

Homework 1: Linear Regression

Writeup due 23:59 on Friday 6 February 2015

You will do this assignment individually and submit your answers as a PDF via the Canvas course website. There is a mathematical component and a programming component to this homework.

1. Non-Uniformly Weighted Data [7pts]

Consider a data set in which each data point t_n is associated with a weighting factor $r_n > 0$, so that the sum-of-squares error function becomes

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

Find an expression for the solution \mathbf{w}^* that minimizes this error function.

Solution 1: Take the gradient and set to zero:

$$\frac{\partial}{\partial \mathbf{w}} E_D(\mathbf{w}) = - \sum_{n=1}^N r_n \{t_n - \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n)\} \boldsymbol{\phi}(\mathbf{x}_n) = 0.$$

Solve for \mathbf{w} :

$$\begin{aligned} \sum_{n=1}^N r_n t_n \boldsymbol{\phi}(\mathbf{x}_n) &= \left(\sum_{n=1}^N r_n \boldsymbol{\phi}(\mathbf{x}_n) \boldsymbol{\phi}(\mathbf{x}_n)^\top \right) \mathbf{w} \\ \mathbf{w} &= \left(\sum_{n=1}^N r_n \boldsymbol{\phi}(\mathbf{x}_n) \boldsymbol{\phi}(\mathbf{x}_n)^\top \right)^{-1} \left(\sum_{n=1}^N r_n t_n \boldsymbol{\phi}(\mathbf{x}_n) \right) \end{aligned}$$

Solution 2 (Matrix form): We rewrite the error function in terms of matrix products:

$$\begin{aligned} E_D(\mathbf{w}) &= \frac{1}{2} \sum_{n=1}^N r_n \{t_n - \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n)\}^2 \\ &= \frac{1}{2} (\boldsymbol{\Phi} \mathbf{w} - \mathbf{t})^\top \mathbf{R} (\boldsymbol{\Phi} \mathbf{w} - \mathbf{t}) \\ &= \frac{1}{2} (\mathbf{w}^\top \boldsymbol{\Phi}^\top \mathbf{R} \boldsymbol{\Phi} \mathbf{w} - \mathbf{w}^\top \boldsymbol{\Phi}^\top \mathbf{R} \mathbf{t} - \mathbf{t}^\top \mathbf{R} \boldsymbol{\Phi} \mathbf{w} + \mathbf{t}^\top \mathbf{R} \mathbf{t}) \\ &= \frac{1}{2} (\mathbf{w}^\top \boldsymbol{\Phi}^\top \mathbf{R} \boldsymbol{\Phi} \mathbf{w} - 2 \mathbf{t}^\top \mathbf{R} \boldsymbol{\Phi} \mathbf{w} + \mathbf{t}^\top \mathbf{R} \mathbf{t}) \end{aligned}$$

Where we have defined $\mathbf{R} = \text{diag}(r_1, \dots, r_N)$.

Taking the gradient of the error function:

$$\begin{aligned} \nabla E_D(\mathbf{w}) &= \boldsymbol{\Phi}^\top \mathbf{R} \boldsymbol{\Phi} \mathbf{w} - \mathbf{t}^\top \mathbf{R} \boldsymbol{\Phi} \\ \mathbf{w} &= (\boldsymbol{\Phi}^\top \mathbf{R} \boldsymbol{\Phi})^{-1} \mathbf{t}^\top \mathbf{R} \boldsymbol{\Phi} \\ &= (\boldsymbol{\Phi}^\top \mathbf{R} \boldsymbol{\Phi})^{-1} \boldsymbol{\Phi}^\top \mathbf{R} \mathbf{t} \end{aligned}$$

2. Priors and Regularization [7pts]

Consider the Bayesian linear regression model given in Bishop 3.3.1. The prior is

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

where α is the precision parameter that controls the variance of the Gaussian prior. The likelihood can be written as

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

Using the fact that the posterior is the product of the prior and the likelihood (up to a normalization constant), show that maximizing the log posterior (i.e., $\ln p(\mathbf{w} | \mathbf{t}) = \ln p(\mathbf{w} | \alpha) + \ln p(\mathbf{t} | \mathbf{w})$) is equivalent to minimizing the regularized error term given by $E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$ with

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

Do this by writing $\ln p(\mathbf{w} | \mathbf{t})$ as a function of $E_D(\mathbf{w})$ and $E_W(\mathbf{w})$, dropping constant terms if necessary. Conclude that maximizing this posterior is equivalent to minimizing the regularized error term given by $E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$. (Hint: take $\lambda = \alpha / \beta$)

Solution: $p(\mathbf{w} | \alpha)$ is a multivariate normal distribution. Suppose the dimension of \mathbf{w} is $D \times 1$. Plug the mean $\mathbf{0}$ and covariance matrix $\alpha^{-1} \mathbf{I}$ into the PDF of multivariate normal distribution:

$$p(\mathbf{w} | \alpha) = \frac{1}{\sqrt{(2\pi)^D \det(\alpha^{-1} \mathbf{I})}} \exp\left(-\frac{1}{2} \mathbf{w}^\top (\alpha^{-1} \mathbf{I})^{-1} \mathbf{w}\right)$$
$$\ln p(\mathbf{w} | \alpha) = \frac{1}{2} \mathbf{w}^\top \mathbf{w} \alpha + \text{constant} = -\alpha E_W(\mathbf{w}) + \text{constant}$$

Similarly,

$$p(\mathbf{t} | \mathbf{w}) = \frac{1}{\sqrt{2\beta^{-1}\pi}} \prod_{n=1}^N \exp\left(-\frac{(t_n - \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n))^2}{2\beta^{-1}}\right)$$
$$\ln p(\mathbf{t} | \mathbf{w}) = \frac{1}{2} \sum_{n=1}^N (t_n - \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n))^2 \beta + \text{constant} = -\beta E_D(\mathbf{w}) + \text{constant}$$

Therefore, maximizing $\ln p(\mathbf{w} | \alpha) + \ln p(\mathbf{t} | \mathbf{w})$ is equivalent to maximizing $-\beta E_D(\mathbf{w}) - \alpha E_W(\mathbf{w})$. Hence it is equal to minimizing $E_D(\mathbf{w}) + \lambda E_W(\mathbf{w})$, where $\lambda = \alpha / \beta$. (Note that $\beta > 0$ because it is a variance.)

3. Modeling Motorcycle Helmet Forces [10pts]

The objective of this problem is to learn about linear regression with basis functions by modeling the g-forces associated with motorcycle helmet impacts. Download the file `motorcycle.csv` from the course website. It has two columns. The first one is the number of milliseconds since impact and the second is the g-force on the head. The data file looks like this:

```
"time since impact (ms)", "g force"  
2.4,0  
2.6,-1.3  
3.2,-2.7  
3.6,0  
4,-2.7  
6.2,-2.7  
6.6,-2.7  
6.8,-1.3  
...
```

and you can see a plot of the data in Figure 1.

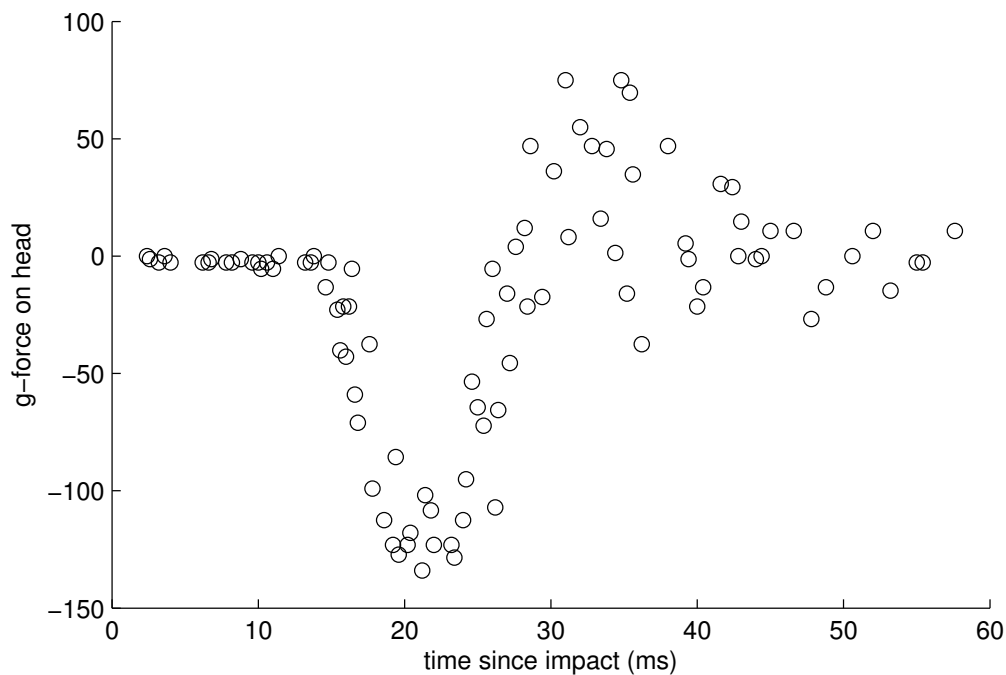


Figure 1: Motorcycle crash helmet data. The horizontal axis is time since impact and the vertical axis is force on the head.

Implement basis function regression with ordinary least squares.¹ Some sample Python code is provided in `linreg.py`. Plot the data and regression lines for the simple linear case, and for each of the following sets of basis functions:

(a) $\phi_j(x) = \exp\{-(x - 10j)/5)^2\}$ for $j = 0, \dots, 6$

(b) $\phi_j(x) = \exp\{-(x - 10j)/10)^2\}$ for $j = 0, \dots, 6$

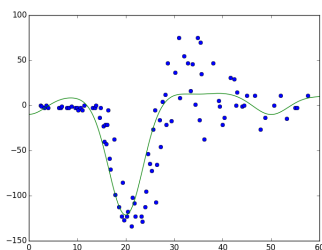
(c) $\phi_j(x) = \exp\{-(x - 10j)/25)^2\}$ for $j = 0, \dots, 6$

(d) $\phi_j(x) = x^j$ for $j = 0, \dots, 10$

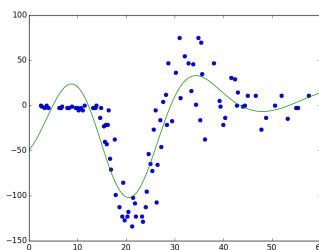
(e) $\phi_j(x) = \sin\{x/j\}$ for $j = 1, \dots, 20$

In addition to the plots, provide one or two sentences for each, explaining whether you think it is fitting well, overfitting or underfitting. If it does not fit well, provide a sentence explaining why.

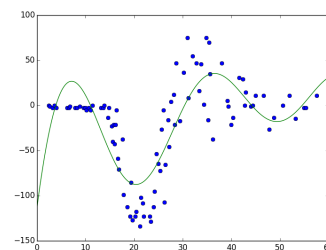
Solution: See the posted solution code.



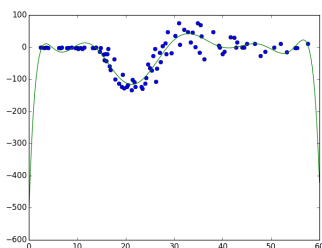
(a)



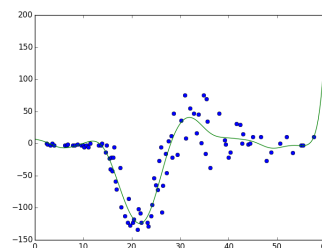
(b)



(c)



(d)



(e)

(a) fits relatively well. (b) and (c) exhibit the problem of underfitting: the flat tail on the left is not properly fit in these two plots. If we plot the basis functions for (b) and (c), we can see that they are wider than (a), and cannot capture the complexity in the data

¹Note that the data clearly don't have fixed variance! There is obviously less variance on the left of the plot. Modeling such *heteroscedastic* data is beyond the scope of the course.

pattern as well. (d) and (e) are overfitting: they fit the given data well, but also have captured the noise. The steep rising / falling tails are not likely patterns of the original data. This is because we fit the data to high-degree basis functions (10th and 20th degree respectively).

Calibration [1pt]

Approximately how long did this homework take you to complete?

Changelog

- **v1.0** – 28 January 2015 at 13:00
- **v1.1** – 28 January 2015 at 22:30. Removed extra λ in $E_W(w)$, minor formatting edits