

CSCI 5521: Fall'19

Introduction To Machine Learning

Homework 3

ARVIND RENGANATHAN renga016@umn.edu

3.

Error rates for MyLogisticReg2 with Boston50						
Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean	S.D
0.16	0.26	0.23	0.23	0.22	0.21	0.03

Error rates for MyLogisticReg2 with Boston75						
Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean	S.D
0.29	0.19	0.23	0.20	0.22	0.22	0.03

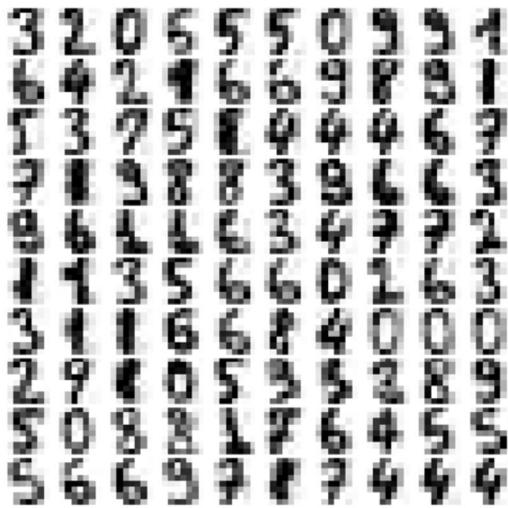
Error rates for Logistic regression with Boston50						
Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean	S.D
0.14	0.07	0.12	0.06	0.10	0.10	0.03

Error rates for Logistic regression with Boston50						
Fold 1	Fold 2	Fold 3	Fold 4	Fold 5	Mean	S.D
0.07	0.10	0.10	0.08	0.08	0.09	0.01

4.

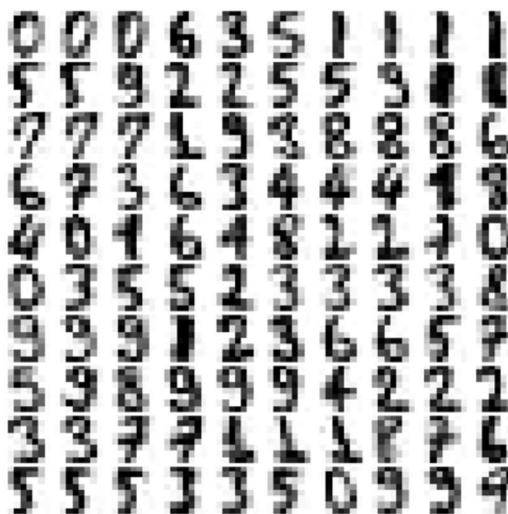
We found that for $\alpha = 90$, we found $d=21$ to be sufficient and for $\alpha = 99$ we found 41 to be sufficient which was expected. For $\alpha = 90$

```
In [11]: digits_new90 = np.dot(data_new90, proj_matrix90.T)
plot_digits(digits_new90)
```



For $\alpha = 99$

```
In [16]: digits_new99 = np.dot(data_new99, proj_matrix99.T)
plot_digits(digits_new99)
```



2. (ii) Let EM Algorithm for a Gaussian mixture model with Σ_K being diagonal, with diagonal elements of σ_{kj}^2

The algorithm alternates between 'E' & 'M' steps

E Step: Using current values of the parameters, compute the responsibilities of the components for data items, by applying Bayes' Rule

$$\pi_{ik} = P(\text{data item } i \text{ came from component } k | x_i)$$

$$= \frac{\pi_k N(x_i | \mu_k, \Sigma_k)}{\sum_{k'} \pi_{k'} N(x_i | \mu_{k'}, \Sigma_{k'})}$$

M Step: Using the current responsibilities, re-estimate the parameters, using weighted averages, with weights given by the responsibilities

$$\pi_k = \frac{1}{N} \sum_i \pi_{ik}, \quad \mu_k = \frac{\sum \pi_{ik} x_i}{\sum_i \pi_{ik}}, \quad \sigma_k^2 = \frac{\sum \pi_{ik} (x_i - \mu_k)^2}{\sum_i \pi_{ik}}$$

We start with initial guess at parameter values (or random), or perhaps with some initial guess at responsibilities (in which case we start with M). We alternate until little change.

b) M step.

$$\pi_{t_n} = \frac{1}{N} \sum_{n=1}^N \frac{N_n}{N}$$
$$M_n = \frac{1}{N_n} \times \cancel{\sum_i p(g_{t_n} | x_i)}$$

$$\pi_{t_n} = \frac{1}{N} \sum_i p(g_{t_n} | x_i)$$

$$M_n = \sum_i p(g_{t_n} | x_i) \cdot x_i / \sum_i p(g_{t_n} | x_i)$$

$$\xi_n = \sum_i p(g_{t_n} | x_i) \cdot (x_i - M_n)^2 / \sum_i p(g_{t_n} | x_i)$$

c)

$$p(g_{t_n} | x_i) = \frac{\pi_{t_n} N(x_i | M_n, \xi_n)}{\sum_{t_n'=1}^K \pi_{t_n'} N(x_i | M_{t_n'}, \xi_{t_n'})}$$

$$1. f(\omega) = \frac{1}{n} \sum_{i=1}^n \left\{ -y_i \omega^T x_i + \log(1 + \exp(\omega^T x_i))^2 \right\} + \frac{\lambda}{2} \|\omega\|^2$$

$$f' = \frac{\partial f}{\partial \omega} = \frac{1}{n} \sum_{i=1}^n \left(-y_i x_i + \frac{e^{\omega^T x_i}}{1 + e^{\omega^T x_i}} x_i \right) + \lambda \omega$$

$$\begin{aligned} f'' = \frac{\partial^2 f}{\partial \omega^2} &= \frac{1}{n} \sum_{i=1}^n \left\{ \frac{(1 + e^{\omega^T x_i})(e^{\omega^T x_i})x_i - e^{\omega^T x_i}e^{\omega^T x_i}x_i^2}{(1 + e^{\omega^T x_i})^2} \right\} + \lambda \\ &= \frac{1}{n} \sum_{i=1}^n \left(\frac{e^{\omega^T x_i}}{(1 + e^{\omega^T x_i})^2} x_i^2 \right) + \lambda \end{aligned}$$

(i) Projected gradient descent \rightarrow

$$\omega_{t+1} = \omega_t - \eta \frac{\partial f(\omega_t)}{\partial \omega_t} \quad \begin{matrix} t \rightarrow t^{\text{th}} \text{ epoch} \\ t+1 \rightarrow (t+1)^{\text{th}} \text{ epoch} \end{matrix}$$

$$\therefore \omega_{t+1} = \begin{cases} \omega_{t+1} & \text{if } \|\omega_{t+1}\|_2 \leq c \\ \frac{c}{\|\omega_t\|_2} \omega_{t+1} & \text{if } \|\omega_{t+1}\|_2 > c \end{cases}$$

The above two conditions are for keeping $\|\omega\| \leq c$
in each epoch.

Gradient descent: iterate until change $< \epsilon$

(b) To show that LR objective is convex, we consider the partial derivatives.

We have to minimize $-\log P(\vec{y} | \vec{x}, \vec{w}) = \sum \log \frac{(1 + \exp(-\vec{y}_i \vec{w}^T \vec{x}_i))}{(1 - \vec{y}_i \vec{w}^T \vec{x}_i)}$

$$\text{Define } g(z) = \frac{1}{1 + e^{-z}} \quad 1 - g(z) = \frac{e^{-z}}{1 + e^{-z}}$$

$$\frac{\partial g(z)}{\partial z} = -g(z)(1 - g(z))$$

$$\frac{\partial \log P(\vec{y} | \vec{x}, \vec{w})}{\partial w_j} = - \sum_{i=1}^n y_i x_{ij} (1 - g(y_i \vec{w}^T \vec{x}_i))$$

$$\frac{\partial^2 \log P(\vec{y} | \vec{x}, \vec{w})}{\partial w_j \partial w_k} = \sum_{i=1}^n y_i^2 x_{ij} x_{ik} g(y_i \vec{w}^T \vec{x}_i)(1 - g(y_i \vec{w}^T \vec{x}_i))$$

To show objective function is convex, we first show hessian is positive semi definite.

A Matrix M is PSD iff $\vec{a}^T M \vec{a} \geq 0$ for all vectors \vec{a} .

Let ∇^2 be hessian for our objective.

Define $P_{ij} = g(y_i \vec{w}^T \vec{x}_i)(1 - g(y_i \vec{w}^T \vec{x}_i))$ & $P_{ii} = P_{ij} \sqrt{p_i}$

Then

$$\vec{a}^T \nabla^2 \vec{a} = \sum_{i=1}^n \sum_{j=1}^d \sum_{k=1}^d a_j a_k x_{ij} x_{ik} p_i,$$
$$= \sum_{i=1}^n \vec{a}^T \vec{p}_i \vec{p}_i^T \vec{a} \geq 0$$

$$\text{as } \vec{a}^T \vec{p}_i \vec{p}_i^T \vec{a} = (\vec{a}^T \vec{p}_i)^2 \geq 0$$

hence, the hessian is PSD. so objective function is convex.

on adding L_2 regularizer $\frac{\lambda}{2} (||w||^2)$ or $\frac{\lambda}{2} \vec{w}^T \vec{w}$ to the objective, the hessian is positive definite & hence objective function is strictly convex. also. it is strictly convex & differentiable

d) The bound is defined as

$$f(w_T) - f(w^*) \leq \frac{\beta}{2} \exp\left(\frac{-4T}{\frac{\beta}{2} + 1}\right) ||w_1 - w^*||^2$$

Here β as is in 1.c) It is β smooth upper bound.
 ~~λ is in 1.b) It is~~ λ is a strongly convex lower bound.

T is step number T of the projected gradient descent algorithm

w_1 is the w initialized at the start of projected gradient descent algorithm

w^* is the optimum value for w

c) A function $g: \mathbb{R}^p \rightarrow \mathbb{R}$ is L -smooth if and only if it is differentiable & is L -Lipschitz continuous

$$\forall \theta_1, \theta_2 \in \mathbb{R}^p \quad \|g'(\theta_1) - g'(\theta_2)\| \leq L \|\theta_1 - \theta_2\|$$

$$f(w) \leq f(w') + \langle \nabla f(w'), w - w' \rangle + \frac{\beta}{2} \|w - w'\|^2$$

$$\forall w, w' \in W$$

All convex quadratic & logistic is smooth.

for smoothness $\gamma_i = 0$

it is enough to check that

$$g: \mathbb{R} \rightarrow \mathbb{R} \quad g(+)=\log(1+\exp(+))$$

has Lipschitz gradient & it does because its second derivative is bounded

$$\nabla f(\beta) = -g'(h(\beta))x_i^T, \quad h(\beta) = x_i^\top \beta$$

& it is twice differentiable

$$\therefore 0 \leq f''(\theta) \leq \beta I,$$