

L12 Logistic Regression

Prof. Xun Jiao

Review

- Logistic Regression
 - Hypothesis

$h_{\theta}(x)$ should give $p(y = 1 \mid x; \theta)$

$$h_{\theta}(x) = \frac{1}{1 + e^{-\theta^T x}}$$

- Decision boundary

Outline

- Cost function
- Gradient descent
- Multi-class Classification
- Scikit Learn Library
- Mini-Project

Cost Function

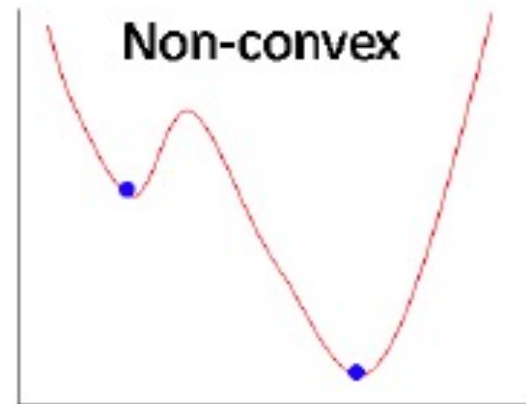
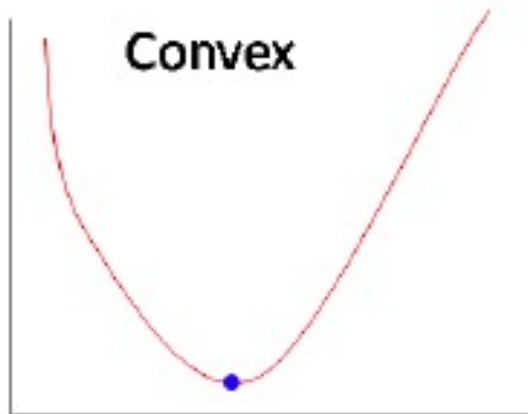
Can't just use squared loss as in linear regression:

$$J(\boldsymbol{\theta}) = \frac{1}{2n} \sum_{i=1}^n \left(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) - y^{(i)} \right)^2$$

– Using the logistic regression model

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^T \mathbf{x}}}$$

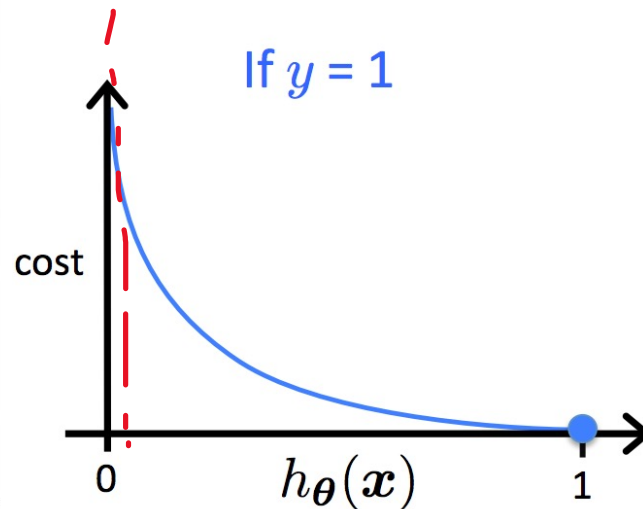
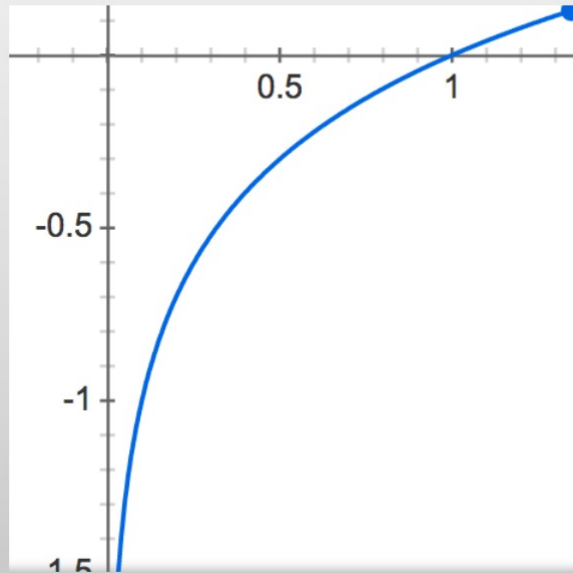
results in a non-convex optimization



Cost function

$$\text{cost}(h_{\theta}(\mathbf{x}), y) = \begin{cases} -\log(h_{\theta}(\mathbf{x})) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(\mathbf{x})) & \text{if } y = 0 \end{cases}$$

Aside: Recall the plot of $\log(z)$



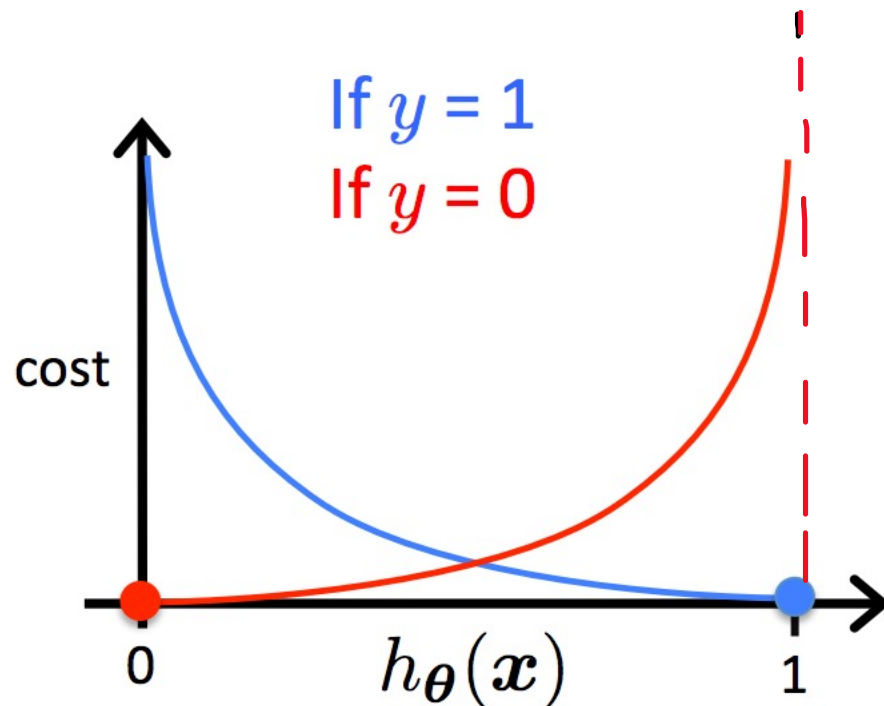
If $y = 1$

- Cost = 0 if prediction is correct
- As $h_{\theta}(\mathbf{x}) \rightarrow 0$, cost $\rightarrow \infty$
- Captures intuition that larger mistakes should get larger penalties
 - e.g., predict $h_{\theta}(\mathbf{x}) = 0$, but $y = 1$

$$\text{cost}(h_{\theta}(\mathbf{x}), y) = \begin{cases} -\log(h_{\theta}(\mathbf{x})) & \text{if } y = 1 \\ -\log(1 - h_{\theta}(\mathbf{x})) & \text{if } y = 0 \end{cases}$$

If $y = 0$

- Cost = 0 if prediction is correct
- As $(1 - h_{\theta}(\mathbf{x})) \rightarrow 0$, $\text{cost} \rightarrow \infty$
- Captures intuition that larger mistakes should get larger penalties



Cost function

$$\text{cost}(h_{\boldsymbol{\theta}}(\mathbf{x}), y) = \begin{cases} -\log(h_{\boldsymbol{\theta}}(\mathbf{x})) & \text{if } y = 1 \\ -\log(1 - h_{\boldsymbol{\theta}}(\mathbf{x})) & \text{if } y = 0 \end{cases}$$



$$J(\boldsymbol{\theta}) = -\sum_{i=1}^n \left[y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right]$$



$$J(\boldsymbol{\theta}) = \sum_{i=1}^n \text{cost}(h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}), y^{(i)})$$

Overall Look

Logistic regression objective:

$$\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$$

$$J(\boldsymbol{\theta}) = - \sum_{i=1}^n \left[y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right]$$

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = g(\boldsymbol{\theta}^{\top} \mathbf{x})$$

$$g(z) = \frac{1}{1 + e^{-z}}$$

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = \frac{1}{1 + e^{-\boldsymbol{\theta}^{\top} \mathbf{x}}}$$

$$J(\boldsymbol{\theta}) = - \sum_{i=1}^n \left[y^{(i)} \log h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)}) + (1 - y^{(i)}) \log (1 - h_{\boldsymbol{\theta}}(\mathbf{x}^{(i)})) \right]$$

Want $\min_{\boldsymbol{\theta}} J(\boldsymbol{\theta})$

- Initialize $\boldsymbol{\theta}$
- Repeat until convergence

$$\theta_j \leftarrow \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\boldsymbol{\theta})$$

simultaneous update
for $j = 0 \dots d$

- Initialize θ
- Repeat until convergence (simultaneous update for $j = 0 \dots d$)

$$\theta_0 \leftarrow \theta_0 - \alpha \sum_{i=1}^n \left(h_{\theta} \left(\mathbf{x}^{(i)} \right) - y^{(i)} \right)$$

$$\theta_j \leftarrow \theta_j - \alpha \left[\sum_{i=1}^n \left(h_{\theta} \left(\mathbf{x}^{(i)} \right) - y^{(i)} \right) x_j^{(i)} \right]$$



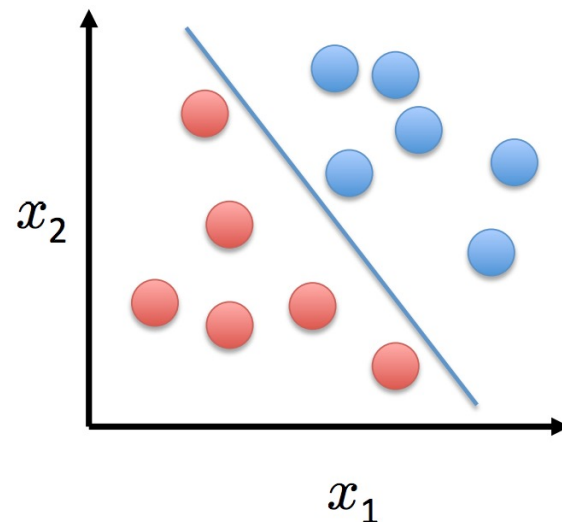
Does this look familiar to you?

Multiclassification

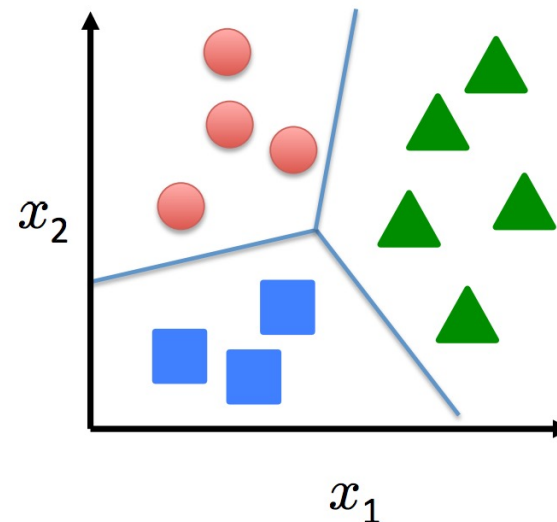
Disease diagnosis: healthy / cold / flu / pneumonia

Object classification: desk / chair / monitor / bookcase

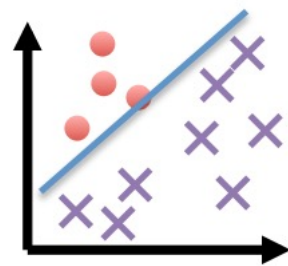
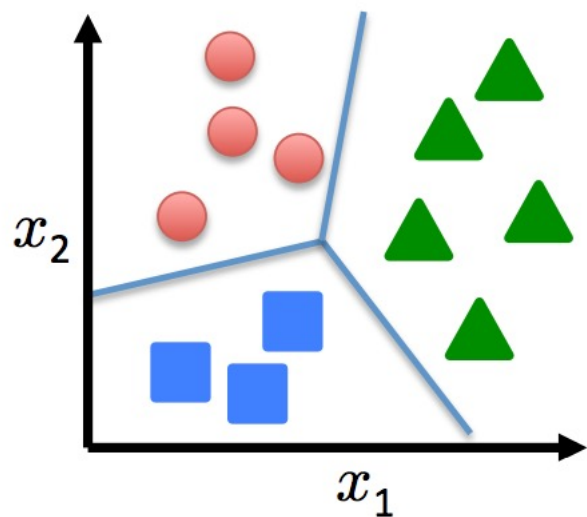
Binary classification:



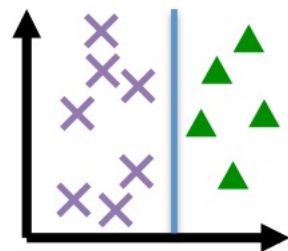
Multi-class classification:



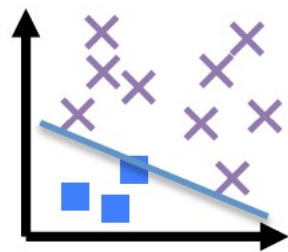
Split into One vs Rest:



$$p(y = 1 \mid \mathbf{x}; \boldsymbol{\theta})$$



$$p(y = \textcolor{red}{1} \mid \mathbf{x}; \boldsymbol{\theta})$$



$$p(y = \textcolor{red}{3} \mid \mathbf{x}; \boldsymbol{\theta})$$

$$h_{\boldsymbol{\theta}}(\mathbf{x}) = p(y = c \mid \mathbf{x}; \boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_C)$$

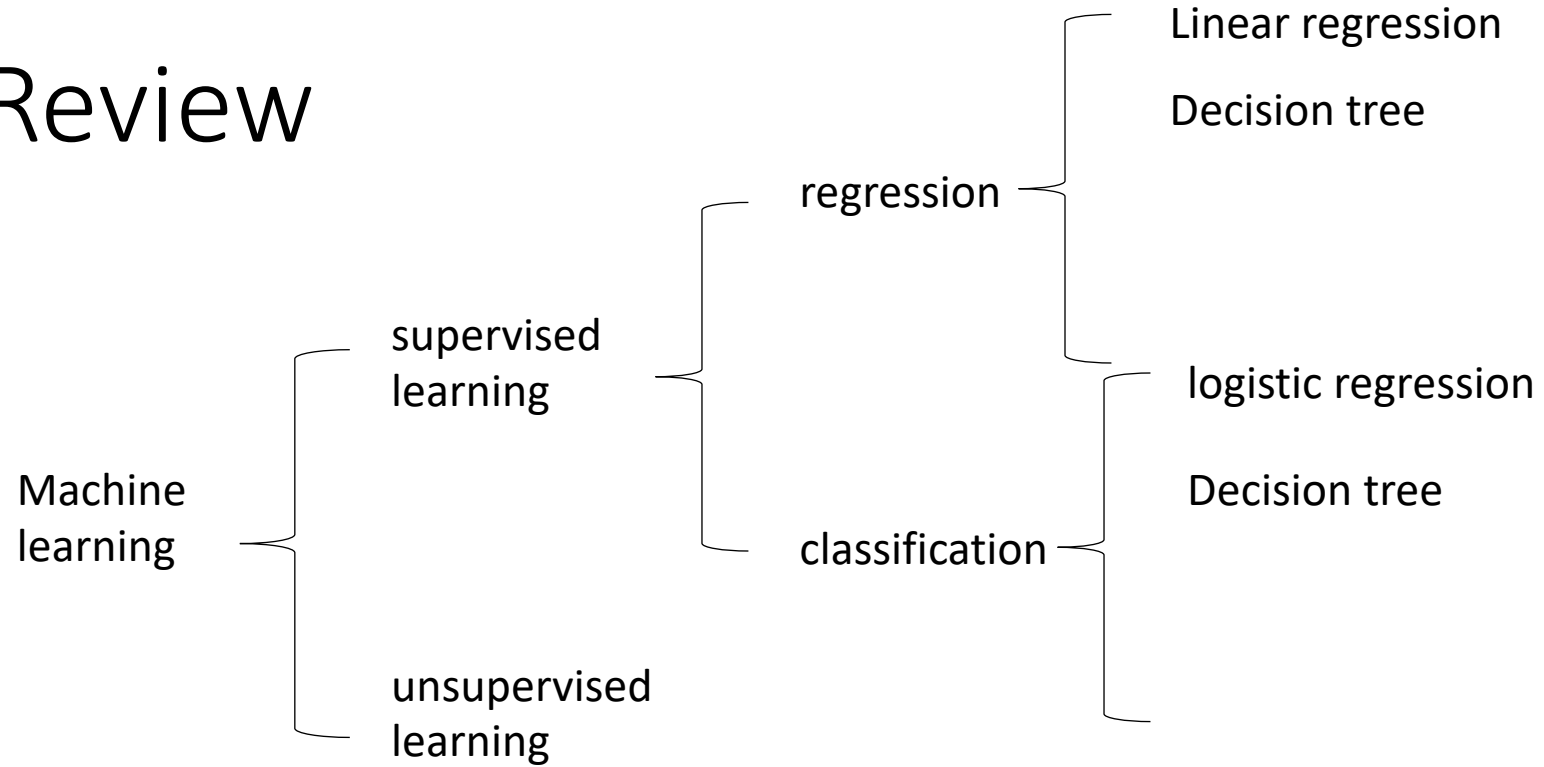
Train a logistic regression classifier for each class i
to predict the probability that $y = i$

Predict class label as the most probable label

$$\max_c h_c(\boldsymbol{x})$$

Decision Tree

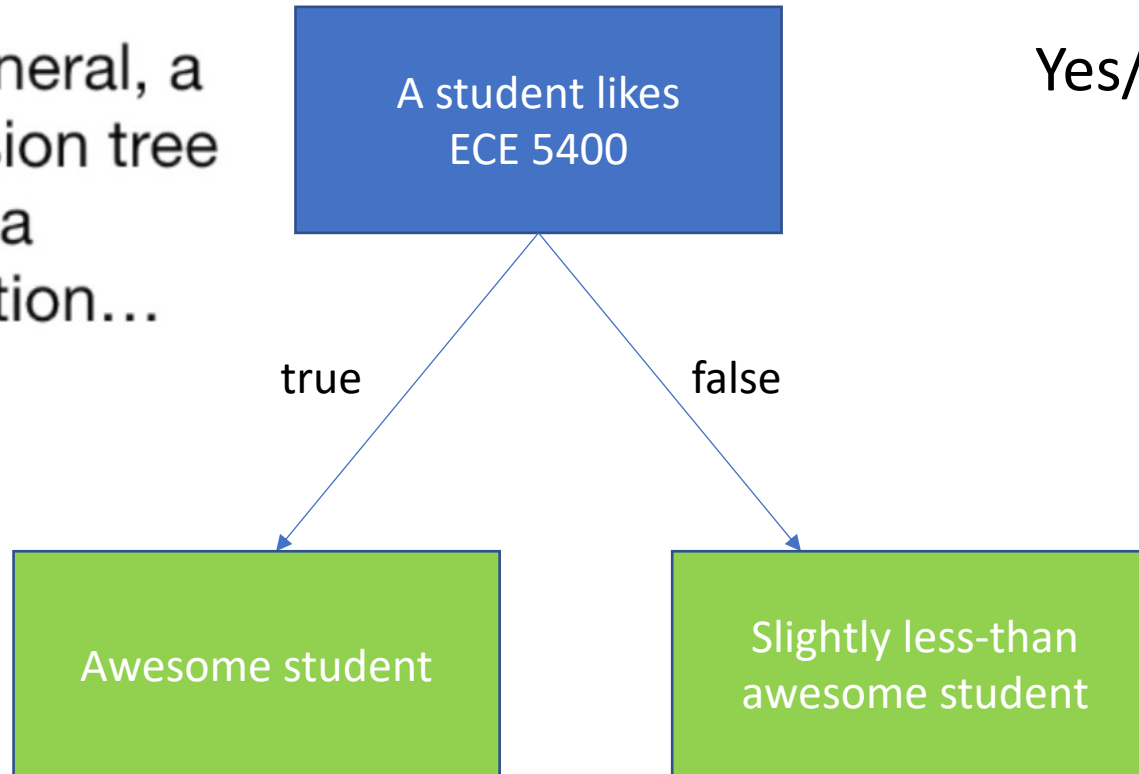
Review



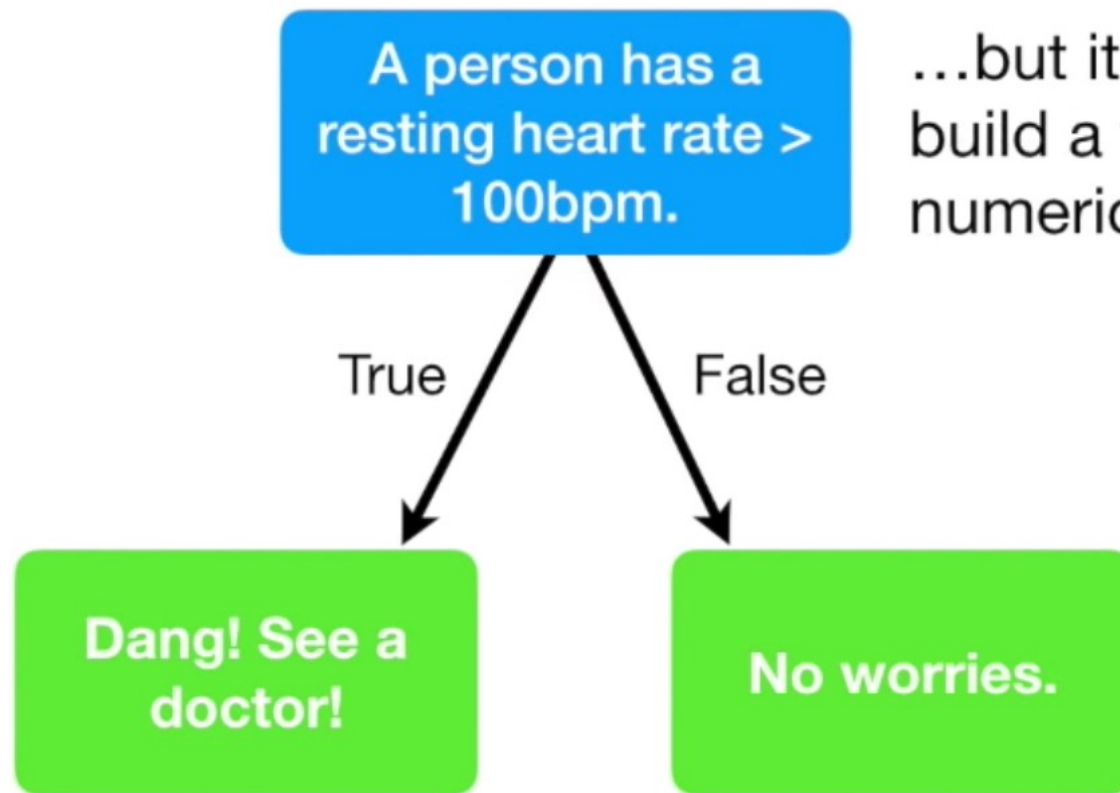
Hypothesis

In general, a decision tree asks a question...

...and then classifies the person based on the answer.

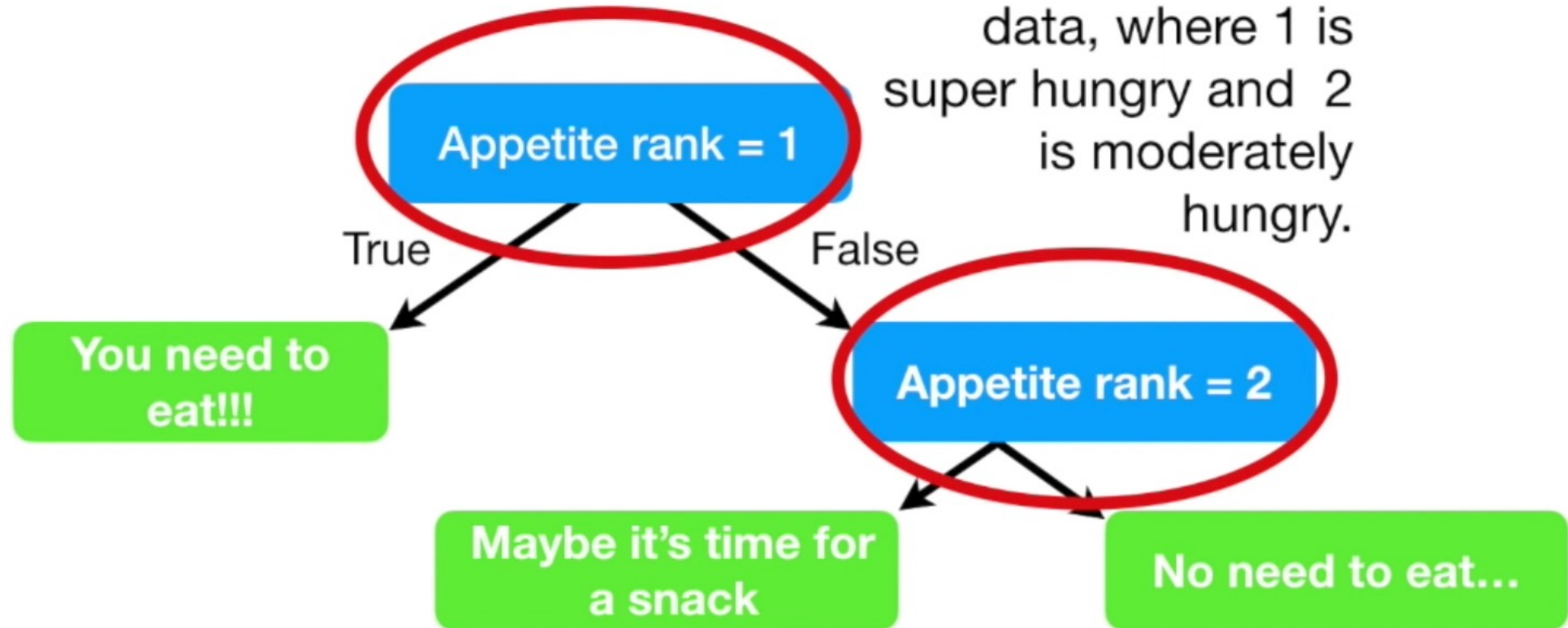


Yes/No question!

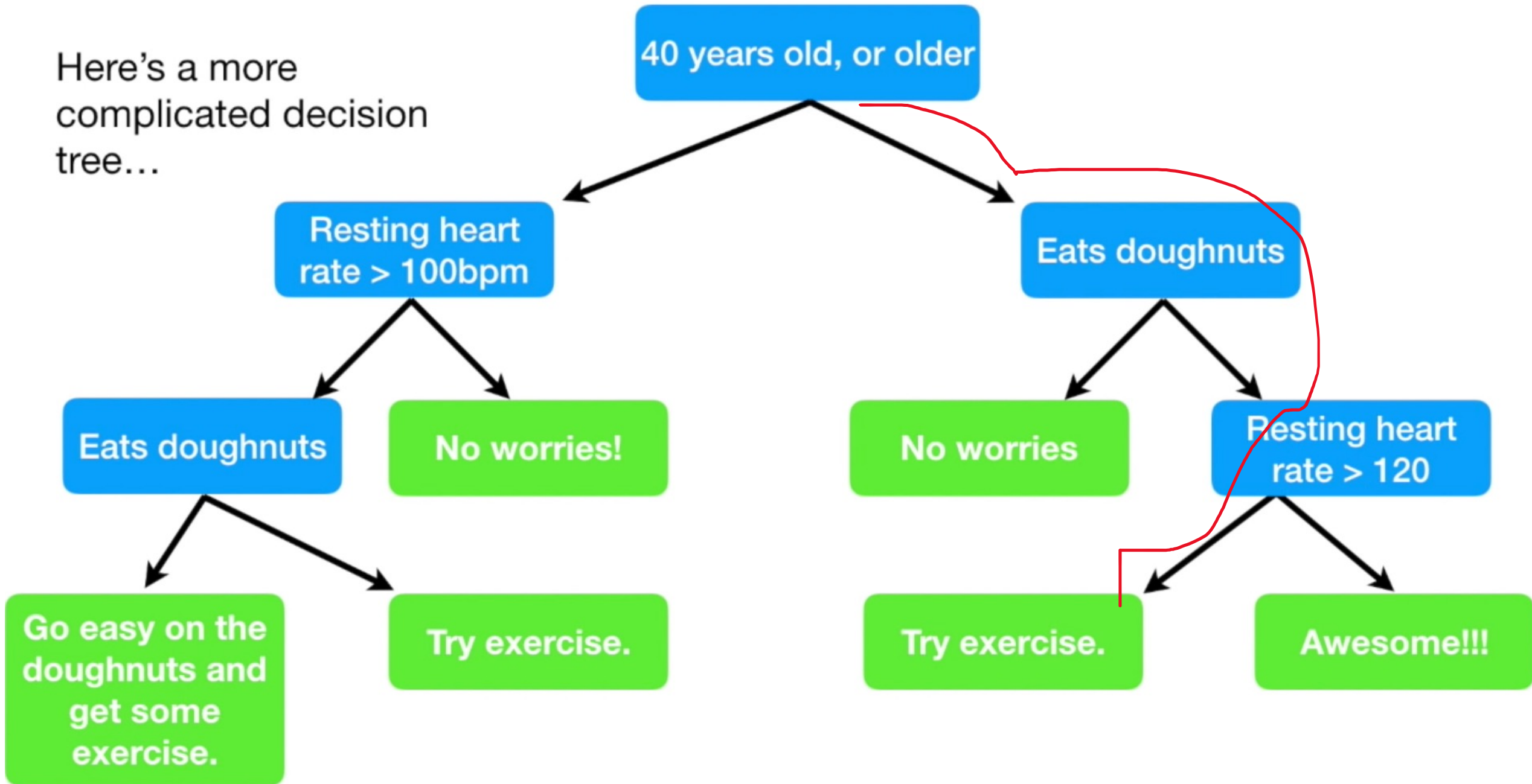


...but it is just as easy to build a tree from numeric data.

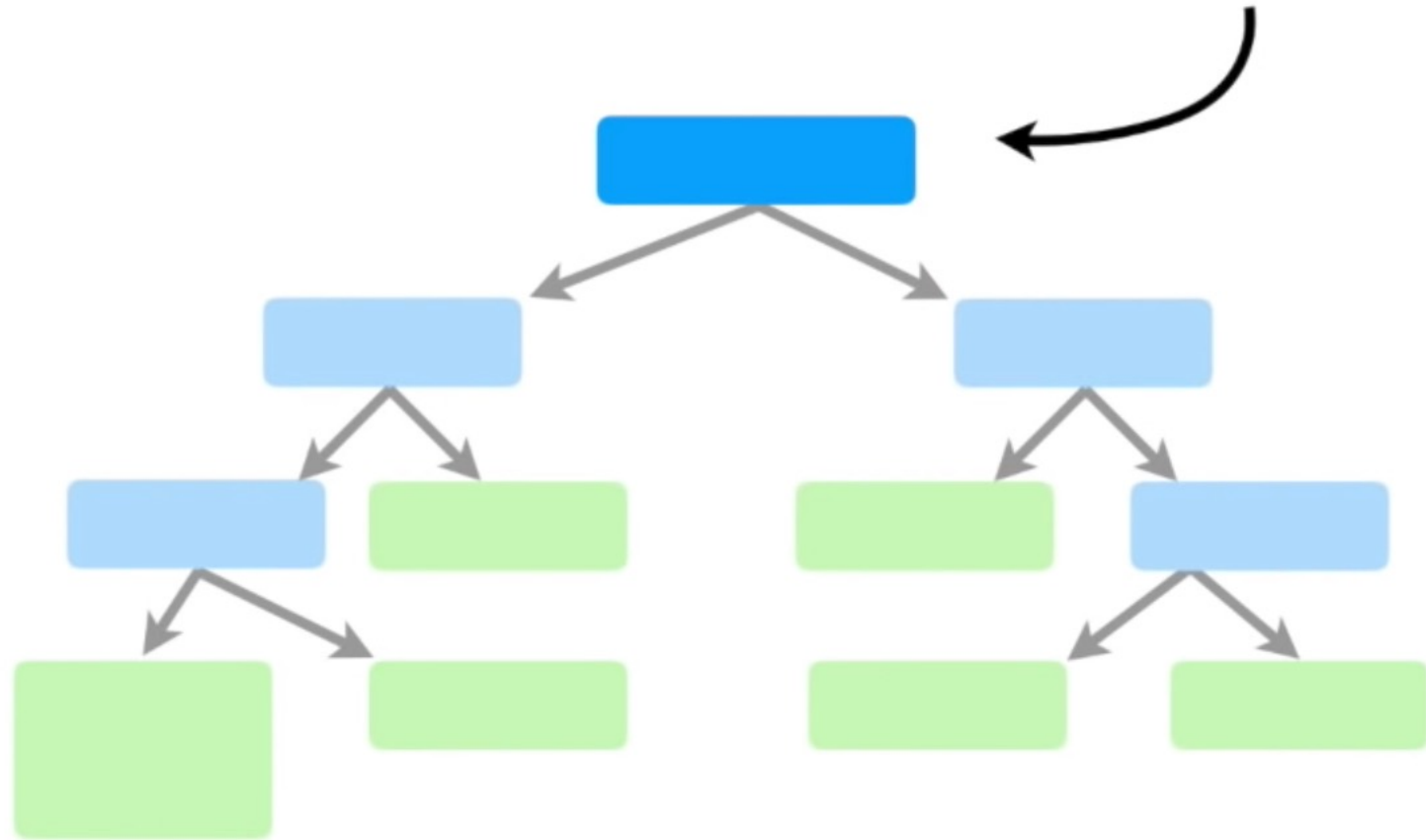
This decision tree is based on **ranked** data, where 1 is super hungry and 2 is moderately hungry.



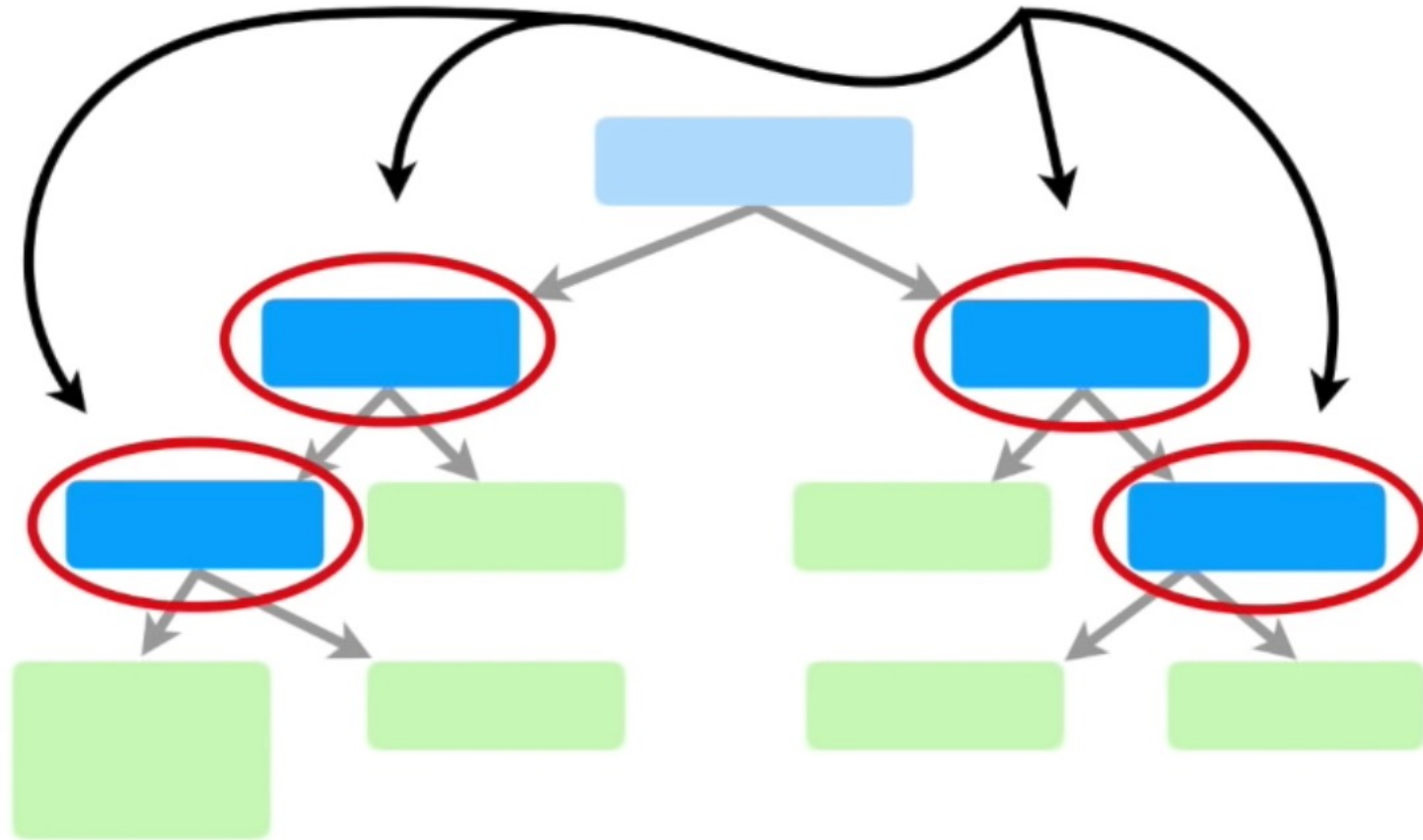
Here's a more complicated decision tree...



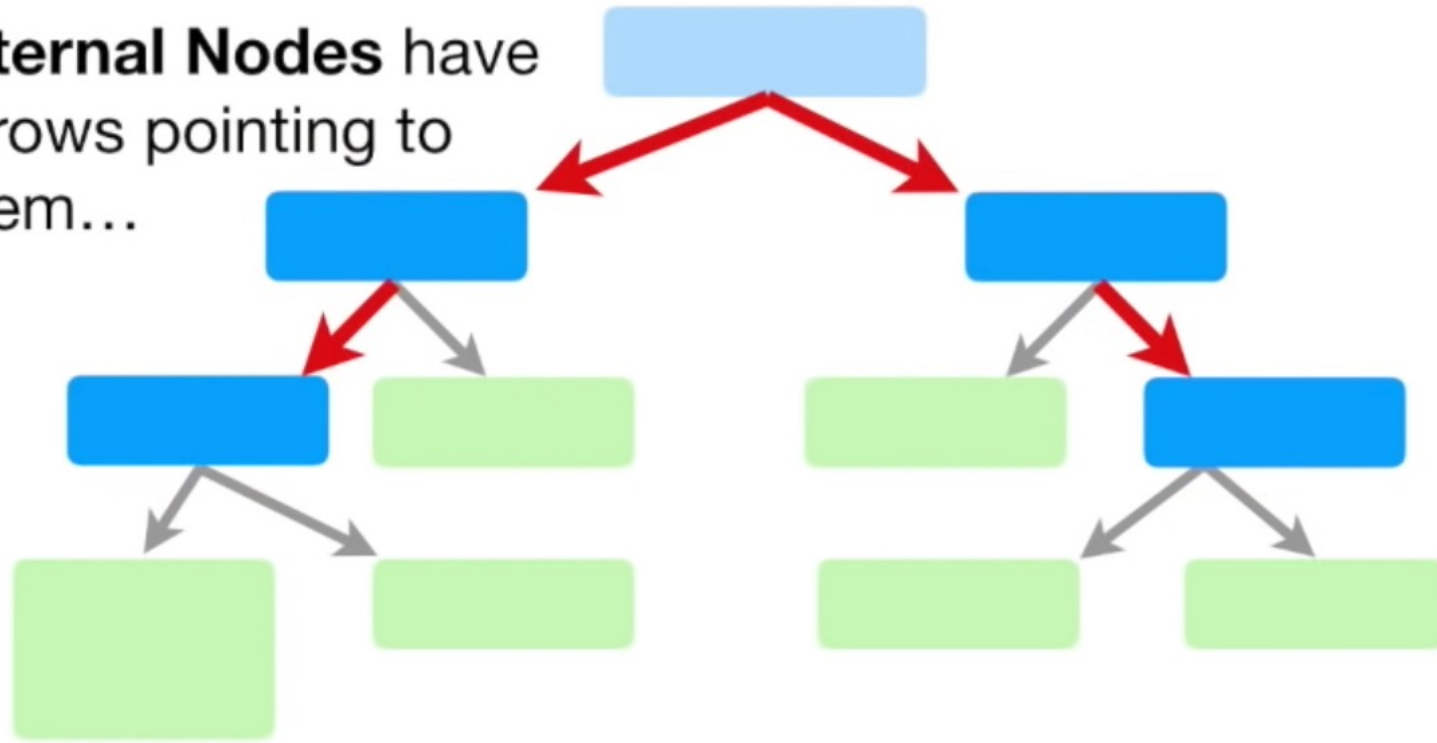
The very top of the tree is called the “**Root Node**” or just “**The Root**”

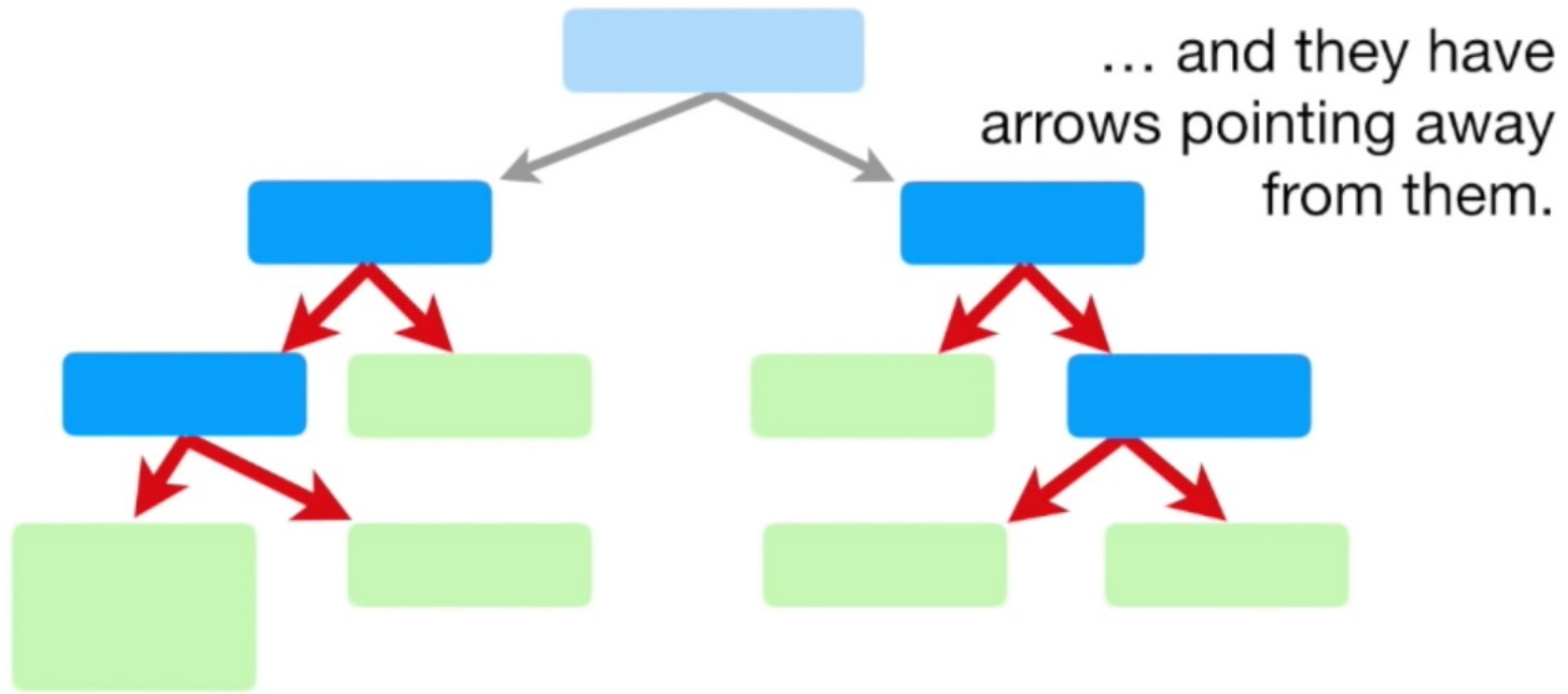


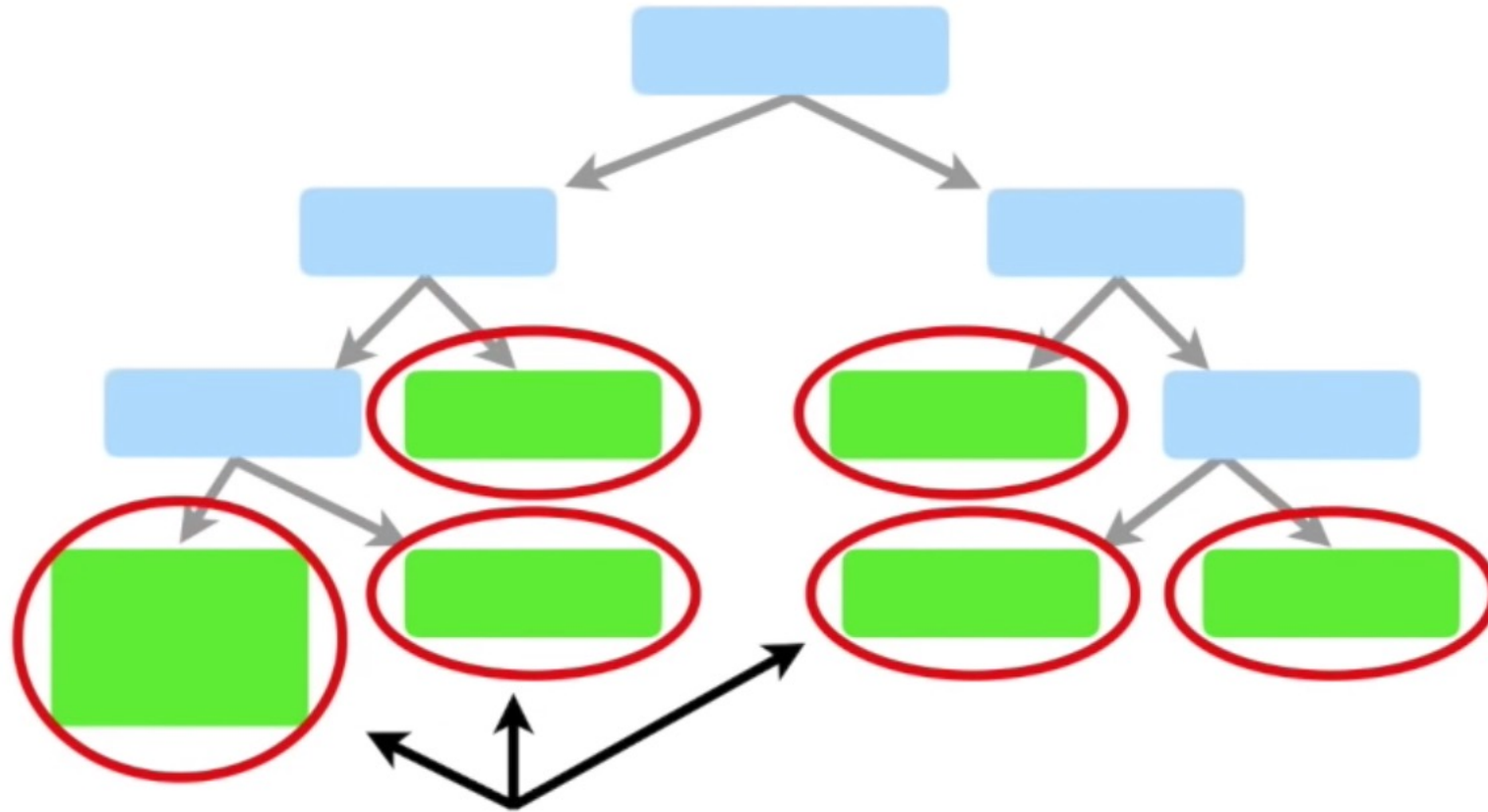
These are called “**Internal Nodes**”, or just “**Nodes**”.



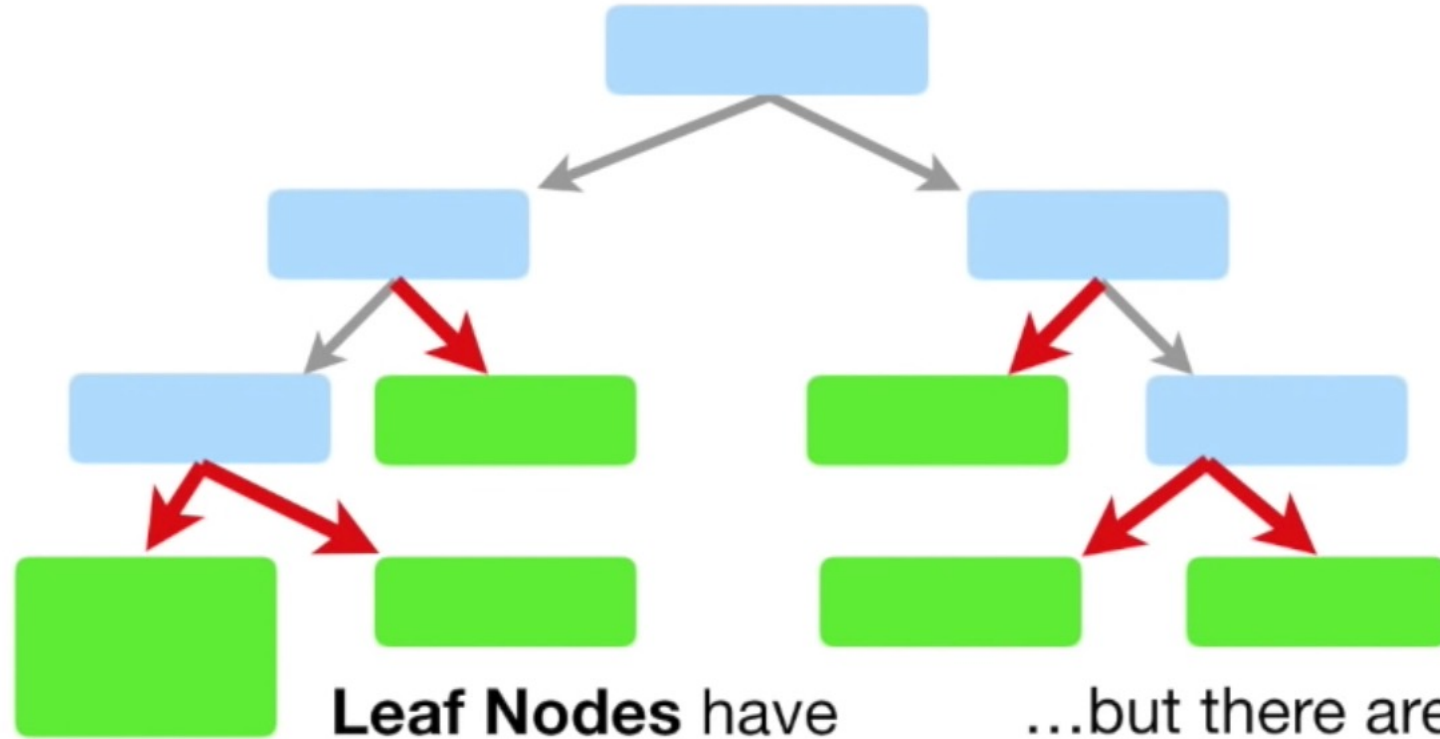
Internal Nodes have
arrows pointing to
them...







Lastly, these are called “**Leaf Nodes**”, or just “**Leaves**”



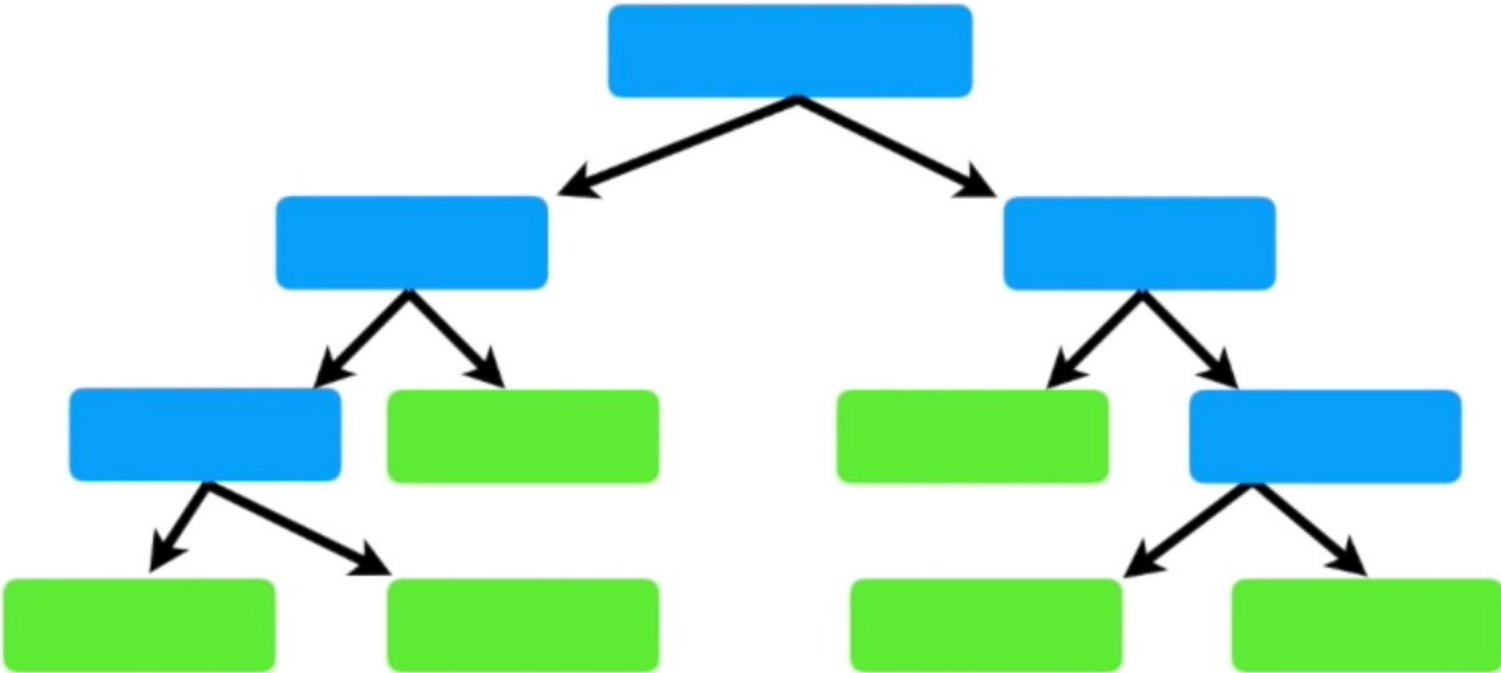
Leaf Nodes have
arrows pointing to
them...

...but there are no
arrows pointing
away from them.

Now we are ready to talk about how to go from a raw table of data...

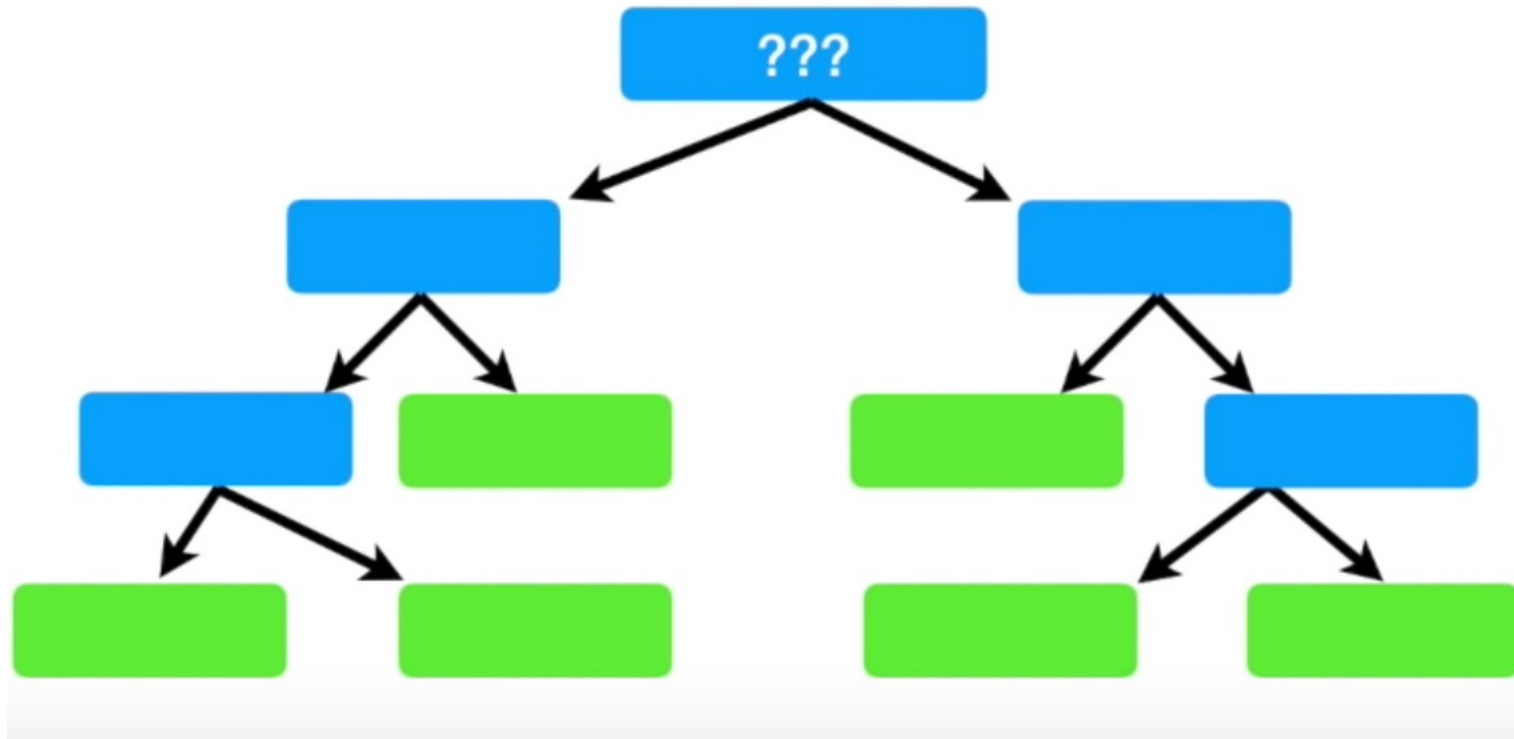
Chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
No	No	No	No
Yes	Yes	Yes	Yes
Yes	Yes	No	No
Yes	No	???	Yes
etc...	etc...	etc...	etc...

...to a decision tree!!!



First Step: who would be the top (root)?

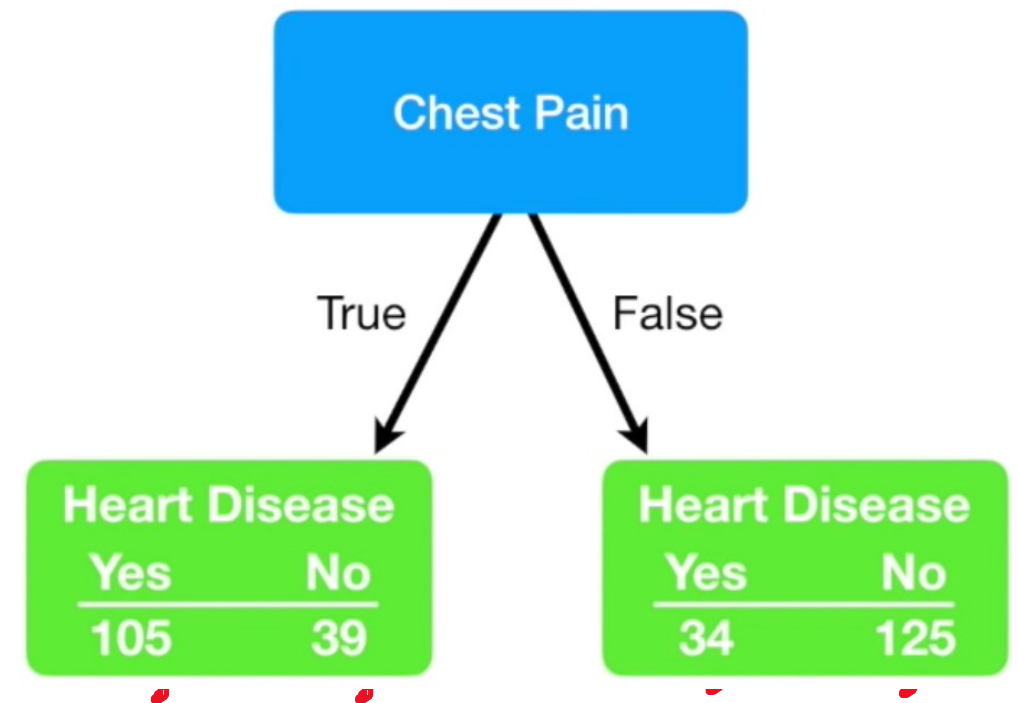
- Chest pain? Blood circulation? Blocked Arteries?
-



If we only use chest pain?

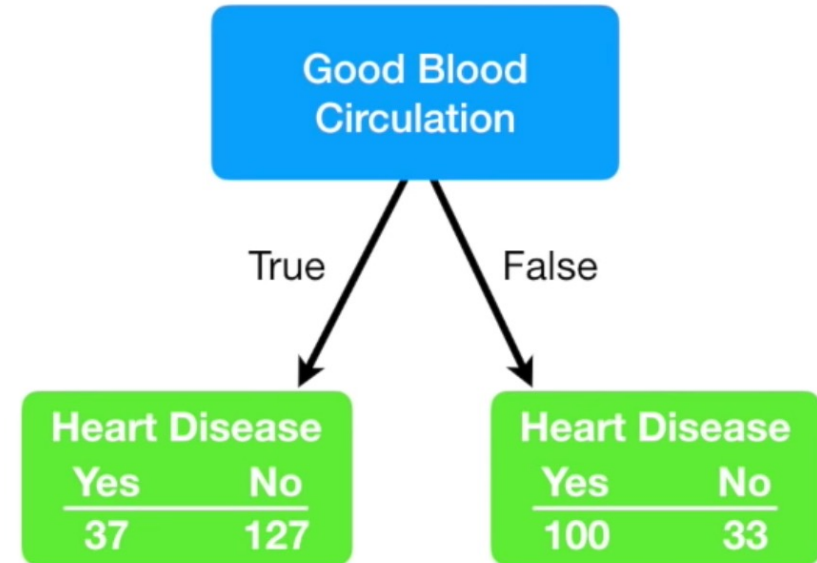
Chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
No	No	No	No
Yes	Yes	Yes	Yes
Yes	Yes	No	No
Yes	No	???	Yes
etc...	etc...	etc...	etc...

Here's a little tree that only takes chest pain into account.



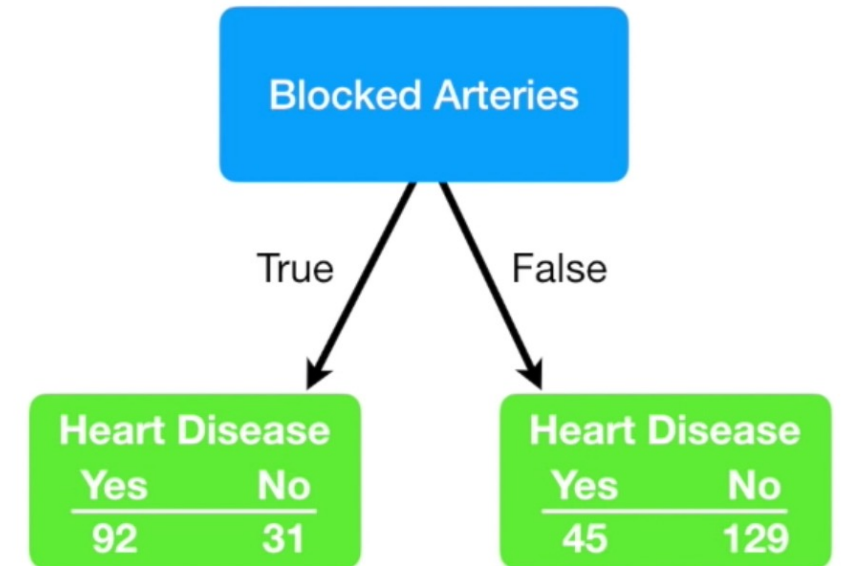
If we only use blood circulation

Chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
No	No	No	No
Yes	Yes	Yes	Yes
Yes	Yes	No	No
Yes	No	???	Yes
etc...	etc...	etc...	etc...

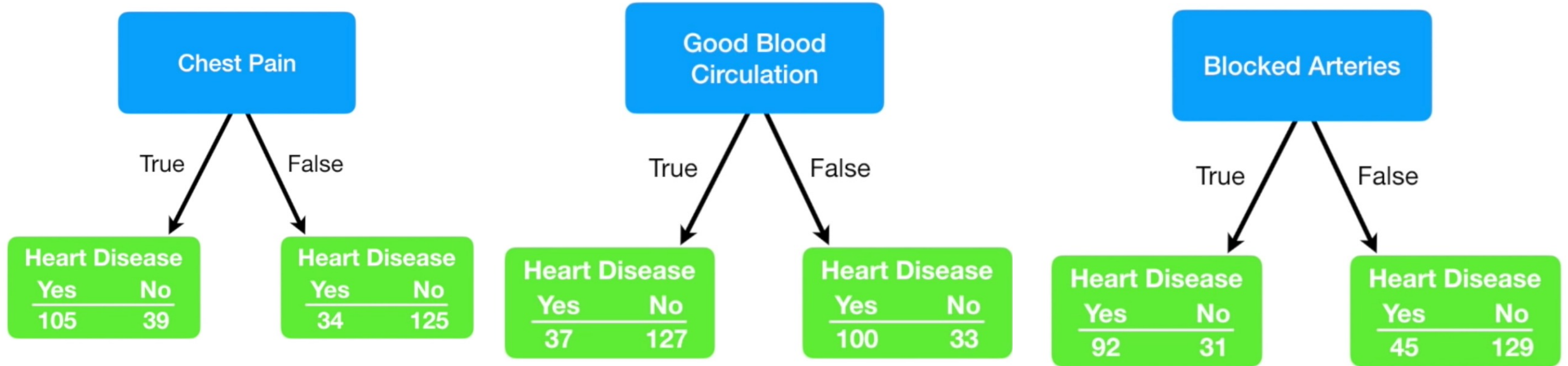


If we only use blocked arteries?

Chest Pain	Good Blood Circulation	Blocked Arteries	Heart Disease
No	No	No	No
Yes	Yes	Yes	Yes
Yes	Yes	No	No
Yes	No	???	Yes
etc...	etc...	etc...	etc...



Which one has best separation ability?



Good, but not perfect:
Can mostly separate, but still
some errors

Good, but not perfect

Good, but not perfect

Think about it: what can be the worst situation???