# Machine Learning and AI in Stata

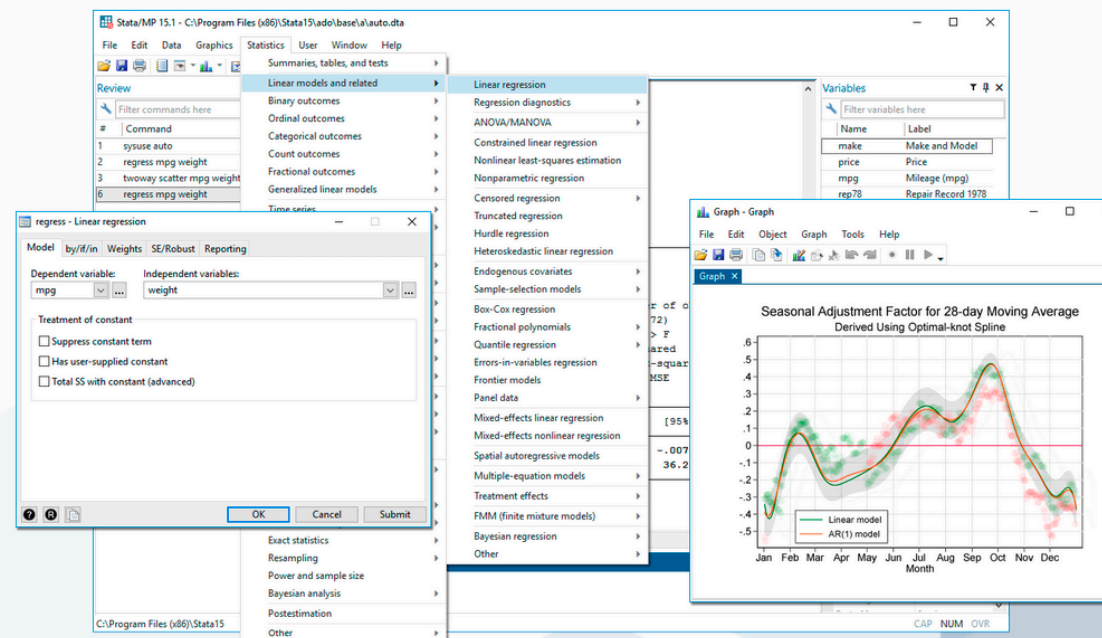**Kiel.AI**

# What is stata?

- Statistical Program like SPSS or Matlab

# Where to get it

- ## [stata.com](stata.com)

- ## Not for free

## New purchases
## Business single-user

release 15

ALL PRICES IN USD

| | Stata/IC | Stata/SE | Stata/MP 2-core | Stata/MP 4-core | Stata/MP >4 cores |
|---|---|---|---|---|---|
| | For mid-sized datasets. | For large datasets. | Fast & for the largest datasets. | Faster. | Even faster. |
| | Perpetual | Perpetual | Perpetual | Perpetual | Select cores |
| | $1,195 USD/perpetual Buy | $1,695 USD/perpetual Buy | $1,995 USD/perpetual Buy | $2,295 USD/perpetual Buy | |

| Product features | Stata/IC | Stata/SE | Stata/MP | | |
|---|---|---|---|---|---|
| Maximum number of variables | 2,048 | 32,767 | 120,000 | | |
| Maximum number of observations | 2.14 billion | 2.14 billion | Up to 20 billion | | |
| Maximum number of independent variables | 798 | 10,998 | 10,998 | | |
| Multicore support | 1-core | 1-core | 2-core | 4-core | 4+ |

## Student pricing

Students currently enrolled at degree-granting institutions may purchase Stata at the prices listed below. Proof of student status (i.e., copy of your university ID card) is required.

**Looking for Small Stata?**

release 15

ALL PRICES IN USD

| | Stata/IC | Stata/SE | Stata/MP 2-core | Stata/MP 4-core |
|---|---|---|---|---|
| | For mid-sized datasets. | For large datasets. | Fast & for the largest datasets. | Faster. |
| Perpetual: | $198 USD Buy | $395 USD Buy | $695 USD Buy | $995 USD Buy |
| Annual: | $89 USD Buy | $235 USD Buy | $395 USD Buy | $545 USD Buy |
| 6 months: | $45 USD Buy | $125 USD Buy | | |

| Product features | Stata/IC | Stata/SE | Stata/MP |
|---|---|---|---|
| Maximum number of variables | 2,048 | 32,767 | 120,000 |
| Maximum number of observations | 2.14 billion | 2.14 billion | Up to 20 billion |
| Maximum number of independent variables | 798 | 10,998 | 10,998 |

# Interface



Past commands appear here

Results are displayed here

Variable list appears here

Data properties appear here

Current working directory appears here

Commands are typed here

Current log status appears here

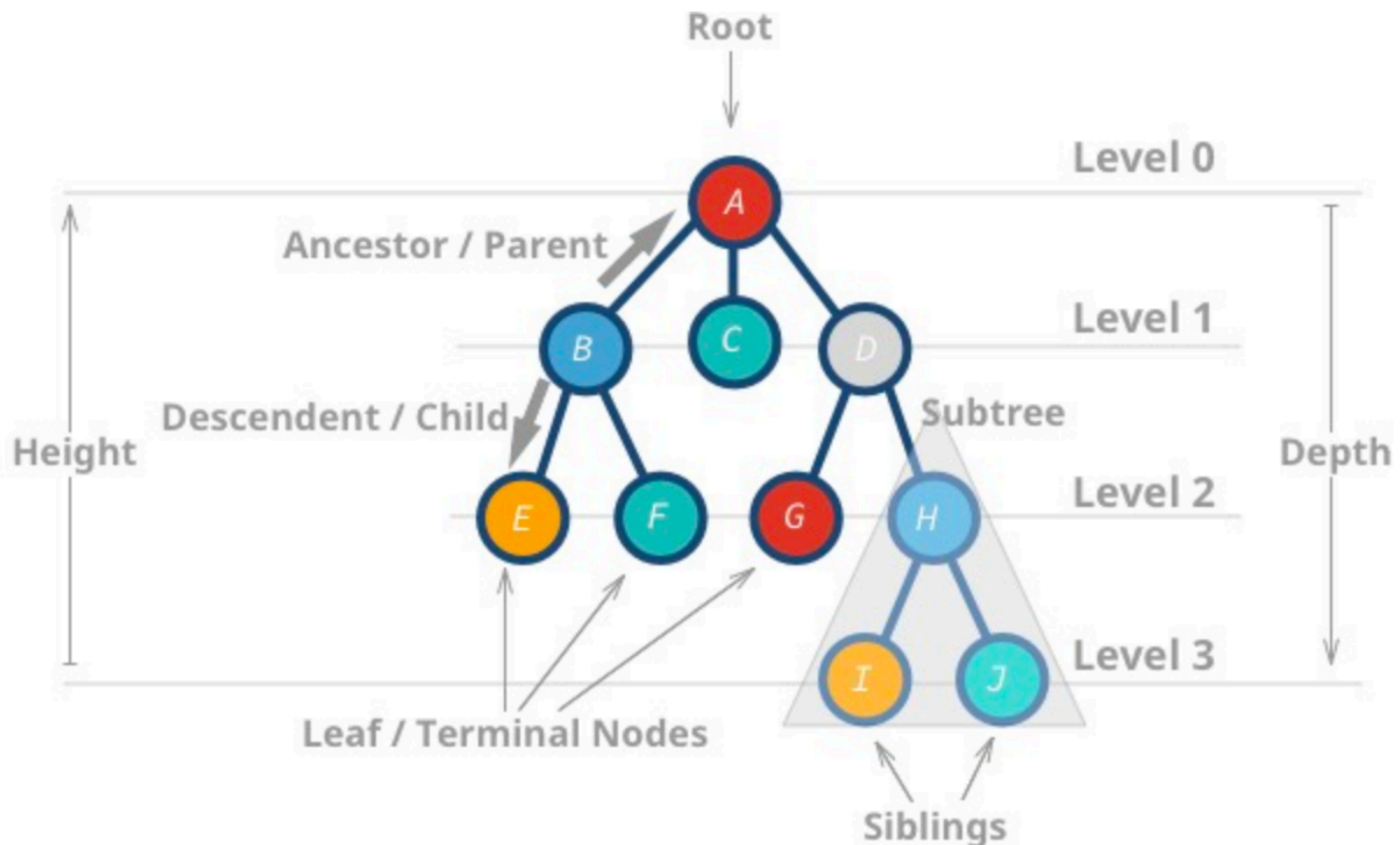Command log status appears here

# Interface

# Differences

- Hybrid between Syntax and Click based (Do-Files)

- Syntax based (R, Python, Matlab)

- Click based (SPSS)

# AI and ML in Stata

- i.e. module CART (Classification and regression trees):
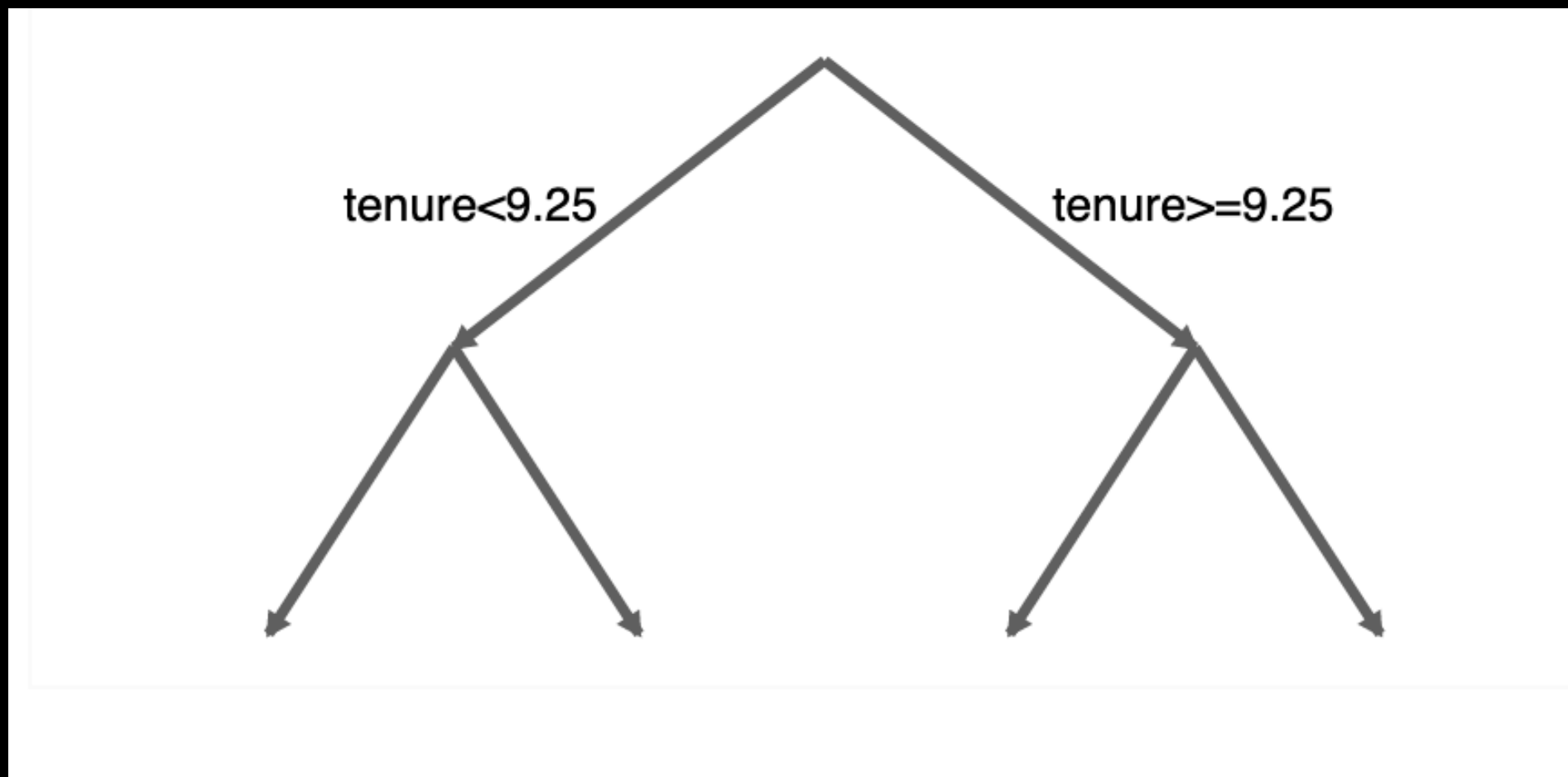
- Install by typing "ssc install cart"

# Trees in Stata

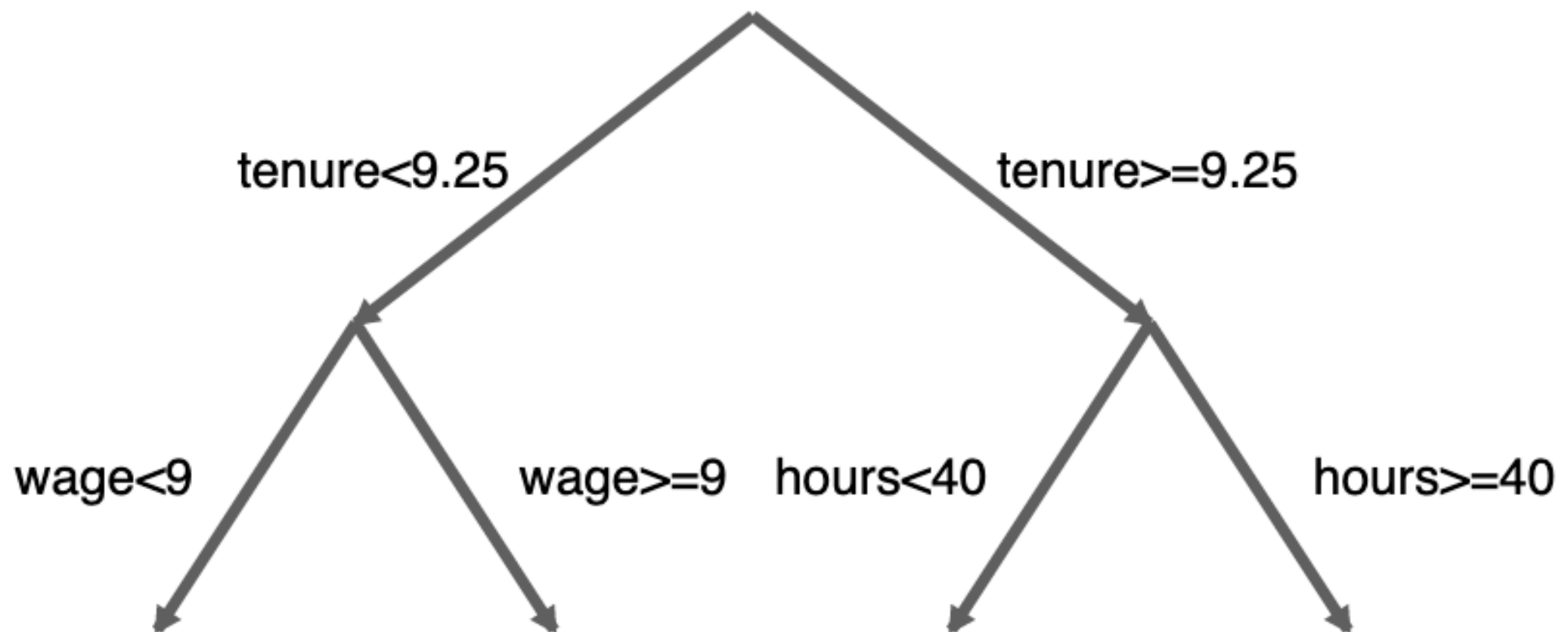- Root (top node), Nodes, Leafs (terminal nodes), …

# Trees in Stata

- Classification in two different groups at each node

# Trees in Stata

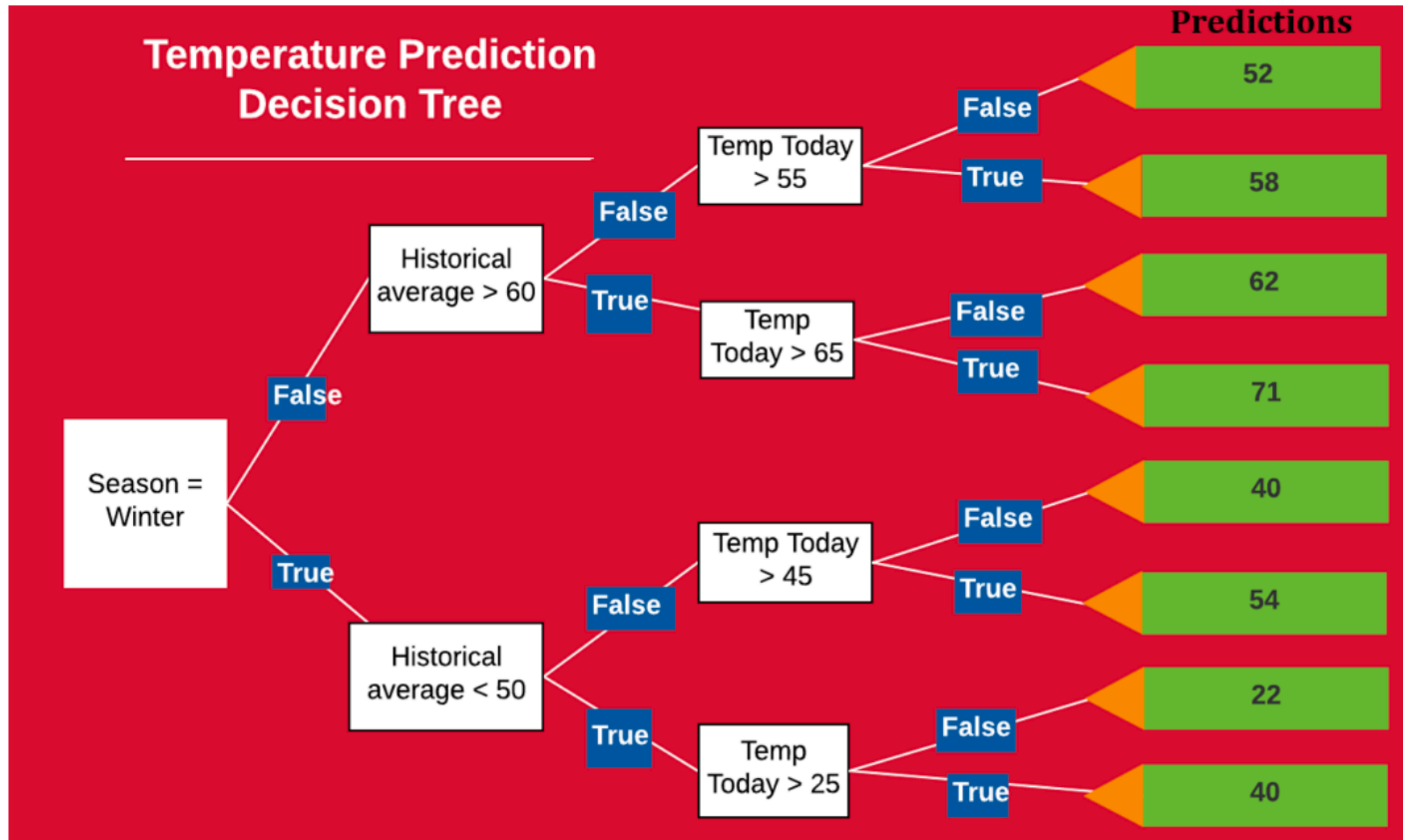- Gets complex very fast: 10 levels of binary splits give you 2^10 = 1024 terminal nodes (leaves)

# Trees in Stata

- Splits at the nodes are by default based on SSR, for binary outcomes it is the number of misclassified observations

# Random Forest

- Collection of Decision Trees (default: 500) with a conclusion or final result.

- Random Forest has a very small MSE (Mean Squared Error) in general

- ML has to be able to react well to new data

- 1. Random sampling of training data

- 2. Random subset of features when splitting nodes

# Temperature prediction

# Conclusion

- Hybrid between click and syntax (intermediate level)

- R or Python are better for more advanced tasks

- But easy tasks like subsets are difficult or not possible