

Linear Models for Regression

Yan Yan

Shenzhen Institute of Advanced Technology
Chinese Academy of Sciences

yan.yan@siat.ac.cn

September 30, 2016

1 Linear Basis Function Models

- Maximum likelihood and least squares
- Geometry of least squares
- Sequential learning

2 Bias-Variance Decomposition

3 Bayesian Linear Regression

- Parameter distribution
- Predictive distribution
- Equivalent kernel

4 Bayesian Model Comparison

5 Evidence Approximation

- Evaluation of Approximation Function
- Maximizing the evidence function

Linear Basis Function Models

The simplest linear model for regression (often simply known as *linear regression*):

$$y(\mathbf{x}, \mathbf{w}) = \omega_0 + \omega_1 x_1 + \cdots + \omega_D x_D = \sum_{j=0}^{M-1} \omega_j \phi_j(\mathbf{x}) = \boldsymbol{\omega}^T \boldsymbol{\phi}(\mathbf{x}) \quad (1)$$

where $\boldsymbol{\omega} = (\omega_0, \dots, \omega_{M-1})^T$ and $\boldsymbol{\phi} = (\phi_0, \dots, \phi_{M-1})^T$. This kinds of models are called linear models.

When ϕ has different types, which means different kinds of *basis function*, we have other kinds of modes, like:

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

or

$$\phi_j(x) = \sigma\left(\frac{x - \mu_j}{s}\right) \quad \text{and} \quad \sigma(a) = \frac{1}{1 + \exp(-a)} \quad (3)$$

Examples of basis function

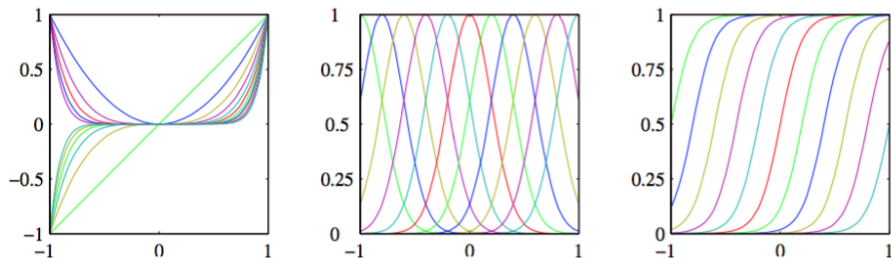


Figure: Examples of basis functions.

Maximum Likelihood and Least Squares I

The target variable t is given by a deterministic function $y(\mathbf{x}, \mathbf{w})$ with additive Gaussian noise so that

$$t = y(\mathbf{x}, \mathbf{w}) + \epsilon \quad (4)$$

where ϵ is a zero mean Gaussian random variable with precision β .

So, we can write

$$p(t|\mathbf{x}, \mathbf{w}, \beta) = \mathcal{N}(t|y(\mathbf{x}, \mathbf{w}), \beta^{-1}) \quad (5)$$

Consider a data set of inputs X with corresponding target value vector \mathbf{t} , make the assumption that these data points are drawn independently from the distribution. So the likelihood function is:

$$\ln p(\mathbf{t}|\mathbf{w}, \beta) = \sum_{n=1}^N \ln \mathcal{N}(t_n|\mathbf{w}^T \phi(\mathbf{x}_n), \beta^{-1}) \quad (6)$$

Maximum Likelihood and Least Squares II

The sum-of-squares error function is defined by:

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \boldsymbol{\phi}(\mathbf{x}_n)\}^2 \quad (7)$$

Let the gradient of the log likelihood function equals to 0. Solving for \mathbf{w} we obtain (which are known as the *normal equations* for the least squares problem):

$$\mathbf{w}_{ML} = (\boldsymbol{\Phi}^T \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^T \mathbf{t} \quad (8)$$

Here $\boldsymbol{\Phi}$ is an $N \times M$ matrix (the *design matrix*):

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

Geometry of Least Squares

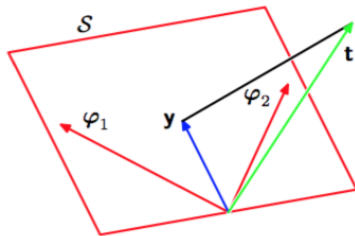


Figure: Examples of basis functions.

In an N -dimensional space whose axes are the value of t_1, \dots, t_N . 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.

Sequential Learning

Batch techniques, such as the maximum likelihood solution which can process large dataset in one go, if the dataset is sufficiently large, it may be worthwhile to use sequential algorithms, known as *on-line* algorithm. Sequential learning is also appropriate for real-time applications.

Stochastic gradient descent is applied. If the error function comprises a sum over data points $E = \sum_n E_n$, the update rule:

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta \nabla E_n \quad (9)$$

For the case of the sum-of-squares error function, this gives:

$$\mathbf{w}^{\tau+1} = \mathbf{w}^{\tau} - \eta (t_n - \mathbf{w}^{(\tau)T} \phi_n) \phi_n \quad (10)$$

Regularized Least Squares I

The regularization terms are added to the error functions. One of the simplest forms of regularizer is given by the sum-of-squares of the weight vector elements:

$$E_W(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w} \quad (11)$$

Consider the sum-of-squares error function with regularizer (quadratic regularizer):

$$E = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(\mathbf{x}_n)\}^2 + \frac{\lambda}{2} \mathbf{w}^T \mathbf{w} \quad (12)$$

This choice of regularizer is known as *weight decay*, because in sequential learning algorithms, it encourages weight values to decay towards zero.

Solving for \mathbf{w} as the gradient of (12) equal to zero, we obtain:

$$\mathbf{w} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t} \quad (13)$$

This represents a simple extension of the least-squares solution (8).

Regularized Least Squares II

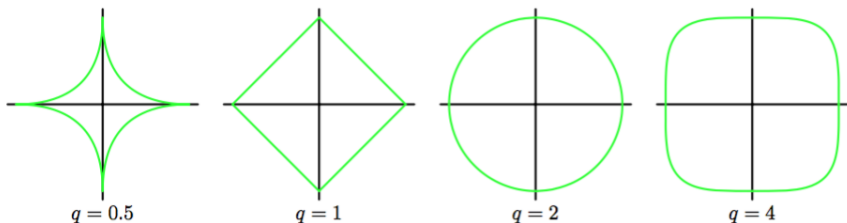


Figure: Contours of the regularization term for various values of the parameter q .

A more general regularizer is

$$\frac{\lambda}{2} \sum_{j=1}^M |w_j|^q \quad (14)$$

when $q=2$ corresponds to the quadratic regularizer in (12).

Regularized Least Squares III

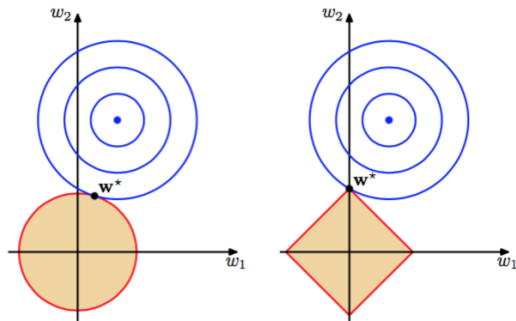


Figure: The error function without regularization and the constraint region for quadratic regularizer (left) and lasso regularizer (right). We can see the w in lasso equals to 0 make a sparse solution.

The case of $q=1$ is known as the *lasso* which leading to a *sparse* model in which the corresponding basis functions play no role. Regularization allows complex models to be trained on data sets of limited size without severe over-fitting, essentially by limiting the effective model complexity.

However, the problem of determining the optimal model complexity is then shifted from one of finding the regularization coefficient λ .

Multiple Outputs

If the target value t turns to a vector \mathbf{t} which means multiple outputs. This could be done via multiple, independent regression problems.

However, usually we use the same set of basis functions to model all components of the target value vector. Suppose we take the conditional distribution of the target vector to be an isotropic Gaussian of the form

$$p(\mathbf{t}|\mathbf{x}, \mathbf{W}, \beta) = \mathcal{N}(\mathbf{t}|\mathbf{W}\phi(\mathbf{x}), \beta^{-1}\mathbf{I}) \quad (15)$$

The log likelihood function is then given by:

$$\ln p(\mathbf{T}|\mathbf{X}, \mathbf{W}, \beta) = \sum_{n=1}^N \mathcal{N}(\mathbf{t}_n|\mathbf{W}^T\phi(\mathbf{x}_n), \beta^{-1}\mathbf{I}) \quad (16)$$

and maximize this function with respect to \mathbf{W} .

Bias-Variance trade-off

The phenomenon of overfitting is an unfortunate property of maximum likelihood. Consider the Bayesian view of model complexity in depth. We consider a frequentist viewpoint of the model complexity issue, known as the bias-variance trade-off.

Various loss functions leads to different corresponding optimal prediction, once we are given the conditional distribution $p(t|\mathbf{x})$. A popular choice is the square loss function, for which the optimal prediction is given by the conditional expectation – we denote by $h(\mathbf{x})$ and is given by:

$$h(\mathbf{x}) = \mathbb{E}(t|\mathbf{x}) = \int tp(t|\mathbf{x})dt \quad (17)$$

Loss Function for Regression - Chap 1.5.5

The decision stage consists of choosing a $y(\mathbf{x})$ of the value of t for each input \mathbf{x} . Suppose we incur a loss $L(t, y(\mathbf{x}))$, the average or expected loss is given by:

$$\mathbb{E}[L] = \int \int L(t, y(\mathbf{x})) p(\mathbf{x}, t) d\mathbf{x} dt \quad (18)$$

Our goal is to choose $y(\mathbf{x})$ so as to minimize $\mathbb{E}[L]$, assume a flexible function $y(\mathbf{x})$, do this using the calculus of variations to give:

$$\frac{\partial \mathbb{E}[L]}{\partial y(\mathbf{x})} = 2 \int \{y(\mathbf{x}) - t\} p(\mathbf{x}, t) dt = 0 \quad (19)$$

Solving for $y(\mathbf{x})$ and using the sum and product rules of probability, we obtain:

$$y(\mathbf{x}) = \frac{\int t p(\mathbf{x}, t) dt}{p(\mathbf{x})} = \int t p(t|\mathbf{x}) dt = \mathbb{E}_t[t|\mathbf{x}] \quad (20)$$

which is the conditional average of t conditioned on \mathbf{x} and is known as the regression function.

Loss Function for Regression - Chap 1.5.5

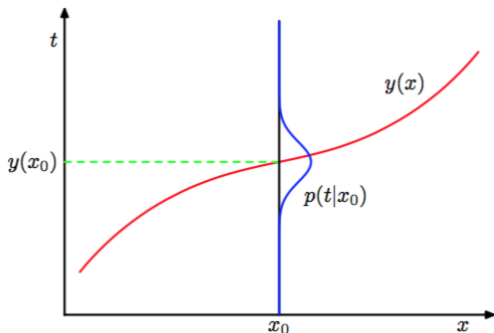


Figure: The regression function $y(x)$, which minimizes the expected squared loss, is given by the mean of the conditional distribution $p(t|x)$

- It is easy to be extended into multiple target variables represented by the vector \mathbf{t} , in which case the optimal solution is the conditional average $y(\mathbf{x}) = \mathbb{E}[\mathbf{t}|\mathbf{x}]$.

Loss Function for Regression - Chap 1.5.5

We can expand the square term as :

$$\{y(\mathbf{x}) - t\}^2 = \{y(\mathbf{x} - \mathbb{E}[t|\mathbf{x}] + \mathbb{E}[t|\mathbf{x}]) - t\}^2 \quad (21)$$

We use $\mathbb{E}[t|\mathbf{x}]$ to denote $\mathbb{E}_t[t|\mathbf{x}]$, and then performing the integral over t , the cross-term in (21) will vanish and then we have:

$$E[L] = \int \{y(\mathbf{x}) - t\}^2 p(\mathbf{x}) d\mathbf{x} + \int \text{var}[t|\mathbf{x}] p(\mathbf{x}) d\mathbf{x} \quad (22)$$

The function we seek to determine enters only in the first term, which will be minimized when $y(\mathbf{x})$ is equal to $\mathbb{E}[t|\mathbf{x}]$, in which case this term will vanish. This is simply the result that we derived previously and that show that the optimal least squares predictor is given by the conditional mean. The second term is the variance of the distribution of t , averaged over \mathbf{x} .

Expected squared loss

So we can write the expected squared loss:

$$E[L] = \int \{y(\mathbf{x}) - h(\mathbf{x})\}^2 p(\mathbf{x}) d\mathbf{x} + \int \int \{h(\mathbf{x} - t)^2\} p(\mathbf{x}) d\mathbf{x} \quad (23)$$

The second term is independent of $y(\mathbf{x})$, arises from the intrinsic noise on the data and represents the minimum achievable value of the expected loss. For the first term, the smallest value is zero. If we had an unlimited supply of data, we might find the regression function $h(\mathbf{x})$ to any desired degree of accuracy, which represents the optimal choice for $y(\mathbf{x})$.

However, in practice we have limited dataset of \mathcal{D} containing N data points, and consequently we do not know the regression function $h(\mathbf{x})$ exactly.

Algorithm Assessment

If we want to model the $h(\mathbf{x})$ using a parametric function $y(\mathbf{x}, \mathbf{w})$ governed by vector \mathbf{w} , then from a Bayesian perspective the uncertainty in our model is expressed through a posterior distribution over \mathbf{w} .

From the frequentist treatment, involves making a point estimate of \mathbf{w} based on the data set \mathcal{D} , and tries instead to interpret the uncertainty of this estimate through the thought experiment:

- 1 suppose we had a large data sets with size N , and each draw independently from $p(t, \mathbf{x})$.
- 2 for the given dataset \mathcal{D} , we use learning algorithm obtain a prediction function $y(\mathbf{x}; \mathcal{D})$.
- 3 different dataset, we have different prediction function.
- 4 then we have different values of squared loss.
- 5 the particular algorithm is then assessed by the average over the ensemble of datasets.

Bias-Variance in Expectation

Consider the integrand of the first term of $\mathbb{E}[L]$, for a particular data set \mathcal{D} takes the form

$$\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}^2 \quad (24)$$

we take its average over the ensemble of data sets. Because this quantity will be dependent on the particular data set \mathcal{D} .

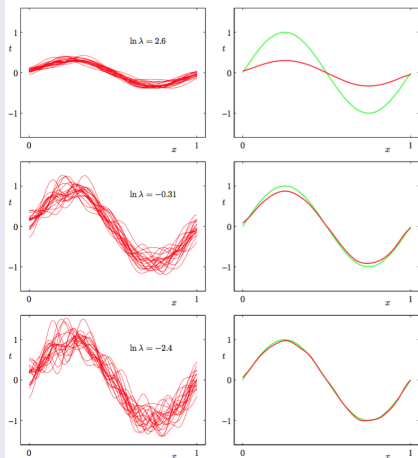
We now take the expectation of this expression with respect to \mathcal{D} and note that the final term will vanish , giving:

$$\begin{aligned} \mathbb{E}_{\mathcal{D}}[\{y(\mathbf{x}; \mathcal{D}) - h(\mathbf{x})\}] &= \underbrace{\{\mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})] - h(\mathbf{x})\}^2}_{(bias)^2} + \\ &\quad \underbrace{\mathbb{E}_{\mathcal{D}}[\{y(\mathbf{x}; \mathcal{D}) - \mathbb{E}_{\mathcal{D}}[y(\mathbf{x}; \mathcal{D})]\}^2]}_{variance} \end{aligned} \quad (25)$$

We see that the expected squared difference between $y(\mathbf{x}, \mathcal{D})$ and the regression function $h(\mathbf{x})$ can be expressed as the sum of two terms.

Bias-Variance in Expectation

Various data points.



Comparisons

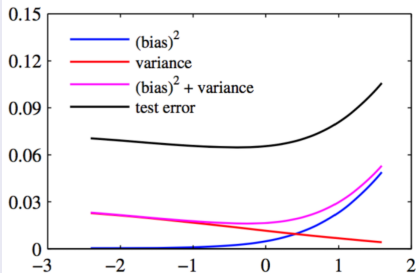


Illustration of the dependence of bias and variance on model complexity governed by a regularization parameter λ and the results for $(bias)^2 + variance$.

Bias-Variance in Expectation

From the comparisons curve we can check the variations for different λ . When λ is large (in the top row of last page left graph), we can see the red curves plot look similar in left column– low variance, but in the right column two curves are quite different – high bias.

In the bottom row which is quite different which λ is small and we have large variance – shown by the high variability between the red curves in the left plot, and low bias shown by the good fit to the regression function.

Note that the result of averaging many solutions for the complex model with $M = 25$ is a very good fit to the regression function.

Indeed, a weighted averaging of multiple solutions lies at the heart of a Bayesian approach, with respect to the posterior distribution of parameters (not with respect to multiple data sets).

Parameter Distribution

Predictive Distribution

Equivalent Kernel

Bayesian Model Comparison

Evaluation of Approximation Function

Maximizeing the Evidence Function

- Lorem ipsum dolor sit amet, consectetur adipiscing elit
- Aliquam blandit faucibus nisi, sit amet dapibus enim tempus eu
- Nulla commodo, erat quis gravida posuere, elit lacus lobortis est, quis porttitor odio mauris at libero
- Nam cursus est eget velit posuere pellentesque
- Vestibulum faucibus velit a augue condimentum quis convallis nulla gravida

Blocks of Highlighted Text

Block 1

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Integer lectus nisl, ultricies in feugiat rutrum, porttitor sit amet augue. Aliquam ut tortor mauris. Sed volutpat ante purus, quis accumsan dolor.

Block 2

Pellentesque sed tellus purus. Class aptent taciti sociosqu ad litora torquent per conubia nostra, per inceptos himenaeos. Vestibulum quis magna at risus dictum tempor eu vitae velit.

Block 3

Suspendisse tincidunt sagittis gravida. Curabitur condimentum, enim sed venenatis rutrum, ipsum neque consectetur orci, sed blandit justo nisi ac lacus.

Heading

- 1 Statement
- 2 Explanation
- 3 Example

Lorem ipsum dolor sit amet, consectetur adipiscing elit. Integer lectus nisl, ultricies in feugiat rutrum, porttitor sit amet augue. Aliquam ut tortor mauris. Sed volutpat ante purus, quis accumsan dolor.

Table

Treatments	Response 1	Response 2
Treatment 1	0.0003262	0.562
Treatment 2	0.0015681	0.910
Treatment 3	0.0009271	0.296

Table: Table caption

Theorem

Theorem (Mass–energy equivalence)

$$E = mc^2$$

Example (Theorem Slide Code)

```
\begin{frame}  
\frametitle{Theorem}  
\begin{theorem}[Mass--energy equivalence]  
$E = mc^2$  
\end{theorem}  
\end{frame}
```

Figure

Uncomment the code on this slide to include your own image from the same directory as the template .TeX file.

An example of the `\cite` command to cite within the presentation:

This statement requires citation [Smith, 2012].



John Smith (2012)

Title of the publication

Journal Name 12(3), 45 – 678.

The End