## 1. Bagging

- k bootstrap sample $(D_1, D_2, D_3 \cdots D_k)$
- 복원추출
- Train distinct model on each $D_i$
- Test a new instance by majority vote or average

① multiple Data sets w/2222

② Build multiple models

3. Combine model

average or vote

→ 문제해결을위해
여러전문가에게문의!

→ model diversity

## 2. Random Forest

- Decision tree 모델을 bagging 하여
  full-grown tree 로 low bias 의 여러 DT를 뽑음

- 다양한 DT가 초여면서 high variance (overfiting)을
  줄임

$$Var(\bar{x}) = \frac{Var(x)}{n}$$   OK

t=1    t=2    t=k-1    t=k

예시 1 ~

$$x = \begin{pmatrix} 1 & 5 & 0 & -12 \\ 0 & 3 & 9 & 1 & -3 \\ 2 & 8 & 9 & 0 & 3 \\ 3 & -1 & 0 & -23 \end{pmatrix} \qquad Y = \begin{pmatrix} G \\ 1 \\ 0 \\ 0 \end{pmatrix}$$

$$\left[\begin{matrix}1\ 3 \ -1\ 0 \ -2 \ 3\end{matrix}\right] x$$
$$\left[0\right]$$

$$x^{(1)} = \begin{pmatrix} 1 & 0 & 2 \\ 2 & 9 & 3 \\ 3 & 0 & 3 \end{pmatrix} \qquad y^{(1)} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \qquad x^{(2)} = \begin{pmatrix} 5 & -1 & 2 \\ 3 & -1 & 3 \\ 8 & 0 & 3 \end{pmatrix}$$

$t_1$ $t_2$

$$y^{(2)} = \begin{pmatrix} 0 \\ 1 \\ 1 \end{pmatrix}$$

① Bootstrapping trainning set $(V)$
+
features $(U)$

overfitting 가능성

Low bias (full-growed tree)
but hig variance

‖ ‖

Low variance

② Feature Importance

$$Importance(X_1) = \frac{1}{M} \sum_{m=1}^{M} \Delta I(X_1, T_m)$$

- $X$ = feature
- $M$ = 모델 (Tree)개수
- $T_m$ = $m$번째 tree
- $\Delta I = \dfrac{X_{左} \quad X_{右}}{T_m}$   $\Delta I$ = error Gap

⇒ $\Delta \sum$은 가장 크게되는 feature 순으로 중요도 나레이가 계산됨