

# Lecture 13: Gaussian Processes

## CS480/680 Intro to Machine Learning

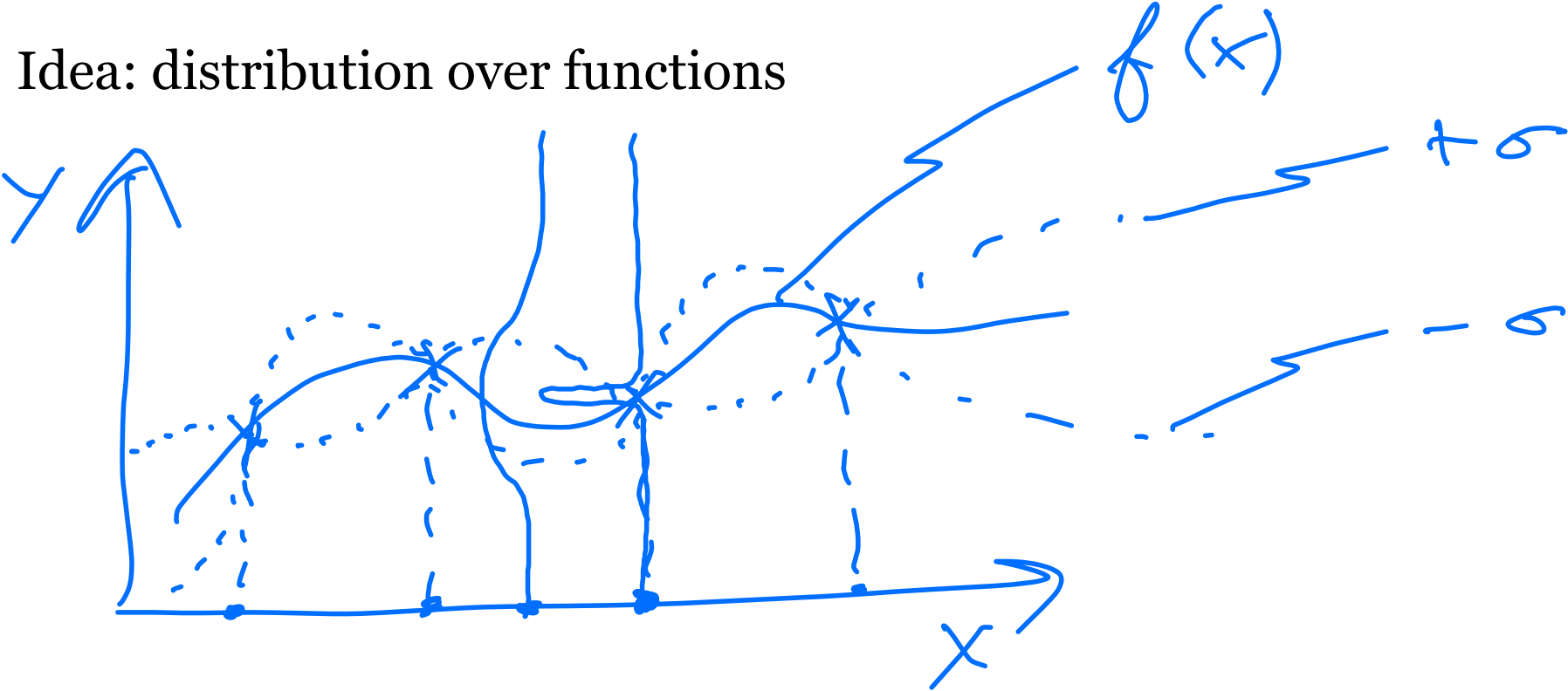
2023-2-28

Pascal Poupart  
David R. Cheriton School of Computer Science



# Gaussian Process Regression

- Idea: distribution over functions





# Bayesian Linear Regression

- Setting:  $f(\mathbf{x}) = \mathbf{w}^T \phi(\mathbf{x})$  and  $y = f(\mathbf{x}) + \epsilon$ 

$\downarrow$   
 unknown

$\downarrow$   
 $N(0, \sigma^2)$

- Function space view:

- Prior:  $\Pr(f(\mathbf{x}_*)) = \Pr(\mathbf{w}^T \phi(\mathbf{x}_*))$ 

$\downarrow$   
 Gaussian

$\xleftarrow{\text{change of variable}} \Pr(\mathbf{w})$   
 $\downarrow$   
 Gaussian

- Posterior:  $\Pr(f(\mathbf{x}_*) | \mathbf{X}, \mathbf{y}) = \Pr(\mathbf{w}^T \phi(\mathbf{x}_*) | \mathbf{X}, \mathbf{y})$ 

$\downarrow$   
 Gaussian

$\xleftarrow{\text{change of variable}} \Pr(\mathbf{w}^T | \mathbf{X}, \mathbf{y})$   
 $\downarrow$   
 Gaussian

# Gaussian Process

- According to the function view, there is a Gaussian at  $f(\mathbf{x}_*)$  for every  $\mathbf{x}_*$ . Those Gaussians are correlated through  $\mathbf{w}$ .
- What is the general form of  $\Pr(f)$  (i.e., distribution over functions)?
- Answer: **Gaussian Process** (infinite dimensional Gaussian distribution)

# Gaussian Process

- Distribution over functions:

$$f(\mathbf{x}) \sim GP(m(\mathbf{x}), k(\mathbf{x}, \mathbf{x}')) \forall \mathbf{x}, \mathbf{x}'$$

- Where  $m(\mathbf{x}) = E(f(\mathbf{x}))$  is the mean

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

is the kernel covariance function

# Mean function $m(\mathbf{x})$

- Compute the mean function  $m(\mathbf{x})$  as follows:
- Let  $f(\mathbf{x}) = \phi(\mathbf{x})^T \mathbf{w}$   
with  $\mathbf{w} \sim N(\mathbf{0}, \alpha^{-1} \mathbf{I})$
- Then  $m(\mathbf{x}) = E(f(\mathbf{x}))$   
 $= E(\mathbf{w})^T \phi(\mathbf{x})$   
 $= \mathbf{0}$

# Kernel covariance function $k(\mathbf{x}, \mathbf{x}')$

- Compute kernel covariance  $k(\mathbf{x}, \mathbf{x}')$  as follows:

- $$\begin{aligned} k(\mathbf{x}, \mathbf{x}') &= E(f(\mathbf{x})f(\mathbf{x}')) \\ &= \phi(\mathbf{x})^T E(\mathbf{w}\mathbf{w}^T)\phi(\mathbf{x}') \\ &= \phi(\mathbf{x})^T \frac{I}{\alpha} \phi(\mathbf{x}') \\ &= \frac{\phi(\mathbf{x})^T \phi(\mathbf{x}')}{\alpha} \end{aligned}$$

- In some cases we can use domain knowledge to specify  $k$  directly.

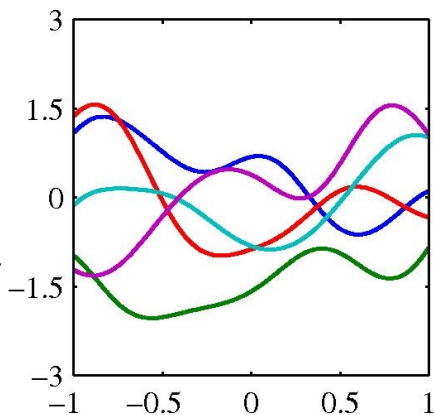


# Examples

- Sampled functions from a Gaussian Process

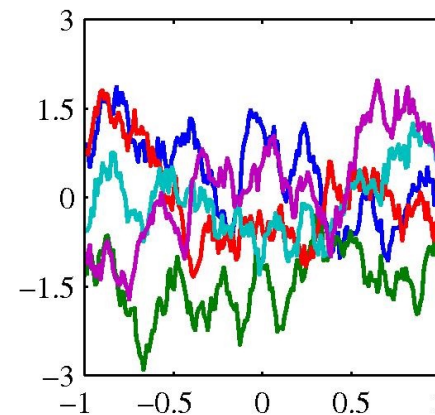
Gaussian kernel

$$k(\mathbf{x}, \mathbf{x}') = e^{-\frac{\|\mathbf{x} - \mathbf{x}'\|^2}{2\sigma^2}}$$



Exponential kernel  
(Brownian motion)

$$k(\mathbf{x}, \mathbf{x}') = e^{-\theta|\mathbf{x} - \mathbf{x}'|}$$



# Gaussian Process Regression

- Gaussian Process Regression corresponds to kernelized Bayesian Linear Regression
- Bayesian Linear Regression:
  - Weight space view
  - Goal:  $\Pr(\mathbf{w}|\mathbf{X}, \mathbf{y})$  (posterior over  $\mathbf{w}$ )
  - Complexity: cubic in # of basis functions
- Gaussian Process Regression:
  - Function space view
  - Goal:  $\Pr(f|\mathbf{X}, \mathbf{y})$  (posterior over  $f$ )
  - Complexity: cubic in # of training points

# Recap: Bayesian Linear Regression

- Prior:  $\Pr(\mathbf{w}) = N(\mathbf{0}, \Sigma)$
- Likelihood:  $\Pr(\mathbf{y}|\mathbf{X}, \mathbf{w}) = N(\mathbf{w}^T \Phi, \sigma^2 I)$
- Posterior:  $\Pr(\mathbf{w}|\mathbf{X}, \mathbf{y}) = N(\bar{\mathbf{w}}, \mathbf{A}^{-1})$   
where  $\bar{\mathbf{w}} = \sigma^{-2} \mathbf{A}^{-1} \Phi \mathbf{y}$  and  $\mathbf{A} = \sigma^{-2} \Phi \Phi^T + \Sigma^{-1}$
- Prediction:  
 $\Pr(y_* | \mathbf{x}_*, \mathbf{X}, \mathbf{y}) = N(\sigma^{-2} \phi(\mathbf{x}_*)^T \mathbf{A}^{-1} \Phi \mathbf{y}, \sigma^2 + \phi(\mathbf{x}_*)^T \mathbf{A}^{-1} \phi(\mathbf{x}_*))$
- Complexity: inversion of  $\mathbf{A}$  is cubic in # of basis functions

# Gaussian Process Regression

- Prior:  $\Pr(f(\cdot)) = N(m(\cdot), k(\cdot, \cdot))$
- Likelihood:  $\Pr(\mathbf{y}|\mathbf{X}, f) = N(f(\mathbf{X}), \sigma^2 \mathbf{I})$
- Posterior:  $\Pr(f(\cdot)|\mathbf{X}, \mathbf{y}) = N(\bar{f}(\cdot), k'(\cdot, \cdot))$   
where  $\bar{f}(\cdot) = k(\cdot, \mathbf{X})(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$  and  
 $k'(\cdot, \cdot) = k(\cdot, \cdot) + \sigma^2 \mathbf{I} - k(\cdot, \mathbf{X})(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} k(\mathbf{X}, \cdot)$
- Prediction:  $\Pr(y_*|\mathbf{x}_*, \mathbf{X}, \mathbf{y}) = N(\bar{f}(\mathbf{x}_*), k'(\mathbf{x}_*, \mathbf{x}_*))$
- Complexity: inversion of  $\mathbf{K} + \sigma^2 \mathbf{I}$  is cubic in # of training points

# Infinite Neural Networks

- Recall: neural networks with a single hidden layer (that contains sufficiently many hidden units) can approximate any function arbitrarily closely
  
- Neal 94: The limit of an infinite single hidden layer neural network is a Gaussian Process

# Bayesian Neural Networks

- Consider neural network with  $J$  hidden units and single identity output unit  $y_k$ :

$$y_k = f(\mathbf{x}; \mathbf{w}) = \sum_{j=1}^J w_{kj} h(\sum_i w_{ji} x_i + w_{j0}) + w_{k0}$$

- Bayesian learning: express prior over the weights
  - Weight space view:

$$\Pr(w_{kj}) \text{ where } E(w_{kj}) = 0, \text{Var}(w_{kj}) = \frac{\alpha}{J} \quad \forall kj,$$

$$\Pr(w_{k0}) \text{ where } E(w_{k0}) = 0, \text{Var}(w_{k0}) = \sigma^2 \quad \forall k$$

- Function space view: when  $J \rightarrow \infty$ , by the central limit theorem, an infinite sum of i.i.d. (identically and independently distributed) variables yields a Gaussian

$$\Pr(f(\mathbf{x})) = N(f(\mathbf{x}) | 0, \alpha E[h(\mathbf{x})h(\mathbf{x}')] + \sigma^2)$$

# Mean Derivation

- Calculation of the mean function:
  - $E[f(\mathbf{x})] = \sum_{j=1}^J E[w_{kj}h(\mathbf{x})] + E[w_{k0}]$ 
    - $= \sum_{j=1}^J E[w_{kj}]E[h(\mathbf{x})] + E[w_{k0}]$
    - $= \sum_{j=1}^J 0 E[h(\mathbf{x})] + 0$
    - $= 0$

# Covariance Derivation

- $$\begin{aligned} & Cov[f(\mathbf{x}), f(\mathbf{x}')] \\ &= E[f(\mathbf{x})f(\mathbf{x}')] - E[f(\mathbf{x})]E[f(\mathbf{x}')] \\ &= E[f(\mathbf{x})f(\mathbf{x}')] \\ &= E\left[\left(\sum_j w_{kj}h_j(\mathbf{x}) + w_{k0}\right)\left(\sum_j w_{kj}h_j(\mathbf{x}') + w_{k0}\right)\right] \\ &= \sum_{j=1}^J E[w_{kj}h_j(\mathbf{x})w_{kj}h_j(\mathbf{x}')] + E[w_{k0}w_{k0}] \\ &= \sum_{j=1}^J E[w_{kj}^2]E[h_j(\mathbf{x})h_j(\mathbf{x}')] + E[w_{k0}^2] \\ &= \sum_{j=1}^J Var[w_{kj}]E[h(\mathbf{x})h(\mathbf{x}')] + Var(w_{k0}) \\ &= \sum_{j=1}^J \frac{\alpha}{J} E[h(\mathbf{x})h(\mathbf{x}')] + \sigma^2 \\ &= \alpha E[h(\mathbf{x})h(\mathbf{x}')] + \sigma^2 \end{aligned}$$



# Bayesian Neural Networks

- When # of hidden units  $J \rightarrow \infty$ , then Bayesian neural net is equivalent to a Gaussian Process

$$\Pr(f(\cdot)) = GP(f(\cdot) | 0, \alpha E[h(\cdot)h(\cdot)] + \sigma^2)$$

- Note: this works for
  - Any activation function  $h$
  - Any i.i.d. prior over the weights with mean 0

# Case Study: AIBO Gait Optimization



# Gait Optimization

- Problem: find best parameter setting of the gait controller to maximize walking speed
  - Why?: Fast robots have a better chance of winning in robotic soccer
- Solutions:
  - Stochastic hill climbing
  - **Gaussian Processes**
    - Lizotte, Wang, Bowling, Schuurmans (2007) Automatic Gait Optimization with Gaussian Processes, *International Joint Conferences on Artificial Intelligence (IJCAI)*.

# Search Problem

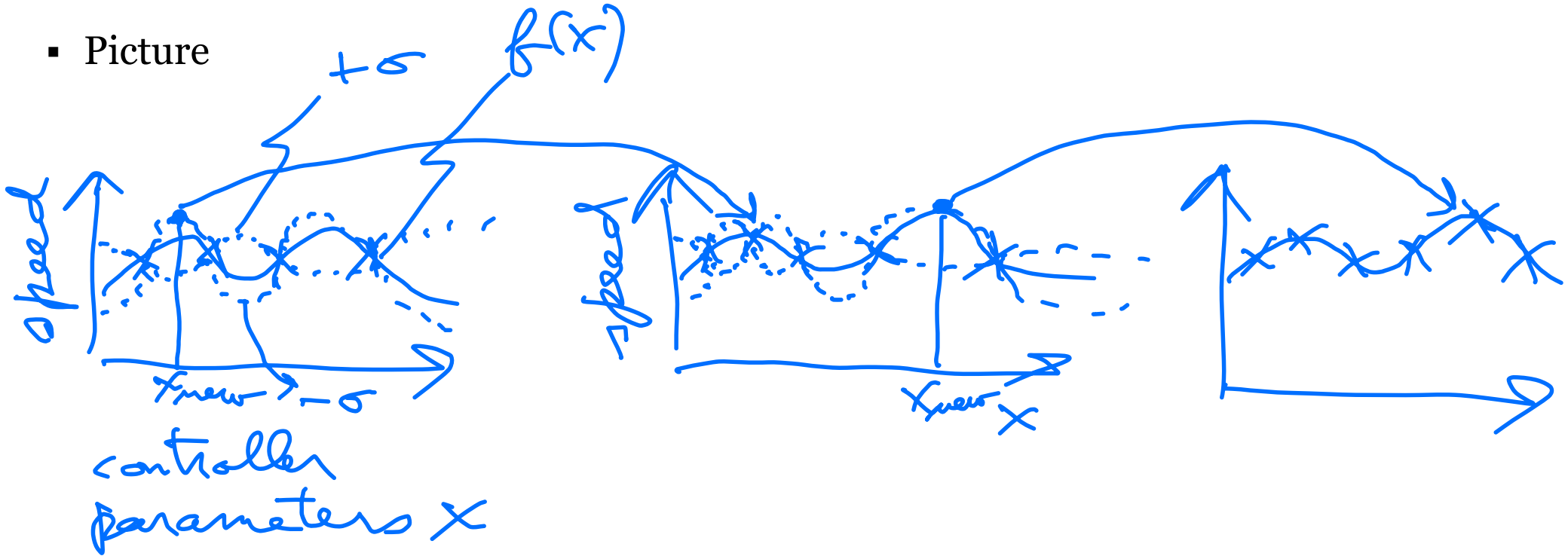
- Let  $\mathbf{x} \in \mathfrak{R}^{15}$ , be a vector of 15 parameters that defines a controller for gait
- Let  $f: \mathbf{x} \rightarrow \mathfrak{R}$  be a mapping from controller parameters to gait speed
- Problem: find parameters  $\mathbf{x}^*$  that yield highest speed.

$$\mathbf{x}^* \leftarrow \operatorname{argmax}_{\mathbf{x}} f(\mathbf{x})$$

But  $f$  is unknown...

# Approach

- Picture



# Approach

- Initialize  $f(\cdot) \sim GP(m(\cdot), k(\cdot, \cdot))$

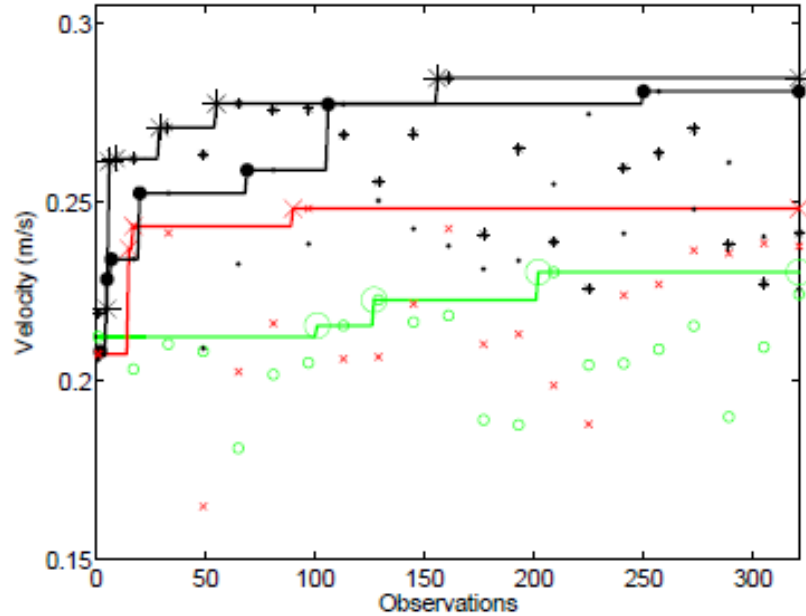
- Repeat:

- Select new  $\mathbf{x}$ :

$$\mathbf{x}_{new} \leftarrow \operatorname{argmax}_{\mathbf{x}} \frac{k(\mathbf{x}, \mathbf{x})}{\max_{\mathbf{x}' \in X} f(\mathbf{x}') - m(\mathbf{x})}$$

- Evaluate  $f(\mathbf{x}_{new})$  by observing speed of robot with parameters set to  $\mathbf{x}_{new}$
  - Update Gaussian process:
    - $\mathbf{X} \leftarrow \mathbf{X} \cup \{\mathbf{x}_{new}\}$  and  $\mathbf{y} \leftarrow \mathbf{y} \cup f(\mathbf{x}_{new})$
    - $m(\cdot) \leftarrow k(\cdot, \mathbf{X})(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} \mathbf{y}$
    - $k(\cdot, \cdot) \leftarrow k(\cdot, \cdot) + \sigma^2 \mathbf{I} - k(\cdot, \mathbf{X})(\mathbf{K} + \sigma^2 \mathbf{I})^{-1} k(\mathbf{X}, \cdot)$

# Results



(●) GP w/MPI	(*) GP w/MPI	(×) H.C1mb	(○) U.Rand
0.281 m/s	0.285 m/s	0.248 m/s	0.230 m/s
$\sigma_f^2 = 0.06$	$\sigma_f^2 = 0.6$		

Gaussian kernel:

$$k(\mathbf{x}, \mathbf{x}') = \sigma_f^2 e^{-\frac{1}{2}(\mathbf{x}-\mathbf{x}')^T S(\mathbf{x}-\mathbf{x}')}$$