

## COMS 4772

## Homework Set 4

1. You may use the fact that *expectation is a linear operator*.

(a) For a random variable  $X$ , let  $EX$  denote its expected value. Show that

$$E((X - EX)(X - EX)^T) = E(XX^T) - EX(EX)^T.$$

The quantity on the left hand side is the variance-covariance matrix for  $X$ , which we will call  $V(X)$ .

Solution:

$$\begin{aligned} V(X) &= E\left[(X - E(X))(X - E(X))^T\right] \\ &= E\left[XX^T - E(X)X^T - X[E(X)]^T + E(X)[E(X)]^T\right] \\ &= E[XX^T] - E[E(X)X^T] - E[X[E(X)]^T] + E[E(X)E(X)^T] \\ &= E[XX^T] - E(X)[E(X)]^T \end{aligned}$$

(b) Show that, for any (appropriately sized) matrix  $A$  we have

$$V(AX) = A(V(X))A^T.$$

Solution:

$$\begin{aligned} V(AX) &= E\left[(AX - E(AX))(AX - E(AX))^T\right] \\ &= E\left[AXX^T A^T - E(AX)X^T A^T - AX E(X^T A^T) + E(AX)E(X^T A^T)\right] \\ &= E(AXX^T A^T) + E(AX)E(X^T A^T) \\ &= AE(XX^T)A^T + AE(X)E(X^T)A^T \\ &= A[V(X)]A^T \end{aligned}$$

(c) Show that

$$E(\|X\|^2) = \text{trace}(V(X)) + \|EX\|^2.$$

Solution:

$$\begin{aligned} E(\|X\|^2) &= E(X^\top X) \\ \|E(X)\|^2 &= [E(X)]^\top E(X) \\ \text{tr}(V(X)) &= \text{tr}(E(XX^\top) - E(X)E(X)^\top) = E(X^\top X) - [E(X)]^\top E(X) \end{aligned}$$

Clearly, by doing some addition, the above statement is true.

(d) Solve the stochastic optimization problem

$$\min_y E\|X - y\|_2^2,$$

where  $X$  is a random vector, and the expectation is taken with respect to  $X$ . What is the minimizer? What's the minimum value?

$$\min_y E\|X - y\|_2^2 = \min_y E_x[X^\top X - 2X^\top y + y^\top y] = \min_y E[X^\top X] - 2y^\top E[X] + y^\top y$$

The minimizer is the mean of  $X$ .

$$\frac{\partial}{\partial y} \min_y E\|X - y\|_2^2 = 0 \Rightarrow y^* = E[X].$$

Let's substitute  $y^*$  for  $y$ . The minimum value of the cost function is the variance of  $X$ .

$$E\|X - y^*\|_2^2 = E[X^\top X] - [E(X)]^2 = \text{Var}(X).$$

2. Frobenius norm estimation. Suppose we want to estimate

$$\|A\|_F^2 = \text{trace}(A^\top A)$$

of a large matrix  $A$ . One way to do this is to hit  $A$  by random vectors  $w$ , and then measure the resulting norm.

(a) Find a sufficient conditions on a random vector  $w$  that ensures

$$E\|Aw\|^2 = \|A\|_F^2.$$

Prove that your condition works.

$$E\|Aw\|^2 = \text{tr}(V(Aw)) + \|E(Aw)\|^2 = E[Aww^\top A^\top] - E[Aw] \cdot E[Aw] = E[Aww^\top A^\top] = \text{tr}(A^\top A)$$

Since  $\text{tr}(A^\top A) = \text{tr}(AA^\top)$ , we have  $A \cdot E[ww^\top]A^\top = \text{tr}(AA^\top)$ .

$$\text{Therefore, } E[ww^\top] = I_k \Rightarrow \text{tr}(E[ww^\top]) = \text{tr}(I_k) = k \Rightarrow E[w^\top w] = k$$

A sufficient condition is that the squared norm of  $w$  be  $k$ , where  $k$  is the dimensionality of the vector.

- (b) What's a simple example of a distribution that satisfies the condition you derived above?

We can sample from a uniform distribution on the unit  $n$ -ball such that  $w$  is a unit vector.

- (c) Explain how you can put the relationship you found to practical use to estimate  $\|A\|_F^2$  for a large  $A$ . In particular, you must explain how to estimate  $\|A\|_F^2$  more or less accurately, depending on the need.

Sample some  $w$ 's uniformly on the unit  $n$ -ball and then compute  $\frac{1}{N} \sum_{i=1}^N k \|Aw\|^2$  where  $k$  is the dimensionality of vector  $w$ . Increasing the  $N$  will improve the accuracy of the estimation.

- (d) Test out the idea in Matlab. Generate a random matrix  $A$ , maybe  $500 \times 1000$ . Compute its Frobenius norm using `norm(A, 'fro')` command. Compare this to the result of your approach. Are they close? Is your approach faster?

With a large enough  $A$ , the estimate is very close. In one example, for  $N = 10000$ , the estimated squared Frobenius norm is 165,278.7 while the actual value is 166,551, leading to a .7639% error. The estimation approach is not faster than calling the Matlab Frobenius norm, but this method appears to scale quite nicely with very large matrices.

3. Consider again the logistic regression problem. Included with this homework is the covtype dataset (500K examples, 54 features).

Consider again the logistic regression formulation:

$$\min_{\theta} \frac{1}{N} \sum_{i=1}^N \log(1 + \exp(\tilde{x}_i^T \theta)) + \lambda \|\theta\|_2$$

where  $\tilde{x}_i = -y_i x_i$  and you can take  $\lambda = 0.01$  (small regularization).

Implement a stochastic gradient method for this problem.

Use the following options for step length:

- (a) Pre-specified constant
- (b) Decreasing with the rule  $\alpha(k) \propto \frac{1}{k}$  (with some initialization)
- (c) Decreasing with rule  $\alpha(k) \propto \frac{1}{k^{0.6}}$  (with some initialization)

Divide covtype into two datasets, 90% training and 10% testing. Tune each of the three previous step size routines (i.e. adjust the constant or the constant initialization) until you are happy each one performs reasonably well. Make a graph showing the value of the *test likelihood* as a function of the iterates for each of the three strategies.

4. (BONUS)

- (a) Change the counting in the previous problem to be as a function of *effective passes through the data*, rather than iterations. For example, five iterations with batch size 1 should be no different than one iteration with batch size 5 in this metric.

- (b) For the pre-specified constant step length strategy, compare test likelihood as a function of effective passes through the data for different random batch sizes, e.g. 1, 10, and 100.
- (c) Again for pre-specified constant step length strategy, implement a growing batch size strategy, where the size of the batch increases with iterations. Can this strategy beat the fixed batch size strategy, with respect to effective passes through the data?

```
''' Logistic regression using stochastic gradient descent with options for ste
'''
```

```
from scipy.io import loadmat
import numpy as np
import pandas as pd
import random
from numpy import log
from numpy import exp
from numpy import dot
from numpy.linalg import norm

data = loadmat('covtype.mat')
y = data['y']
X = data['X']
N, K = X.shape
X_tilde = compute_X_tilde(X, y)
X_train, X_test = data_split(X_tilde, p, N, K) #X_train and X_test are X_tilde

#class LogisticRegression:

def data_split(X, p, N, K):
    # check to see that 0 < p < 1
    n_train = int(float(p) * N)
    row_idx = random.shuffle([i for i in xrange(n_train)])
    X_train = X[row_idx[:n_train], :]
    X_test = X[row_idx[n_train:], :]
    return X_train, X_test

def sigmoid(x):
    return 1.0/(1 + exp(-x))

def compute_X_tilde(X, y):
    return -y*X

def SGD_step(x_tilde, theta, lamb):
    grad = np.zeros(K)
    LL = 0
    grad += x_tilde.dot(sigmoid(-x_tilde.dot(theta))) + lamb
    LL += log(1 + exp(x_tilde.dot(theta))) + lamb*norm(theta,2)
    return grad, LL
```

```

def SGD(X, theta, lamb, alpha, N, K):
    row_idx = random.shuffle([i for i in xrange(N)])
    for idx in row_idx:
        grad, LL = SGD_step(X[idx,:], theta, lamb)

def not_converged(a, b, eps):
    if max(abs(a-b)) > eps:
        return True
    return False

def main():
    # randomly initialize thetas
    theta_old = np.zeros(K)
    LL_list = list()

    while not_converged(theta_old, theta):
        grad, LL = SGD(X_tilde, theta_old, lamb, alpha, N, K)
        theta = theta_old + grad

```