

Assignment 1: Solutions

1 Probabilistic Modeling (12 marks)

1. (4 marks) *What is the type of distribution that describes this situation? What are the parameters μ of this distribution? (See PRML Appendix B)*

Appendix B / Sec. 2.2 in PRML, on the multinomial distribution, is a useful read.

The parameter is a vector that determines the probability of each party winning. For example, this could be represented as a 4-dimensional vector $\mu \in \mathbb{R}^4$, where μ_i is the probability that party i wins.

2. (2 marks) *What would be the value of the parameters μ for an election where the outcome is an equal chance of any party winning?*

$$\mu = (1/4, 1/4, 1/4, 1/4)$$

3. (2 marks) *What would be the value of the parameters μ for an election that is completely “rigged”? E.g. the party currently in power is definitely going to win.*

$$\mu = (0, 0, 1, 0)$$

4. (4 marks) *Suppose my prior is that the Green Party has completely rigged the election. Assume I see a set of polls where the NDP has the largest share of the vote in each poll. What would be my posterior probability on the parameters μ ?*

Recall that the posterior is proportional to the likelihood times the prior. In this case, the poll results would be data, and one could compute the likelihood of the data.

If I truly believe the Green Party has perfectly rigged the election, the observed data are inconsistent with this. If I had a model of polls that allowed a party (e.g. the NDP) to win the poll even though they have no chance of winning the election, then this would be fine, and my posterior probability on the parameters μ would remain a Dirac delta function on $\mu = (0, 0, 0, 1)$.

2 Regularized Least-Squares Linear Regression (15 marks)

Show that the minimizer for least-squares linear regression with L_2 regularization is $\mathbf{w} = (\lambda \mathbf{I} + \Phi^T \Phi)^{-1} \Phi^T \mathbf{t}$.

The regularized error function we wish to minimize is:

$$E(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{y(x_n, \mathbf{w}) - t_n\}^2 + \frac{\lambda}{2} \|\mathbf{w}\|^2 \quad (1)$$

We take its derivatives and set them to zero:

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

$$\nabla E(\mathbf{w}) = 0 \quad (3)$$

$$\Leftrightarrow \sum_{n=1}^N t_n \phi(\mathbf{x}_n)^T = \left[\sum_{n=1}^N \mathbf{w}^T \phi(\mathbf{x}_n) \phi(\mathbf{x}_n)^T \right] + \lambda \mathbf{w}^T \quad (4)$$

$$\mathbf{t}^T \Phi = \mathbf{w}^T \Phi^T \Phi + \lambda \mathbf{w}^T \quad (5)$$

$$\mathbf{t}^T \Phi = \mathbf{w}^T (\Phi^T \Phi + \lambda \mathbf{I}) \quad (6)$$

$$\Phi^T \mathbf{t} = (\Phi^T \Phi + \lambda \mathbf{I})^T \mathbf{w} \quad (7)$$

$$\Phi^T \mathbf{t} = (\Phi^T \Phi + \lambda \mathbf{I}) \mathbf{w} \quad (8)$$

$$\mathbf{w} = (\Phi^T \Phi + \lambda \mathbf{I})^{-1} \Phi^T \mathbf{t} \quad (9)$$

3 Training vs. Test Error (12 marks)

1. (4 marks) *Suppose we perform unregularized regression on a dataset. Is the **training error** with a degree 10 polynomial always lower than or equal to that using a degree 9 polynomial? Explain.*

Yes. Degree 10 polynomials contain degree 9 polynomials. Unregularized regression leads to the optimal solution (caveats possible here). At worst, training error should be the same using a degree 10 polynomial as that using a degree 9 polynomial.

2. (4 marks) *Suppose we perform unregularized regression on a dataset. Is the **testing error** with a degree 10 polynomial always lower than or equal to that using a degree 9 polynomial? Explain.*

No. As discussed in lecture, more complex models can over-fit to a given training dataset, and therefore it is possible that a more complex model (degree 10 polynomial) could have higher testing error compared to a simpler model (degree 9 polynomial).

3. (4 marks) *Suppose we perform unregularized regression on a dataset. Is the **training error** always lower than the **testing error**? Explain.*

No. There are no guarantees on the relationship between training error and testing error. While it is often the case that training error is lower than testing error, it is always possible that due to a particular testing set choice, the testing data points lie perfectly on the learned curve, while some of the training data points do not.

4 Regression (40 marks)

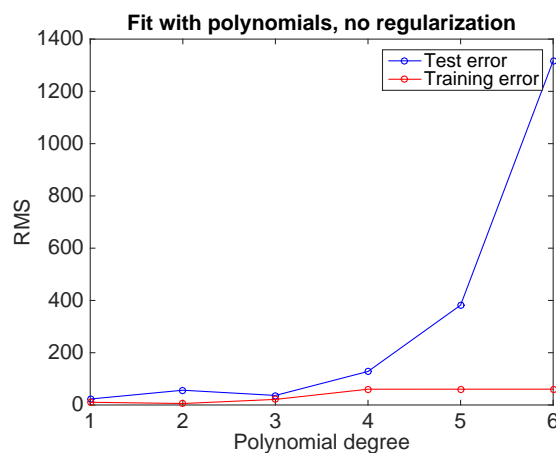
4.1 Getting started

1. (2 marks) *Which country had the highest child mortality rate in 1990? What was the rate?*
Niger, 313.7.

2. (2 marks) Which country had the highest child mortality rate in 2011? What was the rate?
Sierra Leone, 185.3.
3. (2 marks) Some countries are missing some features (see original `.xlsx/.csv` spreadsheet). How is this handled in the function `loadUnicefData.m`?
The median value of the other countries' values is used for this feature.

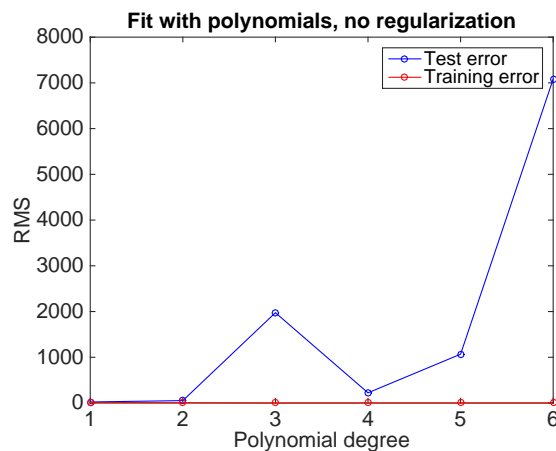
4.2 Polynomial Regression

1. Un-normalized data results in the following training and test errors.

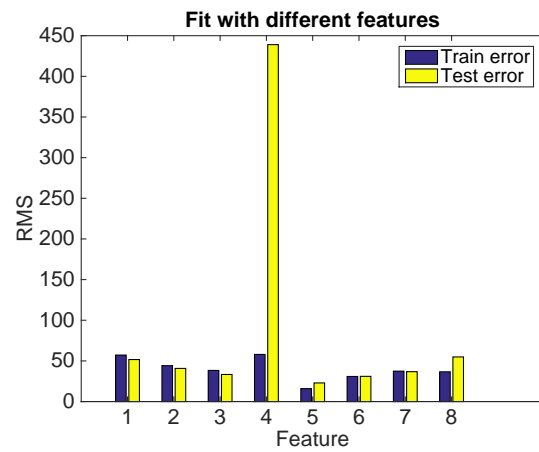


Note that training error actually increases with larger degree. This is due to numerical instabilities, due to large ranges in the values of inputs.

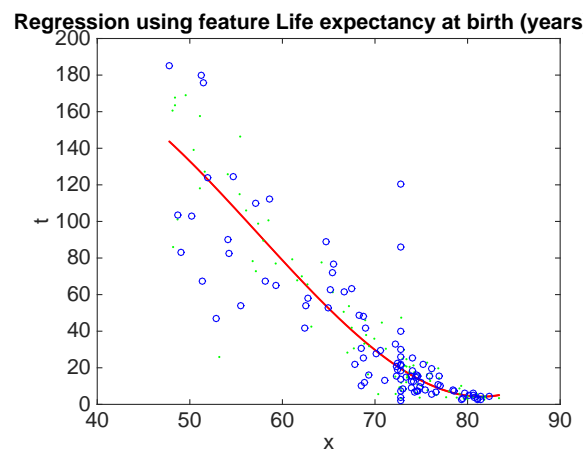
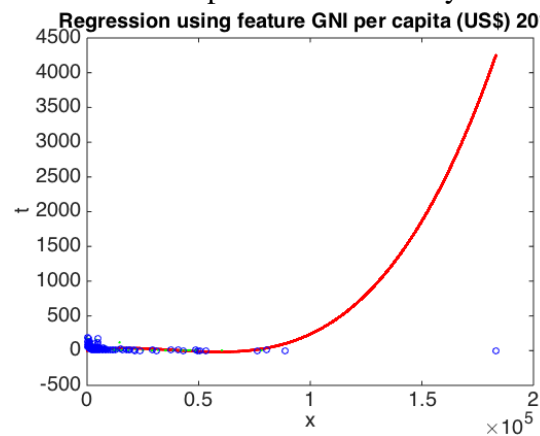
Results with normalized data are below.

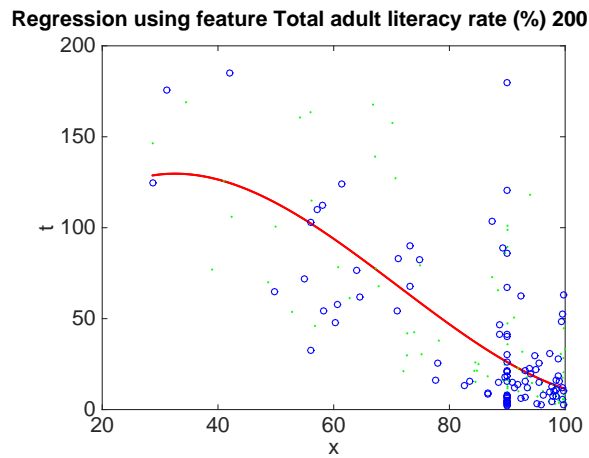


2. Single feature regression. Summary bar chart.



Fits for different features. Note the problems caused by outliers with large values of GNI.

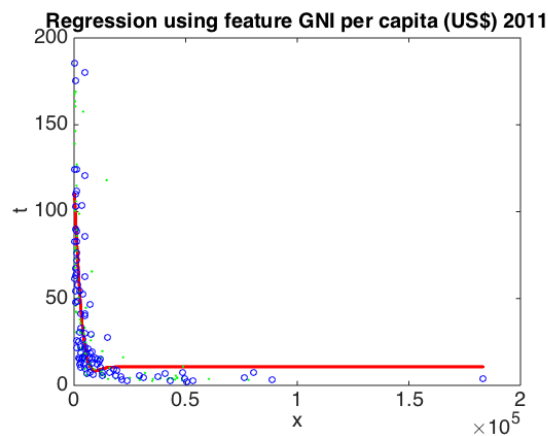




The training error for feature 11 (GNI per capita) is very high. To see what happened, produce plots of the training data points, learned polynomial, and test data points. The code `visualize_1d.m` may be useful.

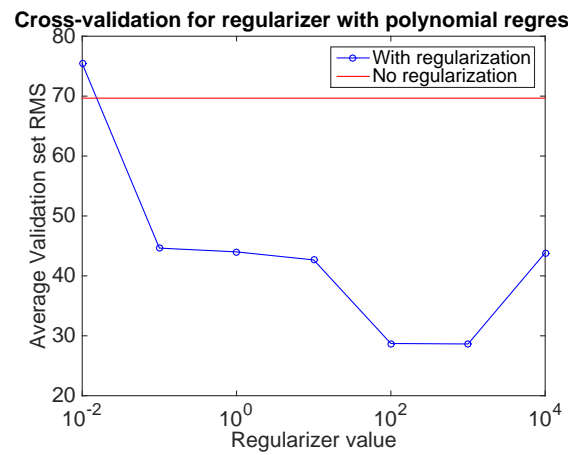
In your report, include plots of the fits for degree 3 polynomials for features 8 (Life expectancy), 11, 12 (Life expectancy), 13 (literacy).

4.3 Sigmoid Basis Functions



Training error is 27.4302, testing error is 30.6816.

4.4 Regularized Polynomial Regression



The value of cross-validation error for $\lambda = 1000$ is lowest, at 28.64.