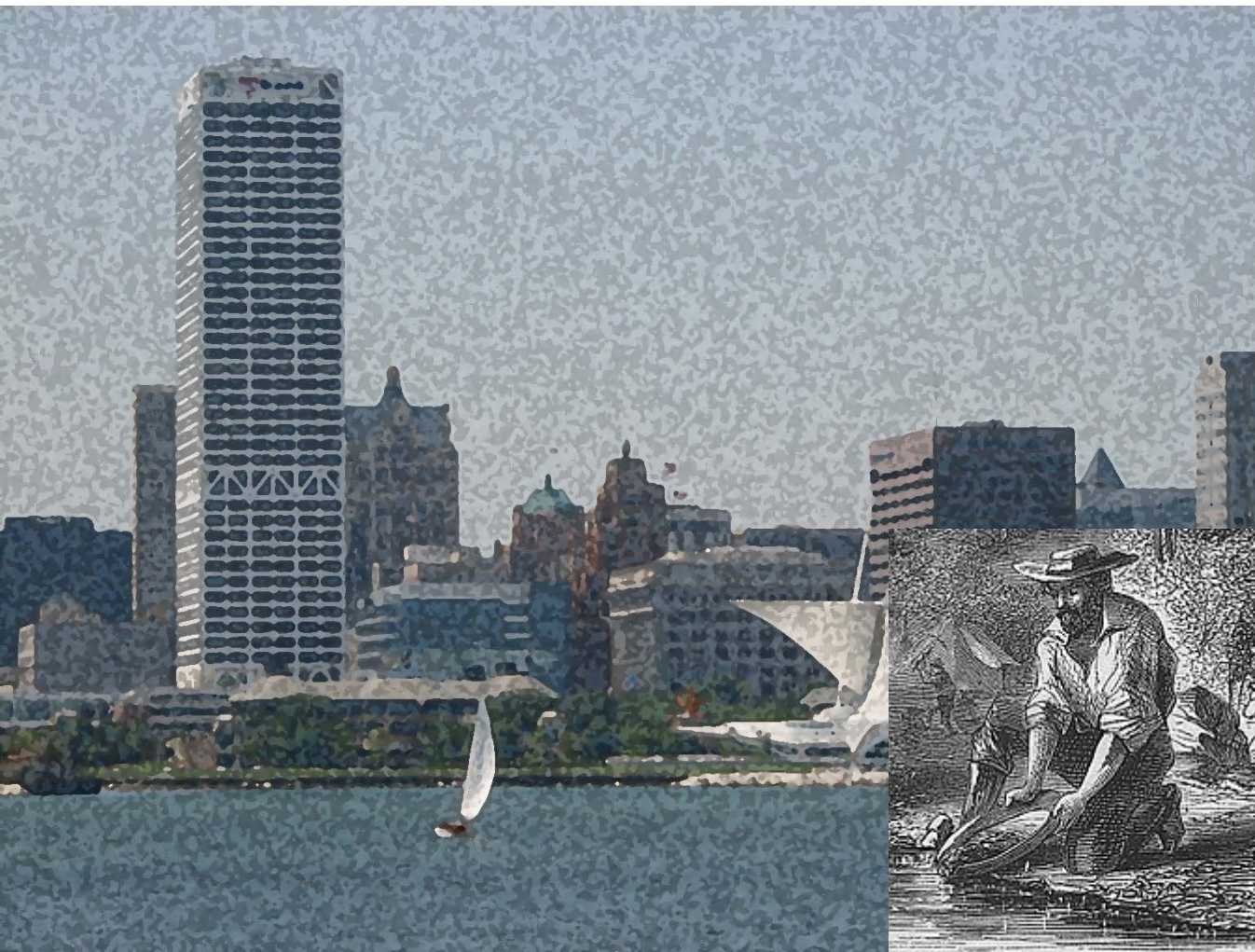# Finding the Gold in Your Data
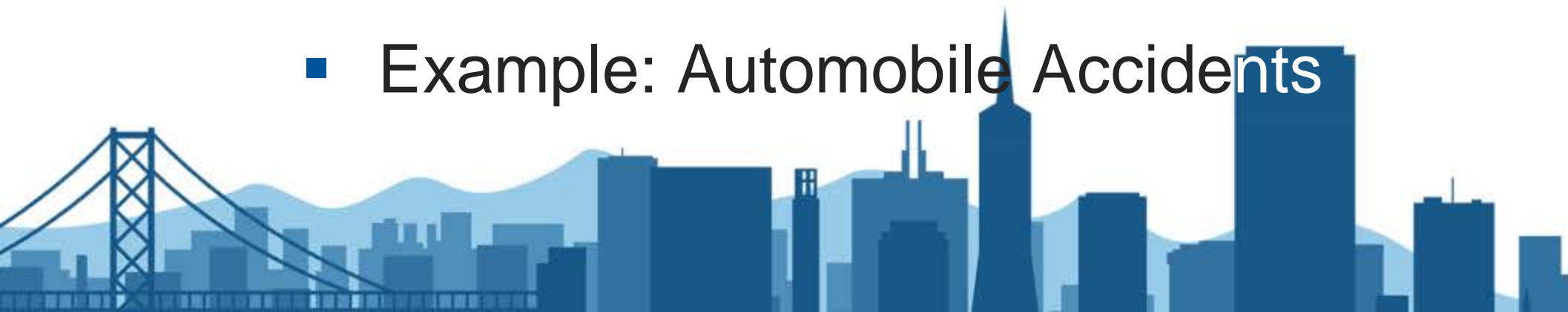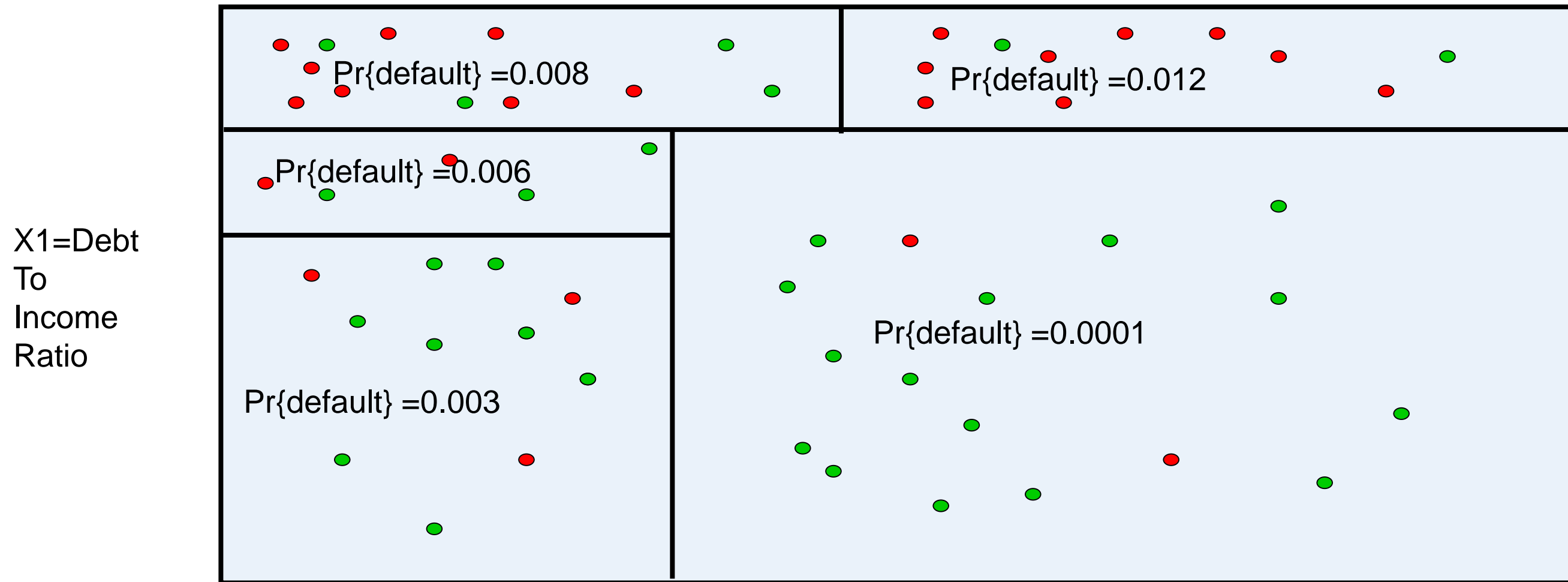
- An introduction to Data Mining
- Originally presented @ SAS Global Forum

# Decision Trees

- A "divisive" method (splits)

- Start with "root node" – all in one group

- Get splitting rules

- Response often binary

- Result is a "tree"

- Example: Loan Defaults

- Example: Framingham Heart Study
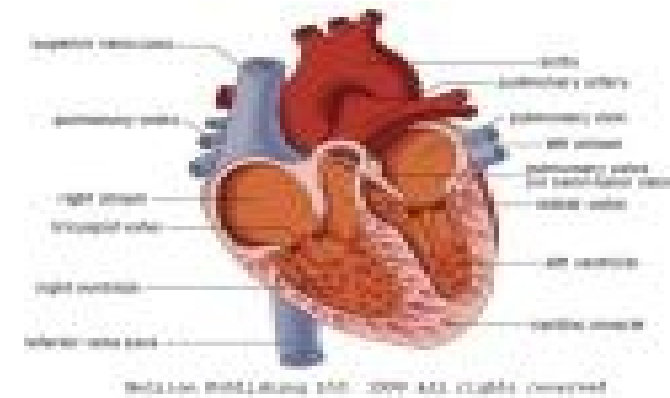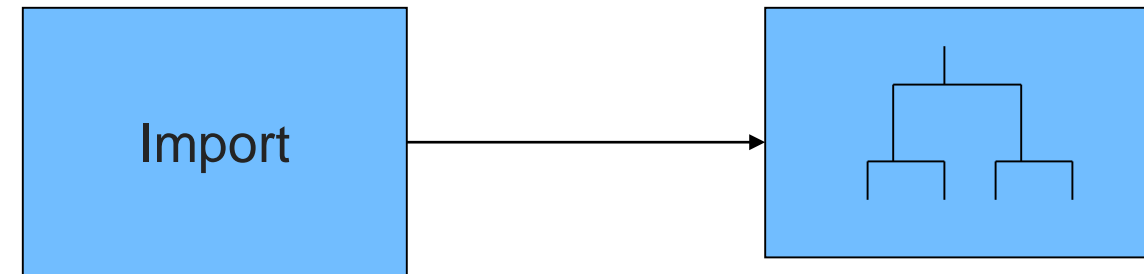
- Example: Automobile Accidents

# Recursive Splitting



X1=Debt To Income Ratio

Pr{default} =0.008

Pr{default} =0.012

Pr{default} =0.006

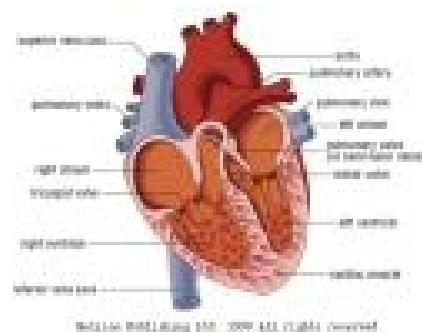Pr{default} =0.003

Pr{default} =0.0001

X2 = Age

- No default
- Default
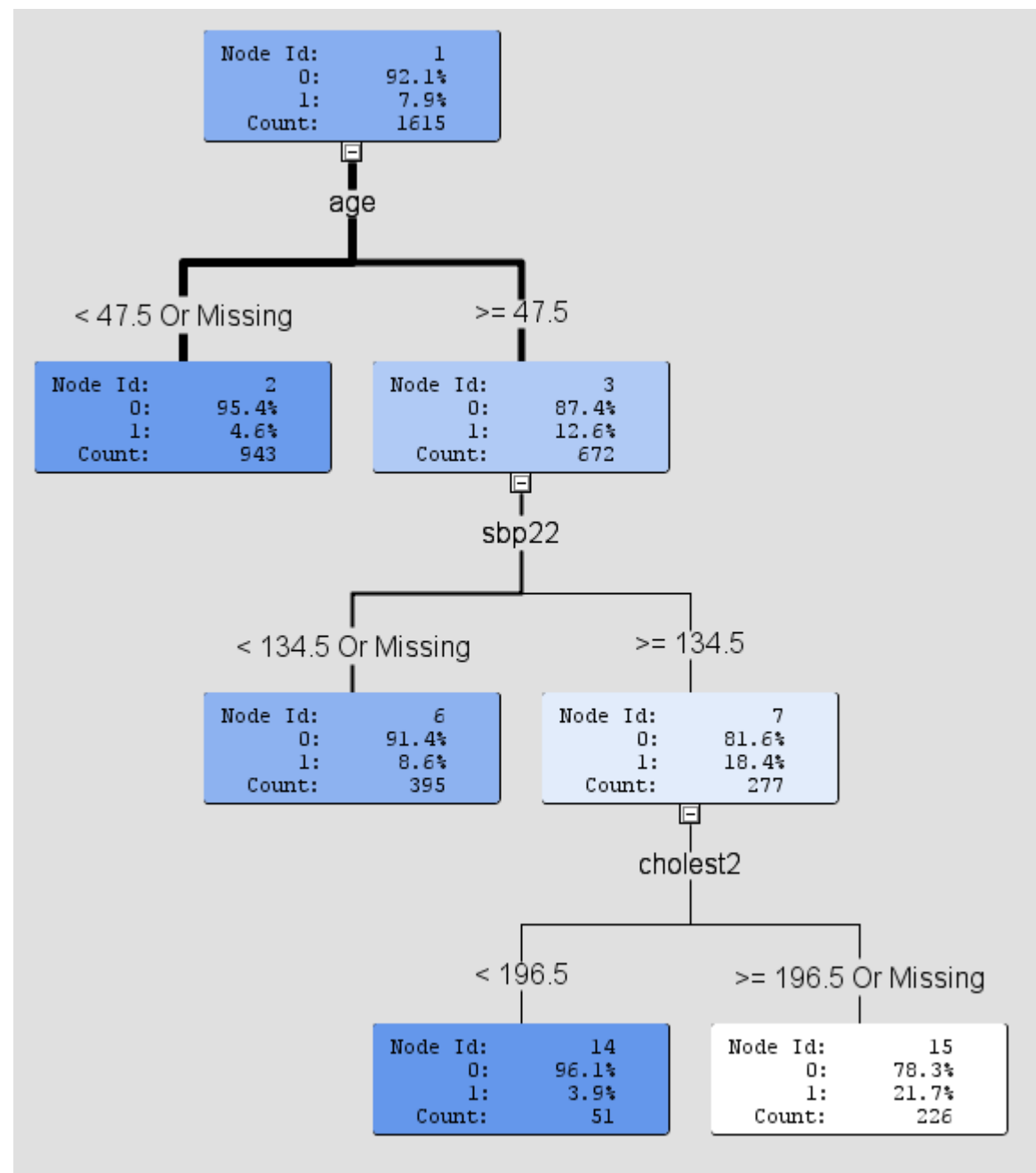
# Some Actual Data

- **Framingham Heart Study**

- **First Stage Coronary Heart Disease**

  - P{CHD} = Function of:

    » Age  - no drug yet!  ☹

    » Cholesterol

    » Systolic BP

Example of a "tree" →

# How to make splits?

- Contingency tables

### DEPENDENT (effect)

| | Heart Disease No | Heart Disease Yes | |
|---|---|---|---|
| **Low BP** | 95 | 5 | 100 |
| **High BP** | 55 | 45 | 100 |
| | 150 | 50 | |

180 ?
240?

### INDEPENDENT (no effect)

| | Heart Disease No | Heart Disease Yes | |
|---|---|---|---|
| | 75 | 25 | 100 |
| | 75 | 25 | 100 |
| | 150 | 50 | |

# How to make splits?

- Contingency tables



Heart Disease

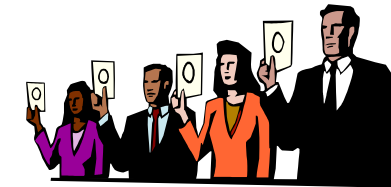|  | No | Yes |  |
|---|---|---|---|
| Low BP | 95 | 5 | 100 |
|  | 75 | 25 |  |
| High BP | 55 | 45 | 100 |
|  | 75 | 25 |  |
|  | 150 | 50 |  |

180 ?
240?

DEPENDENT (effect)

$$\chi^2 = \sum_{allcells} \frac{(\ \textbf{Observed - Expected}\ )^2}{\textbf{Expected}} =$$

$$2(400/75) + 2(400/25) = 42.67$$

Compare to tables – Significant!

(Why "Significant" ???)

SAS.GLOBALFORUM

$H_0$: 💙
$H_1$: ✖

$H_0$: Innocence
$H_1$: Guilt

Beyond reasonable
doubt
$P < 0.05$

| | |
|---|---|
| 95 | 5 |
| 75 | 25 |
| 55 | 45 |
| 75 | 25 |

$H_0$: No association
$H_{1:}$ BP and heart disease
_are_ associated

$P = 0.00000000064$

Framingham Conclusion:  Sufficient evidence _against_ the (null) hypothesis of no relationship.
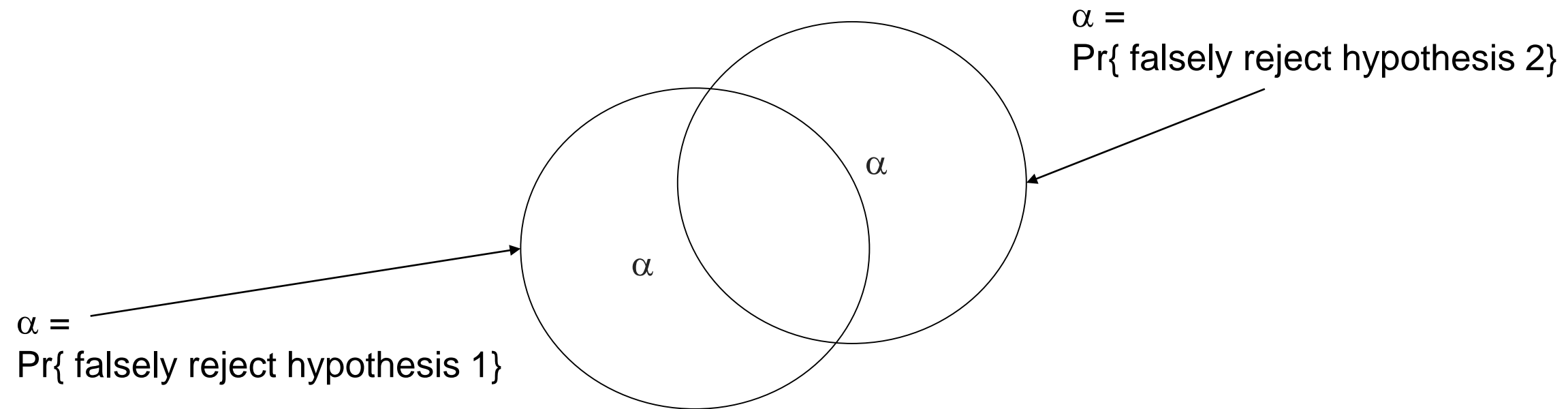
# How to make splits?

- Which variable to use?

- Where to split?
  - Cholesterol > ____
  - Systolic BP > _____

- Idea – Pick BP cutoff  to minimize p-value  for $\chi^2$

- Split point data-derived!

- What does "significance" mean now?

# Multiple testing

$\alpha =$
Pr{ falsely reject hypothesis 2}

$\alpha$

$\alpha$

$\alpha =$
Pr{ falsely reject hypothesis 1}

Pr{ falsely reject one or the other} < $2\alpha$
Desired:  0.05 probability or less
Solution: Compare 2(p-value) to 0.05

SAS.**GLOBAL**FORUM

# Other Sp    lit Criteria

- **Gini Diversity Index**
    - (1)      { A A A A B A B B C B}
    - Pick 2, Pr{different} = 1-Pr{AA}-Pr{BB}-Pr{CC}
        - » 1-[10+6+0]/45=29/45=0.64
    - (2)      { A A B C B A A B C C }
        - » 1-[6+3+3]/45 = 33/45 = 0.73 → (2) IS MORE DIVERSE, LESS PURE

- **Shannon Entropy**
    - Larger → more diverse (less pure)
    - $-\Sigma_i \, p_i \, \log_2(p_i)$

{0.5, 0.4, 0.1} → 1.36
{0.4, 0.2, 0.3} → 1.51    (more diverse)
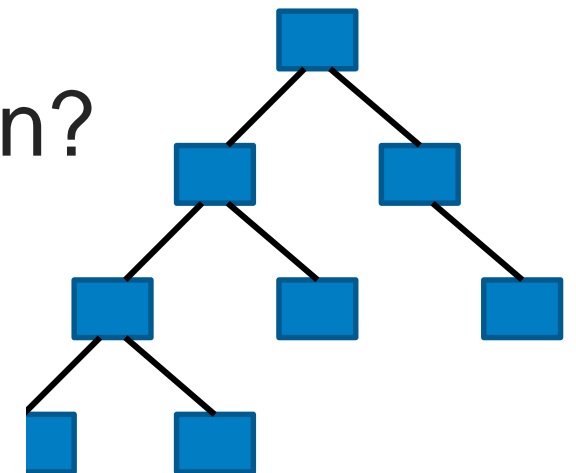
SAS. GLOBAL FORUM

# Validation

- Traditional stats – small dataset, need all observations to estimate parameters of interest.

- Data mining – loads of data, can afford "holdout sample"

- Variation: n-fold cross validation
  - Randomly divide data into n sets
  - Estimate on n-1, validate on 1
  - Repeat n times, using each set as holdout.

# Pruning

- Grow bushy tree on the "fit data"

- Classify validation (holdout) data

- Likely farthest out branches do not improve, possibly hurt fit on validation data

- Prune non-helpful branches.

- What is "helpful"?  What is good discriminator criterion?

# Goals

- <u>Split</u> if diversity in parent "node" > summed diversities in child nodes

- <u>Prune</u> to optimize
  - Estimates
  - Decisions
  - Ranking

- in validation data

# Accounting for Costs

- Pardon me (sir, ma'am) can you spare some change?

- Say "sir" to male +$2.00

- Say "ma'am" to female +$5.00

- Say "sir" to female -$1.00 (balm for slapped face)

- Say "ma'am" to male -$10.00 (nose splint)

# Including Probabilities

Leaf has Pr(**M**)=**.7**, Pr(**F**)=**.3**        You say:

|  | Sir | Ma'am |
|---|---|---|
| True Gender | | |
| **M** | **0.7** (2) | **0.7** (-10) |
| **F** | **0.3** (-1) | **0.3** (5) |
|  | +$1.10 | −$5.50 |

Expected profit is
2(0.7)-1(0.3) = $1.10
if **I** say "sir"

Expected profit is
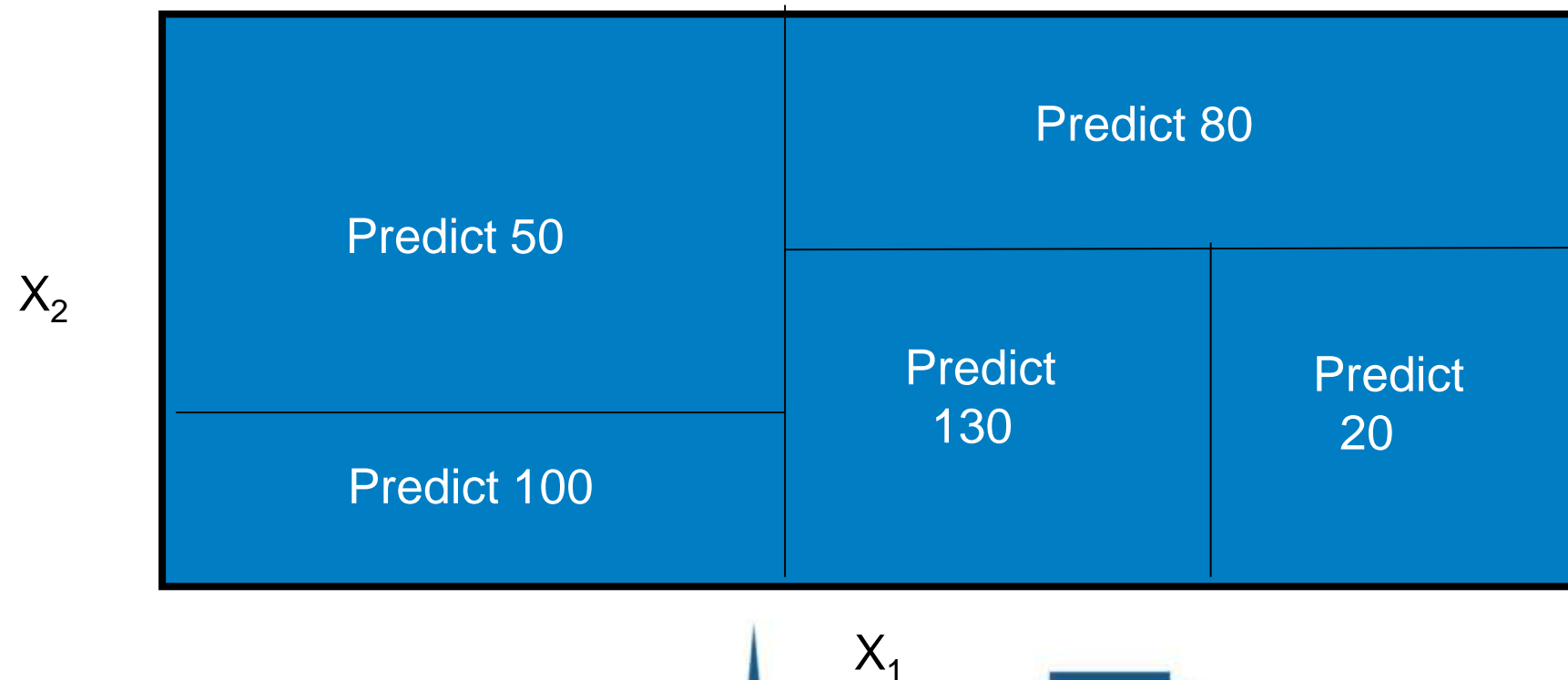-7+1.5 = -$5.50 (a loss)
if **I** say "Ma'am"

Weight leaf profits by leaf size (# obsns.) and sum.
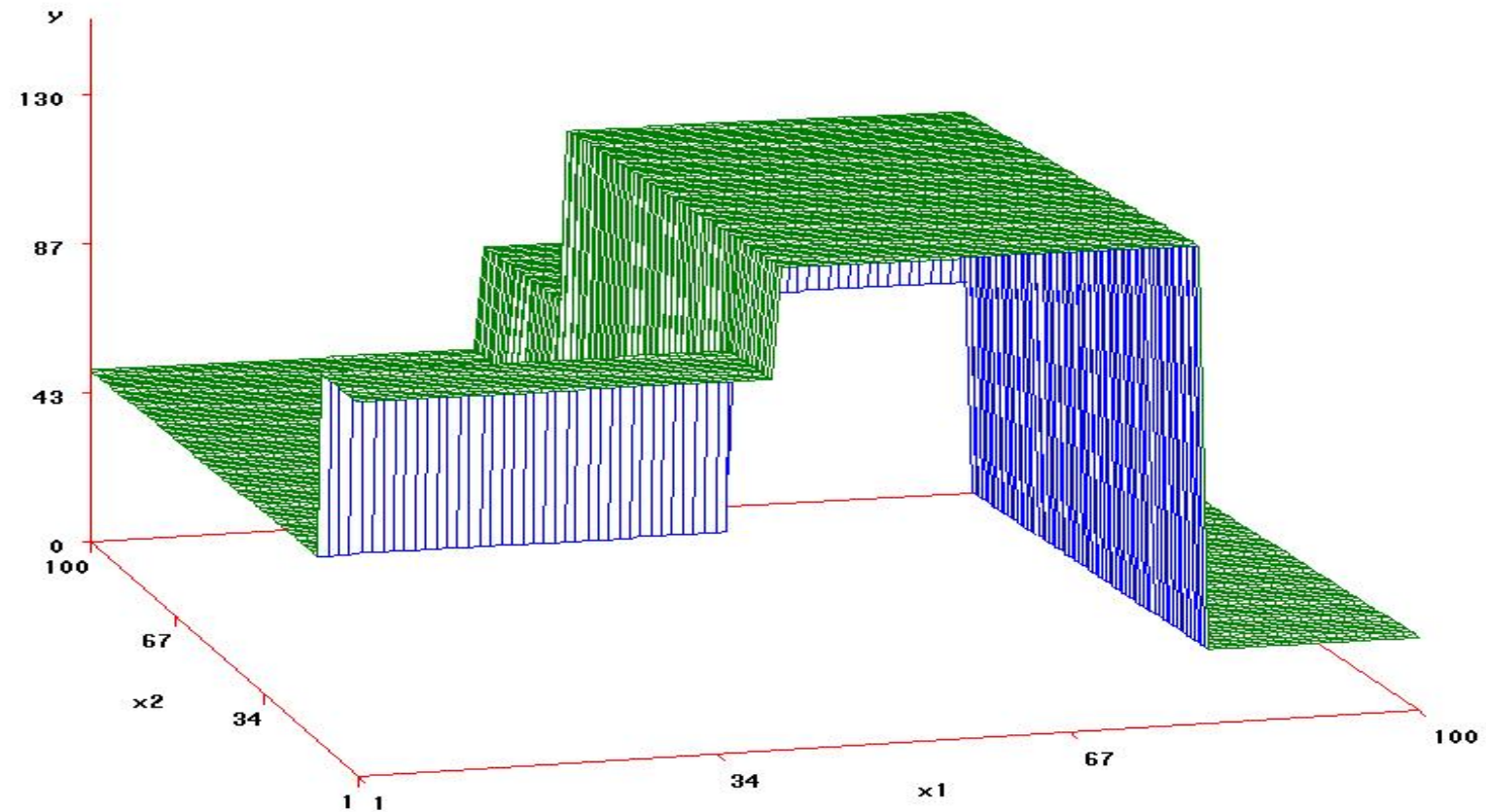
Prune (and split) to maximize profits.

# Regression Trees

- *Continuous* response Y
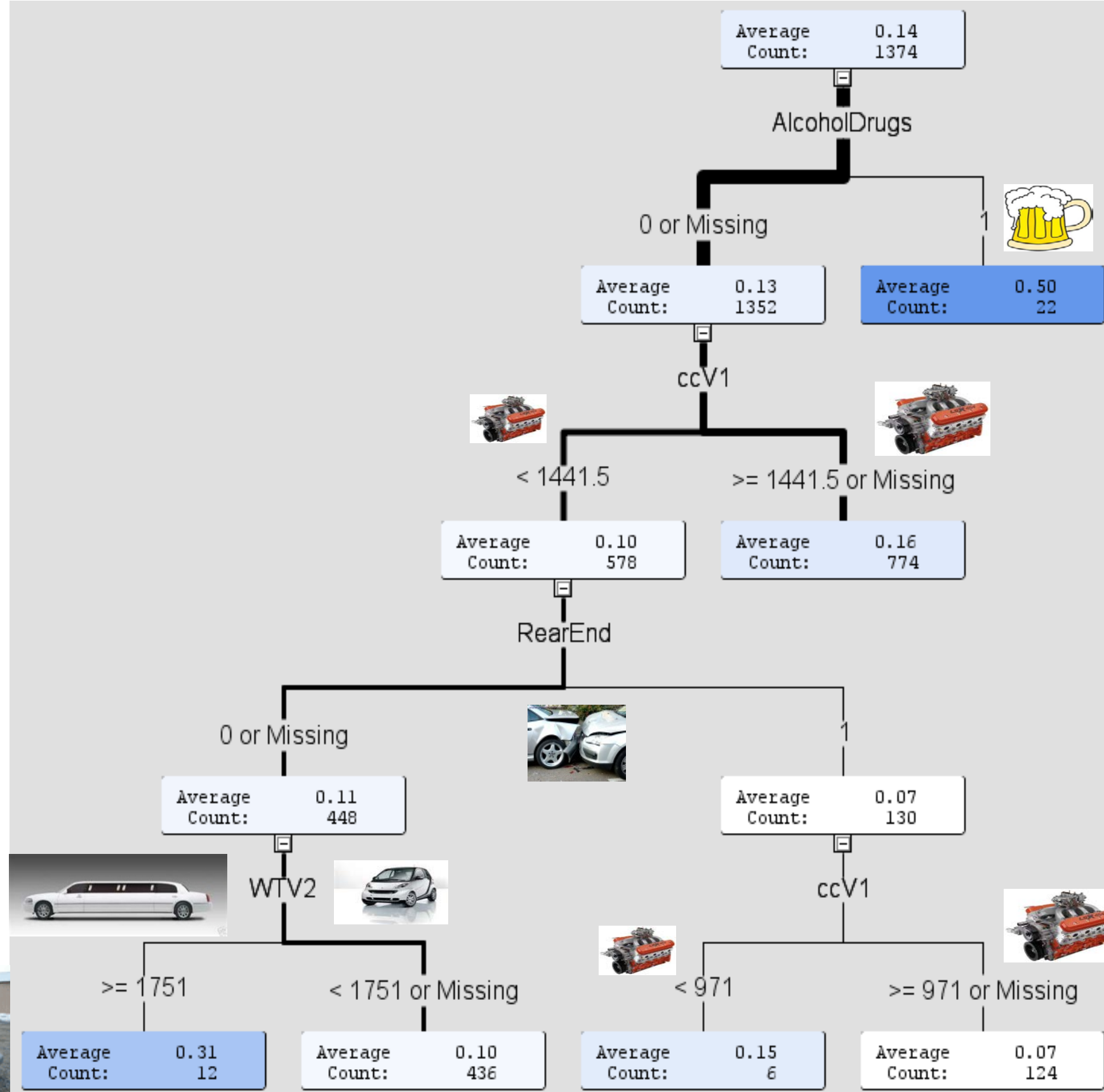- Predicted response $P_i$ constant in regions $i=1, \ldots, 5$

# Regression Trees

- Predict $P_i$ in cell i.

- $Y_{ij}$ $j^{th}$ response in cell i.

- Split to minimize $\Sigma_i \Sigma_j (Y_{ij} - P_i)^2$

Real data
example:
Traffic accidents
in Portugal*

Y = injury
induced "cost to
society"

Help - I ran
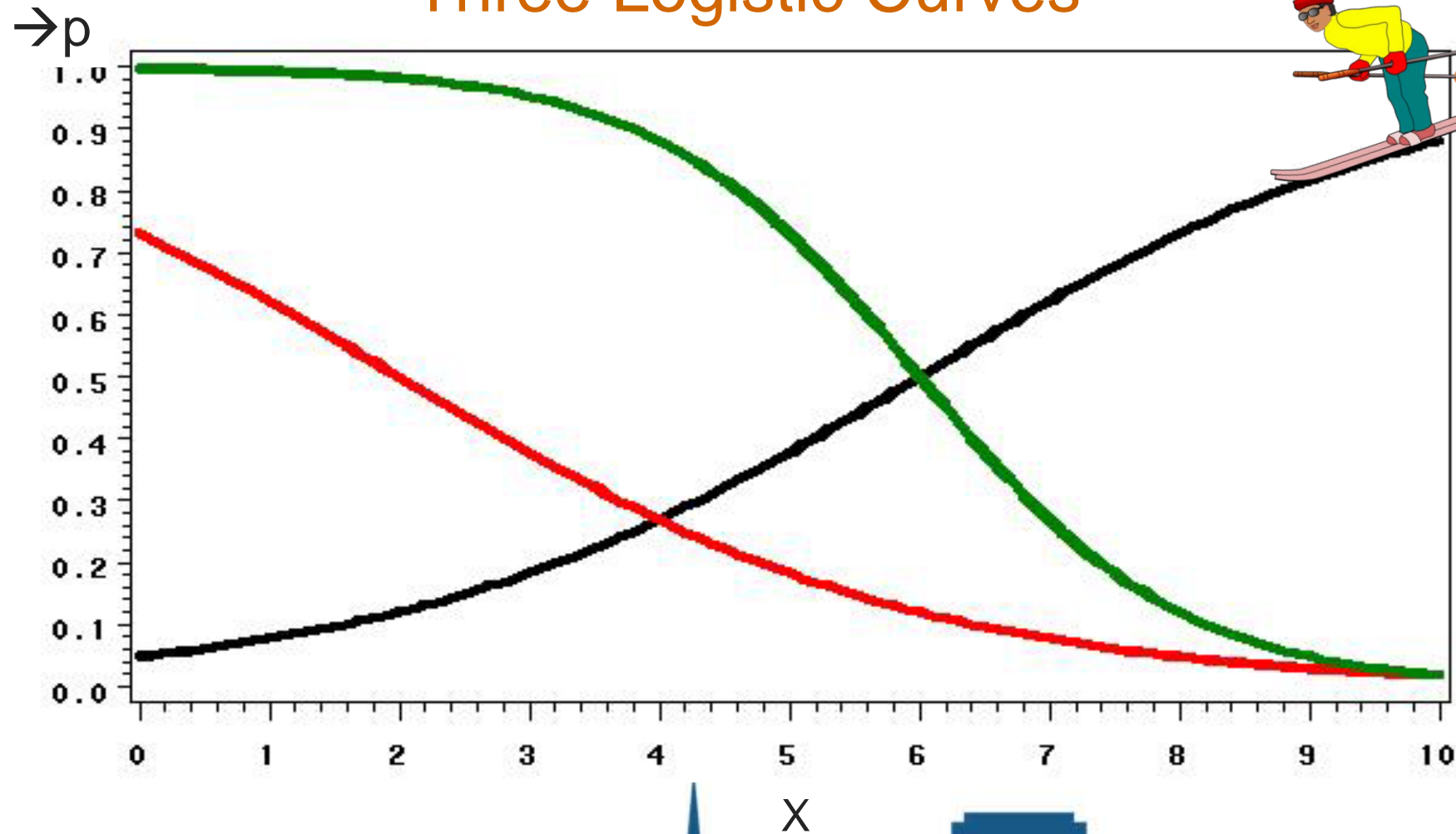Into  a "tree"

Help - I ran
Into  a "tree"

# Logistic Regression

- Logistic – another classifier

- Older – "tried & true" method

- Predict probability of response from input variables ("Features")

- Linear regression gives infinite range of predictions

- 0 < probability < 1 so not linear regression.

# Logistic Regression

$$\frac{e^{a+bX}}{(1+e^{a+bX})} \rightarrow p$$

Three Logistic Curves



X
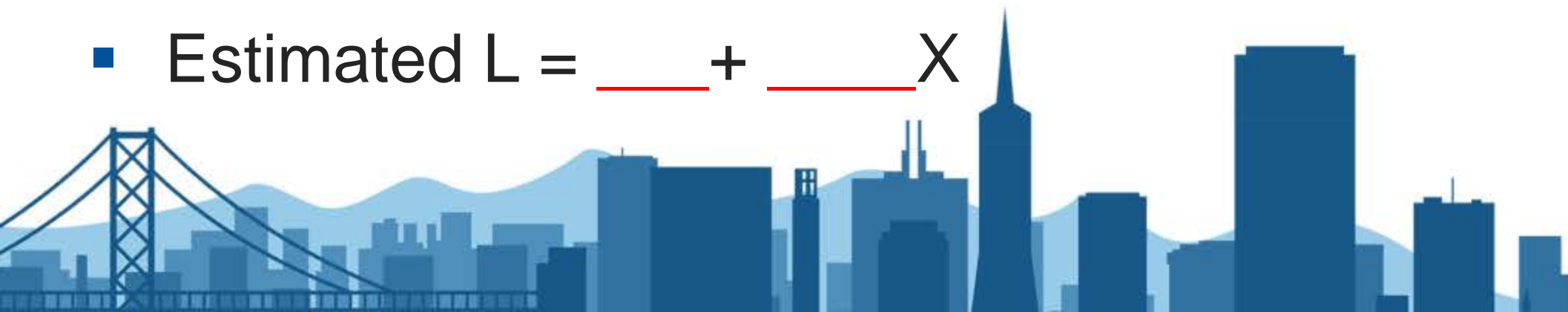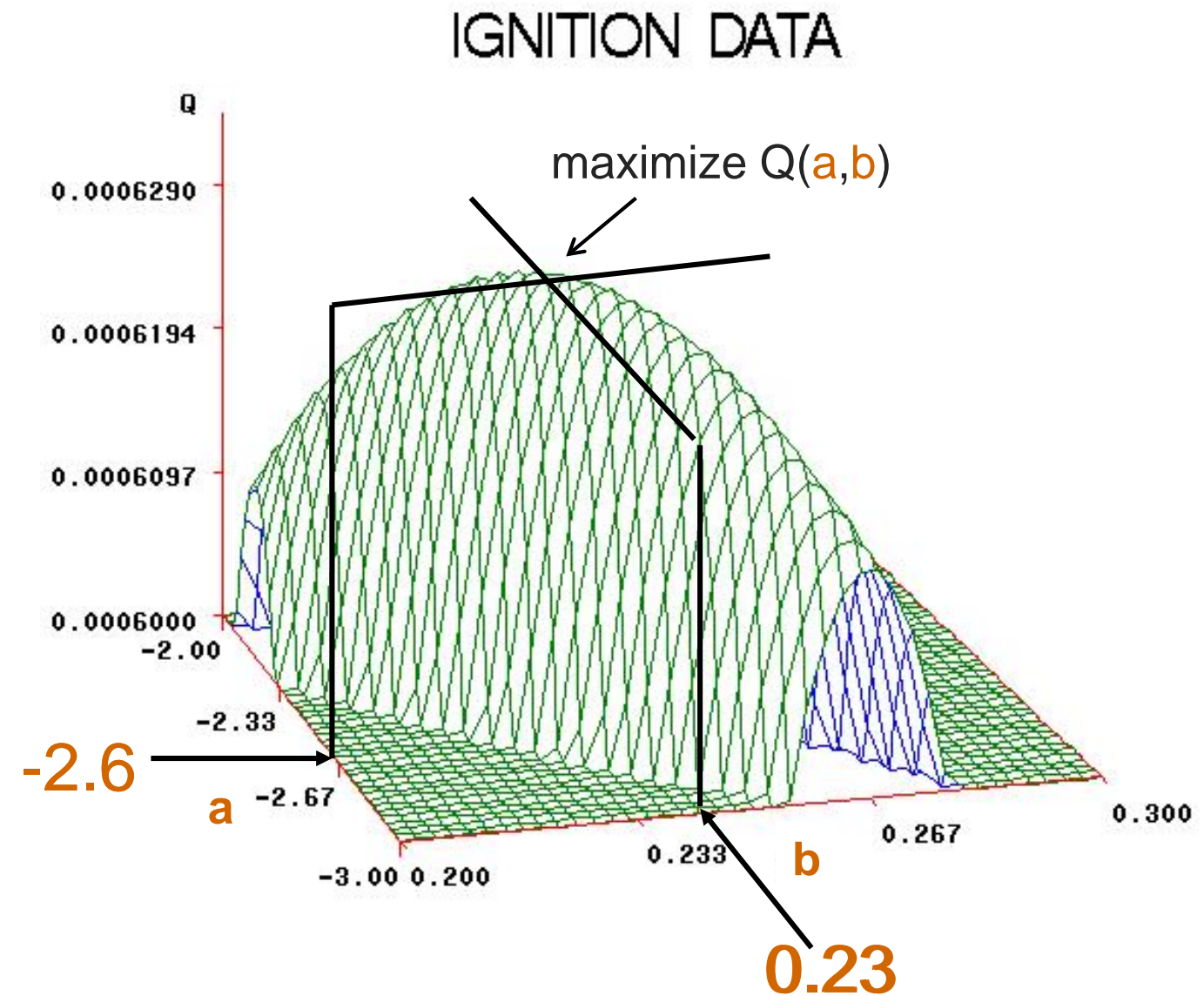
# Example: Seat Fabric Ignition

- Flame exposure time = X

- Y=1 → ignited,     Y=0→ did not ignite
  - Y=0, X= 3, 5, 9 10 ,                    🔥 13,      🔥 16          🔥
  - Y=1, X =                    11, 12      14, 15,          17, 25, 30

- $Q=(1-p_1)(1-p_2)(1-p_3)(1-p_4)p_5 p_6 (1-p_7)p_8 p_9 (1-p_{10})p_{11}p_{12}p_{13}$

- p's all different  $p_i = f(a+bX_i) = e^{a+bX_i}/(1+e^{a+bX_i})$

- Find a,b to maximize Q(a,b)

- Logistic idea:

- Given temperature X, compute $L(x) = a + bX$ then $p = e^L/(1+e^L)$

- $p(i) = e^{a+bX_i}/(1+e^{a+bX_i})$

- Write p(i) if response, 1-p(i) if not

- Multiply all n of these together, find a,b to maximize this "likelihood"

- Estimated L = ___ + ____X

IGNITION DATA



maximize Q(a,b)

Q

0.0006290

0.0006194

0.0006097

0.0006000
-2.00

-2.33

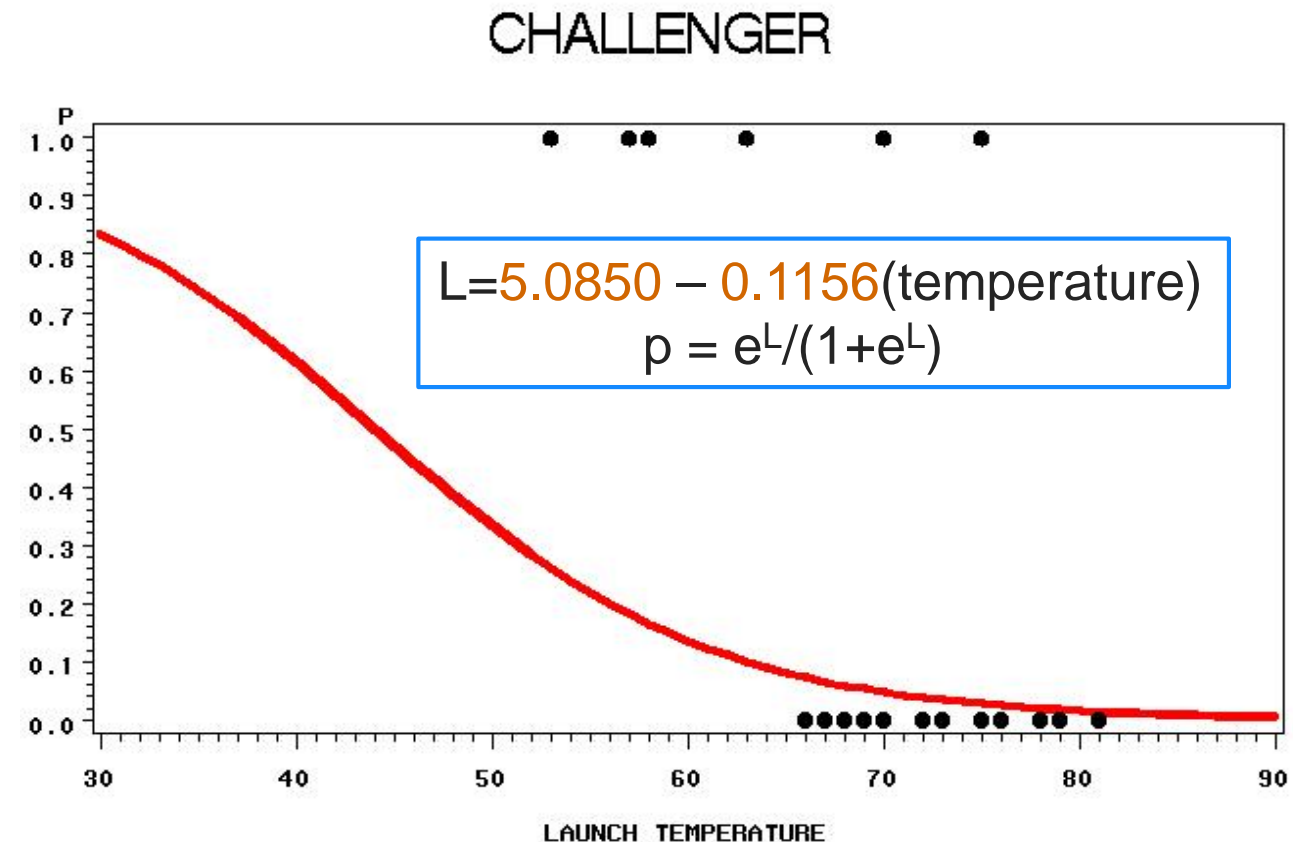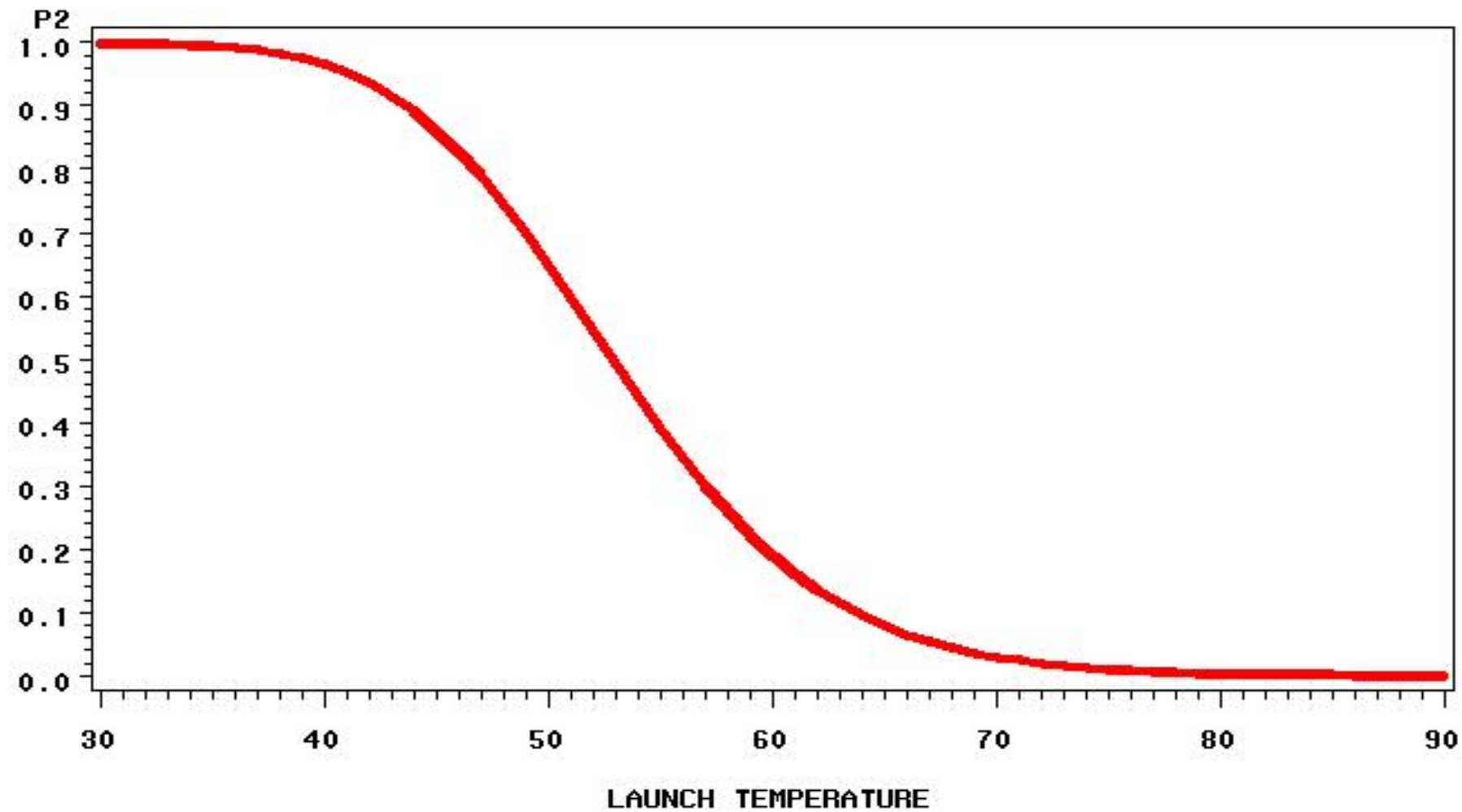-2.6

a  -2.67

-3.00 0.200

0.233

b

0.267

0.300

0.23

# Example: Shuttle Missions

- O-rings failed in Challenger disaster

- Prior flights "erosion" and "blowby" in O-rings (6 per mission)

- Feature: Temperature at liftoff

- Target: (1) - erosion or blowby vs. no problem (0)

CHALLENGER

$$L = 5.0850 - 0.1156(\text{temperature})$$
$$p = e^L/(1+e^L)$$

LAUNCH TEMPERATURE

Pr{2 OR MORE FAILURES}

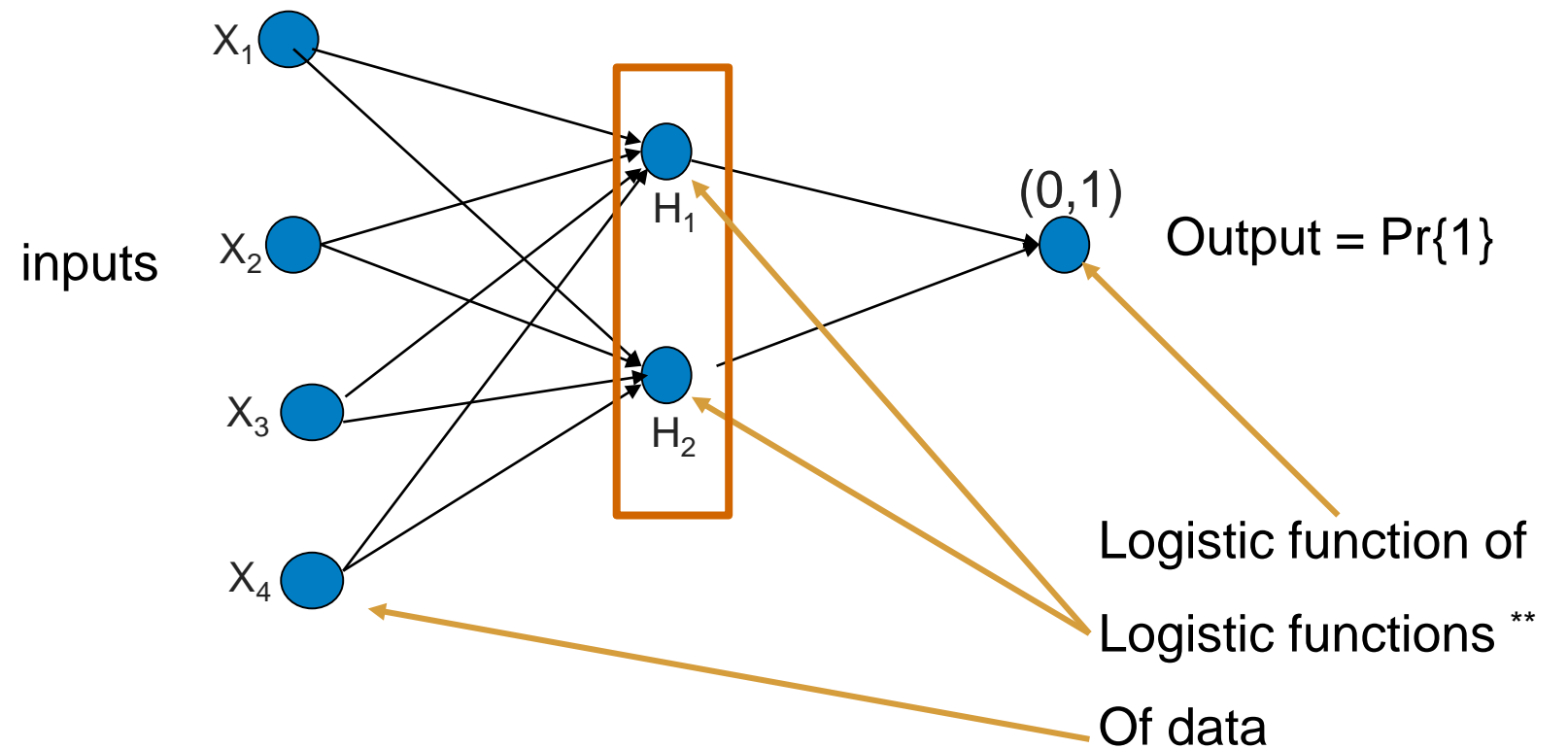$$\Pr\{2 \ or \ more\} = 1 - p_X{}^0(1 - p_X)^6 - 6p_X(1 - p_X)^5$$

# Neural Networks

- Very flexible functions
- "Hidden Layers"
- "Multilayer Perceptron"

inputs

$X_1$  $X_2$  $X_3$  $X_4$

$H_1$  $H_2$

$(0,1)$  Output = Pr{1}

Logistic function of

Logistic functions **

Of data

** (note: Hyperbolic tangent functions are just reparameterized logistic functions)

SAS. GLOBAL FORUM

Three Logistic Curves

$p(i) = \exp(L(i))/(1+\exp(L(i)))$ where $L(1) = -3+.5X$, $L(2)=1-0.5X$, and $L(3)=6-X$

Example:
$Y = a + b1\ H1 + b2\ H2 + b3\ H3$
$Y = 4 + 1\ H1 + 2\ H2 - 4\ H3$

"bias"

"weights"

$H_1$    $b_1$

$H_2$    $b_2$         Y

$H_3$    $b_3$

Arrows represent linear combinations of "basis functions," e.g. logistic curves (hyperbolic tangents)

Combining for Neural Network
$4 + p(1) + 2\ p(2) - 4\ p(3)$

# A Complex Neural Network Surface



("biases")

* Cumulative Lift Chart
  - Go from leaf of most to least predicted response.
  - Lift is

    proportion responding in first p%
    overall population response rate

Lift
3.3→

1→

Lift Chart

response

Predicted pct response

high ←-------------------------------→ low

# A Combined Example

Cell Phone Texting Locations

Black circle: ○
    Phone moved > 50 feet in first two minutes of texting.

Green dot: ●
    Phone moved < 50 feet. .



**Dots & Circles**

**Tree**     **Neural Net**     **Logistic Regression** ← Three Models

← Training Data
Lift Charts

← Validation Data
Lift Charts

← Resulting Surfaces

# Association Analysis is just elementary probability with new names



A: Purchase Milk

B: Purchase Cereal

0.3+0.2+0.1+0.4 = 1.0

Cereal=> Milk

Rule   B=> A  "people who buy B will buy A"

Support:
    Support= Pr{A and B} =  **0.2**

Independence means that Pr{A|B} = Pr{A} =  **0.5**
Pr{A} = **0.5** = Expected confidence if there is no
relation to B..

Confidence:
    Confidence = Pr{A|B}=Pr{A and B}/Pr{B}=**2/3**
??-  Is the confidence in B=>A the same as the
confidence in A=>B??  (yes, no)

Lift:
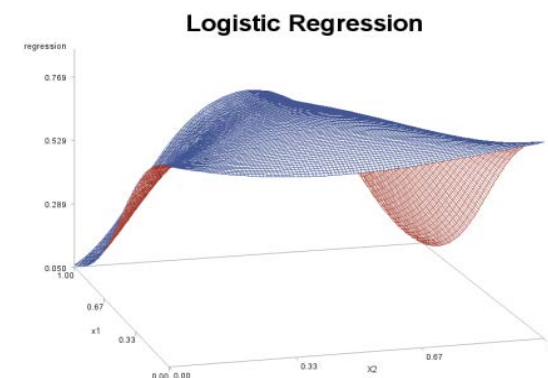    Lift = confidence / E{confidence} = (2/3) / (1/2) = **1.33**
    Gain = 33%

A
0.3

B
0.2  ←0.1

0.4

B

Marketing A to the 30%
of people who buy **B** will
result in 33% better sales
than marketing to a random
30% of the people.

SAS.**GLOBAL**FORUM

# Unsupervised Learning

- We have the "features" (predictors)

- We do NOT have the response even on a training data set (_UN_supervised)

- Another name for clustering

- EM

  - Large number of clusters with k-means (k clusters)
  - Ward's method to combine (less clusters)
  - One more k means

# Text Mining

Hypothetical collection of news releases  ("corpus")  :

release 1: Did the NCAA investigate the basketball scores and
                vote for sanctions?
release 2: Republicans voted for and Democrats voted against
                it for the win.
   (etc.)

Compute word counts:

|  | NCAA | basketball | score | vote | Republican | Democrat | win |
|---|---|---|---|---|---|---|---|
| Release 1 | 1 | 1 | 1 | 1 | 0 | 0 | 0 |
| Release 2 | 0 | 0 | 0 | 2 | 1 | 1 | 1 |

# Text Mining Mini-Example: Word counts in 16 e-mails

←------------------------------words------------------------------→

| document | Election | President | Republican | Basketball | Democrat | Voters | NCAA | Liar | Tournament | Speech | Wins | Score_V | Score_N |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 20 | 8 | 10 | 12 | 6 | 0 | 1 | 5 | 3 | 8 | 18 | 15 | 21 |
| 2 | 5 | 6 | 9 | 5 | 4 | 2 | 0 | 9 | 0 | 12 | 12 | 9 | 0 |
| 3 | 0 | 2 | 0 | 14 | 0 | 2 | 12 | 0 | 16 | 4 | 24 | 19 | 30 |
| 4 | 8 | 9 | 7 | 0 | 12 | 14 | 2 | 12 | 3 | 15 | 22 | 8 | 2 |
| 5 | 0 | 0 | 4 | 16 | 0 | 0 | 15 | 2 | 17 | 3 | 9 | 0 | 1 |
| 6 | 10 | 6 | 9 | 5 | 5 | 19 | 5 | 20 | 0 | 18 | 13 | 9 | 14 |
| 7 | 2 | 3 | 1 | 13 | 0 | 1 | 12 | 13 | 20 | 0 | 0 | 1 | 6 |
| 8 | 4 | 1 | 4 | 16 | 2 | 4 | 9 | 0 | 12 | 9 | 3 | 0 | 0 |
| 9 | 26 | 13 | 9 | 2 | 16 | 20 | 6 | 24 | 4 | 30 | 9 | 10 | 14 |
| 10 | 19 | 22 | 10 | 11 | 9 | 12 | 0 | 14 | 10 | 22 | 3 | 1 | 0 |
| 11 | 2 | 0 | 0 | 14 | 1 | 3 | 12 | 0 | 16 | 12 | 17 | 23 | 8 |
| 12 | 16 | 19 | 21 | 0 | 13 | 9 | 0 | 16 | 4 | 12 | 0 | 0 | 2 |
| 13 | 14 | 17 | 12 | 0 | 20 | 19 | 0 | 12 | 5 | 9 | 6 | 1 | 4 |
| 14 | 1 | 0 | 4 | 21 | 3 | 6 | 9 | 3 | 8 | 0 | 3 | 10 | 20 |

## Eigenvalues of the Correlation Matrix

|   | Eigenvalue | Difference | Proportion | Cumulative |
|---|------------|------------|------------|------------|
| 1 | 7.10954264 | 4.80499109 | 0.5469 | **0.5469** |
| 2 | 2.30455155 | 1.30162837 | 0.1773 | 0.7242 |
| 3 | 1.00292318 | 0.23404351 | 0.0771 | 0.8013 |
| 4 | 0.76887967 | 0.21070080 | 0.0591 | 0.8605 |
| 5 | 0.55817886 | 0.10084923 | 0.0429 | 0.9034 |
|   |            | (more)     |        |        |
| 13 | 0.0008758 |            | 0.0001 | 1.0000 |

| Variable | Prin1 |
|----------|-------|
| Basketball | -.320074 |
| NCAA | -.314093 |
| Tournament | -.277484 |
| Score_V | -.134625 |
| Score_N | -.120083 |
| Wins | -.080110 |
| | |
| Speech | 0.273525 |
| Voters | 0.294129 |
| Liar | 0.309145 |
| Election | 0.315647 |
| Republican | 0.318973 |
| President | 0.333439 |
| Democrat | 0.336873 |

**Plotting 2 Words Only**



Prin 2

Prin 1

55% of the variation in these 13-dimensional vectors occurs in one dimension.

SAS.**GLOBAL**FORUM

## Eigenvalues of the Correlation Matrix

| | Eigenvalue | Difference | Proportion | Cumulative |
|---|---|---|---|---|
| 1 | 7.10954264 | 4.80499109 | 0.5469 | 0.5469 |
| 2 | 2.30455155 | 1.30162837 | 0.1773 | 0.7242 |
| 3 | 1.00292318 | 0.23404351 | 0.0771 | 0.8013 |
| 4 | 0.76887967 | 0.21070080 | 0.0591 | 0.8605 |
| 5 | 0.55817886 | 0.10084923 | 0.0429 | 0.9034 |
| | | (more) | | |
| 13 | 0.000875 | | 0.0001 | 1.0000 |

55% of the variation in these 13-dimensional vectors occurs in one dimension.



Plotting 2 Words Only

| Variable | Prin1 |
|---|---|
| Basketball | -.320074 |
| NCAA | -.314093 |
| Tournament | -.277484 |
| Score_V | -.134625 |
| Score_N | -.120083 |
| Wins | -.080110 |
| | |
| Speech | 0.273525 |
| Voters | 0.294129 |
| Liar | 0.309145 |
| Election | 0.315647 |
| Republican | 0.318973 |
| President | 0.333439 |
| Democrat | 0.336873 |

| document | CLUSTER | Print1 | Election | President | Republican | Basketball | Democrat | Voters | NCAA | Liar | Tournament | Speech | Wins | Score $\overline{V}$ | Score $\overline{N}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 3 | 1 | -3.63815 | 0 | 2 | 0 | 14 | 0 | 2 | 12 | 0 | 16 | 4 | 24 | 19 | 30 |
| 11 | 1 | -3.02803 | 2 | 0 | 0 | 14 | 1 | 3 | 12 | 0 | 16 | 12 | 17 | 23 | 8 |
| 5 | 1 | -2.98347 | 0 | 0 | 4 | 16 | 0 | 0 | 15 | 2 | 17 | 3 | 9 | 0 | 1 |
| 14 | 1 | -2.48381 | 1 | 0 | 4 | 21 | 3 | 6 | 9 | 3 | 8 | 0 | 3 | 10 | 20 |
| 7 | 1 | -2.37638 | 2 | 3 | 1 | 13 | 0 | 1 | 12 | 13 | 20 | 0 | 0 | 1 | 6 |
| 8 | 1 | -1.79370 | 4 | 1 | 4 | 16 | 2 | 4 | 9 | 0 | 12 | 9 | 3 | 0 | 0 |

(biggest gap)

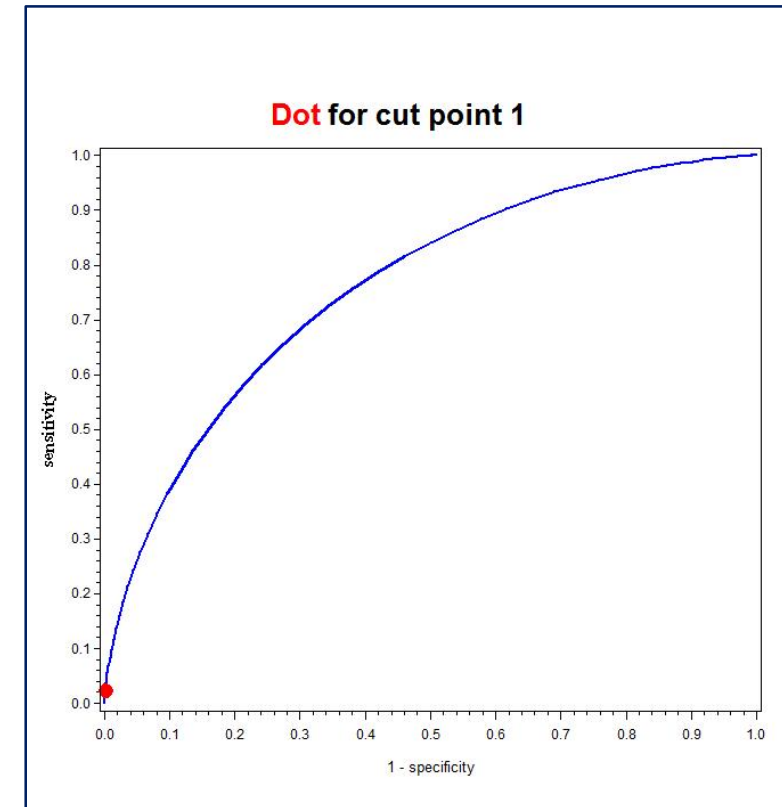| document | CLUSTER | Print1 | Election | President | Republican | Basketball | Democrat | Voters | NCAA | Liar | Tournament | Speech | Wins | Score $\overline{V}$ | Score $\overline{N}$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|---|
| 1 | 2 | -0.00738 | 20 | 8 | 10 | 12 | 6 | 0 | 1 | 5 | 3 | 8 | 18 | 15 | 21 |
| 2 | 2 | 0.48514 | 5 | 6 | 9 | 5 | 4 | 2 | 0 | 9 | 0 | 12 | 12 | 9 | 0 |
| 6 | 2 | 1.54559 | 10 | 6 | 9 | 5 | 5 | 19 | 5 | 20 | 0 | 18 | 13 | 9 | 14 |
| 4 | 2 | 1.59833 | 8 | 9 | 7 | 0 | 12 | 14 | 2 | 12 | 3 | 15 | 22 | 8 | 2 |
| 10 | 2 | 2.49069 | 19 | 22 | 10 | 11 | 9 | 12 | 0 | 14 | 10 | 22 | 3 | 1 | 0 |
| 13 | 2 | 3.16620 | 14 | 17 | 12 | 0 | 20 | 19 | 0 | 12 | 5 | 9 | 6 | 1 | 4 |
| 12 | 2 | 3.48420 | 16 | 19 | 21 | 0 | 13 | 9 | 0 | 16 | 4 | 12 | 0 | 0 | 2 |
| 9 | 2 | 3.54077 | 26 | 13 | 9 | 2 | 16 | 20 | 6 | 24 | 4 | 30 | 9 | 10 | 14 |

Sports Documents

Politics Documents

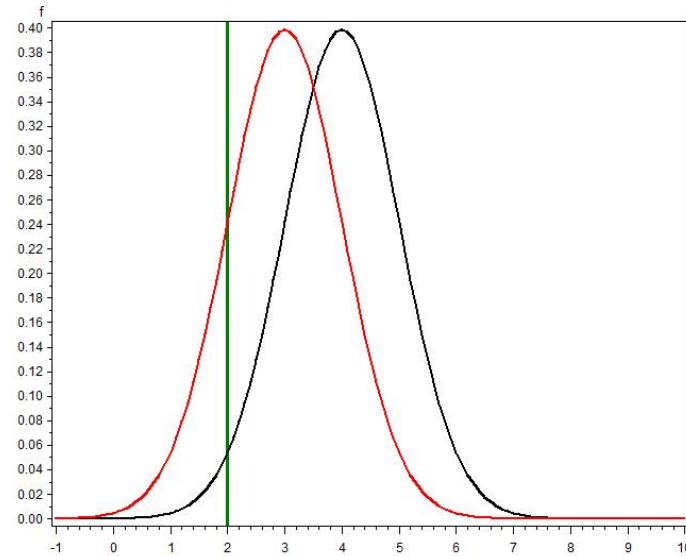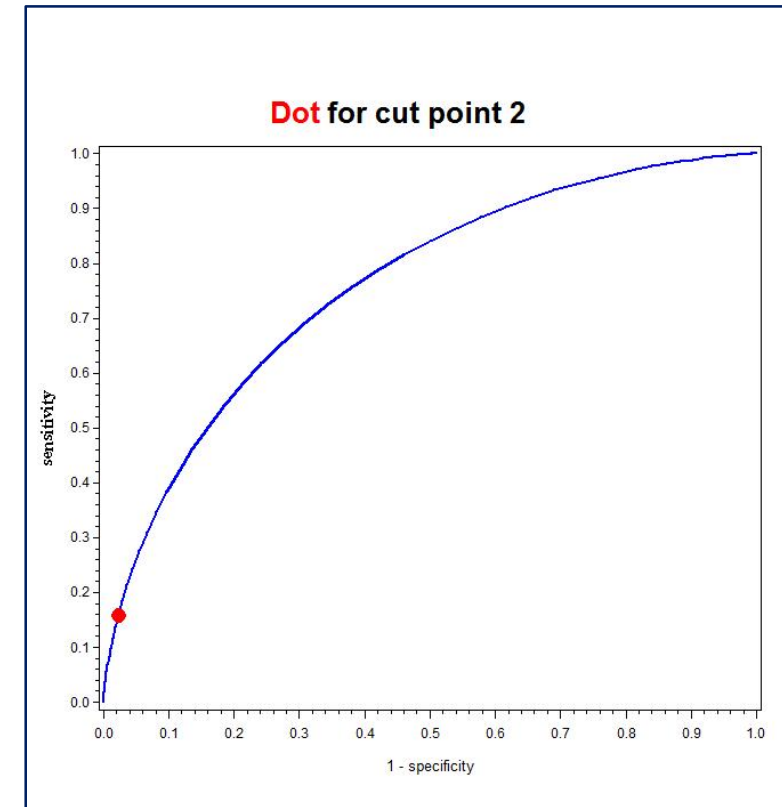# Receiver Operating Characteristic Curve

Cut point   1



Logits of 1s
red

Logits of 0s
black
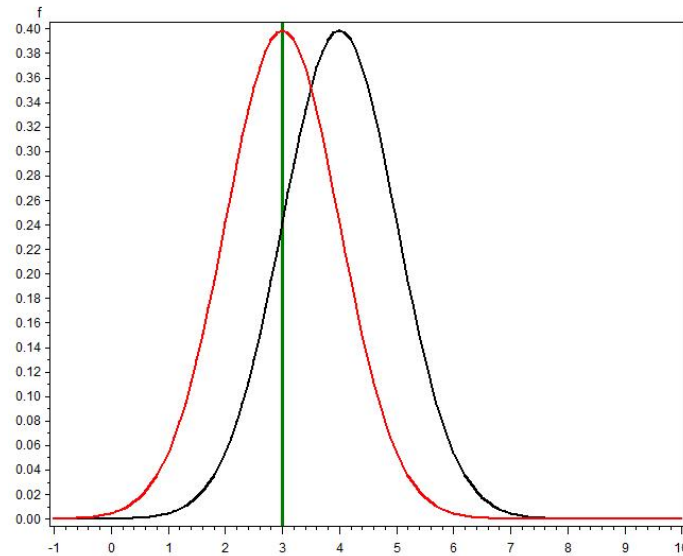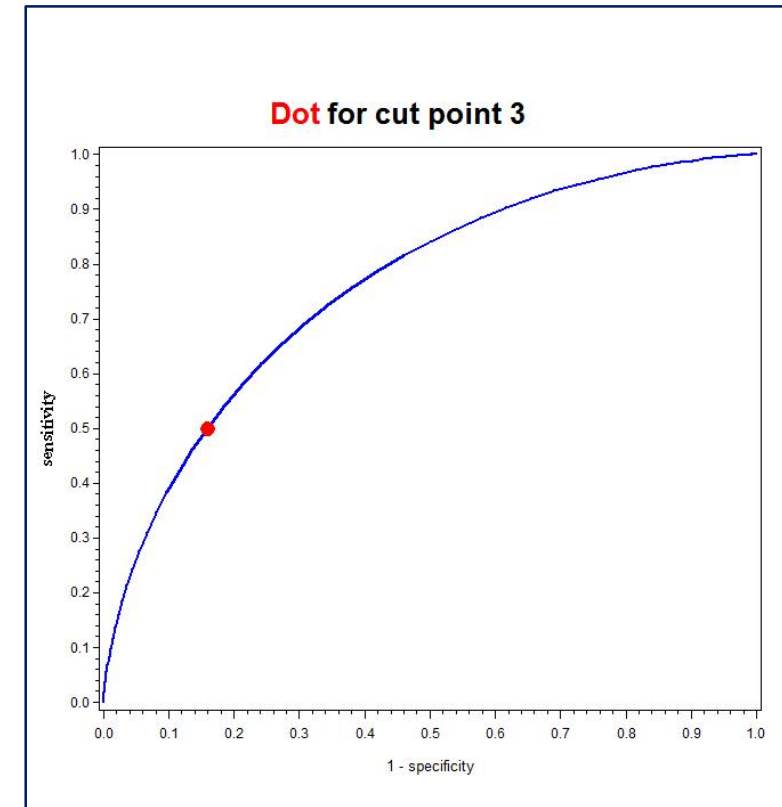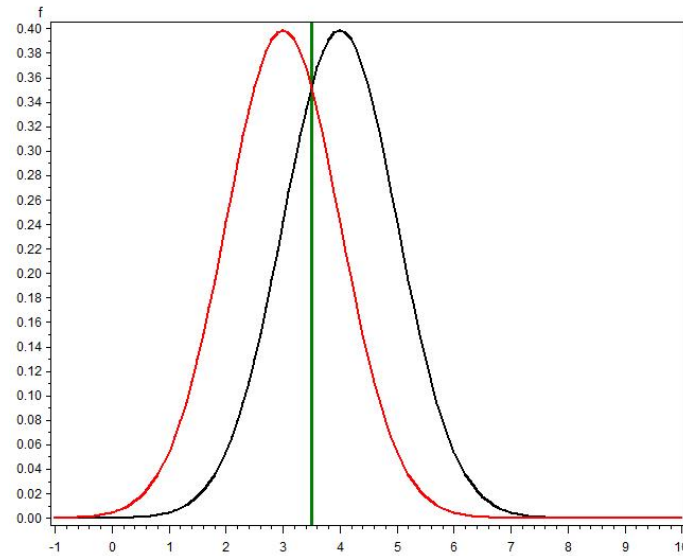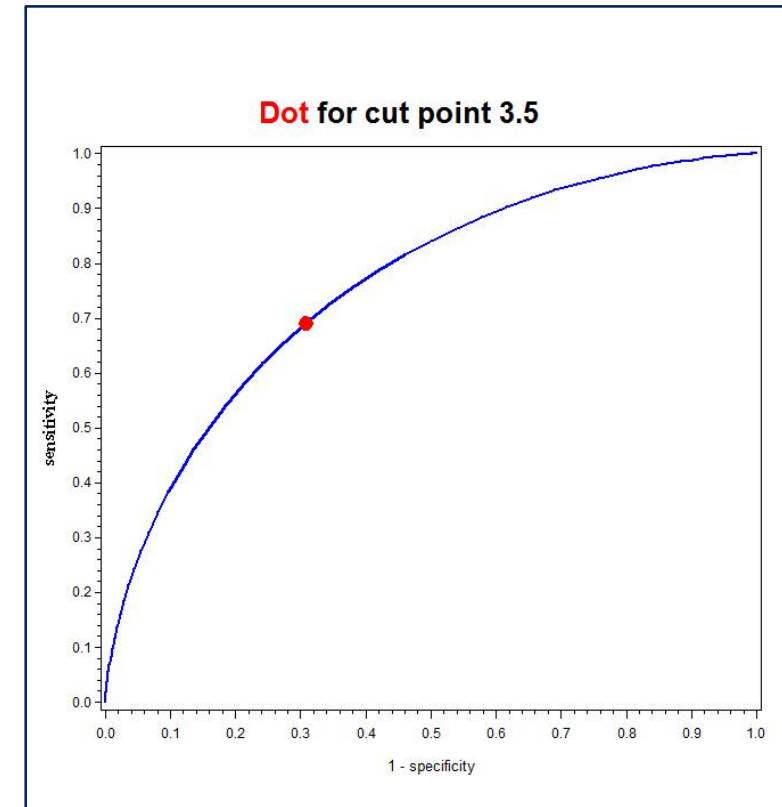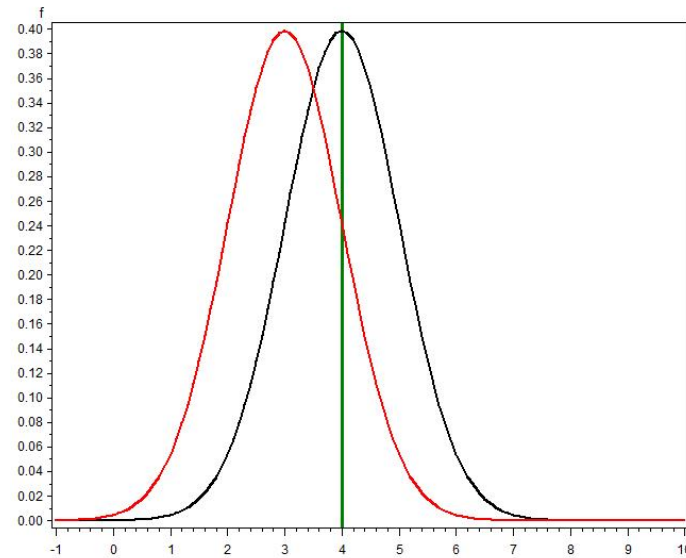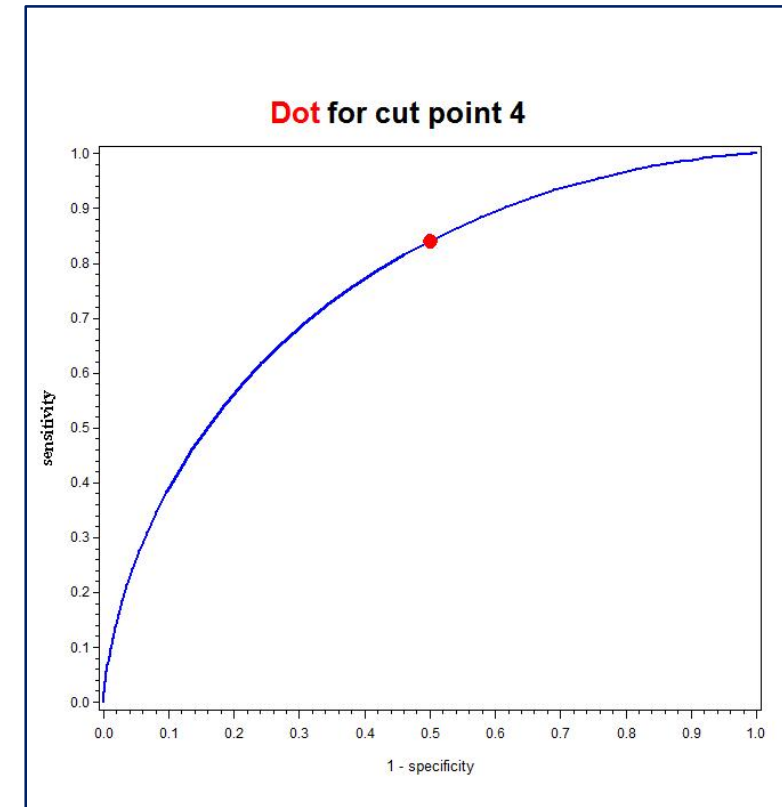
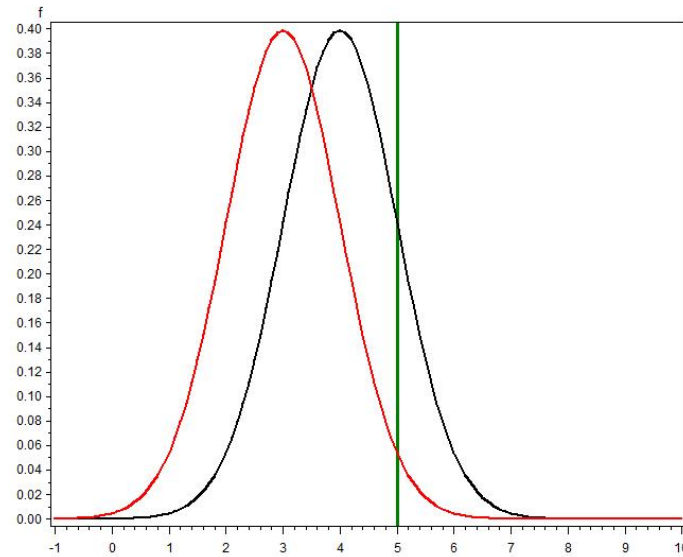# Receiver Operating Characteristic Curve

Cut point   2



Dot for cut point 2

Logits of 1s
red

Logits of 0s
black

# Receiver Operating Characteristic Curve

Cut point   3



Dot for cut point 3

Logits of 1s
red

Logits of 0s
black

SAS.**GLOBAL**FORUM

# Receiver Operating Characteristic Curve

Cut point   3.5



Logits of 1s
red

Logits of 0s
black



Dot for cut point 3.5

# Receiver Operating Characteristic Curve

Cut point   4



Dot for cut point 4

Logits of 1s
red

Logits of 0s
black

# Receiver Operating Characteristic Curve

Cut point   5



Logits of 1s
red

Logits of 0s
black



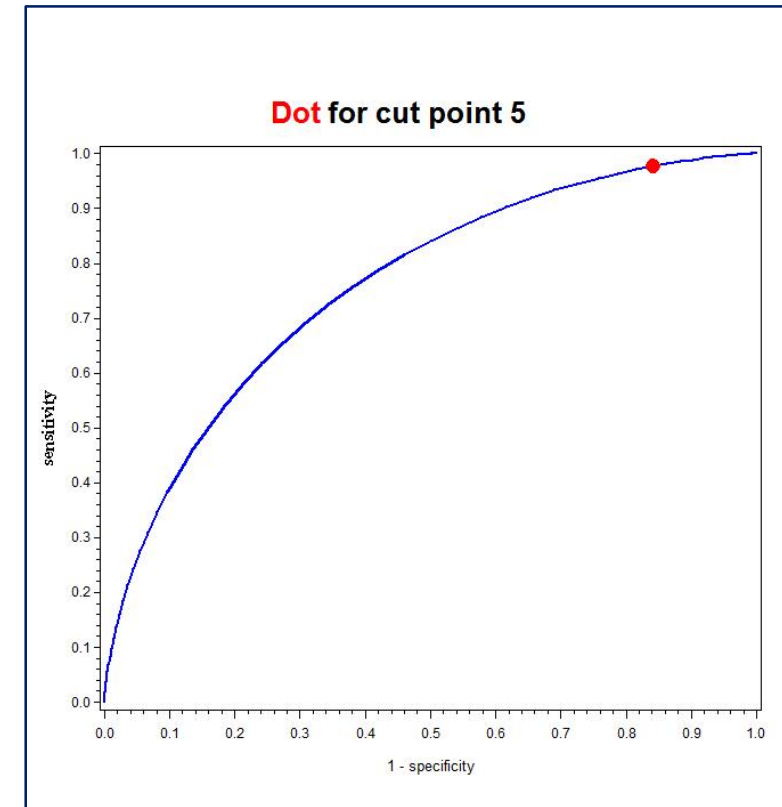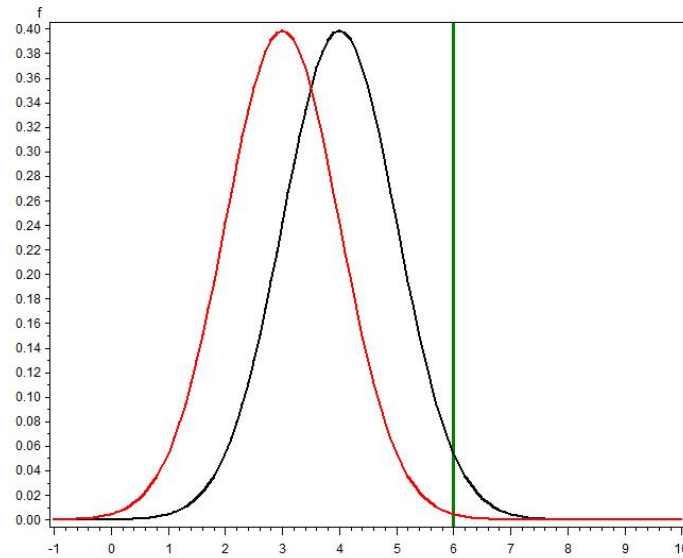SAS.**GLOBAL**FORUM

# Receiver Operating Characteristic Curve

Cut point   6



Logits of 1s
red

Logits of 0s
black



Dot for cut point 6

SAS.**GLOBAL**FORUM

"LOOKS LIKE MY TIME IS ALMOST UP"