# SDS 385: Homework 4

G. Paulon

October 2, 2016

**Problem 1.  Improving SGD for logistic regression**

(A) We first implement the stochastic gradient descent coupled with backtracking for logistic regression. This algorithm uses minibatches of the data in order to pick the correct step size for a certain number of iterations. The step size is refreshed every $k$ iterations, which leads to a natural decay over time. The code is reported in Listing A.1. The results of this algorithm applied to the cancer dataset are not as promising as the results obtained with the adaptive gradient descent. For this reason, in the following we will discuss an efficient implementation and a big-data application only of the second algorithm.

(B) The Adaptive Gradient Descent has been implemented. Let us describe the updating rules for this algorithm, which uses a diagonal approximation of the Hessian at every iteration. Let us first recall that the negative log-likelihood of the $i^{th}$ data point is given by

$$l_i(\alpha, \boldsymbol{\beta}) \propto -y_i \log(w_i) - (m_i - y_i) \log(1 - w_i) - \sum_{j=1}^{p} |\beta_j|$$

$$= -m_i \log(1 - w_i) - y_i \log\left(\frac{w_i}{1 - w_i}\right) - \sum_{j=1}^{p} |\beta_j|$$

$$= m_i \log(1 + \exp(\alpha + \boldsymbol{x}_i^T \boldsymbol{\beta})) - y_i(\alpha + \boldsymbol{x}_i^T \boldsymbol{\beta}) - \sum_{j=1}^{p} |\beta_j|$$

where we added a penalization term that forces sparsity in the problem. Let us remark that we split the intercept and the other covariates. This is done for notation purposes, as the intercept is not penalized in the last term of the sum.

The gradient with respect to $\alpha$ of this contribution is

$$\nabla_\alpha l_i(\alpha, \boldsymbol{\beta}) \propto m_i \frac{e^{\alpha + \boldsymbol{x}_i^T \boldsymbol{\beta}}}{1 + e^{\alpha + \boldsymbol{x}_i^T \boldsymbol{\beta}}} - y_i$$

$$= \hat{y}_i - y_i,$$

where

$$\hat{y}_i = m_i \frac{e^{\alpha + \boldsymbol{x}_i^T \boldsymbol{\beta}}}{1 + e^{\alpha + \boldsymbol{x}_i^T \boldsymbol{\beta}}}.$$

The gradient with respect to $\boldsymbol{\beta}$ of the individual negative log-likelihood is

$$\nabla_{\beta_j} l_i(\alpha, \boldsymbol{\beta}) \propto m_i \frac{x_{ij} e^{\alpha + \boldsymbol{x}_i^T \boldsymbol{\beta}}}{1 + e^{\alpha + \boldsymbol{x}_i^T \boldsymbol{\beta}}} - x_{ij} y_i - \lambda \frac{\beta_j}{|\beta_j|}$$

$$= x_{ij}(\hat{y}_i - y_i) - \lambda \frac{\beta_j}{|\beta_j|}.$$

The code is illustrated in A.2.

The algorithm has been applied to the usual dataset containing a cancer classification problem. The speed-up is remarkable if compared to the same version of the algorithm implemented in R language. The convergence results are displayed in Figure 1.
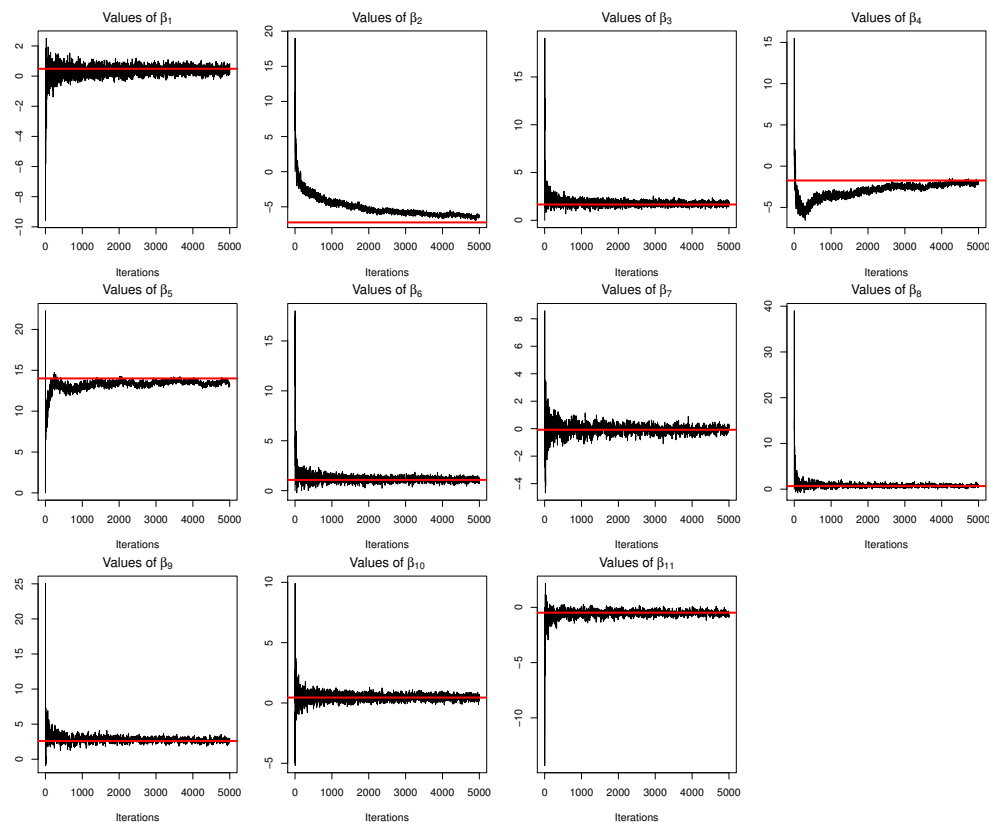
*Figure 1: Results of the Adagrad method. The red lines denote the results obtained with Newton method.*

## Problem 2.   Putting it all together on some biggish data

The algorithm has been improved by adding support for sparse matrices and it has been applied to a large (and sparse) dataset whose goal is to detect malicious URLs in web traffic.

The optimal value of the regularizer parameter $\lambda$ has been chosen by 5-fold cross validation. The estimate of the test error (prediction accuracy) in this case is around $98.0\%$, which is a remarkable result.

| | $\lambda = 0$ | $\lambda = 10^{-10}$ | $\lambda = 10^{-9}$ | $\lambda = 10^{-8}$ | $\lambda = 10^{-7}$ | $\lambda = 10^{-6}$ | $\lambda = 10^{-5}$ | $\lambda = 10^{-4}$ |
|---|---|---|---|---|---|---|---|---|
| Mean | 0.9856 | 0.9856 | 0.9855 | 0.9853 | 0.9841 | 0.9787 | 0.9724 | 0.9626 |
| Std. dev. | 0.0006 | 0.0005 | 0.0005 | 0.0006 | 0.0008 | 0.0014 | 0.0019 | 0.0028 |

*Table 1: Cross-validation test error estimates for several values of the regularization parameter $\lambda$.*

# Appendix A

# R code

```r
SGD.linesearch <- function(y, X, mi, beta0, maxiter = 100000, tol = 1E-8){
  # Function for the gradient descent of the logit model coupled with backtracking
  # -------------------------------------------------------------------------------
  # Args:
  #   - y: response vector (length n)
  #   - X: matrix of the features (n*p)
  #   - mi: vector of the number of trials, always 1 in the logit framework (length n)
  #   - beta0: initial regression parameters (length p)
  #   - maxiter: number of maximum iterations the algorithm will perform if not converging
  #   - tol: tolerance threshold to convergence
  # Returns:
  #   - ll: values of the log-likelihood for every iteration
  #   - beta: final optimal regression parameters
  #   - alpha: selected step size for every iteration
  # ------------------------------------------------
  N <- dim(X)[1]
  P <- length(beta0)

  # Select the updating time and the size of the minibatches
  train.epoch <- floor(maxiter/1000)
  size.minibatch <- floor(N/10)

  # Sample the data points to calculate the gradient
  idx <- sample(1:N, maxiter, replace = TRUE)
  # Sample the indexes for the minibatch
  minibatches <- matrix(sample(1:N, 1000*size.minibatch, replace = TRUE), 1000, size.minibatch)
  idx.minibatch <- minibatches[1, ]

  betas <- array(NA, dim=c(maxiter, length(beta0)))
  betas[1, ] <- beta0 # Initial guess

  av.ll <- array(NA, dim = maxiter)
  ll <- log.lik(betas[1, ], y[idx[1]], matrix(X[idx[1],], nrow=1), mi[idx[1]])
  av.ll[1] <- ll

  for (iter in 2:maxiter){

    # Choose the direction according to the gradient of the single individual idx[i]
    grad <- grad.loglik(betas[iter-1,], y[idx[iter]], matrix(X[idx[iter], ],nrow=1), mi[idx[iter
        ]])
```

3

```
40        dir <- - grad
41        # if we need to update the optimal step
42        if ((iter == 2) | (iter %% train.epoch == 0)){
43          idx.minibatch <- minibatches[ceiling(iter/train.epoch), ]
44          grad.minibatch <- grad.loglik(betas[iter-1, ], y[idx.minibatch], X[idx.minibatch, ], mi[idx
                  .minibatch]) / length(idx.minibatch)
45          alpha <- optimal.step(betas[iter-1, ], -grad.minibatch, grad.minibatch, y[idx[iter]],
                  matrix(X[idx[iter], ], nrow=1), mi[idx[iter]])
46        }
47
48        # Update Beta parameters
49        betas[iter, ] <- betas[iter-1, ] + alpha * dir
50        # Compute log-likelihood and average log-likelihood
51        ll <- log.lik(betas[iter, ], y[idx[iter]], matrix(X[idx[iter],], nrow=1), mi[idx[iter]])
52        av.ll[iter] <- ((av.ll[iter-1])*(iter-1) + ll)/iter
53
54        # Convergence check
55        if (abs(av.ll[iter-1] - av.ll[iter]) / (av.ll[iter-1] + 1E-10) < tol){
56          av.ll <- av.ll[1:iter]
57          betas <- betas[1:iter, ]
58          break;
59        }
60
61        else if (iter == maxiter & abs(av.ll[iter-1] - av.ll[iter])/(av.ll[iter-1] + 1E-10) >= tol){
62          break;
63        }
64    }
65    return(list("ll" = av.ll, "beta" = betas))
66 }
```

*Listing A.1: Function implementing the SGD coupled with backtracking.*

```
1  #define ARMA_64BIT_WORD
2  // [[Rcpp::depends(RcppEigen)]]
3  #include <RcppEigen.h>
4  #include <iostream>
5  #include <algorithm>
6  #include <string>
7
8  using namespace Rcpp;
9  using Eigen::VectorXd;
10 using Eigen::SparseVector;
11
12 typedef Eigen::MappedSparseMatrix<double>  MapMatd;
13
14 // Function to compute the sign of a generic type
15 template <typename T> int sgn(T val) {
16   return (T(0) < val) - (val < T(0));
17 }
18
19 // [[Rcpp::export]]
20 SEXP Ada_Grad(MapMatd& X, VectorXd& y, VectorXd& m, VectorXd& beta0, double eta = 1.0, unsigned
        int npass = 1, double lambda = 0.0, double discount = 0.01){
21   // X is the design matrix stored in sparse column-major format
22   // i.e. with features for case i stores in column i
23   // y is the response vector
24   // m is the vector of sample sizes
```

```
25
26    unsigned int N = X.cols();
27    unsigned int P = X.rows();
28
29    // x is the sparse vector that, at each iteration, will contain the columns of X
30    SparseVector<double> x(P);
31
32    // Initialize parameters
33    double psi0, epsi, yhat, delta, h;
34    double this_grad = 0.0;
35    double mu, gammatilde;
36
37    // Initialize parameters alpha and beta
38    double alpha = 1.0;
39    VectorXd beta(P);
40
41    // Initialize historical gradients (cumulative)
42    double hist_int = 0.0; // historical gradient of the intercept term
43    VectorXd hist_grad(P); // historical gradient of the beta parameters
44    for (int j = 0; j < P; j++){
45      hist_grad(j) = 1e-3;
46      beta(j) = beta0(j);
47    }
48
49    // Vector denoting how long has it been since the last update of each feature
50    NumericVector last_update(P, 0.0);
51
52    // negative log likelihood for assessing fit
53    double nll_avg = 0.0;
54    NumericVector nll_tracker(npass*N, 0.0);
55
56    unsigned int k = 0; // global interation counter
57    unsigned int j; // j is an index that will denote the positions of the nonzero elements in x
58    for(unsigned int pass = 0; pass < npass; pass++) {
59
60      // Loop over each observation (columns of X)
61      for(unsigned int i = 0; i < N; i++) {
62
63        // Form linear predictor and E(Y[i]) from features
64        x = X.innerVector(i);
65        psi0 = alpha + x.dot(beta);
66        epsi = exp(psi0);
67        yhat = m[i] * epsi/(1.0 + epsi);
68
69        // (1) Update nll average
70        nll_avg = (1.0 - discount) * nll_avg + discount * (m[i] * log(1 + epsi) - y[i] * psi0);
71        nll_tracker[k] = nll_avg;
72
73        // (2) Update intercept
74        delta = y[i] - yhat; // gradient with respect to the intercept
75        hist_int += delta * delta; // update historical gradient of the intercept
76        alpha += (eta/sqrt(hist_int)) * delta;
77
78        // (3) Update beta: iterate over the active features for this instance
79        for (SparseVector<double>::InnerIterator it(x); it; ++it) {
80
81          // Which feature is this?
```

```
82          j = it.index();
83
84          // STEP (a): aggregate all the penalty-only updates since the last time we updated this
                feature.
85          // This is a form of lazy updating in which we approximate all the "penalty-only" updates
                at once.
86          double skip = k - last_update(j);
87          h = sqrt(hist_grad(j));
88          gammatilde = skip * eta/h;
89          beta(j) = sgn(beta(j)) * fmax(0.0, fabs(beta(j)) - gammatilde*lambda);
90
91          // Update the last-update vector
92          last_update(j) = k;
93
94          // STEP (b): Now we compute the update for this observation.
95          // gradient of negative log likelihood
96          this_grad = delta*it.value();
97
98          // update adaGrad scaling for this feature
99          hist_grad(j) += this_grad*this_grad;
100
101         // scaled stepsize
102         h = sqrt(hist_grad(j));
103       gammatilde = eta/h;
104       mu = beta(j) + gammatilde * this_grad;
105         beta(j) = sgn(mu) * fmax(0.0, fabs(mu) - gammatilde * lambda);
106       }
107       k++; // increment global counter
108     }
109   }
110
111   // (4) Update last penalties:
112   // At the very end, apply the accumulated penalty for the variables we have not touched
          recently
113   for (int j = 0; j < P; j++) {
114     double skip = k - last_update(j);
115     h = sqrt(hist_grad(j));
116     gammatilde = skip*eta/h;
117     beta(j) = sgn(beta(j))*fmax(0.0, fabs(beta(j)) - gammatilde*lambda);
118   }
119
120   return List::create(Named("alpha") = alpha,
121                       Named("beta") = beta);
122 }
```

*Listing A.2: Function implementing Adagrad algorithm exploiting C++'s library Eigen.*