

ROB313 Assignment #2

1. Assignment Objectives

The purpose of this assignment is to evaluate the effectiveness of Generalized Linear Models with various permutations in the regression estimation of the standard data sets provided in the class material. The different permutations of the algorithm each have their own function calls, and include parameter passing to allow for the optimization and parameter sweep for values such as the regularization constants or basis function values. These functions accept a training set with feature vectors (x) and corresponding labels (y), a testing/validation set of feature vectors, and the various parameters that can be modified and optimized. First, a non-kernelized version is tested against the Mauna Loa set, and an optimal regularization parameter is found. Next, a kernelized GLM, utilizing Cholesky factorization, is created and tested against the Rosenbrock set. The computational cost and translational properties of the kernel are evaluated. Finally, radial basis function is used as the GLM kernel, and is evaluated on all 3 regression data sets. In addition to the coding portion, two proofs are solved on paper, scanned and included to the report.

2. Algorithm Overview

The initial Generalized Linear Model builds on the linear model from the first assignment. First, the ϕ matrices for the training and testing sets are constructed using the selected basis functions: linear, quadratic, and of the form $x \sin(x)$ and $x \cos(x)$, which were chosen based on the shape of the graph. Then, the matrix undergoes the GLM algorithm, inverted using SVD from the numpy library.

A kernelized GLM was used to abstract away the feature space and put the computation in terms of inner products through the Gram Matrix K . In conjunction, a Cholesky factorization was utilized, and the runtime of this algorithm was compared to the previous. In addition, the kernel was analyzed in terms of translational invariance by creating parameterized plots of the kernel output values.

3. Results

3.1 Non Kernelized General Linear Model

Based on an initial analysis of the Mauna Loa set, it was concluded that the data approximates some combination of a sinusoid, a linear/quadratic function, and a function of the form $x \sin(x)$ & $x \cos(x)$. Thus, these functions were selected as the basis function for the GLM. In addition, the frequency of the sinusoidal components was analyzed. It was estimated to have a period of 0.055, when the training set was plotted with respect to its 1D feature vector. The lambda value parameter sweep was done using the validation set, while the final result was plotted with the testing set

Optimal Lambda Value: 5
Testing RMSE: 0.085866

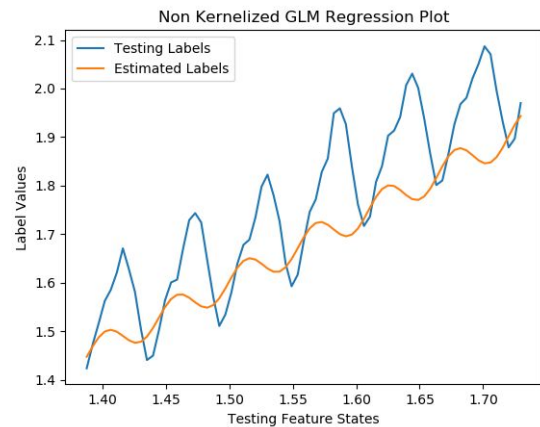


Fig 1: Result of the Mauna Loa regression

3.2 Kernelized General Linear Model

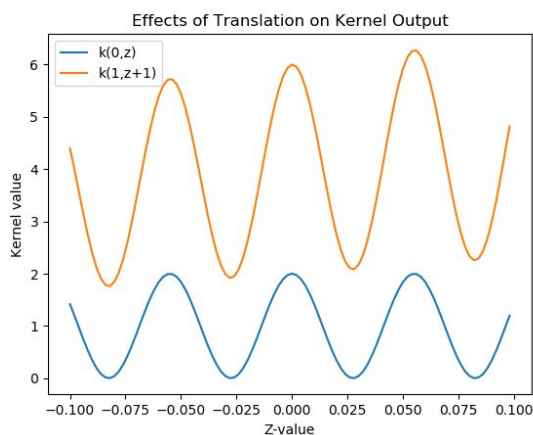


Fig 2: Kernel Output Analysis

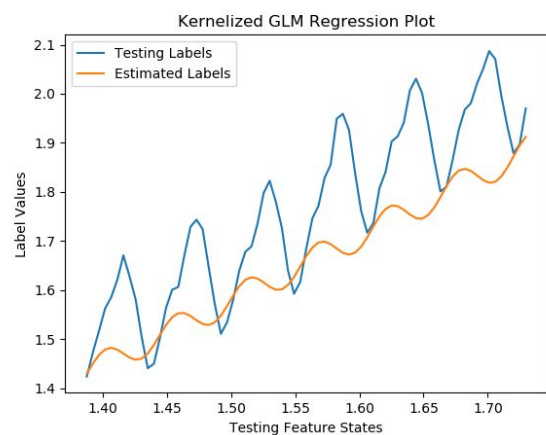


Fig 3: Result of the Mauna Loa regression

Next, the base algorithm was modified and its performance is reevaluated. The basis function space was collapsed using a kernelized implementation, and the weights were calculated using inner products in the dual representation using the Gram Matrix. The kernel was found by calculating the inner product of the basis function vectors for two arbitrary feature inputs. The kernel was then simplified by completing the square and scaling the linear inputs to create a single polynomial kernel term, and the double angle theorem was used to simplify the sinusoid function multiplications, yielding a final kernelized inner product expression:

$$k(x, z) = (1 + x^T z)^2 + (1 + x^T z) \cos(\omega (x^T z))$$

In addition to kernelization, the Singular Value Decomposition was replaced with a Cholesky factorization. In the end, the memory and runtime complexity were significantly poorer with this iteration of the algorithm. On the one hand, with a small feature space, kernelizing the algorithm will actually lead to calculations in a much larger vector space as $N \gg M$. For this reason as well, the N^3 term in the Cholesky's runtime complexity will dominate over the $11N(D+1)^3$ term in the SVD runtime.

Note: The algorithm had exactly the same RMSE error and optimal regression parameter as the previous algorithm

3.3 Radial Basis Function Kernel GLM

Finally, the kernelized version was modified, using a radial basis function Gaussian kernel rather than the previously derived kernel. In this case the kernel parameter as well as the regression parameter were swept, and the algorithm was applied to 2 regression sets and 1 classification set. The optimal parameters were found using the validation set, the RMSE value was calculated on the estimate for the testing set, and the results are summarized in the table below.

	Optimal λ	Optimal θ	Testing RMSE or Accuracy
Mauna Loa Data Set	0.001	1	0.124479
Rosenbrock Data Set	0.001	2	0.193240
Iris Data Set	1	0.5	0.870968

Fig 4: Table summarizing optimal parameters and associated success metric for 3 evaluated data sets

3.4 & 3.5 Proof Questions

Rob313 Assignment #2

Samuel Looper (100304418)

$$④ \quad \operatorname{argmin}_{\mathbf{w} \in \mathbb{R}^D} \underbrace{\|\mathbf{y} - \mathbf{X}\mathbf{w}\|_2^2 + \mathbf{w}^T \mathbf{\Gamma} \mathbf{w}}_{= \otimes} = \operatorname{argmin}$$

$$\otimes = (\mathbf{y} - \mathbf{X}\mathbf{w})^T (\mathbf{y} - \mathbf{X}\mathbf{w}) + \mathbf{w}^T \mathbf{\Gamma} \mathbf{w} = \mathbf{y}^T \mathbf{y} + \mathbf{X}^T \mathbf{w}^T \mathbf{X} \mathbf{w} - 2\mathbf{y}^T \mathbf{X} \mathbf{w} + \mathbf{w}^T \mathbf{\Gamma} \mathbf{w}$$

$$\frac{d}{d\mathbf{w}} \otimes = [(\mathbf{X}^T \mathbf{X})^T + (\mathbf{X}^T \mathbf{X})] \mathbf{w} - 2\mathbf{y}^T \mathbf{X} + (\mathbf{\Gamma}^T + \mathbf{\Gamma}) \mathbf{w} = 0$$

$$= 2\mathbf{X}^T \mathbf{X} \mathbf{w} - 2\mathbf{X}^T \mathbf{y} + 2\mathbf{\Gamma} \mathbf{w} = 0 = \mathbf{X}^T \mathbf{X} \mathbf{w} - \mathbf{X}^T \mathbf{y} + \mathbf{\Gamma} \mathbf{w} = 0$$

$$\therefore \operatorname{argmin}_{\mathbf{w}} \underbrace{\mathbf{X}^T (\mathbf{X} \mathbf{w} - \mathbf{y}) + \mathbf{\Gamma} \mathbf{w}}_{\text{minimizer}} = 0 \quad \text{where} \quad \mathbf{X}^T \mathbf{y} = (\mathbf{X}^T \mathbf{X} + \mathbf{\Gamma}) \mathbf{w}$$

$$\mathbf{w} = (\mathbf{X}^T \mathbf{X} + \mathbf{\Gamma})^{-1} \mathbf{X}^T \mathbf{y}$$

$$⑤ \quad \operatorname{argmin}_{\alpha \in \mathbb{R}^n} \underbrace{\sum_{i=1}^n (y^{(i)} - f(x^{(i)}, \alpha))^2 + \lambda \sum_{i=1}^n \alpha_i^2}_{\otimes} = \operatorname{argmin} \|\mathbf{y} - \hat{\mathbf{y}}\|_2^2 + \lambda \alpha^T \alpha$$

$$= \operatorname{argmin} (\mathbf{y} - \hat{\mathbf{y}})^T (\mathbf{y} - \hat{\mathbf{y}}) + \lambda \alpha^T \alpha$$

$$\hat{\mathbf{y}} = [f(x^{(1)}, \alpha) \dots f(x^{(n)}, \alpha)]^T = [\mathbf{k}(x^{(1)})^T \alpha \dots \mathbf{k}(x^{(n)})^T \alpha]^T$$

$$= [\mathbf{k}(x^{(1)}) \dots \mathbf{k}(x^{(n)})]^T \alpha = \mathbf{K}^T \alpha = \mathbf{K} \alpha$$

$$\therefore \text{we need to } \operatorname{argmin}_{\alpha \in \mathbb{R}^n} [\mathbf{y} - \mathbf{K} \alpha]^T [\mathbf{y} - \mathbf{K} \alpha] + \lambda \alpha^T \alpha$$

$$\otimes = \mathbf{y}^T \mathbf{y} + \alpha^T \mathbf{K}^T \mathbf{K} \alpha - 2\mathbf{y}^T \mathbf{K} \alpha + \lambda \alpha^T \alpha$$

$$\therefore \frac{d}{d\alpha} \otimes = 0 = 2\mathbf{K}^T \mathbf{K} \alpha - 2\mathbf{K}^T \mathbf{y} + 2\lambda \alpha$$

$$\mathbf{K} \mathbf{y} = (\mathbf{K}^T \mathbf{K} + \lambda \mathbf{I}) \alpha \quad \alpha = (\mathbf{K}^T \mathbf{K} + \lambda \mathbf{I})^{-1} \mathbf{K}^T \mathbf{y}$$

Where \mathbf{K} is calculated from training set, \mathbf{y} is from training set, λ can be found by k -fold cross validation & $(\mathbf{K}^T \mathbf{K} + \lambda \mathbf{I})^{-1}$ by SVD or QR factorization iff $\mathbf{K}^T \mathbf{K}$ is full rank

different to regularization w.r.t. \mathbf{w} since $\mathbf{w}^T \mathbf{w} = \alpha^T \mathbf{K} \alpha$ & the extra \mathbf{K} term modifies the equation