

# Intro to Bayesian Linear Regression

"Bayesian" means taking into account prior information.

$$p(x) = \sum_y p(x|y) p(y)$$

↑      ↑      ↑  
posterior   likelihood   prior

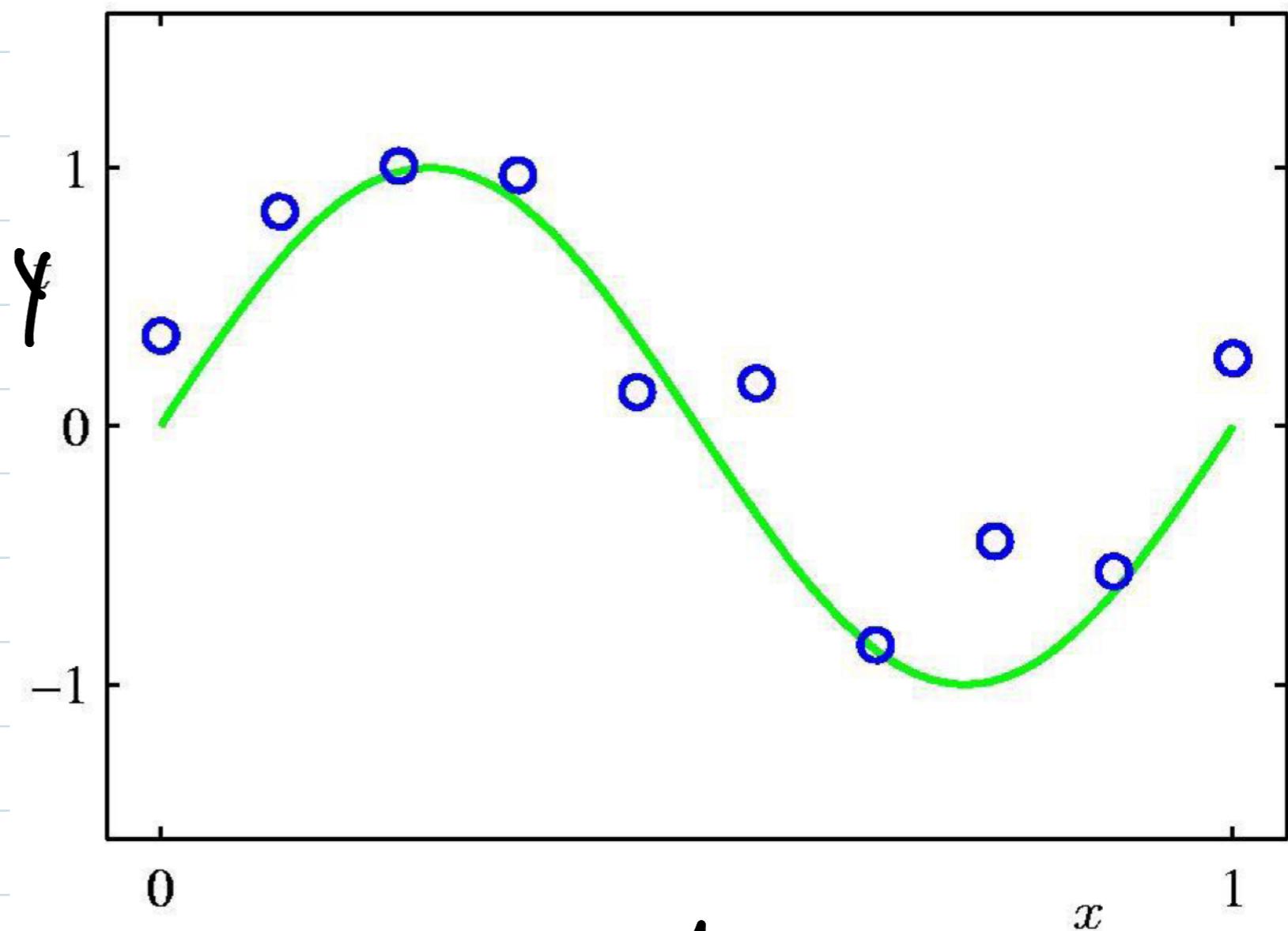
We are going to do regression where we combine data we gather with prior information to compute a predictive distribution.

To get there

- linear basis functions
- Gaussian distributions
- Maximum likelihood estimation
- Bayesian linear regression
- Pendulum example

# Linear Basis Function Models

Example: Polynomial Curve  
Fitting



$$y(x, w) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$
$$w = [w_0 \ w_1 \ w_2 \ \dots \ w_M]$$

Generally  $y(x, w) = \sum_{j=0}^{M-1} w_j \phi_j(x) = w^T \phi(x)$

$x$  vector of inputs

$y$  output

$M$  # of basis functions

$w \in R^M$  weights for basis functions (parameters)

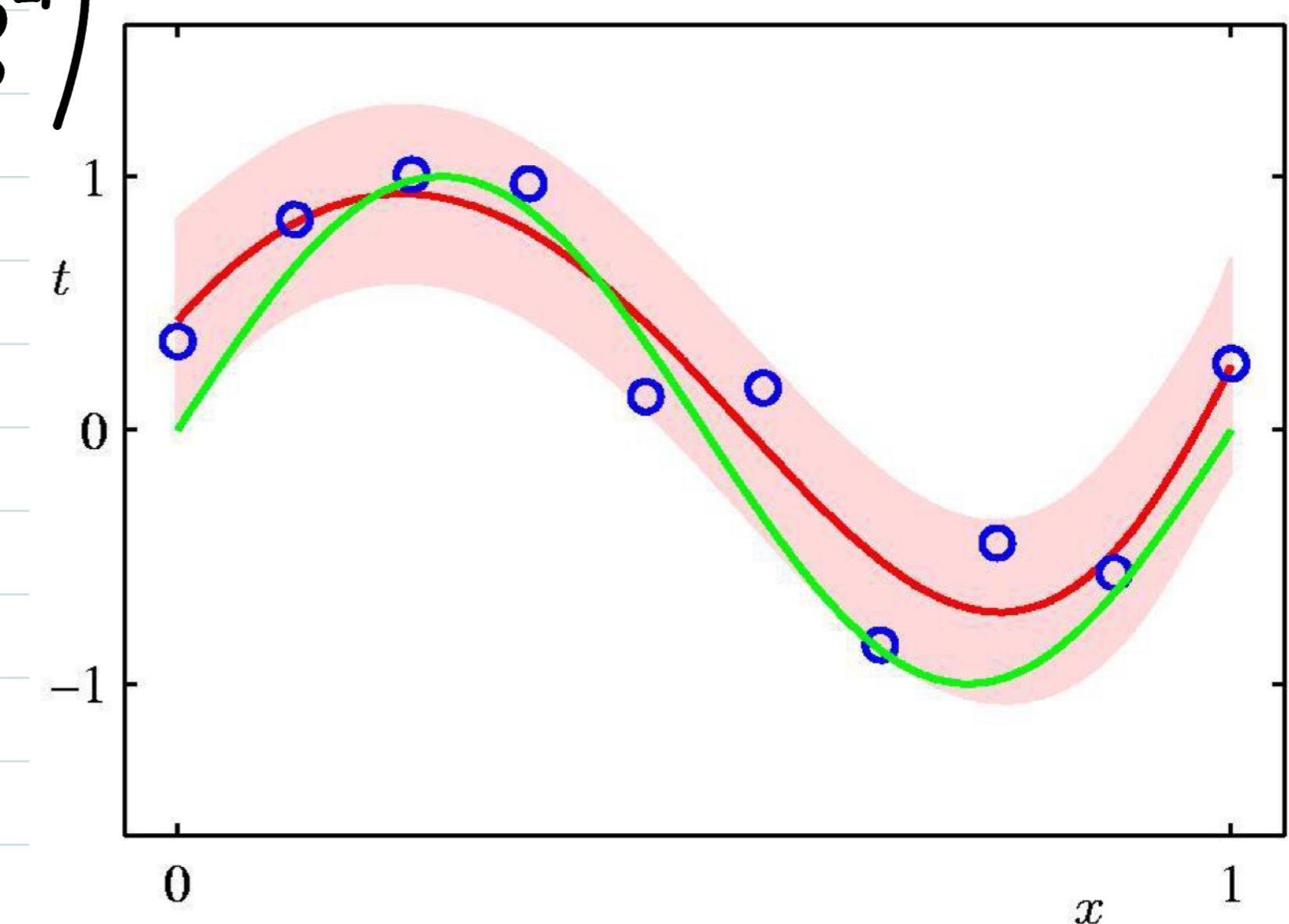
$\phi_j(x)$  is the  $j$ th basis function

We want to find

- $p(w)$  predictive distribution of weights
- $p(y|x, w)$  predictive distribution of outputs

## Predictive Distribution

$$p(t | x, w, \beta) = \mathcal{N}(t | y(x, w), \beta^{-1})$$



What do these predictive distributions look like?

We will assume they are Gaussians

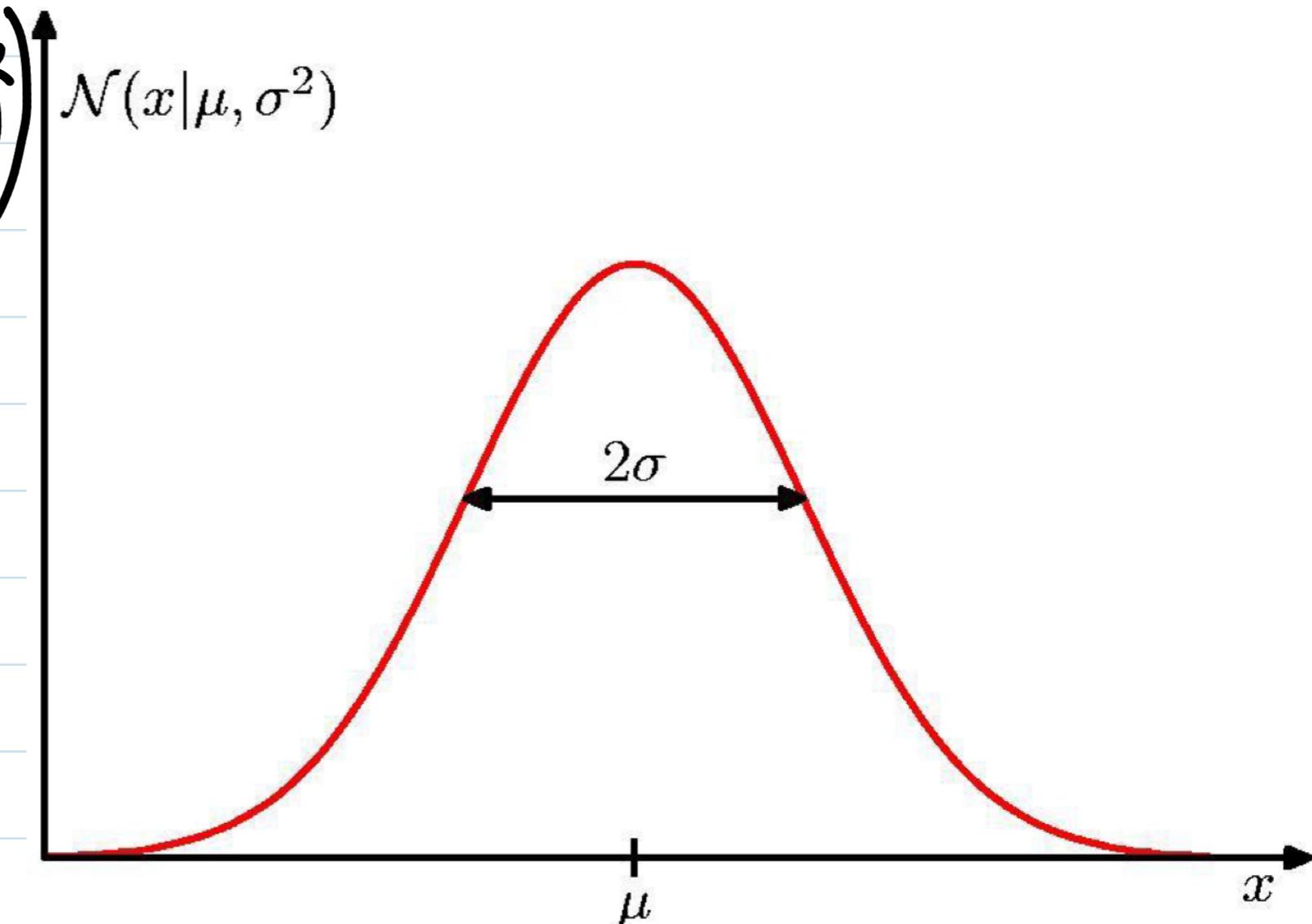
$$N(x|\mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} \exp\left(-\frac{1}{2\sigma^2}(x-\mu)^2\right)$$

probability that a random variable takes the value  $x$  given its mean is  $\mu$  and variance  $\sigma^2$

Notes:

$$N(x|\mu, \sigma^2) > 0$$

$$\int_{-\infty}^{\infty} N(x|\mu, \sigma^2) dx = 1$$



Matlab function  
randn()

$$E[X] = \int_{-\infty}^{\infty} N(x|\mu, \sigma^2) x dx = \mu$$

↑

expected  
value of  $X$

$$E[X^2] = \int_{-\infty}^{\infty} N(x|\mu, \sigma^2) x^2 dx = \mu^2 + \sigma^2$$

$$\text{Var}[X] = E[X^2] - E[X]^2 = \sigma^2$$

$$E[f] = \sum_x p(x) f(x) = \int p(x) f(x) dx$$

$$\text{Var}[f] = E[(f(x) - E[f(x)])^2] = E[f(x)^2] - E[f(x)]^2$$

$$\begin{aligned}\text{COV}[X, Y] &= E_{x,y}[\{x - E[x]\}\{y - E[y]\}] \\ &= E_{x,y}[xy] - E[x]E[y]\end{aligned}$$

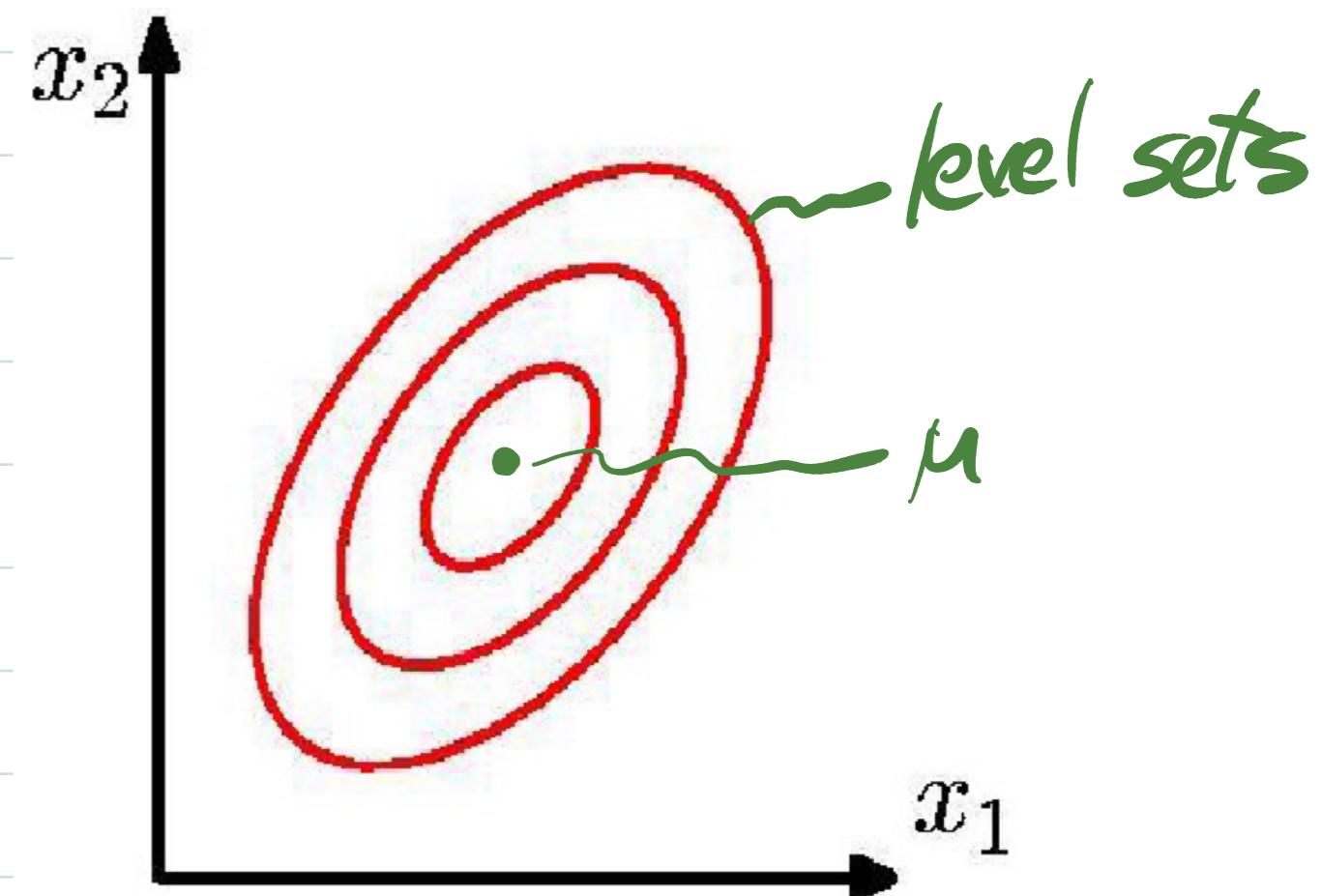
Multivariate Gaussian

$$N(x|\mu, \Sigma) = \frac{1}{(2\pi)^{d/2}} \frac{1}{|\Sigma|^{1/2}} \dots$$

$$\exp\left[-\frac{1}{2} (x-\mu)^T \Sigma^{-1} (x-\mu)\right]$$

$\mu \in \mathbb{R}^d$  mean vector  
 $d$  dimension of  $x$

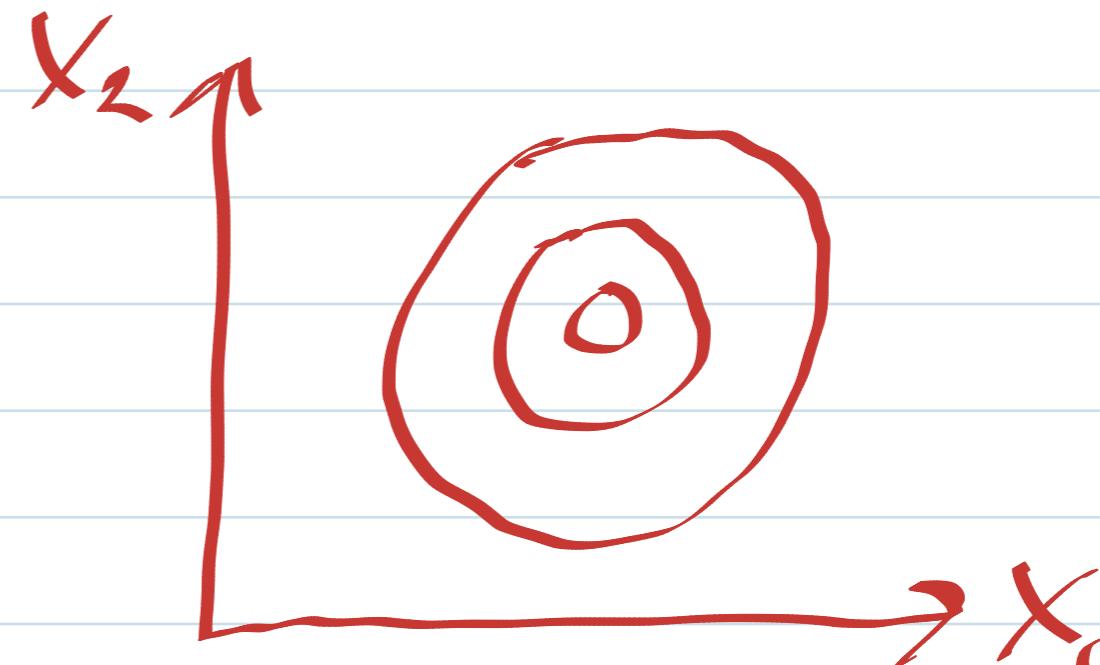
$x \in \mathbb{R}^d$  random variable  
 $\Sigma \in \mathbb{R}^{d \times d}$  covariance matrix



$\Sigma$ , the covariance matrix determines the orientation  
of the level sets

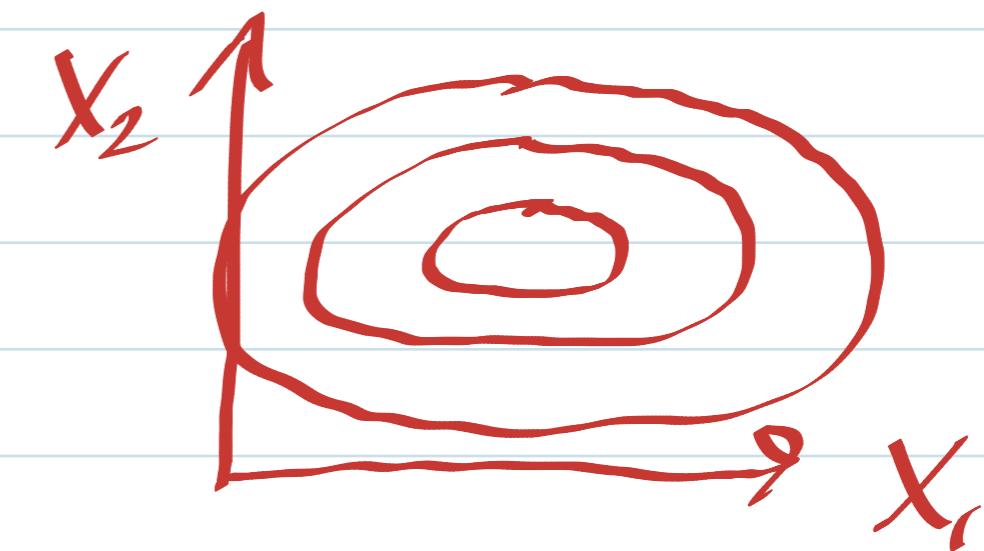
e.g.

$$\Sigma = \begin{bmatrix} \sigma^2 & 0 \\ 0 & \sigma^2 \end{bmatrix}$$

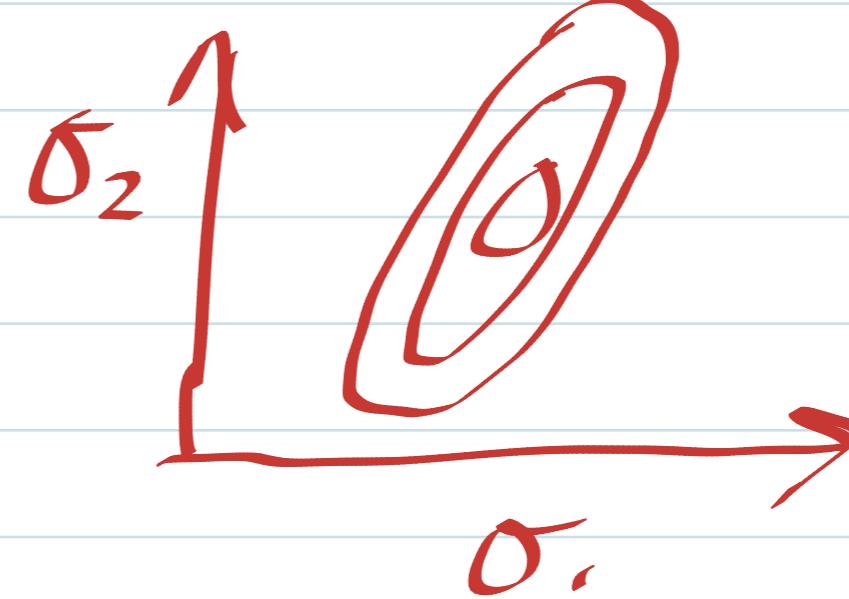


$$\Sigma = \begin{bmatrix} \sigma_1^2 & 0 \\ 0 & \sigma_2^2 \end{bmatrix}$$

$$\sigma_1^2 > \sigma_2^2$$



$$\Sigma = \begin{bmatrix} \sigma_1^2 & \text{cov}(\sigma_1, \sigma_2) \\ \text{cov}(\sigma_1, \sigma_2) & \sigma_2^2 \end{bmatrix}$$



# Maximum Likelihood and Least Squares

assume we make observations from a deterministic function  
with added Gaussian noise

$$t = y(x, w) + \epsilon \quad p(\epsilon | \beta) = N(\epsilon | 0, \beta^{-1})$$

$\beta$  is the inverse of variance  
or a.k.a the confidence

$$y(x, w) = w^T \phi(x)$$

$$p(t | x, w, \beta) = N(t | y(x, w), \beta^{-1})$$

Given observed inputs  $X \in \mathbb{R}^{d \times N}$   
 $x \in \mathbb{R}^d$  input vector

$N$  # of observations

and given observed targets

$$t = [t_1, t_2, \dots, t_n]^T$$

The likelihood function is

$$p(t | X, w, \beta) = \prod_{n=1}^N N(t_n | w^T \phi(x_n), \beta^{-1})$$

The probability of observing the set of targets (noisy observations of  $y$ ) given inputs  $X$  and parameters  $w$  and  $\beta$ .

Goal : Maximize the likelihood over  $w$  and  $\beta$

maximum likelihood estimation (MLE)

Take the log of the likelihood

$$\begin{aligned}\ln p(t|w, \beta) &= \sum_{n=1}^N \ln N(t_n | w^T \phi(x_n), \beta^{-1}) \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln 2\pi - \frac{1}{2} \beta \sum_{n=1}^N [t_n - w^T \phi(x_n)]^2\end{aligned}$$

*Sum of squared errors*

MLE is the same as minimizing the

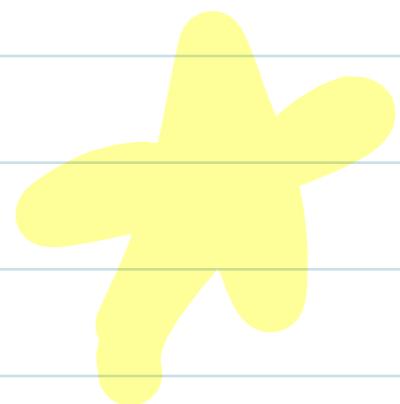
ss error w.r.t  $w, \beta$

To maximize the likelihood take the gradient

$$\nabla_w \ln p(t|w, \beta) = \beta \sum_{n=1}^N [t_n - w^T \phi(x_n)] \phi(x_n)^T = 0$$

Solve for  $W$

$$W_{ML} = (\Phi^T \Phi)^{-1} \Phi^T t$$



where

More-Penrose pseudo-inverse of  $\Phi$

$$\Phi = \begin{bmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \dots & \phi_{M-1}(x_2) \\ \vdots & & & \\ \phi_0(x_N) & \phi_1(x_N) & \dots & \phi_{M-1}(x_N) \end{bmatrix}$$

$\Phi \in \mathbb{R}^{N \times M}$

in Matlab

`pinv(Φ)`

## Multiple Outputs

$$\begin{aligned} p(\bar{\epsilon} | x, W, \beta) &= N(t | y(W, x), \beta^{-1} I) \\ &= N(t | W^T \phi(x), \beta^{-1} I) \end{aligned}$$

$q$  is the dimension of output vector  $\bar{T}$

$M$  is the number of basis function

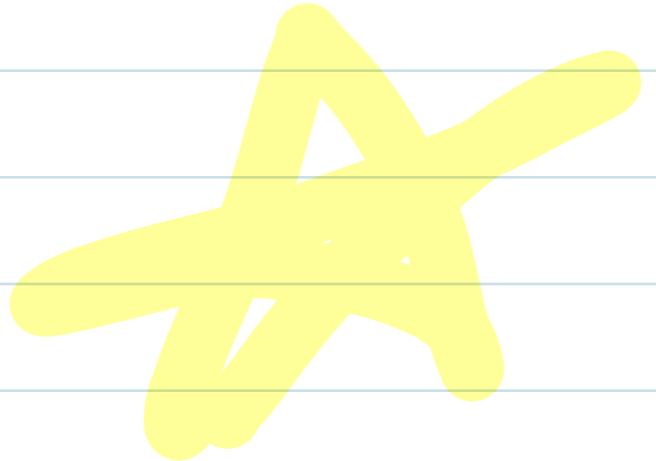
$W \in R^{M \times q}$  matrix of parameters

For outputs  $T = [\bar{\epsilon}_1, \bar{\epsilon}_2, \dots, \bar{\epsilon}_N]^T$  the log likelihood is

$$\begin{aligned} \ln p(T | X, W, \beta) &= \sum_{n=1}^N \ln N(\bar{\epsilon}_n | W^T \phi(x_n), \beta^{-1} I) \\ &= \frac{Nq}{2} \ln \left( \frac{\beta}{2\pi} \right) - \frac{\beta}{2} \sum_{n=1}^N \| t_n - W^T \phi(x_n) \|^2 \end{aligned}$$

maximize w.r.t.  $W$

$$W_{ML} = (\mathcal{I}^T \mathcal{I})^{-1} \mathcal{I}^T T$$



## Bayesian Linear Regression

Define a prior over  $W$

$$p(W) = N(W | m_0, S_0)$$

$m_0$  prior mean of parameters

$S_0$  prior Variance

After a bunch of Gaussian math the posterior

distribution for  $w$  is

$$p(w|t) = N(w|m_N, S_N)$$

★  $m_N = S_N^{-1} (S_0^{-1} m_0 + \beta \mathbb{I}^T t)$  posterior mean

confidence in prior prior mean like MLE estimate  
inverse of observation noise

$$S_N^{-1} = S_0^{-1} + \beta \mathbb{I}^T \mathbb{I}$$
 posterior confidence

↑  
prior confidence  
confidence gained by observing data

A common choice for a prior is

$$p(w) = N(w|0, \alpha^{-1}I)$$

Zero mean with  
confidence  $\propto$

then posterior becomes

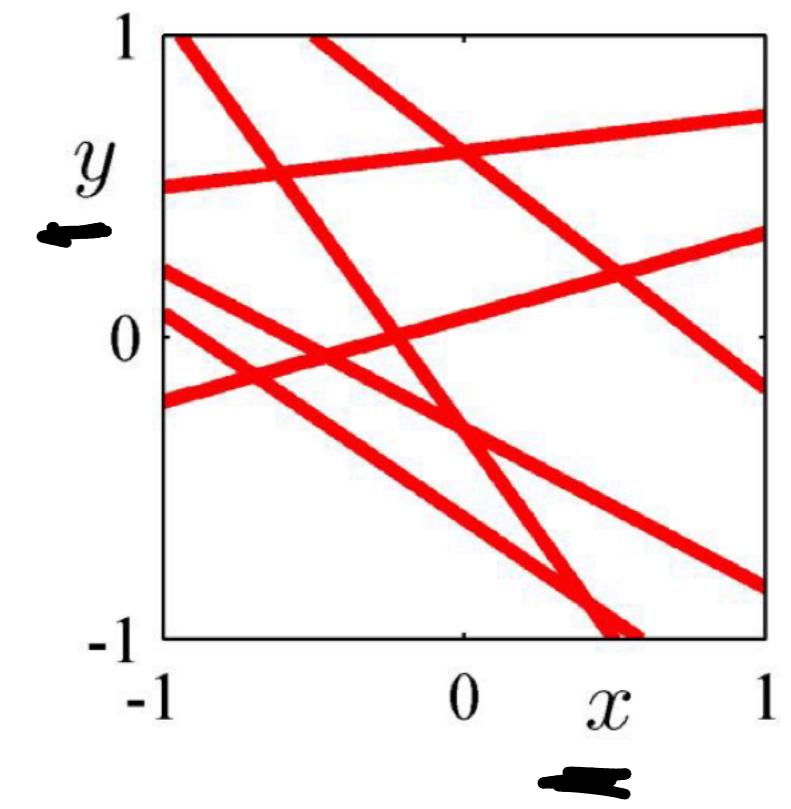
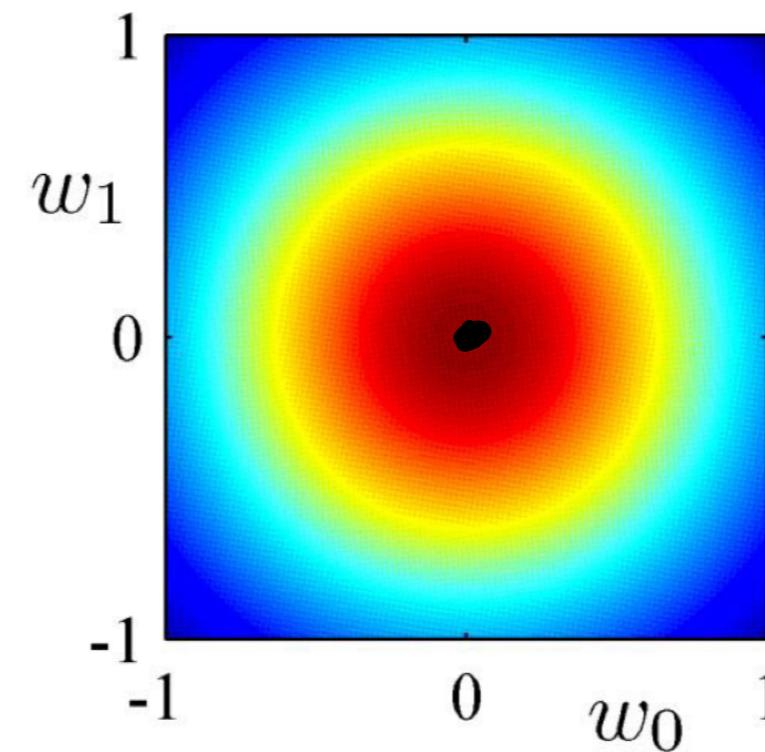
$$m_N = \beta S_N^{-1} \Phi^T t$$

$$S_N^{-1} = \alpha I + \beta \Phi^T \Phi$$

if  $\alpha$  is small (low confidence in prior)

$$m_N = \beta \beta^{-1} (\Phi^T \Phi)^{-1} \Phi^T t \Rightarrow MLE$$

# Examples of Bayesian Regression



function  $y = w_0 + w_1 x$

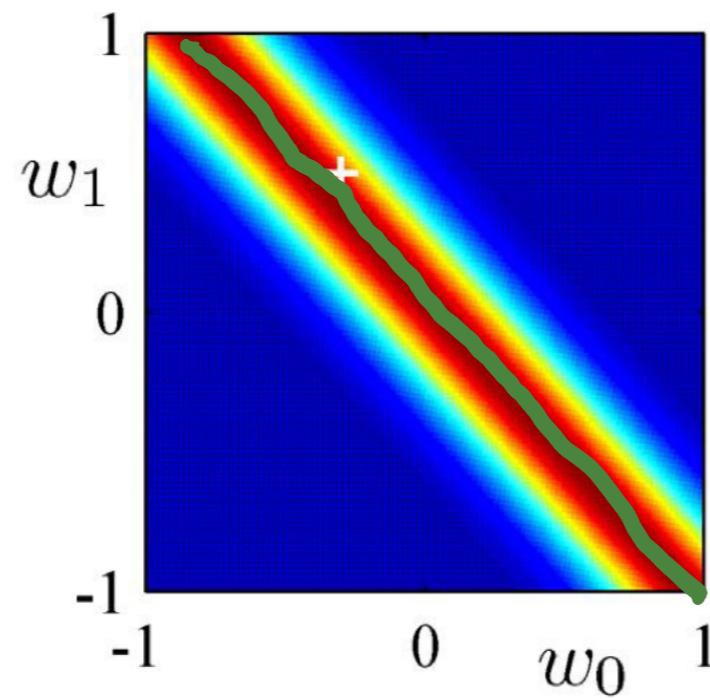
↑  
prior distribution  
of parameters

$$p(w) = N(w | 0, \Sigma_0^{-1})$$

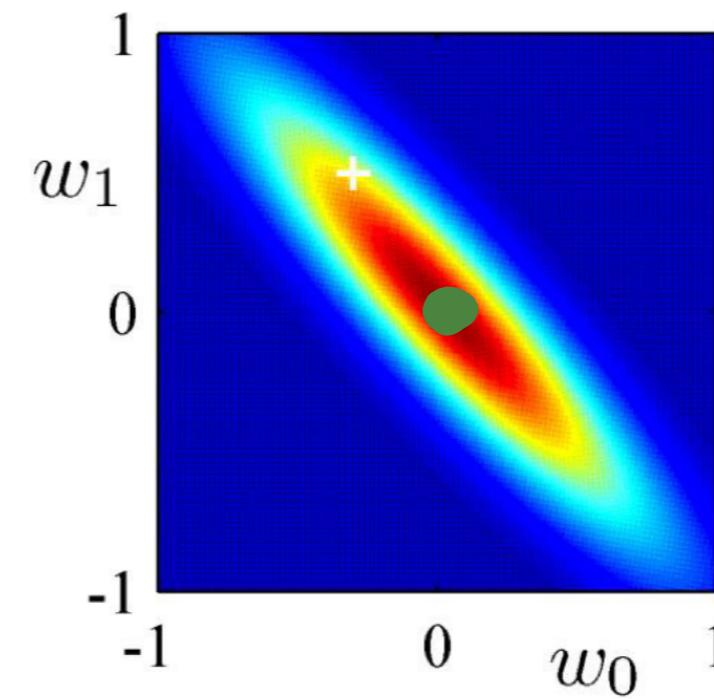
$$m_0 = [0, 0]$$

$$\Sigma_0^{-1} = \begin{bmatrix} \sigma_0^2 & 0 \\ 0 & \sigma_0^2 \end{bmatrix}$$

Observe 1 data point



likelihood



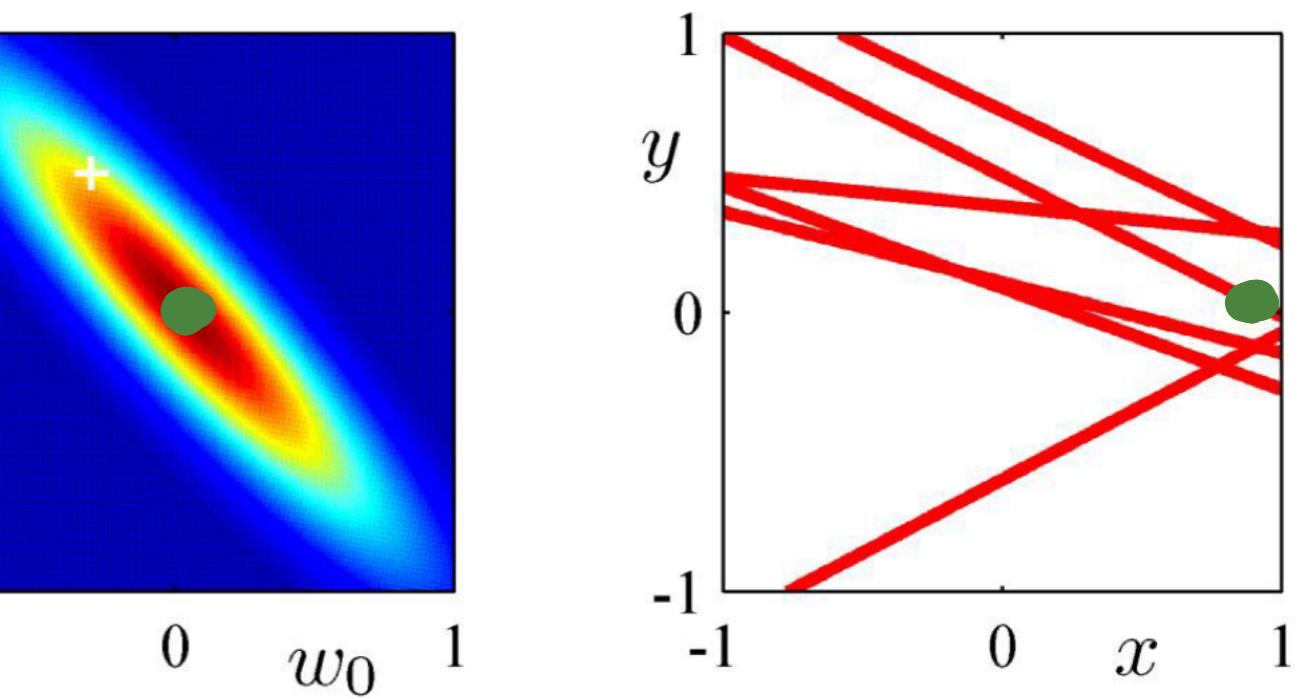
posterior

$$x_1 \approx 0.9$$

$$y_1 \approx 0.1$$

many optimal  
solutions

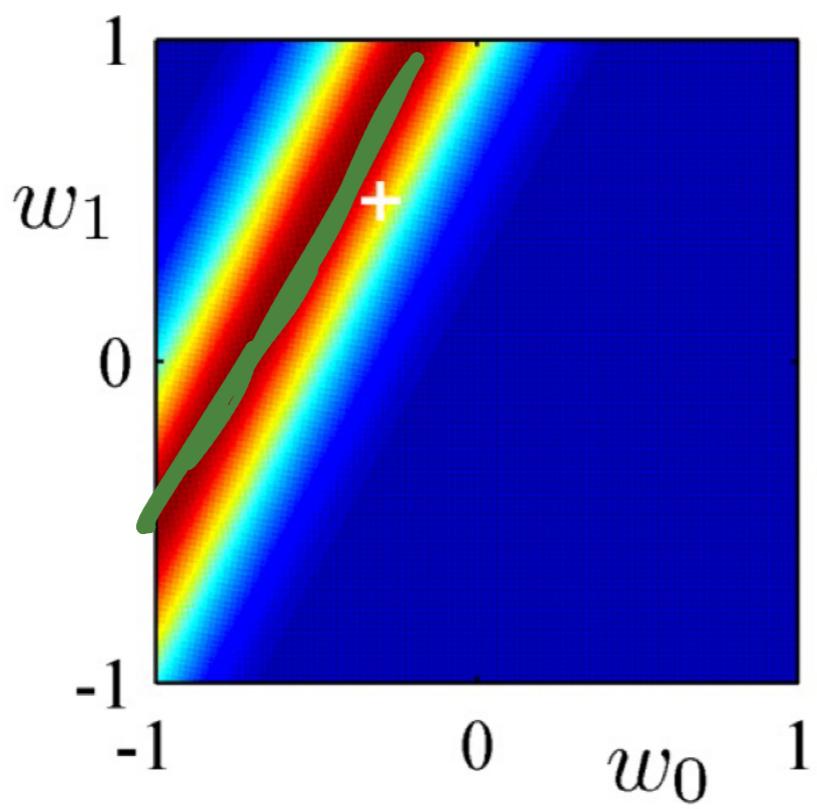
MLE



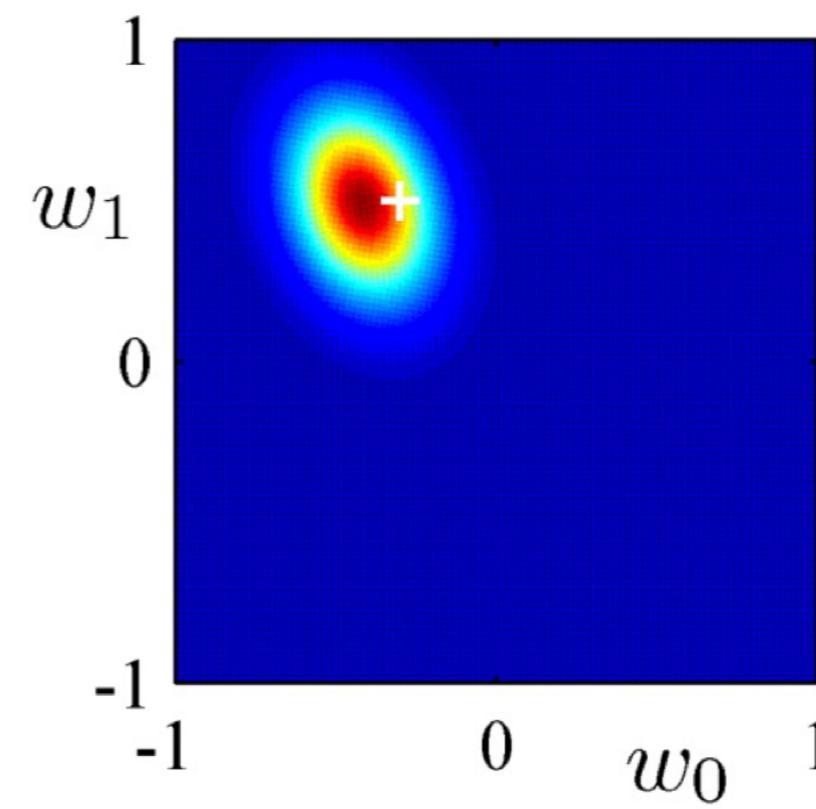
one optimal  
solution

Bayesian

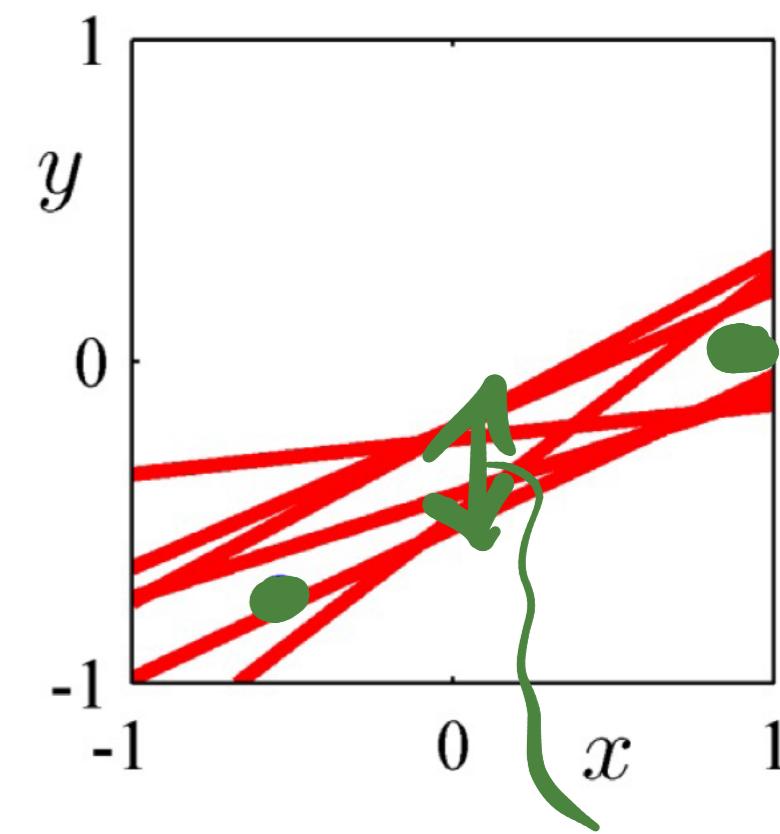
2 data points



likelihood

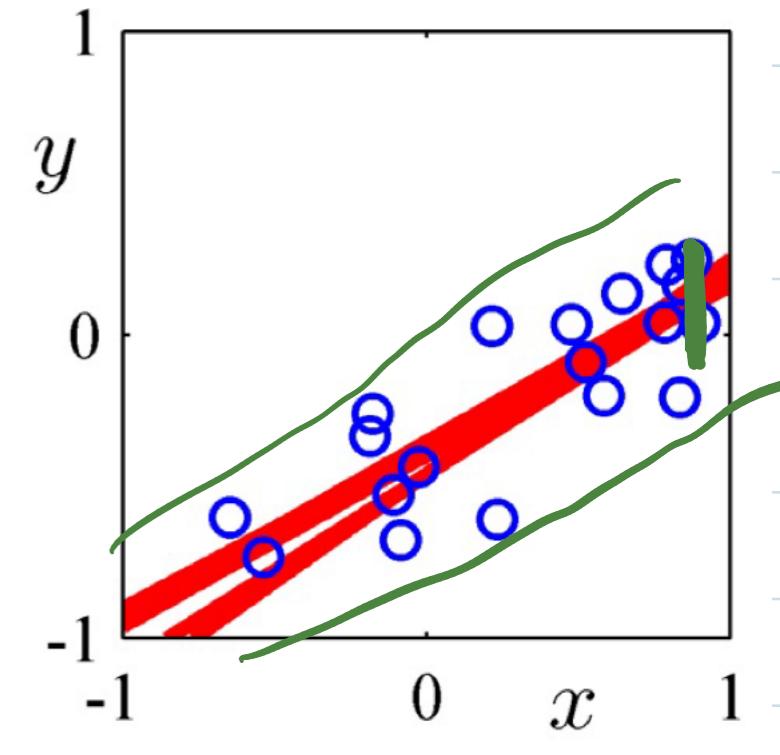
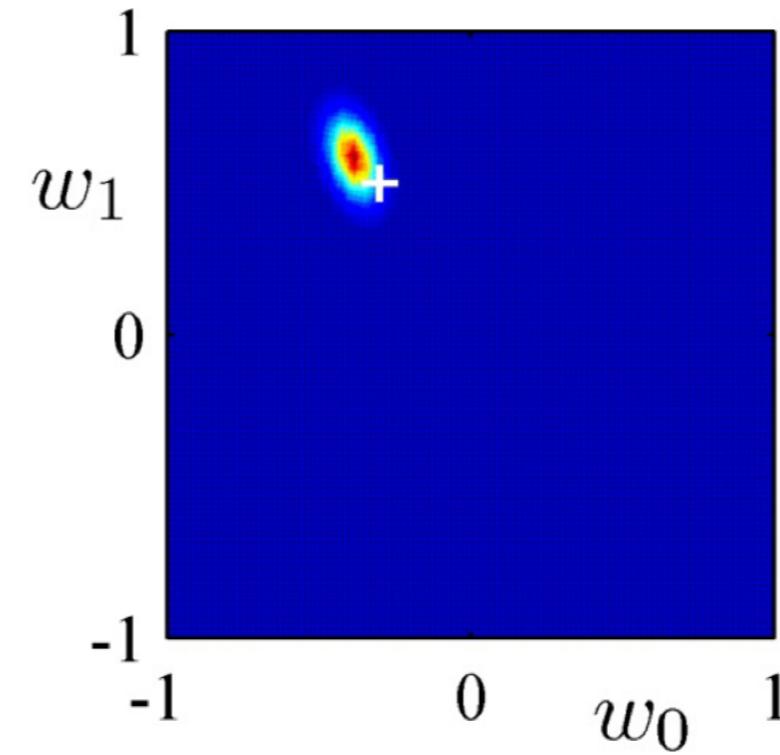
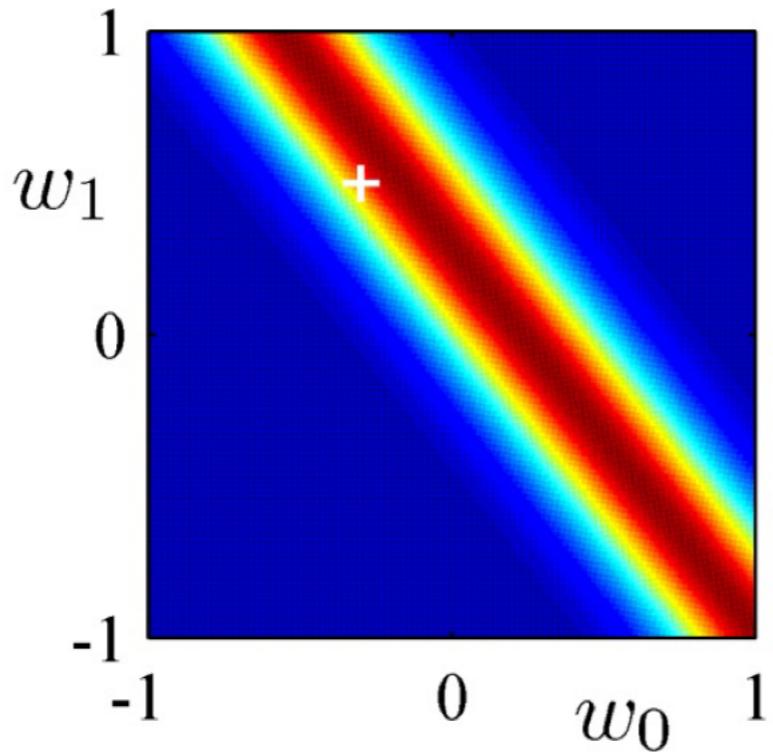


posterior



noise?

20 data  
points



# Predictive Distribution

We want to predict  $t$  for new values of  $X$

$$p(t | \bar{t}, \alpha, \beta) = \int p(t | w, \beta) p(w | \bar{t}, \alpha, \beta) dw$$

new output      vector of previous outputs      confidence in prior      sensor conf  
new  $t$  given parameters      probability of parameters given observed outputs

$$\sigma_N^2(x) = \frac{1}{\beta} + \phi(x)^T S_N^{-1} \phi(x)$$

Sensor Variance      parameter Variance

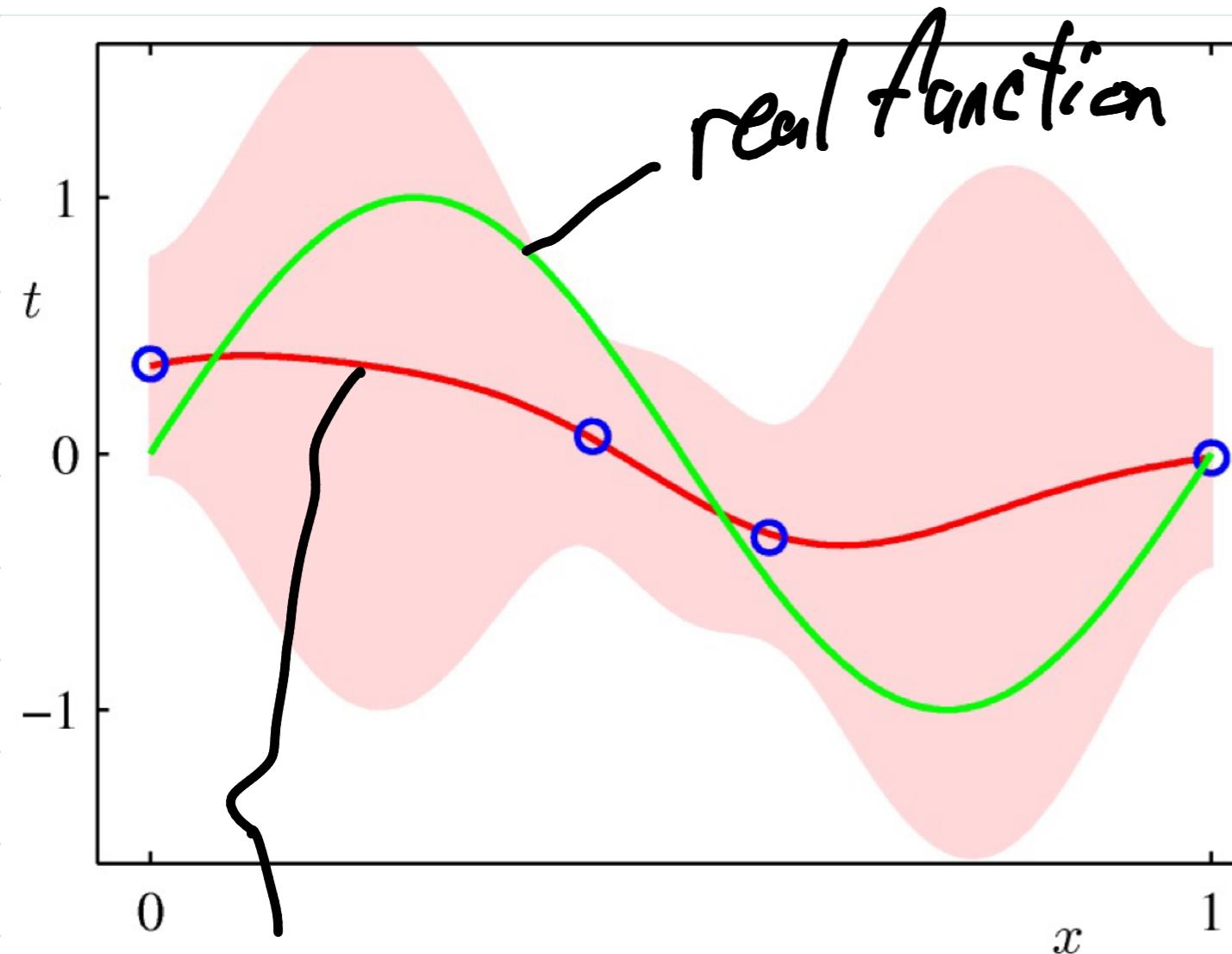
$$= N(t | m_N^T \phi(x), \sigma_N^2(x))$$

mean values of parameters      basis functions

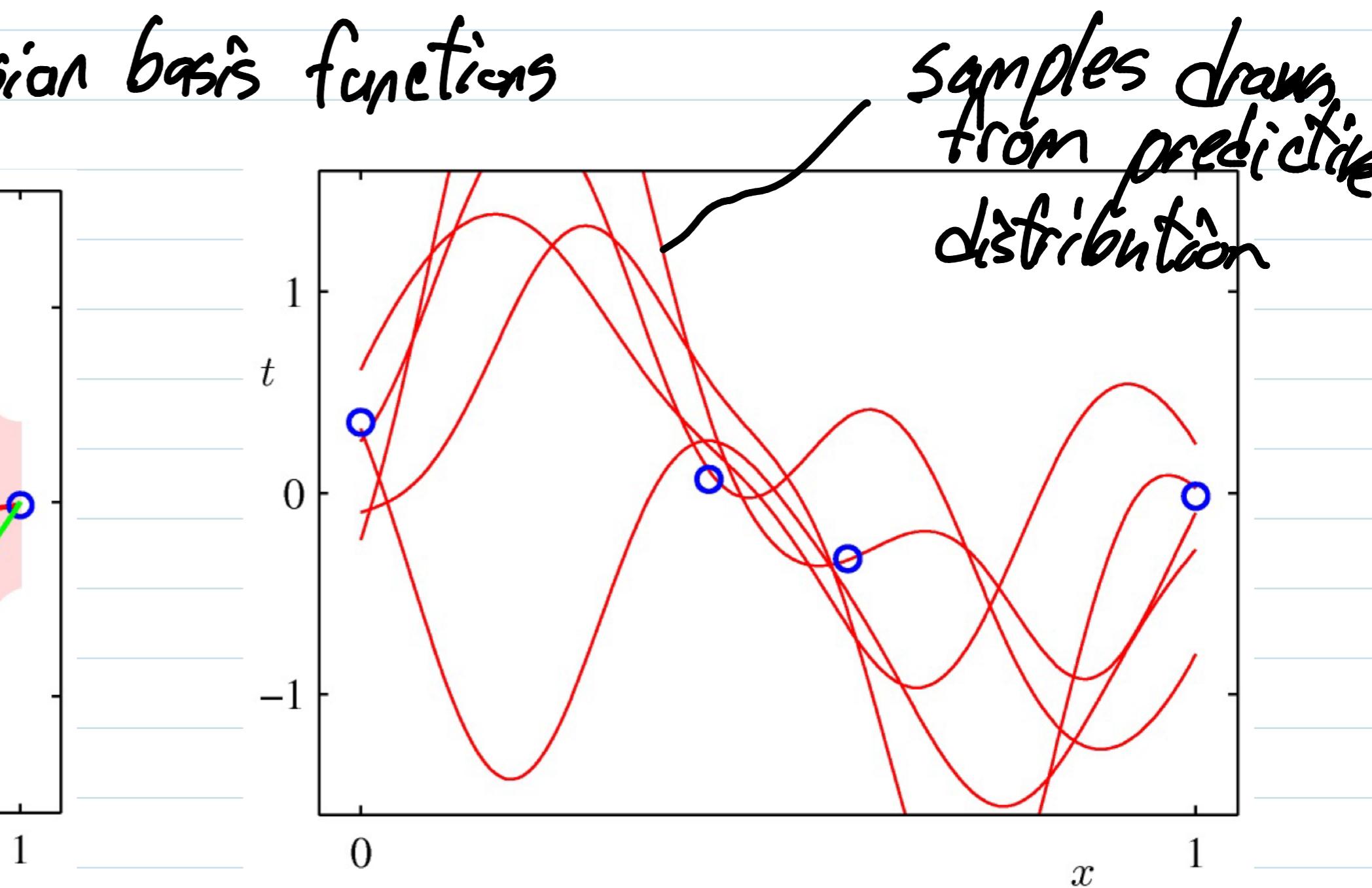
Example:

underlying process is sinusoidal data

4 data points, 9 Gaussian basis functions



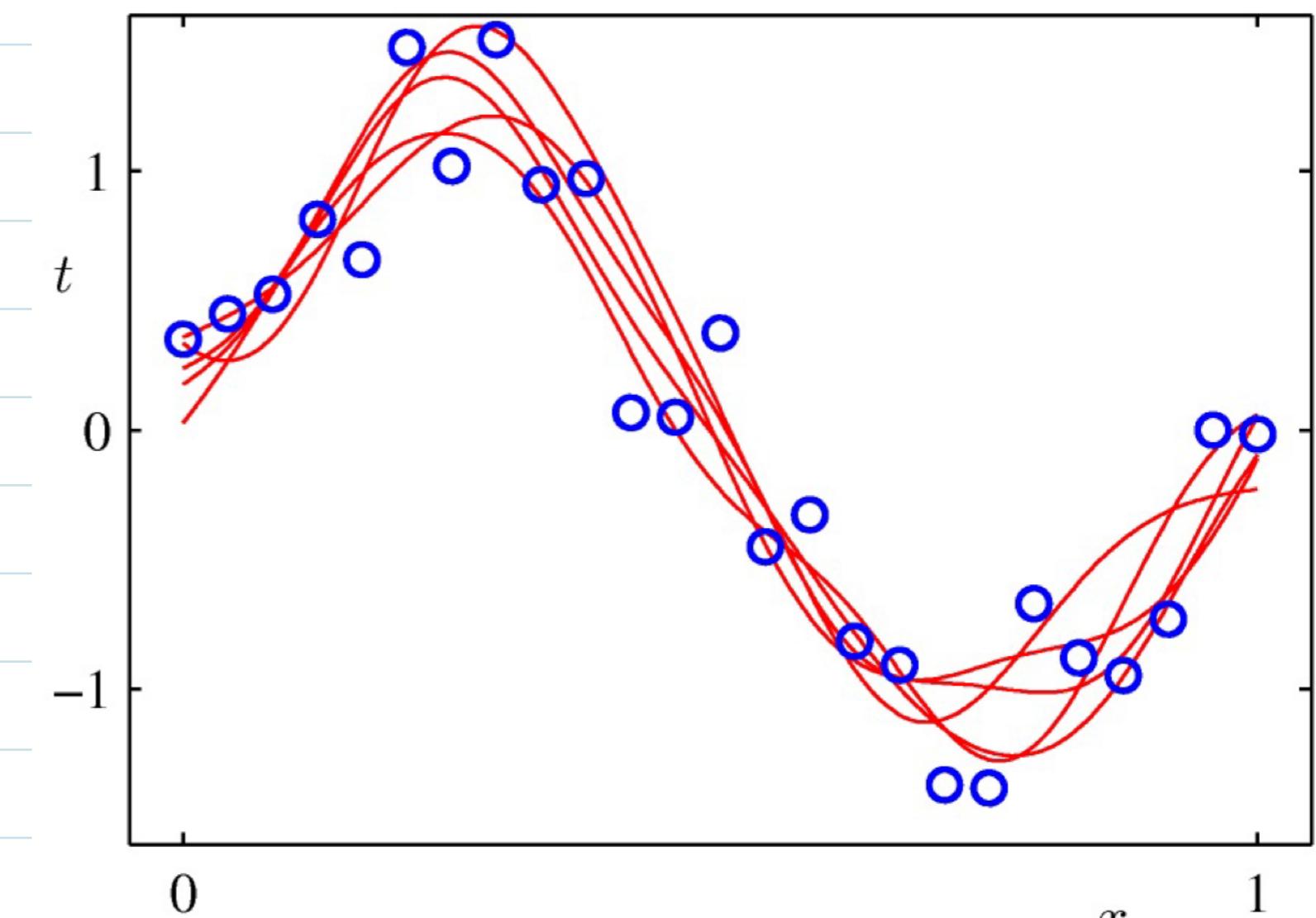
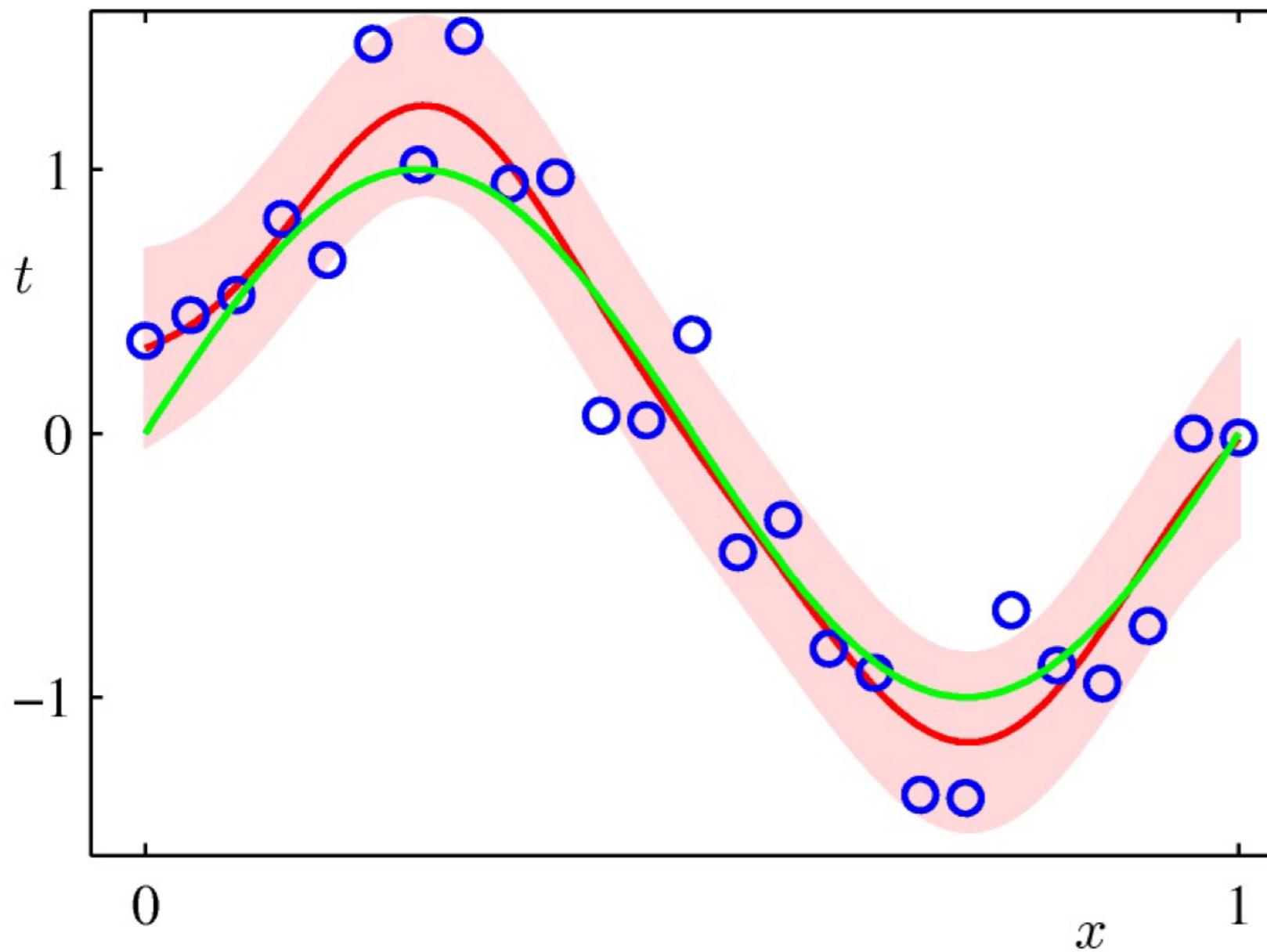
mean.  
of predictive  
distribution



$$t =$$

$$\sum_{i=1}^9 c_i \exp\left(-\frac{(x - \mu_i)^2}{\sigma_i^2}\right)$$

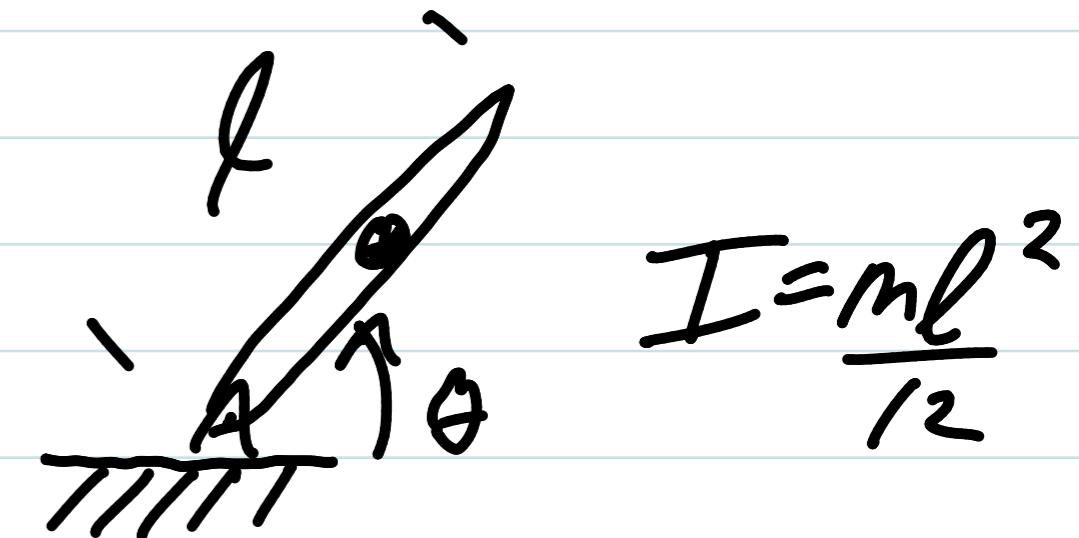
Now with 25 data points



# Pendulum Example

Equation of motion

$$\frac{ml^2}{3}\ddot{\theta} + \frac{mgl\cos\theta}{2} = \tau$$



Goal: find  $m$  and  $l$  given we can measure  $\theta, \dot{\theta}, \tau$

pick  $\ddot{\theta}$  as the output and  $\tau$  and  $\theta$  as input

$$x = [\tau \ \theta]^T \quad y = \ddot{\theta}$$

$$\ddot{\theta} = \frac{3\tau}{ml^2} - \frac{3g\cos\theta}{2l}$$

Write model in the form  $y = w^T \phi(x)$

$$w = \left[ \frac{1}{ml^2} \quad \frac{1}{l} \right]^T$$

$$\phi(x) = \left[ 3t \quad -\frac{3g \cos \theta}{2} \right]^T$$

$$y = w^T \phi(x) = \left[ \frac{1}{ml^2} \quad \frac{1}{l} \right] \begin{bmatrix} 3t \\ -\frac{3g \cos \theta}{2} \end{bmatrix} = \frac{3t}{ml^2} - \frac{3g \cos \theta}{2l}$$

To get started with Bayesian regression

- $\beta^{-1}$  sensor noise
- $m_0$  prior mean of parameters
- $s_0$  prior covariance of parameters

Suppose we move the pendulum around by commanding torques and measure  $\theta$  and  $\ddot{\theta}$ .  $N$  samples

$$\tau = [\ddot{\theta}_1, \ddot{\theta}_2, \dots, \ddot{\theta}_N]^T$$

$$\underline{I} = \begin{bmatrix} 3\tau_1 & -3g\cos\theta_1/2 \\ 3\tau_2 & -3g\cos\theta_2/2 \\ \vdots & \vdots \\ 3\tau_N & -3g\cos\theta_N/2 \end{bmatrix}$$

posterior confidence

$$S_N^{-1} = S_0^{-1} + \beta \underline{I}^T \underline{I}$$

posterior mean of parameters

$$m_N = S_N (S_0^{-1} m_0 + \beta \underline{I}^T \underline{I})$$

$$= \left[ \frac{1}{ml^2}, \frac{1}{l} \right]^T$$

the answer

or a predictive distribution of  $\ddot{\theta}$

$$p(\ddot{\theta} | t, \alpha, \beta) = N(\ddot{\theta} | m_N^\top \phi(x), \sigma_N^2(x))$$

$$\sigma_N^2(x) = \frac{1}{\beta} + \phi^\top(x) S_N \phi(x)$$

Other things to do!

- 1) predict  $\varepsilon$  given  $\theta$  and  $\ddot{\theta}$  by writing the EoM with different basis functions (inverse dynamics)
- 2) write discrete version of EoM

e.g.  $\dot{\theta}_{i+1} = \dot{\theta}_i + \ddot{\theta}_i \Delta t$

$$\dot{\theta}_{i+1} = \dot{\theta}_i + \underbrace{\ddot{\theta}_i}_{\text{function of Parameters}} \Delta t$$

function of Parameters

$$W = \begin{bmatrix} \frac{1}{ml^2} & \frac{1}{l} \\ \frac{1}{l} & 1 \end{bmatrix} = \begin{bmatrix} 1.63 & 1.66 \end{bmatrix}$$