

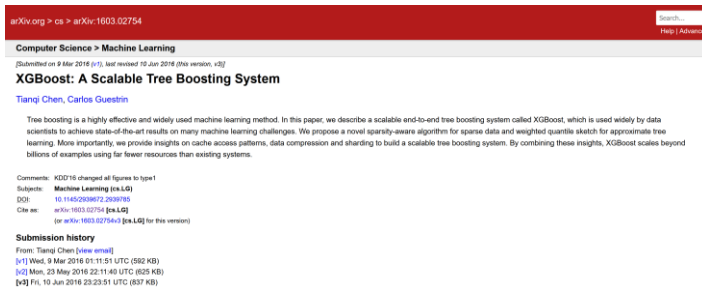
# eXtreme Gradient Boosting

**Xiao Lang**

Department of Mechanics and Maritime Sciences  
Division of Marine Technology  
Gothenburg, SWEDEN

# What is XGBoost ?

- **eXtreme Gradient Boosting = XGBoost**
- **XGBoost is machine learning library like numpy, tensorflow, pytorch**
- **XGBoost has dominated machine learning hackathons and competitions**



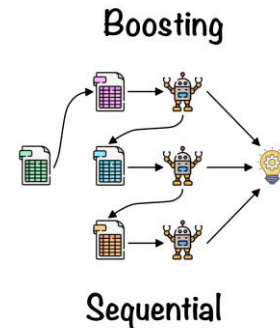
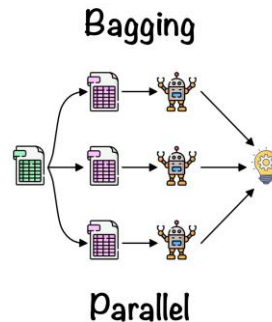
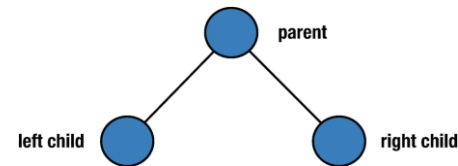
# Evolution of XGBoost

• DT → Boosting → GBDT → XGBoost

Time	Model	Original article
1986	DT	Induction of <b>Decision Trees</b>
1995	Boosting	A Decision-Theoretic Generalization of On-line learning and an Application to <b>Boosting</b>
2001	GBDT	Greedy Function Approximation: A <b>Gradient Boosting</b> Machine
2016	XGBoost	<b>XGBoost</b> : A Scalable Tree Boosting System

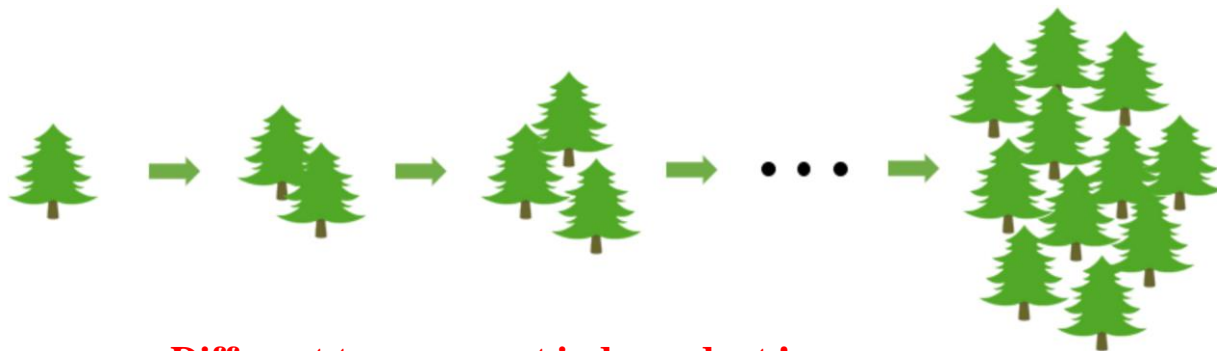
# Evolution of XGBoost

- Decision tree: ID3, C4.5, **CART (binary tree)**
- Boosting: one of ensemble learning method
- Gradient boosting decision tree
- XGBoost: **CART (binary tree)**
  - classification
  - regression



# GBDT & XGBoost similarity

- **Boosting: establish tree (weak evaluator) one by one and accumulation**



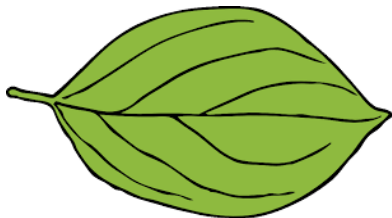
**Different trees are not independent !**

# GBDT & XGBoost difference

- Decision tree

## Regression

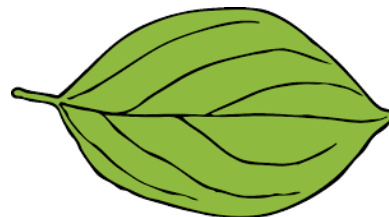
Leaf node: **average**



Sample	Actual
1	0.3
2	0.2
3	1.5
4	0.8
5	0.6
Prediction	0.68

## Classification

Leaf node: **majority**



Sample	Actual
1	0
2	1
3	0
4	0
5	1
Prediction	0

# GBDT & XGBoost difference

- Gradient boosting decision tree

$$\hat{y}_i^{(k)} = \sum_k^K \gamma_k \boxed{h_k(x_i)} \quad \text{Average or majority}$$

- XGBoost

$$\hat{y}_i^{(k)} = \sum_k^K \boxed{f_k(x_i)} \quad \text{Prediction score / leaf weight}$$

# XGBoost parameter

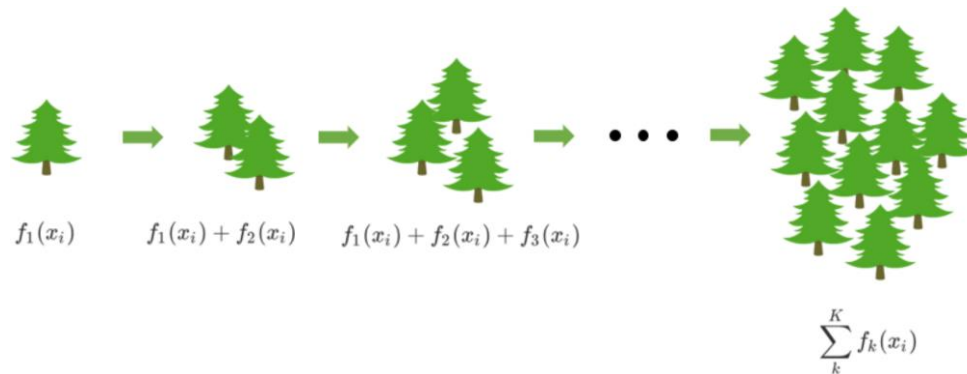
```
class xgboost.XGBRegressor (max_depth=3, learning_rate=0.1, n_estimators=100, silent=True, objective='reg:linear',  
booster='gbtree', n_jobs=1, nthread=None, gamma=0, min_child_weight=1, max_delta_step=0, subsample=1,  
colsample_bytree=1, colsample_bylevel=1, reg_alpha=0, reg_lambda=1, scale_pos_weight=1, base_score=0.5,  
random_state=0, seed=None, missing=None, importance_type='gain', **kwargs)
```

- **Ensemble:** *n\_estimators*, *learning\_rate* (eta)...
- **Weak evaluator:** *max\_depth*, *gamma*, *reg\_alpha*, *reg\_lambda*...
- **Application process:** *n\_jobs*...



# XGBoost parameter - - ensemble

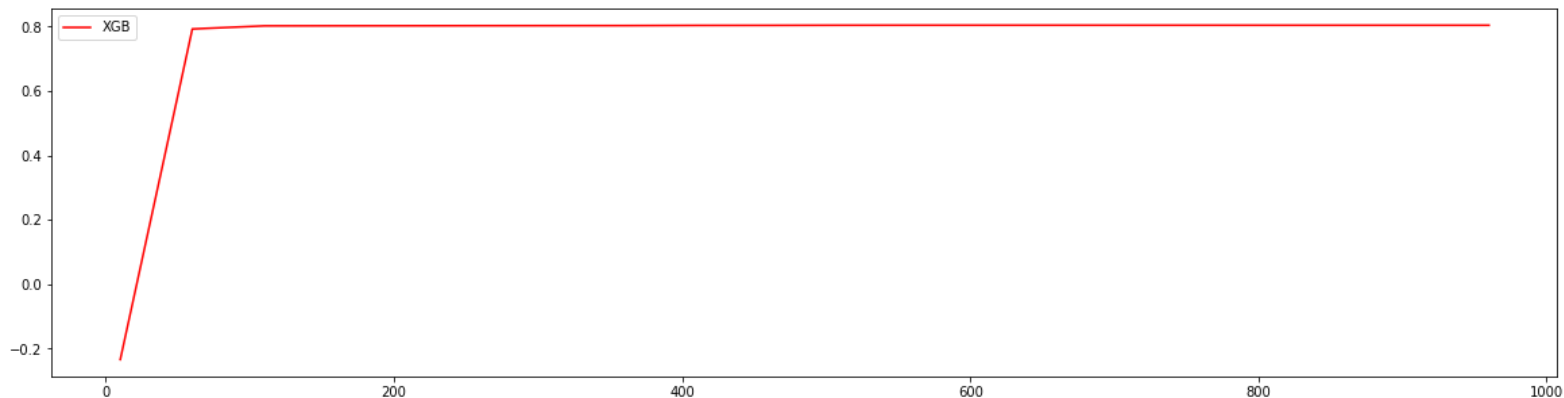
- **n\_estimator**: number of trees to be grown



- **Influence of n\_estimator on the XGBoost model ?**

# XGBoost parameter - - ensemble

- Influence of `n_estimator` on the XGBoost model ?



# XGBoost parameter - - ensemble

- **learning\_rate (eta):** step size shrinkage used in tree grow

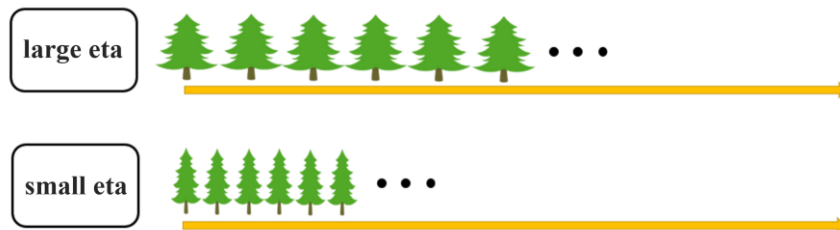
**Logistic regression**  $\Rightarrow$  best fit  $\Rightarrow$  objective function  $\Rightarrow \theta_{k+1} = \theta_k - \alpha * d_{ki}$

**Gradient boosting**  $\Rightarrow$  best prediction  $\Rightarrow$  objective function  $\Rightarrow \hat{y}_i^{(k+1)} = \hat{y}_i^{(k)} + f_{k+1}(x_i)$

# XGBoost parameter - - ensemble

- **learning\_rate (eta):** step size shrinkage used in tree grow

$$\hat{y}_i^{(k+1)} = \hat{y}_i^{(k)} + \eta f_{k+1}(x_i)$$



- **Influence of learning\_rate on the XGBoost model ?**

# Objective function

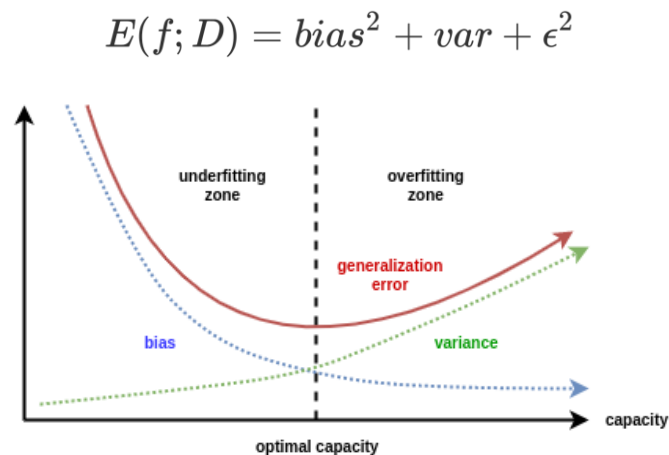
- Logistic regression & SVM - - fixed
- Ensemble model - - optional: **differentiable & can be optimized**
  - RMSE, error, log\_loss...
  - only measures model's generalization ability
- XGBoost: **model performance + computing speed**

# XGBoost objective function

- Obj = **loss function** + **model complexity**

$$Obj = \sum_{i=1}^m l(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k)$$

$$\hat{y}_i^{(t)} = \sum_k^t f_k(x_i) = \sum_k^{t-1} f_k(x_i) + f_t(x_i)$$



# XGBoost objective function

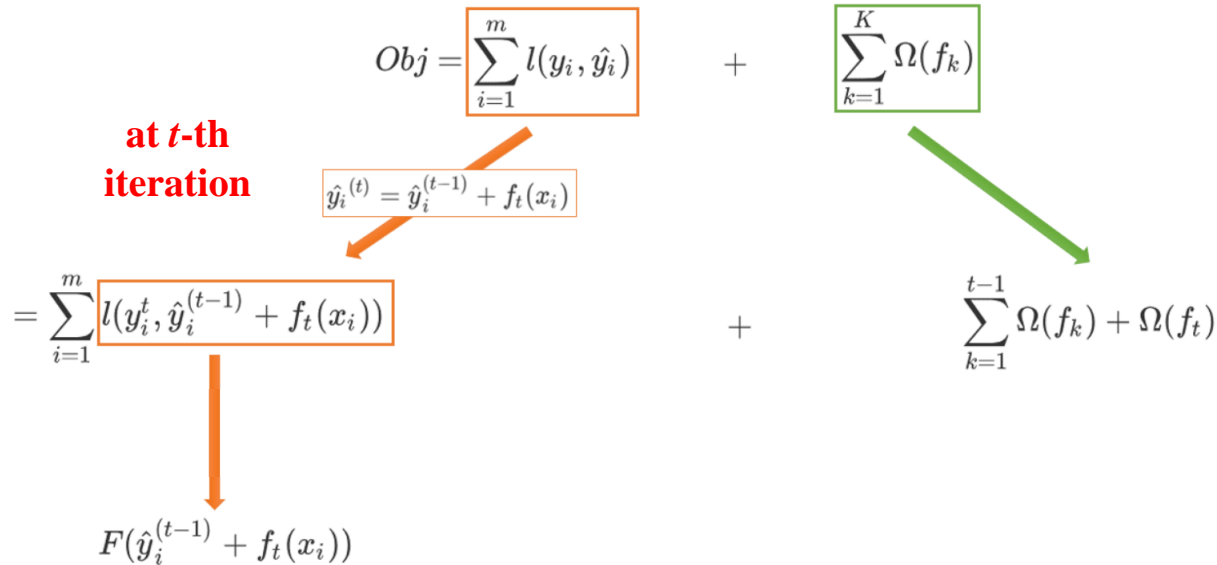
- Solve the objective function

**at  $t$ -th iteration**

$$Obj = \sum_{i=1}^m l(y_i, \hat{y}_i) + \sum_{k=1}^K \Omega(f_k)$$

$$\hat{y}_i^{(t)} = \hat{y}_i^{(t-1)} + f_t(x_i)$$

$$= \sum_{i=1}^m l(y_i^t, \hat{y}_i^{(t-1)} + f_t(x_i)) + \sum_{k=1}^{t-1} \Omega(f_k) + \Omega(f_t)$$

$$F(\hat{y}_i^{(t-1)} + f_t(x_i))$$


# XGBoost objective function

- Taylor Expansion

$$f(x) = \frac{f(c)}{0!} + \frac{f'(c)}{1!}(x-c) + \frac{f''(c)}{2!}(x-c)^2 + \frac{f'''(c)}{3!}(x-c)^3 + \dots$$

assume  $c$  close to  $x$

$$\begin{aligned} f(x + x - c) &\approx \frac{f(c)}{0!} + \frac{f'(c)}{1!}(x-c) + \frac{f''(c)}{2!}(x-c)^2 + \frac{f'''(c)}{3!}(x-c)^3 + \dots \\ &\approx \frac{f(c)}{0!} + \frac{f'(c)}{1!}(x-c) + \frac{f''(c)}{2!}(x-c)^2 \\ &\approx f(c) + f'(c)(x-c) + \frac{f''(c)}{2}(x-c)^2 \end{aligned}$$

$x_1 = x$ ,  $x_2 = x - c$

$$f(x_1 + x_2) \approx f(x_1) + f'(x_1) * x_2 + \frac{f''(x_1)}{2} * x_2^2$$



# XGBoost objective function

- Solve the objective function

$$\mathbf{x}_1 = \mathbf{x}, \mathbf{x}_2 = \mathbf{x} - \mathbf{c}$$

$$f(x_1 + x_2) \approx f(x_1) + x_2 * f'(x_1) + \frac{1}{2}(x_2)^2 * f''(x_1)$$

$$F(\hat{y}_i^{(t-1)} + f_t(x_i)) \approx F(\hat{y}_i^{(t-1)}) + f_t(x_i) * \frac{\partial F(\hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)}} + \frac{1}{2}(f_t(x_i))^2 * \frac{\partial^2 F(\hat{y}_i^{(t-1)})}{\partial (\hat{y}_i^{(t-1)})^2}$$

at  $t$ -th  
iteration

$$\approx l(y_i^t, \hat{y}_i^{(t-1)}) + f_t(x_i) * \frac{\partial l(y_i^t, \hat{y}_i^{(t-1)})}{\partial \hat{y}_i^{(t-1)}} + \frac{1}{2}(f_t(x_i))^2 * \frac{\partial^2 l(y_i^t, \hat{y}_i^{(t-1)})}{\partial (\hat{y}_i^{(t-1)})^2}$$

$$\approx l(y_i^t, \hat{y}_i^{(t-1)}) + f_t(x_i) * g_i + \frac{1}{2}(f_t(x_i))^2 * h_i$$

# XGBoost objective function

- Solve the objective function

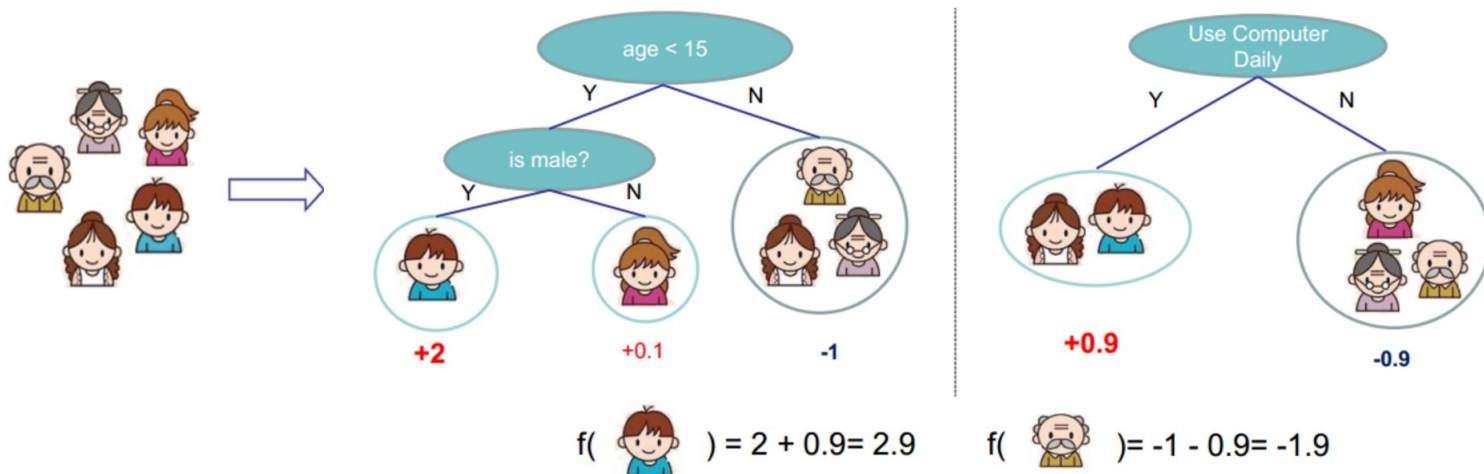
**at  $t$ -th iteration**

$$= \sum_{i=1}^m [l(y_i^t, \hat{y}_i^{(t-1)}) + f_t(x_i)g_i + \frac{1}{2}(f_t(x_i))^2 h_i] + \sum_{k=1}^{t-1} \Omega(f_k) + \Omega(f_t)$$

$$Obj = \sum_{i=1}^m [f_t(x_i)g_i + \frac{1}{2}(f_t(x_i))^2 h_i] + \Omega(f_t)$$

# XGBoost parameter - - weak evaluator

- alpha ( $L1$  regularization) & lambda ( $L2$  regularization)



# XGBoost parameter - - weak evaluator

- alpha (*L1* regularization) & lambda (*L2* regularization)

$$f_t(x_i) = w_{q(x_i)}$$

$$\Omega(f) = \gamma T + \textit{Regularization}$$

$$\begin{aligned} &= \gamma T + \frac{1}{2} \alpha |w| \\ &= \gamma T + \frac{1}{2} \alpha \sum_{j=1}^T |w_j| \end{aligned}$$

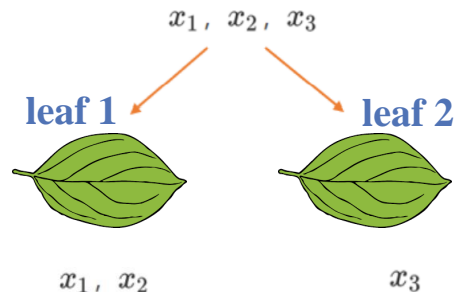
$$= \gamma T + \frac{1}{2} \alpha \sum_{j=1}^T |w_j| + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2$$

$$\begin{aligned} &= \gamma T + \frac{1}{2} \lambda ||w||^2 \\ &= \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2 \end{aligned}$$

# XGBoost tree structures

•  $\omega$  &  $T$

$$\begin{aligned}
 & \sum_{i=1}^m [f_t(x_i)g_i + \frac{1}{2}(f_t(x_i))^2 h_i] + \Omega(f_t) \\
 &= \sum_{i=1}^m [w_{q(x_i)}g_i + \frac{1}{2}w_{q(x_i)}^2 h_i] + \gamma T + \frac{1}{2}\lambda \sum_{j=1}^T w_j^2 \\
 &= \boxed{\sum_{i=1}^m w_{q(x_i)}g_i} + \boxed{\sum_{i=1}^m \frac{1}{2}w_{q(x_i)}^2 h_i} + \gamma T + \frac{1}{2}\lambda \sum_{j=1}^T w_j^2
 \end{aligned}$$



$$\begin{aligned}
 & w_{q(x_1)} * g_1 \\
 & w_{q(x_2)} * g_2
 \end{aligned}$$

$$w_{q(x_3)} * g_3$$

$$w_{q(x_1)} = w_{q(x_2)} = w_1 \quad w_{q(x_3)} = w_2$$

$$\sum_{i=1}^m w_{q(x_i)} * g_i = w_{q(x_1)} * g_1 + w_{q(x_2)} * g_2 + w_{q(x_3)} * g_3$$

$$= w_1(g_1 + g_2) + w_2 * g_3$$

$$= \sum_{j=1}^T (w_j \sum_{i \in I_j} g_i)$$

# XGBoost tree structures

•  $\omega$  &  $T$

$$\begin{aligned}
 & \sum_{j=1}^T (w_j * \sum_{i \in I_j} g_i) + \frac{1}{2} \sum_{j=1}^T (w_j^2 * \sum_{i \in I_j} h_i) + \gamma T + \frac{1}{2} \lambda \sum_{j=1}^T w_j^2 \\
 &= \sum_{j=1}^T \left[ w_j \sum_{i \in I_j} g_i + \frac{1}{2} w_j^2 (\sum_{i \in I_j} h_i + \lambda) \right] + \gamma T
 \end{aligned}$$

$$G_j = \sum_{i \in I_j} g_i, \quad H_j = \sum_{i \in I_j} h_i \quad \longrightarrow \quad = \sum_{j=1}^T \left[ w_j G_j + \frac{1}{2} w_j^2 (H_j + \lambda) \right] + \gamma T$$

$$\frac{\partial F^*(w_j)}{\partial w_j} = G_j + w_j (H_j + \lambda)$$

$$0 = G_j + w_j (H_j + \lambda)$$

$$w_j = - \frac{G_j}{H_j + \lambda}$$

$$F^*(w_j) = w_j G_j + \frac{1}{2} w_j^2 (H_j + \lambda)$$

# XGBoost tree structures






•  $\omega$  &  $T$

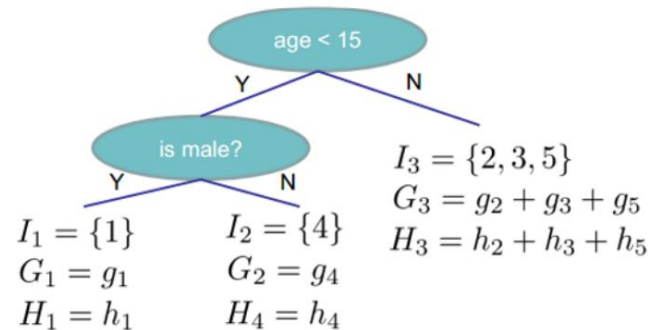
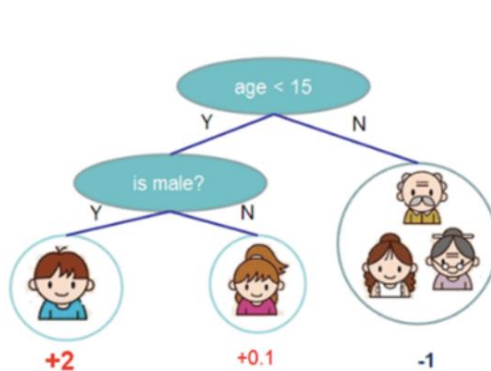
$$\begin{aligned} &= \sum_{j=1}^T \left[ -\frac{G_j}{H_j + \lambda} * G_j + \frac{1}{2} \left( -\frac{G_j}{H_j + \lambda} \right)^2 (H_j + \lambda) \right] + \gamma T \\ &= \sum_{j=1}^T \left[ -\frac{G_j^2}{H_j + \lambda} + \frac{1}{2} * \frac{G_j^2}{H_j + \lambda} \right] + \gamma T \\ &= -\frac{1}{2} \sum_{j=1}^T \frac{G_j^2}{H_j + \lambda} + \gamma T \end{aligned}$$

**structure score**

# XGBoost tree structure

## • Structure score

- 1  g1, h1
- 2  g2, h2
- 3  g3, h3
- 4  g4, h4
- 5  g5, h5



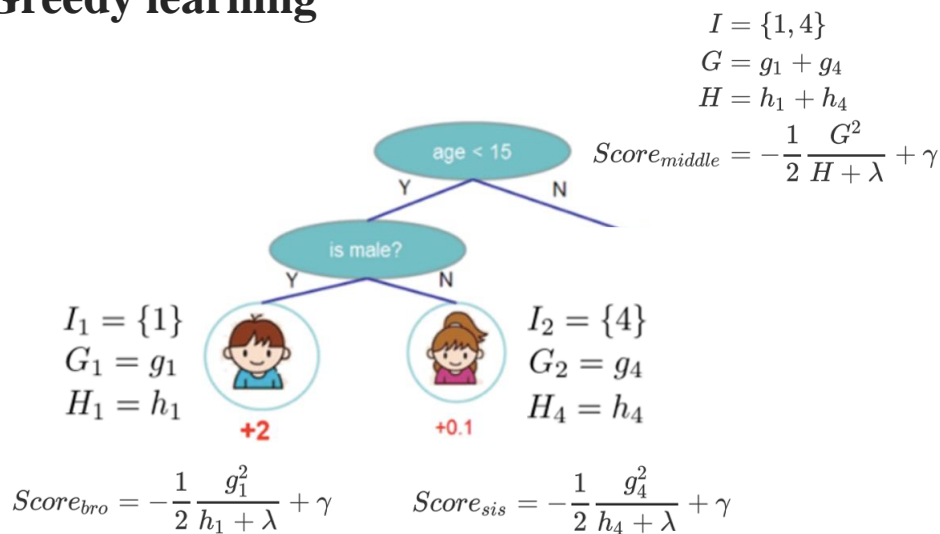
$$Obj = - \sum_j \frac{G_j^2}{H_j + \lambda} + 3\gamma$$

$$Obj = - \left( \frac{g_1^2}{h_1 + \lambda} + \frac{g_4^2}{h_4 + \lambda} + \frac{(g_2 + g_3 + g_5)^2}{h_2 + h_3 + h_5 + \lambda} \right) + 3\gamma$$



# XGBoost tree structure

## • Greedy learning



$$\begin{aligned}
 Gain &= Score_{sis} + Score_{bro} - Score_{middle} \\
 &= -\frac{1}{2} \frac{g_4^2}{h_4 + \lambda} + \gamma - \frac{1}{2} \frac{g_1^2}{h_1 + \lambda} + \gamma - \left( -\frac{1}{2} \frac{G^2}{H + \lambda} + \gamma \right) \\
 &= -\frac{1}{2} \frac{g_4^2}{h_4 + \lambda} + \gamma - \frac{1}{2} \frac{g_1^2}{h_1 + \lambda} + \gamma + \frac{1}{2} \frac{G^2}{H + \lambda} - \gamma \\
 &= -\frac{1}{2} \left[ \frac{g_4^2}{h_4 + \lambda} + \frac{g_1^2}{h_1 + \lambda} - \frac{G^2}{H + \lambda} \right] + \gamma \\
 &= -\frac{1}{2} \left[ \frac{g_4^2}{h_4 + \lambda} + \frac{g_1^2}{h_1 + \lambda} - \frac{(g_1 + g_4)^2}{(h_1 + h_4) + \lambda} \right] + \gamma \\
 Gain &= \frac{1}{2} \left[ \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] - \gamma
 \end{aligned}$$

# XGBoost parameter - - weak evaluator

•  $\gamma$       $Gain = \frac{1}{2} \left[ \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] - \gamma$

$$\frac{1}{2} \left[ \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] - \gamma > 0$$

$$\frac{1}{2} \left[ \frac{G_L^2}{H_L + \lambda} + \frac{G_R^2}{H_R + \lambda} - \frac{(G_L + G_R)^2}{H_L + H_R + \lambda} \right] > \gamma$$

# Summary

- **Efficiency**
  - **support parallel processing implementation**
- **Accuracy**
  - **in-built regularization to reduce overfitting**
  - **more effective tree pruning**
- **Feasibility**
  - **customized objective and evaluation**
  - **tunable parameters**



**CHALMERS**  
UNIVERSITY OF TECHNOLOGY