

LEC004 Demand Forecasting

VG441 SS2020

Cong Shi
Industrial & Operations Engineering
University of Michigan

Ensemble Learning

“The wisdom of the crowd is the collective opinion of a group of individuals rather than that of a single expert.”

“A group of predictors is called an ensemble. Therefore this Machine Learning technique is known as Ensemble Learning. Voilá!”

“Ensemble methods work best when the predictors are as independent of one another as possible. One way to get diverse classifiers is to train them using very different algorithms. This increases the chance that they will make very different types of errors, improving the ensemble’s accuracy.”

Ensemble Learning Techniques

- 些 weak prediction .

- Hard voting classifier (for classification)

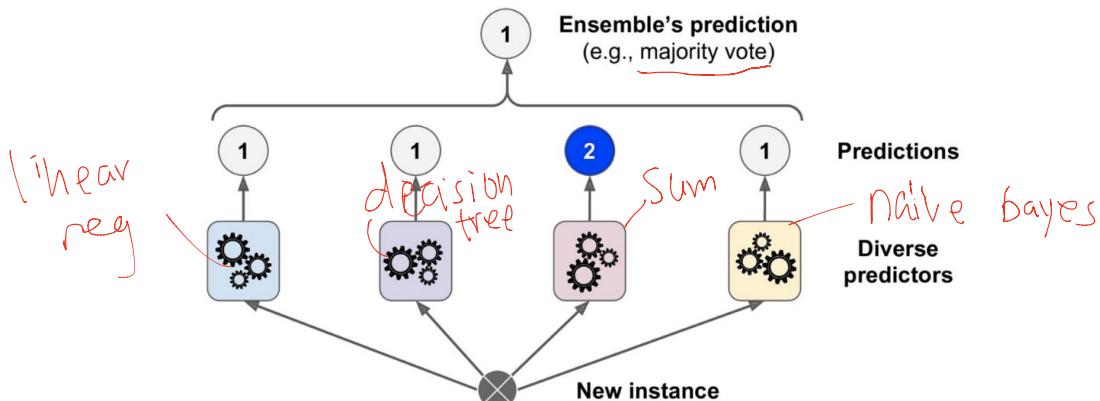


Figure 7-2. Hard voting classifier predictions

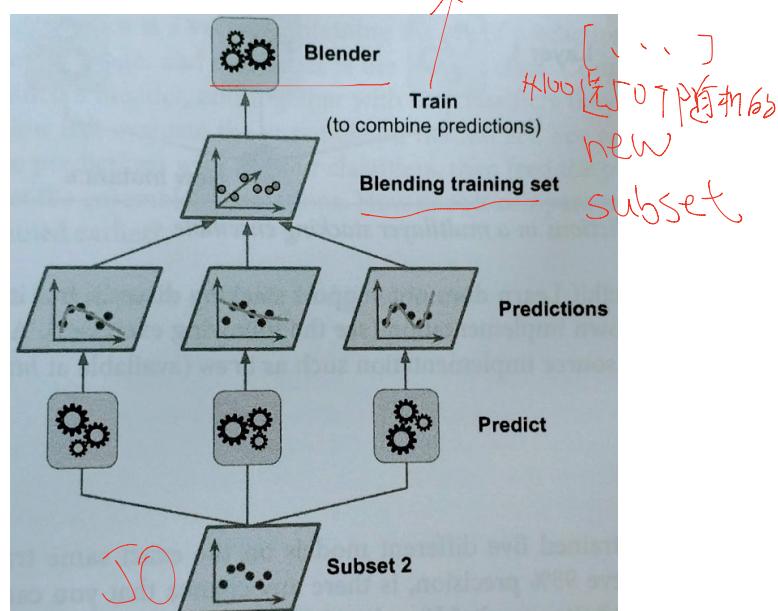
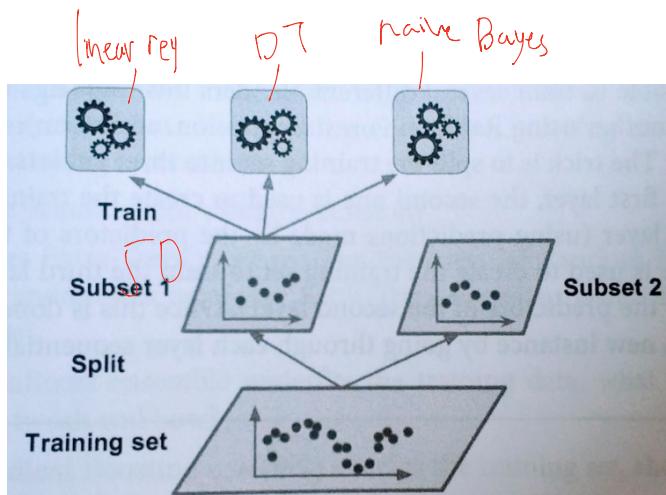
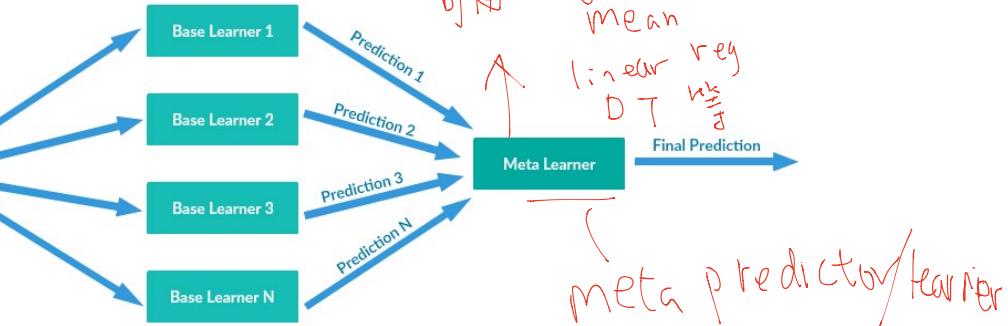
- Averaging or weighted averaged (for regression)

直接取平均 233

Ensemble Learning Techniques

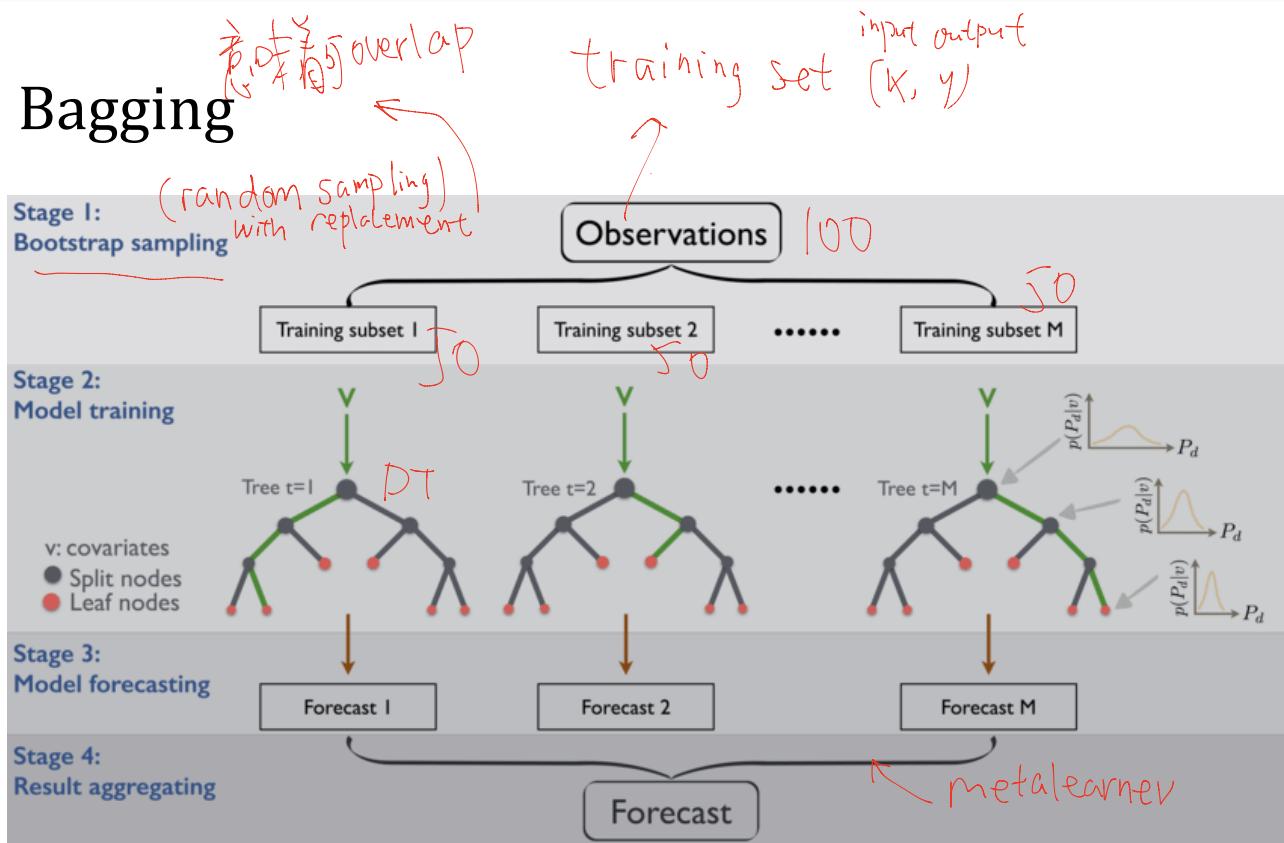
- Stacking

100 train
half half
Data 50 - 50



Ensemble Learning Techniques

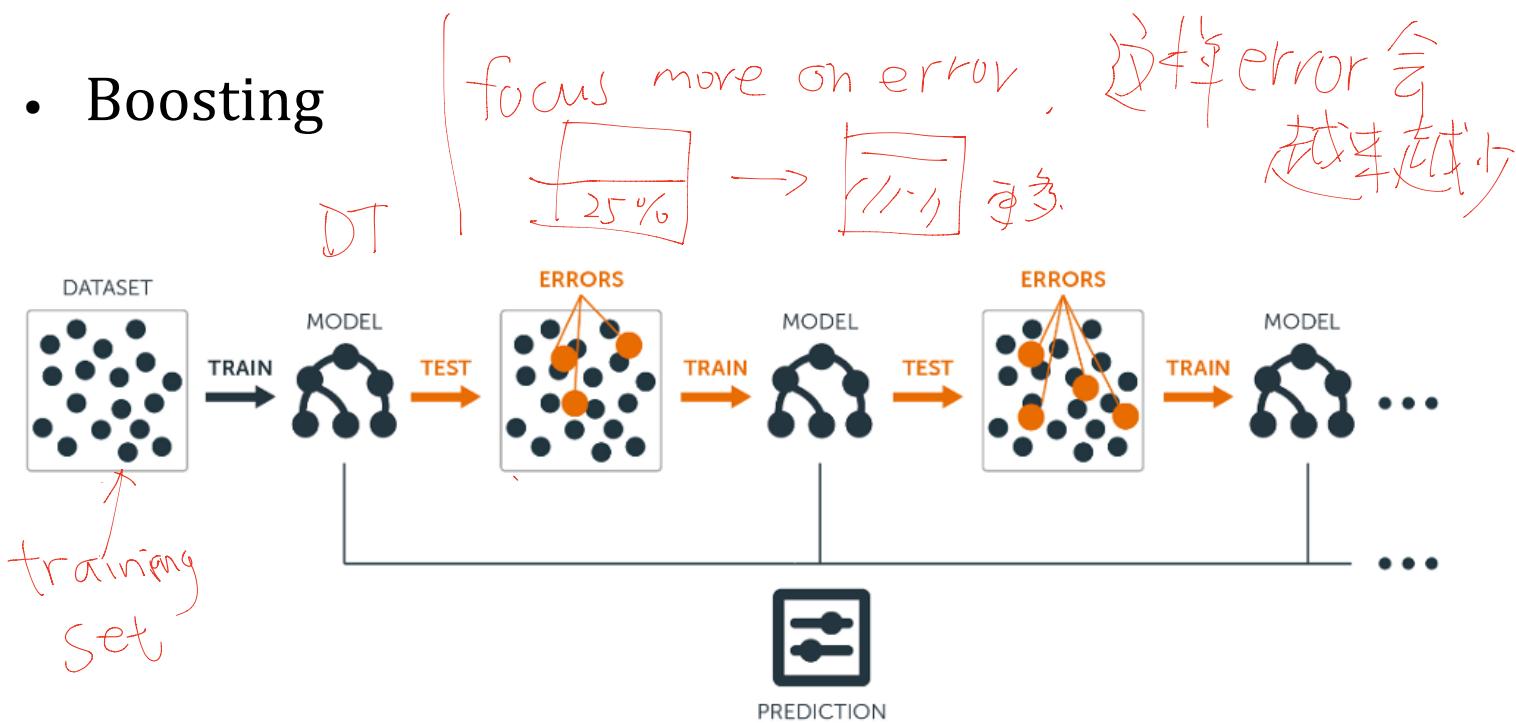
• Bagging



Random Forest

Ensemble Learning Techniques

- Boosting



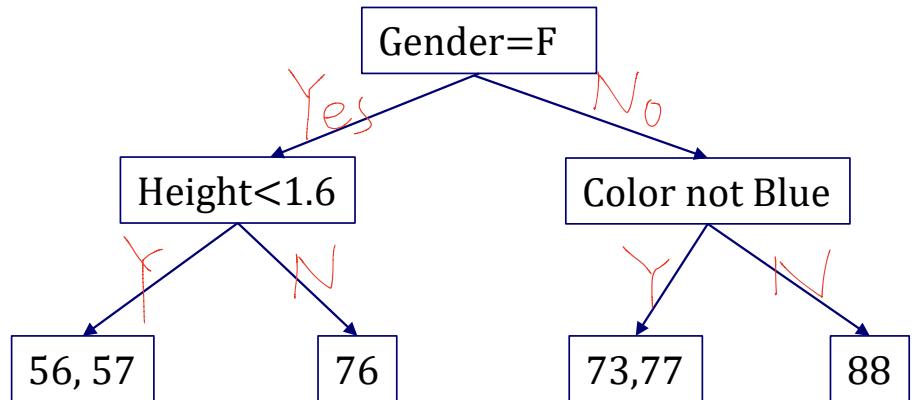
AdaBoost, Gradient Boosting, XGBoost

逐個 error 會 converge
逐個 focus

Decision Tree

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56
1.8	Red	Male	73
1.5	Green	Male	77
1.4	Blue	Female	57

features X Y



Question 1: How do we determine the next node (starting from root)?

把分类法都试一下，取最小，试 threshold

Question 2: Should we split at the current node?

看 deviance 是不是小了，

How to determine and split a node?

impurity
size

Measure of impurity (for regression) is deviance

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56
1.8	Red	Male	73
1.5	Green	Male	77
1.4	Blue	Female	57

88, 76, 56, 73, 77, 57

$$\text{Deviance} = 774.83$$

越低越好

$$\sum(y_i - \bar{y})^2$$

Gender=F

76, 56, 57

88, 73, 77

$$\text{Deviance} = 254 + 120.67 = 374.67$$

Height < 1.6

可以用 Height < 1.6

不用用 Height < 1.5

56, 77, 57

88, 76, 73

Not Blue

88, 56, 57

$$\text{Deviance} = 280.67 + 126 = 406.67$$

$$\text{Deviance} = 8.67 + 662 = 670.67$$

目标是要 Short tree → robust → small variance
large bias 8

Gradient Boosting

F_0 = Initial Model = Taking the mean

Height (m)	Favorite Color	Gender	Weight (kg)
1.6	Blue	Male	88
1.6	Green	Female	76
1.5	Blue	Female	56
1.8	Red	Male	73
1.5	Green	Male	77
1.4	Blue	Female	57



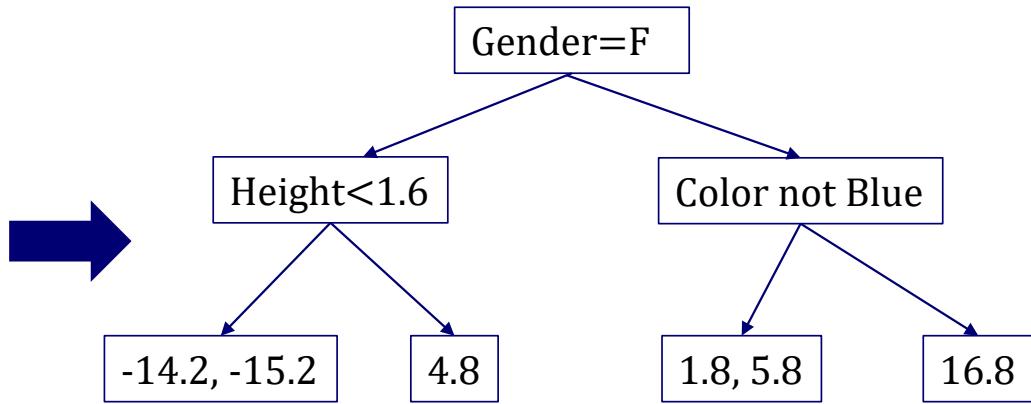
Height (m)	Favorite Color	Gender	Weight (kg)	F_0	PR
1.6	Blue	Male	88	71.2	16.8
1.6	Green	Female	76	71.2	4.8
1.5	Blue	Female	56	71.2	-15.2
1.8	Red	Male	73	71.2	1.8
1.5	Green	Male	77	71.2	5.8
1.4	Blue	Female	57	71.2	-14.2

Pseudo Residual (PR) = True Value - Predicted Value

Gradient Boosting

Fit PR₀ into a decision tree (up to four leaves)

Height (m)	Favorite Color	Gender	PR ₀
1.6	Blue	Male	16.8
1.6	Green	Female	4.8
1.5	Blue	Female	-15.2
1.8	Red	Male	1.8
1.5	Green	Male	5.8
1.4	Blue	Female	-14.2

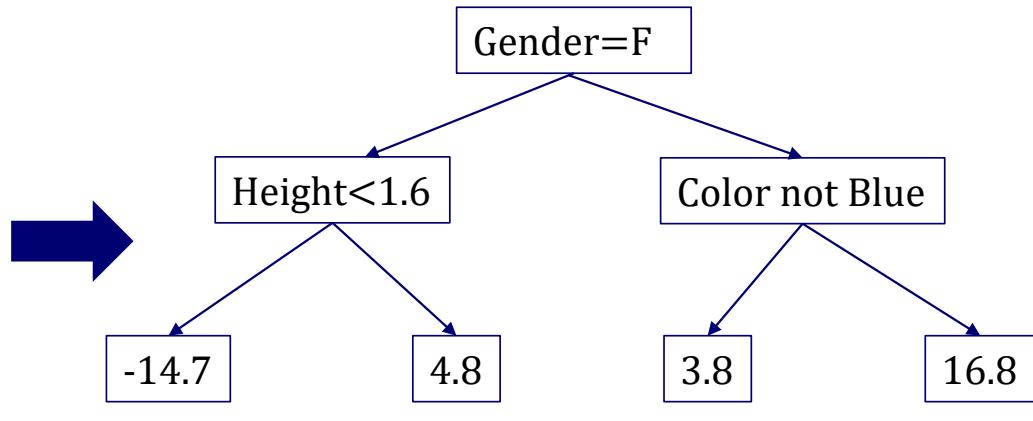


Pseudo Residual (PR) = True Value – Predicted Value

Gradient Boosting

Fit PR0 into a decision tree (up to four leaves)

Height (m)	Favorite Color	Gender	PR0
1.6	Blue	Male	16.8
1.6	Green	Female	4.8
1.5	Blue	Female	-15.2
1.8	Red	Male	1.8
1.5	Green	Male	5.8
1.4	Blue	Female	-14.2



Averaging the residuals on each leaf...

Gradient Boosting

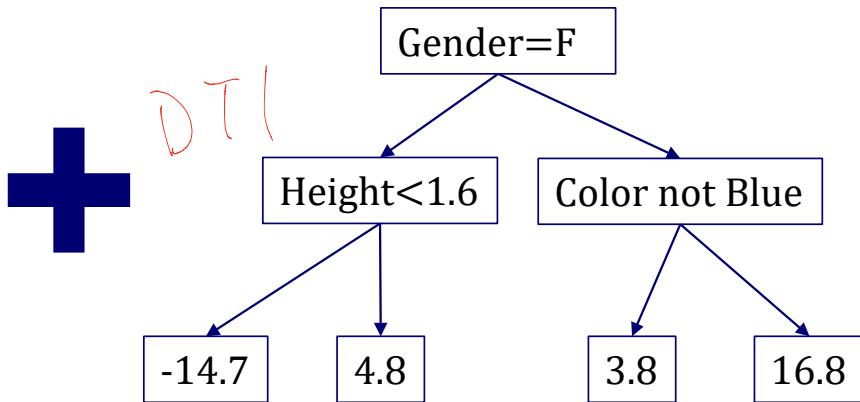
Step-site

next Mean

$$F_1(x) = F_0(x) + \gamma_1 \times \text{Output of DT}(x)$$

Learning rate = 0.1

Height (m)	Favorite Color	Gender	Weight (kg)	F0
1.6	Blue	Male	88	71.2
1.6	Green	Female	76	71.2
1.5	Blue	Female	56	71.2
1.8	Red	Male	73	71.2
1.5	Green	Male	77	71.2
1.4	Blue	Female	57	71.2



$$\left\{
 \begin{aligned}
 F_1((1.6, \text{Blue}, \text{Male})) &= 71.2 + 0.1 \times 16.8 = 72.9 \\
 F_1((1.6, \text{Green}, \text{Female})) &= 71.2 + 0.1 \times 4.8 = 71.7 \\
 F_1((1.5, \text{Blue}, \text{Female})) &= 71.2 + 0.1 \times -14.7 = 69.7 \\
 F_1((1.8, \text{Red}, \text{Male})) &= 71.2 + 0.1 \times 3.8 = 71.6 \\
 F_1((1.5, \text{Green}, \text{Male})) &= 71.2 + 0.1 \times 3.8 = 71.6 \\
 F_1((1.4, \text{Blue}, \text{Female})) &= 71.2 + 0.1 \times -14.7 = 69.7
 \end{aligned}
 \right.$$

Gradient Boosting

So after building the first DT, we obtain...

Height (m)	Favorite Color	Gender	Weight (kg)	F0	PRO	F1	PR1
1.6	Blue	Male	88	71.2	16.8	72.9	15.1
1.6	Green	Female	76	71.2	4.8	71.7	4.3
1.5	Blue	Female	56	71.2	-15.2	69.7	-13.7
1.8	Red	Male	73	71.2	1.8	71.6	1.4
1.5	Green	Male	77	71.2	5.8	71.6	5.4
1.4	Blue	Female	57	71.2	-14.2	69.7	-12.7

$$PR_1 = \text{Weight} - F_1$$

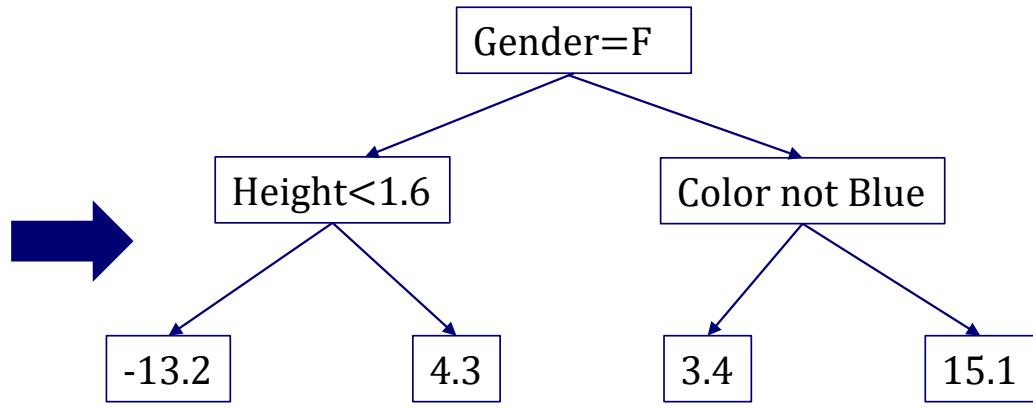
第一波
第二波



Gradient Boosting

Fit PR1 into a decision tree (up to four leaves)

Height (m)	Favorite Color	Gender	PR1
1.6	Blue	Male	15.1
1.6	Green	Female	4.3
1.5	Blue	Female	-13.7
1.8	Red	Male	1.4
1.5	Green	Male	5.4
1.4	Blue	Female	-12.7



新的樹 需要重新算
有可能不同

Gradient Boosting

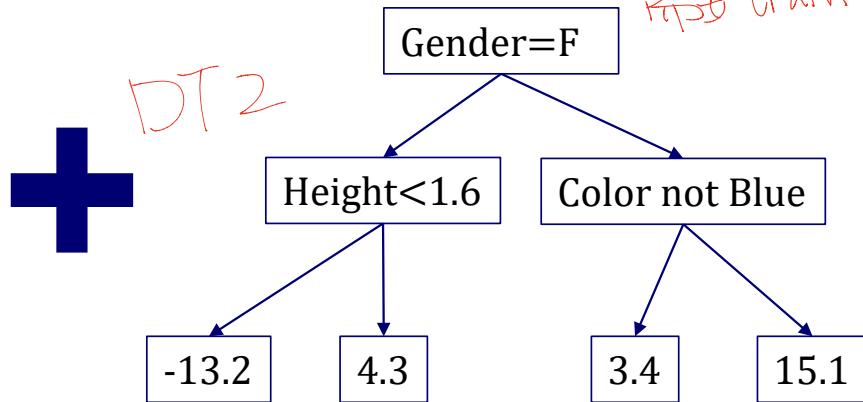
Learning rate = 0.1

$$F_2(x) = F_1(x) + \gamma_2 \times \text{Output of DT}(x)$$

hyper parameter

DT train

Height (m)	Favorite Color	Gender	Weight (kg)	F1
1.6	Blue	Male	88	72.9
1.6	Green	Female	76	71.7
1.5	Blue	Female	56	69.7
1.8	Red	Male	73	71.6
1.5	Green	Male	77	71.6
1.4	Blue	Female	57	69.7



$$F_2((1.6, \text{Blue}, \text{Male})) = 72.9 + 0.1 \times 15.1 = 74.4$$

$$F_2((1.6, \text{Green}, \text{Female})) = 71.7 + 0.1 \times 4.3 = 72.1$$

$$F_2((1.5, \text{Blue}, \text{Female})) = 69.7 + 0.1 \times -13.2 = 68.4$$

$$F_2((1.8, \text{Red}, \text{Male})) = 71.6 + 0.1 \times 3.4 = 71.9$$

$$F_2((1.5, \text{Green}, \text{Male})) = 71.6 + 0.1 \times 3.4 = 71.9$$

$$F_2((1.4, \text{Blue}, \text{Female})) = 69.7 + 0.1 \times -13.2 = 68.4$$

Gradient Boosting

typically会 iterate 100次

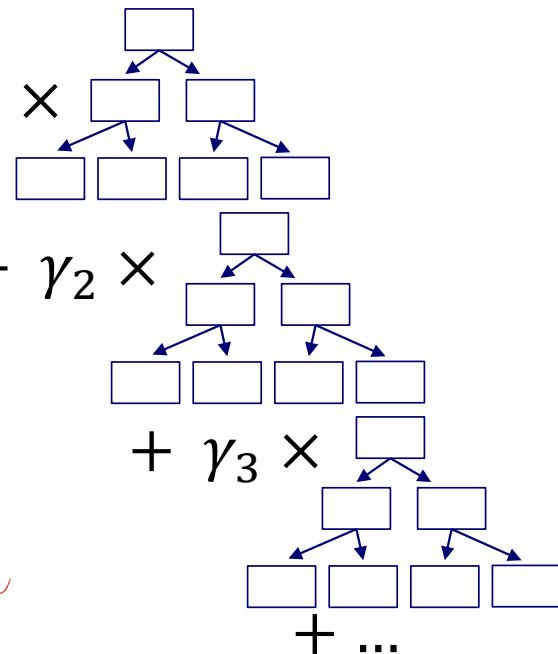
So after building the second DT, we obtain...

Height (m)	Favorite Color	Gender	Weight (kg)	F0	PR0	F1	PR1	F2	PR2
1.6	Blue	Male	88	71.2	16.8	72.9	15.1	74.4	13.6
1.6	Green	Female	76	71.2	4.8	71.7	4.3	72.1	3.9
1.5	Blue	Female	56	71.2	-15.2	69.7	-13.7	68.4	-12.4
1.8	Red	Male	73	71.2	1.8	71.6	1.4	71.9	1.1
1.5	Green	Male	77	71.2	5.8	71.6	5.4	71.9	5.1
1.4	Blue	Female	57	71.2	-14.2	69.7	-12.7	68.4	-11.4

Notice the PR's are shrinking: Small steps towards the right direction!



Gradient Boosting

$$F_m = F_0 + \gamma_1 \times \text{Decision Tree}_1 + \gamma_2 \times \text{Decision Tree}_2 + \gamma_3 \times \text{Decision Tree}_3 + \dots$$


Fit the new PR into DT

Mid term

合計

Stop until the pre-specified #DTs or the PR stops improving!

(e.g. 100)

等到 |PR| < 停下來

Python Time!

- `from sklearn import ensemble`



Some Mathematics...

gradient 就是 PR.

$$\frac{1}{2} \sum (y - F(x))^2$$

Input: training set $\{(x_i, y_i)\}_{i=1}^n$, a differentiable loss function $L(y, F(x))$, number of iterations M .

Algorithm:

1. Initialize model with a constant value:

$$F_0(x) = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, \gamma).$$

least square loss function

2. For $m = 1$ to M :

1. Compute so-called *pseudo-residuals*:

$$r_{im} = - \left[\frac{\partial L(y_i, F(x_i))}{\partial F(x_i)} \right]_{F(x)=F_{m-1}(x)} \quad \text{for } i = 1, \dots, n.$$

Lec 004 #

2. Fit a base learner (or weak learner, e.g. tree) $h_m(x)$ to pseudo-residuals, i.e. train it using the training set $\{(x_i, r_{im})\}_{i=1}^n$.

3. Compute multiplier γ_m by solving the following *one-dimensional optimization* problem:

$$\gamma_m = \arg \min_{\gamma} \sum_{i=1}^n L(y_i, F_{m-1}(x_i) + \gamma h_m(x_i)).$$

4. Update the model:

$$F_m(x) = F_{m-1}(x) + \gamma_m h_m(x).$$

5. Output $F_M(x)$.

learning rate

把 PRs in leaf 平均

Gradient Boosting

- Works exceptionally well in practice
- Won a series of Kaggle competitions
- More robust and explainable than DNN

Data science → Kaggle
• com

u.32n