
ECE 273 Project: Coordinate Descent, Dykstra's and LASSO

Sanjeev Anthia Ganesh
A53308351
santhiag@eng.ucsd.edu

Abstract

1 This project aims to study Dykstra's projection algorithm and block co-ordinate de-
2 scent to solve certain types of optimization problems. It also studies the application
3 of these algorithms to solve a problem called LASSO.

4 1 Co-ordinate Descent Algorithm

5 Input: We are given a convex differentiable function

$$f(x) : R^n \rightarrow R$$

6 whose minimizer needs to be found. Let's assume we have found the global minimizer of this convex
7 function. This statement translates to the fact that the minimizer is achieved along all the co-ordinate
8 axes. In the Co-ordinate Algorithm, we shall explore the fact about the types of functions that
9 Co-ordinate Descent Algorithm works on as well as convergence criteria need for these functions.

10 1.1 Algorithm

11 In Block Coordinate Descent algorithm, we first choose a point

$$x^0 = (x_1^0, x_2^0, \dots, x_n^0) \in \text{dom}(f)$$

12 . In each iteration, we choose a single coordinate or a block of coordinates and minimize the function
13 along each of these coordinates. At iteration $r+1, r \geq 0$, given, $x^r = (x_1^r, x_2^r, \dots, x_n^r), x^r \in \text{dom}(f)$
14 choose an index $s = 1, 2, 3, \dots, n$ and compute a new iterate

$$x^{r+1} = (x_1^{r+1}, x_2^{r+1}, \dots, x_n^{r+1}) \in \text{dom}(f)$$

15 that satisfies the following equation,

$$x_s^{r+1} \in \arg \max_{x_s} f(x_1^r, x_2^r, \dots, x_s, x_{s+1}^r, \dots, x_n^r)$$

16

$$x_j^{r+1} = x_j^r, \forall j \neq s$$

17 This algorithm is repeated across all the indexes and multiple iterations till convergence to the global
18 minimizer is achieved.

19 The regularized regression problem can be stated as follows.

$$\min_{\beta \in R^p} (1/2) * ||y - X\beta||_2^2 + \sum_{k=1}^N h_i(w_i)$$

where $\beta \in R^{p \times 1}$ is a feature vector

$X \in R^{N \times p}$ is a transformation matrix

$y \in R^{N \times 1}$ is a observation vector.

(1)

20 On an overall level, the co-ordinate descent algorithm for a regularized regression problem given
 21 above can be represented as

$$\beta_i^{(k)} = \arg \min_{\beta \in R^p} (1/2) * ||y - \sum_{j < i} X_j \beta_j - \sum_{j > i} X_j \beta_j||_2^2 + h_i(w_i) \quad (2)$$

$i = 1, 2, \dots, d$

22 1.2 Compatible function types and Convergence of the Algorithm

23 The characterization of the objective function for convergence has been explored by Tseng and we
 24 shall use Lemma 3.1 and Theorem 4.1a of his paper to prove convergence of the Co-ordinate Descent
 25 Algorithm.

26 For Block Coordinate Descent algorithm to be compatible with the function whose minimizer needs
 27 to be found, the function needs to have certain properties. We see that if the function is convex
 28 differentiable, then co-ordinate descent would tend to achieve a minimum. If the function is convex but
 29 not differentiable, there is a question whether block-coordinate descent could achieve the minimizer
 30 since the function is not smooth. Tseng, 2001 has provided several conditions needed by a function
 31 that is convex and non-differentiable to compatible with the Co-ordinate Descent Algorithm. Consider
 32 a function

$$f(x_1, x_2, \dots, x_n) = f_0(x_1, x_2, \dots, x_n) + \sum_{k=1}^N f_k(x_k)$$

33 Here, in this function, we separated the differentiable part f_0 from the non-differentiable part f_k .
 34 Assume that the level set $X^0 = \{x : f(x) \leq f(x^0)\}$ is compact and f is continuous on X^0 where x^0
 35 is the initialization point taken. If this is the case, then the function f is continuous and bounded.
 36 Moreover, if the function $f(x_1, x_2, \dots, x_n)$ is pseudoconvex in (x_k, x_i) for every $i, k \in [1, 2, \dots, N]$, and
 37 if the function is regular at every $x \in X^0$ then every cluster point x^r is a stationary point of f . A
 38 function f is regular when the lower directional derivative $f'(x, d) \geq 0 \forall d \in d_1, \dots, d_k, \dots, d_N$. For a
 39 function f to be regular, the $\text{dom}(f_0)$ should be open and must be Gateaux differentiable. Here, we
 40 imply the condition for regularity by taking into account of the smoothness of f_0 . Even if the other
 41 separable non-differentiable functions $f_k, k = 1, 2, \dots, N$ are not smooth, the claim for regularity holds
 42 valid.

43 1.3 Examples on convergence

44 For example, consider the function on Fig. 1. Here, the function is convex and differentiable, therefore
 45 at the global minimizer, I know that the gradient of the function with respect to all the coordinates
 46 would be zero. So, Coordinate descent computes a point that minimizes the function along each
 47 coordinate axis. Therefore, the coordinate descent Algorithm shall converge to the global minima.

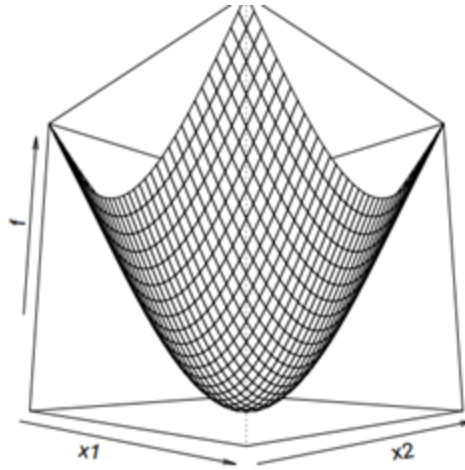


Figure 1: A convex differentiable function

48 Consider the function represented on the Fig. 2. Here, the function is convex , non-smooth and my
 49 minimizer is denoted by the blue dot. Let's say that my algorithm is at the point of intersection of red
 50 dotted lines. According to Coordinate Descent Algorithm, in an iteration, I would try to minimize
 51 along one of the coordinate axes. But here, when I am stuck at a non-stationary point minimizing the
 52 function along either of the axis returns the same stationary point which is obviously not a minimizer.
 53 So, the Coordinate Descent Algorithm gets stuck at this non-stationary point. One solution would
 54 be to minimize along both the axis simultaneously to find out the minimizer. This means that the
 55 non-differentiable part of the convex function is not separable. However, the function represented in
 56 3 has a separable non-differentiable part. Coordinate descent algorithm works in the function to find
 57 the minimization point

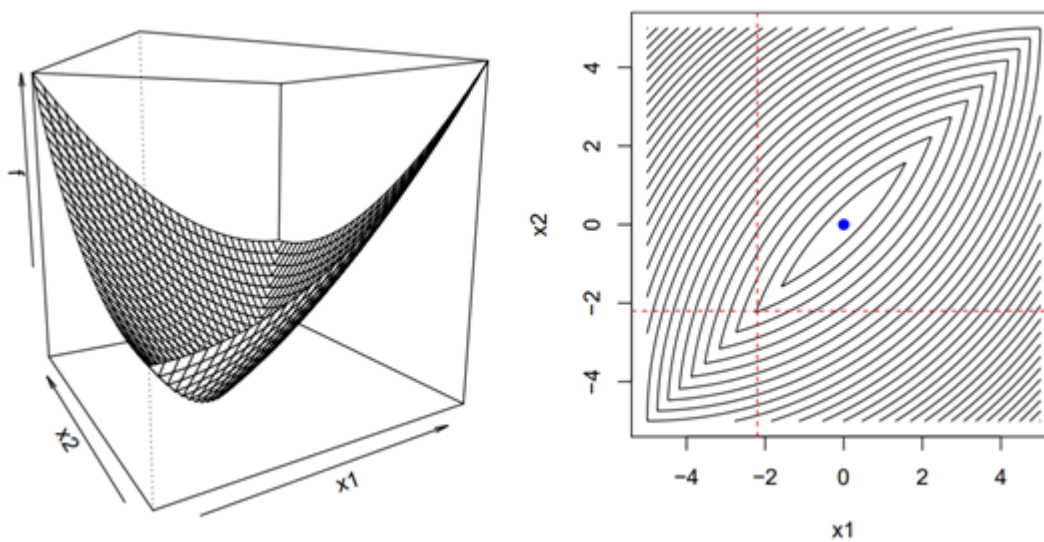


Figure 2: A convex non-differentiable function

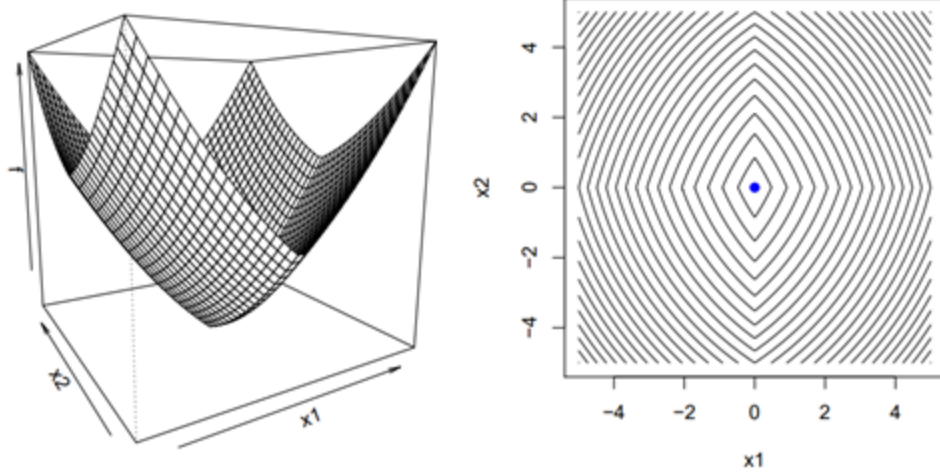


Figure 3: A convex separable non-differentiable function

2 Dykstra's Projection Algorithm

Dykstra's Projection Algorithm is a method for solving the best approximation problem given in Eqn. 3 where $C_1 \cap C_2 \dots \cap C_n$ denote the intersection of convex sets. Initially one is unaware of the sets of points contained in the intersection of these convex sets. So, there is an iterative approach followed where the point x is projected from one convex set to another convex set till the minimum is achieved. This algorithm differs from the Alternating Projections in a way that this algorithm tracks the residual after each projection onto $C_1 \cap C_2 \dots \cap C_n$. Meanwhile, the Alternating projection's provides a point from the feasible set of points lying in the intersection, the Dykstra's Algorithm not only provides a point from the feasible set but it also provides the solution in the intersection of convex sets which is near to the point in question.

$$\min_x \|y - x\|_2^2 \text{ given } x \in C_1 \cap C_2 \dots \cap C_n \quad (3)$$

2.1 Algorithm

First, we initialize $u(0)^d = y, z(0)^1 = z(0)^d = 0$, and then repeat, for $k = 1, 2, 3, \dots$

$$\begin{aligned} u_0^{(k)} &= u_d^{(k-1)} \\ u_i^{(k)} &= P_{C_i}(u_{i-1}^{(k)} + z_i^{(k-1)}) \\ z_i^{(k)} &= u_{i-1}^{(k)} + z_i^{(k-1)} - u_i^{(k)} \\ &\forall i = 1, 2, \dots, d \end{aligned} \quad (4)$$

where $P_{C_i}(x)$ denotes the convex projection of x onto the closed convex set C_i . Now, the residuals along with the projections are tracked in a cyclic order until the best approximation problem reaches its minimum.

3 LASSO

3.1 Primal Solution of LASSO

Consider the problem of Regularized regression problem in Eqn. 5 where the penalty is the L1 norm of the feature vector beta. Then the problem becomes a LASSO problem. LASSO is widely used because of the sparse solutions they return. The goal of the penalty function here, with respect to the

hyperparameter λ , is to minimize the L1 norm of the solution. So, trying to minimize the L1 norm, one encounters a sparse vectors. L1 norm has this capability because of its diamond shape.

$$\min_{\beta \in R^p} (1/2) * ||y - X\beta||_2^2 + \lambda \sum_{k=1}^N ||\beta||_1$$

where $\beta \in R^{p \times 1}$ is a feature vector

$X \in R^{N \times p}$ is a transformation matrix

$y \in R^{N \times 1}$ is a observation vector.

3.2 Dual of LASSO

The dual of the LASSO can be represented as the best approximation problem where the original objective function is represented as follows:

$$\max_{\theta} (1/2) * (||y||_2^2 - ||y - \theta||_2^2) \text{ subject to } ||X_i^T \theta|| \leq \lambda$$

where i varies from 1 to p .

This problem can be converted since I know that $||y||_2^2$ is a constant term and the negative sign can be treated as minimize the objective function.

$$\min_{\theta} (1/2) * (||y - \theta||_2^2) \text{ subject to } ||X_i^T \theta|| \leq \lambda$$

where i varies from 1 to p .

3.3 Relationship between primal and dual solutions

Suppose that $h_i, i = 1, \dots, n$ are seminorms, which we can express in the general form $h_i(v) = \max_{D_i} \langle v, d_i \rangle$, where $D_i \subseteq R^{p_i}$ is a closed, convex set containing 0, for $i = 1, \dots, d$. Suppose also that $C_i = (X_i^T)^{-1}(D_i) = \{v \in R^n : X_i^T v \in D_i\}$, the inverse image of D_i under the linear mapping X_i^T , for $i = 1, \dots, d$.

Equation wise, the abover statement could be written in the following form where we assume that $X_i \in R^{n \times p_i}$ is a full column rank matrix and $h_i(v) = \max_{d \in D_i} \langle d, v \rangle$ for a closed convex set containing 0. Then for $C_i = (X_i^T)^{-1}(D_i) \subseteq R^n$ and any $b \in R^n$,

$$\hat{w}_i = \arg \min_{w_i \in R^{p_i}} (1/2) * ||b - X_i w_i||_2^2 + h_i(w_i)$$

$$\Leftrightarrow \hat{w}_i = (Id - P_{C_i})(b)$$

The best approximation problem and regularized regression problem are dual to each other which satisfy $\hat{u} = y - X\hat{w}$. The Dykstra's algorithm and Coordinate Descent Algorithm are equivalent and satisfy each other at all iterations $k = 1, 2, 3, \dots$ via duality. From Eqn. 2 and 4, these iterations can be equated as

$$z_i^{(k)} = X_i \beta_i^{(k)} \text{ and } u_i^{(k)} = y - \sum_{j < i} X_j \beta_j^{(k-1)} - \sum_{j > i} X_j \beta_j^{(k-1)}$$

where $i = 1, 2, 3, \dots, d$

3.4 Convergence on LASSO

From Eqn. 5, we know that the sequences defined by the Coordinate Descent Algorithm are defined and bounded. If the function f is regular, then every cluster point from each of the iteration is a stationary point. Now, for f to be a regular function, the $\text{dom}(f_0)$ should be open and f_0 should be gateaux differentiable. The f_0 in this problem is differentiable and the $\text{dom}(f_0) \in R^n$ is open. Therefore, the Theorem 4.1a proposed by Tseng is satisfied and the LASSO converges to a co-ordinate wise minimum point.

104 3.5 Solution to primal problem using Coordinate Descent

105 For the Co-ordinate descent algorithm to provide stationary points at each iteration, we note that the
 106 non-differentiable L1 norm is separable from the differentiable convex function. Therefore, we can
 107 perform minimization with respect to each of the coordinate axis.

108 Minimizing over β_i , with $\beta_j, j \neq i$:

$$0 = X_i^T(X_i\beta_i) + X_i^T(X_{-i}\beta_{-i} - y) + \lambda s_i$$

109 where $s_i \in \partial|\beta_i|$. The solution is given by soft-thresholding

110 where

$$\beta_i = S_{\lambda/\|X_i\|_2^2}(X_i^T(y - X_{-i}\beta_{-i})/(X_i^T * X_i)) \quad (10)$$

111 The algorithm here is repeated for $i = 1, 2, \dots, p, 1, 2, \dots$

112 The soft-thresholding operator is denoted by

$$\begin{aligned} S_\lambda(x) &= x - \lambda \text{ if } x > \lambda \\ &= x + \lambda \text{ if } x < -\lambda \\ &= 0 \text{ if } x = 0 \end{aligned} \quad (11)$$

113 The soft-thresholding operator is invoked here because of the non-differentiability of the L1-norm at
 114 β equals to zero. The sub-differential concept where the derivative of L1 norm can take any value
 115 between -1 to 1 at the non-stationary point of L1 norm is assumed and respectively used as the
 116 soft-thresholding operator.

117 3.6 Solution to dual problem using Dykstra's Projection Algorithm

118 The dual problem is solved by Dykstra's projection algorithm from Eqn. 4 where each of
 119 the closed convex set is denoted by $\|X_i^T\|_1 < \lambda$ where $i = 1, 2, \dots, p$ (number of features). The
 120 projection of a point towards a convex set is solved by a convex solver. The relationship obtained
 121 between the Coordinate Descent Algorithm and the Dykstra's Algorithm is that

$$\theta = y - X\beta$$

122 where θ is obtained from the Dykstra's Algorithm and β is obtained from the Coordinate Descent.

123 4 Numerical Results

124 To evaluate the performance of the Algorithm, three separate cases are tested - Variation with respect
 125 to the number of observations (N), variation with respect to the sparsity (s), Variation with respect
 126 to the number of features (p). The test parameters here are - mean squared error and the support
 127 recovery performance. Mean square error denotes the square of the difference between the original
 128 feature vector β and the predicted $\hat{\beta}$. Support recovery performance denotes the accuracy with
 129 which the predicted $\hat{\beta}$ vector has predicted the features of least importance (the features with zero
 130 weights).

131 4.1 Variation of N

132 In this experiment, the Value of p is fixed at 100 and the number of non-zero elements in the sparse
 133 vector s is fixed at 4. The number of observations is varied from 20 to 200 in steps of 10. The error
 134 is relatively small considering which suggests that for a number of observations as small as 20, the
 135 sparse vector is recovered sufficiently as seen in Fig. 4. Therefore, in the next experiment, N was
 136 varied from 2 to 24 in steps of 2 as seen in the section 4.1.2. Here, we see that for a number of
 137 observations as small as the number of sparse elements, the error is high. Relative to the error is the
 138 support recovery performance. Here, the predicted feature vector should predict the features of least
 139 importance. As we see from Fig. 5 and Fig. 7, the support recovery performance is greater than 98
 140 percent with 100 percent recovery for N values larger than 20.

141 **4.1.1 Larger Values of N**

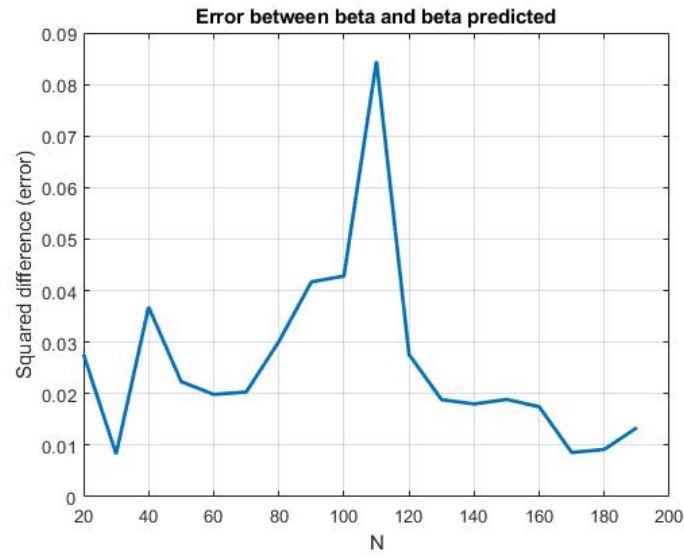


Figure 4: L2 error between the predicted feature vector and original feature vector

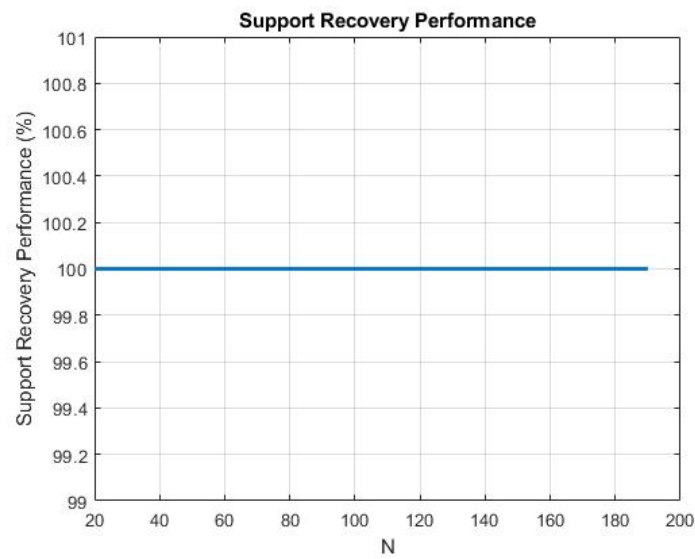


Figure 5: Support Recovery Performance for larger N

142 4.1.2 Smaller values of N

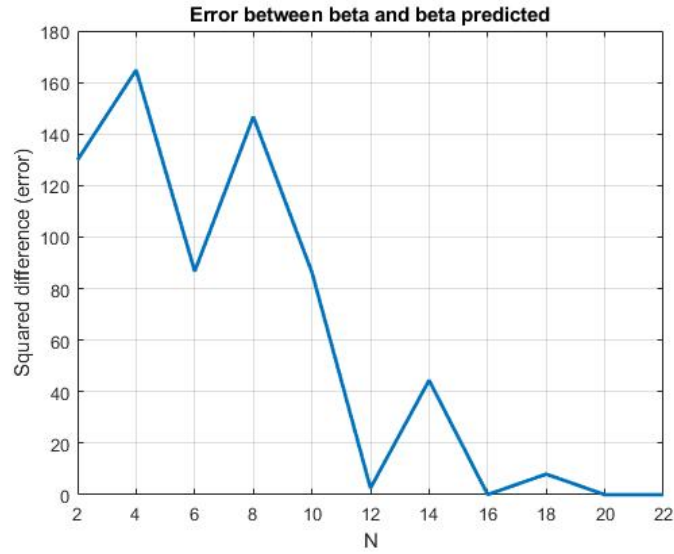


Figure 6: Mean squared error between the predicted feature vector and original feature vector

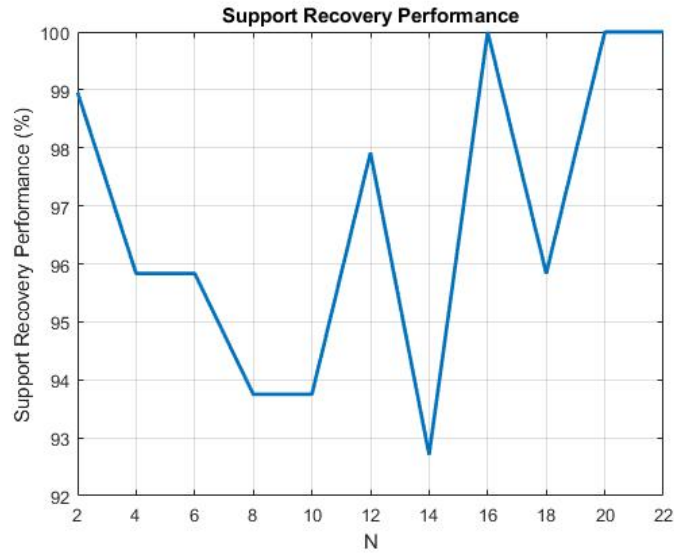


Figure 7: Support Recovery Performance for smaller N

143 4.1.3 Runtime Performance

144 One of the information that I can get is how much Runtime is required by the algorithm for processing
 145 a large number of observations. The CVXPY with highly optimized python module performs 1000
 146 times better than the current implementation of Co-ordinate Descent Algorithm. The solver used by
 147 the CVXPY module is QSQP solver. 8 shows the runtime of Co-ordinate Descent Algorithm across
 148 various values of N and 9 shows the runtime of equivalent problem executed by CVXPY module.

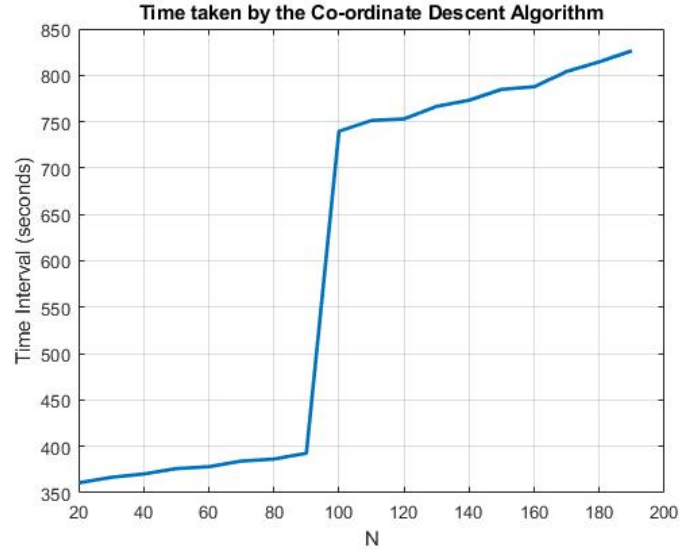


Figure 8: Runtime evaluation of Co-ordinate Descent Algorithm

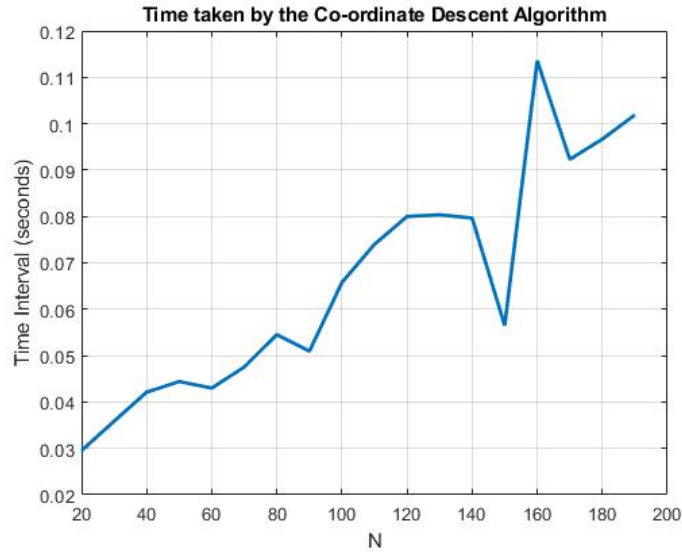


Figure 9: Runtime evaluation of CVXPY module

149 4.2 Variation of sparsity N

150 In this experiment, N is fixed at 50 and p is fixed at 100. The number of non-zero elements is varied
 151 from 5 to 90 in the steps of 5. The whole point in this experiment is to figure out, given the setting, the
 152 features and the observations, we are trying to get an idea on how the support recovery performance
 153 is based on the sparsity.

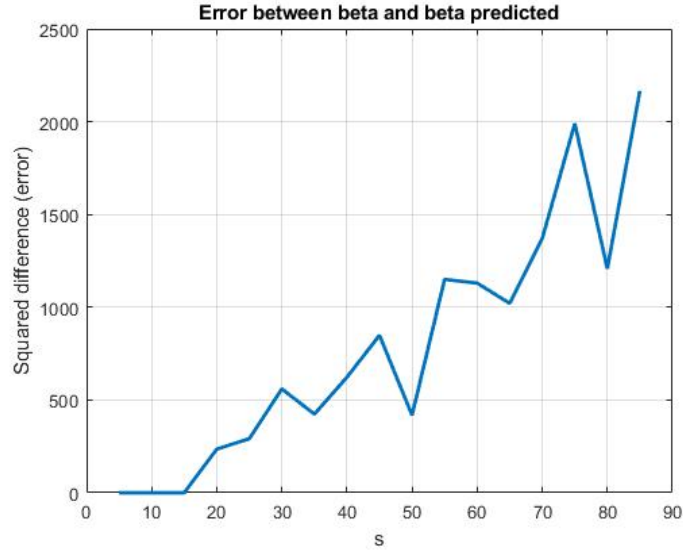


Figure 10: L2 Error between beta and beta predicted

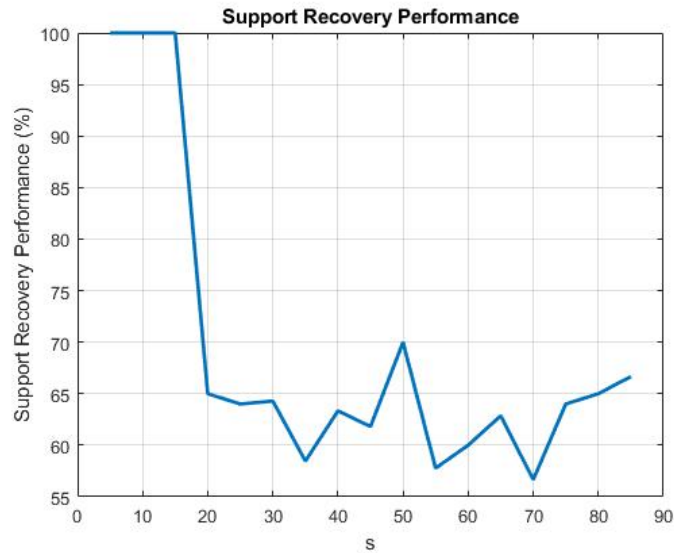


Figure 11: Support Recovery Performance in terms of sparsity

154 4.3 Variation of number of features p

155 In this experiment, p is varied from 100 to 300 whereas N is fixed at 50 and the number of non-
 156 elements in the sparse feature vector is fixed at 8. The support recovery performance as p increases is
 157 around 90 % whereas the error performance is fixed across all the elements over p and no trend is
 158 seemed to be observed. Fig. 12 shows the L2 error wrt the number of features and Fig. 13 shows the
 159 support recovery performance.

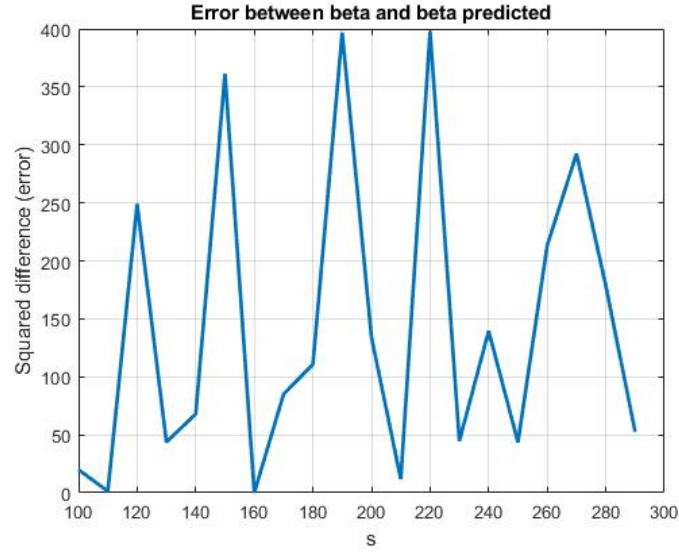


Figure 12: L2 Error between beta and beta predicted in terms of number of features

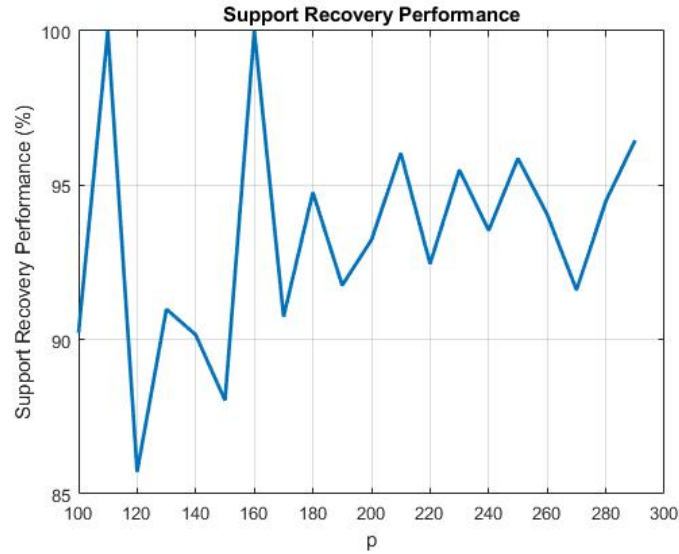


Figure 13: Support Recovery Performance in terms of number of features

160 4.4 Variation of N and s where s is varied linearly with N

161 In this experiment, the value of N is varied from 20 to 200 as well as the sparsity simultaneously
 162 from 4 to 40. The trend we see in this experiment is that with increasing number of observations, we
 163 find the the number of non-zero elements as well increases. This can be seen in the support recovery
 164 performance where we have mapped N to a sparsity such that for 20 number of observations, my
 165 sparsity feature vector would have 4 sparse elements and at 200 number of observations, my sparse
 166 vector would have 40 non-zero elements. It is seen that even with the increase of number of non-zero
 167 elements, my support recovery performance doesn't reduce with increase in number of observations
 168 as seen in Fig. 15.

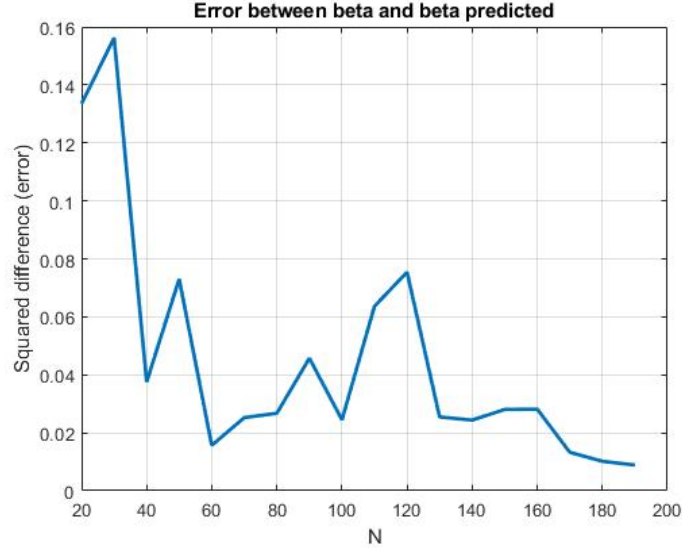


Figure 14: L2 Error between beta and beta predicted in terms of number of features

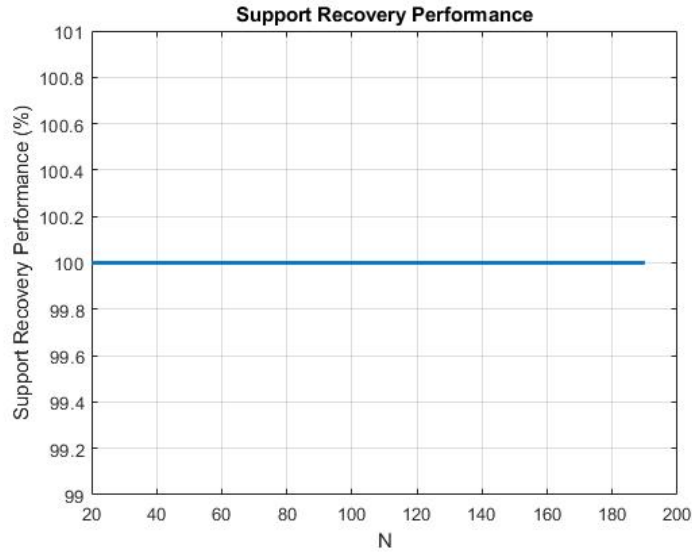


Figure 15: Support Recovery Performance in terms of number of observations

5 Applications to Co-ordinate Descent Algorithm

One of the areas where Co-ordinate Descent could be applied in wireless communications are computing channel capacity and Antenna Array synthesis. In this project, a basic Antenna array synthesis using Co-ordinate Descent is observed. So, in the Antenna Array Synthesis problem, we are given some constraints on the radiation pattern and number of antennas. We need to figure out the complex weights associated with each of the Antenna element to achieve the specification put forth by the radiation pattern. Now, the specification can be location of nulls, the level of side lobes and the direction of main lobes. So, in this project, a simple linear array synthesis based on Dolph-Chebyshev array weights is tried. The constraint here is that the number of Antenna Elements is 8 and the minimum side lobe level achieved is 30dB lesser than the main lobe. Fig. .16 shows such radiation pattern.

180 Now, I have radiation pattern with respect to θ from -90 to 90 degrees. Given the condition of the
 181 number of Antenna Elements and spacing between the elements, I would be able to find out the
 182 transformation matrix H which translates the weight vector to the radiation pattern. Each element in
 183 the matrix H is given by

$$H(\theta, n) = \exp(1i * (2 * \pi * \text{spacing} / \lambda) * \sin(\theta) * n)$$

184 where λ is the resonant wavelength of the Antenna, n is the antenna element and θ is the
 185 angle in radians.

186 We can see that H is a full row rank matrix. Therefore, the problem can be treated as a least squares
 187 problem where the Coordinate Descent iterations are same as LASSO iterations except that λ
 188 is zero here and all the elements are in complex space. For as few iterations as λ , the weight
 189 vector predicted by the Coordinate Descent Iterations matched the radiation pattern in the constraints.

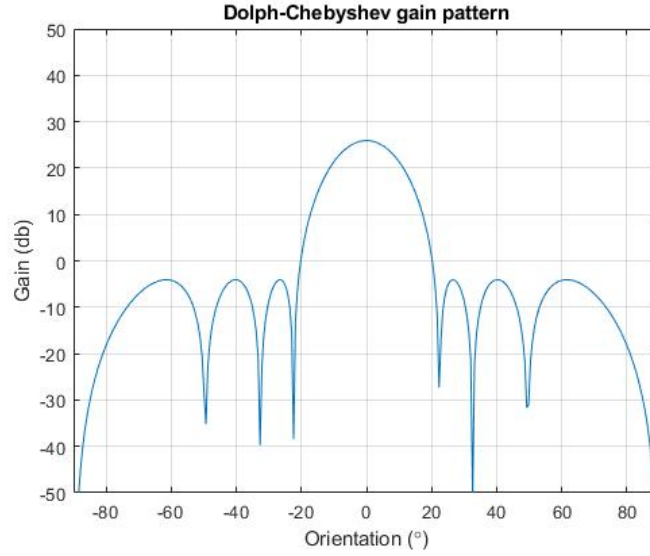


Figure 16: Dolph-Chebyshev weights based radiation pattern of Linear Array

190 References

- 191 [1] [https://stats.stackexchange.com/questions/123672/coordinate-descent-soft-thresholding-update-operator-for-](https://stats.stackexchange.com/questions/123672/coordinate-descent-soft-thresholding-update-operator-for-lasso)
 192 [lasso](https://stats.stackexchange.com/questions/123672/coordinate-descent-soft-thresholding-update-operator-for-lasso)
- 193 [2] <https://www.cs.cmu.edu/~ggordon/10725-F12/slides/25-coord-desc.pdf>
- 194 [3] Tibshirani, R. J. (2017). Dykstra's algorithm, admm, and coordinate descent: Connections, insights, and
 195 extensions.
- 196 [4] Tseng, P. (2001). Convergence of a block coordinate descent method for non differentiable minimization.
- 197 [5] <https://blog.mlreview.com/l1-norm-regularization-and-sparsity-explained-for-dummies-5b0e4be3938a>
- 198 [6] Wang, Y. , He, X. , Wang, J. , Berezin, S. and Mathis, W. (2015) Antenna Array Pattern Synthesis
 199 via Coordinate Descent Method. Journal of Electromagnetic Analysis and Applications, 7, 168-177. doi:
 200 10.4236/jemaa.2015.75018.