

Final Project

Evan Ott

December 5, 2016

Stochastic Gradient Descent (Logistic Regression)

Let's load in the data and Dr. Scott's algorithm:

```
library(ggplot2)

## Warning: package 'ggplot2' was built under R version 3.3.2

library(latex2exp)
library(Matrix)
library(Rcpp)
library(RcppEigen)
library(reshape2)

Rcpp::sourceCpp("scott_sgdlogit.cpp")

raw_data = read.csv("master.csv", header = TRUE)

# Read specific columns from the data

kegg = as.character(raw_data[, 1])
gene = as.character(raw_data[, 2])

data = raw_data[, 3:ncol(raw_data)]

N = ncol(data)
P = nrow(data)
M = rep(1, N) # Used in logistic regression

# For C++ code, it's easiest to use a sparse matrix
# TODO: might want to consider scaling the counts, either in the way DESeq2 does
# it or something similar.
X = Matrix(as.matrix(unname(data)), sparse = TRUE)

# X[i, j] == data[i, j] is non-negative, so this trims out the genes with no counts
X_nzero = X[which(rowSums(data) > 0), ]

# Convert the column names to 0 = control, 1 = treatment
Ynames = names(data)
Y = stringi::stri_endswith(Ynames, fixed = "T") * 1
```

Now, let's create a transformed version of X where the values are on the log scale.

```
# Copy it as to not damage the original
X_nzero_log = X_nzero

# Oooh I figured out a sweet hack. So the dgCMatrx class could be used in the
# following way:
# X@i are the row indices - 1 (that have ncol sections in non-decreasing order)
```

```

# X@p are the indices - 1 in X@i where we switch to a new column
#
# OR
# X@x is the (non-zero) data itself, so why not just edit that? =D
# Here, add a small value because log(1) == 0, and we want to maintain a value there.
# #oneliner
X_nozero_log@x = log2(X_nozero_log@x + 0.1)

```

Now, let's try to use cross-validation (on the limited number of samples we have) to determine a good range for λ , the LASSO penalty factor in the stochastic gradient decent logistic regression code.

```

# Initial guess for the betas.
beta_init = rep(0, P)
# For the weighting in the skip scenario for a feature
# TODO: algorithm is very sensitive to this value
eta = 0.0002
# passes through the data
nIter = 100
# exponentially-weighted moving average factor
# TODO: could explore other values.
discount = 0.01

# X_cv = X_nozero
X_cv = X_nozero_log

folds = 3
col_breaks = cut(1:length(Y), folds, labels = FALSE)
abs_errors = c()
zero_one_errors = c()
lambdas = 10 ^ seq(-1, 4, 0.1)
for (lambda in lambdas) {
  absFoldErrors = c()
  zeroOneFoldErrors = c()
  for (fold in 1:folds) {
    test_ind = which(col_breaks == fold)
    train_ind = which(col_breaks != fold)
    X_nozero_train = X_cv[, train_ind]
    Y_train = Y[train_ind]
    X_nozero_test = X_cv[, test_ind]
    Y_test = Y[test_ind]

    train_result = sparsesgd_logit(X_nozero_train, Y_train, M[train_ind], eta, nIter, beta_init, lambda)

    intercept = train_result$alpha
    beta = train_result$beta
    prediction_test = 1/(1 + exp(-(intercept + t(X_nozero_test) %*% beta)))

    # metrics to choose from
    zero_one_error_test = sum(Y_test != round(prediction_test))
    abs_error_test = sum(abs(Y_test - prediction_test))

    absFoldErrors = c(absFoldErrors, abs_error_test)
    zeroOneFoldErrors = c(zeroOneFoldErrors, zero_one_error_test)
  }
}

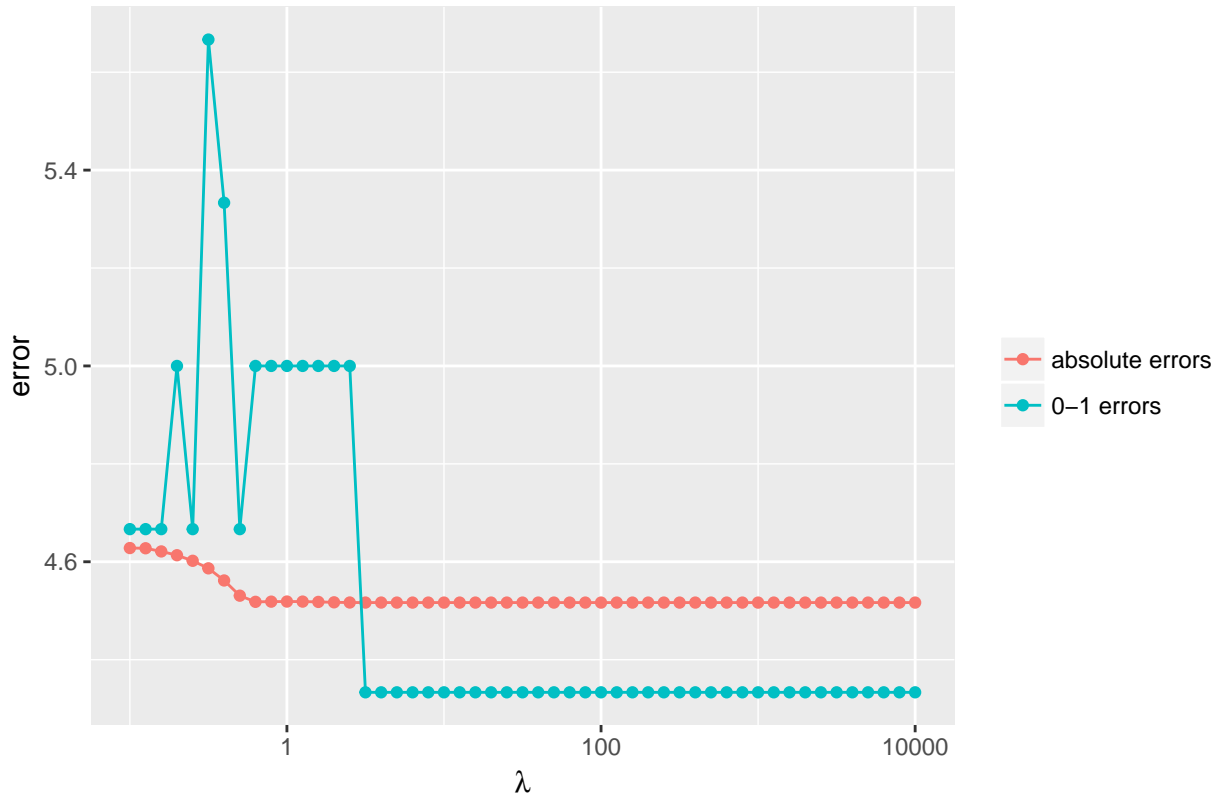
```

```

abs_errors = c(abs_errors, mean(absFoldErrors))
zero_one_errors = c(zero_one_errors, mean(zeroOneFoldErrors))
}

```

3-fold cross-validation out-of-sample errors



So based on 3-fold cross-validation, the absolute error $\sum_i |y_i - \hat{y}_i|$ seems to reach its best (minimal) value when $\lambda = 10^{-0.2}$ while the 0-1 error $\sum_i \delta(y_i, \text{round}(\hat{y}_i))$ (number of incorrect predictions) exhibits the same behavior around $\lambda = 10^{0.5}$.

Let's take the latter, since the absolute error shouldn't be different.

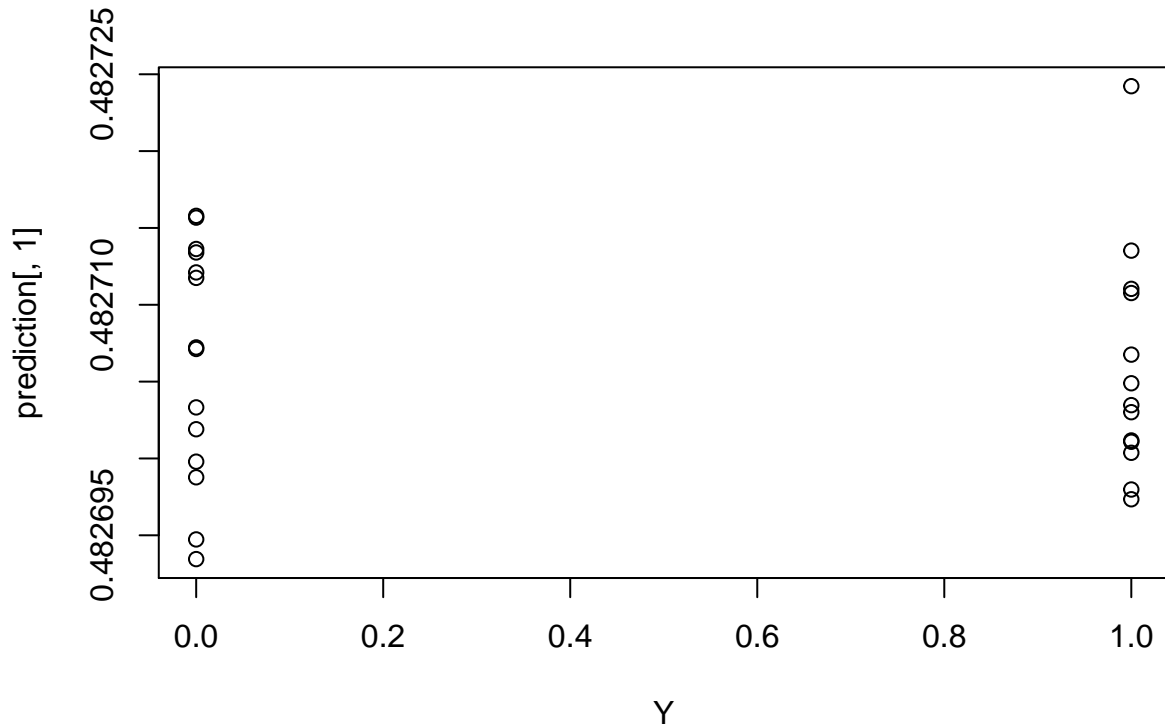
```

# Initial guess for the betas.
beta_init = rep(0, P)
# For the weighting in the skip scenario for a feature
# TODO: algorithm is very sensitive to this value
eta = 0.0002
# passes through the data
nIter = 100
# L1 regularization
lambda = 10^0.5
# exponentially-weighted moving average factor
# TODO: could explore other values.
discount = 0.01

# Run the algorithm
result = sparsesgd_logit(X_nozero_log, Y, M, eta, nIter, beta_init, lambda, discount)
intercept = result$alpha
beta = result$beta
prediction = 1/(1 + exp(-(intercept + t(X_nozero_log) %*% beta)))

```

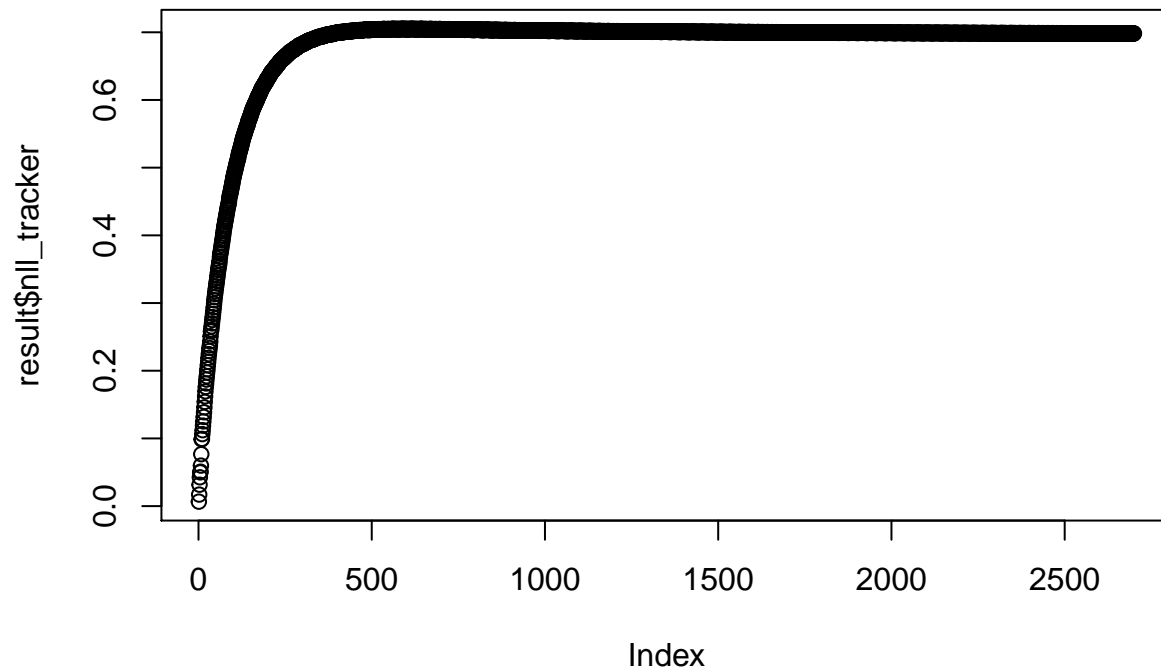
```
plot(prediction[, 1] ~ Y)
```



So, these predictions are, to use a technical term, “totally garbage.” Every single prediction is just slightly below $\frac{1}{2}$ so the algorithm is effectively saying, “I dunno, but maybe this data is control not treatment” so while the results below are potentially a starting point, it’d be rather foolish to grant them much authority.

```
## How many genes are in the final model? 25
## Which genes?
## adef_c_5_1802
## aot180_c_13_1642
## bdim_c_51_1620
## belk_c_262_3655
## belk_c_623_6138
## blon2249_c_1_1767
## ecas2549_c_2_485
## ecas2549_c_2_487
## ecas2549_c_2_492
## ecas2549_c_3_766
## efae1080_c_1_429
## gmor_c_6_936
## lfus_c_18_3454
## mneo_c_1_1955
## mneo_c_1_4619
## NCBIABWK_c_4_1154
## NCBIACFE_c_1_366
## NCBIACFU_c_47_2923
## not020_c_81_1706
## nsic_c_7_1124
## pmir_c_10_1585
## pot302_c_4_620
## psal_c_109_2279
## ptan_c_1_266
```

```
## rp1c_c_1_2351
```



DESeq2

Now, let's use a pre-fabricated solution instead.

```
library(DESeq2)
#library(BiocParallel)
# Allow for parallelization of DESeq2 code.
#register(MulticoreParam(4))

get_results = function(d) {
  nms = rownames(d)
  log2change = d[, 2]
  padj = d[, 6]
  idx = order(nms)
  return(data.frame(names = nms[idx], log2change = log2change[idx], padj = padj[idx]))
}

# Get the counts
count_data_raw = raw_data[, 3:ncol(raw_data)]

#####
# Format the data in the way DESeq2 expects
#####
new_colnames = rep(NULL, N)
condition = rep(NULL, N)
control = 0
treatment = 0
for (i in 1:N) {
  base = c("control", "treatment")[Y[i] + 1]
```

```

val = 0
if (base == "control") {
  control = control + 1
  val = control
} else {
  treatment = treatment + 1
  val = treatment
}
new_colnames[i] = paste0(base, val)
condition[i] = base
}
count_data = count_data_raw
colnames(count_data) = new_colnames
rownames(count_data) = gene

col_data = data.frame(condition)
rownames(col_data) = new_colnames

# remove the unobserved genes from the get-go
count_data = count_data[which(rowSums(data) > 0), ]

# Run DESeq2
dds = DESeqDataSetFromMatrix(countData = count_data,
                             colData = col_data,
                             design = ~ condition)
dds = DESeq(dds)#, parallel = TRUE)
res = results(dds)#, parallel = TRUE)

# plot(sort(res$pvalue))
# points(sort(res$padj), col = "red")
# P - sum(is.na(res$pvalue))
# P - sum(is.na(res$padj))

```

```
## How many genes are in the final model? 61
```

```
## [1] "Which genes in default DESeq2?"
```

```
##          names log2change      padj
## 1   abau_c_1_3356  -3.276361 0.0765527545
## 2   aisr_c_20_1998  -3.623824 0.0809405192
## 3   amas2385_c_4_872 -3.989862 0.0795164174
## 4   aot170_c_17_1854 -3.798250 0.0663167998
## 5   aot171_c_140_1750 -3.501483 0.0999237768
## 6   aot171_c_3_100  -3.701837 0.0894612894
## 7   aot172_c_99_1663 -3.288316 0.0729303040
## 8   aot448_c_6_964   3.611067 0.0458968210
## 9   apre_c_1_700    -3.880782 0.0999237768
## 10  cdip_c_1_1509   -4.746225 0.0263185564
## 11  chom_c_27_1386  -4.444855 0.0729303040
## 12  chom_c_55_2086  -4.119120 0.0729303040
## 13  cper_c_2_296    3.468789 0.0618240866
## 14  cure_c_1_1578   -3.149004 0.0729303040
## 15  dpig_c_1_8      5.053061 0.0164361041
## 16  efae1080_c_1_913 -3.923278 0.0992735671
## 17  esak_c_1_3901   -3.614598 0.0729303040
```

```
## 18 fnucp_c_3_1442 3.400389 0.0319580926
## 19 fper2555_c_2_652 4.693934 0.0319580926
## 20 gmor_c_19_1592 -4.400610 0.0263185564
## 21 lcat_c_18_1871 -3.454624 0.0802019105
## 22 lgas_c_1_25 -3.187085 0.0729303040
## 23 lgoo_c_14_911 -3.837013 0.0999237768
## 24 lmon_c_1_308 -4.066528 0.0994749580
## 25 lot107_c_9_1346 4.593022 0.0795164174
## 26 mlot_c_1_5349 -3.915125 0.0458968210
## 27 mneo_c_1_1226 -3.903096 0.0795164174
## 28 mneo_c_1_1611 -3.646940 0.0555936452
## 29 mot186_c_3_1980 -3.871212 0.0319580926
## 30 mtub_c_1_788 4.436387 0.0458968210
## 31 nbac_c_3_267 -4.703797 0.0319580926
## 32 NCBIABIX_c_1_1199 -3.672483 0.0729303040
## 33 opro_c_21_2166 -4.104369 0.0999237768
## 34 pend_c_1_176 -4.634184 0.0466055374
## 35 peno_c_75_2748 -4.137255 0.0729303040
## 36 pmel_c_20_1104 -4.833471 0.0458968210
## 37 pmuls_c_6_923 -4.824002 0.0729303040
## 38 pot786_c_28_1621 -4.365366 0.0729303040
## 39 pple_c_7_1120 -4.617908 0.0729303040
## 40 pstu_c_1_2331 -4.815178 0.0663167998
## 41 pver_c_10_1377 -4.920895 0.0812257002
## 42 raer_c_1_313 -3.464626 0.0729303040
## 43 raer_c_17_1833 -4.833198 0.0319580926
## 44 raer_c_2_422 -3.354042 0.0988788613
## 45 raer_c_5_1016 -4.696354 0.0729303040
## 46 raer_c_9_1389 -5.312673 0.0007006765
## 47 rden_c_1_560 -4.300782 0.0865728904
## 48 rden1994_c_1_356 -4.165709 0.0765527545
## 49 sked_c_1_1556 -3.935474 0.0988788613
## 50 sked_c_1_1767 -4.088478 0.0729303040
## 51 smal_c_1_1545 -4.671608 0.0729303040
## 52 smit_c_1_1759 -4.502563 0.0729303040
## 53 smut_c_1_1468 5.209984 0.0231500070
## 54 smut_c_1_1548 3.969355 0.0999237768
## 55 smut_c_1_1794 4.633750 0.0555936452
## 56 smut_c_1_371 3.440298 0.0663167998
## 57 smut_c_1_836 3.724677 0.0598894208
## 58 sot138_c_21_1472 -4.072691 0.0812257002
## 59 sot149_c_17_1639 -2.654332 0.0618240866
## 60 ssal_c_18_1887 -3.574628 0.0466055374
## 61 tmed_c_12_2247 -4.149022 0.0729303040
```

```
# Not exactly sure what this plot does either.
```

```
# plotMA(res, main="DESeq2", ylim=c(-3,3))
```

```
# Benjamini-Hochberg by hand (DESeq2 does something similar...)
```

```
alpha = 0.1
```

```
sorted_pval = sort(res$pvalue, na.last = TRUE)
```

```
numNA = sum(is.na(res$pvalue))
```

```
bh = sapply(1:P, function(i){ sorted_pval[i] <= alpha * i / (P - numNA)})
```

```
threshold = sorted_pval[max(which(bh))]
```

```

cat(paste("How many genes in customized Benjamini-Hochberg procedure",
          length(which(res$pvalue < threshold)), "\n", sep = "\n")) # 1 gene here

## How many genes in customized Benjamini-Hochberg procedure
## 1
print("Which genes in customized Benjamini-Hochberg procedure")

## [1] "Which genes in customized Benjamini-Hochberg procedure"
get_results(res[which(res$pvalue < threshold), ])

##           names log2change          padj
## 1 raer_c_9_1389  -5.312673 0.0007006765
# These two are equal -- throws out genes that were never observed section 1.5.3 of DESeq2 paper
print(paste(sum(rowSums(count_data) == 0),
            sum(is.na(res$pvalue))))

## [1] "0 0"
# 118687 genes are thrown out by the "independent filtering" for having a low
# mean normalized count
sum(is.na(res$padj)) - sum(is.na(res$pvalue))

## [1] 118687
# Not entirely sure what this plot means.
# plot(metadata(res)$filterNumRej,
#       type = "b", ylab = "number of rejections",
#       xlab = "quantiles of filter")
# lines(metadata(res)$lo.fit, col = "red")
# abline(v = metadata(res)$filterTheta)

#nofilter
# If you turn off the filtering, only two of the adjusted p-values are less than 0.1
# where the adjustment basically just accounts for Benjamini-Hochberg, I think
resNoFilt <- results(dds, independentFiltering = FALSE)
addmargins(table(filtering = (res$padj < .1),
                  noFiltering = (resNoFilt$padj < .1)))

##           noFiltering
## filtering FALSE  TRUE   Sum
##   FALSE 33700     0 33700
##   TRUE   59      2   61
##   Sum   33759     2 33761

## How many genes without independent filtering
## 2

## [1] "How many genes without independent filtering"

##           names log2change          padj
## 1   dpig_c_1_8   5.053061 0.074217328
## 2 raer_c_9_1389 -5.312673 0.003163909

```