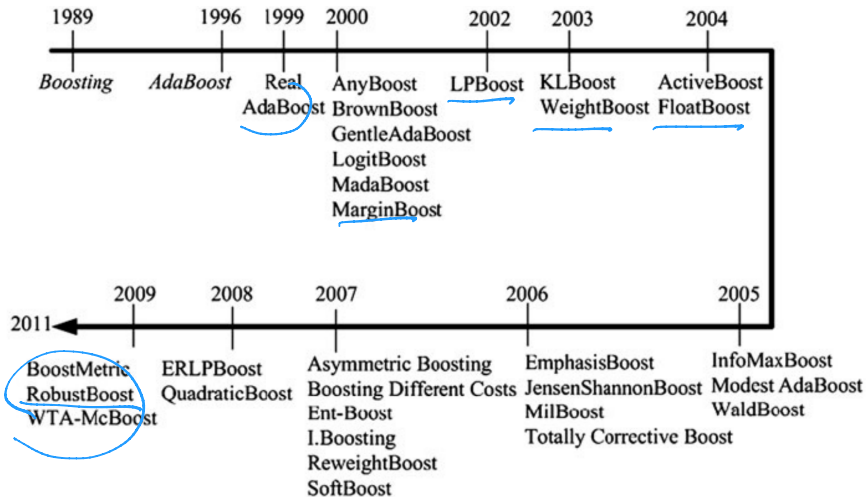  - Boosting too much will eventually increase variance



GradientBoostingClassifier

https://ml-course.github.io/master/notebooks/05%20-%20Ensemble%20Learning.html

# History of Boosting



Image from A.J. Ferreira and M.A.T. Figueiredo, in Ensemble Machine Learning: Methods and Applications. 2012.

*consider the update of the learning rate $\alpha$.*

$X \cdot Y$     $X := (x_1, \ldots, x_n)$
$Y := (y_1, \ldots, y_n)$     $y_i \in \{-1, +1\}$
↳ classification.

- AdaBoost (a special case of gradient boosting) uses the exponential loss, i.e.,

*strong model.*     *outputs :*
$H(x_1) \cdots H(x_n)$.

$$\mathcal{L}(H) = \sum_{i=1}^{n} e^{-y_i H(\mathbf{x}_i)}$$

and learns $\alpha$ adaptively.

- The gradient is then $q_i = \frac{\partial \mathcal{L}}{\partial H(\mathbf{x}_i)} = -y_i e^{-y_i H(\mathbf{x}_i)}$.

$H(x_i) = 0.9$      $y_i = 1$      $e^{-0.9}$

$H(x_j) = 0.1$      $y_j = 1.$     $e^{-0.1}$

- AdaBoost (a special case of gradient boosting) uses the exponential loss, i.e.,

$$\mathcal{L}(H) = \sum_{i=1}^{n} e^{-y_i H(\mathbf{x}_i)}$$

and learns $\alpha$ adaptively.

- The gradient is then $q_i = \frac{\partial \mathcal{L}}{\partial H(\mathbf{x}_i)} = -y_i e^{-y_i H(\mathbf{x}_i)}$.

  *sample weight.*

- For convenience, let $\boxed{w_i = \frac{1}{Z} e^{-y_i H(\mathbf{x}_i)}}$, where $Z = \sum_{i=1}^{n} e^{-y_i H(\mathbf{x}_i)}$ (such that $\sum_{i=1}^{n} w_i = 1$). Then $w_i$ is the relative contribution of the training point $(\mathbf{x}_i, y_i)$ to the overall loss.

- We here consider a binary classification task, i.e., $y \in \{-1, +1\}$, and let $h(\mathbf{x}) \in \{-1, +1\}, \forall \mathbf{x}$.

  Q: How to get h based on $\{q_i\}_{i=1}^{n}$ ?

- Note that $\mathcal{L}(H) = \sum_{i=1}^{n} e^{-y_i H(\mathbf{x}_i)}$, $q_i = \frac{\partial \mathcal{L}}{\partial H(\mathbf{x}_i)} = -y_i e^{-y_i H(\mathbf{x}_i)}$,
  $w_i = \frac{1}{Z} e^{-y_i H(\mathbf{x}_i)}$, $Z = \sum_{i=1}^{n} e^{-y_i H(\mathbf{x}_i)}$ (a constant w.r.t $h$),
  $\sum_{i=1}^{n} w_i = 1$, $y \in \{-1, +1\}$, and $h(\mathbf{x}) \in \{-1, +1\}$.

- We have

*Taylor approximation.*

$$h = \operatorname{argmin}_{h \in \mathbb{H}} \sum_{i=1}^{n} q_i h(\mathbf{x}_i)$$

$q_i = -y_i w_i \cdot Z$

$$= \operatorname{argmin}_{h \in \mathbb{H}} -\sum_{i=1}^{n} w_i y_i h(\mathbf{x}_i) \qquad \begin{cases} 1 \cdot 1 = 1 \\ 1 \cdot -1 = -1 \\ -1 \cdot 1 = -1 \\ -1 \cdot -1 = 1 \end{cases} \tag{7}$$

$$= \operatorname{argmin}_{h \in \mathbb{H}} \sum_{i: h(\mathbf{x}_i) \neq y_i} w_i - \sum_{i: h(\mathbf{x}_i) = y_i} w_i \qquad (h(x) \in \{-1, +1\})$$

$(2) \cdot (\mathbf{w}) - 1 \le$

$\sum_{i=1}^{n} w_i = 1 \implies$

$$= \operatorname{argmin}_{h \in \mathbb{H}} \sum_{i: h(\mathbf{x}_i) \neq y_i} w_i$$

$\sum_{i: h(x_i) = y_i} w_i = 1 - \sum_{i: h(x_i) \neq y_i} w_i$

$(*)$

* The last equality holds because $\sum_{i=1}^{n} w_i = 1$.

$$w_i = \frac{1}{z} e^{-\overset{-1}{\underset{-1}{\tilde{y_i}}} \overset{+1}{\underset{-1}{\tilde{H(i)}}}} = \frac{1}{z} e^{1}$$

$$= \frac{1}{z} e^{-1}$$

- We have got

$$h = \mathrm{argmin}_{h \in \mathbb{H}} \sum_{i: h(\mathbf{x}_i) \neq y_i} w_i$$

$\epsilon := \sum_{i: h(\mathbf{x}_i) y_i = -1} w_i$ can be regarded as a weighted classification error, where $w_i = \frac{1}{z} e^{-y_i H(\mathbf{x}_i)}$.

Note that a miss-classified point by $H$ gets a larger weight.

## Adaptive Boosting (AdaBoost)

- Given $h$, we find $\alpha$ via
$$\alpha = \operatorname{argmin}_\alpha \mathcal{L}(H + \alpha h) = \operatorname{argmin}_\alpha \sum_{i=1}^n e^{-y_i(H(\mathbf{x}_i) + \alpha h(\mathbf{x}_i))} \qquad (8)$$

- Differentiate w.r.t $\alpha$ and equate with zero:
$$\sum_{i=1}^n y_i h(\mathbf{x}_i) e^{-(y_i H(\mathbf{x}_i) + \alpha y_i h(\mathbf{x}_i))} = 0.$$

# Adaptive Boosting (AdaBoost)

- Given $h$, we find $\alpha$ via

$$\alpha = \text{argmin}_\alpha \mathcal{L}(H + \alpha h) = \text{argmin}_\alpha \sum_{i=1}^n e^{-y_i(H(\mathbf{x}_i) + \alpha h(\mathbf{x}_i))} \quad (8)$$

- Differentiate w.r.t $\alpha$ and equate with zero:

$$\sum_{i=1}^n y_i h(\mathbf{x}_i) e^{-(y_i H(\mathbf{x}_i) + \alpha y_i h(\mathbf{x}_i))} = 0.$$

$$\begin{cases} w_i = \frac{1}{z} e^{-y_i h(x_i)} \\ \varepsilon := \sum_{i: h(x_i) \neq y_i} w_i \end{cases}$$

It follows that

$$\sum_{i: h(\mathbf{x}_i) y_i = 1} e^{-(y_i H(\mathbf{x}_i) + \alpha y_i h(\mathbf{x}_i))} - \sum_{i: h(\mathbf{x}_i) y_i = -1} e^{-(y_i H(\mathbf{x}_i) + \alpha y_i h(\mathbf{x}_i))} = 0$$
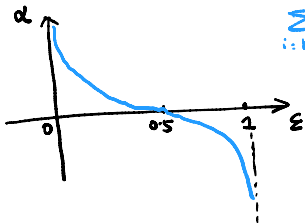
$$w_i \cdot z$$

$$\sum_{i: h(\mathbf{x}_i) y_i = 1} w_i e^{-\alpha} - \sum_{i: h(\mathbf{x}_i) y_i = -1} w_i e^{\alpha} = 0$$

We have $(1 - \epsilon)e^{-\alpha} - \epsilon e^{\alpha} = 0$, $e^{2\alpha} = \frac{1 - \epsilon}{\epsilon}$, and get

$$\alpha = \frac{1}{2} \ln \frac{1 - \epsilon}{\epsilon}. \quad (9)$$

$$\alpha = \frac{1}{2} \log \frac{1-\varepsilon}{\varepsilon}$$



$1 - \varepsilon :=$

$\sum_{i \; : \; h(x_i) = y_i} w_i$

$C$: the number of correctly predicted labels

$\varepsilon > 0.5 :=$ the new model is worse then random guessing.

## Adaptive Boosting (AdaBoost)

- Pseudo code of AdaBoost for binary classification

**Require:** $\ell, \{(\mathbf{x}_i, y_i)\}_{i=1}^{n}, \mathbb{H}$
1: $H_0 = 0$
2: $w_i = 1/n, \ i = 1, 2, , \ldots, n$
3: **for** $t = 0 : T - 1$ **do**
4:      $h_{t+1} = \mathrm{argmin}_{h \in \mathbb{H}} \sum_{i:h(\mathbf{x}_i) \neq y_i} w_i$
5:      $\epsilon = \sum_{i:h(\mathbf{x}_i) \neq y_i} w_i$
6:      **if** $\epsilon < 1/2$ **then**
7:          $\alpha = \frac{1}{2} \ln \frac{1-\epsilon}{\epsilon}$
8:          $H_{t+1} \leftarrow H_t + \alpha h_{t+1}$
9:          $w_i = \frac{1}{Z} \exp\left(-y_i H_{t+1}(\mathbf{x}_i)\right), \ i = 1, 2, \ldots, n$
10:      **else**
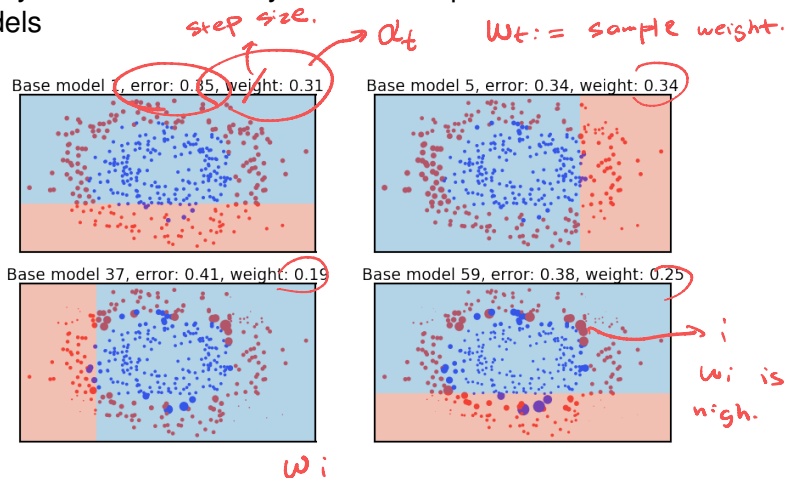11:          Return $H_t$
12:      **end if**
13: **end for**
**Ensure:** $H_T$

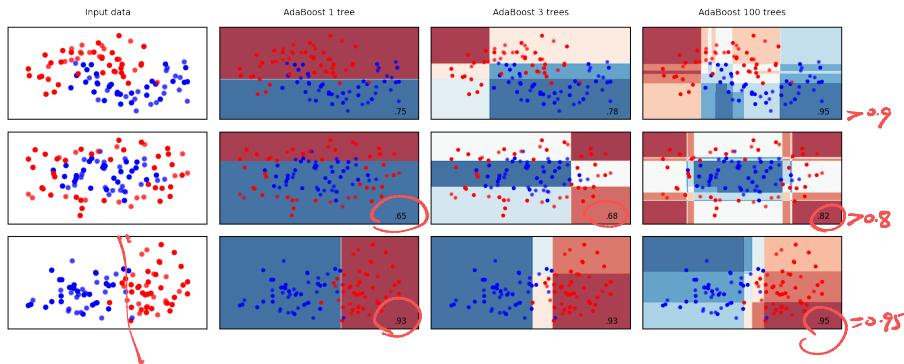- Binary classification on synthetic data: performance of the base models



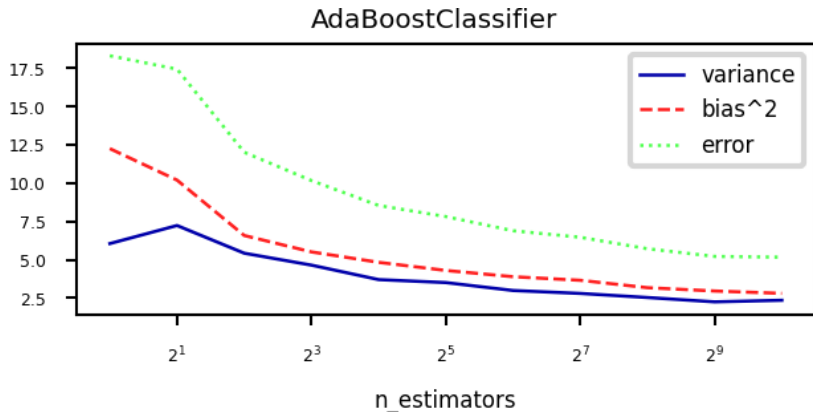The size of the dot indicates the weight of the data point.

https://ml-course.github.io/master/notebooks/05%20-%20Ensemble%20Learning.html

# Numerical example of AdaBoost

- Binary classification on synthetic data: performance of the ensemble classifier



https://ml-course.github.io/master/notebooks/05%20-%20Ensemble%20Learning.html

# Numerical example of AdaBoost

- Bias-Variance analysis
  - AdaBoost reduces bias (and a little variance)
  - Boosting too much will eventually increase variance



AdaBoostClassifier

https://ml-course.github.io/master/notebooks/05%20-%20Ensemble%20Learning.html

*special instance.*

- AdaBoost minimizes the exponential loss function that can make the algorithm sensitive to outliers. With Gradient Boosting, any differentiable loss function can be utilized. Gradient Boosting is more robust to outliers than AdaBoost.

- AdaBoost up-weights observations that were misclassified before. Gradient boosting identifies difficult observations by large residuals computed in the previous iterations. *$\ell_2$ - norm loss*

- AdaBoost was mainly designed for binary classification problems and can be utilized to boost the performance of decision trees. Gradient Boosting is used to solve the differentiable loss function problem. The technique can be used for both classification and regression problems.

- AdaBoost can be regarded as a special case of Gradient Boosting.

.

## Discussion on boosting

- Boosting is a great way to turn a week learner into a strong leaner.
- It defines a whole family of algorithms, including Gradient Boosting, AdaBoost, LogitBoost, and many others.
- Gradient boosted decision tree is very useful for learning to rank (to be introduced in future).
- AdaBoost is an extremely powerful algorithm, that turns any weak learner that can classify any weighted version of the training set with below 0.5 error into a strong learner whose training error decreases exponentially.

A few slides are adapted from https:
//www.cs.cornell.edu/courses/cs4780/2022fa/lectures/lecturenote19.html

"Scalable, Portable and Distributed Gradient Boosting (GBDT, GBRT or GBM) Library, for Python, R, Java, Scala, C++ and more. Runs on single machine, Hadoop, Spark, Dask, Flink, and DataFlow."

- A faster version of gradient boosting.
- Normal regression trees: split to minimize squared loss of leaf predictions.
- XGBoost trees only fit residuals: split so that residuals in leaf are more similar.
- Gradient descent sped up by using the second derivative of the loss function.
- Strong regularization by pre-pruning the trees.
- Column and row are randomly subsampled when computing splits.
- . . .

---

# LightGBM[34](optional)

"A fast, distributed, high performance gradient boosting (GBT, GBDT, GBRT, GBM or MART) framework based on decision tree algorithms, used for ranking, classification and many other machine learning tasks."

- Uses gradient-based sampling
- Use all instances with large gradients/residuals
- Randomly sample instances with small gradients, ignore the rest
- Intuition: samples with small gradients are already well-trained.
- Requires adapted information gain criterion
- Does smarter encoding of categorical features

---

[3]https://lightgbm.readthedocs.io/en/v3.3.2/
[4]https://github.com/microsoft/LightGBM

# CatBoost[56](optional)

"A fast, scalable, high performance Gradient Boosting on Decision Trees library, used for ranking, classification, regression and other machine learning tasks for Python, R, Java, C++. Supports computation on CPU and GPU."

- Another fast boosting technique
- Optimized for categorical variables
- Uses bagged and smoothed version of target encoding
- Uses symmetric trees: same split for all nodes on a given level aka
- Allows monotonicity constraints for numeric features
- Model must be be a non-decreasing function of these features
- Lots of tooling (e.g. GPU training)
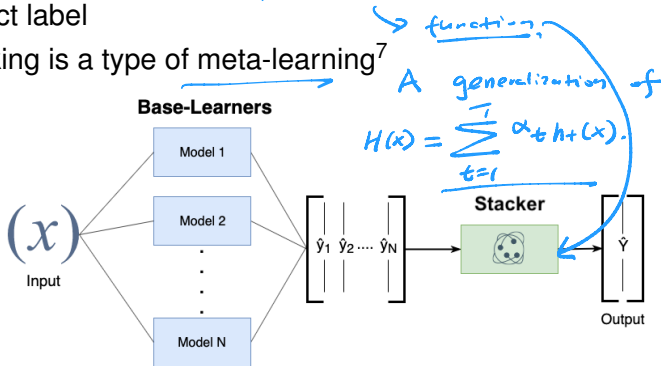- Focusing on optimizing decision trees for categorical variables

---

[5] https://catboost.ai/
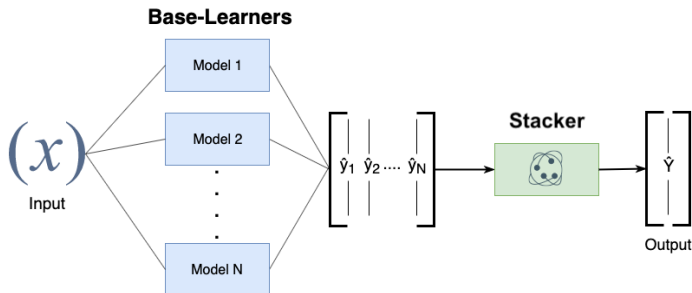[6] https://github.com/catboost/catboost

# Stacking

- Choose *M* different base-models, generate predictions
- Stacker (meta-model) learns mapping between predictions and correct label
- Stacking is a type of meta-learning[7]



*handwritten annotations:* function; A generalization of $H(x) = \sum_{t=1}^{T} \alpha_t h_t(x).$

**Base-Learners**

Model 1
Model 2
Model N

$(x)$ Input

$\hat{y}_1 \ \hat{y}_2 \cdots \hat{y}_N$

**Stacker**

$\hat{y}$

Output

# Stacking



**Base-Learners**

Model 1

Model 2

Model N

$\hat{y}_1$ $\hat{y}_2$ .... $\hat{y}_N$

$(x)$

Input

**Stacker**

$\hat{Y}$

Output

- Stacking can be repeated: multi-level stacking
- Popular stackers: linear models (fast) and gradient boosting (accurate)
- Models need to be sufficiently different, be experts at different parts of the data
- Can be very accurate, but also very slow to predict

# Summary for ensemble learning

- Bagging
  - Vanilla bagging
  - Random forest
- Boosting
  - Gradient boosting
  - AdaBoost
  - . . .
- Stacking
- Other ensemble learning techniques
  - Bayes optimal classifier
  - Bayesian model averaging
  - Any combination of different ensembling techniques
  - Mixture of experts
  - . . .

# Learning outcomes

- Understand the main ideas of bagging, boosting, and stacking
- Understand the algorithms of gradient boosting and AdaBoost
- Know the connection and difference between gradient boosting and AdaBoost
- Be able to utilize some packages of ensemble learning such as XGBoost to solve real problems