10-301/601: Introduction to Machine Learning Lecture 27 – Gaussian Processes

Matt Gormley & Henry Chai 4/23/25

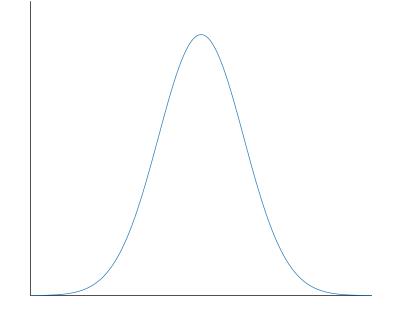
Front Matter

- Announcements
 - HW9 released 4/17, due 4/24 (Thursday) at 11:59 PM
 - You may only use at most 2 late days on HW9
 - Exam 3 on 5/1 from 1 PM to 3 PM
 - We will not use the full 3-hour window
 - All topics from Lectures 17 to 25 (inclusive) are inscope, excluding the MLE/MAP portion of Lecture 17
 - Exam 1 and 2 content may be referenced but will not be the primary focus of any question
 - Please watch Piazza carefully for your room and seat assignments
 - You are allowed to bring one letter-size sheet of notes; you may put whatever you want on both sides

Gaussians

(Univariate) Gaussians:

$$x \sim \mathcal{N}(x; \mu = 0, \sigma^2 = 1)$$

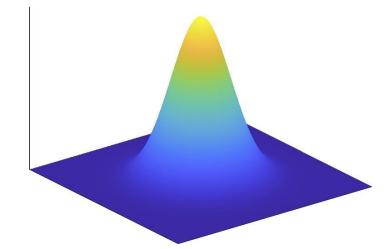


Multivariate Gaussians:

$$\mathbf{x} = [x_1, ..., x_D]^T \mathcal{R}^D$$

$$\sim \mathcal{N}(\mathbf{x}; \boldsymbol{\mu} = \mathbf{0}_D, \boldsymbol{\Sigma} = I_D)$$

$$\mathcal{R}^D \mathcal{R}^D$$
Covariance matrix



Some fun facts about Gaussians

Closure under linear transformations:

Tf
$$\propto N(\overline{p}, Z)$$

 $+ \overline{b} \sim N(A\overline{p} + \overline{b})$ $+ ZA^{T}$

• Closure under addition follows

If
$$\overrightarrow{x} \sim N(\overrightarrow{p}, \overrightarrow{z}) \rightarrow \overrightarrow{y} \sim N(\overrightarrow{m}, S)$$

then $\overrightarrow{x} + \overrightarrow{y} \sim N(\overrightarrow{p} + \overrightarrow{m}, z + S)$

Closure under conditioning:

If
$$\begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \sim N(\begin{bmatrix} N_1 \\ N_2 \end{bmatrix}, \begin{bmatrix} Z_1 & Z_{12} \\ Z_{21} & Z_{22} \end{bmatrix})$$

 $+ \sum_{i=1}^{N} |x_i| \leq C \sim N(y_1 + Z_{12} \sum_{i=2}^{N} (C - y_2),$
 $- \sum_{i=1}^{N} \sum_{i=2}^{N} |Z_{21}|$

Some old friends

Gaussian process =

Bayesian linear regression + Kernels

Some old friends

Gaussian process =

Bayesian linear regression + Kernels

Recall: MAP for Linear Regression

If we assume a linear model with additive Gaussian noise

$$y = Xw + \epsilon$$
 where $\epsilon \sim N(\mathbf{0}_N, \sigma^2 I_N) \rightarrow y \sim N(Xw, \sigma^2 I_N)$

and independent identical Gaussian priors on the weights...

$$\mathbf{w} \sim N\left(\mathbf{w}_{D+1}, \frac{\sigma^2}{\lambda} I_{D+1}\right) \rightarrow \mathbf{p}(\mathbf{w}) \propto \exp\left(-\frac{1}{2\sigma^2}(\lambda \mathbf{w}^T \mathbf{w})\right)$$

• ... then, the MAP of **w** is the ridge regression solution!

$$\mathbf{w}_{MAP} = \underset{\mathbf{w}}{\operatorname{argmin}} (X\mathbf{w} - \mathbf{y})^{T} (X\mathbf{w} - \mathbf{y}) + \lambda \mathbf{w}^{T} \mathbf{w}$$
$$= (X^{T}X + \lambda I_{D+1})^{-1} X^{T} \mathbf{y}$$

 Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$y = Xw + \epsilon$$
 where $\epsilon \sim N(\mathbf{0}_N, \sigma^2 I_N)$ and $w \sim N(\mathbf{0}_{D+1}, \Sigma)$

then,

$$\mathbf{y} \sim N(X\mathbf{0}_{D+1} + \mathbf{0}_N = \mathbf{0}_N, X\Sigma X^T + \sigma^2 I_N)$$

 Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$y = Xw + \epsilon \text{ where } \epsilon \sim N(\mathbf{0}_{N}, \sigma^{2}I_{N}) \text{ and } w \sim N(\mathbf{0}_{D+1}, \Sigma)$$
then,
$$\begin{bmatrix} w \\ y \end{bmatrix} \sim N\left(\begin{bmatrix} \mathbf{0}_{D+1} \\ \mathbf{0}_{N} \end{bmatrix}, \begin{bmatrix} \Sigma \\ ???? \\ X\Sigma X^{T} + \sigma^{2}I_{N} \end{bmatrix}\right)$$

$$= \chi \text{ cov}\left(\overrightarrow{u}, \overrightarrow{\chi}\overrightarrow{u}\right)$$

$$= \chi \text{ cov}\left(\overrightarrow{u}, \overrightarrow{u}\right)$$

 Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$y=Xw+\epsilon \text{ where }\epsilon \sim N(\mathbf{0}_N,\sigma^2I_N) \text{ and } w\sim N(\mathbf{0}_{D+1},\Sigma)$$
 then,

$$\begin{bmatrix} \boldsymbol{w} \\ \boldsymbol{y} \end{bmatrix} \sim N \begin{pmatrix} \begin{bmatrix} \mathbf{0}_{D+1} \\ \mathbf{0}_N \end{bmatrix}, \begin{bmatrix} \Sigma & X\Sigma \\ X\Sigma & X\Sigma X^T + \sigma^2 I_N \end{bmatrix} \end{pmatrix}$$

 Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$y = Xw + \epsilon$$
 where $\epsilon \sim N(\mathbf{0}_N, \sigma^2 I_N)$ and $w \sim N(\mathbf{0}_{D+1}, \Sigma)$

then,

$$\boldsymbol{w} \mid \boldsymbol{y} \sim N(\boldsymbol{\mu}_{POST}, \boldsymbol{\Sigma}_{POST})$$

where

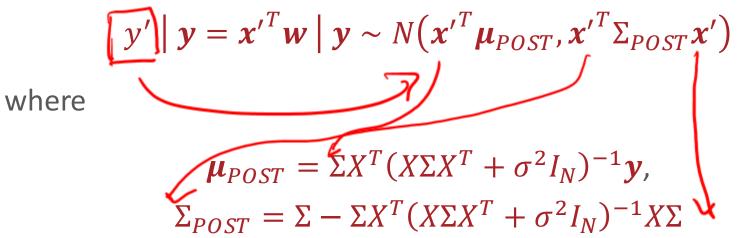
$$\boldsymbol{\mu}_{POST} = \Sigma X^{T} (X \Sigma X^{T} + \sigma^{2} I_{N})^{-1} \boldsymbol{y},$$

$$\Sigma_{POST} = \Sigma - \Sigma X^{T} (X \Sigma X^{T} + \sigma^{2} I_{N})^{-1} X \Sigma$$

 Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$y = Xw + \epsilon$$
 where $\epsilon \sim N(\mathbf{0}_N, \sigma^2 I_N)$ and $w \sim N(\mathbf{0}_{D+1}, \Sigma)$

then given a new test data point x', the prediction is



$$\vec{y} = \begin{bmatrix} y(1) \\ y(2) \\ \vdots \\ y(N) \end{bmatrix}$$

$$\chi = \begin{bmatrix} 1 & \frac{1}{x^2(2)T} \\ 1 & \frac{1}{x^2(2)T} \\ \vdots \\ 1 & \frac{1}{x^2(N)T} \end{bmatrix}$$

 Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$\overrightarrow{y} = Xw + \epsilon$$
 where $\epsilon \sim N(\mathbf{0}_N, \sigma^2 I_N)$ and $w \sim N(\mathbf{0}_{D+1}, \Sigma)$

then given a new test data point x', the prediction is

$$y' \mid y = x'^T w \mid y \sim N(\mu_{PRED}, \Sigma_{PRED})$$

where

$$\mu_{PRED} = \mathbf{x'}^T \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} \mathbf{y},$$

$$\Sigma_{PRED} = \mathbf{x'}^T \Sigma \mathbf{x'} - \mathbf{x'}^T \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} X \Sigma \mathbf{x'}$$

Some old friends

Gaussian process =

Bayesian linear regression + Kernels

Some new friends

Gaussian process =

Bayesian linear regression + Kernels

 Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$y = Xw + \epsilon$$
 where $\epsilon \sim N(\mathbf{0}_N, \sigma^2 I_N)$ and $w \sim N(\mathbf{0}_{D+1}, \Sigma)$

then given a new test data point x', the prediction is

$$y' \mid y = x'^T w \mid y \sim N(\mu_{PRED}, \Sigma_{PRED})$$

where

$$\mu_{PRED} = \mathbf{x'}^T \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} \mathbf{y},$$

$$\Sigma_{PRED} = \mathbf{x'}^T \Sigma \mathbf{x'} - \mathbf{x'}^T \Sigma X^T (X \Sigma X^T + \sigma^2 I_N)^{-1} X \Sigma \mathbf{x'}$$

Bayesian Non-linear Regression...

$$\Phi = \begin{bmatrix} 1 & \phi(x^{(1)})^T \\ 1 & \phi(x^{(2)})^T \\ \vdots & \vdots \\ 1 & \phi(x^{(N)})^T \end{bmatrix}$$

 Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$y = \Phi \omega + \epsilon$$
 where $\epsilon \sim N(\mathbf{0}_N, \sigma^2 I_N)$ and $\omega \sim N(\mathbf{0}_{D'+1}, \Sigma)$

then given a new test data point x', the prediction is

$$y' \mid y = \phi(x')^T \omega \mid y \sim N(\mu_{PRED}, \Sigma_{PRED})$$

where

$$\mu_{PRED} = \phi(\mathbf{x}')^T \Sigma \Phi^T (\Phi \Sigma \Phi^T + \sigma^2 I_N)^{-1} \mathbf{y},$$

$$\Sigma_{PRED}$$

$$= \phi(\mathbf{x}')^T \Sigma \phi(\mathbf{x}') - \phi(\mathbf{x}')^T \Sigma \Phi^T (\Phi \Sigma \Phi^T + \sigma^2 I_N)^{-1} \Phi \Sigma \phi(\mathbf{x}')$$

Bayesian Non-linear Regression can be "kernelized"

$$\Phi = \begin{bmatrix} 1 & \phi(\boldsymbol{x}^{(1)})^T \\ 1 & \phi(\boldsymbol{x}^{(2)})^T \\ \vdots & \vdots \\ 1 & \phi(\boldsymbol{x}^{(N)})^T \end{bmatrix}$$

 Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$y = \Phi \omega + \epsilon$$
 where $\epsilon \sim N(\mathbf{0}_N, \sigma^2 I_N)$ and $\omega \sim N(\mathbf{0}_{D'+1}, \Sigma)$

then given a new test data point x', the prediction is

$$y' \mid y = \phi(x')^T \boldsymbol{\omega} \mid y \sim N(\boldsymbol{\mu}_{PRED}, \Sigma_{PRED})$$

where

$$\mu_{PRED} = \phi(\mathbf{x}')^T \Sigma \Phi^T (\Phi \Sigma \Phi^T + \sigma^2 I_N)^{-1} \mathbf{y},$$

 Σ_{PRED}

$$= \phi(\mathbf{x}')^T \Sigma \phi(\mathbf{x}') - \phi(\mathbf{x}')^T \Sigma \Phi^T (\Phi \Sigma \Phi^T + \sigma^2 I_N)^{-1} \Phi \Sigma \phi(\mathbf{x}')$$

• Define a **kernel function** as

$$K(X,X) = \Phi \Sigma \Phi^T$$

$$K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \Sigma \phi(\mathbf{x}')$$

Bayesian Linear Regression can be kernelized!

 Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$y = \Phi \omega + \epsilon$$
 where $\epsilon \sim N(\mathbf{0}_N, \sigma^2 I_N)$ and $\omega \sim N(\mathbf{0}_{D'+1}, \Sigma)$

then given a new test data point x', the prediction is

$$y' \mid y = \phi(x')^T \boldsymbol{\omega} \mid y \sim N(\boldsymbol{\mu}_{PRED}, \Sigma_{PRED})$$

where

$$\mu_{PRED} = K(x', X)(K(X, X) + \sigma^2 I_N)^{-1} y,$$

$$\Sigma_{PRED} = K(\mathbf{x}', \mathbf{x}') - K(\mathbf{x}', X)(K(X, X) + \sigma^2 I_N) \mathcal{E} K(X, \mathbf{x})$$

Define the kernel function to be

$$K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \Sigma \phi(\mathbf{x}')$$

Wait, what happened to the weights?

 Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$y = \Phi \omega + \epsilon$$
 where $\epsilon \sim N(\mathbf{0}_N, \sigma^2 I_N)$ and $\omega \sim N(\mathbf{0}_{D'+1}, \Sigma)$

then given a new test data point x', the prediction is

$$y' \mid y = \phi(x')^T \boldsymbol{\omega} \mid y \sim N(\boldsymbol{\mu}_{PRED}, \Sigma_{PRED})$$

where

$$\boldsymbol{\mu}_{PRED} = K(\boldsymbol{x}', X)(K(X, X) + \sigma^2 I_N)^{-1} \boldsymbol{y},$$

$$\Sigma_{PRED} = K(\mathbf{x}', \mathbf{x}') - K(\mathbf{x}', X)(K(X, X) + \sigma^2 I_N)^{-1}K(X, \mathbf{x})$$

Define the kernel function to be

$$K(\mathbf{x}, \mathbf{x}') = \phi(\mathbf{x})^T \Sigma \phi(\mathbf{x}')$$

Some old friends

Gaussian process =

Bayesian linear regression + Kernels

A new perspective

Gaussian process =

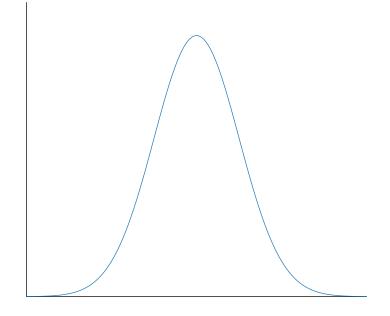
The extension of a Gaussian

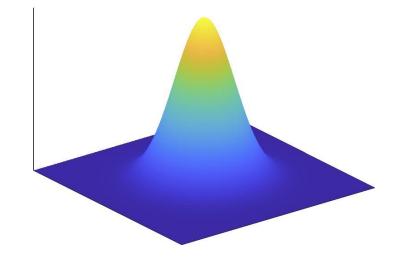
distribution to functions

Gaussians

• (Univariate) Gaussians:

$$x \sim \mathcal{N}(x; \mu = 0, \sigma^2 = 1)$$





Multivariate Gaussians:

$$\mathbf{x} = [x_1, ..., x_D]^T$$
$$\sim \mathcal{N}(\mathbf{x}; \boldsymbol{\mu} = \mathbf{0}_D, \boldsymbol{\Sigma} = I_D)$$

$$f: \mathbb{R}^p \mapsto \mathbb{R} \sim \mathcal{GP}(f; \mu(x))$$
, $\Sigma(x, x')$

Gaussian Process (GP)

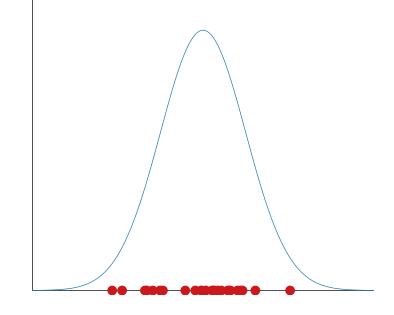
X

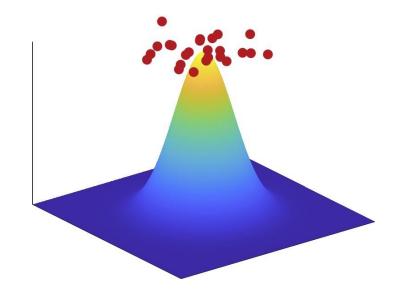
$$f \sim \mathcal{GP}(\mu, \Sigma) \rightarrow f(x) \sim \mathcal{N}(\mu(x), \Sigma(x, x))$$

Gaussians

• (Univariate) Gaussians:

$$x \sim \mathcal{N}(x; \mu = 0, \sigma^2 = 1)$$

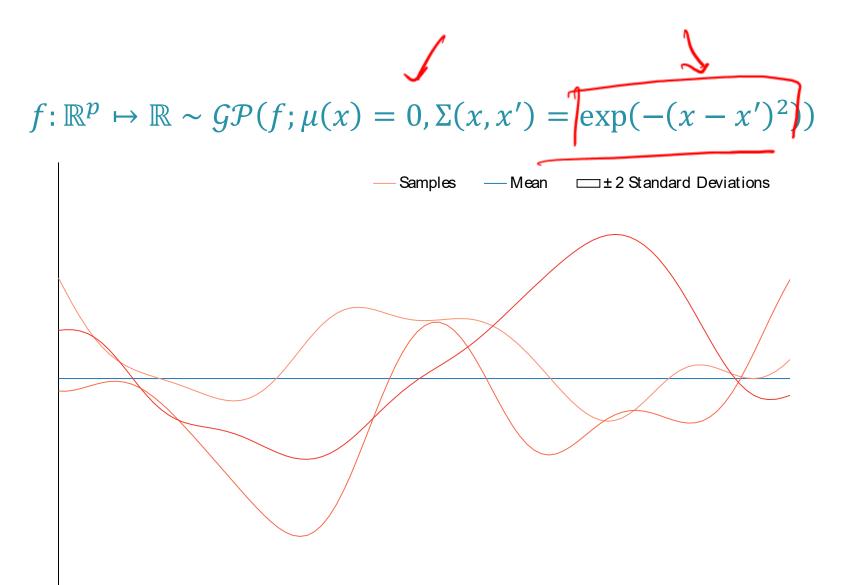




Multivariate Gaussians:

$$\mathbf{x} = [x_1, ..., x_D]^T$$
$$\sim \mathcal{N}(\mathbf{x}; \boldsymbol{\mu} = \mathbf{0}_D, \boldsymbol{\Sigma} = I_D)$$

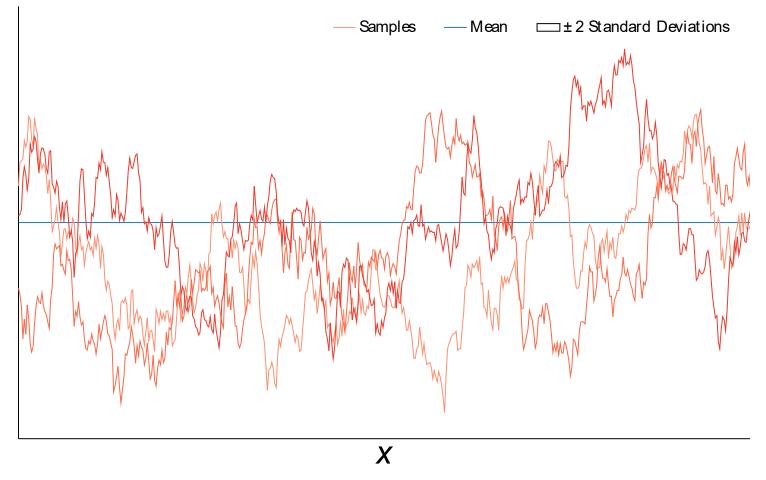
Gaussian Process (GP)



$$f \sim \mathcal{GP}(\mu, \Sigma) \rightarrow f(x) \sim \mathcal{N}(\mu(x), \Sigma(x, x))$$

Gaussian Process (GP)

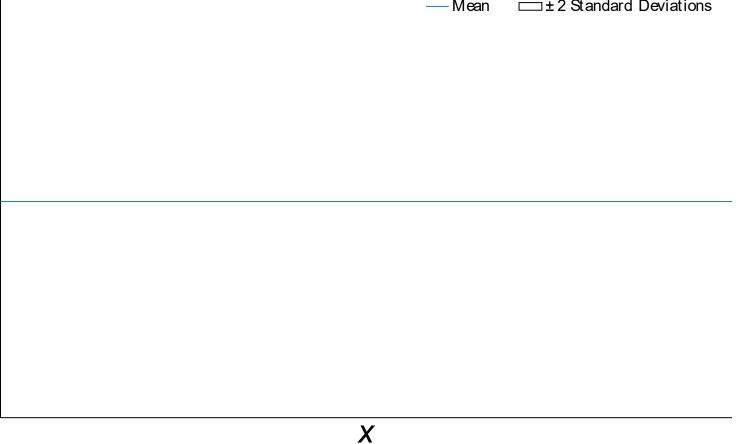
$$f: \mathbb{R}^p \mapsto \mathbb{R} \sim \mathcal{GP}(f; \mu(x) = 0, \Sigma(x, x') = \exp(-|x - x'|))$$



$$f \sim \mathcal{GP}(\mu, \Sigma) \rightarrow f(x) \sim \mathcal{N}(\mu(x), \Sigma(x, x))$$

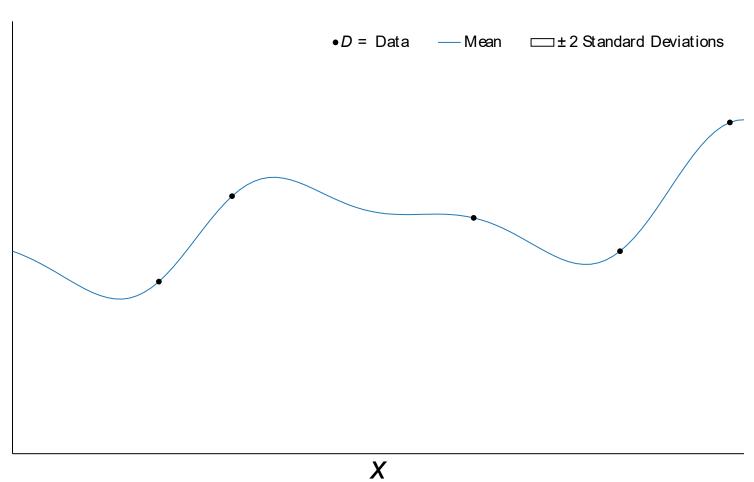
$f: \mathbb{R}^p \mapsto \mathbb{R} \sim \mathcal{GP}(f; \mu(x) = 0, \Sigma(x, x') = \exp(-(x - x')^2))$

GP Prior



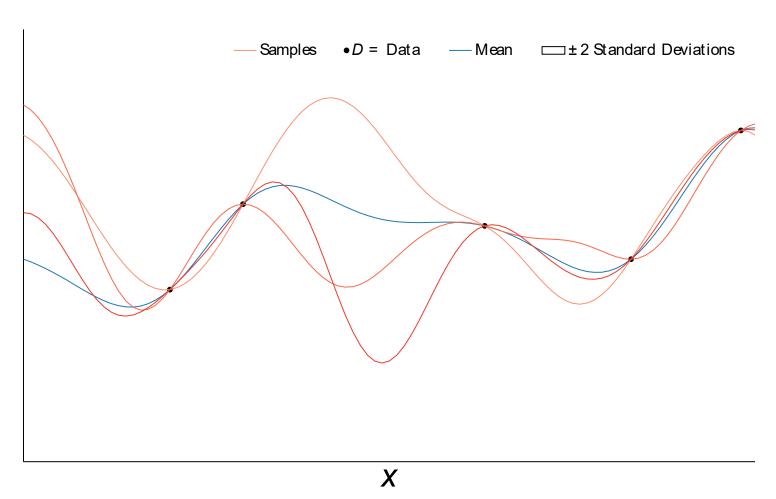
GP Posterior





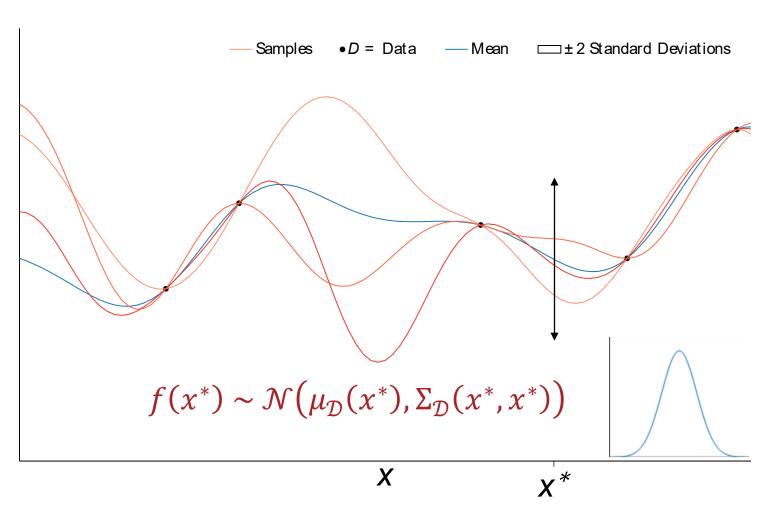
GP Posterior

$f \mid \mathcal{D} \sim \mathcal{GP}(f; \mu_{\mathcal{D}}, \Sigma_{\mathcal{D}})$

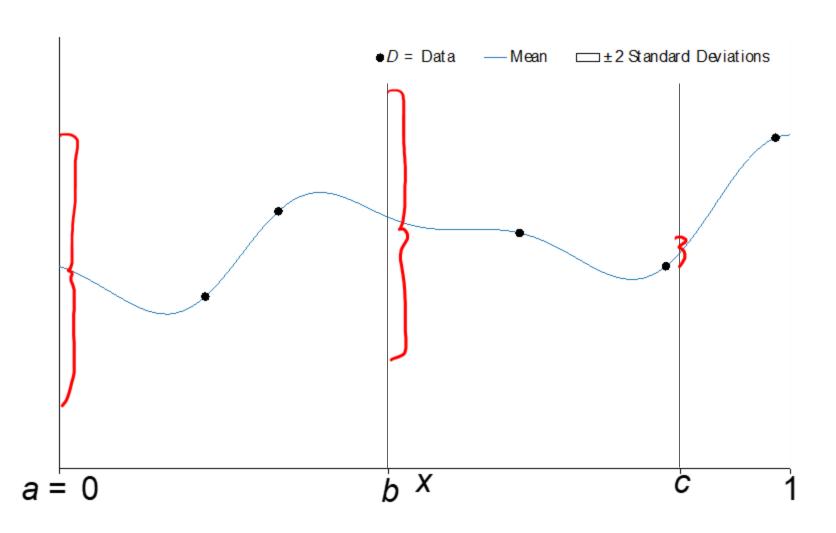


GP Posterior

$f \mid \mathcal{D} \sim \mathcal{GP}(f; \mu_{\mathcal{D}}, \Sigma_{\mathcal{D}})$

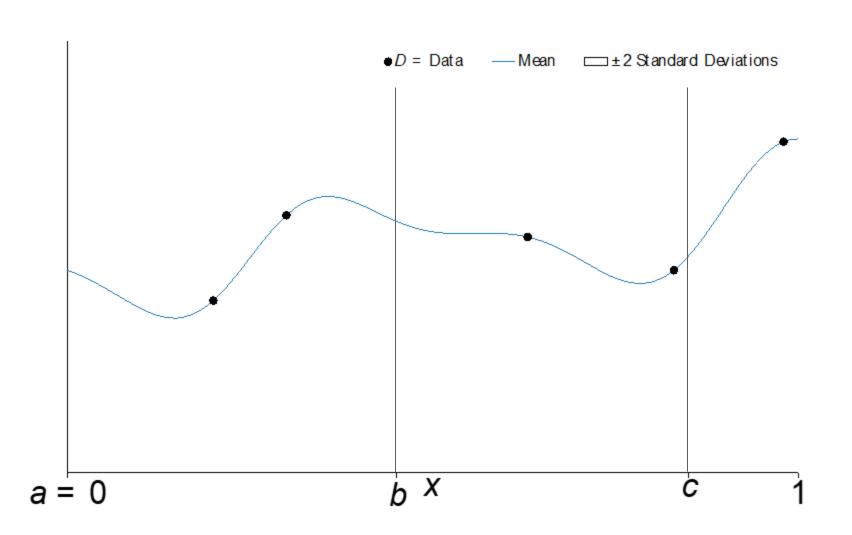


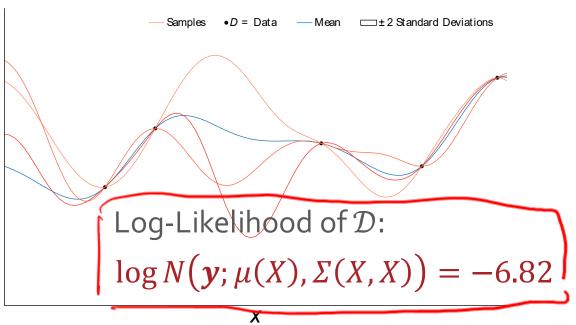
Active Learning



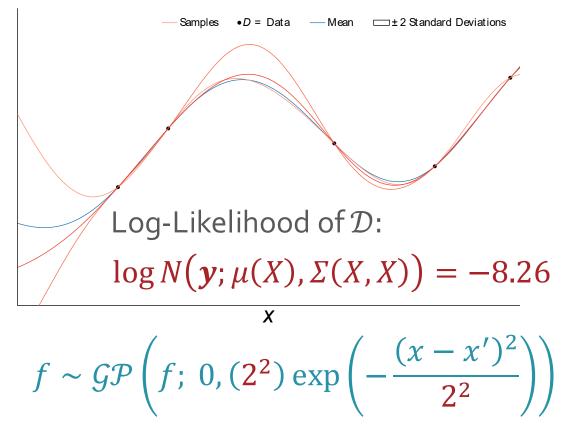
Suppose you can add one data point to your training dataset.

Which value of *x* would you add and why?



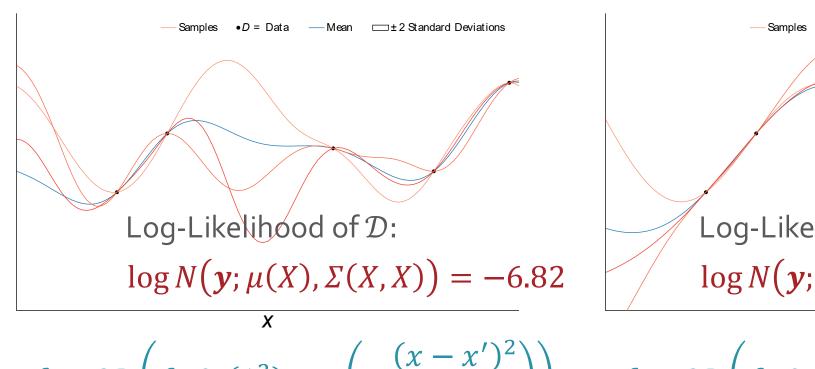


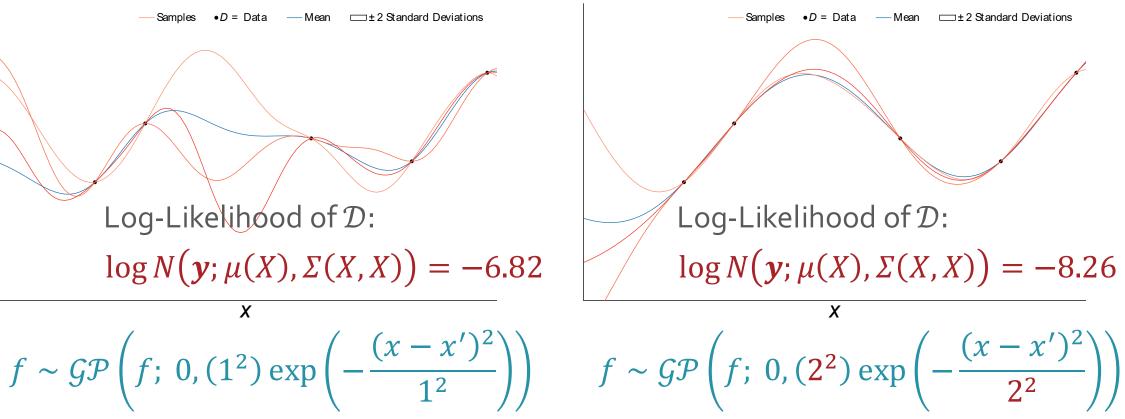
$$f \sim \mathcal{GP}\left(f; \ 0, (1^2) \exp\left(-\frac{(x-x')^2}{1^2}\right)\right)$$



Kernel Hyperparameters

- Can be set via MLE
- As long as μ and Σ are differentiable, the log-likelihood is differentiable with respect to the kernel hyperparameters





Noise

 Assume a linear model with additive Gaussian noise and a zero-mean Gaussian prior on the weights:

$$y = \Phi \omega + \epsilon$$
 where $\epsilon \sim N(\mathbf{0}_N, \sigma^2 I_N)$ and $\omega \sim N(\mathbf{0}_{D'+1}, \Sigma)$

then given a new test data point x', the prediction is

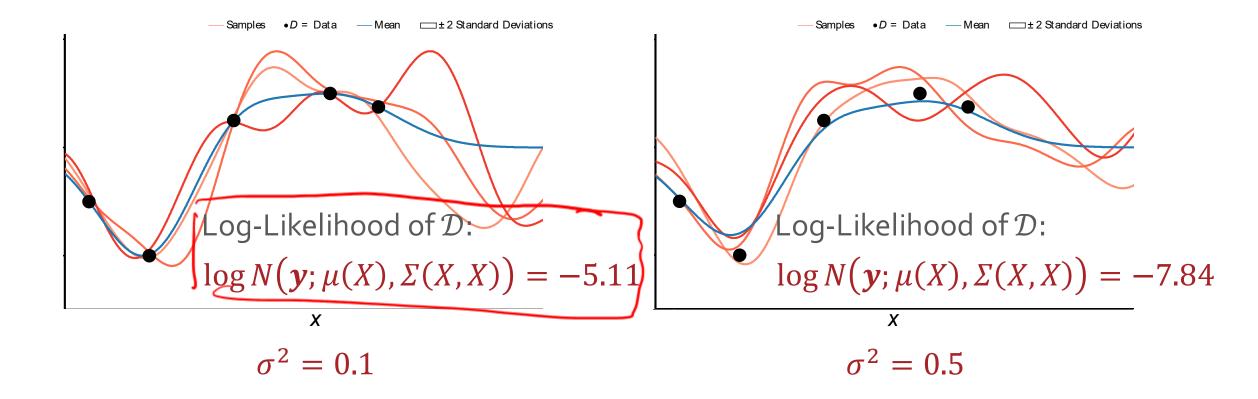
$$y' \mid y = \phi(x')^T \omega \mid y \sim N(\mu_{PRED}, \Sigma_{PRED})$$

where

$$\mu_{PRED} = K(\mathbf{x}', X)(K(X, X) + \sigma^2 I_N)^{-1} \mathbf{y},$$

$$\Sigma_{PRED} = K(x', x') - K(x', X)(K(X, X) + \sigma^2 I_N)^{-1} K(X, x)$$

- σ^2 is another hyperparameter we can tune
 - $\sigma^2=0$ is a noiseless fit: the mean will always pass through the observations exactly; $\sigma^2>0$ allows for deviations



Noise