

# Lab 4: Linear Regression

## 1 Overview

Linear regression is a fundamental tool for learning the relationship between a collection of independent variables and corresponding dependent variables. For simplicity, in this lab, we consider modeling a function  $\mathbb{R} \rightarrow \mathbb{R}$  by finding the least squared error fit to the data. For a collection of observations  $(x_n, t_n)$  and basis functions  $\phi_j(x)$ , the penalized least squares estimate is a linear combination  $\sum_j w_j \phi_j(x)$  which minimizes

$$\min_{\mathbf{w}} \sum_{n=1}^N \frac{1}{2} \left( \sum_j w_j \phi_j(x_n) - t_n \right)^2 + \frac{1}{2\mu^2} \sum_j w_j^2.$$

One gets to this model from a probabilistic perspective by assuming a Gaussian noise model with a Gaussian prior on the parameters. The nature of this problem makes it simple to solve for the minimum (also known as the maximum likelihood estimate) using standard linear algebra methods.

In this lab you are provided with code to work with two types of basis functions  $\phi_j$ : Gaussian curves and hat functions. After writing a function which performs this least squares fit, you will apply the model to a small synthetic data set and study the effect of the prior (a.k.a. penalty) parameter  $\mu$ .

## 2 Provided Resources

- `lsefit.m` - Function to be completed which calculates a maximum likelihood fit to observations  $(x, t) \in \mathbb{R}^2$  in a least squared error sense, using a provided collection of basis functions and the specified penalty  $\mu$ .
- `simple.mat` - A small synthetic collection of observations  $(x, t) \in \mathbb{R}^2$  with  $x \in [0, 2\pi]$  sampled from a simple function (sin) with a small amount of noise.
- `test.mat` - Another collection of observations of the same function as in `simple.mat`, to be used for model selection. The contained variables are named `test_x` and `test_t`.
- `gauss_basis.m` - Function which generates a basis of Gaussian curve functions over the specified interval. The output is parameters for these basis functions, to be used by `func_gauss.m`. Each column of the output contains the mean and standard deviation for a Gaussian function in one dimension.
- `hat_basis.m` - Function which generates a basis of hat functions over the specified interval. The output is parameters for these basis functions, to be used by `func_hat.m`. Each columns of the output contains the start and end of the support interval for a hat function.
- `func_gauss.m` - Evaluates a Gaussian curve in one dimension with the given parameters, at the specified locations.
- `func_hat.m` - Evaluates a hat function in one dimension with the given parameters, at the specified locations.
- `eval_basis.m` - Constructs the design matrix  $\Phi$  given basis functions, their parameters, and data locations  $x$ .

### 3 Guide

1. Complete the function `lsefit.m`. This function should find a maximum likelihood estimate for the coefficients  $w_j$  in the model mentioned previously. To do this, let the design matrix  $\Phi$  be such that  $\Phi_{i,j} = \phi_j(x_i)$ . Differentiating the energy and setting the gradient equal to zero, we arrive at (you should be able to derive this):

$$\mathbf{w} = \left( \Phi^\top \Phi + \frac{1}{\mu^2} I \right)^{-1} \left( \Phi^\top \mathbf{t} \right).$$

You will need to employ `eval_basis.m` to construct  $\Phi$  using the basis function and parameters passed as an argument to `lsefit.m`. For example, the value of the second basis function at the location  $x = 5$  is calculated using `func(5, params(:,2))`.

2. Train the model on the data in `simple.mat` using ten hat functions and  $\mu = 10^5$ . Plot and turn in the learned model (the function fit to the data) on the interval  $[0, 2\pi]$ .
3. Do the same for other values of the hyperparameter such as  $\mu = 10$  and  $\mu = 1$ .
4. What, if anything, is the hyperparameter controlling? How does it impact the solution to the problem?
5. Train the model with ten Gaussian basis functions and  $\mu = 10^5$ , then plot and turn in the learned model evaluated (the function fit to the data) on the interval  $[0, 2\pi]$ .
6. As before, explore other values of the hyperparameter such as  $\mu = 10$  and  $\mu = 1$ .
7. What, if anything, is the hyperparameter controlling in this case? How does the dependence of the solution on the hyperparameter differ for this basis when compared to the hat basis?
8. Fit the model with integer values of  $\mu = 1$  to 100, and using a Gaussian basis of  $M = 10$  elements on the data in `simple.mat`. For each model, calculate the squared error for the observations in `test.mat` (therefore testing the model). Generate and turn in a plot with  $\mu$  on the x-axis and the total squared model error on the test data along the y-axis.
9. What value of  $\mu$ , when trained on the data in `simple.mat`, performs best on the data in `test.mat`? How do you know? Explain the shape of the plot you generated in the previous step.
10. Repeat the process, now fixing  $\mu = 13$  and varying the number of basis elements from  $M = 1$  to 100. Generate and turn in a plot with the number of basis elements on the x-axis and the error for the test data on the y-axis.

### 4 Extra

The Gaussian prior placed on the coefficients is only one possibility. In another situation, we may want to use an exponential prior resulting in the model

$$\sum_n \left( \sum_j w_j \phi_j(x_n) - t_n \right)^2 + \frac{1}{\mu} \sum_j |w_j|.$$

1. How do you think this model will perform differently than the Gaussian prior?
2. What is the optimality condition in this case (set the gradient equal to zero, but be careful about the absolute value)? How could you solve this problem computationally?