

ROBUST LOCATION AND SCALE ESTIMATION  
WITH CENSORED OUTCOMES

Jerome H. Friedman

Stanford University

# MACHINE LEARNING

$$y = F(\mathbf{x}, \mathbf{z})$$

$y$  = outcome variable

$\mathbf{x} = (x_1 \cdots, x_p)$  observed predictor variables

$\mathbf{z} = (z_1, z_2, \cdots)$  other variables

Goal: estimate  $E[y | \mathbf{x}]$  given data  $\{y_i, \mathbf{x}_i\}_{i=1}^N$

## STATISTICAL MODEL

$$y = f(\mathbf{x}) + s(\mathbf{x}) \cdot \epsilon$$

$f(\mathbf{x}) = E[y | \mathbf{x}]$  location function

$s(\mathbf{x}) > 0$  scale function

$\epsilon =$  random variable,  $E[\epsilon | \mathbf{x}] = 0$

Prediction:  $\hat{y} = f(\mathbf{x})$

$s(\mathbf{x}) \cdot \epsilon =$  “irreducible error” (unavoidable)

## REDUCIBLE ERROR

$$r(\mathbf{x}) = E | f(\mathbf{x}) - \hat{f}(\mathbf{x}) |$$

$f(\mathbf{x})$  = optimal location (target) function

$\hat{f}(\mathbf{x})$  = estimate based on training data & ML method

ML goal: methods to reduce  $r(\mathbf{x})$

Statistics goal: methods to estimate  $r(\mathbf{x})$

Prediction error ( $y$ ) = Reducible + Irreducible

Usually: Irreducible  $s(\mathbf{x}) \gg$  Reducible  $r(\mathbf{x})$

## USUAL ASSUMPTIONS

$s(\mathbf{x}) = s = \text{constant}$  (homoscedasticity)

$\epsilon \sim N(0, 1)$  (normality)

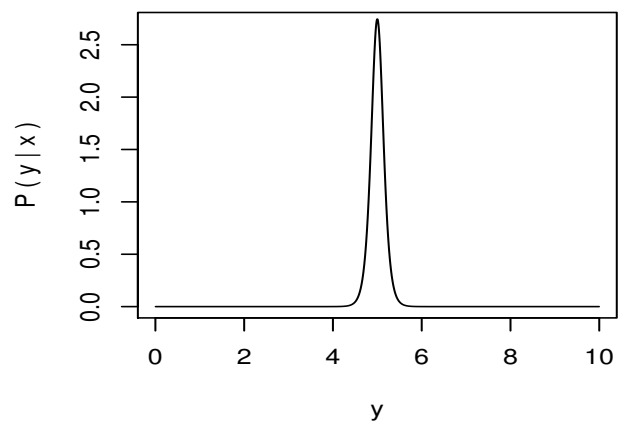
## HOMOSCEDASTICITY

$$F(\mathbf{x}, \mathbf{z}) = f(\mathbf{x}) + g(\mathbf{z}) \quad \text{additive}$$

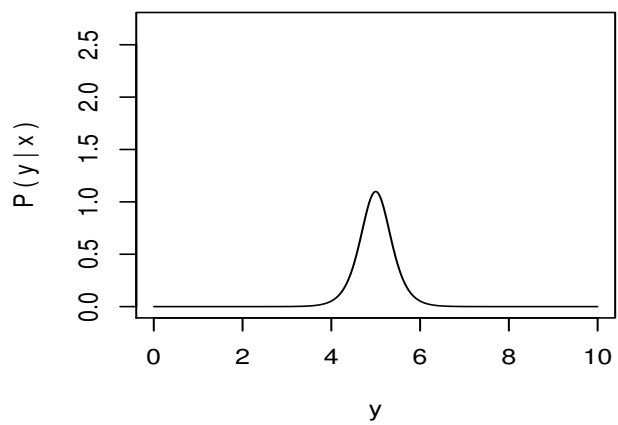
$$p(\mathbf{x}, \mathbf{z}) \implies \text{scale}[g(\mathbf{z}) \mid \mathbf{x}] = \text{constant}$$

Not very likely

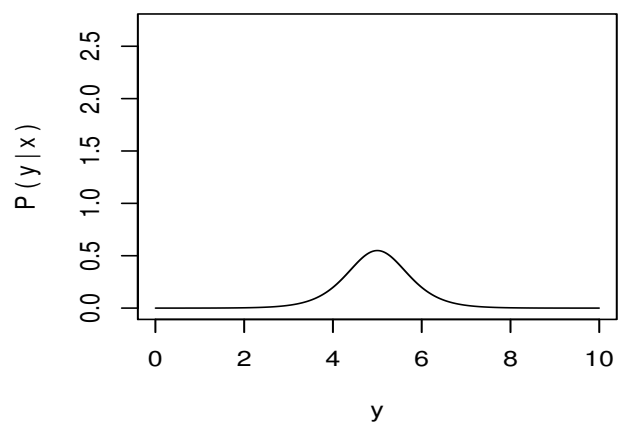
**loc = 5, scale = 0.1**



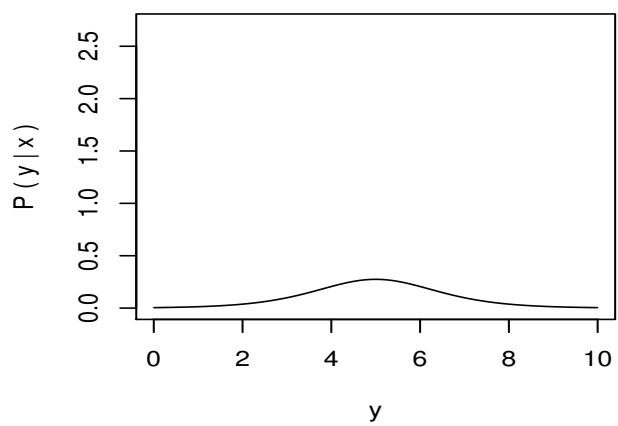
**loc = 5, scale = 0.25**



**loc = 5, scale = 0.5**



**loc = 5, scale = 1**



NORMALITY - not very likely either

Tukey:

“small residuals  $\simeq$  normal, larger have heavier tails.”

Heterodistributionality



## Heterodistributionality

Robustness:

Choose compromise  $\bar{p}(\epsilon)$

good properties for others

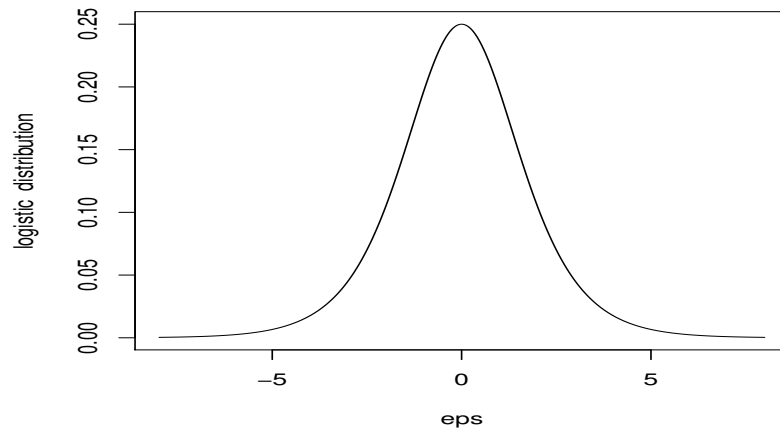
$\bar{p}(\epsilon) = \text{normal, not good!}$

## LOGISTIC DISTRIBUTION

$$\epsilon | \mathbf{x} = (y - f(\mathbf{x})) / s(\mathbf{x})$$

$$\bar{p}(\epsilon) = \frac{e^{-\epsilon}}{s(1+e^{-\epsilon})^2}$$

small  $|\epsilon| \sim$  normal, large  $|\epsilon| \sim$  exponential



Prediction:  $\hat{y} = \hat{f}(\mathbf{x})$

$$\hat{f}(\mathbf{x}) = \arg \min_{f \in F} \sum_{i=1}^N [\varepsilon_i + 2 \log(1 + e^{-\varepsilon_i})]$$

$$\varepsilon_i = (y_i - f(\mathbf{x}_i)) / s(\mathbf{x}_i)$$

minimized at  $f(\mathbf{x}_i) = y_i$  indep  $s(\mathbf{x}_i)$

$1/s(\mathbf{x}_i) \sim$  “weight” for obs  $i$

controls relative influence of  $i$  to fit

Using incorrect  $s(\mathbf{x})$  to estimate  $f(\mathbf{x})$

increases variance, not bias

assume  $s(\mathbf{x}) = \text{constant}$  usually not too bad

## ESTIMATE $\hat{s}(\mathbf{x})$

(1) Improve  $\hat{f}(\mathbf{x})$  in high variance settings.

(2) Important inferential statistic:

(a) prediction interval  $\sim$  accuracy of  $\hat{y}$ -prediction:

$$\text{logistic: } IQR[y | f(\mathbf{x})] = 2 s(\mathbf{x}) / \log(3)$$

(b) can affect decision

(3) Crucial with censoring

CENSORING ( $y$ -value partially known)

Data:  $\{y_i, \mathbf{x}_i\}_1^N \rightarrow \{a_i, b_i, \mathbf{x}_i\}_1^N$

$$a_i \leq y_i \leq b_i$$

$a_i = b_i = y_i \Rightarrow y$ -value known

$a_i = -\infty \Rightarrow$  censored below  $b_i$

$b_i = \infty \Rightarrow$  censored above  $a_i$

Otherwise: interval censored  $[a_i, b_i]$

## Special Case

$\{a_i, b_i\} \rightarrow K$  disjoint intervals (bins):

$K = 2 \Rightarrow$  usual binary logistic regression

$K > 2 \Rightarrow$  ordered multiclass logistic regression

## LIKELIHOOD

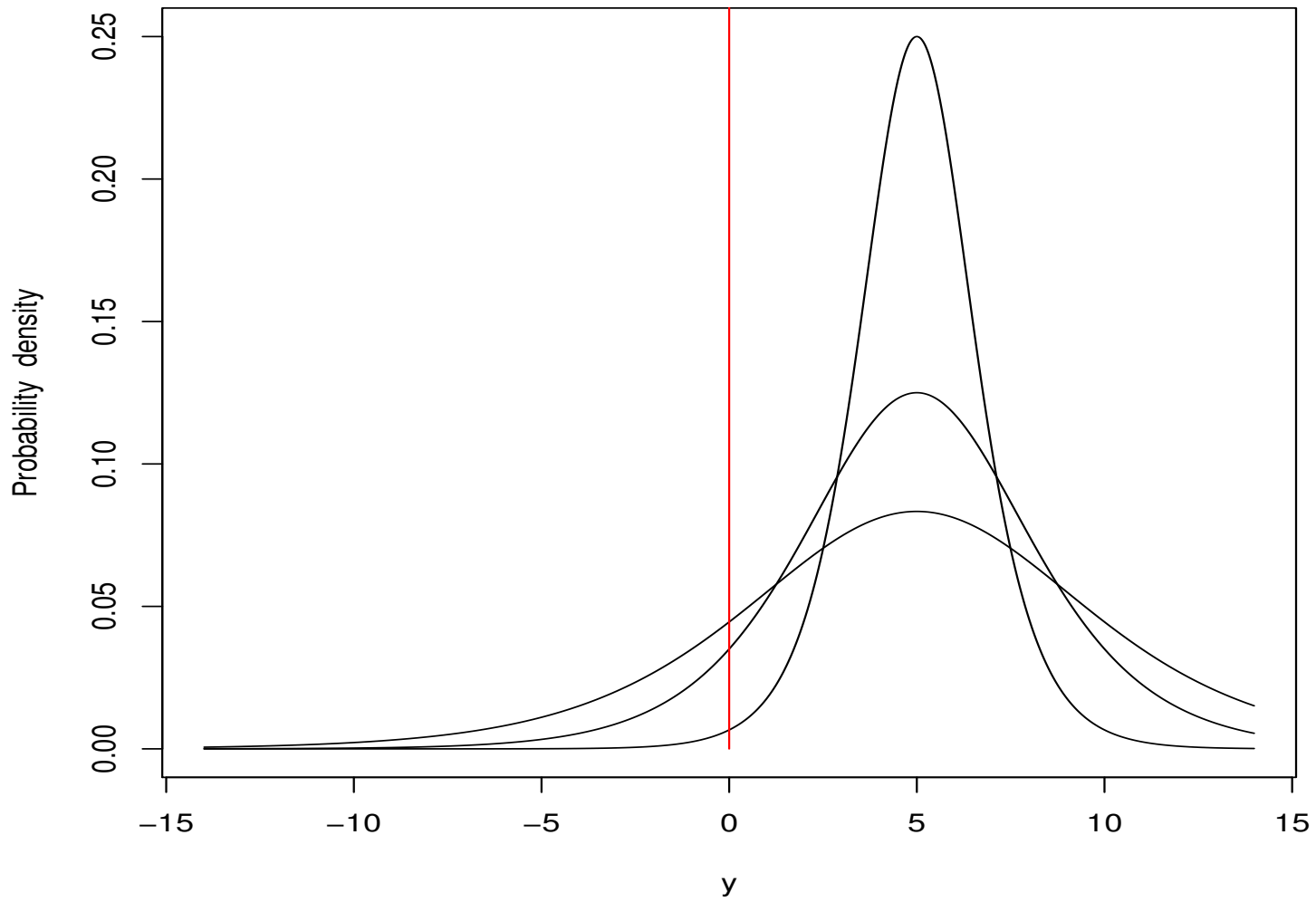
$$\Pr(a \leq y \leq b) = \frac{1}{1+e^{-(b-f)/s}} - \frac{1}{1+e^{-(a-f)/s}}$$

Depends strongly on *both*  $f$  and  $s$

Need to estimate *both*  $f(\mathbf{x})$  and  $s(\mathbf{x})$



**Logistic distribution:  $f = 5$**



## EXERCISE

$$[\hat{f}(\mathbf{x}), \hat{s}(\mathbf{x})] = \arg \min_{(f,s) \in F} \sum_{i=1}^N L[a_i, b_i, f(\mathbf{x}_i), s(\mathbf{x}_i)]$$

$$L(a, b, f, s) = -\log \left( \frac{1}{1+e^{-(b-f)/s}} - \frac{1}{1+e^{-(a-f)/s}} \right)$$

## PROBLEM

$L(a, b, f, s)$  NOT convex in  $s$

IS convex in  $t = 1/s \Rightarrow$  solve for  $t$

Constraint  $t > 0 \Rightarrow$  solve for  $\log(t) = -\log(s)$

## GRADIENT BOOSTED TREE ENSEMBLES

Ann. Statist, **29**. 1189 – 1232 (2001)

$$\hat{f}(\mathbf{x}) = \sum_{k=1}^{K_f} T_k^{(f)}(\mathbf{x})$$

$$\log(\hat{s}(\mathbf{x})) = \sum_{k=1}^{K_s} T_k^{(s)}(\mathbf{x})$$

$$T_k(\mathbf{x}) = \text{CART-tree}(\mathbf{x})$$

## ITERATIVE GRADIENT BOOSTING

Start:  $\hat{s}(\mathbf{x}) = \text{constant}$

Loop {

$$\hat{f}(\mathbf{x}) = \text{tree-boost } f(\mathbf{x}) \text{ given } \hat{s}(\mathbf{x})$$

$$\log(\hat{s}(\mathbf{x})) = \text{tree-boost } \log(s(\mathbf{x})) \text{ given } \hat{f}(\mathbf{x})$$

}

Until no change

## DIAGNOSTICS

$$(1) \textit{median} [y | f(\mathbf{x})] = f(\mathbf{x})$$

$$(2) \textit{median} [|y - f(\mathbf{x})| | s(\mathbf{x})] = s(\mathbf{x}) \cdot \log(3)$$

$$(3) \# (y_i \in [u, v] | f_i \in [g, h]) =$$

$$\sum_{f_i \in [g, h]} \left( \frac{1}{1 + e^{-(v - f_i)/s_i}} - \frac{1}{1 + e^{-(u - f_i)/s_i}} \right)$$

$$(f_i = \hat{f}(\mathbf{x}_i), \quad s_i = \hat{s}(\mathbf{x}_i))$$

## California Housing Price Data (STATLIB Repository)

$N = 20460$  CA neighborhoods (1990 census block groups)

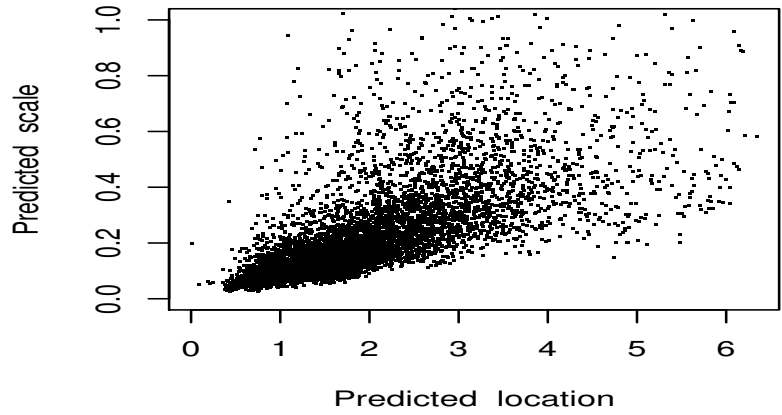
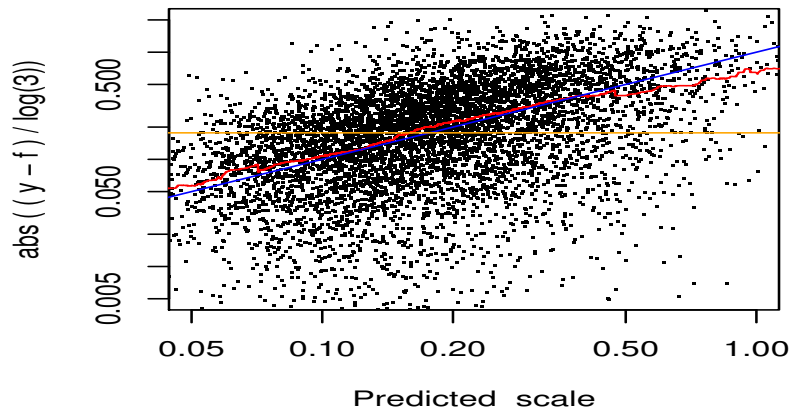
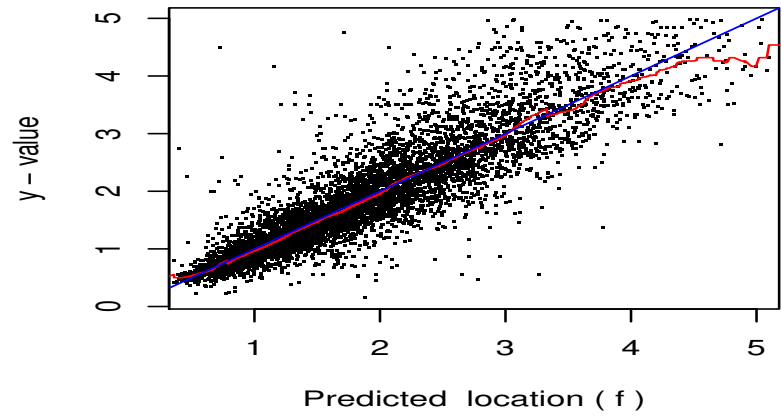
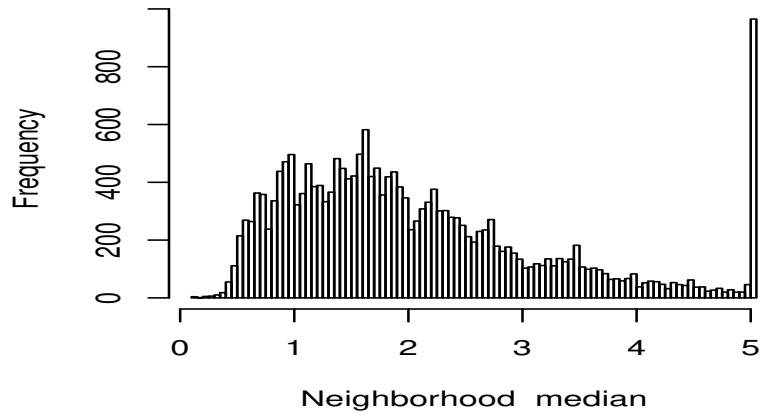
$y =$  Median House Value

$\mathbf{x} =$  (Median Income, Housing Median Age,

Ave No Rooms, Ave No Bedrooms,

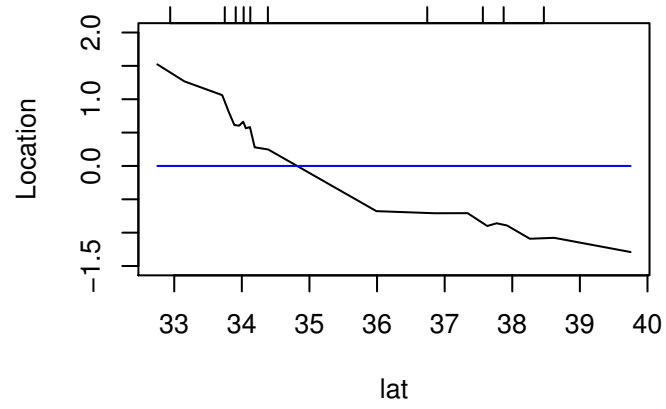
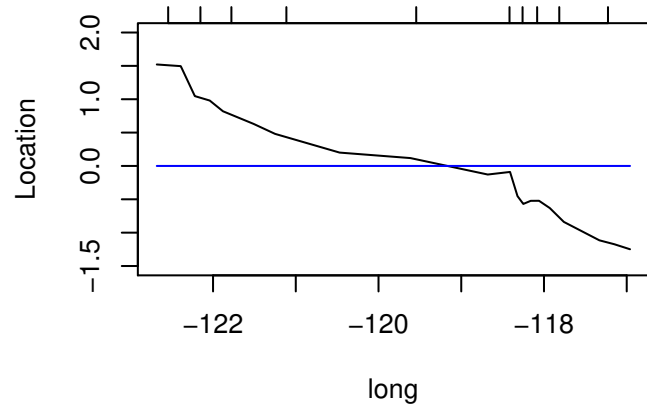
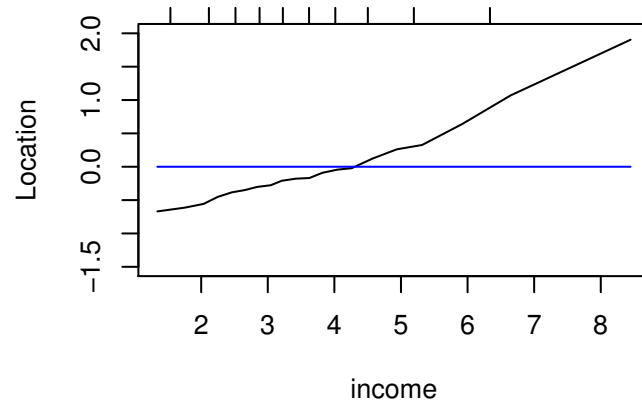
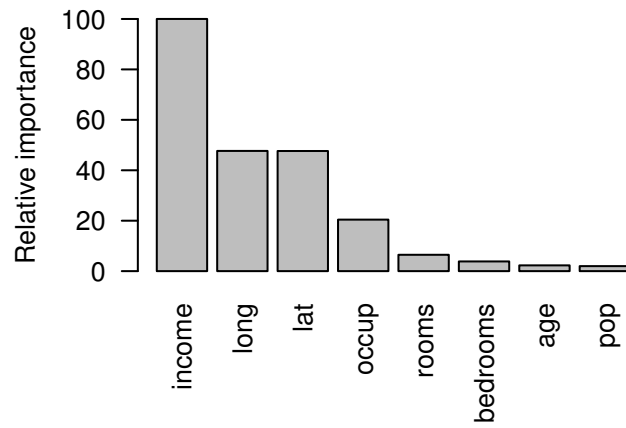
Population, Ave Occupancy, Latitude, Longitude)

### CA housing prices

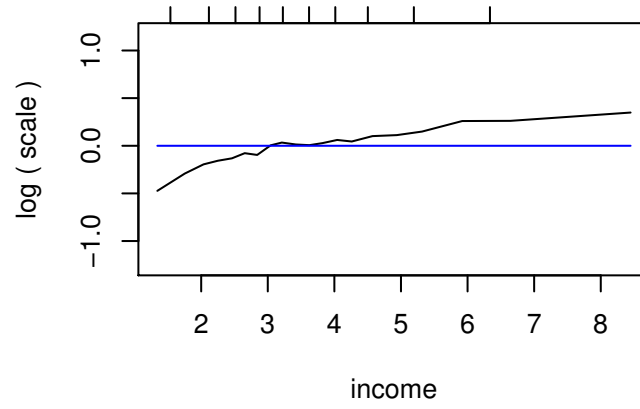
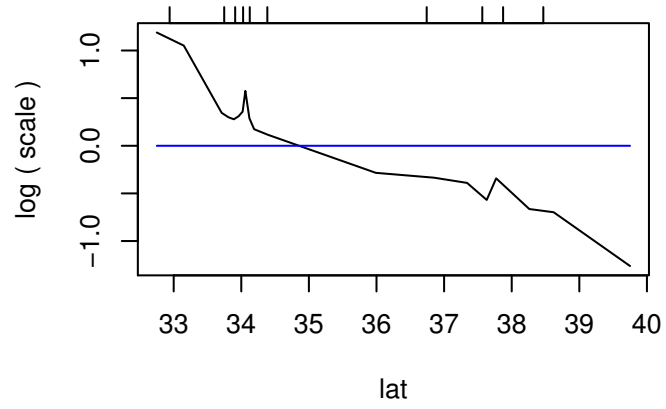
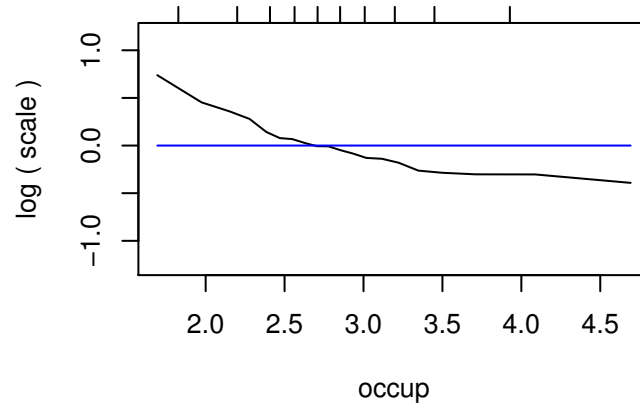
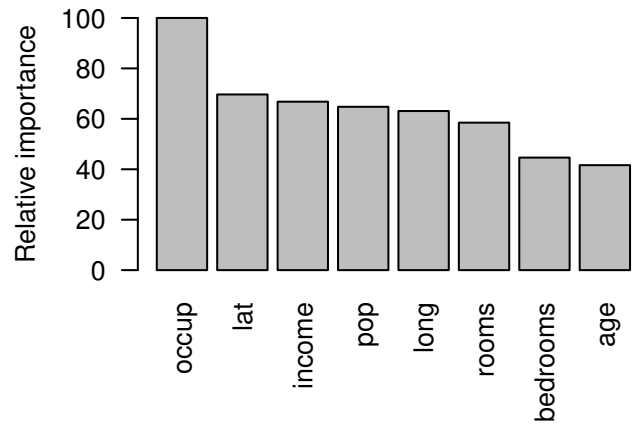




### CA housing : location model



### CA housing : log ( scale ) model



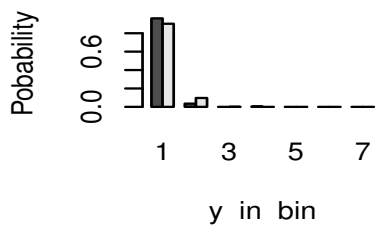
## QUESTIONNAIRE DATA

$$N = 8857, \quad p = 13$$

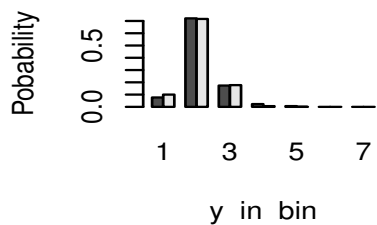
$$y = AGE \in \left\{ \begin{array}{l|l} 14 & 17 \\ 18 & 24 \\ 25 & 34 \\ 35 & 44 \\ 45 & 54 \\ 55 & 64 \\ 65 & \infty \end{array} \right.$$

$\mathbf{x} =$  (Occupation, Type of Home, Sex,  
Marital Status, Education, Income,  
Lived in BA, Dual Incomes, Persons in  
Household, Persons in Household  $< 18$ ,  
Householder Status, Ethnicity, Language)

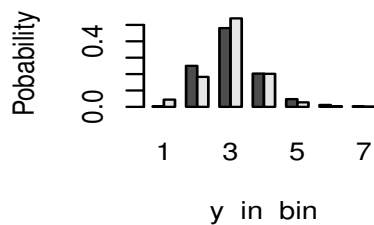
**f in bin 1 : 285**



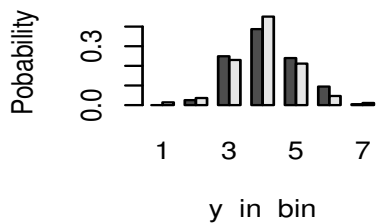
**f in bin 2 : 484**



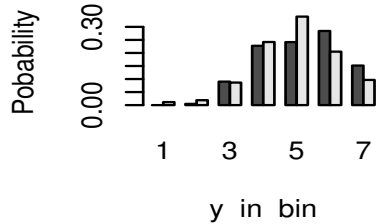
**f in bin 3 : 847**



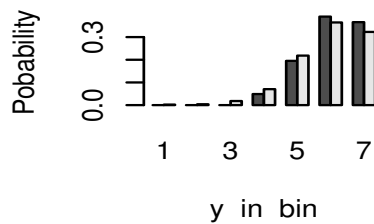
**f in bin 4 : 868**



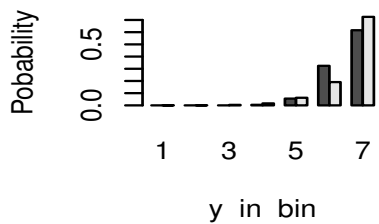
**f in bin 5 : 211**



**f in bin 6 : 41**

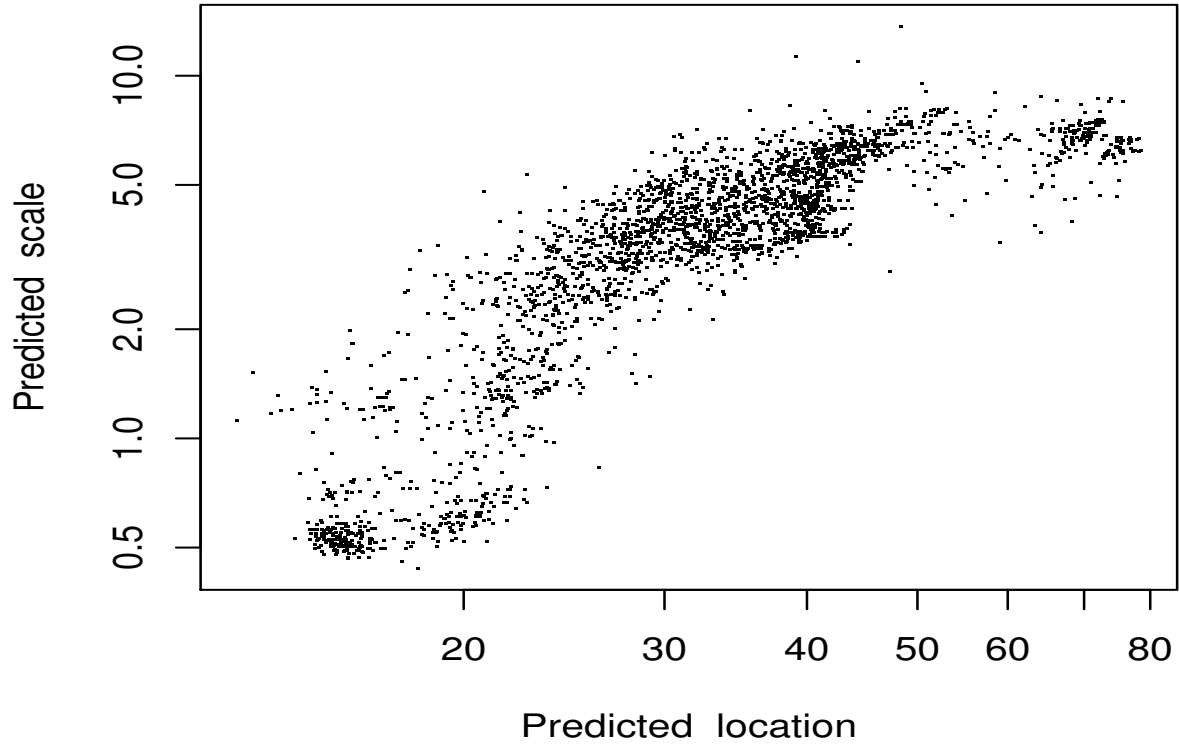


**f in bin 7 : 216**

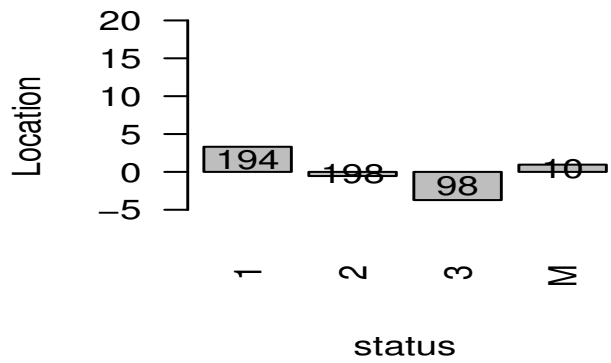
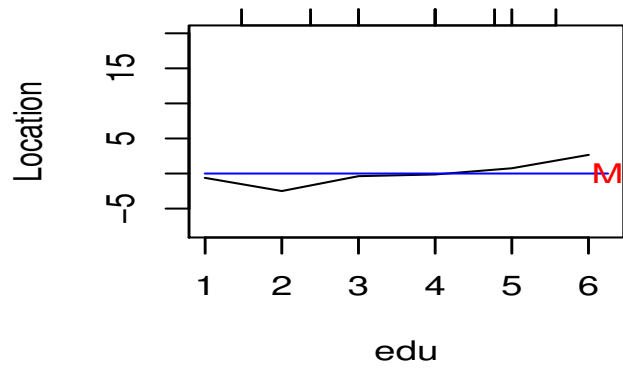
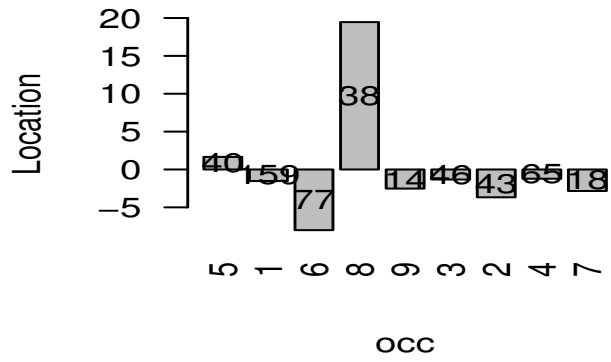
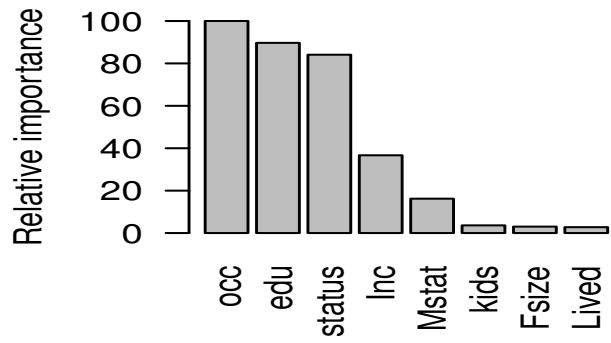


AGE predictions

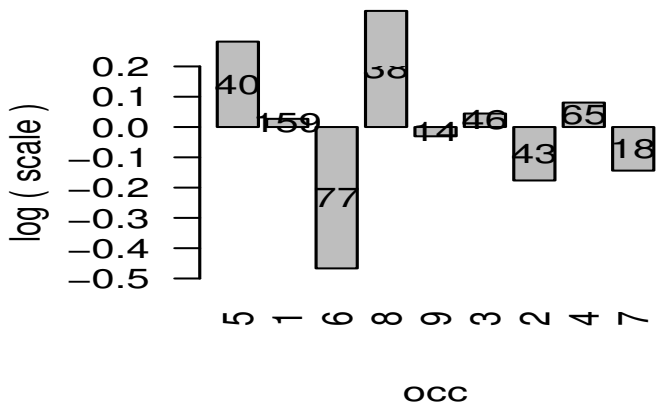
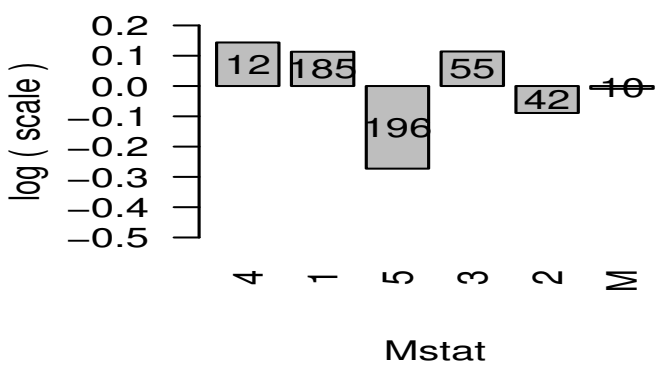
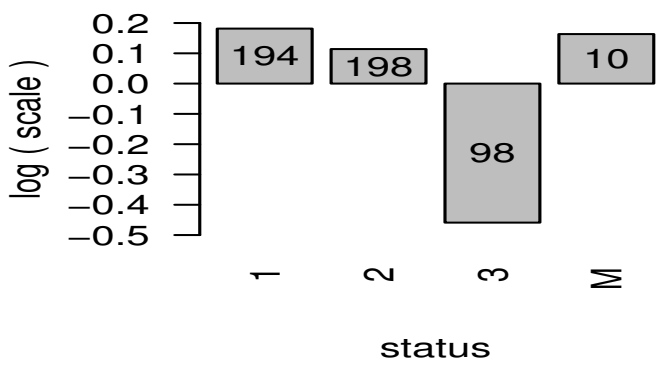
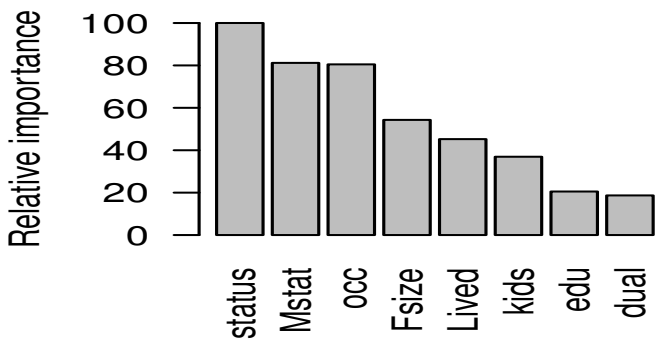
# AGE



### Location Model



### Scale Model





## Wine Quality Data (Irvine Repository)

$N = 6497$  samples of Portuguese "Vinho Verde"

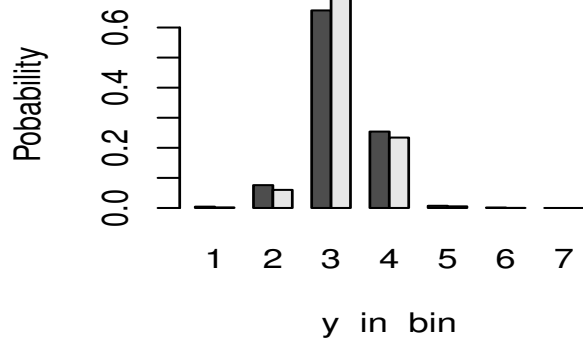
$\tilde{y} =$  Quality: integer (1, 2, ..., 10)

median of at least 3 expert evaluations

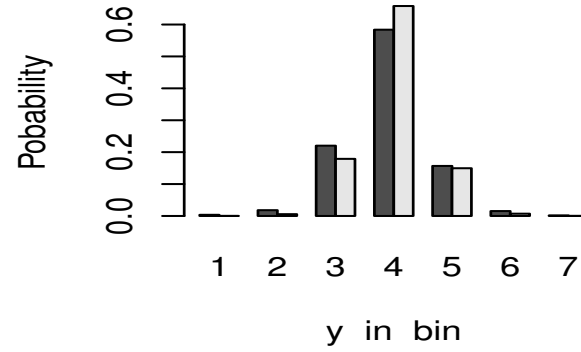
$\tilde{y} = k \Rightarrow y \in [k - 1/2, k + 1/2]$

$x =$ (Fixed acidity, Volatile acidity, Citric acid,  
Residual sugar, Chlorides, Free sulfur dioxide  
Total sulfur dioxide, Density, pH, Sulfates,  
Alcohol)

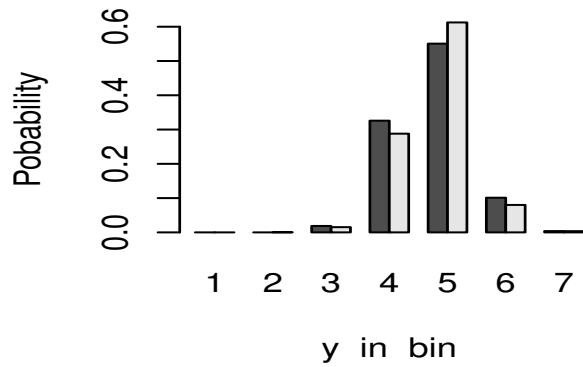
**f in bin 3 : 685**



**f in bin 4 : 1103**

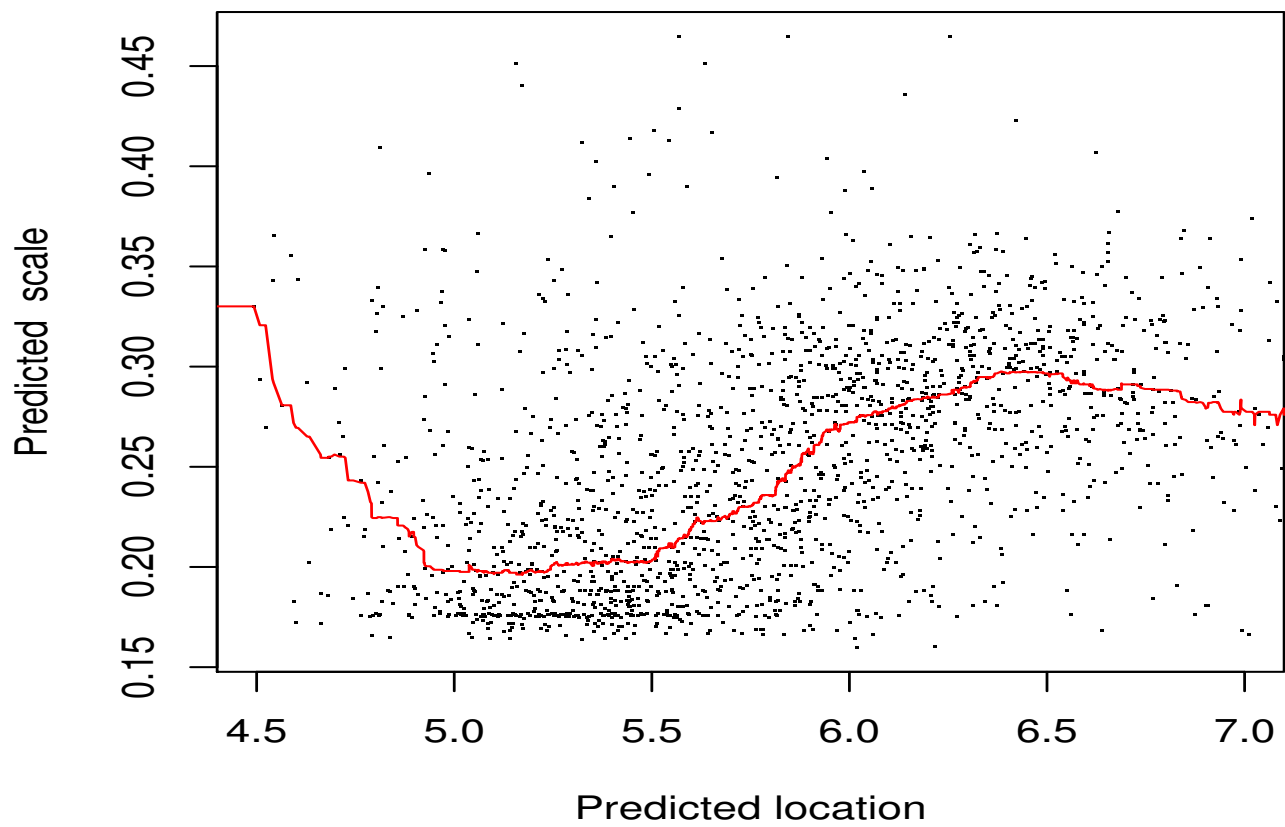


**f in bin 5 : 267**

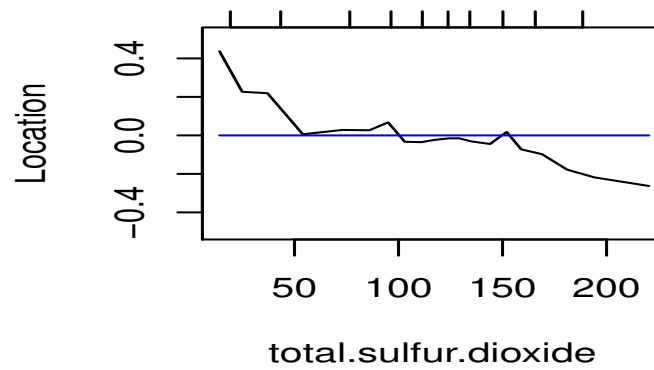
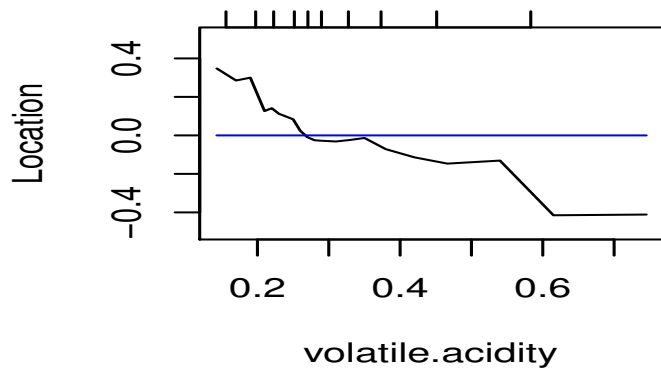
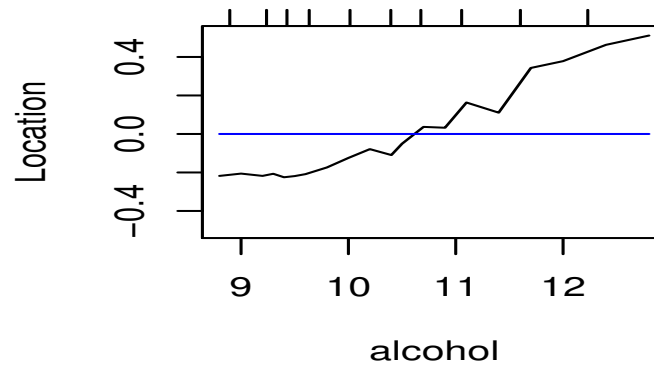
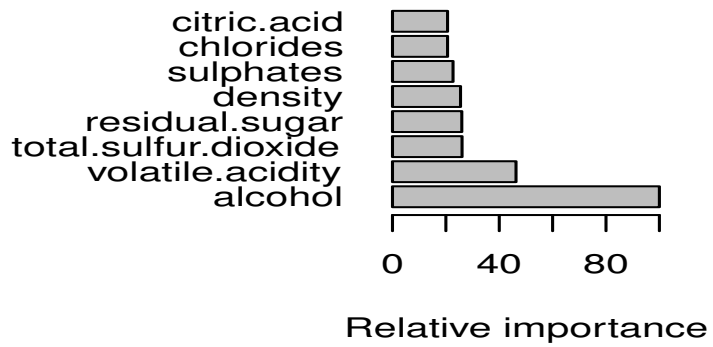


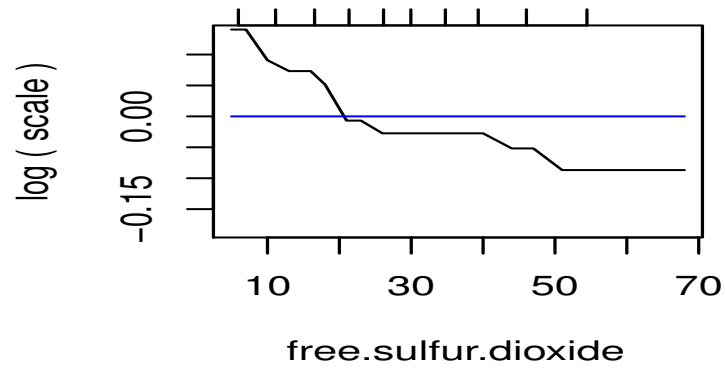
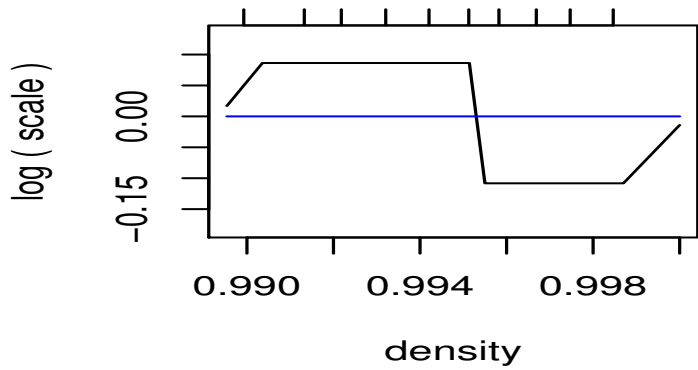
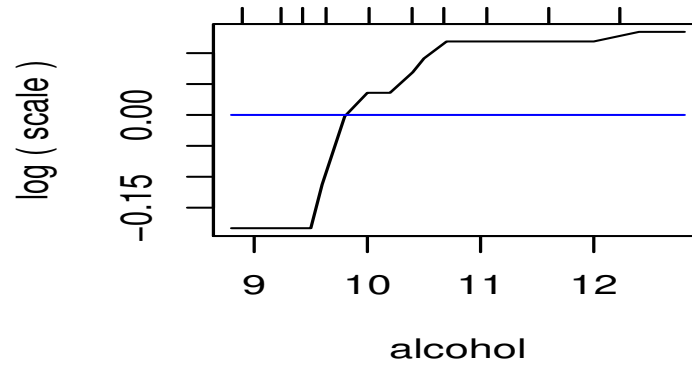
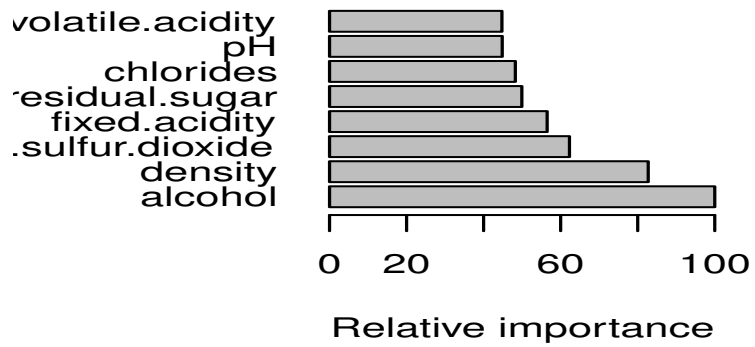
Wine quality data

## Wine quality data



### Wine: location





## ORDERED MULTICLASS LOGISTIC REGRESSION

$$y_i \in \{C_1 < C_2, \dots, C_{K-1} < C_K\}$$

Interval censored:

$\{a_i, b_i\} \rightarrow K$  disjoint intervals (bins):

$$\{b_0, b_1, \dots, b_K\} \quad b_0 = -\infty, \quad b_K = \infty$$

bins  $\sim$  classes with separating boundaries

$$\mathbf{b} = \{b_1, b_2, \dots, b_{K-1}\} \quad \text{unknown}$$

(overall location & scale arbitrary)

## OPTIMAL SCALING (aka ACE)

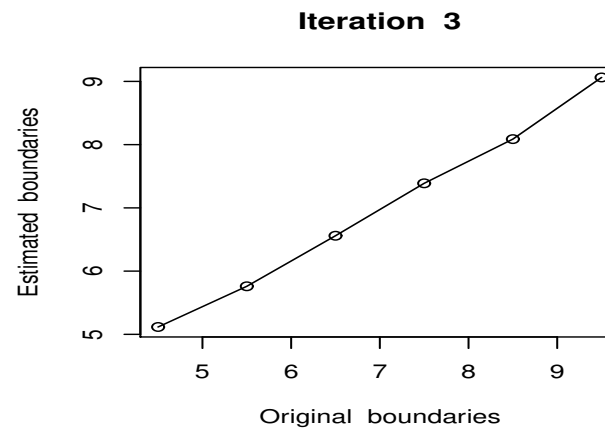
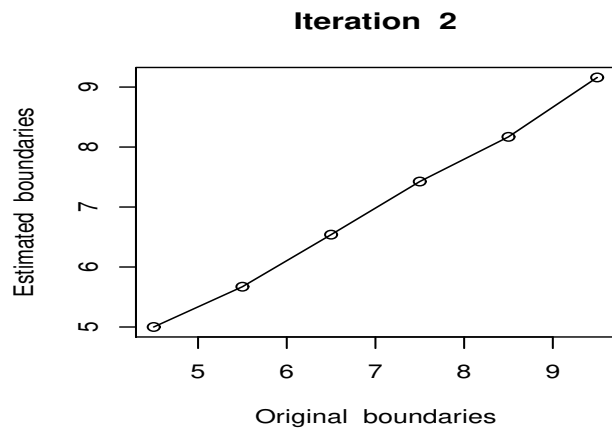
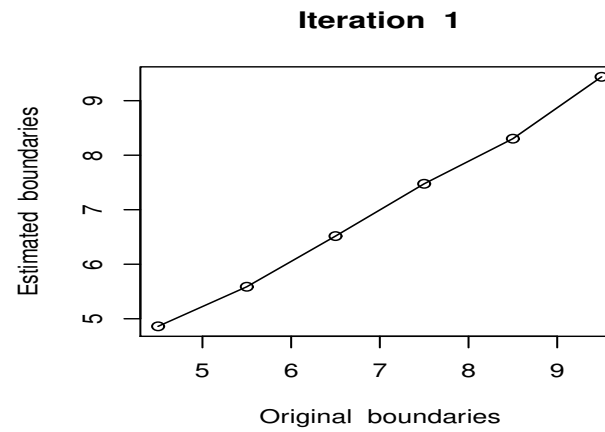
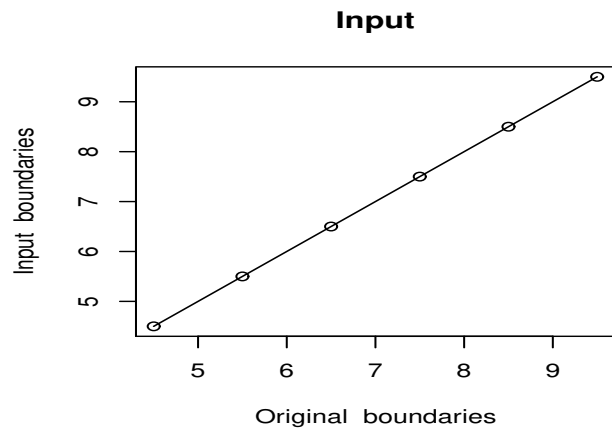
$$[\hat{\mathbf{b}}, \hat{f}(\mathbf{x}), \hat{s}(\mathbf{x})] = \arg \min_{\mathbf{b}, (f,s) \in F} \sum_{i=1}^N L[b_{k(i)-1}, b_{k(i)}, f(\mathbf{x}_i), s(\mathbf{x}_i)]$$

$$L(u, v, f, s) = -\log \left( \frac{1}{1+e^{(f-v)/s}} - \frac{1}{1+e^{(f-u)/s}} \right)$$

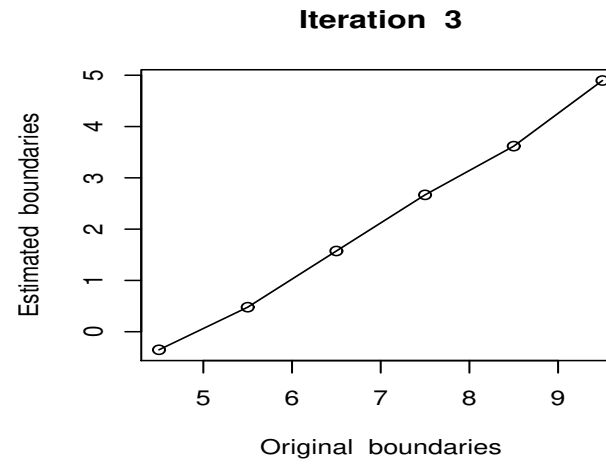
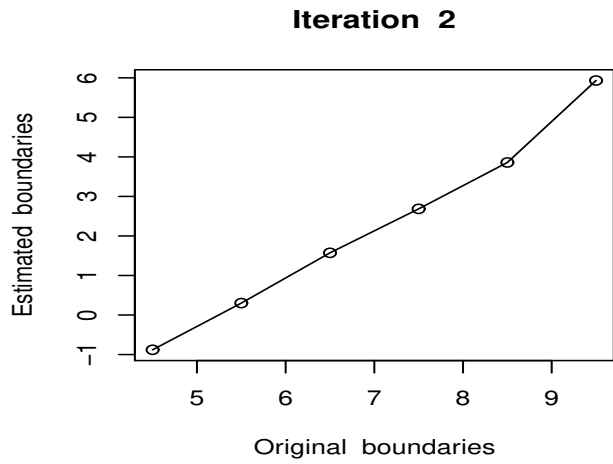
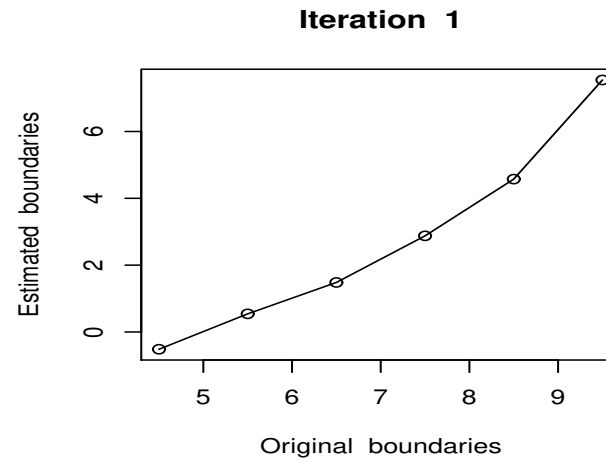
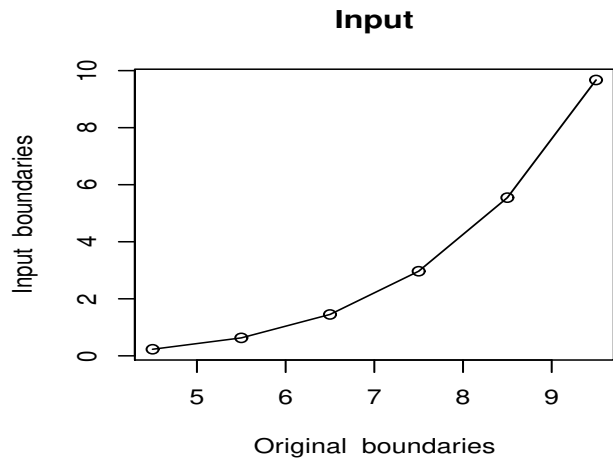
Alternating optimization:

$$\mathbf{b} : \sum_{i=1}^N \left( \frac{1}{1 + e^{(f(\mathbf{x}_i) - b_k)/s(\mathbf{x}_i)}} - \sum_{j=1}^k I(y_i = c_j) \right) = 0$$





Wine quality data – optimal scaling



Wine quality data – optimal scaling

## ASYMMETRIC ERRORS

$$y | \mathbf{x} = f(\mathbf{x}) + \begin{cases} s_l(\mathbf{x}) \cdot \varepsilon & \varepsilon \leq 0 \\ s_u(\mathbf{x}) \cdot \varepsilon & \varepsilon > 0 \end{cases}$$

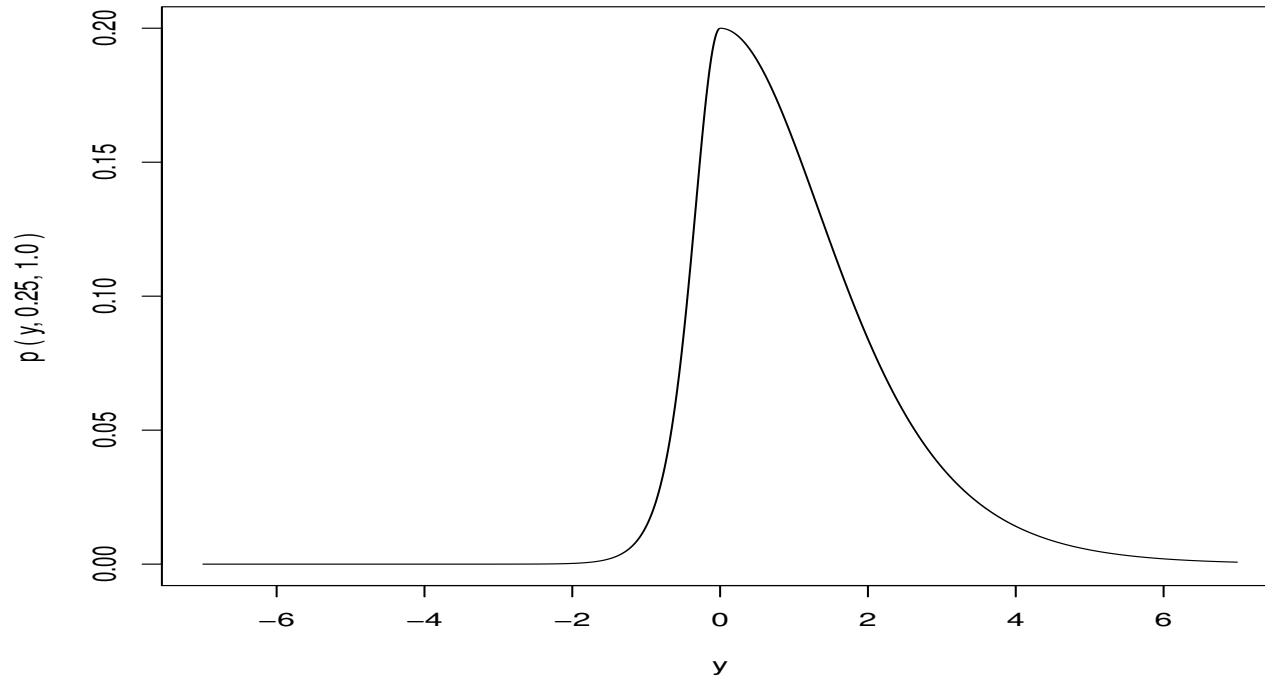
$f(\mathbf{x}) = \text{mode at } \mathbf{x}$

$s_l(\mathbf{x}) = \text{lower scale at } \mathbf{x}$

$s_u(\mathbf{x}) = \text{upper scale at } \mathbf{x}$

$\varepsilon \sim \text{standard logistic distribution}$

### Asymmetric logistic



$$f = 0, \quad s_l = 1/4, \quad s_u = 1$$

## EXERCISE (no censoring)

$$[\hat{f}(\mathbf{x}), \hat{s}_l(\mathbf{x}), \hat{s}_u(\mathbf{x})]$$

$$= \arg \min_{(f, s_l, s_u) \in F} \sum_{i=1}^N L[y_i, f(\mathbf{x}_i), s_l(\mathbf{x}_i), s_u(\mathbf{x}_i)]$$

$$L[y, f, s_l, s_u] =$$

$$L[y, f, s_l] \cdot I[y - f \leq 0] + L[y, f, s_u] \cdot I[y - f > 0]$$

$$L(y, f, s) = \log(s) + (y - f)/s + 2 \log(1 + e^{-(y-f)/s})$$

Iterative gradient boosting

## ASYMMETRIC DIAGNOSTICS

$$(1) \text{ median } [y \mid f(\mathbf{x}), s_l(\mathbf{x}), s_u(\mathbf{x})] = f(\mathbf{x})$$

$$+ \begin{cases} s_u(\mathbf{x}) \log \left( \frac{3s_u(\mathbf{x}) - s_l(\mathbf{x})}{s_u(\mathbf{x}) + s_l(\mathbf{x})} \right) & s_l(\mathbf{x}) \leq s_u(\mathbf{x}) \\ -s_l(\mathbf{x}) \log \left( \frac{3s_l(\mathbf{x}) - s_u(\mathbf{x})}{s_u(\mathbf{x}) + s_l(\mathbf{x})} \right) & s_l(\mathbf{x}) > s_u(\mathbf{x}) \end{cases}$$

$$(2) \text{ median}_{y \leq f(\mathbf{x})} [ |y - f(\mathbf{x})| \mid s_l(\mathbf{x}) ] / \log(3) = s_l(\mathbf{x})$$

$$(3) \text{ median}_{y > f(\mathbf{x})} [ |y - f(\mathbf{x})| \mid s_u(\mathbf{x}) ] / \log(3) = s_u(\mathbf{x})$$

(1/2) – Million Song Dataset (Irvine Repository)

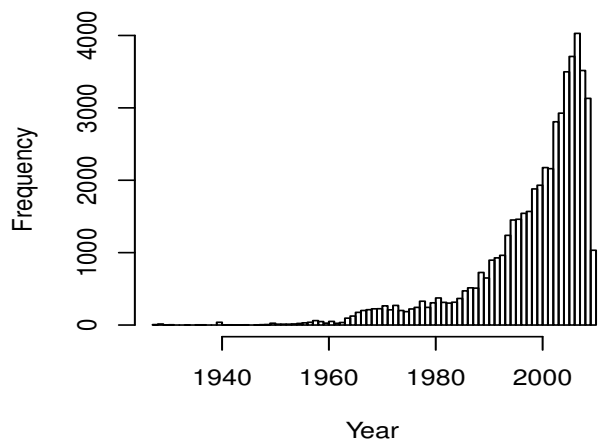
$N = 515345$  songs (463715 train, 51630 test)

$y =$  year released (1922 – 2011)

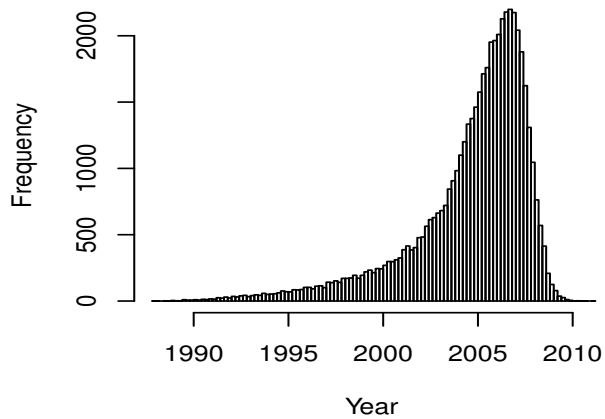
$\mathbf{x} = 90$  attributes ( Echo Nest API):

12 timbre average, 78 = timbre covariance

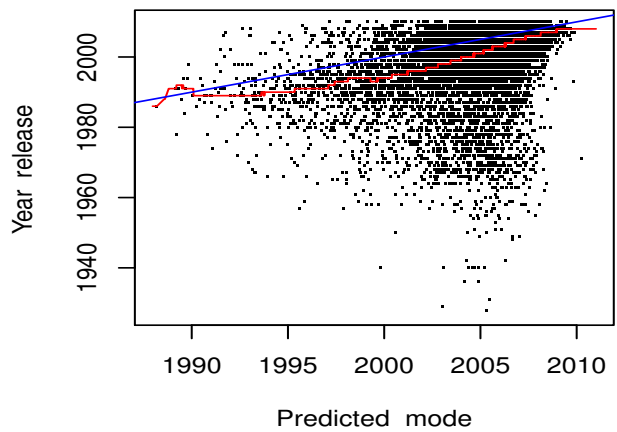
**MSD song release**



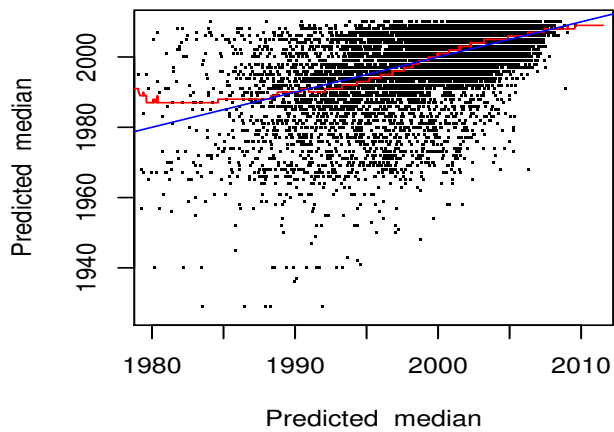
**Predicted mode year | x**



**Mode**

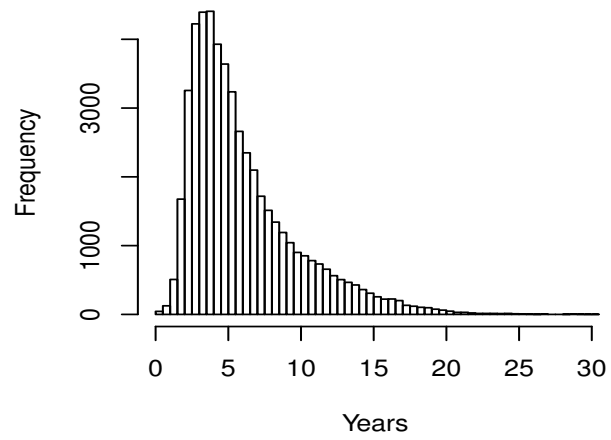


**Median**

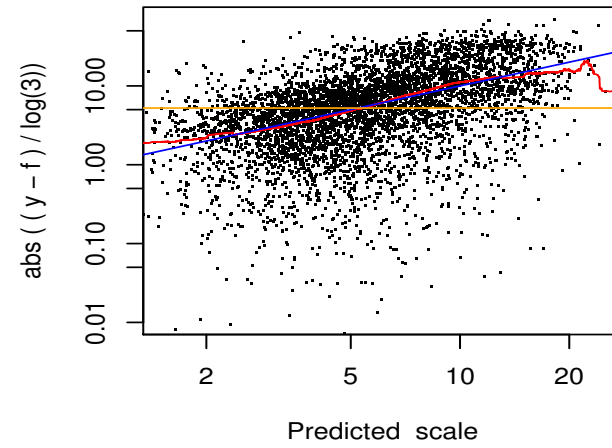




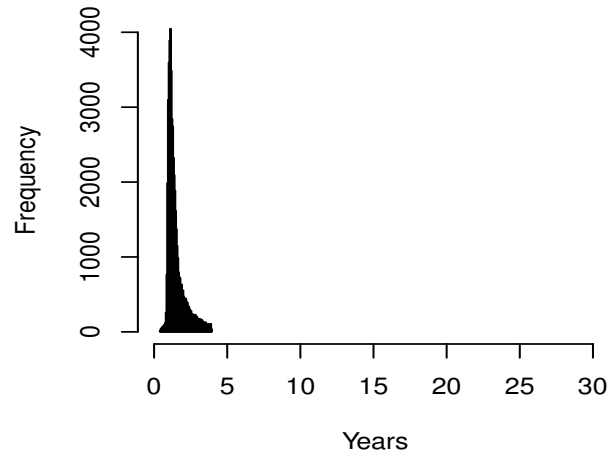
**Lower scale**



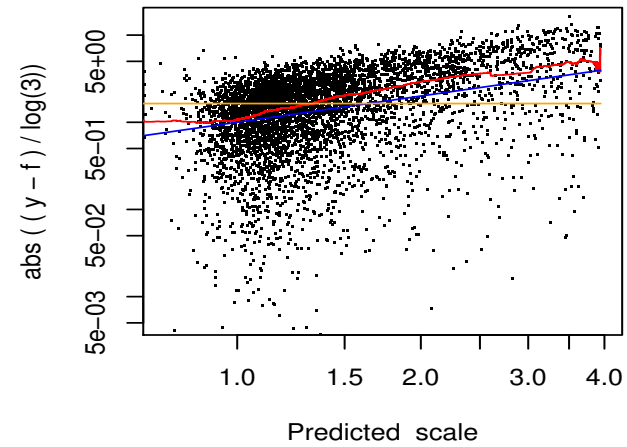
**Lower scale**

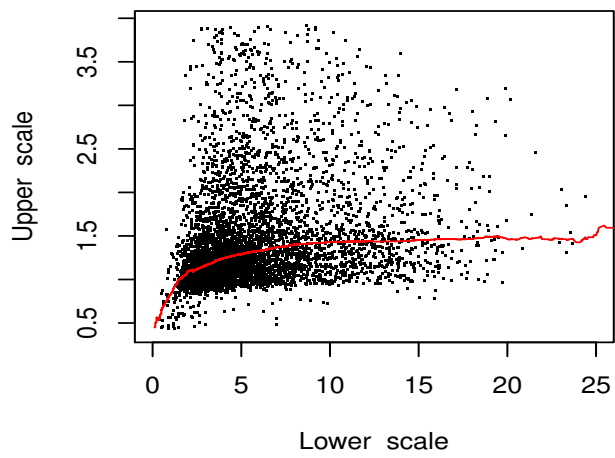
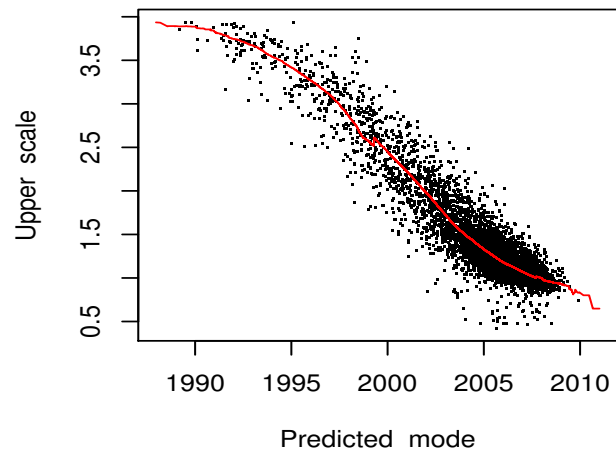
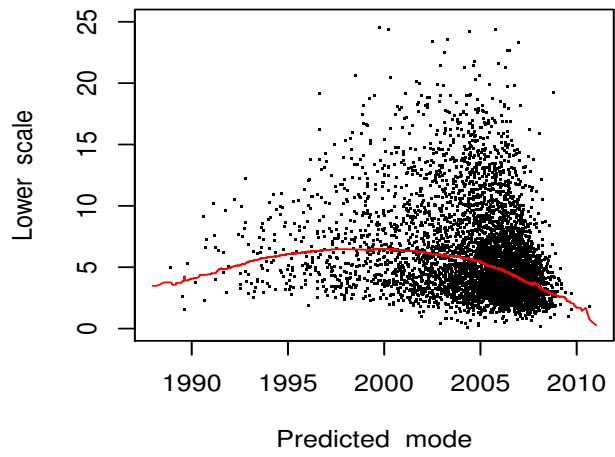


**Upper scale**

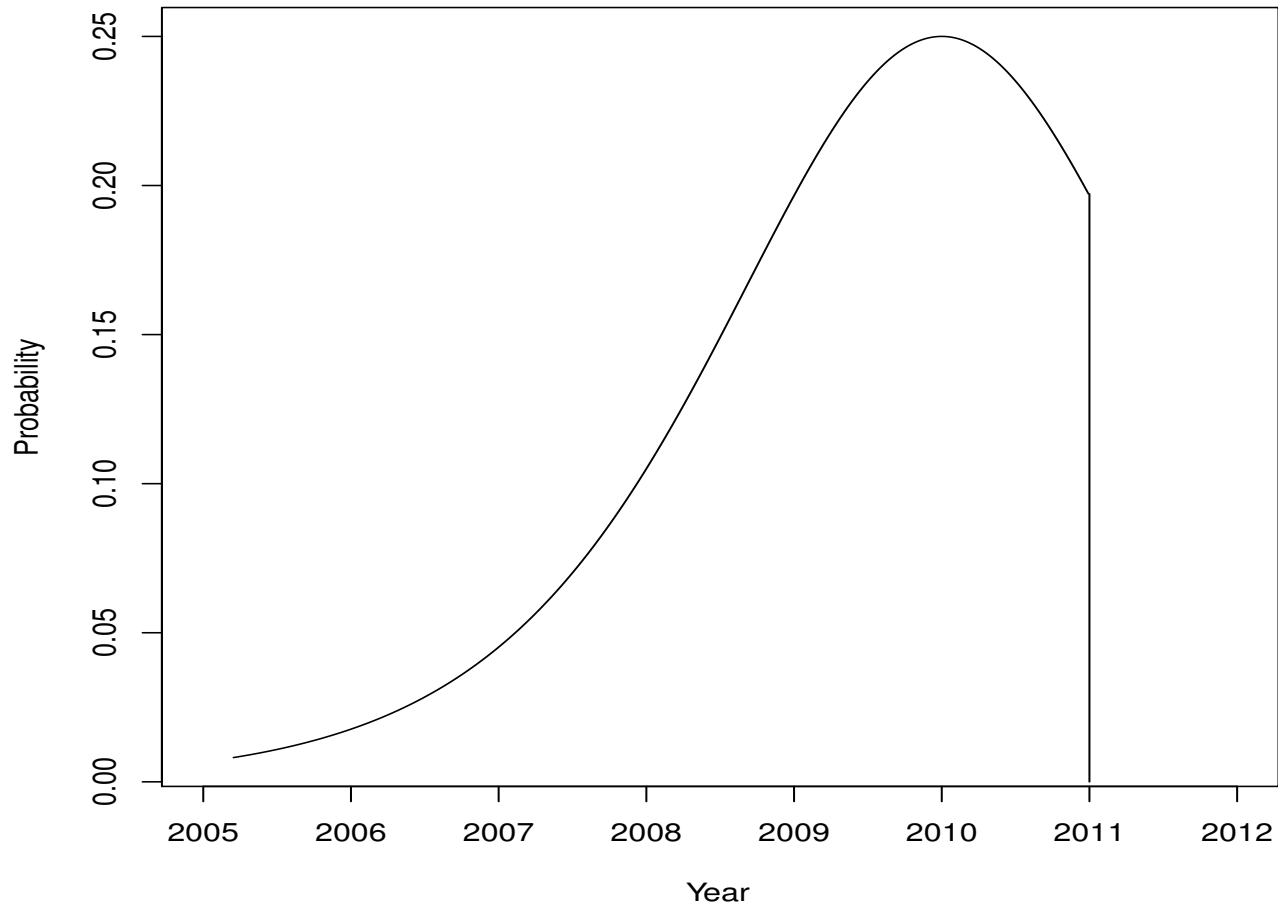


**Upper scale**





### Truncated logistic



Predict  $y \mid \mathbf{x}$ , settle for  $p(y \mid \mathbf{x})$

1. ROBUST (logistic) estimate: loc  $f(\mathbf{x})$  & scale  $s(\mathbf{x})$
2. General censoring:  $a_i \leq y_i \leq b_i$  (ubiquitous)
3. Graphical diagnostics
4. Ordered multi-class classification
5. Asymmetric  $p(y \mid \mathbf{x})$  :  $f(\mathbf{x})$ ,  $s_l(\mathbf{x})$ ,  $s_u(\mathbf{x})$

## REFERENCES

Gradient boosting: *Ann. Statist.*, **29**. 1189 – 1232 (2001)

Optimal scaling:

ACE: *J. Amer. Statist. Assoc.* **80**. 580 – 598 (1985)

GIFI: *Nonlinear Multivariate Analysis*. Wiley, N.Y. (1990)

Slides: <http://statweb.stanford.edu/~jhf/talks/kdd.pdf>