A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# A New Approach to the Parallel Coordinates Method for Large Data Sets

Norm Matloff and Yingkang Xie
University of California at Davis

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

JSM 2014
Boston, MA USA
August 5, 2014

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# What Is Parallel Coordinates Visualization?

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# What Is Parallel Coordinates Visualization?

- If have k variables, draw k vertical axes.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# What Is Parallel Coordinates Visualization?

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

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# What Is Parallel Coordinates Visualization?

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

- Very old idea (late 1800s!).

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# What Is Parallel Coordinates Visualization?

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

- Very old idea (late 1800s!).

- But only popularized 100 years later.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# What Is Parallel Coordinates Visualization?

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

- Very old idea (late 1800s!).

- But only popularized 100 years later.

- Nice math theory, using affine geometry, aiding practical interpretation;

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# What Is Parallel Coordinates Visualization?

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

- Very old idea (late 1800s!).

- But only popularized 100 years later.

- Nice math theory, using affine geometry, aiding practical interpretation; e.g. *Parallel Coordinates: Visual Multidimensional Geometry and Its Applications*, Alfred Inselberg, Springer, 2009.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# What Is Parallel Coordinates Visualization?

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

- Very old idea (late 1800s!).

- But only popularized 100 years later.

- Nice math theory, using affine geometry, aiding practical interpretation; e.g. *Parallel Coordinates: Visual Multidimensional Geometry and Its Applications*, Alfred Inselberg, Springer, 2009.

- Available in **lattice**, **MASS**, **GGally** etc. —

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# What Is Parallel Coordinates Visualization?

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

- Very old idea (late 1800s!).

- But only popularized 100 years later.

- Nice math theory, using affine geometry, aiding practical interpretation; e.g. *Parallel Coordinates: Visual Multidimensional Geometry and Its Applications*, Alfred Inselberg, Springer, 2009.

- Available in **lattice**, **MASS**, **GGally** etc. — but use our approach instead. :-)

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Example

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Example

Example: Height/weight/age data.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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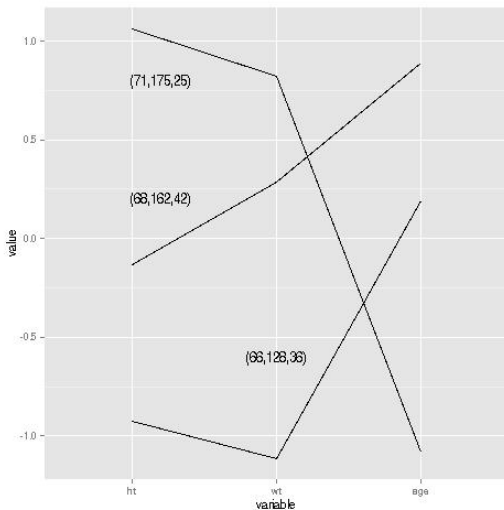mat-
loff.wordpress.com

# Example

Example: Height/weight/age data.

```
> d
  ht  wt age
1 71 175  25
2 66 128  36
3 68 162  42
> library(GGally)
> p <-
+  ggparcoord(d,...
> p <- p + annotate(..
...
```

Vertical axes
use centered,
scaled values.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Problems with Parallel Coordinates

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

- Various solutions; see reviews in Heinrich and Weiskopf (IEEE VIS 2014), Zhou *et al* (IEEE-VGTC, 2008);

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

- Various solutions; see reviews in Heinrich and Weiskopf (IEEE VIS 2014), Zhou *et al* (IEEE-VGTC, 2008); e.g. making the lines fainter, or combining them.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

- Various solutions; see reviews in Heinrich and Weiskopf (IEEE VIS 2014), Zhou *et al* (IEEE-VGTC, 2008); e.g. making the lines fainter, or combining them.

- **But the larger $n$, the less effective these solutions are**,

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Problems with Parallel Coordinates

- Highly cluttered, "black screen" problem.

- Various solutions; see reviews in Heinrich and Weiskopf (IEEE VIS 2014), Zhou *et al* (IEEE-VGTC, 2008); e.g. making the lines fainter, or combining them.

- **But the larger $n$, the less effective these solutions are**, especially with large $p$.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# A New Approach

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# A New Approach

- Our solution:

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# A New Approach

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

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# A New Approach

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

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# A New Approach

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

- Clutter does NOT increase with $n$.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# A New Approach

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

- Clutter does NOT increase with $n$.

- Very versatile.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# A New Approach

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

- Clutter does NOT increase with *n*.

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

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# A New Approach

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

- Clutter does NOT increase with *n*.

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

  - What ht/wt/age combinations are typical overall? **(General analysis.)**

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# A New Approach

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

- Clutter does NOT increase with *n*.

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

  - What ht/wt/age combinations are typical overall? **(General analysis.)**
  - What ht/wt/age combinations are typical within groups?

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# A New Approach

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

- Clutter does NOT increase with *n*.

- Very versatile. E.g., height/weight/age:
  - What ht/wt/age combinations are typical overall? **(General analysis.)**
  - What ht/wt/age combinations are typical within groups? **(Group comparison.)**

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# A New Approach

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

- Clutter does NOT increase with *n*.

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

  - What ht/wt/age combinations are typical overall?
    **(General analysis.)**
  - What ht/wt/age combinations are typical within groups?
    **(Group comparison.)**
  - What ht/wt/age combinations are rare?

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# A New Approach

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

- Clutter does NOT increase with *n*.

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

    - What ht/wt/age combinations are typical overall?
      **(General analysis.)**
    - What ht/wt/age combinations are typical within groups?
      **(Group comparison.)**
    - What ht/wt/age combinations are rare? **(Outlier hunting.)**

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# A New Approach

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

- Clutter does NOT increase with *n*.

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

  - What ht/wt/age combinations are typical overall? **(General analysis.)**
  - What ht/wt/age combinations are typical within groups? **(Group comparison.)**
  - What ht/wt/age combinations are rare? **(Outlier hunting.)**
  - What ht/wt/age combinations are "locally typical"?

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# A New Approach

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

- Clutter does NOT increase with *n*.

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

  - What ht/wt/age combinations are typical overall?
    **(General analysis.)**
  - What ht/wt/age combinations are typical within groups?
    **(Group comparison.)**
  - What ht/wt/age combinations are rare? **(Outlier hunting.)**
  - What ht/wt/age combinations are "locally typical"?
    **(Cluster hunting.)**

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# A New Approach

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

- Clutter does NOT increase with *n*.

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

  - What ht/wt/age combinations are typical overall? **(General analysis.)**
  - What ht/wt/age combinations are typical within groups? **(Group comparison.)**
  - What ht/wt/age combinations are rare? **(Outlier hunting.)**
  - What ht/wt/age combinations are "locally typical"? **(Cluster hunting.)**
  - Bonus: Regression diagnostics.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# A New Approach

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

- Clutter does NOT increase with *n*.

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

  - What ht/wt/age combinations are typical overall? **(General analysis.)**
  - What ht/wt/age combinations are typical within groups? **(Group comparison.)**
  - What ht/wt/age combinations are rare? **(Outlier hunting.)**
  - What ht/wt/age combinations are "locally typical"? **(Cluster hunting.)**
  - Bonus: Regression diagnostics.

- Implemented in a package **freqparcoord** on CRAN.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Example: Taxi Data

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Example: Taxi Data

Example: Taxi data, http://www.theatlantic.com/video/
index/253385/taxi-data-visualization/.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Example: Taxi Data

Example: Taxi data, `http://www.theatlantic.com/video/index/253385/taxi-data-visualization/`.

- We used a 100K subsample.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Example: Taxi Data

Example: Taxi data, `http://www.theatlantic.com/video/index/253385/taxi-data-visualization/`.

- We used a 100K subsample.

- Consists of **data** and **fare** portions, different variables:

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Example: Taxi Data

Example: Taxi data, http://www.theatlantic.com/video/
index/253385/taxi-data-visualization/.

- We used a 100K subsample.

- Consists of **data** and **fare** portions, different variables:

  - **data:** passenger_count, trip_time_in_secs, trip_distance,
    pickup_longitude, pickup_latitude, dropoff_longitude,
    dropoff_latitude, pickuptime

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Example: Taxi Data

Example: Taxi data, http://www.theatlantic.com/video/
index/253385/taxi-data-visualization/.

- We used a 100K subsample.

- Consists of **data** and **fare** portions, different variables:

    - **data:** passenger_count, trip_time_in_secs, trip_distance,
      pickup_longitude, pickup_latitude, dropoff_longitude,
      dropoff_latitude, pickuptime

    - **fare:** fare_amount, surcharge, mta_tax, tip_amount,
      tolls_amount, total_amount, cmt, crd (paid with credit
      card), tippc, booltip (tip, yes or no), pickuptime, daytime

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Outlier Hunting First

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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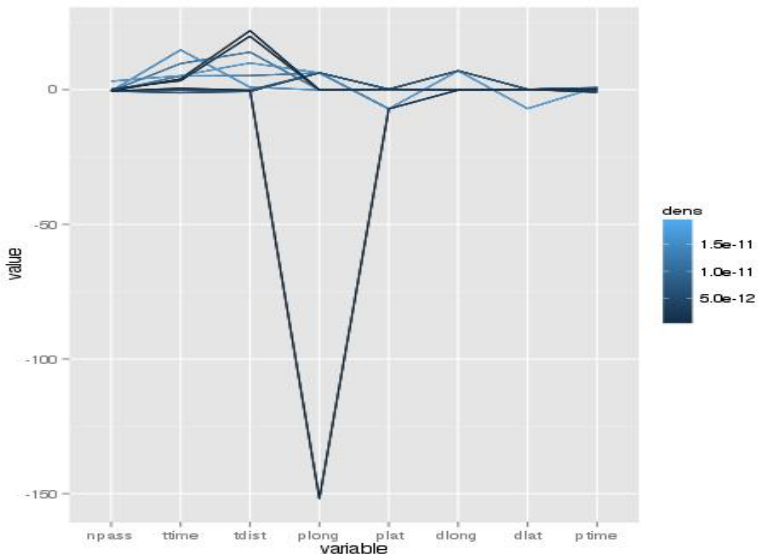*R/stat blog:*
mat-
loff.wordpress.com

# Outlier Hunting First

```
p <- freqparcoord(d100,-10,c(8:15),keepidxs=8)
```

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

## Outlier Hunting First

p <- freqparcoord(d100,-10,c(8:15),keepidxs=8)

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Outliers, cont'd.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Outliers, cont'd.

```
> p$xdisp[,11:14]
       plong       plat       dlong      dlat
   -74.00399   40.742107   -73.94696   40.81335
     0.00000    0.000000   -73.96590   40.80481
   -74.00748   40.703709   -74.07885   40.43142
     0.00000   40.783333     0.00000   40.79044
     0.00000   40.835121     0.00000   40.84693
     0.00000   40.733334     0.00000   40.74148
   -73.88925   40.769035   -73.94363   40.75264
 -1837.04530    0.041667   -73.96226   40.76774
   -73.98628   40.752365   -73.77634   40.64601
     0.00000    0.000000     0.00000    0.00000
```

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Outliers, cont'd.

```
> p$xdisp[,11:14]
        plong       plat      dlong      dlat
   −74.00399   40.742107  −73.94696   40.81335
     0.00000    0.000000  −73.96590   40.80481
   −74.00748   40.703709  −74.07885   40.43142
     0.00000   40.783333    0.00000   40.79044
     0.00000   40.835121    0.00000   40.84693
     0.00000   40.733334    0.00000   40.74148
   −73.88925   40.769035  −73.94363   40.75264
 −1837.04530    0.041667  −73.96226   40.76774
   −73.98628   40.752365  −73.77634   40.64601
     0.00000    0.000000    0.00000    0.00000
```

Bad cases (-1800, 0s) removed (IDs in **p$xdisp** but not shown here).

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Outliers, cont'd.

```
> p$xdisp[,11:14]
         plong        plat       dlong       dlat
    −74.00399    40.742107   −73.94696    40.81335
      0.00000     0.000000   −73.96590    40.80481
    −74.00748    40.703709   −74.07885    40.43142
      0.00000    40.783333     0.00000    40.79044
      0.00000    40.835121     0.00000    40.84693
      0.00000    40.733334     0.00000    40.74148
    −73.88925    40.769035   −73.94363    40.75264
  −1837.04530     0.041667   −73.96226    40.76774
    −73.98628    40.752365   −73.77634    40.64601
      0.00000     0.000000     0.00000     0.00000
```

Bad cases (-1800, 0s) removed (IDs in **p$xdisp** but not shown
here). Trip from Altoona, PA to NYC <u>not</u> removed.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Outliers, cont'd.

```
> p$xdisp[,11:14]
       plong        plat       dlong       dlat
   -74.00399   40.742107   -73.94696   40.81335
     0.00000    0.000000   -73.96590   40.80481
   -74.00748   40.703709   -74.07885   40.43142
     0.00000   40.783333     0.00000   40.79044
     0.00000   40.835121     0.00000   40.84693
     0.00000   40.733334     0.00000   40.74148
   -73.88925   40.769035   -73.94363   40.75264
 -1837.04530    0.041667   -73.96226   40.76774
   -73.98628   40.752365   -73.77634   40.64601
     0.00000    0.000000     0.00000    0.00000
```

Bad cases (-1800, 0s) removed (IDs in **p$xdisp** but not shown
here). Trip from Altoona, PA to NYC <u>not</u> removed.
**Illustrates another advantage of displaying just a few
"typical" cases.**

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# General Analysis

A New
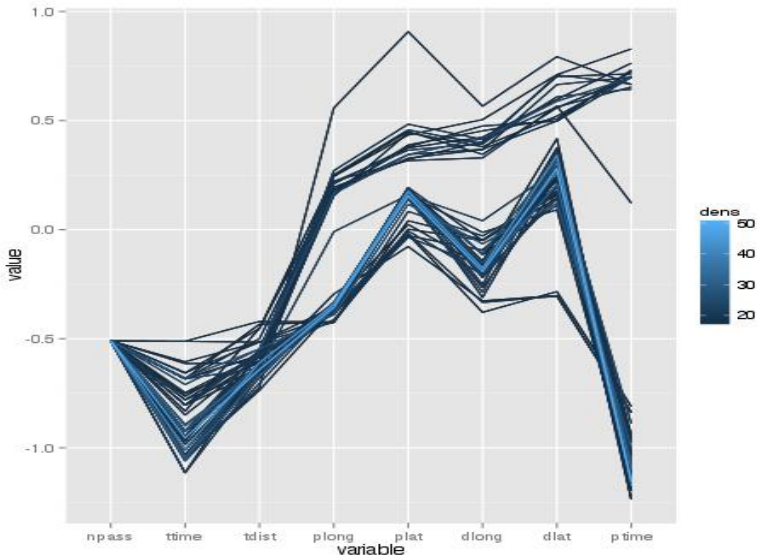Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# General Analysis

```
p <- freqparcoord(d100,50,c(8:15),keepidxs=8)
```

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

## General Analysis

```
p <- freqparcoord(d100,50,c(8:15),keepidxs=8)
```

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# General Analysis, cont'd.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# General Analysis, cont'd.

What do we see?

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# General Analysis, cont'd.

What do we see?

- Already see at least two clusters, largely differing on pickup/dropoff location and time of day.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# General Analysis, cont'd.

What do we see?

- Already see at least two clusters, largely differing on pickup/dropoff location and time of day.

- Note there is much more variation in trip time than in trip distance—due to variation in traffic.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Cluster Analysis

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Cluster Analysis

```
p <- freqparcoord(d100,1,c(8:15),method="locmax",klm=1000,
    cls=cl4,keepidxs=15)
```

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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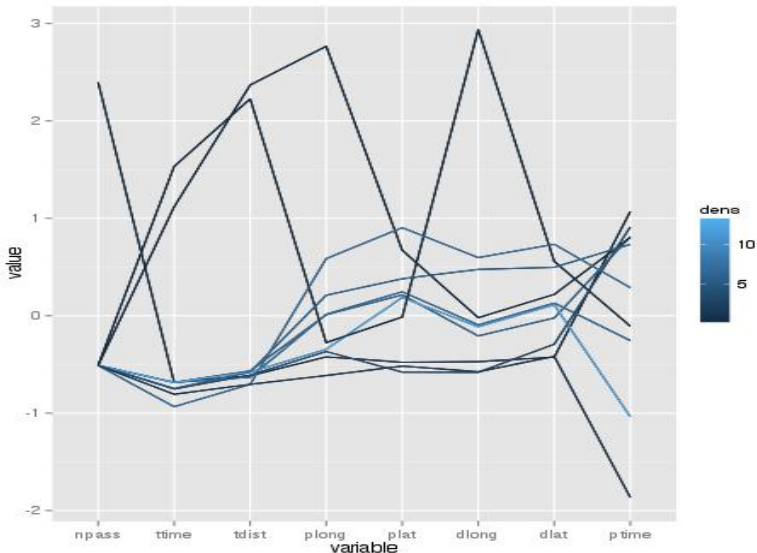R/stat blog:
mat-
loff.wordpress.com

## Cluster Analysis

```
p <- freqparcoord(d100,1,c(8:15),method="locmax",klm=1000,
    cls=cl4,keepidxs=15)
```

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Clustering, cont.d

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Clustering, cont.d

- We see perhaps 8-9 clusters.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Clustering, cont.d

- We see perhaps 8-9 clusters.

- Varying in short vs. long trip distance, pickup/dropoff location, time of day.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Clustering, cont.d

- We see perhaps 8-9 clusters.

- Varying in short vs. long trip distance, pickup/dropoff location, time of day.

- "Changing of the guard," 2 top lines:

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Clustering, cont.d

- We see perhaps 8-9 clusters.

- Varying in short vs. long trip distance, pickup/dropoff location, time of day.

- "Changing of the guard," 2 top lines:
  - Around 1:45 p.m., mid-Manhattan $\rightarrow$ La Guardia Airport.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Clustering, cont.d

- We see perhaps 8-9 clusters.

- Varying in short vs. long trip distance, pickup/dropoff location, time of day.

- "Changing of the guard," 2 top lines:
  - Around 1:45 p.m., mid-Manhattan $\rightarrow$ La Guardia Airport.
  - Around 7:30 p.m., La Guardia Airport $\rightarrow$ mid-Manhattan.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Clustering, cont.d

- We see perhaps 8-9 clusters.

- Varying in short vs. long trip distance, pickup/dropoff location, time of day.

- "Changing of the guard," 2 top lines:
  - Around 1:45 p.m., mid-Manhattan $\rightarrow$ La Guardia Airport.
  - Around 7:30 p.m., La Guardia Airport $\rightarrow$ mid-Manhattan.
  - **Good example of the use of viewing variables together, rather than individually.**

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Group by # of Passengers

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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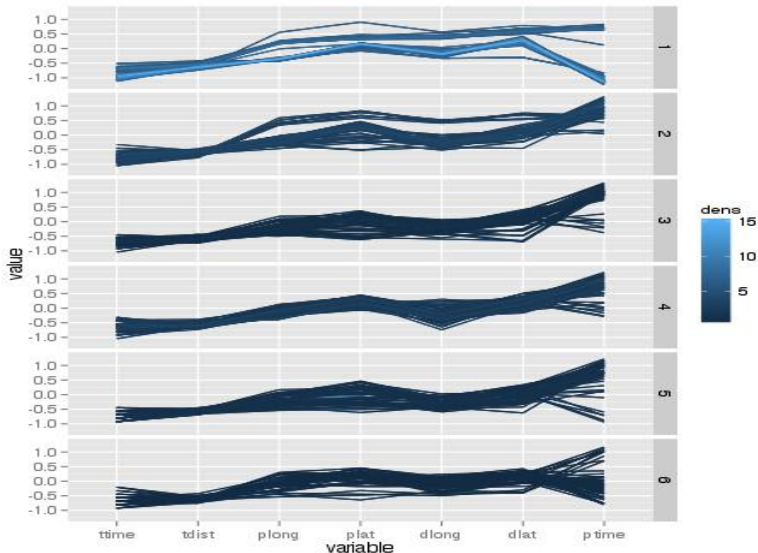*R/stat blog:*
mat-
loff.wordpress.com

# Group by # of Passengers

```
p <- freqparcoord(d100,50,c(9:15),grpvar=8)
```

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Group by # of Passengers

```
p <- freqparcoord(d100,50,c(9:15),grpvar=8)
```

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# # of Passengers, cont'd.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# # of Passengers, cont'd.

- The 1-passenger trips tend to be earlier in the day, some late.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# # of Passengers, cont'd.

- The 1-passenger trips tend to be earlier in the day, some late.
- The 2-4-passenger trips tend to be later in the day.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# # of Passengers, cont'd.

- The 1-passenger trips tend to be earlier in the day, some late.

- The 2-4-passenger trips tend to be later in the day.

- The 5-6 passenger trips (families?) more diverse in time.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Credit Card vs. Cash

A New
Approach to
the Parallel
Coordinates
Method for
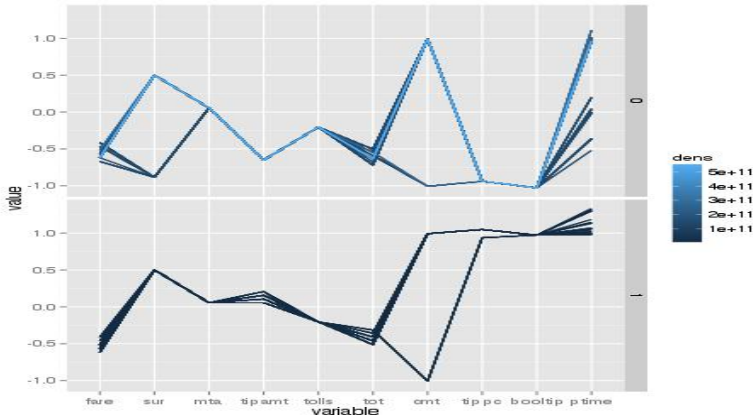Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Credit Card vs. Cash

```
p <- freqparcoord(fare100,10,c(6,7,9:12,14:17),grpvar=13)
```

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Credit Card vs. Cash

```
p <- freqparcoord(fare100,10,c(6,7,9:12,14:17),grpvar=13)
```
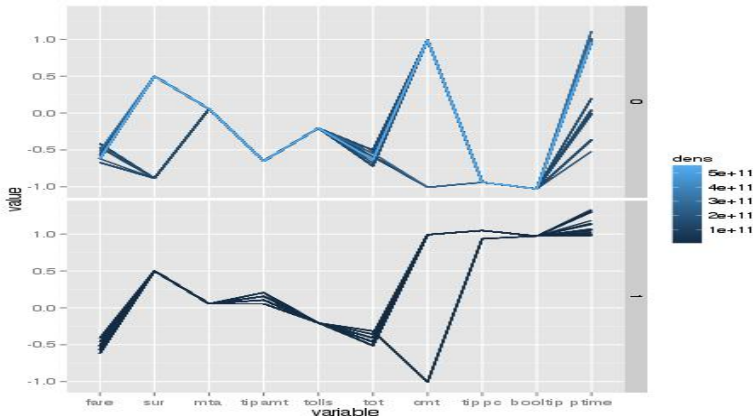


Not much difference, e.g. in base fare.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Credit Card vs. Cash

```
p <- freqparcoord(fare100,10,c(6,7,9:12,14:17),grpvar=13)
```



Not much difference, e.g. in base fare. Some difference in time
of day.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Credit Card vs. Cash

```
p <- freqparcoord(fare100,10,c(6,7,9:12,14:17),grpvar=13)
```
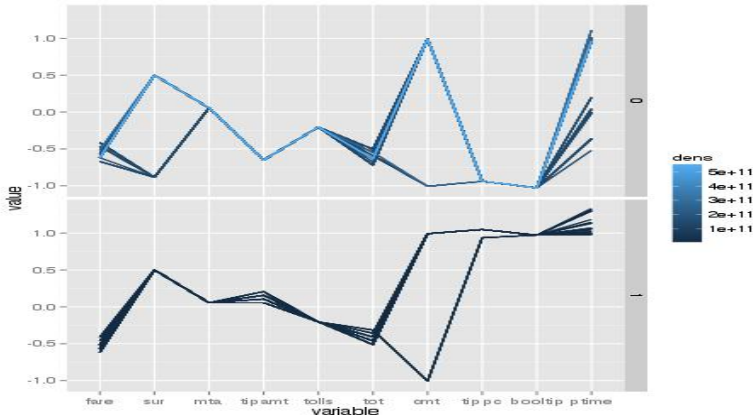


Not much difference, e.g. in base fare. Some difference in time
of day. But stark difference in tips!

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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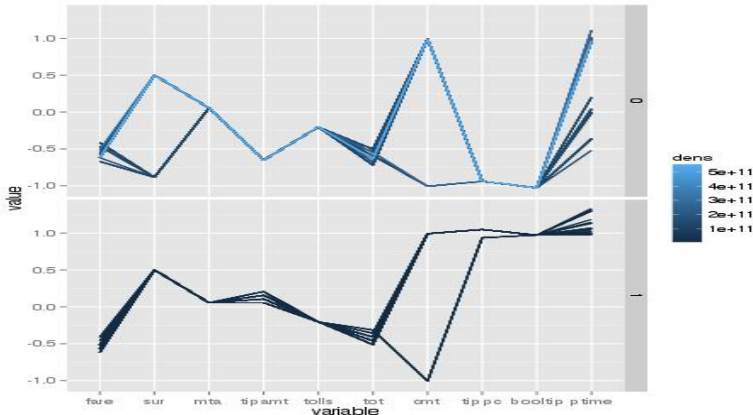ykxie@ucdavis.edu
*R/stat blog:*
mat-
loff.wordpress.com

# Application: Regression
# Diagnostics

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Application: Regression Diagnostics

- Compute *divergences* (<u>not</u> residuals):

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Application: Regression
# Diagnostics

- Compute *divergences* (<u>not</u> residuals):

  $div_i = param\_est_i$ - $nonparam\_est_i$

- Use **freqparcoord()** on the divergences,

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Application: Regression
# Diagnostics

- Compute *divergences* (<u>not</u> residuals):

  $div_i = param\_est_i$ - $nonparam\_est_i$

- Use **freqparcoord()** on the divergences, to identify regions
  of predictor space in which there is systematic over- or
  underestimation of the true regression function.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Application: Regression Diagnostics

- Compute *divergences* (<u>not</u> residuals):

    $div_i = param\_est_i$ - $nonparam\_est_i$

- Use **freqparcoord()** on the divergences, to identify regions of predictor space in which there is systematic over- or underestimation of the true regression function.

- See our **useR! 2014** slides, at `http://heather.cs.ucdavis.edu/freqparcoord/UseR2014Slides.pdf`.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Conclusions

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Conclusions

- A new approach to parallel coordinates.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Conclusions

- A new approach to parallel coordinates.
- Key point: Plots only a few "typical" lines.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Conclusions

- A new approach to parallel coordinates.
- Key point: Plots only a few "typical" lines.
- This reduces clutter—no increase in clutter as $n$ grows!

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Conclusions

- A new approach to parallel coordinates.
- Key point: Plots only a few "typical" lines.
- This reduces clutter—no increase in clutter as $n$ grows!
- Uses: general analysis; group comparison; cluster detection; outlier hunting.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Conclusions

- A new approach to parallel coordinates.

- Key point: Plots only a few "typical" lines.

- This reduces clutter—no increase in clutter as $n$ grows!

- Uses: general analysis; group comparison; cluster detection; outlier hunting.

- Bonus: Regression diagnostics.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Conclusions

- A new approach to parallel coordinates.

- Key point: Plots only a few "typical" lines.

- This reduces clutter—no increase in clutter as *n* grows!

- Uses: general analysis; group comparison; cluster detection; outlier hunting.

- Bonus: Regression diagnostics.

- Package **freqparcoord** on CRAN.

A New
Approach to
the Parallel
Coordinates
Method for
Large Data
Sets

Norm Matloff
and Yingkang
Xie
University of
California at
Davis

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

# Conclusions

- A new approach to parallel coordinates.

- Key point: Plots only a few "typical" lines.

- This reduces clutter—no increase in clutter as $n$ grows!

- Uses: general analysis; group comparison; cluster detection; outlier hunting.

- Bonus: Regression diagnostics.

- Package **freqparcoord** on CRAN.

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