

## Statistical Models and Computing Methods, Problem Set 1

October 15, 2020

Due 10/29/2020

### Problem 1.

(1) Show that  $X \sim \mathcal{N}(0, 1)$  is the maximum entropy distribution such that  $\mathbb{E}X = 0$  and  $\mathbb{E}X^2 = 1$ .

(2) Generalize the result in (1) for the maximum entropy distribution given the first  $k$  moments, *i.e.*,  $\mathbb{E}X^i = m_i$ ,  $i = 1, \dots, k$ .

### Problem 2.

Let  $Y_1, \dots, Y_n$  be a set of independent random variables with the following pdfs

$$p(y_i|\theta_i) = \exp(y_i b(\theta_i) + c(\theta_i) + d(y_i)), \quad i = 1, \dots, n$$

Let  $\mathbb{E}(Y_i) = \mu_i(\theta_i)$ ,  $g(\mu_i) = x_i^T \beta$ , where  $g$  is the link function and  $\beta \in \mathbb{R}^d$  is the vector of model parameters.

(1) Denote  $g(\mu_i)$  as  $\eta_i$ , and let  $s$  be the score function of  $\beta$ . Show that

$$s_j = \sum_{i=1}^n \frac{(y_i - \mu_i) x_{ij}}{\text{Var}(Y_i)} \frac{\partial \mu_i}{\partial \eta_i}, \quad j = 1, \dots, d$$

(2) Let  $\mathcal{I}$  be the Fisher information matrix. Show that

$$\mathcal{I}_{jk} = \mathbb{E}(s_j s_k) = \sum_{i=1}^n \frac{x_{ij} x_{ik}}{\text{Var}(Y_i)} \left( \frac{\partial \mu_i}{\partial \eta_i} \right)^2, \quad \forall 1 \leq j, k \leq d.$$

### Problem 3.

Use the following code to generate covariate matrices  $X$

```
1 import numpy as np
2 np.random.seed(1234)
3
4 n = 100
5 X = np.random.normal(size=(n, 2))
```

(1) Generate  $n = 100$  observations  $Y$  following the logistic regression model with true parameter  $\beta_0 = (-2, 1)$ .

(2) Find the MLE using the iteratively reweighted least square algorithm.

(3) Repeat (1) and (2) for 100 instances. Compare the MLEs with the asymptotical distribution  $\hat{\beta} \sim \mathcal{N}(\beta_0, \mathcal{I}^{-1}(\beta_0))$ . Present your result with a scatter plot for MLEs with contours for the pdf of the asymptotical distribution.

(4) Try the same for  $n = 10000$ . Does the asymptotical distribution provide a better fit to the MLEs? You can use the empirical covariance matrix of the MLEs for comparison.

**Problem 4.**

Consider the probit regression model

$$Y|X, \beta \sim \text{Bernoulli}(p), \quad p = \Phi(X\beta)$$

where  $\Phi$  is the cumulative distribution function of the standard normal distribution. Similarly as in Problem 3, generate a large covariate matrix  $X$  with 100000 instances and 100 features, and response  $Y$  with true parameter  $\beta_0$

```
1 import numpy as np
2 np.random.seed(1234)
3
4 n, d = 100000, 100
5 X = np.random.normal(size=(n,d))
6 beta_0 = np.random.normal(size=d)
```

- (1) Compare gradient descent and nesterov's accelerated gradient descent.
- (2) Compare vanilla stochastic gradient descent with different adaptive stochastic gradient descent methods, including AdaGrad, RMSprop, and Adam. Using minibatch sizes 32, 64, 128.
- (3) Bonus question. Generate a random mask matrix  $M$  as follows and use it to sparsify the covariance matrix  $X$

```
1 np.random.seed(1234)
2
3 sparse_rate = 0.3
4 M = np.random.uniform(size=(n,d)) < sparse_rate
5 X[M] = 0.
```

Repeat your experiments in (2), and compare with the results for the full covariance matrix.