

High-Dim Project

name1

Columbia University
uni1@columbia.edu

name2

Columbia University
uni2@columbia.edu

name3

Columbia University
uni3@columbia.edu

1 Introduction

2 Background

2.1 Multiclass Classification with Group Lasso

The task of multiclass classification involves the prediction of a class label l where the number of possible labels is $k > 2$. More often than not, the original problem is transformed into k binary classification problems, i.e. 1-vs.-all classification and positive prediction with the highest confidence is selected as the label. This approach has the disadvantage of having to train k different models.

An alternative formulation, direct multiclass classification, tackles this problem directly by solving the following argmax problem:

$$y_i = \operatorname{argmax}_c W_{:c}^T x_i$$

where $W \in \mathcal{R}^{p \times k}$ is a weight matrix, with W_{ij} corresponding to the i -th feature of class j . In this paper, we refer to features as elements in the instance data x . A feature in x is associated with k weights in W , one for each class.

The decision function above suggests a max-margin style loss function. More specifically, we use the squared hinge loss:

$$l(W) = \sum_{i=1}^n \sum_{r \neq y_i}^k \max(1 - (W_{:y_i}^T x_i - W_{:r}^T x_i), 0)^2$$

The minimization of l directly will lead to a minimizer W^* that is dense. Sparse solutions are often explicitly sought, with model compactness leading to fast prediction at test time. In order to obtain a sparse W^* , a regularization term $r(W)$ is often applied, yielding the objective function:

$$\min_W l(W) + r(W).$$

Many choices are available for the regularizer r . In (ref ???), they use the group lasso, where each row in W is a group. The associated regularizer then is $r(W) = \lambda \sum_j \|W_{j:}\|_2$ where λ is a parameter that adjusts the strength of the regularization. This has the effect of producing a few rows of non-zero values in W ; since each row corresponds to an individual feature, the optimal sparse W^* yields a fast-evaluating decision function, i.e. most features are ignored at test time.

To minimize this multiclass classification group lasso objective, ??? use coordiante descent, iteratively solving a sub-problem with respect to a single group. Figure ? shows a general outline of algorithm that involves computing the partial gradient with respect to the current group j , the prox operator of the L2 norm, and a final line search to identify an appropriate step size for the current update.

```
for  $i \leftarrow 1, \dots, \max \text{iters}$  do
  for  $j \leftarrow 1, \dots, p$  do
    Compute gradient  $l'(W)_{j:}$ 
    Choose  $\mathcal{L}_j$ 
    Compute
       $V_j = W_{j:} - \frac{1}{\mathcal{L}_j} l'(W)_{j:}$ 
       $W_{j:}^* = \operatorname{Prox}_{\frac{\lambda}{\mathcal{L}_j} \|\cdot\|_2}(V_j)$ 
       $\delta = W_{j:}^* - W_{j:}$ 
    Choose  $\alpha$ 
     $W_{j:} \leftarrow W_{j:} + \alpha \delta$ 
  end
end
```

Efficient computation of this objective is possible by storing current loss for each data point. Let A be an $n \times k$ matrix where the i, r -th element corresponds to $(1 - (W_{:y_i}^T x_i - W_{:r}^T x_i))$. The gradient can then be calculated as $l'(W)_{j:} = \frac{2}{n} \sum_{i=1}^n \sum_{r \neq y_i} \max(A_{ir}, 0)(x_{ij}e_{y_i} - x_{ir}e_r)$ where e_r is a k dimensional vector with

zeroes everywhere except for a 1 at the r -th position. We only have to examine elements in A for which the corresponding x_{ij} is non-zero. When x_i is sparse, more often than not x_{ij} is zero and can be ignored.

2.2 Latent Group Lasso

One limitation of group lasso is that it assumes that group assignments are non-overlapping. In some domains, this can be too restrictive an assumption. For example, in document classification, individual words are used as features. If we were to construct groupings of these features, we might run into a case where one word could reasonably be added to several groups. The overlapping or latent group lasso was introduced to handle such cases.

??? develop a theoretical justification for the latent group lasso, as well as its equivalence to a regular group lasso in a higher dimensional space. Let \mathcal{G} be the set of (possibly overlapping) groups, where $g \in \mathcal{G}$ is a set of indices of covariates associated with that group. Let our data consist of vectors x_i in p dimensions, and let w be the corresponding weight vector in p dimensions that we would like to learn. Finally, define $\text{supp}(v)$ to be the support of v , i.e. the indices of the non-zero elements in v .

For each group $g \in \mathcal{G}$ we associate a latent vector $v^g \in \mathcal{R}^p$ where $\text{supp}(v^g) = g$, i.e. the nonzero elements in the v^g correspond to the indices in the group g . The original weight vector w can be interpreted as a sum of the latent vectors, or $w = \sum_{g \in \mathcal{G}} v^g$. ??? arrive at the following minimization problem

$$\begin{aligned} \min_{w, v^g} l(w) + \lambda \sum_{g \in \mathcal{G}} d_g \|v^g\|_2 \\ \text{s.t. } w = \sum_g v^g \end{aligned}$$

??? show that when the original problem is regression, $w^T x = \left(\sum_g v^g \right)^T x = \hat{v}^T \hat{x}$ where $\hat{v} = (v^g)_{g \in \mathcal{G}}$ and $\hat{x} = \bigoplus_{g \in \mathcal{G}} (x_i)_{i \in g}$, i.e. \hat{x} is the restrictions of each g stacked on top of each other. \hat{x}, \hat{v} have dimension $\sum_{g \in \mathcal{G}} |g|$. In this formulation, the optimal \hat{v}^* can be found using regular non-overlapping group lasso.

3 Our Model

Given n training vectors $x_i \in \mathcal{R}_d$ and their class labels $y_i \in \{1, \dots, m\}$, our goal is to compute

W such that it maximizes the accuracy of our prediction and it is group-wise sparse.

In our model, we minimize the following objective function :

$$\begin{aligned} \min_{W \in \mathcal{R}^{d \times m}} F(W) = \\ \frac{1}{n} \sum_{i=1}^n \sum_{r \neq y_i} \max(1 - (W_{:y_i}^T \cdot x_i - W_{:r}^T \cdot x_i), 0)^2 \\ + \lambda \sum_{g \in \mathcal{G}} \sum_{r=1}^d \|W_{g,r}\|_2 \end{aligned}$$

The first term is the multiclass squared hinge loss function. We want the dot product of an instance and its feature vector to be as large as possible, and the dot product of this instance and the rest feature vectors to be as small as possible. And as long as their difference is greater than a margin (1 in this case), we won't penalize it. In the second term, $W_{g,r}$ means a block of weights in group g and class m . The L2-norm regularization is computed and sum up for each block. The $\lambda > 0$ is a parameter controls the trade-off between the hinge loss and the L2-norm regularization.

4 Data

4.1 Newsgroup Data

4.1.1 Group Identification

4.2 Artificial Data

For the datasets described above, we can't tell with 100 percent confidence that the datasets follow the assumptions of the group structures for the features. And even if they are indeed structured that way, we maybe wrong with the method of coming up with the groups. These issues make it difficult to access our model.

To get rid of all these problems and validate the effectiveness of our model, we created artificial data that followed the underlying assumptions of the model. First, we generate a sparse weight matrix W to represent the relationship between features and classes. The weight matrix W has an internal structure in which features are grouped together. And also, only a small number of groups

have non-zero weights. This makes the matrix sparse.

```
[[ 0.    0.248 -0.766  0.    0. ]
 [ 0.   -0.386  0.754  0.    0. ]
 [ 0.   -0.282  0.488  0.    0. ]
 [ 0.    0.656 -0.616  0.    0. ]
 [ 0.   -0.118 -0.516  0.    0. ]
 [ 0.    0.886  0.    0.   -0.56 ]
 [ 0.   -0.96  0.    0.    0.972]
 [ 0.    0.332  0.    0.   -0.192]
 [ 0.    0.138  0.    0.    0.436]
 [ 0.   -0.134  0.    0.    0.154]
 [ 0.042  0.    0.   -0.216  0. ]
 [ 0.772  0.    0.    0.166  0. ]
 [-0.376  0.    0.    0.01  0. ]
 [-0.344  0.    0.    0.394  0. ]
 [-0.376  0.    0.    0.178  0. ]
 [ 0.    0.    0.276  0.    0. ]
 [ 0.    0.    0.556  0.    0. ]
 [ 0.    0.    0.712  0.    0. ]
 [ 0.    0.   -0.5  0.    0. ]
 [ 0.    0.   -0.18  0.    0. ]
 [ 0.    0.   -0.648  0.    0. ]
 [ 0.    0.   -0.04  0.    0. ]
 [ 0.    0.   -0.824  0.    0. ]
 [ 0.    0.   -0.044  0.    0. ]
 [ 0.    0.   -0.314  0.    0. ]]
```

Figure 1: Group-wise sparse weight matrix generated: 5 classes, 25 features in 5 groups

Then we generate random vectors, each of which has a length of the number of all features, and calculate dot product with the weight matrix W to get the class assignments for these random vectors. The random vectors X and the class assignments Y make up the training data set.

Our goal is to infer this weight matrix W from X and Y using our model. By generating the data set using this method, we can test the effectiveness of our model on a noiseless dataset with right underlying assumptions.

5 Results

5.1 Newsgroup Data

5.2 Artificial Data

Shape Recovery. One of the main indicator of the effectiveness of our model is to see whether the calculated weight matrix is sparse group-wise. Our experiments show so. The following figures shows in a typical trial, the generated target weight matrix and the recovered weight matrix by our model. By comparing them side by side, we can tell that their sparsity patterns are similar.

The different weights between the calculated matrix and the target matrix can be due to many factors. First, sample coverage is a major factor. In our simulation data, the sample size is small. Limited by the time complexity of the algorithm, it's difficult to complete the computation for a very large sample size in a reasonable time. Also,

```
[[ 0.    0.634  0.    0.    0. ]
 [ 0.   -0.472  0.    0.    0. ]
 [ 0.   -0.494  0.    0.    0. ]
 [ 0.   -0.332  0.    0.    0. ]
 [ 0.    0.036  0.    0.    0. ]
 [ 0.    0.   -0.328 -0.362  0. ]
 [ 0.    0.   -0.758 -0.124  0. ]
 [ 0.    0.   -0.048  0.764  0. ]
 [ 0.    0.   -0.62  -0.904  0. ]
 [ 0.    0.    0.908 -0.628  0. ]
 [ 0.    0.    0.    0.    0. ]
 [ 0.    0.    0.    0.    0. ]
 [ 0.    0.    0.    0.    0. ]
 [ 0.    0.    0.    0.    0. ]
 [ 0.    0.    0.    0.    0. ]
 [ 0.    0.    0.    0.    0. ]
 [ 0.    0.    0.    0.    0. ]
 [ 0.    0.    0.    0.    0. ]
 [ 0.    0.    0.    0.    0. ]
 [ 0.206 -0.212  0.    0.    0. ]
 [ 0.678 -0.496  0.    0.    0. ]
 [-0.076  0.184  0.    0.    0. ]
 [-0.366  0.392  0.    0.    0. ]
 [-0.04  0.278  0.    0.    0. ]]
```

Figure 2: Target group-wise sparse weight matrix generated: 5 classes, 25 features in 5 groups

the sampling is random. There is no guaranteed that the inferred weights leads to the target weights.

Accuracy. In our experiments, the generator algorithm was configured to produce 150 random vectors from the underlying model where it consists of 5 classes and 25 features in 5 groups. The accuracy achieved was about 60% to 70%.

6 Conclusion

```

[[ 0.      1.816e-01 0.      0.      0.      ]
 [ 0.      1.828e-02 0.      0.      0.      ]
 [ 0.      1.898e-02 0.      0.      0.      ]
 [ 0.     -5.964e-03 0.      0.      0.      ]
 [ 0.      3.866e-02 0.      0.      0.      ]
 [ 0.      0.     -5.741e-02 -1.187e-01 0.      ]
 [ 0.      0.     -9.831e-02 -9.476e-02 0.      ]
 [ 0.      0.     -5.357e-02 -5.227e-02 0.      ]
 [ 0.      0.     -6.184e-02 -1.363e-01 0.      ]
 [ 0.      0.     -9.737e-05 -1.237e-01 0.      ]
 [ 1.843e-01 0.      0.      0.      0.      ]
 [ 8.137e-02 0.      0.      0.      0.      ]
 [ 8.823e-02 0.      0.      0.      0.      ]
 [ 8.286e-02 0.      0.      0.      0.      ]
 [ 5.945e-02 0.      0.      0.      0.      ]
 [ 4.898e-03 0.      0.      0.      0.      ]
 [ 2.967e-03 0.      0.      0.      0.      ]
 [ 4.154e-03 0.      0.      0.      0.      ]
 [ 4.843e-03 0.      0.      0.      0.      ]
 [ 3.573e-03 0.      0.      0.      0.      ]
 [ 2.809e-01 0.      0.      0.      0.      ]
 [ 4.841e-01 0.      0.      0.      0.      ]
 [ 1.149e-01 0.      0.      0.      0.      ]
 [ -6.285e-02 0.      0.      0.      0.      ]
 [ 9.619e-02 0.      0.      0.      0.      ]]

```

Figure 3: Weight matrix calculated from the training data. It's similar to the one generated.