# Decision Tree

It can be applied for both kind of problems

1. Regression technique
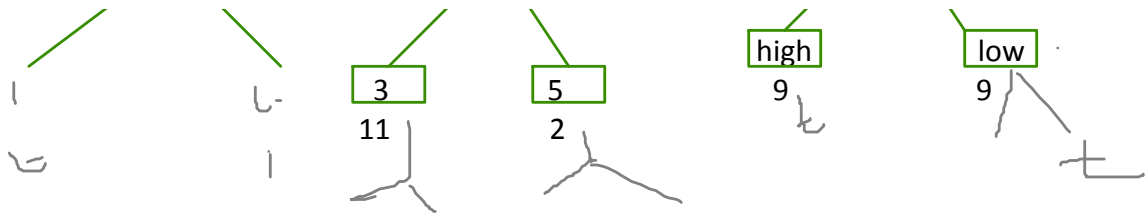2. Classification technique

```
                    ┌──────┐
                    │ CEO  │   --> Root node
                    └──────┘
              ┌─────────────────────┐
         ┌────────┐             ┌────────┐
         │ Sr. Mgr│             │ Sr. Mgr│
         └────────┘             └────────┘
          ┌──────┐             ┌──────────────┐
      ┌──────┐ ┌──────┐    ┌──────┐        ┌──────┐  Intermediate nodes
      └──────┘ └──────┘    └──────┘        └──────┘
                         ┌──────┐    ┌──────┐
                         └──────┘    └──────┘
                      ┌──┐
                      └──┘  ┌──────┐
                            └──────┘   Terminal Nodes
```

To identify the Root variables, we have many techniques

1. Miss classification Error (R only)
2. Gini Index (R and Python)  --> Classification
3. Information gain / Entropy ( Python)  --> Classification
4. Mean square error ( Python)  --> Regression
5. ID3
6. CHAID
7. Variance reduction technique
8. C4.5

-------------------------------------------------------------------------------------------------------------------------------------

```
                        credit (40)
          ┌───────────────┼───────────────┐
     ┌──────────┐    ┌──────────┐    ┌──────────┐
     │ Excellent│    │ fair     │    │ poor     │
     └──────────┘    └──────────┘    └──────────┘
          9              13              18
       ┌─────┐        ┌─────┐        ┌──────────┐
                      ┌───┐ ┌───┐  ┌──────┐ ┌──────┐
                      │ 3 │ │ 5 │  │ high │ │ low  │
                      └───┘ └───┘  └──────┘ └──────┘
                                      9        9
```

```
                                          high        low
                  3        5              9           9
                 11        2
```

----------------------------------------------------------------------------------------

Gini Index

Gini =

$$\sum_{i=1}^{C} (p_i)^2$$

$\Sigma \; Wi(p2 + q2)$

Gini index will be calculated for each of the X variable and will see which X variable contains highest Gini index that variable will become Root Variable.

**Steps to Calculate Gini for a split**

1. Calculate Gini for sub-nodes, using formula sum of square of probability for success and failure (p^2+q^2).
2. Calculate Gini for split using weighted Gini score of each node of that split

3. To calculate the exact Gini score by multiplying Total(step 1 * step 2)

We have to verify same calculative method for each of the x variable and finalized with one of the X variable.

```
                    Gender (30)

         Male(20)                    Female(10)

     IX          X              IX            X
     13          7              1             9

  Yes      No
 (70%)    (30%)
```