# Reading Minds is Easy, MRI-te?

1   **Shilpa Kumar**                    **Siddhartha S. Gorti**
2   Department of Computer Science    Department of Computer Science
3   University of Washington          University of Washington
4   Seattle, WA 98105                 Seattle, WA 98105
5   *shilpak@cs.uw.edu*               *sgorti3@cs.uw.edu*

6                                    **Abstract**

7       This project involves discovering how fMRI brain scans can be used to
8       predict what word a person is reading based off activation patterns in the
9       brain. The goal is to first learn 218 sparse linear models, each predicting a
10      semantic feature of the word based on an fMRI input. Using Coordinate
11      Descent and Stochastic Coordinate Descent for Lasso, as well Proximal
12      Gradient Descent with L1 penalty, we figured out the coefficients to predict
13      the value of a semantic feature. Given a new brain scan input, we are able
14      to build a 218-dimensional vector representing the values of semantic
15      features of the word read. Given two candidate words, we use binary
16      classification to choose the word whose semantic feature vector is closer to
17      the predicted one.

18

## 1   Overview of Methods

20  We used multiple methods to find the sparse linear model for each semantic feature, to
21  minimize overfitting. Lasso is a regression analysis method that uses $\lambda\|w\|1$ to penalize the
22  magnitude of the coefficients. Lasso allows coefficients to be 'pushed' to 0; this allows for
23  us to achieve sparse linear models without compromising some important features of the
24  model. Each different lasso method used was tested with multiple L1 penalty tuning
25  parameters, and across several semantic features.
26

### 1.1   Coordinate Descent for LASSO

28  Coordinate descent is used in cases where the gradient cannot be easily computed. Lasso
29  uses $\lambda\|w\|1$ as a penalty on the loss function, with $\lambda$ being the tuning parameter. The absolute
30  value function is not differentiable at $x = 0$, and a way to work around the complicated
31  gradient that results is to use coordinate descent instead of gradient descent.

32



33          Figure 1: General Algorithm for Coordinate Descent for LASSO[1]

---

[1]Murphy

34
## 1.2 Stochastic Coordinate Descent (SCD)

The coordinate descent algorithm involves going through all data points multiple times per iteration, which can be time-consuming for large datasets. To combat this problem, we utilized stochastic coordinate descent, which uses one random index associated with the voxel (one column of input matrix X) per iteration to minimize the loss function L(w, λ).

40

$$
\begin{aligned}
&\textbf{input} \ : X \in R^{n \times p}, Y \in R^n, \text{ and } \lambda \\
&\textbf{output: } \hat{w} = \arg\min_w L(w, \lambda) \\
&\text{Set } \tilde{X} = [X, -X]; \\
&\text{Initialize } \tilde{w}_j = 0 \text{ for } j = 1, \ldots, 2p; \\
&\textbf{while } not\ converged \ \textbf{do} \\
&\quad \text{Choose } j \text{ uniformly at random ;} \\
&\quad \tilde{w}_j \leftarrow \tilde{w}_j + \max\{-\tilde{w}_j, -\nabla_j \tilde{L}(\tilde{w}, \lambda)\} \\
&\textbf{end} \\
&\text{Set } w_j = \tilde{w}_j - \tilde{w}_{j+p} \text{ for } j = 1, \ldots, p.
\end{aligned}
$$

Figure 2: General Algorithm for Stochastic Coordinate Descent[2]

42
## 1.3 Proximal Gradient Descent (PGD)

To compare our previous algorithms with a method that does not use a version of coordinate descent, we used proximal gradient descent with L1 penalty. Proximal methods are useful for large-scale problems where other methods might be too slow.[3]
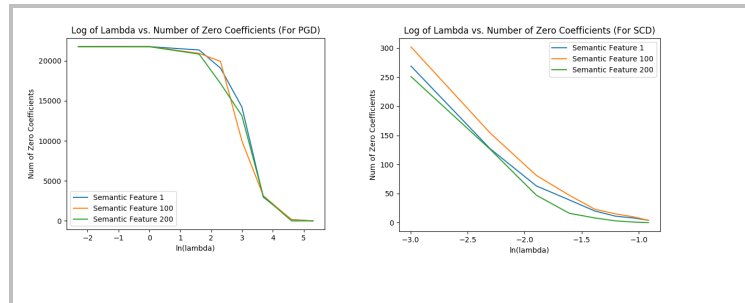
47

$$
\begin{aligned}
\boldsymbol{\theta}_{k+1} &= \operatorname*{argmin}_{\mathbf{z}} \left[ t_k R(\mathbf{z}) + \frac{1}{2} \|\mathbf{z} - \mathbf{u}_k\|_2^2 \right] = \operatorname{prox}_{t_k R}(\mathbf{u}_k) \\
\mathbf{u}_k &= \boldsymbol{\theta}_k - t_k \mathbf{g}_k \\
\mathbf{g}_k &= \nabla L(\boldsymbol{\theta}_k)
\end{aligned}
$$

Figure 2: General Algorithm for Proximal Gradient Descent

49
## 2 Choosing Tuning Parameters

To help us solve overfitting, we introduce the concept of regularization which helps prevent the generation of overfit models. In building these models, we used two different methods to select our tuning parameters, k-fold cross validation and test set validation. While k-fold is more comprehensive, using the test set to find the optimal lambda value was faster and because this computation had to be done for each model, we settled on using the test set to choose our tuning parameter.

57
## 3 Results

59



---

[2] Homework 3, CSE 599
[3] Murphy

60  Figure 4: Log of tuning parameter vs number of nonzero coefficients
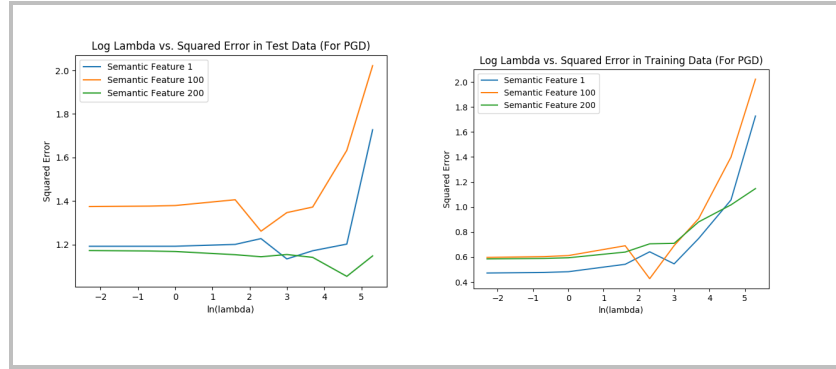
61



62  Figure 5: Log of tuning parameter vs test and training error for PGD
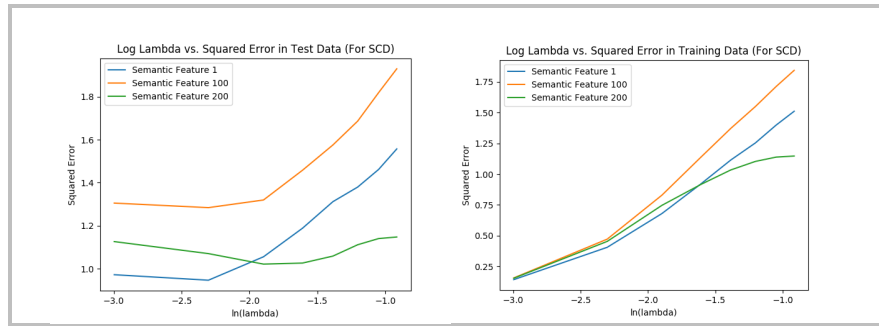
63

64



65  Figure 6: Log of tuning parameter vs test and training error for SCD

66

## 4    Comparison

Coordinate descent for Lasso, on our dataset, took an unreasonably long amount of time. For predicting the model of just one semantic feature, the algorithm took upwards of 2.5 hours. We used coordinate descent simply as a baseline comparison, but did not calculate training error or test error for various lambdas because of the length of time necessary for 218 models and various lambdas.

Both our stochastic coordinate descent (SCD) and proximal gradient descent (PGD) followed the correct trends in training error and test error. Training error exponentially increased for various semantic features when lambda was increased, and test error decreased up to a certain lambda, and increased again. However, the graphs for PGD seemed less smooth than those of SCD, and the training error graph for PGD had a small dip in error near $\ln(\lambda) = 2$. In addition, PGD took more time to complete than SCD for the same number of iterations. For the number of nonzero coefficients $\lambda$, convergence took much longer for PGD.

## 5    Future Work

For the next phase of our project we will be working towards feeding in brain scan input into all these linear models to generate a 218 dimensional vector. Then given a choice between two words, we will be comparing the generated vector (representing the semantic features) between the semantic features vector for both words, picking the word whose semantic features vector most closely resembles the generated one, using 1-NN binary classification.

88   **References**

89    [1] Bradley, Joseph K., Aapo Kyrola, Danny Bickson, and Carlos Geustrin. *Parallel Coordinate*
90   *Descent for L1-Regularized Loss Minimization*. 2011. MS. Carnegie Mellon University, Pittsburgh.

91   [2] Caramanis, and Sanghavi. "Lecture 4 — September 11." EE 381V: Large Scale Optimization.
92   UTexas. Lecture.

93   [3] "Homework 3 -- Midterm." CSE 599 / Stat 592: Machine Learning (Statistics) for Big Data.
94   University of Washington Computer Science. Homework Assignment.

95   [4] Murphy, Kevin P. Machine Learning: A Probabilistic Perspective. Cambridge: The MIT Press,
96   2012. Print.

97   [5] Shalev-Schwartz, Shai, and Ambuj Tewari. *Stochastic Methods for L1 Regularized Loss*
98   *Minimization*. 2011. MS. Toyota Technological Institute at Chicago, Chicago.