Norm Matloff University of California at Davis

# Long Live (Big Data-Fied) Statistics!

Norm Matloff University of California at Davis

Joint Statistical Meetings Montreal, August 4, 2013

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Norm Matloff University of California at Davis



▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

# Prolog

## Where are we with Big Data?

Norm Matloff University of California at Davis



▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

# Prolog

#### Where are we with Big Data?

- Role of statistics?
- Role of parallel computation?
- Interactions between the two?

Norm Matloff University of California at Davis

#### Where are we?

#### Norm Matloff University of California at Davis

#### Where are we?

Attitudes and worries:



#### Where are we?

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Attitudes and worries:

• "With Big Data, you don't need inference methods."



Where are we?

Norm Matloff University of California at Davis

Attitudes and worries:

- "With Big Data, you don't need inference methods."
- "With Machine Learning, you don't need statistics."

#### Where are we?

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Long Live (Big Data-Fied) Statistics!

Norm Matloff University of California at Davis

Attitudes and worries:

- "With Big Data, you don't need inference methods."
- "With Machine Learning, you don't need statistics."
- Stat community left out of the Big Data revolution (e.g. *Amstat News*, June 2013).

Norm Matloff University of California at Davis

#### Yes, stat is still needed!



#### Yes, stat is still needed!

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Norm Matloff University of California at Davis

#### Yes, stat is still needed!

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

Actually, stat is needed more than ever, e.g.:

 Inference is an issue even for big n, once one considers subsets, where n becomes smaller. Same if do not have p << n.</li>

Norm Matloff University of California at Davis

#### Yes, stat is still needed!

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

- Inference is an issue even for big n, once one considers subsets, where n becomes smaller. Same if do not have p << n.</li>
- Almost all machine learning techniques are revivals of old stat methods.

Norm Matloff University of California at Davis

#### Yes, stat is still needed!

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

- Inference is an issue even for big n, once one considers subsets, where n becomes smaller. Same if do not have p << n.</li>
- Almost all machine learning techniques are revivals of old stat methods. And if you don't understand stat, you won't be able to use ML methods effectively.

Norm Matloff University of California at Davis

#### Yes, stat is still needed!

- Inference is an issue even for big n, once one considers subsets, where n becomes smaller. Same if do not have p << n.</li>
- Almost all machine learning techniques are revivals of old stat methods. And if you don't understand stat, you won't be able to use ML methods effectively.
- An old stat technique—nonparametric curve estimation—now more useful than ever, for Big Data Graphics.

Norm Matloff University of California at Davis

#### Yes, stat is still needed!

- Inference is an issue even for big n, once one considers subsets, where n becomes smaller. Same if do not have p << n.</li>
- Almost all machine learning techniques are revivals of old stat methods. And if you don't understand stat, you won't be able to use ML methods effectively.
- An old stat technique—nonparametric curve estimation—now more useful than ever, for Big Data Graphics.
- The Curse of Dimensionality hasn't gone away.

Norm Matloff University of California at Davis

#### Yes, stat is still needed!

Actually, stat is needed more than ever, e.g.:

- Inference is an issue even for big n, once one considers subsets, where n becomes smaller. Same if do not have p << n.</li>
- Almost all machine learning techniques are revivals of old stat methods. And if you don't understand stat, you won't be able to use ML methods effectively.
- An old stat technique—nonparametric curve estimation—now more useful than ever, for Big Data Graphics.
- The Curse of Dimensionality hasn't gone away. Impossible to understand without stat.

Norm Matloff University of California at Davis

## Setting



### Setting

Consider the classical (though not universal) data format:



# Setting

#### Consider the classical (though not universal) data format:

• n observations/cases/instances/...



# Setting

Consider the classical (though not universal) data format:

- n observations/cases/instances/...
- p variables/features/attributes/...

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ

Long Live (Big Data-Fied) Statistics!

Norm Matloff University of California at Davis

Consider the classical (though not universal) data format:

- n observations/cases/instances/...
- p variables/features/attributes/...
- Assumed i.i.d.

Long Live (Big Data-Fied) Statistics!

Norm Matloff University of California at Davis

Consider the classical (though not universal) data format:

- n observations/cases/instances/...
- p variables/features/attributes/...
- Assumed i.i.d.

(Here I'm trying to include terminology from the nonstat communities.)

Long Live (Big Data-Fied) Statistics!

Norm Matloff University of California at Davis

Consider the classical (though not universal) data format:

- n observations/cases/instances/...
- p variables/features/attributes/...
- Assumed i.i.d.

(Here I'm trying to include terminology from the nonstat communities.)

Does "Big" Data mean big n or big p or both?

Long Live (Big Data-Fied) Statistics!

Norm Matloff University of California at Davis

Consider the classical (though not universal) data format:

- n observations/cases/instances/...
- p variables/features/attributes/...
- Assumed i.i.d.

(Here I'm trying to include terminology from the nonstat communities.)

Does "Big" Data mean big n or big p or both?

This talk will contain one "big n" section and two "big p" section.

Norm Matloff University of California at Davis

# Big n

Norm Matloff University of California at Davis



# Part I: Big-n Problem

Norm Matloff University of California at Davis

#### Big-n can be handled

Norm Matloff University of California at Davis

### Big-n can be handled

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Norm Matloff University of California at Davis

# Big-n can be handled

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Big n can generally be handled:

• Many (though certainly not all) computations "additive," thus "embarrassingly parallel."

Norm Matloff University of California at Davis

# Big-n can be handled

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

- Many (though certainly not all) computations "additive," thus "embarrassingly parallel."
- Thus amenable to parallel computation, especially distributed data, MapReduce etc.

Norm Matloff University of California at Davis

# Big-n can be handled

- Many (though certainly not all) computations "additive," thus "embarrassingly parallel."
- Thus amenable to parallel computation, especially distributed data, MapReduce etc.
- "Chunks averaging method" (CAM) (Fan *et al*, 2007; Matloff, 2010; etc.) can turn most statistical computations into embarrassingly parallel.

Norm Matloff University of California at Davis

# Big-n can be handled

- Many (though certainly not all) computations "additive," thus "embarrassingly parallel."
- Thus amenable to parallel computation, especially distributed data, MapReduce etc.
- "Chunks averaging method" (CAM) (Fan *et al*, 2007; Matloff, 2010; etc.) can turn most statistical computations into embarrassingly parallel. ("Software alchemy.")

Norm Matloff University of California at Davis

## Chunk Averaging Method

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Norm Matloff University of California at Davis

# Chunk Averaging Method

• Key point: CAM converts non-embaarrassingly parallel algs to additive ones that are **statistically equivalent** (same standard errors).

Norm Matloff University of California at Davis

# Chunk Averaging Method

- Key point: CAM converts non-embaarrassingly parallel algs to additive ones that are **statistically equivalent** (same standard errors).
- Example: Regression.

Norm Matloff University of California at Davis

# Chunk Averaging Method

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

- Key point: CAM converts non-embaarrassingly parallel algs to additive ones that are **statistically equivalent** (same standard errors).
- Example: Regression.
  - Break observations into chunks.
Norm Matloff University of California at Davis

# Chunk Averaging Method

- Key point: CAM converts non-embaarrassingly parallel algs to additive ones that are **statistically equivalent** (same standard errors).
- Example: Regression.
  - Break observations into chunks.
  - Fit regression equation to each chunk.

Norm Matloff University of California at Davis

# Chunk Averaging Method

- Key point: CAM converts non-embaarrassingly parallel algs to additive ones that are **statistically equivalent** (same standard errors).
- Example: Regression.
  - Break observations into chunks.
  - Fit regression equation to each chunk.
  - Average the results.

Norm Matloff University of California at Davis

# Chunk Averaging Method

- Key point: CAM converts non-embaarrassingly parallel algs to additive ones that are **statistically equivalent** (same standard errors).
- Example: Regression.
  - Break observations into chunks.
  - Fit regression equation to each chunk.
  - Average the results.
- Produces statistically equivalent results for large n.

Norm Matloff University of California at Davis

# Chunk Averaging Method

- Key point: CAM converts non-embaarrassingly parallel algs to additive ones that are **statistically equivalent** (same standard errors).
- Example: Regression.
  - Break observations into chunks.
  - Fit regression equation to each chunk.
  - Average the results.
- Produces statistically equivalent results for large n.
- Essentially and i.i.d.-based method, e.g. quantile regression, hazard function estimation, tree methods, etc.

Norm Matloff University of California at Davis

# Chunk Averaging Method

- Key point: CAM converts non-embaarrassingly parallel algs to additive ones that are **statistically equivalent** (same standard errors).
- Example: Regression.
  - Break observations into chunks.
  - Fit regression equation to each chunk.
  - Average the results.
- Produces statistically equivalent results for large n.
- Essentially and i.i.d.-based method, e.g. quantile regression, hazard function estimation, tree methods, etc.

• Superlinear speedup.

Norm Matloff University of California at Davis

# Chunk Averaging Method

- Key point: CAM converts non-embaarrassingly parallel algs to additive ones that are **statistically equivalent** (same standard errors).
- Example: Regression.
  - Break observations into chunks.
  - Fit regression equation to each chunk.
  - Average the results.
- Produces statistically equivalent results for large n.
- Essentially and i.i.d.-based method, e.g. quantile regression, hazard function estimation, tree methods, etc.
- *Superlinear* speedup. E.g. quantile regression, 5.31X for 4 threads.

Norm Matloff University of California at Davis

# Chunk Averaging Method

- Key point: CAM converts non-embaarrassingly parallel algs to additive ones that are **statistically equivalent** (same standard errors).
- Example: Regression.
  - Break observations into chunks.
  - Fit regression equation to each chunk.
  - Average the results.
- Produces statistically equivalent results for large n.
- Essentially and i.i.d.-based method, e.g. quantile regression, hazard function estimation, tree methods, etc.
- *Superlinear* speedup. E.g. quantile regression, 5.31X for 4 threads. Can be faster even for just one core.

Norm Matloff University of California at Davis

#### Part II

<□ > < @ > < E > < E > E のQ @

Norm Matloff University of California at Davis



▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

# Part II: A New Graphical Approach to Big-p Problem

Norm Matloff University of California at Davis

#### Graphing Lots of Variables

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Norm Matloff University of California at Davis

Issues:

### Graphing Lots of Variables

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?



Norm Matloff University of California at Davis

## Graphing Lots of Variables

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Issues:

• Can deal (a little bit) better with big-p if we can display lots of variables on the same graph.

Norm Matloff University of California at Davis

# Graphing Lots of Variables

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Issues:

- Can deal (a little bit) better with big-p if we can display lots of variables on the same graph.
- Can't display all at once, but try to get at least several.

Norm Matloff University of California at Davis

# Graphing Lots of Variables

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Issues:

- Can deal (a little bit) better with big-p if we can display lots of variables on the same graph.
- Can't display all at once, but try to get at least several.
- Problems:

Norm Matloff University of California at Davis

# Graphing Lots of Variables

#### Issues:

- Can deal (a little bit) better with big-p if we can display lots of variables on the same graph.
- Can't display all at once, but try to get at least several.
- Problems:
  - Displaying > 2 variables on a 2-dimensional device.

Norm Matloff University of California at Davis

# Graphing Lots of Variables

#### Issues:

- Can deal (a little bit) better with big-p if we can display lots of variables on the same graph.
- Can't display all at once, but try to get at least several.
- Problems:
  - Displaying > 2 variables on a 2-dimensional device.
  - "Black screen problem"—with big n, at least parts of the screen become solid black.

Norm Matloff University of California at Davis

#### Graphing Lots of Variables

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Norm Matloff University of California at Davis

#### Graphing Lots of Variables

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Norm Matloff University of California at Davis

### Graphing Lots of Variables

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Some existing methods:

Norm Matloff University of California at Davis

# Graphing Lots of Variables

Some existing methods:

• Grand tours:

Norm Matloff University of California at Davis

# Graphing Lots of Variables

Some existing methods:

• Grand tours: Show sequence of rotations and projections, e.g. **tourr** in R.

Norm Matloff University of California at Davis

# Graphing Lots of Variables

Some existing methods:

 Grand tours: Show sequence of rotations and projections, e.g. tourr in R. Nice, visually appealing.

Norm Matloff University of California at Davis

# Graphing Lots of Variables

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Some existing methods:

• Grand tours: Show sequence of rotations and projections, e.g. **tourr** in R.

Nice, visually appealing. But hard to discern exact relations, and suffers from black-screen problem.

Norm Matloff University of California at Davis

# Graphing Lots of Variables

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Some existing methods:

- Grand tours: Show sequence of rotations and projections, e.g. **tourr** in R.
  - Nice, visually appealing. But hard to discern exact relations, and suffers from black-screen problem.
- Parallel coordinates:

Norm Matloff University of California at Davis

# Graphing Lots of Variables

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Some existing methods:

• Grand tours: Show sequence of rotations and projections, e.g. **tourr** in R.

Nice, visually appealing. But hard to discern exact relations, and suffers from black-screen problem.

• Parallel coordinates: Draw one vertical axis for each variable. Draw a set of connecting lines for each data point.

Norm Matloff University of California at Davis

# Graphing Lots of Variables

Some existing methods:

• Grand tours: Show sequence of rotations and projections, e.g. **tourr** in R.

Nice, visually appealing. But hard to discern exact relations, and suffers from black-screen problem.

• Parallel coordinates: Draw one vertical axis for each variable. Draw a set of connecting lines for each data point.

Hard to understand noncontiguous axes, black screen problem.

Norm Matloff University of California at Davis

#### Statistics to the rescue!

Norm Matloff University of California at Davis

#### Statistics to the rescue!

An obvious solution to the black-screen problem:

Norm Matloff University of California at Davis

#### Statistics to the rescue!

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

An obvious solution to the black-screen problem: nonparametric curve estimation.

Norm Matloff University of California at Davis

#### Statistics to the rescue!

▲ロト ▲周ト ▲ヨト ▲ヨト ヨー のくで

Norm Matloff University of California at Davis

### Statistics to the rescue!

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

An obvious solution to the black-screen problem: nonparametric curve estimation. Example: Scatter plots.

• Ordinary plot would fill the screen.

Norm Matloff University of California at Davis

### Statistics to the rescue!

▲ロト ▲冊 ▶ ▲ ヨ ▶ ▲ ヨ ▶ ● の Q @

- Ordinary plot would fill the screen.
- Solution: Draw the nonpar. 2-dim. density estimate instead.

Norm Matloff University of California at Davis

### Statistics to the rescue!

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

- Ordinary plot would fill the screen.
- Solution: Draw the nonpar. 2-dim. density estimate instead.
- That makes 3 dimensions, but code third dimension (density height) via color.

Norm Matloff University of California at Davis

### Statistics to the rescue!

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

- Ordinary plot would fill the screen.
- Solution: Draw the nonpar. 2-dim. density estimate instead.
- That makes 3 dimensions, but code third dimension (density height) via color.
- E.g. scatterSmooth() in R.

Norm Matloff University of California at Davis

#### Displaying 3 Vars. in 2 Dims.

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Norm Matloff University of California at Davis

# Displaying 3 Vars. in 2 Dims.

The **scatterSmooth()** example actually shows how to display 3 variables in 2 dimensions:
Norm Matloff University of California at Davis

# Displaying 3 Vars. in 2 Dims.

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

The **scatterSmooth()** example actually shows how to display 3 variables in 2 dimensions: Say have variables X, Y, Z.

Norm Matloff University of California at Davis

# Displaying 3 Vars. in 2 Dims.

The **scatterSmooth()** example actually shows how to display 3 variables in 2 dimensions: Say have variables X, Y, Z.

• Plot regression function of Z, **color coded**, against X and Y.

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Norm Matloff University of California at Davis

# Displaying 3 Vars. in 2 Dims.

The **scatterSmooth()** example actually shows how to display 3 variables in 2 dimensions: Say have variables X, Y, Z.

- Plot regression function of Z, **color coded**, against X and Y.
- Regression function:  $m(s,t) = E(Z \mid X = s, Y = t)$  (i.e. general, not assuming param. model).

Norm Matloff University of California at Davis

# Displaying 3 Vars. in 2 Dims.

The **scatterSmooth()** example actually shows how to display 3 variables in 2 dimensions: Say have variables X, Y, Z.

- Plot regression function of Z, **color coded**, against X and Y.
- Regression function:  $m(s,t) = E(Z \mid X = s, Y = t)$  (i.e. general, not assuming param. model).

• Use nonpar. estimation, e.g. nearest-neighbor.

Norm Matloff University of California at Davis

# Displaying 3 Vars. in 2 Dims.

The **scatterSmooth()** example actually shows how to display 3 variables in 2 dimensions: Say have variables X, Y, Z.

- Plot regression function of Z, **color coded**, against X and Y.
- Regression function:  $m(s,t) = E(Z \mid X = s, Y = t)$  (i.e. general, not assuming param. model).

< □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > < □ > <

- Use nonpar. estimation, e.g. nearest-neighbor.
- 3 vars. in 2 dims.!

Norm Matloff University of California at Davis

# Displaying 3 Vars. in 2 Dims.

The **scatterSmooth()** example actually shows how to display 3 variables in 2 dimensions: Say have variables X, Y, Z.

- Plot regression function of Z, **color coded**, against X and Y.
- Regression function:  $m(s,t) = E(Z \mid X = s, Y = t)$  (i.e. general, not assuming param. model).

- Use nonpar. estimation, e.g. nearest-neighbor.
- 3 vars. in 2 dims.! (No perspective plotting.)

Norm Matloff University of California at Davis

### Proposed Boundary Method

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Norm Matloff University of California at Davis

## Proposed Boundary Method

I introduce here a new approach to plotting multiple variables in 2 dims.,

Norm Matloff University of California at Davis

## Proposed Boundary Method

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

I introduce here a new approach to plotting multiple variables in 2 dims., based on "boundary curves."

Norm Matloff University of California at Davis

## Proposed Boundary Method

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

I introduce here a new approach to plotting multiple variables in 2 dims., based on "boundary curves." First, again consider 3 variables, X, Y and Z.

Norm Matloff University of California at Davis

## Proposed Boundary Method

I introduce here a new approach to plotting multiple variables in 2 dims., based on "boundary curves." First, again consider 3 variables, X, Y and Z.

• For user-chosen b, boundary is the set

$$\{(s,t): E(Z|X=s,Y=t)=b\}$$
 (1)

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Norm Matloff University of California at Davis

## Proposed Boundary Method

I introduce here a new approach to plotting multiple variables in 2 dims., based on "boundary curves." First, again consider 3 variables, X, Y and Z.

• For user-chosen b, boundary is the set

$$\{(s,t): E(Z|X=s,Y=t)=b\}$$
 (1)

• User might set b = E(Z) (overall, unconditional mean).

Norm Matloff University of California at Davis

## Proposed Boundary Method

I introduce here a new approach to plotting multiple variables in 2 dims., based on "boundary curves." First, again consider 3 variables, X, Y and Z.

• For user-chosen b, boundary is the set

$$\{(s,t): E(Z|X=s, Y=t) = b\}$$
 (1)

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

- User might set b = E(Z) (overall, unconditional mean).
- Plot estimate of the boundary curve.

Norm Matloff University of California at Davis

## Proposed Boundary Method

I introduce here a new approach to plotting multiple variables in 2 dims., based on "boundary curves." First, again consider 3 variables, X, Y and Z.

• For user-chosen b, boundary is the set

$$\{(s,t): E(Z|X=s, Y=t) = b\}$$
 (1)

- User might set b = E(Z) (overall, unconditional mean).
- Plot estimate of the boundary curve.
- Displaying 3 vars. in 2 dims.

Norm Matloff University of California at Davis

### Boundary Plot Example

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Norm Matloff University of California at Davis

### Boundary Plot Example

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

• Bank account data, UCI repository.

Norm Matloff University of California at Davis

## Boundary Plot Example

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

- Bank account data, UCI repository.
- X = age of customer, Y = current bank account, Z = say Yes to open new type of account

Norm Matloff University of California at Davis

## Boundary Plot Example

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

- Bank account data, UCI repository.
- X = age of customer, Y = current bank account, Z = say Yes to open new type of account

Norm Matloff University of California at Davis

### Bank Example

◆□▶ ◆□▶ ◆∃▶ ◆∃▶ = のへで

#### Norm Matloff University of California at Davis

### Bank Example

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへぐ



Norm Matloff University of California at Davis



## Bank Example

 Above line means, above-avg. prob. sign up for new account.

(日)、(四)、(E)、(E)、(E)

Norm Matloff University of California at Davis



## Bank Example

- Above line means, above-avg. prob. sign up for new account.
- Near retire ⇒
  "hardest sell" !

Norm Matloff University of California at Davis



# Bank Example

- Above line means, above-avg. prob. sign up for new account.
- Near retire ⇒
  "hardest sell" !
- Those around 60 need a large balance before willing to try new account.

イロト 不得 トイヨト イヨト

Norm Matloff University of California at Davis

### More Than 3 Vars. in 2 Dims.

Norm Matloff University of California at Davis

### More Than 3 Vars. in 2 Dims.

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

• Plotting boundaries has been done before.

Norm Matloff University of California at Davis

### More Than 3 Vars. in 2 Dims.

- Plotting boundaries has been done before.
- But the idea here is to display <u>several</u> boundaries at once, so as to display more variables in one 2-dim. graph.

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Norm Matloff University of California at Davis

### Example: Adult Data

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Norm Matloff University of California at Davis

### Example: Adult Data

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

UCI Adult data

Norm Matloff University of California at Davis

### Example: Adult Data

- UCI Adult data
- X = age, Y = education, Z = high income

Norm Matloff University of California at Davis

### Example: Adult Data

- UCI Adult data
- X = age, Y = education, Z = high income
- But now add a  $\underline{4th}$  variable: W = gender

Norm Matloff University of California at Davis

### Example: Adult Data

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

- UCI Adult data
- X = age, Y = education, Z = high income
- But now add a  $\underline{4th}$  variable: W = gender
- Plot 2 boundary curves, one male and one female.

Norm Matloff University of California at Davis

## Example: Adult Data

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

- UCI Adult data
- X = age, Y = education, Z = high income
- But now add a  $\underline{4th}$  variable: W = gender
- Plot 2 boundary curves, one male and one female.
- Thus display <u>4</u> variables in 2 dims.

Norm Matloff University of California at Davis

### Adult Example

(ロ)、(型)、(E)、(E)、 E) の(の)

Norm Matloff

Adult Example

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

University of California at Davis



Norm Matloff University of California at Davis



## Adult Example

 Above line means, higher-thanavg. prob. of high income.

(日)、(四)、(E)、(E)、(E)

Norm Matloff University of California at Davis



## Adult Example

- Above line means, higher-thanavg. prob. of high income.
- Before age 35, not much difference.

(日) (個) (目) (目) (目) (目)
Norm Matloff University of California at Davis



# Adult Example

- Above line means, higher-thanavg. prob. of high income.
- Before age 35, not much difference.
- After age 35, women need much more education than men to likely have high income.

(日)、

э

Norm Matloff University of California at Davis

## Example: Flight Lateness

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Norm Matloff University of California at Davis

# Example: Flight Lateness

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

• Airline lateness data.

Norm Matloff University of California at Davis

# Example: Flight Lateness

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

- Airline lateness data.
- X = departure delay, Y = distance, Z = arrival lateness,
  - $\mathsf{W}=\mathsf{originating}\ \mathsf{airport}\ (\mathsf{here,}\ \mathsf{SFO},\ \mathsf{IAD},\ \mathsf{IAH}),$

Norm Matloff University of California at Davis

# Example: Flight Lateness

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

- Airline lateness data.
- X = departure delay, Y = distance, Z = arrival lateness, W = originating airport (here, SFO, IAD, IAH), so again, displaying 4 variables in 2 dims.
- 3 curves, one for each airport

Norm Matloff University of California at Davis

# Example: Flight Lateness

- Airline lateness data.
- X = departure delay, Y = distance, Z = arrival lateness, W = originating airport (here, SFO, IAD, IAH), so again, displaying 4 variables in 2 dims.
- 3 curves, one for each airport
- Could add V = daytime vs. evening, for 6 curves, thus displaying 5 variables in 2 dims.

Norm Matloff University of California at Davis

# Example: Flight Lateness

- Airline lateness data.
- X = departure delay, Y = distance, Z = arrival lateness, W = originating airport (here, SFO, IAD, IAH), so again, displaying 4 variables in 2 dims.
- 3 curves, one for each airport
- Could add V = daytime vs. evening, for 6 curves, thus displaying 5 variables in 2 dims.
- Could plot straight regressions too, but boundaries always enable us to plot "one more variable."

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

#### Norm Matloff University of California at Davis

## Airline Example

▲□▶ ▲圖▶ ▲圖▶ ▲圖▶ = ● ● ●

Norm Matloff University of California at Davis Airline Example

3000 -SFO distance - 0000 -IAD IAH 1000 -100 dep delay

◆□▶ ◆□▶ ◆豆▶ ◆豆▶ ̄豆 \_ のへぐ

Norm Matloff University of California at Davis



## Airline Example

 Above line means, higher-thanavg. mean delay.

・ロト ・雪 ト ・ ヨ ト ・ ヨ ト

æ

Norm Matloff University of California at Davis



# Airline Example

- Above line means, higher-thanavg. mean delay.
- SFO seems to be doing better.

(日) (個) (目) (目) (目) (目)

Norm Matloff University of California at Davis



# Airline Example

- Above line means, higher-thanavg. mean delay.
- SFO seems to be doing better. Need a very long flight to have above-avg. delay, relative to the others.

Norm Matloff University of California at Davis

## Computation

<□ > < @ > < E > < E > E のQ @



Data-Fied) Statistics! Norm Matloff

#### Norm Matloff University of California at Davis

## Computation

• R's (contour() not used (don't want "islands").

## Computation

Long Live (Big Data-Fied) Statistics!

- R's (contour() not used (don't want "islands").
- Estimate regression (via fast kNN, FNN library).

## Computation

Long Live (Big Data-Fied) Statistics!

- R's (contour() not used (don't want "islands").
- Estimate regression (via fast kNN, FNN library).
- Find "boundary band," all points near the estimate boundary.

## Computation

Long Live (Big Data-Fied) Statistics!

- R's (contour() not used (don't want "islands").
- Estimate regression (via fast kNN, FNN library).
- Find "boundary band," all points near the estimate boundary.
- Smooth the band.

Norm Matloff University of California at Davis

## Parallel Computation

Norm Matloff University of California at Davis

## Parallel Computation

#### Norm Matloff University of California at Davis

## Parallel Computation

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

Computation can be voluminous.

• Parallel processing.

#### Norm Matloff University of California at Davis

# Parallel Computation

- Parallel processing.
- Take advantage of superlinearity from CAM.

#### Norm Matloff University of California at Davis

# Parallel Computation

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

- Parallel processing.
- Take advantage of superlinearity from CAM.
- Break into chunks, but only find near nghbrs. within chunks, not across chunks.

#### Norm Matloff University of California at Davis

# Parallel Computation

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

- Parallel processing.
- Take advantage of superlinearity from CAM.
- Break into chunks, but only find near nghbrs. within chunks, not across chunks.
- The "A" part of CAM comes in the smoothing of the band.

Part III

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Long Live (Big Data-Fied) Statistics!

Norm Matloff University of California at Davis

# Part III: Big p and the Curse of Dimensionality

Exorcizing the Curse of Dimensionality

Part III

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Long Live (Big Data-Fied) Statistics!

Norm Matloff University of California at Davis

# Part III: Big p and the Curse of Dimensionality

Exorcizing the Curse of Dimensionality Some small steps in that direction.

Norm Matloff University of California at Davis

# Big p

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 のへぐ

Norm Matloff University of California at Davis

# Big p

• Theoretical considerations imply that should have  $p < \sqrt{n}$  in regression case (Portnoy, 1968).

Long Live (Big Data-Fied) Statistics!

- Theoretical considerations imply that should have  $p < \sqrt{n}$  in regression case (Portnoy, 1968).
- Yet today p >> n is commonplace.

Statistics! Norm Matloff University of California at Davis

Long Live (Big

Data-Fied)

- Theoretical considerations imply that should have  $p < \sqrt{n}$  in regression case (Portnoy, 1968).
- Yet today p >> n is commonplace.
- This causes "multiple inference" problems (e.g. familywise error rates).

Statistics! Norm Matloff University of California at Davis

Long Live (Big

Data-Fied)

- Theoretical considerations imply that should have  $p < \sqrt{n}$  in regression case (Portnoy, 1968).
- Yet today p >> n is commonplace.
- This causes "multiple inference" problems (e.g. familywise error rates).
- So, e.g., Cl radii 1.96 std.err. $(\hat{\theta})$  might NOT be "essentially 0."

Statistics! Norm Matloff University of California at Davis

Long Live (Big

Data-Fied)

- Theoretical considerations imply that should have  $p < \sqrt{n}$  in regression case (Portnoy, 1968).
- Yet today *p* >> *n* is commonplace.
- This causes "multiple inference" problems (e.g. familywise error rates).
- So, e.g., Cl radii 1.96 std.err.(θ̂) might NOT be "essentially 0." I.e., Big n not big after all.

Data-Fied) Statistics! Norm Matloff University of California at

Davis

Long Live (Big

- Theoretical considerations imply that should have  $p < \sqrt{n}$  in regression case (Portnoy, 1968).
- Yet today p >> n is commonplace.
- This causes "multiple inference" problems (e.g. familywise error rates).
- So, e.g., Cl radii 1.96 std.err.(θ̂) might NOT be "essentially 0." I.e., Big n not big after all.
- And the ever-present Curse of Dimensionality.

Norm Matloff University of California at Davis

## Principle Components Analysis

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?



Norm Matloff University of California at Davis

# Principle Components Analysis

• What sizes of p relative to n might be problematic for PCA?



Norm Matloff University of California at Davis

# Principle Components Analysis

- What sizes of p relative to n might be problematic for PCA?
  - Sample covariance matrix V has p(p-1)/2 distinct entries.



Norm Matloff University of California at Davis

# Principle Components Analysis

- What sizes of p relative to n might be problematic for PCA?
  - Sample covariance matrix V has p(p-1)/2 distinct entries.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

• Data matrix has np entries.
Norm Matloff University of California at Davis

# Principle Components Analysis

- What sizes of p relative to n might be problematic for PCA?
  - Sample covariance matrix V has p(p-1)/2 distinct entries.
  - Data matrix has np entries.
  - So V is completely determined (except roundoff error) if np = p(p-1)/2.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

Norm Matloff University of California at Davis

# Principle Components Analysis

- What sizes of p relative to n might be problematic for PCA?
  - Sample covariance matrix V has p(p-1)/2 distinct entries.
  - Data matrix has np entries.
  - So V is completely determined (except roundoff error) if np = p(p-1)/2.

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

• So, have problem if p > 2n, roughly.

Norm Matloff University of California at Davis

# PCA Experiment



Norm Matloff University of California at Davis

#### Simulation experiment:

# PCA Experiment

Norm Matloff University of California at Davis

#### Simulation experiment:

•  $Y_1$ ,  $Y_2$  indep. N(0,1);  $X_1 = Y_1 + Y_2$ ,  $X_2 = Y_1 - Y_2$ ,  $X_3$ , ...,  $X_p$  iid N(0,1), indep. of  $X_1$ ,  $X_2$ .

**PCA** Experiment

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

#### Norm Matloff University of California at Davis

#### Simulation experiment:

•  $Y_1$ ,  $Y_2$  indep. N(0,1);  $X_1 = Y_1 + Y_2$ ,  $X_2 = Y_1 - Y_2$ ,  $X_3$ , ...,  $X_p$  iid N(0,1), indep. of  $X_1$ ,  $X_2$ .

**PCA** Experiment

• First PC should be (1,0,0,...) or (0,1,0,...).

Norm Matloff University of California at Davis

#### Simulation experiment:

•  $Y_1$ ,  $Y_2$  indep. N(0,1);  $X_1 = Y_1 + Y_2$ ,  $X_2 = Y_1 - Y_2$ ,  $X_3$ , ...,  $X_p$  iid N(0,1), indep. of  $X_1$ ,  $X_2$ .

PCA Experiment

• First PC should be (1,0,0,...) or (0,1,0,...).

> sim
function(n,p) {
 y1 <- rnorm(n); y2 <- rnorm(n);
 x1 <- y1+y2; x2 <- y1-y2; p2 <- p - 2
 x <-</pre>

cbind(x1,x2,matrix(rnorm(n\*p2),ncol=p2))
cvr <- cov(x)
which.max(</pre>

abs(eigen(cvr,symmetric=T)\$vectors[,1]))

Norm Matloff University of California at Davis

# Simulation, cont'd.

Norm Matloff University of California at Davis

### Return value from **sim()** should be 1 or 2. Let's see: > sim(500, 400)[1] 1 > sim(500,800)[1] 1 > sim(500,800)[1] 2 > sim(500, 1200)[1] 439 > sim(500, 1200)[1] 2 > sim(500, 1200)[1] 1 > sim(500, 1200)[1] 905

## Simulation, cont'd.

Norm Matloff University of California at Davis

# Simulation, cont'd.

Norm Matloff University of California at Davis

### Simulation, cont'd.

• When n < p/2—very common in practice!—sometimes right but sometimes get phantom PCs.

Norm Matloff University of California at Davis

# Simulation, cont'd.

- When n < p/2—very common in practice!—sometimes right but sometimes get phantom PCs.
- On the other hand, results of Johnstone (2000) suggest that as long as n > p/2 we might be OK.

Norm Matloff University of California at Davis

# Simulation, cont'd.

- When n < p/2—very common in practice!—sometimes right but sometimes get phantom PCs.
- On the other hand, results of Johnstone (2000) suggest that as long as n > p/2 we might be OK.
- Moreover, in practice the variables are correlated, often very highly so, in regular patterns. I suspect this makes it "more OK."

Norm Matloff University of California at Davis

# Exorcizing the Curse?

▲□▶ ▲圖▶ ▲臣▶ ▲臣▶ 三臣 - のへで

Norm Matloff University of California at Davis

### Exorcizing the Curse?

◆□▶ ◆□▶ ◆臣▶ ◆臣▶ 臣 の�?

• The term *curse of dimensionality* goes back 50 years.

Norm Matloff University of California at Davis

# Exorcizing the Curse?

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のQ@

- The term *curse of dimensionality* goes back 50 years.
- In last 10-15 years, it has gotten scarier:

Norm Matloff University of California at Davis

# Exorcizing the Curse?

- The term *curse of dimensionality* goes back 50 years.
- In last 10-15 years, it has gotten scarier: Berry(1999) proved that any 2 points are approximately the same distance from each other!

Norm Matloff University of California at Davis

# Exorcizing the Curse?

- The term *curse of dimensionality* goes back 50 years.
- In last 10-15 years, it has gotten scarier: Berry(1999) proved that any 2 points are approximately the same distance from each other!
- So, e.g., nearest-neighbor methods look iffy.

Norm Matloff University of California at Davis

# Exorcizing the Curse?

- The term *curse of dimensionality* goes back 50 years.
- In last 10-15 years, it has gotten scarier: Berry(1999) proved that any 2 points are approximately the same distance from each other!
- So, e.g., nearest-neighbor methods look iffy.
- My own rough derivation:

Norm Matloff University of California at Davis

# Exorcizing the Curse?

- The term *curse of dimensionality* goes back 50 years.
- In last 10-15 years, it has gotten scarier: Berry(1999) proved that any 2 points are approximately the same distance from each other!
- So, e.g., nearest-neighbor methods look iffy.
- My own rough derivation:
  - Suppose the p distance components are iid.

Norm Matloff University of California at Davis

# Exorcizing the Curse?

- The term *curse of dimensionality* goes back 50 years.
- In last 10-15 years, it has gotten scarier: Berry(1999) proved that any 2 points are approximately the same distance from each other!
- So, e.g., nearest-neighbor methods look iffy.
- My own rough derivation:
  - Suppose the p distance components are iid.
  - $\sqrt{Var(distance)}/E(distance) \rightarrow 0$  as  $p > \infty$

Norm Matloff University of California at Davis

# Exorcizing the Curse?

- The term *curse of dimensionality* goes back 50 years.
- In last 10-15 years, it has gotten scarier: Berry(1999) proved that any 2 points are approximately the same distance from each other!
- So, e.g., nearest-neighbor methods look iffy.
- My own rough derivation:
  - Suppose the p distance components are iid.
  - $\sqrt{Var(distance)}/E(distance) \rightarrow 0$  as  $p > \infty$
  - So, distances are approximately constant.

Norm Matloff University of California at Davis

### Some Hope

<□ > < @ > < E > < E > E のQ @



#### Norm Matloff University of California at Davis

# Some Hope

### Some Hope:

• But all that involves equally-weighted components in distance.

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Long Live (Big Data-Fied) Statistics!

Norm Matloff University of California at Davis

- But all that involves equally-weighted components in distance.
- Yet, arguably we should have weights, according to importance of the variables.

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Long Live (Big Data-Fied) Statistics!

Norm Matloff University of California at Davis

- But all that involves equally-weighted components in distance.
- Yet, arguably we should have weights, according to importance of the variables.
- Then the above problem goes away. (Coef. of var. does not go to 0.)

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Long Live (Big Data-Fied) Statistics!

Norm Matloff University of California at Davis

- But all that involves equally-weighted components in distance.
- Yet, arguably we should have weights, according to importance of the variables.
- Then the above problem goes away. (Coef. of var. does not go to 0.)
- But how set the weights?

▲ロト ▲帰ト ▲ヨト ▲ヨト 三日 - の々ぐ

Long Live (Big Data-Fied) Statistics!

Norm Matloff University of California at Davis

- But all that involves equally-weighted components in distance.
- Yet, arguably we should have weights, according to importance of the variables.
- Then the above problem goes away. (Coef. of var. does not go to 0.)
- But how set the weights?
- Stay tuned...

Norm Matloff University of California at Davis

### Misc.

<□ > < @ > < E > < E > E のQ @

#### Norm Matloff University of California at Davis

# Misc.

### Online materiasl:

The visualization code is available for your use and comments/suggestions:

http://heather.cs.ucdavis.edu/BigDataVis.html
These slides are there too.

#### **Acknowlegements:**

The author would like to thank Noah Gift, Marnie Dunsmore, Nicholas Lewin-Koh, and David Scott for helpful discussions, and Hao Chen and Bill Hsu for use of their high-performance computing equipment.