# Gaussian Processes for Regression: Models, Algorithms, and Applications

Tamara Broderick

Associate Professor
MIT

http://www.tamarabroderick.com/tutorials.html

# Why Gaussian processes (GPs)?
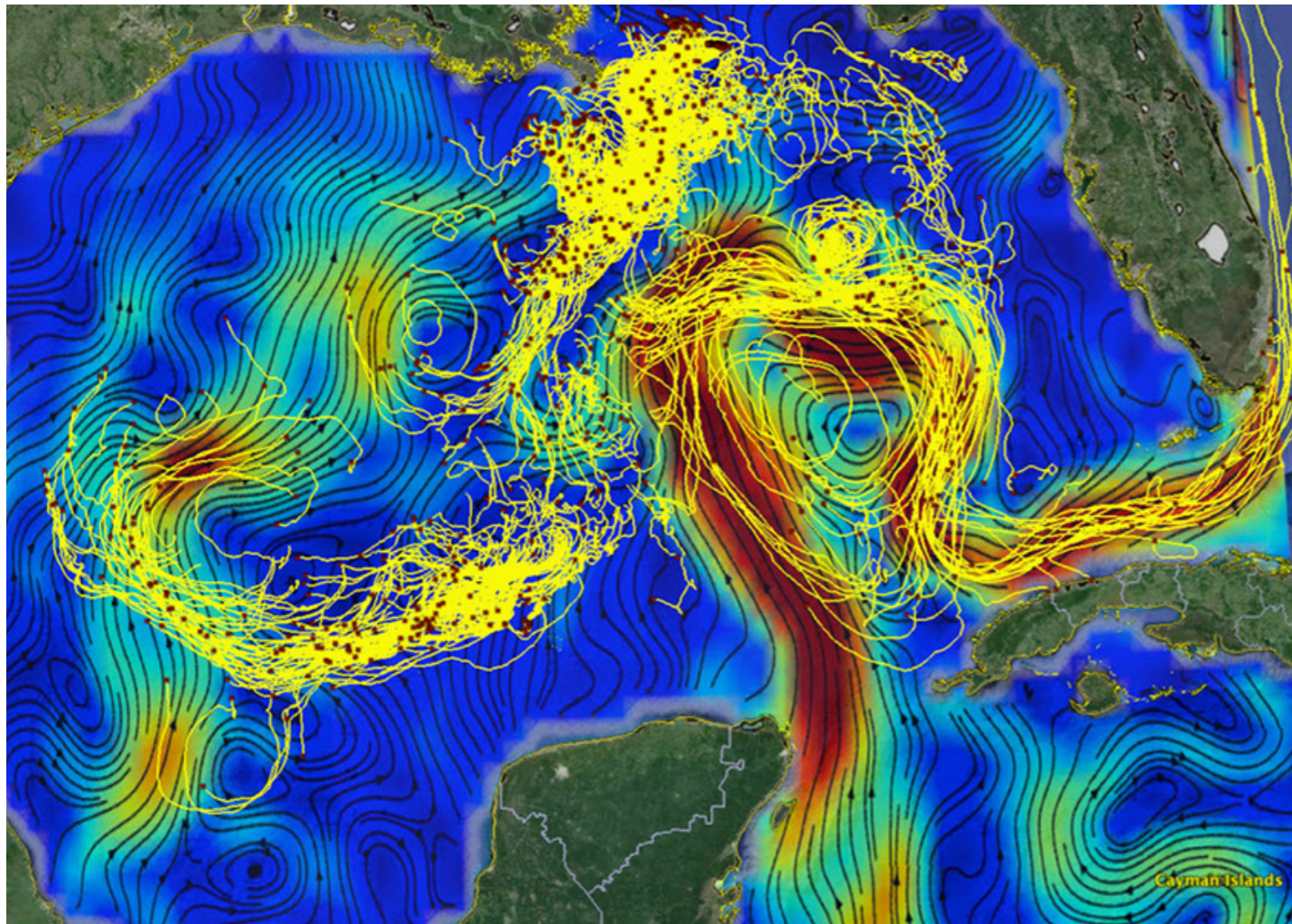
# Why Gaussian processes (GPs)?

- Often want to estimate/"predict" some continuous outcome as a function of certain inputs (*regression*)

# Why Gaussian processes (GPs)?

- Often want to estimate/"predict" some continuous outcome as a function of certain inputs (*regression*)
- GPs are good at certain types of regression problems

1

# Why Gaussian processes (GPs)?

- Often want to estimate/"predict" some continuous outcome as a function of certain inputs (*regression*)
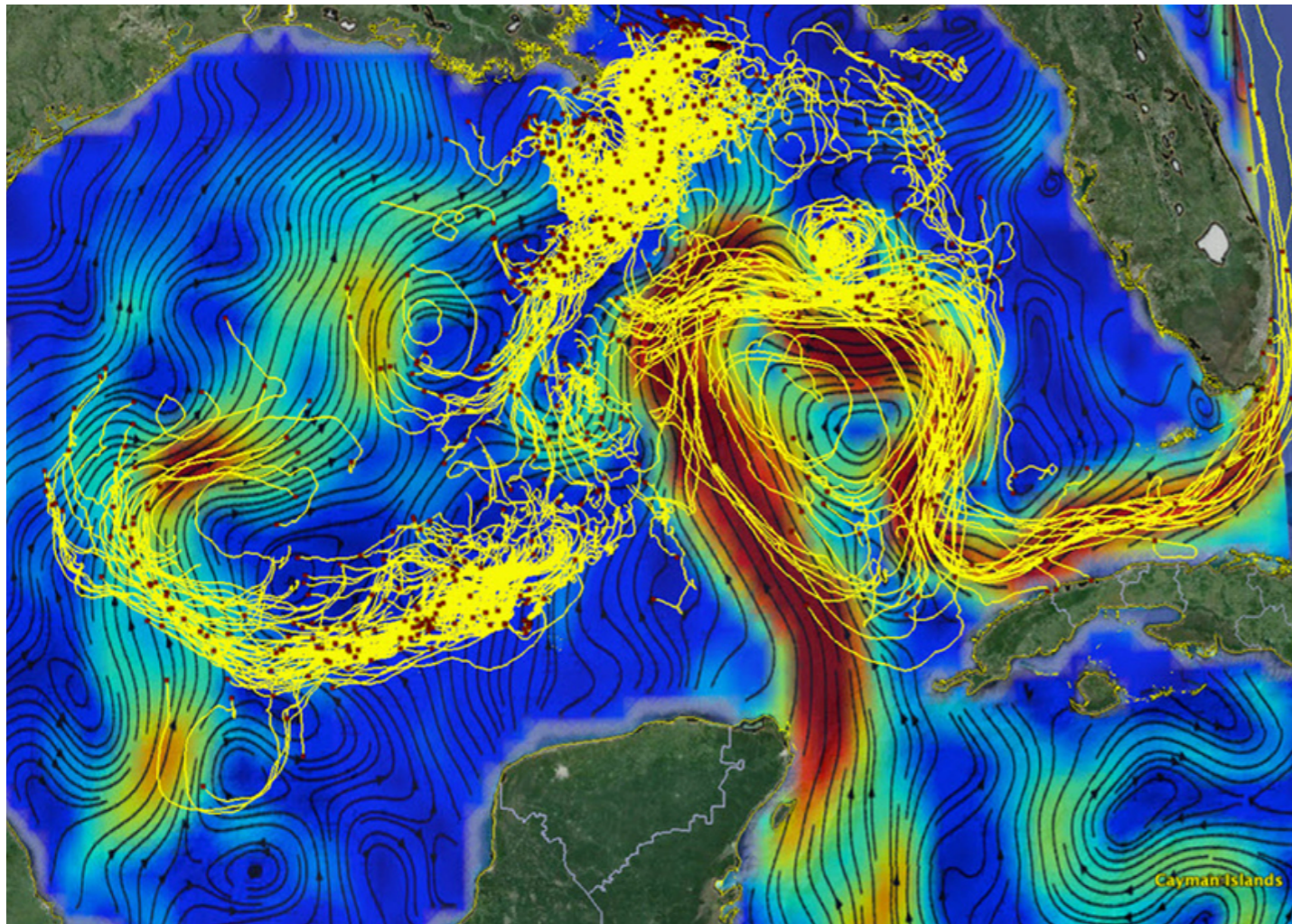- GPs are good at certain types of regression problems

Example:



[Ryan, Özgökmen 2023; Zewe 2023; Gonçalves et al 2019; Lodise et al 2020; Berlinghieri et al 2023]

# Why Gaussian processes (GPs)?

- Often want to estimate/"predict" some continuous outcome as a function of certain inputs (*regression*)
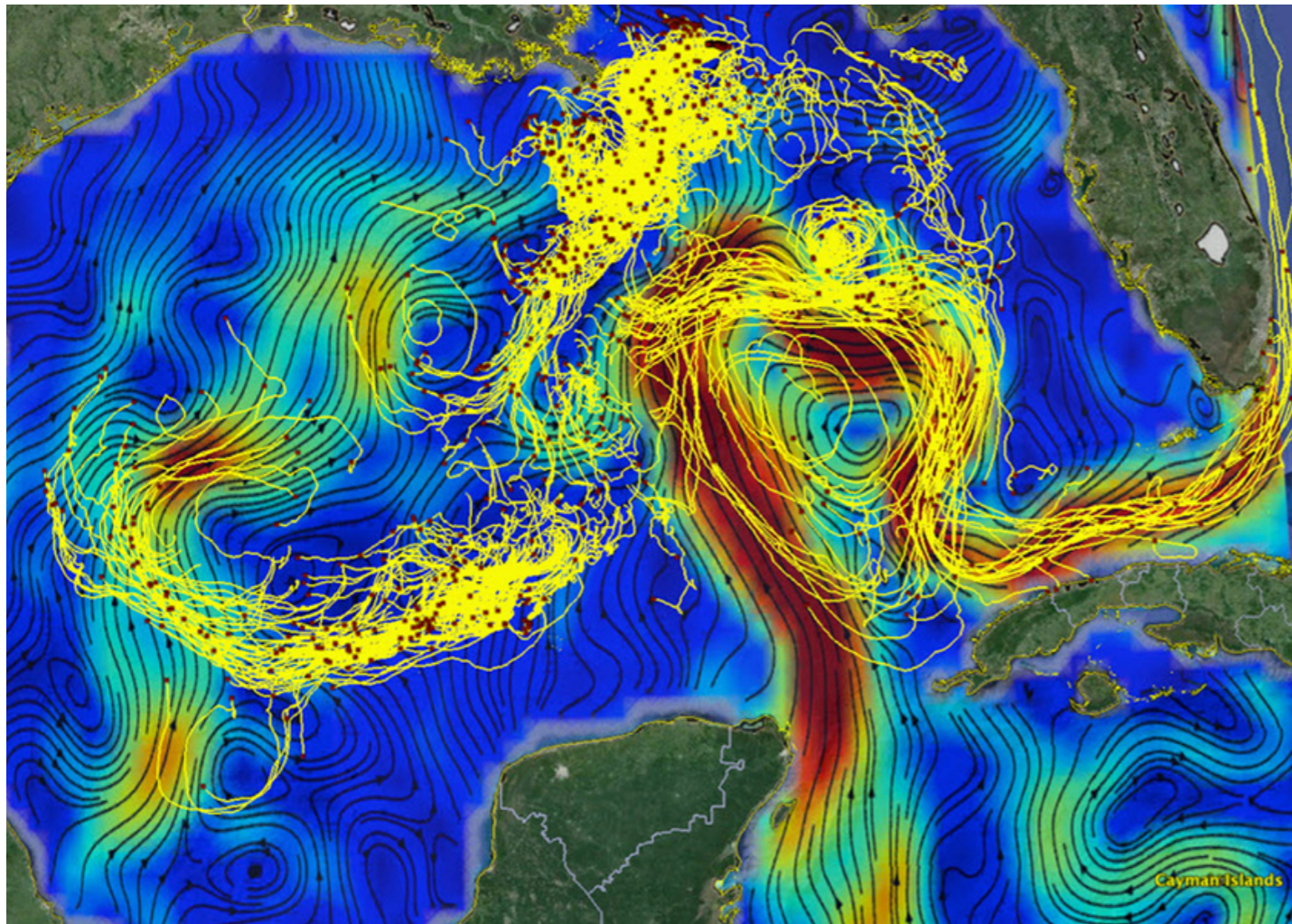- GPs are good at certain types of regression problems



Example:
- The ocean current (velocity vector field) varies by space & time

[Ryan, Özgökmen 2023; Zewe 2023; Gonçalves et al 2019; Lodise et al 2020; Berlinghieri et al 2023]

1

# Why Gaussian processes (GPs)?

- Often want to estimate/"predict" some continuous outcome as a function of certain inputs (*regression*)
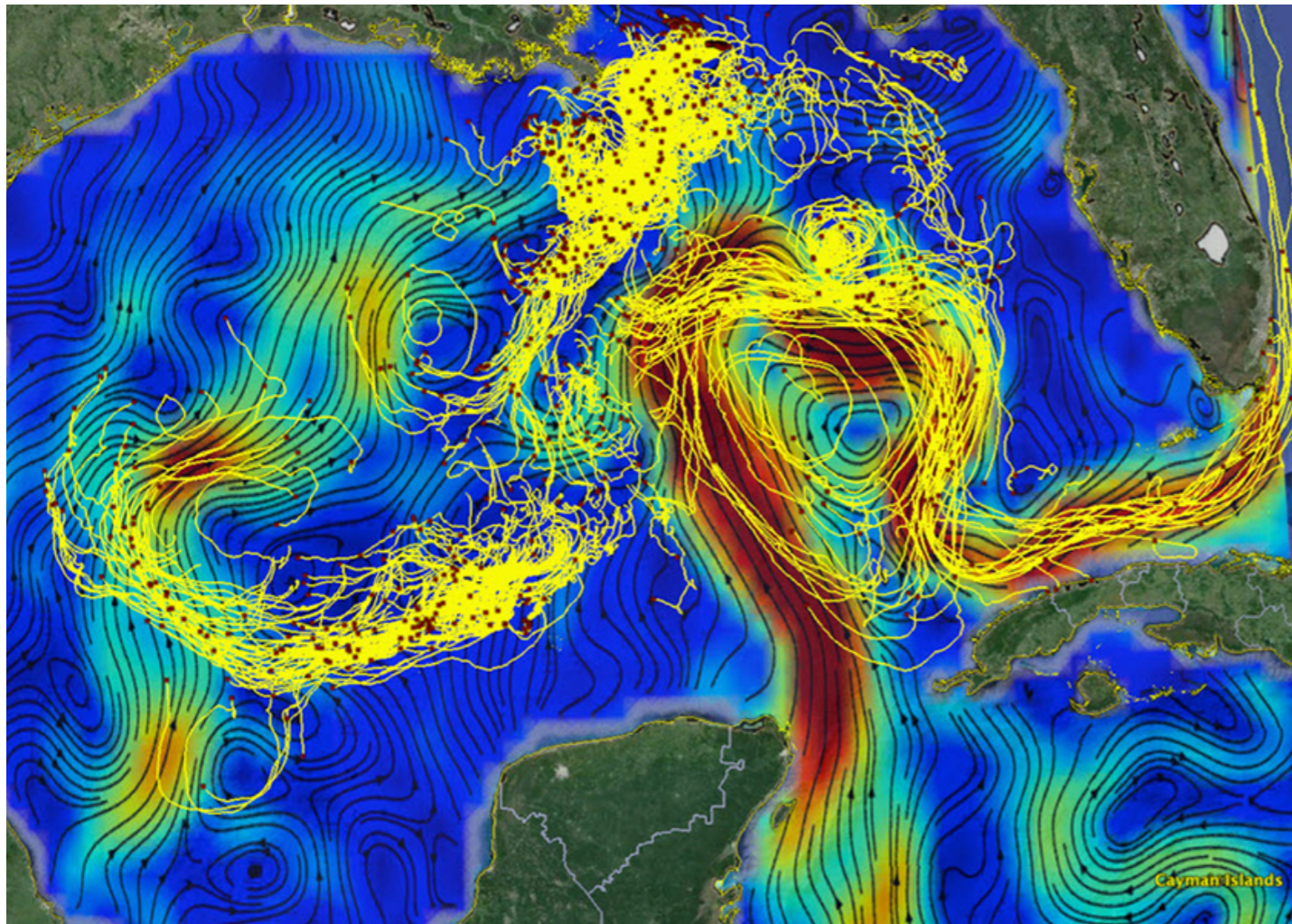- GPs are good at certain types of regression problems



Example:
- The ocean current (velocity vector field) varies by space & time
- Scientists get sparse observations of the current from buoys

[Ryan, Özgökmen 2023; Zewe 2023; Gonçalves et al 2019; Lodise et al 2020; Berlinghieri et al 2023]

1

# Why Gaussian processes (GPs)?

- Often want to estimate/"predict" some continuous outcome as a function of certain inputs (*regression*)
- GPs are good at certain types of regression problems



Example:
- The ocean current (velocity vector field) varies by space & time
- Scientists get sparse observations of the current from buoys
- Goal: estimate the current

[Ryan, Özgökmen 2023; Zewe 2023; Gonçalves et al 2019; Lodise et al 2020; Berlinghieri et al 2023]

1

# Why Gaussian processes (GPs)?

- Often want to estimate/"predict" some continuous outcome as a function of certain inputs (*regression*)
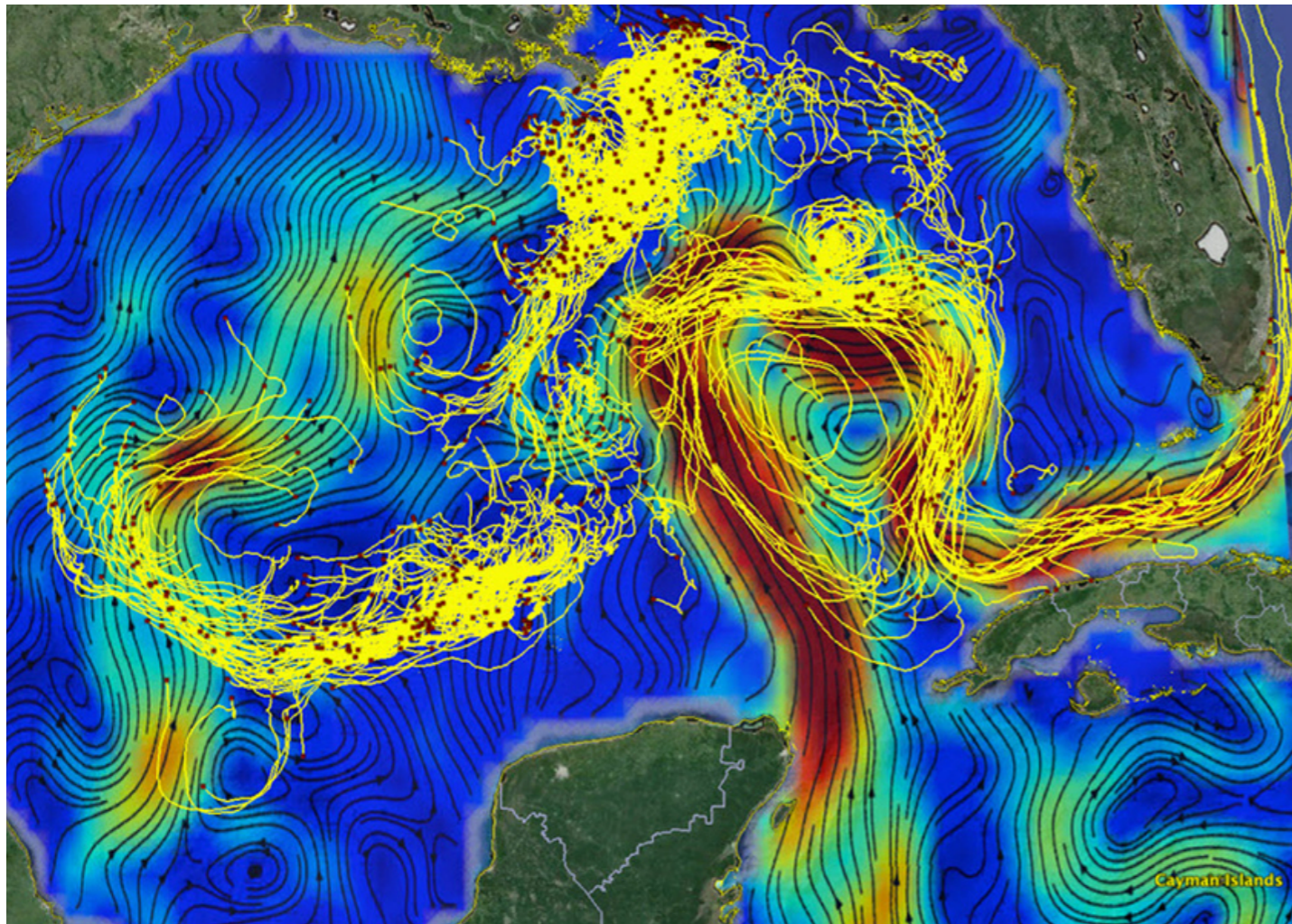- GPs are good at certain types of regression problems



Example:
- The ocean current (velocity vector field) varies by space & time
- Scientists get sparse observations of the current from buoys
- Goal: estimate the current "predict"

[Ryan, Özgökmen 2023; Zewe 2023; Gonçalves et al 2019; Lodise et al 2020; Berlinghieri et al 2023]

# Why Gaussian processes (GPs)?

- Often want to estimate/"predict" some continuous outcome as a function of certain inputs (*regression*)
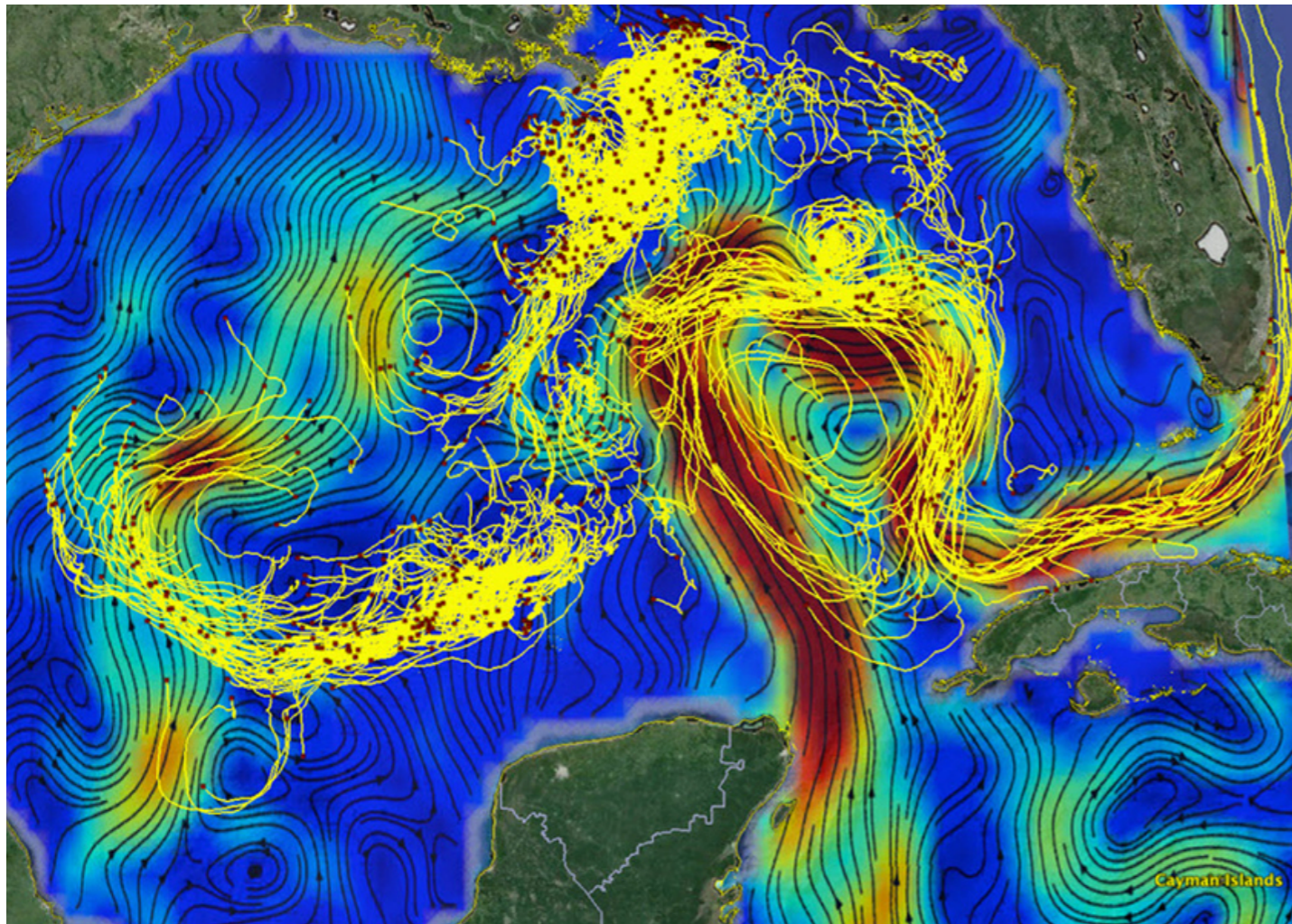- GPs are good at certain types of regression problems



Example:
- The ocean current (velocity vector field) varies by space & time
- Scientists get sparse observations of the current from buoys
- Goal: estimate the current "predict"

[Ryan, Özgökmen 2023; Zewe 2023; Gonçalves et al 2019; Lodise et al 2020; Berlinghieri et al 2023]

1

# Why Gaussian processes (GPs)?

- Often want to estimate/"predict" some continuous outcome as a function of certain inputs (*regression*)
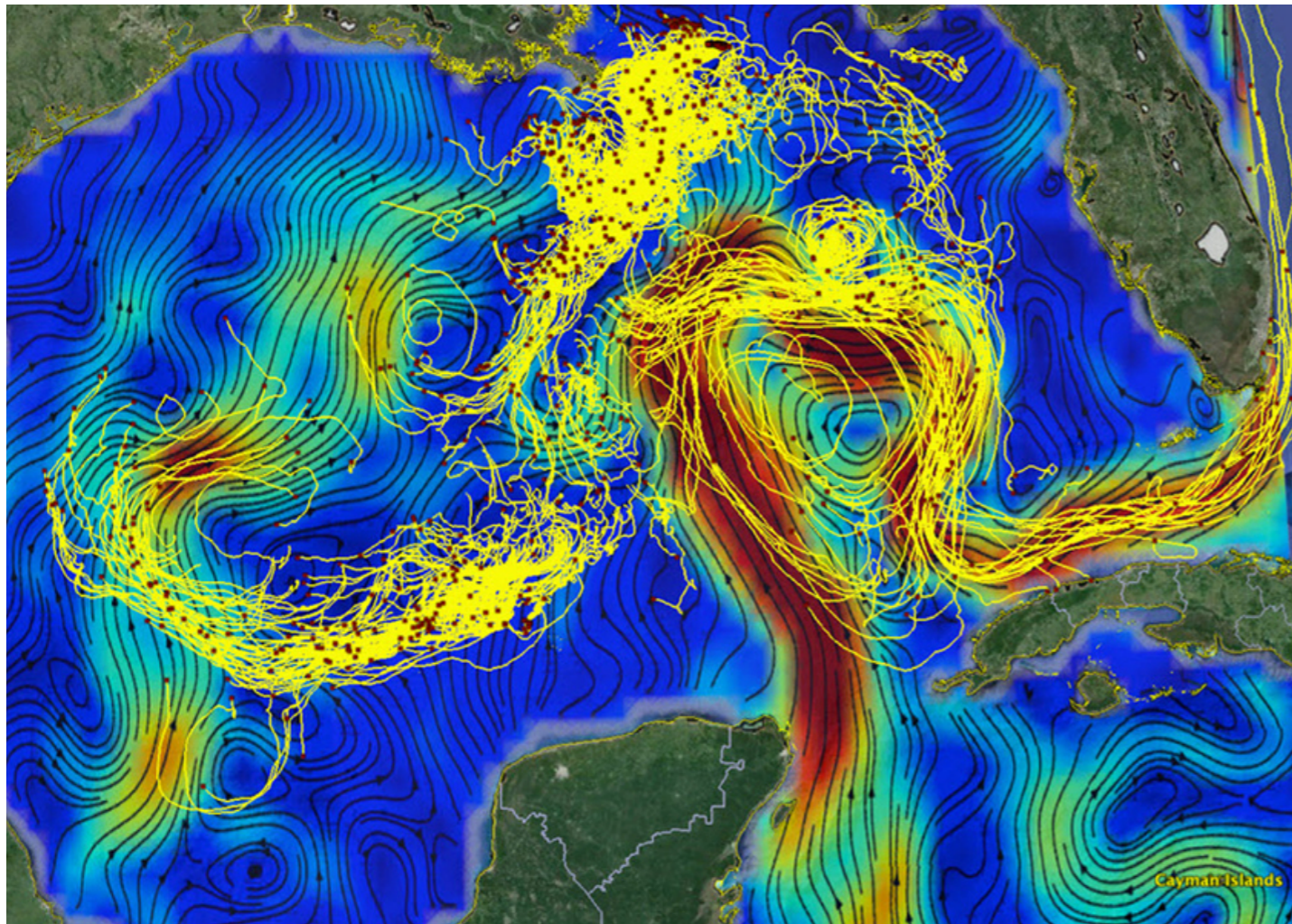- GPs are good at certain types of regression problems



Example:
- The ocean current (velocity vector field) varies by space & time
- Scientists get sparse observations of the current from buoys
- Goal: estimate the current   "predict"

[Ryan, Özgökmen 2023; Zewe 2023; Gonçalves et al 2019; Lodise et al 2020; Berlinghieri et al 2023]

Challenges:

# Why Gaussian processes (GPs)?

- Often want to estimate/"predict" some continuous outcome as a function of certain inputs (*regression*)
- GPs are good at certain types of regression problems
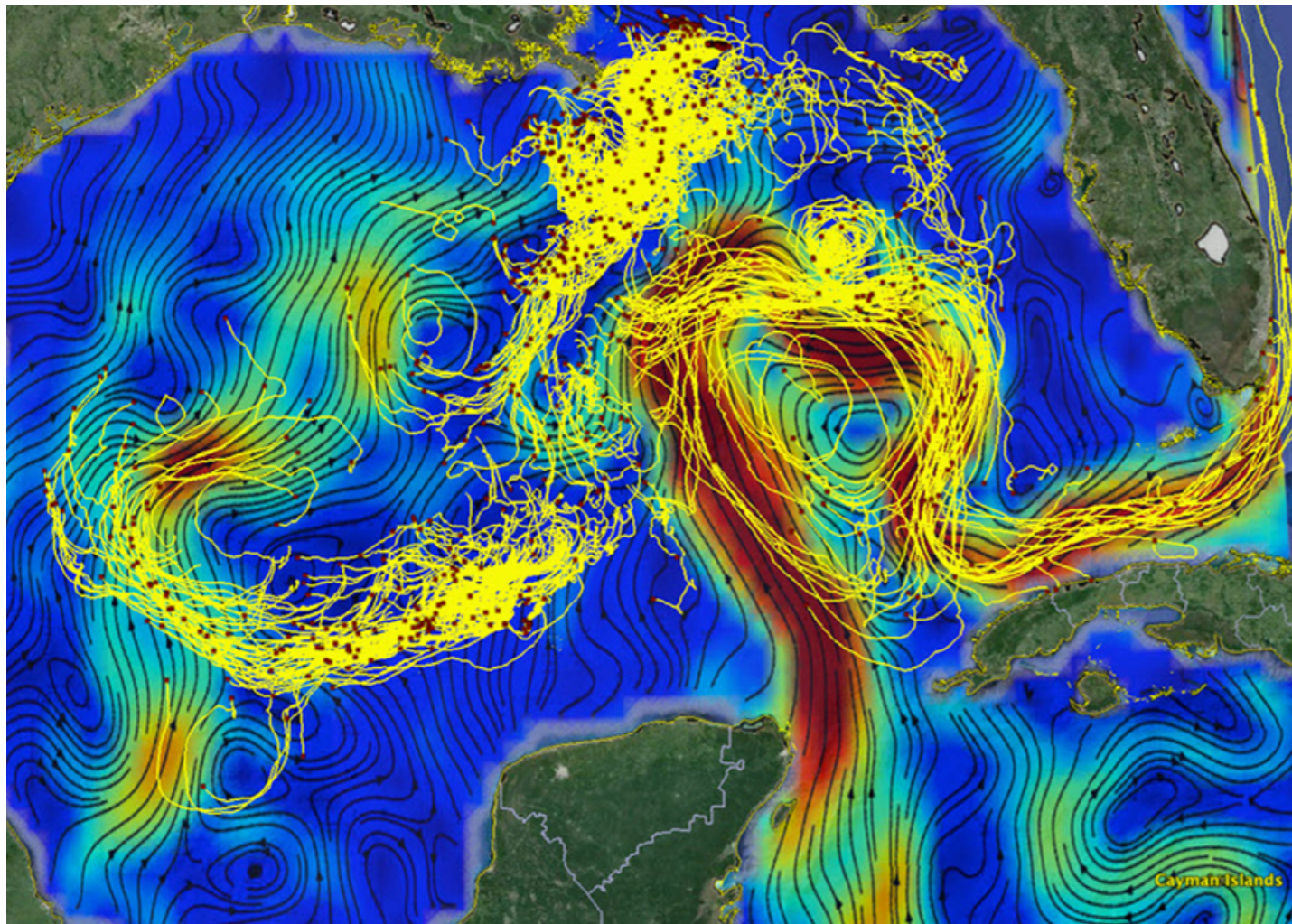


Example:

- The ocean current (velocity vector field) varies by space & time
- Scientists get sparse observations of the current from buoys
- Goal: estimate the current    "predict"

[Ryan, Özgökmen 2023; Zewe 2023; Gonçalves et al 2019; Lodise et al 2020; Berlinghieri et al 2023]

Challenges:  • Sparse & expensive data, not on a grid

# Why Gaussian processes (GPs)?

- Often want to estimate/"predict" some continuous outcome as a function of certain inputs (*regression*)
- GPs are good at certain types of regression problems
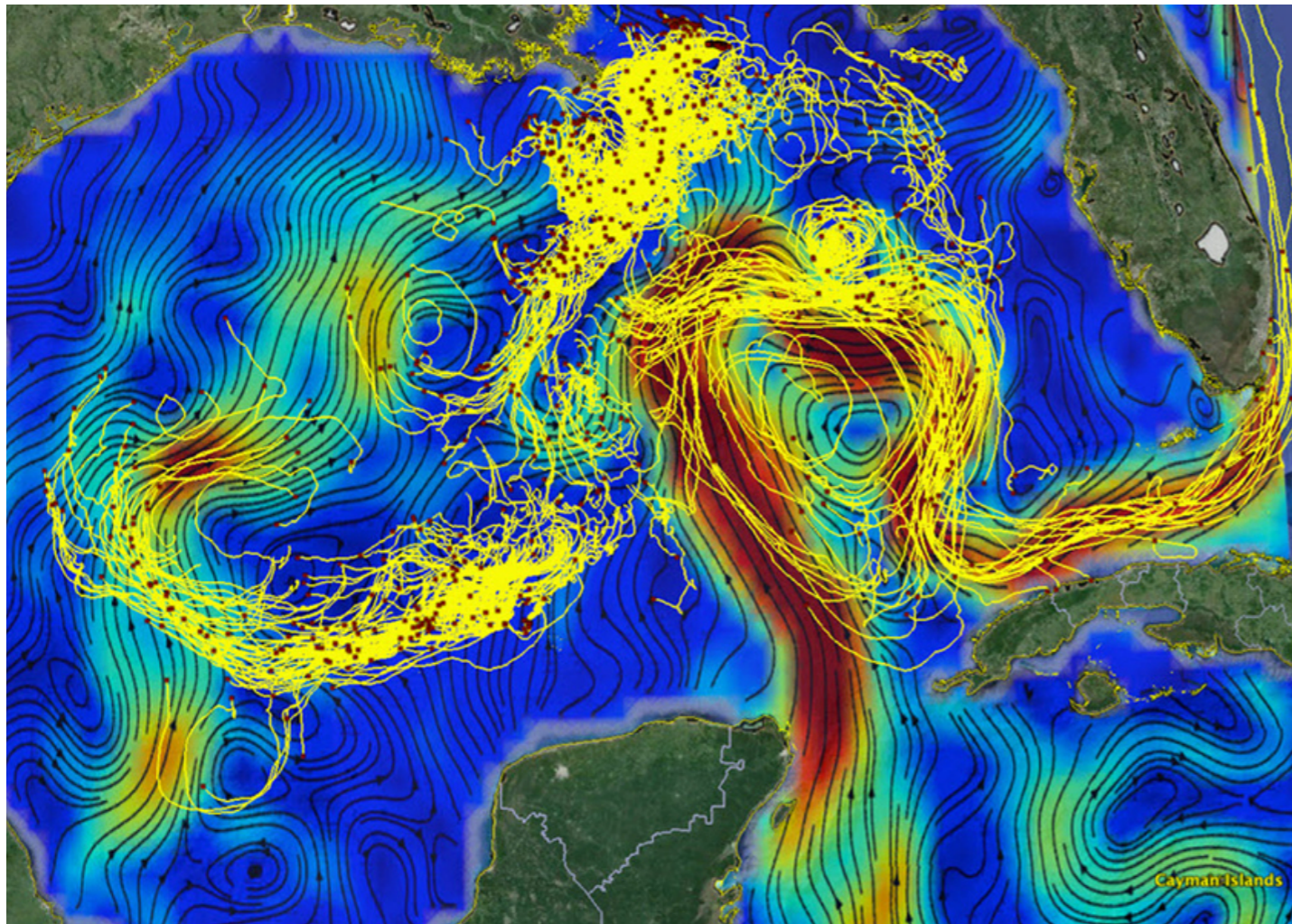


Example:
- The ocean current (velocity vector field) varies by space & time
- Scientists get sparse observations of the current from buoys
- Goal: estimate the current    "predict"

[Ryan, Özgökmen 2023; Zewe 2023; Gonçalves et al 2019; Lodise et al 2020; Berlinghieri et al 2023]

Challenges: • Sparse & expensive data, not on a grid
- Current is highly nonlinear but smooth in space-time

# Why Gaussian processes (GPs)?

- Often want to estimate/"predict" some continuous outcome as a function of certain inputs (*regression*)
- GPs are good at certain types of regression problems



Example:
- The ocean current (velocity vector field) varies by space & time
- Scientists get sparse observations of the current from buoys
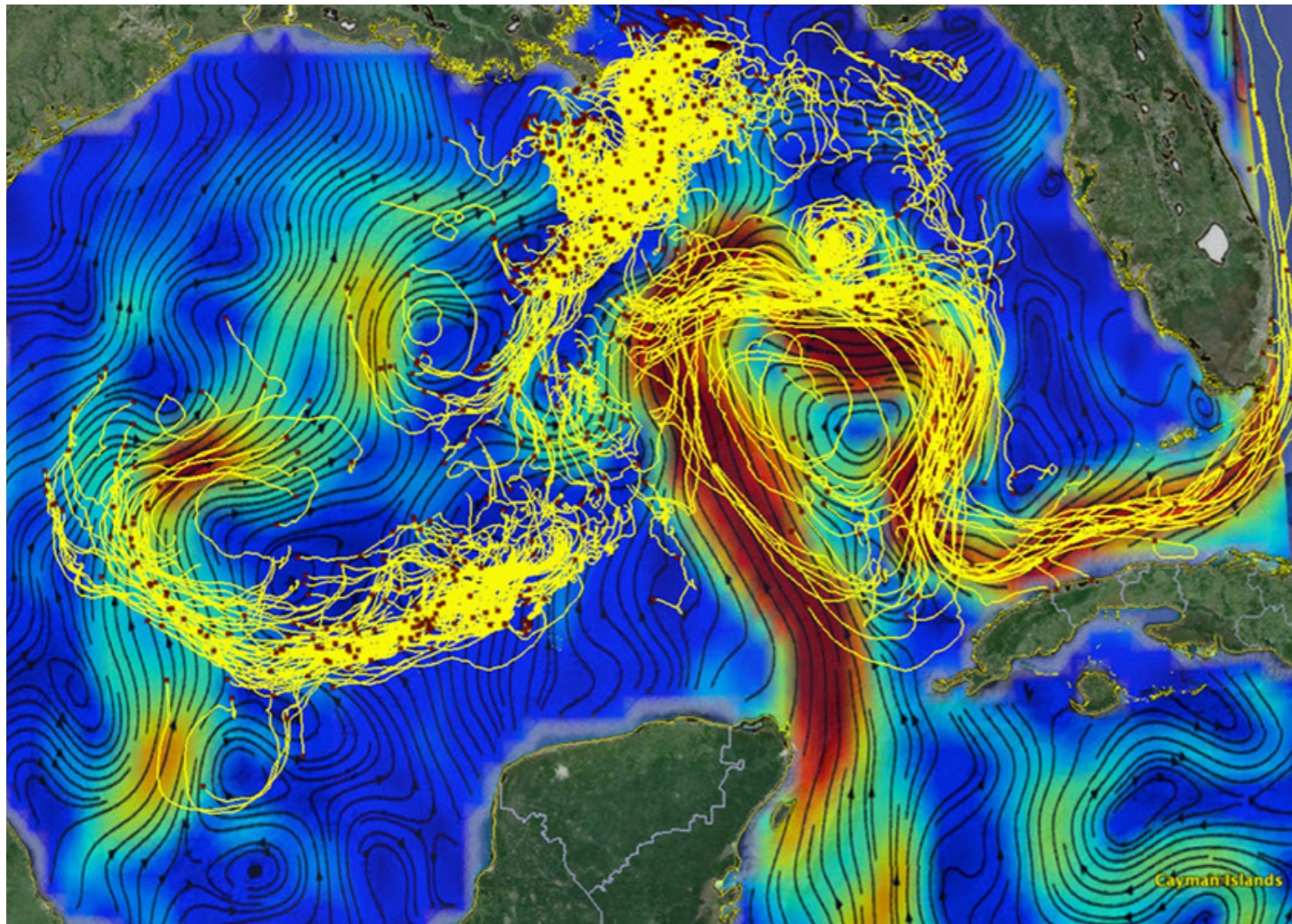- Goal: estimate the current "predict"

[Ryan, Özgökmen 2023; Zewe 2023; Gonçalves et al 2019; Lodise et al 2020; Berlinghieri et al 2023]

Challenges: • Sparse & expensive data, not on a grid
- Current is highly nonlinear but smooth in space-time
- Want uncertainty quantification

# Why Gaussian processes (GPs)?

# Why Gaussian processes (GPs)?
Example: "Surrogate model"



[Gramacy, Lee 2008] [Gramacy 2020]

# Why Gaussian processes (GPs)?

Example: "Surrogate model"
- The lift force of a rocket booster varies as a function of speed at re-entry, angle of attack, and sideslip angle



Lift

[Gramacy, Lee 2008]
[Gramacy 2020]

# Why Gaussian processes (GPs)?

Example: "Surrogate model"

- The lift force of a rocket booster varies as a function of speed at re-entry, angle of attack, and sideslip angle
- Scientists can run expensive simulations at chosen input settings



Lift

[Gramacy, Lee 2008] [Gramacy 2020]

# Why Gaussian processes (GPs)?

Example: "Surrogate model"
- The lift force of a rocket booster varies as a function of speed at re-entry, angle of attack, and sideslip angle
- Scientists can run expensive simulations at chosen input settings
- Goal: estimate how lift varies as a function of inputs



Lift

lift

1.0
0.5
0.0

alpha (angle of attack)
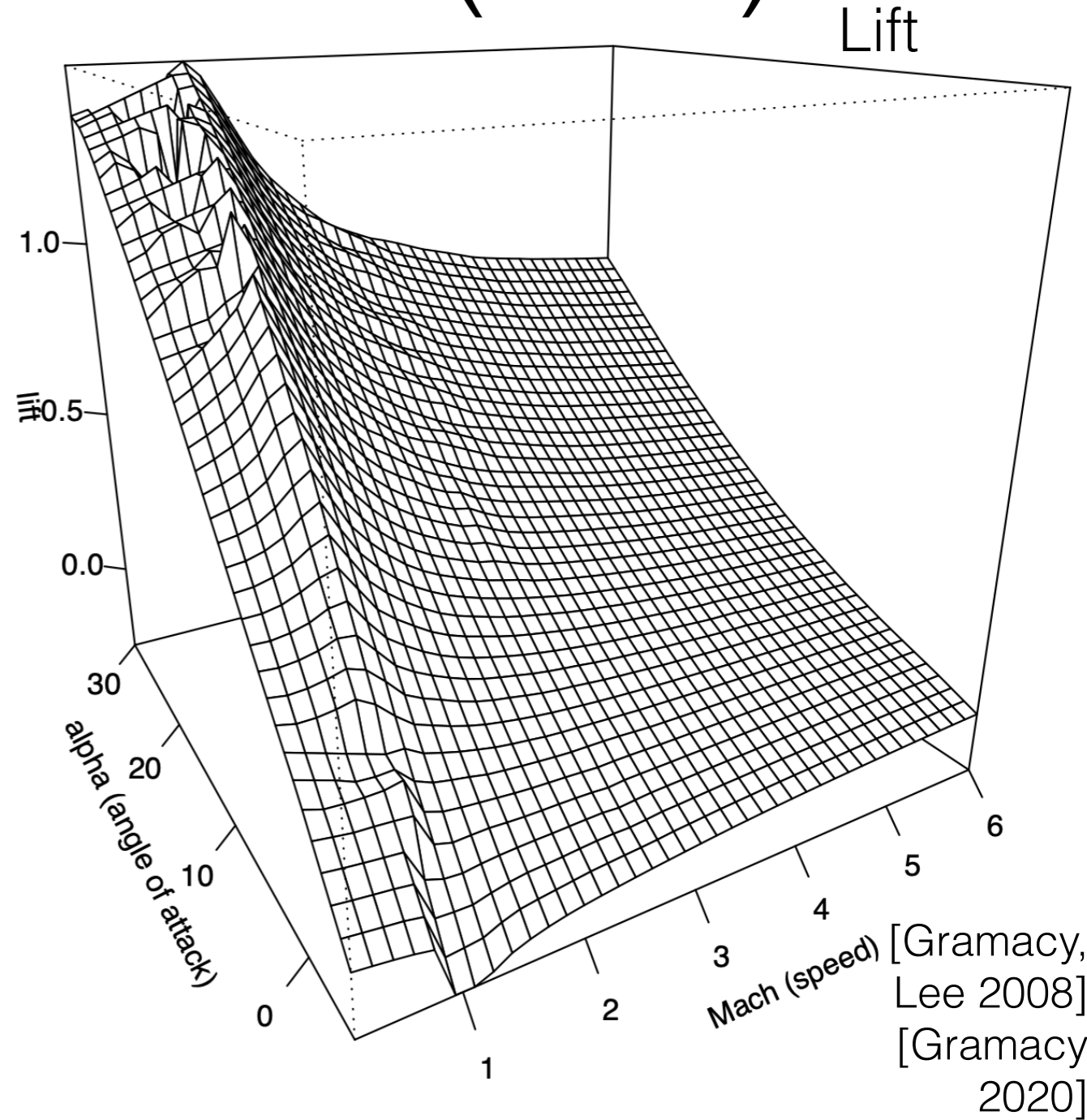30
20
10
0

Mach (speed)
1  2  3  4  5  6

[Gramacy, Lee 2008]
[Gramacy 2020]

2

# Why Gaussian processes (GPs)?

Example: "Surrogate model"

- The lift force of a rocket booster varies as a function of speed at re-entry, angle of attack, and sideslip angle
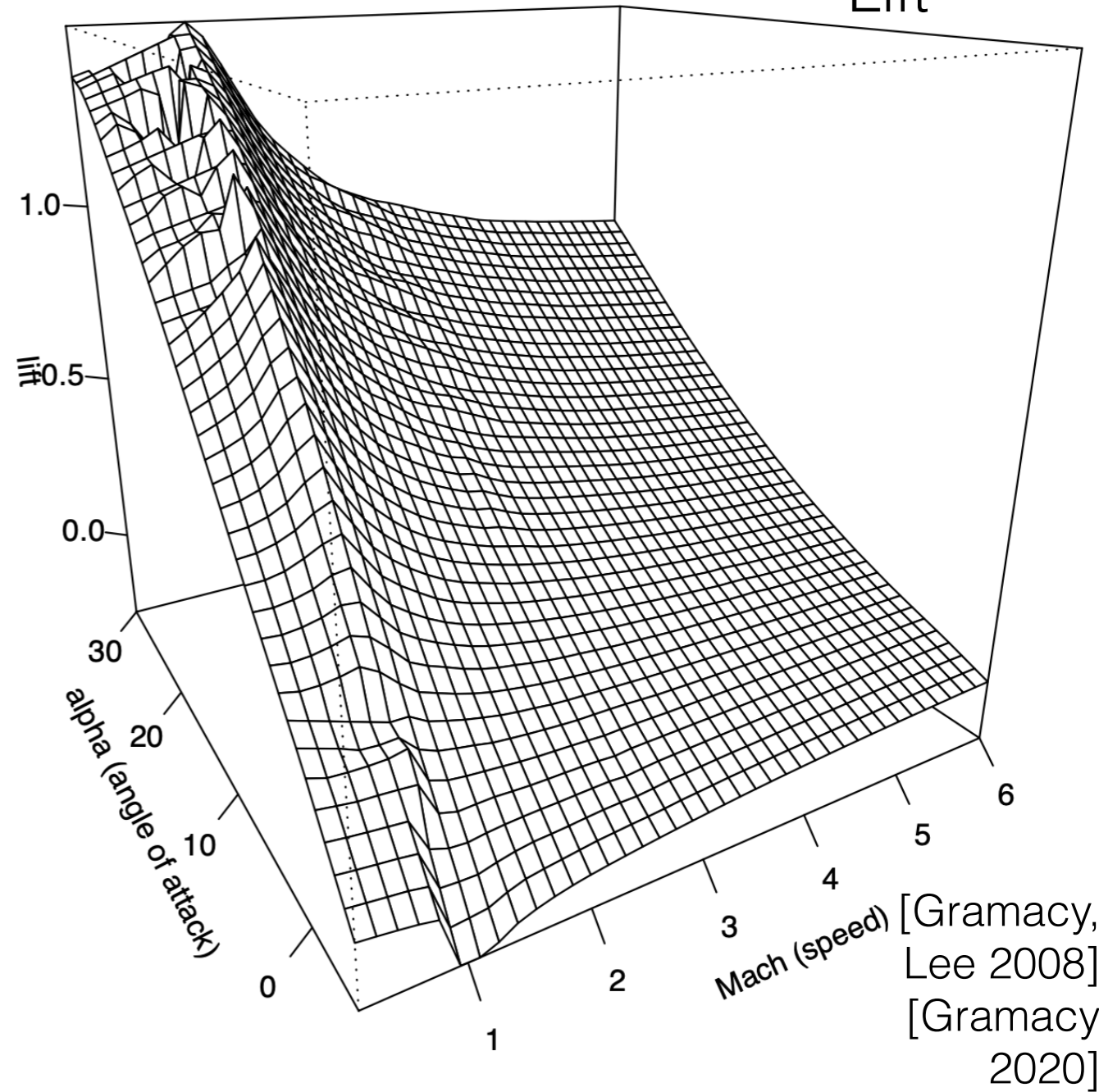- Scientists can run expensive simulations at chosen input settings
- Goal: estimate how lift varies as a function of inputs

Motif:



Lift

lift 1.0 0.5 0.0

alpha (angle of attack) 30 20 10 0

Mach (speed) 1 2 3 4 5 6

[Gramacy, Lee 2008]
[Gramacy 2020]

# Why Gaussian processes (GPs)?

Example: "Surrogate model"
- The lift force of a rocket booster varies as a function of speed at re-entry, angle of attack, and sideslip angle
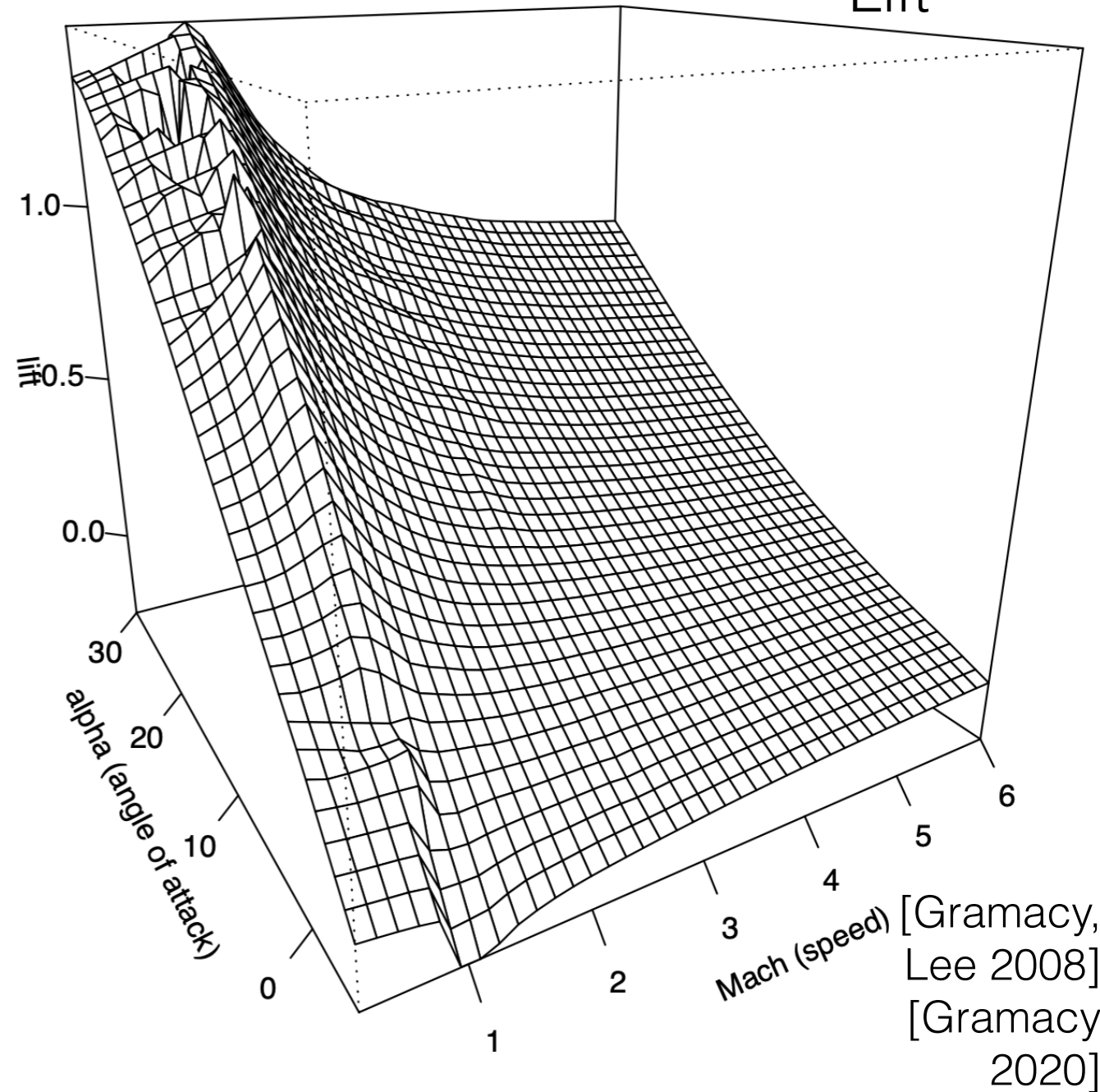- Scientists can run expensive simulations at chosen input settings
- Goal: estimate how lift varies as a function of inputs

Motif: • Sparse, costly data



[Gramacy, Lee 2008]
[Gramacy 2020]

# Why Gaussian processes (GPs)?

Example: "Surrogate model"

- The lift force of a rocket booster varies as a function of speed at re-entry, angle of attack, and sideslip angle
- Scientists can run expensive simulations at chosen input settings
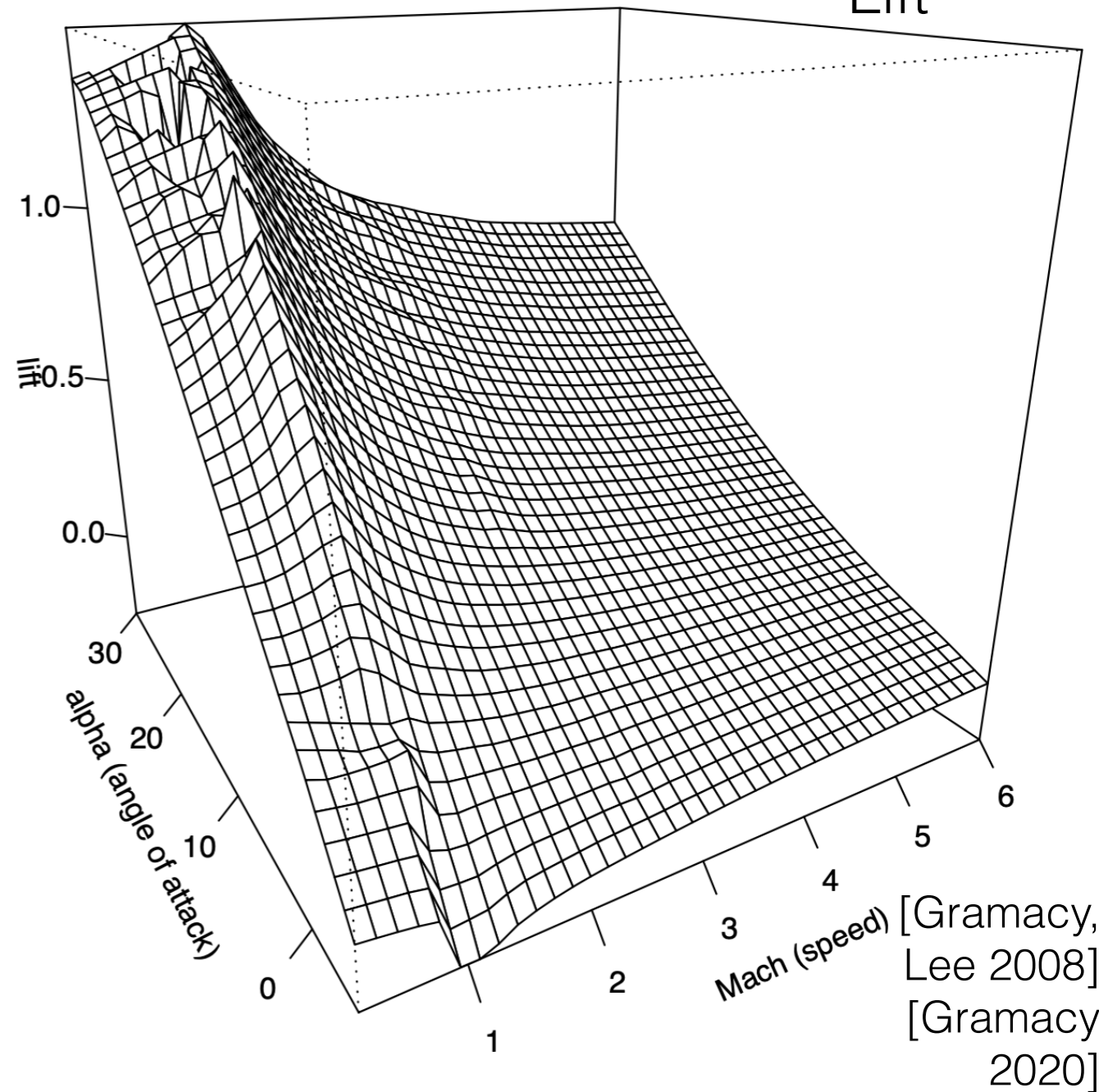- Goal: estimate how lift varies as a function of inputs

Motif:
- Sparse, costly data
- Output may have a nonlinear relationship to the inputs



Lift

[Gramacy, Lee 2008]
[Gramacy 2020]

2

# Why Gaussian processes (GPs)?

Example: "Surrogate model"
- The lift force of a rocket booster varies as a function of speed at re-entry, angle of attack, and sideslip angle
- Scientists can run expensive simulations at chosen input settings
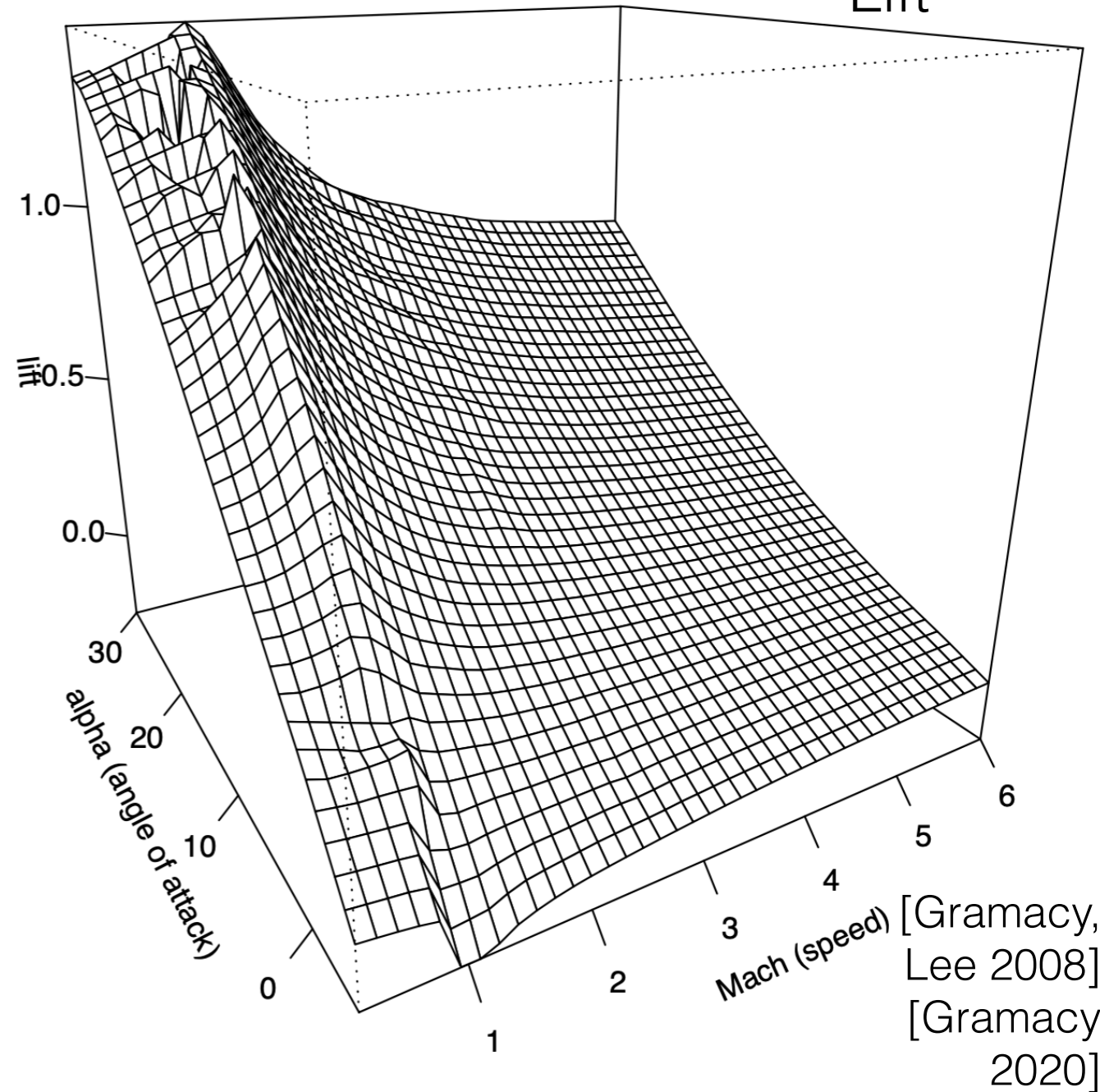- Goal: estimate how lift varies as a function of inputs

Motif:
- Sparse, costly data
- Output may have a nonlinear relationship to the inputs
- Want uncertainty quantification



[Gramacy, Lee 2008]
[Gramacy 2020]

# Why Gaussian processes (GPs)?

Example: "Surrogate model"

- The lift force of a rocket booster varies as a function of speed at re-entry, angle of attack, and sideslip angle
- Scientists can run expensive simulations at chosen input settings
- Goal: estimate how lift varies as a function of inputs
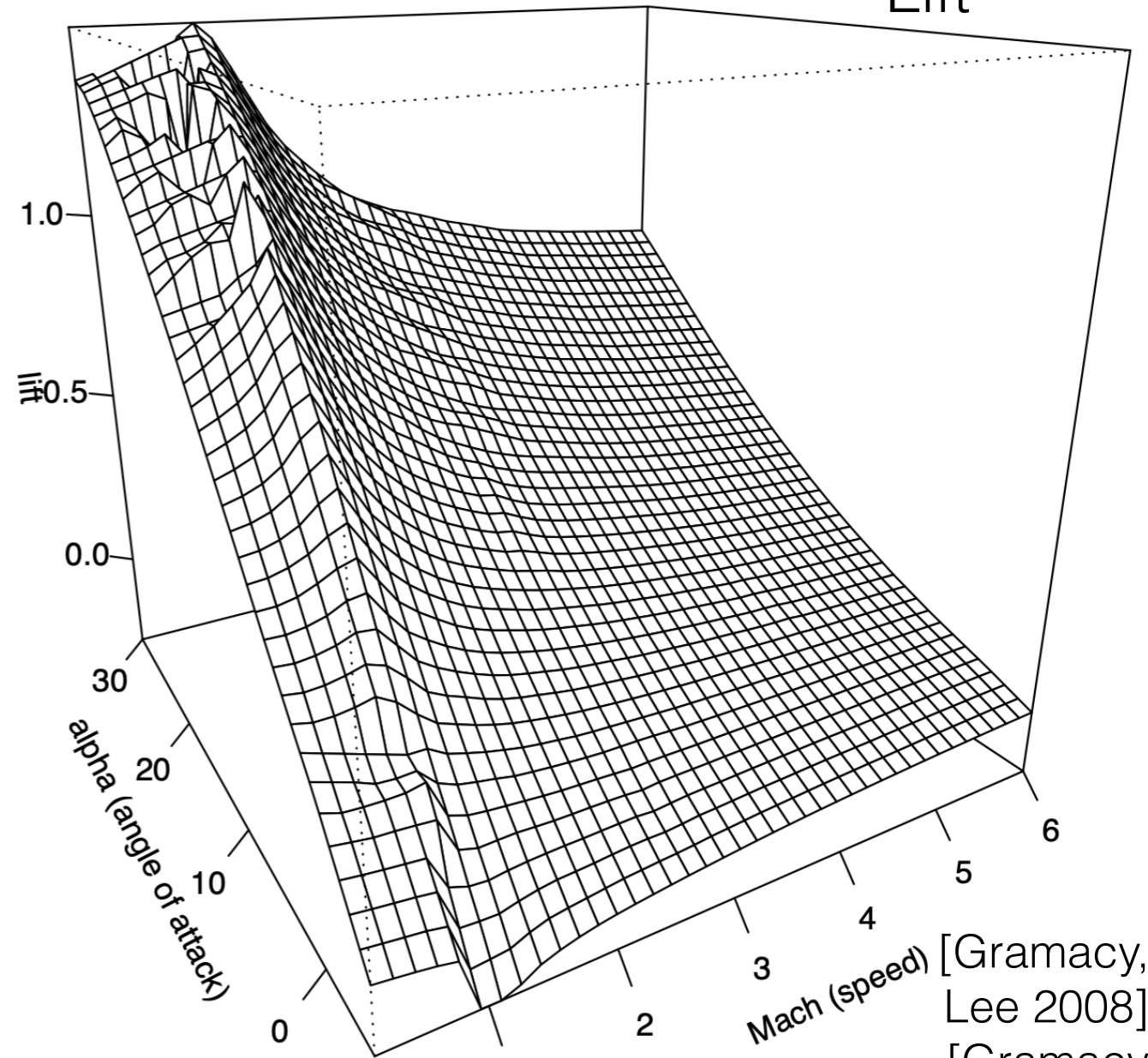
Motif:
- Sparse, costly data
- Output may have a nonlinear relationship to the inputs
- Want uncertainty quantification
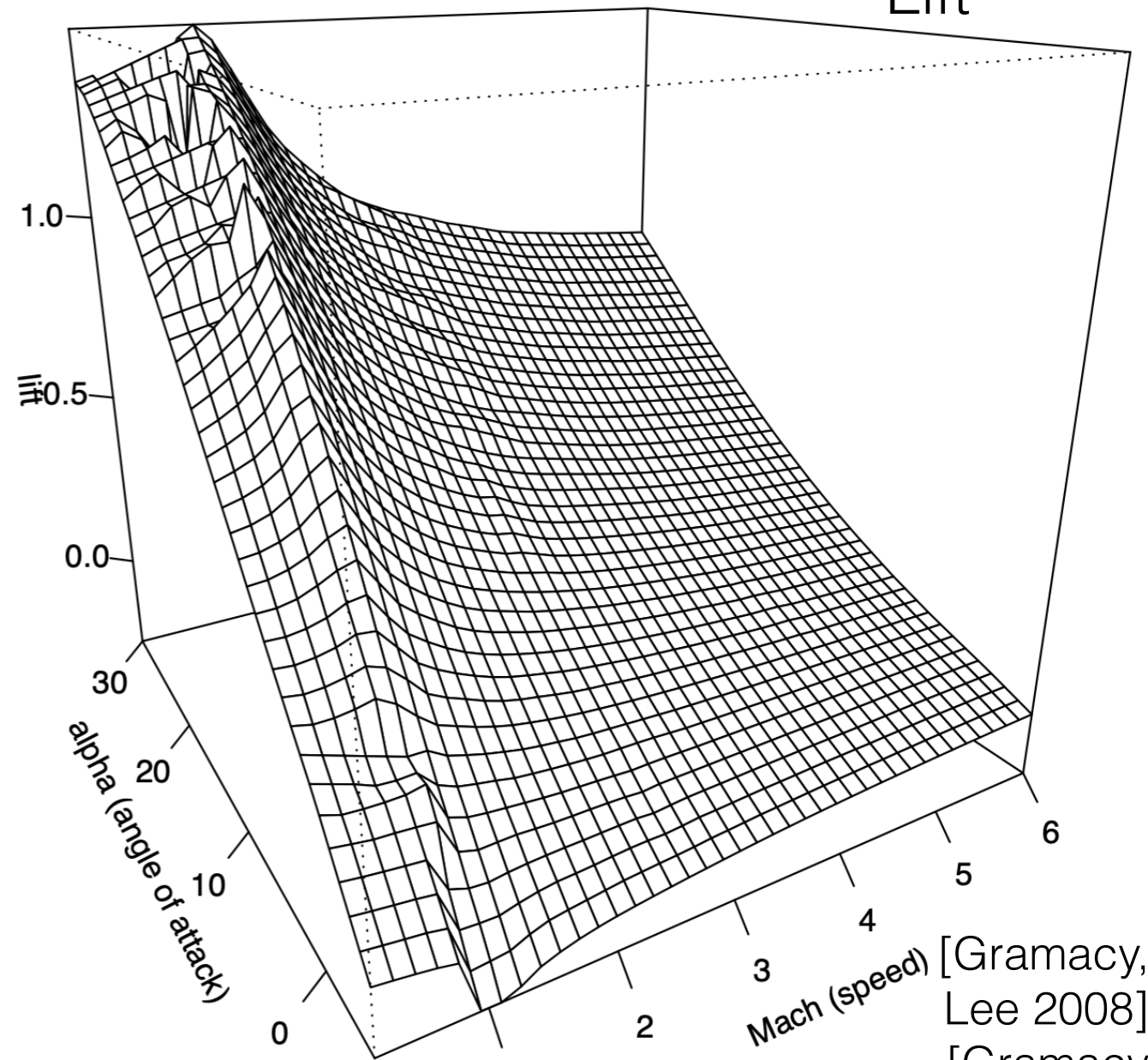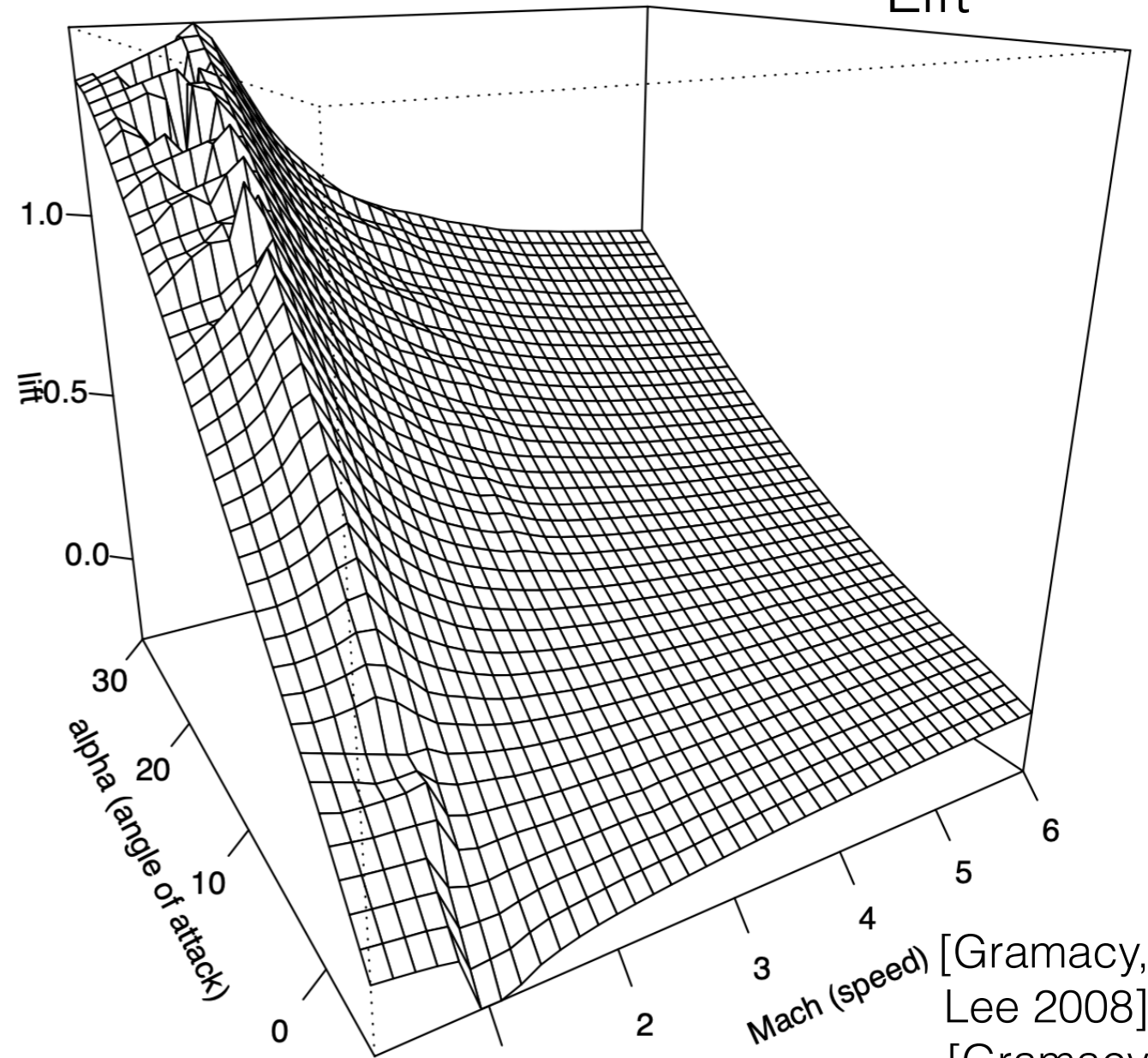- Low-dimensional inputs

Lift

[Gramacy, Lee 2008]
[Gramacy 2020]

# Why Gaussian processes (GPs)?

Example: "Surrogate model"

- The lift force of a rocket booster varies as a function of speed at re-entry, angle of attack, and sideslip angle
- Scientists can run expensive simulations at chosen input settings
- Goal: estimate how lift varies as a function of inputs

Motif:
- Sparse, costly data
  - Output may have a nonlinear relationship[1] to the inputs
  - Want uncertainty quantification
  - Low-dimensional inputs

One more example: learn (& optimize) performance in machine learning as a function of tuning parameters



Lift

[Gramacy, Lee 2008]
[Gramacy 2020]

[Snoek et al 2012, 2015; Garnett 2023]

# Roadmap

# Roadmap

- A Bayesian approach

# Roadmap

- A Bayesian approach

- What is a Gaussian process?

# Roadmap

- A Bayesian approach

- What is a Gaussian process?

  - Popular version using a squared exponential kernel

# Roadmap

- A Bayesian approach

- What is a Gaussian process?

  - Popular version using a squared exponential kernel

- Gaussian process inference

# Roadmap

- A Bayesian approach

- What is a Gaussian process?

  - Popular version using a squared exponential kernel

- Gaussian process inference

  - Prediction & uncertainty quantification

# Roadmap

- A Bayesian approach

- What is a Gaussian process?

  - Popular version using a squared exponential kernel

- Gaussian process inference

  - Prediction & uncertainty quantification


- Goal:
  - Learn the mechanism behind standard GPs to identify benefits and pitfalls

# A Bayesian approach

# A Bayesian approach

- p(unknowns | data)

# A Bayesian approach

- p(unknowns | data)

Given the data we've seen, what do we know about the underlying function?

# A Bayesian approach

- p(unknowns | data)

Given the data we've seen, what do we know about the underlying function?



Legend:
- x  data
- —  $f(x)$ best guess
- ▨  $f(x)$ 95% interval

# A Bayesian approach

- p(unknowns | data)

Given the data we've seen, what do we know about the underlying function?



Legend:
- x data
- $f(x)$ best guess
- $f(x)$ 95% interval

# A Bayesian approach

-

Given the data we've seen, what do we know about the underlying function?



Legend:
- **x** data
- —— $f(x)$ best guess
- $f(x)$ 95% interval

4

# A Bayesian approach

- p(unknowns | data)

Given the data we've seen, what do we know about the underlying function?



Legend:
- **x** data
- — $f(x)$ best guess
- $f(x)$ 95% interval

# A Bayesian approach

- p(unknowns | data) $\propto$ p(data | unknowns) p(unknowns)

Given the data we've seen, what do we know about the underlying function?



x — data
— $f(x)$ best guess
▓ $f(x)$ 95% interval

4

# A Bayesian approach

- p(unknowns | data) $\propto$ p(data | unknowns) p(unknowns)

Given the data we've seen, what do we know about the underlying function?

A (statistical) model that can generate functions and data of interest



Legend:
- **x** data
- — $f(x)$ best guess
- $f(x)$ 95% interval

# A Bayesian approach

- p(unknowns | data) $\propto$ p(data | unknowns) p(unknowns)

Given the data we've seen, what do we know about the underlying function?

A (statistical) model that can generate functions and data of interest



4

# A Bayesian approach

- p(unknowns | data) $\propto$ p(data | unknowns) p(unknowns)

Given the data we've seen, what do we know about the underlying function?

A (statistical) model that can generate functions and data of interest



4

# A Bayesian approach

- p(unknowns | data) $\propto$ p(data | unknowns) p(unknowns)

Given the data we've seen, what do we know about the underlying function?

A (statistical) model that can generate functions and data of interest



Legend:
- **x** data
- — $f(x)$ best guess
- $f(x)$ 95% interval

4

# Multivariate Gaussian

# Multivariate Gaussian

- *M*=2, bivariate Gaussian: $[y^{(1)}, y^{(2)}]^\top \sim \mathcal{N}(\mu, K)$
  - With $\mu = [0, 0]^\top$ and $K = \sigma^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$ correlation

# Multivariate Gaussian

- *M*=2, bivariate Gaussian: $[y^{(1)}, y^{(2)}]^\top \sim \mathcal{N}(\mu, K)$

  - With $\mu = [0, 0]^\top$ and $K = \sigma^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$ correlation

# Multivariate Gaussian using locations

- $M$=2, bivariate Gaussian: $[y^{(1)}, y^{(2)}]^\top \sim \mathcal{N}(\mu, K)$

  - With $\mu = [0, 0]^\top$ and $K = \sigma^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$ correlation

# Multivariate Gaussian using locations

- *M*=2, bivariate Gaussian: $[y^{(1)}, y^{(2)}]^\top \sim \mathcal{N}(\mu, K)$

  - With $\mu = [0,0]^\top$ and $K = \sigma^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$ correlation

# Multivariate Gaussian using locations

- $M{=}2$, bivariate Gaussian: $[y^{(1)}, y^{(2)}]^\top \sim \mathcal{N}(\mu, K)$

  - With $\mu = [0, 0]^\top$ and $K = \sigma^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$ correlation

# Multivariate Gaussian using locations

- $M=2$, bivariate Gaussian: $[y^{(1)}, y^{(2)}]^\top \sim \mathcal{N}(\mu, K)$

  - With $\mu = [0, 0]^\top$ and $K = \sigma^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$ correlation

# Multivariate Gaussian using locations

- $M=2$, bivariate Gaussian: $[y^{(1)}, y^{(2)}]^\top \sim \mathcal{N}(\mu, K)$

  - With $\mu = [0, 0]^\top$ and $K = \sigma^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$ correlation



  - What if we let the correlation depend on the *x*'s?

# Multivariate Gaussian using locations

- $M$=2, bivariate Gaussian: $[y^{(1)}, y^{(2)}]^\top \sim \mathcal{N}(\mu, K)$

  - With $\mu = [0, 0]^\top$ and $K = \sigma^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$ correlation

- What if we let the correlation depend on the $x$'s?
  - Let $\rho = \rho(|x^{(1)} - x^{(2)}|)$

# Multivariate Gaussian using locations

- $M$=2, bivariate Gaussian: $[y^{(1)}, y^{(2)}]^\top \sim \mathcal{N}(\mu, K)$

  - With $\mu = [0, 0]^\top$ and $K = \sigma^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$ correlation

- What if we let the correlation depend on the $x$'s?
  - Let $\rho = \rho(|x^{(1)} - x^{(2)}|)$
    - Where the correlation goes to 1 as the $x$'s get close

# Multivariate Gaussian using locations

- $M=2$, bivariate Gaussian: $[y^{(1)}, y^{(2)}]^\top \sim \mathcal{N}(\mu, K)$

  - With $\mu = [0, 0]^\top$ and $K = \sigma^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$ correlation



- What if we let the correlation depend on the *x*'s?
  - Let $\rho = \rho(|x^{(1)} - x^{(2)}|)$
    - Where the correlation goes to 1 as the *x*'s get close
    - And goes to 0 as the *x*'s get far
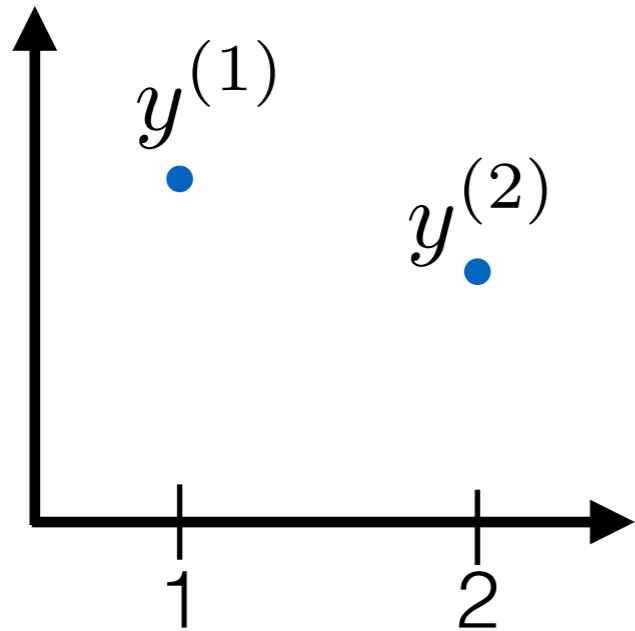
# Multivariate Gaussian using locations

- $M$=2, bivariate Gaussian: $[y^{(1)}, y^{(2)}]^\top \sim \mathcal{N}(\mu, K)$

  - With $\mu = [0, 0]^\top$ and $K = \sigma^2 \begin{bmatrix} 1 & \rho \\ \rho & 1 \end{bmatrix}$ correlation
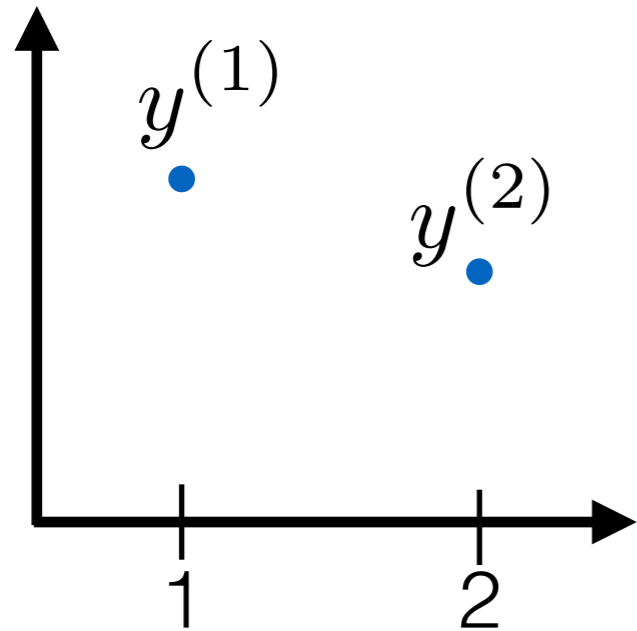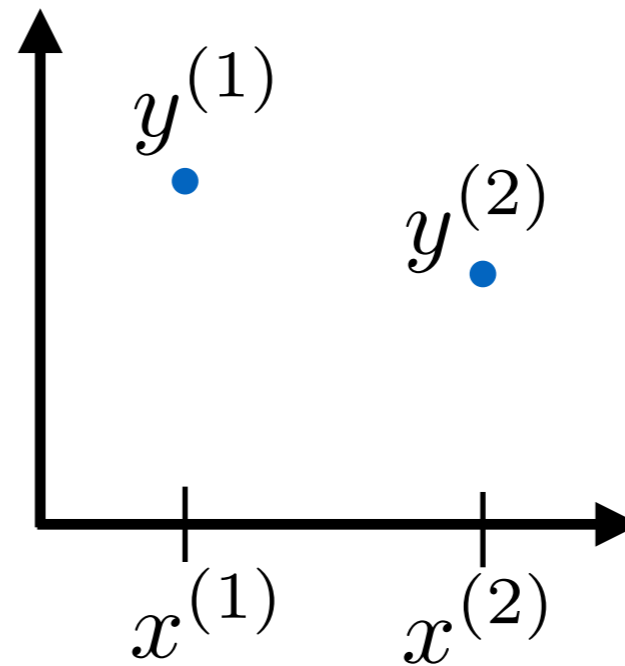


- What if we let the correlation depend on the *x*'s?
  - Let $\rho = \rho(|x^{(1)} - x^{(2)}|)$
    - Where the correlation goes to 1 as the *x*'s get close
    - And goes to 0 as the *x*'s get far

[demo]

# Multivariate Gaussian using locations

- Next: Similar setup but an $M$-long Gaussian instead of just 2-long (bivariate)

# Multivariate Gaussian using locations

- Next: Similar setup but an *M*-long Gaussian instead of just 2-long (bivariate)
  - We have *M* locations (need not be in order):
  $$x^{(m)} \in (-\infty, \infty)$$

# Multivariate Gaussian using locations

- Next: Similar setup but an *M*-long Gaussian instead of just 2-long (bivariate)
  - We have *M* locations (need not be in order):
  $$x^{(m)} \in (-\infty, \infty)$$
  - We're going to generate $[y^{(1)}, \ldots, y^{(M)}]^\top \sim \mathcal{N}(\mu, K)$

# Multivariate Gaussian using locations

- Next: Similar setup but an *M*-long Gaussian instead of just 2-long (bivariate)
  - We have *M* locations (need not be in order):
  $$x^{(m)} \in (-\infty, \infty)$$
  - We're going to generate $[y^{(1)}, \ldots, y^{(M)}]^\top \sim \mathcal{N}(\mu, K)$
    - Take $\mu = \mathbf{0}_M$ and *K* such that
    $$K_{m,m'} = \sigma^2 \rho(|x^{(m)} - x^{(m')}|)$$

# Multivariate Gaussian using locations

- Next: Similar setup but an *M*-long Gaussian instead of just 2-long (bivariate)
  - We have *M* locations (need not be in order):
  $$x^{(m)} \in (-\infty, \infty)$$
  - We're going to generate $[y^{(1)}, \ldots, y^{(M)}]^\top \sim \mathcal{N}(\mu, K)$
    - Take $\mu = \mathbf{0}_M$ and *K* such that
    $$K_{m,m'} = \sigma^2 \rho(|x^{(m)} - x^{(m')}|)$$
    - Let's try $\rho(\Delta) = \exp(-\frac{1}{2}\Delta^2)$

# Multivariate Gaussian using locations

- Next: Similar setup but an *M*-long Gaussian instead of just 2-long (bivariate)
  - We have *M* locations (need not be in order):
  $$x^{(m)} \in (-\infty, \infty)$$
  - We're going to generate $[y^{(1)}, \ldots, y^{(M)}]^\top \sim \mathcal{N}(\mu, K)$
    - Take $\mu = \mathbf{0}_M$ and *K* such that
    $$K_{m,m'} = \sigma^2 \rho(|x^{(m)} - x^{(m')}|)$$
    - Let's try $\rho(\Delta) = \exp(-\frac{1}{2}\Delta^2)$
      - Check:

# Multivariate Gaussian using locations

- Next: Similar setup but an *M*-long Gaussian instead of just 2-long (bivariate)
  - We have *M* locations (need not be in order):
  $$x^{(m)} \in (-\infty, \infty)$$
- We're going to generate $[y^{(1)}, \ldots, y^{(M)}]^\top \sim \mathcal{N}(\mu, K)$
  - Take $\mu = \mathbf{0}_M$ and *K* such that
  $$K_{m,m'} = \sigma^2 \rho(|x^{(m)} - x^{(m')}|)$$
  - Let's try $\rho(\Delta) = \exp(-\frac{1}{2}\Delta^2)$
    - Check: $\rho(0) = $ ?

# Multivariate Gaussian using locations

- Next: Similar setup but an *M*-long Gaussian instead of just 2-long (bivariate)
  - We have *M* locations (need not be in order):
  $$x^{(m)} \in (-\infty, \infty)$$
- We're going to generate $[y^{(1)}, \ldots, y^{(M)}]^\top \sim \mathcal{N}(\mu, K)$
  - Take $\mu = \mathbf{0}_M$ and *K* such that
  $$K_{m,m'} = \sigma^2 \rho(|x^{(m)} - x^{(m')}|)$$
  - Let's try $\rho(\Delta) = \exp(-\frac{1}{2}\Delta^2)$
    - Check: $\rho(0) = \boxed{1}$

# Multivariate Gaussian using locations

- Next: Similar setup but an *M*-long Gaussian instead of just 2-long (bivariate)
  - We have *M* locations (need not be in order):
  $$x^{(m)} \in (-\infty, \infty)$$
- We're going to generate $[y^{(1)}, \ldots, y^{(M)}]^\top \sim \mathcal{N}(\mu, K)$
  - Take $\mu = \mathbf{0}_M$ and *K* such that
  $$K_{m,m'} = \sigma^2 \rho(|x^{(m)} - x^{(m')}|)$$
- Let's try $\rho(\Delta) = \exp(-\frac{1}{2}\Delta^2)$
  - Check: $\rho(0) = 1$ , $\rho(\Delta)$ is ▮?▮ as $\Delta$ increases

# Multivariate Gaussian using locations

- Next: Similar setup but an *M*-long Gaussian instead of just 2-long (bivariate)
  - We have *M* locations (need not be in order):
    $$x^{(m)} \in (-\infty, \infty)$$
  - We're going to generate $[y^{(1)}, \ldots, y^{(M)}]^\top \sim \mathcal{N}(\mu, K)$
    - Take $\mu = \mathbf{0}_M$ and *K* such that
      $$K_{m,m'} = \sigma^2 \rho(|x^{(m)} - x^{(m')}|)$$
    - Let's try $\rho(\Delta) = \exp(-\frac{1}{2}\Delta^2)$
      - Check: $\rho(0) = 1$, $\rho(\Delta)$ is decreasing as $\Delta$ increases

# Multivariate Gaussian using locations

- Next: Similar setup but an *M*-long Gaussian instead of just 2-long (bivariate)
  - We have *M* locations (need not be in order):
  $$x^{(m)} \in (-\infty, \infty)$$
- We're going to generate $[y^{(1)}, \ldots, y^{(M)}]^\top \sim \mathcal{N}(\mu, K)$
  - Take $\mu = \mathbf{0}_M$ and *K* such that
  $$K_{m,m'} = \sigma^2 \rho(|x^{(m)} - x^{(m')}|)$$
- Let's try $\rho(\Delta) = \exp(-\frac{1}{2}\Delta^2)$
  - Check: $\rho(0) = 1$, $\rho(\Delta)$ is decreasing as $\Delta$ increases
  - And $\rho(\Delta) \to \boxed{?}$ as $\Delta \to \infty$

# Multivariate Gaussian using locations

- Next: Similar setup but an *M*-long Gaussian instead of just 2-long (bivariate)
  - We have *M* locations (need not be in order):
  $$x^{(m)} \in (-\infty, \infty)$$
- We're going to generate $[y^{(1)}, \ldots, y^{(M)}]^\top \sim \mathcal{N}(\mu, K)$
  - Take $\mu = \mathbf{0}_M$ and *K* such that
  $$K_{m,m'} = \sigma^2 \rho(|x^{(m)} - x^{(m')}|)$$
- Let's try $\rho(\Delta) = \exp(-\frac{1}{2}\Delta^2)$
  - Check: $\rho(0) = 1$ , $\rho(\Delta)$ is decreasing as $\Delta$ increases
  - And $\rho(\Delta) \rightarrow 0$ as $\Delta \rightarrow \infty$

# Multivariate Gaussian using locations

- Next: Similar setup but an *M*-long Gaussian instead of just 2-long (bivariate)
  - We have *M* locations (need not be in order):
  $$x^{(m)} \in (-\infty, \infty)$$
  - We're going to generate $[y^{(1)}, \ldots, y^{(M)}]^\top \sim \mathcal{N}(\mu, K)$
    - Take $\mu = \mathbf{0}_M$ and *K* such that
    $$K_{m,m'} = \sigma^2 \rho(|x^{(m)} - x^{(m')}|)$$
  - Let's try $\rho(\Delta) = \exp(-\frac{1}{2}\Delta^2)$
    - Check: $\rho(0) = 1$, $\rho(\Delta)$ is decreasing as $\Delta$ increases
    - And $\rho(\Delta) \to 0$ as $\Delta \to \infty$

Note: our previous example was the special case where *M*=2

5

# Multivariate Gaussian using locations

- Next: Similar setup but an *M*-long Gaussian instead of just 2-long (bivariate)
  - We have *M* locations (need not be in order):
$$x^{(m)} \in (-\infty, \infty)$$
  - We're going to generate $[y^{(1)}, \ldots, y^{(M)}]^\top \sim \mathcal{N}(\mu, K)$
    - Take $\mu = \mathbf{0}_M$ and *K* such that
$$K_{m,m'} = \sigma^2 \rho(|x^{(m)} - x^{(m')}|)$$
  - Let's try $\rho(\Delta) = \exp(-\frac{1}{2}\Delta^2)$
    - Check: $\rho(0) = 1$, $\rho(\Delta)$ is decreasing as $\Delta$ increases
    - And $\rho(\Delta) \to 0$ as $\Delta \to \infty$

[demo1, demo2]

Note: our previous example was the special case where *M*=2

# Multivariate Gaussian using locations

- Next: Similar setup but an *M*-long Gaussian instead of just 2-long (bivariate)
  - We have *M* locations (need not be in order):
  $$x^{(m)} \in (-\infty, \infty)$$
  - We're going to generate $[y^{(1)}, \ldots, y^{(M)}]^\top \sim \mathcal{N}(\mu, K)$
    - Take $\mu = \mathbf{0}_M$ and *K* such that
    $$K_{m,m'} = \sigma^2 \rho(|x^{(m)} - x^{(m')}|)$$
  - Let's try $\rho(\Delta) = \exp(-\frac{1}{2}\Delta^2)$
    - Check: $\rho(0) = 1$, $\rho(\Delta)$ is decreasing as $\Delta$ increases
    - And $\rho(\Delta) \to 0$ as $\Delta \to \infty$

Note: our previous example was the special case where *M*=2

[demo1, demo2]

We just drew random functions from a type of "Gaussian process"!

5

# Gaussian processes

# Gaussian processes

- Definition: "A *Gaussian process* is a collection of random variables, any finite number of which have a joint Gaussian distribution." [Rasmussen and Williams 2006; a much much older idea!]

# Gaussian processes

- Definition: "A *Gaussian process* is a collection of random variables, any finite number of which have a joint Gaussian distribution." [Rasmussen and Williams 2006; a much much older idea!]
  - E.g. the function $f(\mathbf{x})$ is a collection indexed by input $\mathbf{x}$

# Gaussian processes

- Definition: "A *Gaussian process* is a collection of random variables, any finite number of which have a joint Gaussian distribution." [Rasmussen and Williams 2006; a much much older idea!]
  - E.g. the function $f(\mathbf{x})$ is a collection indexed by input $\mathbf{x}$
- It is specified by its mean function and covariance function:

$$f \sim \mathcal{GP}(m, k)$$

# Gaussian processes

- Definition: "A *Gaussian process* is a collection of random variables, any finite number of which have a joint Gaussian distribution." [Rasmussen and Williams 2006; a much much older idea!]
  - E.g. the function $f(\mathbf{x})$ is a collection indexed by input $\mathbf{x}$
- It is specified by its mean function and covariance function:

$$f \sim \mathcal{GP}(m, k)$$

  - Mean function $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$

# Gaussian processes

- Definition: "A *Gaussian process* is a collection of random variables, any finite number of which have a joint Gaussian distribution." [Rasmussen and Williams 2006; a much much older idea!]
  - E.g. the function $f(\mathbf{x})$ is a collection indexed by input $\mathbf{x}$
- It is specified by its mean function and covariance function:

$$f \sim \mathcal{GP}(m, k)$$

  - Mean function $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$

  - Covariance function (a.k.a. *kernel*)
    $$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$$

# Gaussian processes

- Definition: "A *Gaussian process* is a collection of random variables, any finite number of which have a joint Gaussian distribution." [Rasmussen and Williams 2006; a much much older idea!]
  - E.g. the function $f(\mathbf{x})$ is a collection indexed by input $\mathbf{x}$
- It is specified by its mean function and covariance function:

$$f \sim \mathcal{GP}(m, k)$$

*x* could be just about anything, but in this tutorial, we'll assume it's a real vector

- Mean function $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$
- Covariance function (a.k.a. *kernel*)

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$$

# Gaussian processes

- Definition: "A *Gaussian process* is a collection of random variables, any finite number of which have a joint Gaussian distribution." [Rasmussen and Williams 2006; a much much older idea!]
  - E.g. the function $f(\mathbf{x})$ is a collection indexed by input $\mathbf{x}$
- It is specified by its mean function and covariance function:

$$f \sim \mathcal{GP}(m, k)$$

*x* could be just about anything, but in this tutorial, we'll assume it's a real vector

  - Mean function $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$
  - Covariance function (a.k.a. *kernel*)
$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$$
- A common default (e.g. in software) is $m(\mathbf{x}) = 0$

# Gaussian processes

- Definition: "A *Gaussian process* is a collection of random variables, any finite number of which have a joint Gaussian distribution." [Rasmussen and Williams 2006; a much much older idea!]
  - E.g. the function $f(\mathbf{x})$ is a collection indexed by input $\mathbf{x}$
- It is specified by its mean function and covariance function:

$$f \sim \mathcal{GP}(m, k)$$

*$\mathbf{x}$ could be just about anything, but in this tutorial, we'll assume it's a real vector*

  - Mean function $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$
  - Covariance function (a.k.a. *kernel*)
  $$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$$

- A common default (e.g. in software) is $m(\mathbf{x}) = 0$
- One very commonly used covariance function is the *squared exponential* or *radial basis function (RBF)*

# Gaussian processes

- Definition: "A *Gaussian process* is a collection of random variables, any finite number of which have a joint Gaussian distribution." [Rasmussen and Williams 2006; a much much older idea!]
  - E.g. the function $f(\mathbf{x})$ is a collection indexed by input $\mathbf{x}$
- It is specified by its mean function and covariance function:

$$f \sim \mathcal{GP}(m, k)$$

*x* could be just about anything, but in this tutorial, we'll assume it's a real vector

  - Mean function $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$
  - Covariance function (a.k.a. *kernel*)

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$$

- A common default (e.g. in software) is $m(\mathbf{x}) = 0$
- One very commonly used covariance function is the *squared exponential* or *radial basis function (RBF)*
  - We'll see a more general form later, but for now we're using: $k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp(-\frac{1}{2}\|\mathbf{x} - \mathbf{x}'\|^2)$

# Gaussian processes

- Definition: "A *Gaussian process* is a collection of random variables, any finite number of which have a joint Gaussian distribution." [Rasmussen and Williams 2006; a much much older idea!]
  - E.g. the function $f(\mathbf{x})$ is a collection indexed by input $\mathbf{x}$
- It is specified by its mean function and covariance function:
$$f \sim \mathcal{GP}(m, k)$$
*x* could be just about anything, but in this tutorial, we'll assume it's a real vector
  - Mean function $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$
  - Covariance function (a.k.a. *kernel*)
$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$$
- A common default (e.g. in software) is $m(\mathbf{x}) = 0$
- One very commonly used covariance function is the *squared exponential* or *radial basis function (RBF)*
  - We'll see a more general form later, but for now we're using: $k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp(-\frac{1}{2}\|\mathbf{x} - \mathbf{x}'\|^2)$
- For now, assume data is observed without noise

# Gaussian processes

- Definition: "A *Gaussian process* is a collection of random variables, any finite number of which have a joint Gaussian distribution." [Rasmussen and Williams 2006; a much much older idea!]
  - E.g. the function $f(\mathbf{x})$ is a collection indexed by input $\mathbf{x}$
- It is specified by its mean function and covariance function:

$$f \sim \mathcal{GP}(m, k)$$

*x* could be just about anything, but in this tutorial, we'll assume it's a real vector

  - Mean function $m(\mathbf{x}) = \mathbb{E}[f(\mathbf{x})]$
  - Covariance function (a.k.a. *kernel*)

$$k(\mathbf{x}, \mathbf{x}') = \mathbb{E}[(f(\mathbf{x}) - m(\mathbf{x}))(f(\mathbf{x}') - m(\mathbf{x}'))]$$

- A common default (e.g. in software) is $m(\mathbf{x}) = 0$
- One very commonly used covariance function is the *squared exponential* or *radial basis function (RBF)*
  - We'll see a more general form later, but for now we're using: $k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp(-\frac{1}{2}\|\mathbf{x} - \mathbf{x}'\|^2)$
- For now, assume data is observed without noise [demo1,2]

# Inference about unknowns given data

# Inference about unknowns given data

- Let *X* collect the *N* "training" data points (indexed 1 to *N*)

  *X: NxD*

# Inference about unknowns given data

- Let *X* collect the *N* "training" data points (indexed 1 to *N*)
- Let *X'* collect the *M* "test" data points
  
  *X: NxD*

# Inference about unknowns given data

- Let *X* collect the *N* "training" data points (indexed 1 to *N*)
- Let *X'* collect the *M* "test" data points          *X: NxD*
  - Where we want to evaluate the function

# Inference about unknowns given data

- Let *X* collect the *N* "training" data points (indexed 1 to *N*)
- Let *X'* collect the *M* "test" data points    *X: NxD*
  - Where we want to evaluate the function
  - Indexed *N*+1 to *N*+*M*

# Inference about unknowns given data

- Let *X* collect the *N* "training" data points (indexed 1 to *N*)
- Let *X'* collect the *M* "test" data points $\qquad$ *X: NxD*
  - Where we want to evaluate the function
  - Indexed *N*+1 to *N*+*M*

$$\begin{array}{c} Nx1 \\ Mx1 \end{array} \begin{bmatrix} f(X) \\ f(X') \end{bmatrix}$$

# Inference about unknowns given data

- Let *X* collect the *N* "training" data points (indexed 1 to *N*)
- Let *X'* collect the *M* "test" data points       *X: NxD*
  - Where we want to evaluate the function
  - Indexed *N+*1 to *N+M*

- Then by our model

$$\begin{matrix} Nx1 \\ Mx1 \end{matrix} \begin{bmatrix} f(X) \\ f(X') \end{bmatrix} \sim \boxed{?}$$

# Inference about unknowns given data

- Let *X* collect the *N* "training" data points (indexed 1 to *N*)
- Let *X'* collect the *M* "test" data points
  *X: NxD*
  - Where we want to evaluate the function
  - Indexed *N+1* to *N+M*

- Then by our model

$$\begin{matrix} Nx1 \\ Mx1 \end{matrix} \begin{bmatrix} f(X) \\ f(X') \end{bmatrix} \sim \mathcal{N}$$

# Inference about unknowns given data

- Let $X$ collect the $N$ "training" data points (indexed 1 to $N$)
- Let $X'$ collect the $M$ "test" data points $\quad$ <span style="color:#6fb7e8">$X: N \times D$</span>
  - Where we want to evaluate the function
  - Indexed $N{+}1$ to $N{+}M$

- Then by our model

$$\begin{array}{c}\color{#6fb7e8}N\times 1\\[8pt]\color{#6fb7e8}M\times 1\end{array}\begin{bmatrix} f(X) \\ f(X') \end{bmatrix} \sim \mathcal{N}\left( \quad ? \quad \right.$$

# Inference about unknowns given data

- Let *X* collect the *N* "training" data points (indexed 1 to *N*)
- Let *X'* collect the *M* "test" data points      *X: NxD*
  - Where we want to evaluate the function
  - Indexed *N+1* to *N+M*

- Then by our model

$$
\begin{array}{c} Nx1 \\ Mx1 \end{array} \begin{bmatrix} f(X) \\ f(X') \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \right.
$$

$$\underbrace{\qquad}_{(N+M)x1}$$

# Inference about unknowns given data

- Let *X* collect the *N* "training" data points (indexed 1 to *N*)
- Let *X'* collect the *M* "test" data points $\quad X: NxD$
  - Where we want to evaluate the function
  - Indexed *N*+1 to *N*+*M*
- $K(X,X')$ is the *NxM* matrix with (*n*,*m*) entry $k(x^{(n)}, x^{(N+m)})$
- Then by our model

$$\begin{array}{l} Nx1 \\ Mx1 \end{array} \begin{bmatrix} f(X) \\ f(X') \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \right.$$

$$\underbrace{\phantom{aaaaaaaa}}_{(N+M)x1}$$

7

# Inference about unknowns given data

- Let *X* collect the *N* "training" data points (indexed 1 to *N*)
- Let *X'* collect the *M* "test" data points      *X: NxD*
  - Where we want to evaluate the function
  - Indexed *N+1* to *N+M*
- *K(X,X')* is the *NxM* matrix with (*n*,*m*) entry $k(x^{(n)}, x^{(N+m)})$
- Then by our model

$$
\begin{array}{l} Nx1 \\ Mx1 \end{array}
\begin{bmatrix} f(X) \\ f(X') \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K(X,X) & K(X,X') \\ K(X',X) & K(X',X') \end{bmatrix}\right)
$$

$\underbrace{\qquad\qquad}_{(N+M)x1}$

# Inference about unknowns given data

- Let *X* collect the *N* "training" data points (indexed 1 to *N*)
- Let *X'* collect the *M* "test" data points     *X: NxD*
  - Where we want to evaluate the function
  - Indexed *N+1* to *N+M*
- *K(X,X')* is the *NxM* matrix with (*n,m*) entry $k(x^{(n)}, x^{(N+m)})$
- Then by our model

$$
\begin{array}{c} Nx1 \\ Mx1 \end{array}
\begin{bmatrix} f(X) \\ f(X') \end{bmatrix}
\sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} K(X,X) & K(X,X') \\ K(X',X) & K(X',X') \end{bmatrix} \right)
$$

$\underbrace{\phantom{xxxxx}}_{(N+M)x1}$   $\underbrace{\phantom{xxxxxxxxxxxxxxxxx}}_{(N+M)x(N+M)}$

# Inference about unknowns given data

- Let *X* collect the *N* "training" data points (indexed 1 to *N*)
- Let *X'* collect the *M* "test" data points $\quad$ *X: NxD*
  - Where we want to evaluate the function
  - Indexed *N+*1 to *N+M*
- *K(X,X')* is the *NxM* matrix with (*n,m*) entry $k(x^{(n)}, x^{(N+m)})$
- Then by our model

$$
\begin{array}{c} Nx1 \\ Mx1 \end{array}
\begin{bmatrix} f(X) \\ f(X') \end{bmatrix}
\sim \mathcal{N}\left( \mathbf{0}, \begin{bmatrix} K(X,X) & K(X,X') \\ K(X',X) & K(X',X') \end{bmatrix} \right)
$$

$$\underbrace{\phantom{0}}_{(N+M)x1} \qquad \underbrace{\phantom{K(X',X')}}_{(N+M)x(N+M)}$$

- The conditional satisfies $f(X')|f(X), X, X' \sim$

# Inference about unknowns given data

- Let *X* collect the *N* "training" data points (indexed 1 to *N*)
- Let *X'* collect the *M* "test" data points          *X: NxD*
  - Where we want to evaluate the function
  - Indexed *N*+1 to *N*+*M*
- *K*(*X*,*X'*) is the *NxM* matrix with (*n*,*m*) entry $k(x^{(n)}, x^{(N+m)})$
- Then by our model

$$\begin{array}{c} N\text{x}1 \\ M\text{x}1 \end{array} \begin{bmatrix} f(X) \\ f(X') \end{bmatrix} \sim \mathcal{N}\left( \mathbf{0}, \begin{bmatrix} K(X,X) & K(X,X') \\ K(X',X) & K(X',X') \end{bmatrix} \right)$$

$\underbrace{\qquad}_{(N+M)\text{x}1} \underbrace{\qquad\qquad}_{(N+M)\text{x}(N+M)}$

- The conditional satisfies $f(X')|f(X), X, X' \sim \mathcal{N}$

7

# Inference about unknowns given data

- Let $X$ collect the $N$ "training" data points (indexed 1 to $N$)
- Let $X'$ collect the $M$ "test" data points $\qquad$ <span style="color:#4a90d9">$X$: $N$x$D$</span>
  - Where we want to evaluate the function
  - Indexed $N+1$ to $N+M$
- $K(X,X')$ is the $N$x$M$ matrix with ($n,m$) entry $k(x^{(n)}, x^{(N+m)})$
- Then by our model

$$\begin{matrix} Nx1 \\ Mx1 \end{matrix} \begin{bmatrix} f(X) \\ f(X') \end{bmatrix} \sim \mathcal{N}\left(\mathbf{0}, \begin{bmatrix} K(X,X) & K(X,X') \\ K(X',X) & K(X',X') \end{bmatrix}\right)$$

<span style="color:#4a90d9">$(N+M)$x1 $\qquad$ $(N+M)$x$(N+M)$</span>

- The conditional satisfies $f(X')|f(X), X, X' \sim \mathcal{N}$
  - Can compute mean & covariance in closed form with Gaussian facts $\quad$ <span style="color:#4a90d9">$M$x1 $\qquad$ $M$x$M$</span>

7

# Inference about unknowns given data

- Let *X* collect the *N* "training" data points (indexed 1 to *N*)
- Let *X'* collect the *M* "test" data points     $X: N \times D$
  - Where we want to evaluate the function
  - Indexed *N*+1 to *N*+*M*
- *K*(*X*,*X'*) is the *N*x*M* matrix with (*n*,*m*) entry $k(x^{(n)}, x^{(N+m)})$
- Then by our model

$$N\times1 \begin{bmatrix} f(X) \\ f(X') \end{bmatrix} \sim \mathcal{N} \left( \mathbf{0}, \begin{bmatrix} K(X,X) & K(X,X') \\ K(X',X) & K(X',X') \end{bmatrix} \right)$$

$(N+M)\times1 \qquad (N+M)\times(N+M)$

- The conditional satisfies $f(X')|f(X), X, X' \sim \mathcal{N}$
  - Can compute mean & covariance in closed form with Gaussian facts     $M\times1 \qquad M\times M$

- We'll infer *f*(*X'*) given our simulated data; recall we're using
$$k(x, x') = \sigma^2 \exp(-\tfrac{1}{2}(x - x')^2), \sigma = 1$$

# Inference about unknowns given data

- Let *X* collect the *N* "training" data points (indexed 1 to *N*)
- Let *X'* collect the *M* "test" data points $\qquad$ <span style="color:#6cb4e4">*X: NxD*</span>
  - Where we want to evaluate the function
  - Indexed *N+1* to *N+M*
- *K(X,X')* is the *NxM* matrix with (*n,m*) entry $k(x^{(n)}, x^{(N+m)})$
- Then by our model

$$
\begin{array}{c} {\color{#6cb4e4} N\text{x}1} \\ {\color{#6cb4e4} M\text{x}1} \end{array}
\begin{bmatrix} f(X) \\ f(X') \end{bmatrix} \sim \mathcal{N}\left( \mathbf{0}, \begin{bmatrix} K(X,X) & K(X,X') \\ K(X',X) & K(X',X') \end{bmatrix} \right)
$$

<span style="color:#6cb4e4">*(N+M)x1*</span> $\qquad$ <span style="color:#6cb4e4">*(N+M)x(N+M)*</span>

- The conditional satisfies $f(X')|f(X), X, X' \sim \mathcal{N}$
  - Can compute mean & covariance in closed form with Gaussian facts $\quad$ <span style="color:#6cb4e4">*Mx1*</span> $\qquad$ <span style="color:#6cb4e4">*MxM*</span>

- We'll infer *f(X')* given our simulated data; recall we're using
$$ k(x,x') = \sigma^2 \exp(-\tfrac{1}{2}(x-x')^2), \sigma = 1 $$

7 $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad$ [demo1,2]

# Squared exponential kernel revisited

# Squared exponential kernel revisited

- What if we happened to measure our data on a different scale?                [demo1]

# Squared exponential kernel revisited

- What if we happened to measure our data on a different scale? [demo1]
- We've been using this particular kernel:
$$k(x, x') = \sigma^2 \exp(-\tfrac{1}{2}(x - x')^2), \sigma = 1$$

# Squared exponential kernel revisited

- What if we happened to measure our data on a different scale?                         [demo1]
- We've been using this particular kernel:

$$k(x, x') = \sigma^2 \exp(-\tfrac{1}{2}(x - x')^2), \sigma = 1$$

  - What do we expect from the scale of *f*(*x*) a priori?

8

# Squared exponential kernel revisited

- What if we happened to measure our data on a different scale?                    [demo1]

- We've been using this particular kernel:
  $$k(x, x') = \sigma^2 \exp(-\tfrac{1}{2}(x - x')^2), \sigma = 1$$

  - What do we expect from the scale of *f*(*x*) a priori?

    - At one *x*, with ~95% probability a priori, $f(x) \in$    ?

# Squared exponential kernel revisited

- What if we happened to measure our data on a different scale?                                            [demo1]
- We've been using this particular kernel:
$$k(x, x') = \sigma^2 \exp(-\tfrac{1}{2}(x - x')^2), \sigma = 1$$

  - What do we expect from the scale of *f*(*x*) a priori?
    - At one *x*, with ~95% probability a priori, $f(x) \in (-2, 2)$

# Squared exponential kernel revisited

- What if we happened to measure our data on a different scale?                                   [demo1]
- We've been using this particular kernel:
  $$k(x, x') = \sigma^2 \exp(-\tfrac{1}{2}(x - x')^2), \sigma = 1$$
  - What do we expect from the scale of *f*(*x*) a priori?
    - At one *x*, with ~95% probability a priori, $f(x) \in (-2, 2)$
    - Marginal variance cannot increase with data

# Squared exponential kernel revisited

- What if we happened to measure our data on a different scale?  [demo1, demo2]
- We've been using this particular kernel:
$$k(x, x') = \sigma^2 \exp(-\tfrac{1}{2}(x - x')^2), \sigma = 1$$
  - What do we expect from the scale of *f*(*x*) a priori?
    - At one *x*, with ~95% probability a priori, $f(x) \in (-2, 2)$
    - Marginal variance cannot increase with data

# Squared exponential kernel revisited

- What if we happened to measure our data on a different scale?  [demo1, demo2]
- We've been using this particular kernel:
$$k(x, x') = \sigma^2 \exp(-\tfrac{1}{2}(x - x')^2), \sigma = 1$$
  - What do we expect from the scale of *f*(*x*) a priori?
    - At one *x*, with ~95% probability a priori, $f(x) \in (-2, 2)$
    - Marginal variance cannot increase with data
  - What counts as "close" in *x*?

# Squared exponential kernel revisited

- What if we happened to measure our data on a different scale?                    [demo1, demo2]
- We've been using this particular kernel:
$$k(x, x') = \sigma^2 \boxed{\exp(-\tfrac{1}{2}(x - x')^2)}, \sigma = 1$$
  - What do we expect from the scale of *f*(*x*) a priori?
    - At one *x*, with ~95% probability a priori, $f(x) \in (-2, 2)$
    - Marginal variance cannot increase with data
  - What counts as "close" in *x*?

# Squared exponential kernel revisited

- What if we happened to measure our data on a different scale? [demo1, demo2]
- We've been using this particular kernel:
$$k(x, x') = \sigma^2 \exp(-\tfrac{1}{2}(x - x')^2), \sigma = 1$$
  - What do we expect from the scale of *f*(*x*) a priori?
    - At one *x*, with ~95% probability a priori, $f(x) \in (-2, 2)$
    - Marginal variance cannot increase with data
  - What counts as "close" in *x*?

$$\exp(-\tfrac{1}{2}2^2) \approx 0.14 \qquad \exp(-\tfrac{1}{2}3^2) \approx 0.011 \qquad \exp(-\tfrac{1}{2}4^2) \approx 0.00034$$

# Squared exponential kernel revisited

- What if we happened to measure our data on a different scale? [demo1, demo2]
- We've been using this particular kernel:
$$k(x, x') = \sigma^2 \exp(-\tfrac{1}{2}(x - x')^2), \sigma = 1$$
  - What do we expect from the scale of *f*(*x*) a priori?
    - At one *x*, with ~95% probability a priori, $f(x) \in (-2, 2)$
    - Marginal variance cannot increase with data
  - What counts as "close" in *x*?

$$\exp(-\tfrac{1}{2}2^2) \approx 0.14 \qquad \exp(-\tfrac{1}{2}3^2) \approx 0.011 \qquad \exp(-\tfrac{1}{2}4^2) \approx 0.00034$$

- What can we do to handle different *x* and *f*(*x*) scales?

# Squared exponential kernel revisited

- What if we happened to measure our data on a different scale? [demo1, demo2]
- We've been using this particular kernel:
$$k(x, x') = \sigma^2 \exp(-\tfrac{1}{2}(x - x')^2), \sigma = 1$$
  - What do we expect from the scale of *f*(*x*) a priori?
    - At one *x*, with ~95% probability a priori, $f(x) \in (-2, 2)$
    - Marginal variance cannot increase with data
  - What counts as "close" in *x*?

$$\exp(-\tfrac{1}{2}2^2) \approx 0.14 \quad \exp(-\tfrac{1}{2}3^2) \approx 0.011 \quad \exp(-\tfrac{1}{2}4^2) \approx 0.00034$$

- What can we do to handle different *x* and *f*(*x*) scales?
  - Normalization in *y* can help; in *x,* can still be hiccups
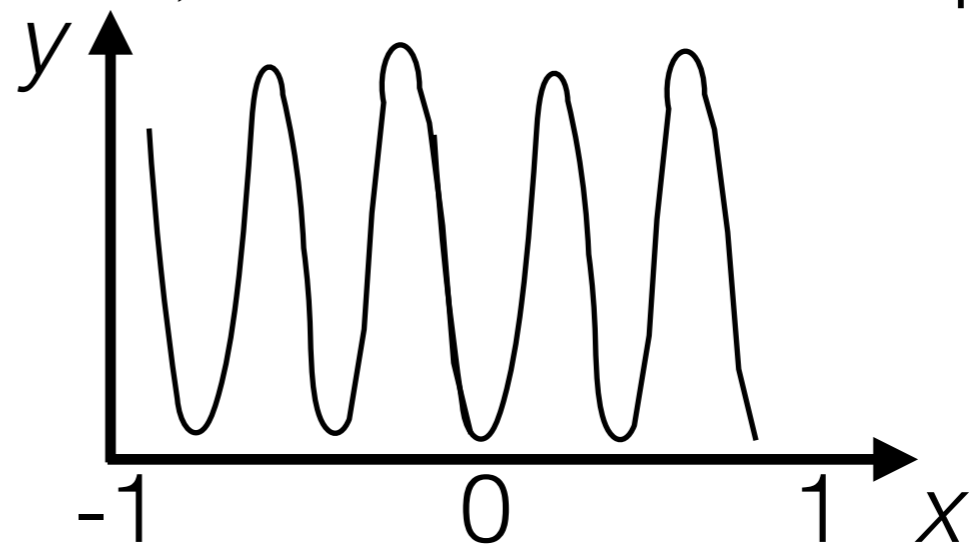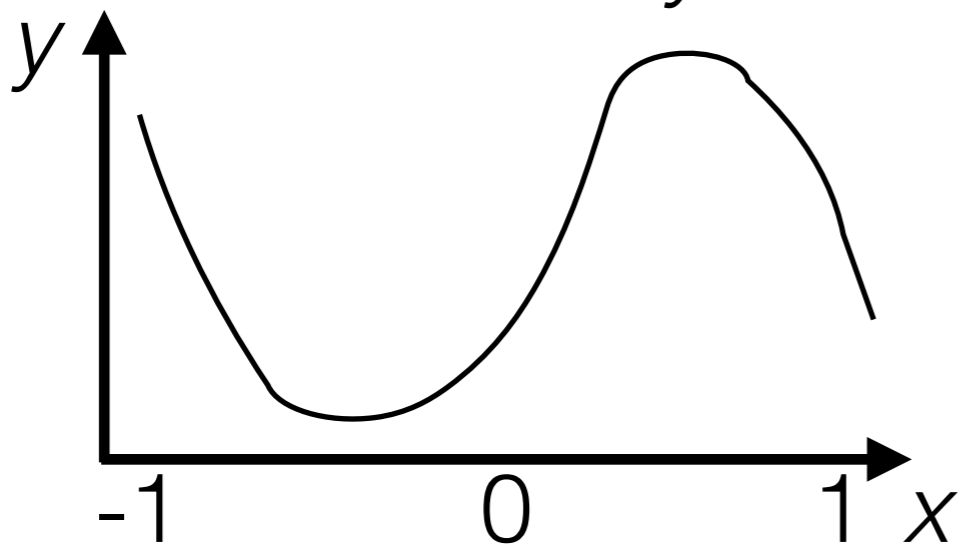
# Squared exponential kernel revisited

- What if we happened to measure our data on a different scale? [demo1, demo2]
- We've been using this particular kernel:
$$k(x, x') = \sigma^2 \exp(-\tfrac{1}{2}(x - x')^2), \sigma = 1$$
  - What do we expect from the scale of *f*(*x*) a priori?
    - At one *x*, with ~95% probability a priori, $f(x) \in (-2, 2)$
    - Marginal variance cannot increase with data
  - What counts as "close" in *x*?

$$\exp(-\tfrac{1}{2}2^2) \approx 0.14 \qquad \exp(-\tfrac{1}{2}3^2) \approx 0.011 \qquad \exp(-\tfrac{1}{2}4^2) \approx 0.00034$$

- What can we do to handle different *x* and *f*(*x*) scales?
  - Normalization in *y* can help; in *x*, can still be hiccups

# Squared exponential kernel revisited

- A common option in practice and in software is to fit the *hyperparameters* of a more general *squared exponential kernel* from data

# Squared exponential kernel revisited

- A common option in practice and in software is to fit the *hyperparameters* of a more general *squared exponential kernel* from data
  - More general form of the squared exponential:

$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp\left(-\frac{1}{2} \sum_{d=1}^{D} \frac{(x_d - x'_d)^2}{\ell_d^2}\right)$$

# Squared exponential kernel revisited

- A common option in practice and in software is to fit the *hyperparameters* of a more general *squared exponential kernel* from data
  - More general form of the squared exponential:

$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp\left(-\frac{1}{2} \sum_{d=1}^{D} \frac{(x_d - x'_d)^2}{\ell_d^2}\right)$$

signal variance          lengthscales

# Squared exponential kernel revisited

- A common option in practice and in software is to fit the *hyperparameters* of a more general *squared exponential kernel* from data
  - More general form of the squared exponential:

$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp(-\frac{1}{2} \sum_{d=1}^{D} \frac{(x_d - x'_d)^2}{\ell_d^2})$$

signal variance    lengthscales

  - *Parameters* (here, *f*) parametrize the distribution of the data. If we knew them, we could generate the data.

# Squared exponential kernel revisited

- A common option in practice and in software is to fit the *hyperparameters* of a more general *squared exponential kernel* from data
  - More general form of the squared exponential:

$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp(-\frac{1}{2} \sum_{d=1}^{D} \frac{(x_d - x'_d)^2}{\ell_d^2})$$

signal variance                 lengthscales

  - *Parameters* (here, *f*) parametrize the distribution of the data. If we knew them, we could generate the data.
    - GPs: *nonparametric* model: infinite # of latent params

# Squared exponential kernel revisited

- A common option in practice and in software is to fit the *hyperparameters* of a more general *squared exponential kernel* from data
  - More general form of the squared exponential:

$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp(-\frac{1}{2} \sum_{d=1}^{D} \frac{(x_d - x'_d)^2}{\ell_d^2})$$

signal variance          lengthscales

  - *Parameters* (here, *f*) parametrize the distribution of the data. If we knew them, we could generate the data.
    - GPs: *nonparametric* model: infinite # of latent params
  - *Hyperparameters* parametrize the distribution of the parameters. If known, we could generate the parameters.

# Squared exponential kernel revisited

- A common option in practice and in software is to fit the *hyperparameters* of a more general *squared exponential kernel* from data
  - More general form of the squared exponential:

$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp(-\frac{1}{2}\sum_{d=1}^{D}\frac{(x_d - x'_d)^2}{\ell_d^2})$$

signal variance          lengthscales

  - *Parameters* (here, *f*) parametrize the distribution of the data. If we knew them, we could generate the data.
    - GPs: *nonparametric* model: infinite # of latent params
  - *Hyperparameters* parametrize the distribution of the parameters. If known, we could generate the parameters.
- Algorithm:

# Squared exponential kernel revisited

- A common option in practice and in software is to fit the *hyperparameters* of a more general *squared exponential kernel* from data
  - More general form of the squared exponential:

$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp\left(-\frac{1}{2}\sum_{d=1}^{D}\frac{(x_d - x'_d)^2}{\ell_d^2}\right)$$

signal variance              lengthscales

  - *Parameters* (here, *f*) parametrize the distribution of the data. If we knew them, we could generate the data.
    - GPs: *nonparametric* model: infinite # of latent params
  - *Hyperparameters* parametrize the distribution of the parameters. If known, we could generate the parameters.
- Algorithm:
  - Fit a value for the hyperparameters using the data.

# Squared exponential kernel revisited

- A common option in practice and in software is to fit the *hyperparameters* of a more general *squared exponential kernel* from data
  - More general form of the squared exponential:

$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp(-\frac{1}{2}\sum_{d=1}^{D}\frac{(x_d - x'_d)^2}{\ell_d^2})$$

signal variance        lengthscales

  - *Parameters* (here, *f*) parametrize the distribution of the data. If we knew them, we could generate the data.
    - GPs: *nonparametric* model: infinite # of latent params
  - *Hyperparameters* parametrize the distribution of the parameters. If known, we could generate the parameters.
- Algorithm:
  - Fit a value for the hyperparameters using the data.
  - Given those values, now compute and report the mean and uncertainty intervals.

# Squared exponential kernel revisited

- A common option in practice and in software is to fit the *hyperparameters* of a more general *squared exponential kernel* from data
  - More general form of the squared exponential:

$$k(\mathbf{x}, \mathbf{x}') = \sigma^2 \exp(-\frac{1}{2} \sum_{d=1}^{D} \frac{(x_d - x_d')^2}{\ell_d^2})$$

signal variance                lengthscales

  - *Parameters* (here, *f*) parametrize the distribution of the data. If we knew them, we could generate the data.
    - GPs: *nonparametric* model: infinite # of latent params
  - *Hyperparameters* parametrize the distribution of the parameters. If known, we could generate the parameters.
- Algorithm:
  - Fit a value for the hyperparameters using the data.
  - Given those values, now compute and report the mean and uncertainty intervals.        [demo1,2,3]

# Roadmap

# Roadmap

- A Bayesian approach

# Roadmap

- A Bayesian approach
- What is a Gaussian process?

# Roadmap

- A Bayesian approach
- What is a Gaussian process?
  - Popular version using a squared exponential kernel

# Roadmap

- A Bayesian approach
- What is a Gaussian process?
  - Popular version using a squared exponential kernel
    - Hyperparameters: length scale(s), signal variance, observation-noise variance

# Roadmap

- A Bayesian approach
- What is a Gaussian process?
  - Popular version using a squared exponential kernel
    - Hyperparameters: length scale(s), signal variance, observation-noise variance
- Gaussian process inference

# Roadmap

- A Bayesian approach
- What is a Gaussian process?
  - Popular version using a squared exponential kernel
    - Hyperparameters: length scale(s), signal variance, observation-noise variance
- Gaussian process inference
  - Prediction & uncertainty quantification

# Roadmap

- A Bayesian approach
- What is a Gaussian process?
  - Popular version using a squared exponential kernel
    - Hyperparameters: length scale(s), signal variance, observation-noise variance
- Gaussian process inference
  - Prediction & uncertainty quantification
  - Limitations (not just for GPs):

# Roadmap

- A Bayesian approach
- What is a Gaussian process?
  - Popular version using a squared exponential kernel
    - Hyperparameters: length scale(s), signal variance, observation-noise variance
- Gaussian process inference
  - Prediction & uncertainty quantification
  - Limitations (not just for GPs):
    - E.g. extrapolation

# Roadmap

- A Bayesian approach
- What is a Gaussian process?
  - Popular version using a squared exponential kernel
    - Hyperparameters: length scale(s), signal variance, observation-noise variance
- Gaussian process inference
- Prediction & uncertainty quantification
- Limitations (not just for GPs):
  - E.g. extrapolation
  - High-dimensional inputs

# Roadmap

- A Bayesian approach
- What is a Gaussian process?
  - Popular version using a squared exponential kernel
    - Hyperparameters: length scale(s), signal variance, observation-noise variance
- Gaussian process inference
  - Prediction & uncertainty quantification
  - Limitations (not just for GPs):
    - E.g. extrapolation
    - High-dimensional inputs
- Always ask: What uncertainty are we quantifying?

# Roadmap

- A Bayesian approach
- What is a Gaussian process?
  - Popular version using a squared exponential kernel
    - Hyperparameters: length scale(s), signal variance, observation-noise variance
- Gaussian process inference
  - Prediction & uncertainty quantification
  - Limitations (not just for GPs):
    - E.g. extrapolation
    - High-dimensional inputs
- Always ask: What uncertainty are we quantifying?

Goal:
  - Learn the mechanism behind standard GPs to identify benefits and pitfalls

# Resources

- Rasmussen and Williams 2006. *Gaussian Processes for Machine Learning*. https://gaussianprocess.org/gpml/

  - Chapters 1, 2, 4, 5

- Gramacy 2020. *Surrogates: Gaussian process modeling, design and optimization for the applied sciences*. https://bookdown.org/rbg/surrogates/

- Garnett 2023. *Bayesian Optimization*. https://bayesoptbook.com/

- Software options include:

  - scikit-learn, GPy, GPflow, GPyTorch

- My setup for this tutorial: pip install X

  - X = jupyterlab, notebook, numpy, matplotlib, scikit-learn

# References (1/1)

Belkhiri, L., Tiri, A., & Mouni, L. (2020). Spatial distribution of the groundwater quality using kriging and Co-kriging interpolations. *Groundwater for Sustainable Development*, 11, 100473.

Berlinghieri, R., et al. (2023). Gaussian processes at the Helm(holtz): A more fluid model for ocean currents. *ICML*.

Binois, M., & Wycoff, N. (2022). A survey on high-dimensional Gaussian process modeling with application to Bayesian optimization. ACM Transactions on Evolutionary Learning and Optimization, 2(2), 1-26.

Garnett, R. (2023). *Bayesian Optimization*. Cambridge University Press.

Gramacy, R. B., & Lee, H. K. H. (2008). Bayesian treed Gaussian process models with an application to computer modeling. Journal of the American Statistical Association, 103(483), 1119-1130.

Gramacy, R. B. (2020). *Surrogates: Gaussian process modeling, design, and optimization for the applied sciences*. Chapman and Hall/CRC.

Ryan, E., & Özgökmen, T. Image credit.

Snoek, J., Larochelle, H., & Adams, R. P. (2012). Practical Bayesian optimization of machine learning algorithms. *NeurIPS.*

Snoek, J., et al. (2015). Scalable Bayesian optimization using deep neural networks. *ICML* (pp. 2171-2180). PMLR.

Williams, C. K., & Rasmussen, C. E. (2006). *Gaussian Processes for Machine Learning*. MIT Press.

Zewe, A. "A better way to study ocean currents." MIT News. May 17, 2023.