

Machine Learning

Gaussian Processes

Carl Henrik Ek - carlhenrik.ek@bristol.ac.uk

October 16, 2018

<http://www.carlhenrik.com>

Introduction

$$k(\mathbf{x}_i, \mathbf{x}_j) = \phi(\mathbf{x}_i)^T \phi(\mathbf{x}_j)$$

- kernel functions describe inner-products in an induced representation
- induced representation lives in what is called a Hilbert Space
- importantly the space is metric

Euclidean Distance

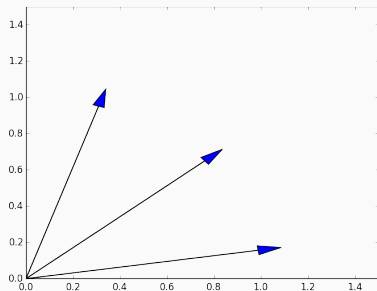
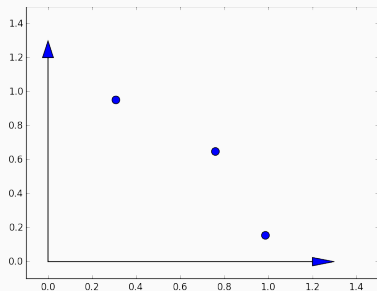
$$\|\mathbf{x}_i - \mathbf{x}_j\|^2 = \mathbf{x}_i^T \mathbf{x}_i - 2\mathbf{x}_i^T \mathbf{x}_j + \mathbf{x}_j^T \mathbf{x}_j$$

Kernelised Euclidean Distance

$$\|\phi(\mathbf{x}_i) - \phi(\mathbf{x}_j)\|^2 = k(\mathbf{x}_i, \mathbf{x}_i) - 2k(\mathbf{x}_i, \mathbf{x}_j) + k(\mathbf{x}_j, \mathbf{x}_j)$$

$$\begin{aligned}\sigma(\mathbf{X}, \mathbf{Y}) &= \mathbb{E} [(\mathbf{X} - \mathbb{E}[\mathbf{X}])^T (\mathbf{Y} - \mathbb{E}[\mathbf{Y}])] = \\ &= \mathbb{E}[\mathbf{X}^T \mathbf{Y}] - \mathbb{E}[\mathbf{X}]^T \mathbb{E}[\mathbf{Y}] = \{\mathbb{E}[\mathbf{X}] = \mathbb{E}[\mathbf{Y}] = \mathbf{0}\} = \\ &= \mathbb{E}[\mathbf{X}^T \mathbf{Y}]\end{aligned}$$

Kernels



$$\begin{aligned}\sigma(\mathbf{X}, \mathbf{Y}) &= \begin{bmatrix} x_{11} & x_{21} & x_{31} \\ x_{12} & x_{22} & x_{32} \end{bmatrix} \begin{bmatrix} y_{11} & y_{12} \\ y_{21} & y_{22} \\ y_{31} & y_{32} \end{bmatrix} = \\ &= \begin{bmatrix} x_{11}y_{11} + x_{21}y_{21} + x_{31}y_{31} & x_{11}y_{12} + x_{21}y_{22} + x_{31}y_{32} \\ x_{12}y_{11} + x_{22}y_{21} + x_{32}y_{31} & x_{12}y_{12} + x_{22}y_{22} + x_{32}y_{32} \end{bmatrix}\end{aligned}$$

$$\begin{aligned}\sigma(\mathbf{X}, \mathbf{Y}) &= \begin{bmatrix} x_{11} & x_{21} & x_{31} \\ x_{12} & x_{22} & x_{32} \end{bmatrix} \begin{bmatrix} y_{11} & y_{12} \\ y_{21} & y_{22} \\ y_{31} & y_{32} \end{bmatrix} = \\ &= \begin{bmatrix} x_{11}y_{11} + x_{21}y_{21} + x_{31}y_{31} & x_{11}y_{12} + x_{21}y_{22} + x_{31}y_{32} \\ x_{12}y_{11} + x_{22}y_{21} + x_{32}y_{31} & x_{12}y_{12} + x_{22}y_{22} + x_{32}y_{32} \end{bmatrix}\end{aligned}$$

$$\begin{aligned}\sigma(\mathbf{X}^T, \mathbf{Y}^T) &= \begin{bmatrix} x_{11} & x_{12} \\ x_{21} & x_{22} \\ x_{31} & x_{32} \end{bmatrix} \begin{bmatrix} y_{11} & y_{21} & y_{31} \\ y_{12} & y_{22} & y_{32} \end{bmatrix} = \\ &= \begin{bmatrix} x_{11}y_{11} + x_{12}y_{12} & x_{11}y_{21} + x_{12}y_{22} & x_{11}y_{31} + x_{12}y_{32} \\ x_{21}y_{11} + x_{22}y_{12} & x_{21}y_{21} + x_{22}y_{22} & x_{21}y_{31} + x_{22}y_{32} \\ x_{31}y_{11} + x_{32}y_{12} & x_{31}y_{21} + x_{32}y_{22} & x_{31}y_{31} + x_{32}y_{32} \end{bmatrix}\end{aligned}$$

Kernels and Covariances

- Covariance between columns: $\mathbf{X}^T \mathbf{Y}$ (data-dimensions)
- Covariance between rows: $\mathbf{X} \mathbf{Y}^T$ (data-points)
- Kernels: $k(\mathbf{x}, \mathbf{y}) = \phi(\mathbf{x})^T \phi(\mathbf{y})$
- A kernel function describes the co-variance of the data points

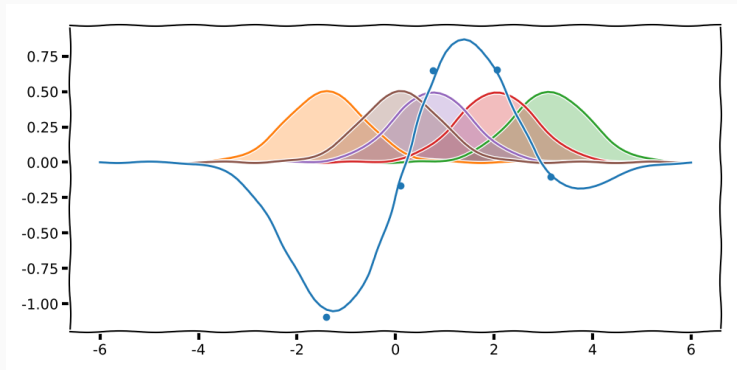
$$k(\mathbf{x}_i, \mathbf{x}_j) = \sigma^2 e^{-\frac{1}{2\ell^2}(\mathbf{x}_i - \mathbf{x}_j)^T(\mathbf{x}_i - \mathbf{x}_j)}$$

Exponentiated Quadratic

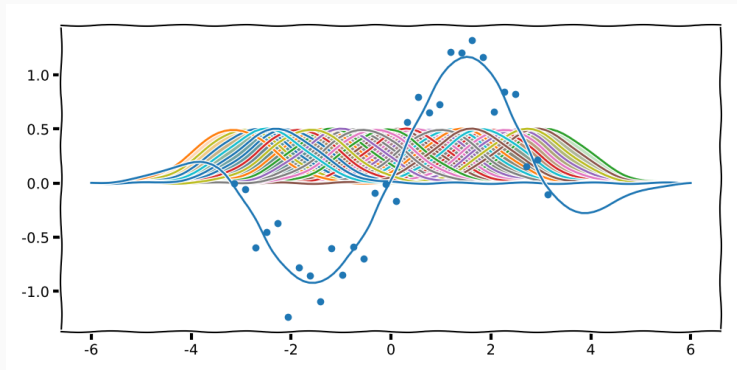
- How does the data vary along the dimensions spanned by the data
- RBF, Squared Exponential, Exponentiated Quadratic
- Co-variance smoothly decays with distance
- You can build new kernels out of other kernels [1] p. 296

$$\begin{aligned}y(\mathbf{x}_*) &= \mathbf{w}^T \mathbf{x}_* = \mathbf{a}^T \mathbf{x} \mathbf{x}_* = \mathbf{a}^T k(\mathbf{x}, \mathbf{x}_*) = \\&= ((\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{t})^T k(\mathbf{x}, \mathbf{x}_*) = k(\mathbf{x}_*, \mathbf{x})(\mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{t} \\p(\mathbf{t} | \mathbf{w}, \mathbf{x}) &= \prod_n^N p(t_n | \mathbf{w}, \mathbf{x}) = \prod_n^N \mathcal{N}(t_n | \mathbf{w}^T \mathbf{x}_n, \sigma^2 \mathbf{I}) \\p(\mathbf{w}) &= \mathcal{N}(\mathbf{0}, \tau^2 \mathbf{I})\end{aligned}$$

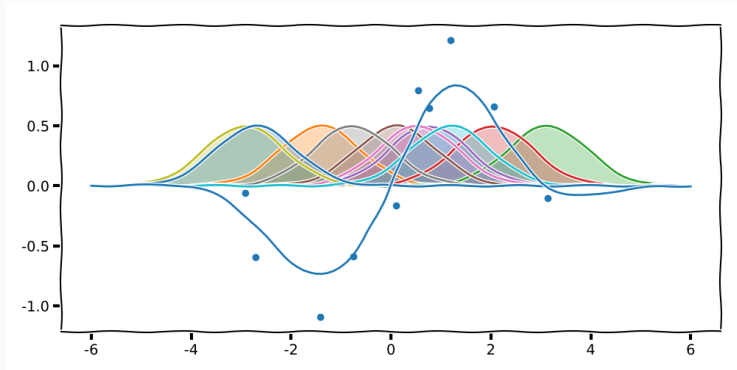
Kernel Regression



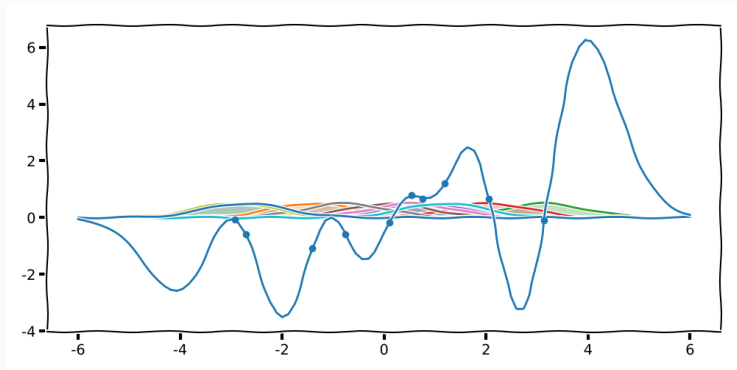
Kernel Regression



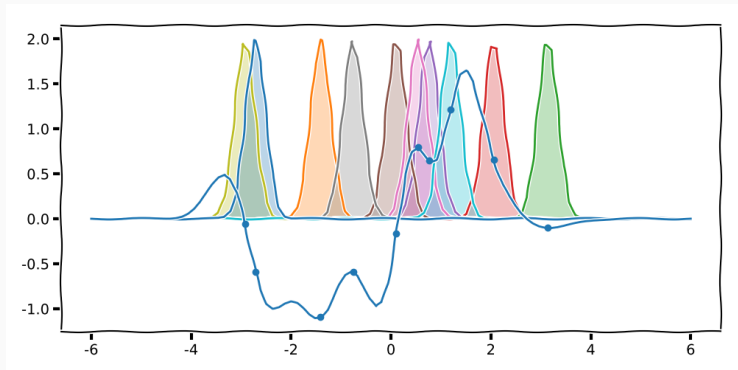
Kernel Regression



Kernel Regression



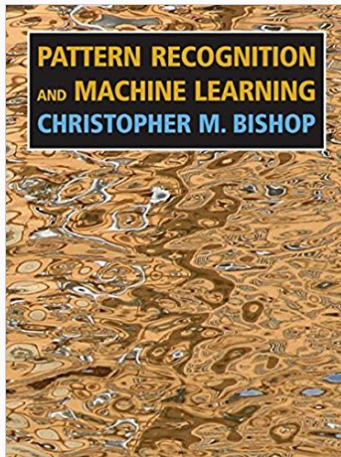
Kernel Regression



Question

- *How do you know which basis function to choose?*
- *Would you call this overfitting?*

Gaussian Processes





IUDICIUM POSTERIUM DISCIPULUS EST PRIORIS

¹<http://gpss.cc>

- We have uncertainty in our observed outputs

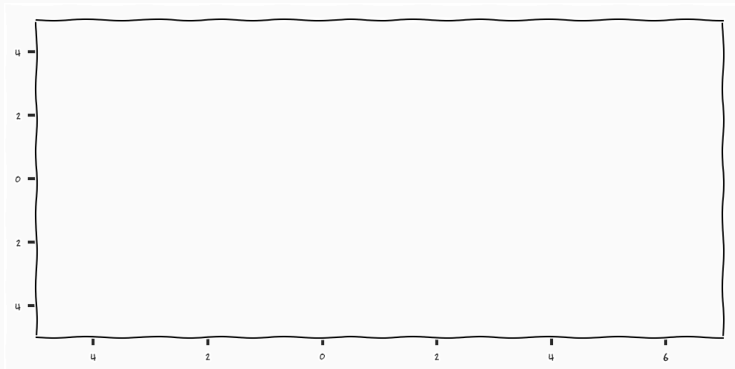
- We have uncertainty in our observed outputs
- We have no uncertainty in our mapping

- We have uncertainty in our observed outputs
- We have no uncertainty in our mapping
 - Linear, it is a line

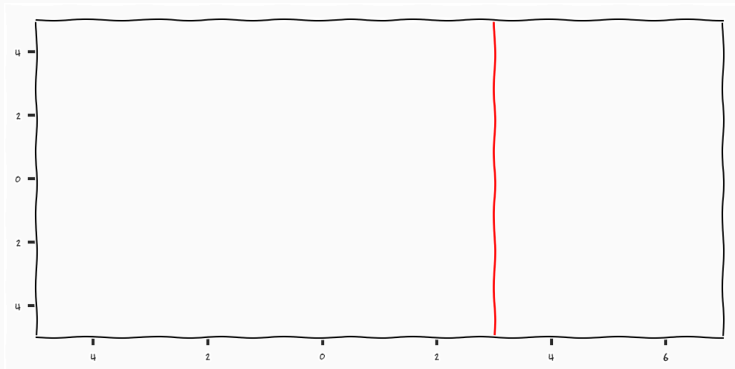
- We have uncertainty in our observed outputs
- We have no uncertainty in our mapping
 - Linear, it is a line
 - Kernels, it is this specific basis function

- We have uncertainty in our observed outputs
- We have no uncertainty in our mapping
 - Linear, it is a line
 - Kernels, it is this specific basis function
- *need a prior assumption over the space of functions*

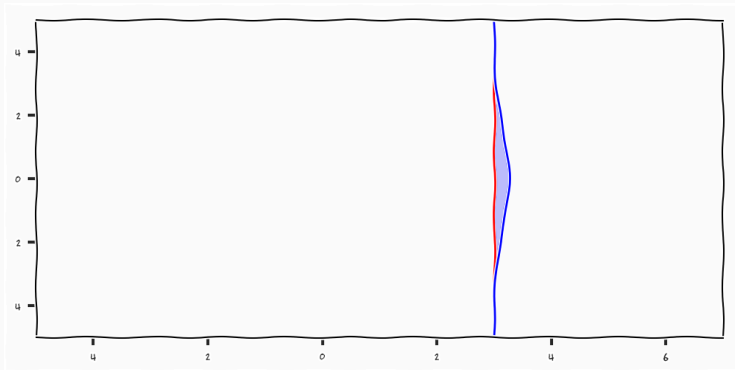
Gaussian Processes



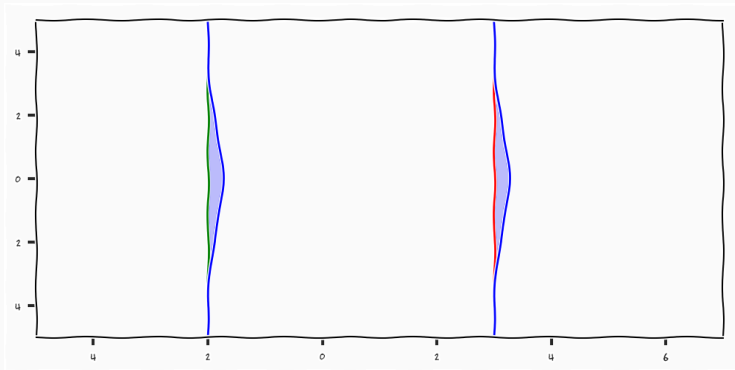
Gaussian Processes



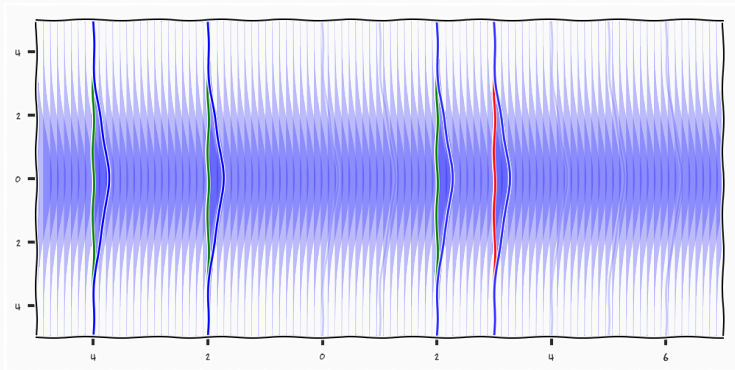
Gaussian Processes



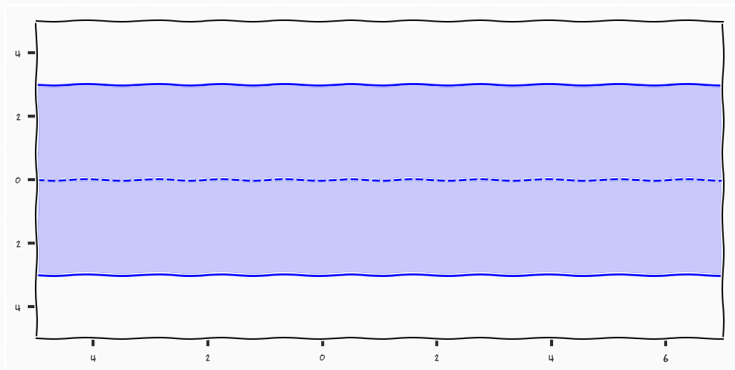
Gaussian Processes



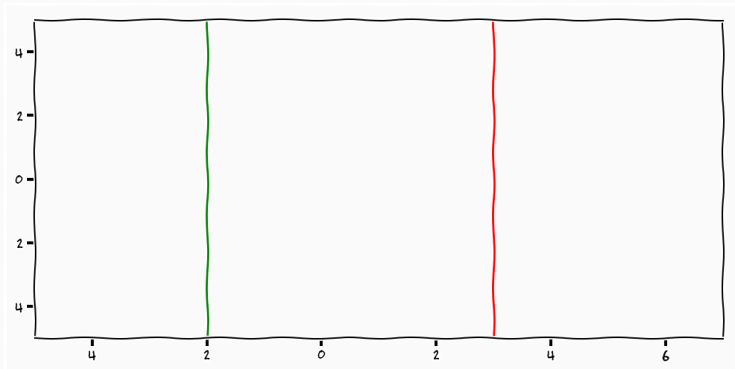
Gaussian Processes



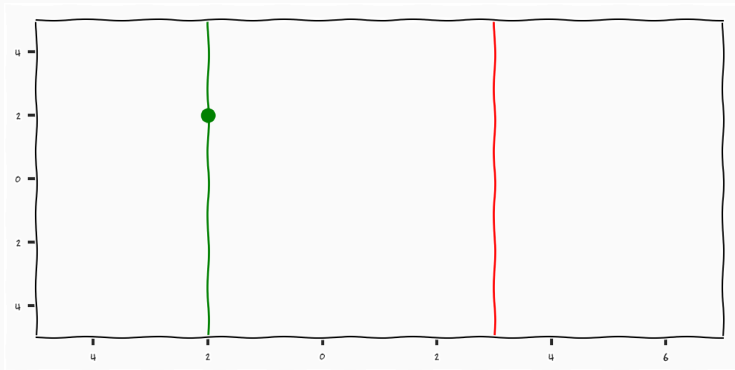
Gaussian Processes



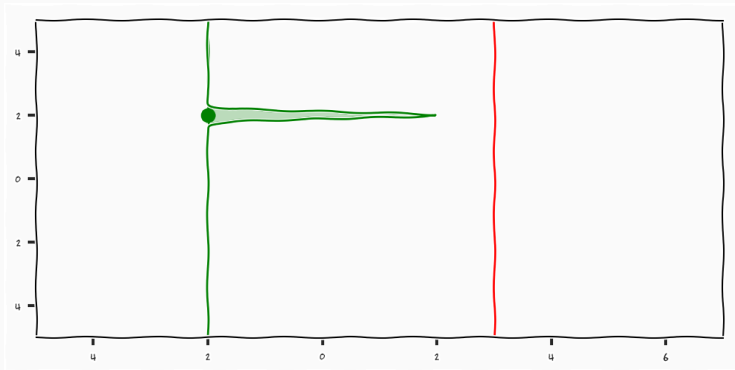
Gaussian Processes



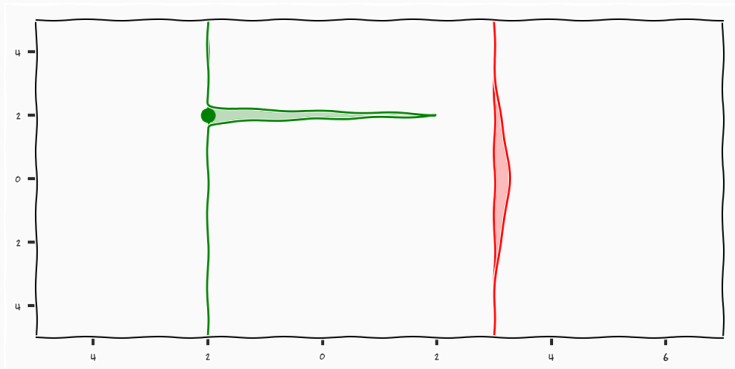
Gaussian Processes



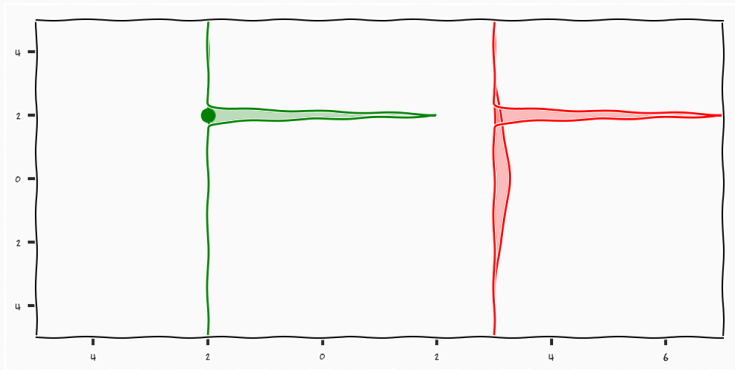
Gaussian Processes



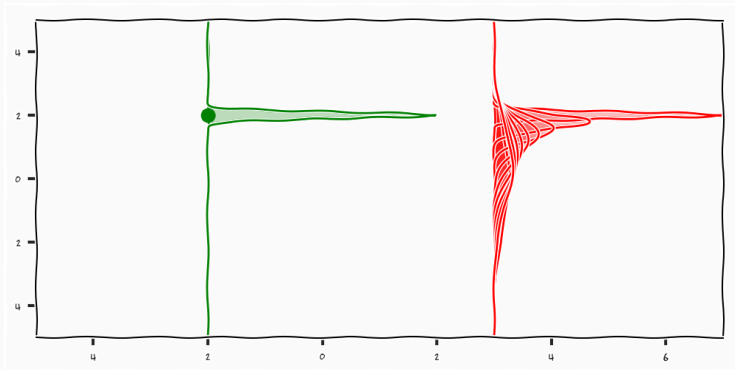
Gaussian Processes



Gaussian Processes

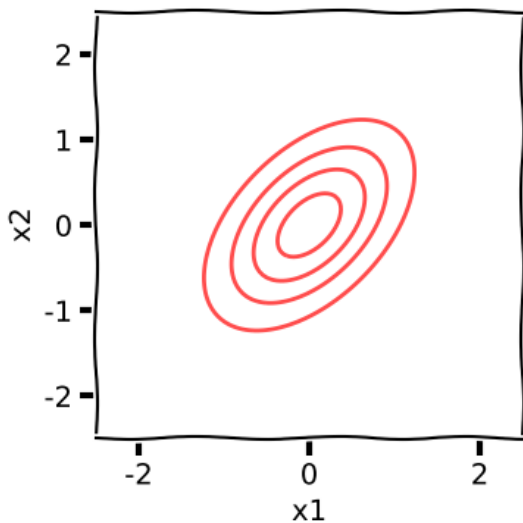


Gaussian Processes

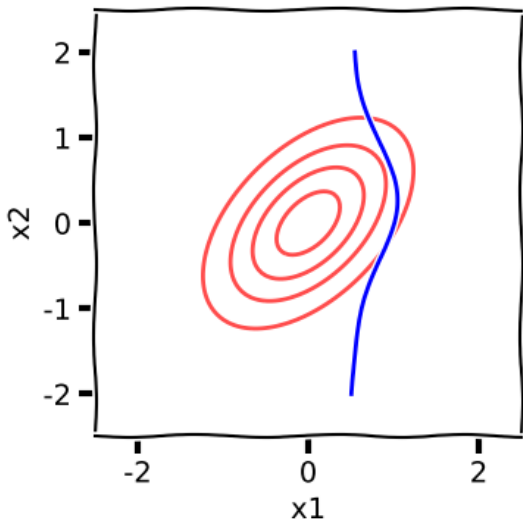


$$\mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}\right)$$

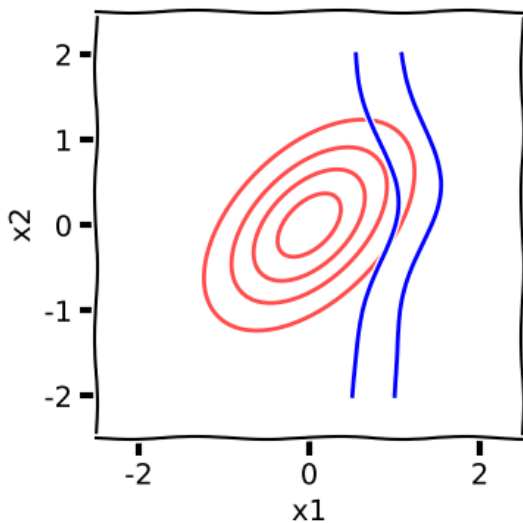
Conditional Gaussians



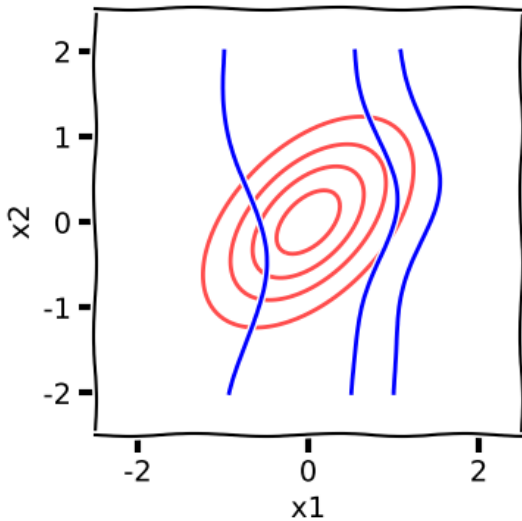
Conditional Gaussians



Conditional Gaussians

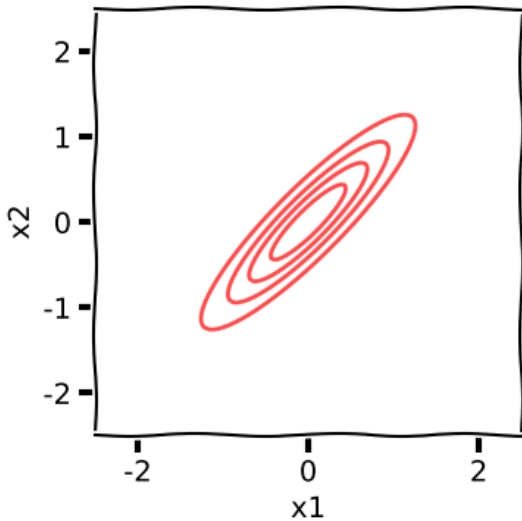


Conditional Gaussians

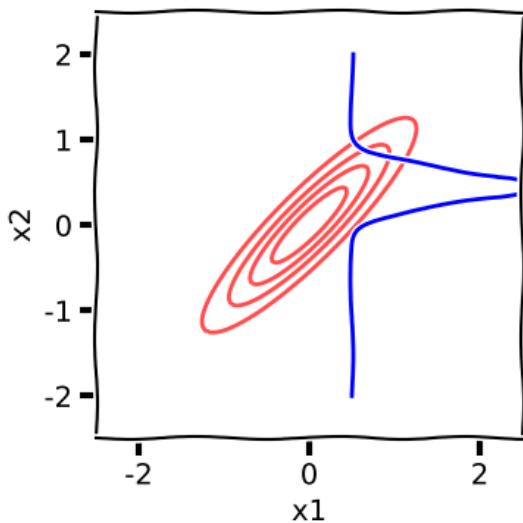


$$\mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}\right)$$

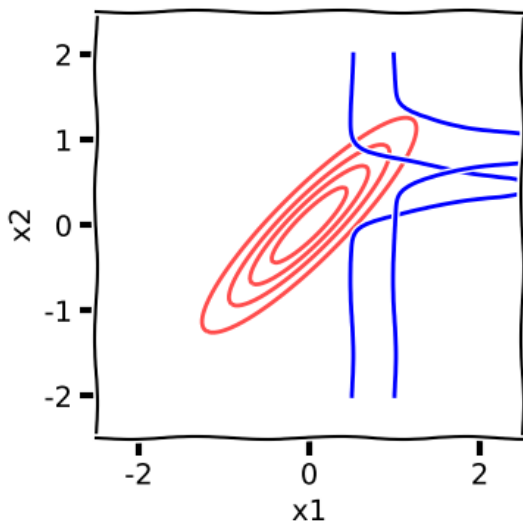
Conditional Gaussians



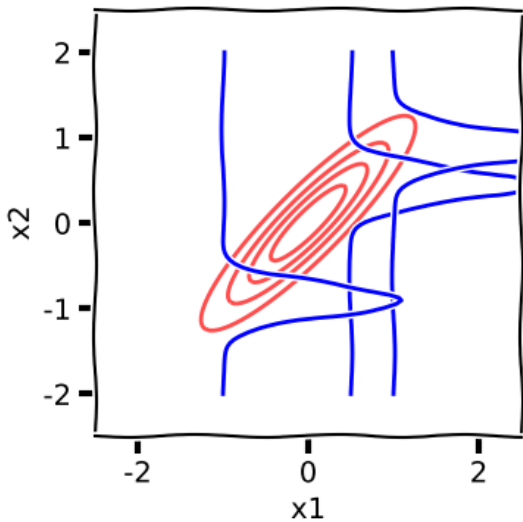
Conditional Gaussians



Conditional Gaussians

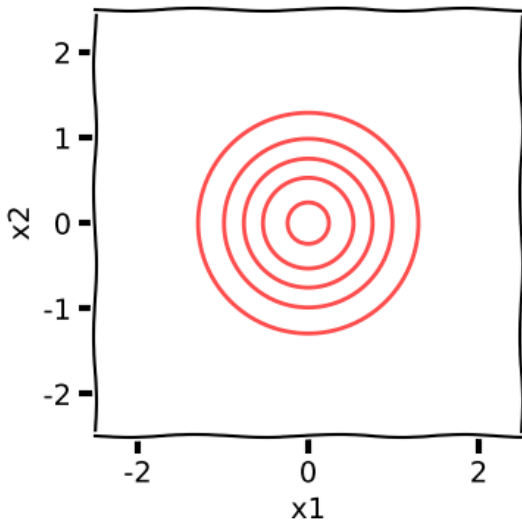


Conditional Gaussians

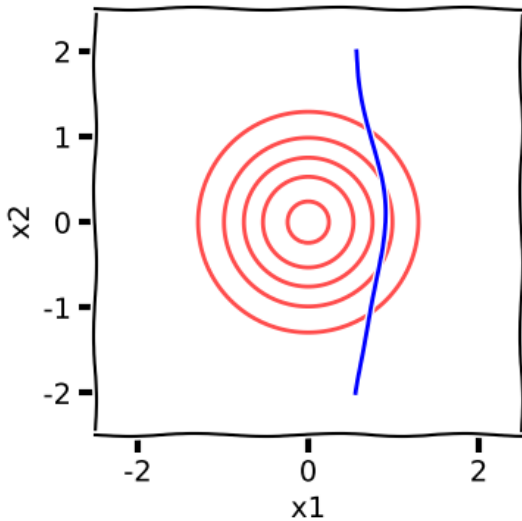


$$\mathcal{N}\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$$

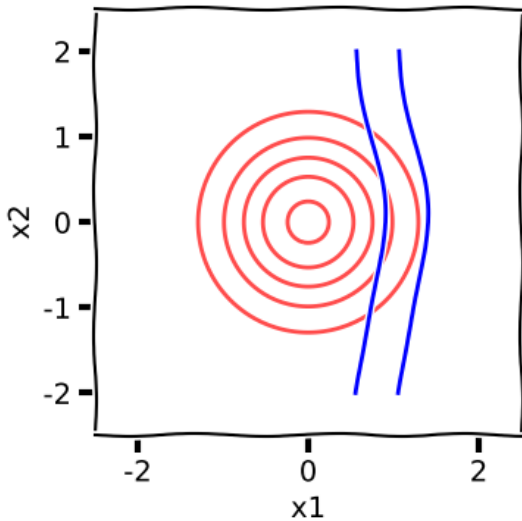
Conditional Gaussians



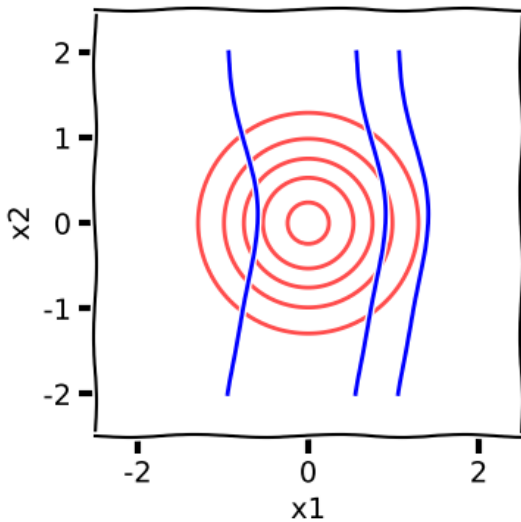
Conditional Gaussians



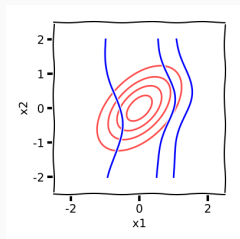
Conditional Gaussians



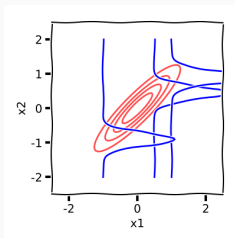
Conditional Gaussians



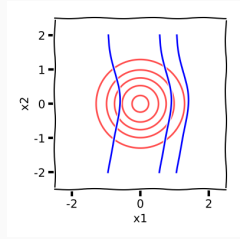
Conditional Gaussians



$$N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.5 \\ 0.5 & 1 \end{bmatrix}\right)$$



$$N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0.9 \\ 0.9 & 1 \end{bmatrix}\right)$$



$$N\left(\begin{bmatrix} 0 \\ 0 \end{bmatrix}, \begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}\right)$$

Uncertainty over functions

- Regression model,

$$\mathbf{y}_i = f(\mathbf{x}_i) + \epsilon$$

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

- Introduce f_i as *instansiation* of function,

$$f_i = f(\mathbf{x}_i),$$

- as a new random variable.

Uncertainty over functions

- Regression model,

$$\begin{aligned}y_i &= f(\mathbf{x}_i) + \epsilon \\ \epsilon &\sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})\end{aligned}$$

- Introduce f_i as *instansiation* of function,

$$f_i = f(\mathbf{x}_i),$$

- as a new random variable.
- now we have a "handle" to specify our assumptions over

Uncertainty over functions

Model,

$$p(\mathbf{y}, \mathbf{f}, \boldsymbol{\theta}) = p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{x}, \boldsymbol{\theta})p(\boldsymbol{\theta})$$

Want to "push" \mathbf{x} through a mapping f of which we are uncertain,

$$p(\mathbf{f}|\mathbf{x}, \boldsymbol{\theta}),$$

prior over instantiations of function.

$$p(\mathbf{f}|\mathbf{x}, \boldsymbol{\theta}) \sim \mathcal{GP}(\mu(\mathbf{x}), k(\mathbf{x}, \mathbf{x}))$$

Definition (Gaussian Process)

A Gaussian Process is an infinite collection of random variables where **any** subset is jointly gaussian. The process is specified by a mean function $\mu(\cdot)$ and a co-variance function $k(\cdot, \cdot)$

Gaussian Marginal

$$p(f_1, f_2, \dots, f_N, \dots | \mathbf{x}, \boldsymbol{\theta}) = \mathcal{GP}(\mu(\mathbf{x}), k(\mathbf{x}, \mathbf{x}))$$

$$= \mathcal{N} \left(\begin{bmatrix} \mu(x_1) \\ \mu(x_2) \\ \vdots \\ \mu(x_N) \\ \vdots \end{bmatrix}, \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \cdots & k(x_1, x_N) & \cdots \\ k(x_2, x_1) & k(x_2, x_2) & \cdots & k(x_2, x_N) & \cdots \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ k(x_N, x_1) & k(x_N, x_2) & \cdots & k(x_N, x_N) & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix} \right)$$

- Remember the Gaussian marginal distribution,

$$p(f_1, f_2) = \mathcal{N} \left(\begin{bmatrix} \mu(x_1) \\ \mu(x_2) \end{bmatrix}, \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) \\ k(x_2, x_1) & k(x_2, x_2) \end{bmatrix} \right)$$

$$p(\mathbf{f}|\mathbf{x}, \boldsymbol{\theta}) \sim \mathcal{GP}(\mu(\mathbf{x}), k(\mathbf{x}, \mathbf{x}))$$

$$\mathbf{y}_i = f_i + \epsilon$$

$$\epsilon \sim \mathcal{N}(\mathbf{0}, \sigma^2 \mathbf{I})$$

$$p(\mathbf{y}|\mathbf{x}, \boldsymbol{\theta}) = \int p(\mathbf{y}|\mathbf{f})p(\mathbf{f}|\mathbf{x}, \boldsymbol{\theta})d\mathbf{f}$$

\mathcal{GP} is infinite, but we only observe finite amount of data. This means conditioning on a subset of the data, the \mathcal{GP} is a just a Gaussian distribution, which is self-conjugate and we know how to do everything

The Mean Function

- Function of only the input location
- What do I expect the function value to be **only** accounting for the input location

The Covariance Function

- Function of **two** input locations
- How should the information from other locations with **known** function value observations effect my estimate

The Mean Function

- Function of only the input location
- What do I expect the function value to be **only** accounting for the input location
- We will assume this to be constant

The Covariance Function

- Function of **two** input locations
- How should the information from other locations with **known** function value observations effect my estimate

The Mean Function

- Function of only the input location
- What do I expect the function value to be **only** accounting for the input location
- We will assume this to be constant

The Covariance Function

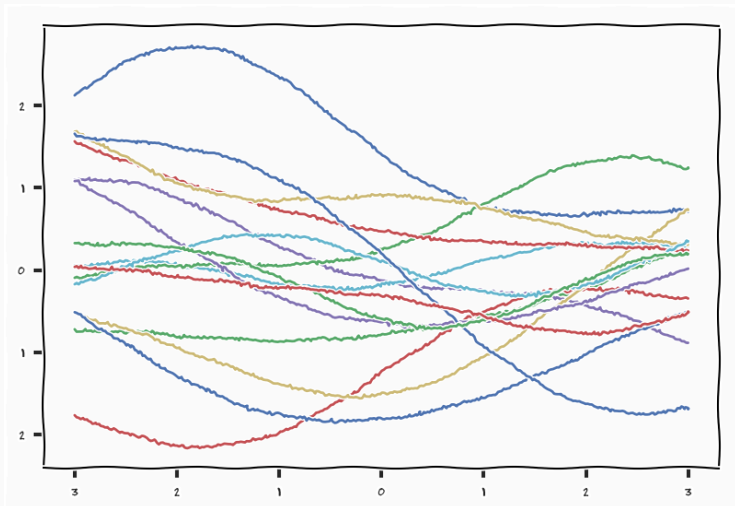
- Function of **two** input locations
- How should the information from other locations with **known** function value observations effect my estimate
- Encodes the behavior of the function

Gaussian Process: Samples

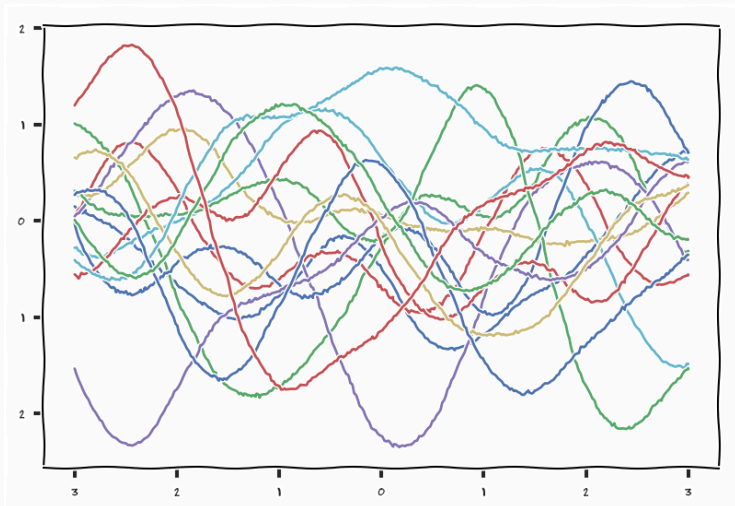
$$p(f_1, f_2, \dots, f_N, \dots | \mathbf{x}, \boldsymbol{\theta}) = \mathcal{GP}(\mu(\mathbf{x}), k(\mathbf{x}, \mathbf{x}))$$

$$= \mathcal{N} \left(\begin{bmatrix} \mu(x_1) \\ \mu(x_2) \\ \vdots \\ \mu(x_N) \\ \vdots \end{bmatrix}, \begin{bmatrix} k(x_1, x_1) & k(x_1, x_2) & \cdots & k(x_1, x_N) & \cdots \\ k(x_2, x_1) & k(x_2, x_2) & \cdots & k(x_2, x_N) & \cdots \\ \vdots & \vdots & \ddots & \vdots & \vdots \\ k(x_N, x_1) & k(x_N, x_2) & \cdots & k(x_N, x_N) & \cdots \\ \vdots & \vdots & \vdots & \vdots & \ddots \end{bmatrix} \right)$$

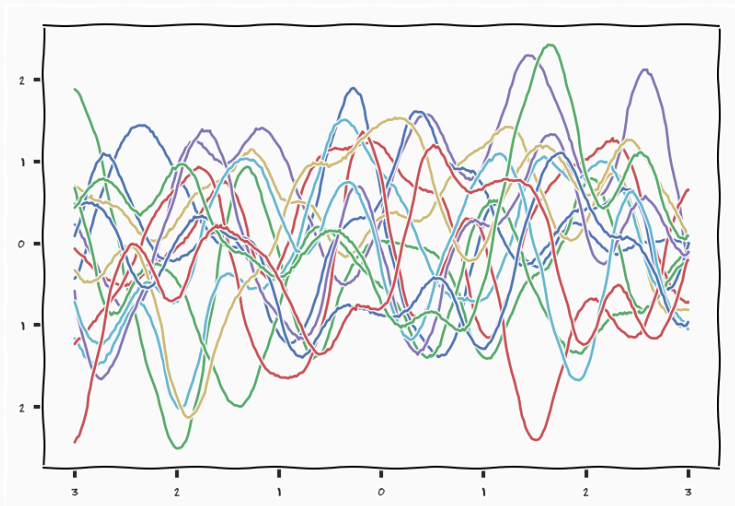
Gaussian Processes: Samples



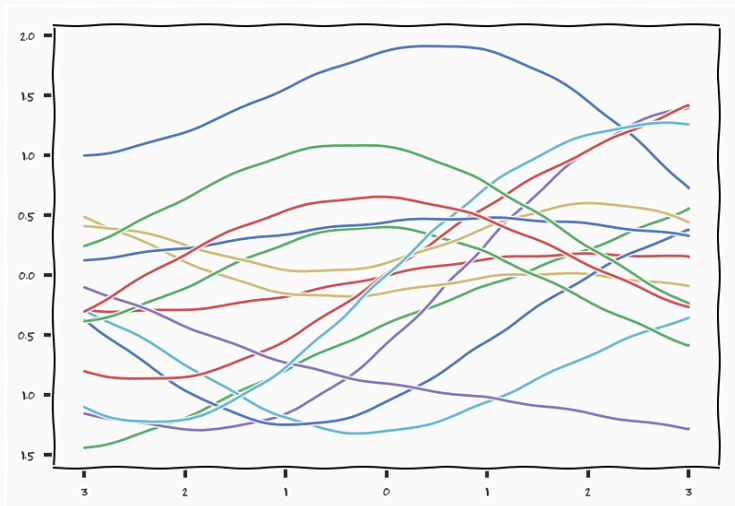
Gaussian Processes: Samples



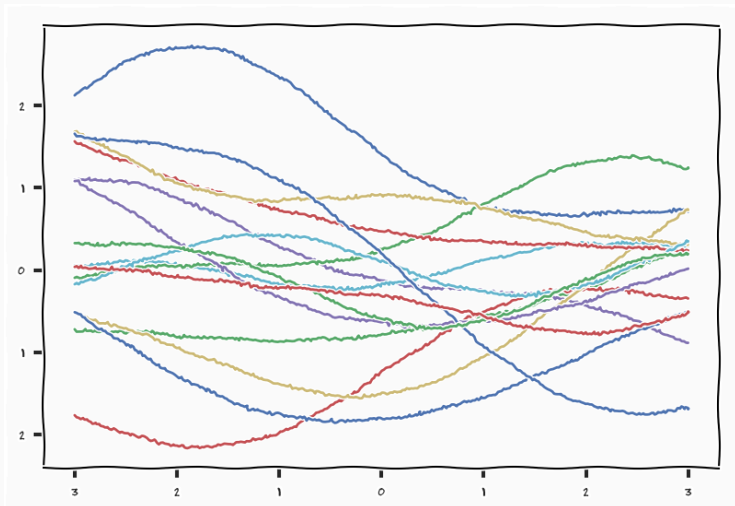
Gaussian Processes: Samples



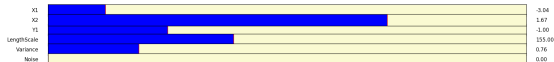
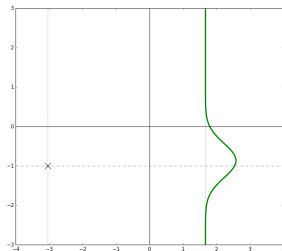
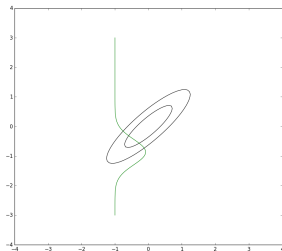
Gaussian Processes: Samples



Gaussian Processes: Samples



Demo



Reset

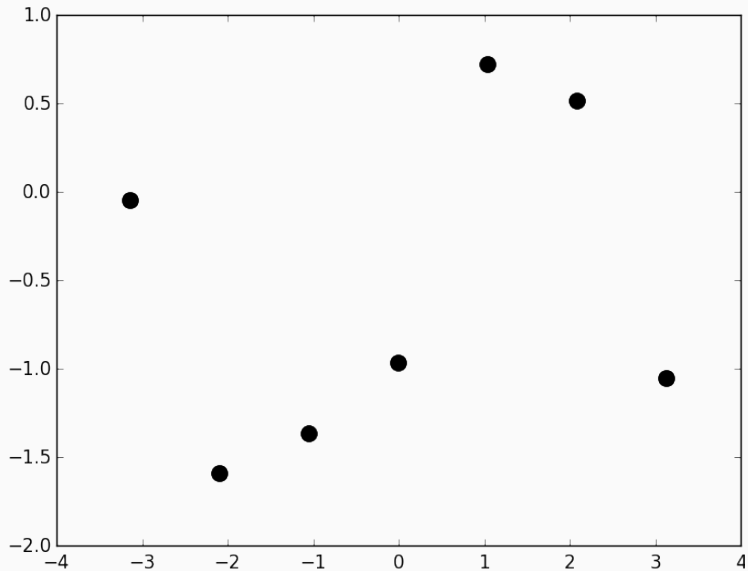
- All instantiations are jointly Gaussian

$$\begin{bmatrix} \mathbf{f} \\ f_* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ 0 \end{bmatrix}, \begin{bmatrix} k(\mathbf{X}, \mathbf{X}) & k(\mathbf{X}, \mathbf{x}_*) \\ k(\mathbf{x}_*, \mathbf{X}) & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right)$$

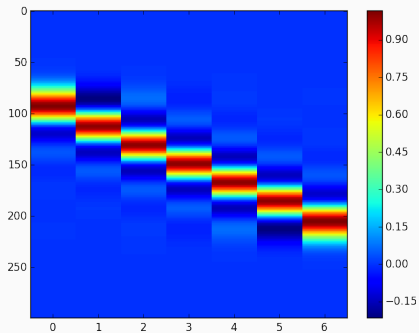
- Conditional Gaussian (same as always)

$$p(f_* | \mathbf{x}_*, \mathbf{X}, \mathbf{f}) = \mathcal{N}(k(\mathbf{x}_*, \mathbf{X})^\top k(\mathbf{X}, \mathbf{X})^{-1} \mathbf{f}, \\ k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X})^\top k(\mathbf{X}, \mathbf{X})^{-1} k(\mathbf{X}, \mathbf{x}_*))$$

Gaussian Process: Posterior

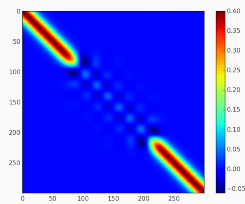
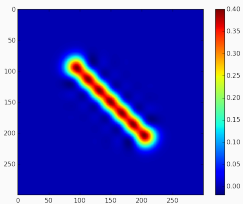
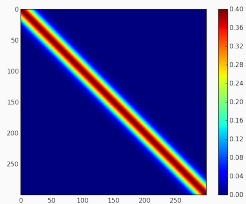


Does it make sense



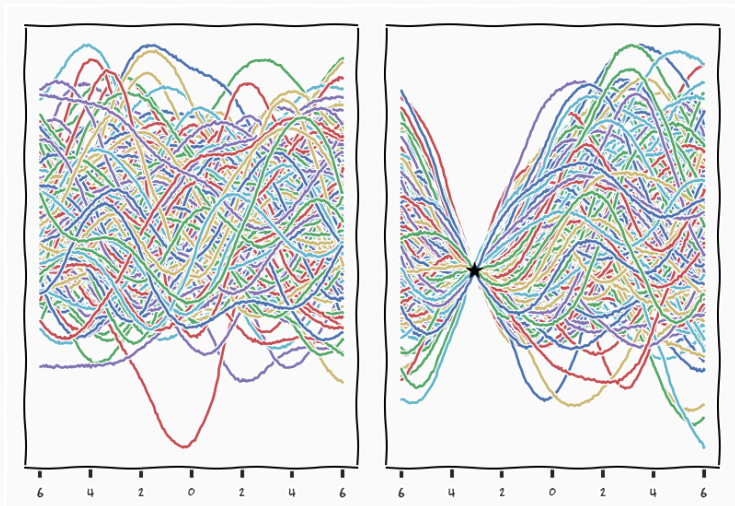
$$k(\mathbf{x}_*, \mathbf{X})^T k(\mathbf{X}, \mathbf{X})^{-1} \mathbf{f}$$

Does it make sense

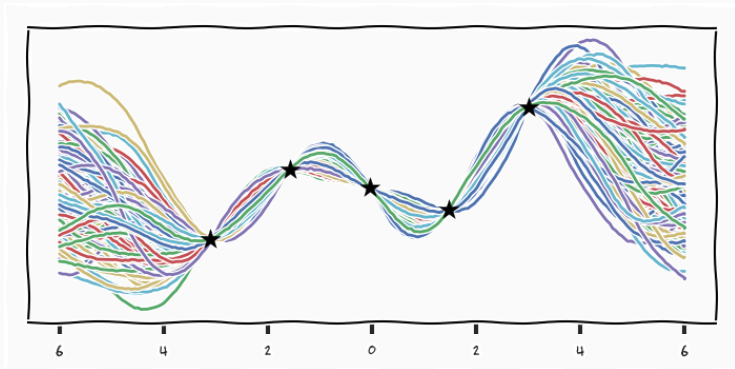


$$k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{X})^T k(\mathbf{X}, \mathbf{X})^{-1} k(\mathbf{X}, \mathbf{x}_*)$$

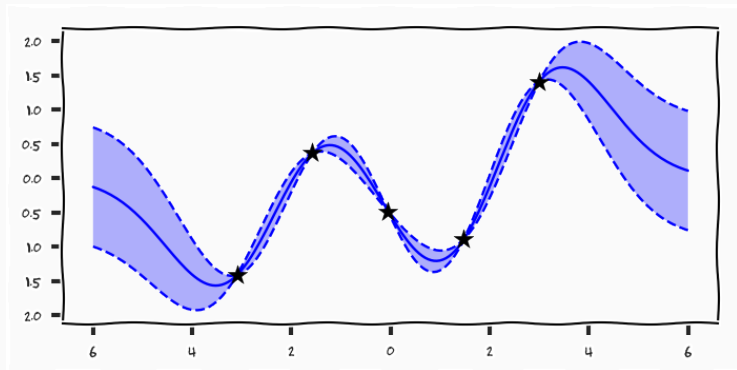
Gaussian Processes: Posterior Samples



Gaussian Processes Posterior



Gaussian Processes Posterior



Gaussian Processes: Noisy observations

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I} & k(\mathbf{X}, \mathbf{x}_*) \\ k(\mathbf{x}_*, \mathbf{X}) & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right)$$

$$p(f_* | \mathbf{x}_*, \mathbf{x}, \mathbf{f}, \boldsymbol{\theta}) = \mathcal{N}(k(\mathbf{x}_*, \mathbf{x})^T (K(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1} \mathbf{f}, \\ k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{x})^T (K(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1} K(\mathbf{x}, \mathbf{x}_*))$$

- Add noise to observations

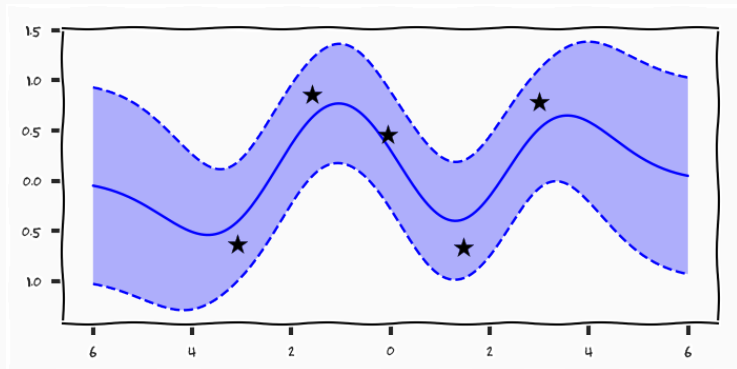
Gaussian Processes: Noisy observations

$$\begin{bmatrix} \mathbf{y} \\ \mathbf{f}_* \end{bmatrix} \sim \mathcal{N} \left(\begin{bmatrix} \mathbf{0} \\ \mathbf{0} \end{bmatrix}, \begin{bmatrix} k(\mathbf{X}, \mathbf{X}) + \sigma^2 \mathbf{I} & k(\mathbf{X}, \mathbf{x}_*) \\ k(\mathbf{x}_*, \mathbf{X}) & k(\mathbf{x}_*, \mathbf{x}_*) \end{bmatrix} \right)$$

$$p(f_* | \mathbf{x}_*, \mathbf{x}, \mathbf{f}, \boldsymbol{\theta}) = \mathcal{N}(k(\mathbf{x}_*, \mathbf{x})^T (K(\mathbf{x}, \mathbf{x} + \sigma^2 \mathbf{I}))^{-1} \mathbf{f}, \\ k(\mathbf{x}_*, \mathbf{x}_*) - k(\mathbf{x}_*, \mathbf{x})^T (K(\mathbf{x}, \mathbf{x}) + \sigma^2 \mathbf{I})^{-1} K(\mathbf{x}, \mathbf{x}_*))$$

- Add noise to observations
- *Do you recognise the mean?*

Gaussian Processes



Question 1-14

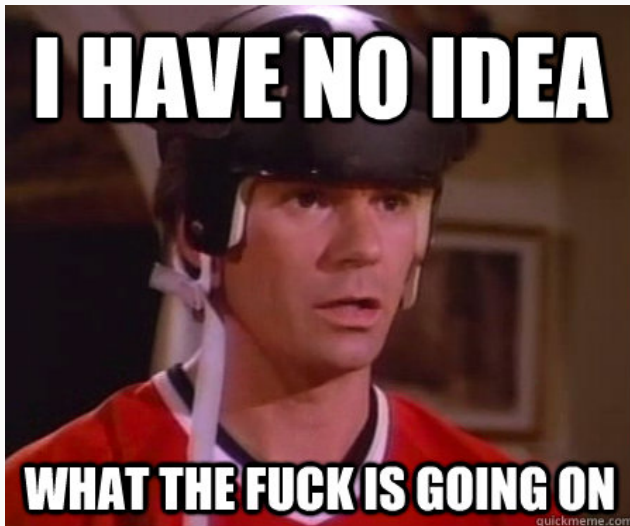
Summary

Summary

- Repeat of the machine learning procedure
 - assumption + data + compute \rightarrow updated assumption
 - don't worry it will become clear eventually
- Gaussian processes
 - infinite generalisation of Gaussian distribution
 - prior over the space of functions
 - contains **all** functions

eof

Part II



I don't think so

$$p(\text{ML}|\text{COMS30007}) = \frac{p(\text{COMS30007}|\text{ML})p(\text{ML})}{p(\text{COMS30007})}$$

- Why is Machine Learning Hard
 - we don't learn how to solve a specific task
 - we learn how to learn how to solve every task
 - its meta knowledge

- Why is Machine Learning Hard
 - we don't learn how to solve a specific task
 - we learn how to learn how to solve every task
 - its meta knowledge
- Why is Machine Learning Easy
 - you have examples of learning from anything
 - its the one thing that humans are good at

$$(\alpha \mathbf{I} + \beta \phi(\mathbf{X})^T \phi(\mathbf{X}))^{-1} \phi(\mathbf{X})^T \mathbf{t}$$

Laplace Demon [2]

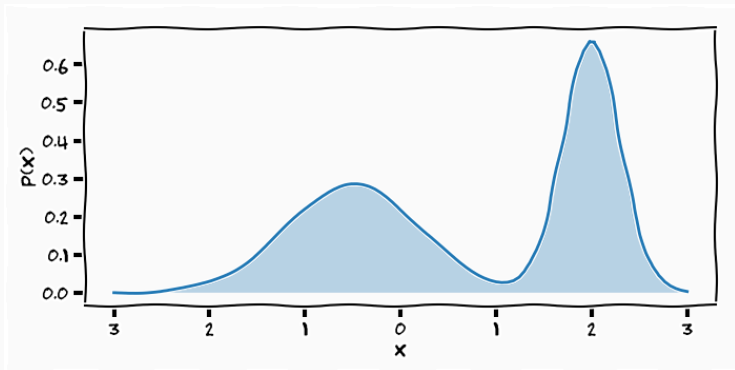


Laplace's Demon [2]

An intelligence which at a given instant knew all the forces acting in nature and the position of every object in the universe - if endowed with a brain sufficiently vast to make all necessary calculations - could describe with a single formula the motions of the largest astronomical bodies and those of the smallest atoms. To such an intelligence, nothing would be uncertain; the future, like the past, would be an open book.

All these efforts in the search for truth tend to lead the mind continuously towards the intelligence we have just mentioned, although it will always remain infinitely distant from this intelligence.

Uncertainty



Variables

```
def f(x):  
    if x == 1:  
        return 2  
    else:  
        return 1  
  
x = 1  
print(f(x))
```

Definition (Variable)

In elementary mathematics, a variable is an alphabetic character representing a number, called the value of the variable, which is either arbitrary, not fully specified, or unknown.

Random Variables

```
import numpy as np
```

```
def f(x):  
    if x > 3.0:  
        return 2  
    else:  
        return 1
```

```
x = np.random.normal(10.0,2.0,1)
```

Definition (Random Variable)

In probability and statistics, a random variable, random quantity, aleatory variable, or stochastic variable is a variable whose possible values are numerical outcomes of a random phenomenon.²

²https://en.wikipedia.org/wiki/Random_variable

Random Variables

```
import numpy as np
```

```
def f(x):  
    if x > 3.0:  
        return 2  
    else:  
        return 1
```

```
x = np.random.normal(10.0,2.0,1)
```

$$p(x) = \mathcal{N}(x|10.0, 2.0)$$

²https://en.wikipedia.org/wiki/Random_variable

Random Variables

```
import numpy as np

def f(x):
    return np.random.normal(x,1.0,1)

x = np.random.normal(10.0,2.0,1)
```

- x is random

$$p(x) = \mathcal{N}(x|10.0, 2.0)$$

- f is random

$$p(f|x) = \mathcal{N}(f|x, 1.0)$$

"I understand that the likelihood function is a normal distribution with mean Wx and variance $\hat{\Sigma}^{-1}$, but where do we actually find the W values to calculate it? Do we sample all of the W combinations from our prior distribution and check to see if lines close to the point would be drawn? Do we choose them from our prior based on which ones are more likely? Do we use calculus to estimate them?"

"- How the likelihood function is derived for a single point, and how it is iteratively updated during regression"

"- How the mean and covariance of the posterior are iteratively calculated during regression, making the link between them, the basis functions and the design matrix"

Conditional Distributions

$$p(\mathbf{w}|m_0, S_0) = \mathcal{N}(\mathbf{w}|m_0, S_0) = \frac{1}{((2\pi)^D |S_0|)^{\frac{1}{2}}} e^{\frac{1}{2}(\mathbf{w}-\mathbf{m}_0)^T S^{-1}(\mathbf{w}-\mathbf{m}_0)}$$

- The above is a function of \mathbf{w}
- The function has "parameters" m_0 and S_0
- In order to evaluate the function the parameters needs to be set

Interpreting Bayesian Probabilities³

"Our prior is our assumption - when we say updated assumption, is this correct to say this is our posterior? Except this updated assumption isn't really an assumption, because it is a function over y and we therefore can't use it as a prior in another equation."

³reddit URL

Interpreting Bayesian Probabilities³

"When we state m_n and s_n as the parameters of our posterior, does this mean 'after we've multiplied the prior with n amounts of likelihoods from n data points'?"

Interpreting Bayesian Probabilities³

"Given the marginal likelihood can be interpreted as the distribution/probability of observing our data given our model, does this mean that the predicted posterior is simply our marginal likelihood? Or is the predicted posterior simply an entirely different model (using our previous results from the model learning w)?"

"Suggestion: go back over the linear regression lecture, but about 5x slower, with a view to making it a little more obvious how to tackle Q12 in the coursework"

"Do we need to really know the derivation for this? And when it says that in this form the prediction is now made using the training set.. does it use both the x and t values from the training set or just the target value ' t '?"

⁴reddit URL

$$p(t_*|\mathbf{t}, \mathbf{x}_*, \mathbf{X}, \alpha, \beta) = \int p(t_*|\mathbf{x}_*, \mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \alpha, \beta) d\mathbf{w}$$

- Likelihood

$$p(t_*|\mathbf{x}_*, \mathbf{w}, \beta) = \mathcal{N}(t_*|\mathbf{w}^T \phi(\mathbf{x}_*), \beta^{-1}) = f_1(t_*)$$

- Posterior

$$\begin{aligned} p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \alpha, \beta) &= \mathcal{N}(\mathbf{w}|\beta (\alpha \mathbf{I} + \beta \phi(\mathbf{X})^T \phi(\mathbf{X}))^{-1} \phi(\mathbf{X})^T \mathbf{t}, \\ &= (\alpha \mathbf{I} + \beta \phi(\mathbf{X})^T \phi(\mathbf{X}))^{-1}) = f_2(\mathbf{w}) \end{aligned}$$

⁵reddit URL

- Likelihood

$$p(t_*|\mathbf{x}_*, \mathbf{w}, \beta) = \mathcal{N}(t_*|\mathbf{w}^T\phi(\mathbf{x}_*), \beta^{-1}) = f_1(t_*)$$

- Posterior

$$\begin{aligned} p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \alpha, \beta) &= \mathcal{N}(\mathbf{w}|\beta(\alpha\mathbf{I} + \beta\phi(\mathbf{X})^T\phi(\mathbf{X}))^{-1}\phi(\mathbf{X})^T\mathbf{t}, \\ &\quad (\alpha\mathbf{I} + \beta\phi(\mathbf{X})^T\phi(\mathbf{X}))^{-1}) = f_2(\mathbf{w}) \end{aligned}$$

$$\begin{aligned} p(t_*|\mathbf{t}, \mathbf{x}_*, \mathbf{X}, \alpha, \beta) &= \int p(t_*|\mathbf{x}_*, \mathbf{w}, \beta)p(\mathbf{w}|\mathbf{t}, \mathbf{X}, \alpha, \beta)d\mathbf{w} \\ &= \int f_1(t_*)f_2(\mathbf{w})d\mathbf{w} = \int g(t_*, \mathbf{w})d\mathbf{w} = h(t_*) \end{aligned}$$

⁵reddit URL

I don't think so

$$p(\text{ML}|\text{COMS30007}) = \frac{p(\text{COMS30007}|\text{ML})p(\text{ML})}{p(\text{COMS30007})}$$

References



Christopher M. Bishop.

Pattern Recognition and Machine Learning (Information Science and Statistics).

Springer-Verlag New York, Inc., Secaucus, NJ, USA, 2006.



Pierre Simon Laplace.

A philosophical essay on probabilities, 1814.