

CPSC 340 Assignment 6 (due Friday November 29 at 11:55pm)

1 Robust PCA

The function `example_RPCA` loads a dataset X where each row contains the pixels from a single frame of a video of a highway. The demo applies PCA to this dataset and then uses this to reconstruct the original image. It then shows the following 3 images for each frame of the first 50 frames (pausing for a tenth of a second on each frame):

1. The original frame.
2. The reconstruction based on PCA.
3. A binary image showing locations where the reconstruction error is non-trivial.

Recently, latent-factor models have been proposed as a strategy for “background subtraction”: trying to separate objects from their background. In this case, the background is the highway and the objects are the cars on the highway. In this demo, we see that PCA does an ok job of identifying the cars on the highway in that it does tend to identify the locations of cars. However, the results aren’t great as it identifies quite a few irrelevant parts of the image as objects.

Robust PCA is a variation on PCA where we replace the L2-norm with the L1-norm,

$$f(Z, W) = \sum_{i=1}^n \sum_{j=1}^d |w_j^T z_i - x_{ij}|,$$

and it has recently been proposed as a more effective model for background subtraction. [Write a new function, `robustPCA`, that uses the Huber loss to approximate the absolute value to implement robust PCA. Hand in your code.](#)

Hint: most of the work has been done for you in the function `PCA_gradient`. This function implements an alternating minimization approach to minimizing the PCA objective (without enforcing orthogonality). This gradient-based approach to PCA can be modified to use the Huber loss. A reasonable value of the hyper-parameter ϵ in the Huber loss might be 0.01 for this problem. You may want to use a smaller version of the dataset when debugging your objective/gradient code.

```
File Edit View Run Kernel Tabs Settings Help
Untitled-1.jupyter X IE example_RPCA.j X IE PCA.j X
79 # s = s + w = X
80 # s = zeros(size(W))
81 for j in 1:d
82     for i in 1:n
83         show_size(W[i,:])
84         #
85         r = dot(Z[i,1:M],Z[i,:]) - X[i]
86         if abs(r) > 0.01
87             r = r + 0.01 * (abs(r) - 0.01 * 0.5)
88         else
89             r = r + 0.5*r
90         end
91         # s = s + sign(r) * min(abs(r),1) * transpose(Z[i,:])
92     end
93 end
94 @print("Iteration %d, loss = %f", iter, f/length(X))
95 if (fold - f)/length(X) < 1e-2
96     break
97 end
98 end
99
100 # we didn't enforce that W was orthogonal so we need to optimize to find Z
101 compress(Xhat) = compress_gradientDescent(Xhat,W,Mu)
102 expand(Z) = expandfunc(Z,W,Mu)
103 return CompressModel(compress,expand,W)
104 end
105
106 function compress_gradientDescent(Xhat,W,Mu)
107     (T,K) = size(Xhat)
108     k = size(W,1)
109     Xcentered = Xhat - repeat(Mu,T,1)
110     Z = zeros(T,k)
111     func(f) = pcAGD(f,Z,Xcentered,W)
112     [f] = findMin(func,Z,[],verbose=false)
113     return Z
114 end
115
116 function pcAGD(f,Z,X,W)
117     # Scale vector of parameters into matrix
118     n = size(X,1)
119     k = size(W,1)
120     Z = reshape(Z,n,k)
121     #
122     # s = zeros(size(Z))
123     for j in 1:d
124         for i in 1:n
125             r = dot(W[i,:],Z[i,:]) - X[i]
126             if abs(r) > 0.01
127                 r = r + 0.01 * (abs(r) - 0.01 * 0.5)
128             else
129                 r = r + 0.5*r
130             end
131             show_size(f[i,:])
132             show_size(transpose(W[i,:]))
133             # s[i,:] = s[i,:] + sign(r) * min(abs(r),0.01) * (W[i,:])'
134         end
135     end
136     # Return function and gradient vector
137     return (f,g[])
138 end
139
140 function pcAGD(W,K,K,2)
141     # Scale vector of parameters into matrix
142     d = size(X,2)
143     n = size(Z,2)
144     W = reshape(W,n,d)
145     #
146     # Compute function value
147     # f = 2*W * X
148     # f = (1/2)sum(W.*Z)
149     #
150     # Compute derivative with respect to each residual
151     # df = R
152     #
153     # Multiply by Z' to get elements of gradient
154     # g = Z'*df
155     #
156     # s = zeros(size(W))
157     for j in 1:d
158         for i in 1:n
159             r = dot(Z[i,1:M],W[i,:]) - X[i]
160             if abs(r) > 0.01
161                 r = r + 0.01 * (abs(r) - 0.01 * 0.5)
162             else
163                 r = r + 0.5*r
164             end
165             # s[i,:] = s[i,:] + sign(r) * min(abs(r),0.01) * (Z[i,:])'
166         end
167     end
168     # Return function and gradient vector
169     return (f,g[])
170 end
171
172
173
174
175
176
177
178
179
```

Answer:

2 Multi-Dimensional Scaling

The function `example_MDS` loads the animals dataset and then applies gradient descent to minimize the following multi-dimensional scaling (MDS) objective (starting from the PCA solution):

$$f(Z) = \frac{1}{2} \sum_{i=1}^n \sum_{j=i+1}^n (\|z_i - z_j\| - \|x_i - x_j\|)^2. \quad (1)$$

The result of applying MDS is shown below. Although this visualization isn't perfect (with "gorilla" being placed close to the dogs and "otter" being placed close to two types of bears), this visualization does organize the animals in a mostly-logical way.

2.1 ISOMAP

Euclidean distances between very different animals are unlikely to be particularly meaningful. However, since related animals tend to share similar traits we might expect the animals to live on a low-dimensional manifold. This suggests that ISOMAP may give a better visualization. Make a new function `ISOMAP` that computes the approximate geodesic distance (shortest path through a graph where the edges are only between nodes that are k -nearest neighbour) between each pair of points, and then fits a standard MDS model (1) using gradient descent. [Hand in your code and the plot of the result when using the 3-nearest neighbours.](#)

Hint: the function `dijkstra` (in `misc.jl`) can be used to compute the shortest (weighted) distance between two points in a weighted graph. This function requires an n by n matrix giving the weights on each edge (use ∞ as the weight for absent edges). Note that ISOMAP uses an undirected graph, while the k -nearest neighbour graph might be asymmetric. You can use the usual heuristics to turn this into an undirected graph of including an edge i to j if i is a KNN of j or if j is a KNN of i . (Also, be careful not to include the point itself in the KNN list).

2.2 ISOMAP with Disconnected Graph

An issue with measuring distances on graphs is that the graph may not be connected. For example, if you run your ISOMAP code with 2-nearest neighbours then some of the distances are infinite. One heuristic to address this is to set these infinite distances to the maximum distance in the graph (i.e., the maximum geodesic distance between any two points that are connected), which will encourage non-connected points to be far apart. Modify your ISOMAP function to implement this heuristic. [Hand in your code and the plot of the result when using the 2-nearest neighbours.](#)

3 Neural Networks

3.1 Neural Networks by Hand

Suppose that we train a neural network with sigmoid activations and one hidden layer and obtain the following parameters (assume that we don't use any bias variables):

$$W = \begin{bmatrix} -2 & 2 & -1 \\ 1 & -2 & 0 \end{bmatrix}, v = \begin{bmatrix} 3 \\ 1 \end{bmatrix}.$$

Assuming that we are doing regression, [for a training example with features \$x_i^T = \[-3 \ -2 \ 2\]\$ what are the values in this network of \$z_i\$, \$h\(z_i\)\$, and \$\hat{y}_i\$?](#)

3.2 Neural Network Tuning - Regression

The file `example_nnet.jl` runs a stochastic gradient method to train a neural network on the *basisData* dataset from a previous assignment. However, in its current form it doesn't fit the data very well. Modify the training procedure to improve the performance of the neural network. [Hand in your plot after changing the code to have better performance, and list the changes you made.](#)

Hint: there are many possible strategies you could take to improve performance. Below are some suggestions, but note that the some will be more effective than others:

- Changing the network structure (*nHidden* is a vector giving the number of hidden units in each layer).
- Changing the training procedure (you can change the stochastic gradient step-size, use mini-batches, run it for more iterations, add momentum, switch to *findMin*, and so on).
- Transform the data by standardizing the features, standardizing the targets, and so on.
- Add regularization (L2-regularization, L1-regularization, dropout, and so on).
- Add bias variables within the hidden layers.
- Change the loss function or the non-linearities (right now it uses squared error and tanh to introduce non-linearity).
- Use mini-batches of data, possibly with batch normalization.

3.3 Neural Network Tuning - Classification

The file `example_usps.jl` runs a stochastic gradient method to train a neural network on a set of images of digits. Modify the training procedure to improve the performance of the neural network. [List the changes you made and the best test performance that you were able to obtain.](#)

4 Very-Short Answer Questions

1. [Is the NMF loss function convex? What is an optimization method you could use to try to minimize it?](#)
2. [Consider fitting a linear latent-factor model, using L1-regularization of the \$w_c\$ values and L2-regularization of the \$z_i\$ values,](#)

$$f(Z, W) = \frac{1}{2} \|ZW - X\|_F^2 + \lambda_W \sum_{c=1}^k [\|w_c\|_1] + \frac{\lambda_Z}{2} \sum_{i=1}^n [\|z_i\|^2],$$

- (a) [What is the effect of \$\lambda_Z\$ on the two parts of the fundamental trade-off in machine learning?](#)
 - (b) [What is the effect of \$k\$ on the two parts?](#)
 - (c) [Would either of answers to the previous two questions change if \$\lambda_W = 0\$?](#)
3. [Which is better for recommending movies to a new user, collaborative filtering or content-based filtering? Briefly justify your answer.](#)
 4. [Is ISOMAP mainly used for supervised or unsupervised learning? Is it parametric or non-parametric?](#)
 5. [What is the difference between Euclidean distance and geodesic distance?](#)
 6. [Are neural networks mainly used for supervised or unsupervised learning? Are they parametric or nonparametric?](#)

7. The loss for a neural network is typically non-convex. Give one set of hyperparameters for which the loss is actually convex.
8. Assuming we have some procedure that returns a random global minimum of a neural network objective, how does the depth of a neural network affect the fundamental trade-off? For a convolutional network, how would the width of the convolutions affect the fundamental trade-off?
9. What is the “vanishing gradient” problem with neural networks based on sigmoid non-linearities?
10. List 3 forms of regularization we use to prevent overfitting in neural networks.
11. Convolutional networks seem like a pain... why not just use regular (“fully connected”) neural networks for image classification?