Weiyi Chen
wec427@g.harvard.edu
CS181-S16

Collaborators: N.A.

# Homework 1: Linear Regression

You should submit your answers as a PDF via the Canvas course website. There is a mathematical component and a programming component to this homework. You may collaborate with others, but are expected to list collaborators, and write up your problem sets individually.

Please type your solutions after the corresponding problems using this LaTeX template, and start each problem on a new page.

---

**Problem 1** (Centering and Ridge Regression, 7pts)

Consider a data set in which each data input vector $x \in \mathbb{R}^n$ is centered, meaning $\forall x, \sum_i x_i = 0$. Let $X \in \mathbb{R}^{n \times m}$ be the input matrix, the columns of which are the input vectors. Let $\lambda$ be a positive constant. We define:

$$J(w, w_0) = (y - Xw - w_0 \mathbf{1})^T (y - Xw - w_0 \mathbf{1}) + \lambda w^T w$$

(a) Compute the gradient of $J(w, w_0)$ with respect to $w_0$. Simplify as much as you can for full credit.

(b) Compute the gradient of $J(w, w_0)$ with respect to $w$. Simplify as much as you can for full credit. Make sure to give your answer in matrix form.

(c) Suppose that $\lambda > 0$. Knowing that $J$ is a convex function of its arguments, conclude that a global optimizer of $J(w, w_0)$ is

$$w_0 = \frac{1}{n} \sum_i y_i \tag{1}$$

$$w = (X^T X + \lambda I)^{-1} X^T y \tag{2}$$

Before taking the inverse of a matrix, prove that it is invertible.

---

**Solution**

(a) The gradient of $J(w, w_0)$ with respect to $w_0$ is

$$
\begin{aligned}
\nabla_{w_0} J(w, w_0) &= 2(y - Xw - w_0 \mathbf{1})^T \cdot (-\mathbf{1}) \\
&= -2(y - Xw - w_0 \mathbf{1})^T \mathbf{1} \\
&= -2(y - Xw)^T \mathbf{1} + 2 w_0 n \\
&= -2 y^T \cdot \mathbf{1} + w^T X^T \mathbf{1} + 2 w_0 n \\
&= -2 \sum_i y_i + w^T \mathbf{0} + 2 w_0 n \\
&= -2 \sum_i y_i + 2 w_0 n
\end{aligned}
$$

(b) The gradient of $J(w, w_0)$ with respect to $w$ is

$$
\begin{aligned}
\nabla_w J(w, w_0) &= 2(y - Xw - w_0\mathbf{1})^T \cdot (-X) + 2\lambda w^T \\
&= -2(y - Xw - w_0\mathbf{1})^T X + 2\lambda w^T \\
&= 2w^T(X^TX + \lambda I) - 2y^TX + 2w_0\mathbf{1}^TX \\
&= 2w^T(X^TX + \lambda I) - 2y^TX + 2w_0\mathbf{0} \\
&= 2w^T(X^TX + \lambda I) - 2y^TX
\end{aligned}
$$

(c) A global optimizer of $J(w, w_0)$ is derived by letting the above gradients as zero, i.e.

$$
\nabla_{w_0} J(w, w_0) = -2\sum_i y_i + 2w_0 n = 0 \Rightarrow w_0 = \frac{1}{n}\sum_i y_i
$$

$$
\nabla_w J(w, w_0) = 2w^T(X^TX + \lambda I) - 2y^TX = \mathbf{0} \Rightarrow w = (X^TX + \lambda I)^{-1}X^Ty
$$

In order to prove the inverse of $X^TX + \lambda I$ exists, we can verify its determinant is non-zero, i.e.

$$
det(X^TX + \lambda I) = |X|^2 + \lambda > 0
$$

given $|X|^2 \geq 0$ and $\lambda > 0$, therefore $X^TX + \lambda I$ is invertible.

**Problem 2** (Priors and Regularization,7pts)

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

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

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

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

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., $\ln p(w \mid t) = \ln p(w|\alpha) + \ln p(t \mid w)$) is equivalent to minimizing the regularized error term given by $E_D(w) + \lambda E_W(w)$ with

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

$$E_W(w) = \frac{1}{2} w^\mathsf{T} w$$

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

**Solution**

Write $\ln p(w \mid t)$ as a function of $E_D(w)$ and $E_W(w)$, i.e.

$$\ln p(w \mid t) = \ln p(w \mid \alpha) + \ln p(t \mid w)$$

$$= \ln \mathcal{N}(w \mid 0, \alpha^{-1}I) + \ln \prod_{n=1}^{N} \mathcal{N}(t_n \mid w^\mathsf{T} \phi(x_n), \beta^{-1})$$

$$= \ln \mathcal{N}(w \mid 0, \alpha^{-1}I) + \sum_{n=1}^{N} \ln \mathcal{N}(t_n \mid w^\mathsf{T} \phi(x_n), \beta^{-1})$$

where the multivariate normal distribution formula is

$$\mathcal{N}(x \mid \mu, \Sigma) = \frac{1}{\sqrt{(2\pi)^n |\Sigma|}} \exp\left(-\frac{1}{2}(x - \mu)^\mathsf{T} \Sigma^{-1}(x - \mu)\right)$$

via substituting, we derive

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

Dropping the constant terms, we have

$$\ln p(w \mid t) - C = -\frac{1}{2}\alpha w^T w - \frac{1}{2}\beta \sum_{n=1}^{N} (t_n - w^\mathsf{T} \phi(x_n))^2$$

$$= -\alpha E_W(w) - \beta E_D(w)$$

Taking $\lambda = \alpha/\beta$,

$$\frac{1}{\beta}\left(\ln p(\boldsymbol{w}\,|\,\boldsymbol{t}) - C\right) = -\lambda E_W(\boldsymbol{w}) - E_D(\boldsymbol{w})$$

Therefore, maximizing the posterior is equivalent to maximize $-\lambda E_W(\boldsymbol{w}) - E_D(\boldsymbol{w})$, or in other words, minimize the regularized error terms, i.e. $E_D(\boldsymbol{w}) + \lambda E_W(\boldsymbol{w})$.

## 3. Modeling Changes in Congress [10pts]

The objective of this problem is to learn about linear regression with basis functions by modeling the average age of the US Congress. The file `congress-ages.csv` contains the data you will use for this problem. It has two columns. The first one is an integer that indicates the Congress number. Currently, the 114th Congress is in session. The second is the average age of that members of that Congress. The data file looks like this:

```
congress,average_age
80,52.4959
81,52.6415
82,53.2328
83,53.1657
84,53.4142
85,54.1689
86,53.1581
87,53.5886
```
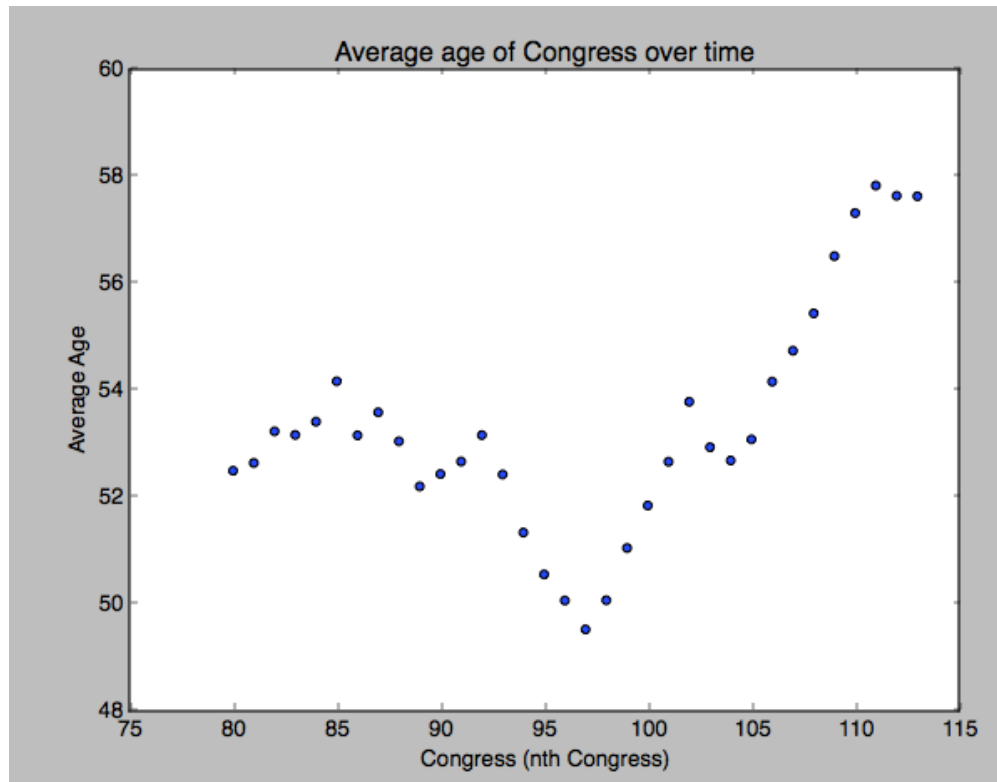
and you can see a plot of the data in Figure 1.



Figure 1: Average age of Congress. The horizontal axis is the Congress number, and the vertical axis is the average age of the congressmen.

**Problem 3** (Modeling Changes in Congress, 10pts)

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

  (a) $\phi_j(x) = x^j$ for $j = 1, \ldots, 7$

  (b) $\phi_j(x) = x^j$ for $j = 1, \ldots, 3$

  (c) $\phi_j(x) = \sin\{x/j\}$ for $j = 1, \ldots, 4$

  (d) $\phi_j(x) = \sin\{x/j\}$ for $j = 1, \ldots, 7$

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

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. A good fit should capture the most important trends in the data.

**Solution**

Basis function regression is implemented in `linreg.py`. Below are the plots for each case -
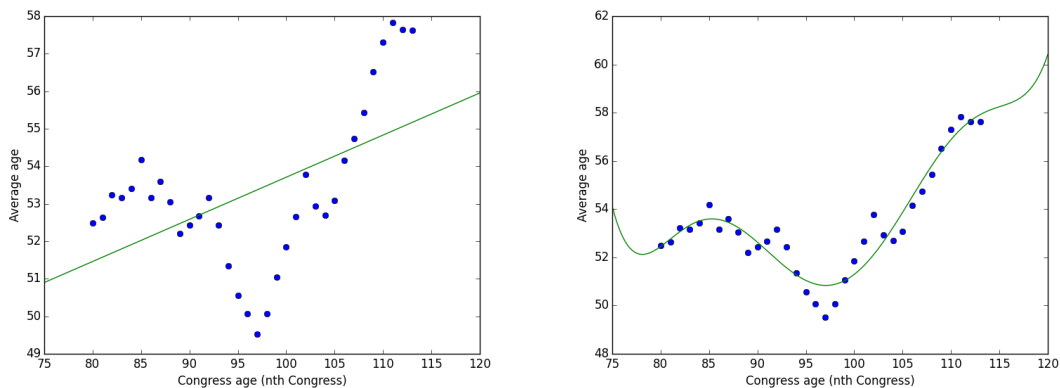


Figure 2: The data and regression line for the simple linear case and (a)

Following is my thinking of fitting well or not for each case -

- Fitting well: case (a), case (b), case (c), case (d)

- Overfitting: case (e)

- Underfitting: the simple linear case

The simple linear case does not fit well because it cannot capture the most important trends in the data, most of points are off from the line far away and the straight line cannot capture the curve indicated by the points either.

Case (e) is overfitting, though it captures every point exactly, it also fit to the error of each point, which is not expected. When the regression line goes out of the x range of the data points, the curve is going off far away to negative values very quickly, which is clearly unreasonable from our perspective.
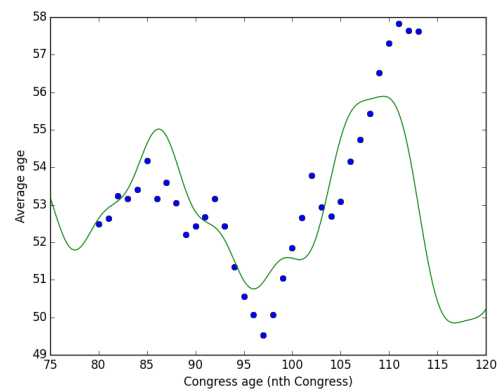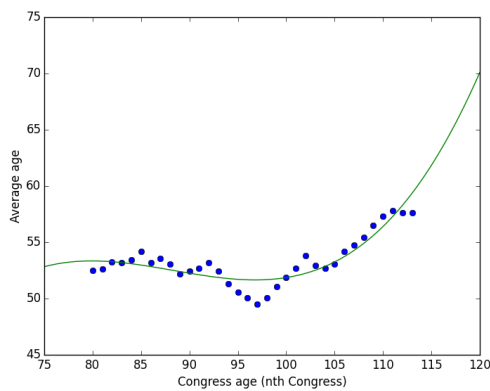
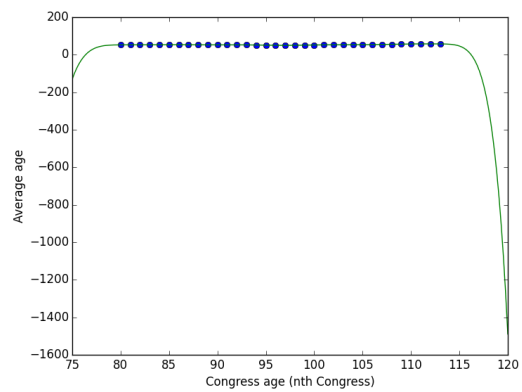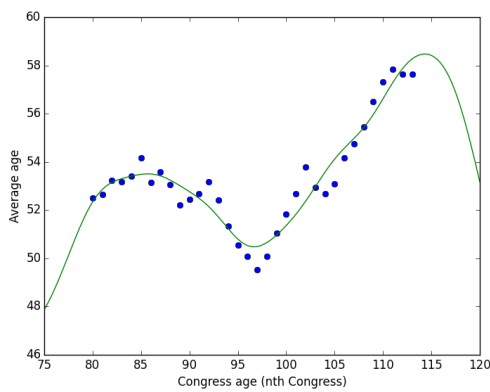Figure 3: The data and regression line for cases (b) and (c)



Figure 4: The data and regression line for cases (d) and (e)

**Problem 4** (Calibration, 1pt)

Approximately how long did this homework take you to complete?

**Answer:** 3 hours. Roughly 1 hour for each question, including investigation time on unfamiliar knowledge.