

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 L^AT_EX 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 $x \in \mathbb{R}^m$. Let $X \in \mathbb{R}^{n \times m}$ be the input matrix, the rows of which are the input vectors, and the columns of which are centered at 0. Let λ be a positive constant. We define:

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

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

$$w_0 = \frac{1}{n} \sum_i y_i \quad (1)$$

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

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

Solution

First we simplify:

$$J(w, w_0) = (y - Xw - w_0 \mathbf{1})^T (y - Xw - w_0 \mathbf{1}) + \lambda w^T w \quad (3)$$

$$J(w, w_0) = (y^T - w^T X^T - \mathbf{1}^T w_0) (y - Xw - w_0 \mathbf{1}) + \lambda w^T w \quad (4)$$

$$J(w, w_0) = y^T y - y^T Xw - y^T \mathbf{1} w_0 - w^T X^T y + w^T X^T Xw + w^T X^T \mathbf{1} w_0 - \mathbf{1}^T w_0 y + \mathbf{1}^T w_0 Xw + \mathbf{1}^T \mathbf{1} w_0^2 + \lambda w^T w \quad (5)$$

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

Using the expanded form we take the derivative with respect to w_0

$$\nabla_{w_0} = -y^T \mathbf{1} + w^T X^T \mathbf{1} - y^T \mathbf{1} + w^T X^T \mathbf{1} + 2(\mathbf{1}^T \mathbf{1}) w_0 \quad (6)$$

Now we can simplify $\mathbf{1}^T \mathbf{1} = n$ because this is the dot product of two 1-vectors of length n . Furthermore, because the columns of our matrix are centered at 0, $\mathbf{w}^T \mathbf{X}^T \mathbf{1} = 0$, and thus our final answer is:

$$\nabla_{w_0} = 2(nw_0 - \mathbf{y}^T \mathbf{1}) \quad (7)$$

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

Using the expanded form we take the derivative with respect to \mathbf{w}

$$\nabla_{\mathbf{w}} = -\mathbf{y}^T \mathbf{X} - \mathbf{y}^T \mathbf{X} + 2(\mathbf{w}^T \mathbf{X}^T \mathbf{X}) + \mathbf{1}^T \mathbf{X} w_0 + \mathbf{1}^T \mathbf{X} w_0 + 2\lambda \mathbf{w}^T \quad (8)$$

Once again, because the columns of our matrix are centered at 0, $\mathbf{1}^T \mathbf{X} w_0 = 0$, and thus our final answer is:

$$\nabla_{\mathbf{w}} = 2(\lambda \mathbf{w}^T - \mathbf{y}^T \mathbf{X} + \mathbf{w}^T \mathbf{X}^T \mathbf{X}) \quad (9)$$

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

$$w_0 = \frac{1}{n} \sum_i y_i \quad (10)$$

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (11)$$

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

Since J is convex, we only need to see where the gradient is 0. Thus we set our values in lines (7) and (9) equal to 0:

$$0 = 2(nw_0 - \mathbf{y}^T \mathbf{1}) \quad (12)$$

$$nw_0 = \mathbf{y}^T \mathbf{1} \quad (13)$$

$$w_0 = \frac{1}{n} \sum_i y_i \quad (14)$$

$$(15)$$

$$0 = 2(\lambda \mathbf{w}^T - \mathbf{y}^T \mathbf{X} + \mathbf{w}^T \mathbf{X}^T \mathbf{X}) \quad (16)$$

$$\left[\mathbf{w}^T \mathbf{X}^T \mathbf{X} + \lambda \mathbf{w}^T \right]^T = \left[\mathbf{y}^T \mathbf{X} \right]^T \quad (17)$$

$$\mathbf{X}^T \mathbf{X} \mathbf{w} + \lambda \mathbf{w} = \mathbf{X}^T \mathbf{y} \quad (18)$$

$$(\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{w} = \mathbf{X}^T \mathbf{y} \quad (19)$$

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I})^{-1} \mathbf{X}^T \mathbf{y} \quad (20)$$

Now to prove that the matrix is invertible, we will show that the matrix is Positive-definite. A positive-definite matrix $A \implies$ eigenvalues are positive, $\implies \det A > 0$, $\implies A$ is invertible.

Positive-definite matrix $A \implies$ positive eigenvalues,
 $\implies \det A > 0$
 $\implies A$ is invertible.

Now to show that a matrix is positive-definite we must show that

$$\mathbf{v}^T \mathbf{A} \mathbf{v} > 0 \quad \forall \mathbf{v} \neq 0$$

.

Looking at our matrix $\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}$:

$$\mathbf{v}^T (\mathbf{X}^T \mathbf{X} + \lambda \mathbf{I}) \mathbf{v} > 0 \tag{21}$$

$$\mathbf{v}^T \mathbf{X}^T \mathbf{X} \mathbf{v} + \mathbf{v}^T \lambda \mathbf{I} \mathbf{v} > 0 \tag{22}$$

$$\|\mathbf{X} \mathbf{v}\|^2 + \lambda \|\mathbf{v}\|^2 > 0 \tag{23}$$

Since these are both lengths, we know that they must both be greater than or equal to 0 in all cases. Moreover, since the definition defines $\mathbf{v} \neq 0$, we know that the final term must be greater than 0 proving that the matrix is in fact positive definite and thus invertible.

Problem 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

First let us consider the log-likelihood for a standard normal, $\mathcal{N}(x | \mu, \sigma^2)$.

$$p(x | \mu, \sigma^2) = \prod_{i=1}^N \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{1}{2\sigma^2}(x_i - \mu)^2}$$

$$\ln(p(x | \mu, \sigma^2)) = \sum_{i=1}^N \ln\left(\frac{1}{\sqrt{2\pi\sigma^2}}\right) + \sum_{i=1}^N -\frac{1}{2\sigma^2}(x_i - \mu)^2$$

$$\ln(p(x | \mu, \sigma^2)) = N \left(\ln\left(2\pi\sigma^2\right)^{-\frac{1}{2}} \right) - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2$$

$$\ln(p(x | \mu, \sigma^2)) = -\frac{N}{2} \ln(2\pi) - \frac{N}{2} \ln(\sigma^2) - \frac{1}{2\sigma^2} \sum_{i=1}^N (x_i - \mu)^2$$

Since both of our values are normally distributed, we can throw them into this key, dropping additive constants and adding them together (because we have taken the log, multiplicative constants are now additive). This yields:

$$\ln(p(\mathbf{w} \mid \alpha) * p(\mathbf{t} \mid \mathbf{w})) = \ln(\mathcal{N}(\mathbf{w} \mid \mathbf{0}, \alpha^{-1} \mathbf{I}) + \ln \prod_{n=1}^N \mathcal{N}(t_n \mid \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n), \beta^{-1})$$

$$\ln(p(\mathbf{w} \mid \alpha) * p(\mathbf{t} \mid \mathbf{w})) = -\frac{N}{2} \ln(2\pi) - \frac{N}{2} \ln(\alpha^{-1} \mathbf{I}) - \frac{\alpha}{2\mathbf{I}} \sum_i (w_i - 0)^2 - \frac{N}{2} \ln(2\pi) - \frac{N}{2} \ln(\beta^{-1}) - \frac{\beta}{2} \sum_n (t_n - \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n))^2$$

$$\ln(p(\mathbf{w} \mid \alpha) * p(\mathbf{t} \mid \mathbf{w})) = -\left(\frac{\alpha}{2\mathbf{I}} \sum_i (w_i)^2 + \frac{\beta}{2} \sum_n (t_n - \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n))^2 \right)$$

$$\ln(p(\mathbf{w} \mid \alpha) * p(\mathbf{t} \mid \mathbf{w})) = -\left(\frac{\alpha}{2\beta} \mathbf{w}^T \mathbf{w} + \frac{1}{2} \sum_n (t_n - \mathbf{w}^\top \boldsymbol{\phi}(\mathbf{x}_n))^2 \right)$$

$$\ln(p(\mathbf{w} \mid \alpha) * p(\mathbf{t} \mid \mathbf{w})) = -(\lambda E_W(\mathbf{w}) + E_D(\mathbf{w})) \text{ If } \lambda = \frac{\alpha}{\beta}$$

Since this is negative, and we are maximizing this value, we clearly want to minimize the positive of this value.

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.

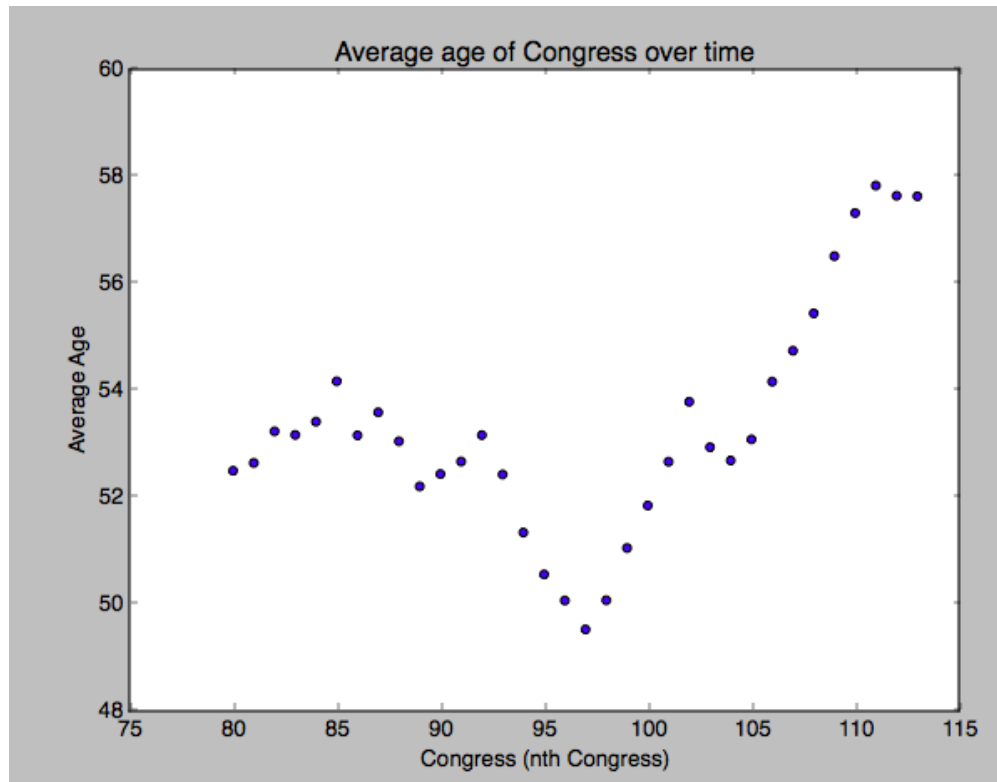


Figure 1: Average age of Congress. The horizontal axis is the Congress number, and the vertical axis is the average age of the congressmen.

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 $j = 1, \dots, 7$
- (b) $\phi_j(x) = x^j$ for $j = 1, \dots, 3$
- (c) $\phi_j(x) = \sin\{x/j\}$ for $j = 1, \dots, 4$
- (d) $\phi_j(x) = \sin\{x/j\}$ for $j = 1, \dots, 7$
- (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. A good fit should capture the most important trends in the data.

Solution

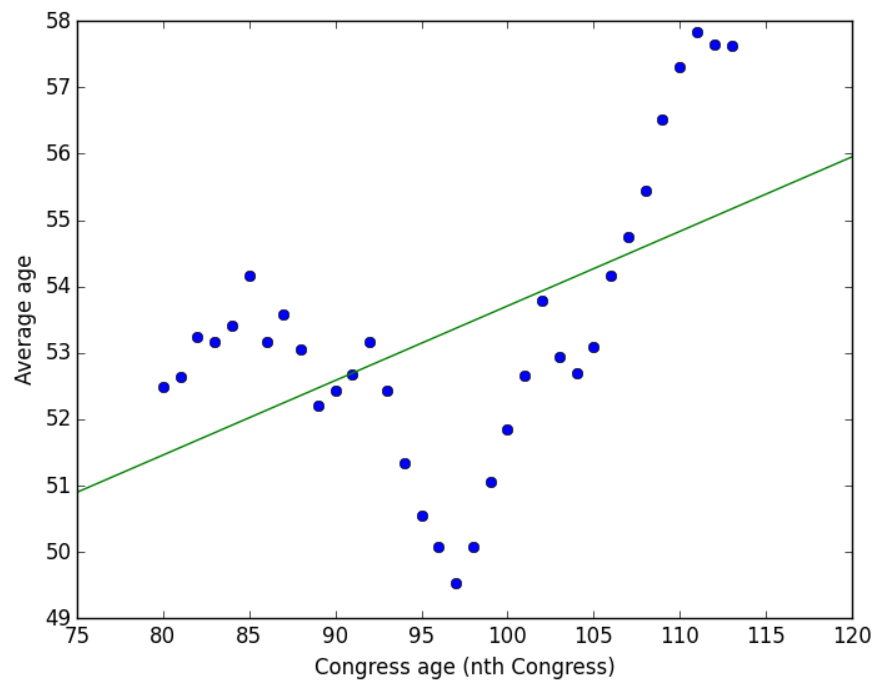


Figure 2: This is a simple linear regression, and clearly looks to be underfitting because it is missing the whole downward spike in the middle.

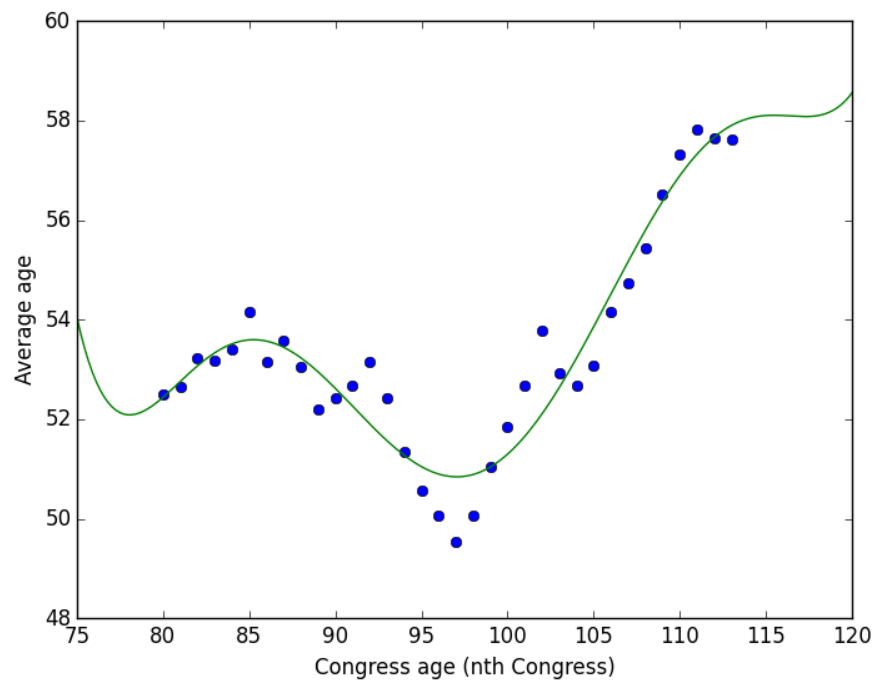


Figure 3: This example seems to be fitting the data better at the trend of a dropping of the congressional age. However since this uses terms to the 7th power, it may not be good at predicting future ages of Congress.

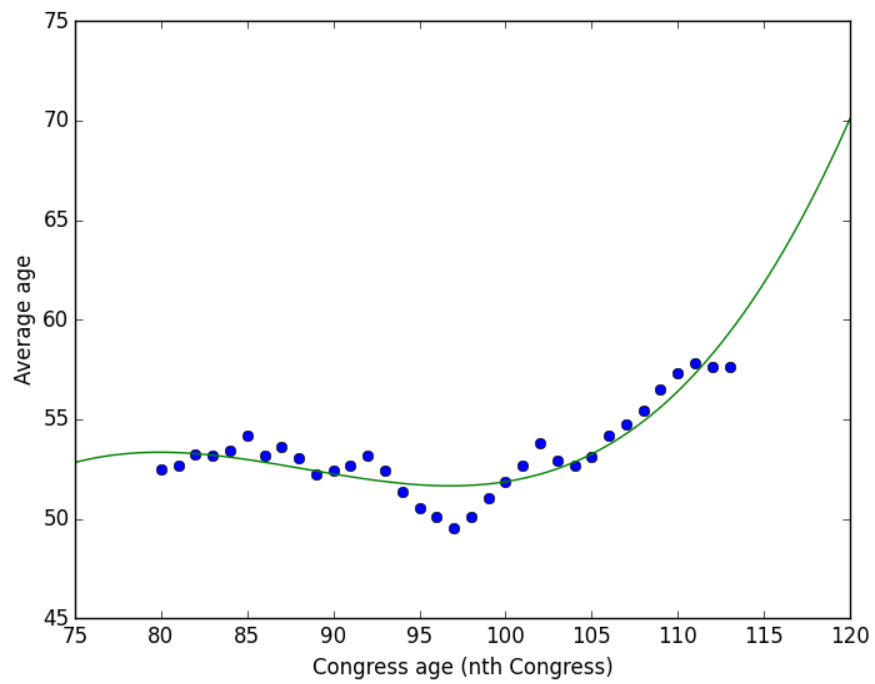


Figure 4: While this generally fits the previous data, the fit grows exponentially in the future, which likely won't model future Congressional compositions.

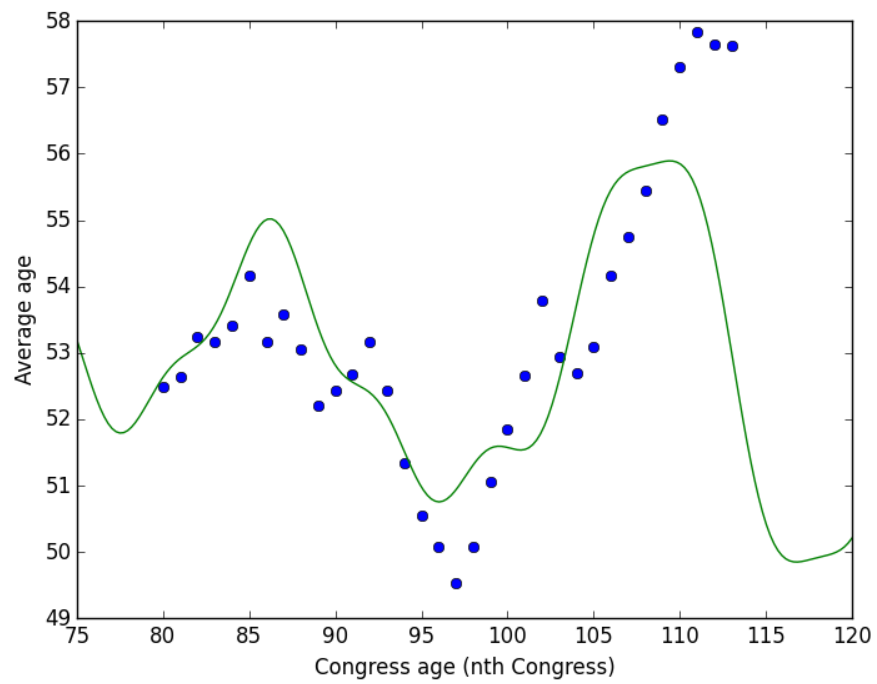


Figure 5: This is underfitting the data, missing the drop of age around the 97th congress and overinflating a spike after the 86th.

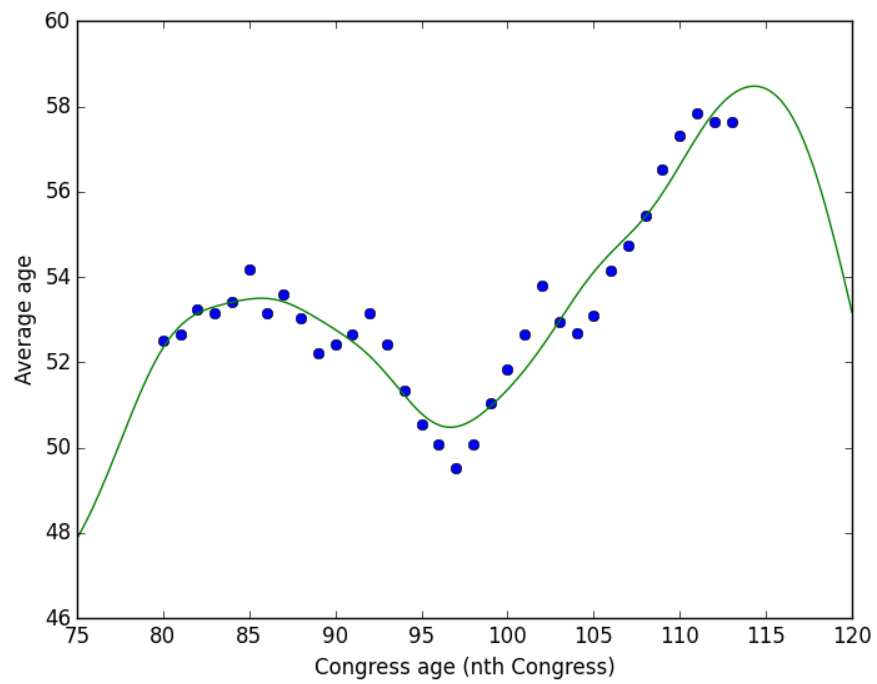


Figure 6: This fits the data very well, not overvaluing any parts. A word of caution is that it quickly expects the age of Congress to decrease in the future.

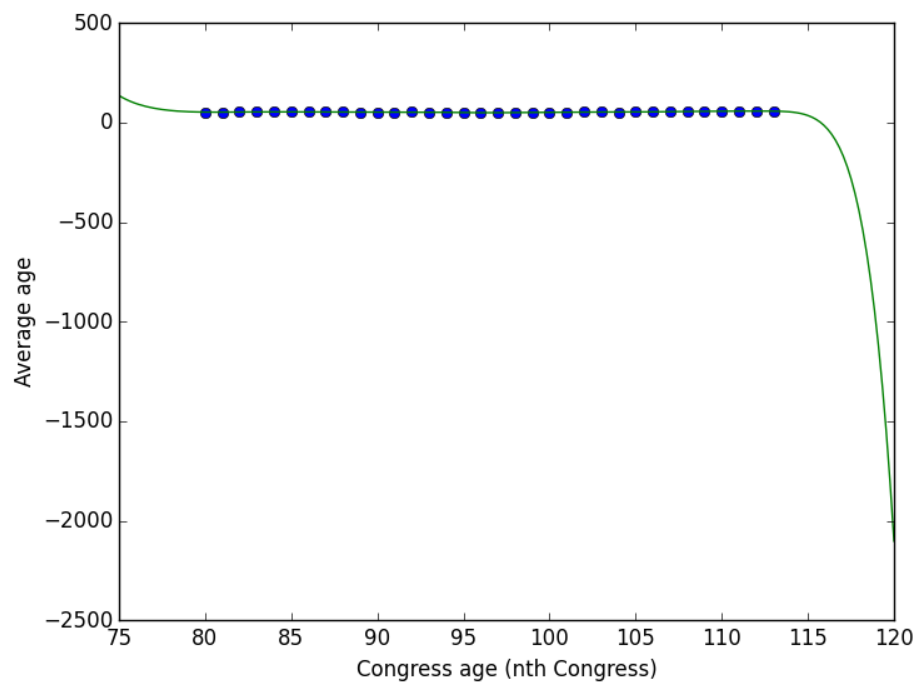


Figure 7: This figure super overfits the data. The fact that there are essentially as many terms as datapoints allows the program to draw the line directly through every point. This has 0 predictive value because this model is so overfit to past data, it doesn't try to distinguish the signal from the noise.

Problem 4 (Calibration, 1pt)

Approximately how long did this homework take you to complete?

Answer: 12