# Spatial Data and Analysis

## Discussion 2

**Felipe González**

UC Berkeley

September 11$th$

# Admin stuff

1. Check the discussion syllabus in the website

2. Remember to use "[spatial]" in the subject header

3. Check discussion slides, there is useful stuff

# Outline

# Miscellaneous #1: vector and locations

- For simplicity think of a vector as a $n \times 1$ matrix: $\vec{a} = \begin{bmatrix} a_1 \\ \vdots \\ a_n \end{bmatrix}$

- $N$ locations can be described by 2 vectors of size $N$:

$$(\vec{x}, \vec{y}) = \begin{bmatrix} x_1 & y_1 \\ \vdots & \vdots \\ x_n & y_n \end{bmatrix}$$

  where $(x_i, y_i)$ is a single location with $i = 1, \ldots, N$

- Random locations: $(\vec{x}, \vec{y})$ is `[randn(N, 1) randn(N,1)]`

# Miscellaneous #2: matlab v/s stata

▶ Matlab reads the code `plot(A,B)` and plots *A* in the *x*-axis and *B* in the *y*-axis ∼ `plot(x,y)`

▶ Stata reads the code `scatter A B` and plots *A* in the *y*-axis and *B* in the *x*-axis ∼ `scatter y x`

▶ In a way, matlab thinks more mathematically and stata more statistically.

▶ When plotting locations in matlab: `plot(lon,lat)`, <u>not</u> the perhaps more intuitive `plot(lat,lon)`

# Logical operators

▶ Basic logical operators we will use repeatedly:

$$
\begin{array}{rcl}
\text{EQUAL} & : & = \\
\text{GREATER} & : & > \\
\text{GREATER OR EQUAL} & : & >= \\
\text{AND} & : & \& \\
\text{OR} & : & | \\
\text{IS NOT} & : & \sim=
\end{array}
$$

▶ Useful to create indicator variables and others

# Example

- Examples of logical operators:

```
1   % Parameters
2       A = [1 0 1];
3       B = [0 0 1];
4       C = [0 4 0];
5
6   % Operations
7       C = (A==1 & B==1) ;
8       D = (A==1 | B==1) ;
9       E = A(A >1)        ;
10      F = C(C >1)        ;
```

# Example (cont.)

▶ You very rarely have to loop through a whole matrix to get only specific parts of it

▶ Exercise: find an approximate value for the mean of negative numbers in a normal distribution with 1,000 values

```
1       X = randn(1000,1) ;
2       mean(X(X<0));
```

# if

- Sometimes `if` statements can be useful:

```matlab
1  % Parameters
2      N = 100                 ;
3      A = ceil(randn(N,1));
4      B = NaN(N,1)            ;
5
6  % If statement
7      if   A(10, 1) == 1
8           disp('Hi everyone')
9      else disp('Its hot')
10     end
```

# Comments

▶ Note that the command == is used when using a statement, and = when assigning values

▶ Short circuit behavior:
$$A \ \& \ B \ : \ A \text{ and } B \text{ are evaluated}$$
$$A \ \&\& \ B \ : \ B \text{ only evaluated if } A \text{ is true}$$

▶ Multiple statements need to be separated by parentheses:

```
1  % Statement 1
2      (x == 3 & x > 4) | x < 4
3  % Statement 2
4      x == 3 & (x > 4 | x < 4)
```

# Loops

▶ Loops allow you to execute a block of code many times. For example, create a vector $x$ with ten ones using a loop:

```
1  for i = 1 : 10
2      x(i, 1) = 1 ;
3  end
```

▶ However, this is obviously inefficient since we can create the same vector simply typing `x=ones(10,1)`.

▶ Knowing when to use a loop is an important part of being efficient at coding. Indentation is important for readability

# Example

- Create a fake dataset with panel structure where variables are 100 ID's and 10 years:

| ID | Year |
|----|------|
| 1 | 1 |
| ⋮ | ⋮ |
| 1 | 10 |
| ⋮ | ⋮ |
| 100 | 1 |
| ⋮ | ⋮ |
| 100 | 10 |

# Solution

```
1   % Parameters
2       ID = 100          ;
3       t  = 10           ;
4       x  = NaN(ID*t, 2) ;
5
6   % Loop
7       for i = 1:ID
8           x((i-1)*10+1:i*10, 1) = i     ;
9           x((i-1)*10+1:i*10, 2) = 1:10  ;
10      end
```

# Monte Carlo

- Monte Carlo is an algorithm that uses random sampling to obtain numerical results

- This method is useful in applied statistics to obtain quick answers to a wide range of questions

- Consider $y_i = 5 + 0.5x_i + \varepsilon_i$. Then:
    1. what is the distribution of an OLS estimator using $(y, x)$?
    2. what are the consequences of classical measurement error in $x$?

- Assume $x \sim N(0, 1)$, $\varepsilon \sim N(0, 1)$, $x \perp \varepsilon$ and different magnitudes for the measurement error

# Structure and parameters

▶ Think about the following structure to answer these questions: define parameters, procedure, output

```
1   % Parameters
2       clear ; clc
3       s = 1000    ;
4       N = 100     ;
5       b = NaN(s,1);
6       g = NaN(s,1);
```

▶ Always useful to think about what is the output of the process you're writing
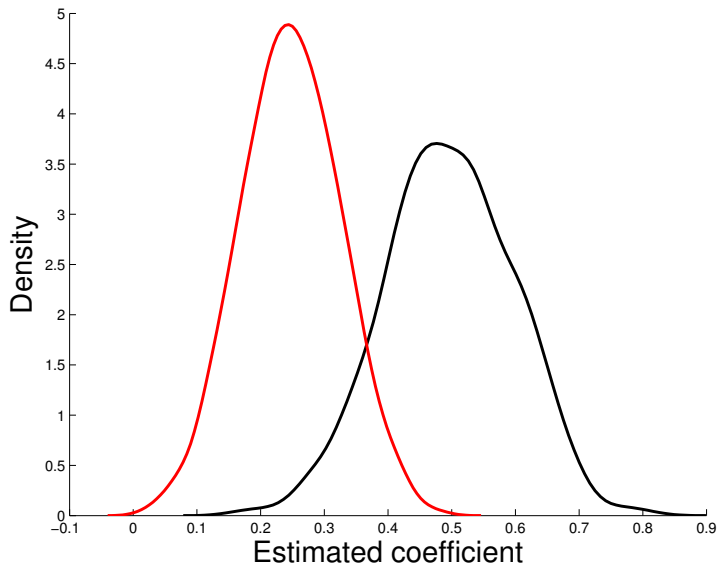
```matlab
1   % Loop
2       for i = 1:s
3       % Variables
4           x = randn(N, 1) ;
5           z = x - randn(N,1) ;
6           y = 5 + .5 * x + randn(N,1);
7       % OLS coefficients
8           X = [ones(N,1) x] ;
9           beta = (X' * X) \ (X' * y) ;
10          b(i,1) = beta(2, 1);
11      % Attenuated OLS coefficients
12          Z = [ones(N,1) z] ;
13          gamma = (Z' * Z) \ (Z' * y) ;
14          g(i,1) = gamma(2, 1);
15      end
```

# Output

```
1   % Figure
2       [a,b] = ksdensity(b) ;
3       [c,d] = ksdensity(g) ;
4
5       hold on
6       plot(b, a, 'k', 'LineWidth', 2)
7       plot(d, c, 'r', 'LineWidth', 2)
8           ylabel('Density'    , 'FontSize' , 20)
9           xlabel('Coefficient', 'FontSize' , 20)
10          box off
11          set(gcf, 'Color', 'w')
12      hold off
13
14      export_fig 'ols.pdf'
```

Figure: *Distribution of OLS coefficient and attenuation bias*
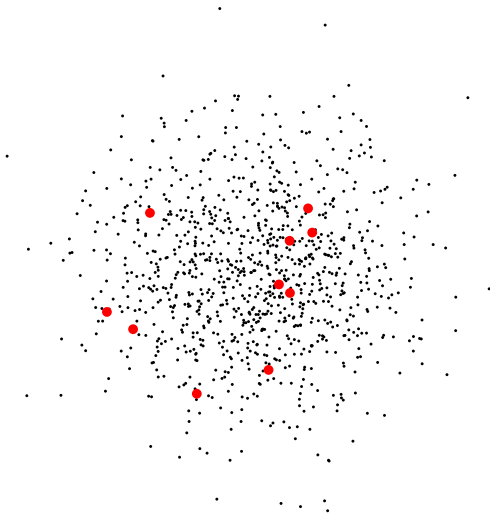
# Application: school choice

- Consider a simple model of school choice in which students enroll in the school that is closer to their home

- As researcher you are worried that students' locations could be measured with error

- This is crucial if we want to predict how many students will enroll in each school (e.g., construction of schools)

- Use point processes, distances, Monte Carlo and the model of school choice to explore the consequences of measurement error

```matlab
1  % Parameters
2      N     = 1000    ; % students
3      S     = 10      ; % schools
4      x     = 0:.01:1 ; % size of noise
5      ERROR = NaN(N, size(x, 2)) ;
6  % Location of schools
7      lat_s = randn(S,1)    ;
8      lon_s = randn(S,1)    ;
9      LOC_s = [lon_s lat_s] ;
10 % Location of students
11     lat_n = randn(N,1)    ;
12     lon_n = randn(N,1)    ;
13     LOC_n = [lon_n lat_n] ;
14 % Distance and school choice
15     D = pdist2(LOC_n, LOC_s, 'euclidean') ;
16     [d, c] = min(D,[],2) ;
```

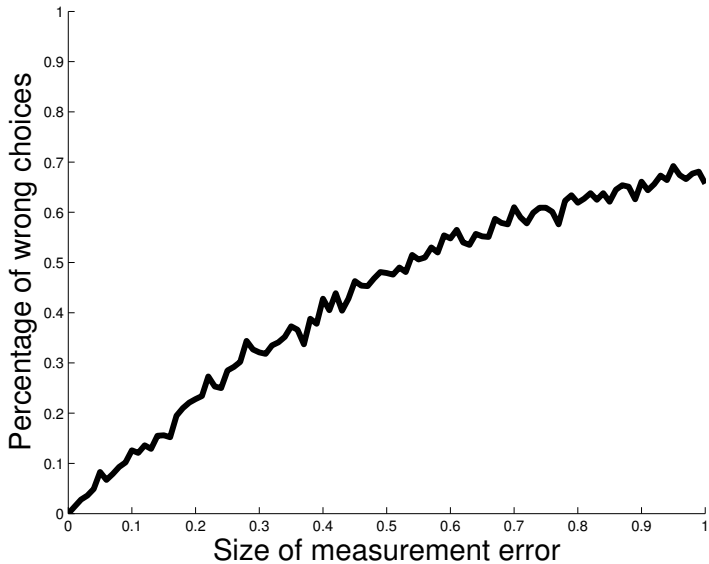Figure: *Schools and students in space*

```matlab
% Loop for different types of noise
    for i = 1:size(x,2)
    % Noisy location of students
        lat_n_n = lat_n + x(1,i)*randn(N,1) ;
        lon_n_n = lon_n + x(1,i)*randn(N,1) ;
        LOC_n_n = [lon_n_n lat_n_n]   ;
    % Noisy distance from students to schools
        D_n = pdist2(LOC_n_n, LOC_s, 'euclidean') ;
    % Noisy school choice
        [d_n, c_n] = min(D_n,[],2) ;
    % Actual and noisy school choice
        ERROR(:, i) = (c ~= c_n) ;
    end
% Percentage of errors in choices
    result = sum(ERROR, 1)./N ;
```

# Output

```
1  % Figure
2  plot(x, result, 'k', 'LineWidth', 4)
3    ylabel('Percentage of wrong choices', 'FontSize', 20)
4    xlabel('Size of measurement error'  , 'FontSize', 20)
5    axis([0 1 0 1])
6    box off
7    set(gcf, 'Color', 'w')
8
9    export_fig 'error_1.pdf'
```

Figure: *Mistakes due to measurement error*

# Incorporating uncertainty

```matlab
1   % Loop for different location of students
2       for j = 1:M
3       % Location of students
4           lat_n = randn(N,1)     ;
5           lon_n = randn(N,1)     ;
6           LOC_n = [lon_n lat_n] ;
7       % Distance and school choice
8           D = pdist2(LOC_n, LOC_s, 'euclidean') ;
9           [d, c] = min(D,[],2) ;
10      % Loop for different types of noise
11          for i = 1 : size(x, 2)
12          % Noisy location of students
13              lat_n_n = lat_n + x(1,i)*randn(N,1) ;
```
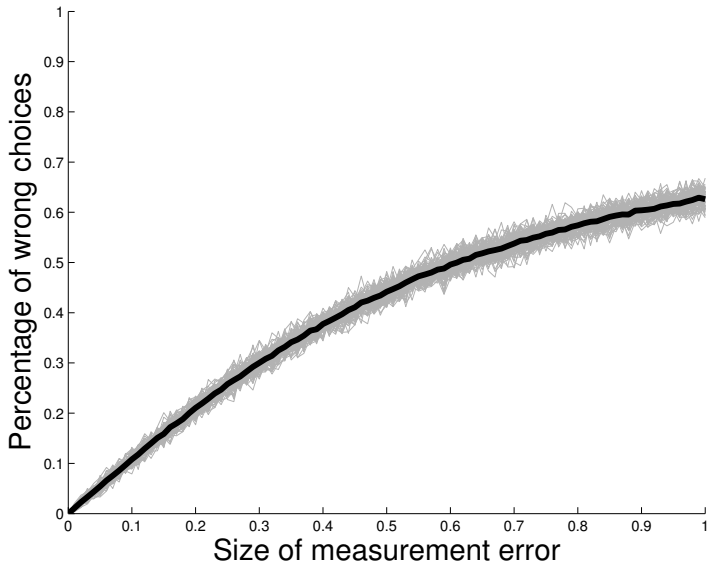
# Incorporating uncertainty

```matlab
14                    lon_n_n = lon_n + x(1,i)*randn(N,1) ;
15                    LOC_n_n = [lon_n_n lat_n_n]    ;
16            % Noisy distance and noisy school choice
17                    D_n = pdist2(LOC_n_n, LOC_s, 'euclidean') ;
18                    [d_n, c_n] = min(D_n,[],2) ;
19            % Actual and noisy school choice
20                    ERROR(:, i) = (c ~= c_n) ;
21            end
22        % Percentage of errors in choices
23            result(j,:) = sum(ERROR, 1)./N ;
24        end
```

# Output

```
1  % Plot error
2  hold on
3  plot(x, result, 'Color', [.7 .7 .7], 'LineWidth', .5)
4  plot(x, mean(result), 'k', 'LineWidth', 4 )
5   ylabel('Percentage of wrong choices', 'FontSize' , 20)
6   xlabel('Size of measurement error'  , 'FontSize' , 20)
7   axis([0 1 0 1])
8   box off
9   set(gcf, 'Color', 'w')
10  hold off
11
12  export_fig 'error_2.pdf'
```

Figure: *Incorporating uncertainty using Montecarlo*

# Additional resources

1. Explore `export_fig`

2. Explore `mcode`

3. Explore options of commands we have used (e.g, `plot`, `pdist2`)