

CS60020: Foundations of Algorithm Design and Machine Learning

Sourangshu Bhattacharya

**Some slides are taken from Christopher Bishop and
Geoffrey Hinton's courses**

OVERVIEW

What is Machine Learning?

It is **very hard** to write programs that solve problems like **recognizing a face**.

We don't know what program to write because we don't know **how our brain does it**.

Even if we had a good idea about how to do it, the program might be **horrendously complicated**.

Instead of writing a program by hand, we collect **lots of examples** that **specify the correct output** for a given input.

A machine learning algorithm then takes these examples and produces a program that does the job.

The program produced by the learning algorithm may look very different from a typical hand-written program. It may contain millions of numbers.

If we do it right, the program works for **new cases** as well as the ones we trained it on.

A classic example of a task that requires machine learning: It is very hard to say what makes a 2

0 0 0 1 1 1 1 1 1 2

2 2 2 2 2 2 2 3 3 3

3 4 4 4 4 4 5 5 5 5

6 6 7 7 7 7 7 8 8 8

— 9 9 9 9 9 9 9 9 —

Some more examples of tasks that are best solved by using a learning algorithm

Recognizing patterns:

- Facial identities or facial expressions

- Handwritten or spoken words

- Medical images

Generating patterns:

- Generating images or motion sequences

Recognizing anomalies:

- Unusual sequences of credit card transactions

- Unusual patterns of sensor readings in a nuclear power plant or unusual sound in your car engine.

Prediction:

- Future stock prices or currency exchange rates

Some web-based examples of machine learning

The web contains a lot of data. Tasks with very big datasets often use machine learning

especially if the data is noisy or non-stationary.

Spam filtering, fraud detection:

The enemy adapts so we must adapt too.

Recommendation systems:

Lots of noisy data. Million dollar prize!

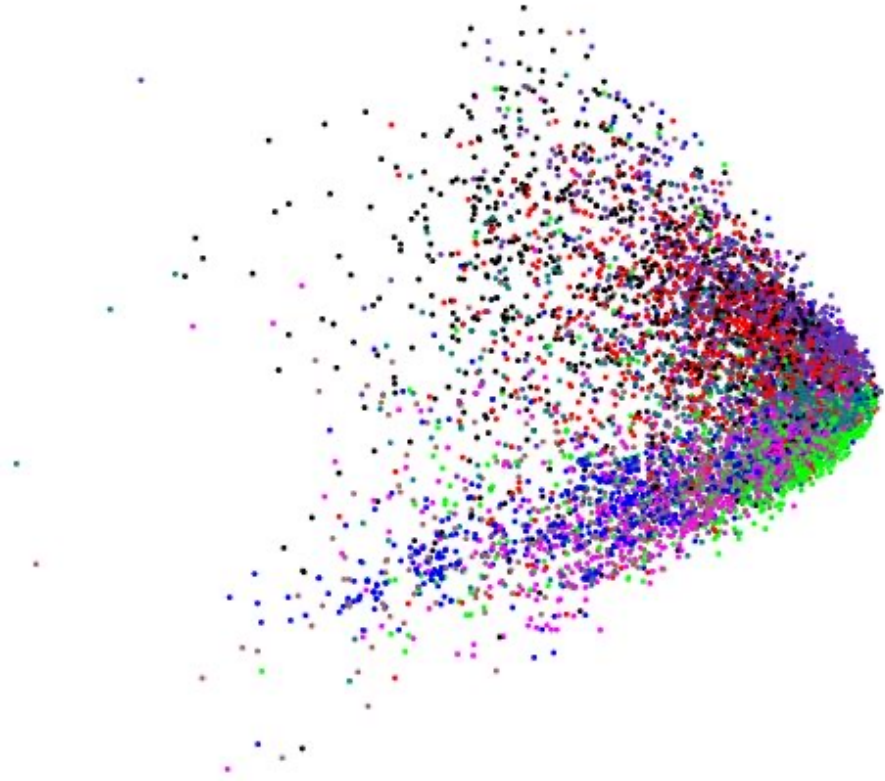
Information retrieval:

Find documents or images with similar content.

Data Visualization:

Display a huge database in a revealing way

Displaying the structure of a set of documents using Latent Semantic Analysis (a form of PCA)



Each document is converted to a vector of word counts. This vector is then mapped to two coordinates and displayed as a colored dot. The colors represent the hand-labeled classes.

When the documents are laid out in 2-D, the classes are not used. So we can judge how good the algorithm is by seeing if the classes are separated.

Machine Learning & Symbolic AI

Knowledge Representation works with facts/assertions and develops rules of logical inference. The rules can handle quantifiers. Learning and uncertainty are usually ignored.

Expert Systems used logical rules or conditional probabilities provided by “experts” for specific domains.

Graphical Models treat uncertainty properly and allow learning (but they often ignore quantifiers and use a fixed set of variables)

Set of logical assertions → values of a subset of the variables and local models of the probabilistic interactions between variables.

Logical inference → probability distributions over subsets of the unobserved variables (or individual ones)

Learning = refining the local models of the interactions.

Machine Learning & Statistics

A lot of machine learning is just a rediscovery of things that statisticians already knew.

But the emphasis is very different:

A good piece of statistics: Clever proof that a relatively simple estimation procedure is asymptotically unbiased.

A good piece of machine learning: Demonstration that a complicated algorithm produces impressive results on a specific task.

Data-mining: Using very simple machine learning techniques on very large databases because computers are too slow to do anything more interesting with ten billion examples.

A spectrum of machine learning tasks

Statistics-----Artificial Intelligence

Low-dimensional data (e.g. less than 100 dimensions)

Lots of noise in the data

There is not much structure in the data, and what structure there is, can be represented by a fairly simple model.

The main problem is distinguishing true structure from noise.

High-dimensional data (e.g. more than 100 dimensions)

The noise is not sufficient to obscure the structure in the data if we process it right.

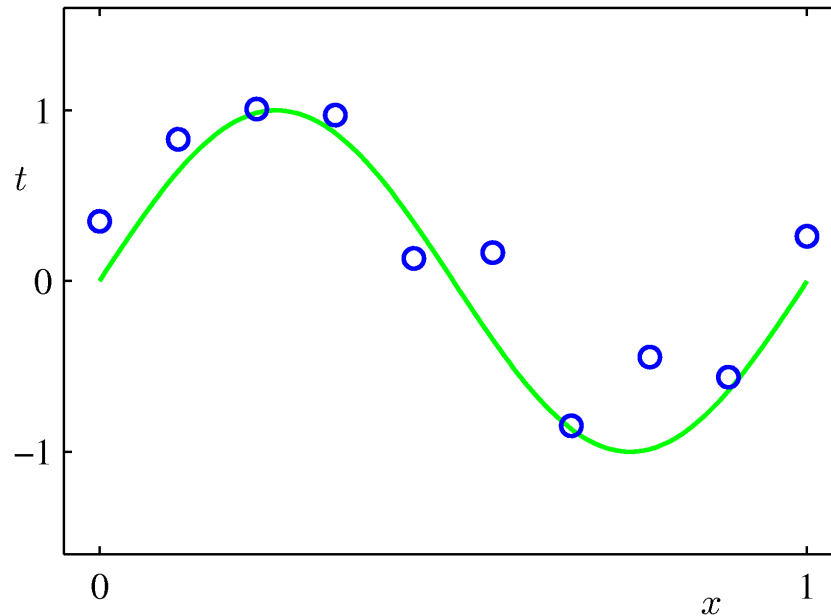
There is a huge amount of structure in the data, but the structure is too complicated to be represented by a simple model.

The main problem is figuring out a way to represent the complicated structure that allows it to be learned.

REGRESSION

Linear Basis Function Models (1)

Example: Polynomial Curve Fitting



$$y(x, \mathbf{w}) = w_0 + w_1x + w_2x^2 + \dots + w_Mx^M = \sum_{j=0}^M w_jx^j$$

Linear Basis Function Models (2)

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.

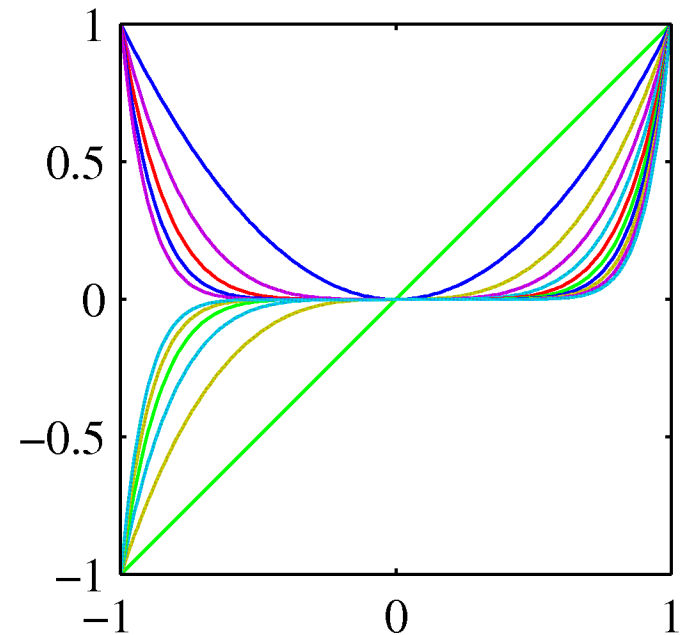
In the simplest case, we use linear basis functions : $\phi_d(\mathbf{x}) = x_d$.

Linear Basis Function Models (3)

Polynomial basis functions:

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

These are global; a small change in x affect all basis functions.

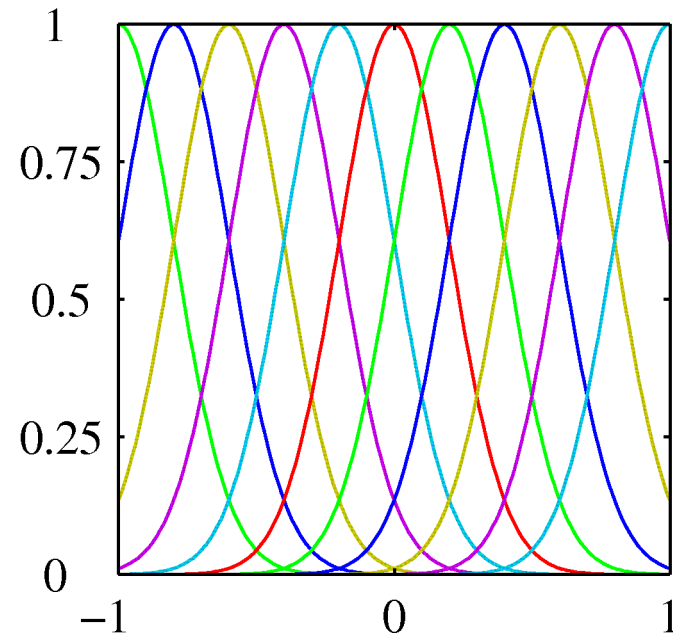


Linear Basis Function Models (4)

Gaussian basis functions:

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

These are local; a small change in x only affect nearby basis functions. μ_j and s control location and scale (width).



Linear Basis Function Models (5)

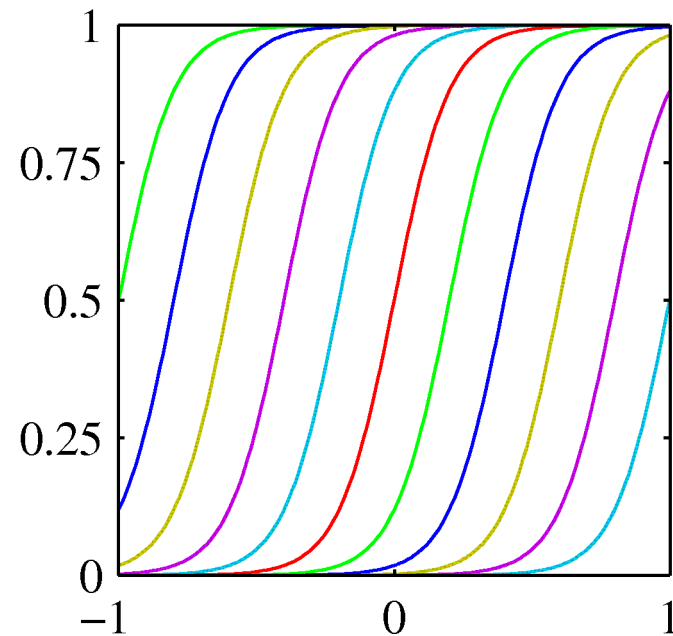
Sigmoidal basis functions:

$$\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right)$$

where

$$\sigma(a) = \frac{1}{1 + \exp(-a)}.$$

Also these are local; a small change in x only affect nearby basis functions. μ_j and s control location and scale (slope).

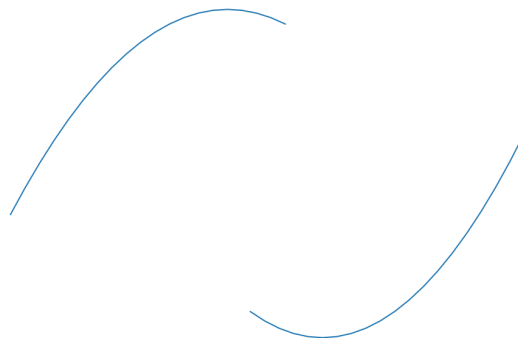


Least Squares Estimation

A polynomial curve is represented by the parameters w .

$$f(x) = x - x^2$$

$$f(x) = x + x^2$$



Error (loss) function for a given parameter:

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

Estimate $w^* = \min_w E(w)$

Maximum Likelihood and Least Squares (1)

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,

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

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 Least Squares (2)

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 (3)

Computing the gradient and setting it to zero yields

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

Solving for \mathbf{w} , we get

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

The Moore-Penrose
pseudo-inverse, Φ^\dagger .

where

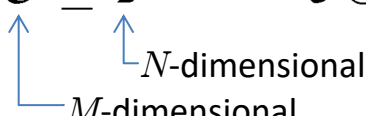
$$\Phi = \begin{pmatrix} \phi_0(\mathbf{x}_1) & \phi_1(\mathbf{x}_1) & \cdots & \phi_{M-1}(\mathbf{x}_1) \\ \phi_0(\mathbf{x}_2) & \phi_1(\mathbf{x}_2) & \cdots & \phi_{M-1}(\mathbf{x}_2) \\ \vdots & \vdots & \ddots & \vdots \\ \phi_0(\mathbf{x}_N) & \phi_1(\mathbf{x}_N) & \cdots & \phi_{M-1}(\mathbf{x}_N)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

Geometry of Least Squares

Consider

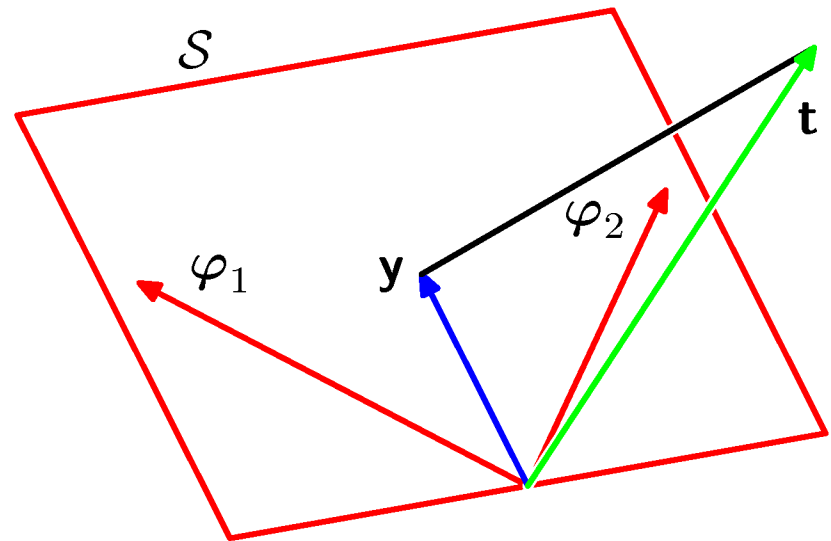
$$\mathbf{y} = \Phi \mathbf{w}_{\text{ML}} = [\varphi_1, \dots, \varphi_M] \mathbf{w}_{\text{ML}}.$$

$$\mathbf{y} \in \mathcal{S} \subseteq \mathcal{T} \quad \mathbf{t} \in \mathcal{T}$$


M-dimensional N-dimensional

\mathcal{S} is spanned by $\varphi_1, \dots, \varphi_M$.

\mathbf{w}_{ML} minimizes the distance between \mathbf{t} and its orthogonal projection on \mathcal{S} , i.e. \mathbf{y} .



Normal Equations

$$\underset{\substack{p \times p \quad p \times 1}}{(\mathbf{A}^T \mathbf{A})} \hat{\boldsymbol{\beta}} = \underset{p \times 1}{\mathbf{A}^T \mathbf{Y}}$$

If $(\mathbf{A}^T \mathbf{A})$ is invertible,

$$\hat{\boldsymbol{\beta}} = (\mathbf{A}^T \mathbf{A})^{-1} \mathbf{A}^T \mathbf{Y} \qquad \hat{f}_n^L(X) = X \hat{\boldsymbol{\beta}}$$

When is $(\mathbf{A}^T \mathbf{A})$ invertible ?

Recall: **Full rank matrices are invertible.**

What if $(\mathbf{A}^T \mathbf{A})$ is not invertible ?

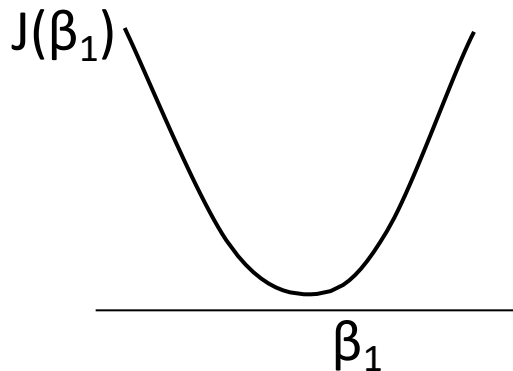
Gradient Descent

Even when $(\mathbf{A}^T \mathbf{A})$ is invertible, might be computationally expensive if \mathbf{A} is huge.

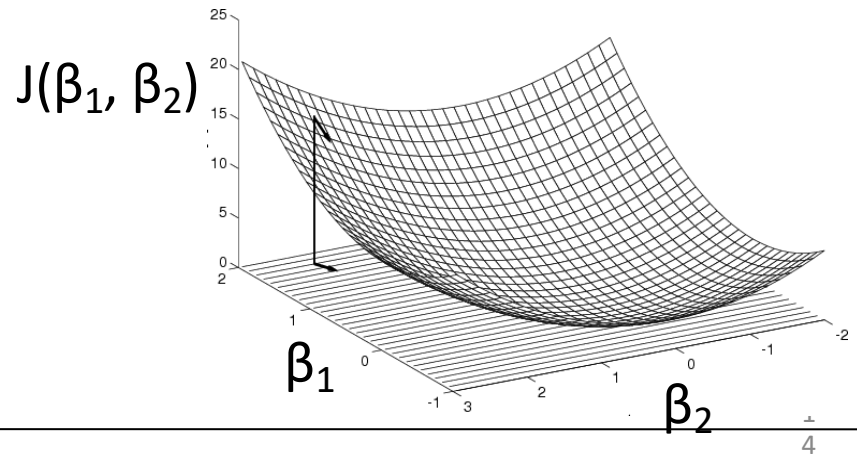
$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) = \arg \min_{\beta} J(\beta)$$

Treat as optimization problem

Observation: $J(\beta)$ is convex in β .



How to find the minimizer?



Gradient Descent

Even when $(\mathbf{A}^T \mathbf{A})$ is invertible, might be computationally expensive if \mathbf{A} is huge.

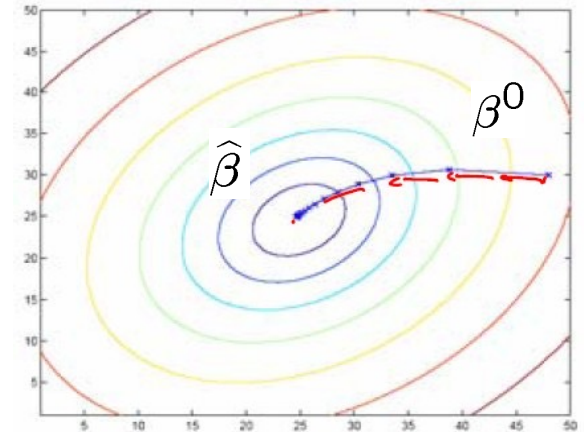
$$\hat{\beta} = \arg \min_{\beta} \frac{1}{n} (\mathbf{A}\beta - \mathbf{Y})^T (\mathbf{A}\beta - \mathbf{Y}) = \arg \min_{\beta} J(\beta)$$

Since $J(\beta)$ is convex, move along negative of gradient

Initialize: β^0

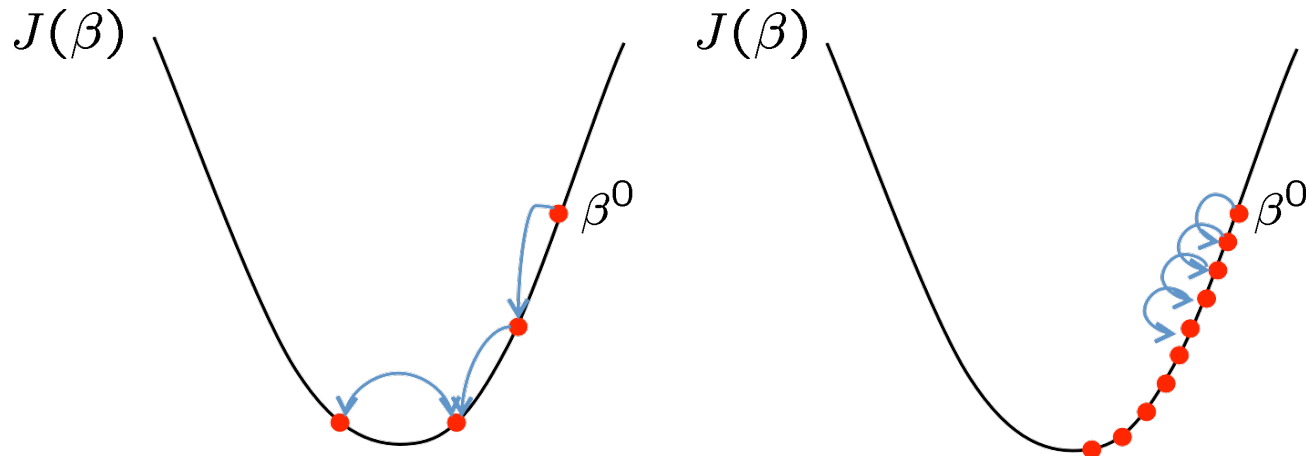
step size

Update:
$$\beta^{t+1} = \beta^t - \frac{\alpha}{2} \frac{\partial J(\beta)}{\partial \beta} \bigg|_t$$
$$= \beta^t - \alpha \underbrace{\mathbf{A}^T (\mathbf{A}\beta^t - \mathbf{Y})}_{0 \text{ if } \hat{\beta} = \beta^t}$$



Stop: when some criterion met e.g. fixed # iterations, or $\left. \frac{\partial J(\beta)}{\partial \beta} \right|_{\beta^t} < \epsilon$.

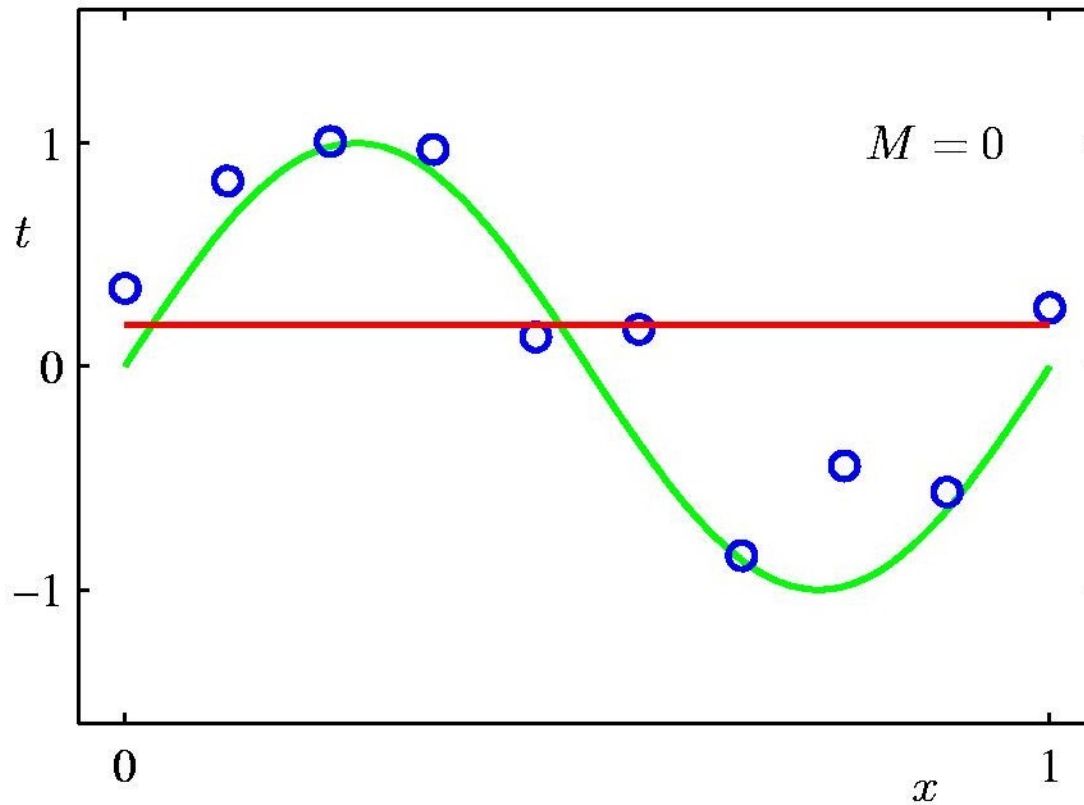
Effect of step-size α



Large $\alpha \Rightarrow$ Fast convergence but larger residual error
Also possible oscillations

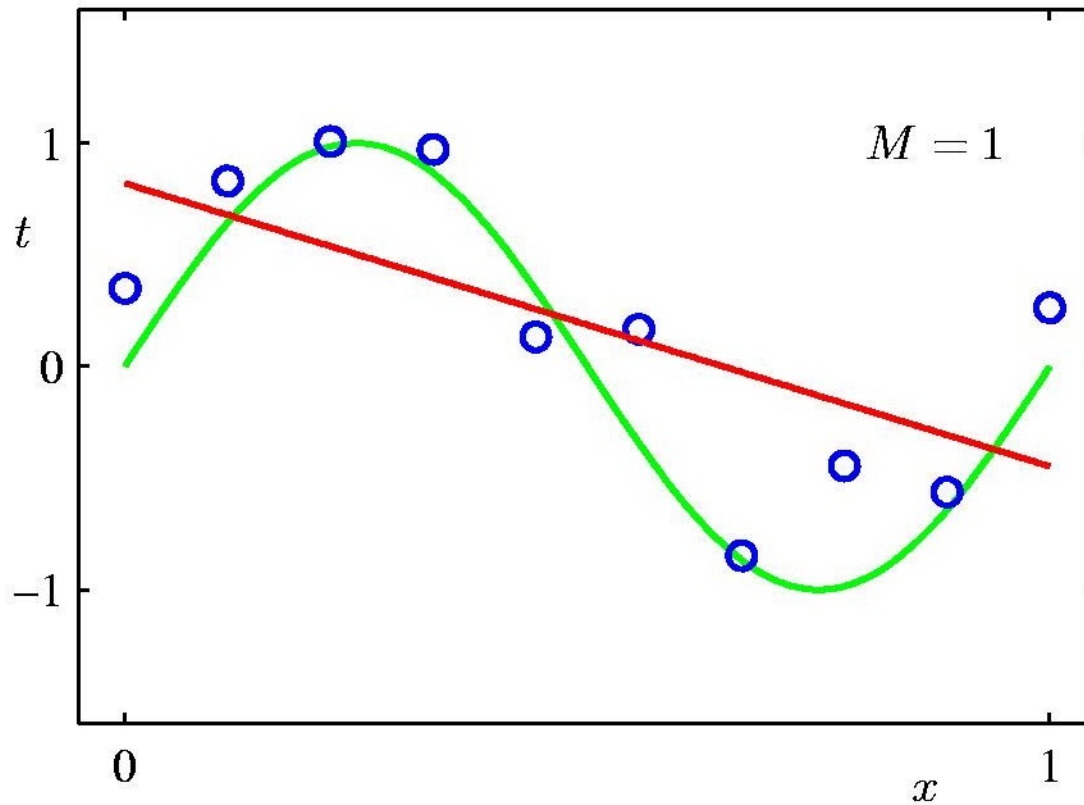
Small $\alpha \Rightarrow$ Slow convergence but small residual error

0th Order Polynomial

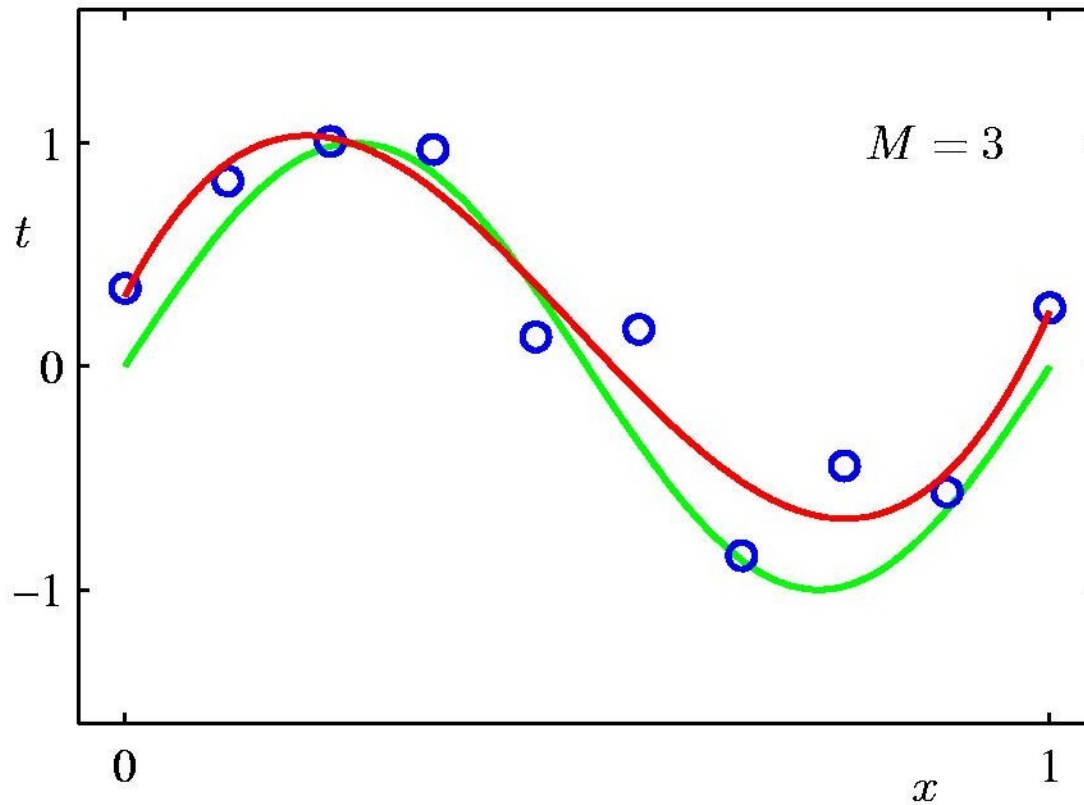


$n=10$

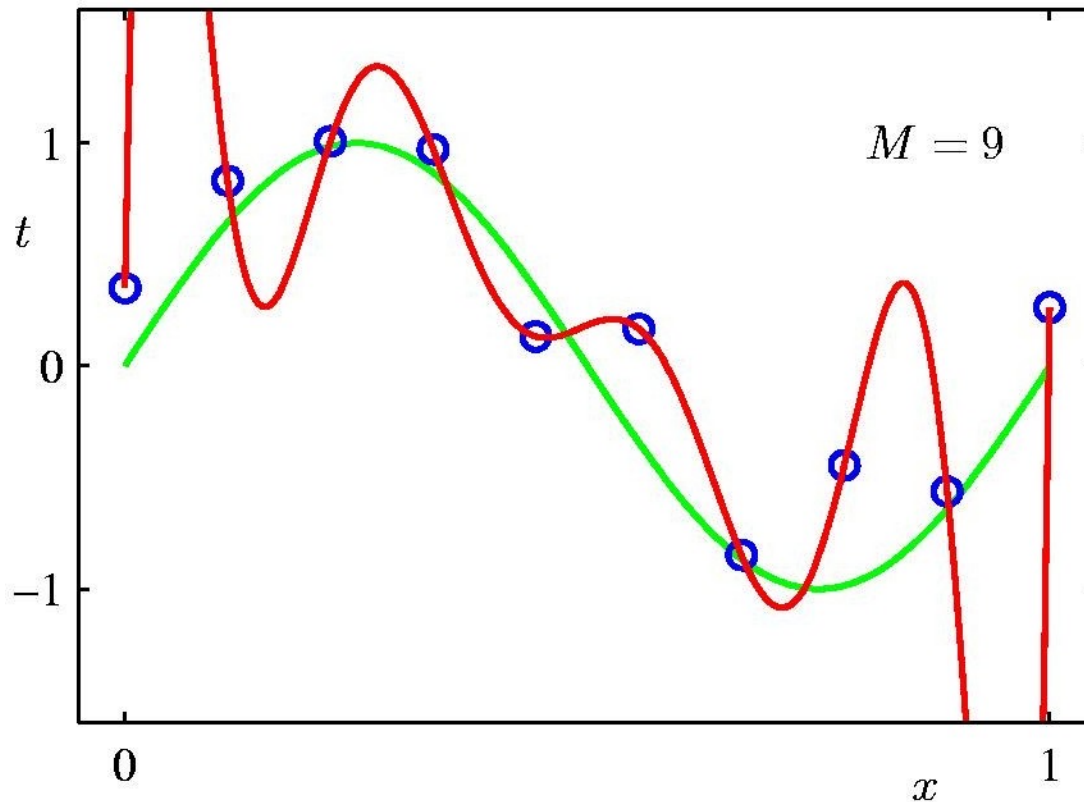
1st Order Polynomial



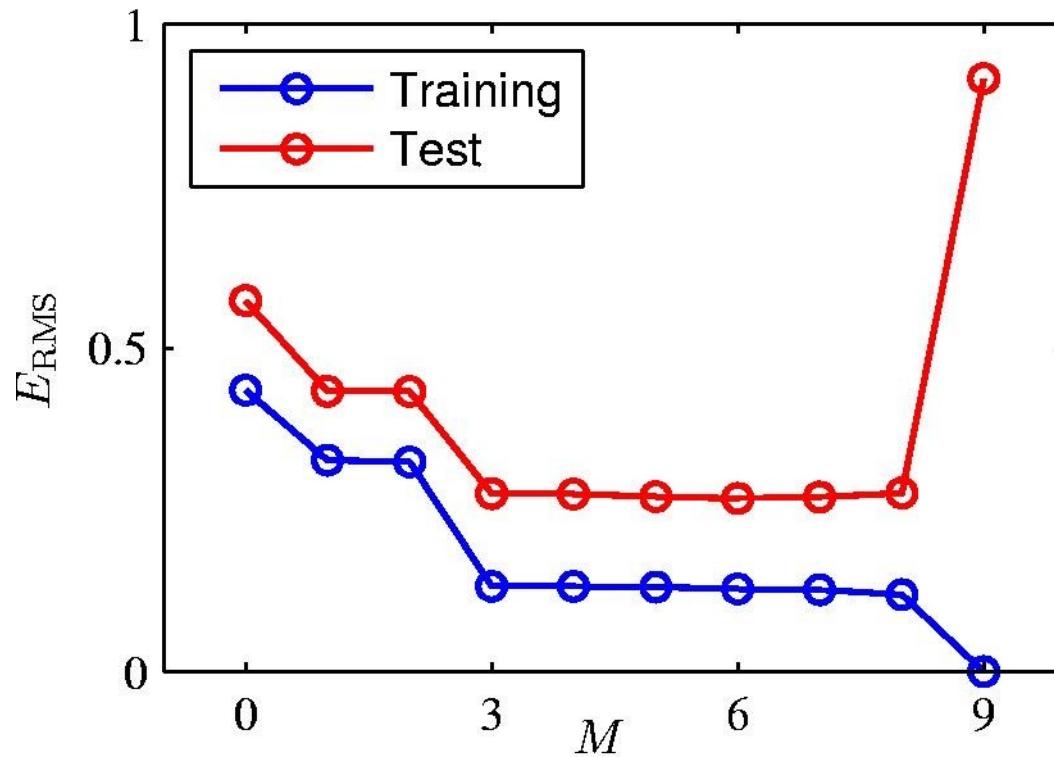
3rd Order Polynomial



9th Order Polynomial



Over-fitting



Root-Mean-Square (RMS) Error

Polynomial Coefficients

	$M = 0$	$M = 1$	$M = 3$	$M = 9$
w_0^*	0.19	0.82	0.31	0.35
w_1^*		-1.27	7.99	232.37
w_2^*			-25.43	-5321.83
w_3^*			17.37	48568.31
w_4^*				-231639.30
w_5^*				640042.26
w_6^*				-1061800.52
w_7^*				1042400.18
w_8^*				-557682.99
w_9^*				125201.43

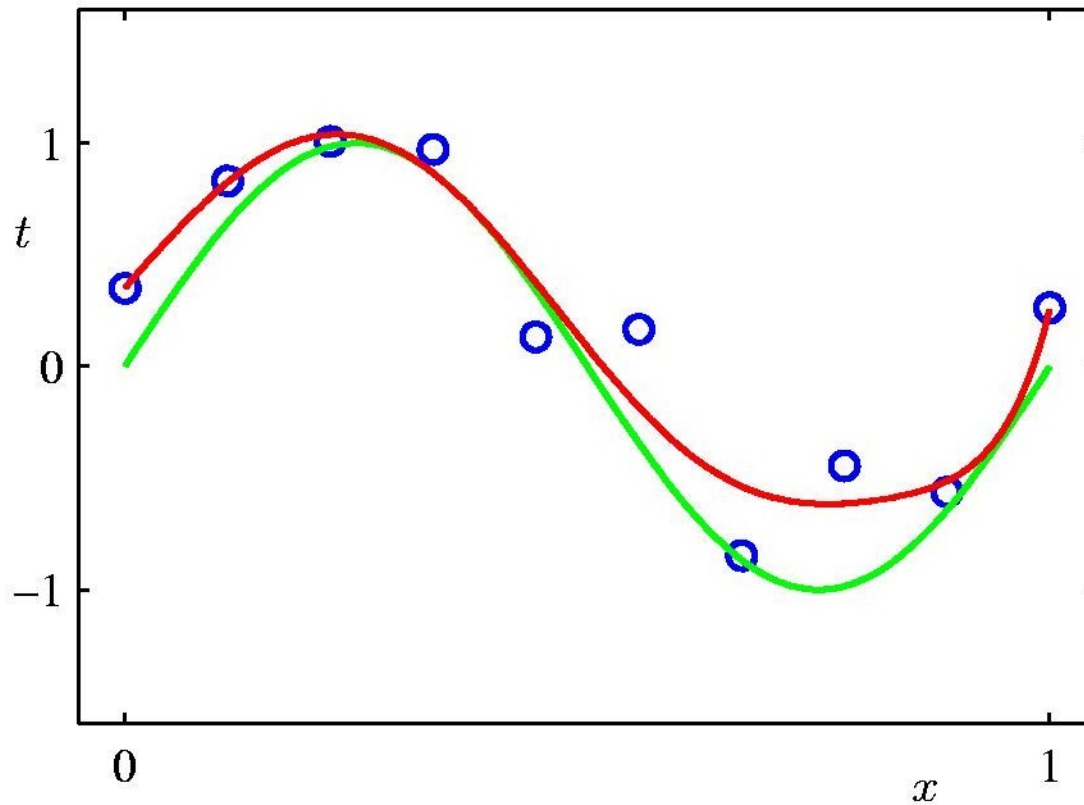
Regularization

Penalize large coefficient values

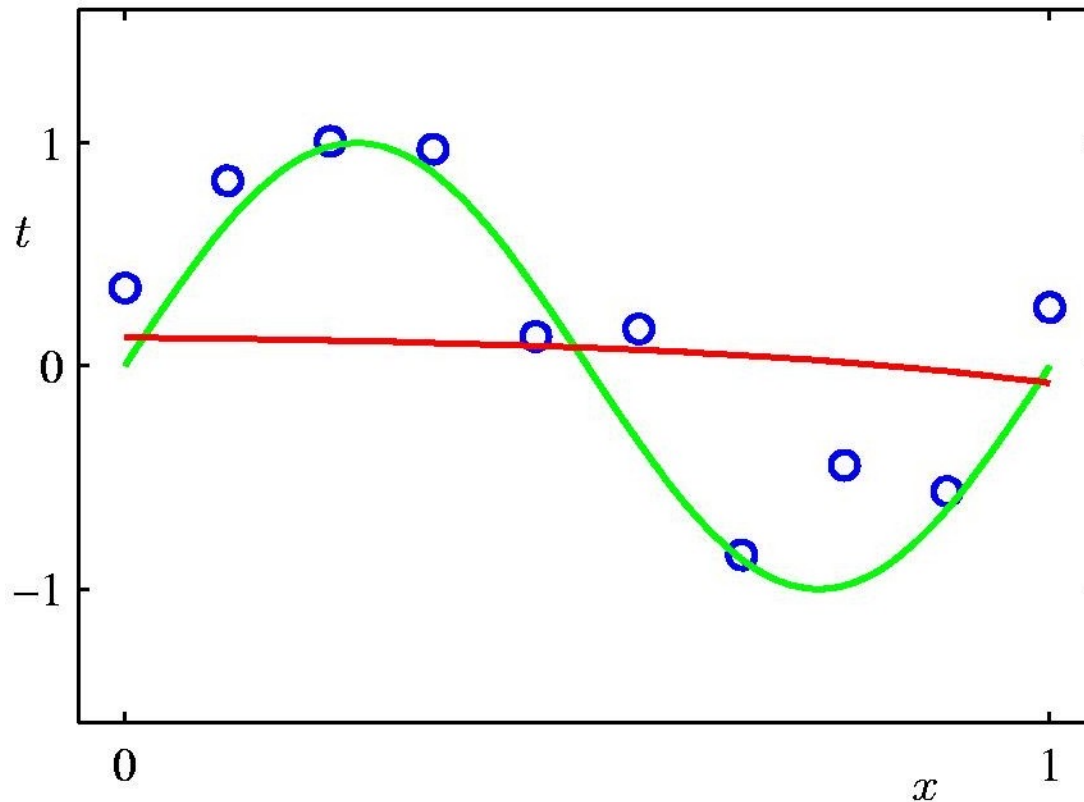
$$J_{\mathbf{x},\mathbf{y}}(\mathbf{w}) = \frac{1}{2} \sum_i \left(y^i - \sum_j w_j \phi_j(\mathbf{x}^i) \right)^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2$$

Regularization:

$$\ln \lambda = -18$$

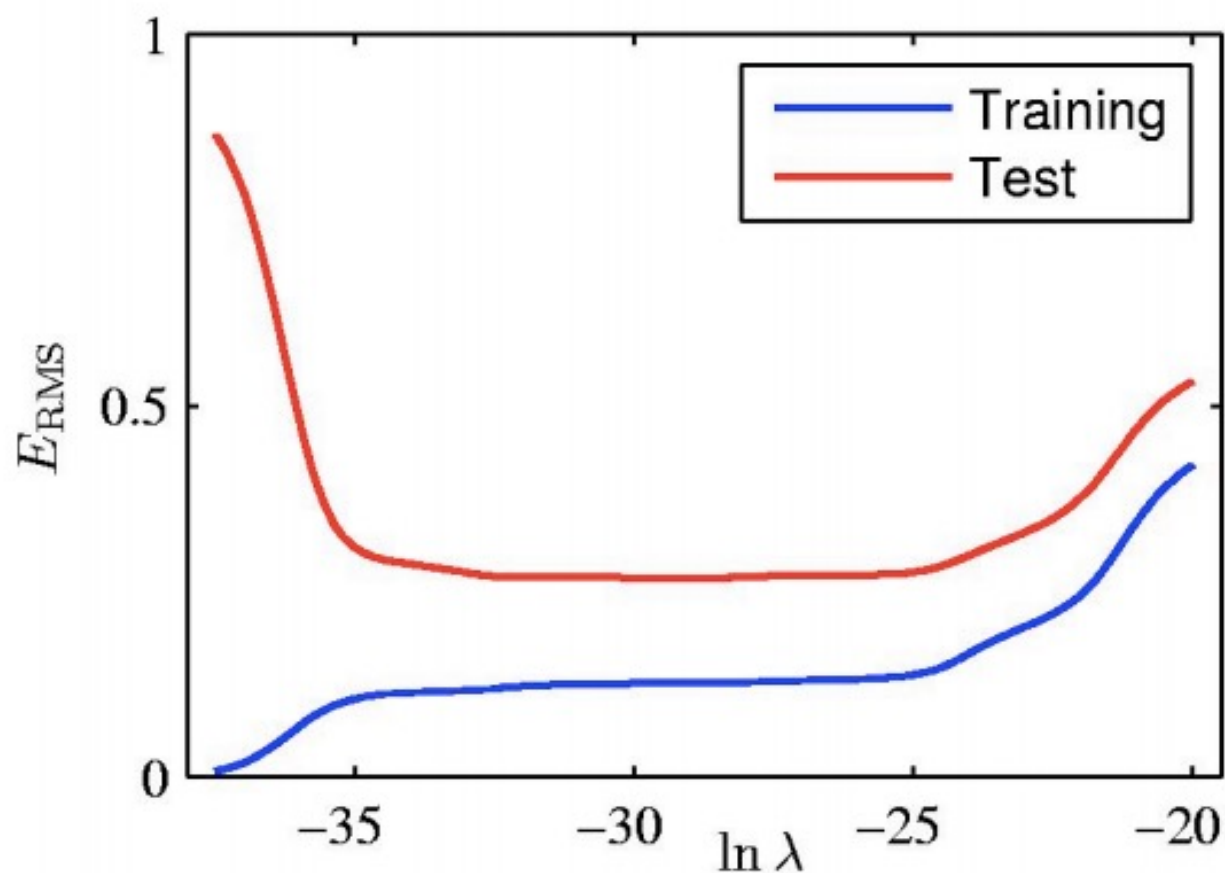


Over Regularization



Regularization

9th Order Polynomial



Regularized Least Squares (1)

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

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

which is minimized by

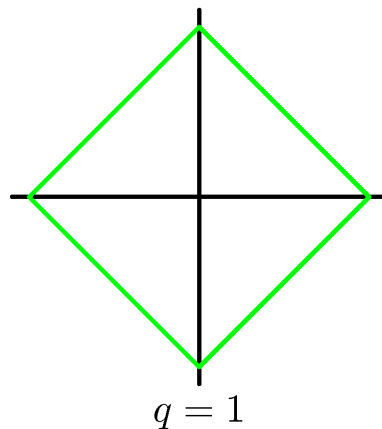
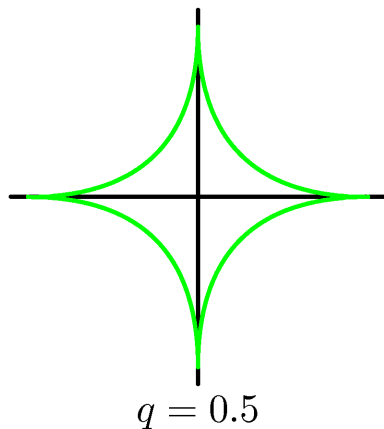
$$\mathbf{w} = \left(\lambda \mathbf{I} + \Phi^T \Phi \right)^{-1} \Phi^T \mathbf{t}.$$

λ is called the regularization coefficient.

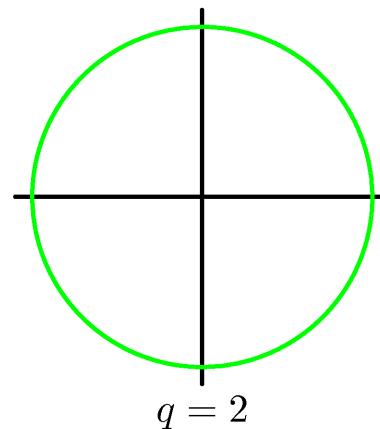
Regularized Least Squares (2)

With a more general regularizer, we have

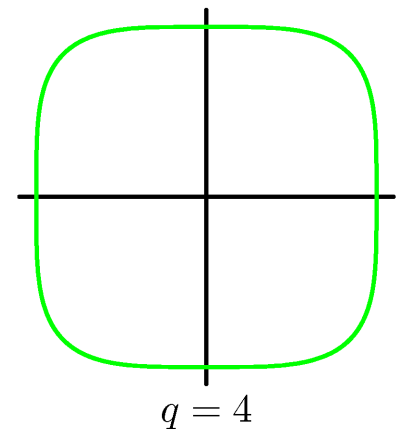
$$\frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \sum_{j=1}^M |w_j|^q$$



Lasso

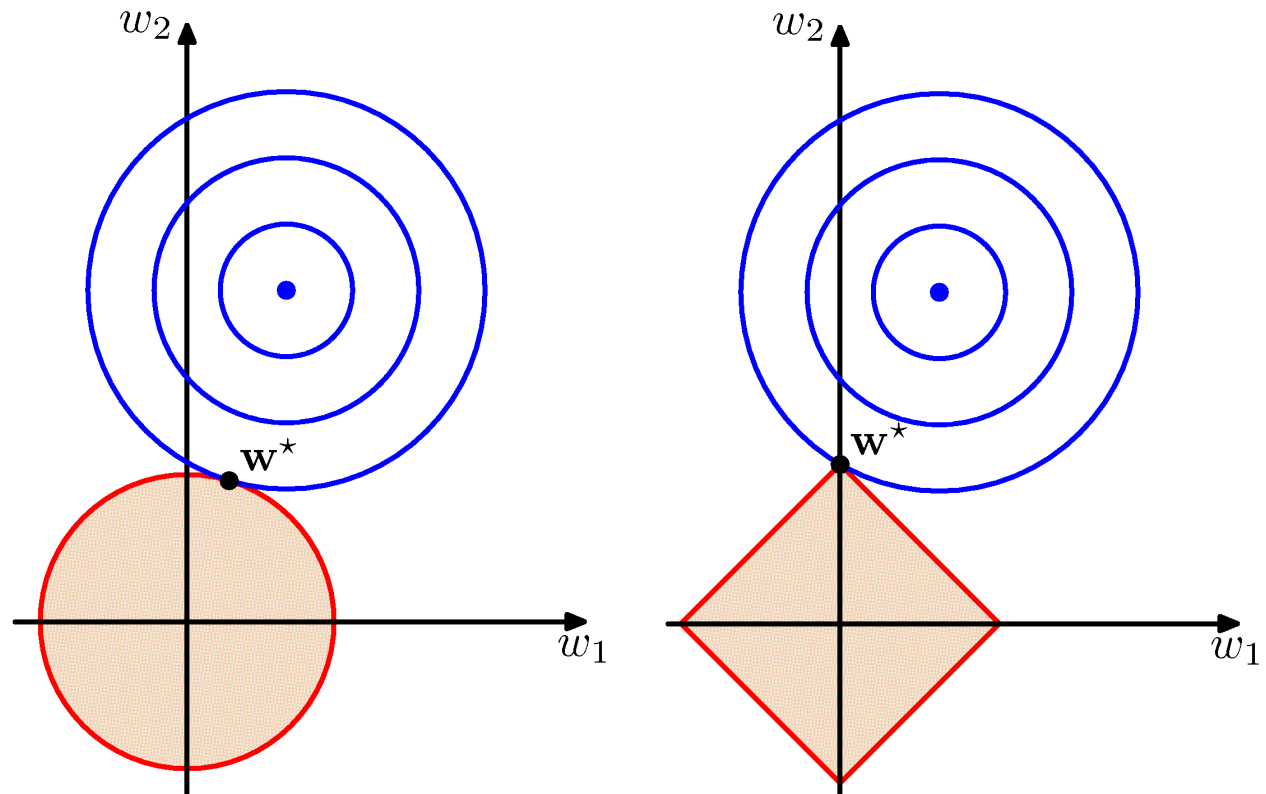


Quadratic



Regularized Least Squares (3)

Lasso tends to generate sparser solutions than a quadratic regularizer.



CLASSIFICATION

Discrete and Continuous Labels

Classification

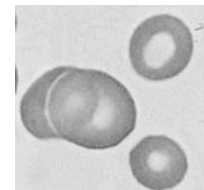
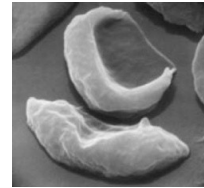


X = Document



Sports
Science
News

Y = Topic



X = Cell Image



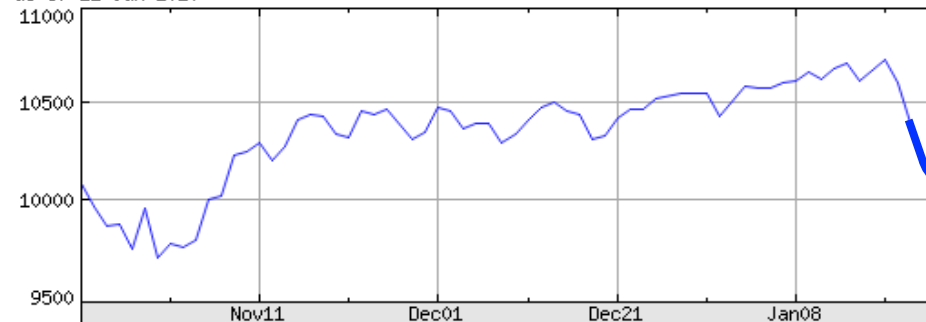
Anemic cell
Healthy cell

Y = Diagnosis

Regression

Stock Market
Prediction

DJ INDU AVERAGE (DOW JONES & CO)
as of 22-Jan-2010



X = Feb01

Y = ?

Copyright 2010 Yahoo! Inc.

Prev Cls: ----

<http://finance.yahoo.com/>

An example application

An emergency room in a hospital measures 17 variables (e.g., blood pressure, age, etc) of newly admitted patients.

A decision is needed: whether to put a new patient in an intensive-care unit.

Due to the high cost of ICU, those patients who may survive less than a month are given higher priority.

Problem: to predict **high-risk patients** and discriminate them from **low-risk patients**.

Another application

A credit card company receives thousands of applications for new cards. Each application contains information about an applicant,

age

Marital status

annual salary

outstanding debts

credit rating

etc.

Problem: to decide whether an application should approved, or to classify applications into two categories, **approved** and **not approved**.

The data and the goal

Data: A set of data records (also called examples, instances or cases) described by

k attributes: $A_1, A_2, \dots A_k$.

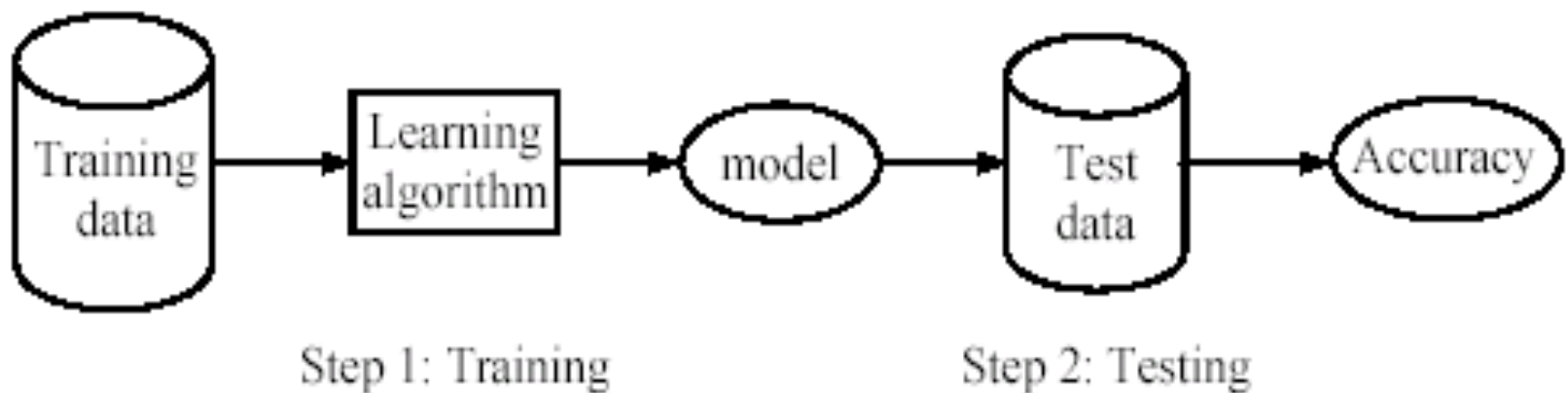
a class: Each example is labelled with a pre-defined class.

Goal: To learn a **classification model** from the data that can be used to predict the classes of new (future, or test) cases/instances.

Supervised learning process: two steps

- **Learning (training):** Learn a model using the training data
- **Testing:** Test the model using **unseen test data** to assess the model accuracy

$$Accuracy = \frac{\text{Number of correct classifications}}{\text{Total number of test cases}},$$



Least squares classification

Binary classification.

Each class is described by it's own linear model:

$$y(x) = w^T x + w_0$$

Compactly written as:

$$y(\mathbf{x}) = \mathbf{W}^T \mathbf{x}$$

\mathbf{W} is $[w \ w_0]$.

$$E_D(\mathbf{W}) = 1/2 (\mathbf{XW} - \mathbf{t})^T (\mathbf{XW} - \mathbf{t})$$

n^{th} row of \mathbf{X} is x_n , the n^{th} datapoint.

\mathbf{t} is vector of +1, -1.

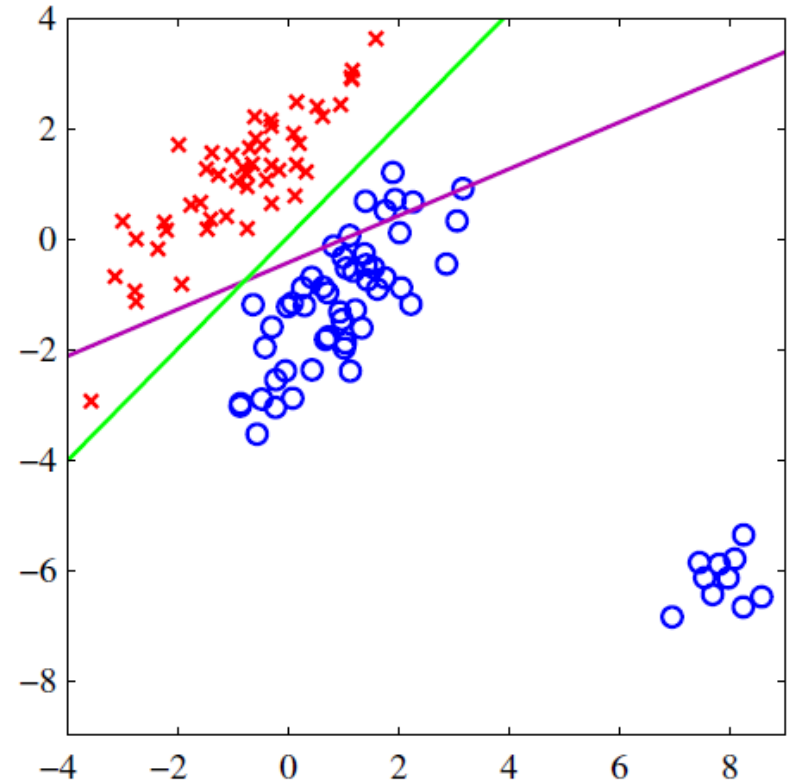
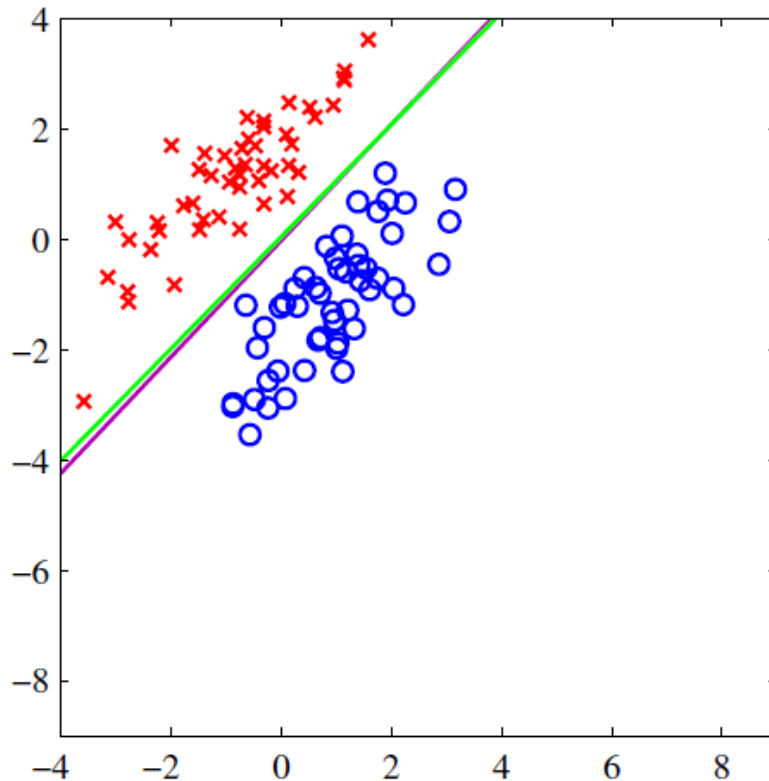
Least squares classification

Least squares \mathbf{W} is:

$$\mathbf{W} = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T \mathbf{t}$$

Problem is affected by outliers.

Least squares classification



From Linear to Logistic Regression

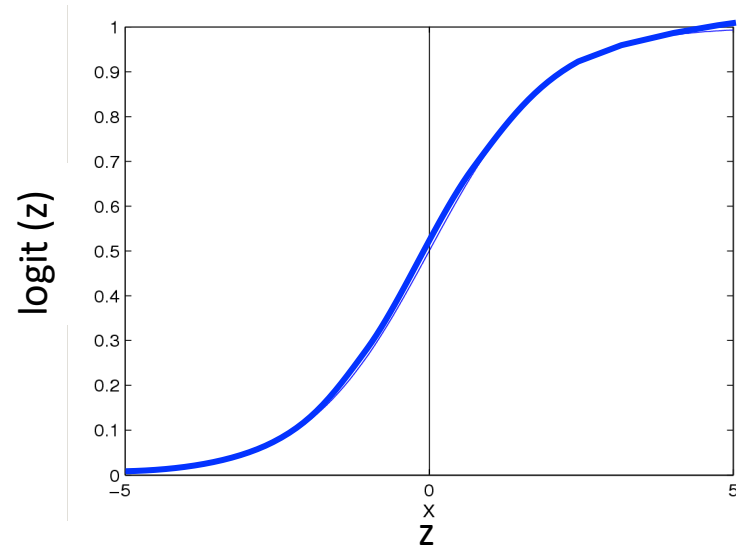
Assumes the following functional form for $P(Y|X)$:

$$P(Y = 1|X) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

Logistic function applied to a linear function of the data

**Logistic function
(or Sigmoid):**

$$\frac{1}{1 + \exp(-z)}$$



Features can be discrete or continuous!

Logistic Regression is a Linear Classifier!

Assumes the following functional form for $P(Y|X)$:

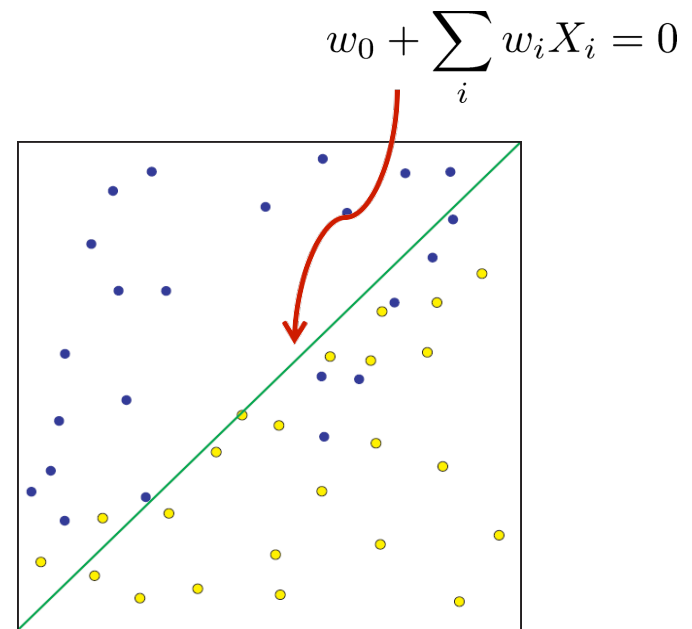
$$P(Y = 1|X) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

Decision boundary:

$$P(Y = 0|X) \stackrel{0}{\geq} P(Y = 1|X) \stackrel{1}{1}$$

$$w_0 + \sum_i w_i X_i \stackrel{0}{\geq} 0 \stackrel{1}{1}$$

(Linear Decision Boundary)



Logistic Regression is a Linear Classifier!

Assumes the following functional form for $P(Y|X)$:

$$P(Y = 1|X) = \frac{1}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$\Rightarrow P(Y = 0|X) = \frac{\exp(w_0 + \sum_i w_i X_i)}{1 + \exp(w_0 + \sum_i w_i X_i)}$$

$$\Rightarrow \frac{P(Y = 0|X)}{P(Y = 1|X)} = \exp(w_0 + \sum_i w_i X_i) \stackrel{0}{\geq} \stackrel{1}{1}$$

$$\Rightarrow w_0 + \sum_i w_i X_i \stackrel{0}{\geq} \stackrel{1}{0}$$

Logistic Regression

Label $t \in \{+1, -1\}$ modeled as:

$$P(t = 1|x, w) = \sigma(w^T x)$$

$$P(y|x, w) = \sigma(yw^T x), y \in \{+1, -1\}$$

Given a set of parameters w , the probability or likelihood of a datapoint (x, t) :

$$P(t|x, w) = \sigma(tw^T x)$$

Logistic Regression

Given a training dataset $\{(x_1, t_1), \dots, (x_N, t_N)\}$,
log likelihood of a model w is given by:

$$L(w) = \sum_n \ln(P(t_n | x_n, w))$$

Using principle of maximum likelihood, the
best w is given by:

$$w^* = \arg \max_w L(w)$$

Logistic Regression

Final Problem:

$$\max_w \sum_{i=1}^n -\log(1 + \exp(-t_n w^T x_n))$$

Or, $\min_w \sum_{i=1}^n \log(1 + \exp(-t_n w^T x_n))$

Error function:

$$E(w) = \sum_{i=1}^n \log(1 + \exp(-t_n w^T x_n))$$

$E(w)$ is convex.

Logistic Regression

Final Problem:

$$\max_w \sum_{i=1}^n -\log(1 + \exp(-t_n w^T x_n))$$

Regularized Version:

$$\max_w \sum_{i=1}^n -\log(1 + \exp(-t_n w^T x_n)) - \lambda w^T w$$

Or,
$$\min_w \sum_{i=1}^n \log(1 + \exp(-t_n w^T x_n)) + \lambda ||w||^2$$

Properties of Error function

Derivatives:

$$\nabla E(w) = \sum_{i=1}^n -(1 - \sigma(t_i w^T x_i))(t_i x_i)$$

$$\nabla E(w) = \sum_{i=1}^n (\sigma(w^T x_i) - t_i) x_i$$

$$\nabla^2 E(w) = \sum_{i=1}^n \sigma(t_i w^T x_i) (1 - \sigma(t_i w^T x_i)) x_i x_i^T$$

Gradient Descent

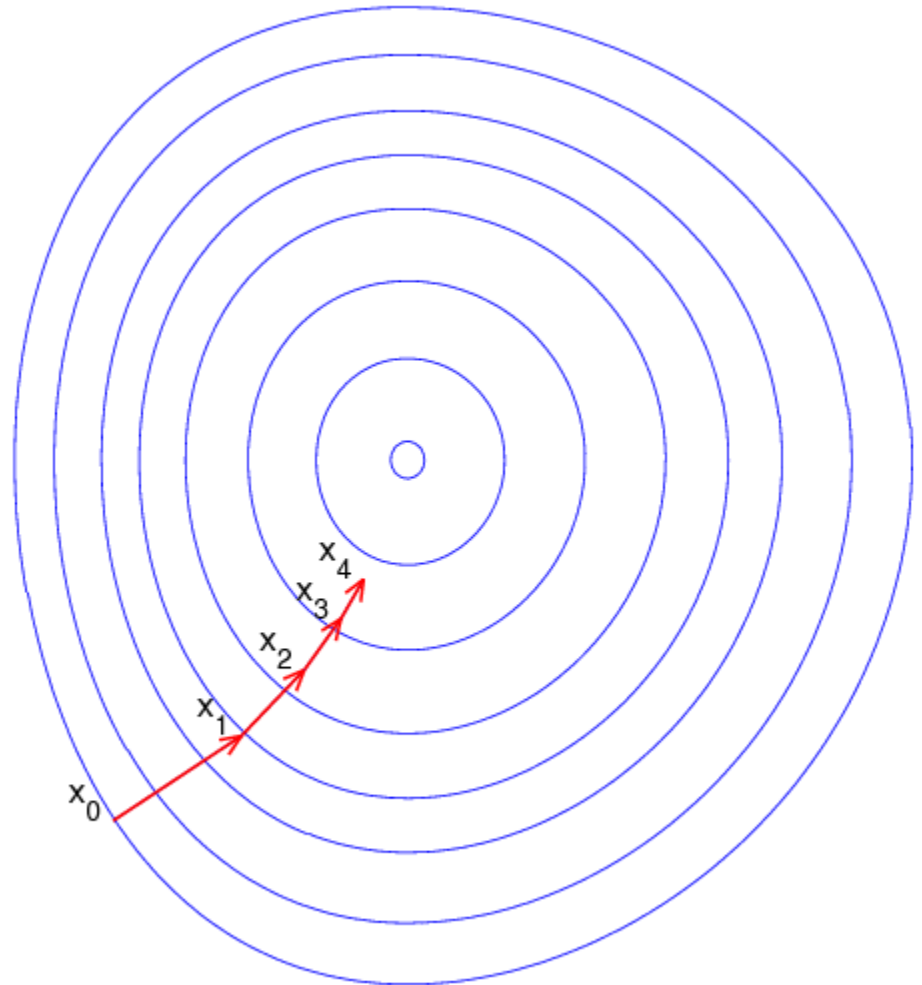
Problem: $\min f(x)$

$f(x)$: differentiable

$g(x)$: gradient of $f(x)$

Negative gradient is
steepest descent
direction.

At each step move in
the gradient direction
so that there is
“sufficient decrease”



Gradient Descent

input : Function f , Gradient ∇f

output: Optimal solution w^*

Initialize $w_0 \leftarrow 0, k \leftarrow 0$

while $|\nabla f_k| > \epsilon$ **do**

 Compute $\alpha_k \leftarrow \text{linesearch}(f, -\nabla f_k, w_k)$

 Set $w_{k+1} \leftarrow w_k - \alpha_k \nabla f_k$

 Evaluate ∇f_{k+1}

$k \leftarrow k + 1$

end

$w^* \leftarrow w_k$

Logistic Regression for more than 2 classes

- Logistic regression in more general case, where $Y \in \{y_1, \dots, y_K\}$

for $k < K$

$$P(Y = y_k | X) = \frac{\exp(w_{k0} + \sum_{i=1}^d w_{ki} X_i)}{1 + \sum_{j=1}^{K-1} \exp(w_{j0} + \sum_{i=1}^d w_{ji} X_i)}$$

for $k=K$ (normalization, so no weights for this class)

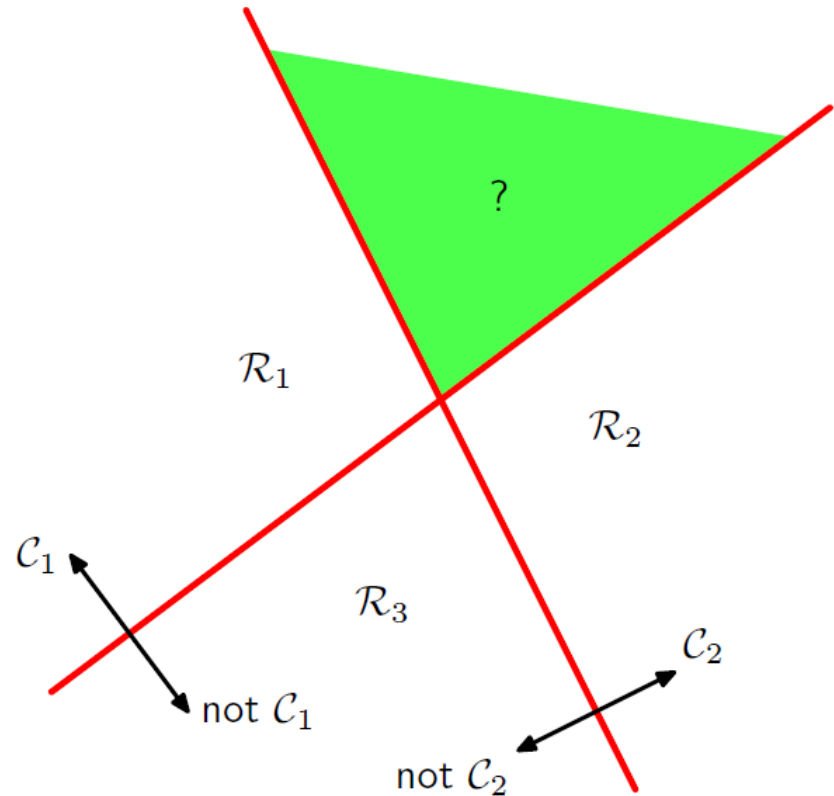
$$P(Y = y_K | X) = \frac{1}{1 + \sum_{j=1}^{K-1} \exp(w_{j0} + \sum_{i=1}^d w_{ji} X_i)}$$

Multiple classes

One-vs-all: $K - 1$ hyperplanes each separating C_1, \dots, C_{K-1} classes from rest.

Otherwise C_K

Low number of classifiers.

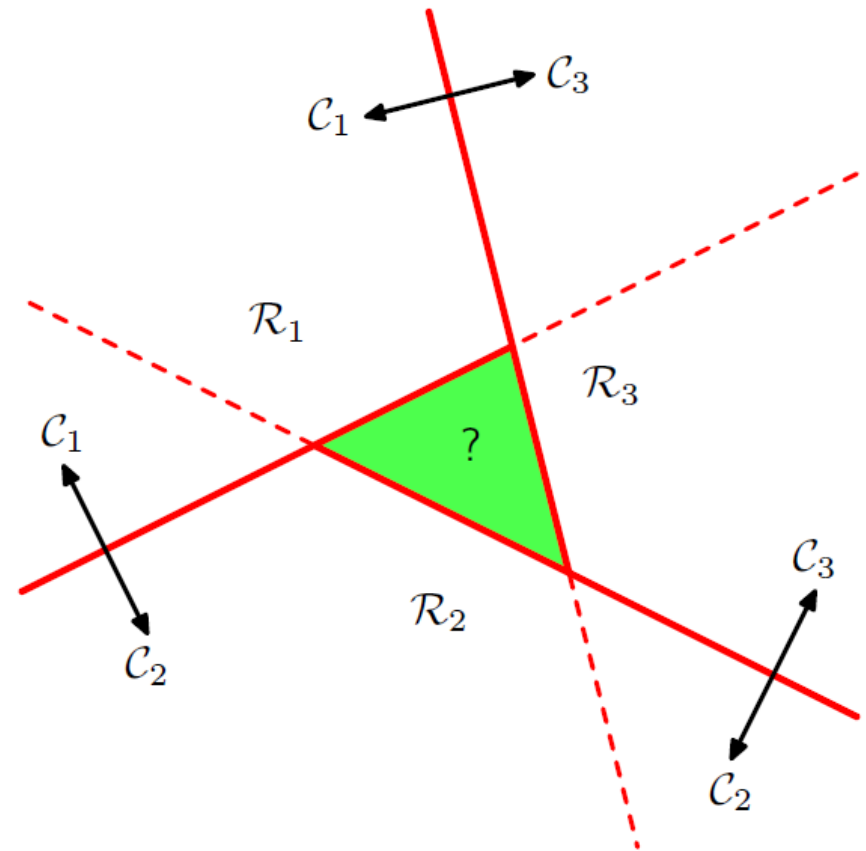


Multiple classes

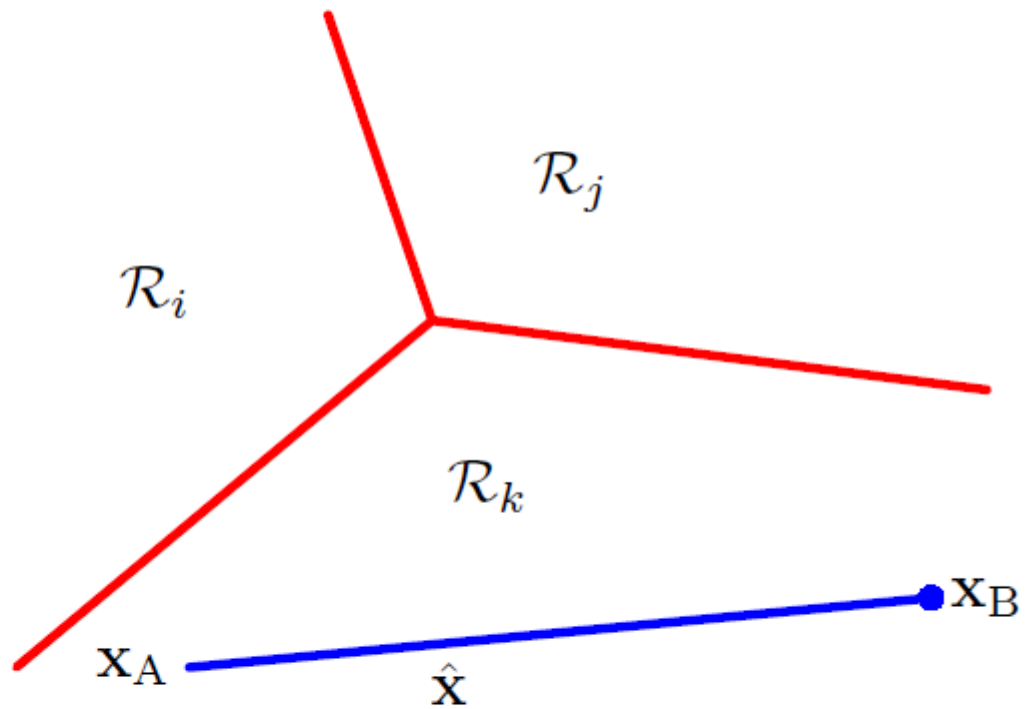
One-vs-one: Every pair $C_i - C_j$ get a boundary.

Final by majority vote.

High number of
classifiers.



Multiple Classes



Multiple classes

K-linear discriminant functions:

$$y_k(x) = w_k^T \mathbf{x} + w_{k0}$$

Assign x to C_k if $y_k(x) \geq y_j(x)$ for all $j \neq k$

Decision boundary:

$$(w_k - w_j)^T \mathbf{x} + (w_{k0} - w_{j0}) = 0$$

Decision region is singly connected:

$$x = \lambda x_A + (1 - \lambda)x_B$$

If x_A and x_B have same label, so does x .

MORE REGRESSION

The Bias-Variance Decomposition (1)

Recall the *expected squared loss*,

$$\mathbb{E}[L] = \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \underbrace{\iint \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) d\mathbf{x} dt}_{\text{noise}}$$

where

$$h(\mathbf{x}) = \mathbb{E}[t|\mathbf{x}] = \int tp(t|\mathbf{x}) dt.$$

The second term of $\mathbb{E}[L]$ corresponds to the noise inherent in the random variable t .

What about the first term?

The Bias-Variance Decomposition (2)

Suppose we were given multiple data sets, each of size N . Any particular data set, \mathcal{D} , will give a particular function $y(\mathbf{x}; \mathcal{D})$. We then have

$$\begin{aligned} & \{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2 \\ &= \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] + \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ &= \{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2 + \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 \\ &\quad + 2\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}. \end{aligned}$$

The Bias-Variance Decomposition (3)

Taking the expectation over \mathcal{D} yields

$$\begin{aligned} \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2] \\ = \underbrace{\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2}_{(\text{bias})^2} + \underbrace{\mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2]}_{\text{variance}}. \end{aligned}$$

The Bias-Variance Decomposition (4)

Thus we can write

$$\text{expected loss} = (\text{bias})^2 + \text{variance} + \text{noise}$$

where

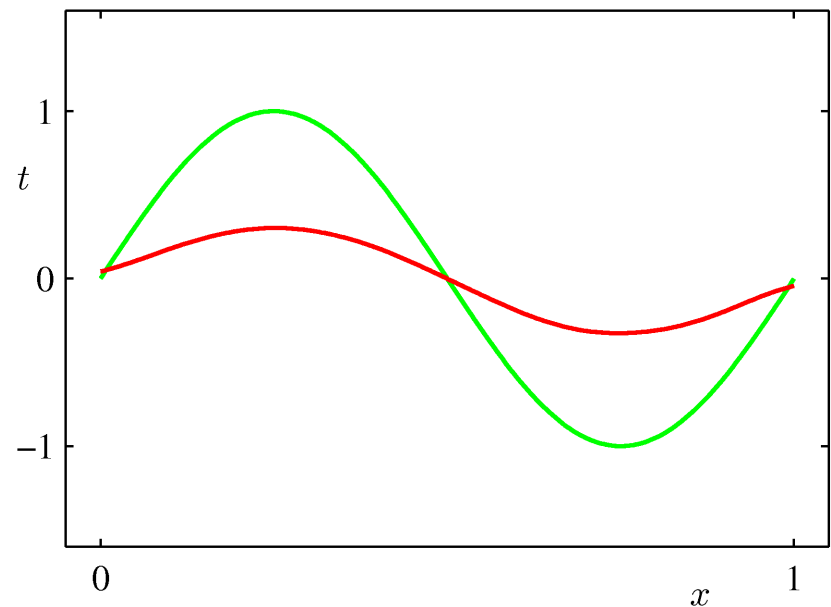
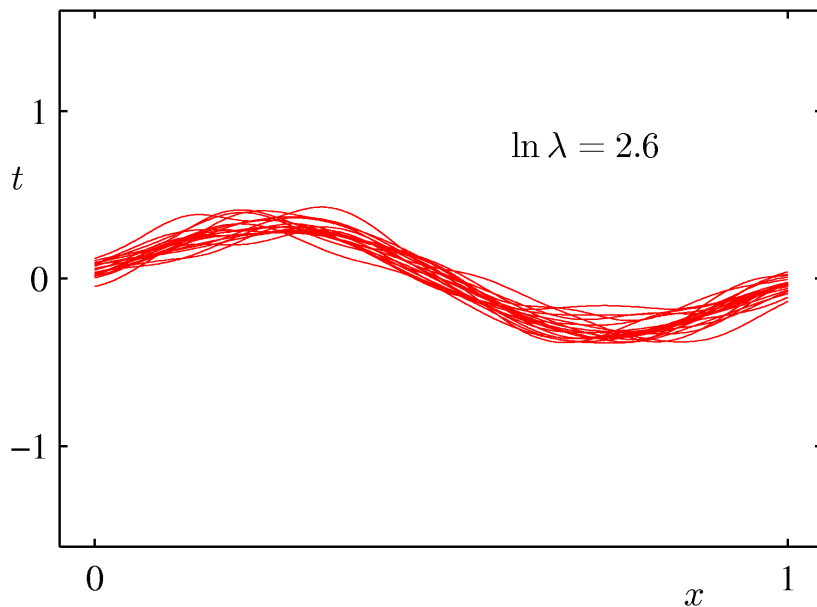
$$(\text{bias})^2 = \int \{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2 p(\mathbf{x}) \, d\mathbf{x}$$

$$\text{variance} = \int \mathbb{E}_{\mathcal{D}} [\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2] p(\mathbf{x}) \, d\mathbf{x}$$

$$\text{noise} = \iint \{h(\mathbf{x}) - t\}^2 p(\mathbf{x}, t) \, d\mathbf{x} \, dt$$

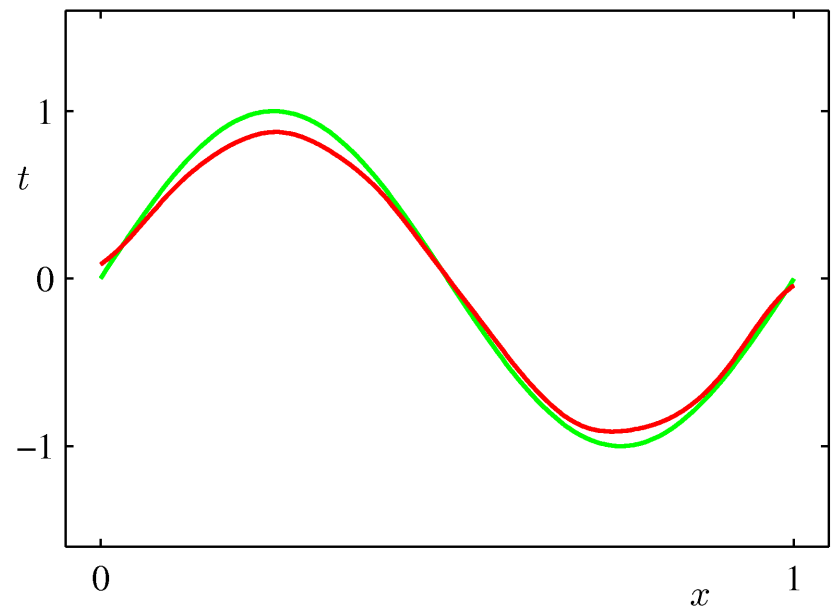
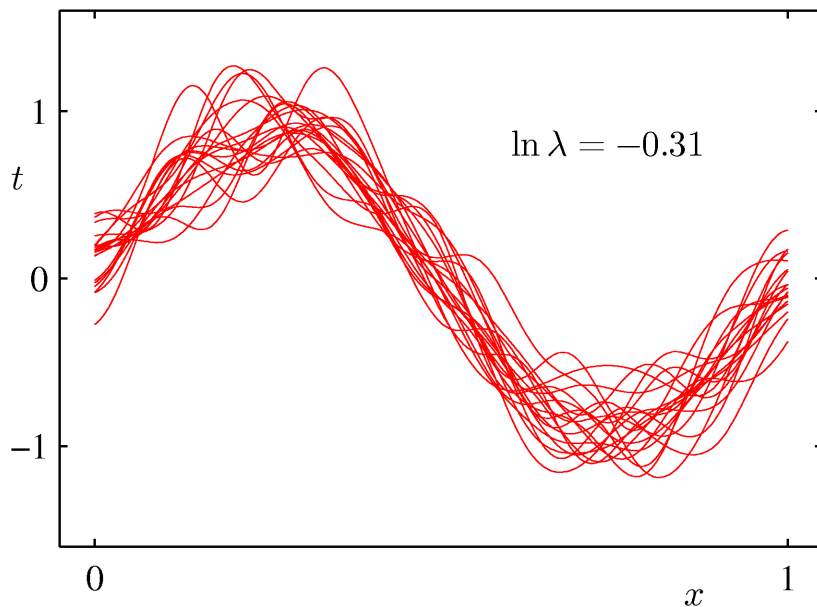
The Bias-Variance Decomposition (5)

Example: 25 data sets from the sinusoidal, varying the degree of regularization, λ .



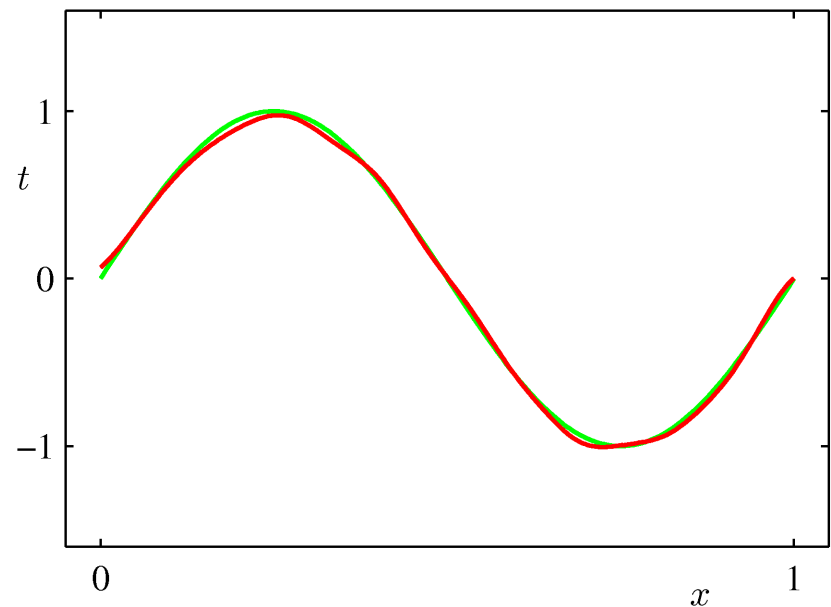
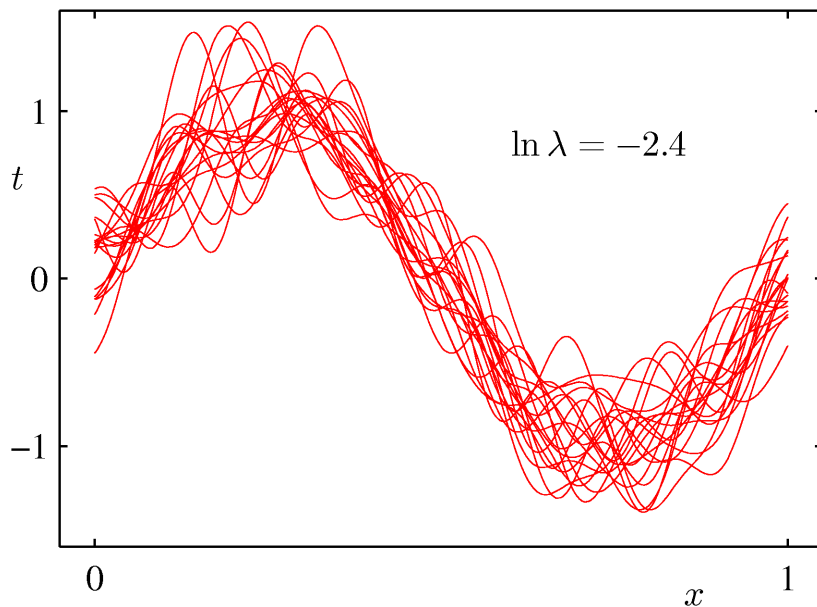
The Bias-Variance Decomposition (6)

Example: 25 data sets from the sinusoidal, varying the degree of regularization, λ .



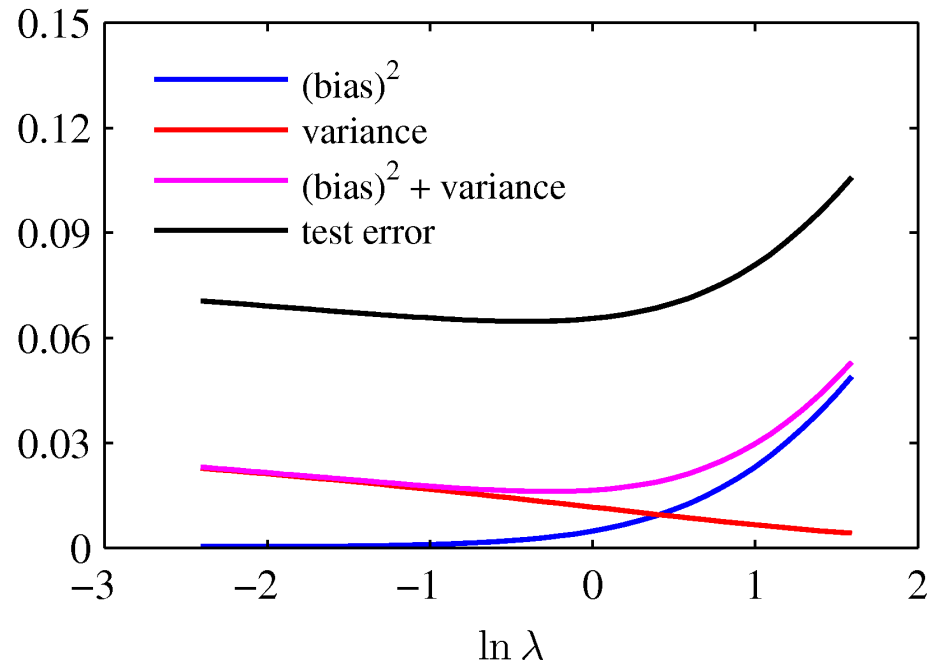
The Bias-Variance Decomposition (7)

Example: 25 data sets from the sinusoidal, varying the degree of regularization, λ .



The Bias-Variance Trade-off

From these plots, we note that an over-regularized model (large λ) will have a high bias, while an under-regularized model (small λ) will have a high variance.



Bayesian Linear Regression (1)

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 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

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

Bayesian Linear Regression (2)

A common choice for the prior is

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

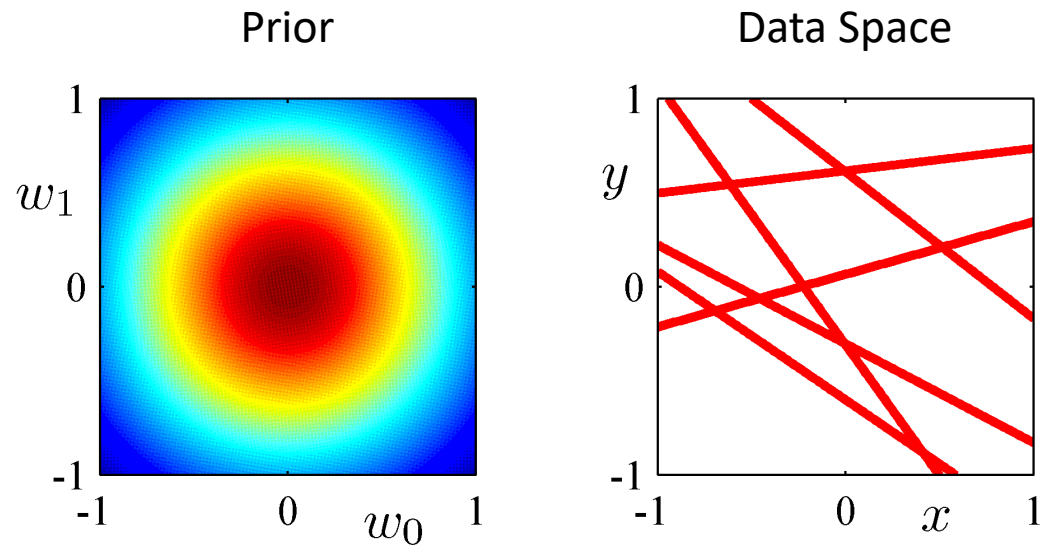
for which

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

Next we consider an example ...

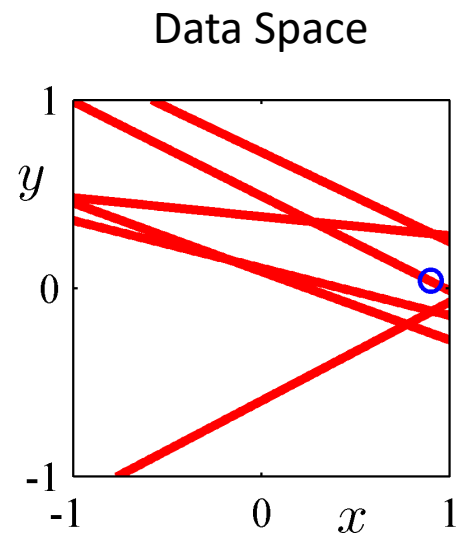
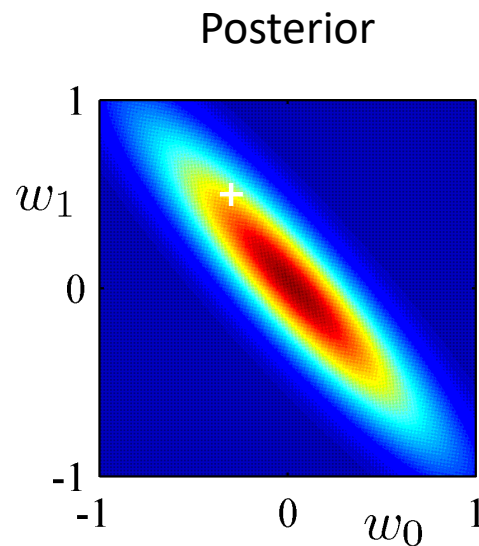
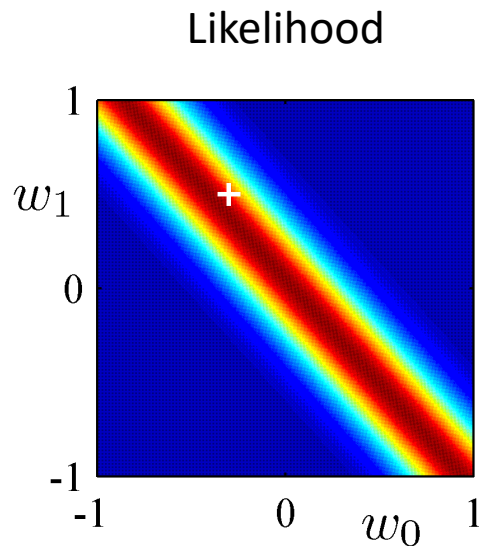
Bayesian Linear Regression (3)

0 data points observed



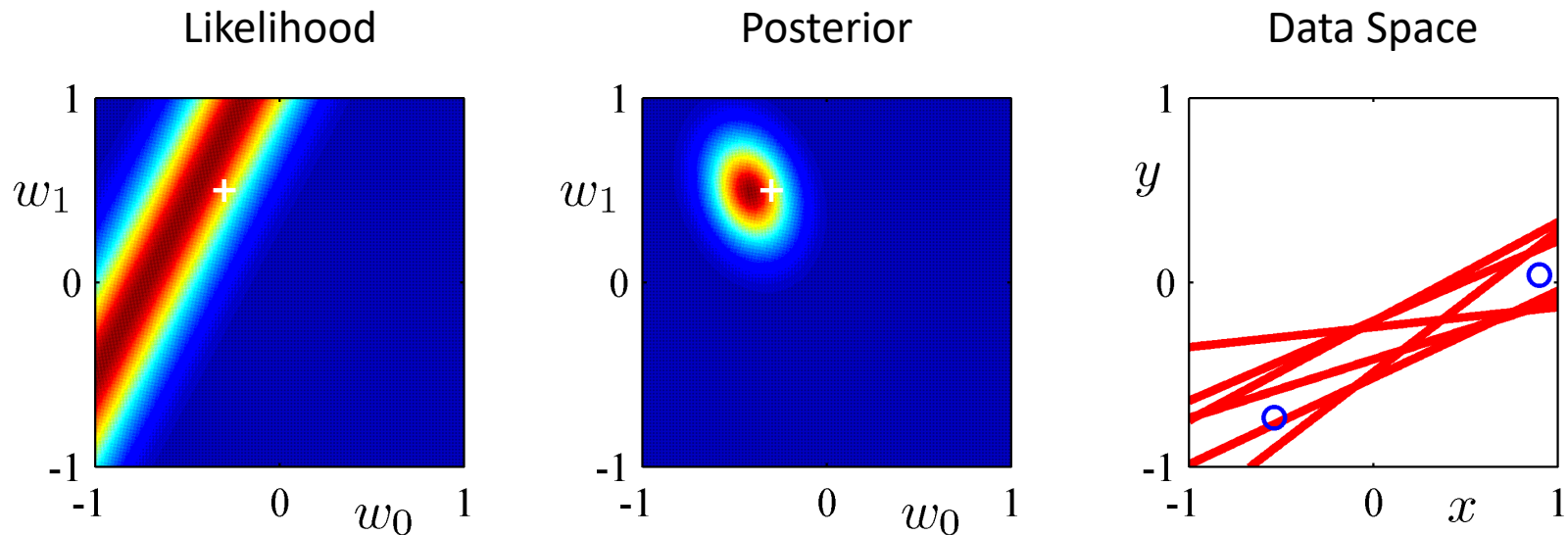
Bayesian Linear Regression (4)

1 data point observed



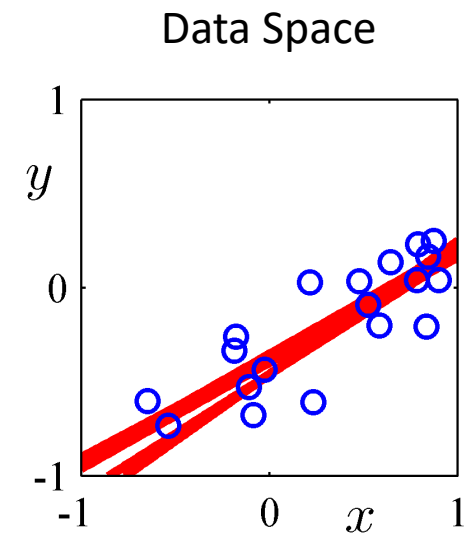
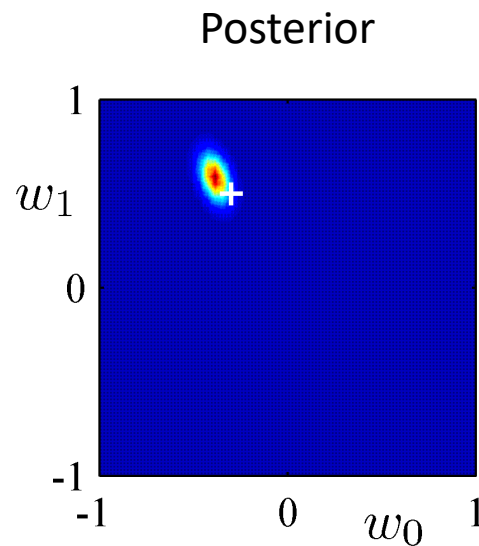
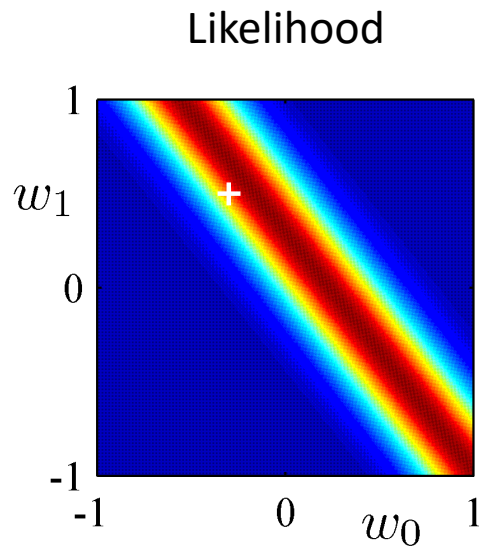
Bayesian Linear Regression (5)

2 data points observed



Bayesian Linear Regression (6)

20 data points observed



Predictive Distribution (1)

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

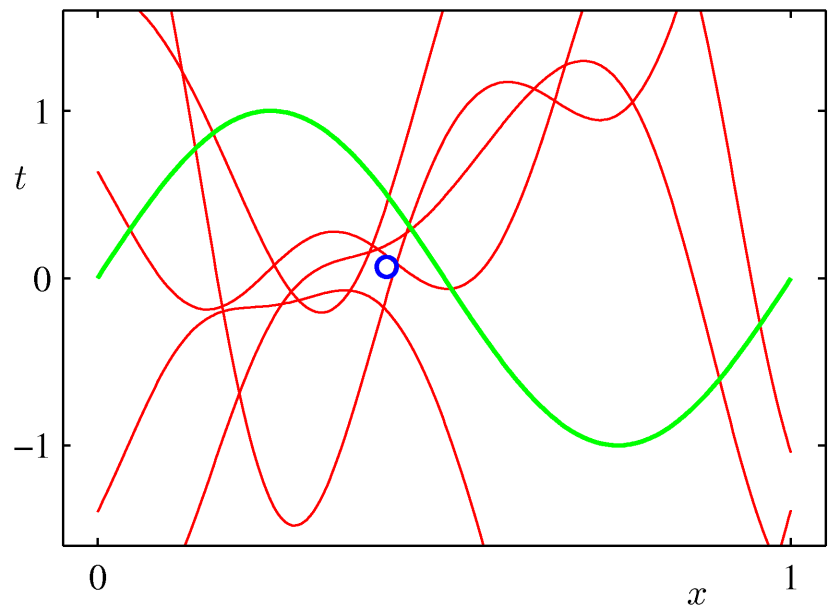
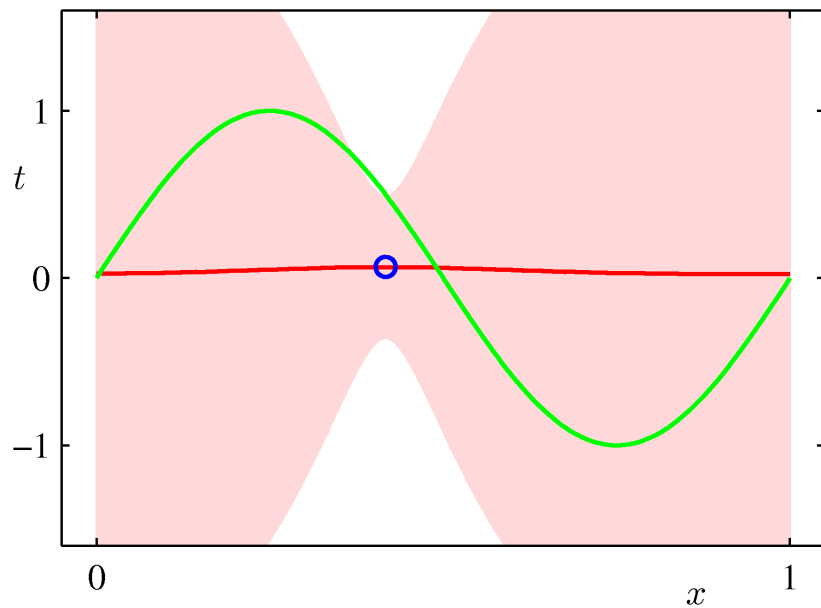
$$\begin{aligned} p(t|\mathbf{t}, \alpha, \beta) &= \int p(t|\mathbf{w}, \beta) p(\mathbf{w}|\mathbf{t}, \alpha, \beta) d\mathbf{w} \\ &= \mathcal{N}(t|\mathbf{m}_N^T \boldsymbol{\phi}(\mathbf{x}), \sigma_N^2(\mathbf{x})) \end{aligned}$$

where

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

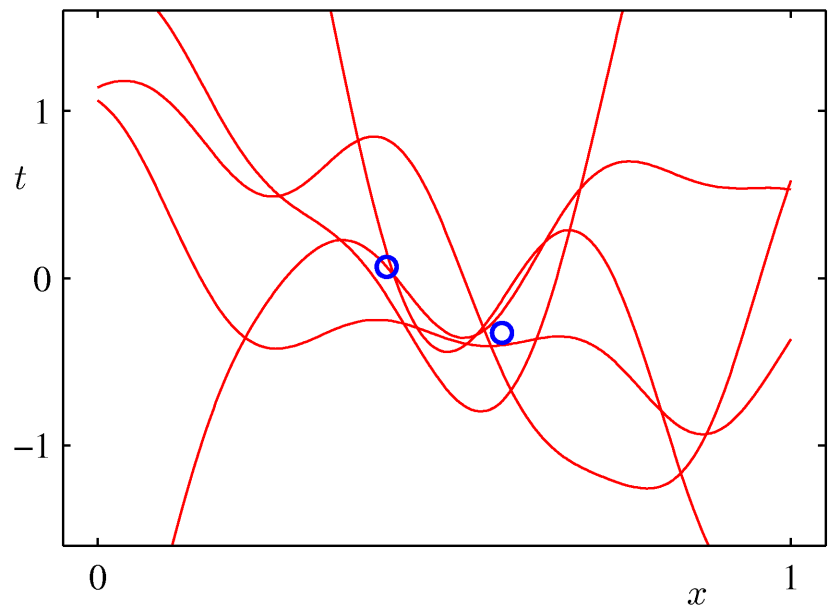
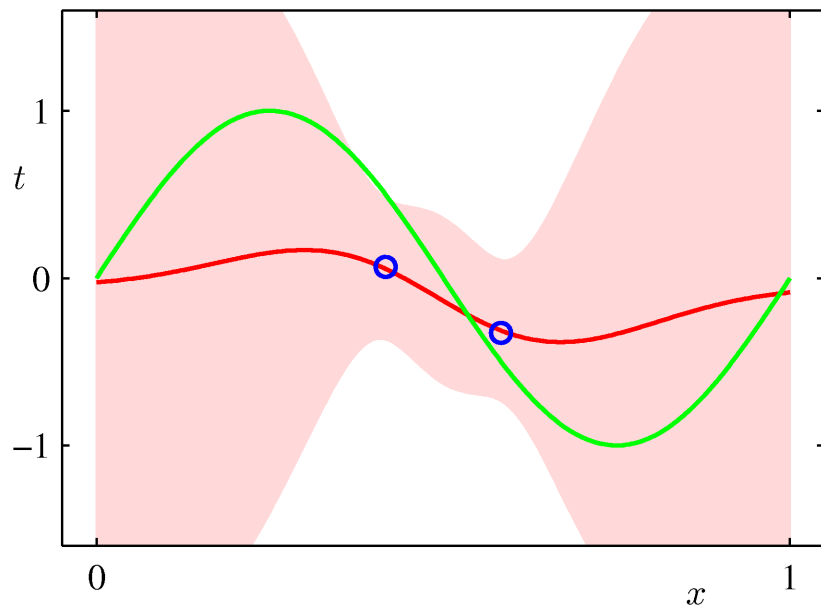
Predictive Distribution (2)

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



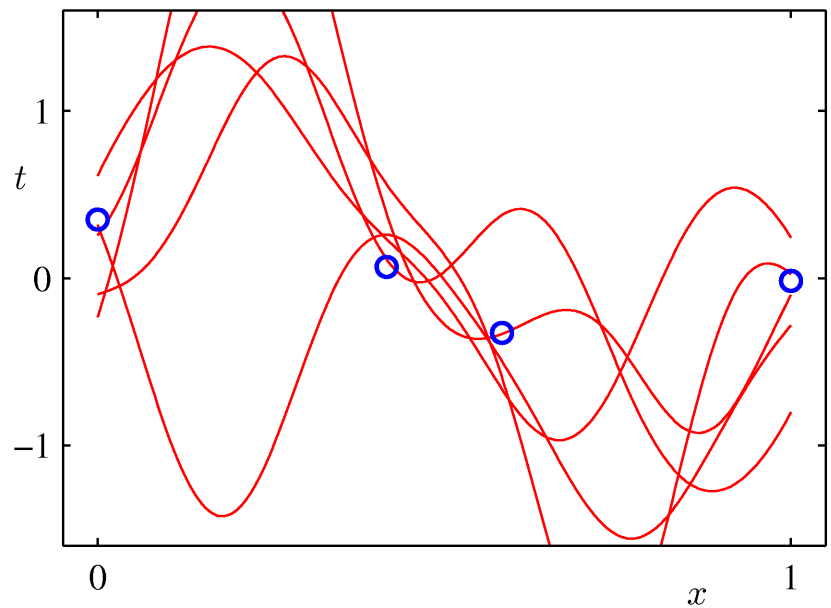
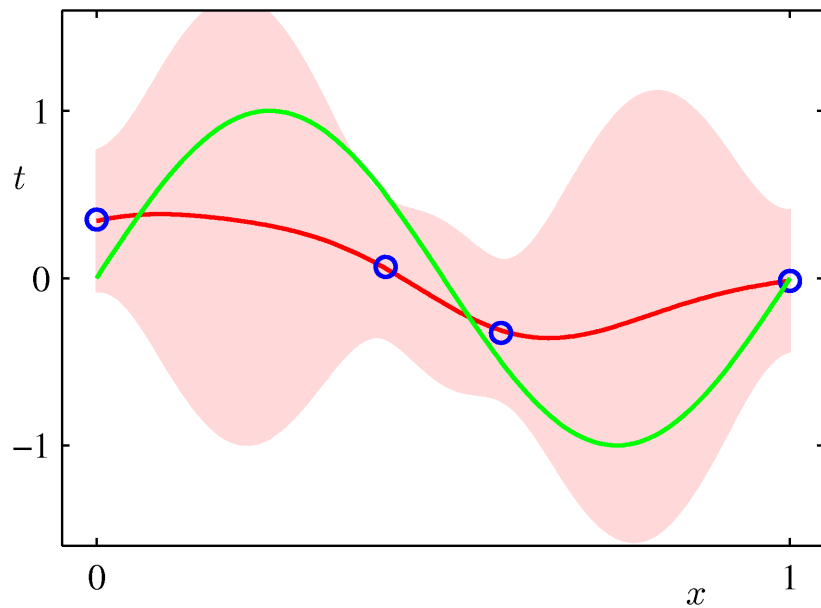
Predictive Distribution (3)

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



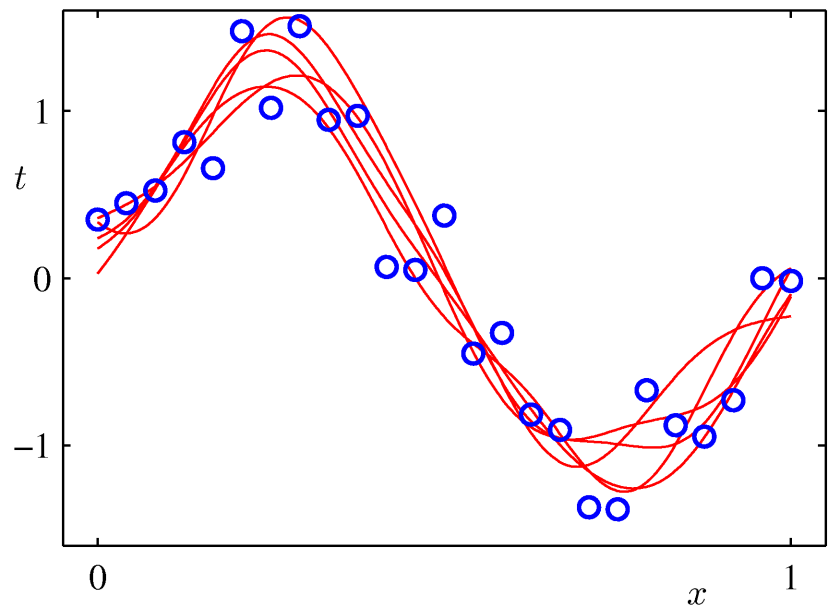
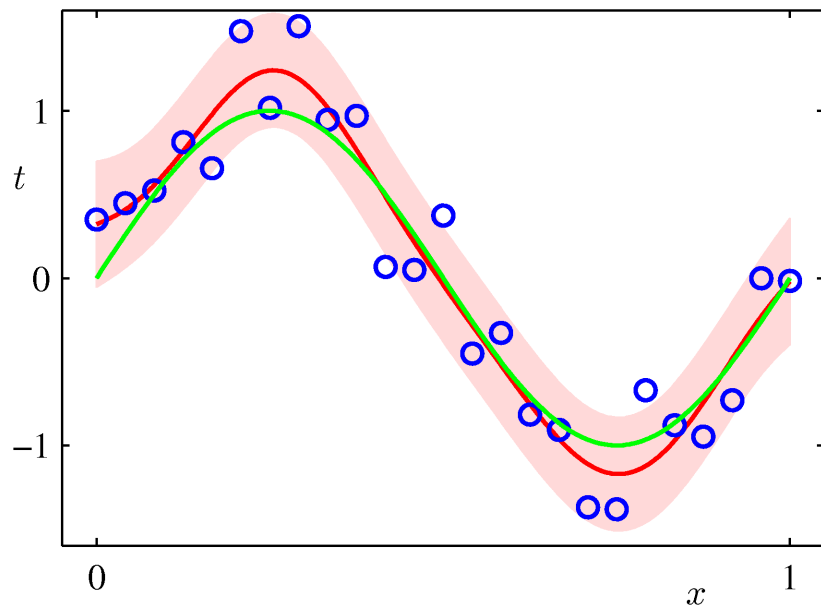
Predictive Distribution (4)

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



Predictive Distribution (5)

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



Multiple Outputs (1)

Analogously to the single output case we have:

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

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

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

Multiple Outputs (2)

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

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

If we consider a single target variable, t_k , we see that

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

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