# Feature importance

# Decision tree: induction



Age > 45?

Age > 65?

Income > 10000?

Degree > 2?

Height > 1.7?

Age > 50?

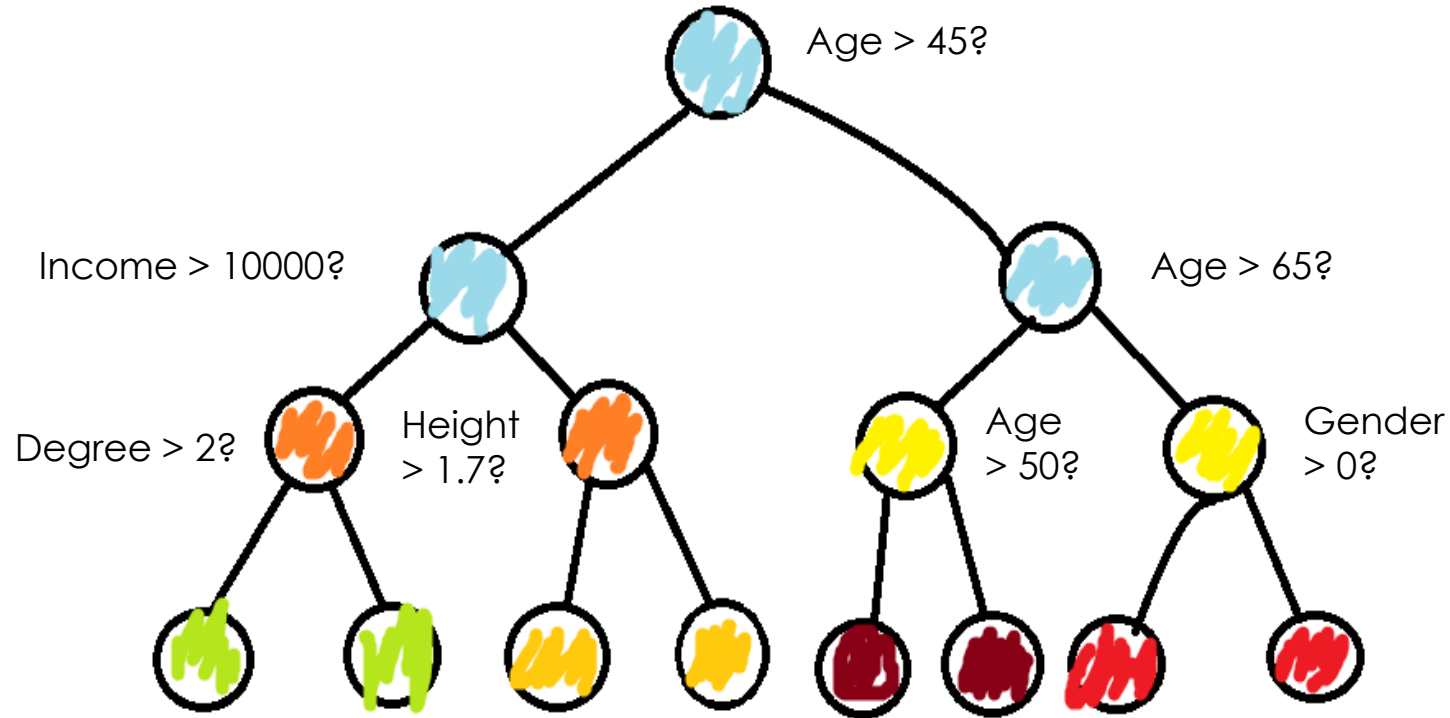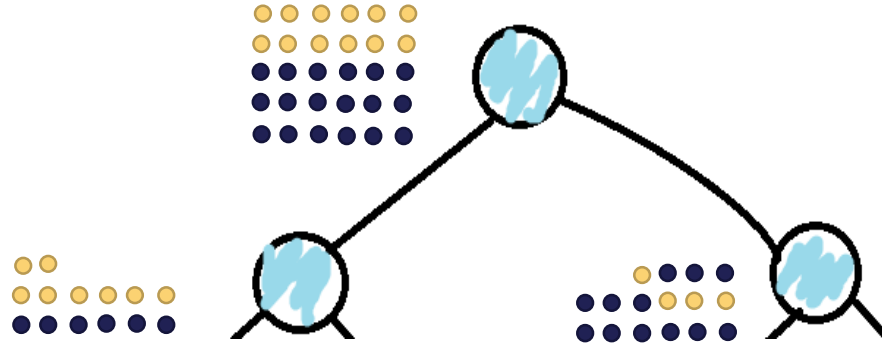Gender > 0?

At each note there is an increase in purity.
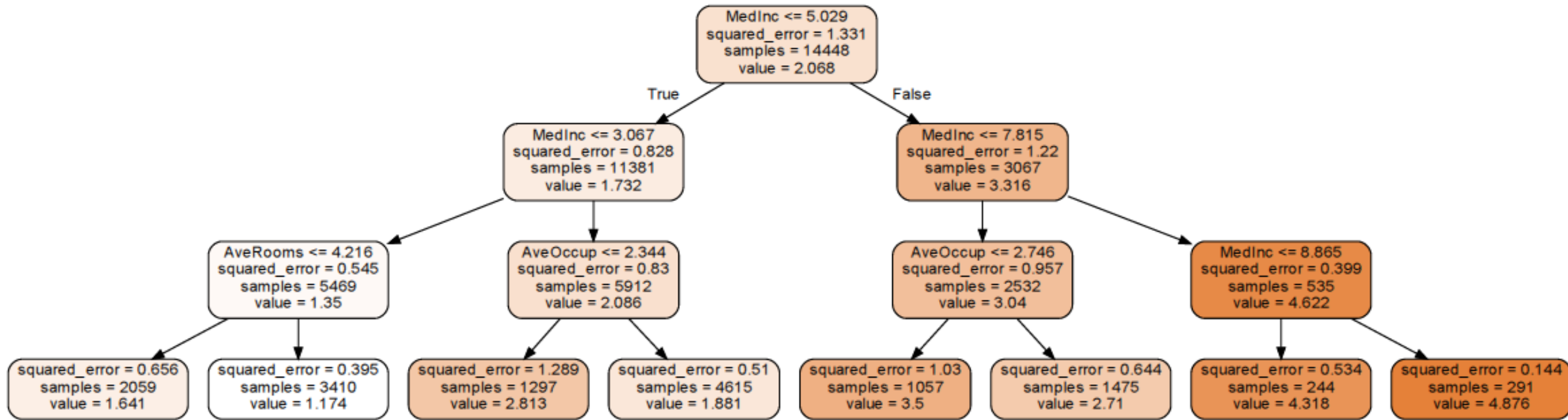
# Feature importance



- How often a feature is chosen.

- How big the increase in purity is.
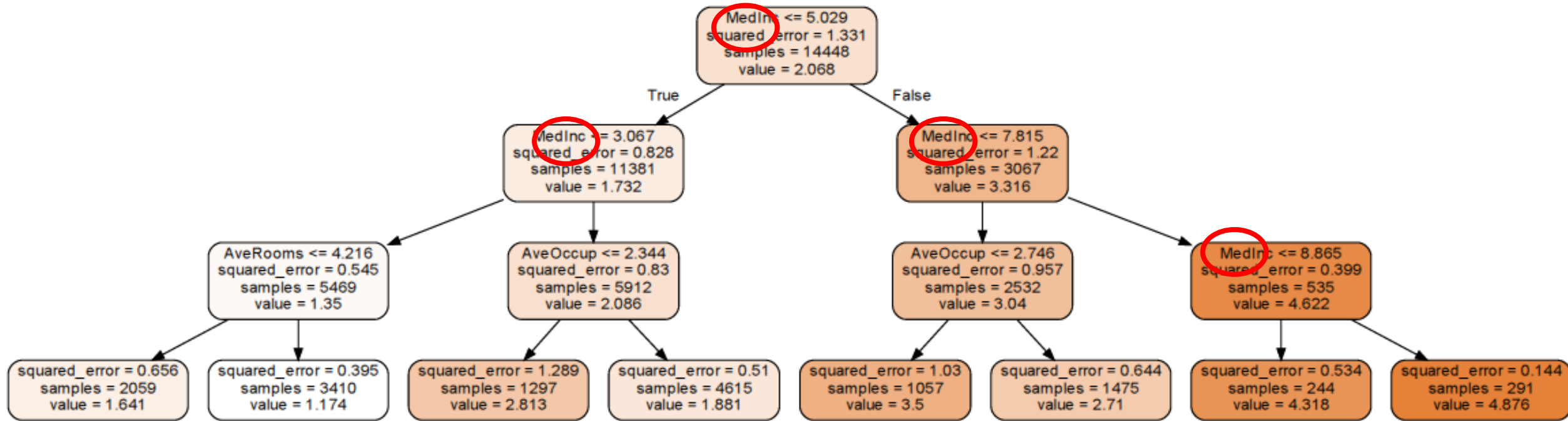
# Feature importance



```
importances[node.feature] +=

weighted_n_node_samples * node.impurity –

(left.weighted_n_node_samples * left.impurity +

right.weighted_n_node_samples * right.impurity)
```
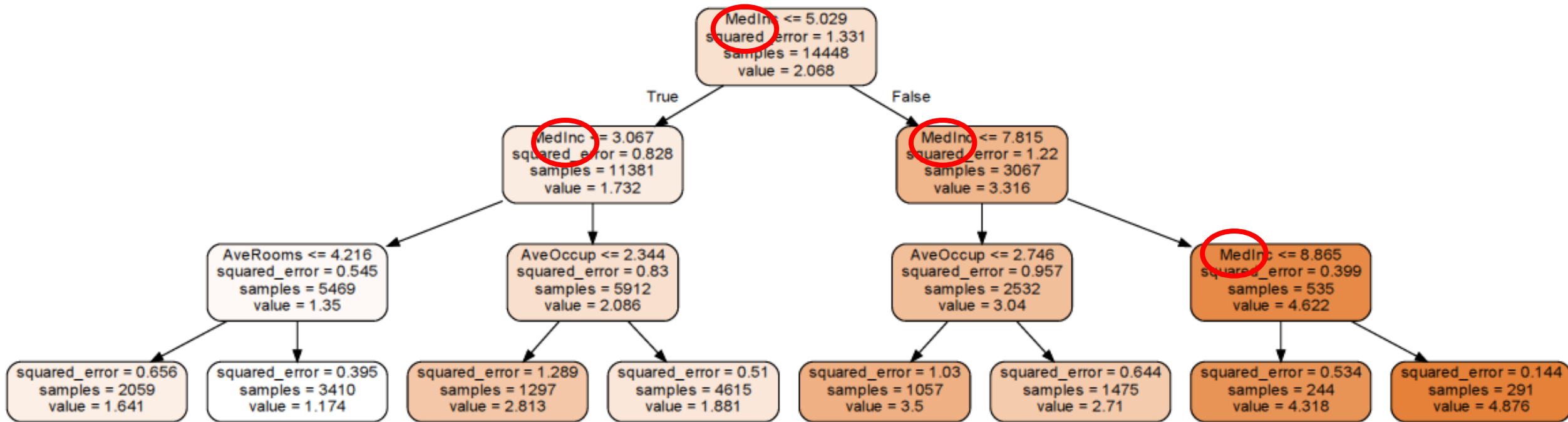
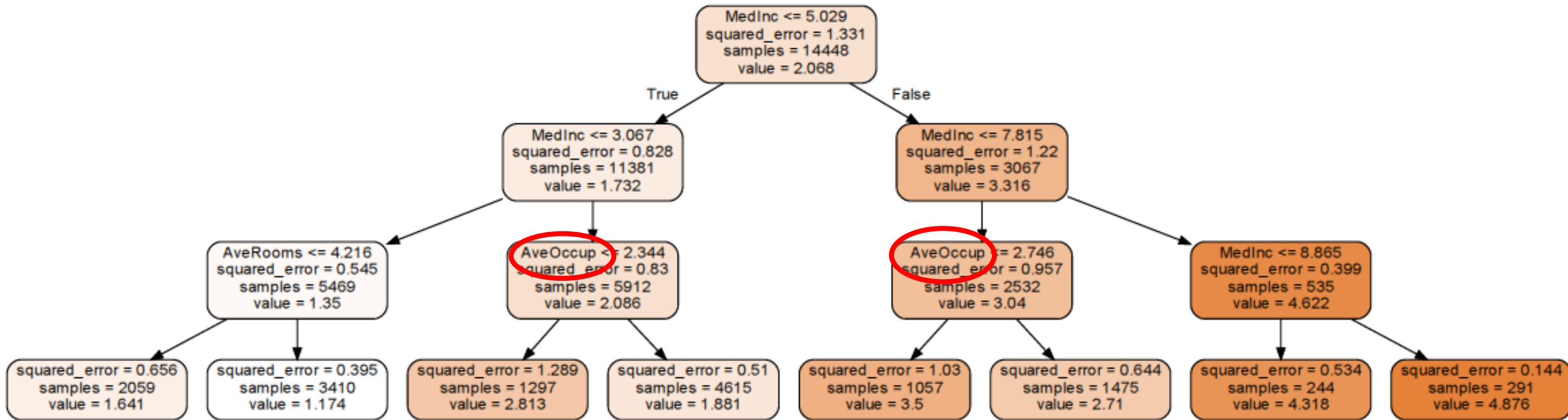# Feature importance

# Feature importance

# Feature importance



```
MedInc[node-1] = 1 * 1.331 – (11381/14448 * 0.828 + 3067/14448 * 1.22)

MedInc[node-2] = 1 * 0.828 – (5469/11381 * 0.545 + 5912/11382 * 0.83)

…

MedInc = MedInc[node-1] + MedInc[node-2] … + MedInc[node-n]
```
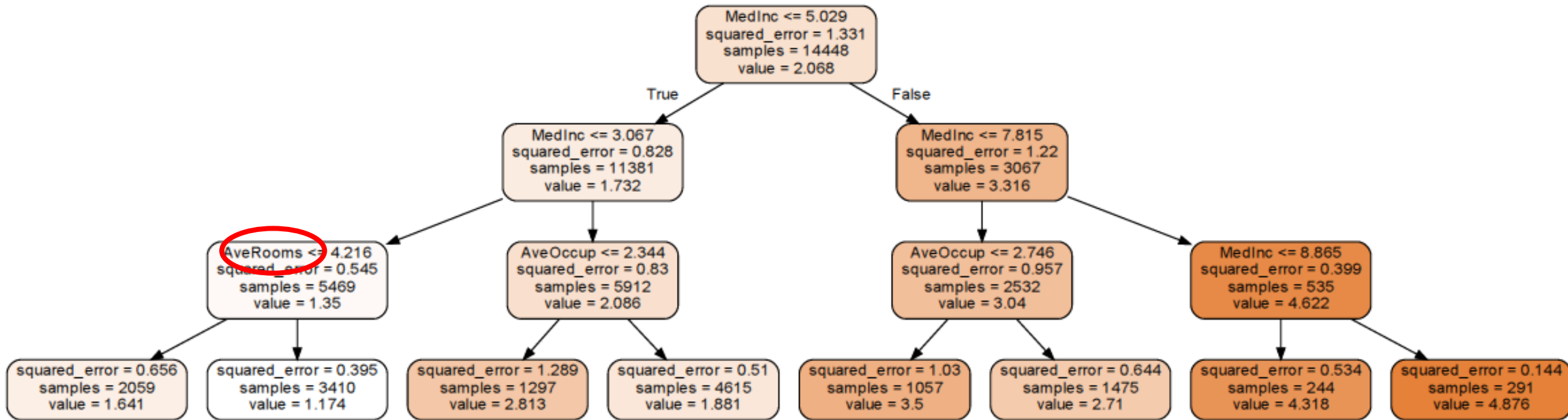
# Feature importance



$$AveOccup = AveOccup[node-1] + AveOccup[node-2]$$

# Feature importance
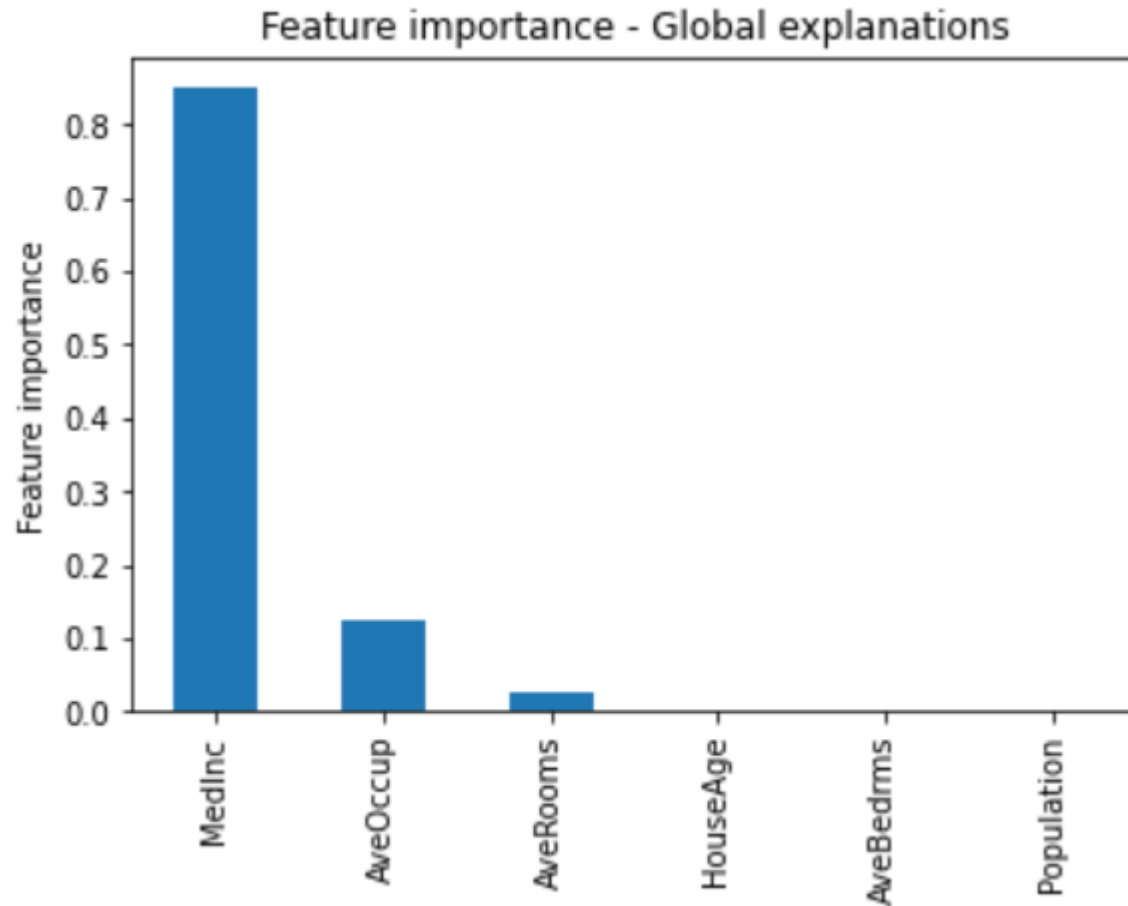


AveRooms = AveRooms[node-1]

# Feature importance

Sum(importance) = Importance[MedInc] + Importance[AveOccup] + Importance[AveRooms]

Importance[MedInc] = Importance[MedInc] / Sum(importance)

Importance[AveOccup] = Importance[AveOccup] / Sum(importance)

Importance[AveRooms] = Importance[AveRooms] / Sum(importance)

# Feature importance


Feature importance - Global explanations

- Features at top nodes have generally greater importance.
  - higher decrease in impurity.

- Features used in multiple nodes are more important.
  - Higher cumulative impurity decrease.