

Bayesian Regression: Basis Functions, MLE & Regularized Least Squares, Multiple Outputs, Inference with σ^2 unknown, Zellner's g-Prior, Uninformative Priors

Prof. Nicholas Zabaras
Center for Informatics and Computational Science

<https://cics.nd.edu/>

University of Notre Dame
Notre Dame, Indiana, USA

Email: nzabaras@gmail.com
URL: <https://www.zabaras.com/>

September 24, 2018

Statistical Computing, University of Notre Dame, Notre Dame, IN, USA (Fall 2018, N. Zabaras)



Contents

- Linear basis function models, Maximum likelihood and least squares,
Geometry of least squares, Convexity of the NLL , Sequential learning,
Robust Linear Regression, Regularized least squares, Multiple Outputs

- Bayesian linear regression, Parameter posterior distribution, A Note on Data Centering, Numerical Example, Predictive distribution, Bayesian inference in linear regression when σ^2 is unknown, Zellner's g-Prior,
Uninformative (Semi-Conjugate) Prior, Evidence Approximation

Following closely:

- Chris Bishop's PRML book, Chapter 3
- Kevin Murphy, Machine Learning: A probabilistic Perspective, Chapter 7
- Regression using parametric discriminative models in pmtk3 (run TutRegr.m in Pmtk3)

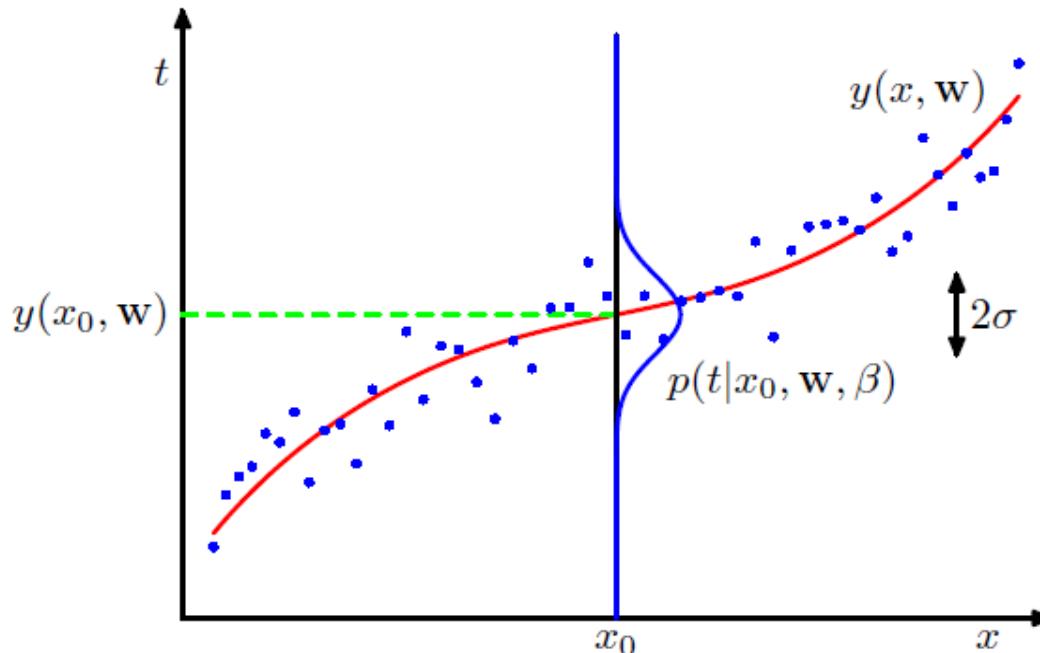


Linear Regression

- We already considered in an earlier lecture an example of linear regression – polynomial curve fitting.
- We are interested in a linear combination – regression – of a ``fixed set'' of nonlinear basis functions.
- Supervised learning: N observations $\{x_n\}$ with corresponding target values $\{t_n\}$ are provided.
- The goal is to predict t for a new value x .
- We construct a function such that $y(x)$ is a prediction of t .
- We follow a Bayesian perspective and model the predictive distribution $p(t|x)$.

Linear Regression

- From the conditional distribution $p(t|x)$, we can make point estimates of t for a given x by minimizing a ‘loss function’.
- For a quadratic loss function, the point estimate is the conditional mean $y(x, w) = \mathbb{E}[t|x]$.



Linear Regression

- The simplest linear model for regression is one that involves a linear combination of the input variables

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1 x_1 + \dots + w_D x_D = \sum_{i=0}^D w_i x_i$$

where $\mathbf{x} = (x_1 \quad x_2 \quad \dots \quad x_D)^T$ and we have defined $x_0 = 1$

- This is often simply known as linear regression. D is the input dimensionality.



Linear Basis Functions Models

- More generally:

$$y(\mathbf{x}, \mathbf{w}) = w_0 + w_1\phi_1(\mathbf{x}) + \dots + w_{M-1}\phi_{M-1}(\mathbf{x}) = \sum_{i=0}^{M-1} w_i\phi_i(\mathbf{x}) = \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}), \text{ where: } \phi_0(\mathbf{x}) = 1$$

where $\phi_i(\mathbf{x})$ are known as basis functions and

$$\boldsymbol{\phi}(\mathbf{x}) = (\phi_0(\mathbf{x}) \quad \phi_1(\mathbf{x}) \quad \dots \quad \phi_{M-1}(\mathbf{x}))^T$$

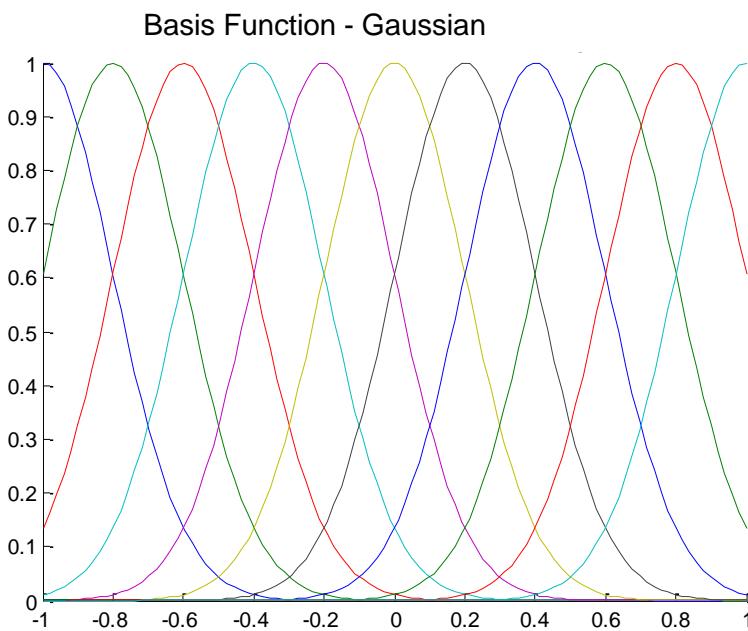
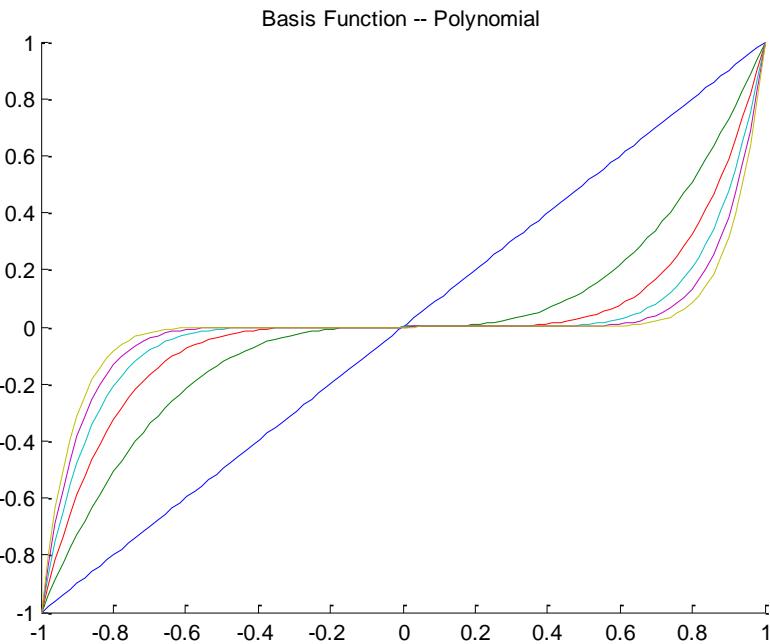
- The parameter w_0 allows for any fixed offset in the data and is called the **bias parameter**.
- For convenience, we define an additional dummy ‘basis function’ so that,

$$\phi_0(\mathbf{x}) = 1$$

- Often $\{\phi_i(\mathbf{x})\}$ represent features extracted from the data \mathbf{x} .



Polynomial and Gaussian Basis Functions



[MatLab code](#)

$$\mu = -1 : 0.2 : 1$$

$$s = 0.2$$

Gaussian basis functions:

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

Polynomial basis functions
(scalar input, global support):

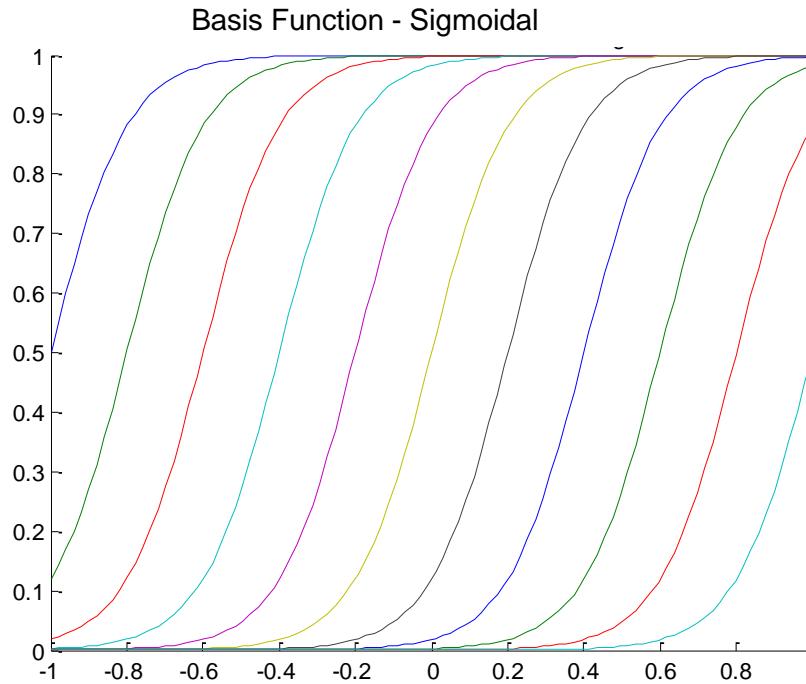
$$\phi_j(x) = x^j$$



Logistic Sigmoidal Basis Functions

$$\mu = -1 : 0.2 : 1$$

$$s = 0.1$$



[MatLab code](#)

Sigmoidal basis functions:

$$\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right), \quad \sigma(a) = \frac{1}{1 + e^{-a}}$$

$$\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}} = 2\sigma(2a) - 1$$

$\sigma(a)$: logistic sigmoid function

Sigmoidal and Tanh Basis Functions

- The sigmoidal and tanh basis functions are related:

$$\tanh(a) = \frac{e^a - e^{-a}}{e^a + e^{-a}} = 2\sigma(2a) - 1 \quad \text{where} \quad \sigma(a) = \frac{1}{1 + e^{-a}}$$

- A general linear combination of logistic sigmoidal functions is equivalent to a linear combination of tanh functions:

$$\begin{aligned} y(x, w) &= w_0 + \sum_{j=1}^M w_j \sigma\left(\frac{x - \mu_j}{s}\right) = w_0 + \sum_{j=1}^M w_j \sigma\left(2 \frac{x - \mu_j}{2s}\right) \\ &= w_0 + \sum_{j=1}^M w_j \frac{1 + \tanh\left(\frac{x - \mu_j}{2s}\right)}{2} = u_0 + \sum_{j=1}^M u_j \tanh\left(\frac{x - \mu_j}{2s}\right) \end{aligned}$$

$$\text{where: } u_0 = w_0 + \frac{1}{2} \sum_{j=1}^M w_j, \quad u_j = \frac{w_j}{2}$$



Choice of Basis Functions

- We are interested in functions of *local support* to explore adaptivity.
- Local support functions comprise a spectrum of different spatial frequencies.
- An example is wavelets that are local both spatially and in frequency.
 - They are however useful only when the input is defined on a lattice.

Maximum Likelihood and Least Squares

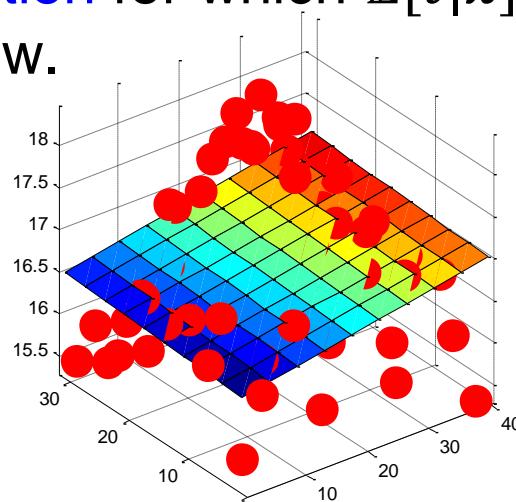
- Assume observations from a deterministic function with added Gaussian noise:

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

which is the same as saying,

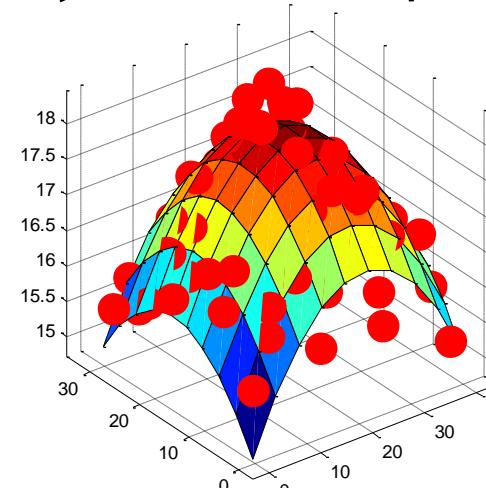
$$p(t | \mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t | y(\mathbf{x}, \mathbf{w}), \beta^{-1})$$

- Here β is the precision. This is based on a squared loss function for which $\mathbb{E}[t|\mathbf{x}] = y(\mathbf{x}, \mathbf{w})$. A 2D example is shown below.



$$\mathbb{E}[t | \mathbf{x}] = w_0 + w_1 x_1 + w_2 x_2$$

Statistical Computing, University of Notre Dame, Notre Dame, IN, USA (Fall 2018, N. Zabaras)



Run [surfaceFitDemo](#) from PMTK3

Maximum Likelihood and Least Squares

- Given observed inputs, $X = \{x_1, \dots, x_N\}$, and targets $t = \{t_1, \dots, t_N\}^T$, we obtain the likelihood function

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

- We often use the *log-likelihood*:

$$\ell(w, \beta) = \log p(\mathcal{D} | \theta) = \sum_{n=1}^N \log \mathcal{N}(t_n | w^T \phi(x_n), \beta^{-1}), \theta = (w, \beta)$$

- Instead of maximizing the log-likelihood, one equivalently can minimize *the negative log-likelihood (NLL)*

$$NLL(w, \beta) = -\sum_{n=1}^N \log \mathcal{N}(t_n | w^T \phi(x_n), \beta^{-1})$$

Maximum Likelihood and Least Squares

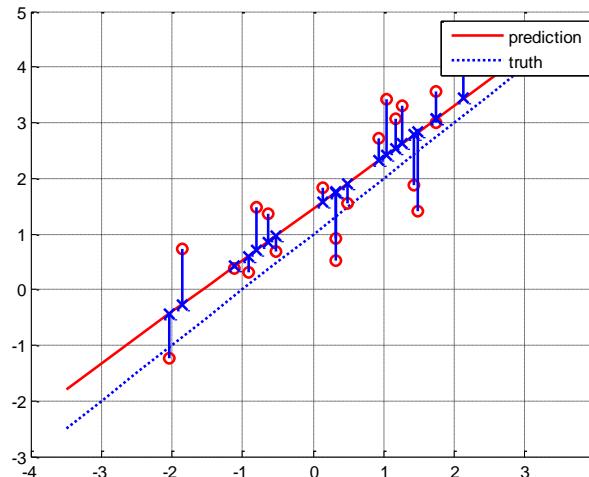
- Taking the log of the likelihood, we obtain:

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

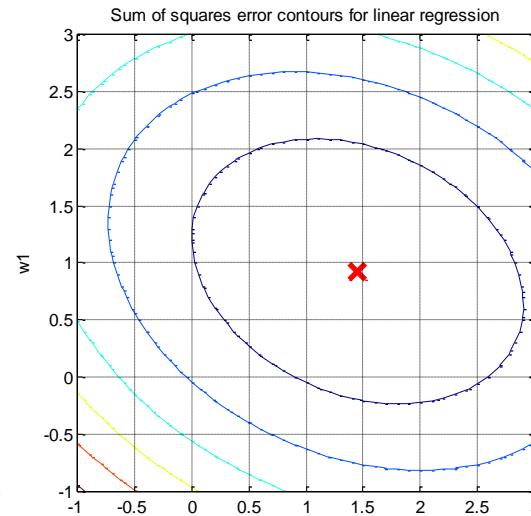
where with $E_D(\mathbf{w})$ we have denoted:

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

- RSS is often known as the *residual sum of squares* or *sum of squared errors (SSE)* and *MSE* is the mean squared error.
- Computing \mathbf{w} via *MLE* is the same as *Least Squares*.



Run [residualsDemo](#) from PMTK3



The NLL is a quadratic bowl with a unique minimum (the MLE estimate)

Run [contoursSSEdemo](#), from PMTK3



Maximum Likelihood and Least Squares

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

- Setting the gradient of the log-likelihood (written here as a row vector) wrt \mathbf{w} equal to zero:

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

- This equation can be solved for \mathbf{w} .



Maximum Likelihood and Least Squares

$$\mathbf{w}^T \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T = \sum_{n=1}^N t_n \phi(\mathbf{x}_n)^T \Rightarrow \mathbf{w}^T \Phi^T \Phi = (\Phi^T \mathbf{t})^T \Rightarrow \Phi^T \Phi \mathbf{w} = \Phi^T \mathbf{t}$$

□ We obtain (normal equation; ordinary least squares solution):

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \equiv \Phi^\dagger \mathbf{t}, \quad \Phi^\dagger : \text{Moore-Penrose pseudo-inverse}$$

where we have defined:

$$\Phi = \begin{pmatrix} \phi^T(\mathbf{x}_1) \\ \phi^T(\mathbf{x}_2) \\ \vdots \\ \phi^T(\mathbf{x}_N) \end{pmatrix} = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & .. & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & .. & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \vdots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & .. & \phi_{M-1}(\mathbf{x}_N) \end{pmatrix}, \quad \Phi^T = (\phi(\mathbf{x}_1) \quad \phi(\mathbf{x}_2) \quad .. \quad \phi(\mathbf{x}_N))$$
$$\Phi^T = (\phi_0 \quad \phi_1 \quad .. \quad \phi_N)$$
$$\phi(\mathbf{x}) = (\phi_0(\mathbf{x}_i) \quad \phi_1(\mathbf{x}_i) \quad .. \quad \phi_{M-1}(\mathbf{x}_i))^T$$

□ Note that indeed:

$$\Phi^T \Phi = \sum_{n=1}^N \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T, \quad \Phi^T \mathbf{t} = \sum_{n=1}^N t_n \phi(\mathbf{x}_n)$$

Maximum Likelihood and Least Squares

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

- Taking now the log of the likelihood with respect to β gives:

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

- So the MLE variance β_{ML} is equal to the residual variance of the target values around the regression function.

Computing the Bias Parameter

- If we make the bias parameter w_0 explicit, then the error function becomes

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \left(t_n - w_0 - \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}_n) \right)^2$$

- Setting the derivative with respect to w_0 equal to zero, and solving for w_0 , we obtain

$$w_0 = \frac{1}{N} \sum_{n=1}^N t_n - \frac{1}{N} \sum_{n=1}^N \sum_{j=1}^{M-1} w_j \phi_j(\mathbf{x}_n) \Rightarrow$$

$$w_0 = \bar{t} - \sum_{j=1}^{M-1} w_j \bar{\phi}_j, \quad \bar{t} = \frac{1}{N} \sum_{n=1}^N t_n, \quad \bar{\phi}_j = \frac{1}{N} \sum_{n=1}^N \phi_j(\mathbf{x}_n)$$

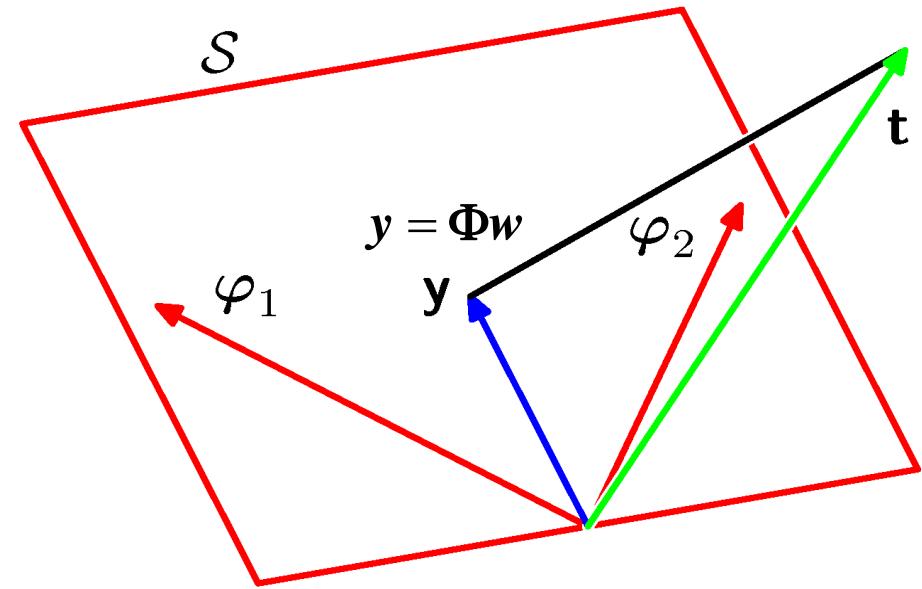
- The bias parameter w_0 compensates for the difference between the averages of the target values and the weighted sum of the averages of the basis function values.



Geometry of Least Squares

- We look for a geometrical interpretation of the least-squares solution in an N -dimensional space. t is a vector in that space with components t_1, \dots, t_N ($N > M$).
- The least-squares regression function is obtained by finding the orthogonal projection of the data vector t onto the subspace spanned by the basis functions $\varphi_j(x)$.
- Note φ_j is here the j -th column of Φ .

$$\varphi_j \equiv (\phi_j(x_1), \dots, \phi_j(x_N))^T$$



$$\Phi = (\varphi_0 \quad \varphi_1 \quad \dots \quad \varphi_{M-1}) = \begin{pmatrix} \phi_0(x_1) & \phi_1(x_1) & \dots & \phi_{M-1}(x_1) \\ \phi_0(x_2) & \phi_1(x_2) & \dots & \phi_{M-1}(x_2) \\ \vdots & \vdots & \vdots & \vdots \\ \phi_0(x_N) & \phi_1(x_N) & \dots & \phi_{M-1}(x_N) \end{pmatrix}$$

Geometry of Least Squares

- We are looking for w such that the projection error

$$t - y = t - \Phi w$$

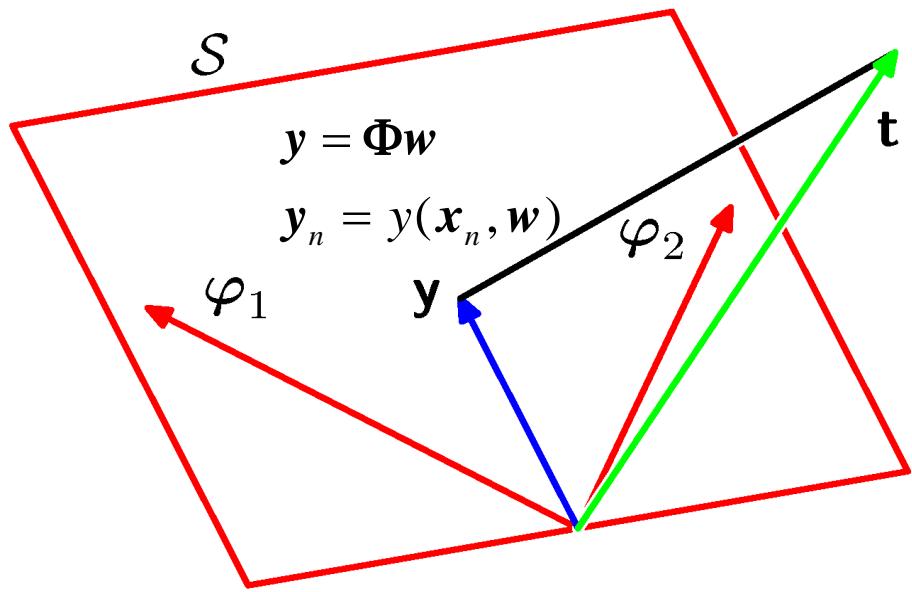
is orthogonal to the basis φ_j , i.e. such that:

$$\Phi^T (t - \Phi w) = 0$$

- These are the normal equations we derived earlier.

$$\Phi^T = \begin{pmatrix} \varphi_0^T \\ \varphi_1^T \\ \vdots \\ \varphi_{M-1}^T \end{pmatrix}$$

$$\varphi_j(x) = (\phi_j(x_1), \dots, \phi_j(x_N))^T$$

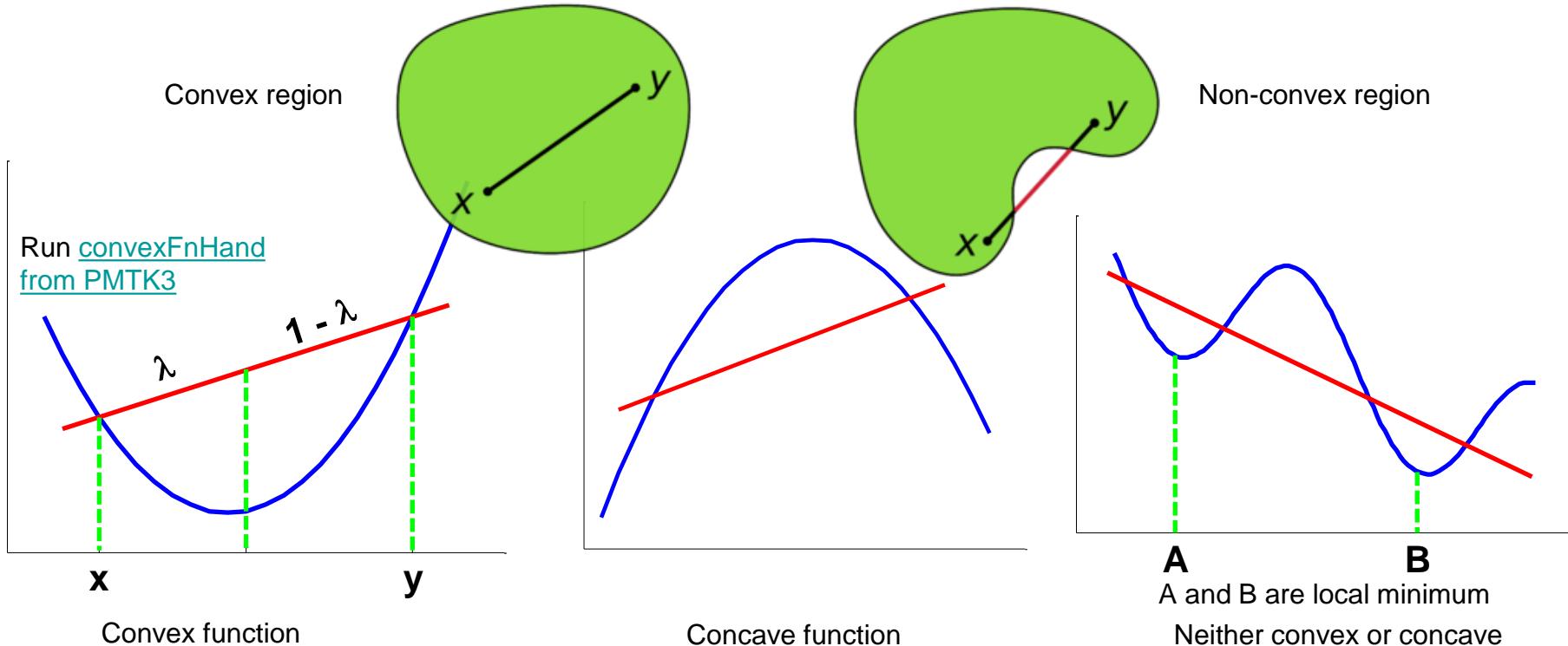


S : M –dimensional subspace spanned by $\varphi_j(x)$

$$\Phi = [\varphi_0(x) \quad \varphi_1(x) \quad \dots \quad \varphi_{M-1}(x)]$$

Convexity of the NLL

- Convexity of the NLL (positive definite Hessian) leads to a unique globally optimal MLE.
- Some models of interest don't have concave likelihoods and locally optimal MLE estimates are needed.



$$\theta^2, e^\theta, \theta \log \theta \ (\theta > 0)$$



Sequential Learning: LMS Algorithm

- If the data set is large, we use sequential (on-line) algorithms
- We apply the technique of **stochastic (sequential) gradient descent.**
- If the error function comprises a sum over data points $E = \sum_n E_n$

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

then after presentation of pattern n , the stochastic gradient descent algorithm updates the parameter vector \mathbf{w} using

$$\mathbf{w}^{(\tau+1)} = \mathbf{w}^{(\tau)} - \eta \nabla E_n = \mathbf{w}^{(\tau)} + \underbrace{\eta \left(t_n - \mathbf{w}^{(\tau)^T} \phi(\mathbf{x}_n) \right) \phi(\mathbf{x}_n)}_{-\nabla E_n}$$

τ is the iteration number & η the learning rate parameter.

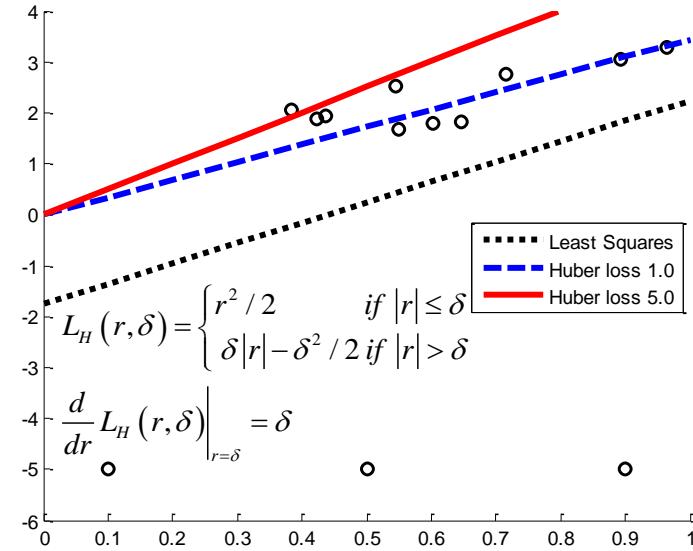
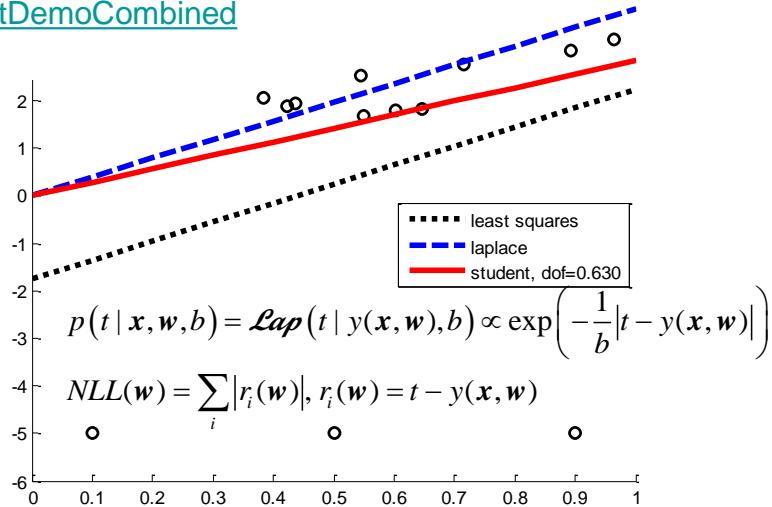
- This is known as **least-mean-squares** or the **LMS algorithm.**



Robust Linear Regression

- Using a Gaussian distribution for the noise,
$$t = y(\mathbf{x}, \mathbf{w}) + \varepsilon, \varepsilon \sim \mathcal{N}(\varepsilon | 0, \beta^{-1})$$
 can result in poor fit especially if we have outliers in the data.
- Squared error penalizes deviations quadratically, so points far from the line have more affect on the fit than points near the line.
- To achieve robustness to outliers one can replace the Gaussian with a distribution that has heavy tails (e.g. the Laplace distribution). Such a distribution assigns higher likelihood to outliers, without having to perturb the regression line to “explain” them.

Run [linregRobustDemoCombined](#)
from PMTK3



Robust Linear Regression

- Using the Laplace distribution leads to L_1 error norm (non-linear objective function) that is difficult to optimize.
- A solution is to transform the problem (by increasing its dimension to $2N + M$) to a *linear programming problem*.

$$r_i \triangleq r_i^+ - r_i^-$$

$$\min_{\mathbf{w}, r_i^+, r_i^-} \sum_i (r_i^+ + r_i^-) \quad s.t. \quad r_i^+ \geq 0, r_i^- \geq 0, \mathbf{w}^T \mathbf{x}_i + r_i^+ - r_i^- = t_i$$

- Note that with our definition above:

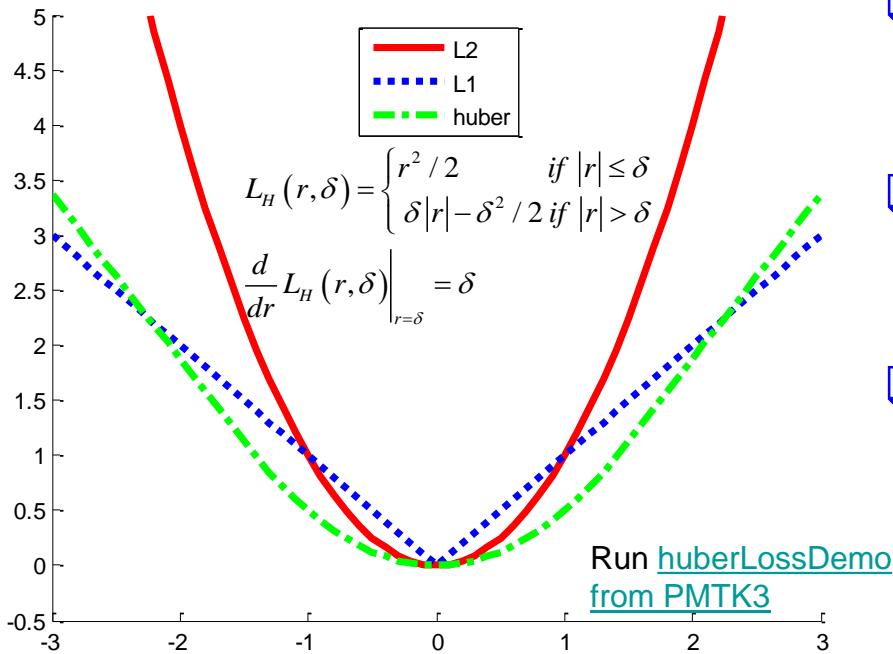
$$r_i^+ = \frac{1}{2}(r_i + |r_i|) = \begin{cases} r_i & \text{if } r_i \geq 0 \\ 0 & \text{if } r_i < 0 \end{cases}, \quad r_i^- = \frac{1}{2}(|r_i| - r_i) = \begin{cases} 0 & \text{if } r_i > 0 \\ -r_i & \text{if } r_i \leq 0 \end{cases}$$

$$|r_i| = r_i^+ + r_i^-$$

- Boyd, S. and L. Vandenberghe (2004). *Convex optimization*. Cambridge



Huber Loss Function



- This is equivalent to L_2 for errors that are smaller than δ , and is equivalent to L_1 for larger errors.
- This loss function is everywhere differentiable, using the fact that $d/dr |r| = sign(r)$ if $r \neq 0$.
- The function is also C_1 continuous, since the gradients of the two parts of the function match at $r = \pm\delta$.

- Optimizing the Huber loss is much faster than using the Laplace likelihood, since we can use standard optimization methods (quasi-Newton) rather than linear programming.
- The Huber method also has a probabilistic interpretation, although it is rather unnatural (Pontil et al. 1998).

- Pontil, M., S. Mukherjee, and F. Girosi (1998). [On the Noise Model of Support Vector Machine Regression](#). Technical report, MIT AI Lab.
- Huber, P. (1964). [Robust estimation of a location parameter](#). *Annals of Statistics* 53, 73a ~A, S101.

Regularized LS - Ridge Regression

- Consider the error function:

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

data term + regularization term

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

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

- Specifically, setting the gradient with respect to \mathbf{w} to zero, and solving for \mathbf{w} as before, we obtain

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

- This is a trivial extension of the least-squares solution we encountered earlier (*Regularized Least Squares – Ridge Regression*)



Regularized Least Squares

- Regularized solution:

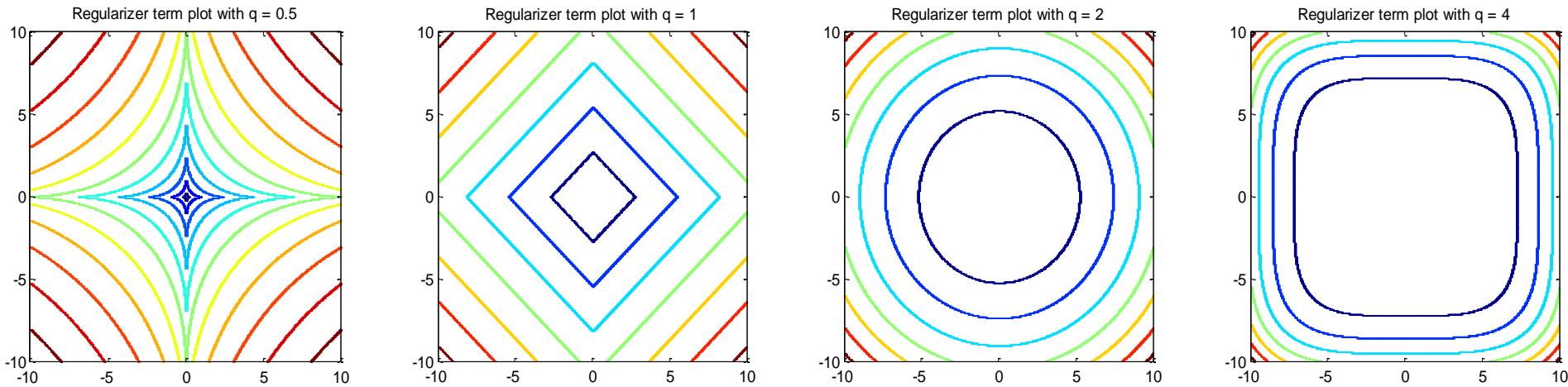
$$w = (\lambda I + \Phi^T \Phi)^{-1} \Phi^T t$$

- Regularization limits the effective model complexity (the appropriate number of basis functions).
- This is replaced with the problem of finding a suitable value of the regularization coefficient λ .
- λ controls how many non-zero w 's (i.e. basis functions) you have.

Regularized Least Squares

- With a more general regularizer, we have

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



[MatLab code](#)

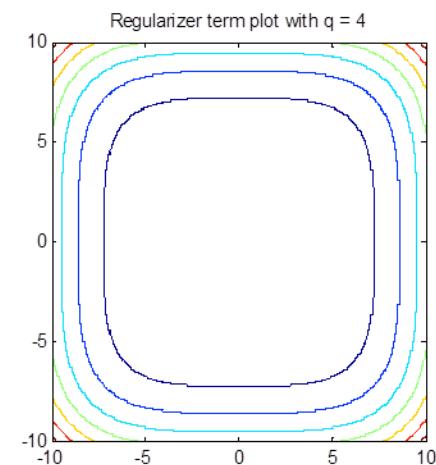
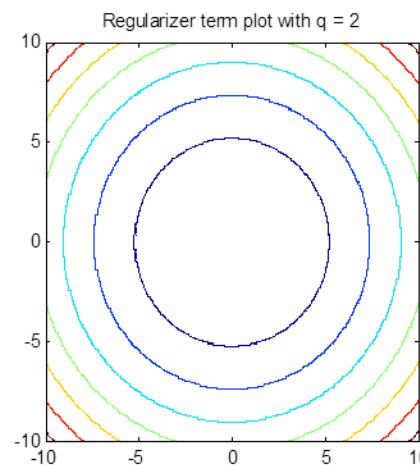
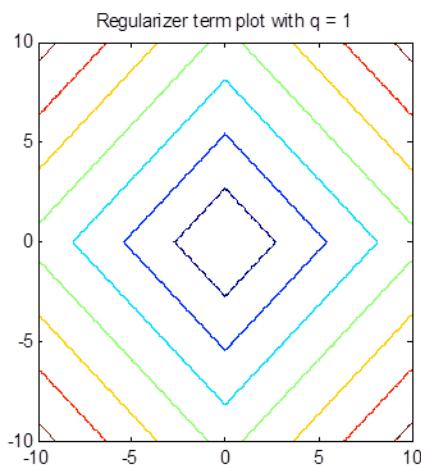
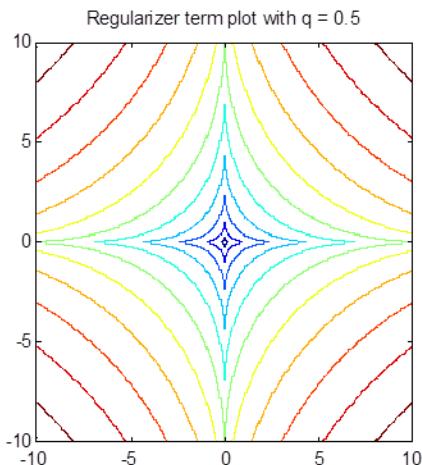
- $q = 1$ is known as the **Lasso regularizer**. These plots show only the regularizer term with $\lambda = 0.7334$.



Regularized Least Squares

- With a more general regularizer, we have

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



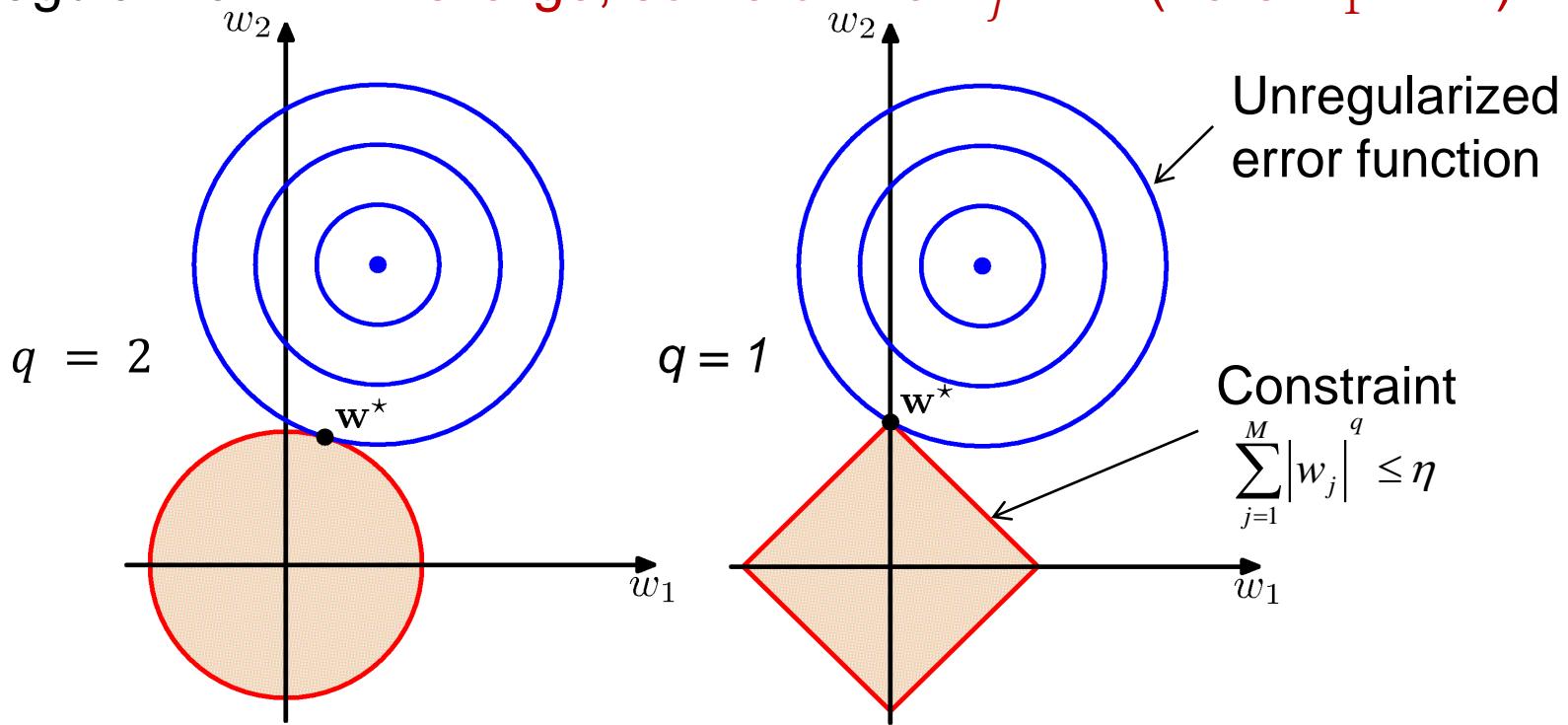
[MatLab code](#)

- $q = 2$ corresponds to the quadratic regularizer.



Regularized Least Squares

- Lasso tends to generate sparser solutions than a quadratic regularizer – if λ is large, some of the $w_j \rightarrow 0$ (here $w_1 = 0$).



- Here, we consider that the regularized least squares solution is equivalent to minimizing the unregularized sum of squares with the constraint shown for some η (see proof next).

Regularized Least Squares

- Let us write the constraint in the equivalent form: $\frac{1}{2} \left(\sum_{j=1}^M |w_j|^q - \eta \right) \leq 0$
- This leads to the following Lagrangian function:

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

- This is identical to our regularized least squares (RLS) in the dependence on \mathbf{w} .

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

- For a particular $\lambda > 0$, let $\mathbf{w}^*(\lambda)$ be the solution of the RLS in (*).
- From the Kuhn-Tucker optimality conditions for $L(\mathbf{w}, \lambda)$ we then see:

$$\eta = \sum_{j=1}^M |w_j^*|^q$$



Kuhn-Tucker Optimality Conditions

- Consider the following constraint **minimization** problem:

$$\min_x f(\mathbf{x}), \text{subject to } g(\mathbf{x}) \geq 0$$

- This is equivalent as the minimization with respect to \mathbf{x} and λ of the following Lagrangian:

$$\min_{\mathbf{x}, \lambda} L(\mathbf{x}, \lambda) = f(\mathbf{x}) - \lambda g(\mathbf{x})$$

subject to the following (Kuhn-Tucker) conditions:

$$\lambda \geq 0, g(\mathbf{x}) \geq 0, \lambda g(\mathbf{x}) = 0$$

- Note for maximization problems, the Lagrangian should be modified as: $L(\mathbf{x}, \lambda) = f(\mathbf{x}) + \lambda g(\mathbf{x})$



Multiple Outputs-Isotropic Covariance

- If we want to predict $K > 1$ target variables, we use the same basis for all components of the target vector):

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

\mathbf{W} is an $M \times K$ matrix and \mathbf{t} is K dimensional.

- 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

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

K-Independent Regression Problems

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

- As before, we can maximize this function with respect to \mathbf{W} , giving

$$\mathbf{W}_{ML} = \underbrace{\left(\boldsymbol{\Phi}^T \boldsymbol{\Phi} \right)^{-1}}_{M \times M} \underbrace{\boldsymbol{\Phi}^T}_{M \times N} \underbrace{\mathbf{T}}_{N \times K}$$

- If we examine this result for each target variable t_k , we have (take the k^{th} column of \mathbf{W} and \mathbf{T}):

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

which is identical with the single output case (so **there is decoupling between the target variables**)

- As expected, we obtain K –**independent regression problems.**



Multiple Outputs - Full Covariance

- Let us repeat the earlier formulation but with covariance matrix Σ . If we want to predict $K > 1$ target variables, we use the same basis for all components of the target vector):

$$p(\mathbf{t} | \mathbf{x}, \mathbf{W}, \Sigma) = \mathcal{N}(\mathbf{t} | \mathbf{y}(\mathbf{x}, \mathbf{W}), \Sigma) = \mathcal{N}(\mathbf{t} | \mathbf{W}^T \phi(\mathbf{x}), \Sigma)$$

where \mathbf{W} is an $M \times K$ matrix of parameters

- 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 $\ln p(\mathbf{T} | \mathbf{X}, \mathbf{W}, \Sigma)$:

$$\sum_{n=1}^N \ln \mathcal{N}(t_n | \mathbf{W}^T \phi(x_n), \Sigma) = -\frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}^T \phi(\mathbf{x}_n))^T \Sigma^{-1} (\mathbf{t}_n - \mathbf{W}^T \phi(\mathbf{x}_n))$$



Multiple Outputs - Full Covariance

$$\ln p(\mathbf{T} | \mathbf{X}, \mathbf{W}, \Sigma) = -\frac{N}{2} \ln |\Sigma| - \frac{1}{2} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}^T \boldsymbol{\phi}(\mathbf{x}_n))^T \Sigma^{-1} (\mathbf{t}_n - \mathbf{W}^T \boldsymbol{\phi}(\mathbf{x}_n))$$

- As before, we maximize this function with respect to \mathbf{W} ,

$$0 = -\sum_{n=1}^N \Sigma^{-1} (\mathbf{t}_n - \mathbf{W}^T \boldsymbol{\phi}(\mathbf{x}_n)) \boldsymbol{\phi}(\mathbf{x}_n)^T \Rightarrow \mathbf{W}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{T}$$

- For the ML estimate for Σ , use the result for the MLE of the covariance of a multivariate Gaussian:

$$\Sigma = \frac{1}{N} \sum_{n=1}^N (\mathbf{t}_n - \mathbf{W}_{ML}^T \boldsymbol{\phi}(\mathbf{x}_n)) (\mathbf{t}_n - \mathbf{W}_{ML}^T \boldsymbol{\phi}(\mathbf{x}_n))^T$$

- Note that each column of \mathbf{W}_{ML} is of the form

$$\mathbf{w}_{ML} = (\Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$$

seen for isotropic noise distribution and is independent of Σ !

Bayesian Linear Regression

- Effective model complexity in MLE is governed by the number of basis functions and is controlled by the size of the data set.
- With regularization, the effective model complexity is controlled mainly by λ and still by the number and form of the basis functions.
- *Thus the model complexity for a particular problem cannot be decided simply by maximizing the likelihood function as this leads to excessively complex models and over-fitting.*
- Independent hold-out data can be used to determine model complexity but this wastes data and it is computationally expensive.

Bayesian Linear Regression

- A Bayesian treatment of linear regression avoids the over-fitting of maximum likelihood.
- Bayesian approaches lead to automatic methods of determining model complexity using the training data alone.



Bayesian Linear Regression

- Assume additive Gaussian noise with known precision β . The likelihood function $p(\mathbf{t}|\mathbf{w})$ is the exponential of a quadratic function of \mathbf{w}

$$p(\mathbf{t} | \mathbf{X}, \mathbf{w}, \beta) = \prod_{n=1}^N \mathcal{N}(t_n | \mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) = \left(\frac{\beta}{2\pi} \right)^{N/2} \exp \left(-\frac{\beta}{2} \sum_{n=1}^N \{ t_n - \mathbf{w}^T \phi(\mathbf{x}_n) \}^2 \right)$$

and its **conjugate prior** is Gaussian:

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

- Combining this with the likelihood and using results for marginal and conditional Gaussian distributions, gives **the posterior**

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

where

$$\mathbf{m}_N = \mathbf{S}_N \left(\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t} \right)$$

$$\mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \Phi^T \Phi$$



Posterior Distribution: Derivation

$$p(\mathbf{w} | \mathbf{m}_0, S_0) \propto \exp \left\{ -\frac{1}{2} (\mathbf{w} - \mathbf{m}_0)^T S_0^{-1} (\mathbf{w} - \mathbf{m}_0) \right\},$$

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) \propto \exp \left\{ -\frac{\beta}{2} \sum_{n=1}^N (\phi(x_n)^T \mathbf{w} - t_n)^2 \right\} \Rightarrow$$

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) \propto \exp \left\{ -\frac{1}{2} \mathbf{w}^T \sum_{n=1}^N \beta \phi(x_n) \phi(x_n)^T \mathbf{w} + \beta \sum_{n=1}^N t_n \phi(x_n) \mathbf{w} \right\}$$

- We now have the product of two Gaussians and the posterior is easily computed as:

$$p(\mathbf{w} | \mathbf{x}, \mathbf{t}, S_0, \beta) \propto$$

$$\exp \left(-\frac{1}{2} \mathbf{w}^T S_0^{-1} \mathbf{w} - \mathbf{w}^T S_0^{-1} \mathbf{m}_0 - \frac{1}{2} \mathbf{w}^T \sum_{n=1}^N \beta \phi(x_n) \phi(x_n)^T \mathbf{w} + \beta \mathbf{w}^T \sum_{n=1}^N t_n \phi(x_n) \right)$$

$$\propto \mathcal{N} \left(\mathbf{w} / \underbrace{S_N (S_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t})}_{\mathbf{m}_N}, S_N \right), S_N^{-1} = S_0^{-1} + \sum_{n=1}^N \beta \phi(x_n) \phi(x_n)^T = S_0^{-1} + \beta \Phi^T \Phi$$

Complete the square in \mathbf{w}



Sequential Posterior Calculation

- Note that because the posterior distribution is Gaussian, its posterior mode coincides with its mean.

$$p(\mathbf{w} | \mathbf{t}) = \mathcal{N}(\mathbf{w} | \mathbf{m}_N, \mathbf{S}_N) \quad \mathbf{m}_N = \mathbf{S}_N^{-1} (\mathbf{S}_0^{-1} \mathbf{m}_0 + \beta \Phi^T \mathbf{t}) \\ \mathbf{S}_N^{-1} = \mathbf{S}_0^{-1} + \beta \Phi^T \Phi$$

$$\mathbf{w}_{MAP} = \mathbf{m}_N$$

- The above expressions for the posterior mean and variance can also be written for a **sequential calculation** (we already have observed N data points and now considering an additional data point $(\mathbf{x}_{N+1}, t_{N+1})$). In this case, we have:

$$p(\mathbf{w} | \mathbf{t}_{N+1}, \mathbf{x}_{N+1}, \mathbf{m}_N, \mathbf{S}_N) = \mathcal{N}(\mathbf{w} | \mathbf{m}_{N+1}, \mathbf{S}_{N+1})$$

$$\mathbf{m}_{N+1} = \mathbf{S}_{N+1}^{-1} (\mathbf{S}_N^{-1} \mathbf{m}_N + \beta \phi_{n+1} t_{n+1}) \\ \mathbf{S}_{N+1}^{-1} = \mathbf{S}_N^{-1} + \beta \phi_{n+1} \phi_{n+1}^T$$



Bayesian Linear Regression

- Let us consider for a prior, a zero-mean isotropic Gaussian governed by a single precision parameter α so that

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

and the corresponding posterior distribution over \mathbf{w} is then given by

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

- The log of the posterior is the sum of the log likelihood and the log of the prior and, as a function of \mathbf{w} , takes the form

$$\ln p(\mathbf{w} | \mathbf{t}) = -\frac{\beta}{2} \sum_{n=1}^N \left\{ t_n - \mathbf{w}^T \phi(x_n) \right\}^2 - \frac{\alpha}{2} \mathbf{w}^T \mathbf{w} + const$$

- Thus the MAP estimate is the same as regularized least squares (Ridge Regression) with $\lambda = \alpha / \beta$.



A Note on Data Centering

- In linear regression, it helps to center the data in a way that does not require us to compute the offset term μ . Write the likelihood as:

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \mu, \beta) \propto \exp\left(-\frac{\beta}{2} (\mathbf{t} - \mu \mathbf{1}_N - \Phi \mathbf{w})^T (\mathbf{t} - \mu \mathbf{1}_N - \Phi \mathbf{w})\right)$$

- Let us assume that the input data are centered in each dimension such that:

$$\sum_{j=1}^N \phi_i(\mathbf{x}_j) = 0 \quad \forall i = 1, \dots, M$$

$$\Phi = \begin{pmatrix} \phi^T(\mathbf{x}_1) \\ \phi^T(\mathbf{x}_2) \\ \vdots \\ \phi^T(\mathbf{x}_N) \end{pmatrix} = \begin{pmatrix} \phi_1(\mathbf{x}_1) & \phi_2(\mathbf{x}_1) & \dots & \phi_M(\mathbf{x}_1) \\ \phi_1(\mathbf{x}_2) & \phi_2(\mathbf{x}_2) & \dots & \phi_M(\mathbf{x}_2) \\ \vdots & \vdots & \vdots & \vdots \\ \phi_1(\mathbf{x}_N) & \phi_2(\mathbf{x}_N) & \dots & \phi_M(\mathbf{x}_N) \end{pmatrix}, \quad \Phi^T = (\phi(\mathbf{x}_1) \quad \phi(\mathbf{x}_2) \quad \dots \quad \phi(\mathbf{x}_N))$$
$$\Phi^T = (\phi_1 \quad \phi_2 \quad \dots \quad \phi_N)$$
$$\phi(\mathbf{x}) = (\phi_0(\mathbf{x}_i) \quad \phi_1(\mathbf{x}_i) \quad \dots \quad \phi_{M-1}(\mathbf{x}_i))^T$$

- The mean of the output is equally likely to be positive or negative. Let us put an improper prior $p(\mu) \propto 1$ and integrate μ out.



A Note on Data Centering

- Introducing, $\bar{t} = \frac{1}{N} \sum_{i=1}^N t_i$, the marginal likelihood becomes :

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) \propto \int \exp \left(-\frac{\beta}{2} \left(\underbrace{\mathbf{t} - \bar{t}\mathbf{1}_N - \Phi\mathbf{w}}_A - (\mu - \bar{t})\mathbf{1}_N \right)^T \left(\mathbf{t} - \bar{t}\mathbf{1}_N - \Phi\mathbf{w} - (\mu - \bar{t})\mathbf{1}_N \right) \right) d\mu$$

- Completing the square in μ gives (use the centering of the input):

$$\begin{aligned} p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) &\propto \int \exp \left(-\frac{\beta}{2} \left((\mu - \bar{t})^2 N - 2(\mu - \bar{t}) \underbrace{\mathbf{A}^T \mathbf{1}_N}_{\substack{\bar{t}N - \bar{t}N - \mathbf{w}^T \Phi^T \mathbf{1}_N \\ 0}} + \mathbf{A}^T \mathbf{A} \right) \right) d\mu \\ &\propto \exp \left(-\frac{\beta}{2} \left(\mathbf{t} - \bar{t}\mathbf{1}_N - \Phi\mathbf{w} \right)^T \left(\mathbf{t} - \bar{t}\mathbf{1}_N - \Phi\mathbf{w} \right) \right) \end{aligned}$$

- Our model is now simplified if instead of \mathbf{t} we use (centered output) $\mathbf{t} \leftarrow \mathbf{t} - \bar{t}\mathbf{1}_N$ and the likelihood is simply written as:

$$p(\mathbf{t} | \mathbf{x}, \mathbf{w}, \beta) \propto \exp \left(-\frac{\beta}{2} (\mathbf{t} - \Phi\mathbf{w})^T (\mathbf{t} - \Phi\mathbf{w}) \right)$$

$\hat{\mu} = \bar{t} - \sum_{j=1}^D \bar{\phi}_j^T \mathbf{w}_j$, $[\bar{\phi}_1, \dots, \bar{\phi}_D]$ is formed by averaging each column of Φ

- Recall that the MLE estimate for μ is:



A Note on Data Centering

- To simplify the earlier notation, consider a linear regression model of the form

$$\mathbb{E}[y | \mathbf{x}] = w_0 + \mathbf{w}^T \mathbf{x}$$

- In the context e.g. of MLE, we need to minimize

$$\min_{w_0, \mathbf{w}} = \sum_{i=1}^N (t_i - w_0 - \mathbf{w}^T \mathbf{x}_i)^2$$

- Minimization wrt w_0 gives:

$$\sum_{i=1}^N (t_i - w_0 - \mathbf{w}^T \mathbf{x}_i) = 0 \Rightarrow w_0 N = \bar{t} N - N \mathbf{w}^T \bar{\mathbf{x}} \Rightarrow w_0 = \bar{t} - \mathbf{w}^T \bar{\mathbf{x}}$$

where:

$$\bar{\mathbf{x}} \equiv \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_D \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^N x_{i1} / N \\ \sum_{i=1}^N x_{i2} / N \\ \vdots \\ \sum_{i=1}^N x_{iD} / N \end{pmatrix}, \quad \bar{t} \equiv \sum_{i=1}^N t_i / N$$

- Thus:

$$\hat{w}_0 = \bar{t} - \bar{\mathbf{x}}^T \mathbf{w}$$



A Note on Data Centering

- Substituting the bias term in our objective function gives:

$$\min_{\mathbf{w}} \sum_{i=1}^N \left(t_i - \bar{t} + \mathbf{w}^T \bar{\mathbf{x}} - \mathbf{w}^T \mathbf{x}_i \right)^2 = \min_{\mathbf{w}} \sum_{i=1}^N \left(t_i - \bar{t} + \mathbf{w}^T (\bar{\mathbf{x}} - \mathbf{x}_i) \right)^2$$

- Minimization wrt \mathbf{w} gives:

$$\sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \hat{\mathbf{w}} = \sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{t}_i - \bar{\mathbf{t}})$$

- We thus first compute the MLE of \mathbf{w} using the centered input and output as follows:

$$\hat{\mathbf{w}} = (\mathbf{X}_c^T \mathbf{X}_c)^{-1} \mathbf{X}_c^T \mathbf{t}_c = \left[\sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{x}_i - \bar{\mathbf{x}})^T \right] \left[\sum_{i=1}^N (\mathbf{x}_i - \bar{\mathbf{x}})(\mathbf{t}_i - \bar{\mathbf{t}}) \right]$$
$$\mathbf{X}_c = \mathbf{X} - \bar{\mathbf{X}} = \mathbf{X} - \mathbf{1}_N \bar{\mathbf{x}}^T = \begin{pmatrix} \mathbf{x}_1^T - \bar{\mathbf{x}}^T \\ \mathbf{x}_2^T - \bar{\mathbf{x}}^T \\ \vdots \\ \mathbf{x}_N^T - \bar{\mathbf{x}}^T \end{pmatrix} = \begin{pmatrix} x_{11} - \bar{x}_1 & x_{12} - \bar{x}_2 & \dots & x_{1D} - \bar{x}_D \\ x_{21} - \bar{x}_1 & x_{22} - \bar{x}_2 & \dots & x_{2D} - \bar{x}_D \\ \vdots & \vdots & \ddots & \vdots \\ x_{N1} - \bar{x}_1 & x_{N2} - \bar{x}_2 & \dots & x_{NN} - \bar{x}_D \end{pmatrix}, \quad \bar{\mathbf{x}} \equiv \begin{pmatrix} \bar{x}_1 \\ \bar{x}_2 \\ \vdots \\ \bar{x}_D \end{pmatrix} = \begin{pmatrix} \sum_{i=1}^N x_{i1} / N \\ \sum_{i=1}^N x_{i2} / N \\ \vdots \\ \sum_{i=1}^N x_{iD} / N \end{pmatrix}, \quad \mathbf{t}_c = \mathbf{t} - \bar{\mathbf{t}} = \mathbf{t} - \mathbf{1}_N \bar{\mathbf{t}}, \quad \bar{\mathbf{t}} \equiv \sum_{i=1}^N t_i / N$$

- We can then estimate the MLE estimate of w_0 as follows:

$$\hat{w}_0 = \bar{t} - \bar{\mathbf{x}}^T \mathbf{w}$$

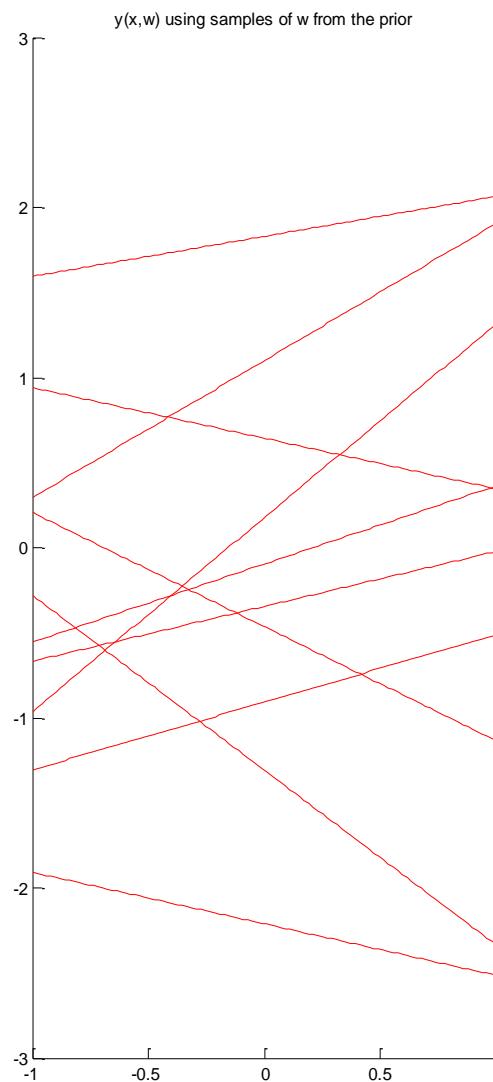
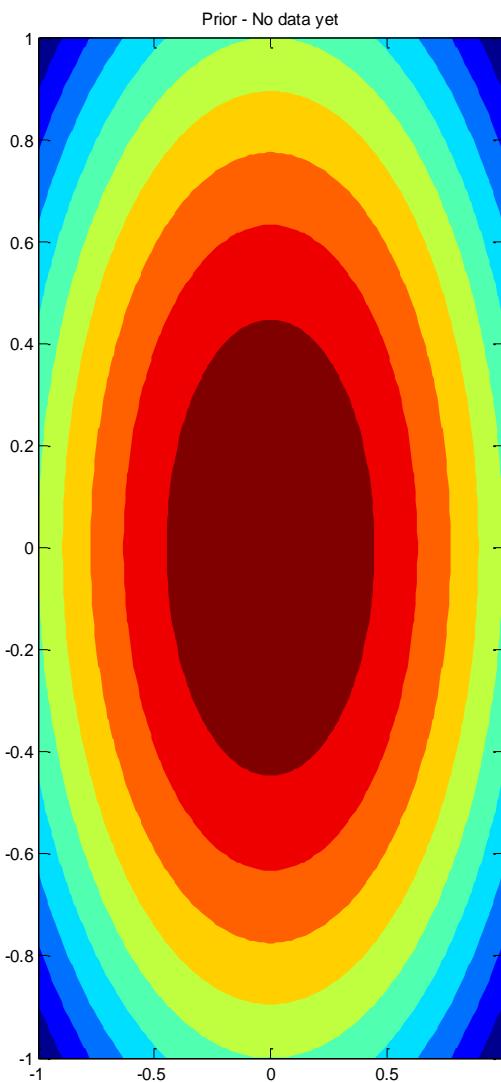


Bayesian Regression: Example

- We generate synthetic data from the function $f(x, a) = a_0 + a_1 x$ with parameter values $a_0 = -0.3$ and $a_1 = 0.5$ by first choosing values of x_n from the uniform distribution $\mathcal{U}(x| -1, 1)$, then evaluating $f(x_n, a)$, and finally adding Gaussian noise with standard deviation of 0.2 to obtain the target values t_n .
- We assume $\beta = (1/0.2)^2 = 25$ and $\alpha = 2.0$.
- We perform Bayesian inference sequentially – one point at a time – so the posterior at each level becomes the new prior.
- We show results after 1, 2 and 22 points have been collected.
- The results include the likelihood contours (for 1 point), the posterior and samples of the regression function from the posterior.



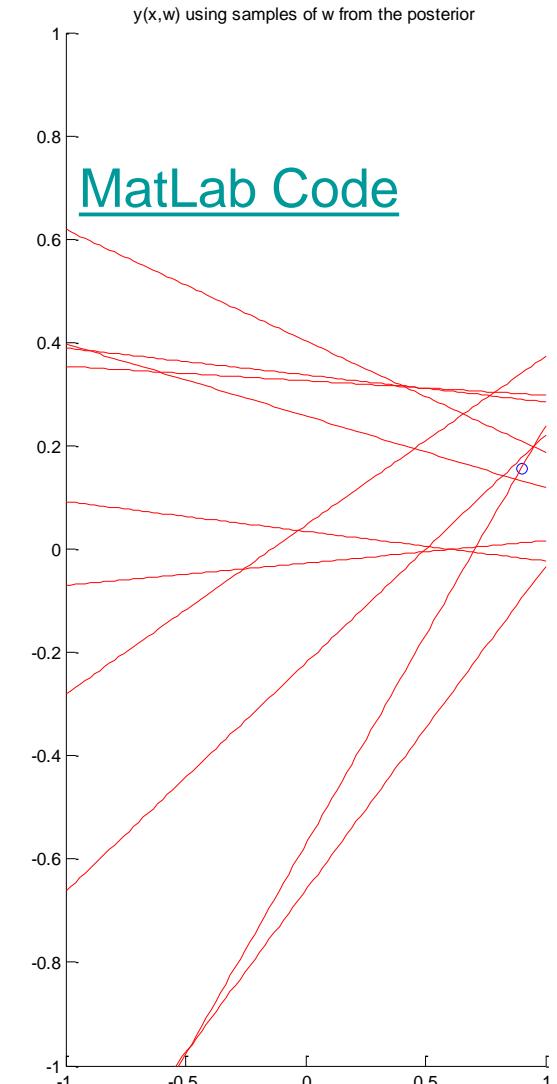
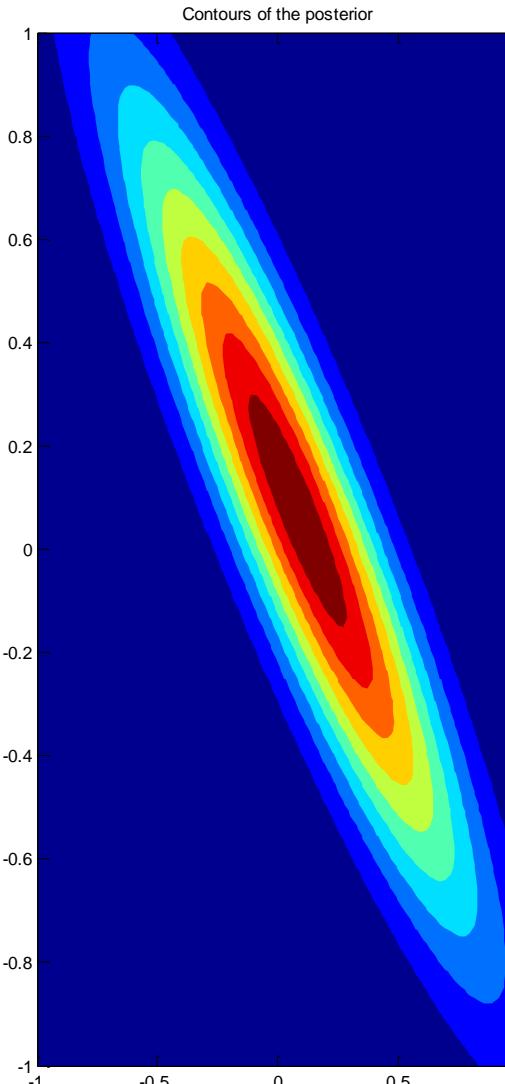
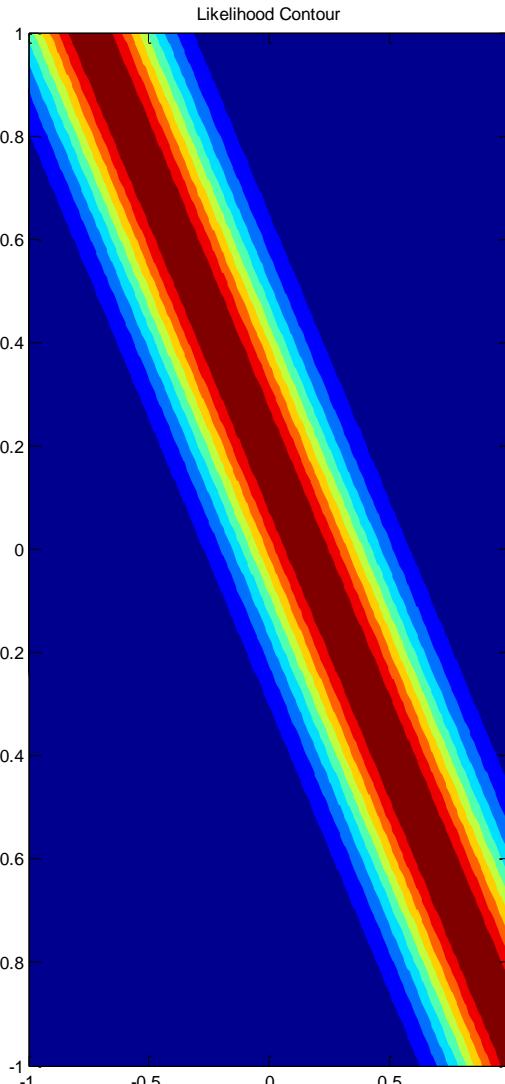
Bayesian Regression: Example



[MatLab Code](#)



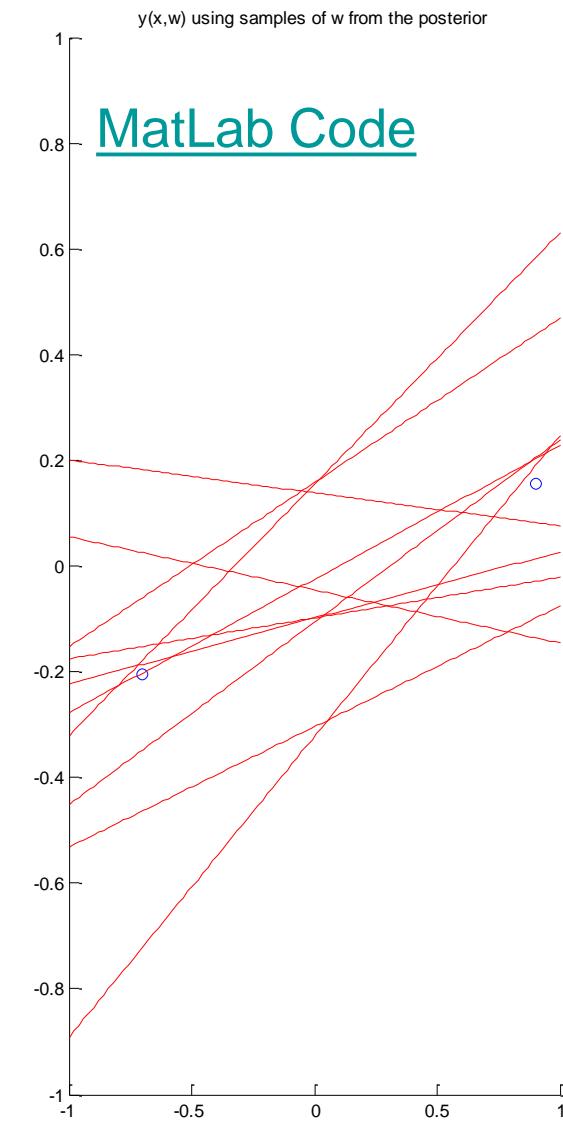
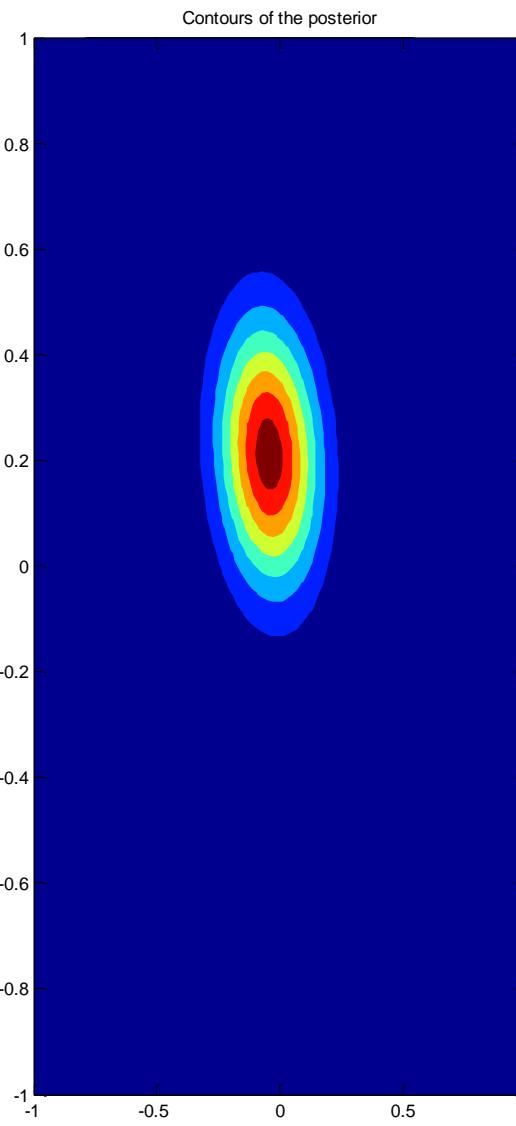
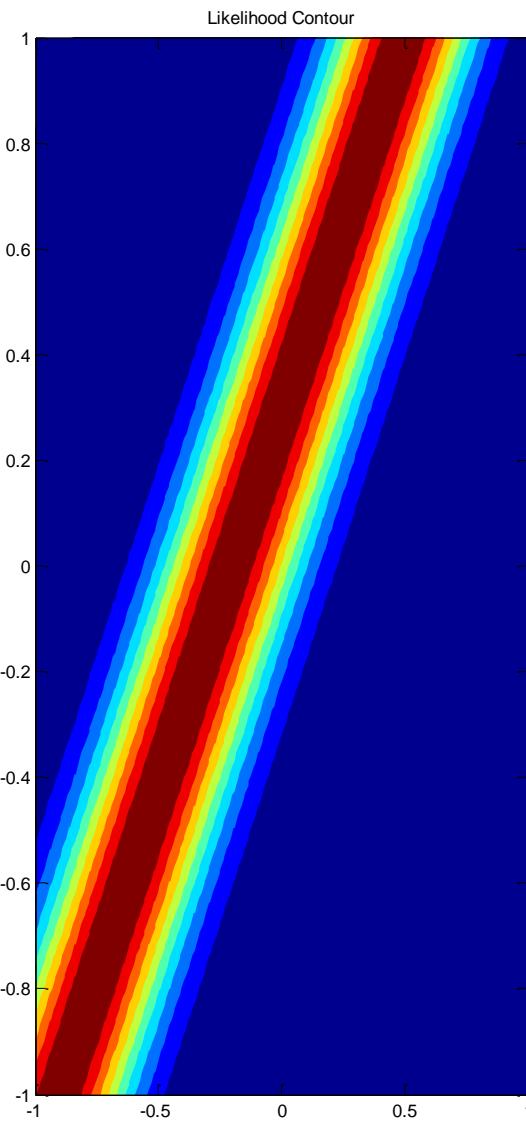
Example: One Data Point Collected



Note that the regression lines pass close to the data point (shown with a circle)



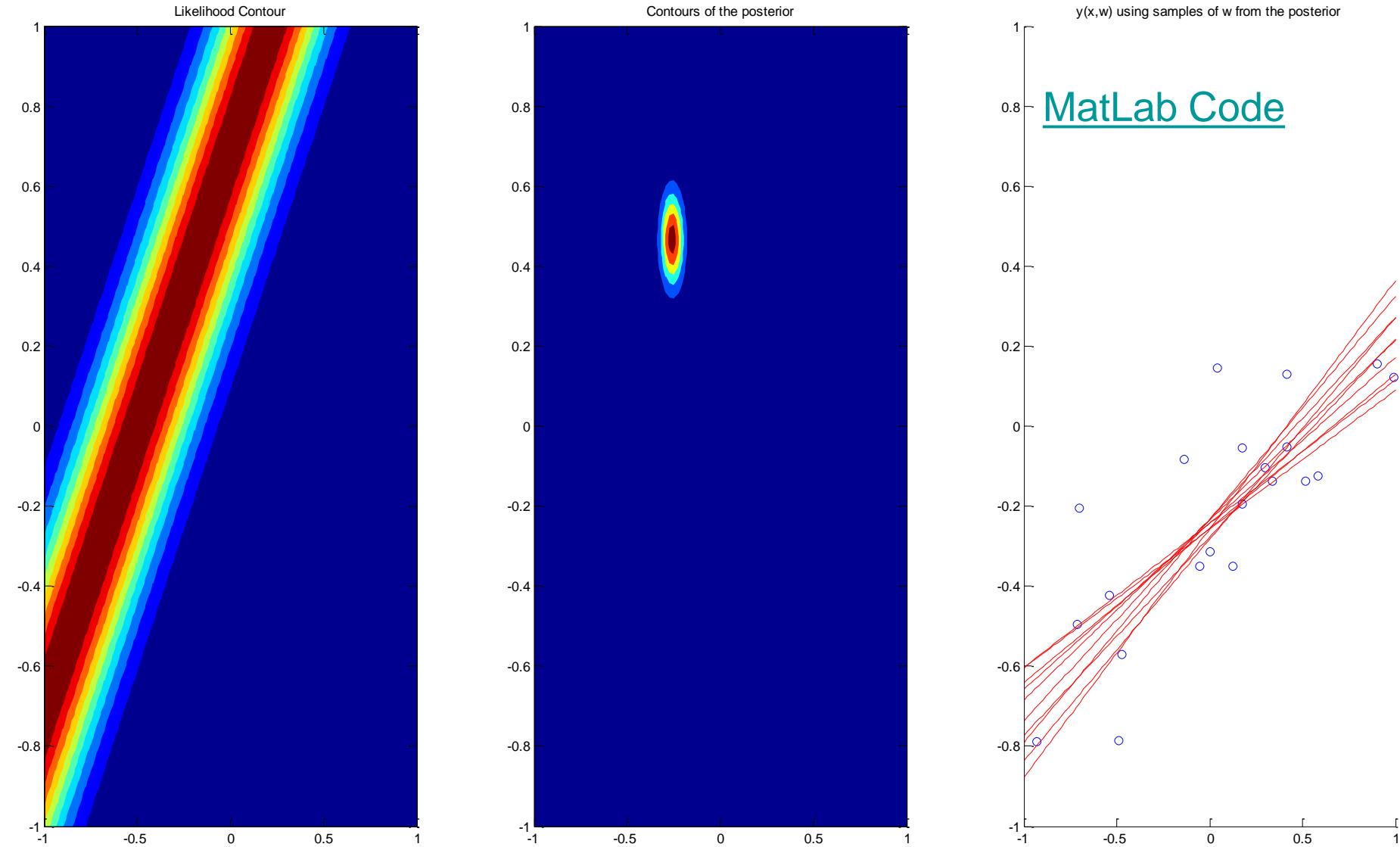
Example: 2nd Data Point Collected



Note that the regression lines now pass close to both data points



Example: 22 Data Points Collected

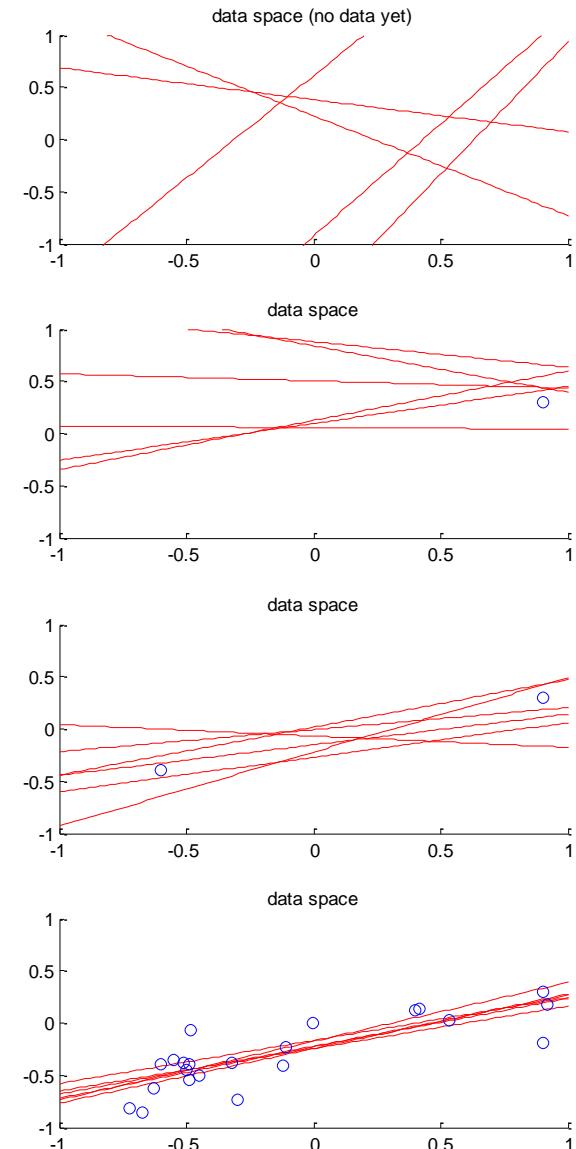
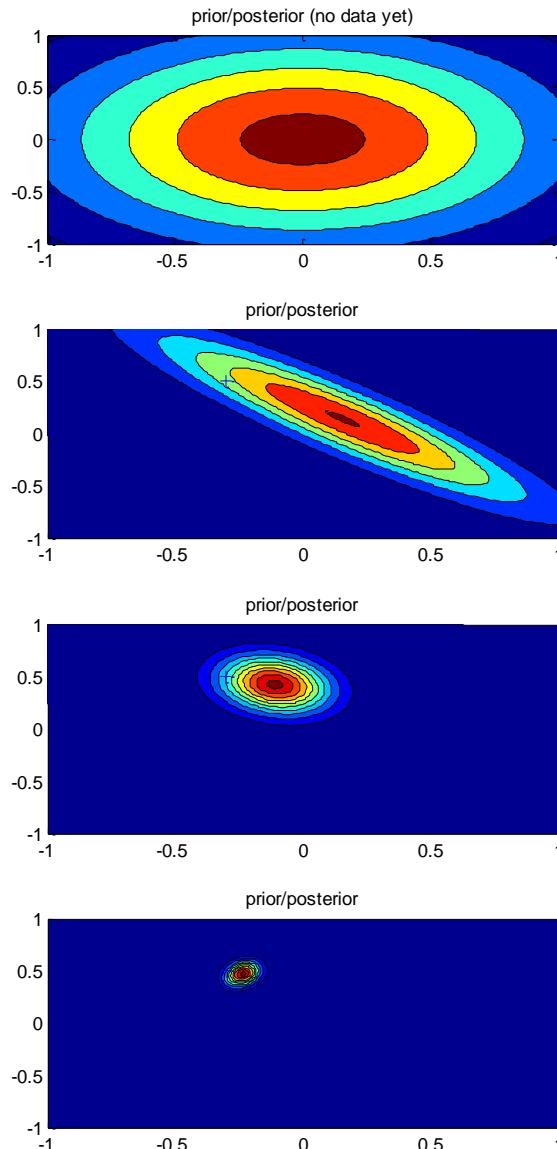
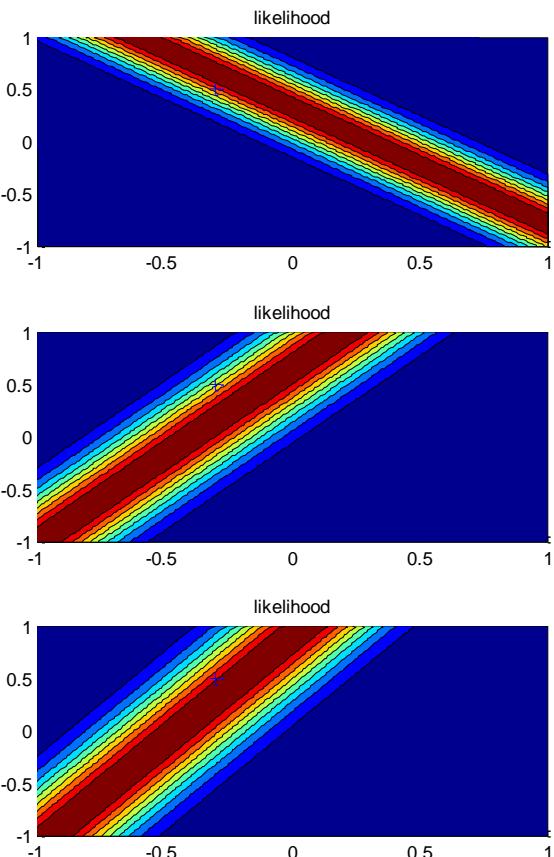


Note that the regression lines after 22 data points have been collected



Summary of Results

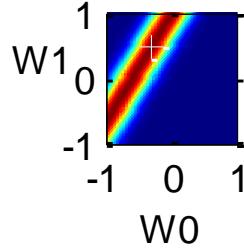
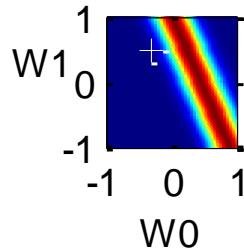
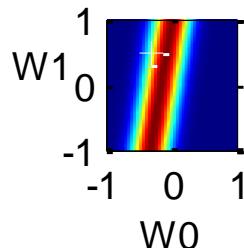
MatLab Code



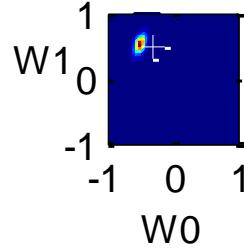
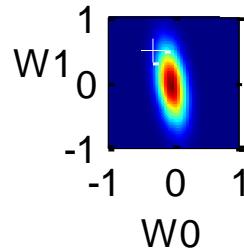
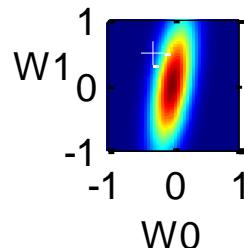
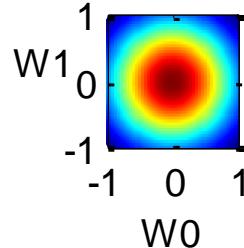
Summary of Results

likelihood

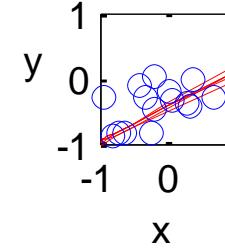
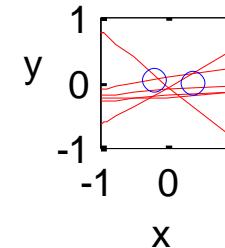
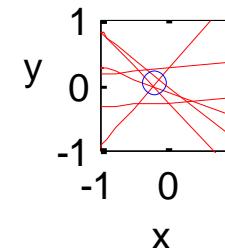
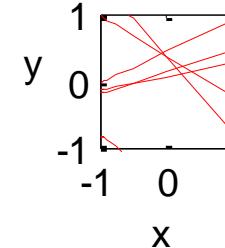
Run [bayesLinRegDemo2d](#)
from PMTK3



prior/posterior



data space



After 20 data points



Predictive Distribution

- We are not interested in w itself but in making predictions of t for new values of x . This requires the predictive distribution

$$p(t | x, \mathbf{x}, \mathbf{t}) = \int p(t | x, w, \beta) p(w | \mathbf{x}, \mathbf{t}, \alpha, \beta) dw = \mathcal{N}\left(t | \mathbf{m}_N^T \phi(x), \sigma_N^2(x)\right)$$

where $\sigma_N^2(x) = \frac{1}{\beta} + \phi(x)^T S_N \phi(x), \quad S_N^{-1} = \alpha \mathbf{I} + \beta \Phi^T \Phi$

- The 1st term represents the noise on the data whereas the 2nd term reflects the uncertainty associated with w .
- Because the noise process and the distribution of w are independent Gaussians, their variances are additive.
- The error bars get larger as we move away from the training points. By contrast, in the plug-in approximation, the error bars are of constant size.
- As additional data points are observed, the posterior distribution becomes narrower.



Predictive Distribution

- In a full Bayesian treatment, we want to compute the predictive distribution, i.e. given the training data \mathbf{x} and \mathbf{t} and a new test point x , we want the distribution:

$$p(t | \mathbf{x}, \mathbf{x}, \mathbf{t}) = \int p(t | \mathbf{x}, \mathbf{w}) p(\mathbf{w} | \mathbf{x}, \mathbf{t}) d\mathbf{w}, \text{ where}$$

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

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

$$\exp\left(-\frac{1}{2}\mathbf{w}^T \alpha \mathbf{I}_{M \times M} \mathbf{w} - \frac{1}{2}\mathbf{w}^T \sum_{n=1}^N \beta \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T \mathbf{w} - \beta \sum_{n=1}^N t_n \boldsymbol{\phi}(x_n) \mathbf{w}\right)$$

$$\propto \mathcal{N}\left(\mathbf{w} | \beta \mathbf{S}_N \sum_{n=1}^N t_n \boldsymbol{\phi}(x_n), \mathbf{S}_N^{-1}\right), \mathbf{S}_N^{-1} = \alpha \mathbf{I}_{M \times M} + \sum_{n=1}^N \beta \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T$$

- To compute the marginal, we use the standard equations for Gaussian Linear Systems from an earlier lecture.



Appendix: Useful Result

$$p(x) = \mathcal{N}(x | \mu, \Lambda^{-1})$$

$$p(y | x) = \mathcal{N}(y | Ax + b, L^{-1})$$

- For the above linear model, we proved in an earlier lecture that the following very useful results about marginal and conditional Gaussian models hold:

$$p(y) = \mathcal{N}(y | A\mu + b, L^{-1} + A\Lambda^{-1}A^T)$$

$$p(x | y) = \mathcal{N}(x | (\Lambda + A^T L A)^{-1} (\Lambda \mu + A^T L (y - b)), (\Lambda + A^T L A)^{-1})$$

Predictive Distribution

$$p(x) = \mathcal{N}(x | \mu, \Lambda^{-1})$$

$$p(y | x) = \mathcal{N}(y | Ax + b, L^{-1})$$

$$p(w | x, t, \alpha, \beta) = \mathcal{N}\left(w | \beta S_N \sum_{n=1}^N t_n \phi(x_n), S_N\right)$$

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

➤ Thus for our problem:

$$x \leftarrow w, \mu = \beta S_N \sum_{n=1}^N t_n \phi(x_n), \Lambda^{-1} = S_N$$

$$y \leftarrow t, A = \phi(x)^T, b = 0, L^{-1} = \beta^{-1}$$

➤ The predictive distribution now takes the form:

$$p(y) = \mathcal{N}(y | A\mu + b, L^{-1} + A\Lambda^{-1}A^T) \Rightarrow$$

$$p(t) = \mathcal{N}\left(t | \phi(x)^T \beta S_N \sum_{n=1}^N t_n \phi(x_n), \beta^{-1} + \phi(x)^T S_N \phi(x)\right)$$

Predictive Distribution

- In a full Bayesian treatment, we want to compute the predictive distribution, i.e. given the training data \mathbf{x} and \mathbf{t} and a new test point x , we want the distribution:

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

$$p(t | \mathbf{x}, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t | m(x), \sigma_N^2(x))$$

where the mean and variance are given by

$$m(x) = \beta \boldsymbol{\phi}(x)^T \mathbf{S}_N \sum_{n=1}^N \boldsymbol{\phi}(x_n) t_n = \boldsymbol{\phi}(x)^T \mathbf{m}_N, \quad \mathbf{m}_N = \beta \mathbf{S}_N \boldsymbol{\Phi}^T \mathbf{t}$$

$$\sigma_N^2(x) = \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S}_N \boldsymbol{\phi}(x) \text{ (*uncertainty in the data + uncertainty in w*)}$$

$$\mathbf{S}_N^{-1} = \alpha \mathbf{I} + \beta \sum_{n=1}^N \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T = \alpha \mathbf{I} + \beta \boldsymbol{\Phi}^T \boldsymbol{\Phi}$$

- Note that: $\sigma_N^2(x) = \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S}_N \boldsymbol{\phi}(x) \xrightarrow{N \rightarrow \infty} \beta^{-1}$ and $\sigma_{N+1}^2(x) \leq \sigma_N^2(x)$



Predictive Distribution

□ It is easy to show:

$$\sigma_N^2(x) = \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S}_N \boldsymbol{\phi}(x) \xrightarrow{N \rightarrow \infty} \beta^{-1} \text{ and } \sigma_{N+1}^2(x) \leq \sigma_N^2(x)$$

□ Note: $\mathbf{S}_{N+1}^{-1} = \alpha \mathbf{I} + \beta \sum_{n=1}^{N+1} \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T = \mathbf{S}_N^{-1} + \beta \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T$

and the identity:

$$(\mathbf{M} + \mathbf{v} \mathbf{v}^T)^{-1} = \mathbf{M}^{-1} - \frac{(\mathbf{M}^{-1} \mathbf{v})(\mathbf{v}^T \mathbf{M}^{-1})}{1 + \mathbf{v}^T \mathbf{M}^{-1} \mathbf{v}}$$

□ Using these results, we can write:

$$\begin{aligned} \sigma_{N+1}^2(x) &= \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S}_{N+1} \boldsymbol{\phi}(x) = \beta^{-1} + \boldsymbol{\phi}(x)^T (\mathbf{S}_N^{-1} + \beta \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T)^{-1} \boldsymbol{\phi}(x) = \beta^{-1} + \\ &\boldsymbol{\phi}(x)^T \left[\mathbf{S}_N - \frac{\beta (\mathbf{S}_N \boldsymbol{\phi}(x_n)) (\boldsymbol{\phi}(x_n)^T \mathbf{S}_N)}{1 + \beta \boldsymbol{\phi}(x_n)^T \mathbf{S}_N \boldsymbol{\phi}(x_n)} \right] \boldsymbol{\phi}(x) = \sigma_N^2(x) - \frac{\beta (\boldsymbol{\phi}(x_n)^T \mathbf{S}_N \boldsymbol{\phi}(x))^2}{1 + \beta \boldsymbol{\phi}(x_n)^T \mathbf{S}_N \boldsymbol{\phi}(x_n)} \leq \sigma_N^2(x) \end{aligned}$$

Predictive Distribution: Summary

- The notation used here is as follows:

$$p(t | x, \mathbf{x}, \mathbf{t}) = \mathcal{N}(t | m(x), \sigma_N^2(x))$$

$$\begin{aligned}m(x) &= \beta \boldsymbol{\phi}(x)^T \mathbf{S}_N \sum_{n=1}^N \boldsymbol{\phi}(x_n) t_n \\ \sigma_N^2(x) &= \beta^{-1} + \boldsymbol{\phi}(x)^T \mathbf{S}_N \boldsymbol{\phi}(x) \\ \mathbf{S}_N^{-1} &= \alpha \mathbf{I} + \beta \sum_{n=1}^N \boldsymbol{\phi}(x_n) \boldsymbol{\phi}(x_n)^T\end{aligned}$$

Note:
Predictive mean and variance are functions of x .

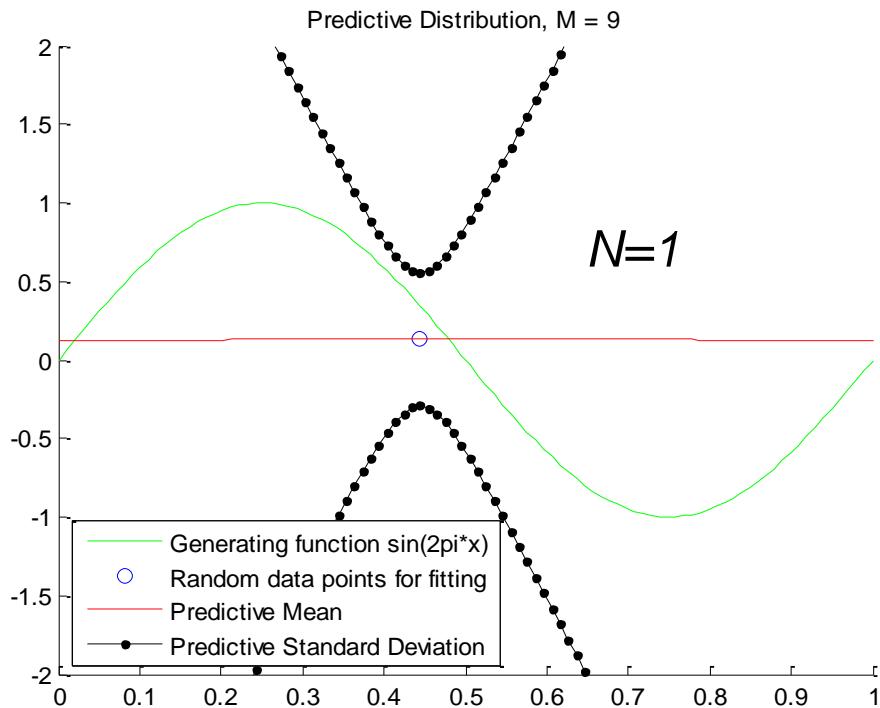
$$\boldsymbol{\phi}(x_n) = \begin{pmatrix} \phi_0(x_n) \\ \phi_1(x_n) \\ \phi_2(x_n) \\ \vdots \\ \phi_{M-1}(x_n) \end{pmatrix}, \boldsymbol{\phi}(x)^T = \{\phi_0(x_n) \quad \phi_1(x_n) \quad \phi_2(x_n) \quad \dots \quad \phi_{M-1}(x_n)\}, \mathbf{I} = \text{unit matrix } M \times M$$

For Polynomial regression:

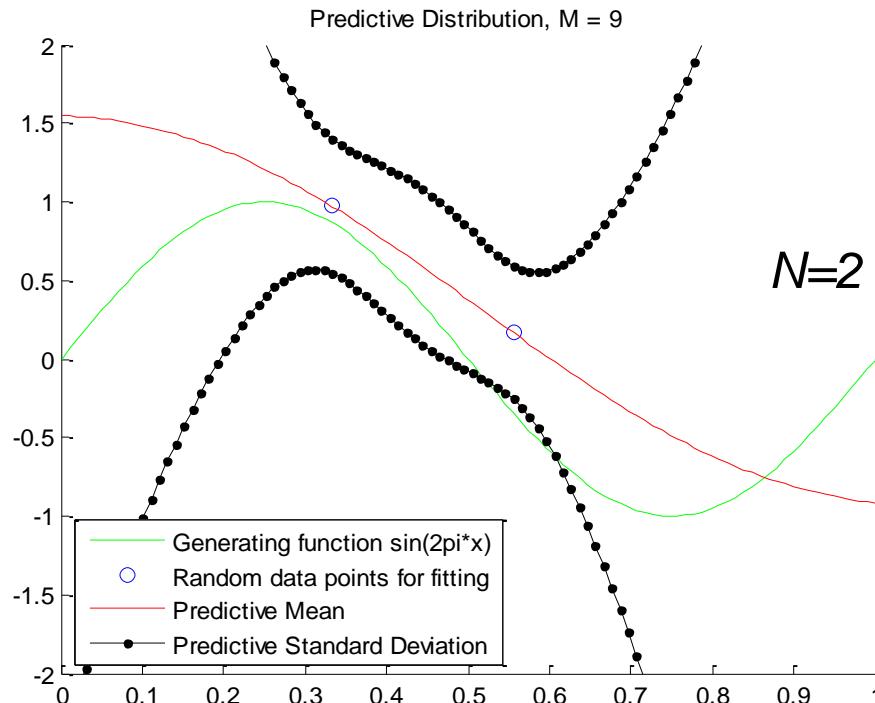
$$\boldsymbol{\phi}(x_n) = \begin{pmatrix} 1 \\ x_n \\ x_n^2 \\ \vdots \\ x_n^{M-1} \end{pmatrix}, \boldsymbol{\phi}(x)^T = \{1 \quad x \quad x^2 \quad \dots \quad x^{M-1}\}$$



Pointwise Uncertainty in the Predictions



$M=9$ Gaussians, 10 parameters
Scale of Gaussians adjusted
with data
 $\alpha = 5 \times 10^{-3}$
 $\beta = 11.1$
Using $N = 1, 2, 4, 10$
Data are given [here](#)

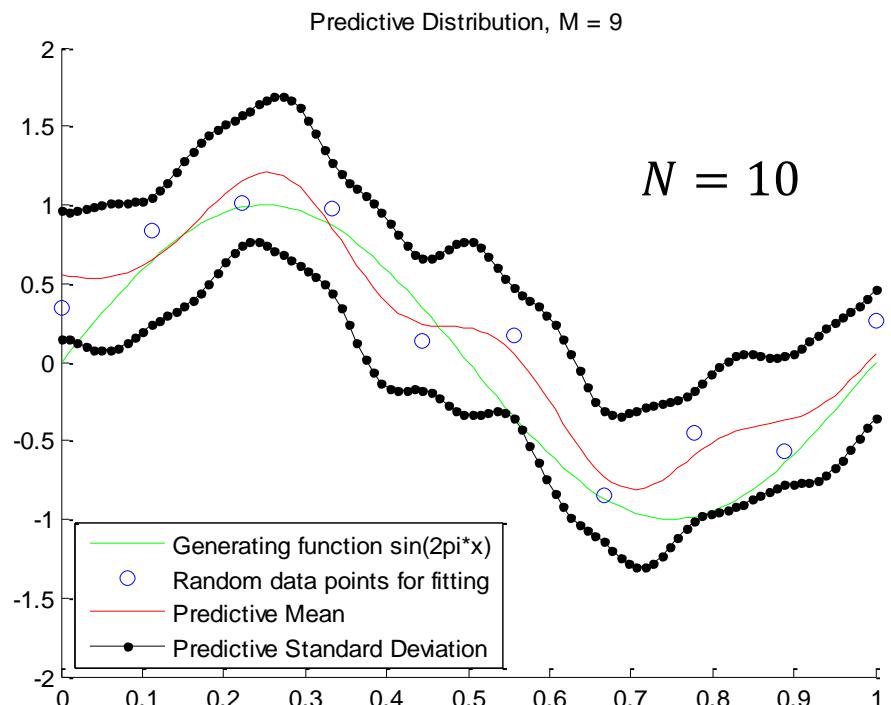
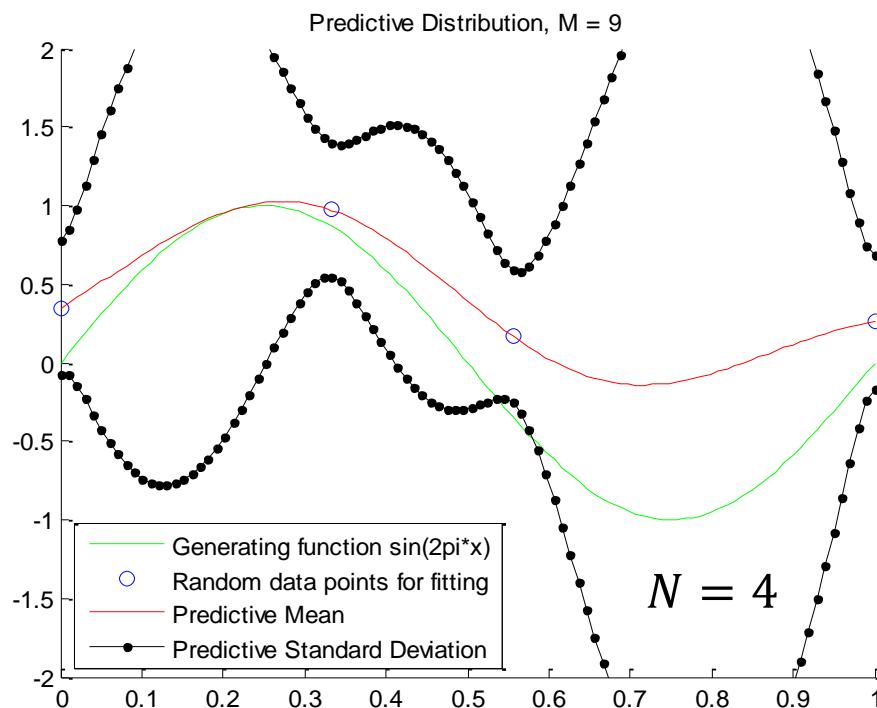


- The predictive uncertainty is smaller near the data.
- The level of uncertainty decreases with N

[MatLab code](#)



Pointwise Uncertainty in the Predictions

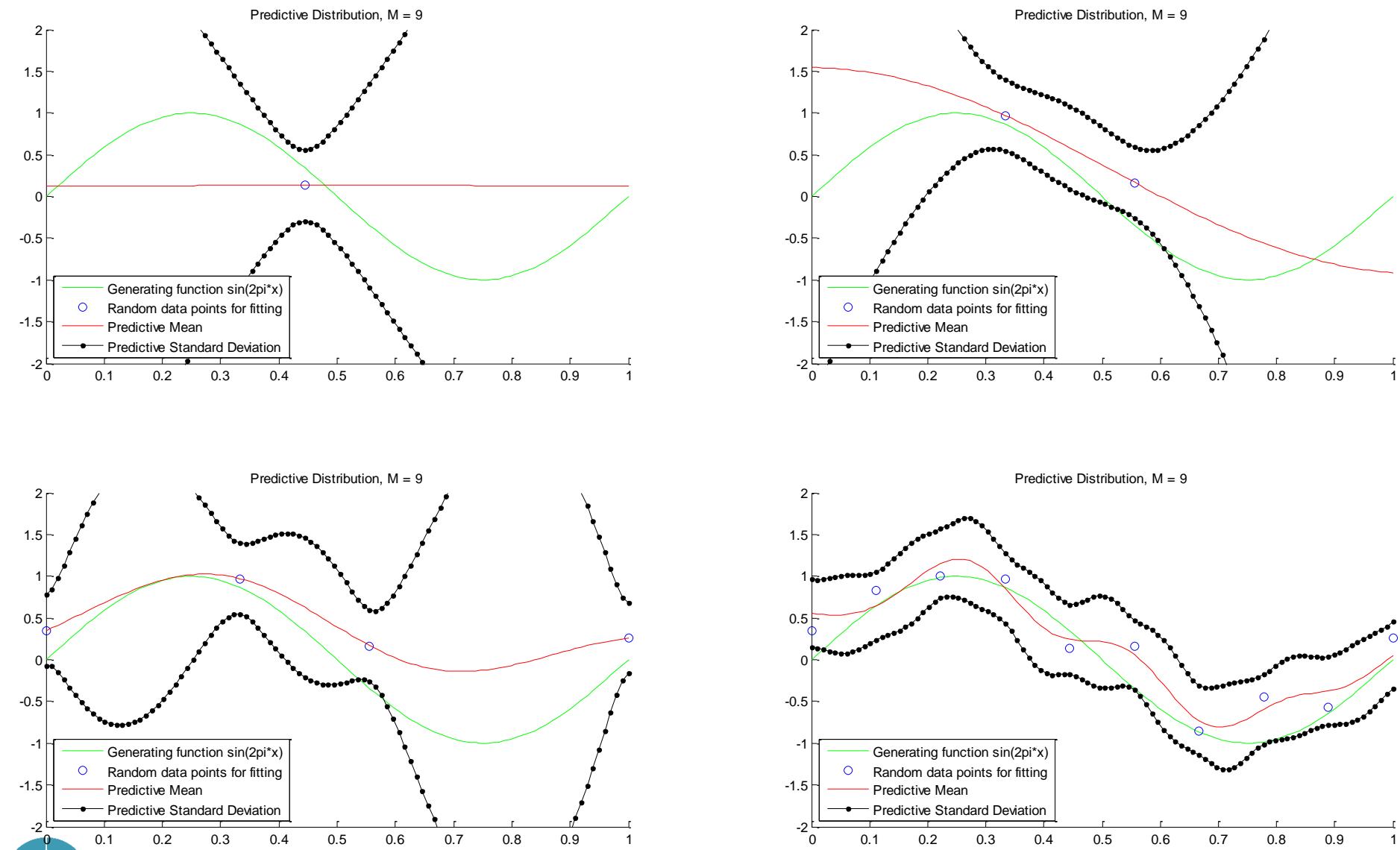


[MatLab code](#)



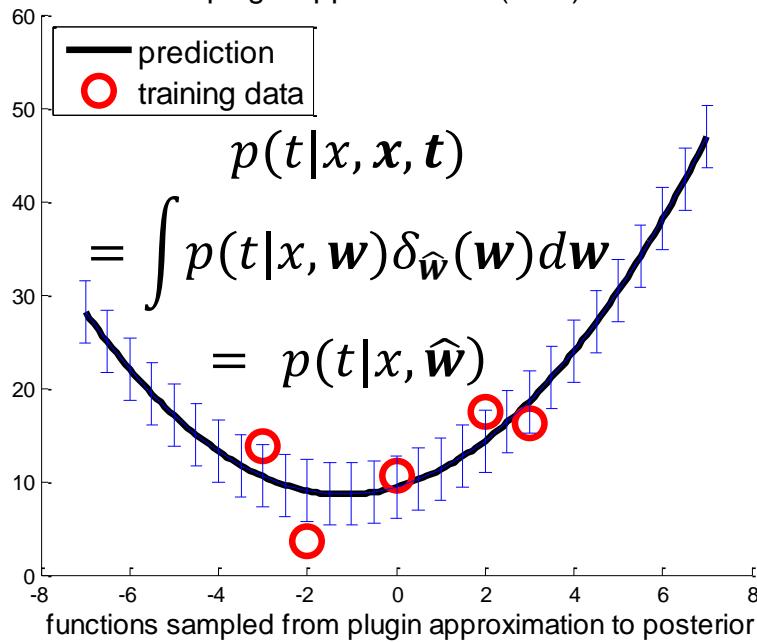
Summary of Results

MatLab code

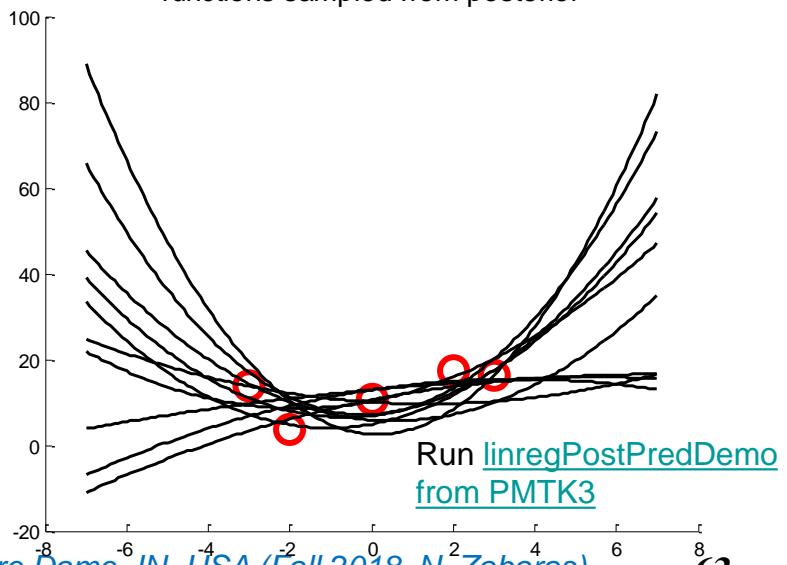
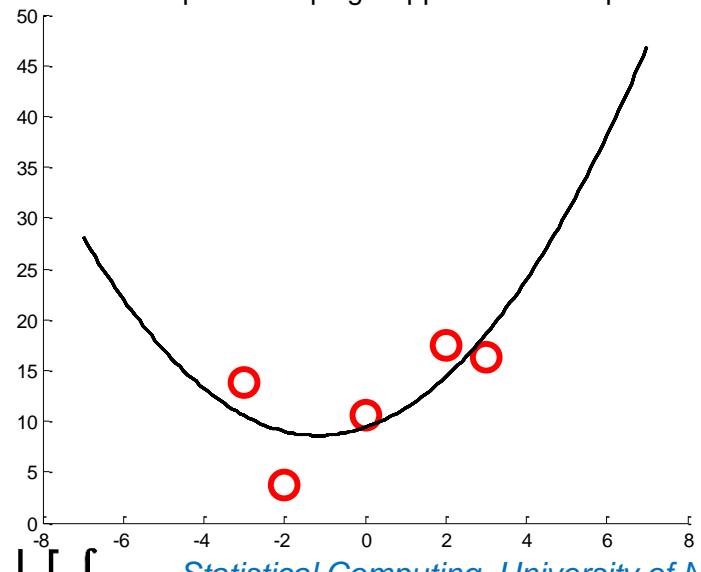
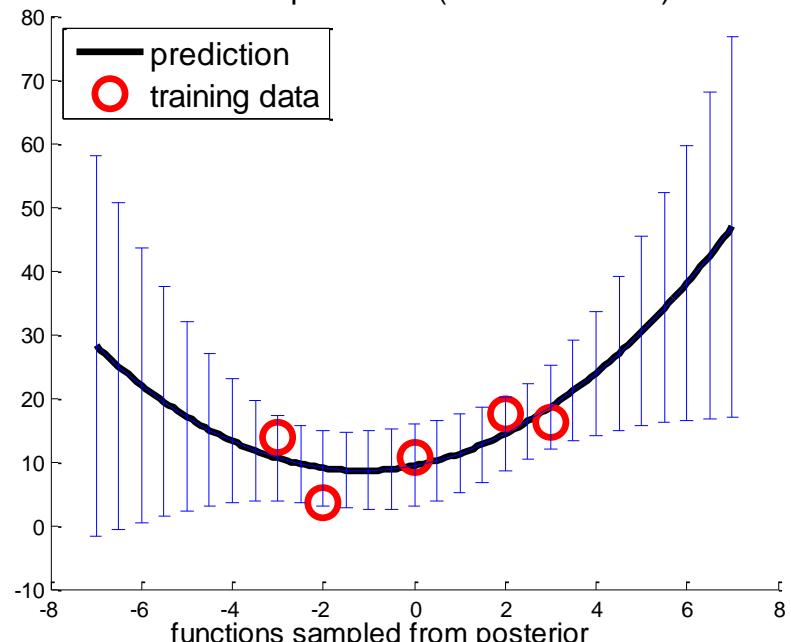


Plugin Approximation

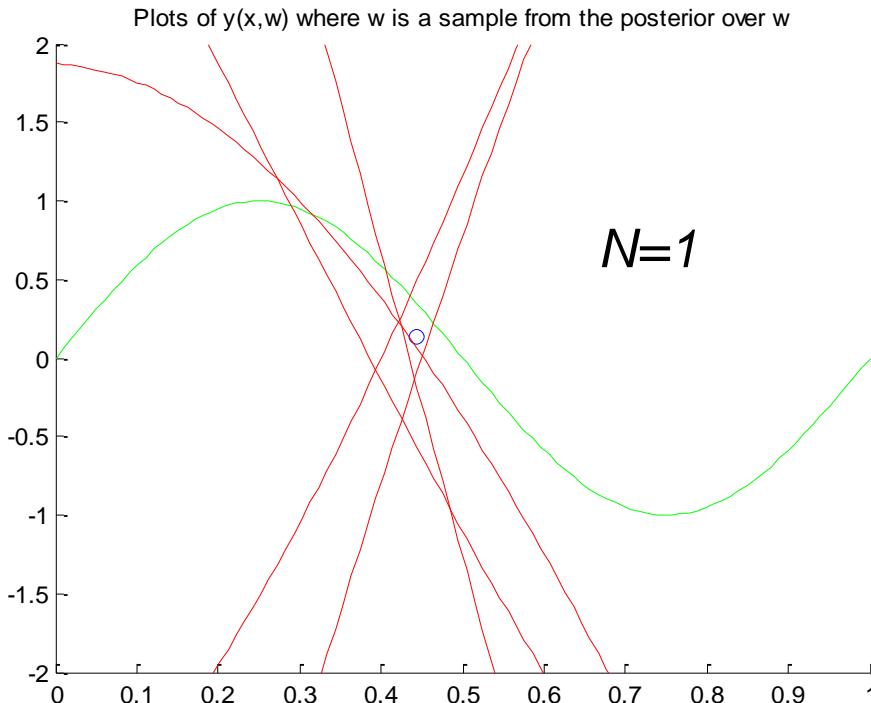
plugin approximation (MLE)



Posterior predictive (known variance)



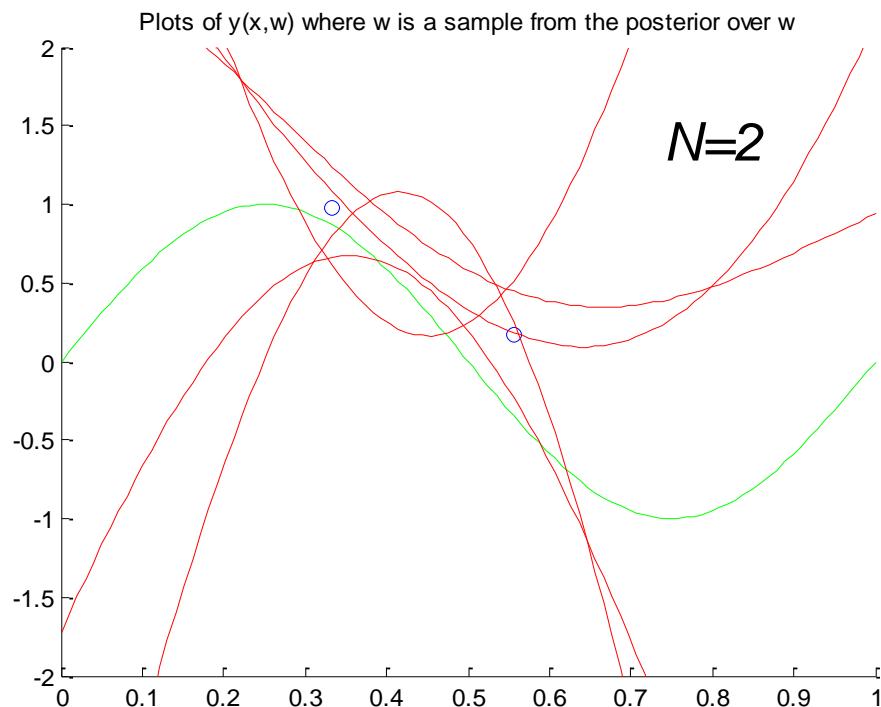
Covariance Between the Predictions



Same data and basis functions
as in the earlier example

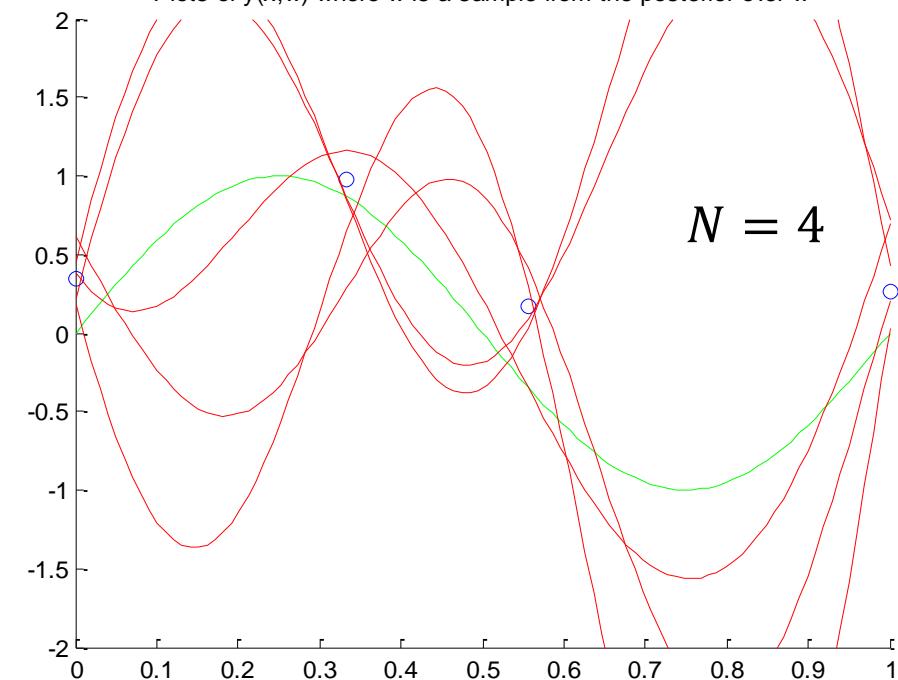
MatLab Code

- Draw samples from the posterior of w and then plot $y(x, w)$. We use the same data as the earlier example.
- We are visualizing the joint uncertainty in the posterior distribution between the y values at two or more x values.



Covariance Between the Predictions

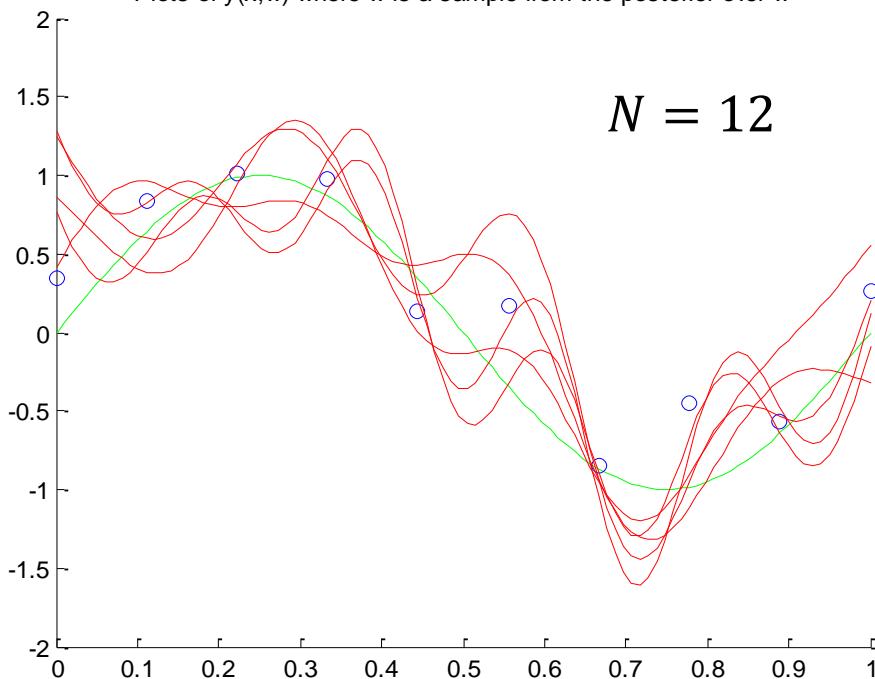
Plots of $y(x, w)$ where w is a sample from the posterior over w



$N = 4$

- Draw samples from the posterior of w and then plot $y(x, w)$
- We are visualizing the joint uncertainty in the posterior distribution between the y values at two or more x values.

Plots of $y(x, w)$ where w is a sample from the posterior over w



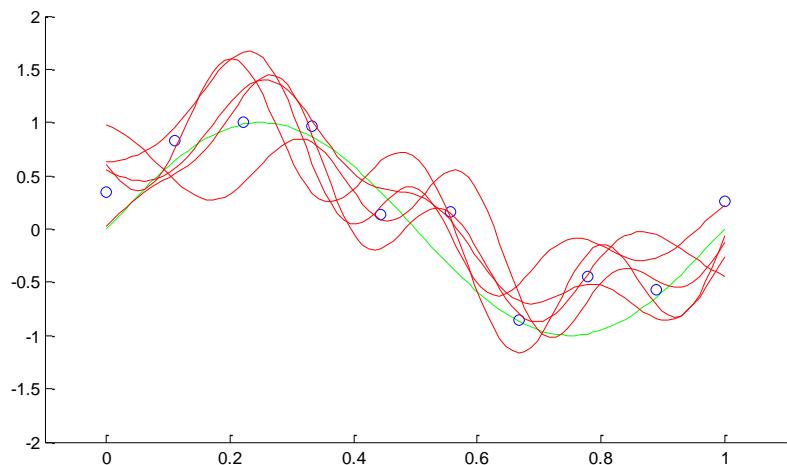
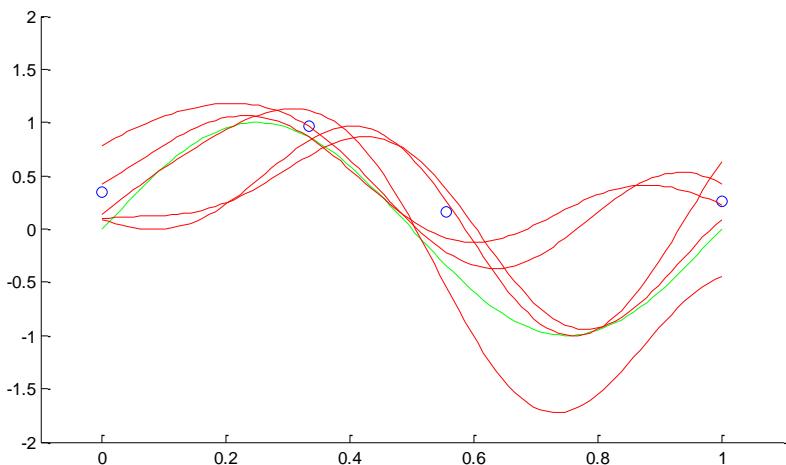
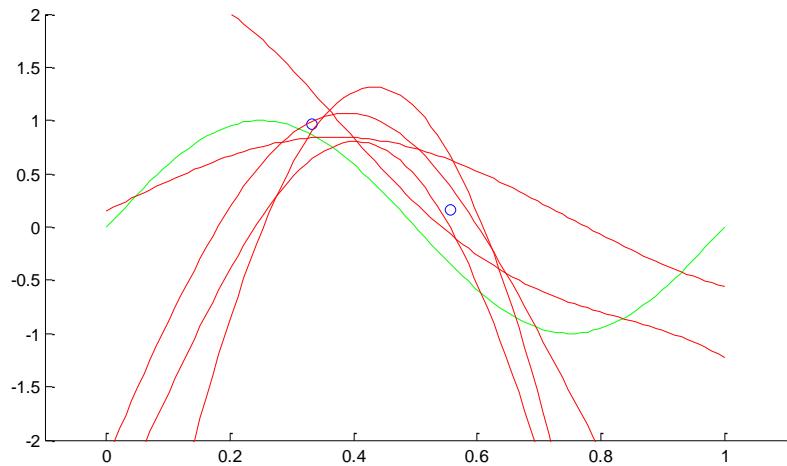
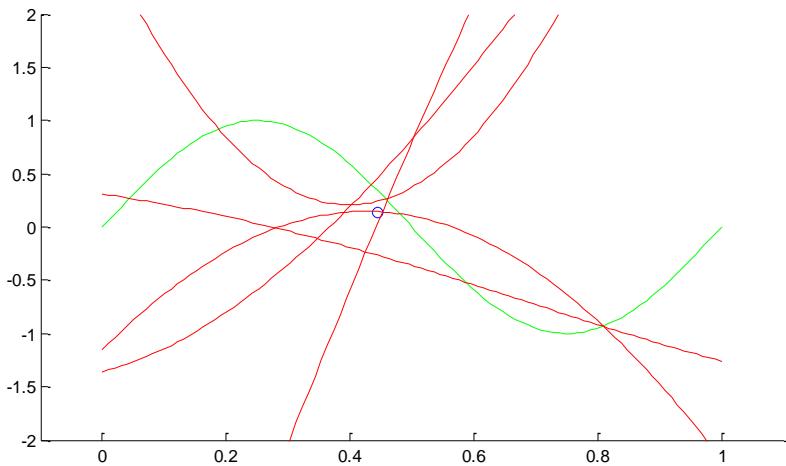
$N = 12$

MatLab Code



Summary of Results

MatLab Code



Gaussian Basis vs. Gaussian Process

- If we use localized basis functions such as Gaussians, then in regions away from the basis function support, the contribution from the second term in the predictive variance will go to zero, leaving only the noise contribution β^{-1} .

$$\sigma_N^2(x) = \beta^{-1} + \phi(x)^T S_N \phi(x) \xrightarrow[\text{support of } \phi(x)]{\text{away from the}} \beta^{-1}$$

- The model becomes very confident in its predictions when extrapolating outside the region occupied by the basis functions. This is an undesirable behavior.
- This problem can be avoided by adopting an alternative Bayesian approach to regression ([Gaussian processes](#)).



Bayesian Inference when σ^2 is Unknown

- Let us extend the previous results for linear regression assuming now that σ^2 is unknown.
- Assume a likelihood of the form:^{*}

$$p(\mathbf{y} | \mathbf{X}, \mathbf{w}, \sigma^2) = \mathcal{N}(\mathbf{y} | \mathbf{X}\mathbf{w}, \sigma^2 \mathbf{I}_N) = \frac{1}{(2\pi)^{N/2}} (\sigma^2)^{-N/2} \exp\left(-\frac{(\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w})}{2\sigma^2}\right)$$

- A conjugate prior has the following form:

$$p(\mathbf{w}, \sigma^2) = \mathcal{NIG}(\mathbf{w}, \sigma^2 | \mathbf{w}_0, \mathbf{V}_0, a_0, b_0) \triangleq \mathcal{N}(\mathbf{w} | \mathbf{w}_0, \sigma^2 \mathbf{V}_0) \mathcal{IG}(\sigma^2 | a_0, b_0) \\ \frac{b_0^{a_0}}{(2\pi)^{D/2} |\mathbf{V}_0|^{1/2} \Gamma(a_0)} (\sigma^2)^{-(a_0 + D/2 + 1)} \exp\left(-\frac{(\mathbf{w} - \mathbf{w}_0)^T \mathbf{V}_0^{-1} (\mathbf{w} - \mathbf{w}_0) + 2b_0}{2\sigma^2}\right)$$

- The posterior is now derived as:

$$p(\mathbf{w}, \sigma^2 | \mathcal{D}) = \frac{b_0^{a_0}}{(2\pi)^{(N+D)/2} |\mathbf{V}_0|^{1/2} \Gamma(a_0)} (\sigma^2)^{-(a_0 + (D+N)/2 + 1)} \exp\left(-\frac{(\mathbf{w} - \mathbf{w}_0)^T \mathbf{V}_0^{-1} (\mathbf{w} - \mathbf{w}_0) + (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + 2b_0}{2\sigma^2}\right)$$

- ❖ In the remaining of this lecture, the response is denoted as \mathbf{y} , the dimensionality of \mathbf{w} as D and the design matrix as \mathbf{X} .



Bayesian Inference when σ^2 is Unknown

$$p(\mathbf{w}, \sigma^2 | \mathcal{D}) = \frac{b_0^{a_0}}{(2\pi)^{(N+D)/2} |\mathbf{V}_0|^{1/2} \Gamma(a_0)} (\sigma^2)^{-(a_0 + (D+N)/2 + 1)} \exp\left(-\frac{(\mathbf{w} - \mathbf{w}_0)^T \mathbf{V}_0^{-1} (\mathbf{w} - \mathbf{w}_0) + (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + 2b_0}{2\sigma^2}\right)$$

□ Let us define the following:

$$\begin{aligned} \mathbf{V}_N &= (\mathbf{V}_0^{-1} + \mathbf{X}^T \mathbf{X})^{-1}, \quad \mathbf{w}_N = \mathbf{V}_N (\mathbf{V}_0^{-1} \mathbf{w}_0 + \mathbf{X}^T \mathbf{y}) \\ a_N &= a_0 + N/2, \quad b_N = b_0 + \frac{1}{2} (\mathbf{w}_0^T \mathbf{V}_0^{-1} \mathbf{w}_0 + \mathbf{y}^T \mathbf{y} - \mathbf{w}_N^T \mathbf{V}_N^{-1} \mathbf{w}_N) \end{aligned}$$

□ With these definitions, one with simple algebra can show:

$$\begin{aligned} p(\mathbf{w}, \sigma^2 | \mathcal{D}) &\propto (\sigma^2)^{-(a_N + D/2 + 1)} \exp\left(-\frac{(\mathbf{w} - \mathbf{w}_0)^T \mathbf{V}_0^{-1} (\mathbf{w} - \mathbf{w}_0) + (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + 2b_N - \mathbf{w}_0^T \mathbf{V}_0^{-1} \mathbf{w}_0 - \mathbf{y}^T \mathbf{y} + \mathbf{w}_N^T \mathbf{V}_N^{-1} \mathbf{w}_N}{2\sigma^2}\right) \\ &\propto (\sigma^2)^{-(a_N + D/2 + 1)} \exp\left(-\frac{(\mathbf{w} - \mathbf{w}_N)^T \mathbf{V}_N^{-1} (\mathbf{w} - \mathbf{w}_N) + 2b_N}{2\sigma^2}\right) \end{aligned}$$

$$p(\mathbf{w}, \sigma^2 | \mathcal{D}) = \mathcal{NIG}(\mathbf{w}, \sigma^2 | \mathbf{w}_N, \mathbf{V}_N, a_N, b_N) \triangleq \mathcal{N}(\mathbf{w} | \mathbf{w}_N, \sigma^2 \mathbf{V}_N) \mathcal{IG}(\sigma^2 | a_N, b_N)$$

□ The posterior marginals can now be derived explicitly:

$$p(\sigma^2 | \mathcal{D}) = \mathcal{IG}(\sigma^2 | a_N, b_N)$$

$$p(\mathbf{w} | \mathcal{D}) = \mathcal{T}_D\left(\mathbf{w}_N, \frac{b_N}{a_N} \mathbf{V}_N, 2a_N\right) \propto \left[1 + \frac{(\mathbf{w} - \mathbf{w}_N)^T \mathbf{V}_N^{-1} (\mathbf{w} - \mathbf{w}_N)}{2b_N}\right]^{-\frac{2a_N + D}{2}}$$

Posterior Marginals

$$p(\mathbf{w} | \mathcal{D}) \propto \int_0^{\infty} (\sigma^2)^{-(a_N + D/2 + 1)} \exp\left(-\frac{(\mathbf{w} - \mathbf{w}_N)^T \mathbf{V}_N^{-1} (\mathbf{w} - \mathbf{w}_N) + 2b_N}{2\sigma^2}\right) d\sigma^2 \propto \left[1 + \frac{(\mathbf{w} - \mathbf{w}_N)^T \mathbf{V}_N^{-1} (\mathbf{w} - \mathbf{w}_N)}{2b_N}\right]^{-\frac{2a_N + D}{2}}$$

- The marginal posterior can be directly written as:

$$p(\mathbf{w} | \mathcal{D}) = \mathcal{I}_D\left(\mathbf{w}_N, \frac{b_N}{a_N} \mathbf{V}_N, 2a_N\right) \propto \left[1 + \frac{(\mathbf{w} - \mathbf{w}_N)^T \mathbf{V}_N^{-1} (\mathbf{w} - \mathbf{w}_N)}{2b_N}\right]^{-\frac{2a_N + D}{2}}$$

- ❖ To compute the integral above, simply set $\lambda = \sigma^{-2}$, $d\sigma^2 = -\lambda^{-2}d\lambda$ and use the normalizing factor of the Gamma distribution $\int_0^{\infty} \lambda^{a-1} e^{-b\lambda} d\lambda = \Gamma(a)b^{-a} \sim b^{-a}$.



Posterior Predictive Distribution

- Consider the posterior predictive for m new test inputs:

$$p(\tilde{\mathbf{y}}|\tilde{\mathbf{X}}, \mathcal{D}) \propto \iint \frac{1}{(2\pi)^{m/2}} (\sigma^2)^{-m/2} \exp\left(-\frac{(\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\mathbf{w})^T(\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\mathbf{w})}{2\sigma^2}\right) (\sigma^2)^{-(a_N + D/2 + 1)} \exp\left(-\frac{(\mathbf{w} - \mathbf{w}_N)^T \mathbf{V}_N^{-1} (\mathbf{w} - \mathbf{w}_N) + 2b_N}{2\sigma^2}\right) d\mathbf{w} d\sigma^2$$

- As a first step, let us integrate in \mathbf{w} by writing:

$$\begin{aligned} & (\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\mathbf{w})^T(\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\mathbf{w}) + (\mathbf{w} - \mathbf{w}_N)^T \mathbf{V}_N^{-1} (\mathbf{w} - \mathbf{w}_N) + 2b_N = \\ & \left(\mathbf{w} - (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \mathbf{V}_N^{-1})^{-1} (\tilde{\mathbf{X}}^T \tilde{\mathbf{y}} + \mathbf{V}_N^{-1} \mathbf{w}_N)\right)^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \mathbf{V}_N^{-1}) \left(\mathbf{w} - (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \mathbf{V}_N^{-1})^{-1} (\tilde{\mathbf{X}}^T \tilde{\mathbf{y}} + \mathbf{V}_N^{-1} \mathbf{w}_N)\right) \\ & - (\tilde{\mathbf{X}}^T \tilde{\mathbf{y}} + \mathbf{V}_N^{-1} \mathbf{w}_N)^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \mathbf{V}_N^{-1})^{-T} (\tilde{\mathbf{X}}^T \tilde{\mathbf{y}} + \mathbf{V}_N^{-1} \mathbf{w}_N) + \mathbf{w}_N^T \mathbf{V}_N^{-1} \mathbf{w}_N + \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} + 2b_N \end{aligned}$$

These terms cancel out from the integration in \mathbf{w}

- Let us denote the last term in the Eq. above as

$$2\beta = -(\tilde{\mathbf{X}}^T \tilde{\mathbf{y}} + \mathbf{V}_N^{-1} \mathbf{w}_N)^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \mathbf{V}_N^{-1})^{-T} (\tilde{\mathbf{X}}^T \tilde{\mathbf{y}} + \mathbf{V}_N^{-1} \mathbf{w}_N) + \mathbf{w}_N^T \mathbf{V}_N^{-1} \mathbf{w}_N + \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} + 2b_N$$

Posterior Predictive Distribution

□ The posterior predictive

$$p(\tilde{\mathbf{y}}|\tilde{\mathbf{X}}, \mathcal{D}) \\ \propto \iint \frac{1}{(2\pi)^{m/2}} (\sigma^2)^{-m/2} \exp\left(-\frac{(\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\mathbf{w})^T(\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\mathbf{w})}{2\sigma^2}\right) (\sigma^2)^{-(a_N + D/2 + 1)} \exp\left(-\frac{(\mathbf{w} - \mathbf{w}_N)^T \mathbf{V}_N^{-1} (\mathbf{w} - \mathbf{w}_N) + 2b_N}{2\sigma^2}\right) d\mathbf{w} d\sigma^2$$

is now simplified using $\lambda = 1/\sigma^2$ and recalling the normalization of the Gamma distribution:

$$p(\tilde{\mathbf{y}}|\tilde{\mathbf{X}}, \mathcal{D}) \propto \int (\lambda)^{m/2+a_N-1} \exp(-\beta\lambda) d\lambda \sim \beta^{-(m/2+a_N)}$$

□ Substituting β and by comparing the two Eqs. one can verify:

$$p(\tilde{\mathbf{y}}|\tilde{\mathbf{X}}, \mathcal{D}) \propto \left(-(\tilde{\mathbf{X}}^T \tilde{\mathbf{y}} + \mathbf{V}_N^{-1} \mathbf{w}_N)^T (\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \mathbf{V}_N^{-1})^{-T} (\tilde{\mathbf{X}}^T \tilde{\mathbf{y}} + \mathbf{V}_N^{-1} \mathbf{w}_N) + \mathbf{w}_N^T \mathbf{V}_N^{-1} \mathbf{w}_N + \tilde{\mathbf{y}}^T \tilde{\mathbf{y}} + 2b_N \right)^{-\left(\frac{m}{2}+a_N\right)}$$
$$\propto \left(1 + \frac{(\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\mathbf{w}_N)^T \left(\frac{b_N}{a_N} (\mathbf{I}_m + \tilde{\mathbf{X}}\mathbf{V}_N \tilde{\mathbf{X}}^T) \right)^{-1} (\tilde{\mathbf{y}} - \tilde{\mathbf{X}}\mathbf{w}_N)}{2a_N} \right)^{-\left(\frac{m}{2}+a_N\right)}$$

← Use the Sherman Morrison Woodbury formula here to show that (symmetry of \mathbf{V}_0 is assumed)

$$(\mathbf{I}_m + \tilde{\mathbf{X}}\mathbf{V}_N \tilde{\mathbf{X}}^T)^{-1} = \mathbf{I}_m - \tilde{\mathbf{X}}(\tilde{\mathbf{X}}^T \tilde{\mathbf{X}} + \mathbf{V}_N^{-1})^{-1} \tilde{\mathbf{X}}^T$$



Bayesian Inference when σ^2 is Unknown

- The posterior predictive is also a Student's \mathcal{T} :

$$p(\tilde{\mathbf{y}}|\tilde{\mathbf{X}}, \mathcal{D}) = \mathcal{T}_m\left(\tilde{\mathbf{y}}|\tilde{\mathbf{X}}\mathbf{w}_N, \frac{b_N}{a_N}(\mathbf{I}_m + \tilde{\mathbf{X}}\mathbf{V}_N\tilde{\mathbf{X}}^T), 2a_N\right)$$

- The predictive variance has two terms

- $\frac{b_N}{a_N}\mathbf{I}_m$ due to the measurement noise
- and $\frac{b_N}{a_N}\tilde{\mathbf{X}}\mathbf{V}_N\tilde{\mathbf{X}}^T$ due to the uncertainty in \mathbf{w} . *The second term depends on how close a test input is to the training data.*



Zellner's G-Prior

$$p(\mathbf{w}, \sigma^2) = \mathcal{NIG}(\mathbf{w}, \sigma^2 | \mathbf{w}_0, \mathbf{V}_0, a_0, b_0) \triangleq \mathcal{N}(\mathbf{w} | \mathbf{w}_0, \sigma^2 \mathbf{V}_0) \mathcal{IG}(\sigma^2 | a_0, b_0)$$

- It is common to set $a_0 = b_0 = 0$, corresponding to an uninformative prior for σ^2 , and to set $\mathbf{w}_0 = \mathbf{0}$ and $\mathbf{V}_0 = g(\mathbf{X}^T \mathbf{X})^{-1}$ for any positive value g .
- This is called Zellner's **g-prior**. Here g plays a role analogous to $1/\lambda$ in ridge regression. However, the prior covariance is proportional to $(\mathbf{X}^T \mathbf{X})^{-1}$ rather than \mathbf{I} .

$$p(\mathbf{w}, \sigma^2) = \mathcal{NIG}(\mathbf{w}, \sigma^2 | 0, g(\mathbf{X}^T \mathbf{X})^{-1}, 0, 0) \triangleq \mathcal{N}(\mathbf{w} | 0, \sigma^2 g(\mathbf{X}^T \mathbf{X})^{-1}) \mathcal{IG}(\sigma^2 | 0, 0)$$

- This ensures that the posterior is invariant to scaling of the inputs.

- Zellner, A. (1986). [On assessing prior distributions and bayesian regression analysis with g-prior distributions](#). In [Bayesian inference and decision techniques](#), Studies of Bayesian and Econometrics and Statistics volume 6. North Holland.
- Minka, T. (2000b). [Bayesian linear regression](#). Technical report, MIT.



Unit Information Prior

$$p(\mathbf{w}, \sigma^2) = \mathcal{NIG}(\mathbf{w}, \sigma^2 | 0, g(\mathbf{X}^T \mathbf{X})^{-1}, 0, 0) \triangleq N(\mathbf{w} | 0, \sigma^2 g(\mathbf{X}^T \mathbf{X})^{-1}) InvGamma(\sigma^2 | 0, 0)$$

- We will see below that if we use an uninformative prior, the posterior precision given N measurements is $\mathbf{V}_N^{-1} = \mathbf{X}^T \mathbf{X}$.
- The **unit information prior** is defined to contain as much information as one sample.
- To create a *unit information prior* for linear regression, we need to use $\mathbf{V}_0^{-1} = \frac{1}{N} \mathbf{X}^T \mathbf{X}$ which is equivalent to the *g-prior* with $g = N$.
- Zellner's prior depends on the data: This is contrary to much of our Bayesian inference discussion!

- Kass, R. and L. Wasserman (1995). A reference bayesian test for nested hypotheses and its relationship to the schwarz criterio. *J. of the Am. Stat. Assoc.* 90(431), 928–934.



Uninformative Prior

- An uninformative prior can be obtained by considering the *uninformative limit of the conjugate g-prior, which corresponds to setting $g = \infty$* . This is equivalent to an improper \mathcal{NIG} prior with $\mathbf{w}_0 = 0$, $\mathbf{V}_0 = \infty\mathbf{I}$, $a_0 = 0$ and $b_0 = 0$, which gives $p(\mathbf{w}, \sigma^2) \propto \sigma^{-(D+2)}$.

$$p(\mathbf{w}, \sigma^2) = \mathcal{NIG}(\mathbf{w}, \sigma^2 | 0, \infty\mathbf{I}, 0, 0) \triangleq \mathcal{N}(\mathbf{w} | 0, \sigma^2 \infty I) \mathcal{IG}(\sigma^2 | 0, 0) \rightarrow \sigma^{-(D+2)}$$

- Alternatively, we can start with the *semi-conjugate prior $p(\mathbf{w}, \sigma^2) = p(\mathbf{w})p(\sigma^2)$, and take each term to its uninformative limit individually*, which gives $p(\mathbf{w}, \sigma^2) \propto \sigma^{-2}$.

This is equivalent to an improper \mathcal{NIG} prior with $\mathbf{w}_0 = 0$, $\mathbf{V} = \infty\mathbf{I}$, $a_0 = -D/2$ and $b_0 = 0$.

$$p(\mathbf{w}, \sigma^2) = \mathcal{NIG}(\mathbf{w}, \sigma^2 | 0, \infty\mathbf{I}, 0, 0) \triangleq \mathcal{N}(\mathbf{w} | 0, \sigma^2 \infty I) \mathcal{IG}\left(\sigma^2 | -\frac{D}{2}, 0\right) \rightarrow \sigma^{-2}$$

Uninformative Prior

- Using the uninformative prior, $p(\mathbf{w}, \sigma^2) \propto \sigma^{-2}$, the corresponding posterior and marginal posteriors are given by

$$p(\mathbf{w}, \sigma^2 | \mathcal{D}) = \mathcal{NIG}(\mathbf{w}, \sigma^2 | \mathbf{w}_N, \mathbf{V}_N, a_N, b_N)$$

$$p(\mathbf{w} | \mathcal{D}) = \mathcal{T}_D \left(\mathbf{w}_N, \frac{b_N}{a_N} \mathbf{V}_N, 2a_N \right) = \mathcal{T}_D \left(\mathbf{w} | \widehat{\mathbf{w}}_{MLE}, \frac{s^2}{N-D} \mathbf{C}, N - D \right)$$

$$\mathbf{V}_N = \mathbf{C} = (\mathbf{V}_0^{-1} + \mathbf{X}^T \mathbf{X})^{-1} \rightarrow (\mathbf{X}^T \mathbf{X})^{-1}, \quad \mathbf{w}_N = \mathbf{V}_N (\mathbf{V}_0^{-1} \mathbf{w}_0 + \mathbf{X}^T \mathbf{y}) \rightarrow (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \widehat{\mathbf{w}}_{MLE}$$
$$a_N = a_0 + N/2 = (N - D)/2,$$

$$\mathbf{b}_N = \mathbf{b}_0 + \frac{1}{2} (\mathbf{w}_0^T \mathbf{V}_0^{-1} \mathbf{w}_0 + \mathbf{y}^T \mathbf{y} - \mathbf{w}_N^T \mathbf{V}_N^{-1} \mathbf{w}_N) = s^2/2, \quad s^2 = (\mathbf{y} - \mathbf{X} \widehat{\mathbf{w}}_{MLE})^T (\mathbf{y} - \mathbf{X} \widehat{\mathbf{w}}_{MLE})$$

$$\mathbf{w}_N = \widehat{\mathbf{w}}_{MLE} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y}$$

- Note in the calculation of s^2 :

$$s^2 = (\mathbf{y} - \mathbf{X} \widehat{\mathbf{w}}_{MLE})^T (\mathbf{y} - \mathbf{X} \widehat{\mathbf{w}}_{MLE}) = (\mathbf{y} - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y})^T (\mathbf{y} - \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y})$$
$$= \mathbf{y}^T \mathbf{y} - \mathbf{y}^T \mathbf{X} (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{y}^T \mathbf{y} - \widehat{\mathbf{w}}_{MLE}^T (\mathbf{X}^T \mathbf{X}) (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{y} = \mathbf{y}^T \mathbf{y} - \widehat{\mathbf{w}}_{MLE}^T \mathbf{V}_N^{-1} \widehat{\mathbf{w}}_{MLE}$$

Frequentist Confidence Interval Vs. Bayesian Marginal Credible Interval

- The use of a (semi-conjugate) uninformative prior is quite interesting since the resulting posterior turns out to be equivalent to the results obtained from frequentist statistics.

$$p(w_j|D) = \mathcal{T}\left(w_j|\widehat{w}_j, \frac{C_{jj}s^2}{N-D}, N-D\right)$$

- This is equivalent to the sampling distribution of the MLE which is given by the following:

$$\frac{w_j - \widehat{w}_j}{s_j} \sim \mathcal{T}_{N-D}, s_j = \sqrt{\frac{C_{jj}s^2}{N-D}}$$

is the standard error of the estimated parameter.

- The frequentist confidence interval and the Bayesian marginal credible interval for the parameters are the same.

- Rice, J. (1995). *Mathematical statistics and data analysis*. Duxbury. 2nd edition (page 542)
- Casella, G. and R. Berger (2002). *Statistical inference*. Duxbury. 2nd edition (page 554)



The Caterpillar Example

- As a worked example of the uninformative prior, consider the [caterpillar dataset](#). We can compute the posterior mean and standard deviation, and the 95% credible intervals (CI) for the regression coefficients.

- The 95% credible intervals are identical to the 95% confidence intervals computed using standard frequentist methods.

	coeff	mean	stddev	95pc CI	sig
w0	10.998	3.06027	[4.652, 17.345]	*	
w1	-0.004	0.00156	[-0.008, -0.001]	*	
w2	-0.054	0.02190	[-0.099, -0.008]	*	
w3	0.068	0.09947	[-0.138, 0.274]		
w4	-1.294	0.56381	[-2.463, -0.124]	*	
w5	0.232	0.10438	[0.015, 0.448]	*	
w6	-0.357	1.56646	[-3.605, 2.892]		
w7	-0.237	1.00601	[-2.324, 1.849]		
w8	0.181	0.23672	[-0.310, 0.672]		
w9	-1.285	0.86485	[-3.079, 0.508]		
w10	-0.433	0.73487	[-1.957, 1.091]		

Run [linregBayesCaterpillar](#) from PMTK3

▪ Marin, J.-M. and C. Robert (2007). [Bayesian Core: a practical approach to computational Bayesian statistics](#). Springer.



The Caterpillar Example

- We can use these marginal posteriors to compute if the coefficients are significantly different from 0 -- check if its 95% CI excludes 0.
- The CIs for coefficients 0, 1, 2, 4, 5 are all significant.
- These results are the same as those produced by a frequentist approach using p-values at the 5% level.
- But note that the MLE does not even exist when $N < D$, so standard frequentist inference theory breaks down in this setting. Bayesian inference theory still works using proper priors.

- Maruyama, Y. and E. George (2008). [A g-prior extension for \$p > n\$](#) . Technical report, U. Tokyo.



Empirical Bayes for Linear Regression

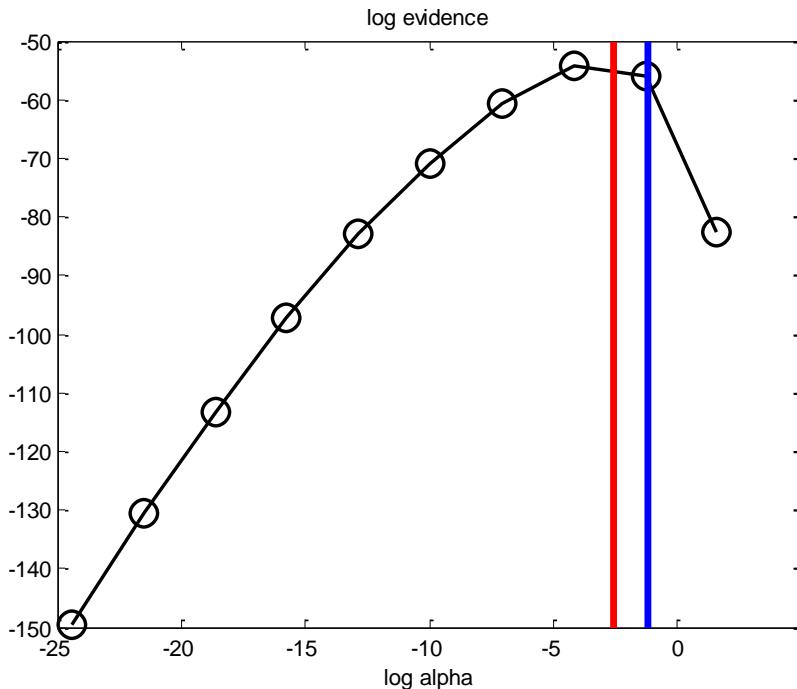
- We describe next an empirical Bayes procedure for picking the hyper-parameters in the prior (we will come back to this and relevance determination in a forthcoming lecture).
- More precisely, we choose $\eta = (\alpha, \lambda)$ to maximize the marginal likelihood, where $\lambda = 1/\sigma^2$ be the precision of the observation noise and α is the precision of the prior,
 $p(\mathbf{w}) = N(\mathbf{w}|\mathbf{0}, \alpha^{-1}\mathbf{I})$.
- This is known as the *evidence procedure*.

- MacKay, D. (1995b). [Probable networks and plausible predictions — a review of practical Bayesian methods for supervised neural networks](#). *Network*.
- Buntine, W. and A. Weigend (1991). [Bayesian backpropagation](#). *Complex Systems* 5, 603–643.
- MacKay, D. (1999). [Comparision of approximate methods for handling hyperparameters](#). *Neural Computation* 11(5), 1035–1068.



Empirical Bayes for Linear Regression

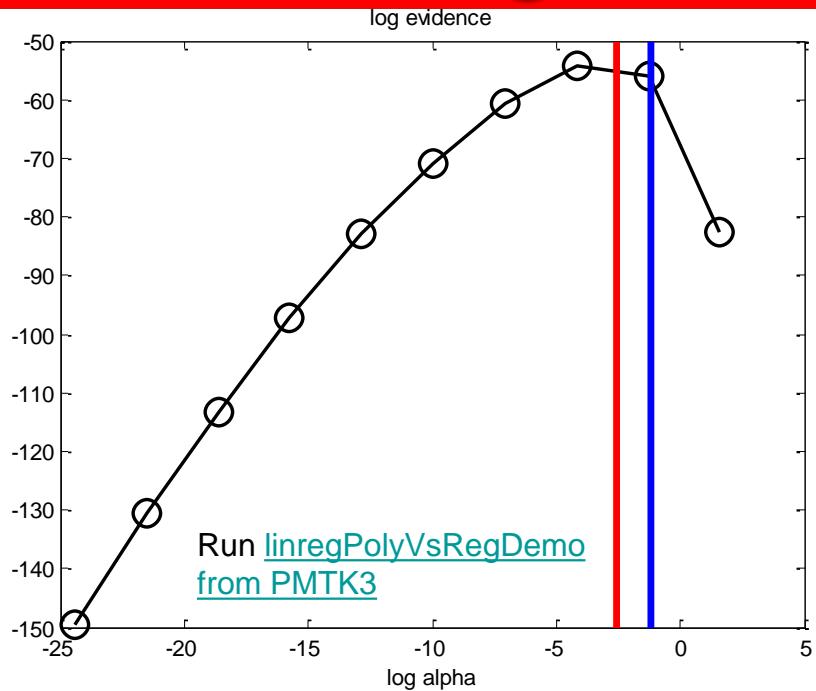
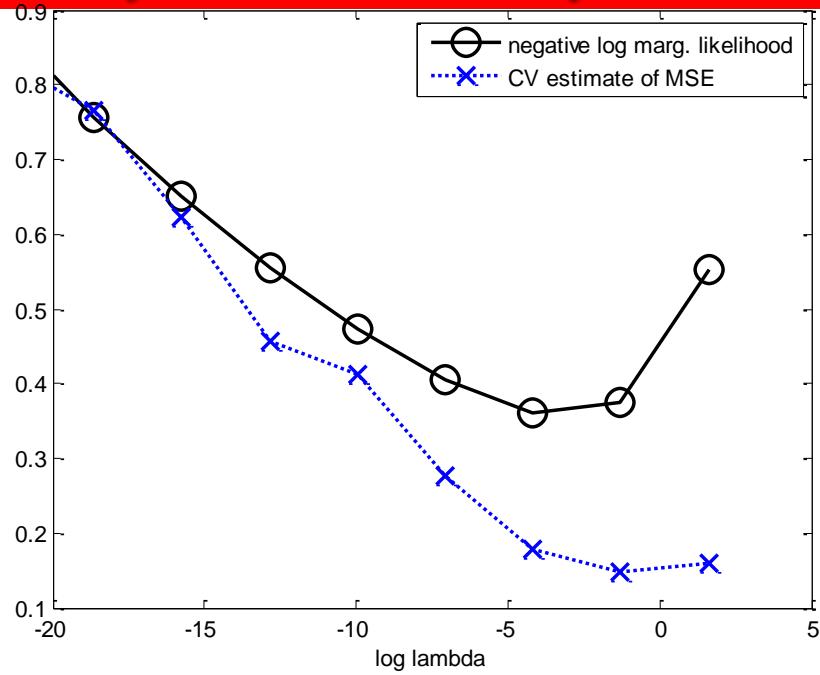
- The evidence procedure provides an alternative to using cross validation.
- In the Figure, the log marginal likelihood is plotted for different values of α , as well as the maximum value found by the optimizer.



Run [linregPolyVsRegDemo](#)
from PMTK3



Empirical Bayes for Linear Regression



- We obtain the same result as 5-CV ($\lambda = 1/\sigma^2$ is fixed in both methods).
- The key advantage of the evidence procedure over CV is that it allows different α_j to be used for every feature.

Automatic Relevancy Determination

- The evidence procedure can be used to perform feature selection (automatic relevancy determination or ARD)
- The evidence procedure is also useful when comparing different kinds of models:

$$\begin{aligned} p(\mathcal{D} | m) &= \iint p(\mathcal{D} | w, m) p(w | m, \eta) p(\eta | m) dw d\eta \\ &\approx \max_{\eta} p(\mathcal{D} | w, m) p(w | m, \eta) p(\eta | m) dw \end{aligned}$$

- It is important to (at least approximately) integrate over η rather than setting it arbitrarily.
- Using variation Bayes models our uncertainty on η rather than computing point estimates.

