

# Statistical Machine Learning

## Lecture 03 Gaussian Process I

Sharif University of Technology  
Spring 2021

Recall:

$D = \{x_1, x_2, \dots, x_n\}$  ← Data

Model:  $p(x|\theta)$  ← likelihood function

$\underline{\theta}$ : parameter Random / unknown

$p(\theta)$  ← prior on  $\theta$

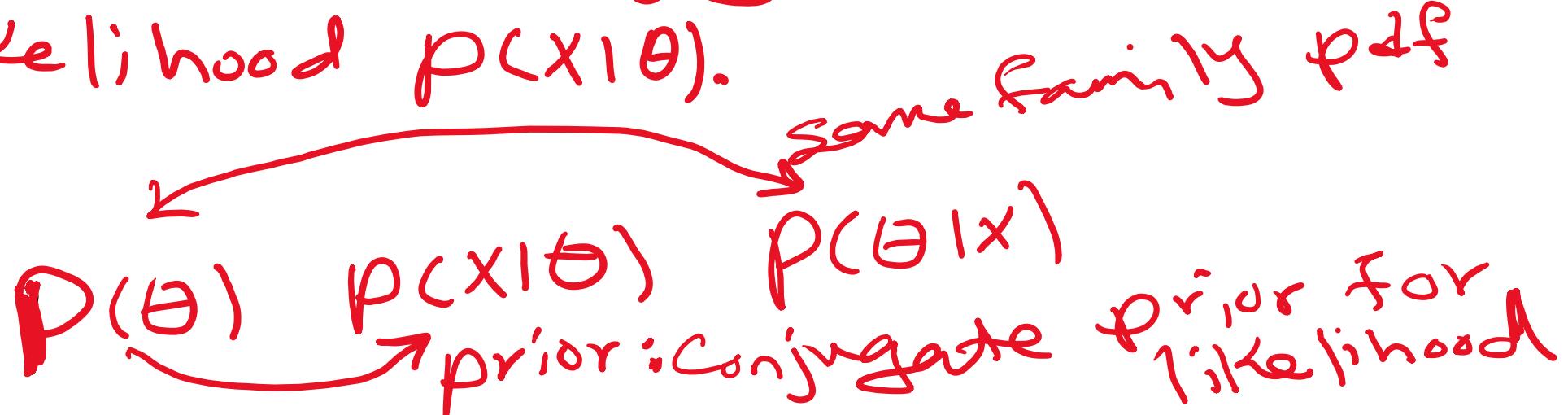
Bayes model:  $p(\theta|x) = \frac{p(x|\theta)p(\theta)}{\text{Posterior } p(x)}$

$\theta$  finite size ← parametric

$\theta$  not finite ← non parametric Bayesian (NPB)

Conjugate prior:

if  $p(x|\theta)$  is in the same family of prob. density fn. as the prior  $p(\theta)$  then prior and posterior called conjugate distributions or prior is the conjugate prior for the likelihood  $p(x|\theta)$ .



Examples of Conjugate prior:

Gaussian family is conjugate to itself.  
(Self-conjugate) with respect to a  
Gaussian likelihood

Normal prior  $\rightarrow$  Normal likelihood (mean)  $\rightarrow$  Normal posterior

Dirichlet dist.  $\rightarrow$  multinomial likelihood  $\rightarrow$  Dirichlet posterior

Prior  $\uparrow$                       Likelihood  $\uparrow$                       Posterior  $\uparrow$

NPB

Gaussian process  
(Regression)

Beta process  
Indian Buffet process

(IBP)  
(latent feature) models

Dirichlet process  
Chinese Restaurant process  
(CRP)  
Latent models  
(clustering)

Gaussian process:

Defines a dist.  $P(f)$  over function  $f$   
where  $f$  is a function mapping some input  
Space  $X$  to  $\underline{\mathbb{R}}$ :  $f: X \rightarrow \underline{\mathbb{R}}$

Let  $f = (f(x_1), f(x_2), \dots, f(x_n))$   
vector n-dim.  $x_i \in X$   
 $f$  is a R.V.

Then  $p(f)$  is a Gaussian process (GP)

if for any finite subset  ~~$\mathcal{S}$~~

$$\{x_1, \dots, x_n\} \subset X$$

the marginal dist. over that

subset  $p(f)$  has multivariate  
Gaussian dist.

GPs are parameterized by a:

$\mu(x)$  ← vector

$C(x, x')$  ← covariance fun

$x, x'$

$$P(f(x), f(x')) \sim N(\mu, \Sigma)$$

$$\mu = \begin{bmatrix} \mu(x) \\ \mu(x') \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} C(x, x) & C(x, x') \\ C(x', x) & C(x', x') \end{bmatrix}$$

$$\rightarrow C(x_i, x_j) = \gamma_0 \exp \left\{ - \frac{|x_i - x_j|}{\lambda} + v_1 + v_2 \delta_{ij} \right\}$$

$$(v_1, v_2, \lambda, \gamma_0)$$

④ How to use GP for nonlinear regression?

Data Set given  $D = \{(x_i, y_i)\}_{i=1}^n \equiv (\underline{x}, \underline{y})$   $f(\cdot)$

Model:  $\underline{y}_i = \underline{f}(\underline{x}_i) + \underline{\epsilon}_i$

$$f \sim GP(\cdot | \Omega, \Sigma)$$

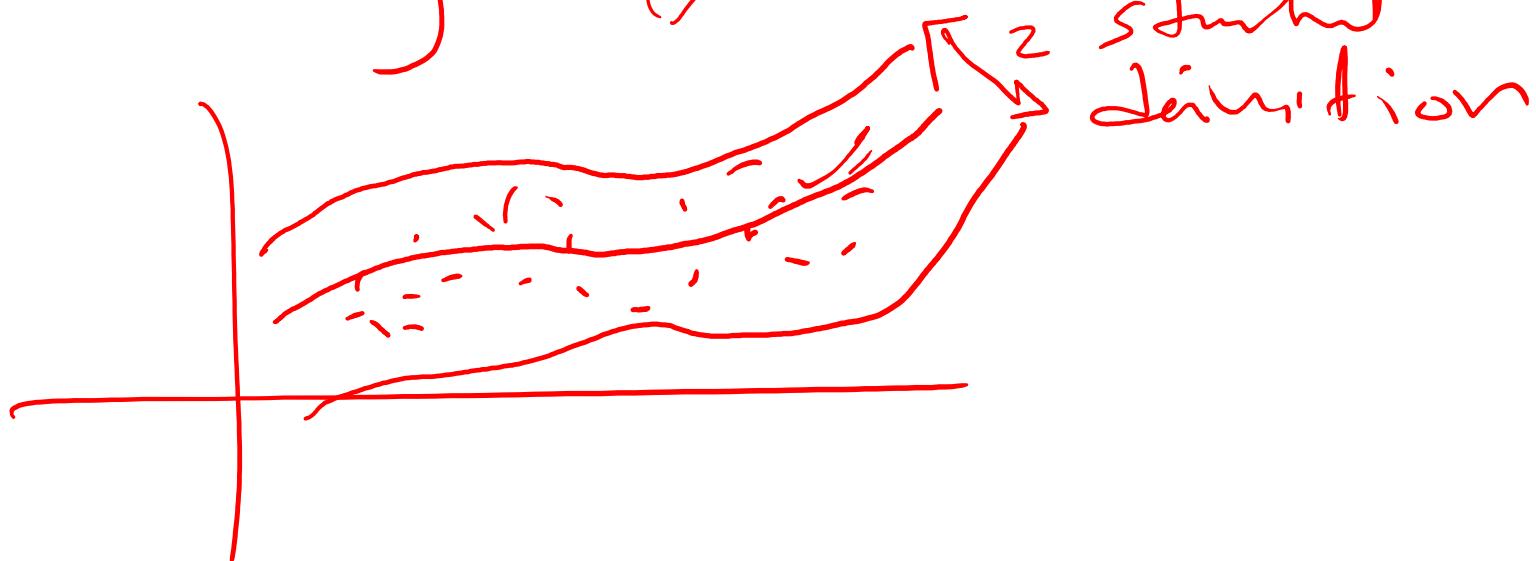
$$\epsilon_i \sim N(\cdot | 0, \sigma^2)$$

Noise

prior of  $f$  is GP  
likelihood is normal  
↓  
posterior is GP

$$p(y'|x; D) = \int p(y'|x, f, D) p(f|D) df$$

$$p(y|x) = \int p(y|f, x) p(f) df$$



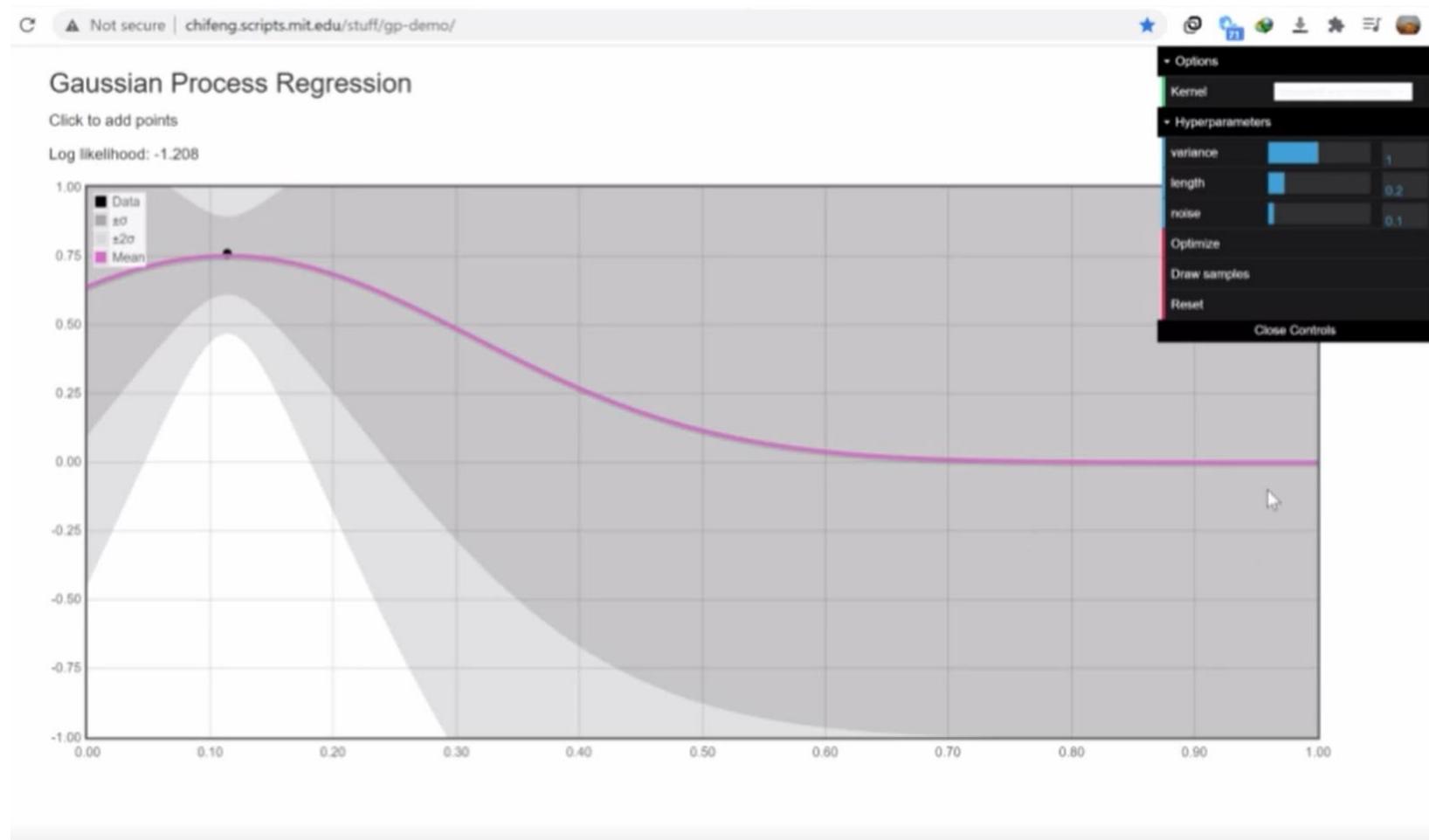
# Demo: Gaussian Process:

<http://chifeng.scripts.mit.edu/stuff/gp-demo/>



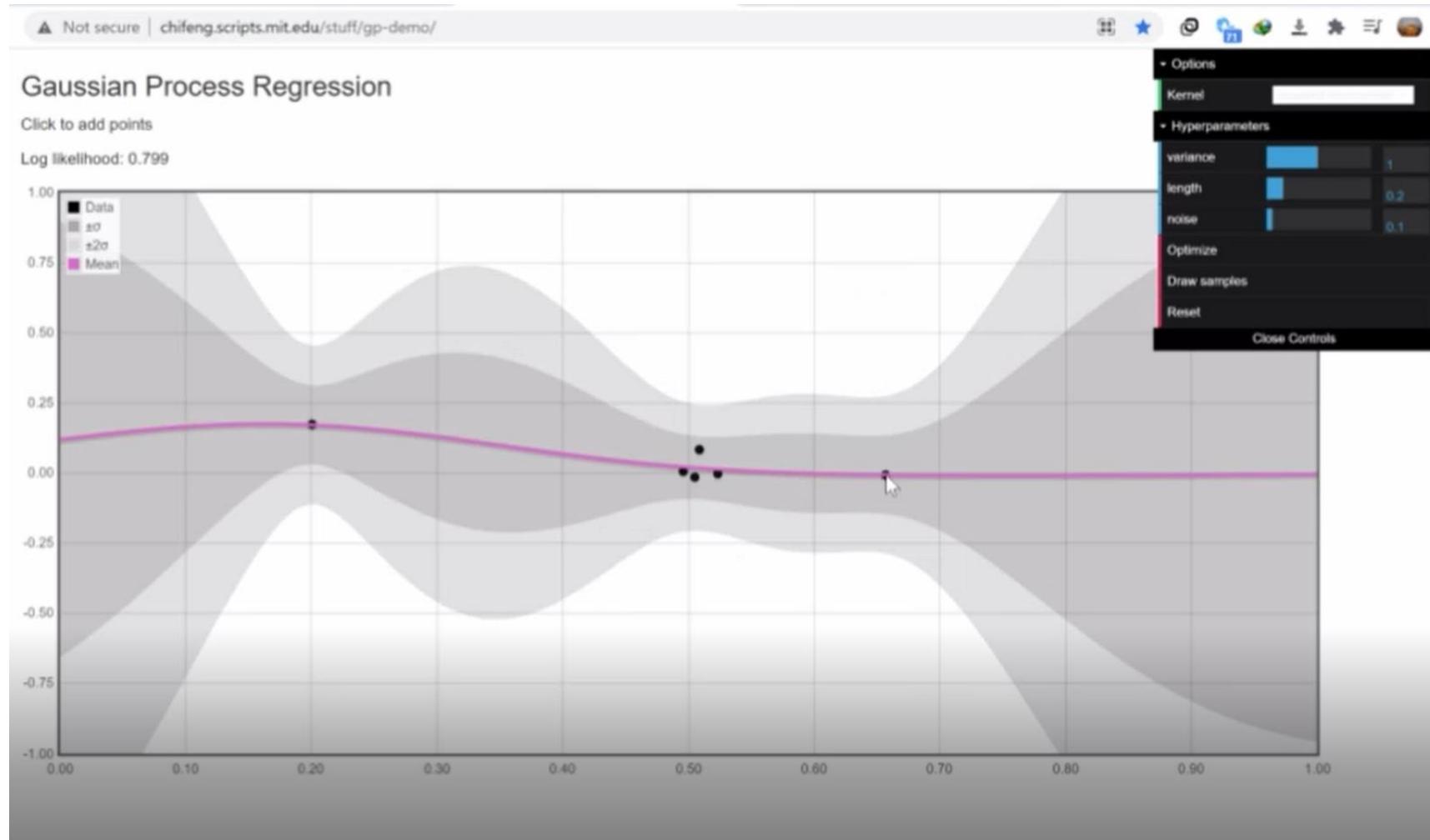
# Demo: Gaussian Process:

<http://chifeng.scripts.mit.edu/stuff/gp-demo/>



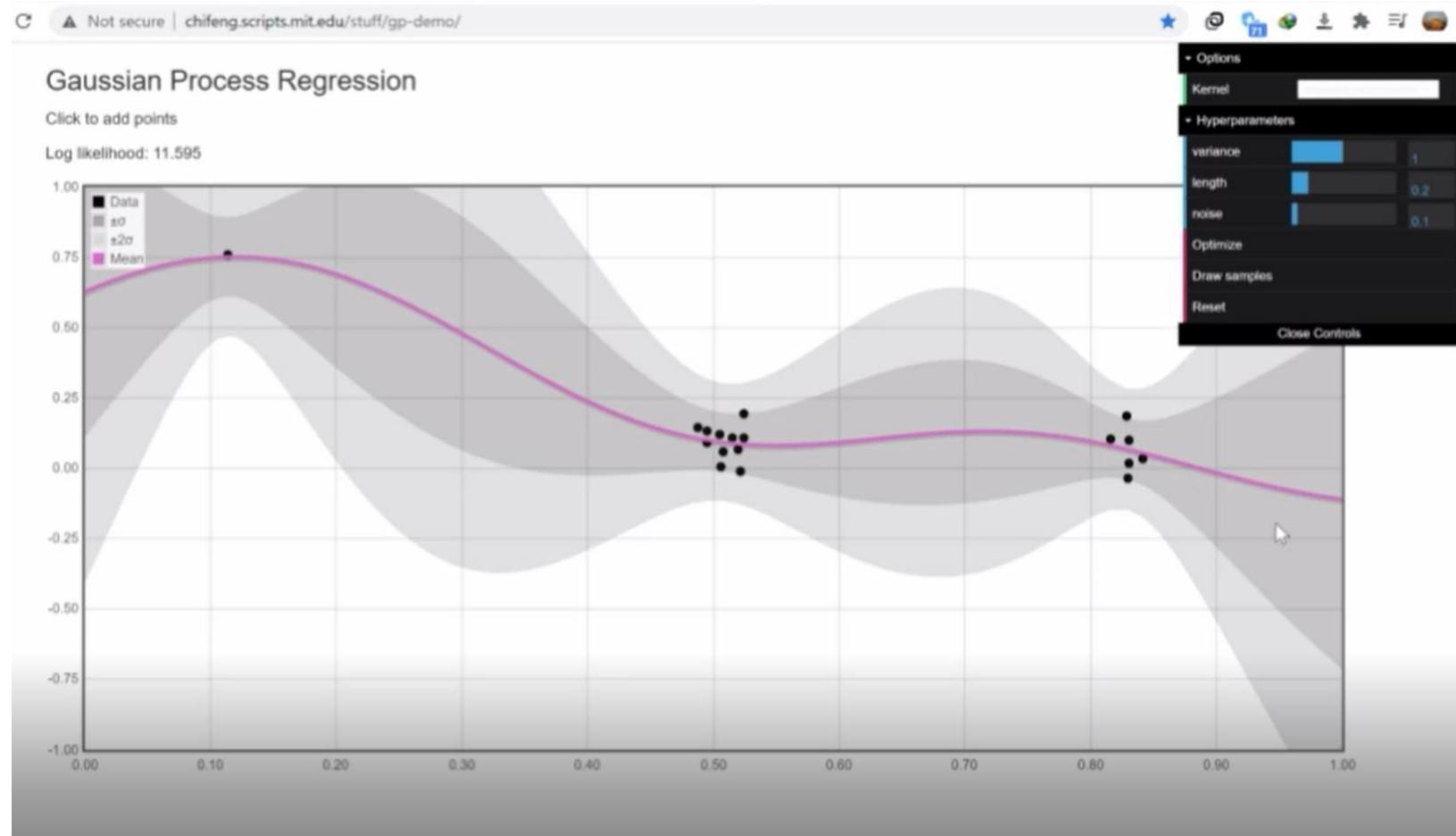
# Demo: Gaussian Process:

<http://chifeng.scripts.mit.edu/stuff/gp-demo/>



# Demo: Gaussian Process:

<http://chifeng.scripts.mit.edu/stuff/gp-demo/>



2. How do you use GP for linear regression?

$x_i$ : input     $y_i$ : output

$$y_i = \beta_0 + \beta_1 x_i + \epsilon_i$$

Basis function  $\phi_l(x_i)$

$$y_i = \sum_{l=1}^K \beta_l \phi_l(x_i) + \epsilon_i$$

← linear reg.  
with  $K$   
basis fn.

$\beta_l \sim N(0, \lambda_l)$

$\epsilon_i \sim N(0, \sigma^2)$

$E[y_i] = 0$

$$\text{Cov}(y_i, y_j) \triangleq c_{ij}$$

$$\triangleq \sum_l \lambda_l \phi_l(x_i) \phi_l(x_j) + \delta_{ij} \sigma^2$$

$\{y_i\}$  is a GP: Cov-f.  $c(x_i, x_j) \triangleq c_{ij}$

Can we use GP for classification? Yes

How? for Binary classification

