# Data Analysis: Statistical Modeling and Computation in Applications

Spatial and Environmental Data: Introduction, Local Correlations

- Environmental data
- Modeling flows
- Short-range spatial correlations
  - intuition
  - 2 variables
  - multiple variables

• Air Quality, Water Quality (pollutants)

- Air Quality, Water Quality (pollutants)
- Weather & climate (temperature, winds, moisture, precipitation, extreme conditions, ...)
- Storms

- Air Quality, Water Quality (pollutants)
- Weather & climate (temperature, winds, moisture, precipitation, extreme conditions, ...)
- Storms
- Ocean dynamics

- Air Quality, Water Quality (pollutants)
- Weather & climate (temperature, winds, moisture, precipitation, extreme conditions, ...)
- Storms
- Ocean dynamics
- Vegetation (forests, algae, ...)

- Air Quality, Water Quality (pollutants)
- Weather & climate (temperature, winds, moisture, precipitation, extreme conditions, ...)
- Storms
- Ocean dynamics
- Vegetation (forests, algae, ...)
- Wildlife monitoring

- Air Quality, Water Quality (pollutants)
- Weather & climate (temperature, winds, moisture, precipitation, extreme conditions, ...)
- Storms
- Ocean dynamics
- Vegetation (forests, algae, ...)
- Wildlife monitoring
- Earthquake magnitudes

- Air Quality, Water Quality (pollutants)
- Weather & climate (temperature, winds, moisture, precipitation, extreme conditions, ...)
- Storms
- Ocean dynamics
- Vegetation (forests, algae, ...)
- Wildlife monitoring
- Earthquake magnitudes

Why do we care? What are questions of interest?

#### • understand underlying processes, changes

• understand underlying processes, changes

(e.g. climate change: "statistics of weather over time")

• impacts on environment, health, economics, society

• understand underlying processes, changes

- impacts on environment, health, economics, society
- policies

• understand underlying processes, changes

- impacts on environment, health, economics, society
- policies
- forecast events, warnings (e.g. community seismic network, storms, ...)

• understand underlying processes, changes

- impacts on environment, health, economics, society
- policies
- forecast events, warnings (e.g. community seismic network, storms, ...)
- resource/energy management, e.g., water, renewable energies

• understand underlying processes, changes

- impacts on environment, health, economics, society
- policies
- forecast events, warnings (e.g. community seismic network, storms, ...)
- resource/energy management, e.g., water, renewable energies
- use in planning, routing, backtracking, control (ships, airplanes, ...)

• understand underlying processes, changes

(e.g. climate change: "statistics of weather over time")

- impacts on environment, health, economics, society
- policies
- forecast events, warnings (e.g. community seismic network, storms, ...)
- resource/energy management, e.g., water, renewable energies
- use in planning, routing, backtracking, control (ships, airplanes, ...)

#### Questions

understand underlying processes, changes

(e.g. climate change: "statistics of weather over time")

- impacts on environment, health, economics, society
- policies
- forecast events, warnings (e.g. community seismic network, storms, ...)
- resource/energy management, e.g., water, renewable energies
- use in planning, routing, backtracking, control (ships, airplanes, ...)

#### Questions

- relationships (correlations, association)
- trends; forecasting

• understand underlying processes, changes

(e.g. climate change: "statistics of weather over time")

- impacts on environment, health, economics, society
- policies
- forecast events, warnings (e.g. community seismic network, storms, ...)
- resource/energy management, e.g., water, renewable energies
- use in planning, routing, backtracking, control (ships, airplanes, ...)

#### Questions

- relationships (correlations, association)
- trends; forecasting
- planning
- quantifying uncertainty, adaptive sensing

#### Environmental Data – What is special?



Image source: H. Kaper, H. Engler. Mathematics & Climate.

#### Environmental Data – What is special?



underlying physical processes; scientific models

Image source: H. Kaper, H. Engler. Mathematics & Climate.

Stefanie Jegelka (and Caroline Uhler)

## Environmental Data - What is special?



## Environmental Data - What is special?



#### spatial and temporal correlation

```
from https://www.epa.gov/outdoor-air-quality-data,
http://public.dep.state.ma.us/MassAir/Pages/MapCurrent.aspx?&ht=1&hi=101
```

#### Environmental Data – What is special?



Figure 1.8. Dipoles in NCEP sea-level data for the period 1948–1967. The color background shows the regions of high activity. The edges represent dipole connections between regions.

#### Environmental Data – What is special?



Figure 1.8. Dipoles in NCEP sea-level data for the period 1948–1967. The color background shows the regions of high activity. The edges represent dipole connections between regions.

#### long-range spatial (anti-)correlations

Source: H. Kaper, H. Engler. Mathematics & Climate.

- correlations in time
- correlations in space (fields)
- scientific models + statistics; simulations
- methodological challenges for statistics

- correlations in time
- correlations in space (fields)
- scientific models + statistics; simulations
- methodological challenges for statistics
  - no controlled studies, only observational data

- correlations in time
- correlations in space (fields)
- scientific models + statistics; simulations
- methodological challenges for statistics
  - no controlled studies, only observational data
  - hypotheses from data

- correlations in time
- correlations in space (fields)
- scientific models + statistics; simulations
- methodological challenges for statistics
  - no controlled studies, only observational data
  - hypotheses from data
  - large data sets

- correlations in time
- correlations in space (fields)
- scientific models + statistics; simulations
- methodological challenges for statistics
  - no controlled studies, only observational data
  - hypotheses from data
  - large data sets
  - heterogeneous data

- correlations in time
- correlations in space (fields)
- scientific models + statistics; simulations
- methodological challenges for statistics
  - no controlled studies, only observational data
  - hypotheses from data
  - large data sets
  - heterogeneous data
  - proxy data

- Environmental data
- Modeling flows
- Short-range spatial correlations
  - intuition
  - 2 variables
  - multiple variables



#### Data for the problem set: flow in space & time

## Flow fields

• V(x, t): flow at location x at time t vector in 2d or 3d

1 7 7 7 1 1 1 1  $\uparrow \uparrow \uparrow$ ≁ V(x, t) at a fixed time t

## Flow fields

- V(x, t): flow at location x at time t vector in 2d or 3d
- velocity: vector length



## Flow fields

- V(x, t): flow at location x at time t vector in 2d or 3d
- velocity: vector length



V(x, t) at a fixed time t What does variation in time mean?

# Working with flow fields

• forward prediction
# Working with flow fields



# Working with flow fields

 forward prediction simulate (propagate) distributions, include variation in time



# Working with flow fields

- forward prediction simulate (propagate) distributions, include variation in time
- hindcasting



8 March 2014: Malaysian Airlines Flight 370 (MH370) disappeared on its way from Kuala Lumpur to Beijing.



source: Wikipedia

#### Airplane pieces found in multiple locations.



image source: BBC





full report:

http://www.geomar.de/fileadmin/content/service/presse/Pressemitteilungen/2016/MH370\_Report\_May2016.pdf

Stefanie Jegelka (and Caroline Uhler)



#### full report:

http://www.geomar.de/fileadmin/content/service/presse/Pressemitteilungen/2016/MH370\_Report\_May2016.pdf

Stefanie Jegelka (and Caroline Uhler)

- Environmental data
- Modeling flows
- Short-range spatial correlations
  - intuition
  - 2 variables
  - multiple variables

### Sensing and correlations in space



- Measure & model correlations in space?
- Estimate temperature / rainfall / gold in other locations?

#### Sensing and correlations in space



- Measure & model correlations in space?
- Estimate temperature / rainfall / gold in other locations?
- Intuition: correlation is a function of distance

Stefanie Jegelka (and Caroline Uhler)

#### Where we are headed





# Intuition: correlated (Gaussian) random variables



weaker correlation

Stefanie Jegelka (and Caroline Uhler)

# Intuition: correlated (Gaussian) random variables



#### strong correlation

Stefanie Jegelka (and Caroline Uhler)

#### Intuition: 50 Gaussian random variables



#### Intuition: 50 Gaussian random variables



#### Intuition: 50 Gaussian random variables



# Multivariate Gaussian distribution

$$y \sim \mathcal{N}(\mu, \Sigma) \qquad p(y) = \frac{1}{(2\pi)^{-d/2} |\Sigma|^{1/2}} \exp\left(-\frac{1}{2}(y-\mu)^{\top} \Sigma^{-1}(y-\mu)\right)$$

$$\begin{array}{l} \mathcal{N}_{\mathcal{N}} \mathcal{N}(\mu, \mathcal{G}^{2}) \qquad p(y) \propto \exp(-\frac{(y-\mu)^{2}}{2\mathcal{G}^{2}}\right)$$

$$\begin{array}{l} \mathcal{N}_{\mathcal{N}} \mathcal{N}(\mu, \mathcal{G}^{2}) \qquad p(y) \propto \exp(-\frac{(y-\mu)^{2}}{2\mathcal{G}^{2}}\right)$$

$$\begin{array}{l} \mathcal{N}_{\mathcal{N}} \mathcal{N}(\mu, \mathcal{G}^{2}) \qquad p(y) \propto \exp(-\frac{(y-\mu)^{2}}{2\mathcal{G}^{2}}\right)$$

$$\begin{array}{l} \mathcal{N}_{\mathcal{N}} \mathcal{N}(\mu, \mathcal{G}^{2}) \qquad p(y) \approx \exp(-\frac{(y-\mu)^{2}}{2\mathcal{G}^{2}}\right)$$

$$\begin{array}{l} \mathcal{N}_{\mathcal{N}} \mathcal{N}(\mu, \mathcal{G}^{2}) \qquad p(y) \approx \exp(-\frac{(y-\mu)^{2}}{2\mathcal{G}^{2}}\right)$$

 $\Sigma_{ij} = \operatorname{cov}(y_i, y_j)$ 

## Multivariate Gaussian distribution

$$y \sim \mathcal{N}(\mu, \Sigma)$$
  $p(y) = rac{1}{(2\pi)^{-d/2} |\Sigma|^{1/2}} \exp\left(-rac{1}{2}(y-\mu)^{\top} \Sigma^{-1}(y-\mu)
ight)$ 



 $\Sigma_{ij} = \operatorname{cov}(y_i, y_j)$ 















- Environmental data
- Modeling flows
- Short-range spatial correlations
  - intuition
  - 2 variables
  - multiple variables

•  $Y_A$ ,  $Y_B$  Gaussian random variables. We observe  $Y_B = y_B$ .

•  $Y_A$ ,  $Y_B$  Gaussian random variables. We observe  $Y_B = y_B$ .

$$\begin{bmatrix} Y_A \\ Y_B \end{bmatrix} \sim \mathcal{N}\left( \begin{bmatrix} \mu_A \\ \mu_B \end{bmatrix}, \begin{bmatrix} \sigma_A^2 & \sigma_{AB} \\ \sigma_{AB} & \sigma_B^2 \end{bmatrix} \right) \quad \longleftarrow$$

• Conditioning:  $p(Y_A | Y_B = y_B)$  is also Gaussian with mean and variance  $p(Y_A | Y_B) = \frac{P(Y_A | Y_B)}{P(Y_B)} = \frac{(-) e \times p(-(-) z^{-1}(-))}{P(Y_B)}$ 

•  $Y_A$ ,  $Y_B$  Gaussian random variables. We observe  $Y_B = y_B$ .

• Conditioning:  $p(Y_A|Y_B = y_B)$  is also Gaussian with mean and variance

$$\mu_{A|B} = \mu_A + \frac{\sigma_{AB}}{\sigma_B^2} (\underline{y_B - \mu_B})$$

•  $Y_A$ ,  $Y_B$  Gaussian random variables. We observe  $Y_B = y_B$ .

$$\left[\begin{array}{c} Y_{A} \\ Y_{B} \end{array}\right] \sim \mathcal{N}\left(\left[\begin{array}{c} \mu_{A} \\ \mu_{B} \end{array}\right], \left[\begin{array}{c} \sigma_{A}^{2} & \sigma_{AB} \\ \sigma_{AB} & \sigma_{B}^{2} \end{array}\right]\right)$$

• Conditioning:  $p(Y_A|Y_B = y_B)$  is also Gaussian with mean and variance

$$\mu_{A|B} = \mu_A + \frac{\sigma_{AB}}{\sigma_B^2} (y_B - \mu_B)$$

$$\sigma_{A|B}^2 = \sigma_A^2 - \sigma_{AB} \sigma_B^{-2} \sigma_{AB}.$$

$$\mu_A \mu_{A|B} \mu_{A|B}$$





$$\sigma_{N*} = [\operatorname{cov}(Y_*, Y_1), \dots, \operatorname{cov}(Y_*, Y_N)]$$



$$\sigma_{N*} = [\operatorname{cov}(Y_*, Y_1), \dots, \operatorname{cov}(Y_*, Y_N)]$$
  
$$\sigma_*^2 = \operatorname{var}(Y_*)$$
  
$$p(Y_*) = \mathcal{N}(O_1 \in \mathcal{C}^2)$$

• Assume  $Y_1, \ldots, Y_N, Y_*$  are jointly Gaussian with  $\mu = 0$  and



$$\sigma_{N*} = [\operatorname{cov}(Y_*, Y_1), \dots, \operatorname{cov}(Y_*, Y_N)]$$
  
$$\sigma_*^2 = \operatorname{var}(Y_*)$$

Observe y<sub>1</sub>,..., y<sub>N</sub>. Conditioning yields a Gaussian p(Y<sub>\*</sub>|y<sub>1</sub>,...y<sub>N</sub>) with parameters

• Assume  $Y_1, \ldots, Y_N, Y_*$  are jointly Gaussian with  $\mu = 0$  and



$$\sigma_{N*} = [\operatorname{cov}(Y_*, Y_1), \dots, \operatorname{cov}(Y_*, Y_N)]$$
  
$$\sigma_*^2 = \operatorname{var}(Y_*)$$

• Observe  $y_1, \ldots, y_N$ . Conditioning yields a Gaussian  $p(Y_*|y_1, \ldots, y_N)$  with parameters

$$\mu_{A|B} = \mu_A + \frac{\sigma_{AB}}{\sigma_B^2} (y_B - \mu_B)$$

• Assume  $Y_1, \ldots, Y_N, Y_*$  are jointly Gaussian with  $\mu = 0$  and



$$\sigma_{N*} = [\operatorname{cov}(Y_*, Y_1), \dots, \operatorname{cov}(Y_*, Y_N)]$$
  
$$\sigma_*^2 = \operatorname{var}(Y_*)$$

• Observe  $y_1, \ldots, y_N$ . Conditioning yields a Gaussian  $p(Y_*|y_1, \ldots, y_N)$  with parameters

$$\mu_{A|B} = \mu_A + \frac{\sigma_{AB}}{\sigma_B^2} (y_B - \mu_B) \implies \mu_{*|1:N} = \mu_* + \sigma_{N*}^\top \Sigma_N^{-1} (y_{1:N} - \mu_{1:N})$$
## Conditioning on multiple variables / partial observation

• Assume  $Y_1, \ldots, Y_N, Y_*$  are jointly Gaussian with  $\mu = 0$  and



$$\sigma_{N*} = [\operatorname{cov}(Y_*, Y_1), \dots, \operatorname{cov}(Y_*, Y_N)]$$
  
$$\sigma_*^2 = \operatorname{var}(Y_*)$$

• Observe  $y_1, \ldots, y_N$ . Conditioning yields a Gaussian  $p(Y_*|y_1, \ldots, y_N)$  with parameters

$$\mu_{A|B} = \mu_A + \frac{\sigma_{AB}}{\sigma_B^2} (y_B - \mu_B) \implies \mu_{*|1:N} = \mu_* + \sigma_{N*}^\top \Sigma_N^{-1} (y_{1:N} - \mu_{1:N})$$
  
$$\sigma_{A|B}^2 = \sigma_A^2 - \sigma_{AB} \sigma_B^{-2} \sigma_{AB}$$

## Conditioning on multiple variables / partial observation

• Assume  $Y_1, \ldots, Y_N, Y_*$  are jointly Gaussian with  $\mu = 0$  and



$$\sigma_{N*} = [\operatorname{cov}(Y_*, Y_1), \dots, \operatorname{cov}(Y_*, Y_N)]$$
  
$$\sigma_*^2 = \operatorname{var}(Y_*)$$

• Observe  $y_1, \ldots, y_N$ . Conditioning yields a Gaussian  $p(Y_*|y_1, \ldots, y_N)$  with parameters

$$\mu_{A|B} = \mu_A + \frac{\sigma_{AB}}{\sigma_B^2} (y_B - \mu_B) \Rightarrow \mu_{*|1:N} = \mu_* + \sigma_{N*}^\top \Sigma_N^{-1} (y_{1:N} - \mu_{1:N})$$
  
$$\sigma_{A|B}^2 = \sigma_A^2 - \sigma_{AB} \sigma_B^{-2} \sigma_{AB} \Rightarrow \sigma_{*|1:N}^2 = \sigma_*^2 - \sigma_{N*}^\top \Sigma_N^{-1} \sigma_{N*}.$$



• prediction: mean/mode  $\hat{y}_* = \mu_{*|1:N} = \mu_* + \sigma_{N*}^\top \Sigma_N^{-1} y_{1:N}$ (for  $\mu = 0$ )

- prediction: mean/mode  $\hat{y}_* = \mu_{*|1:N} = \mu_* + \sigma_{N*}^\top \Sigma_N^{-1} y_{1:N}$ (for  $\mu = 0$ )
- $\bullet$  variance shrinks as we observe more data:  $\sigma_{*|1:\textit{N}}^2 \leq \sigma_*^2$

- prediction: mean/mode  $\hat{y}_* = \mu_{*|1:N} = \mu_* + \sigma_{N*}^\top \Sigma_N^{-1} y_{1:N}$ (for  $\mu = 0$ )
- variance shrinks as we observe more data:  $\sigma^2_{*|1:N} \leq \sigma^2_{*}$
- example of Bayesian Inference.



- prediction: mean/mode  $\hat{y}_* = \mu_{*|1:N} = \mu_* + \sigma_{N*}^\top \Sigma_N^{-1} y_{1:N}$ (for  $\mu = 0$ )
- variance shrinks as we observe more data:  $\sigma^2_{*|1:N} \leq \sigma^2_{*}$
- example of Bayesian Inference.
- Just need the *covariance* ...

- prediction: mean/mode  $\hat{y}_* = \mu_{*|1:N} = \mu_* + \sigma_{N*}^\top \Sigma_N^{-1} y_{1:N}$ (for  $\mu = 0$ )
- variance shrinks as we observe more data:  $\sigma_{*|1:N}^2 \leq \sigma_*^2$
- example of Bayesian Inference.
- Just need the *covariance* ...
- Idea: covariance is a function of distance!

- Specifics of environmental data: spatio-temporal dependencies, physical models, simulations
- Modeling flows via discretization
- Modeling short-range spatial correlations with Gaussians

## Environmental Data, Climate Informatics, Comp. Sustainability:

- H. Kaper, H. Engler. *Mathematics & Climate*. Chapters 1, 8.
- W. Menke, J. Menke. Environmental Data Analysis with Matlab
- Tackling Climate Change with Machine Learning. https://arxiv.org/abs/1906.05433 list of high impact problems
- Computational Sustainability: Computing for a Better World and a Sustainable Future. *Communications of the ACM*, Sep 2019. https://cacm.acm.org/ magazines/2019/9/238970-computational-sustainability/fulltext
- http://www.climateinformatics.org/