

# ORIE 4742 - Info Theory and Bayesian ML

## Gaussian Processes

---

April 28, 2020

Sid Banerjee, ORIE, Cornell

- Bishop- Ch 6 ('Kernel methods')
- Gaussian processes for ML - Rasmussen & Williams

## normal-normal model (Gaussian rv with unknown $\mu$ )

- data  $D = \{X_1, X_2, \dots, X_n\} \in \mathbb{R}^n$
- model  $\mathcal{M}$ :  $X_i$  i.i.d. from  $\mathcal{N}(\mu, \tau)$ , with **unknown**  $\mu$ , **known**  $\tau = 1/\sigma^2$  <sup>precision</sup>

### normal-normal model

- **likelihood**:  $p(D|\mu) \propto \exp\left(-\tau \sum_{i=1}^n (x_i - \mu)^2/2\right)$  <sup>hyperparameters -  $M_\mu, \tau_\mu, \tau$</sup>
- **prior**:  $\mu \sim \mathcal{N}(M_\mu, 1/\tau_\mu) \propto \exp\left(-\tau_\mu(\mu - m_\mu)^2/2\right)$   <sub>$\tau$</sub>
- **posterior**: let  $\bar{x} = \frac{\sum_{i=1}^n x_i}{n}$ ,  $\tau_D = n\tau + \tau_\mu$  and  $m_D = \tau_D^{-1}(n\tau \cdot \bar{x} + \tau_\mu \cdot m_\mu)$

$$p(\mu|D) \sim \mathcal{N}(m_D, \tau_D^{-1})$$

- **posterior predictive distribution**:

$$p(x|D) \sim \mathcal{N}(m_D, \tau^{-1} + \tau_D^{-1})$$

# Bayesian linear regression - fixed basis fns $\phi_0(x)=1, \phi_1(x), \dots, \phi_{M-1}(x)$

- data  $D = \{(t_1, x_1), (t_2, x_2), \dots, (t_N, x_N)\} \in \mathbb{R}^n$
- model  $\mathcal{M}$ :  $t_i = \sum_{j=0}^{M-1} W_j \phi_j(x_i) + \epsilon_i$ , where  $\epsilon_i \sim \mathcal{N}(0, \beta^{-1})$

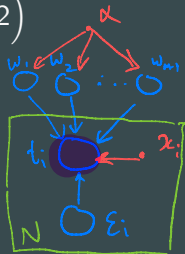
## Bayesian linear regression model

- likelihood:  $p(D|W) \propto \exp\left(-\beta \sum_{i=1}^N (x_i - W^T \phi(x_i))^2 / 2\right)$
- prior:  $W \sim \mathcal{N}(0, \alpha^{-1} I)$
- posterior: let  $m_D = T_D^{-1} \beta \Phi^T t$  and  $T_D = \beta \Phi^T \Phi + \alpha I$

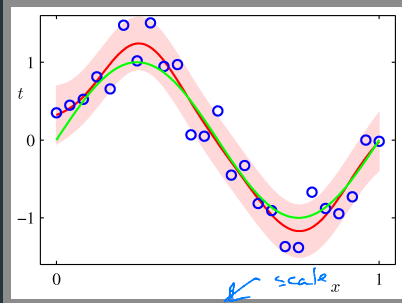
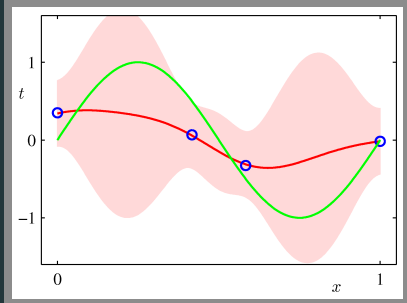
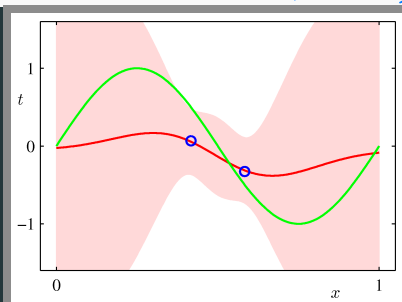
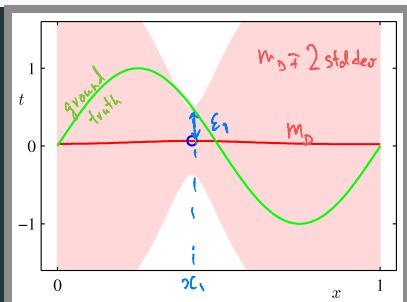
$$p(W|D) \sim \mathcal{N}(m_D, T_D^{-1})$$

- posterior predictive distribution:

$$p(t|x) \sim \mathcal{N}(m_D^T \phi(x), \beta^{-1} + \phi(x)^T T_D^{-1} \phi(x))$$

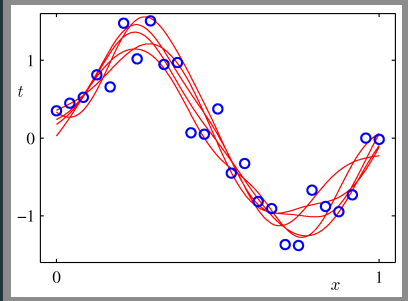
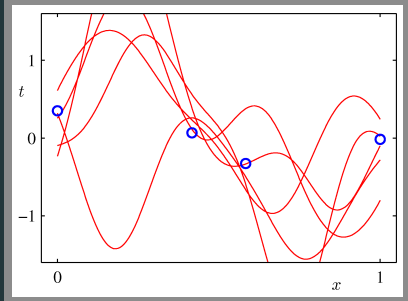
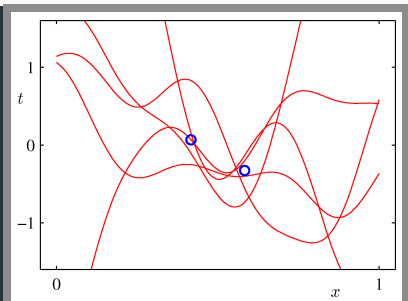
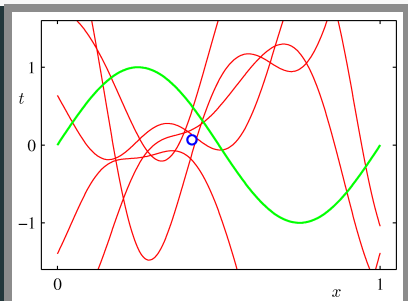


# Bayesian linear regression: posterior prediction (Bishop (h3))



Basis fns = 'Gaussian' -  $\phi(x) = \exp(-\theta_i(x-\mu_i)^2)$

# Bayesian linear regression: posterior sampling



# the 'equivalent' kernel

- data  $D = \{(t_1, x_1), (t_2, x_2), \dots, (t_N, x_N)\} \in \mathbb{R}^n$
- model  $\mathcal{M}$ :  $t_i = \sum_{j=0}^{M-1} W_j \phi(x_i) + \epsilon_i$ , where  $\epsilon_i \sim \mathcal{N}(0, \beta^{-1})$
- **prior**:  $W \sim \mathcal{N}(0, \alpha^{-1}I)$
- **posterior**: let  $m_D = T_D^{-1} \beta \Phi^T t$  and  $T_D = \beta \Phi^T \Phi + \alpha I$ , then

$$t(x|D) = m_D^T \phi(x) + \epsilon_D \sim \mathcal{N}(m_D^T \phi(x), \epsilon_D)$$

where  $\epsilon_D \sim \mathcal{N}(0, \beta^{-1} + \Phi^T T_D^{-1} \Phi)$

*'noise'*      *'uncertainty in params'*

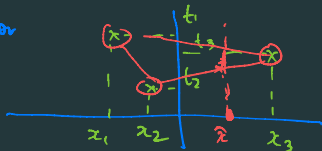
alternately,  $y(x|D) = \sum_{n=1}^N k(x, x_n) t_n$ , where  $k(x, y) = \beta \phi(x)^T T_D^{-1} \phi(y)$

*'equivalent kernel'*

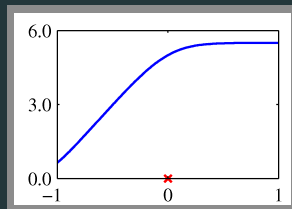
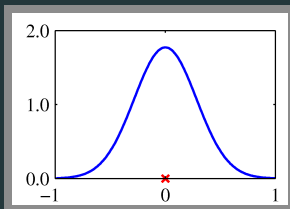
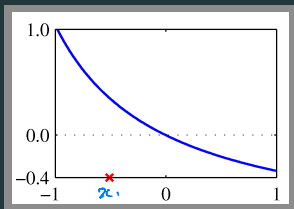
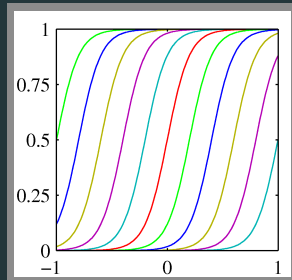
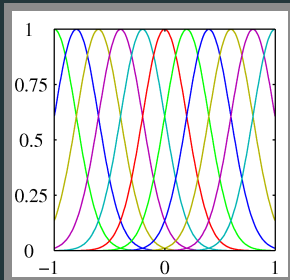
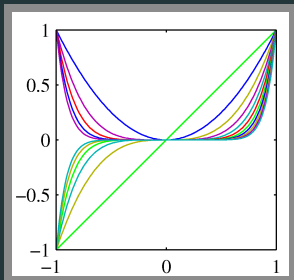
Sum over all data pts

*'weight'*  
= fn of  $x$   
and  $x_n$

$n^{\text{th}}$  observation



# basis functions and equivalent kernels



$$\phi(x) = (1 \ x \ x^2 \ \dots \ x^{M-1}), \quad \phi(x_1)^\top \phi(y) = 1 + x_1 y + x_1^2 y^2 + \dots + x_1^{M-1} y^{M-1}$$

# what are kernel methods? (Ch 6 of Bishop)

- generalized 'nearest-neighbor' methods
- given data  $D = \{(x_1, t_1), \dots, (x_n, t_n)\}$ , the resulting model is

$$\hat{y}(\hat{x}|D) = \sum_{i=1}^n k(x, x_i) t_i + \epsilon_D \sim \mathcal{N}(0, \text{Cov fn. of } \{x_i\})$$

## properties of kernels

function  $k(x, y)$  is a kernel of basis  $\phi(x)$  if  $k_\phi(x, y) = \phi(x)^T \phi(y)$

this is true if  $k$  is

- **symmetric**  $k(x, y) = k(y, x)$
- **positive-definite**  $K = \{k(x_i, x_j)\} \succeq 0$  for all  $\{x_i\}_{i=1}^n, n \in \mathbb{N}$

some special classes of kernels

- **stationary** kernel:  $k(x, y) = \psi(x - y)$
- **homogenous** kernel:  $k(x, y) = \psi(\|x - y\|)$

$$a^T K a \geq 0 \quad \forall a \in n \times 1 \text{ vectors}$$

Combine kernels

$$\begin{aligned} & - c_1 k_1 + c_2 k_2 = k(\phi, \phi) \\ & - \exp(c|k) \end{aligned}$$



# Gaussian process

$G$  is a random fn (Eg.  $G(x) = W_0 + W_1 x$ )

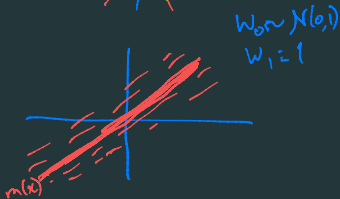
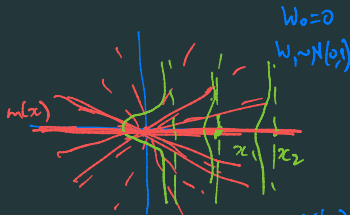
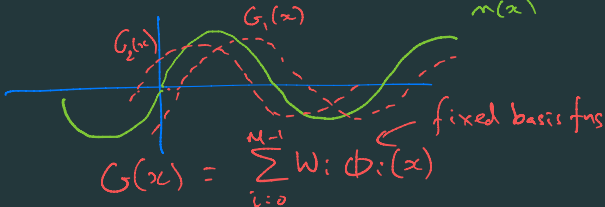
distribution over functions  $G(x)$  such that:

- any finite collection  $(G(x_1), G(x_2), \dots, G(x_n))$  is jointly Gaussian
- specified by mean  $m(x) = \mathbb{E}[G(x)]$  and covariance  $k(x, y) = \mathbb{E}[(G(x) - m(x))(G(y) - m(y))]$  (where  $k$  is a kernel)

example:  $y(x) = w^T \phi(x)$ , with  $w \sim \mathcal{N}(0, \alpha^{-1}I)$

Eg 1 -  $G(x) = W_0 + x$   
 $G(x) = W_1 x$

$G(x) = W_0 + \sin(x)$ ,  $W_0 \sim \mathcal{N}(0, 1)$   
 $m(x)$



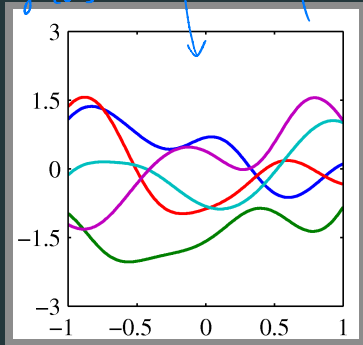
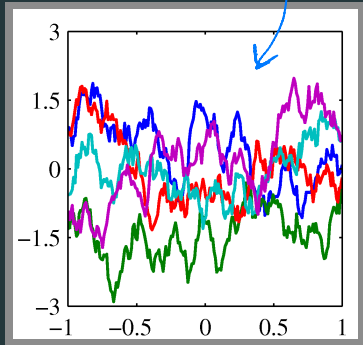
# Gaussian process examples

$$(m(x) = 0 \quad \forall x)$$

distribution over functions  $G(x)$  with jointly Gaussian samples, mean  $m(x) = \mathbb{E}[G(x)]$ , covariance  $k(x, y) = \mathbb{E}[(G(x) - m(x))(G(y) - m(y))]$

examples:  $k(x, y) = \exp(-\theta|x - y|)$ ,  $k(x, y) = \exp(-\theta(x - y)^2)$  (Gaussian kernel, rbf)

stationary, homogeneous



OU (Ornstein-Uhlenbeck process)  
(related to Brownian motion)

# Gaussian process regression (noise-free) $(t_i, x_i \sim t_i = \sum_{j=1}^M w_j \phi(x_i))$

- 'training' data  $D = \{(t_1, x_1), (t_2, x_2), \dots, (t_N, x_N)\} \in \mathbb{R}^n$
  - 'test' data:  $\tilde{x}$
  - model: GP with  $m(x) = 0$ , kernel  $k(x, y)$  input -  $k(x, y) \sim$  rbf
  - prior:  $(t_1, t_2, \dots, t_N, t) \sim \mathcal{N}\left(0, \begin{bmatrix} K_D & k \\ k^T & c \end{bmatrix}\right)$  (hyperparam -  $\theta$ )
 

$k(x_1, x_1)$	$\dots$	$k(x_1, x_N)$	$k(x_1, \tilde{x})$
$k(x_2, x_1)$	$\dots$	$k(x_2, x_N)$	$k(x_2, \tilde{x})$
$k(x_N, x_1)$	$\dots$	$k(x_N, x_N)$	$k(x_N, \tilde{x})$
$k(\tilde{x}, x_1)$	$\dots$	$k(\tilde{x}, x_N)$	$k(\tilde{x}, \tilde{x})$
- where  $K_D = \{k(x_i, x_j)\}$ ,  $k = \{k(\tilde{x}, x_j)\}$ , and  $c = k(\tilde{x}, \tilde{x})$

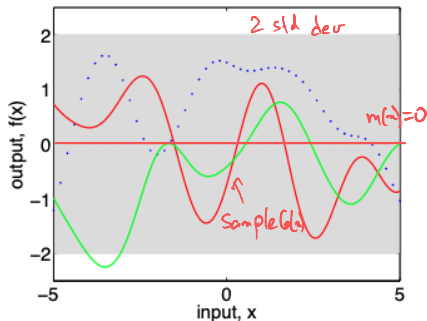
- posterior: conditioning on data  $D$ , we have
 
$$\tilde{t} \sim \mathcal{N}(k^T K_D^{-1} t, c - k^T K_D^{-1} k)$$

Q.  $(x, y, z) \sim \mathcal{N}\left(\begin{pmatrix} m_x \\ m_y \\ m_z \end{pmatrix}, \Sigma\right), P(x|y, z) = ?$

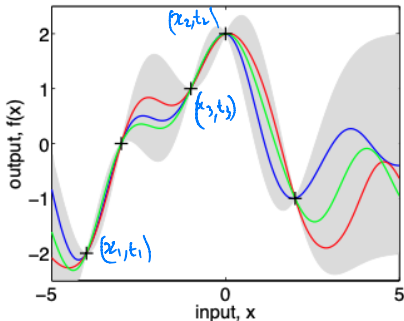
Application - 'Simulation'

# GP regression: example

$m(x)=0$ , rbf



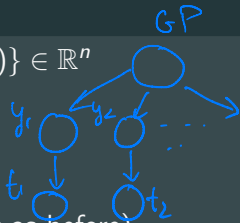
(a), prior



(b), posterior

# Gaussian process regression (with noise)

- 'training' data  $D = \{(t_1, x_1), (t_2, x_2), \dots, (t_N, X_N)\} \in \mathbb{R}^n$
- 'test' data:  $\tilde{x}$
- model:  $(x, y) \sim \text{GP}$  with  $m(x) = 0$ , kernel  $k(x, y)$   
observation  $t_i = y_i + \epsilon_i$  with  $\epsilon_i \sim \mathcal{N}(0, \beta^{-1})$
- prior:  $p(t|y) = \mathcal{N}(y, \beta^{-1}I_{n+1})$  and (with  $K_D, k, c$  as before)



$$\underbrace{(y_1, y_2, \dots, y_N)}_{\text{'training'}} \underbrace{, \tilde{y}}_{\text{'test'}} \sim \mathcal{N} \left( 0, \begin{bmatrix} K_D & k \\ k^T & c \end{bmatrix} \right)$$

$\uparrow m(x)=0$        $\leftarrow k(x,y)$

- posterior: conditioning on data  $D$ , we have

$$\tilde{t} \sim \mathcal{N} \left( k^T (K_D + \beta^{-1}I)^{-1} t, c - k^T (K_D + \beta^{-1}I)^{-1} k \right)$$

# GP noisy regression: example (Bishop) - rbf kernel

