# Using Data Mining Tools of Decision Trees in Quality and Reliability Applications
## *Brief Example on Modern Engineered Wood*

Hyunjoong (June) Kim, Frank Guess and Timothy Young
Yonsei University, Seoul, Korea
University of Tennessee, Knoxville

2006 Joint Research Conference on Statistics
in Quality, Industry and Technology
Knoxville, TN

Paper link: http://stat.bus.utk.edu/techrpts/index.html
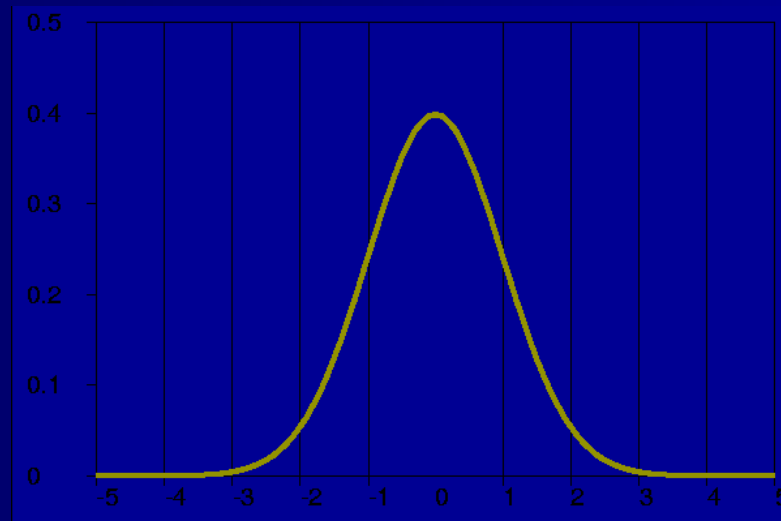
# Data Mining (DM)

## An Important Tool for Industrial Processes and Real-Time Decision-Making

# Data Mining (DM)
# Philosophy of Reducing Variation

- If you cannot quantify variation, how will you reduce variation?



- Once variation is quantified, what are sources of variation?

- DM and Decision Trees powerful methods for identifying sources of variation.

# Decision Trees (DT)

## "Real-time Supervised Classification Learning"

# Philosophy of Data Mining

- DM tools find hidden structures & helpful patterns
- DM provides exploratory data analysis with little human interactions
- Is that good or bad?
- Yes, that is good & it is bad! Both … ! Need balance of both automation & human interactions.
- Win/win - old time Deming strategy and modern process of thinking deeper… .
- DT is aligned well with the DM philosophy

# Data Mining (DM) & Trees

□ Decision trees (classification and regression) are a popular method among DM tools

□ Quick ref: Guidici (2003) *Applied Data Mining: Stat Methods for Bus. & Indus.*

□ We will discuss results of a case study using "GUIDE", one of many DT methods.

# GUIDE

(Generalized, Unbiased, Interaction Detection and Estimation)
Loh (2002)

http://www.stat.wisc.edu/~loh/guide.html

Fit one regression model at each node – multiple regression, stepwise, etc.

1. Use residuals to select split variable (negligible bias)
2. Select split point or split set
3. Prune tree as in CART

See also Kim & Loh (2003) take a "CRUISE" with another DT tool

http://web.utk.edu/~hjkim/  & 2004 work …

# GUIDE
## (Key features: sensitivity to curvature and interactions)

# Case Study
# Medium Density Fiberboard (MDF)

# Case Study
# MDF

- Medium density fiberboard (MDF): a highly used engineered wood composite
- Interested in the tensile strength (psi) or "internal bond" (IB) of MDF
- Destructive testing is performed during the manufacturing process (Goal: maximize & improve product quality and reliability)
- Prevention of unacceptable reliability can result in millions of dollars saving
- Young & Guess (2002) & Guess et al. (2003) - *IJRA*

# Predictors and Response (IB)

☐ Predictors: describe types and manufacturing conditions

- panel density (lbs/ft$^3$)

- panel thickness (inches)

- and others (moisture, line speed, temperature, etc.)

- day of the week (Sunday through Saturday)

- shift (morning, afternoon, night)

- week of the month (first through fourth week)

☐ Response variable is the strength of IB

☐ Next comparing regression only vs. GUIDE?

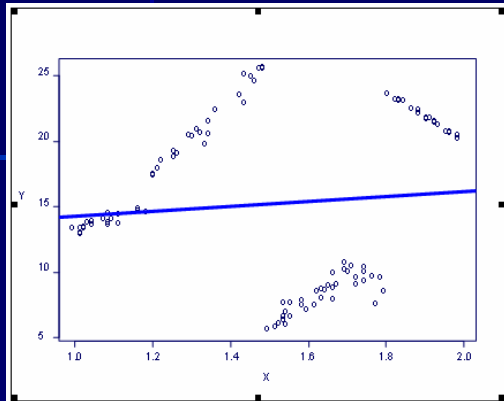# Automated Real-Time Relational Databases



Lab – Destructive Testing
Event Data

Data Warehouse
Real-Time Data

□  Mills collect a lot of data, how much knowledge is gained?

# "Regression Only" Analysis



| Source | | Coef | S.E. | Source | | Coef | S.E. |
|---|---|---|---|---|---|---|---|
| length | | 0.003 | 0.015 | | Sunday | 0.873 | 1.171 |
| density | | 21.371 | 0.476 | | Monday | 5.123 | 1.168 |
| thickness | | 70.701 | 4.317 | | Tuesday | 1.522 | 1.216 |
| width | | -0.063 | 0.106 | Day | Wednesday | -3.264 | 1.201 |
| week | 1 | 3.021 | 0.799 | | Thursday | -1.846 | 1.369 |
| | 2 | -2.841 | 0.797 | | Friday | -4.266 | 1.202 |
| | 3 | 0.736 | 0.735 | | Morning | -1.845 | 0.632 |
| | | | | Shift | Afternoon | 1.535 | 0.639 |

# "Regression Only" Analysis



| Source | | Coef | S.E. | Source | | Coef | S.E. |
|---|---|---|---|---|---|---|---|
| length | | 0.003 | 0.015 | | Sunday | 0.873 | 1.171 |
| density | | 21.371 | 0.476 | | Monday | 5.123 | 1.168 |
| thickness | | 70.701 | 4.317 | | Tuesday | 1.522 | 1.216 |
| width | | -0.063 | 0.106 | Day | Wednesday | -3.264 | 1.201 |
| | 1 | 3.021 | 0.799 | | Thursday | -1.846 | 1.369 |
| week | 2 | -2.841 | 0.797 | | Friday | -4.266 | 1.202 |
| | 3 | 0.736 | 0.735 | | Morning | -1.845 | 0.632 |
| | | | | Shift | Afternoon | 1.535 | 0.639 |

# GUIDE Regression Tree

# Why split, e.g., panel density?



**GUIDE** identifies location of split very easily and quickly

(not simply along pre-defined product type set points)

# Why split, e.g., panel thickness?
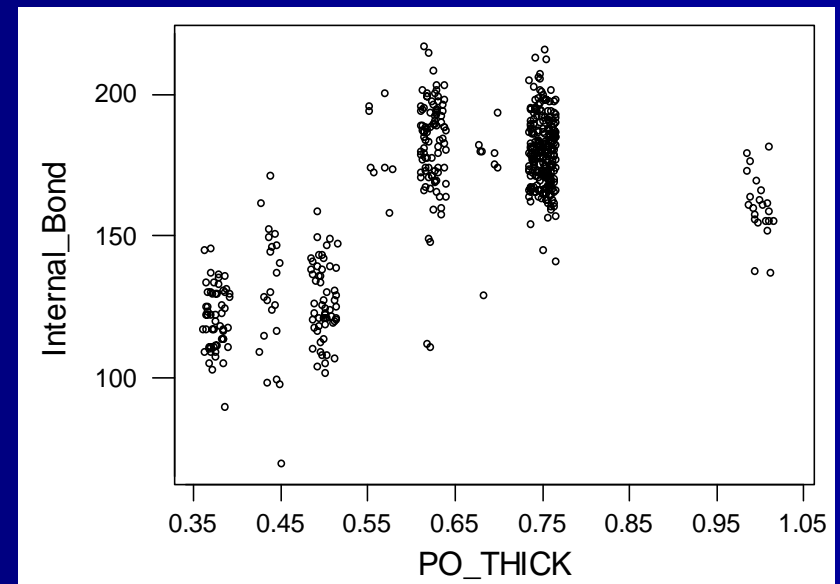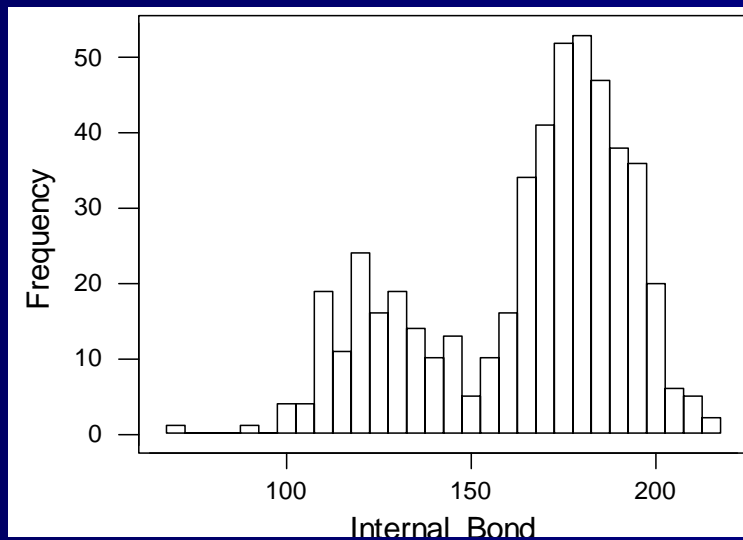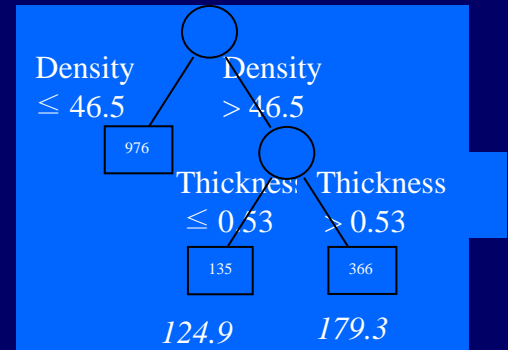


**GUIDE** identifies location of split very easily and quickly

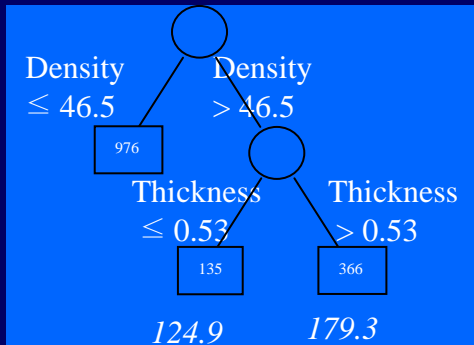(not simply along pre-defined product type set points)

# Density ≤ 46.5

Density ≤ 46.5    Density > 46.5

976

Thickness ≤ 0.53    Thickness > 0.53

135    366

*124.9    179.3*

| Source | Coef | S.E. | Source | | Coef | S.E. |
|---|---|---|---|---|---|---|
| length | 0.009 | 0.012 | | Sunday | 0.037 | 0.903 |
| density | 3.889 | 0.708 | | Monday | -1.770 | 0.971 |
| thickness | -11.063 | 3.765 | | Tuesday | -2.466 | 0.892 |
| width | 0.0362 | 0.087 | Day | Wednesday | 0.806 | 0.997 |
| week | 1 | 1.762 | 0.629 | | Thursday | 0.455 | 1.062 |
| | 2 | -1.516 | 0.625 | | Friday | 1.754 | 0.995 |
| | 3 | 0.286 | 0.569 | | Morning | -1.035 | 0.495 |
| | | | Shift | Afternoon | 0.922 | 0.497 |

Not desirable to have "Day", "Shift" as source of variation!

# Density > 46.5 and Thickness > 0.53



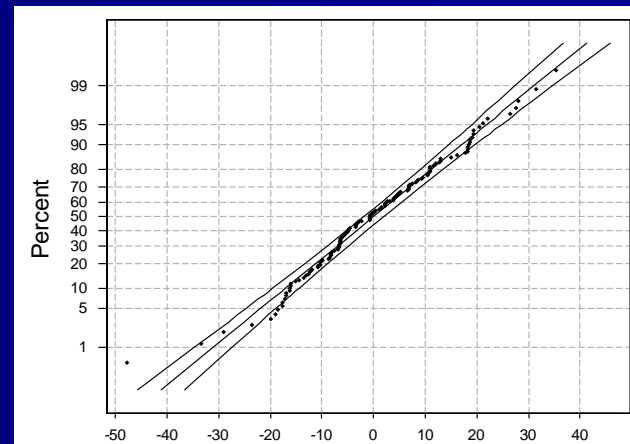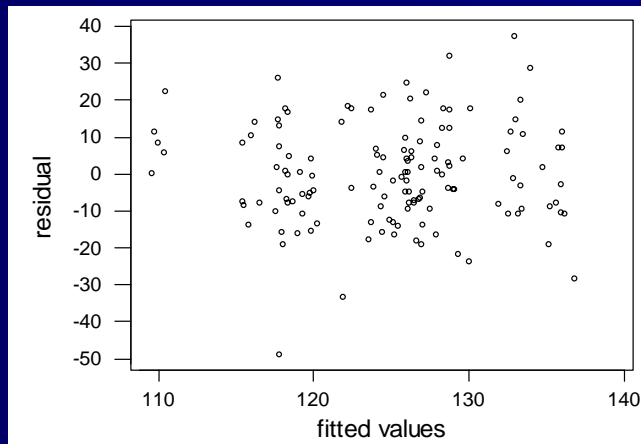| Source | | Coef | S.E. | | Source | Coef | S.E. |
|---|---|---|---|---|---|---|---|
| length | | -0.019 | 0.022 | | Sunday | 0.692 | 1.803 |
| density | | 22.561 | 4.834 | | Monday | -1.866 | 1.829 |
| thickness | | 1.320 | 14.81 | Day | Tuesday | 3.100 | 3.411 |
| width | | 0.089 | 0.181 | | Wednesday | 1.603 | 2.602 |
| week | 1 | 3.614 | 1.245 | | Thursday | 0.269 | 2.808 |
| | 2 | -4.913 | 1.270 | | Friday | -0.331 | 1.928 |
| | 3 | -0.936 | 1.153 | | Morning | -0.777 | 0.990 |
| | | | | Shift | Afternoon | -0.171 | 0.990 |

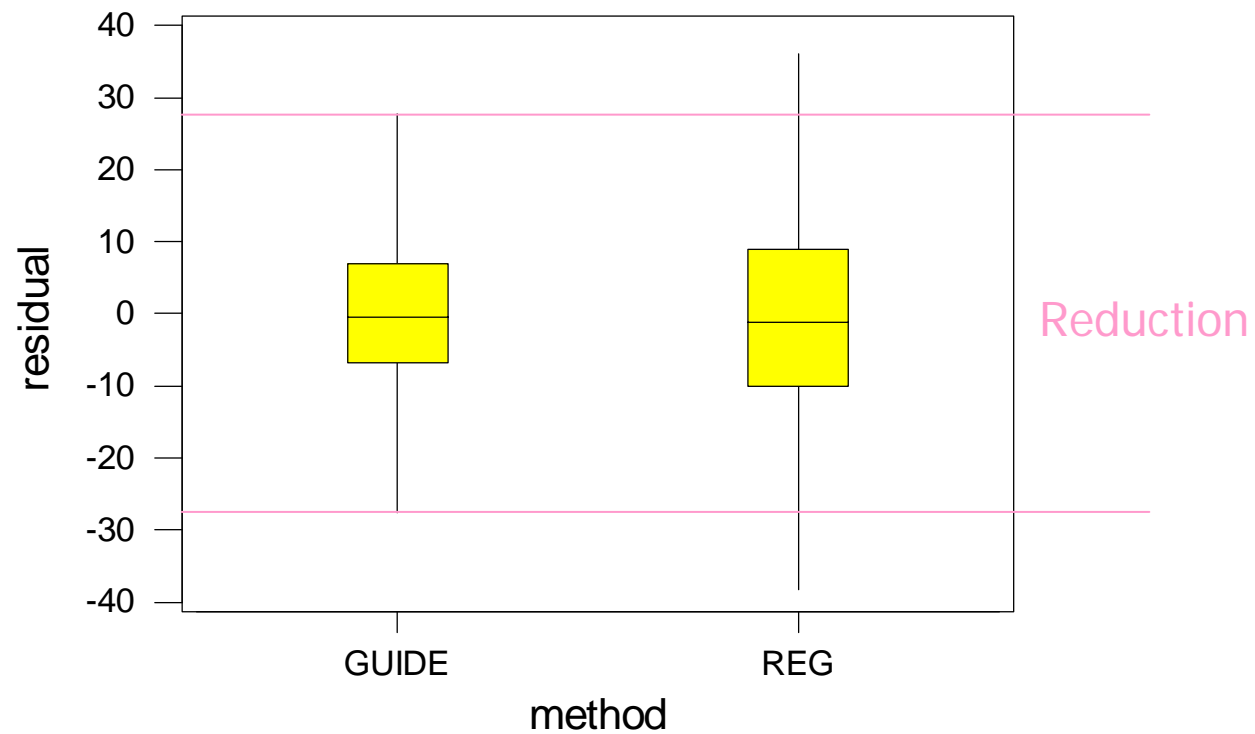Again, not desirable to have "Day" as source of variation!

# GUIDE vs. Regression Only

"Regression Only" Coefficient of Determination: $R^2 = 62.4\%$

"Guide" Coefficient of Determination; $R^2 = 83.0\%$
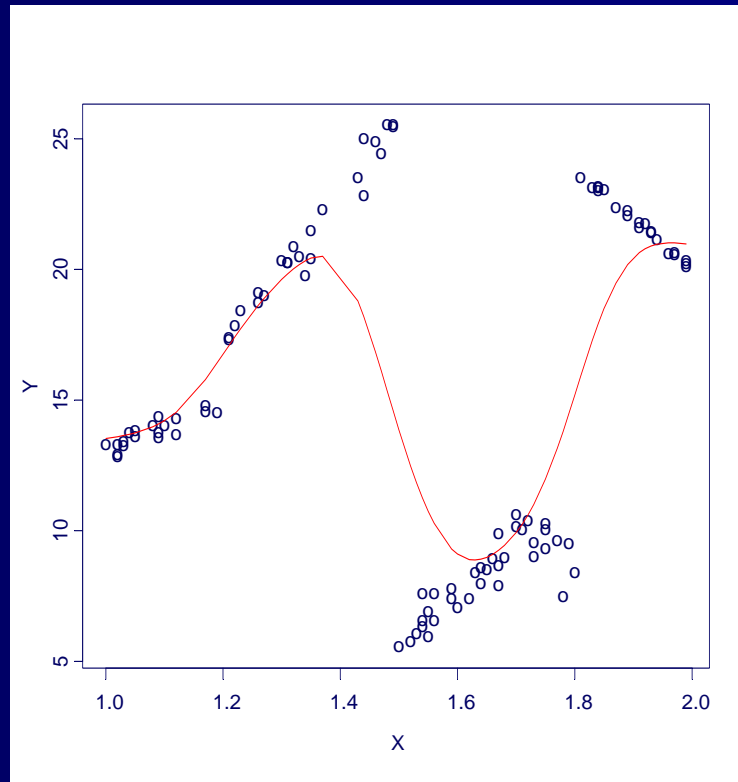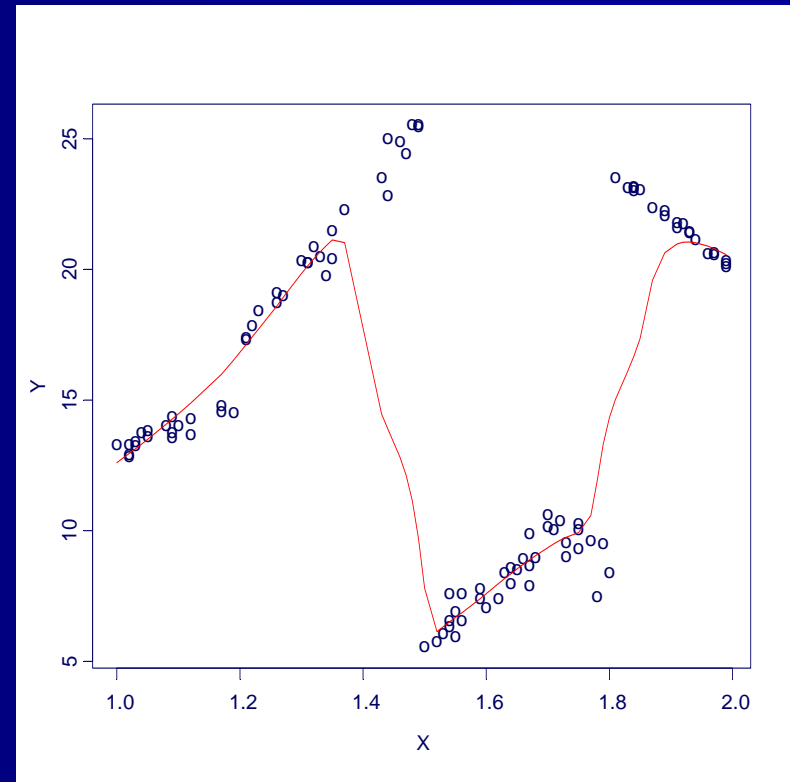
# GUIDE vs. Regression only

# Note

- GUIDE detected useful partitions
- Automatic procedure – little human interaction: wonderful exploratory tool!
- Easy Interpretation
    - follow the tree !
    - important variables appear in the tree
- Better fitted model & better predictions
- Missing data, simpler, interactions, etc.

# Decision Trees may be more helpful in analyzing industrial processes than "Global Modeling"
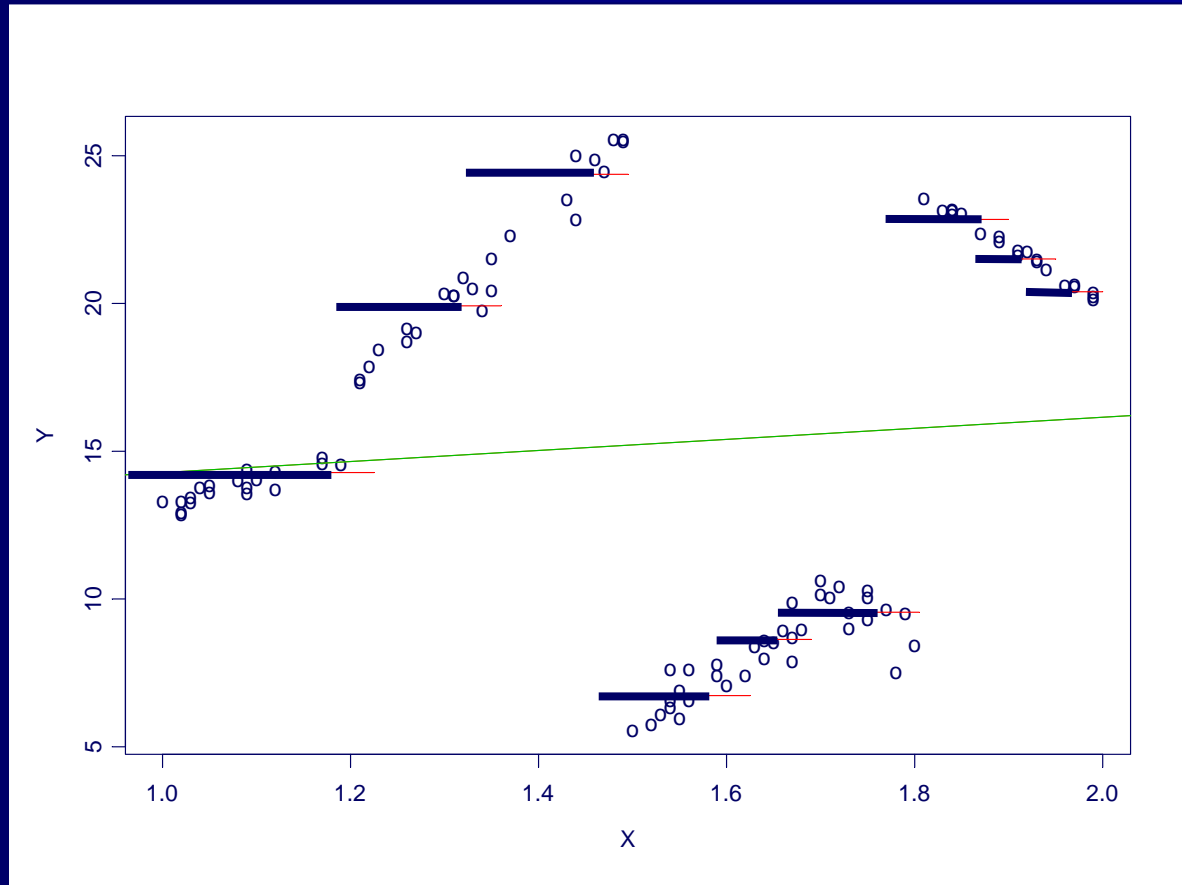
Smoothing Spline

Loess

# CART (Breiman *et al.* 1984) "Classification and Regression Tree"



Some Limitations of CART: mean function fit

# DT Modeling

- Assume different (multiple) correlation structure exists within database
- Tree-based model = data partitioning + statistical modeling + pruning
- Intermediate nodes = partition data space
- Terminal nodes = fit the Node model (example: linear regression, loess, mean)
- Pruning = prevent over-fitting
- "Given a Node model, how to find the partition (heterogeneity) of the data space?"
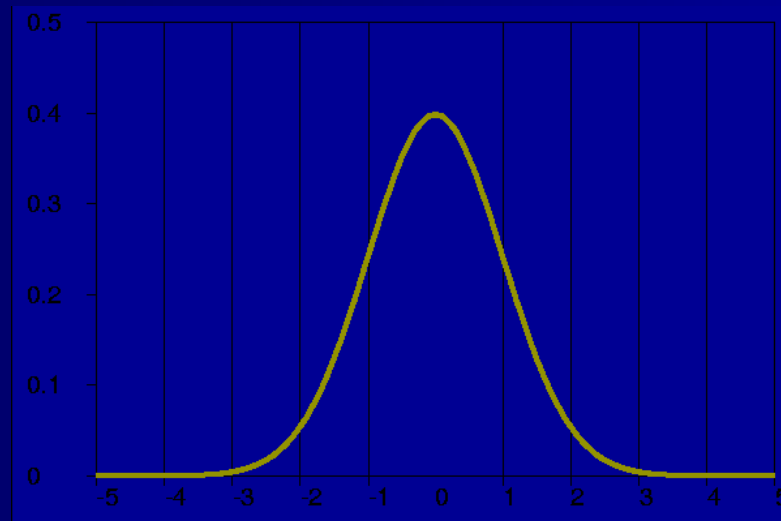
# Conclusions

- Decision Trees, e.g., GUIDE, etc., helpful
- Note there are data mining approaches other than decision trees, such as neural nets, genetic algorithms, etc. (dependent on problem and data heterogeneity)
- Decision Trees are attractive because they show clearly how to reach a decision, and easy to interpret by practitioners

# Decision Trees (DT)
# Philosophy of Reducing Variation

- If you cannot quantify variation, how will you reduce variation?



- Once variation is quantified, what are sources of variation?

- DT is a useful method for identifying sources of variation.

# Links & Future work

GUIDE:

http://www.stat.wisc.edu/~loh/guide.html

CRUISE:

http://web.utk.edu/~hjkim/

Paper:

http://stat.bus.utk.edu/techrpts/index.html