

Predicting Regression Probability Distributions

with Imperfect Data

through Optimal Transformations

Jerome H. Friedman

Stanford University

MACHINE LEARNING / REGRESSION

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

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)

Neither very likely

Tukey:

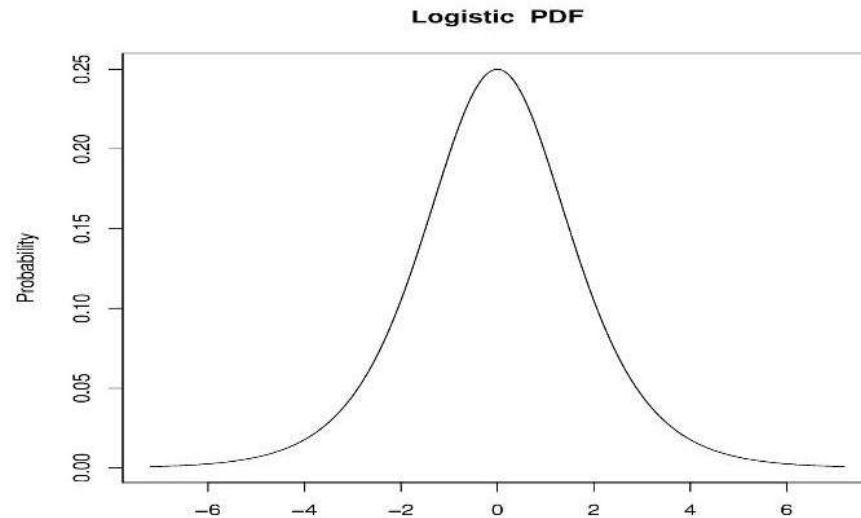
“small residuals \simeq normal, larger have heavier tails.”

LOGISTIC DISTRIBUTION

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

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

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:

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

(b) can affect decision

CENSORING

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

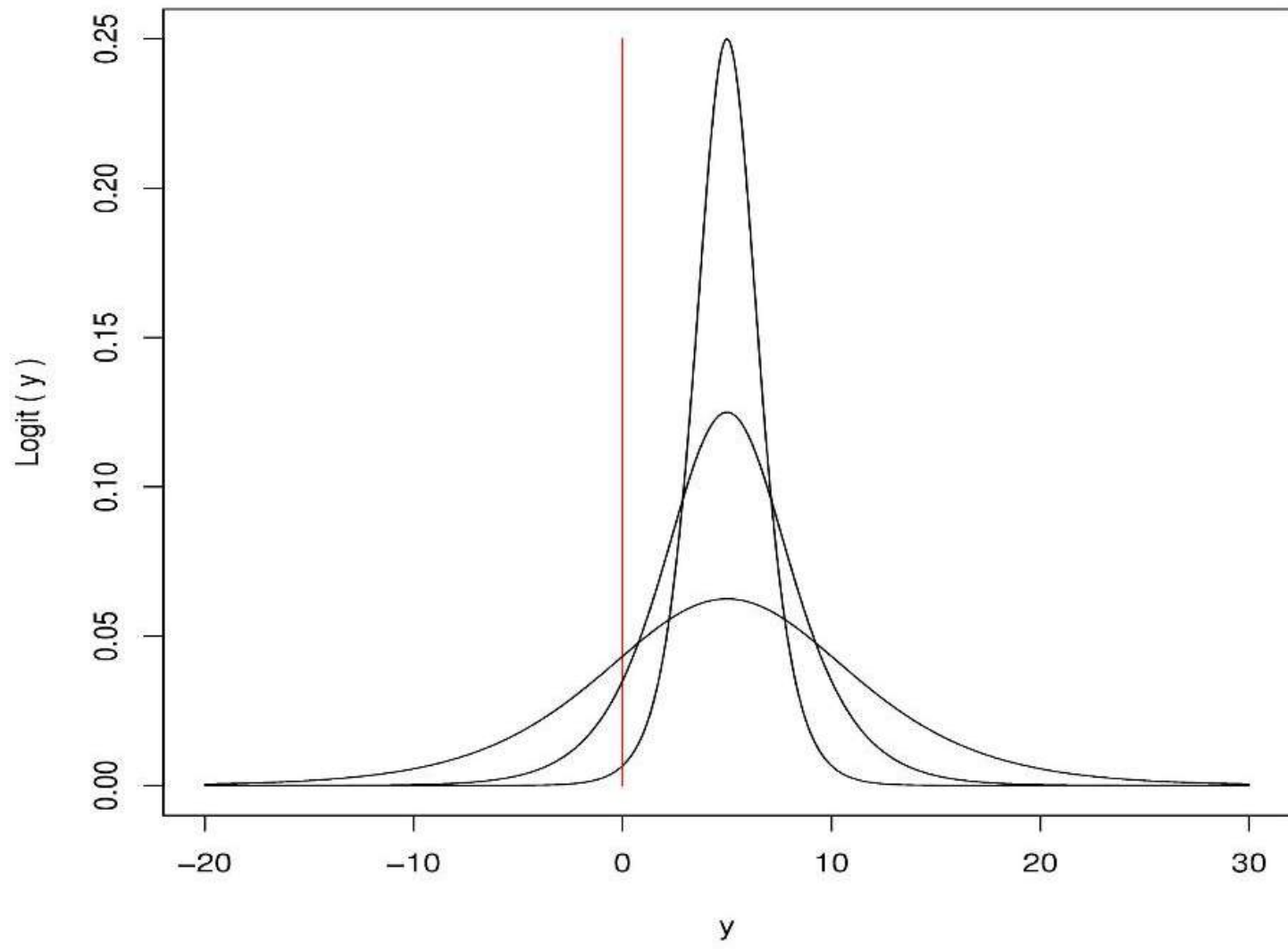
Otherwise: interval censored $[a_i, b_i]$

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 $s(\mathbf{x})$



EXERCISE

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

$$L(a, b, f(\mathbf{x}), s(\mathbf{x})) = \log \left[\frac{1}{1 + \exp((f(\mathbf{x}) - a)/s(\mathbf{x}))} - \frac{1}{1 + \exp((f(\mathbf{x}) - b)/s(\mathbf{x}))} \right]$$

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 {

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

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

}

Until no change

OPTIMAL TRANSFORMATIONS

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

$g(y)$ = unknown monotonic function

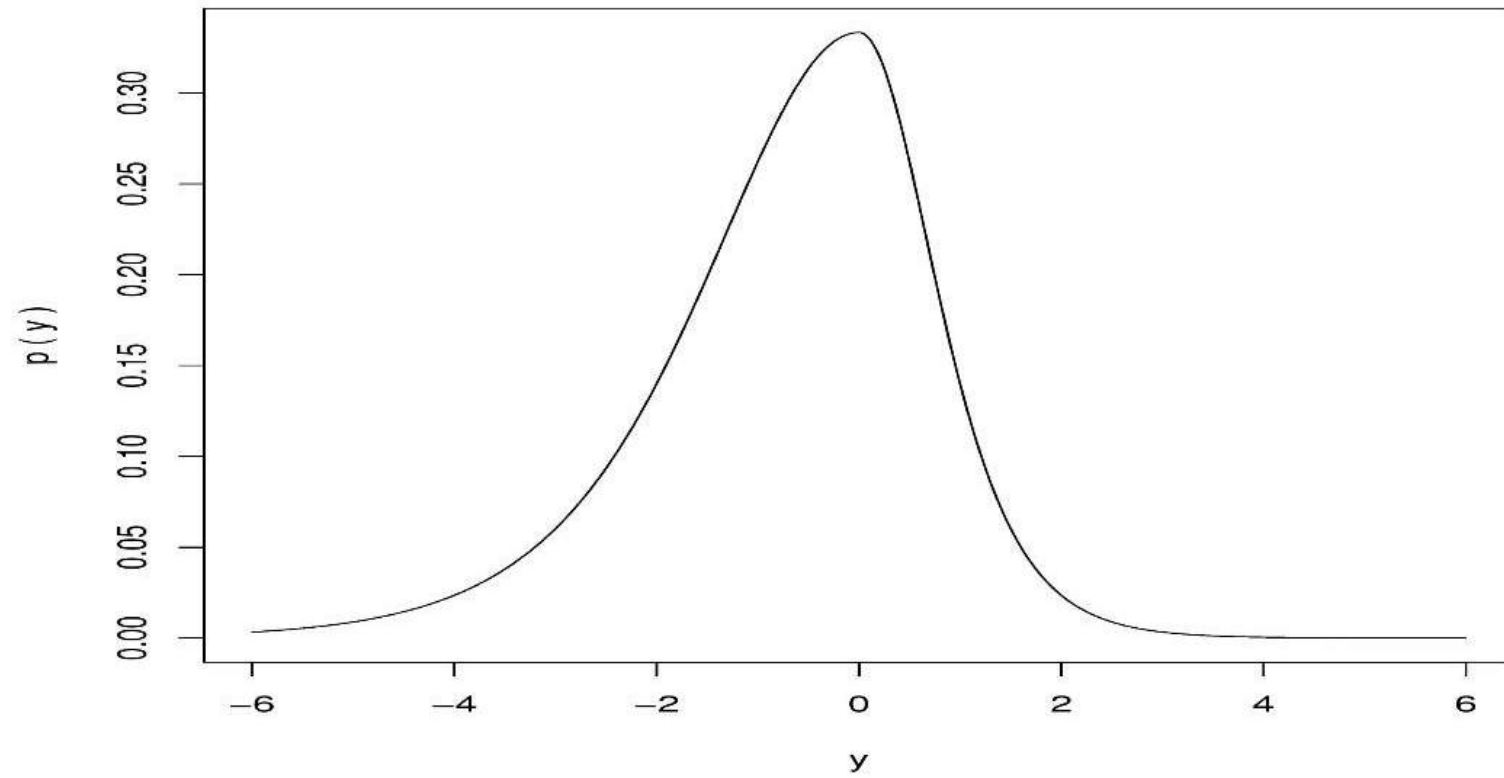
$$[\hat{g}(y), \hat{f}(\mathbf{x}), \hat{s}(\mathbf{x})] = \arg \min_{g, (f, s) \in F} \sum_{i=1}^N L[g(y_i), f(\mathbf{x}_i), s(\mathbf{x}_i)]$$

$$\frac{1}{N} \sum_{i=1}^N CDF(\hat{g}(y), \hat{f}(\mathbf{x}_i), \hat{s}(\mathbf{x}_i)) = CDF_y(y)$$

$$CDF(g(y), f(\mathbf{x}), s(\mathbf{x})) = \frac{1}{1 + \exp((f(\mathbf{x}) - g(y))/s(\mathbf{x}))}$$

ASYMMETRY

Asymmetric Logistic



ASYMMETRIC LOGISTIC DISTRIBUTION

$$p(z | f, s_l, s_u) = \frac{2}{s_l + s_u} \left[\frac{I(z \leq f) \exp((f - z)/s_l)}{(1 + \exp((f - z)/s_l))^2} + \frac{I(z > f) \exp((f - z)/s_u)}{(1 + \exp((f - z)/s_u))^2} \right]$$

$$g(y) = f(\mathbf{x}) + \eta, \quad \begin{cases} \eta = -s_l(\mathbf{x}) \cdot |\varepsilon|, & \text{prob} = s_l(\mathbf{x}) / (s_l(\mathbf{x}) + s_u(\mathbf{x})) \\ \eta = +s_u(\mathbf{x}) \cdot |\varepsilon|, & \text{prob} = s_u(\mathbf{x}) / (s_l(\mathbf{x}) + s_u(\mathbf{x})) \end{cases}$$

f = mode, s_l = lower scale, s_u = upper scale

$$\varepsilon \sim L(0, 1)$$

ASYMMETRIC GRADIENT BOOSTING

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

Loop {

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

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

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

} Until change $<$ threshold.

OPTIMAL TRANSFORMATIONS

$$g(y) \mid \mathbf{x} \sim \text{Logistic}(\hat{f}(\mathbf{x}), \hat{s}_l(\mathbf{x}), \hat{s}_u(\mathbf{x}))$$

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

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

$$\frac{1}{N} \sum_{i=1}^N CDF(\hat{g}(y), \hat{f}(\mathbf{x}_i), \hat{s}_l(\mathbf{x}_i), s_u(\mathbf{x}_i)) = CDF_y(y)$$

$$CDF(g(y), f(\mathbf{x}), s(\mathbf{x})) = \text{messy closed form}$$

DIAGNOSTICS

For $\mathbf{x}_i \in S$:

$$\text{Model: } \Pr(\hat{g}(y) < z) = \frac{1}{|S|} \sum_{i \in S} CDF(z, \hat{f}(\mathbf{x}_i), \hat{s}_l(\mathbf{x}_i), \hat{s}_u(\mathbf{x}_i))$$

$$\text{Empirical: } \widehat{\Pr}(\hat{g}(y) < z) = \frac{1}{|S|} \sum_{i \in S} I(\hat{g}(y_i) \leq z)$$

$$\text{Symmetric: } (g(y_i) - f(x_i))/s(x_i) \sim L(0, 1)$$

$$\text{Asymmetric : } I(-) (\hat{g}(y_i) - \hat{f}(\mathbf{x}_i)) / \hat{s}_l(\mathbf{x}_i)$$

$$+ I(+) (\hat{g}(y_i) - \hat{f}(\mathbf{x}_i)) / \hat{s}_u(\mathbf{x}_i) \sim L(0, 1)$$

DATA SUBSETS S

$$S = \{i \mid r < \hat{f}(\mathbf{x}_i) \leq t \ \& \ u < \hat{s}(\mathbf{x}_i) \leq v\}.$$

$$S = \{i \mid r < \hat{f}(\mathbf{x}_i) \leq t \ \& \ u < \sqrt{\hat{s}_l(\mathbf{x}_i)\hat{s}_u(\mathbf{x}_i)} \leq v\}$$

$$S = \{i \mid r < \hat{f}(\mathbf{x}_i) \leq t \ \& \ u < \hat{s}_l(\mathbf{x}_i) \leq v\}.$$

$$S = \{i \mid r < \hat{f}(\mathbf{x}_i) \leq t \ \& \ u < \hat{s}_{ul}(\mathbf{x}_i) \leq v\}.$$

QUESTIONNAIRE DATA ($N = 8856$)

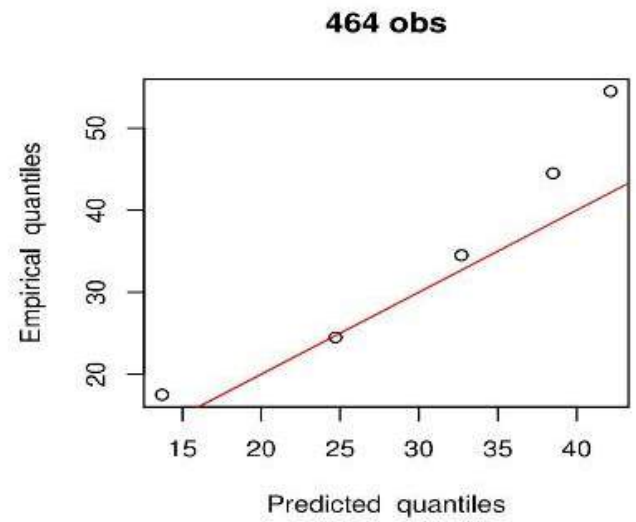
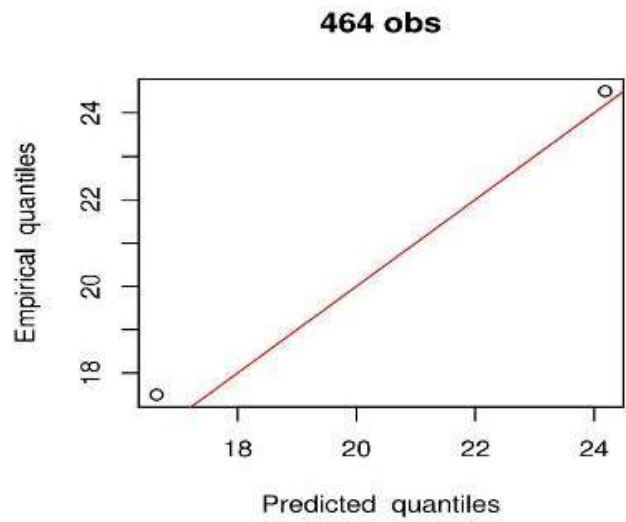
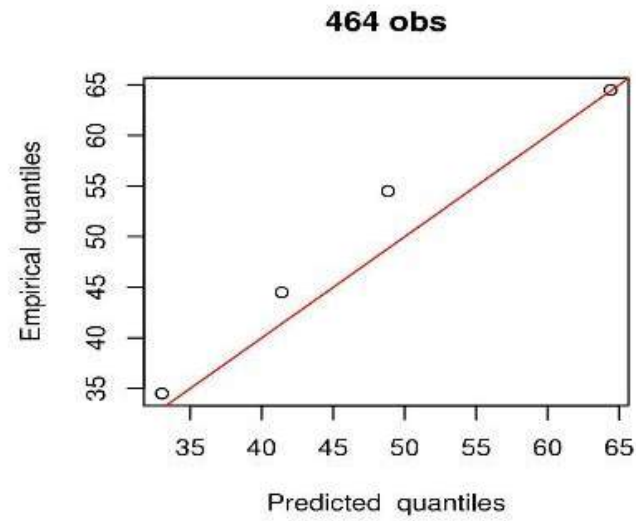
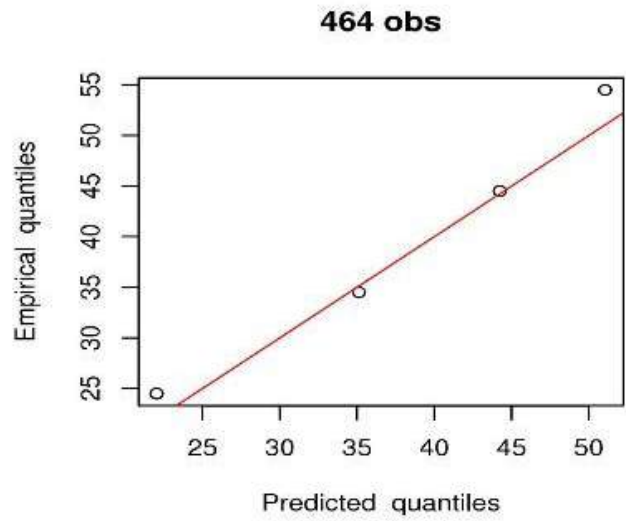
Bay Area Shopping Mall Customers

AGE \in

1	2	3	4	5	6	7
17 & under	18–24	25–34	35–44	45–54	55–64	65 & older

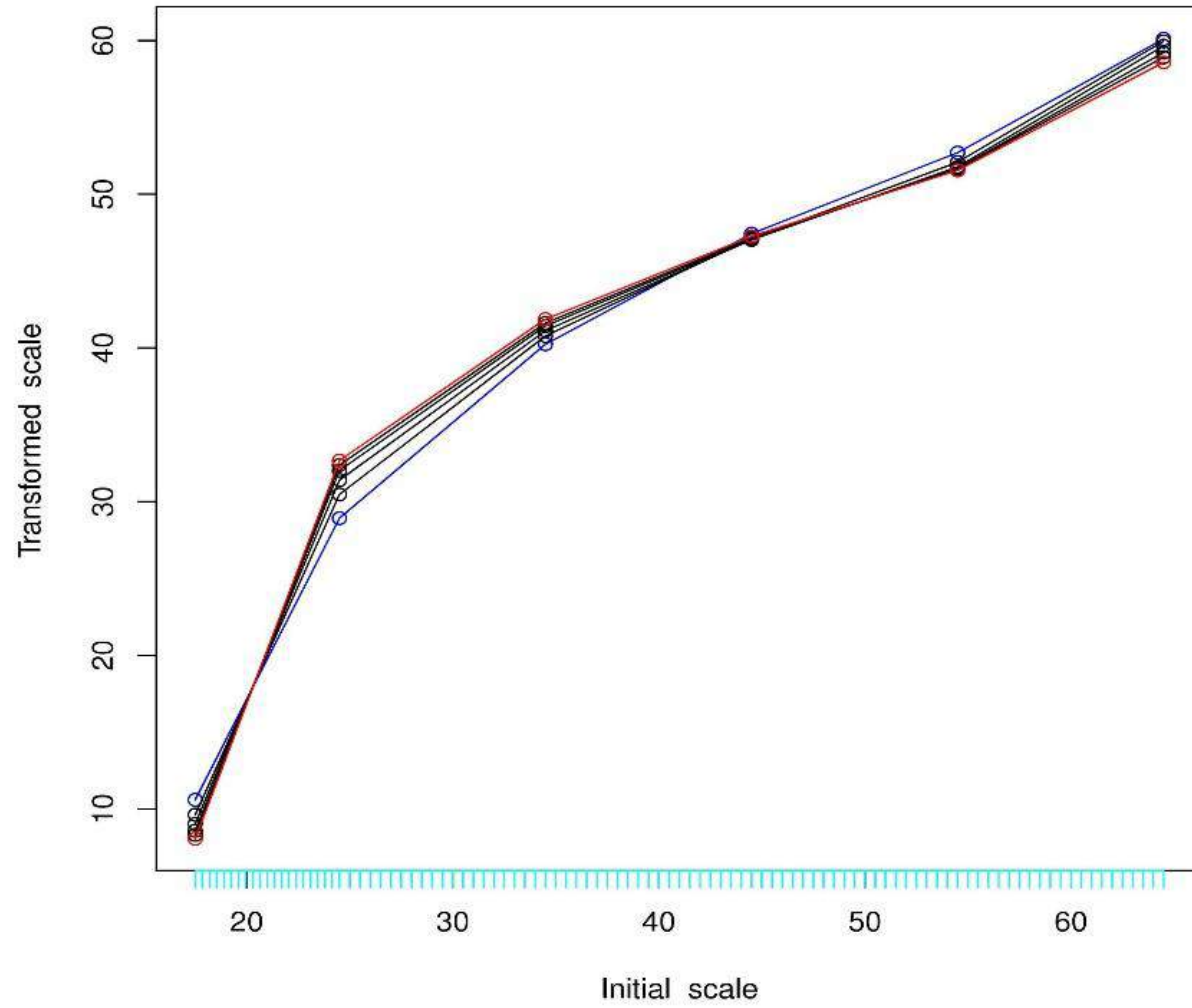
$\mathbf{x} =$

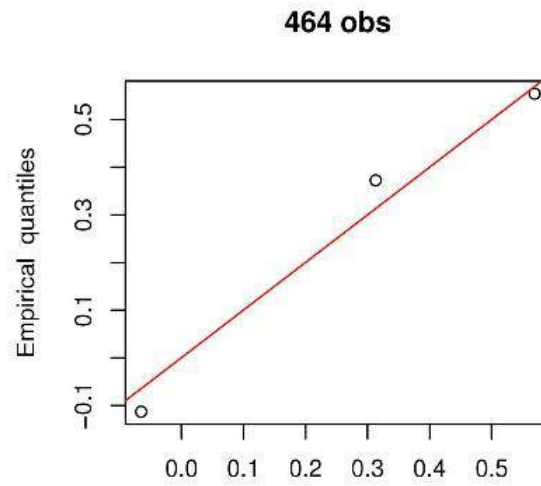
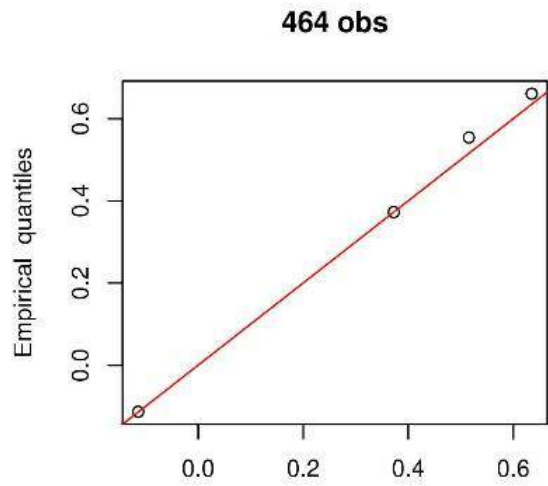
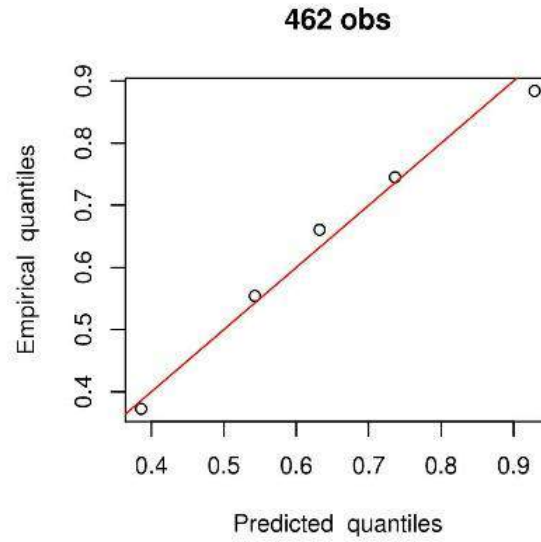
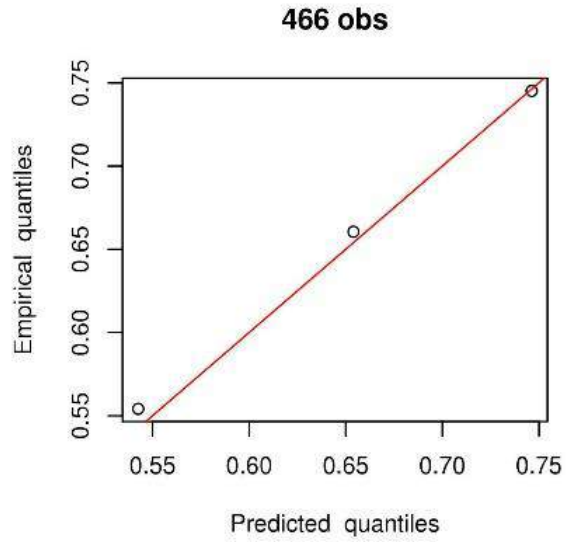
- | | | | |
|---|-------------------|----|-------------------------------|
| 1 | Occupation | 8 | Dual Incomes |
| 2 | Type of home | 9 | Persons in household |
| 3 | Gender | 10 | Persons in household under 18 |
| 4 | Marital status | 11 | Householder status |
| 5 | Education | 12 | Ethnic classification |
| 6 | Annual income | 13 | Language |
| 7 | Lived in Bay Area | | |

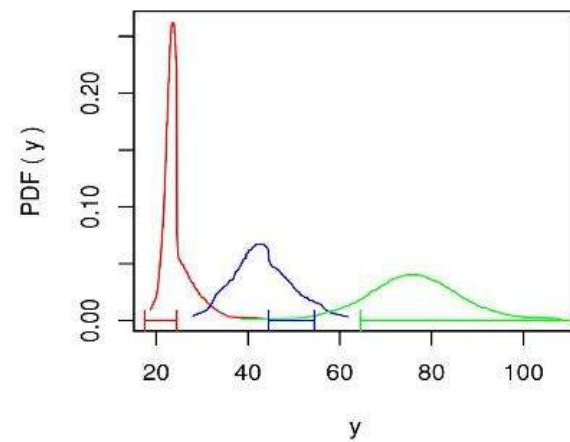
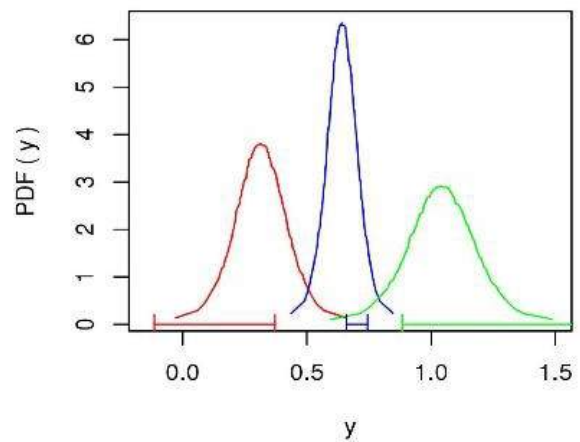
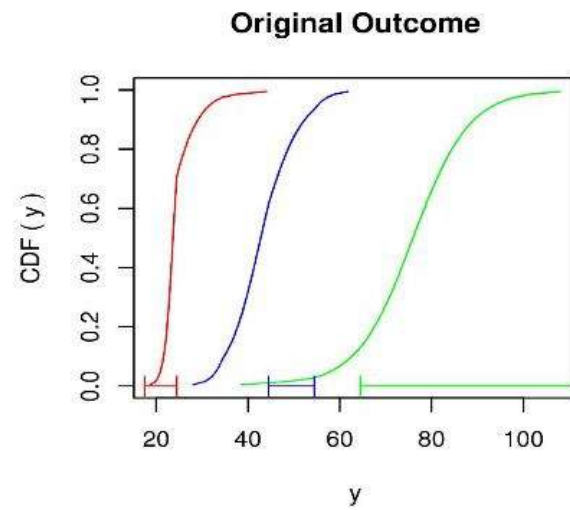
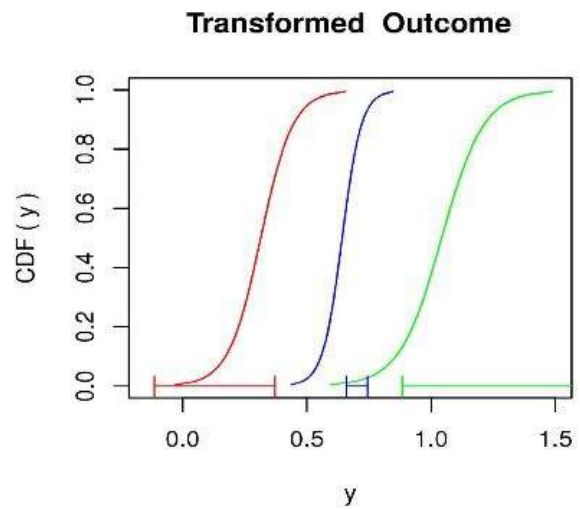


Q-Q plots of empirical versus predicted quantiles for untransformed Omnireg solution.

Age Transformation Sequence







MASHABLE ONLINE NEWS POPULARITY

(Irvine Repository)

$N = 39644$ news articles

25000 train, 7322 mod. sel., 7322 test

$y =$ popularity (number of shares on social media)

$x = 59$ numerical attributes

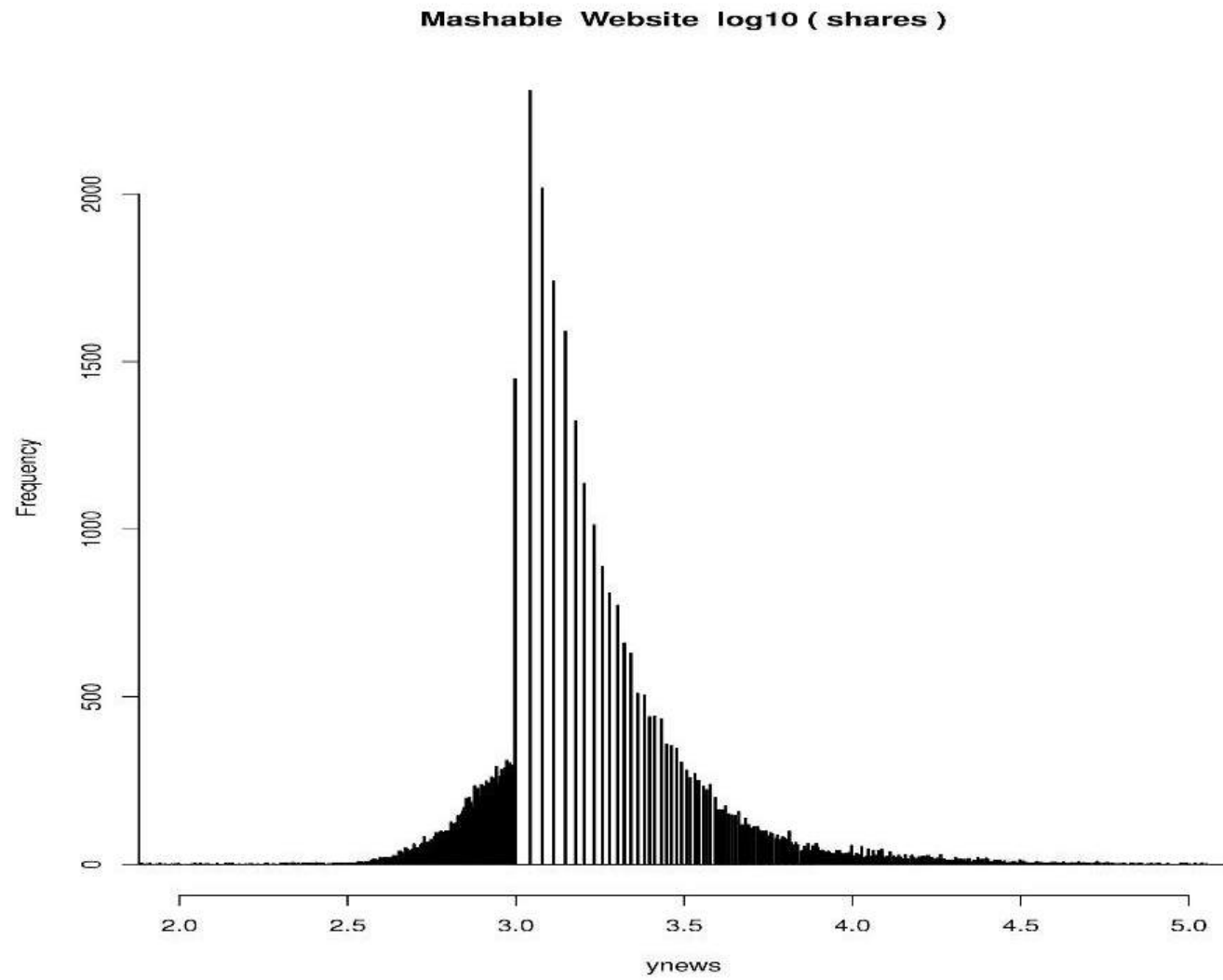
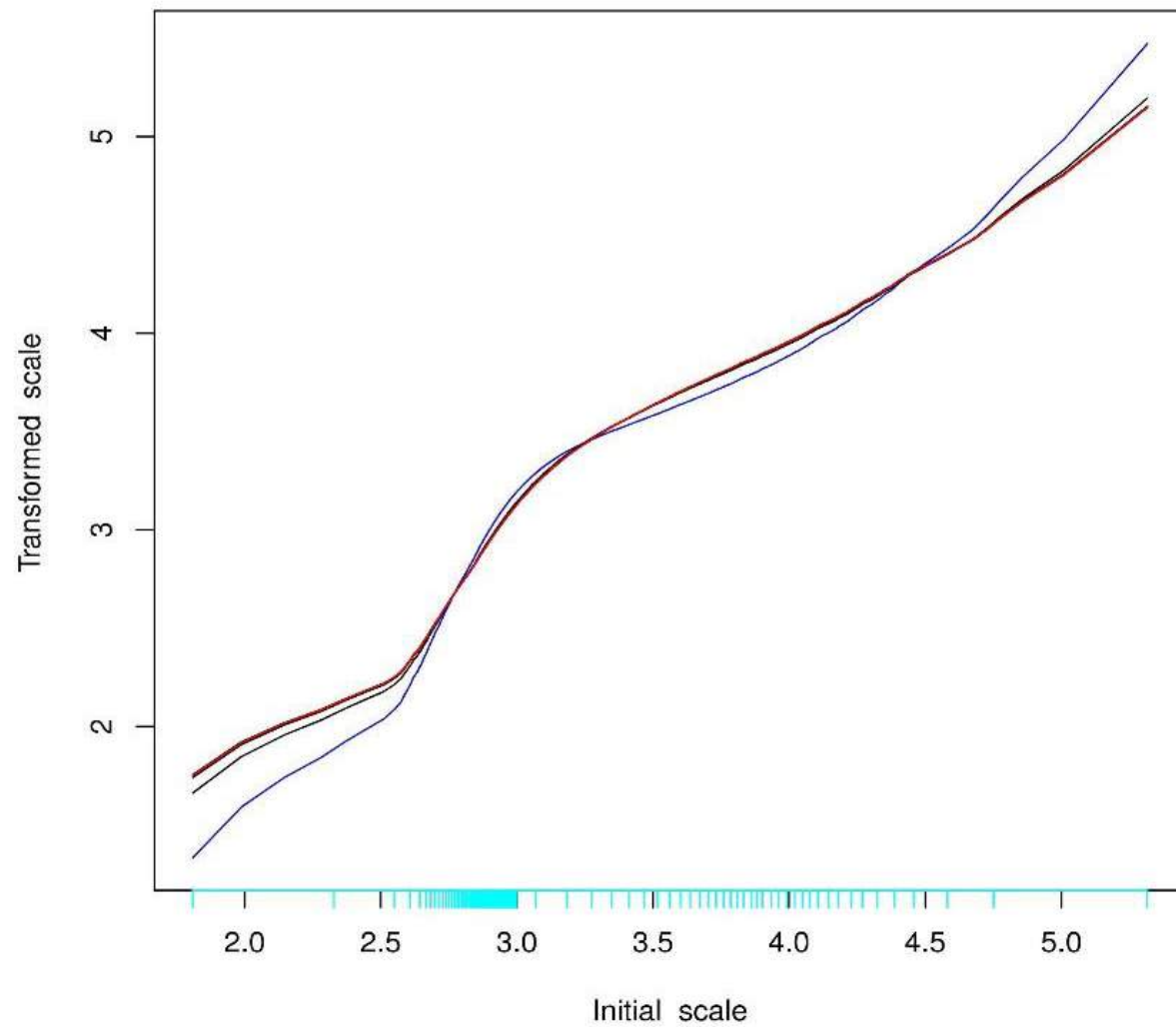
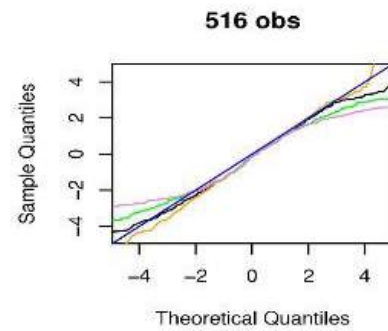
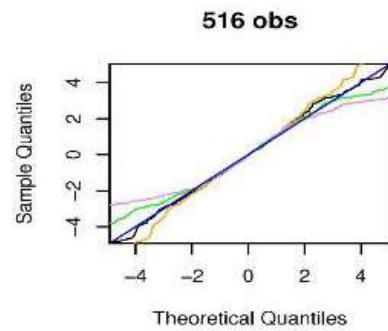
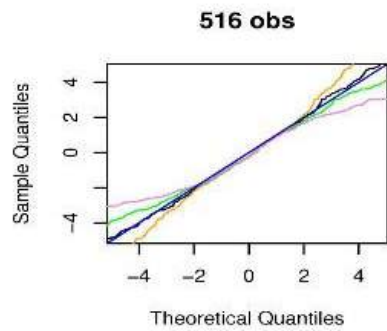
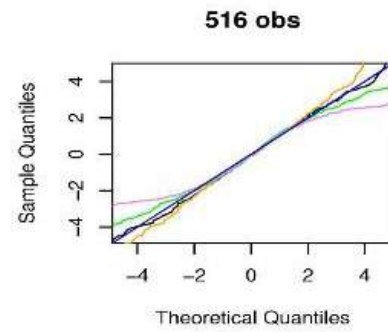
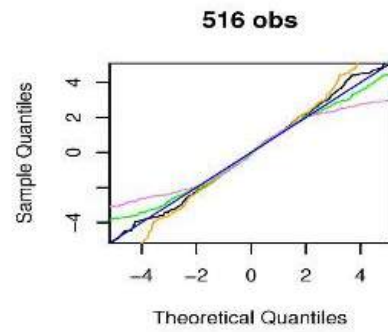
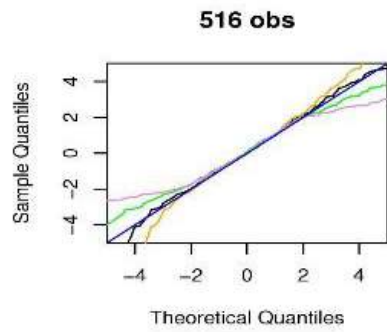
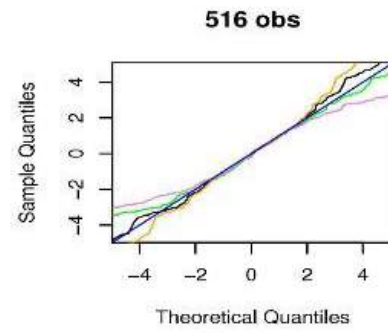
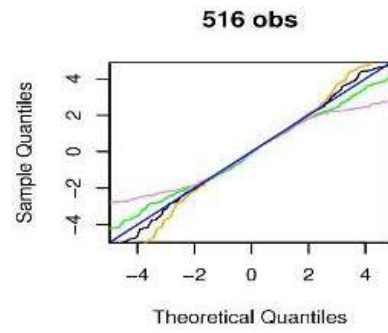
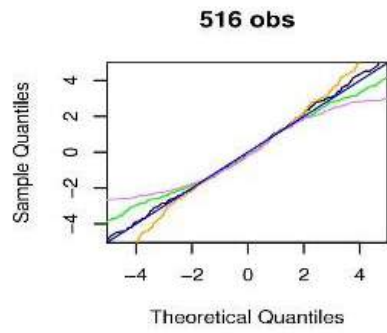


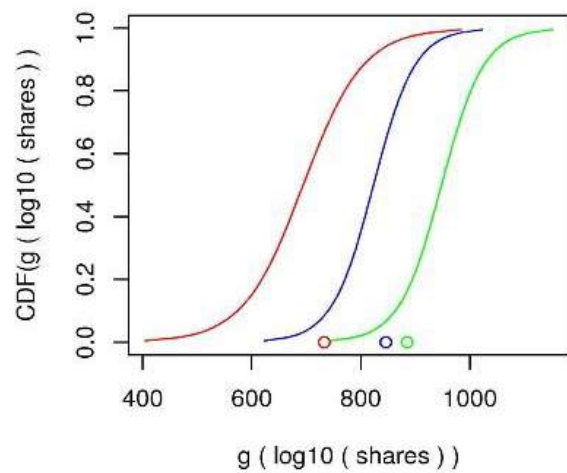
Fig 10

log10 (shares) Transformation Sequence

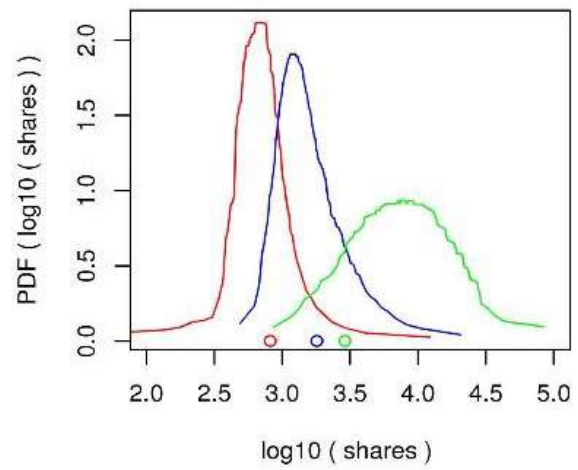
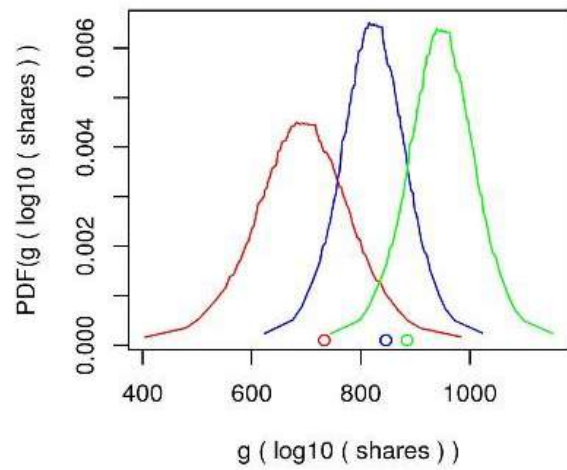
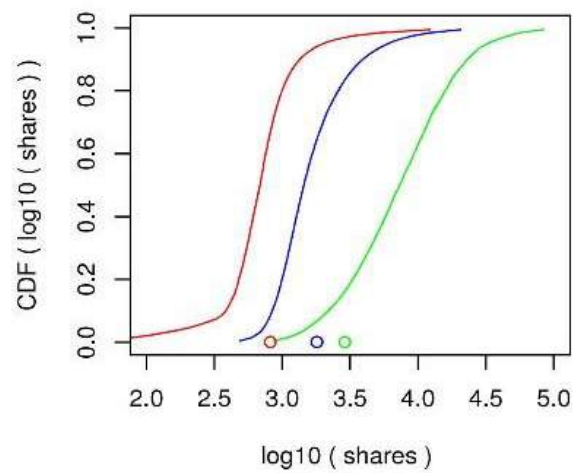




Transformed Outcome



Original Outcome



(1/2) – Million Song Dataset

Song Recordings

(Irvine Repository)

$N(\text{train}) = 463715$ song recordings

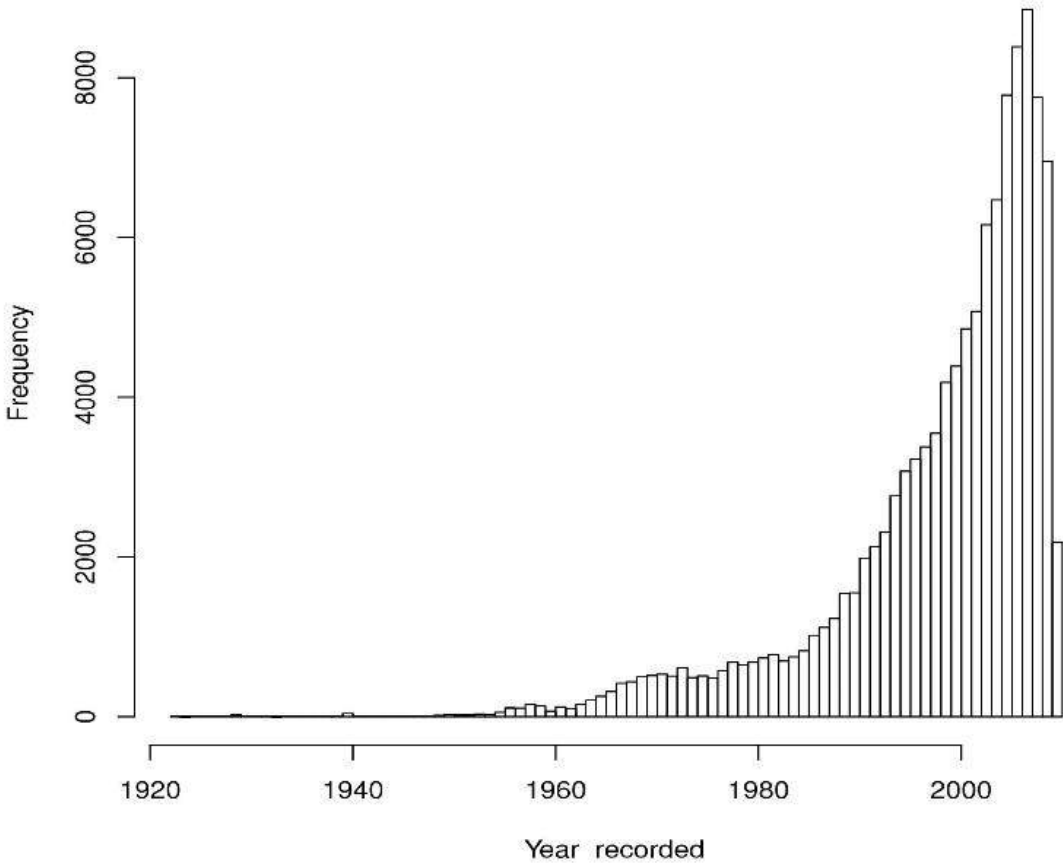
50000 train, 10000 mod. selection

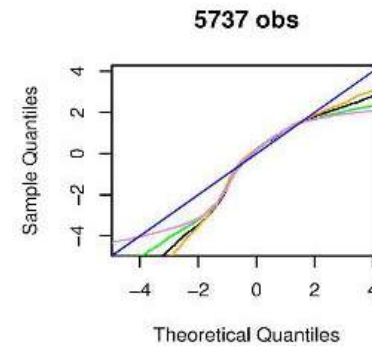
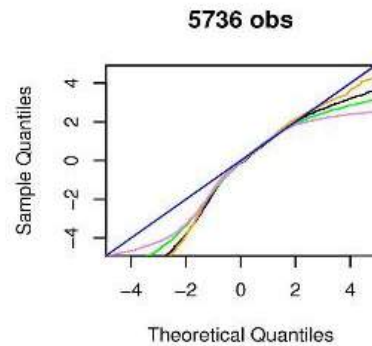
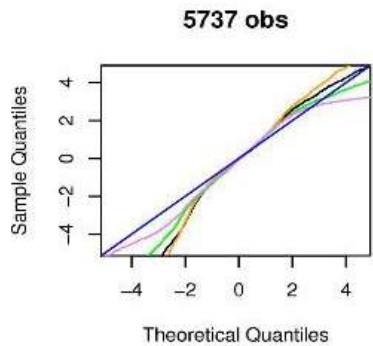
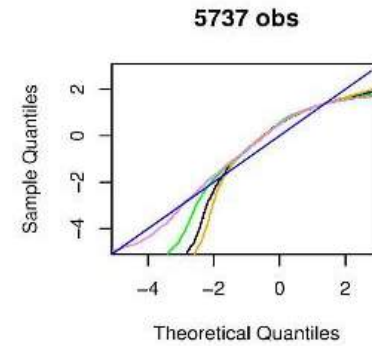
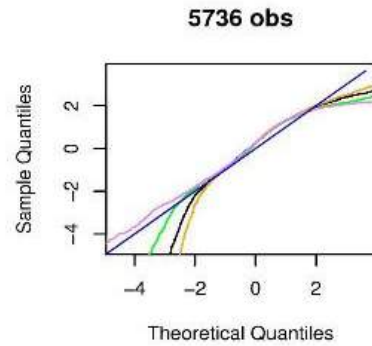
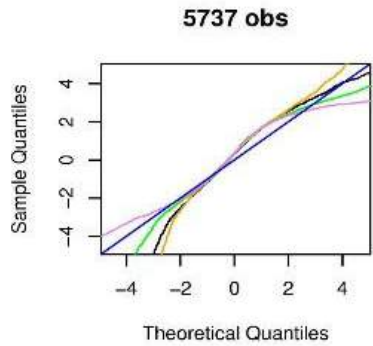
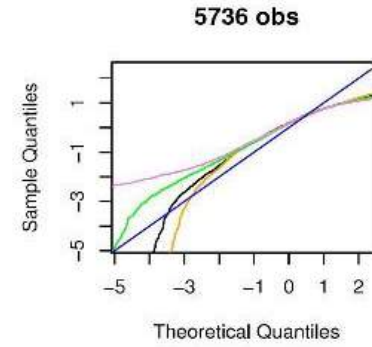
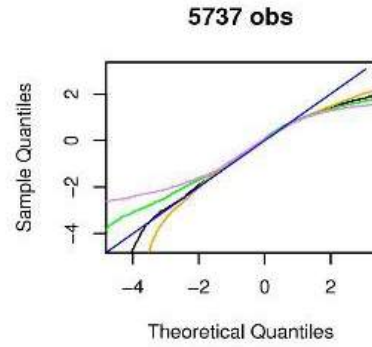
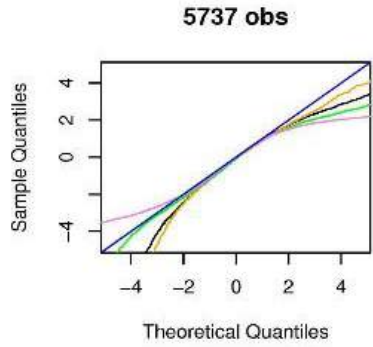
$N(\text{test}) = 51630$

y = year released (1922 – 2011)

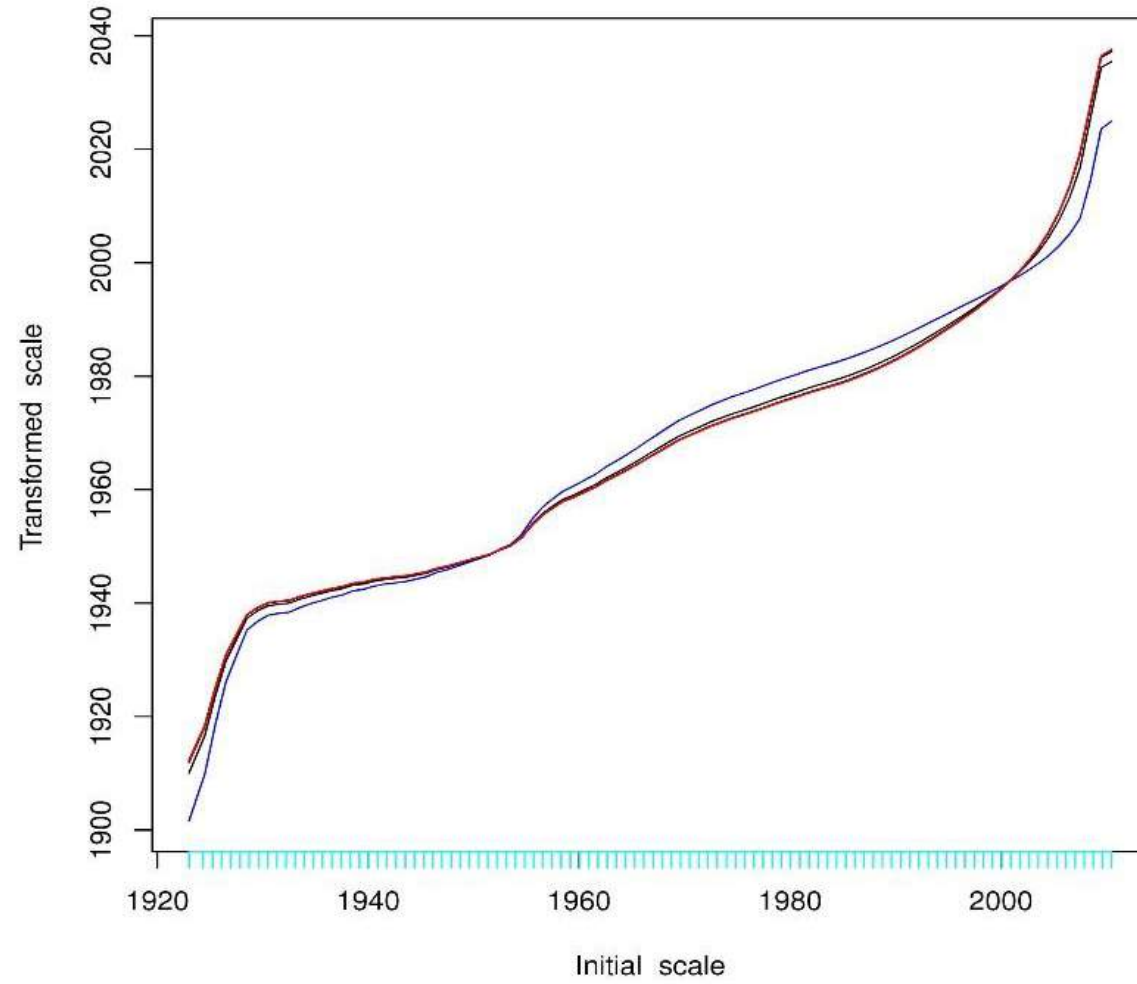
x = 89 acoustic measurements

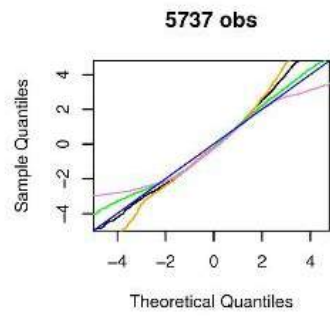
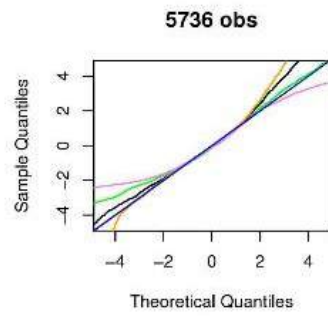
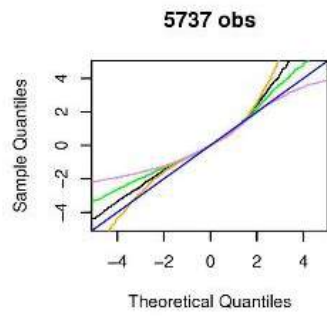
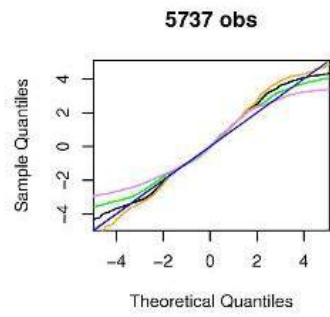
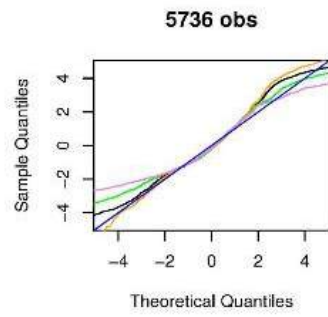
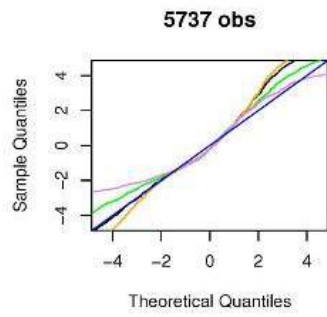
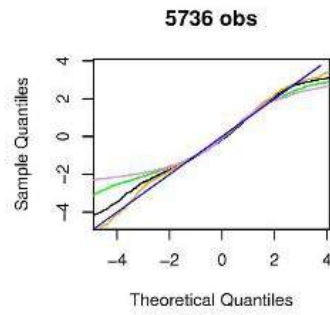
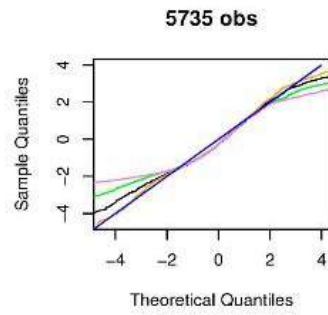
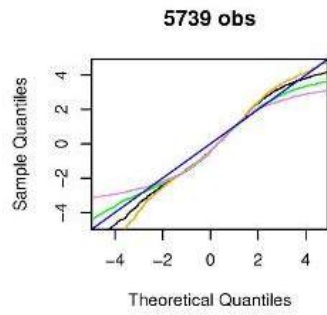
Million Song Data Set

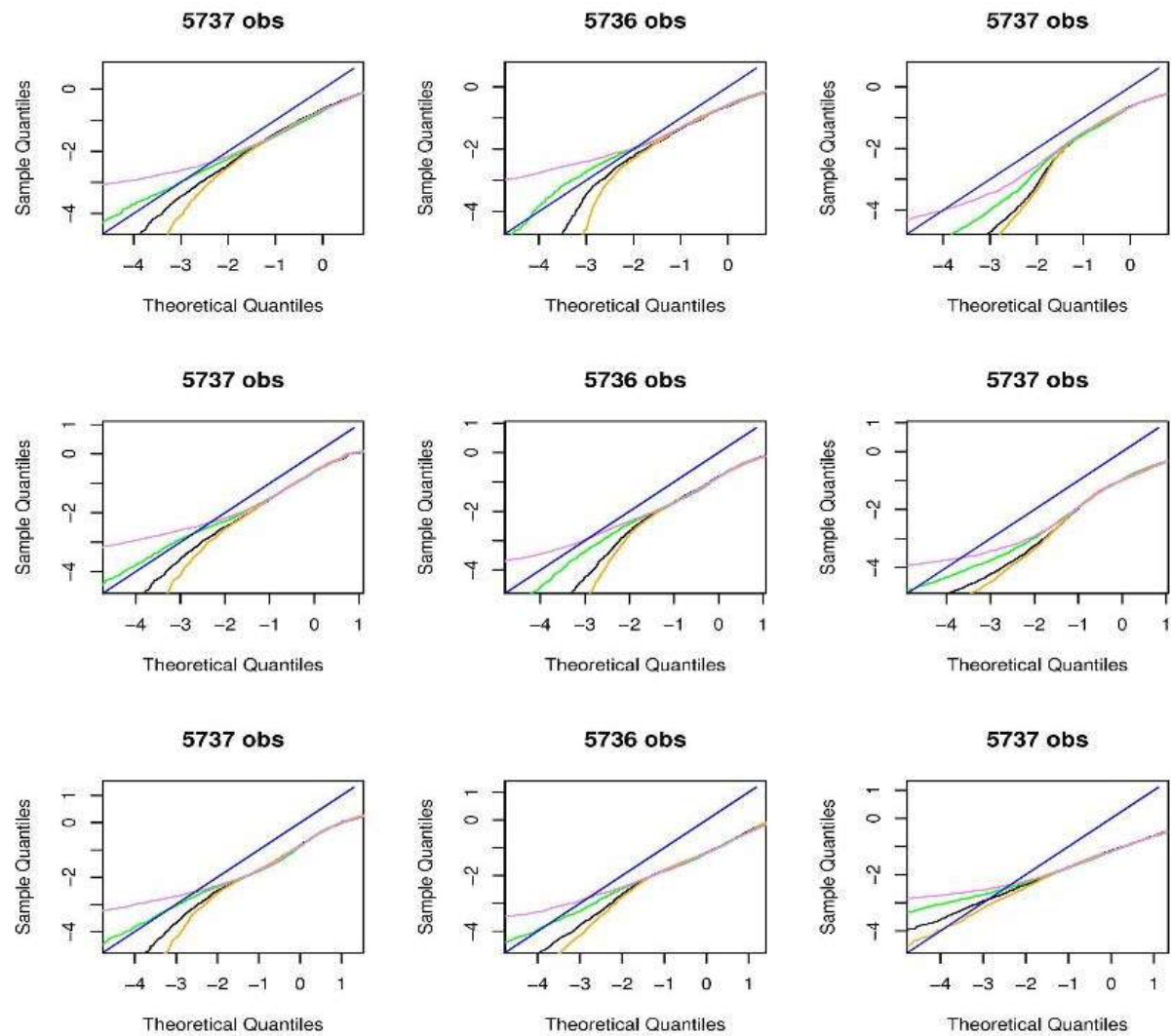




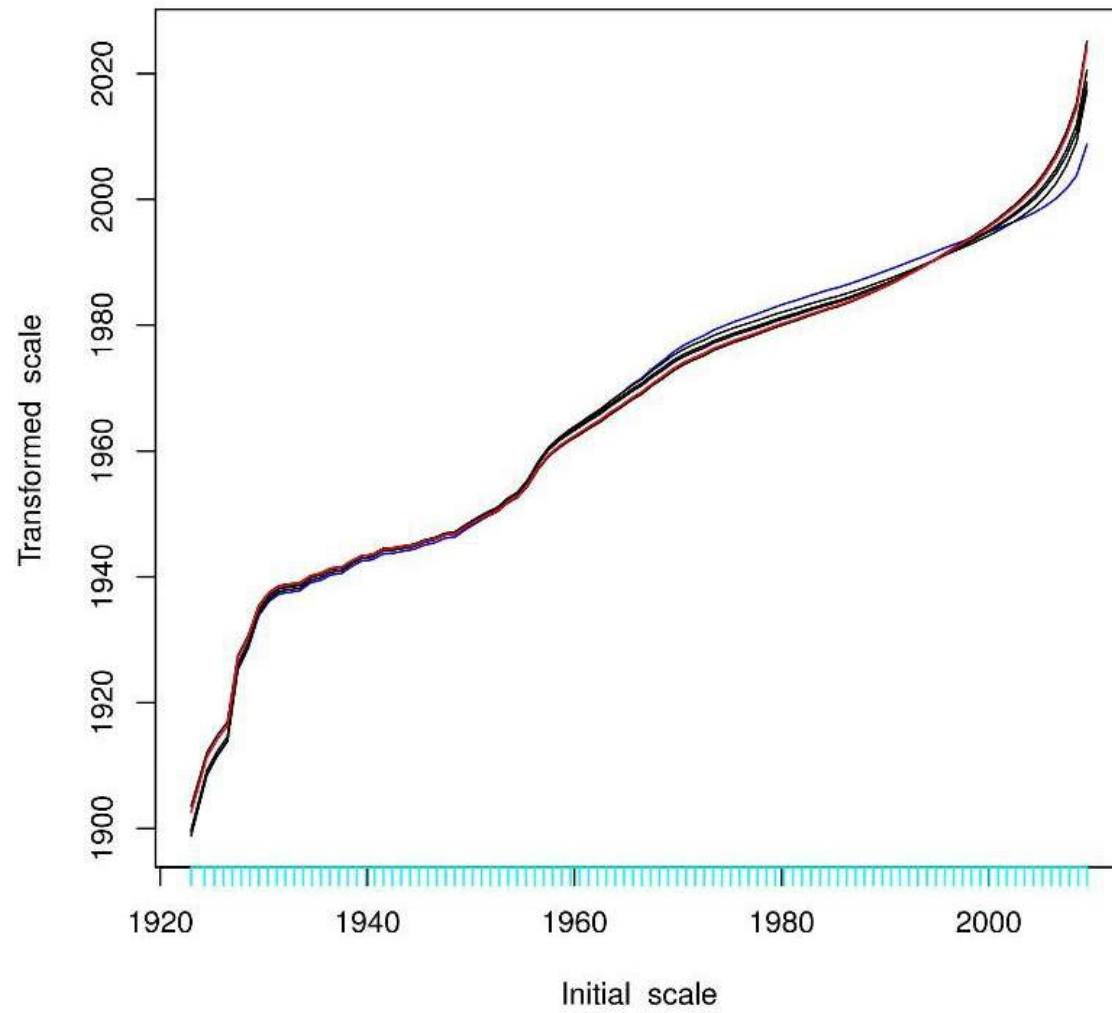
MSD Transformation Iterations

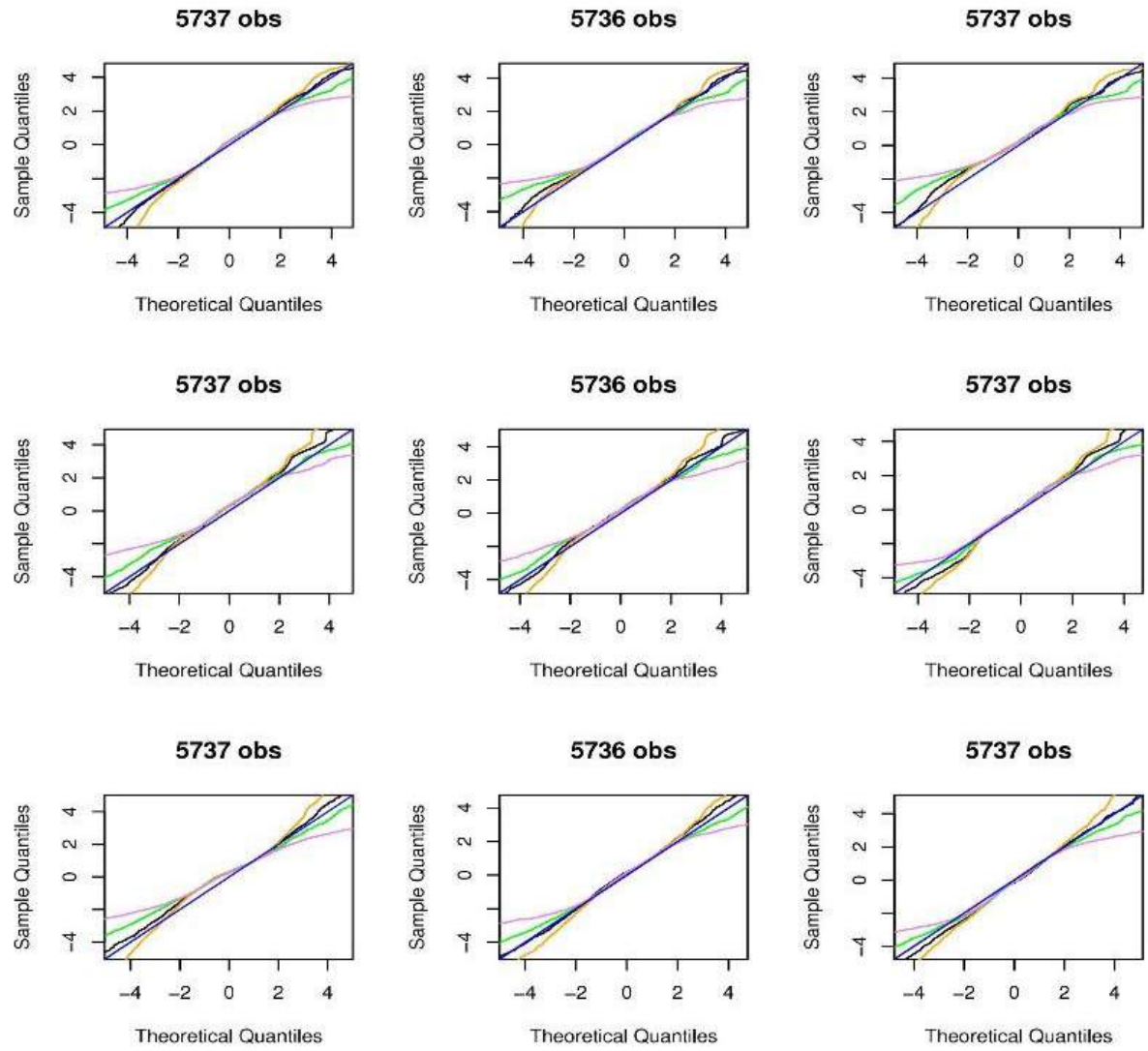




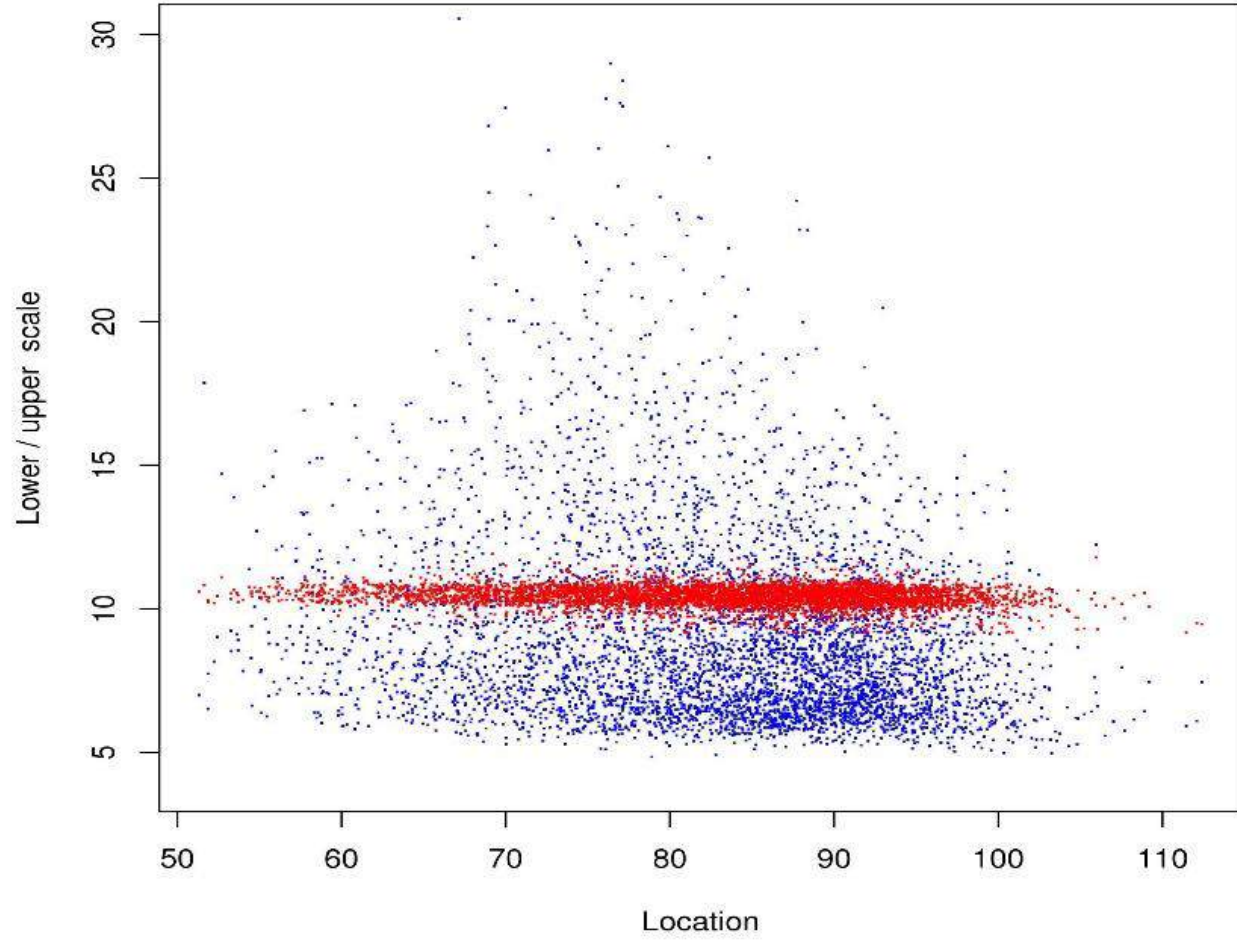


MSD Asymmetric Transformation

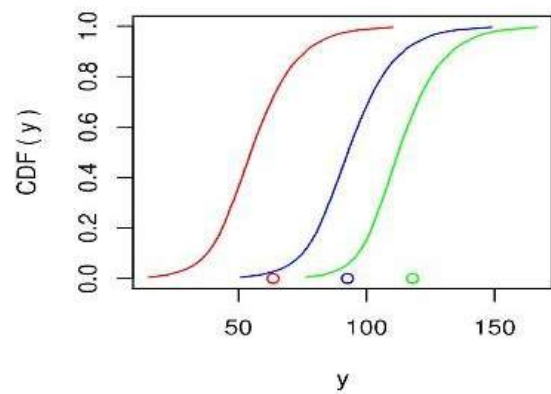




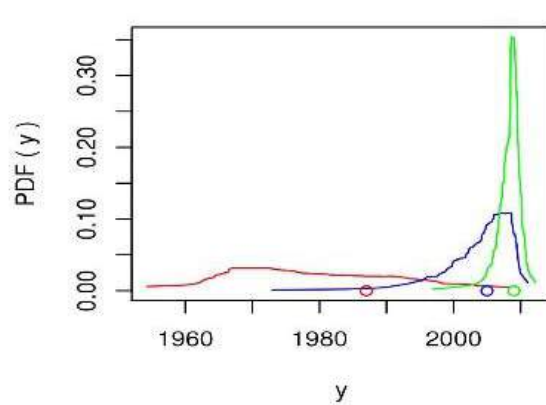
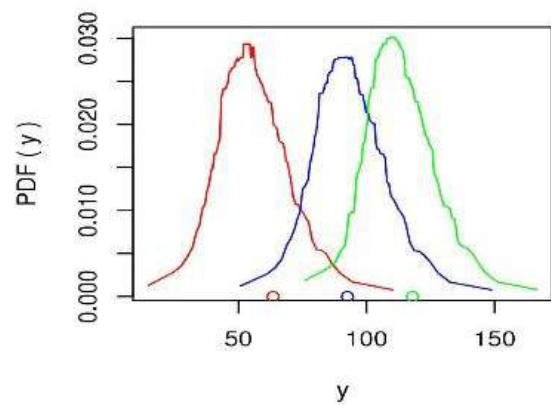
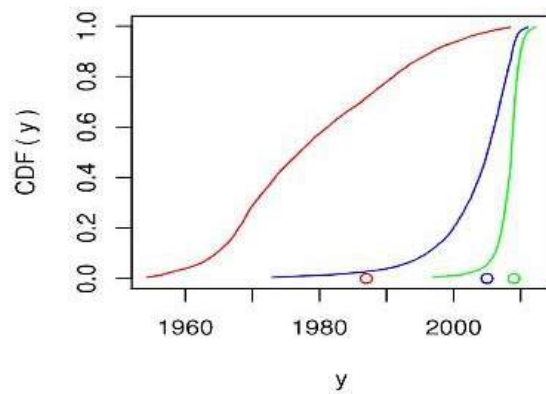
Transformed Asymmetric Error Model



Transformed Outcome



Original Outcome



LESSONS

1. Don't confuse reducible and irreducible error.
2. Tukey was right (normality is rare)
3. Even approximate homoscedasticity is rare
4. Even approximate symmetry is rare
5. Optimal (nonobvious) transformations can help

Slides:

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