# Agglomeration Economies
## Urban Economics

Ignacio Sarmiento-Barbieri

Universidad de los Andes
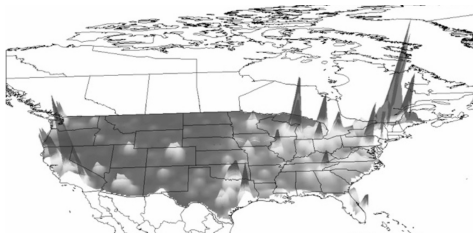
September 23, 2024

# Agenda

# Agglomeration Economies

▶ Why do we see such a remarkable clustering of human activity in a small number of urban areas?
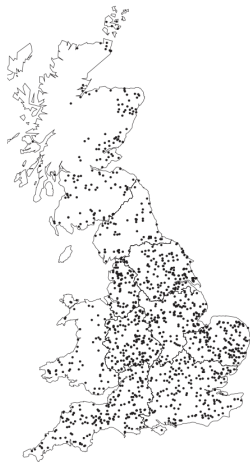
# Evidence of Agglomeration Economies

▶ Three strategies to identify agglomeration economies

1. Show there is too much spatial concentration to be random (Duranton and Overman, 2005)

2. Compare productivity over space (Greenstone et al., 2010)

3. Compare wages and rents across space (Quantitative Spatial Models, Ahlfeldt et al, 2015)

# Agenda

# Spatial Concentration
## Extremes of Localization and Dispersion



(c) Other Agricultural and Forestry Machinery (SIC2932)

(d) Machinery for Textile, Apparel and Leather Production (SIC2954)

# Spatial Concentration
## Ambiguous Cases



(a) Basic Pharmaceuticals
(SIC2441)

(b) Pharmaceutical Preparations
(SIC2442)

# Duranton & Overman Methodology

1. **Select Relevant Establishments:**
   - Choose establishments based on industry and size.
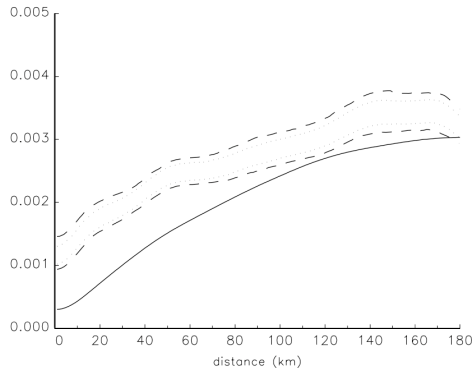   - Consider different thresholds to assess robustness (e.g., include only those contributing to 90% of employment).

2. **Compute Bilateral Distances:**
   - Calculate Euclidean distances between all pairs of establishments.
   - Use Kernel Density Estimation (KDE) to estimate the density of these distances.
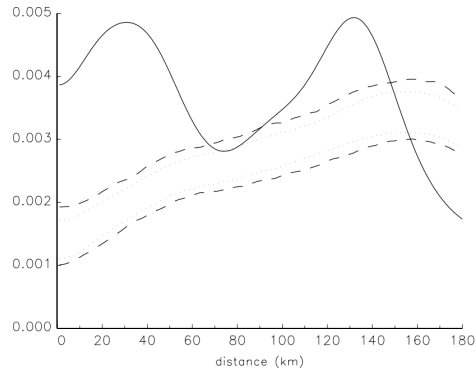
# Spatial Concentration

### K Density Estimates



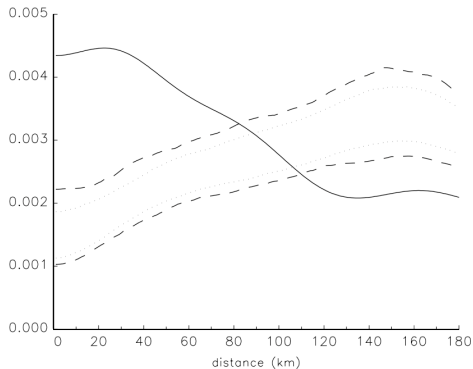(c) Other Agricultural and Forestry Machinery (SIC2932)

(d) Machinery for Textile, Apparel and Leather Production (SIC2954)
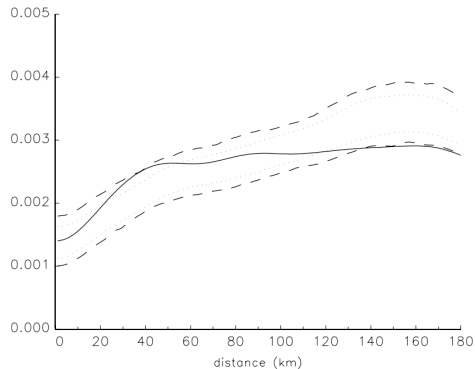
**Figure 2**. *K*-density, local confidence intervals and global confidence bands for four illustrative industries

# Spatial Concentration
## Ambiguous Cases



(a) Basic Pharmaceuticals
(SIC2441)

(b) Pharmaceutical Preparations
(SIC2442)

# Duranton & Overman Methodology

4 **Generate Counterfactuals:**
  - ▶ Randomly assign establishments to locations, maintaining the number of establishments and industrial concentration.
  - ▶ Create 1,000 simulations to construct a baseline for comparison.

5 **Statistical Significance:**
  - ▶ Compare actual densities with simulated counterfactuals to determine if localization is significant.
  - ▶ Use local and global confidence intervals to assess statistical significance of localization.

6 **Localization Metrics:**
  - ▶ Define indices for localization ($\gamma$) and dispersion ($\psi$) at each distance.
  - ▶ Determine global localization or dispersion based on these indices over all distances.

# Agenda

# Agenda

# Univariate density estimation: parametric vs nonparametric methods

▶ Let $f = f(.)$ be the density function of the random variable $X$

▶ Let $x_1, x_2, \ldots x_n$ be a random sample of $X$. Then, $x_i \sim f$ *iid*

▶ How can we estimate the density in a particular point $x_0$, then $f(x_0)$?

# Agenda

# Parametric methods

▶ They assume a particular functional form for $f$.

# Parametric methods

▶ They assume a particular functional form for $f$.

▶ Example:

$$f\left(x_o\right) = \phi\left(x_o\right) = \frac{1}{\sigma\sqrt{2\,\pi}}\exp\left[-\frac{1}{2}\left(\frac{x_o - \mu}{\sigma}\right)^2\right] \tag{1}$$

# Agenda

# Non-parametric methods

▶ They seek to estimate $f(x_0)$) without assuming a particular functional form, only assuming certain regularity conditions of the density (smoothness, differentiability)

# A rudimentary non-parametric estimator: the histogram



▶ The height of the bars is an estimator of the density at any point in the interval

$$\hat{f}(x_o) = \frac{Nber.\,obs\ interval}{n \times w} \tag{2}$$

# Histogram problems: (1) depends on starting point

# Histogram problems: (2) depends on the bandwidth

# Histogram problems: (2) depends on the bandwidth

# Histogram problems: (3) is discontinuous at the ends of the interval

- Note that $\hat{f}(x_1) = \hat{f}(x_2)$
- But $\hat{f}(x_2 + \epsilon) = \frac{1}{4}\hat{f}(x_2)$ for any $\epsilon > 0$

# Another non-parametric estimator: the naive estimator
Discrete case

- X is a discrete RV, *iid*

- Objective: estimate $\Pr(X = x_o) = f(x_o)$

- The naive estimator comes from:

$$\hat{f}(x_o) = \frac{\# \, x_i = x_o}{n} = \frac{1}{n} \sum_{i=1}^{n} I(x_i = x_o) \tag{3}$$

# Another non-parametric estimator: the naive estimator

## Continuous case

- X is a continuous RV ($Pr(X = x_0)$) we evaluate the probability that X is "close" to $x_0$.

- We say that $x_i$ is close to $x_0$ if $x_i$ belongs to the interval $(x_o - h, \ x_o + h)$

# The naive estimator

Notice:

- Every possible point $x_0$ is the center of an interval

- Observations that are within that interval (less than one h away from $x_0$) are weighted by 1/2



Weight of each observation
$$\frac{1}{2}I(x_0 - h < x_i < x_0 + h)$$

$x_0$

# The naive estimator

Notice:

▶ To estimate the density at points a, b and c we construct intervals around them

▶ Then, unlike what happens in the histograms, the intervals of the naive estimator overlap (the initial point no longer matters)

# The naive estimator
Problems

▶ It depends on h: the larger the bandwidth, the more distant observations from $x_0$ are used to estimate $f(x_0)$. The higher h, the smoother the estimated density (h=smoothing parameter)



▶ The weight $1/2I(.)$ is discontinuous at the limits of each interval, generating discontinuities in the estimated density.

▶ The weight treats observations very close to $x_0$ in the same way as others somewhat further away, as long as they belong to the interval of length 2h around $x_0$

# The weighted average estimator or kernel method

▶ The kernel estimator is a generalization of the naive estimator that overcomes some of the deficiencies of the latter.

▶ The weight $1/2I(.)$ of the naive estimator is replaced by a new weight K(.):

$$\hat{f}\ (x_o) = \frac{1}{n \times h} \sum_{i=1}^{n} K\ \left( \frac{x_i - x_0}{h} \right) \tag{4}$$

▶ Rosenblatt-Parzen density estimation approach

# The weighted average estimator or kernel method

▶ The weights or kernels are known functions that satisfy:

1. $K(\phi) \geq 0$
2. $\int_{-\infty}^{\infty} K(\phi) \, d\phi = 1$
3. $K(\phi) = K(-\phi)$
4. $E(K(\phi)) = 0$

# The weighted average estimator or kernel method
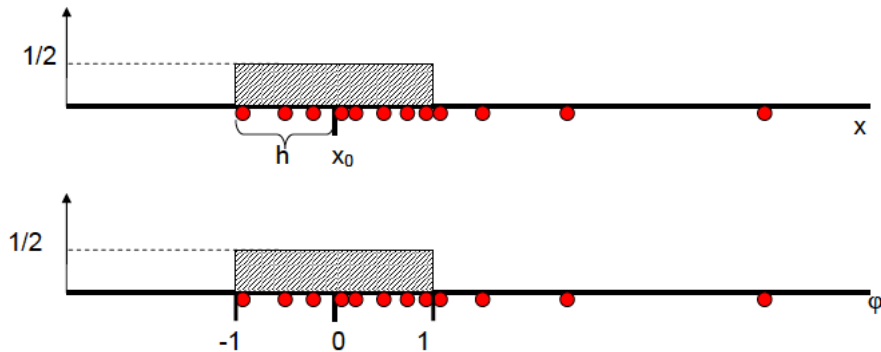Example 1: rectangular or uniform kernel

$$K\left(\phi\right) = \begin{cases} \frac{1}{2} & |\phi| < 1 \\ 0 & o.w. \end{cases} \tag{5}$$

where $\phi = \frac{x_i - x_o}{h}$ In this case, the Kernels estimator matches the naive estimator. Note that the weight is similar to a uniform density function on the interval (-1, 1)

# The weighted average estimator or kernel method

Example 1: rectangular or uniform kernel

# The weighted average estimator or kernel method

Example 2: Gaussian kernel

▶ Gaussian

$$K(\phi) = \frac{1}{\sqrt{2\pi}} \, e^{-\phi^2/2} \tag{6}$$

▶ The kernel function corresponds to the standard normal density function that satisfies the above assumptions.

▶ For this type of kernel, the density estimator is given by:

$$\hat{f}(x_o) = \frac{1}{n \times h} \sum_{i=1}^{n} \frac{1}{\sqrt{2\pi}} \, e^{-\left(\frac{x_i - x_o}{h}\right)^2/2} \tag{7}$$

▶ Important: we don't assume normal distribution for any variable. The functional form is just used to weight the sample observations

# The weighted average estimator or kernel method

Example 2: Gaussian kernel



$$\hat{f}(x_0) = \frac{1}{n \times h} \sum_{i=1}^{n} \frac{1}{\sqrt{2\pi}} e^{-\left(\frac{x_i - x_0}{h}\right)^2 \big/ 2}$$

# The weighted average estimator or kernel method

Example 2: Gaussian kernel (smaller h)



$$\hat{f}(x_0) = \frac{1}{n \times h} \sum_{i=1}^{n} \frac{1}{\sqrt{2\pi}} e^{-\left(\frac{x_i - x_0}{h}\right)^2 / 2}$$

# The weighted average estimator or kernel method
## Other Kernels

| Kernel type | Formula | Support |
|---|---|---|
| Gaussian or normal | $\dfrac{1}{\sqrt{2\pi}}\exp(-\varphi^2/2)$ | $\varphi \in (-\infty, \infty)$ |
| Epanechnikov | $\dfrac{3}{4\sqrt{5}}(1-\dfrac{1}{5}\varphi^2)I(\mid\varphi\mid<\sqrt{5})$ | $\varphi \in \left(-\sqrt{5}, \sqrt{5}\right)$ |
| Epanechnikov modificado | $\dfrac{3}{4}(1-\varphi^2)I(\mid\varphi\mid<1)$ | $\varphi \in (-1, 1)$ |
| Triangular | $(1-\mid\varphi\mid)I(\mid\varphi\mid<1)$ | $\varphi \in (-1, 1)$ |
| Uniform or rectangular | $\dfrac{1}{2}I(\mid\varphi\mid<1)$ | $\varphi \in (-1, 1)$ |

# The weighted average estimator or kernel method
Properties: Bias

▶ The kernel estimator is generally biased.

▶ The approximate expression for the bias is given by:

$$bias[\hat{f}(x_o)] \approx \frac{h^2}{2} f''(x_o) \int_{-\infty}^{\infty} K(\phi)\phi^2 d\phi \tag{8}$$

▶ The approximate expression for the asymptotic variance of the kernel estimator is given by:

$$variance[\hat{f}(x_o)] \approx \frac{1}{n \times h} f(x_o) \int_{-\infty}^{\infty} K^2(\phi) d\phi \tag{9}$$

# The weighted average estimator or kernel method
Bandwidth Selection

$$MSE(\hat{f}(x)) = Var(\hat{f}(x)) + \left(bias(\hat{f}(x))\right)^2 \tag{10}$$

# The weighted average estimator or kernel method
Bandwidth Selection

$$MSE(\hat{f}(x)) = Var(\hat{f}(x)) + \left( bias(\hat{f}(x)) \right)^2 \tag{10}$$

▶ We obtain a bandwidth which globally balances bias and variance by minimizing MSE with respect to h, i.e.,

$$h_{opt} = \left( \frac{\int K^2(z)dz}{(\int z^2 K(z)dz)^2 \int f'(x)^2 dx} \right)^{-1/5} n^{-1/5} \tag{11}$$

# The weighted average estimator or kernel method
Bandwidth Selection

- ▶ There are two popular approaches to bandwidth selection,
    1. Rules-of-thumb,

# The weighted average estimator or kernel method

Rule-of-Thumb

- The rule-of-thumb for choosing the bandwidth makes assumtions about $f$ and $K$

- For example: under a gaussian density and kernel

$$h_{opt} = 1.059\sigma n^{-1/5} \tag{12}$$

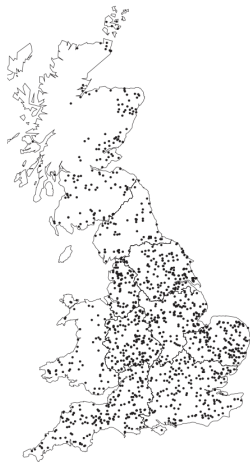# The weighted average estimator or kernel method
## Bandwidth Selection

▶ There are two popular approaches to bandwidth selection,

1. Rules-of-thumb,
2. cross-validation methods,

# Agenda

# Spatial Concentration
## Extremes of Localization and Dispersion

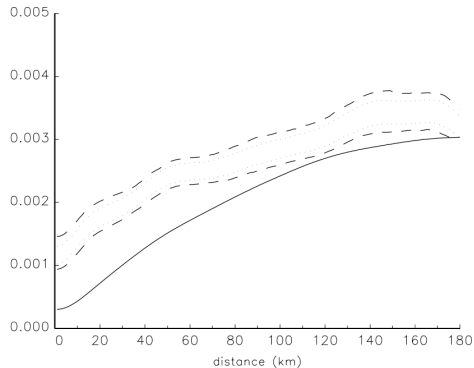

(c) Other Agricultural and Forestry
Machinery (SIC2932)

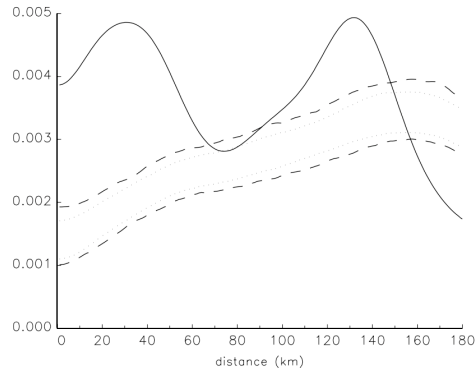(d) Machinery for Textile, Apparel and
Leather Production (SIC2954)

(c) Other Agricultural and Forestry
Machinery (SIC2932)

(d) Machinery for Textile, Apparel and
Leather Production (SIC2954)

**Figure 2**. *K*-density, local confidence intervals and global confidence bands for four illustrative industries
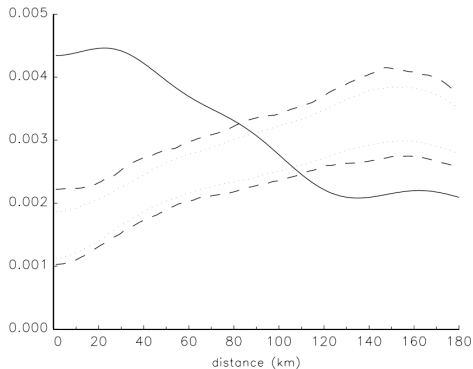
# Spatial Concentration
Ambiguous Cases



(a) Basic Pharmaceuticals
(SIC2441)

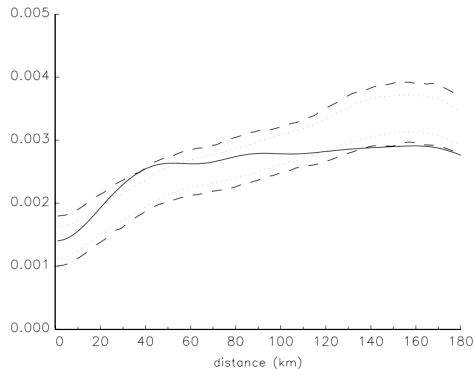(b) Pharmaceutical Preparations
(SIC2442)

# Spatial Concentration

Ambiguous Cases



(a) Basic Pharmaceuticals
(SIC2441)

(b) Pharmaceutical Preparations
(SIC2442)