data2

target = 0

target = 1

data1

target

$$I_G(n) = 1 - \sum_{i=1}^{J} (P_i)^2 \qquad \text{Gini Impurity}$$

$$I_{root} = 1 - \left[ \left(\frac{2}{6}\right)^2 + \left(\frac{4}{6}\right)^2 \right] = 1 - \frac{20}{36} = \frac{4}{9} \qquad \text{without classification line}$$

if target = 1          if target = 0
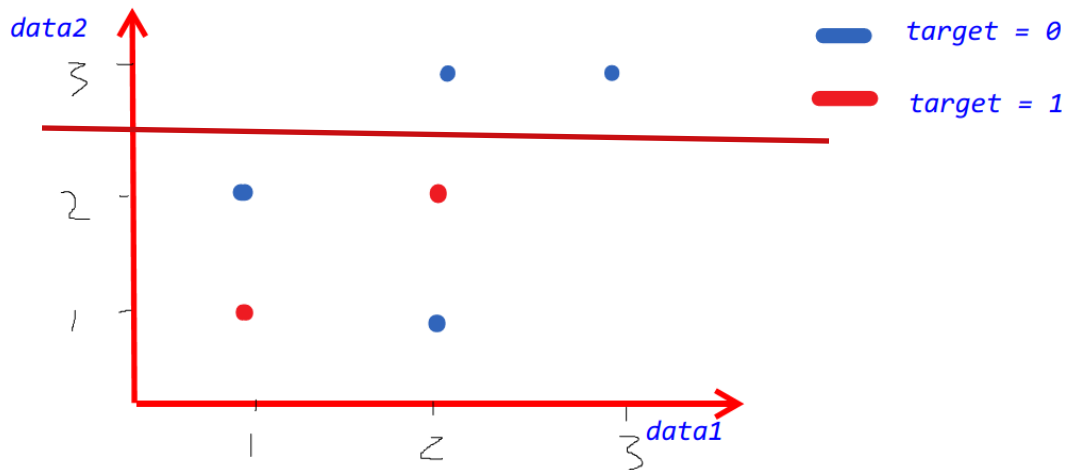
if there's only one type of data in the group, then I(root) = 0. for example: target = 0

target = 1

— target = 0

— target = 1

$$I_{root} = 1 - \left[ \left(\frac{4}{4}\right)^2 + \left(\frac{0}{4}\right)^2 \right]$$

$$= 1 - 1$$

$$= 0$$

data2

3 — (blue dot) (blue dot)

2 — (blue dot) (red dot)

1 — (red dot) (blue dot)

data1
1  2  3

target = 0 (blue)
target = 1 (red)

if we have a classification line(red line)

gini impurity for the upper part                     lower part

$$I = \left(1 - \left[(\tfrac{0}{2})^2 + (\tfrac{2}{2})^2\right]\right) \times \frac{2}{6} + \left(1 - \left[(\tfrac{2}{4})^2 + (\tfrac{2}{4})^2\right]\right) \times \frac{4}{6}$$
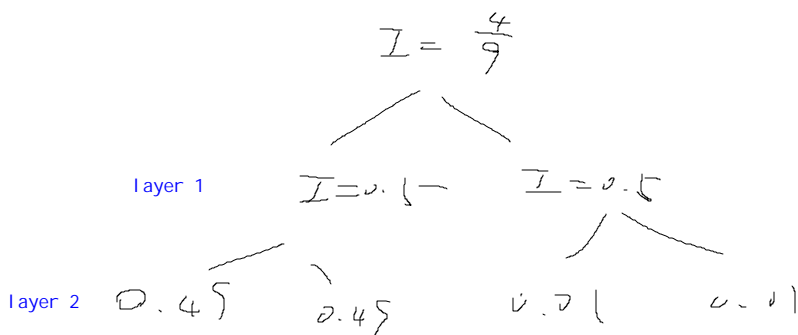
$$= \frac{1}{3} < \frac{4}{9}$$

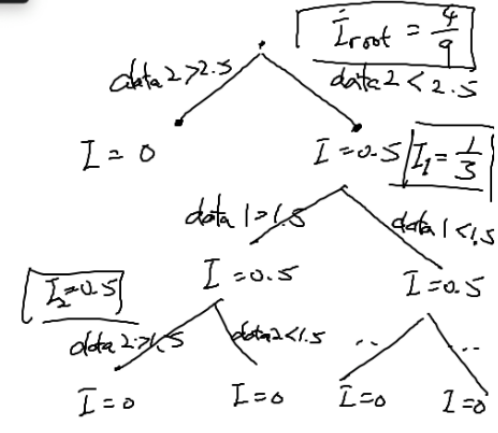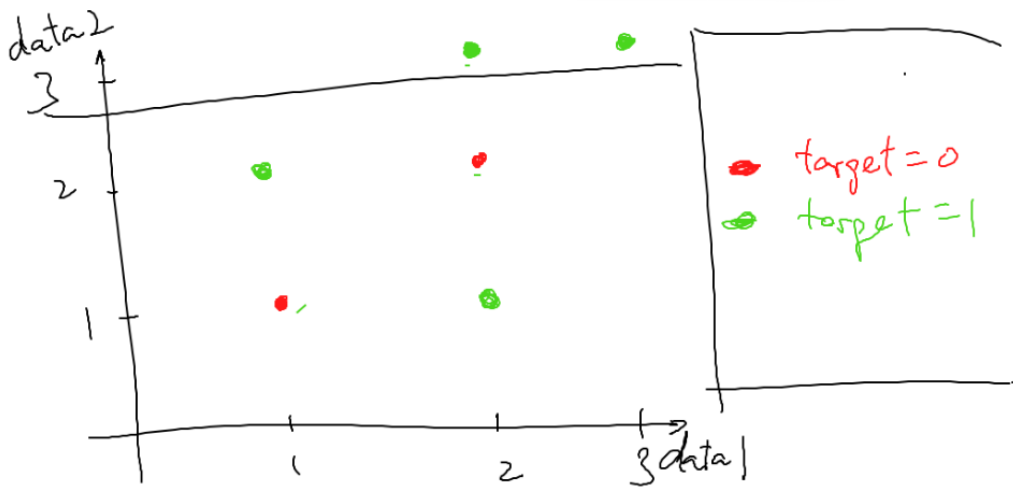so, every time you draw a classification line, it'll be better

Notice, sometimes, we don't need to let I = 0. (overfitting)
if two points are different fundamentally, then it is necessary to add the line,
however, if they are different because other noise, then there's no need to add the line.

In practice, if you find the next level of Gini impurity improves a little, then there's no need to add line any more.

$$I = \frac{4}{9}$$

layer 1      $I = 0.1\sim$      $I = 0.5$

layer 2   $0.45$      $0.45$      $0.7|$      $0.01$

normally, we can limit the number of layer in each tree to solve over fitting problem. ie, if we allow 1 layer, then, layer 2 will not appear.
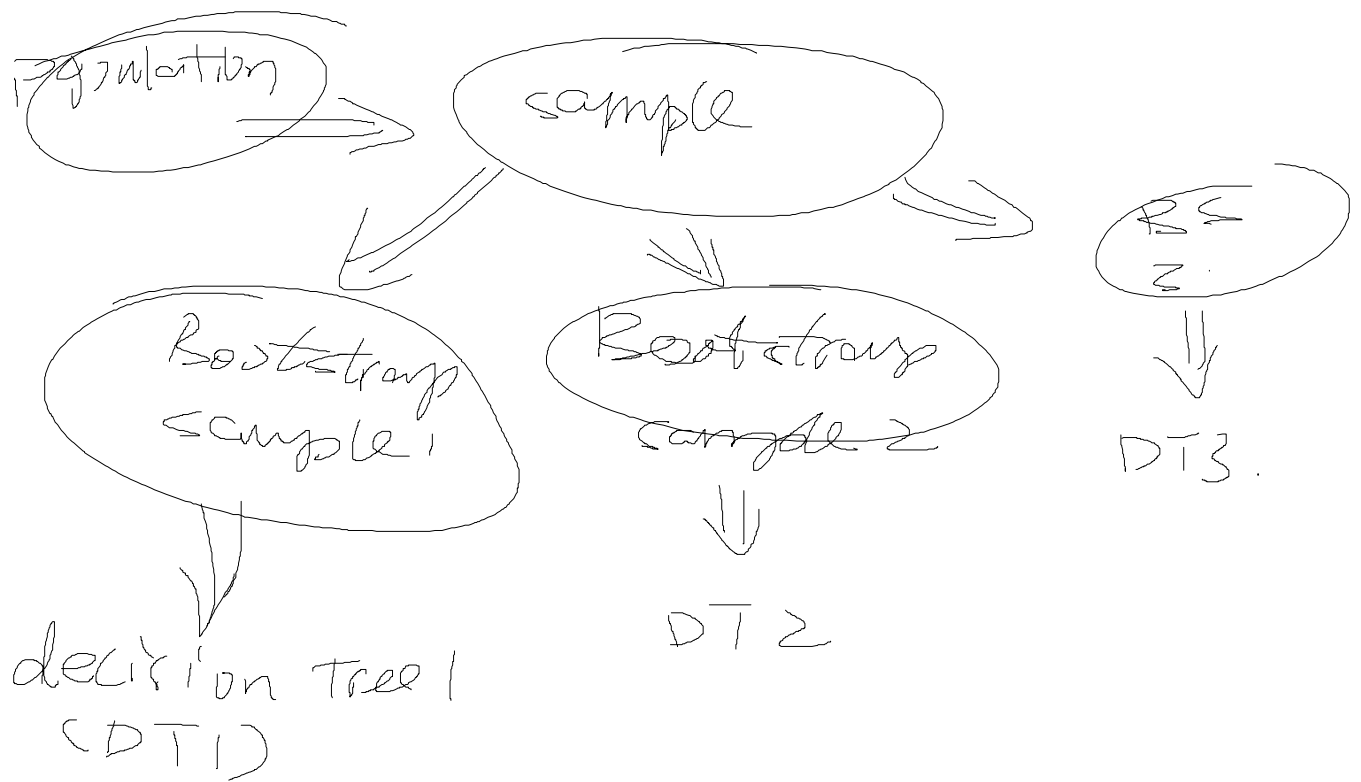
data 2

$$\boxed{\dot{I}_{root} = \frac{4}{9}}$$

data 2 > 2.5 / \ data 2 < 2.5

$I = 0$

$I = 0.5 \mid \boxed{\overline{I_1} = \frac{1}{3}}$

data 1 > 1.5 / \ data 1 < 1.5

$\boxed{I_2 = 0.5}$    $I = 0.5$    $I = 0.5$

data 2 > 1.5 / \ data 2 < 1.5

$I = 0$    $I = 0$    $\hat{I} = 0$    $1 = 0$

target = 0
target = 1

3 data 1

Gini Impurity

$$I_G(n) = 1 - \sum_{i=1}^{J} (P_i)^2$$

$$\overline{I}_{root} = 1 - \left( \left(\frac{2}{6}\right)^2 + \left(\frac{4}{6}\right)^2 \right) = 1 - \frac{20}{36} = 1 - \frac{5}{9} = \boxed{\frac{4}{9}}$$

$$I_1 = \left( 1 - \left( \left(\frac{0}{2}\right)^2 + \left(\frac{2}{2}\right)^2 \right) \right) \cdot \frac{2}{6} + \left( 1 - \left( \left(\frac{2}{4}\right)^2 + \left(\frac{2}{4}\right)^2 \right) \right) \cdot \frac{4}{6}$$

$$= \boxed{\frac{1}{3}}$$

Population

sample

BS 2.

Bootstrap sample 1

Bootstrap sample 2

DT3.

DT 2

decision Tree 1 (DT1)

as bootstrapping is random, each of them have a decision tree. these decision tree are random. this is random forest