

Building statistical models by visualization

Tom Minka
CMU Statistics Dept

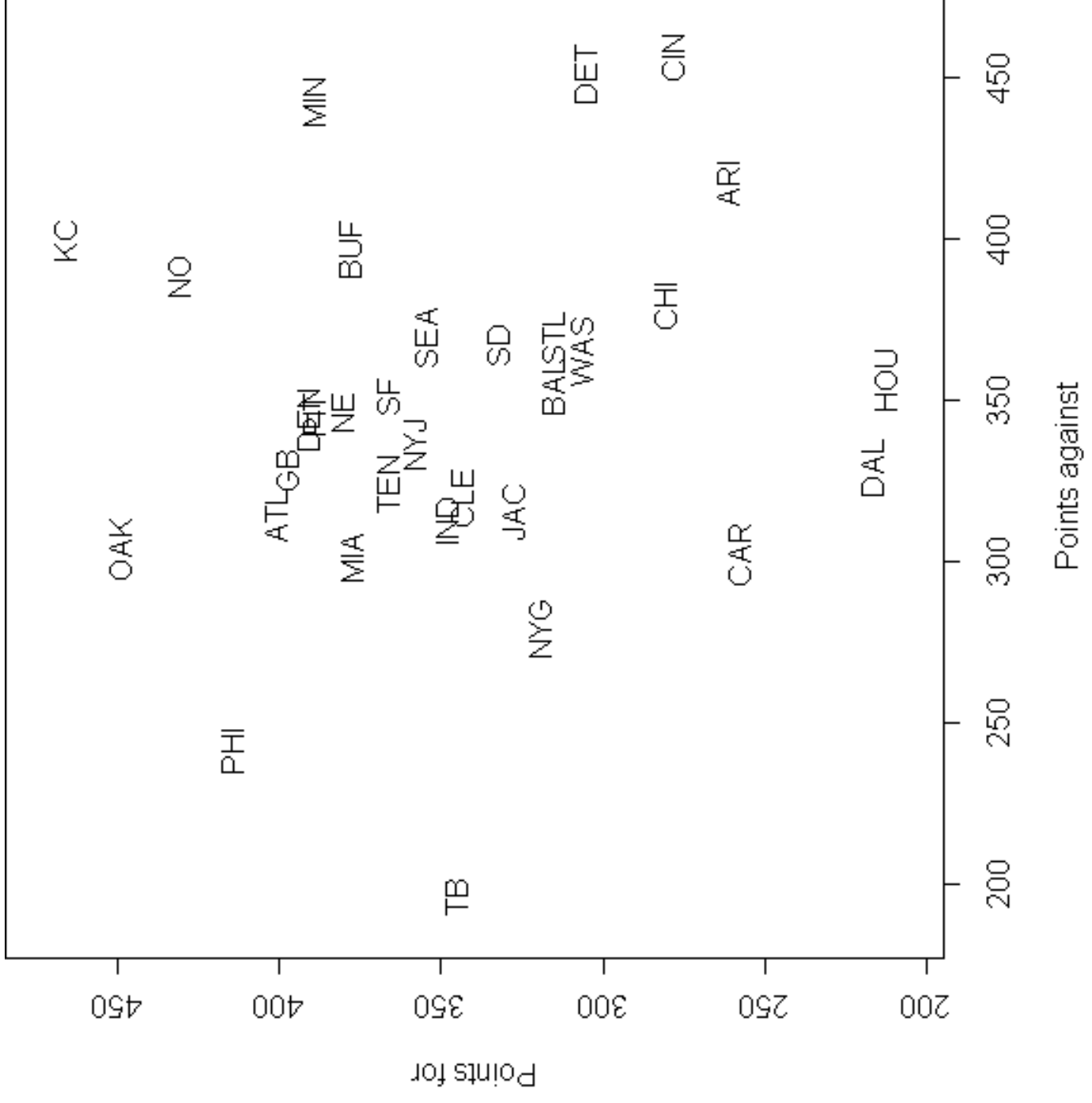
Outline

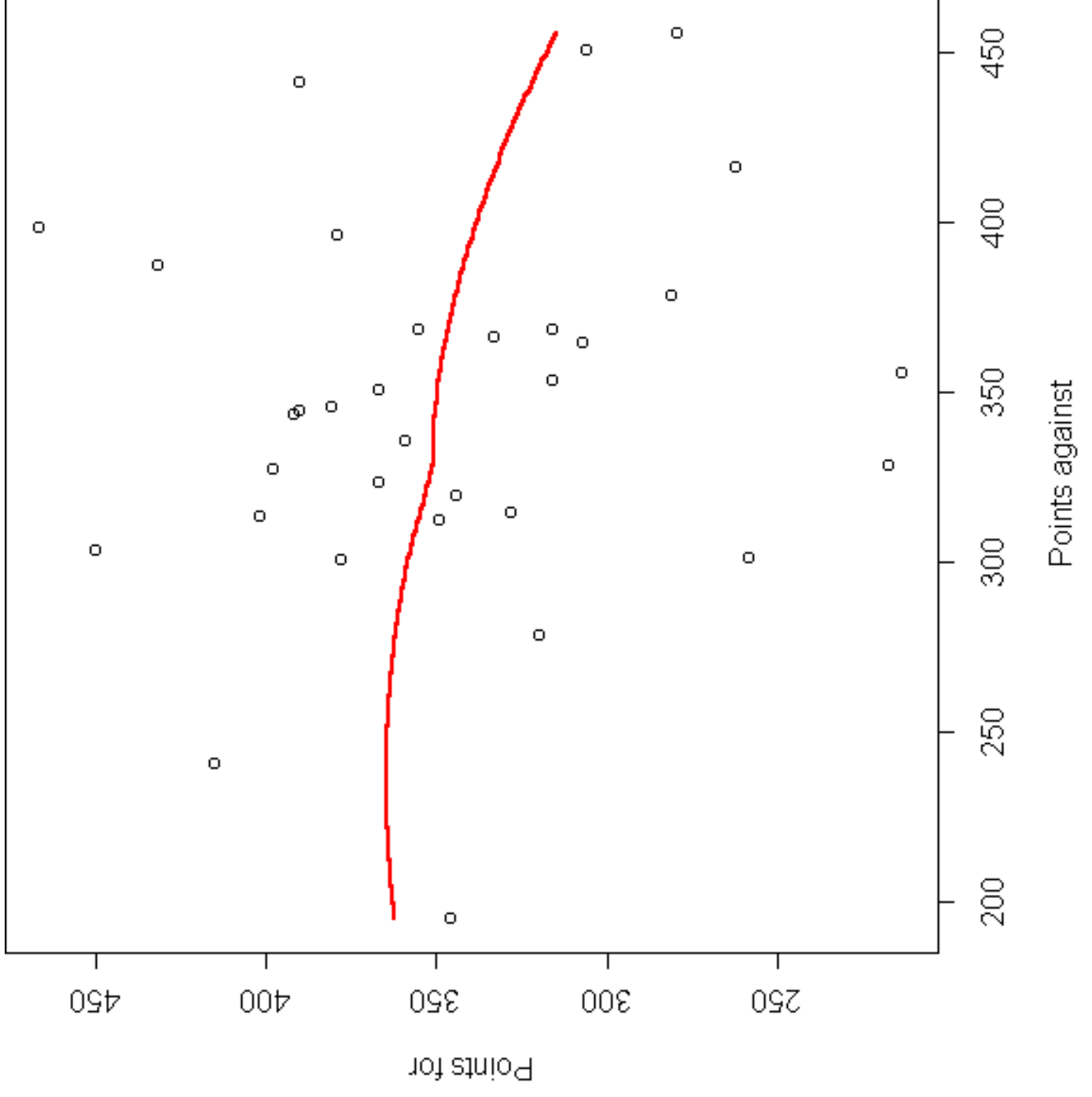
- Scatterplots
 - independence, causality
- QQ plots
 - distribution checking
- Residual plots
 - linearity, outliers
- Projections for regression
 - additivity
- Projections for classification
 - linearity

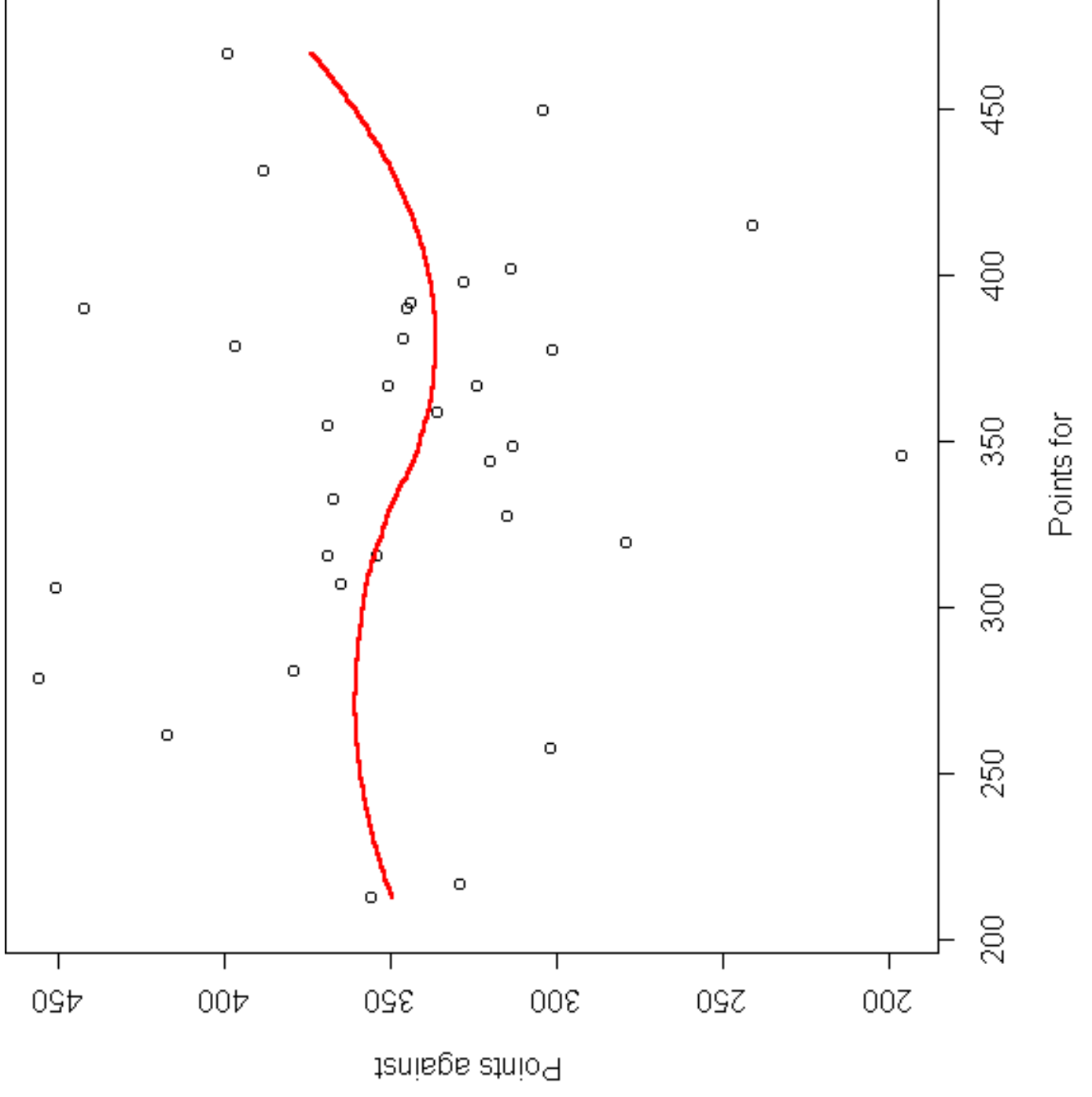
Football statistics

Is there a better representation?

	W	L	T	PF	PA
CIN	2	14	0	279	456
OAK	11	5	0	450	304
GB	12	4	0	398	328
SEA	7	9	0	355	369
JAC	6	10	0	328	315
NO	9	7	0	432	388
KC	8	8	0	467	399
TB	12	4	0	346	196
MIN	6	10	0	390	442
TEN	11	5	0	367	324
CAR	7	9	0	258	302
NYJ	9	7	0	359	336
NE	9	7	0	381	346
STL	7	9	0	316	369
...					

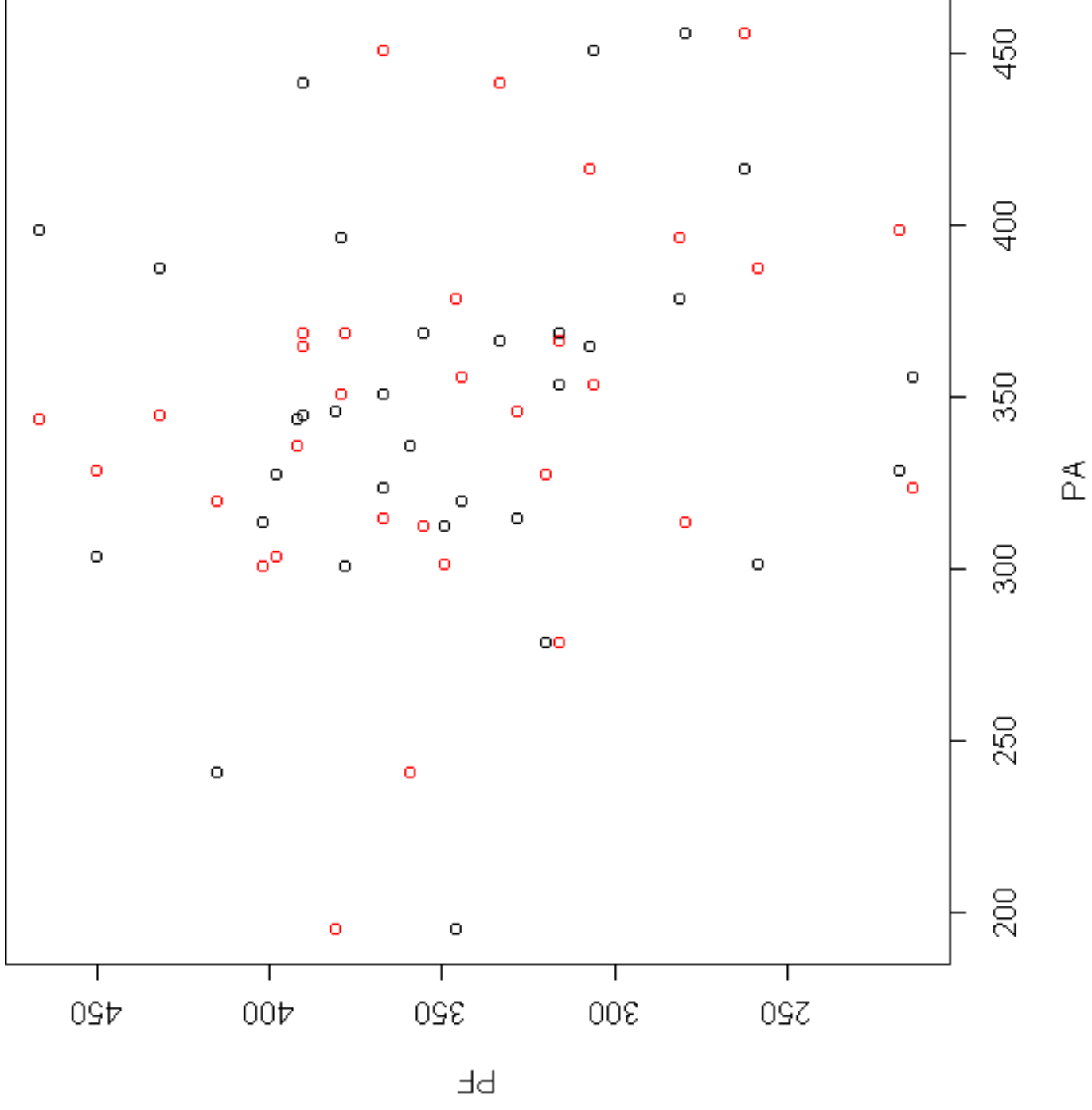


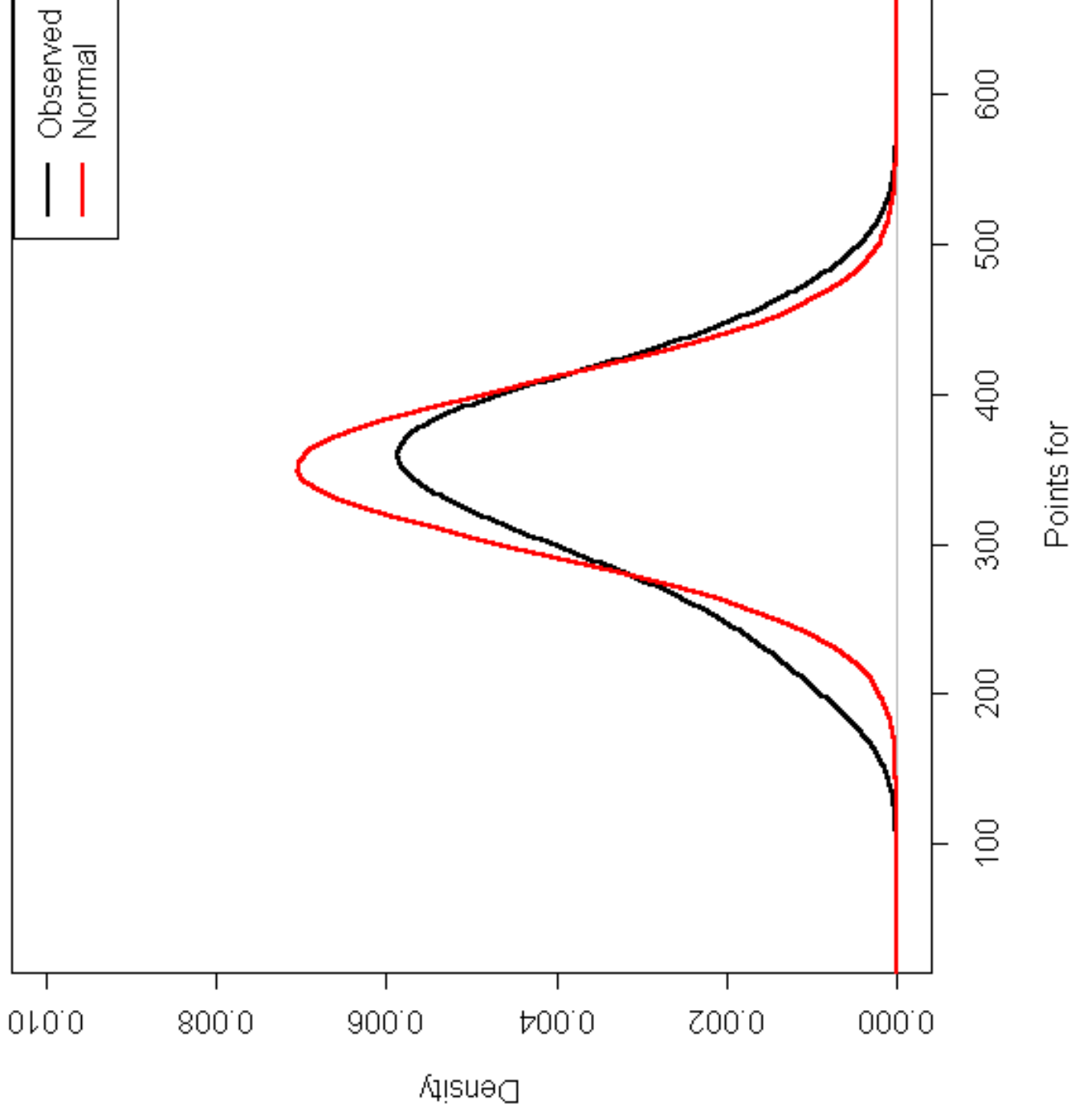




Visual independence test

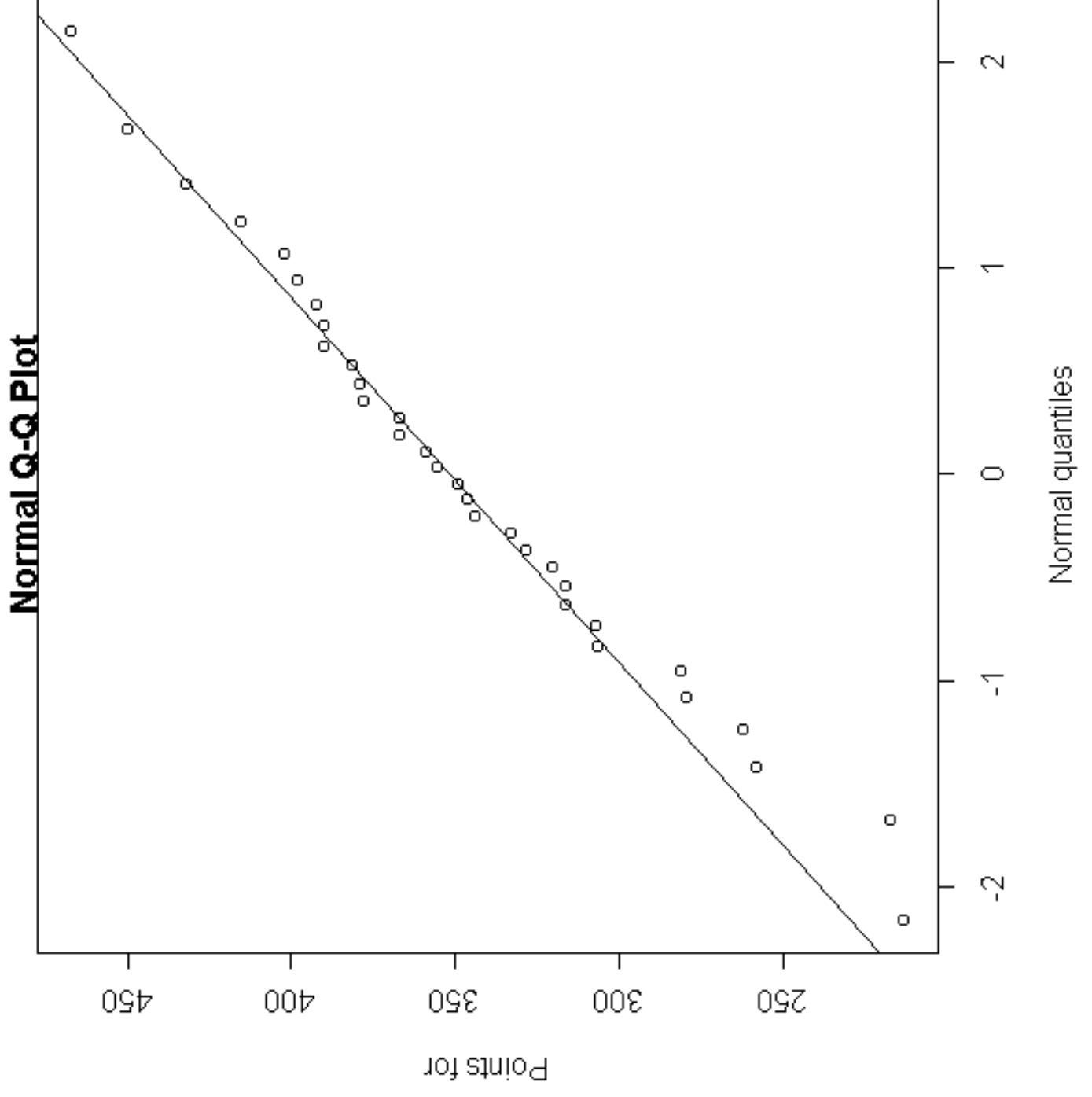
- “Permutation test”
- Randomly pair x values with y values
- If the distribution looks different from the original, the variables are dependent
- No distributional assumptions required





Comparing distributions

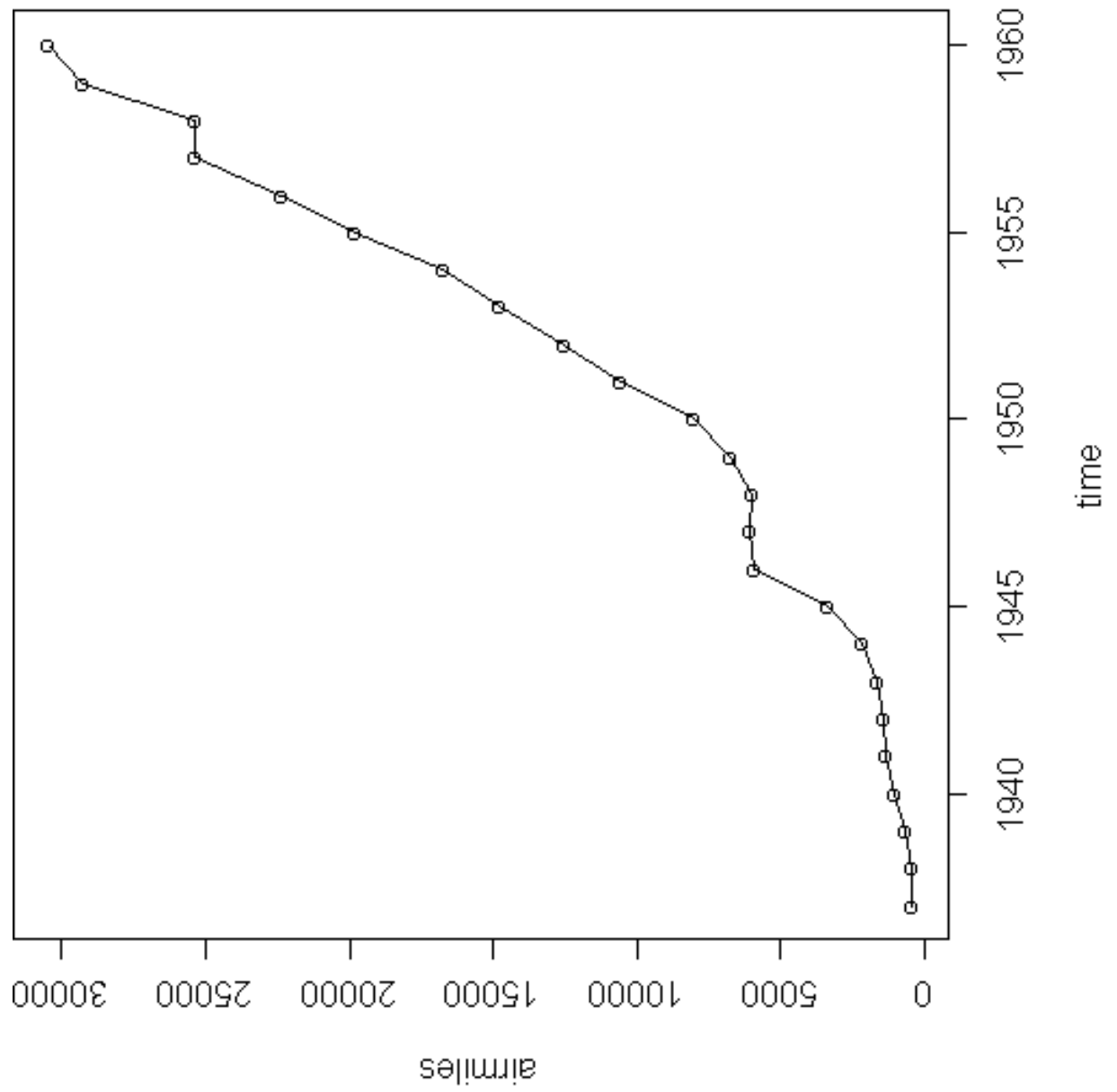
- “Quantile-quantile plot”
- A pseudo-scatterplot for unpaired data
- Quantile of x = fraction of points $< x$
- Plot quantile q in one set against quantile q in the other set, for all q
- Tells you how to transform one variable to have the distribution of the other

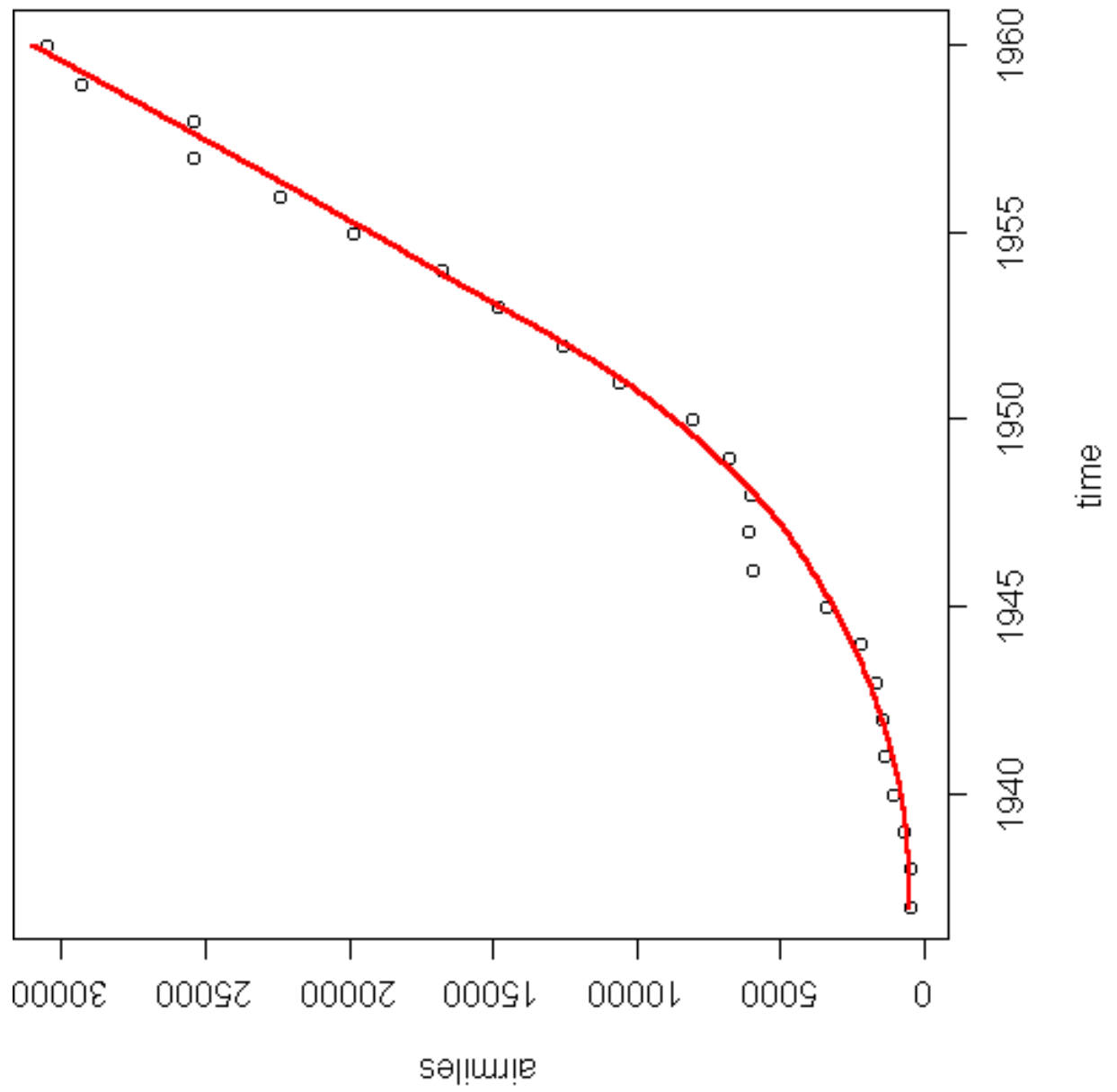


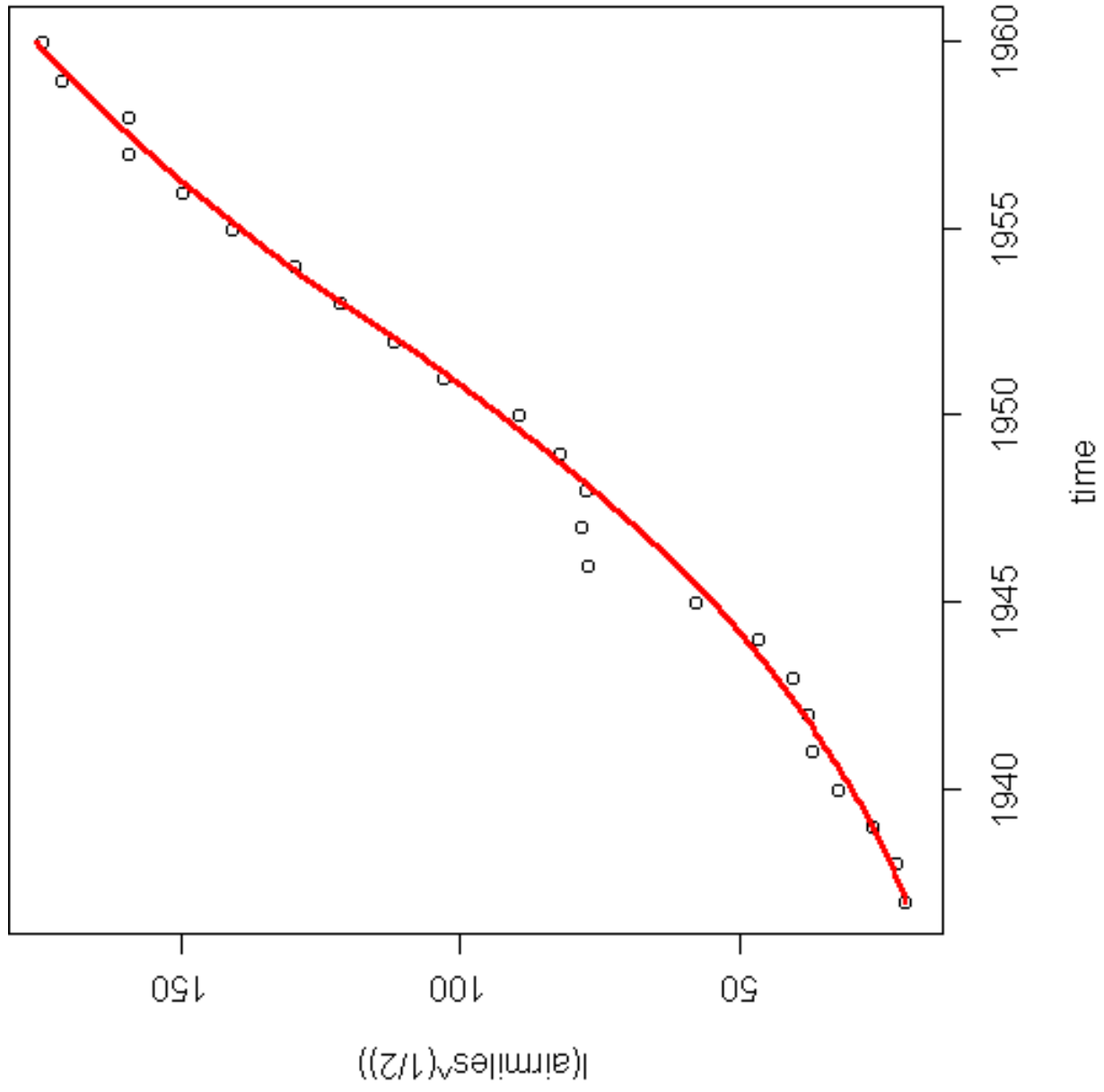
Regression models

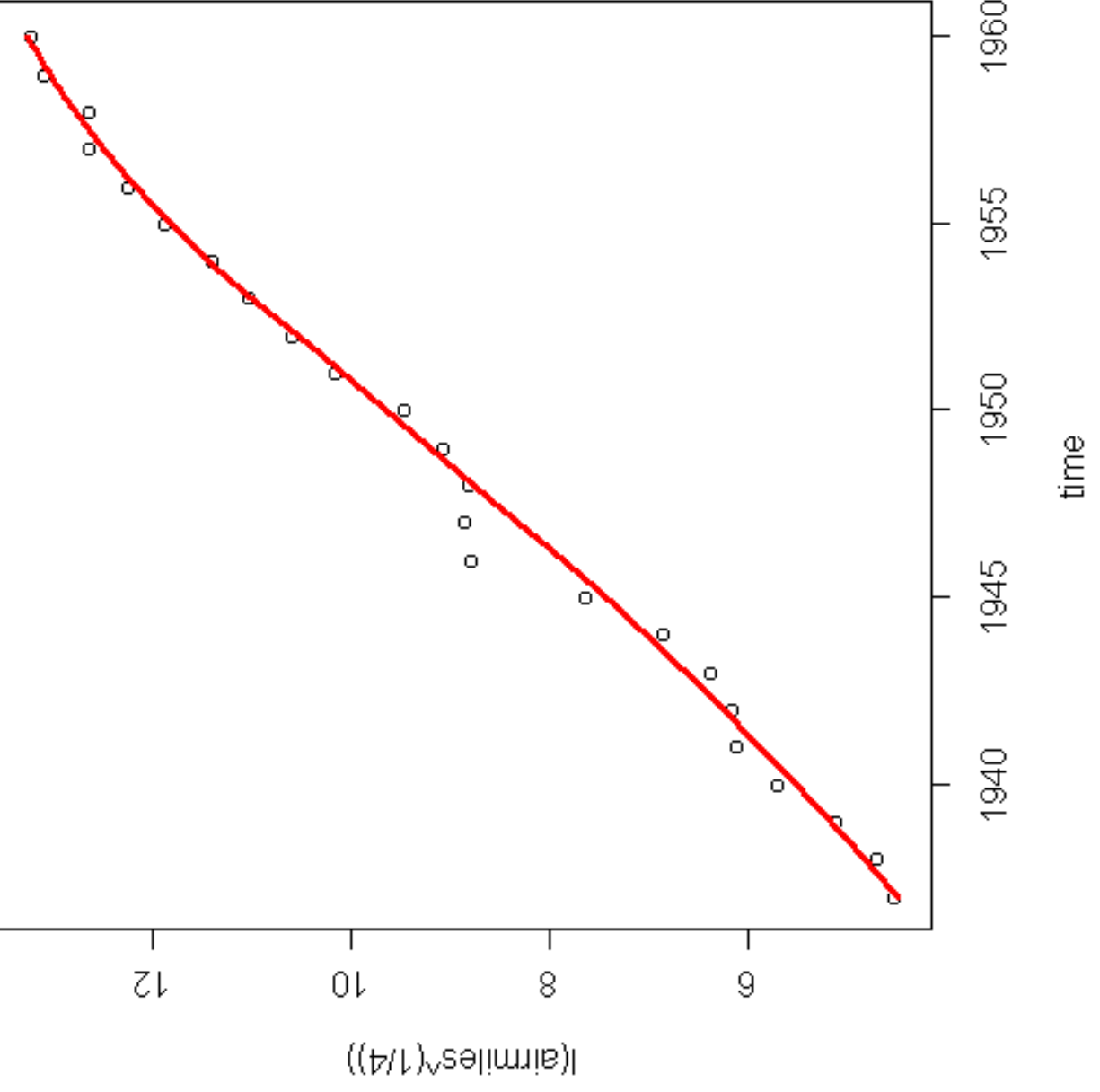
How to do regression visually:

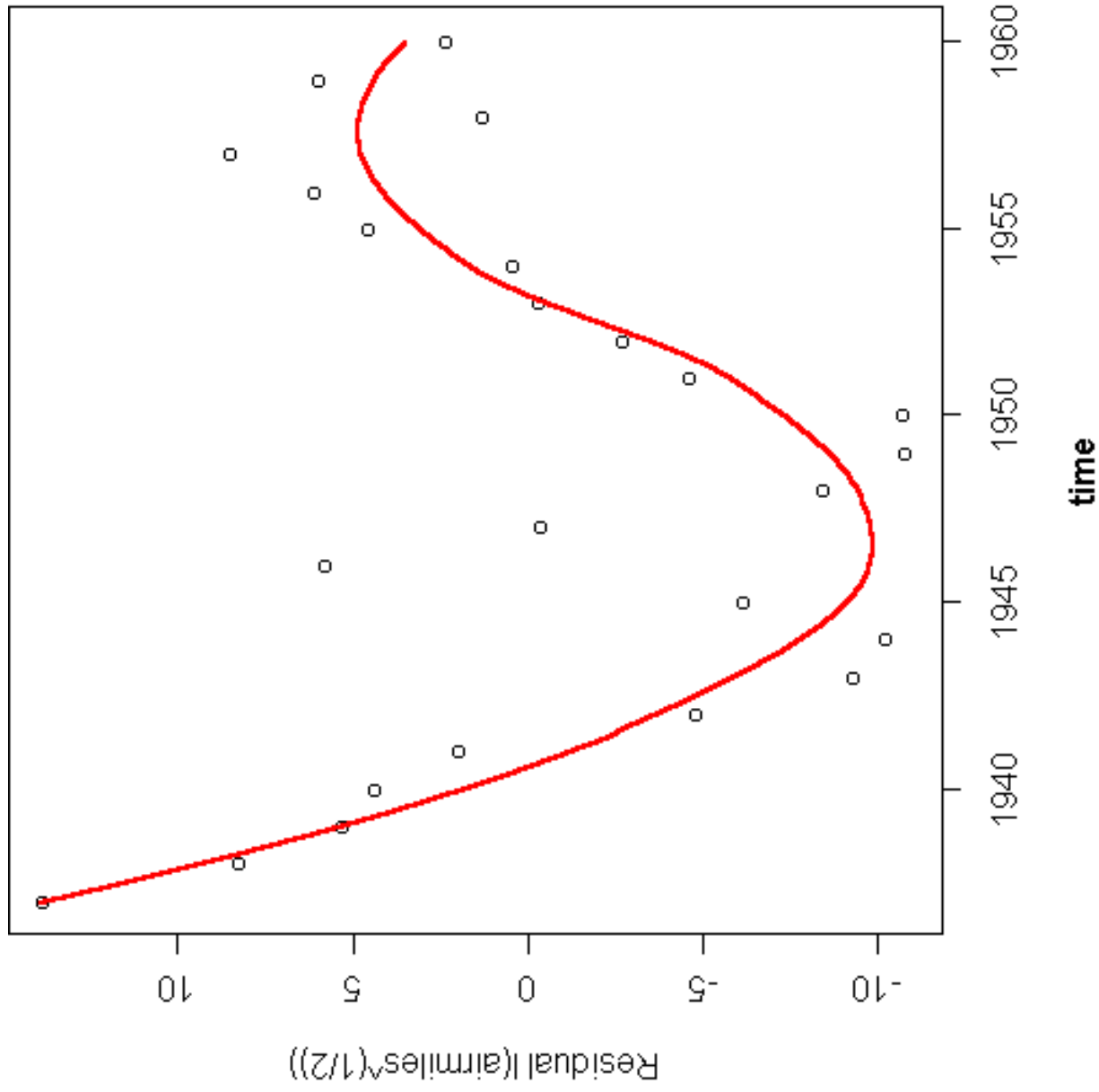
- Transform to make the picture simpler
- Fit a simple model
- Use residuals to suggest more complex models, outliers to remove
- Iterate

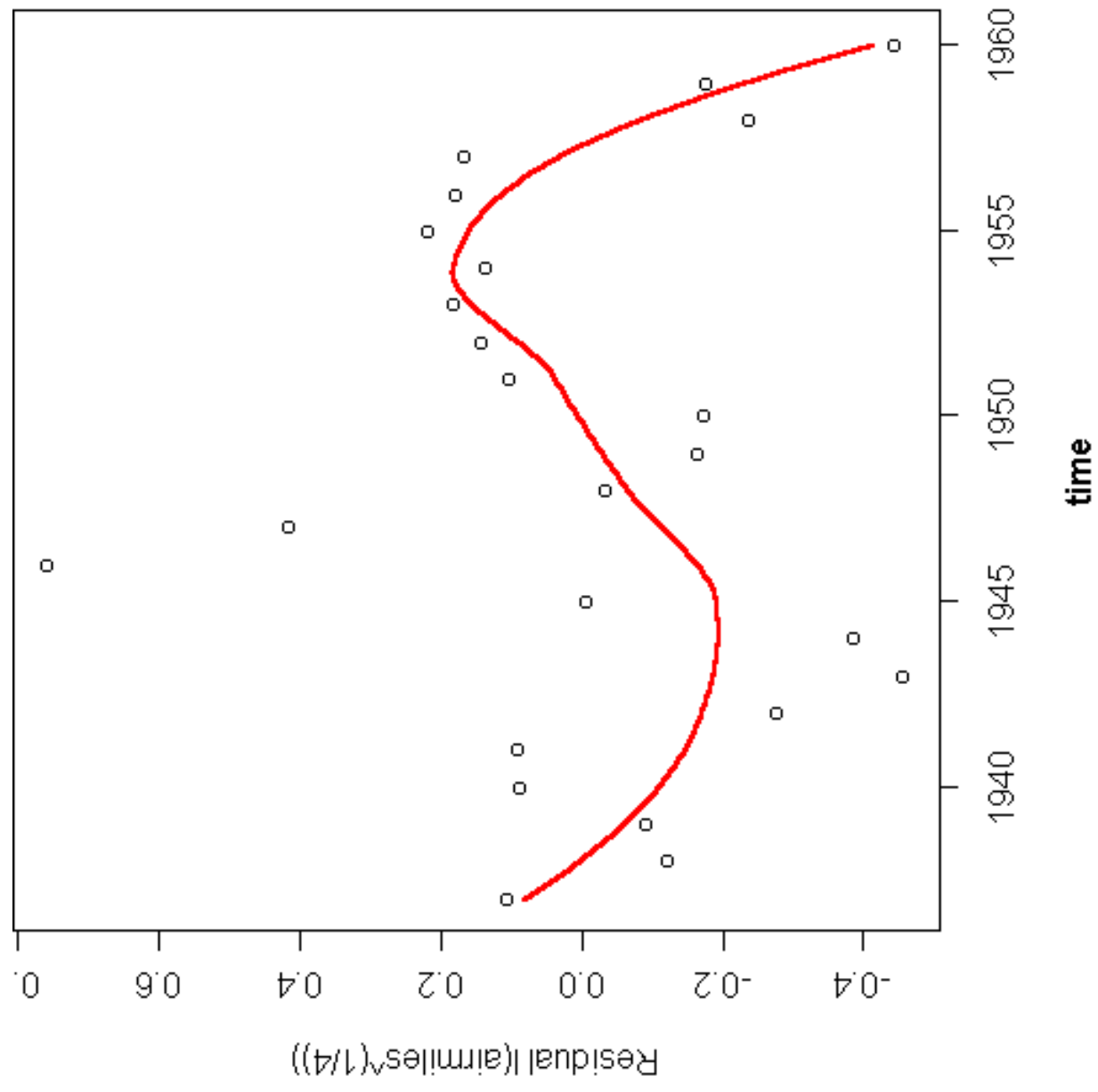


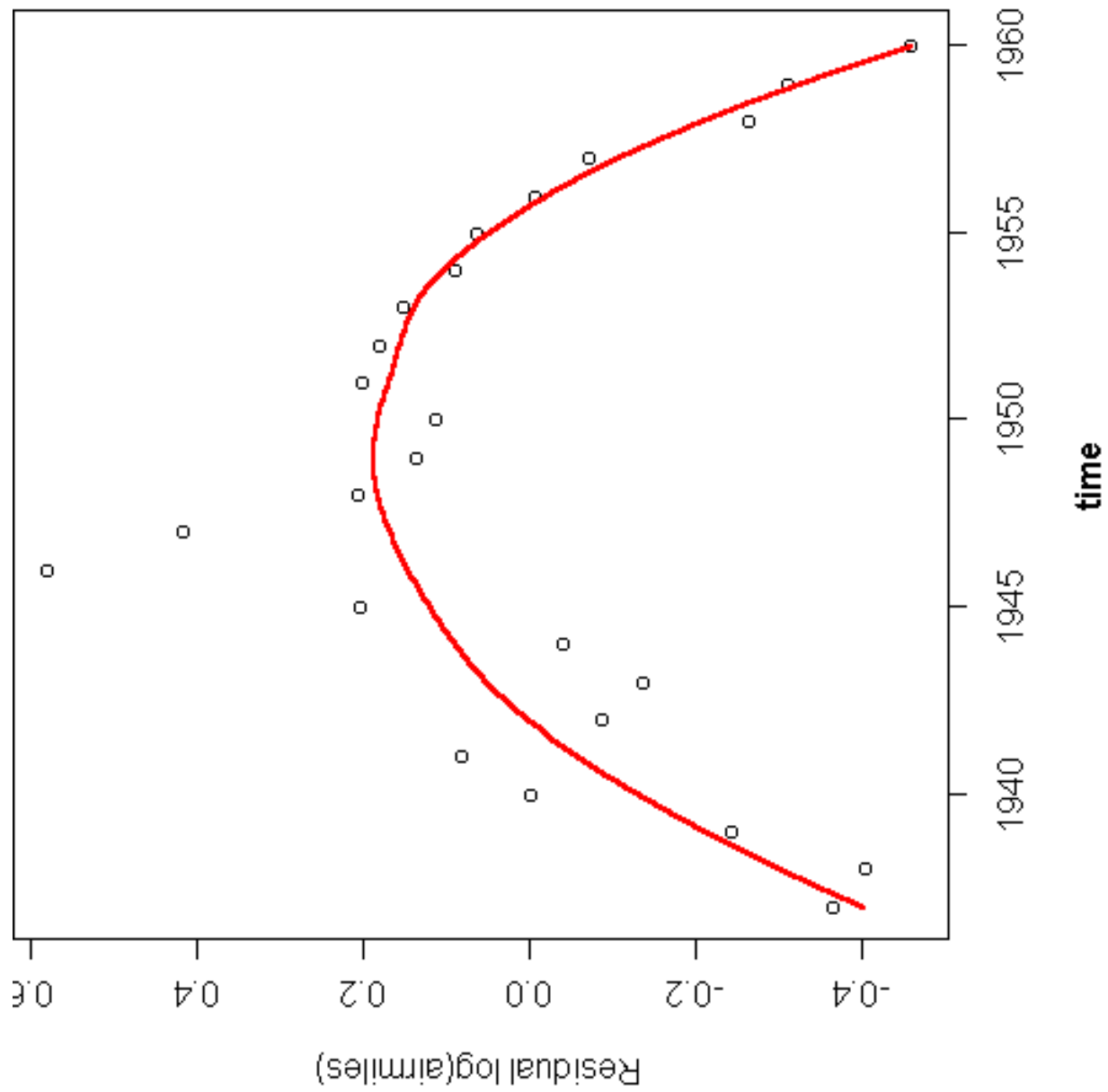


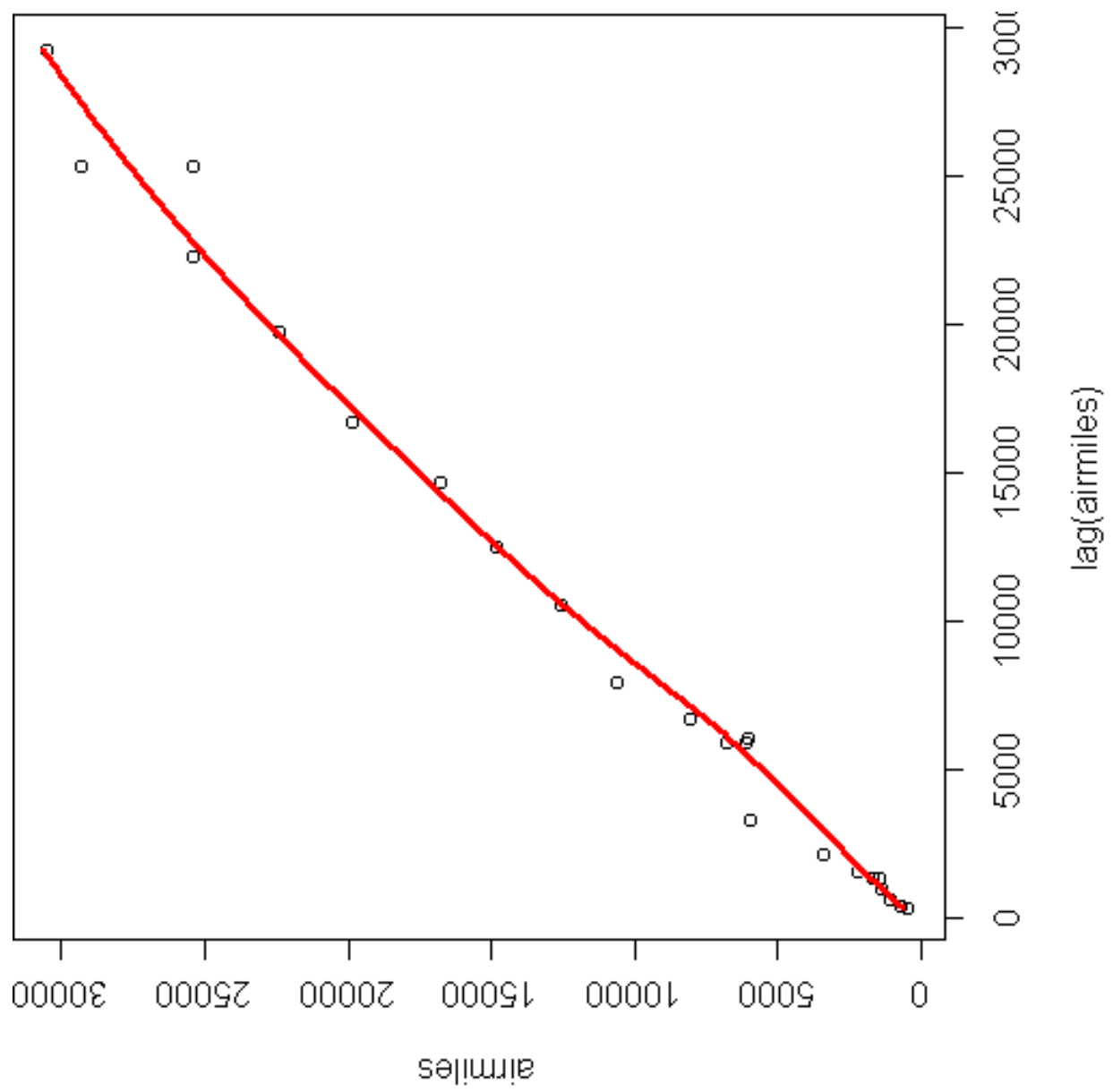


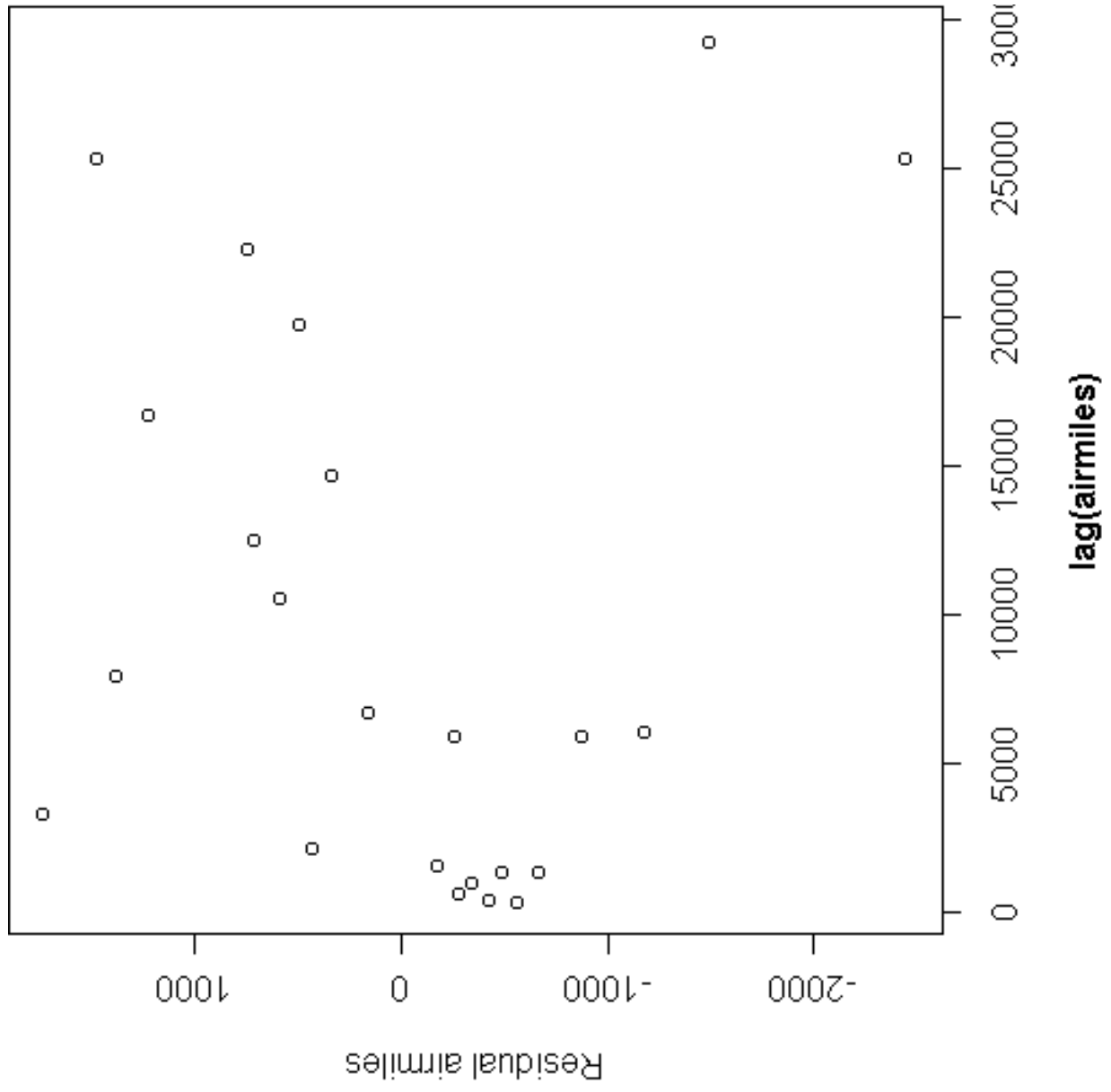


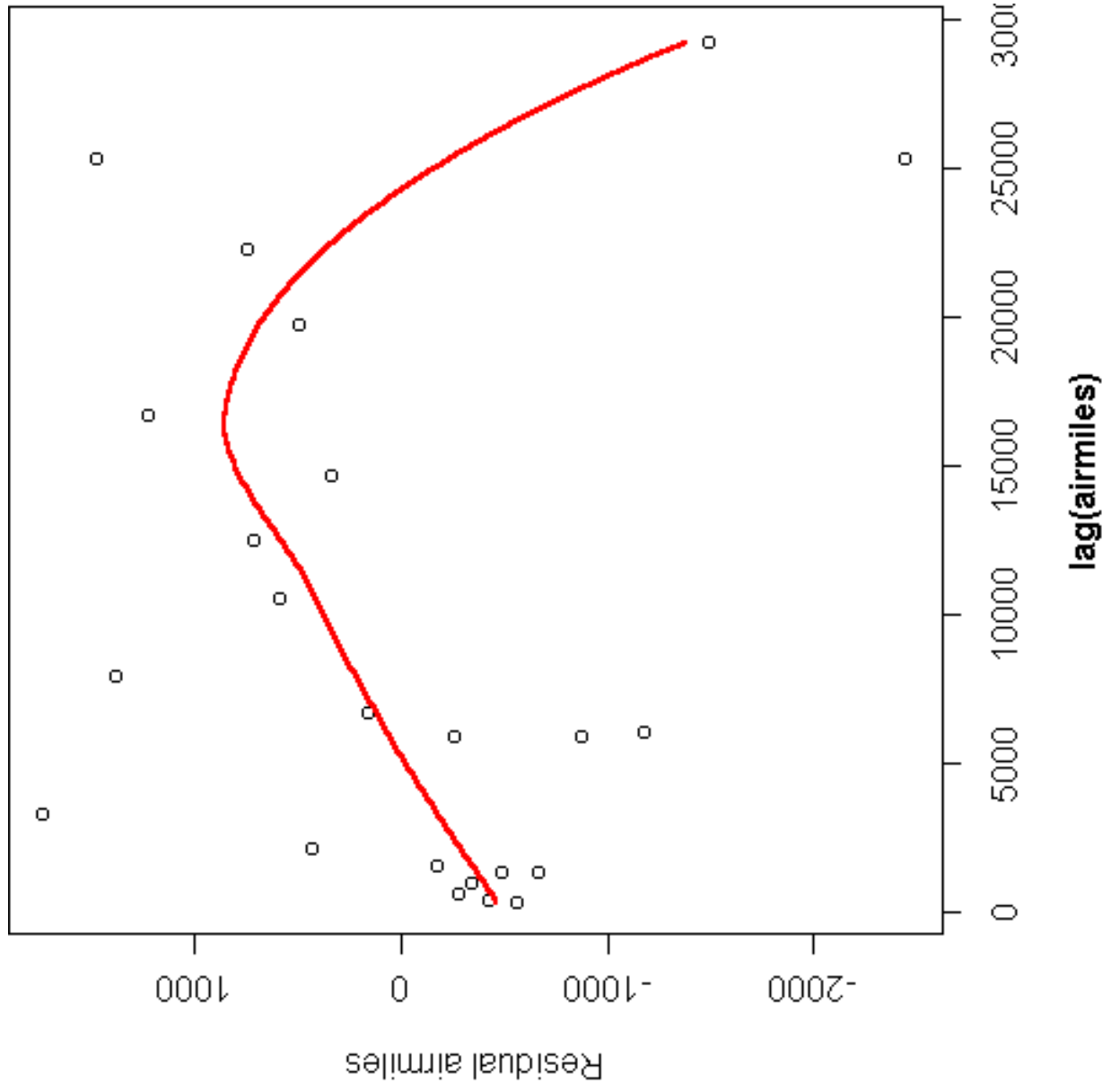


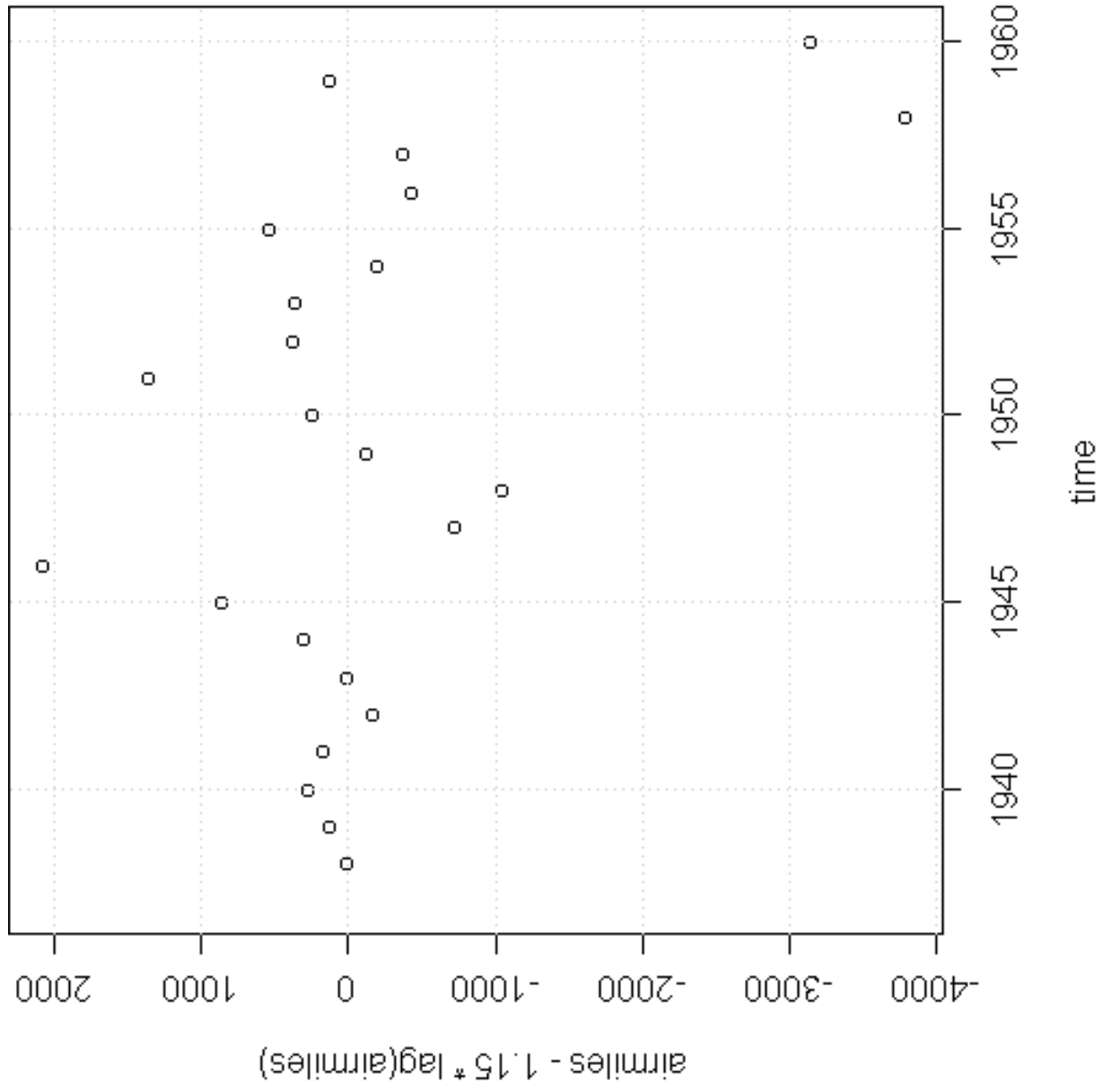


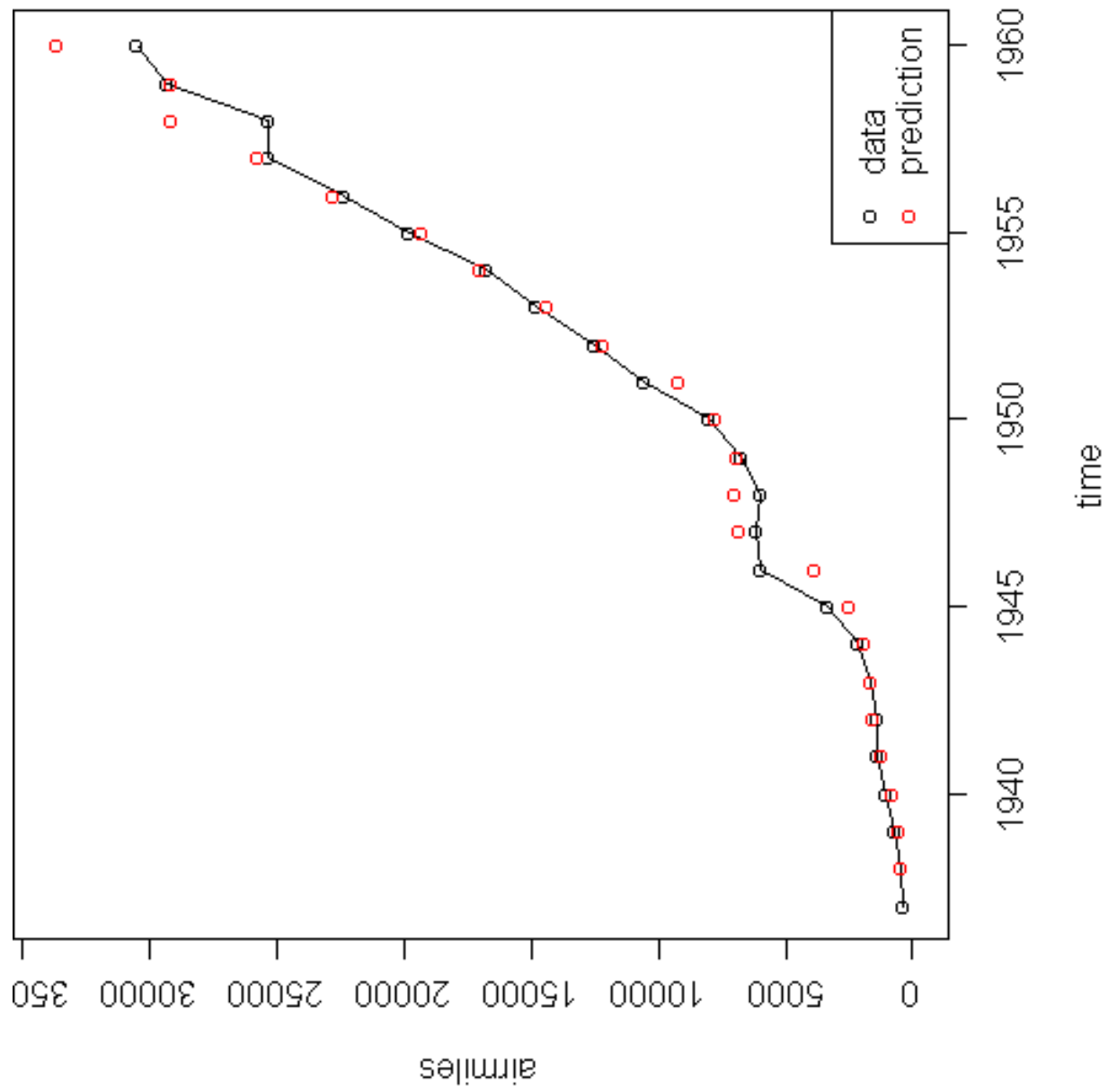












High-dimensional data

- Two basic approaches to visualization
- Many points, few dimensions:
 - Projection
 - Slicing
- Few points, many dimensions:
 - Parallel-coordinates
 - Iconic displays

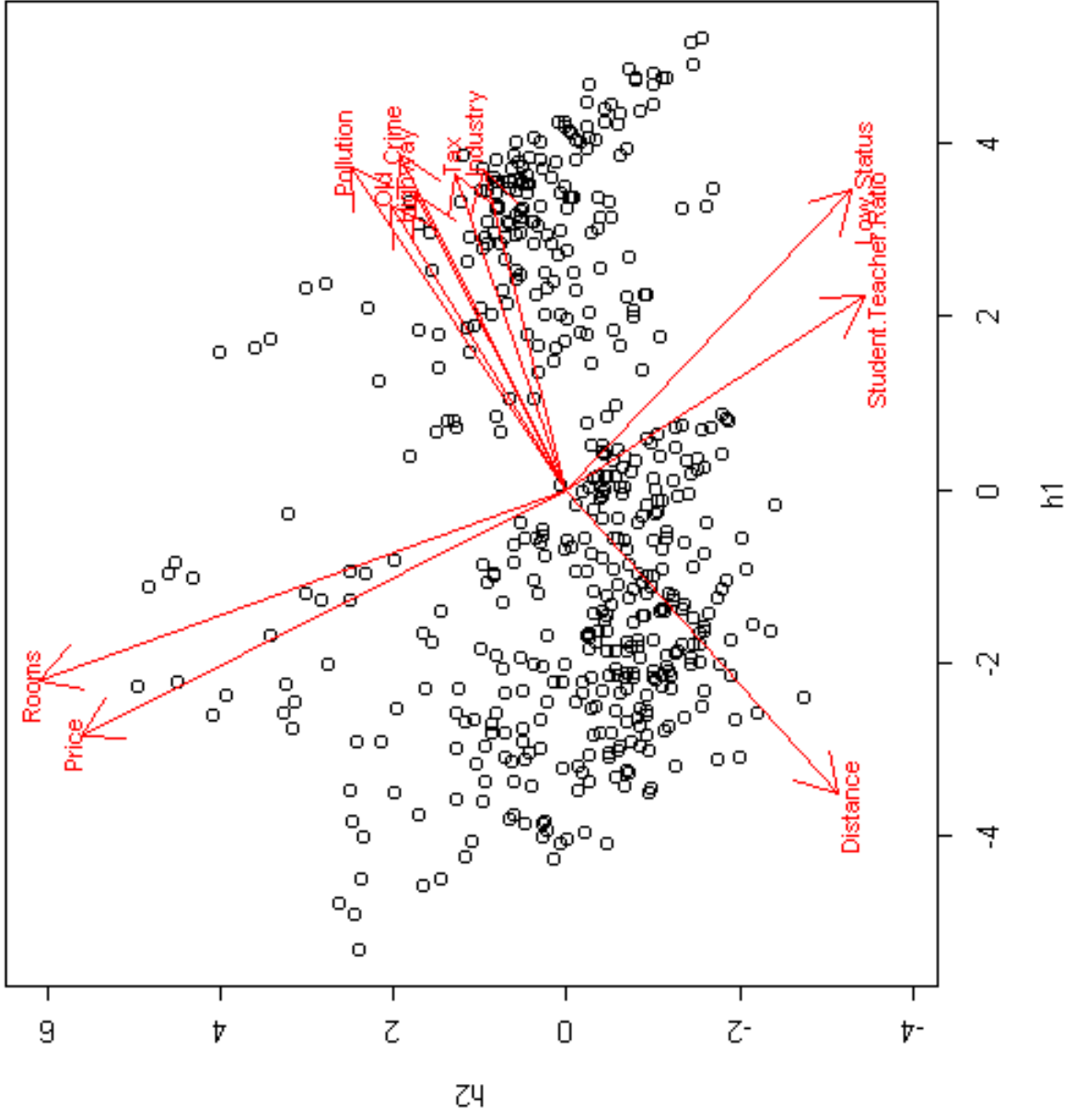
Projections

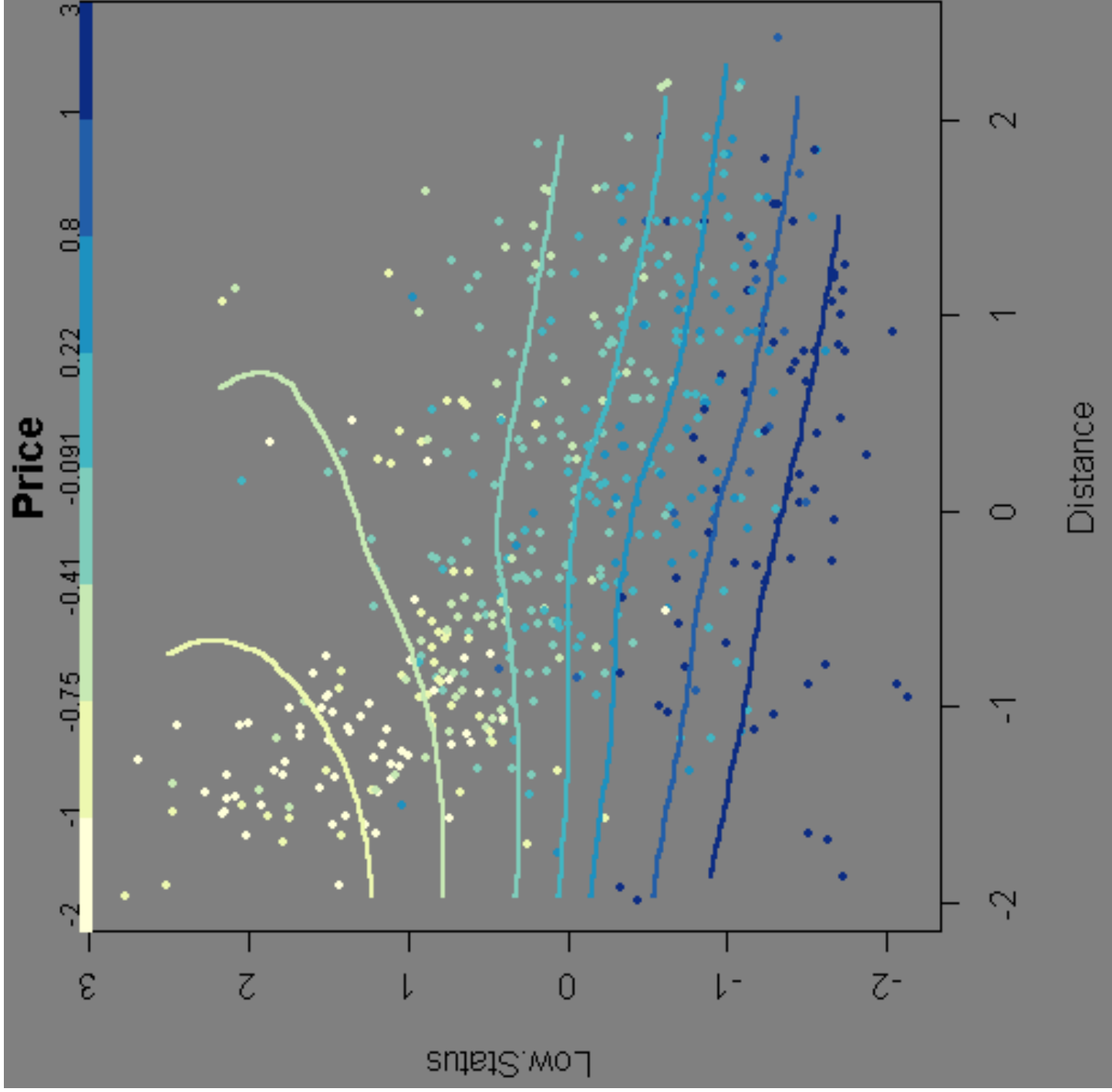
- PCA
 - Maximize the spread of the projected data
- Regression projection
 - Project only the predictors (inputs)
 - Maximize the spread of the response (output)

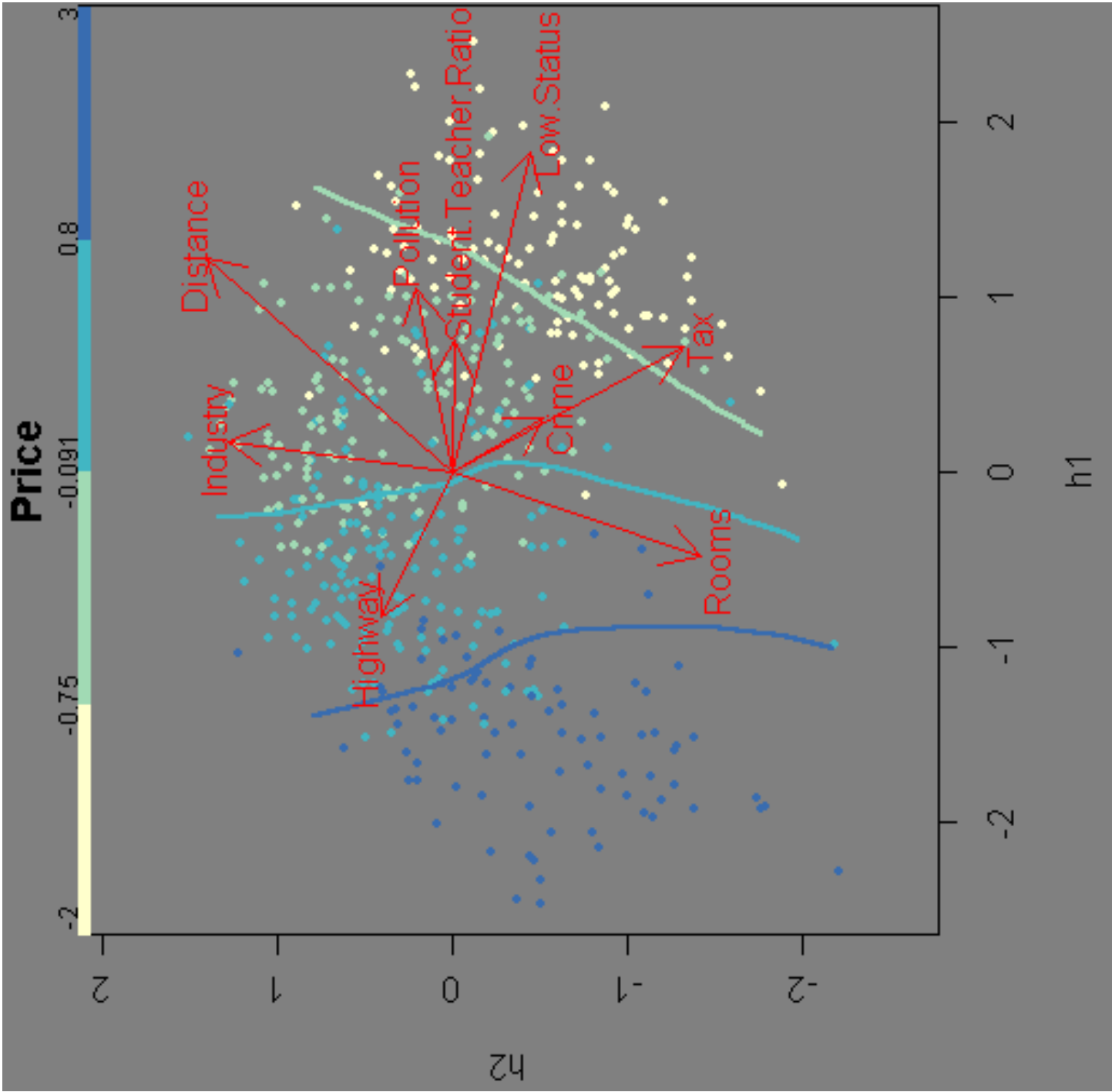
Boston Housing data

- Predict the median house value in Boston census tracts, based on crime, poverty, industry, pollution, etc.
- A regression problem with many predictors
- Is an additive model appropriate?

PCA projection





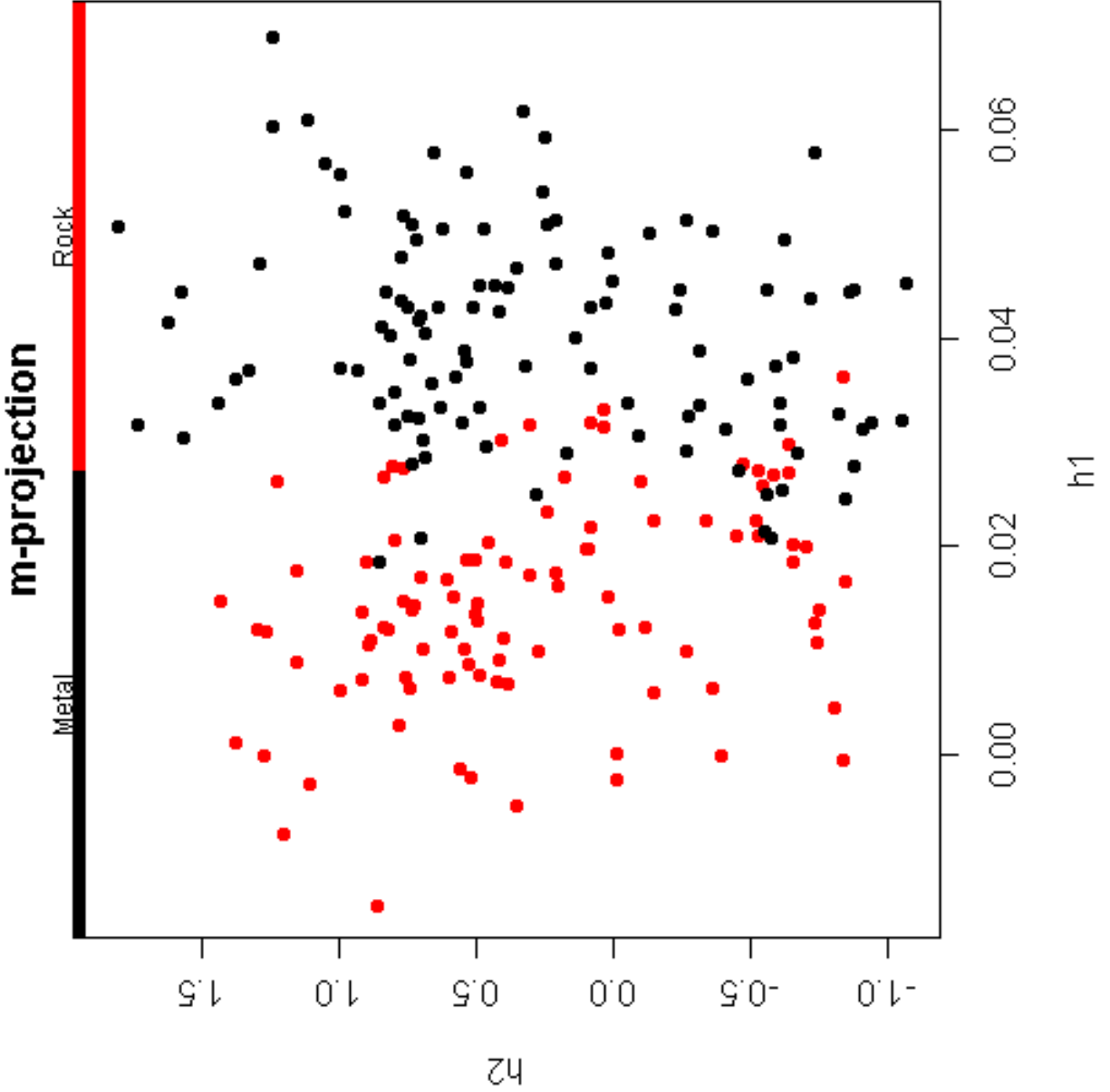


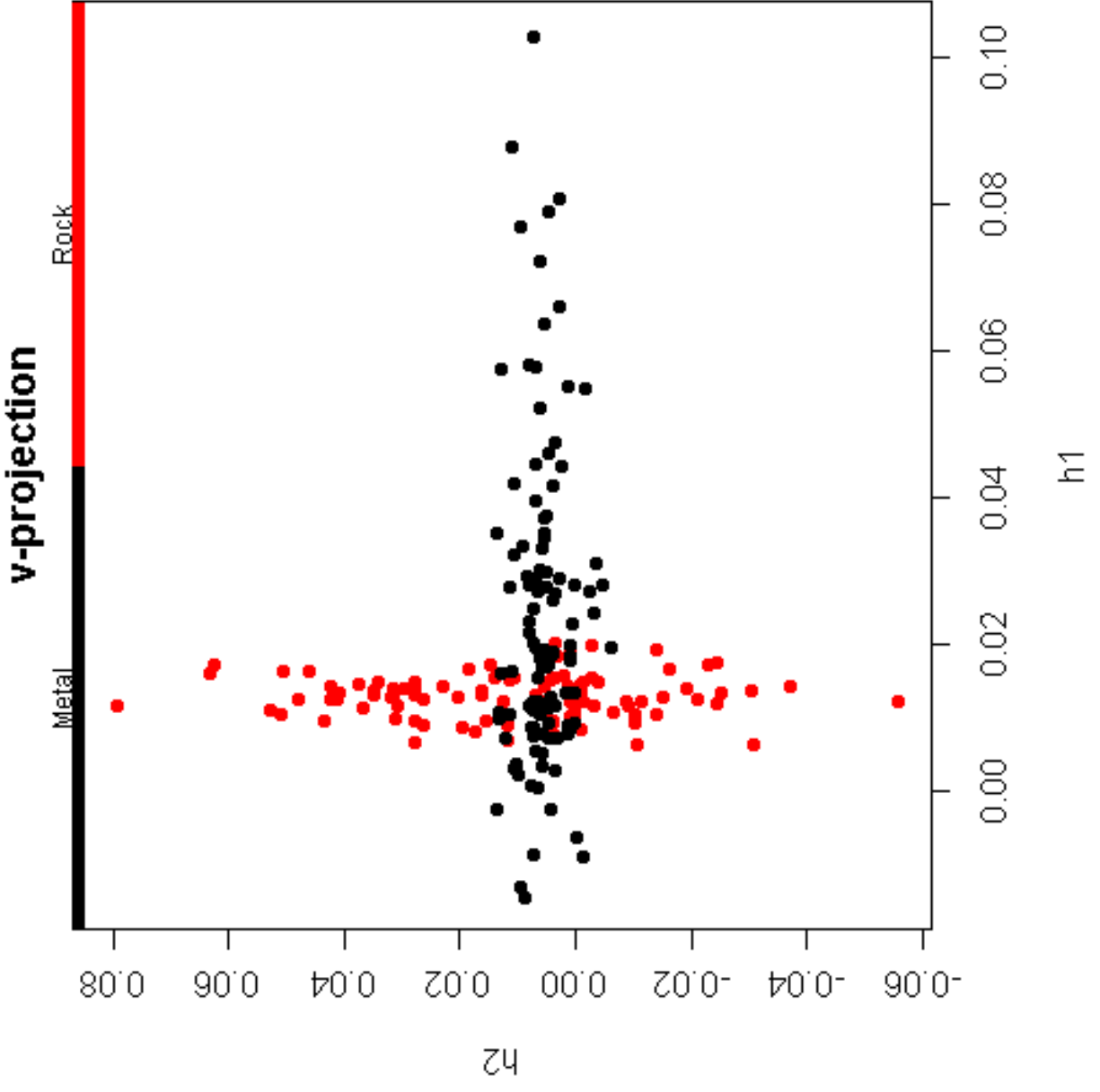
Discriminative projections

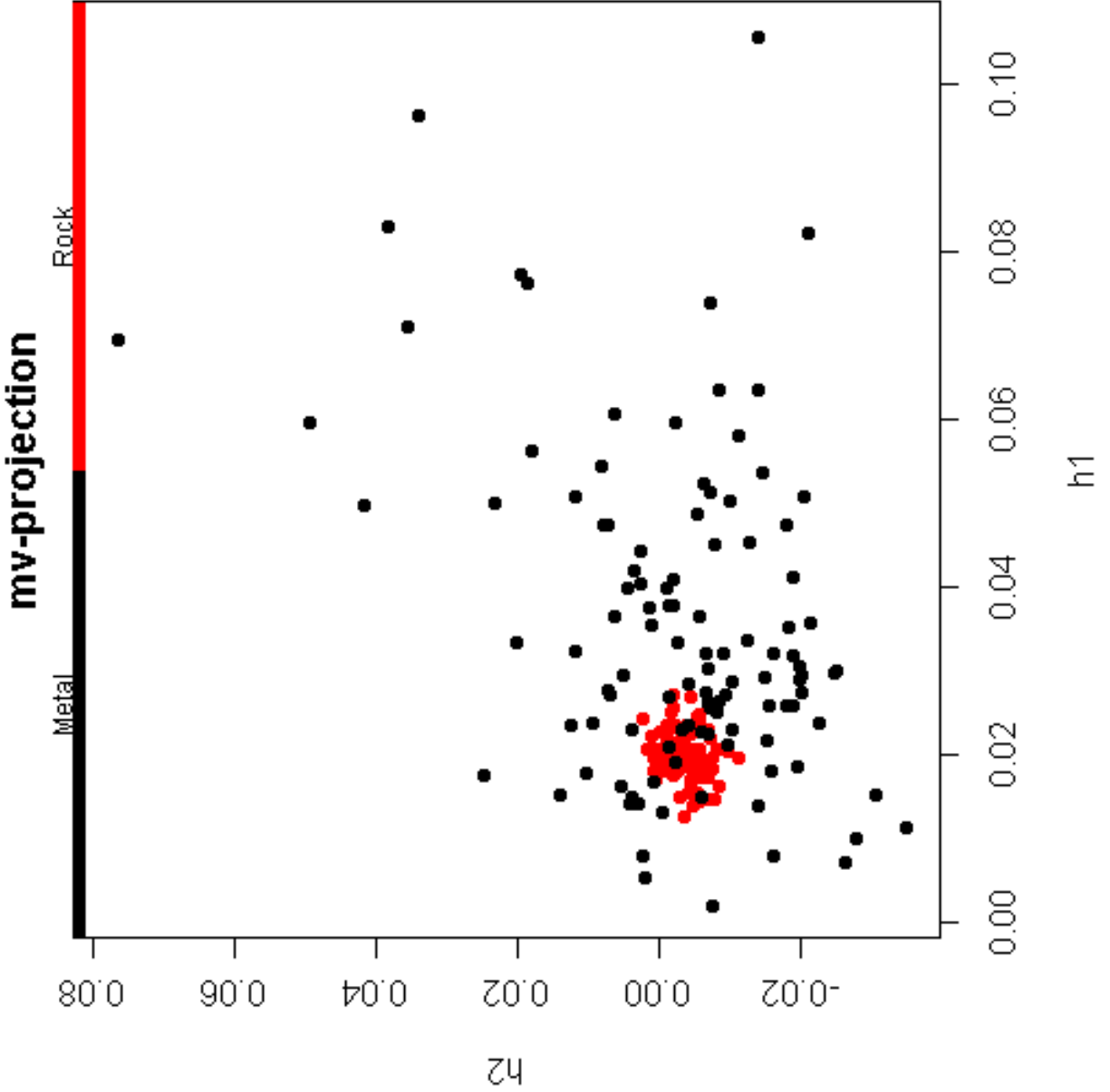
- M-projection (Fisher, LDA)
 - Tries to separate means of classes
- V-projection
 - Tries to separate variances of classes
- MV-projection
 - Maximize KL divergence between Gaussians
 - Separates means and variances

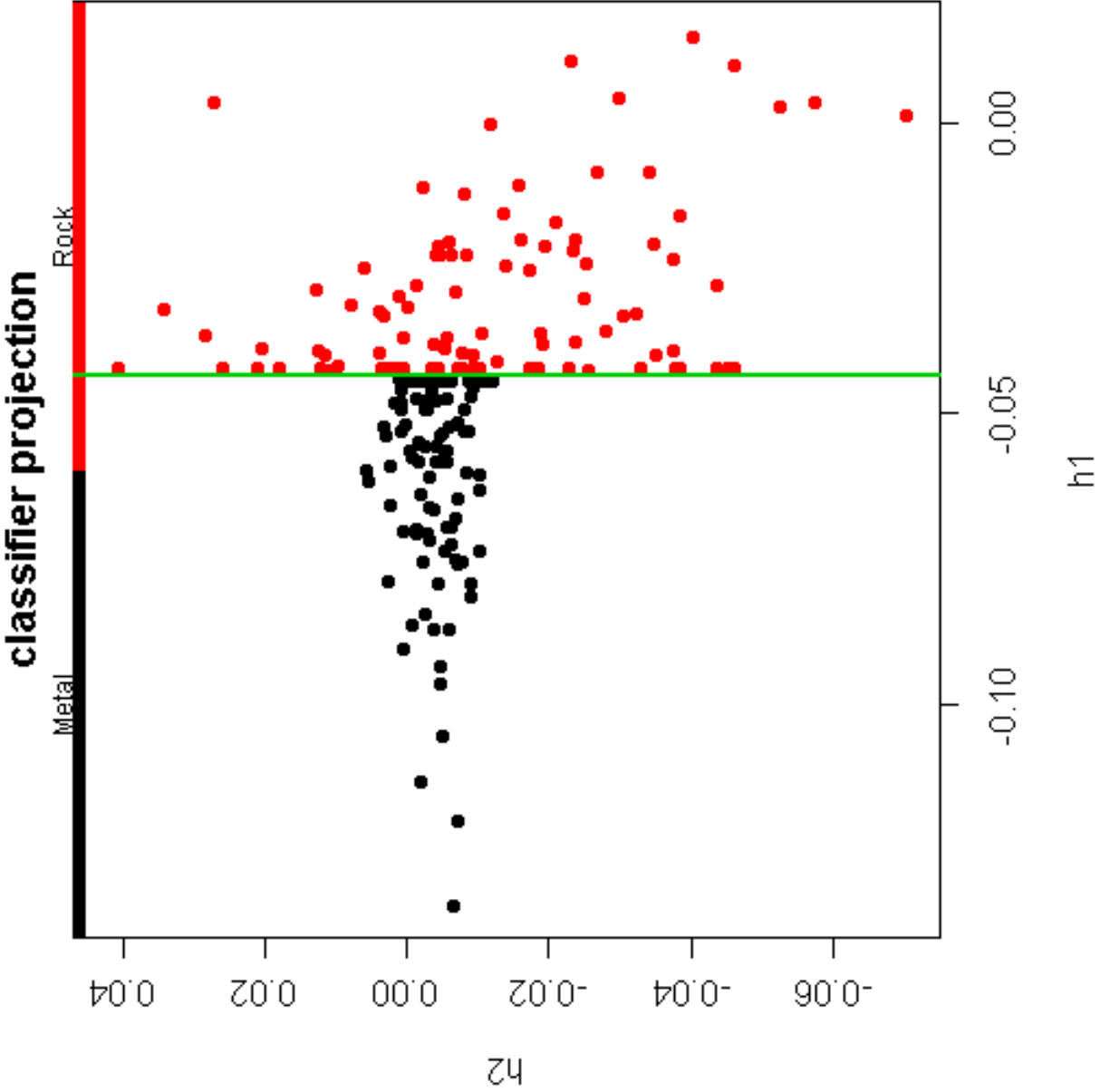
Sonar problem

- Sonar echo is represented by energy in 60 frequency bands
- Mines vs. Rocks
- Dataset is linearly separable, but 1nn consistently beats linear classifiers



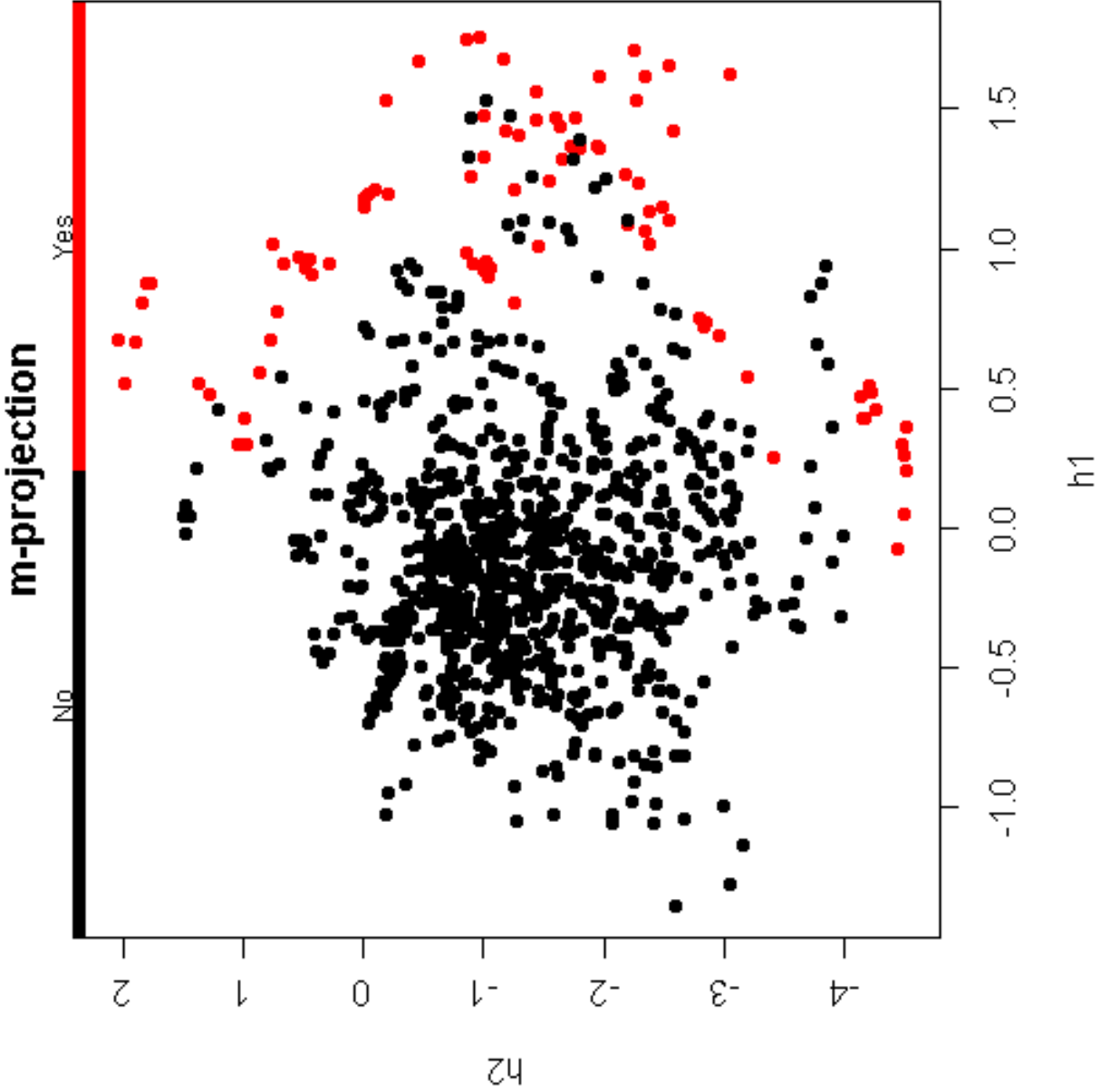


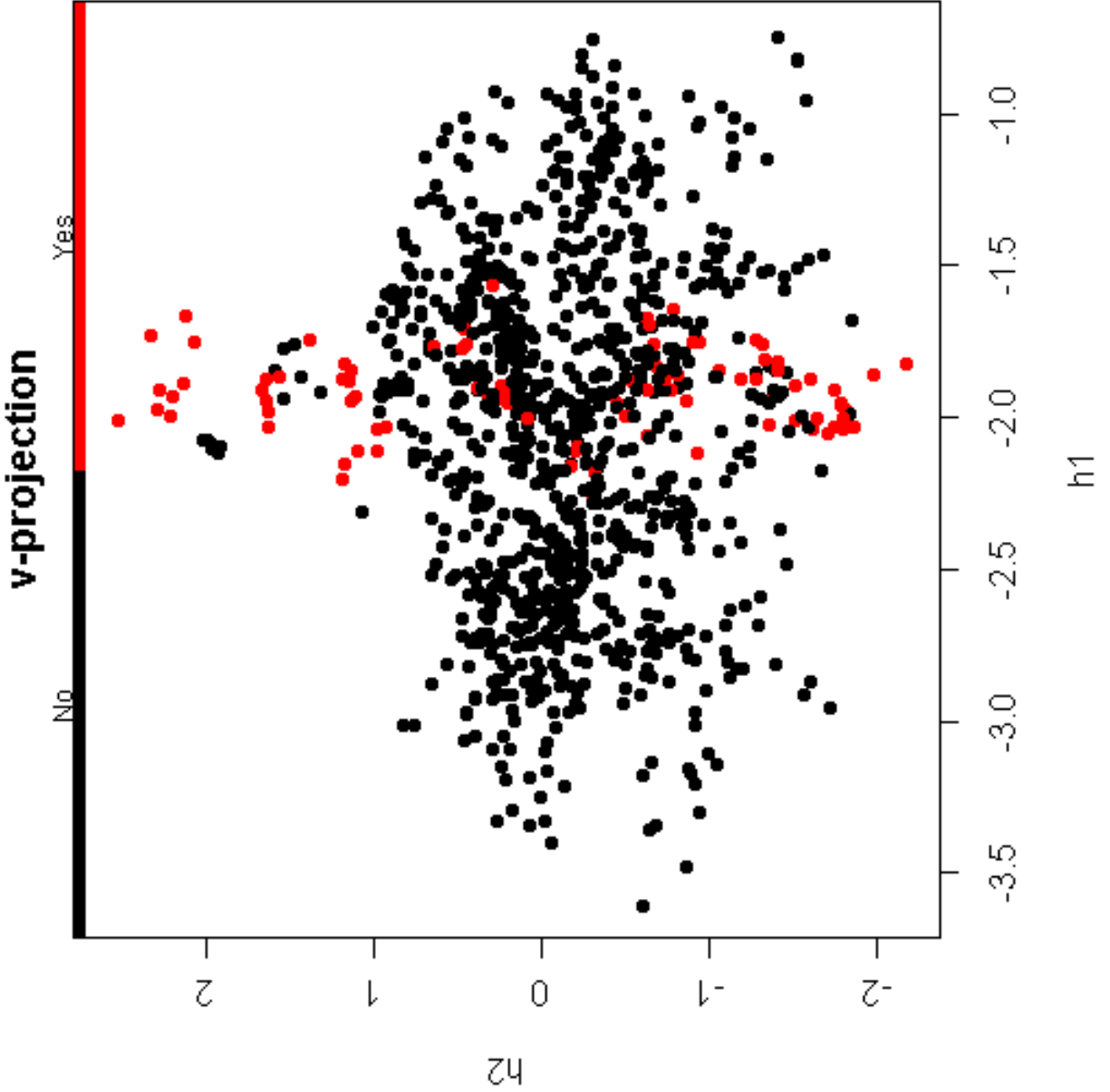


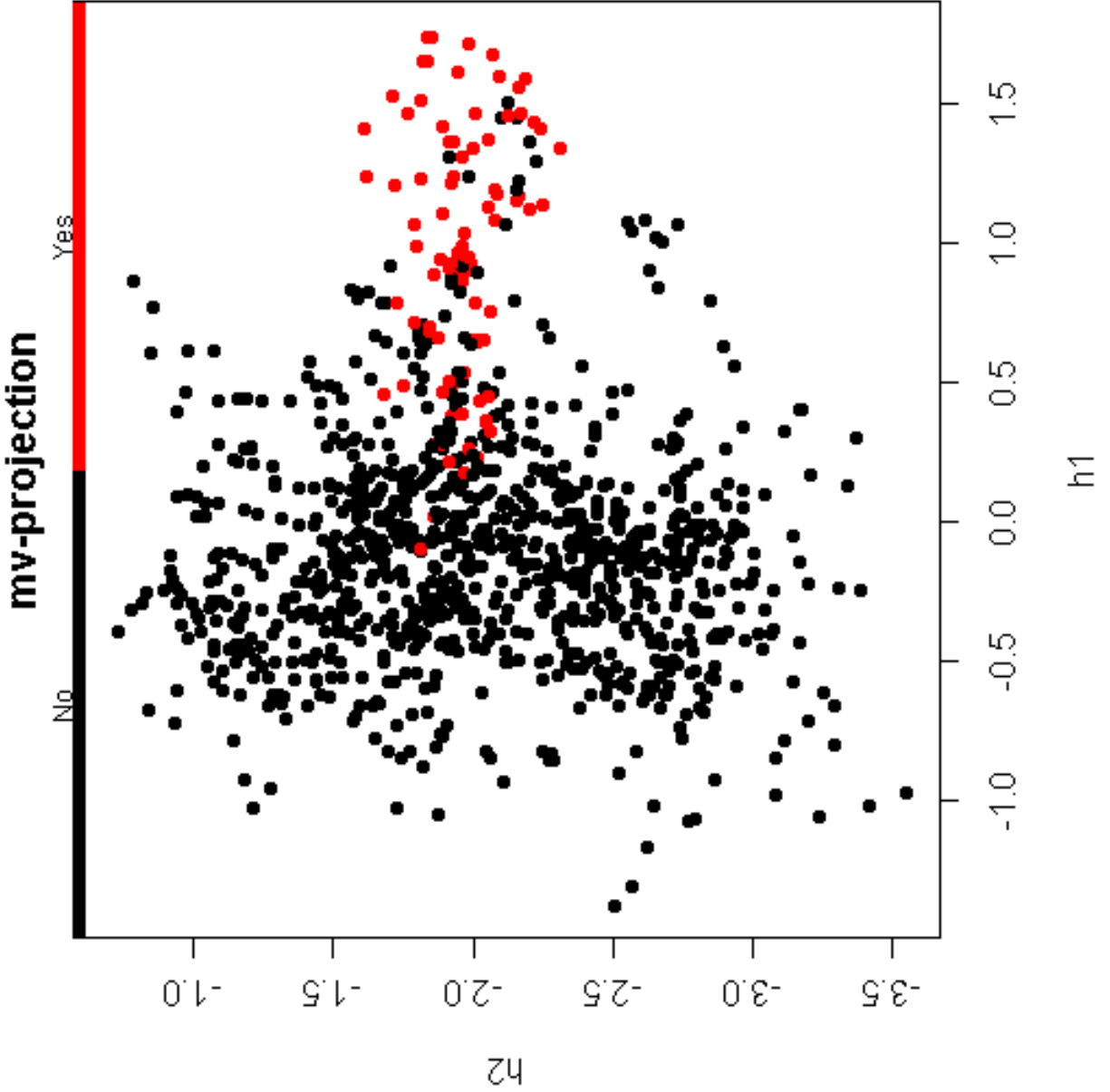


Vowels dataset

- Binary problem: “hid” vs. rest
- Knn and quadratic kernel beat linear

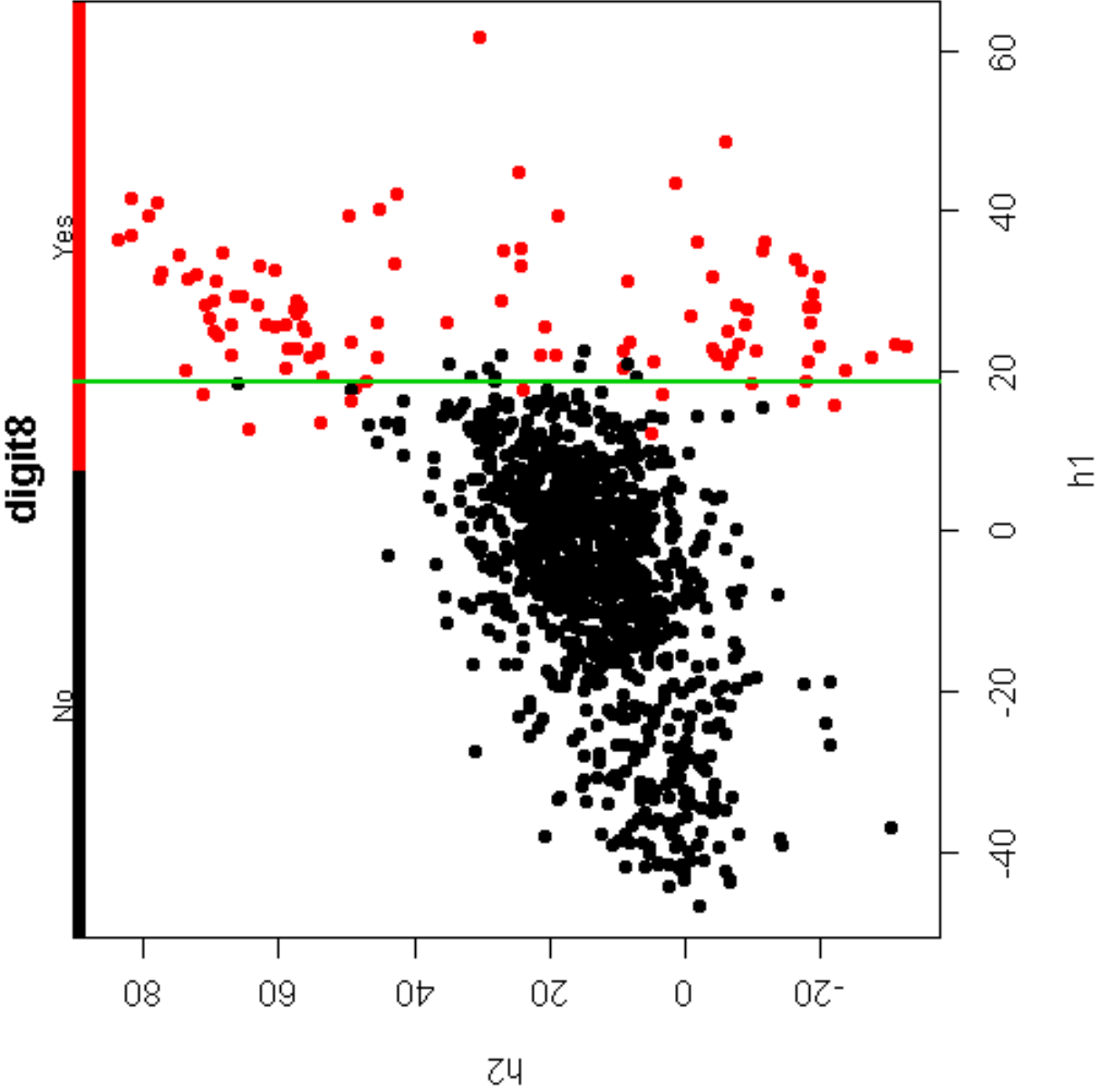


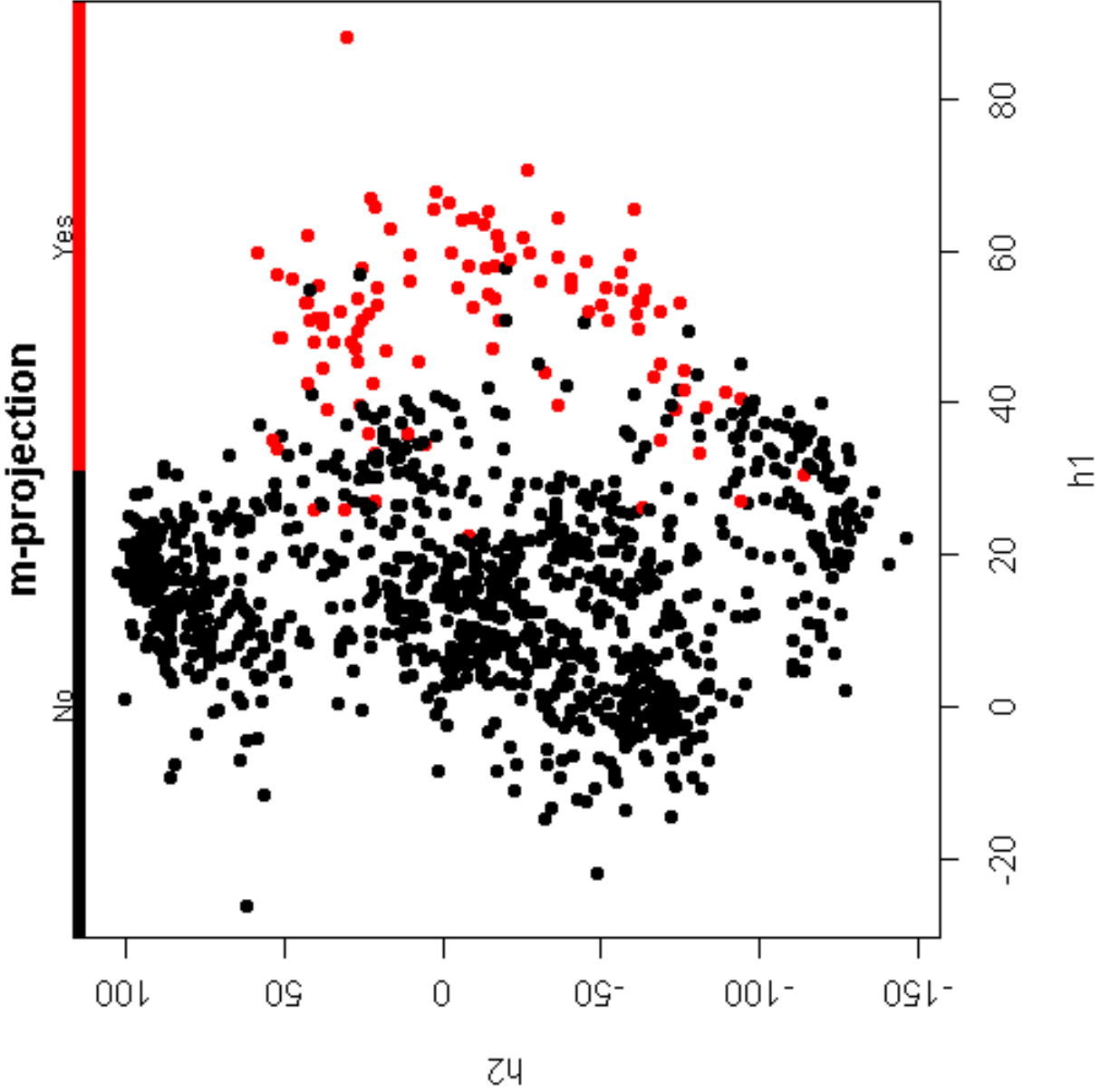


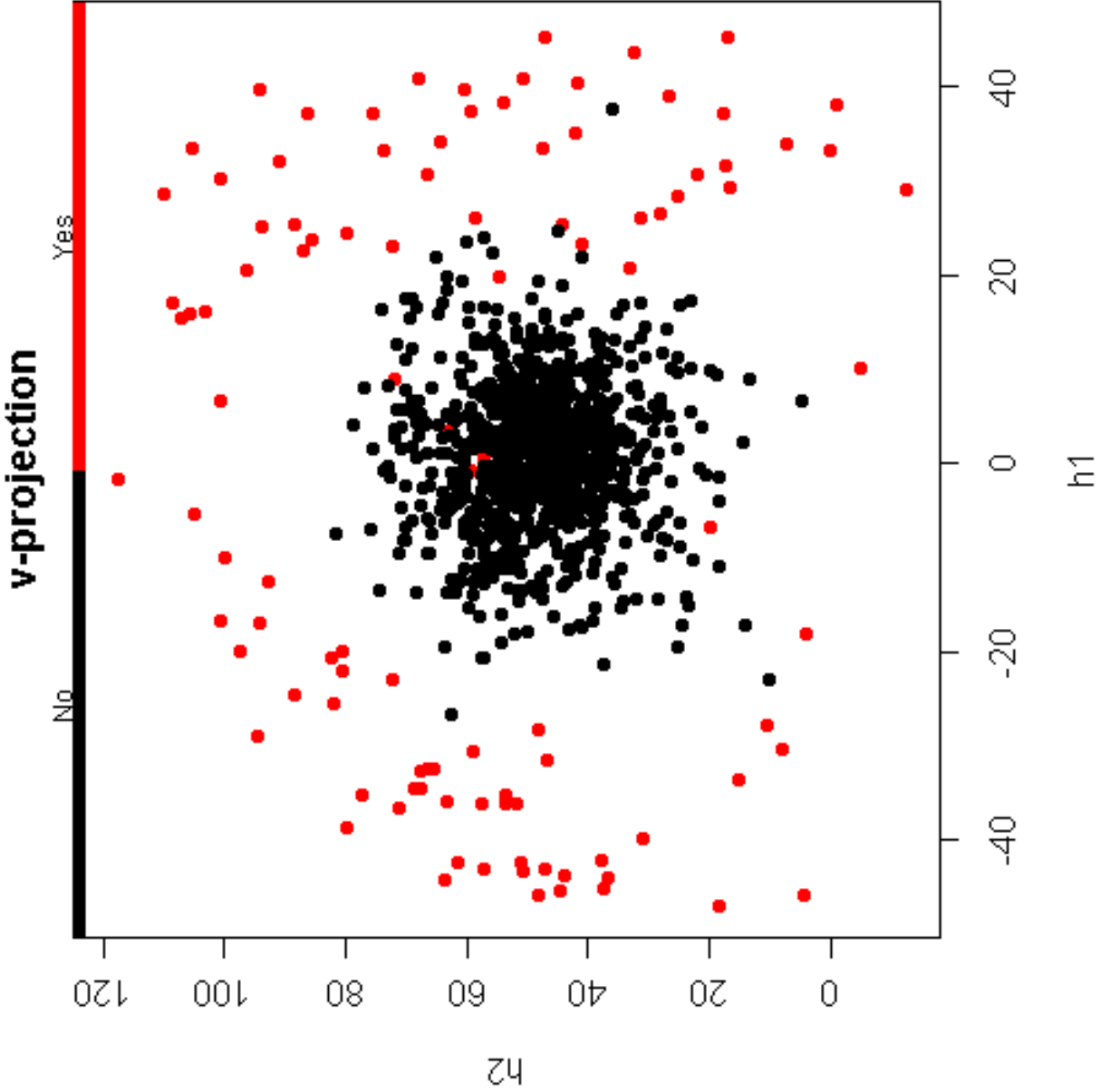


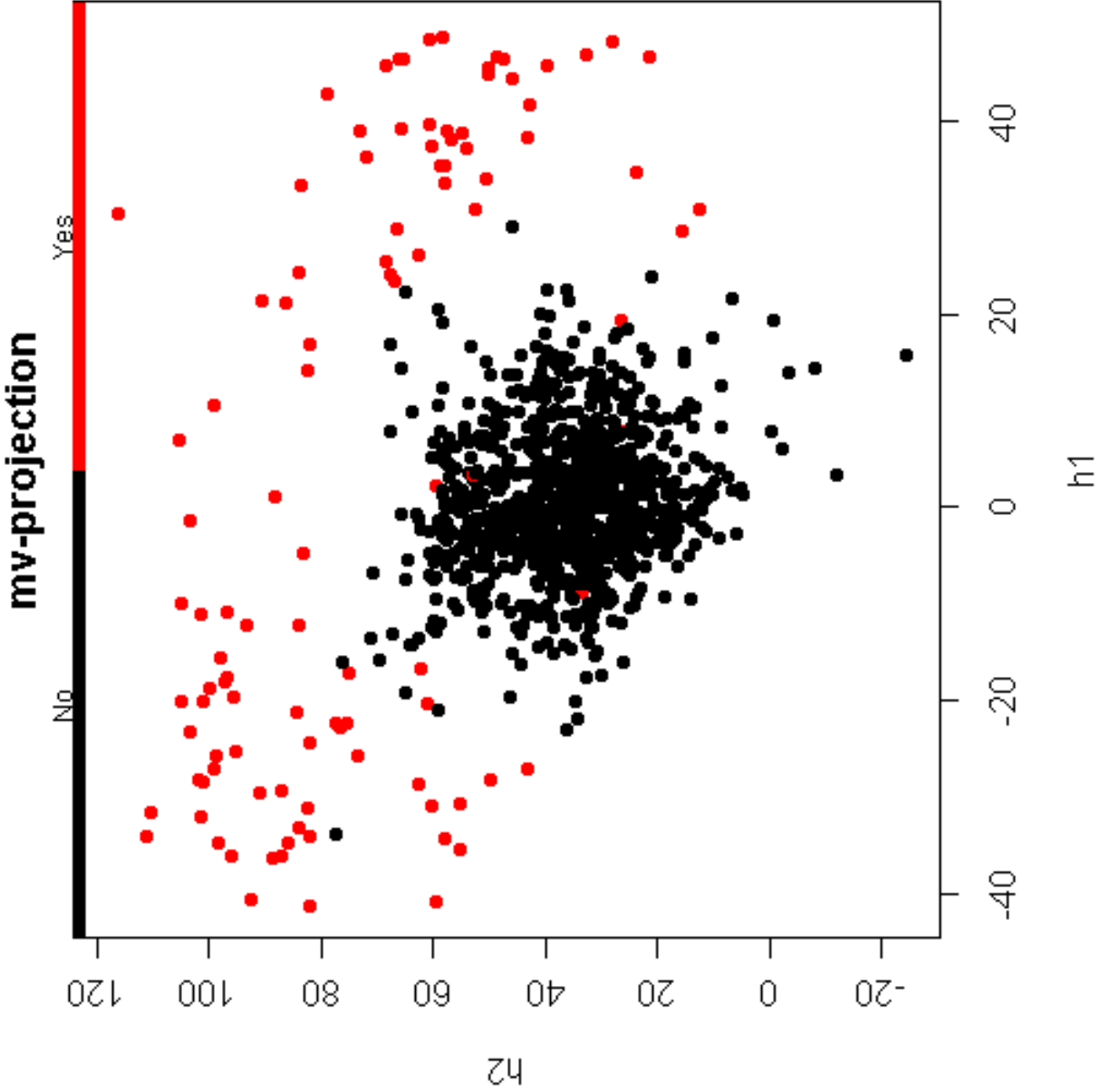
Online digit recognition

- Classify “8” vs. rest
- Knn beats quadratic kernel beats linear









Some good books

- “The Elements of Graphing Data”, William Cleveland, 2nd Ed.
- “Visualizing Data”, William Cleveland
- “The Visual Display of Quantitative Information”, Edward Tufte
- “Exploratory Data Analysis”, John Tukey

Summary

- Visualization is a simple and fast way to check model assumptions and learn about a domain
- Many opportunities still exist to design better graphs, esp. for high dimensions
- Visualization is not “art”, but a well - structured field, worthy of research attention