

Analytics to Support Business Decision in Traditional Industry

25 May 2014
7th R conference (Beijing, China)



Traditional Industry v.s. E-commerce Industry

- NOT big data
- Different time series
- Well defined concentration
- Relationship: Rep Sales play a role
- Without/little online purchasing data
- Machine learning and statistical model
- Market research based on questionnaire
- Time of the year matters (for agriculture)

How Analytics Can Help

➤ Business Opportunities

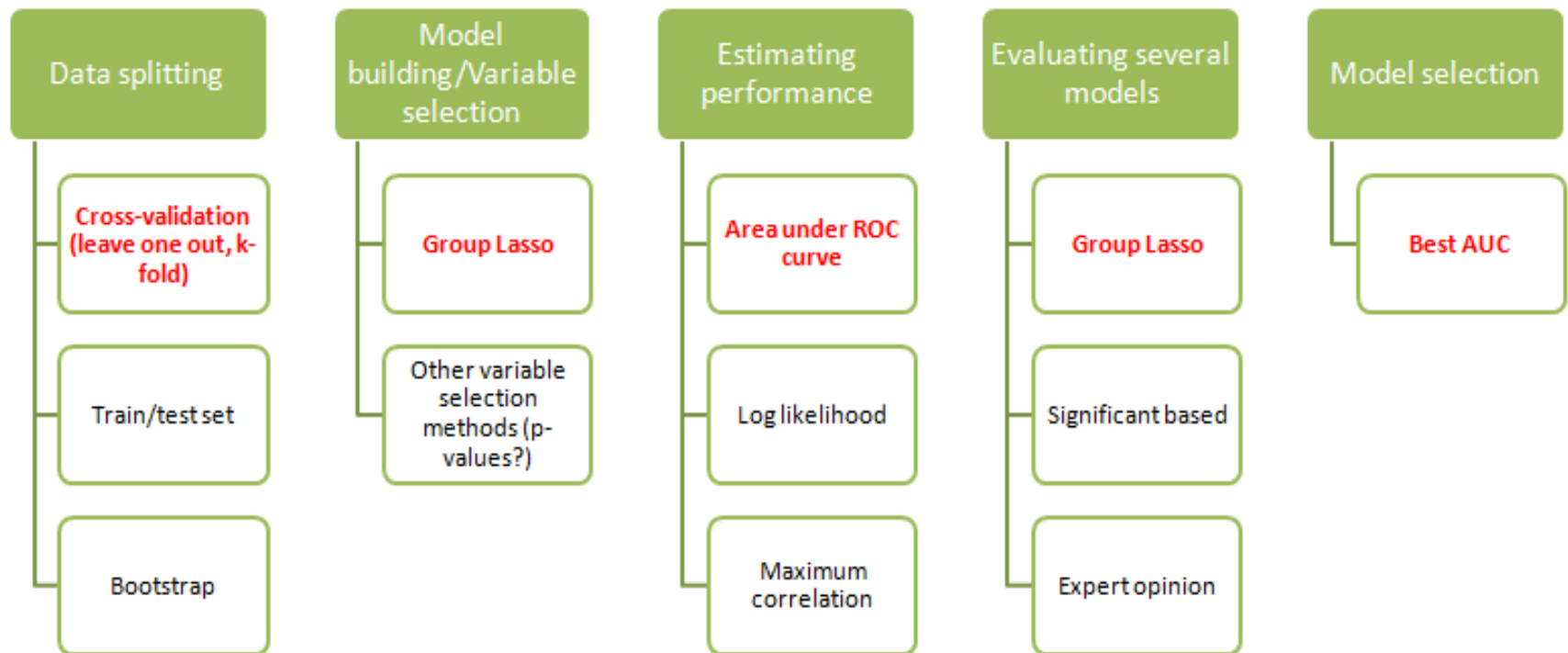
- Cross-selling
- Evaluate marketing campaign
- Target specific population (segmentation)
- R.O.I of different marketing program and services
- Understand customer attrition/conversion/retention
- Predict possible customer reaction before decision implementation
-

Not Academic

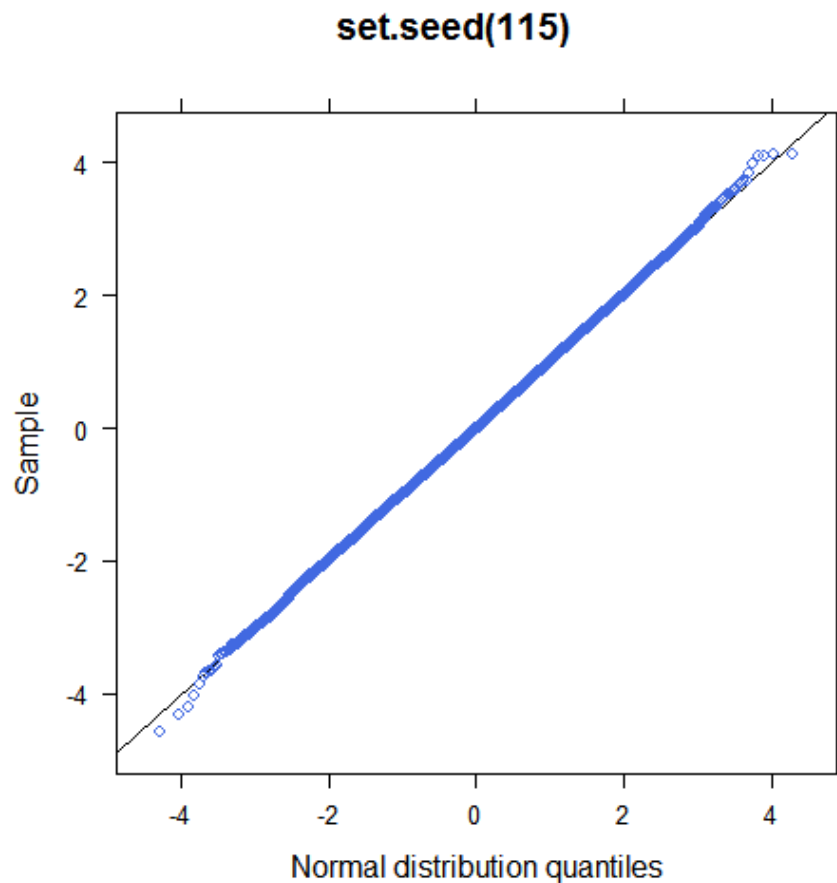
- Dirty data
- Useful v.s. Accurate
- Issues of observational study
 - Hard to define control group
 - Missing value
 - Correlation? Causality?
- Gap between “Implication” and “Implementation”

Scoring System for Customer Retention

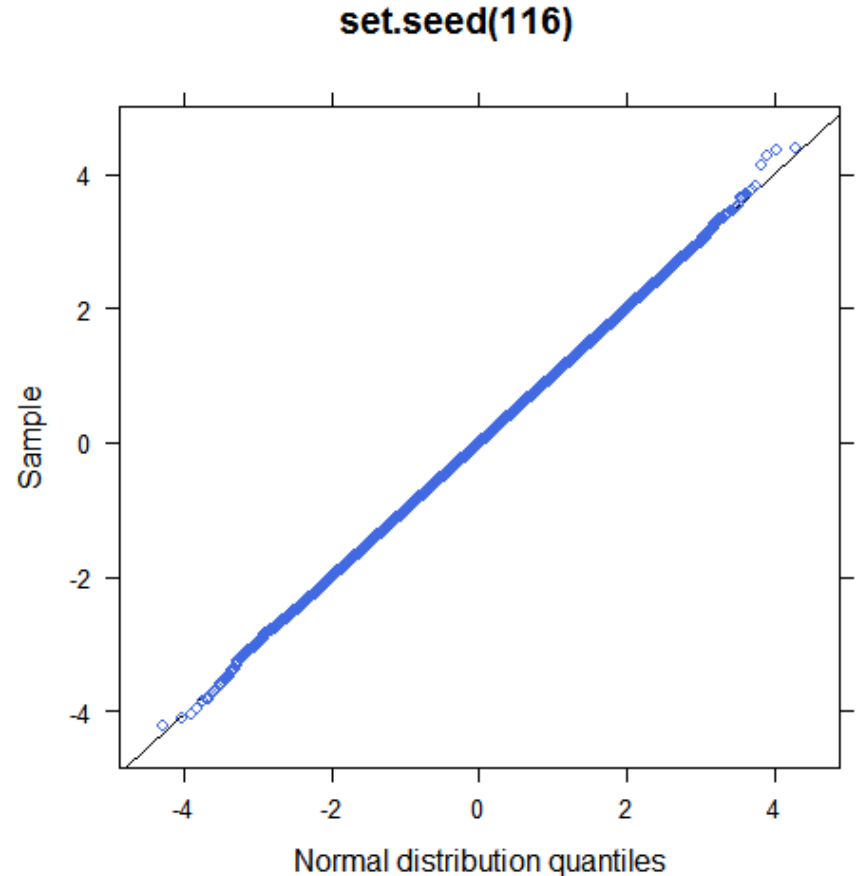
- Objective: score customers by the likelihood of retention
- Response: 0/1
- Predictors: categorical, continuous (different scales, all possible historical transactional records)
- Data Size: 10000 x 200



⁶ 50,000 random t variates, 200 d.f.

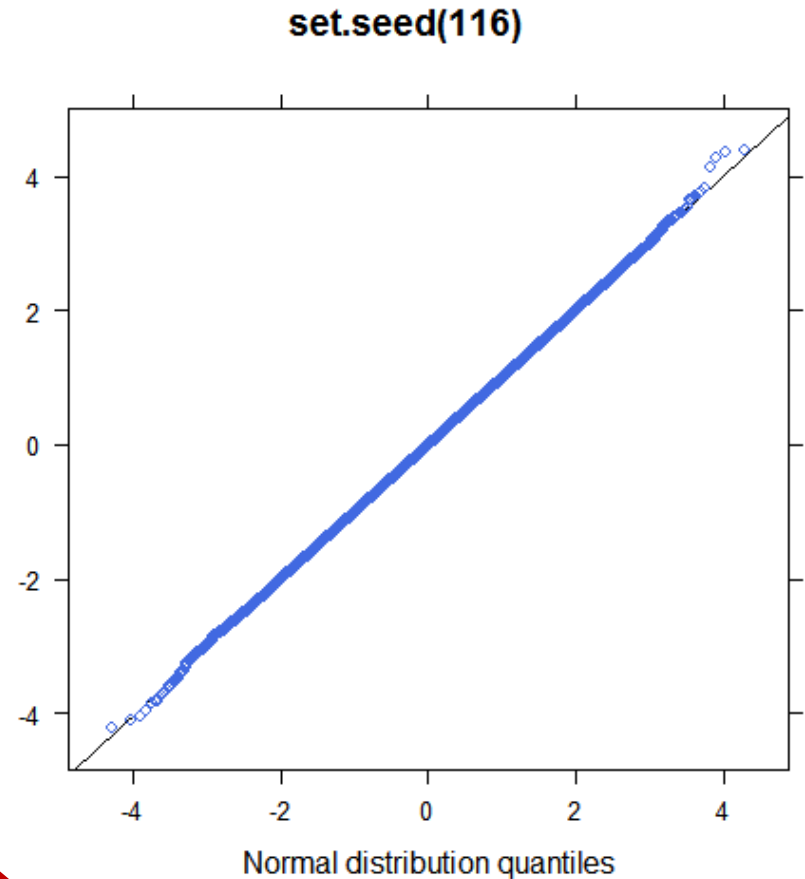
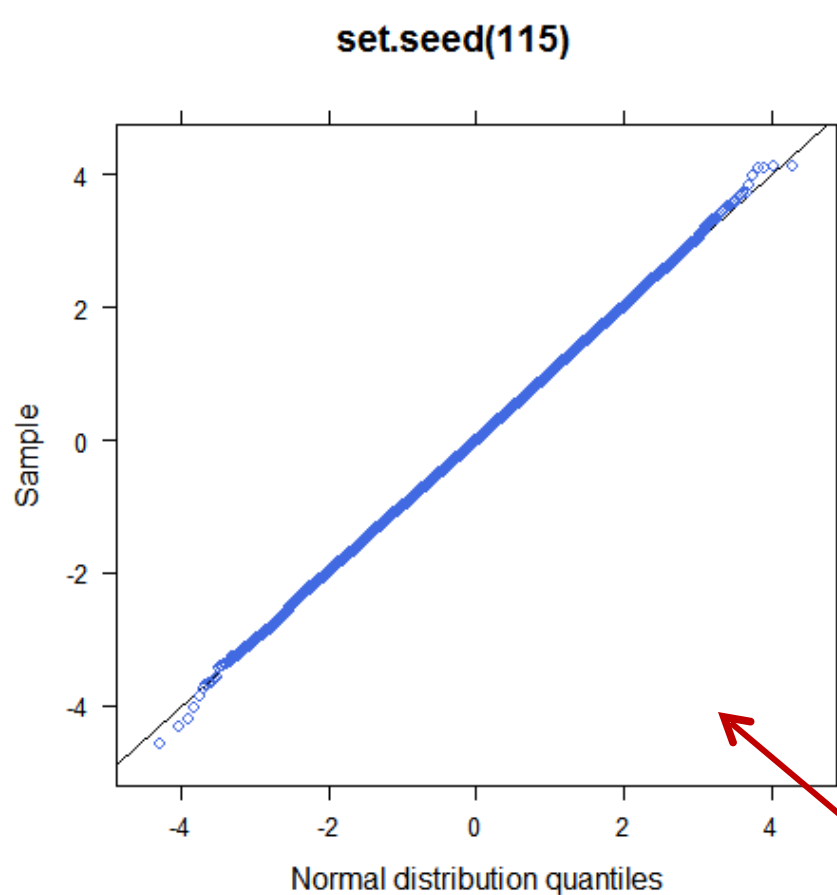


Anderson-Darling normality test
 $A = 1.1849$, $p\text{-value} = 0.004318$



Anderson-Darling normality test
 $A = 0.1731$, $p\text{-value} = 0.9281$

7 50,000 random t variates, 200 d.f.

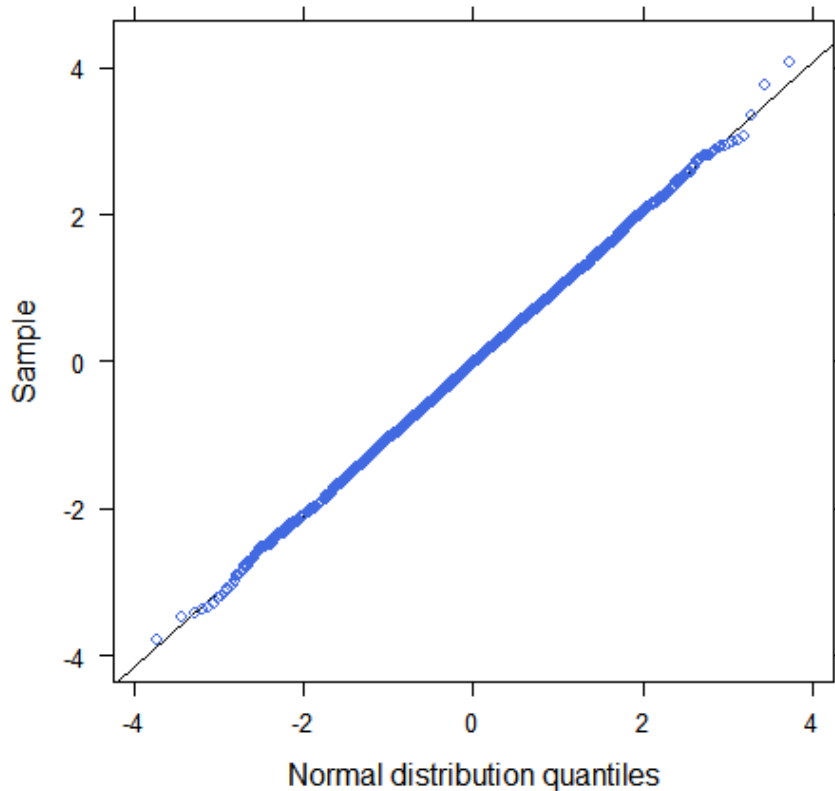


Anderson-Darling normality test
 $A = 1.1849$, $p\text{-value} = 0.004318$

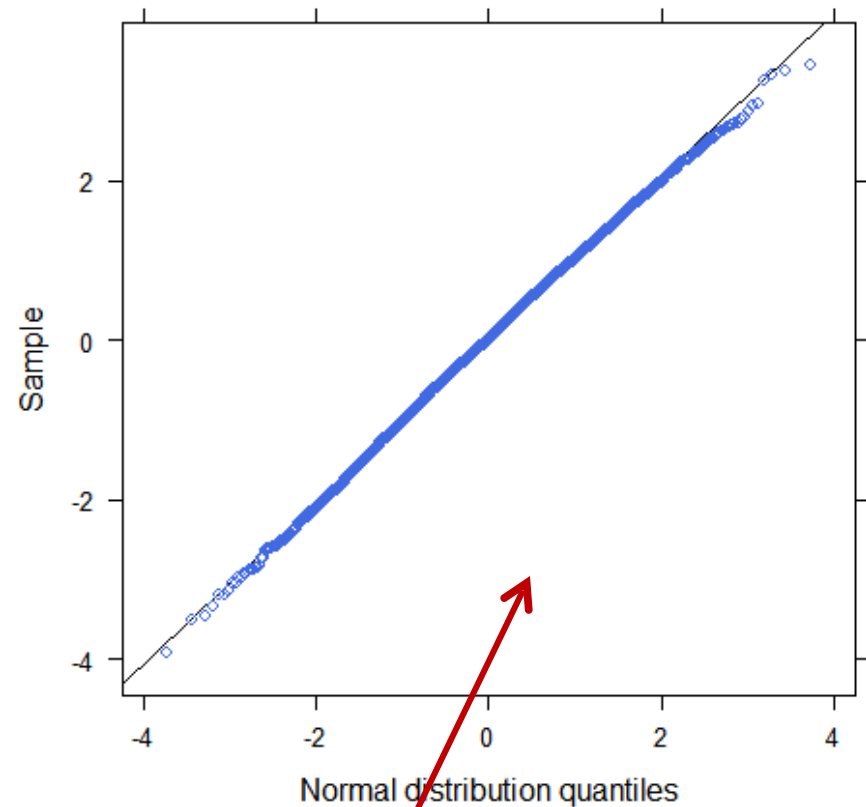
Anderson-Darling normality test
 $A = 0.1731$, $p\text{-value} = 0.9281$

5,000 random t variates, 200 d.f.

set.seed(14)



set.seed(15)

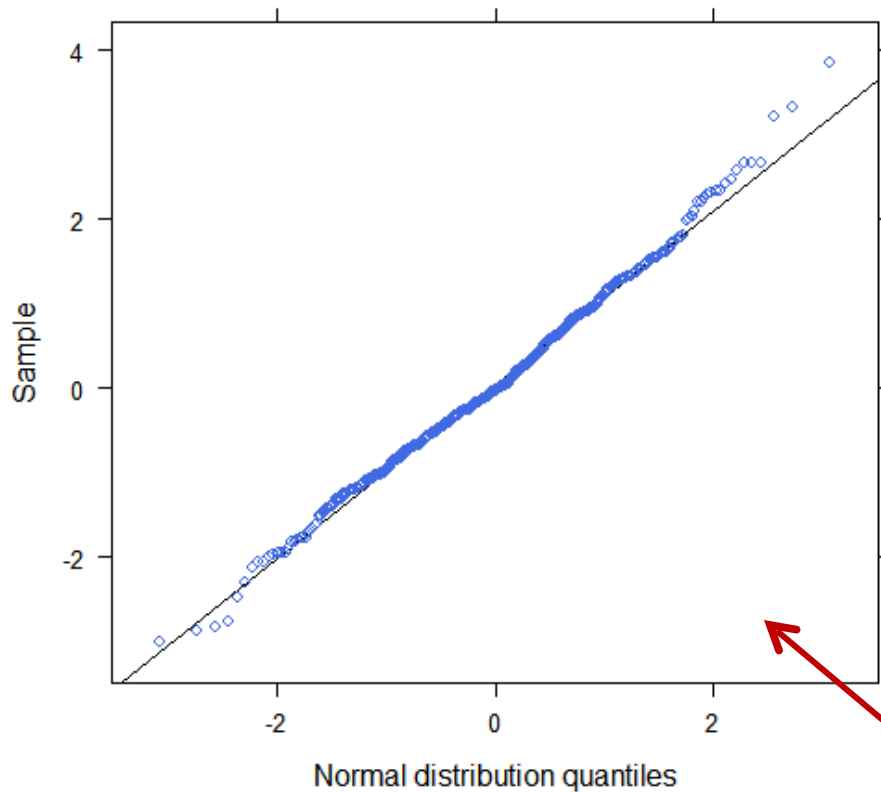


Anderson-Darling normality test
 $A = 0.2676$, $p\text{-value} = 0.6856$

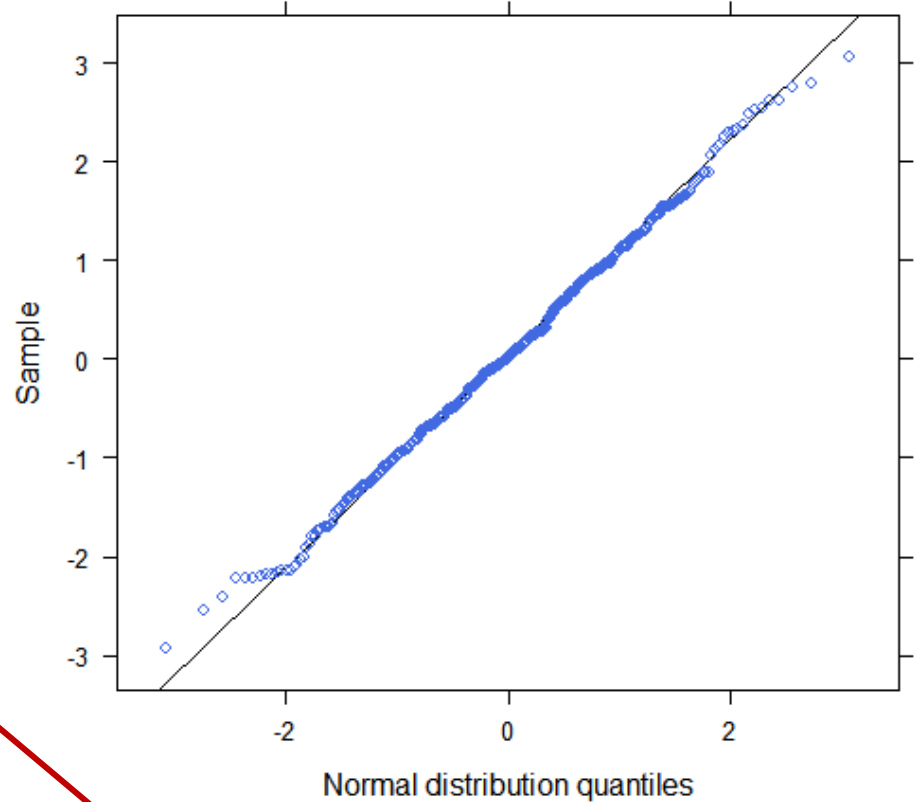
Anderson-Darling normality test
 $A = 1.0479$, $p\text{-value} = 0.009378$

9 500 random t variates, 200 d.f.

set.seed(24)



set.seed(26)

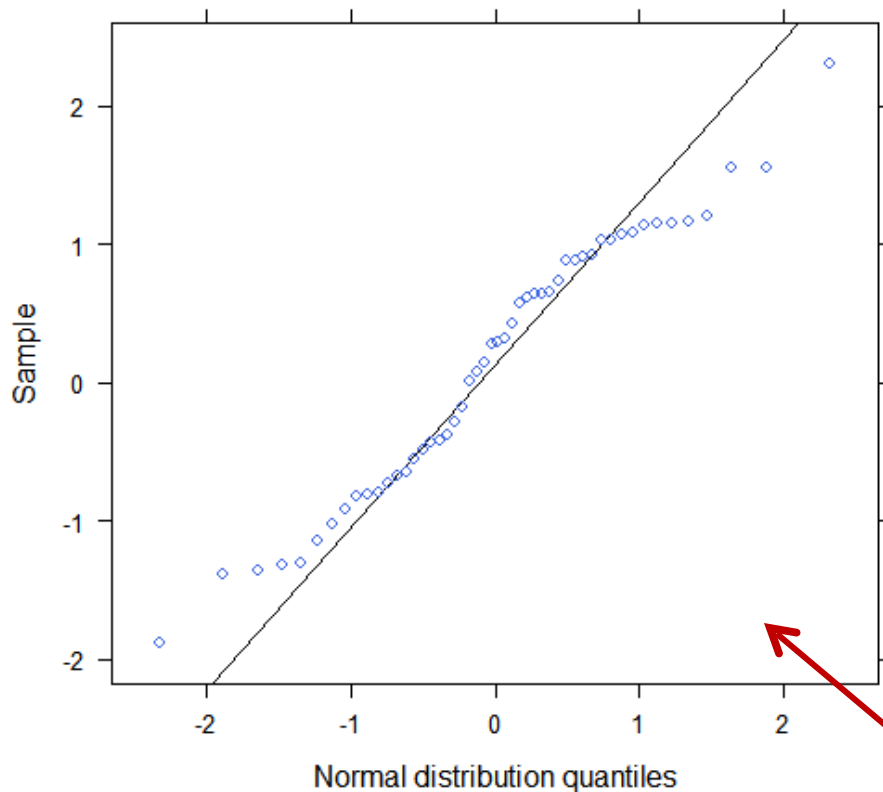


Anderson-Darling normality test
 $A = 0.7488$, $p\text{-value} = 0.05094$

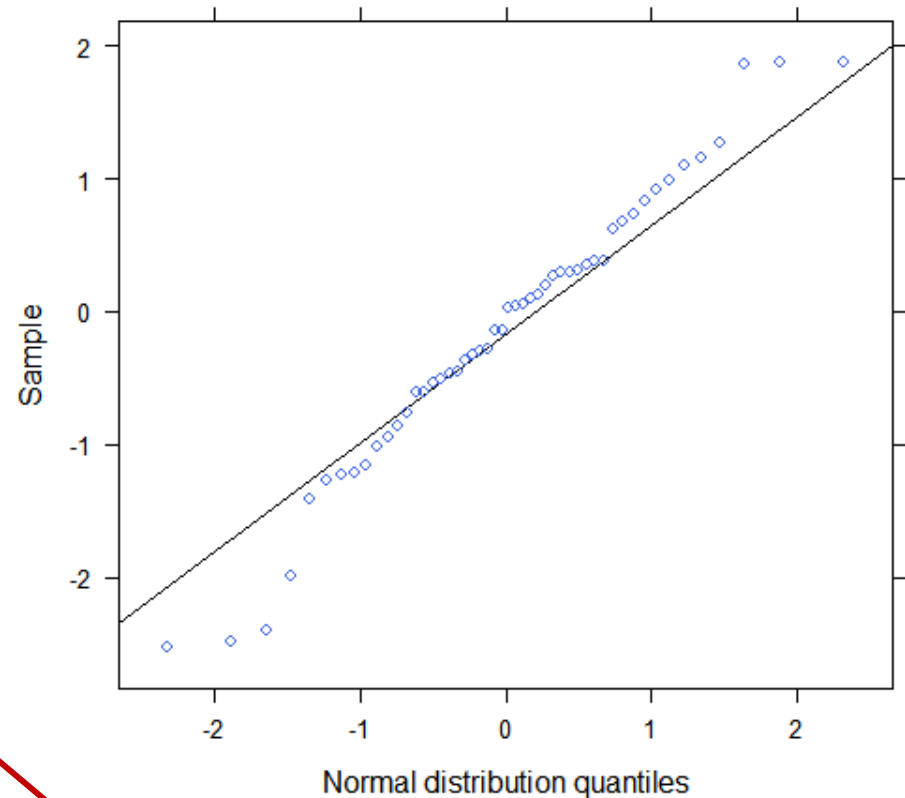
Anderson-Darling normality test
 $A = 0.2427$, $p\text{-value} = 0.7672$

50 random t variates, 200 d.f.

set.seed(183)



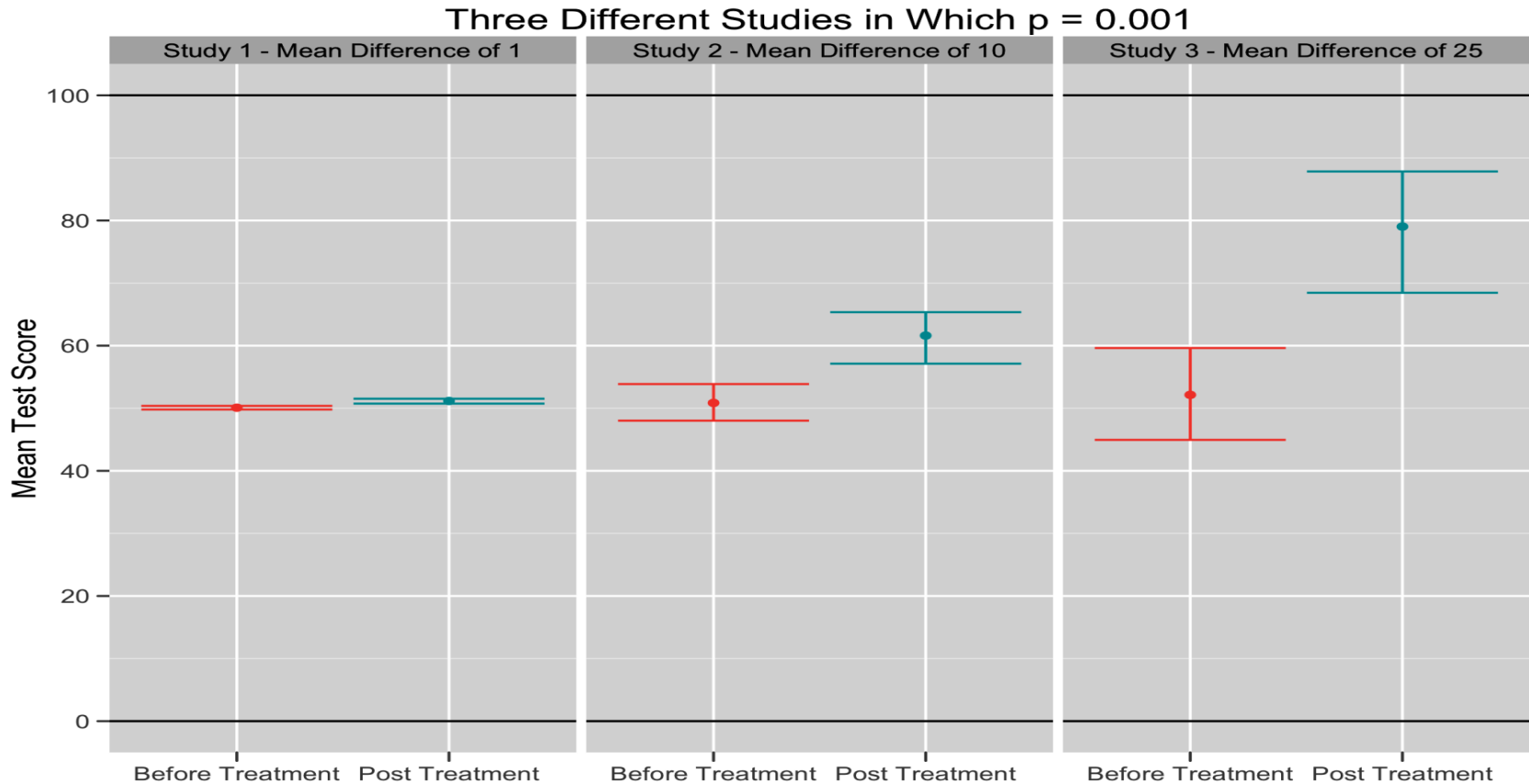
set.seed(140)



Anderson-Darling normality test
 $A = 0.7495$, $p\text{-value} = 0.0477$

Anderson-Darling normality test
 $A = 0.2845$, $p\text{-value} = 0.6156$

Equal p-values are not all equal



Study 1: Non-meaningful difference, precisely measured.

Study 2: Moderate difference, moderate accuracy.

Study 3: Important difference, poorly measured.

P-values collapse two-dimensional data (mean/sd) down to one dimension

12 A Simulation Study

- 800 farms (i.e. $n=800$) and 120 survey questions (i.e. $G=120$) in each dataset
- i^{th} farm is generated from a $Bernoulli(1, p_i)$
- $\ln(\frac{p_i}{1-p_i}) = \beta_0 + \sum_{g=1}^G \mathbf{x}_{i,g}^T \beta_g$, where β_0 is the intercept, $\mathbf{x}_{i,g}$ is a three dimensional indication vector for question answer and β_g
- $\gamma = 0.1, 0.25, 0.5, 1$ and 2 and $\beta_0 = -\frac{40}{3}\gamma$

$$\beta_g = (1, 0, -1) \times \gamma, g = 1, \dots, 40,$$

$$\beta_g = (1, 0, 0) \times \gamma, g = 41, \dots, 80.$$

$$\beta_g = (0, 0, 0) \times \gamma, g = 81, \dots, 120.$$

Simulation Study Result

Simulation study result with various values of coefficient γ . Reported are mean and standard deviation of AUC for both methods and mean difference

Coefficient γ	Group Lasso (mean \pm sd)	Logistic Regression (mean \pm sd)
0.1	0.57 \pm 0.03	0.54 \pm 0.06
0.25	0.71 \pm 0.02	0.64 \pm 0.04
0.5	0.91 \pm 0.03	0.78 \pm 0.03
1	0.92 \pm 0.01	0.82 \pm 0.02
2	0.95 \pm 0.01	0.84 \pm 0.02

Challenges

- Communication
- Approve or change
- 360 degree data
- “Small p” and “Big n”
- Actionable questions
- Short term v.s long term
- Development to implementation



Snedecor Hall

Statistics Department
Statistical Laboratory
Center for Survey Statistics and
Methodology

Donated by Class of 1984



DuPont Pioneer



谢谢！
Thank you!