

# LINEAR REGRESSION

# Credits

2

Probability & Bayesian Inference

- Some of these slides were sourced and/or modified from:
  - Christopher Bishop, Microsoft UK

# Relevant Problems from Murphy

3

Probability & Bayesian Inference

- 7.4, 7.6, 7.7, **7.9**
- Please do 7.9 at least. We will discuss the solution in class.

# Linear Regression Topics

4

Probability & Bayesian Inference

- What is linear regression?
- Example: polynomial curve fitting
- Other basis families
- Solving linear regression problems
- Regularized regression
- Multiple linear regression
- Bayesian linear regression

# What is Linear Regression?

5

Probability & Bayesian Inference

- In classification, we seek to identify the **categorical** class  $C_k$  associate with a given input vector  $x$ .
- In regression, we seek to identify (or **estimate**) a **continuous** variable  $y$  associated with a given input vector  $x$ .
- $y$  is called the **dependent variable**.
- $x$  is called the **independent variable**.
- If  $y$  is a vector, we call this multiple regression.
- We will focus on the case where  $y$  is a scalar.
- Notation:
  - $y$  will denote the continuous model of the dependent variable
  - $t$  will denote discrete noisy observations of the dependent variable (sometimes called the **target variable**).

# Where is the Linear in Linear Regression?

6

Probability & Bayesian Inference

- In regression we assume that  $y$  is a function of  $\mathbf{x}$ .  
The exact nature of this function is governed by an unknown parameter vector  $\mathbf{w}$ :

$$y = y(\mathbf{x}, \mathbf{w})$$

- The regression is linear if  $y$  is linear in  $\mathbf{w}$ . In other words, we can express  $y$  as

$$y = \mathbf{w}^t \phi(\mathbf{x})$$

where

$\phi(\mathbf{x})$  is some (potentially nonlinear) function of  $\mathbf{x}$ .

# Linear Basis Function Models

7

Probability &amp; Bayesian Inference

- Generally

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

- where  $\phi_j(\mathbf{x})$  are known as *basis functions*.
- Typically,  $\Phi_0(\mathbf{x}) = 1$ , so that  $w_0$  acts as a bias.
- In the simplest case, we use linear basis functions :  
 $\Phi_d(\mathbf{x}) = x_d$ .

# Linear Regression Topics

8

Probability & Bayesian Inference

- What is linear regression?
- **Example: polynomial curve fitting**
- Other basis families
- Solving linear regression problems
- Regularized regression
- Multiple linear regression
- Bayesian linear regression

# Example: Polynomial Bases

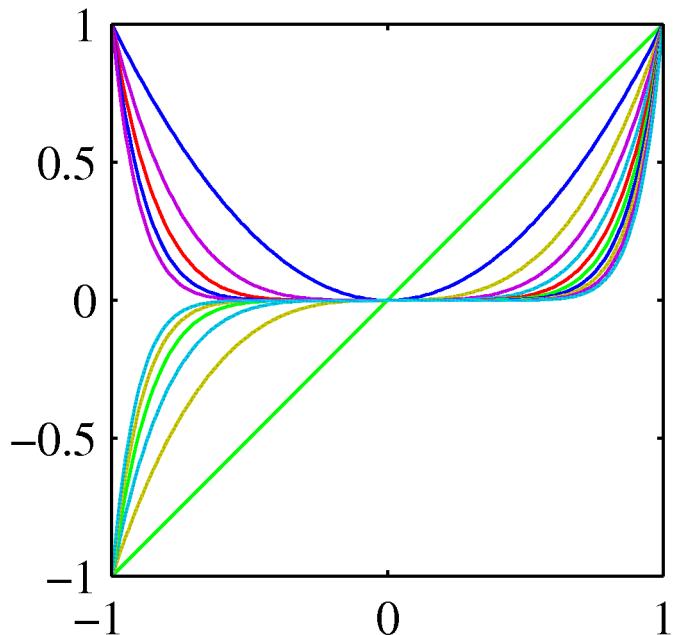
9

Probability & Bayesian Inference

- Polynomial basis functions:

$$\phi_j(x) = x^j.$$

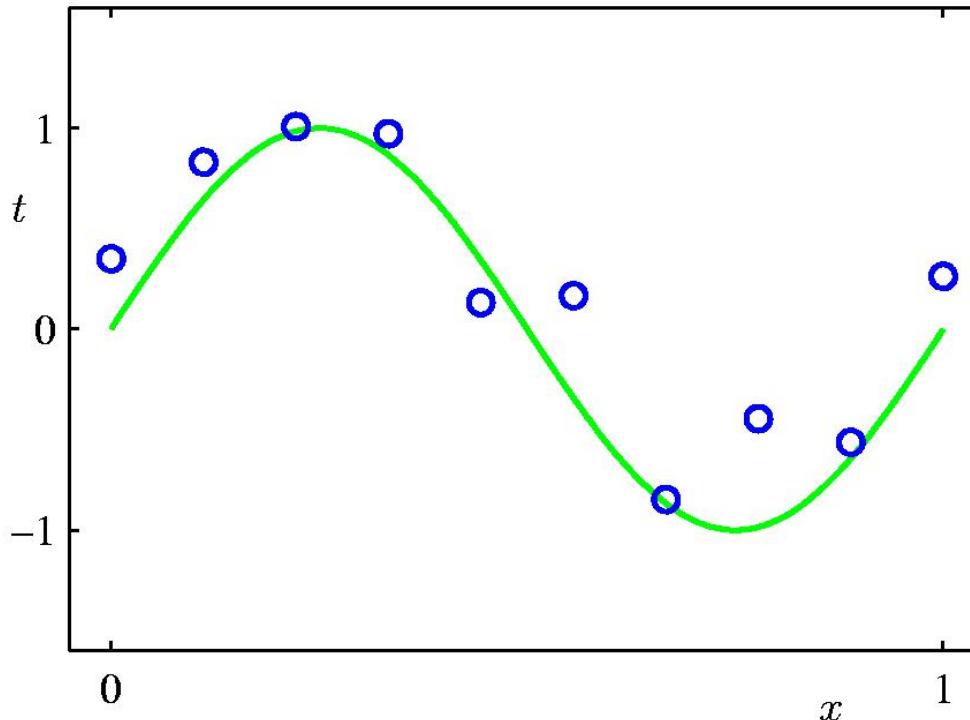
- These are global
  - a small change in  $x$  affects all basis functions.
  - A small change in a basis function affects  $y$  for all  $x$ .



# Example: Polynomial Curve Fitting

10

Probability &amp; Bayesian Inference

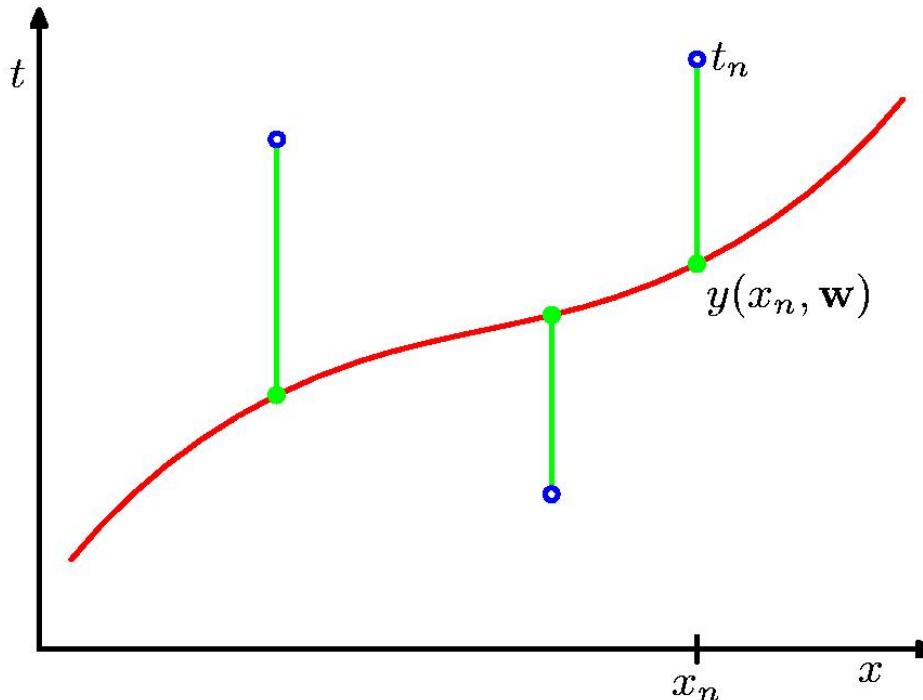


$$y(x, \mathbf{w}) = w_0 + w_1 x + w_2 x^2 + \dots + w_M x^M = \sum_{j=0}^M w_j x^j$$

# Sum-of-Squares Error Function

11

Probability & Bayesian Inference

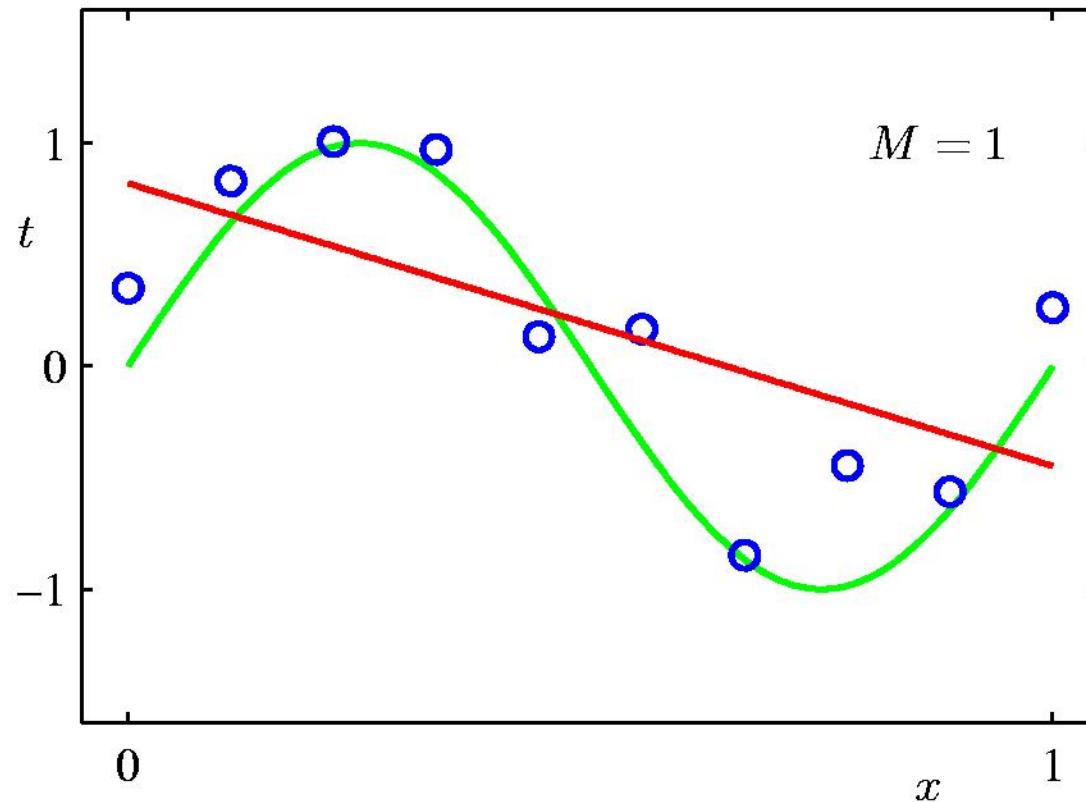


$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2$$

# 1<sup>st</sup> Order Polynomial

12

Probability & Bayesian Inference

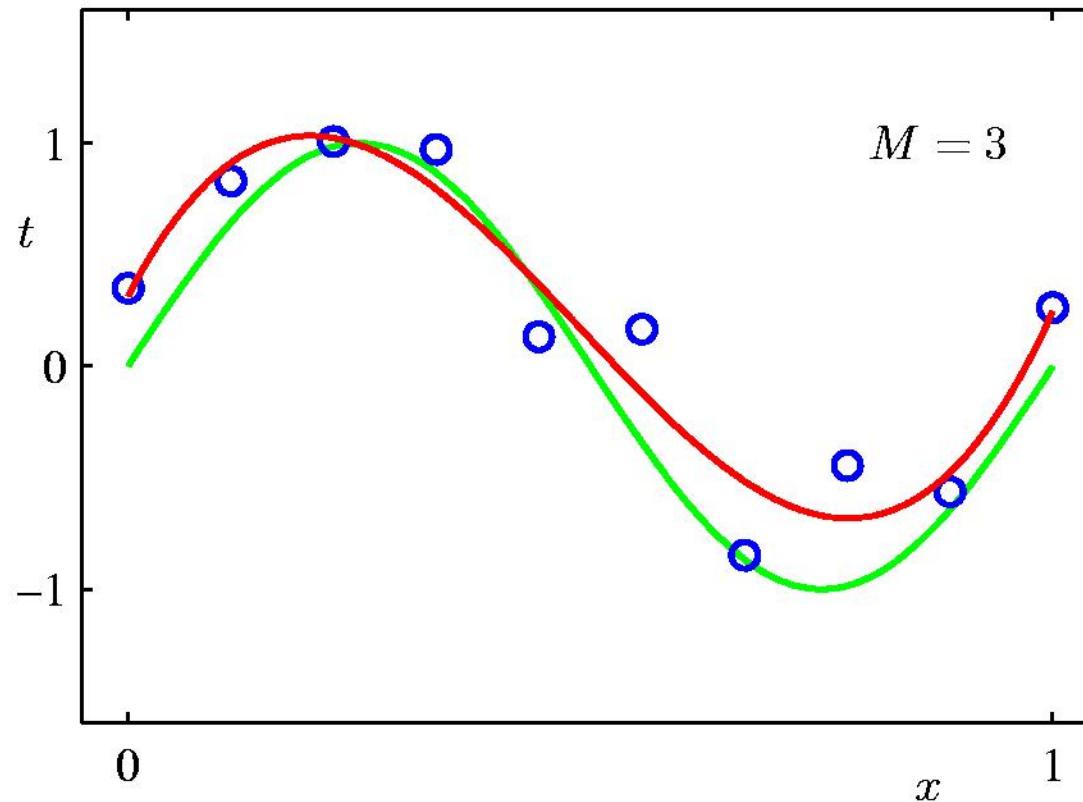


$M = 1$

# 3<sup>rd</sup> Order Polynomial

13

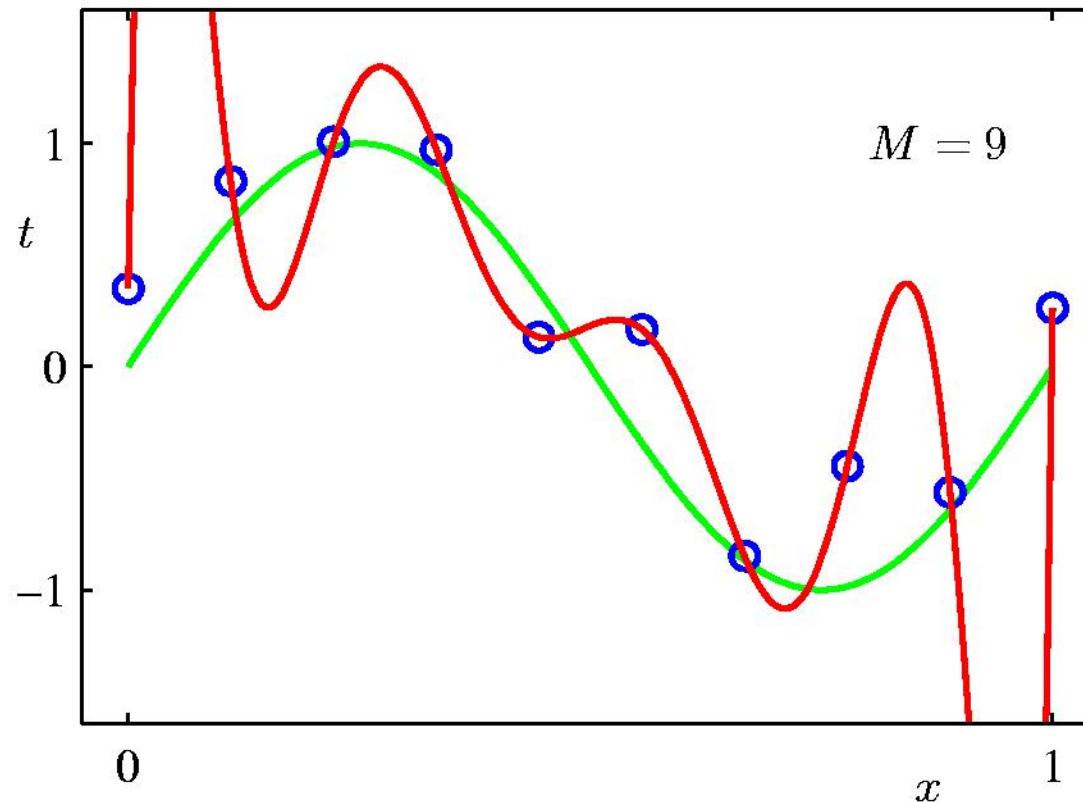
Probability & Bayesian Inference



# 9<sup>th</sup> Order Polynomial

14

Probability & Bayesian Inference



# Regularization

15

Probability & Bayesian Inference

- Penalize large coefficient values

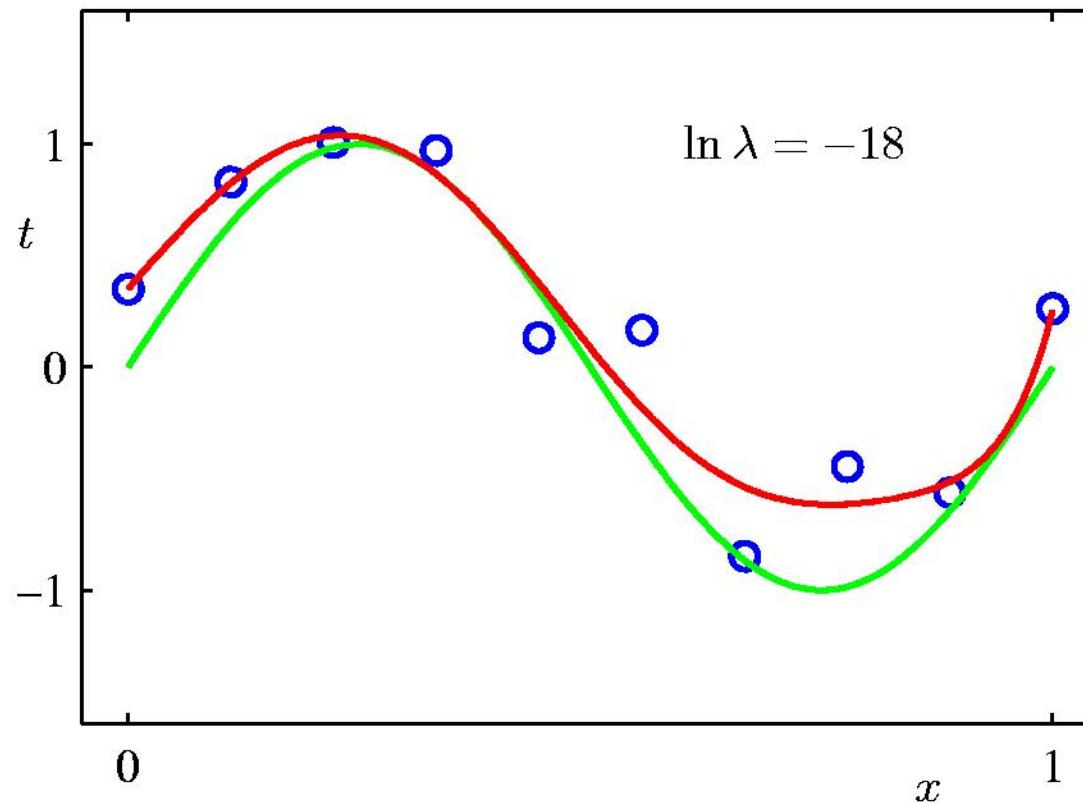
$$\tilde{E}(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

# Regularization

16

Probability & Bayesian Inference

## 9<sup>th</sup> Order Polynomial

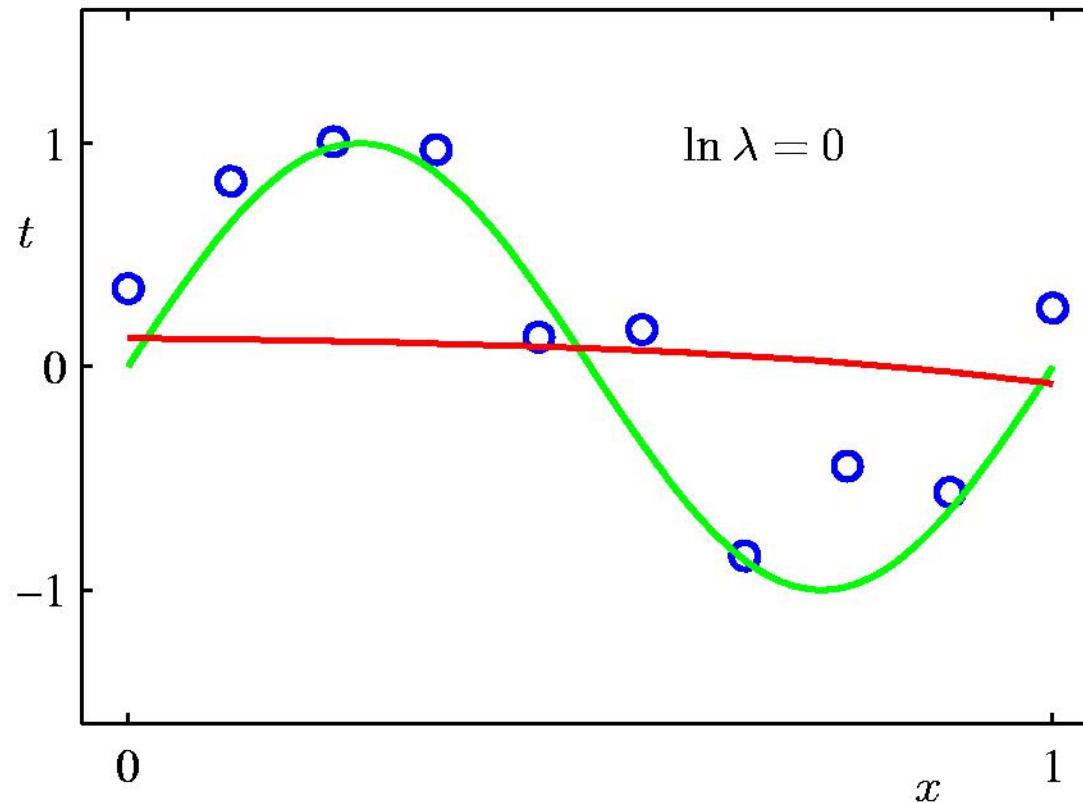


# Regularization

17

Probability &amp; Bayesian Inference

## 9<sup>th</sup> Order Polynomial

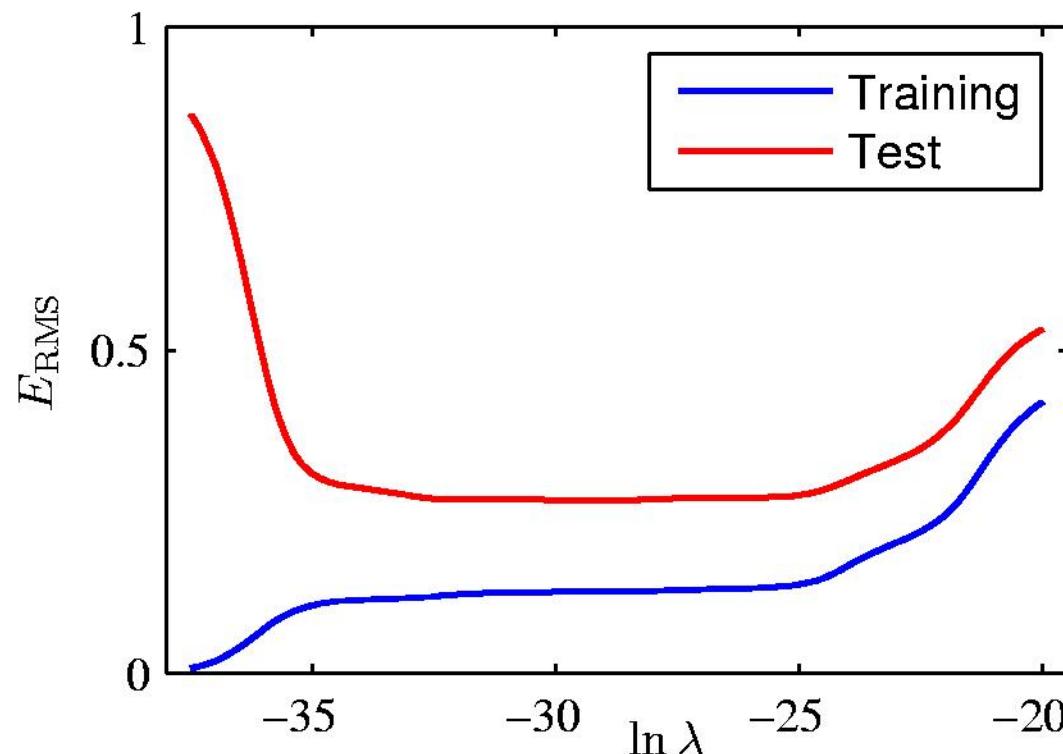


# Regularization

18

Probability & Bayesian Inference

## 9<sup>th</sup> Order Polynomial

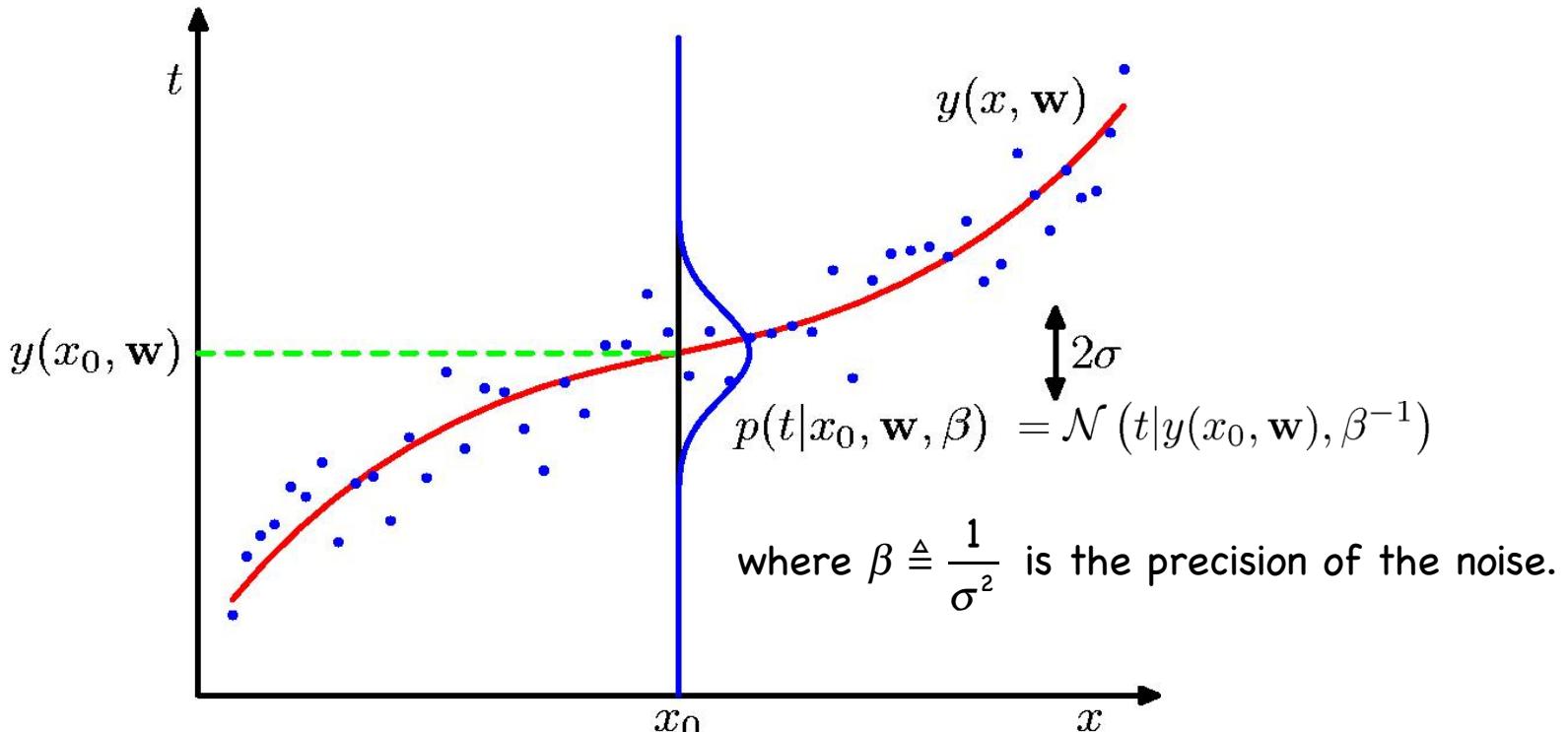


# Probabilistic View of Curve Fitting

19

Probability & Bayesian Inference

- Why least squares?
- Model noise (deviation of data from model) as Gaussian i.i.d.



# Maximum Likelihood

20

Probability &amp; Bayesian Inference

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

- We determine  $\mathbf{w}_{ML}$  by minimizing the squared error  $E(\mathbf{w})$ .

$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \underbrace{\sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2}_{\beta E(\mathbf{w})} + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

- Thus least-squares regression reflects an assumption that the noise is i.i.d. Gaussian.

# Maximum Likelihood

21

Probability &amp; Bayesian Inference

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

- We determine  $\mathbf{w}_{ML}$  by minimizing the squared error  $E(\mathbf{w})$ .

$$\ln p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta) = -\frac{\beta}{2} \underbrace{\sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2}_{\beta E(\mathbf{w})} + \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi)$$

- Now given  $\mathbf{w}_{ML}$ , we can estimate the variance of the noise:

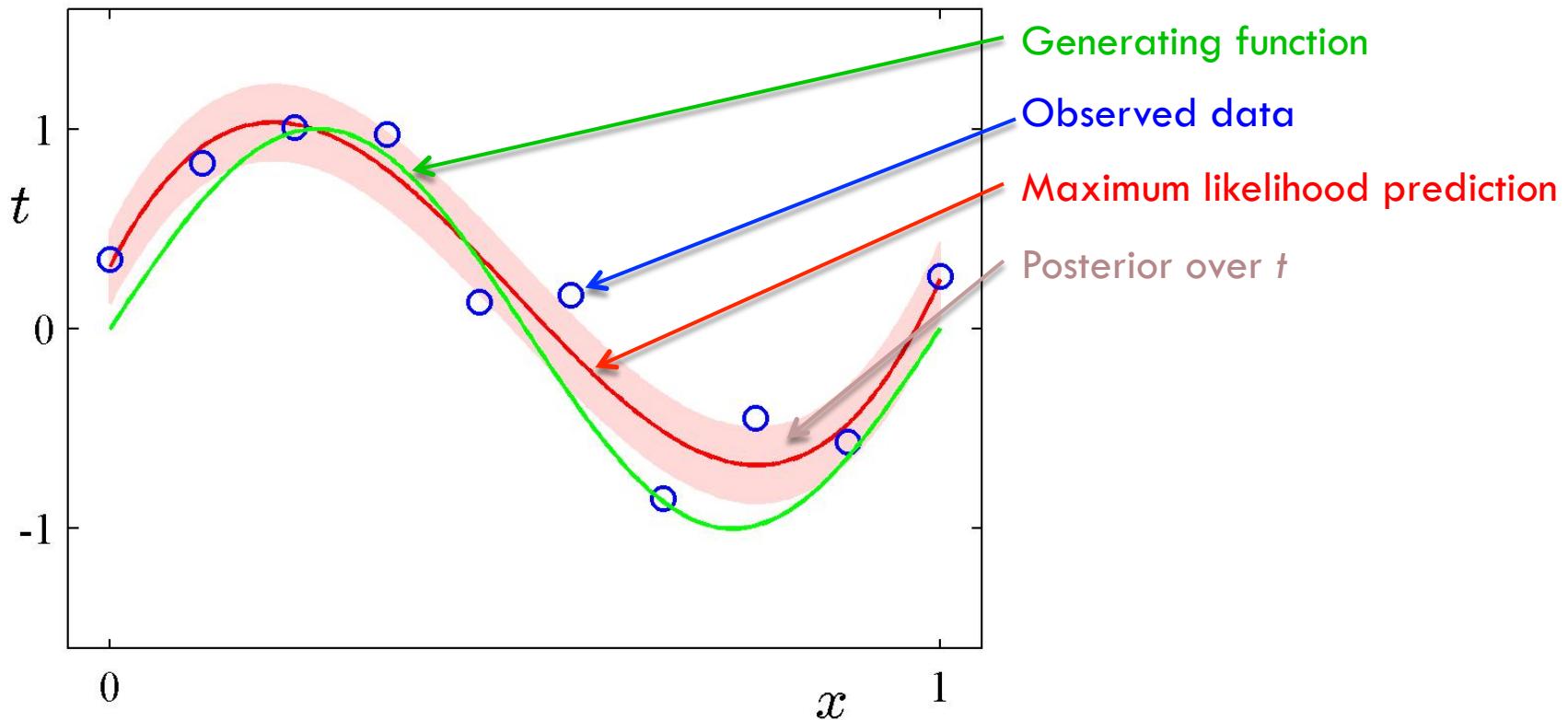
$$\frac{1}{\beta_{ML}} = \frac{1}{N} \sum_{n=1}^N \{y(x_n, \mathbf{w}_{ML}) - t_n\}^2$$

# Predictive Distribution

22

Probability &amp; Bayesian Inference

$$p(t|x, \mathbf{w}_{\text{ML}}, \beta_{\text{ML}}) = \mathcal{N}(t|y(x, \mathbf{w}_{\text{ML}}), \beta_{\text{ML}}^{-1})$$



# MAP: A Step towards Bayes

23

Probability &amp; Bayesian Inference

- Prior knowledge about probable values of  $\mathbf{w}$  can be incorporated into the regression:

$$p(\mathbf{w}|\alpha) = \mathcal{N}(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I}) = \left(\frac{\alpha}{2\pi}\right)^{(M+1)/2} \exp\left\{-\frac{\alpha}{2}\mathbf{w}^T\mathbf{w}\right\}$$

- Now the posterior over  $\mathbf{w}$  is proportional to the product of the likelihood times the prior:

$$p(\mathbf{w}|\mathbf{x}, \mathbf{t}, \alpha, \beta) \propto p(\mathbf{t}|\mathbf{x}, \mathbf{w}, \beta)p(\mathbf{w}|\alpha)$$

- The result is to introduce a new quadratic term in  $\mathbf{w}$  into the error function to be minimized:

$$\beta\tilde{E}(\mathbf{w}) = \frac{\beta}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\alpha}{2}\mathbf{w}^T\mathbf{w}$$

- Thus regularized (ridge) regression reflects a 0-mean isotropic Gaussian prior on the weights.

# Linear Regression Topics

24

Probability & Bayesian Inference

- What is linear regression?
- Example: polynomial curve fitting
- **Other basis families**
- Solving linear regression problems
- Regularized regression
- Multiple linear regression
- Bayesian linear regression

# Gaussian Bases

25

Probability & Bayesian Inference

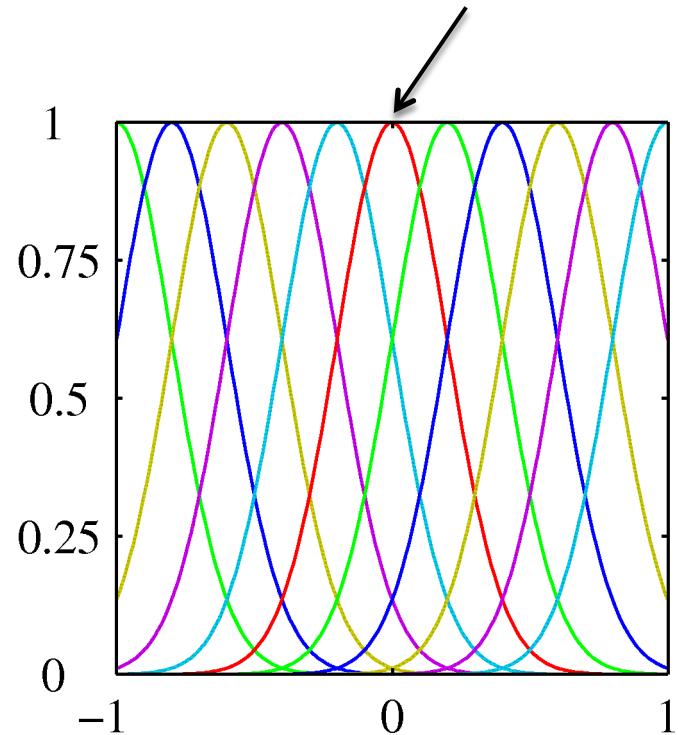
- Gaussian basis functions:

$$\phi_j(x) = \exp \left\{ -\frac{(x - \mu_j)^2}{2s^2} \right\}$$

- These are local:

- a small change in  $x$  affects only nearby basis functions.
- a small change in a basis function affects  $y$  only for nearby  $x$ .
- $\mu_j$  and  $s$  control location and scale (width).

Think of these as interpolation functions.



# Linear Regression Topics

26

Probability & Bayesian Inference

- What is linear regression?
- Example: polynomial curve fitting
- Other basis families
- **Solving linear regression problems**
- Regularized regression
- Multiple linear regression
- Bayesian linear regression

# Maximum Likelihood and Linear Least Squares

27

Probability & Bayesian Inference

- Assume observations from a deterministic function with added Gaussian noise:

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

- which is the same as saying,

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

- where

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

# Maximum Likelihood and Linear Least Squares

28

Probability & Bayesian Inference

- Given observed inputs,  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , and targets,  $\mathbf{t} = [t_1, \dots, t_N]^T$  we obtain the likelihood function

$$p(\mathbf{t}|\mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}).$$

# Maximum Likelihood and Linear Least Squares

29

Probability & Bayesian Inference

- Taking the logarithm, we get

$$\begin{aligned}\ln p(\mathbf{t}|\mathbf{w}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}) \\ &= \frac{N}{2} \ln \beta - \frac{N}{2} \ln(2\pi) - \beta E_D(\mathbf{w})\end{aligned}$$

- where

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2$$

- is the sum-of-squares error.

# Maximum Likelihood and Least Squares

30

Probability &amp; Bayesian Inference

- Computing the gradient and setting it to zero yields

$$\nabla_{\mathbf{w}} \ln p(\mathbf{t}|\mathbf{w}, \beta) = \beta \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\} \boldsymbol{\phi}(\mathbf{x}_n)^T = \mathbf{0}.$$

- Solving for  $\mathbf{w}$ , we get

$$\mathbf{w}_{\text{ML}} = \overbrace{\left( \boldsymbol{\Phi}^T \boldsymbol{\Phi} \right)^{-1} \boldsymbol{\Phi}^T \mathbf{t}}$$

The Moore-Penrose  
pseudo-inverse,  $\boldsymbol{\Phi}^\dagger$ .

- where

$$\boldsymbol{\Phi} = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}.$$

# Linear Regression Topics

31

Probability & Bayesian Inference

- What is linear regression?
- Example: polynomial curve fitting
- Other basis families
- Solving linear regression problems
- **Regularized regression**
- Multiple linear regression
- Bayesian linear regression

# Regularized Least Squares

32

Probability & Bayesian Inference

- Consider the error function:

$$E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$$

Data term + Regularization term

$\lambda$  is called the regularization coefficient.

- With the sum-of-squares error function and a quadratic regularizer, we get

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w}$$

- which is minimized by

$$\mathbf{w} = \left( \lambda \mathbf{I} + \boldsymbol{\Phi}^T \boldsymbol{\Phi} \right)^{-1} \boldsymbol{\Phi}^T \mathbf{t}.$$

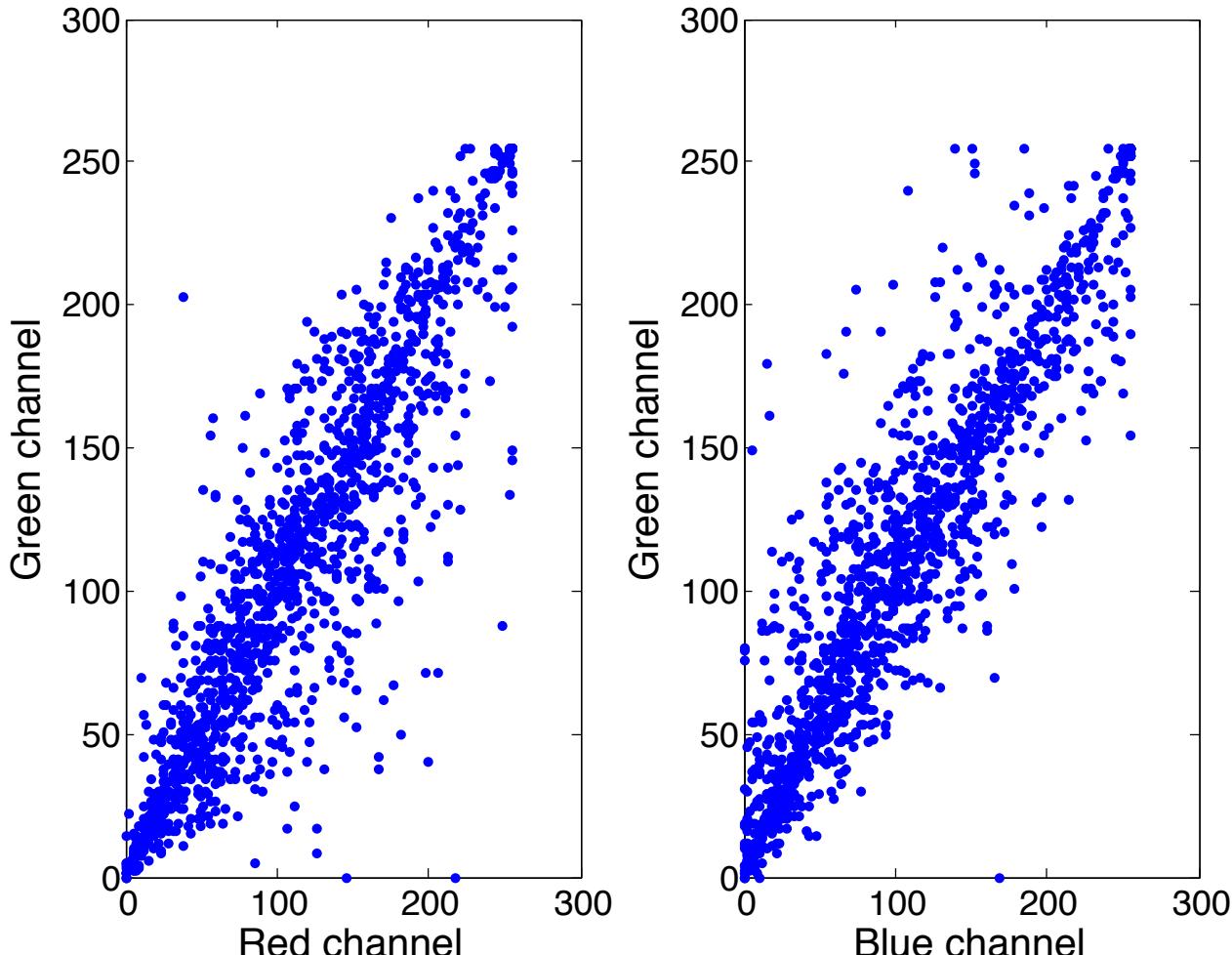
Thus the name ‘ridge regression’



# Application: Colour Restoration

33

Probability & Bayesian Inference



# Application: Colour Restoration

34

Probability & Bayesian Inference

Original Image



Red and Blue Channels Only



Predicted Image



Remove  
Green



Restore  
Green

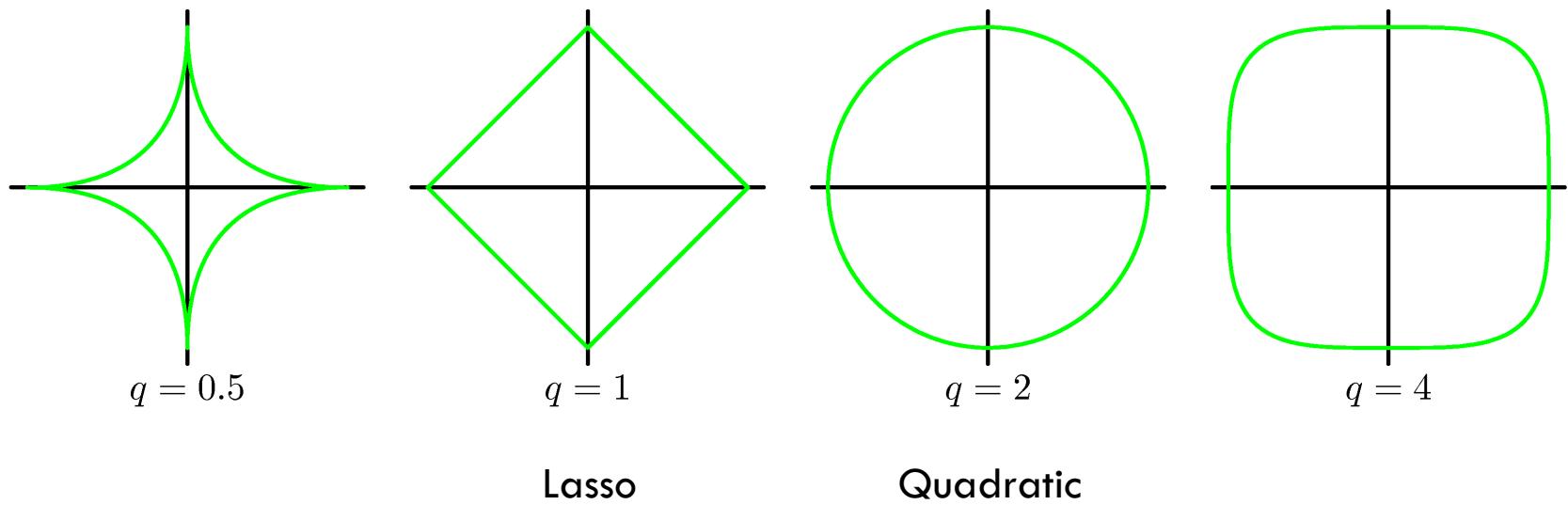
# Regularized Least Squares

35

Probability & Bayesian Inference

- A more general regularizer:

$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$



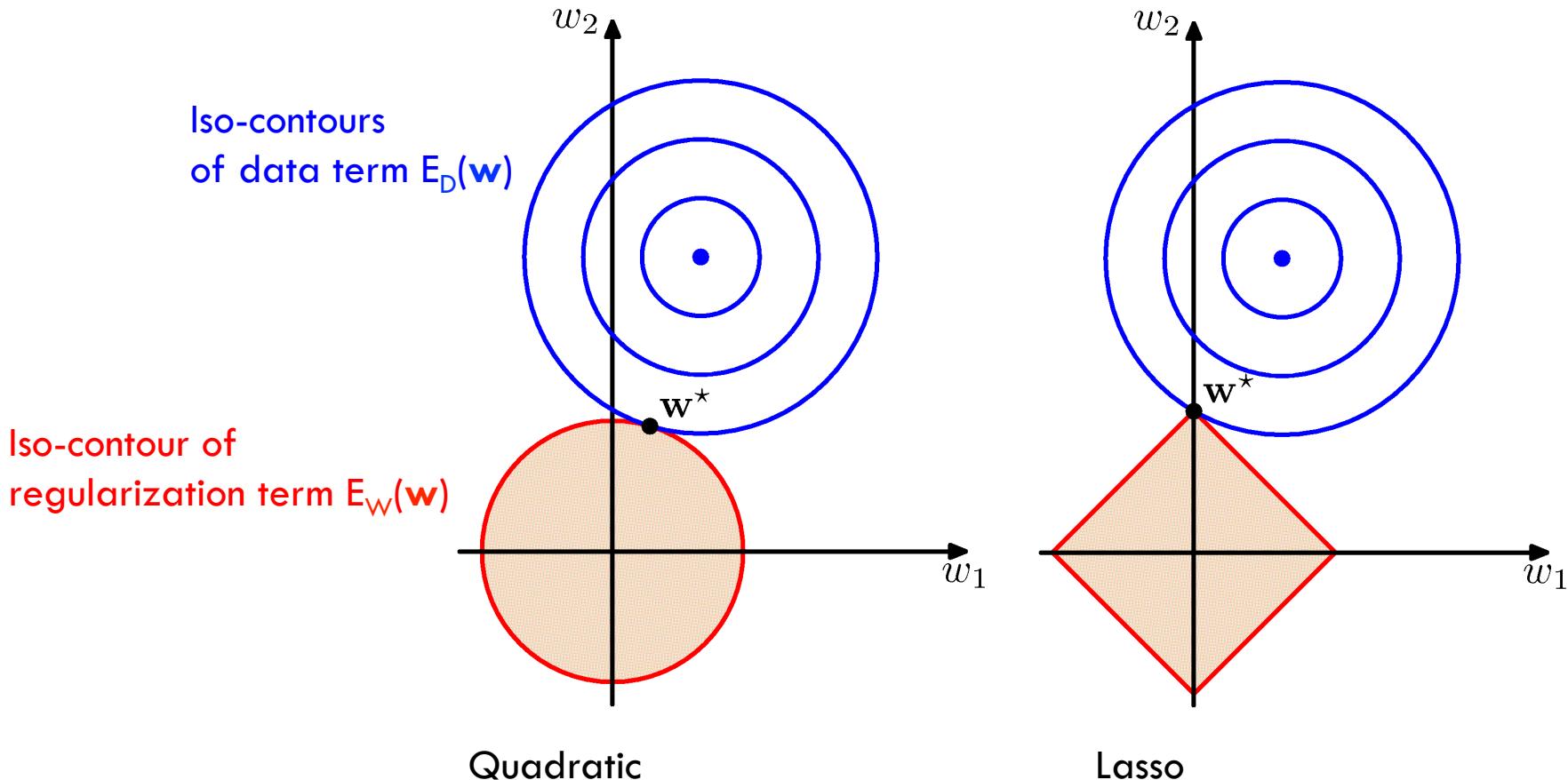
(Least absolute shrinkage and selection operator)

# Regularized Least Squares

36

Probability & Bayesian Inference

- Lasso generates sparse solutions.



# Solving Regularized Systems

37

Probability & Bayesian Inference

- Quadratic regularization has the advantage that the solution is closed form.
- Non-quadratic regularizers generally do not have closed form solutions
- Lasso can be framed as minimizing a quadratic error with linear constraints, and thus represents a convex optimization problem that can be solved by quadratic programming or other convex optimization methods.
- We will discuss quadratic programming when we cover SVMs

# Linear Regression Topics

38

Probability & Bayesian Inference

- What is linear regression?
- Example: polynomial curve fitting
- Other basis families
- Solving linear regression problems
- Regularized regression
- **Multiple linear regression**
- Bayesian linear regression

# Multiple Outputs

39

Probability &amp; Bayesian Inference

- Analogous to the single output case we have:

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

- Given observed inputs  $\mathbf{X} = \{\mathbf{x}_1, \dots, \mathbf{x}_N\}$ , and targets  $\mathbf{T} = [\mathbf{t}_1, \dots, \mathbf{t}_N]^T$

we obtain the log likelihood function

$$\begin{aligned} \ln p(\mathbf{T}|\mathbf{X}, \mathbf{W}, \beta) &= \sum_{n=1}^N \ln \mathcal{N}(\mathbf{t}_n | \mathbf{W}^T \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1}\mathbf{I}) \\ &= \frac{NK}{2} \ln \left( \frac{\beta}{2\pi} \right) - \frac{\beta}{2} \sum_{n=1}^N \|\mathbf{t}_n - \mathbf{W}^T \boldsymbol{\phi}(\mathbf{x}_n)\|^2. \end{aligned}$$

# Multiple Outputs

40

Probability &amp; Bayesian Inference

- Maximizing with respect to  $\mathbf{W}$ , we obtain

$$\mathbf{W}_{\text{ML}} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{T}.$$

- If we consider a single target variable,  $\mathbf{t}_k$ , we see that

$$\mathbf{w}_k = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}_k = \Phi^\dagger \mathbf{t}_k$$

- where  $\mathbf{t}_k = [t_{1k}, \dots, t_{Nk}]^T$ , which is identical with the single output case.

# Some Useful MATLAB Functions

41

Probability & Bayesian Inference

- **polyfit**
  - Least-squares fit of a polynomial of specified order to given data
- **regress**
  - More general function that computes linear weights for least-squares fit

# Linear Regression Topics

42

Probability & Bayesian Inference

- What is linear regression?
- Example: polynomial curve fitting
- Other basis families
- Solving linear regression problems
- Regularized regression
- Multiple linear regression
- **Bayesian linear regression**

# Bayesian Linear Regression



Rev. Thomas Bayes, 1702 - 1761

# Bayesian Linear Regression

44

Probability & Bayesian Inference

- In least-squares, we determine the weights  $w$  that minimize the least squared error between the model and the training data.
- This can result in overlearning!
- Overlearning can be reduced by adding a regularizing term to the error function being minimized.
- Under specific conditions this is equivalent to a Bayesian approach, where we specify a prior distribution over the weight vector.

# Bayesian Linear Regression

45

Probability &amp; Bayesian Inference

- Define a conjugate prior over  $\mathbf{w}$ :

$$p(\mathbf{w}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_0, \mathbf{S}_0).$$

- Combining this with the likelihood function and matching terms, we obtain

$$p(\mathbf{w} | \mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N)$$

- where

$$\begin{aligned}\mathbf{m}_N &= \mathbf{S}_N \left( \mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \boldsymbol{\Phi}^T \mathbf{t} \right) \\ \mathbf{S}_N^{-1} &= \mathbf{S}_0^{-1} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}.\end{aligned}$$

# Bayesian Linear Regression

46

Probability & Bayesian Inference

- A common choice for the prior is

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

- for which

$$\begin{aligned}\mathbf{m}_N &= \beta \mathbf{S}_N \Phi^T \mathbf{t} \\ \mathbf{S}_N^{-1} &= \alpha \mathbf{I} + \beta \Phi^T \Phi.\end{aligned}$$

- Thus  $\mathbf{m}_N$  represents the ridge regression solution with

$$\lambda = \alpha / \beta$$

- Next we consider an example ...

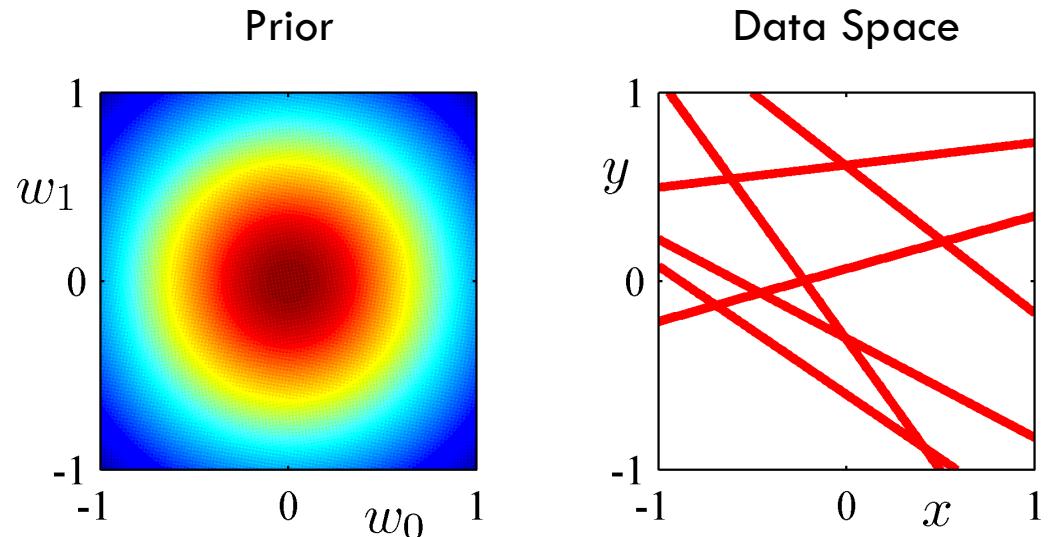
# Bayesian Linear Regression

47

Probability & Bayesian Inference

## Example: fitting a straight line

0 data points observed

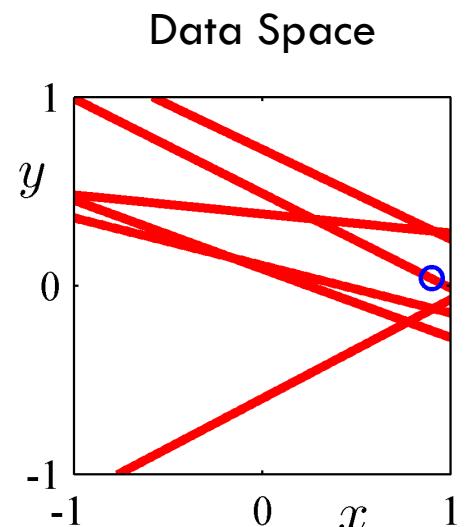
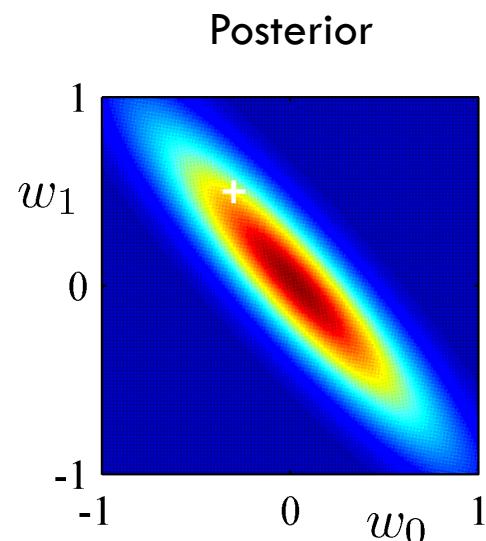
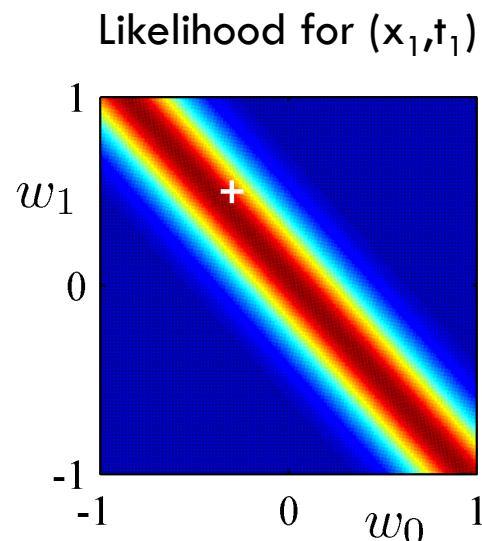


# Bayesian Linear Regression

48

Probability & Bayesian Inference

1 data point observed

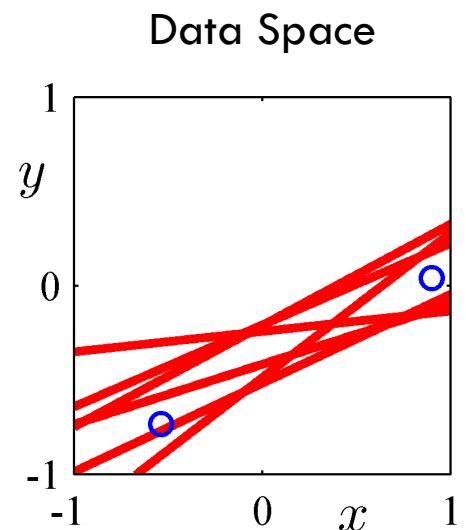
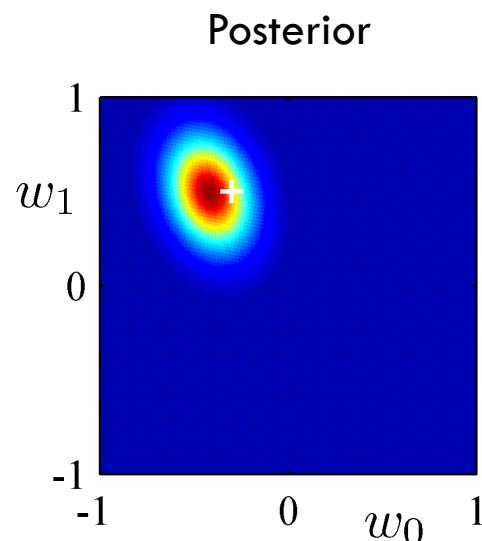
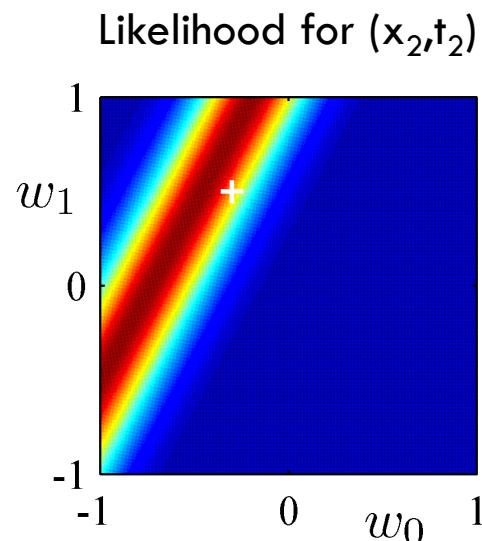


# Bayesian Linear Regression

49

Probability & Bayesian Inference

2 data points observed

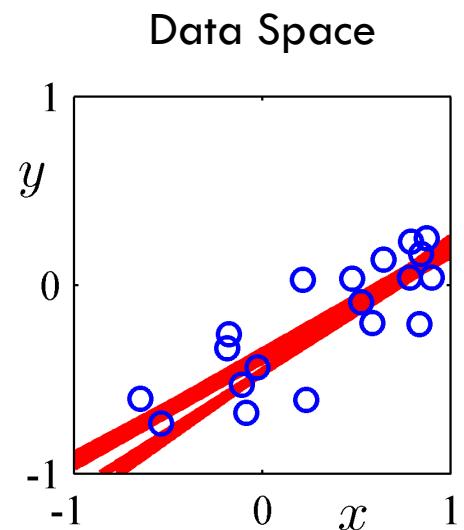
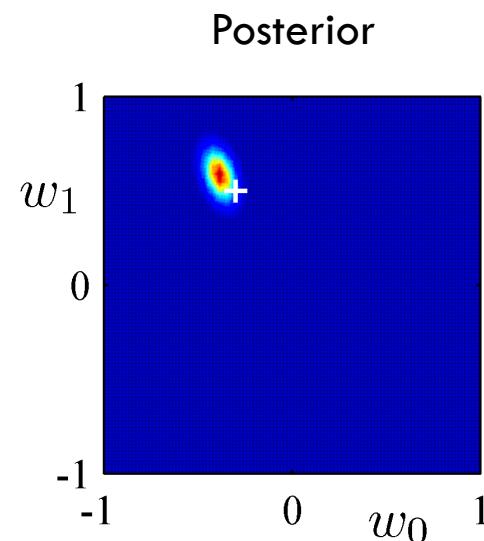
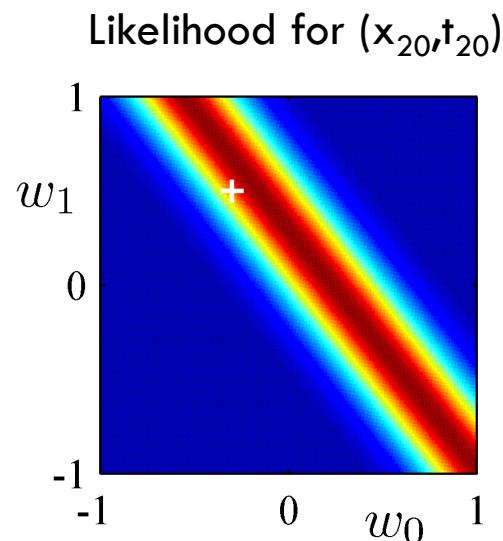


# Bayesian Linear Regression

50

Probability & Bayesian Inference

20 data points observed



# Bayesian Prediction

51

Probability & Bayesian Inference

- In least-squares, or regularized least-squares, we determine specific weights  $w$  that allow us to predict a specific value  $y(x, w)$  for every observed input  $x$ .
- However, our estimate of the weight vector  $w$  will never be perfect! This will introduce error into our prediction.
- In Bayesian prediction, we model the posterior distribution over our predictions, taking into account our uncertainty in model parameters.

# Predictive Distribution

52

Probability &amp; Bayesian Inference

- Predict  $t$  for new values of  $\mathbf{x}$  by integrating over  $\mathbf{w}$ :

$$\begin{aligned} p(t|\mathbf{t}, \alpha, \beta) &= \int p(t|\mathbf{w}, \beta)p(\mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w} \\ &= \mathcal{N}(t|\mathbf{m}_N^T \boldsymbol{\phi}(\mathbf{x}), \sigma_N^2(\mathbf{x})) \end{aligned}$$

- where

$$\sigma_N^2(\mathbf{x}) = \frac{1}{\beta} + \boldsymbol{\phi}(\mathbf{x})^T \mathbf{S}_N \boldsymbol{\phi}(\mathbf{x}).$$

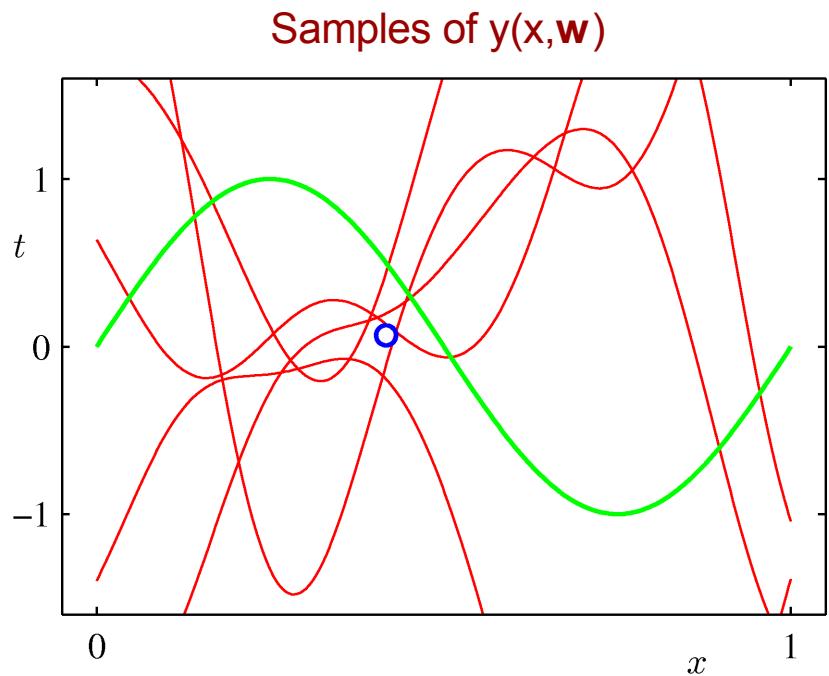
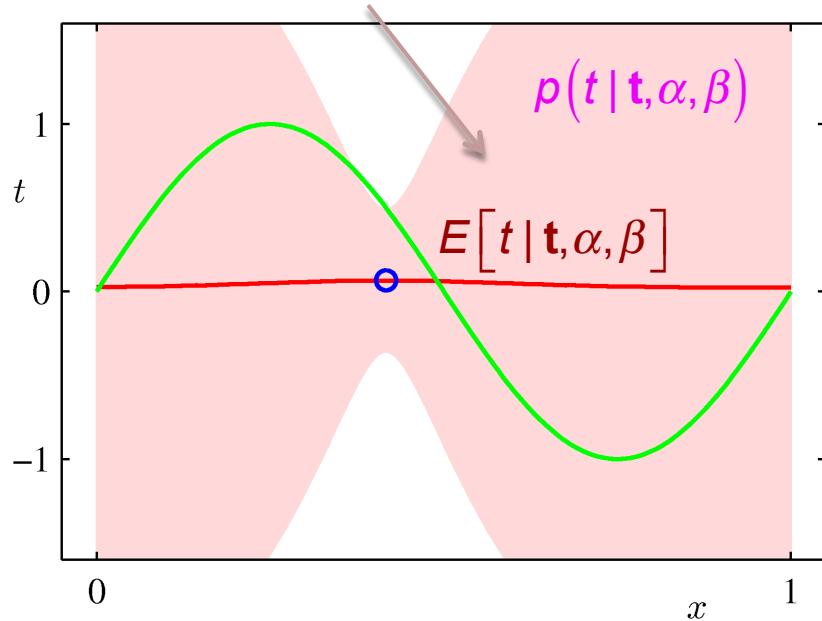
# Predictive Distribution

53

Probability &amp; Bayesian Inference

- Example: Sinusoidal data, 9 Gaussian basis functions, 1 data point

Notice how much bigger our uncertainty is relative to the ML method!!

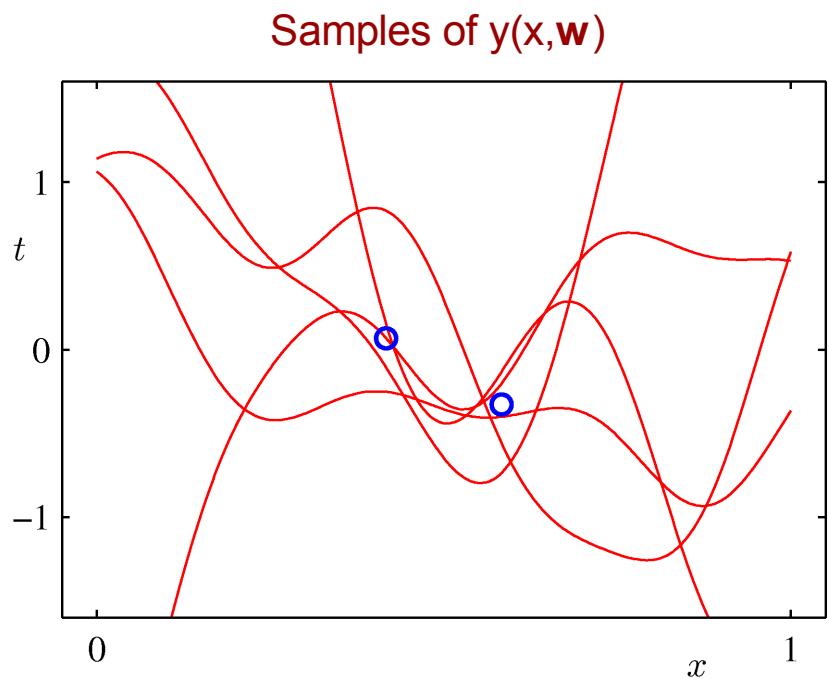
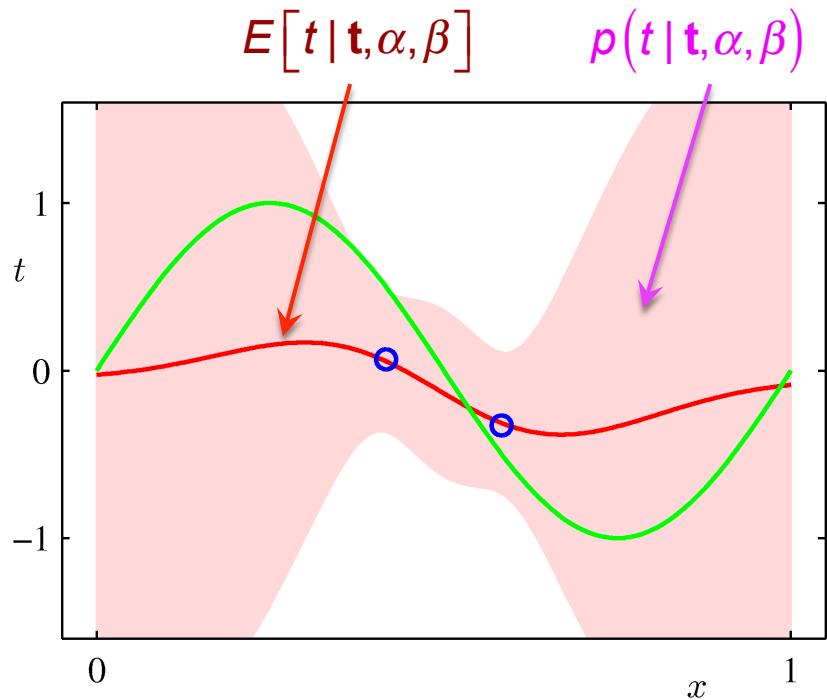


# Predictive Distribution

54

Probability &amp; Bayesian Inference

- Example: Sinusoidal data, 9 Gaussian basis functions, 2 data points

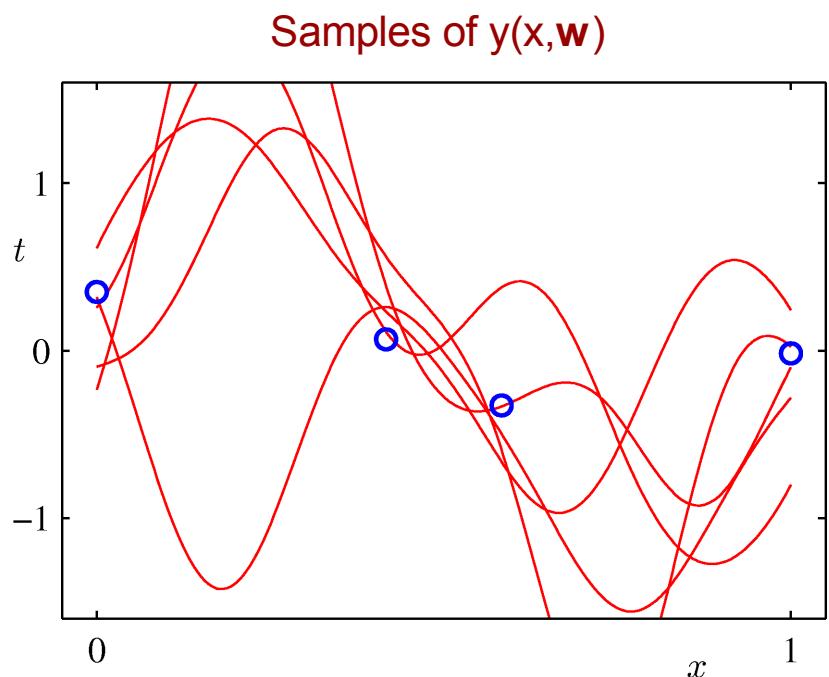
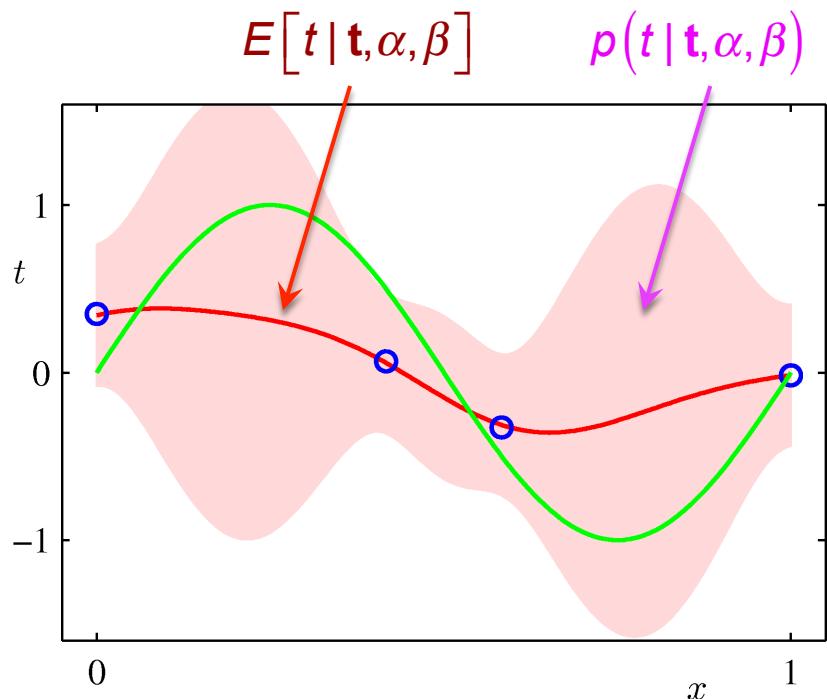


# Predictive Distribution

55

Probability &amp; Bayesian Inference

- Example: Sinusoidal data, 9 Gaussian basis functions, 4 data points

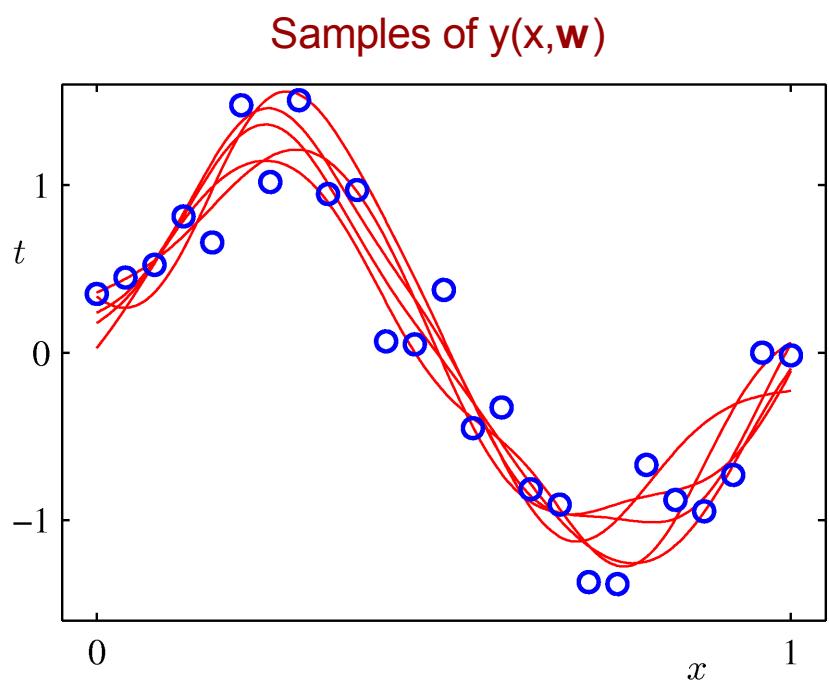
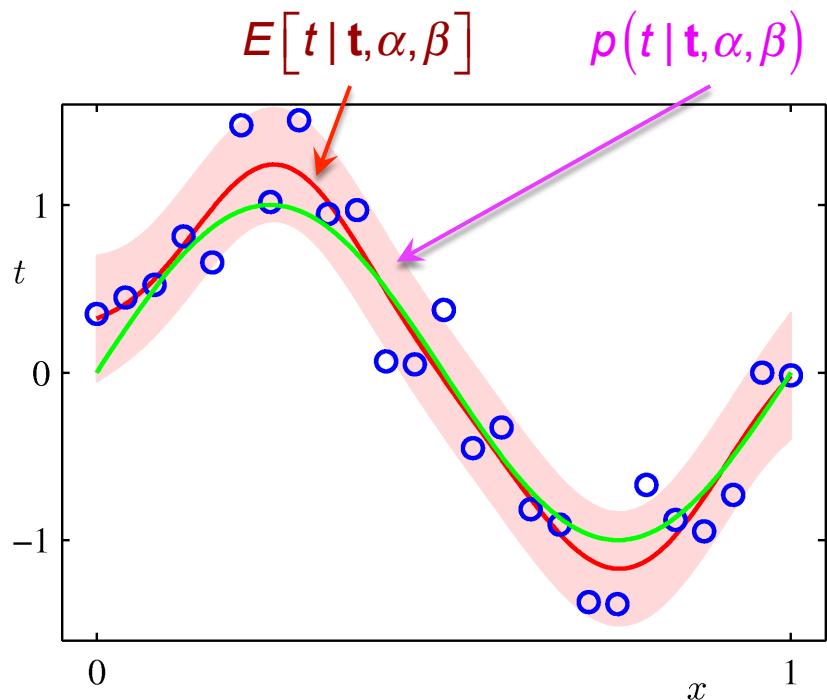


# Predictive Distribution

56

Probability &amp; Bayesian Inference

- Example: Sinusoidal data, 9 Gaussian basis functions, 25 data points



# Linear Regression Topics

57

Probability & Bayesian Inference

- What is linear regression?
- Example: polynomial curve fitting
- Other basis families
- Solving linear regression problems
- Regularized regression
- Multiple linear regression
- Bayesian linear regression