

Lecture 1

Intro to Spatial and Temporal Processes

Dennis Sun
Stats 253

June 23, 2014

Outline of Lecture

① What is Spatial and Temporal Data?

Spatial Data

Temporal Data

Discussion

② Course Logistics

③ Linear Regression

④ Autoregressions

⑤ Recap

Where are we?

① What is Spatial and Temporal Data?

Spatial Data

Temporal Data

Discussion

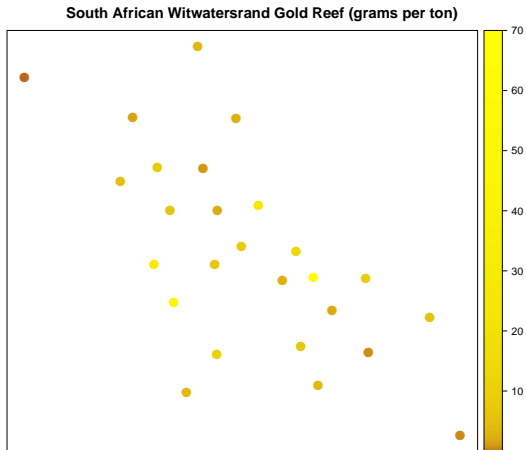
② Course Logistics

③ Linear Regression

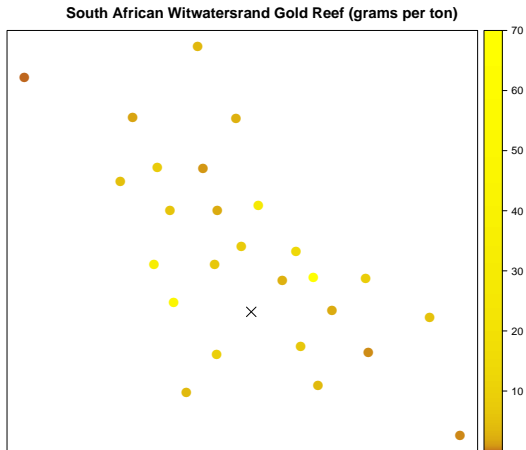
④ Autoregressions

⑤ Recap

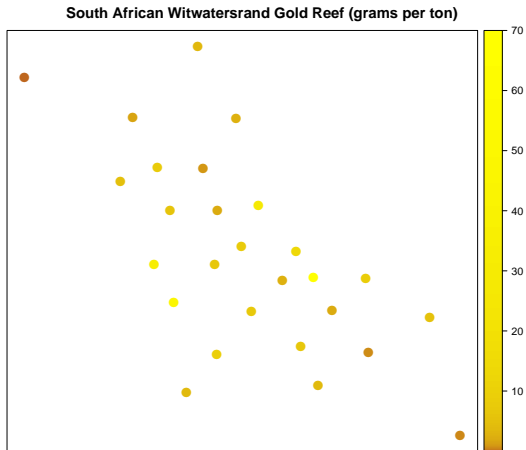
Geostatistics



Geostatistics

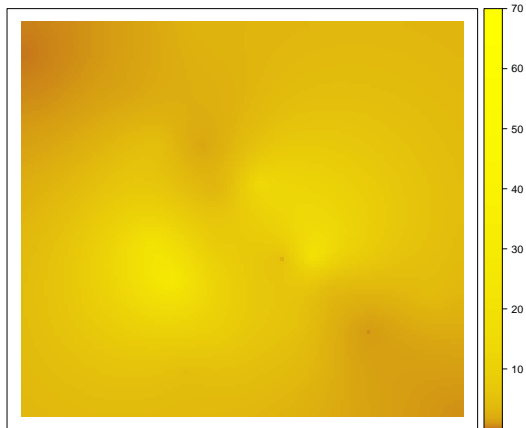


Geostatistics

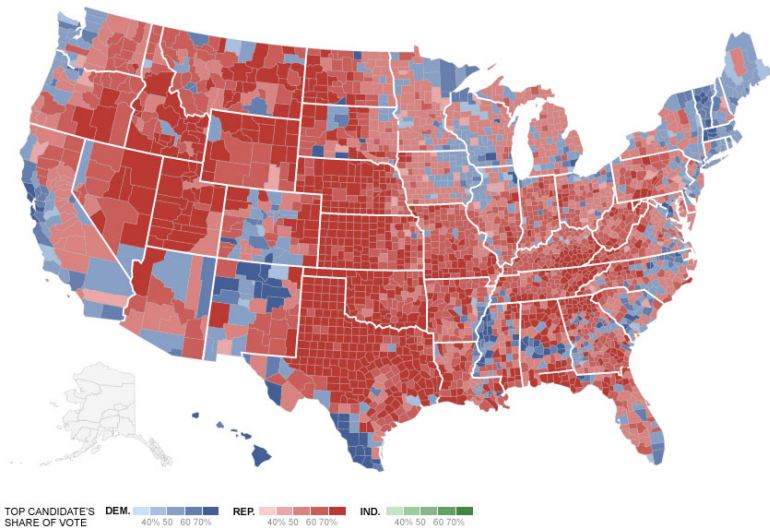


Geostatistics

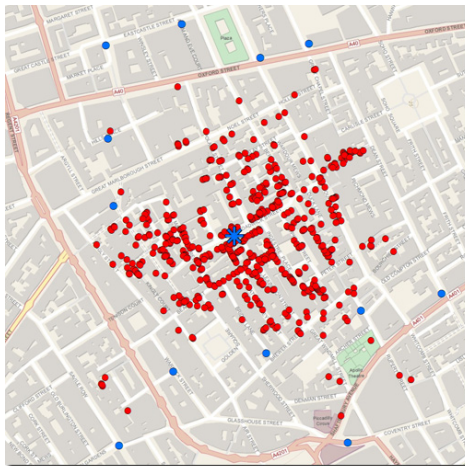
South African Witwatersrand Gold Reef (grams per ton)



Lattice (Areal) Data

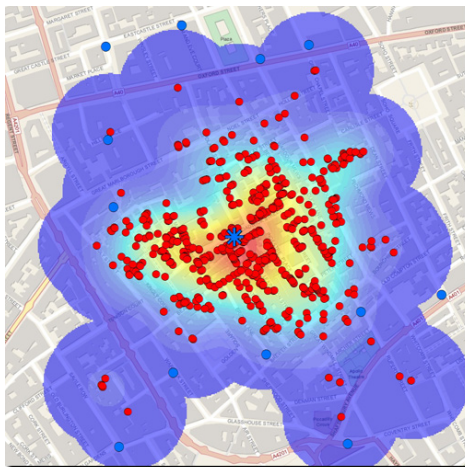


Point Processes



John Snow: 1854 Broad Street Cholera Outbreak

Point Processes



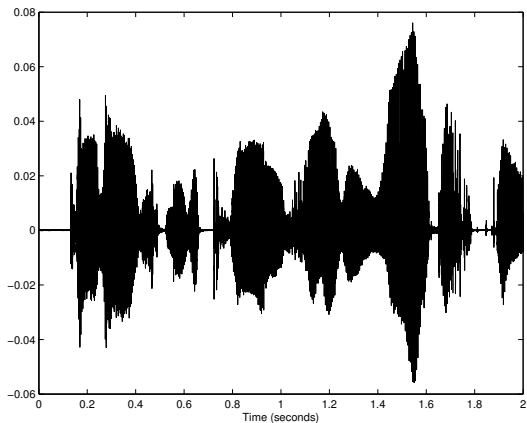
John Snow: 1854 Broad Street Cholera Outbreak

The Division of Spatial Statistics

Cressie (1993) organizes spatial statistics into these three categories.

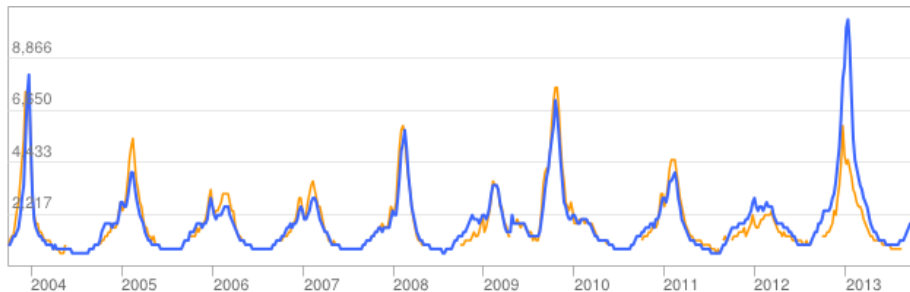
- I Geostatistics: continuous space, labeled observations, goal is prediction
- II Lattice (areal) data: discrete space, labeled observations, goal is inference
- III Point processes: continuous space, unlabeled observations, goal is inference

Time Series Example 1



Human Speech

Time Series Example 2



Google Flu Trends

What do space and time have in common?

The observations y_t (or y_s) are correlated:

$$\text{Cov}(y_t, y_{t'}) \neq 0 \text{ for } t \neq t'$$

Compare with the first assumption you see in most statistics courses:

Let y_i be i.i.d....

How are they different?

Time data is **ordered**, whereas there is no clear ordering for spatial data.

Where are we?

① What is Spatial and Temporal Data?

Spatial Data

Temporal Data

Discussion

② Course Logistics

③ Linear Regression

④ Autoregressions

⑤ Recap

Organization

- Lectures: Mondays and Wednesdays 2:15-3:30pm in Education 334
- Instructors:
 - Dennis Sun
 - Edgar Dobriban
 - Jingshu Wang
- Contact? Office Hours? Sections?
- All information can be found on the course website:
`stats253.stanford.edu`.

Content

- This is a new course.
- The material that we will be covering has not really been synthesized—because it is at the frontiers of statistics!
- We will be loosely following the books
 - Shumway and Stoffer. Time Series Analysis and Applications (with R Applications).
 - Sherman. Spatial Statistics and Spatio-Temporal Data.
- You don't have to purchase these books: they are available for free for Stanford students. (Link on course website.)
- Other useful references:
 - Bivand et al. Applied Spatial Data Analysis with R. (also available free)
 - Cressie and Wikle. Statistics for Spatio-Temporal Data.

Homeworks

- There will be about 4 short homeworks. Each homework will be a case study (data analysis).
- We will provide support for R, but you are free to use any computing environment (e.g., Python, Matlab, C....)
- You may work in pairs. If you do this, please turn in only one copy with both of your names on the front page.
- They will be graded on effort and completion only.

Project

- The goal of this course is to make a small but useful contribution to the world: e.g., a conference publication, open-source code, etc.
- This may sound intimidating, but there's actually a lot of low-hanging fruit in this subject!
- Grading rubric: produce something useful \Rightarrow A+.
- I have no qualms about giving everyone an A+ if you all earn it!

Course Requirements

- The grade will be based on the final project.
- If taking this class CR/NC, the final project is optional, but you are required to complete all homeworks.
- Please enroll in this class if you are able: I promise that you will get much more out of it! The homeworks will be short and instructional.

Where are we?

① What is Spatial and Temporal Data?

Spatial Data

Temporal Data

Discussion

② Course Logistics

③ Linear Regression

④ Autoregressions

⑤ Recap

Review of Linear Regression

Model: $y_i = \beta_1 x_{1i} + \dots + \beta_p x_{pi} + \epsilon_i$, $\epsilon_i \stackrel{iid}{\sim} N(0, \sigma^2)$

- Ordinary least squares: choose β_1, \dots, β_p by solving

Estimation: $\operatorname{argmin}_{\beta_1, \dots, \beta_p} \sum_{i=1}^n (y_i - (\beta_1 x_{1i} + \dots + \beta_p x_{pi}))^2$.

- Write in vector notation as:

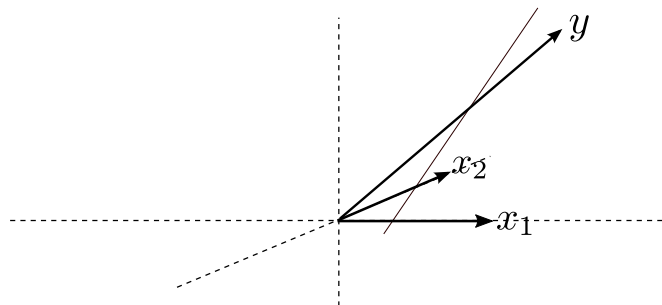
Model: $\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}$, $\boldsymbol{\epsilon} \sim N(0, \sigma^2 I)$

Estimation: $\hat{\boldsymbol{\beta}} = \operatorname{argmin}_{\boldsymbol{\beta}} \|\mathbf{y} - X\boldsymbol{\beta}\|^2$.

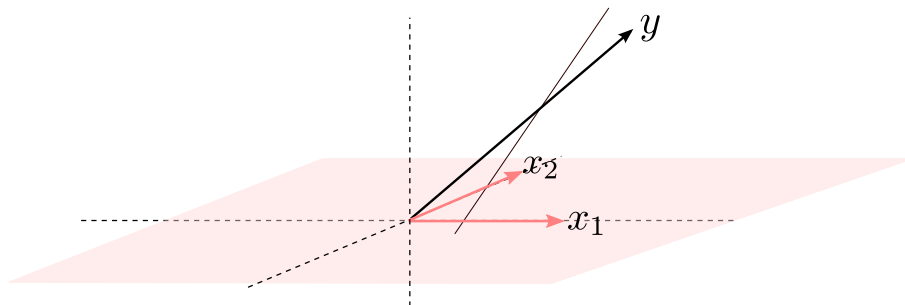
- Solve by differentiating: $\hat{\boldsymbol{\beta}}$ must satisfy

$$2X^T(\mathbf{y} - X\hat{\boldsymbol{\beta}}) = 0 \quad \Rightarrow \quad \hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y}.$$

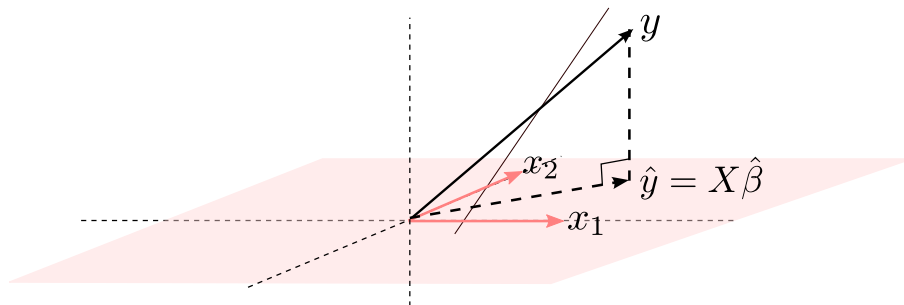
Regression: A Geometric Perspective



Regression: A Geometric Perspective



Regression: A Geometric Perspective



Properties of the Estimator

Recall that the model is

$$\mathbf{y} = X\boldsymbol{\beta} + \boldsymbol{\epsilon}, \quad \boldsymbol{\epsilon} \sim N(0, \sigma^2 I).$$

Is the estimator $\hat{\boldsymbol{\beta}}$ any good for estimating $\boldsymbol{\beta}$?

$$\hat{\boldsymbol{\beta}} = (X^T X)^{-1} X^T \mathbf{y} = \boldsymbol{\beta} + (X^T X)^{-1} X^T \boldsymbol{\epsilon}$$

$$E(\hat{\boldsymbol{\beta}}) = \boldsymbol{\beta}$$

$$\text{Var}(\hat{\boldsymbol{\beta}}) = \sigma^2 (X^T X)^{-1}$$

It is the **best linear unbiased predictor** (i.e., with the smallest variance).

Where are we?

① What is Spatial and Temporal Data?

Spatial Data

Temporal Data

Discussion

② Course Logistics

③ Linear Regression

④ Autoregressions

⑤ Recap

Introducing Dependence

- In linear regression:

$$\text{Cov}(y_i, y_j) = \text{Cov}(\mathbf{x}_i^T \boldsymbol{\beta} + \epsilon_i, \mathbf{x}_j^T \boldsymbol{\beta} + \epsilon_j) = \text{Cov}(\epsilon_i, \epsilon_j) = 0$$

- How can we introduce dependence?

$$y_t = \mathbf{x}_t^T \boldsymbol{\beta} + \phi y_{t-1} + \epsilon_t$$

- How do we estimate ϕ ?

Autoregression

- **Idea:** Write as

$$\underbrace{\begin{bmatrix} y_2 \\ \vdots \\ y_n \end{bmatrix}}_{\tilde{\mathbf{y}}} = \underbrace{\begin{bmatrix} & & y_1 \\ - & X_{2:n} & - \\ & & \vdots \\ & & y_{n-1} \end{bmatrix}}_{\tilde{X}} \begin{bmatrix} \beta \\ \vdots \\ \phi \end{bmatrix} + \epsilon$$

- Regress $\tilde{\mathbf{y}}$ on \tilde{X} to obtain estimates of β and ϕ .
- Hence, we call this an **auto - regressive** (AR) model, meaning “regress on itself.”
- Does this method “work”?

Does it work?

Let's set $X \equiv 0$ for now. So the model is

$$y_t = \phi y_{t-1} + \epsilon_t$$

The least squares estimate is

$$\hat{\phi} = (\mathbf{y}_{1:(n-1)}^T \mathbf{y}_{1:(n-1)})^{-1} \mathbf{y}_{1:(n-1)}^T \mathbf{y}_{2:n}$$

- Is it true that $E(\hat{\phi}) = \phi$? **No.**
- Is it true that $\text{Var}(\hat{\phi}) = \sigma^2 (\mathbf{y}_{1:(n-1)}^T \mathbf{y}_{1:(n-1)})^{-1}$? **No.**
- **However**, it turns out that $\hat{\phi}$ is consistent for ϕ .

$$\hat{\phi} \xrightarrow{p} \phi \text{ as } n \rightarrow \infty.$$

Simulation Example

- Suppose we observe 1000 observations of a random walk:

$$y_t = y_{t-1} + \epsilon_t, \quad \epsilon_t \sim N(0, 1)$$

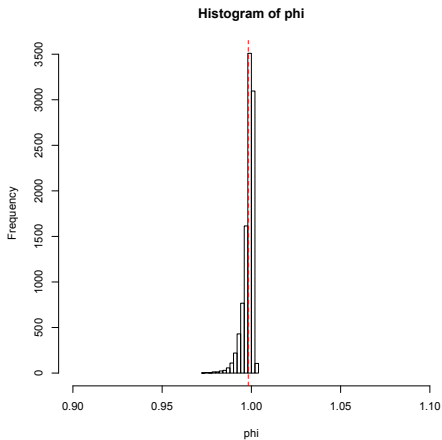
(In this case, $\phi = 1$.)

- Calculate $\hat{\phi}$ by regression.
- R Code:

```
phi <- sapply(1:10000, function(iter) {
  z <- cumsum(rnorm(1000))
  x <- z[1:999]
  y <- z[2:1000]
  return(sum(x*y)/sum(x*x))
})
```


Simulation Example

```
hist(phi, xlim=c(.9,1.1))  
abline(v=mean(phi), col='red', lty=2)
```



$\text{sd}(\text{phi}) = .003$

Simulation Example

- The true standard error of $\hat{\phi}$ is .003.
- Does this agree with what linear regression would tell us?
- R Code:

```
z <- cumsum(rnorm(1000))
x <- z[1:999]
y <- z[2:1000]
model <- lm(y~x-1)
summary(model)
```

Call:
lm(formula = y ~ x - 1)

Residuals:

Min	1Q	Median	3Q	Max
-3.2497	-0.6678	0.0396	0.6699	4.3311

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
x	0.998757	0.001835	544.3	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.9882 on 998 degrees of freedom

Simulation Example

- The true standard error of $\hat{\phi}$ is .003.
- Does this agree with what linear regression would tell us?
- R Code:

```
z <- cumsum(rnorm(1000))
x <- z[1:999]
y <- z[2:1000]
model <- lm(y~x-1)
summary(model)
```

Call:

```
lm(formula = y ~ x - 1)
```

Residuals:

Min	1Q	Median	3Q	Max
-3.2497	-0.6678	0.0396	0.6699	4.3311

Coefficients:

	Estimate	Std. Error	t value	Pr(> t)
x	0.998757	0.001835	544.3	<2e-16 ***

Signif. codes: 0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 0.9882 on 998 degrees of freedom

Good

Simulation Example

- The true standard error of $\hat{\phi}$ is .003.
- Does this agree with what linear regression would tell us?
- R Code:

```
z <- cumsum(rnorm(1000))
x <- z[1:999]
y <- z[2:1000]
model <- lm(y~x-1)
summary(model)
```

```
Call:
lm(formula = y ~ x - 1)
```

```
Residuals:
```

```
   Min       1Q   Median       3Q      Max
-3.2497 -0.6678  0.0396  0.6699  4.3311
```

```
Coefficients:
```

```
   Estimate Std. Error t value Pr(>|t|)
x  0.998757  0.001835    544.3  <2e-16 ***
```

```
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1
```

Good

Bad

```
Residual standard error: 0.9882 on 998 degrees of freedom
```

Conclusions

- Linear regression gives good estimates for the coefficients of an AR process.
- However, it tends to **underestimate** the error (when observations are positively correlated).
- This will make effects look more significant than they really are!
- How can we fix this? *Next lecture.*

Where are we?

① What is Spatial and Temporal Data?

Spatial Data

Temporal Data

Discussion

② Course Logistics

③ Linear Regression

④ Autoregressions

⑤ Recap

What We've Learned

- The similarities and differences between spatial and temporal data.
- Linear regression
- AR processes: the simplest model for correlated data
- The advantages and shortfalls of using regression to estimate parameters in AR processes.

`stats253.stanford.edu`