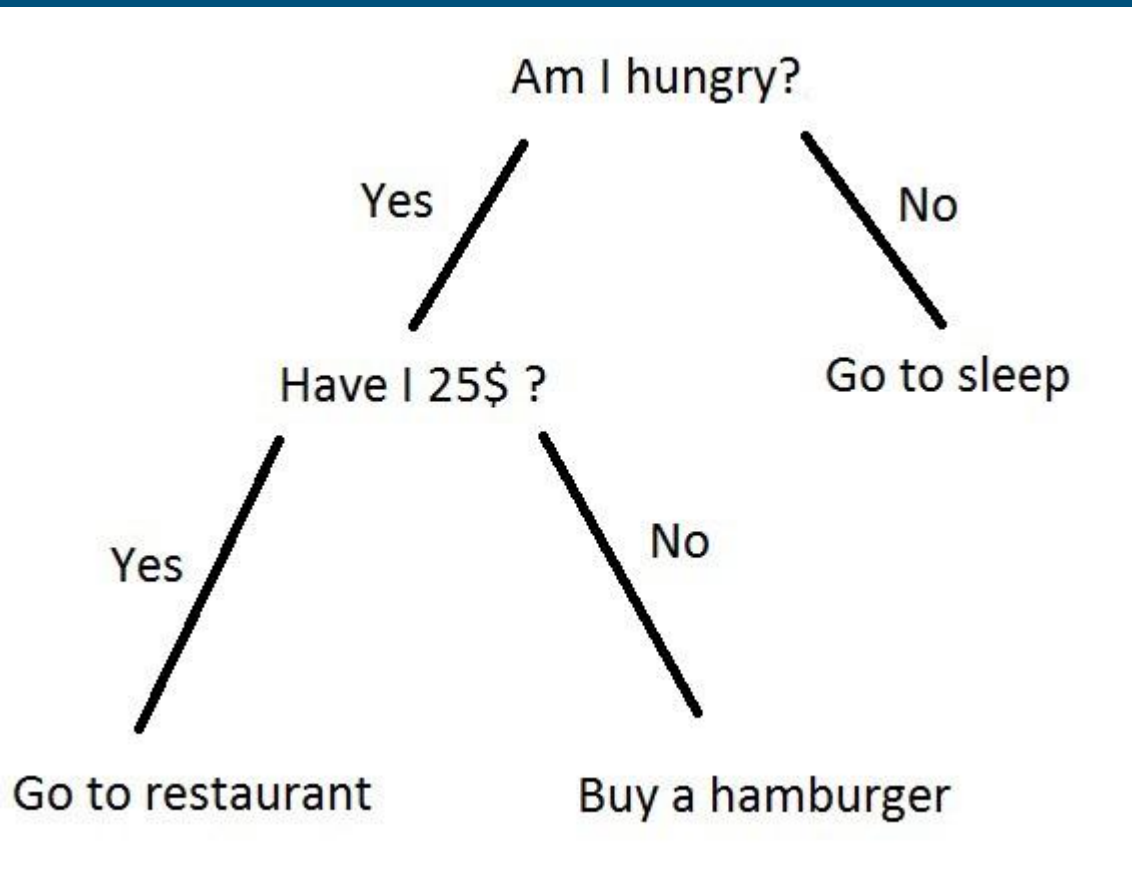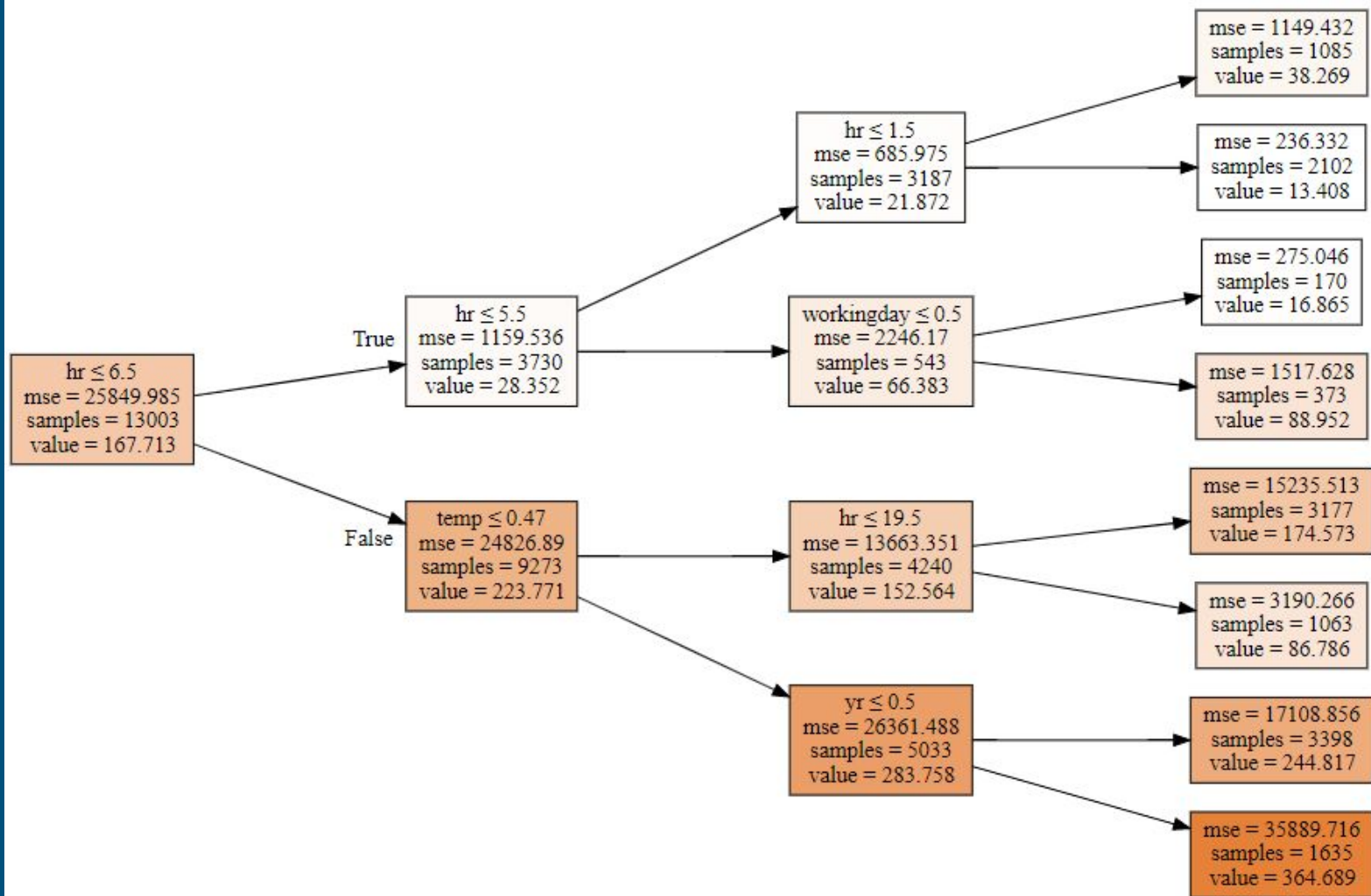# Data Science Course

Lecture 15

# Decision Trees

# Decision Trees

- Classification and Regression Trees or CART.
- Nodes, edges, leafs.
- Sequence of splits.

# Training

- Optimal split: iterate over the predictors and their possible values.
- Target average of resulting sub-groups.
- Calculate MSE.
- Weighted average.

# Evaluation - R²

$$SS_{\text{res}} = \sum_i (y_i - f_i)^2$$

$$SS_{\text{tot}} = \sum_i (y_i - \bar{y})^2,$$

$$R^2 \equiv 1 - \frac{SS_{\text{res}}}{SS_{\text{tot}}}.$$

# Balance

Pros:

- Simple to understand, interpret, visualize.
- Can handle numerical and categorical data.
- Regression and classification.

Cons:

- Over-complex trees that do not generalize the data well.
- Biased trees if some classes dominate.

# Bagging

- Objective: reduce the variance (robustness and accuracy).
- Multiple models not correlated.

# Random Forest

- Small tweak to decorrelate the trees.