

Decision Tree

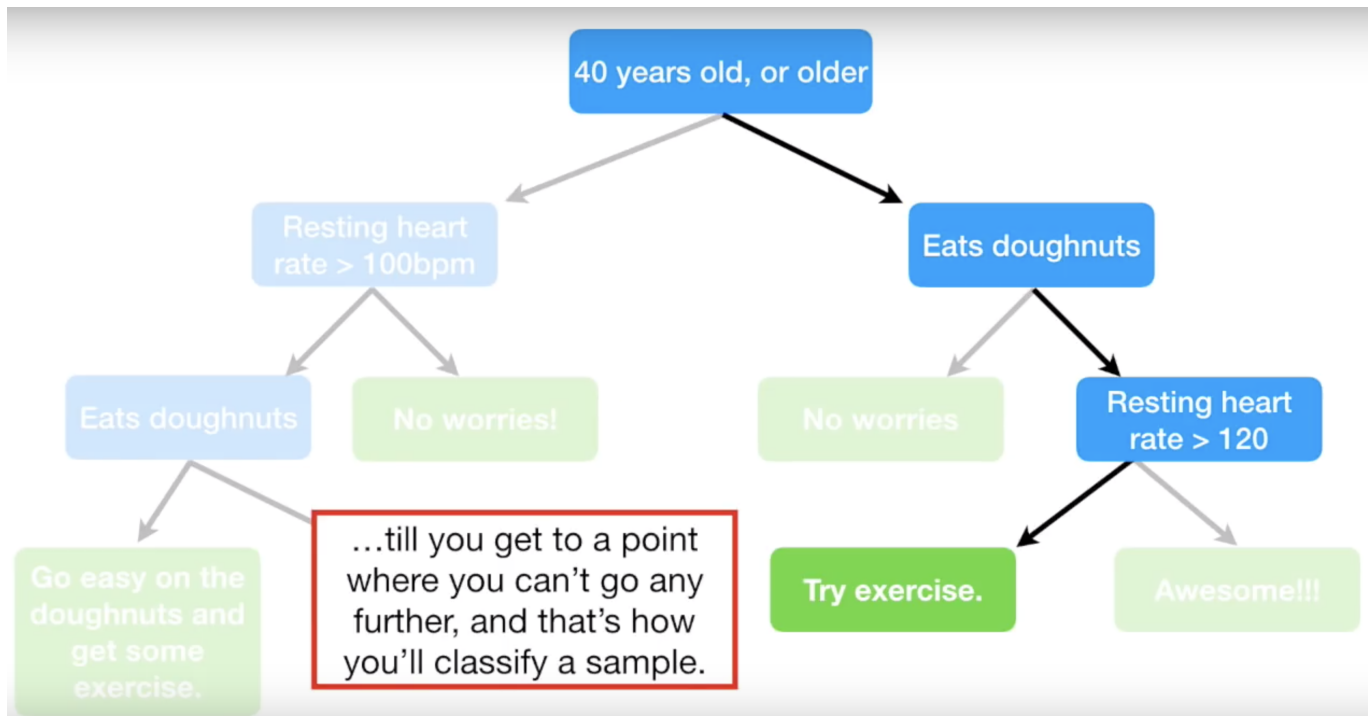
Ein **Decision Tree** stellt eine Frage und klassifiziert dann die Person anhand der Antwort.

Person hat einen hohen Puls - Person sollte zum Arzt

Person hat keinen hohen Puls - Person muss nicht zum Arzt

Klassifizierung können Kategorien oder Zahlen sein!

Man kann auch beides kombinieren!



Namensgebung

Root-Node:

Die **root-node** befindet sich am obersten Ende des Baums. Von dort baut sich der Baum abwärts auf.

Node:

Sind die einzelnen **nodes** des Baumes.

Leaf-nodes:

Bilden den Abschluss eines einzelnen **strang**. Pfeile zeigen zu den jeweiligen **leaf-nodes** aber keine Pfeile gehen von dort weg.

Decision Tree aufbauen

Um zu entscheiden welcher der Faktoren die **root-node** betrachtet man wie gut jeder der einzelnen Faktoren den Outcome (das was man vorhersagen will) beschreibt. Da in den meisten Fällen keiner der Faktoren den Outcome zu 100% richtig beschreibt werden diese **impure** genannt. Diese Unreinheit kann zum Beispiel mit dem **Gini-Koeffizient** berechnet werden.

\$\$\$Für jedes Blatt - Gini Index = 1 - (WS für 'JA')^2 - (WS für 'NEIN')^2\$\$\$

Beispiel: Herzkrankheit vorhersagen

Root-Node: Brustschmerzen Leaf-Node1: (Herzkrank)

Herzkrankheit Ja: 105
Herzkrankheit Nein: 39

Leaf-Node2: (Nicht Herzkrank)

Herzkrankheit Ja: 34
Herzkrankheit Nein: 125

--> Gini-Index für Node1:

$$1 - \left(\frac{105}{105 + 39}\right)^2$$

--> Gini-Index für Node2:

$$1 - \left(\frac{125}{34 + 125}\right)^2$$

Um nun den Gesamten Gini-Index für **Brustschmerzen** als **root-node** zu bestimmen berechnet man die Anzahl aller Beobachtungsmern in dem jeweiligen Node gesehen auf die gesamt Anzahl von Node 1 und Node 2:

$$\left(\frac{144}{144 + 159}\right) * \text{GINI-Index-Node1} + \left(\frac{159}{144 + 159}\right) * \text{GINI-Index-Node2}$$

Das wird jetzt für alle Faktoren gemacht (zum Beispiel noch von geblockten Arterien). Der Faktor mit dem geringsten Gini-Index ist die **root-node**.

Der Baum baut sich nun nach und nach auf - sobald man die Daten nicht mehr trennen kann entsteht eine **leaf-node**.

Random Forests

random forest basieren auf der **decision tree** Methode, kombinieren aber alle Vorteile davon mit einer Lösung für die Nachteile (schlechte Genauigkeit von Decision Trees). Im Grunde ist ein Random forest eine Vielzahl von gebootstrapteten Decision Trees.

Wichtig dabei:

Für einen Random forest wird eine beliebige Anzahl von Parametern aus den Daten gewählt. **Doppelte Ziehung einer Zeile für einen decision tree ist ausdrücklich erlaubt!**

So werden ganz viele unterschiedliche Decision trees erstellt.

Um nun einen Datensatz zu evaluieren wird dieser durch jeden dieser erstellten Bäume gejagt und alle Ergebnisse der einzelnen Bäume (ja - nein) aufsummiert. Die Mehrheit gewinnt und gibt das finale Label an.

Genauigkeit des Random Forest

Die Genauigkeit wird mit den **ausgelassenen** Zeilen des Bootstrappings berechnet da man deren Ausgang kennt. Diese werden ebenfalls einfach durch den random forest gejagt und man kann somit dann eine Genauigkeit berechnen. In der Regel ist diese um einiges besser als bei einem einfachen Entscheidungsbaum!