Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Regression Fit Diagnostics Using freqparcoord

Norm Matloff and Yingkang Xie
University of California at Davis

e-mail: matloff@cs.ucdavis.edu
R/stat blog: matloff.wordpress.com

useR! 2014
UCLA
July 1, 2014

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Intro to freqparcoord

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Intro to freqparcoord

Overview of **freqparcoord**:

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Intro to freqparcoord

Overview of **freqparcoord**:

- Available on CRAN.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Intro to freqparcoord

Overview of **freqparcoord**:

- Available on CRAN.

- New approach to the parallel coordinates data
  visualization method.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Intro to freqparcoord

Overview of **freqparcoord**:

- Available on CRAN.

- New approach to the parallel coordinates data
  visualization method. (Examples presented shortly.)

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Intro to freqparcoord

Overview of **freqparcoord**:

- Available on CRAN.

- New approach to the parallel coordinates data visualization method. (Examples presented shortly.)

- Can also be used for hunting outliers, clusters...

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Intro to freqparcoord

Overview of **freqparcoord**:

- Available on CRAN.

- New approach to the parallel coordinates data visualization method. (Examples presented shortly.)

- Can also be used for hunting outliers, clusters...

- and for regression diagnostics—our topic here.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# What Is Parallel Coordinates
Visualization?

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# What Is Parallel Coordinates
# Visualization?

- Very old idea.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
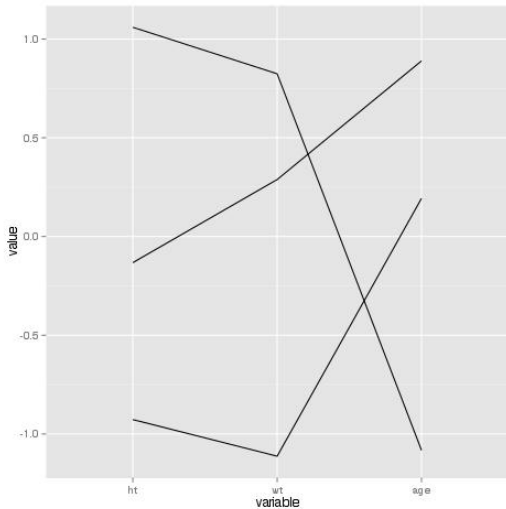loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# What Is Parallel Coordinates Visualization?

- Very old idea.
- If have k variables, draw k vertical axes.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# What Is Parallel Coordinates
# Visualization?

- Very old idea.

- If have k variables, draw k vertical axes. Each data point
  is a polygonal line connecting the value of each variable.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Example

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Example

Example: Height/weight/age data.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Example

Example: Height/weight/age data.

```
> d
  ht  wt age
1 71 175  25
2 66 128  36
3 68 162  42
```

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Problems with Parallel Coordinates

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

- Various solutions, e.g. making the lines fainte, or combining themr

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

- Various solutions, e.g. making the lines fainte, or combining themr

- Our solution: Plot only a few "typical" lines,

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

- Various solutions, e.g. making the lines fainte, or combining themr

- Our solution: Plot only a few "typical" lines, based on estimated multivariate density.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

- Various solutions, e.g. making the lines fainte, or combining themr

- Our solution: Plot only a few "typical" lines, based on estimated multivariate density.

- E.g., height/weight/age:

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

- Various solutions, e.g. making the lines fainte, or combining themr

- Our solution: Plot only a few "typical" lines, based on estimated multivariate density.

- E.g., height/weight/age:

  - What height/weight/age combinations are typical **overall**?

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

- Various solutions, e.g. making the lines fainte, or combining themr

- Our solution: Plot only a few "typical" lines, based on estimated multivariate density.

- E.g., height/weight/age:

  - What height/weight/age combinations are typical **overall**?
  - What height/weight/age combinations are typical within groups?

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

- Various solutions, e.g. making the lines fainte, or combining themr

- Our solution: Plot only a few "typical" lines, based on estimated multivariate density.

- E.g., height/weight/age:

  - What height/weight/age combinations are typical **overall**?
  - What height/weight/age combinations are typical within groups? **Group comparison.**

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis


e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

- Various solutions, e.g. making the lines fainte, or combining themr

- Our solution: Plot only a few "typical" lines, based on estimated multivariate density.

- E.g., height/weight/age:

  - What height/weight/age combinations are typical **overall**?
  - What height/weight/age combinations are typical within groups? **Group comparison.**
  - What height/weight/age combinations are rare?

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

- Various solutions, e.g. making the lines fainte, or combining themr

- Our solution: Plot only a few "typical" lines, based on estimated multivariate density.

- E.g., height/weight/age:

  - What height/weight/age combinations are typical **overall**?
  - What height/weight/age combinations are typical within groups? **Group comparison.**
  - What height/weight/age combinations are rare? **Outlier hunting.**

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

- Various solutions, e.g. making the lines fainte, or combining themr

- Our solution: Plot only a few "typical" lines, based on estimated multivariate density.

- E.g., height/weight/age:

    - What height/weight/age combinations are typical **overall**?
    - What height/weight/age combinations are typical within groups? **Group comparison.**
    - What height/weight/age combinations are rare? **Outlier hunting.**
    - What height/weight/age combinations are "locally typical"?

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

- Various solutions, e.g. making the lines fainte, or combining themr

- Our solution: Plot only a few "typical" lines, based on estimated multivariate density.

- E.g., height/weight/age:

  - What height/weight/age combinations are typical **overall**?
  - What height/weight/age combinations are typical within groups? **Group comparison.**
  - What height/weight/age combinations are rare? **Outlier hunting.**
  - What height/weight/age combinations are "locally typical"? **Cluster hunting.**

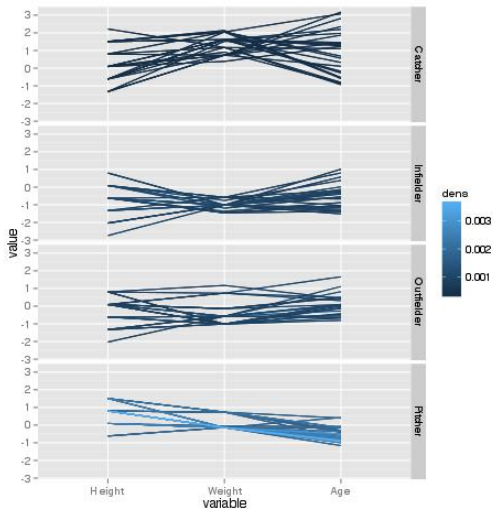Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

- Various solutions, e.g. making the lines fainte, or combining themr

- Our solution: Plot only a few "typical" lines, based on estimated multivariate density.

- E.g., height/weight/age:

  - What height/weight/age combinations are typical **overall**?
  - What height/weight/age combinations are typical within groups? **Group comparison.**
  - What height/weight/age combinations are rare? **Outlier hunting.**
  - What height/weight/age combinations are "locally typical"? **Cluster hunting.**

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
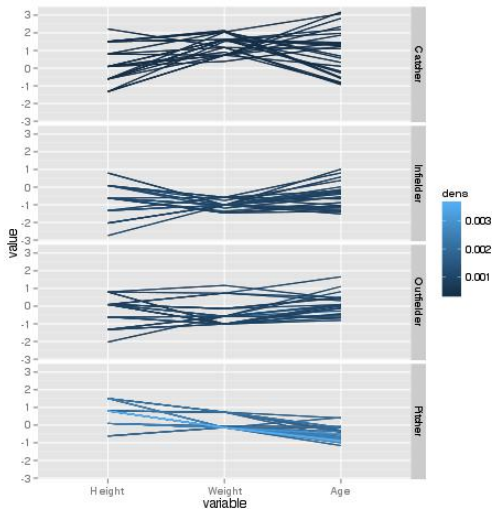mat-
loff.wordpress.com

# UCLA Baseball Player Data

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# UCLA Baseball Player Data

Most typical 25 points for each playing position.



- Catchers heavier,
  vary widely in
  height and age.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com
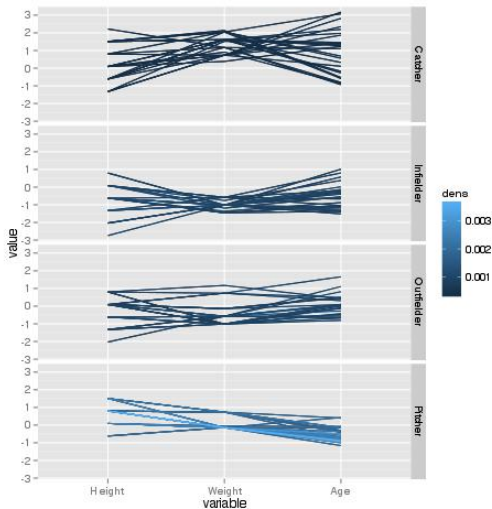
# UCLA Baseball Player Data

Most typical 25 points for each playing position.



- Catchers heavier, vary widely in height and age.

- Pitchers tall, lighter, less variable in age.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# UCLA Baseball Player Data

Most typical 25 points for each playing position.



- Catchers heavier, vary widely in height and age.
- Pitchers tall, lighter, less variable in age.
- Infielders vary considerably in height but not weight.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Application to Regression
# Diagnostics

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Application to Regression
# Diagnostics

Our **freqparcoord** package includes a function **regdiag()**.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Application to Regression Diagnostics

Our **freqparcoord** package includes a function **regdiag()**.

- Focused vertical axis:

```
divergences = fitted parametric model -
              fitted nonparametric model
```

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Application to Regression
# Diagnostics

Our **freqparcoord** package includes a function **regdiag()**.

- Focused vertical axis:

  ```
  divergences = fitted parametric model -
                fitted nonparametric model
  ```

  (Uses k-NN for nonparametric est..)

- The divergences are NOT the residuals (i.e. not actual -
  fitted parametric).

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Application to Regression
# Diagnostics

Our **freqparcoord** package includes a function **regdiag()**.

- Focused vertical axis:

  ```
  divergences = fitted parametric model -
                fitted nonparametric model
  ```

  (Uses k-NN for nonparametric est..)

- The divergences are NOT the residuals (i.e. not actual -
  fitted parametric).

- What **regdiag()** does it look at the typical values among
  the most negative and most positive divergences.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Application to Regression Diagnostics

Our **freqparcoord** package includes a function **regdiag()**.

- Focused vertical axis:

  ```
  divergences = fitted parametric model -
                fitted nonparametric model
  ```

  (Uses k-NN for nonparametric est..)

- The divergences are NOT the residuals (i.e. not actual - fitted parametric).

- What **regdiag()** does it look at the typical values among the most negative and most positive divergences.

- In other words: **regdiag()** asks, "In what region[s] of predictor space is the fit poorer?"

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Example

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Example

Programmers and engineers in Silicon Valley, 2000 Census, 5%
PUMS.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: matloff@cs.ucdavis.edu
R/stat blog:
matloff.wordpress.com

## Example

Programmers and engineers in Silicon Valley, 2000 Census, 5%
PUMS.

```
> data(prgeng)  # fpc. built-in data set
> pg1 <- prgeng
> pg1$ms <- as.integer(pg1$educ == 14)  # MS
> pg1$phd <- as.integer(pg1$educ == 16)  # PhD
> pg1$se <- as.integer(pg1$occ==102)  # s. eng
> l1 <- lm(wageinc ~ age+ms+phd+se+sex, data=pg1
# look at 40% most neg., 40% most pos. divs.
> p <- regdiag(l1, tail=0.40)
> p  # display graph
> p$paramr2  # parametric adj. R2
[1] 0.07027561
> p$nonparamr2  # nonparamr2 R2
[1] 0.1286746
```

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com
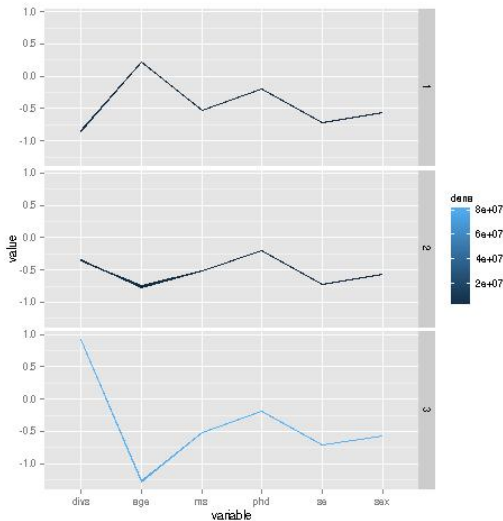
# Analysis of PUMS Data

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Analysis of PUMS Data



- Both $R^2$ values low, but nonpar. 83% higher.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
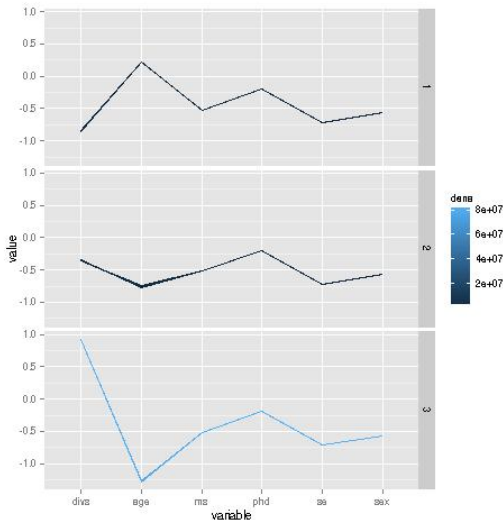R/stat blog:
mat-
loff.wordpress.com

# Analysis of PUMS Data



- Both $R^2$ values low, but nonpar. 83% higher. Room for improvement in param. model!

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
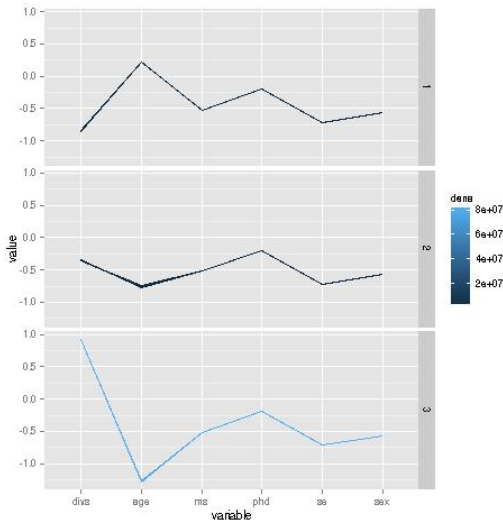mat-
loff.wordpress.com

# Analysis of PUMS Data



- Both $R^2$ values low, but nonpar. 83% higher. Room for improvement in param. model!

- The Age variable seems to be the culprit:

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Analysis of PUMS Data



- Both $R^2$ values low, but nonpar. 83% higher. Room for improvement in param. model!

- The Age variable seems to be the culprit: Overpredict for younger, underpredict for older.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Add Quadratic Term

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Add Quadratic Term

- The "typical divergences" plot suggested adding a quadratic term for Age:

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Add Quadratic Term

- The "typical divergences" plot suggested adding a quadratic term for Age:

  ```
  pg1$age2 <- pg1$age^2
  l2 <- lm(wageinc ~
      age+age2+ms+phd+se+sex, data=pg1)
  ```

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Add Quadratic Term

- The "typical divergences" plot suggested adding a
  quadratic term for Age:

  pg1**$**age2 <− pg1**$**age^2
  l2 <− **lm**( wageinc ~
      age+age2+ms+phd+**se**+sex , **data**=pg1 )

- This brought adj. $R^2$ up from 0.07 to 0.13.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# UCI Adult Data

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# UCI Adult Data

Can use **regdiag()** for generalized linear models too, e.g. logit.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# UCI Adult Data

Can use **regdiag()** for generalized linear models too, e.g. logit.

- Predict a binary
  High Income
  variable, from
  Education, Age,
  Gender, Married.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# UCI Adult Data

Can use **regdiag()** for generalized linear models too, e.g. logit.

- Predict a binary
  High Income
  variable, from
  Education, Age,
  Gender, Married.

- The **regdiag()** plot
  shows younger
  women
  overpredicted, men
  underpredicted.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
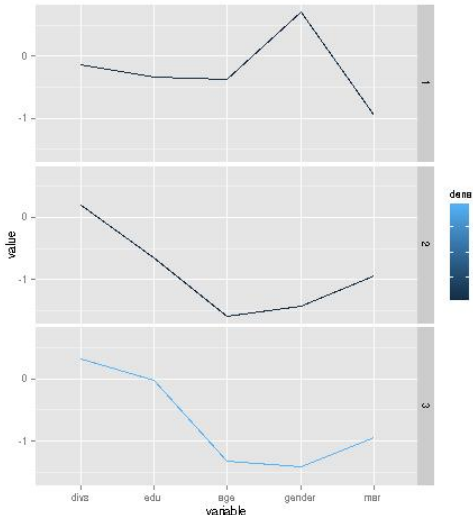loff.wordpress.com

# UCI Adult Data

Can use **regdiag()** for generalized linear models too, e.g. logit.

- Predict a binary
  High Income
  variable, from
  Education, Age,
  Gender, Married.

- The **regdiag()** plot
  shows younger
  women
  overpredicted, men
  underpredicted.

- Thus, might add
  Age $\times$ Gender
  interaction term.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

## UCI Adult Data

Can use **regdiag()** for generalized linear models too, e.g. logit.

- Predict a binary
  High Income
  variable, from
  Education, Age,
  Gender, Married.

- The **regdiag()** plot
  shows younger
  women
  overpredicted, men
  underpredicted.

- Thus, might add
  Age × Gender
  interaction term.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# More on Adult Data

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# More on Adult Data

- Calls:

```
g1 <-
    glm( gt50 ~ edu + age + gender + mar,
    data=newadult, family=binomial)
regdiag(g1)
```

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# More on Adult Data

- Calls:

  ```
  g1 <-
      glm( gt50 ~ edu + age + gender + mar,
      data=newadult, family=binomial)
  regdiag(g1)
  ```

- Addition of interaction term:

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# More on Adult Data

- Calls:

  ```
  g1 <-
      glm( gt50 ~ edu + age + gender + mar,
      data=newadult , family=binomial )
  regdiag ( g1 )
  ```

- Addition of interaction term:

  - Did NOT improve correct-classification rate (81%).

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

## More on Adult Data

- Calls:

  ```
  g1 <-
      glm( gt50 ~ edu + age + gender + mar,
      data=newadult, family=binomial)
  regdiag(g1)
  ```

- Addition of interaction term:
  - Did NOT improve correct-classification rate (81%).
  - BUT changed $\widehat{\beta}_{Gender}$ a lot, from 0.351 to 0.610.
    Interaction term -0.006.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

## More on Adult Data

- Calls:

  ```
  g1 <−
      glm( gt50 ~ edu + age + gender + mar,
      data=newadult, family=binomial)
  regdiag(g1)
  ```

- Addition of interaction term:
  - Did NOT improve correct-classification rate (81%).
  - BUT changed $\widehat{\beta}_{Gender}$ a lot, from 0.351 to 0.610.
    Interaction term -0.006. Male "advantage" in log-odds
    ratio now becomes, e.g. 0.46 at age 25, only 0.28 at age
    55.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Summary

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Summary

- The **freqparcoord** package plots only "typical" lines, thus avoding clutter.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Summary

- The **freqparcoord** package plots only "typical" lines, thus avoding clutter. Can be used for group comparison, outlier hunting, clusters hunting.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Summary

- The **freqparcoord** package plots only "typical" lines, thus avoding clutter. Can be used for group comparison, outlier hunting, clusters hunting.

- The package includes a function **regdiag()** that applies these ideas to regression model diagnostics.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Summary

- The **freqparcoord** package plots only "typical" lines, thus avoding clutter. Can be used for group comparison, outlier hunting, clusters hunting.

- The package includes a function **regdiag()** that applies these ideas to regression model diagnostics.

  - Computes "divergences," i.e. par. fit - nonpar. fit.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Summary

- The **freqparcoord** package plots only "typical" lines, thus avoding clutter. Can be used for group comparison, outlier hunting, clusters hunting.

- The package includes a function **regdiag()** that applies these ideas to regression model diagnostics.

  - Computes "divergences," i.e. par. fit - nonpar. fit.
  - Applies **freqparcoord** to find the most typical divergences, among the most neg. and most pos.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Summary

- The **freqparcoord** package plots only "typical" lines, thus avoding clutter. Can be used for group comparison, outlier hunting, clusters hunting.

- The package includes a function **regdiag()** that applies these ideas to regression model diagnostics.

  - Computes "divergences," i.e. par. fit - nonpar. fit.
  - Applies **freqparcoord** to find the most typical divergences, among the most neg. and most pos.
  - Also reports par., nonpar. $R^2$ values to see whether par. model is "leaving money on the table."

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Summary

- The **freqparcoord** package plots only "typical" lines, thus avoding clutter. Can be used for group comparison, outlier hunting, clusters hunting.

- The package includes a function **regdiag()** that applies these ideas to regression model diagnostics.

  - Computes "divergences," i.e. par. fit - nonpar. fit.
  - Applies **freqparcoord** to find the most typical divergences, among the most neg. and most pos.
  - Also reports par., nonpar. $R^2$ values to see whether par. model is "leaving money on the table."
  - Plots suggest quad., interaction terms to add.

Regression Fit
Diagnostics
Using
freqparcoord

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

e-mail: mat-
loff@cs.ucdavis.edu
R/stat blog:
mat-
loff.wordpress.com

# Summary

- The **freqparcoord** package plots only "typical" lines, thus avoiding clutter. Can be used for group comparison, outlier hunting, clusters hunting.

- The package includes a function **regdiag()** that applies these ideas to regression model diagnostics.

  - Computes "divergences," i.e. par. fit - nonpar. fit.
  - Applies **freqparcoord** to find the most typical divergences, among the most neg. and most pos.
  - Also reports par., nonpar. $R^2$ values to see whether par. model is "leaving money on the table."
  - Plots suggest quad., interaction terms to add.

- Location of these slides:
  http://heather.cs.ucdavis.edu/freqparcoord/
  Slides.pdf