

Homework 5, Fall 2024

Mohammad Alshurbaji

11/27/2024

Problem 1: Vanilla Logistic Regression for multi-classification:

- k classes. - The input $x \in \mathbb{R}^d$

- Prob. to each class is:

$$P(Y=k | X=x) = \frac{\exp(w_k^T x)}{1 + \sum_{l=1}^{k-1} \exp(w_l^T x)}, \text{ for } k=1, 2, \dots, k$$

$$P(Y=k | X=x) = \frac{1}{1 + \sum_{l=1}^{k-1} \exp(w_l^T x)}$$

If we define $w_k = 0$, Then:

$$P(Y=k | X=x) = \frac{\exp(w_k^T x)}{1 + \sum_{l=1}^{k-1} \exp(w_l^T x)}, k=1, \dots, k$$

Question 1: What and how many para. are there to be opt.?

1. Weight vectors (w_k)

→ For each class k , there is $w_k \in \mathbb{R}^d$.

→ assuming $w_k = 0$.

→ means: The # of weight para.:

$$d \times (k-1)$$

2. Bias Parameters: b_k

→ one bias for each class: $(k-1)$

∴ Total Parameters: $d \times (k-1) + (k-1)$

$$= (d+1)(k-1) \#$$

Question 2 $L(w_1, \dots, w_{k-1}) = \sum_{i=1}^n \ln P(Y=y_i | X=x_i)$

$$= \sum_{i=1}^n \ln \left(\frac{\exp(w_k^T x)}{1 + \sum_{l=1}^{k-1} \exp(w_l^T x)} \right)$$

$$= \sum_{i=1}^n \left(\ln(\exp(w_k^T x)) - \ln \left(1 + \sum_{l=1}^{k-1} \exp(w_l^T x) \right) \right)$$

$$= \underbrace{\sum_{i=1}^n w_k^T x}_I - \underbrace{\sum_{i=1}^n \ln \left(1 + \sum_{l=1}^{k-1} \exp(w_l^T x) \right)}_{II}$$

Question 3 The gradient of L w.r.t. w_k :

$$\frac{\partial L}{\partial w_k} = \sum_{i=1}^n$$

I: $\frac{\partial I}{\partial w_k} = x$

II: $\frac{\partial II}{\partial w_k} = \frac{\exp(w_k^T x) x}{1 + \sum_{l=1}^{k-1} \exp(w_l^T x)}$

I + II: $\frac{\partial L}{\partial w_k} = \sum_{i=1}^n \left(x - \left(\frac{\exp(w_k^T x) x}{1 + \sum_{l=1}^{k-1} \exp(w_l^T x)} \right) \right)$

$$= \sum_{i=1}^n \left(x - P(Y=k | X=x_i) x \right)$$

$$= \sum_{i=1}^n \left(\left(\frac{1}{Y=k} P(Y=k | X=x_i) \right) x \right) \quad \# \text{ (In terms of Probability)}$$

Question 4: With adding the regularization term:

$$f(w_1, \dots, w_{k-1}) = \underbrace{L(w_1, \dots, w_{k-1})}_I - \underbrace{\frac{\lambda}{2} \sum_{l=1}^{k-1} \|w_l\|_2^2}_II$$

$$\frac{\partial f}{\partial w_k} :$$

$$I : \frac{\partial I}{\partial w_k} = \sum_{i=1}^n \left((1 - P(Y=k | X=x_i)) X \right)$$

$$II : \frac{\partial II}{\partial w_k} = -\lambda w_k$$

$$\therefore I + II : \frac{\partial f}{\partial w_k} = \sum_{i=1}^n \left((1 - P(Y=k | X=x_i)) X \right) - \lambda w_k$$

#

Question 5: USPS Handwritten Recognition digit dataset.

→ Image size: 16 x 16

→ For each digit (i.e. 0, 1, ..., 9) there are 600 training samples. 500 testing ones.

(a)

→ Continue in the next pages.

Question 4:

Log_grad.m code:”

```
function G=log_grad(y, X, B)
```

```
    [n,d] = size(X);%n: number of samples, d: number of features
```

```
    K = size(B,2) + 1; %Total number of classes
```

```
%compute gradient
```

```
    XB = X * B;
```

```
    expXB = exp(XB);
```

```
    prob = expXB ./ (1 + sum(expXB, 2));
```

```
    prob = [prob, 1 - sum(prob, 2)];
```

```
    G = zeros(d,K-1);
```

```
    for k = 1:K-1
```

```
        indicator = (y == k); % Indicator vector for class k
```

```
        G(:, k) = X' * (indicator - prob(:, k)); % Gradient for class k
```

```
    end
```

```
end
```

```
“
```

```

function G=log_grad(y, X, B)
    [n,d] = size(X);%n: number of samples, d: number of features
    K = size(B,2) + 1; %Total number of classes

    %compute gradient
    XB = X * B;
    expXB = exp(XB);
    prob = expXB ./ (1 + sum(expXB, 2));

    prob = [prob, 1 - sum(prob, 2)];

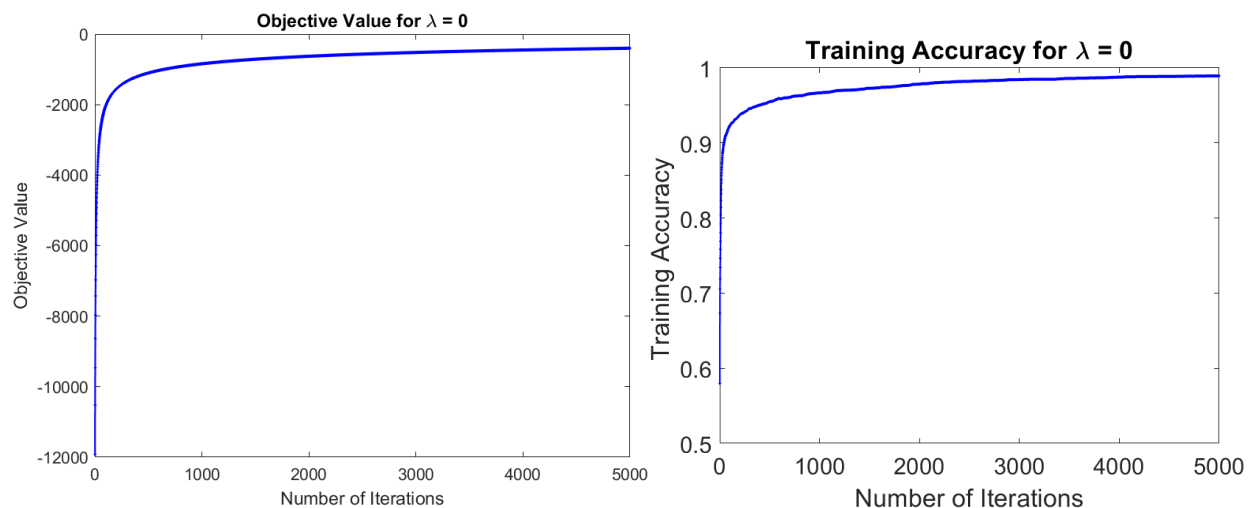
    G = zeros(d,K-1);
    for k = 1:K-1
        indicator = (y == k); % Indicator vector for class k
        G(:, k) = X' * (indicator - prob(:, k)); % Gradient for class k
    end

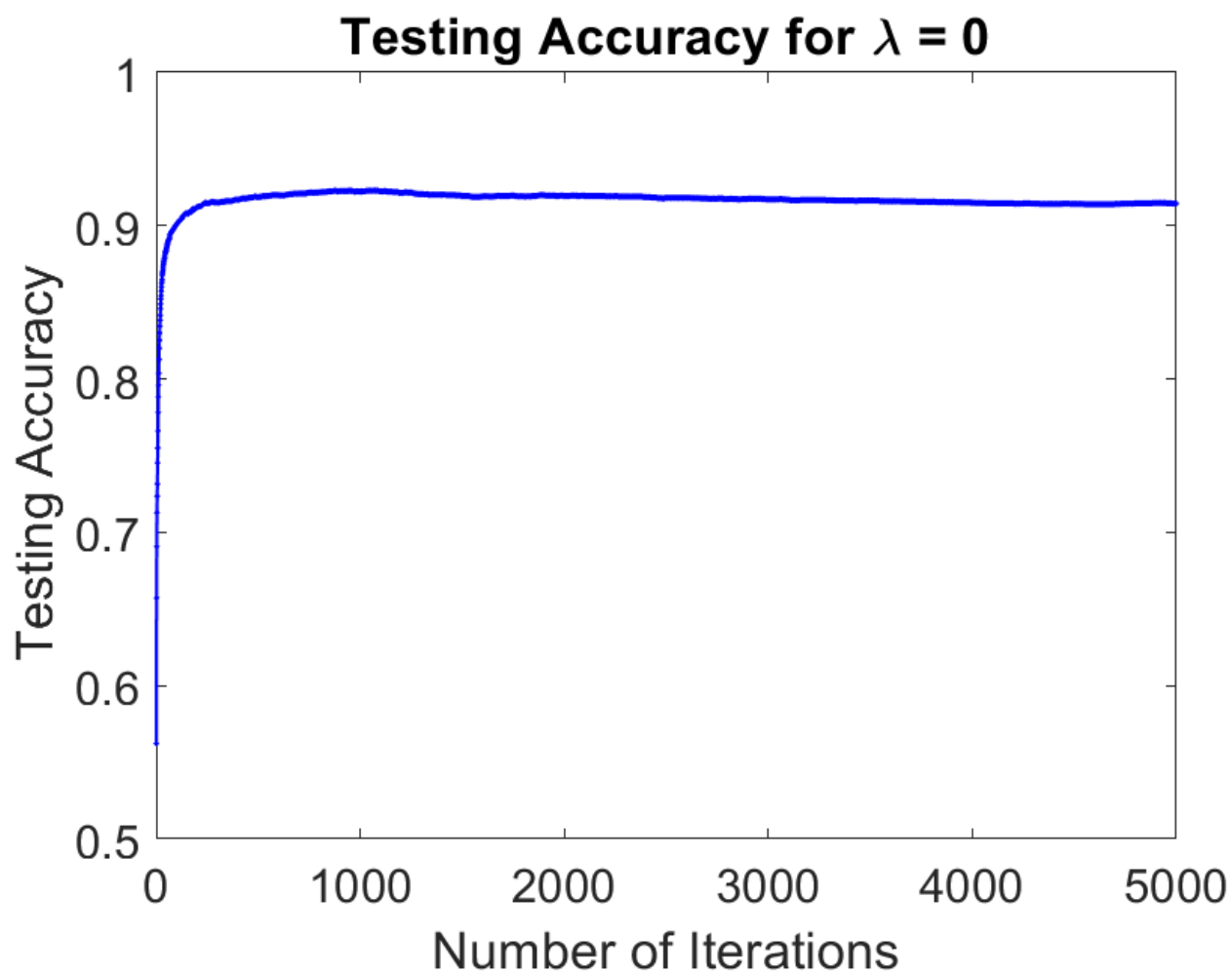
end

```

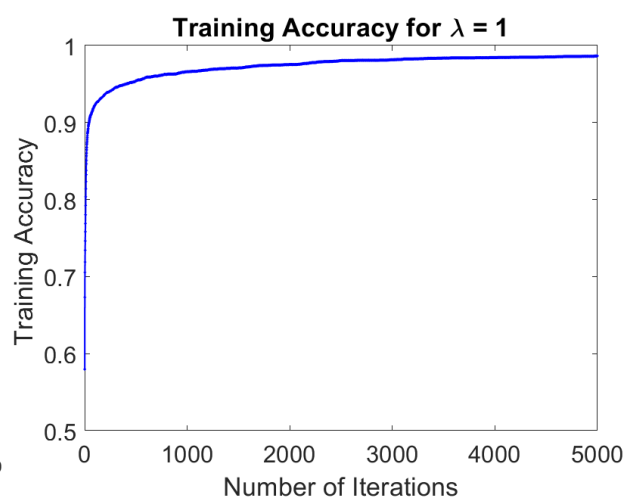
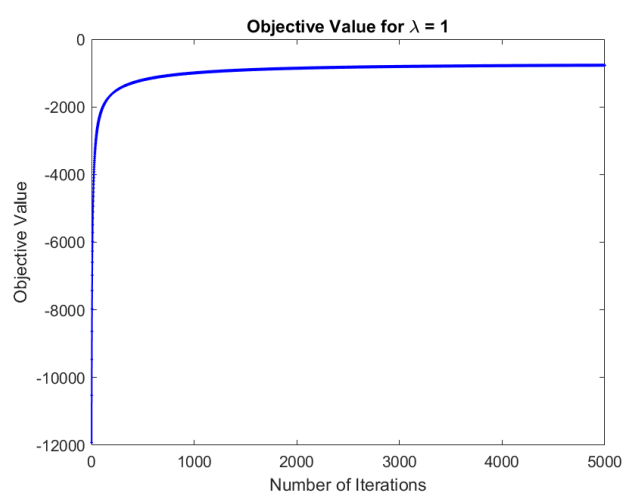
Following this code and implementing the logistic_classify .m and generate these graphs:

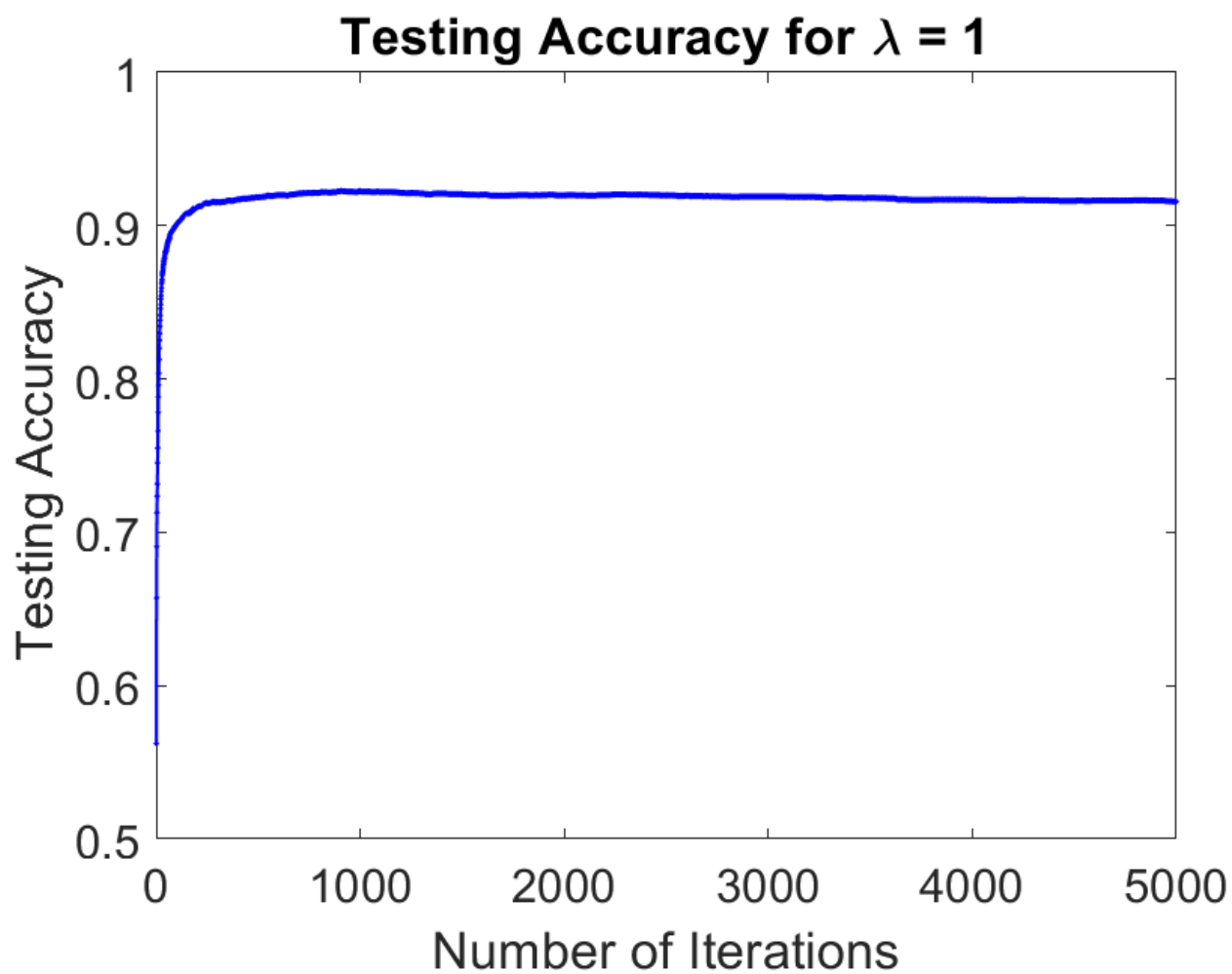
When lambda = 0



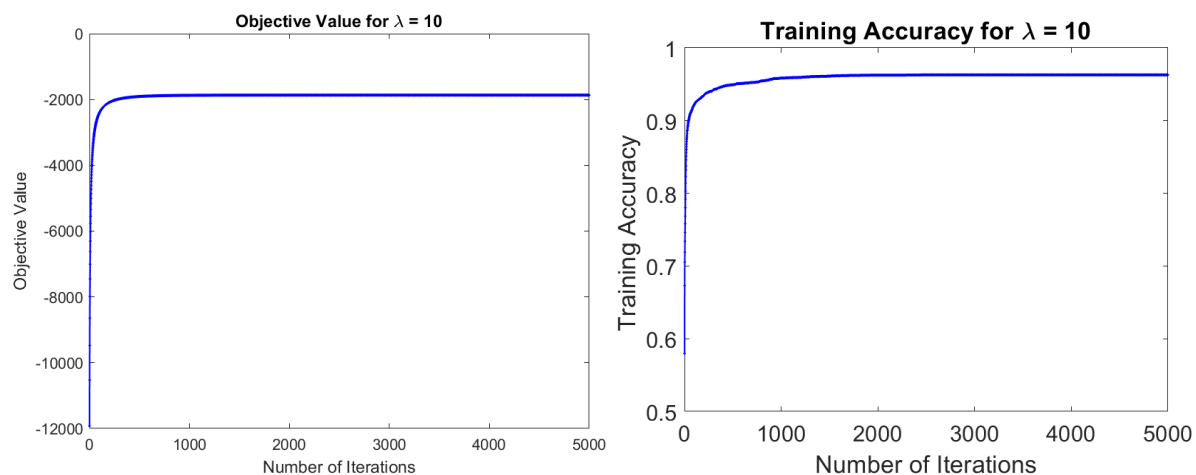


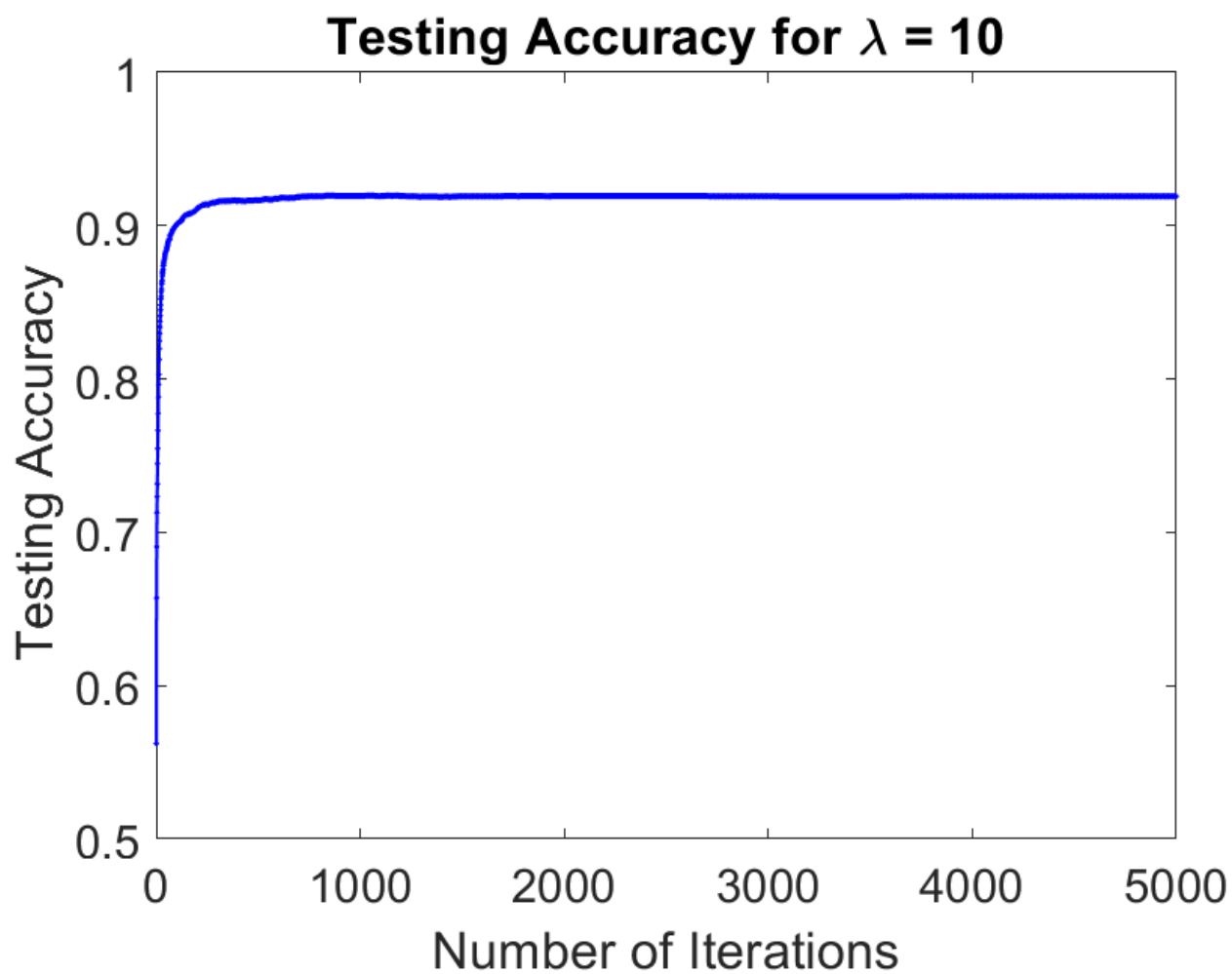
When lambda = 1



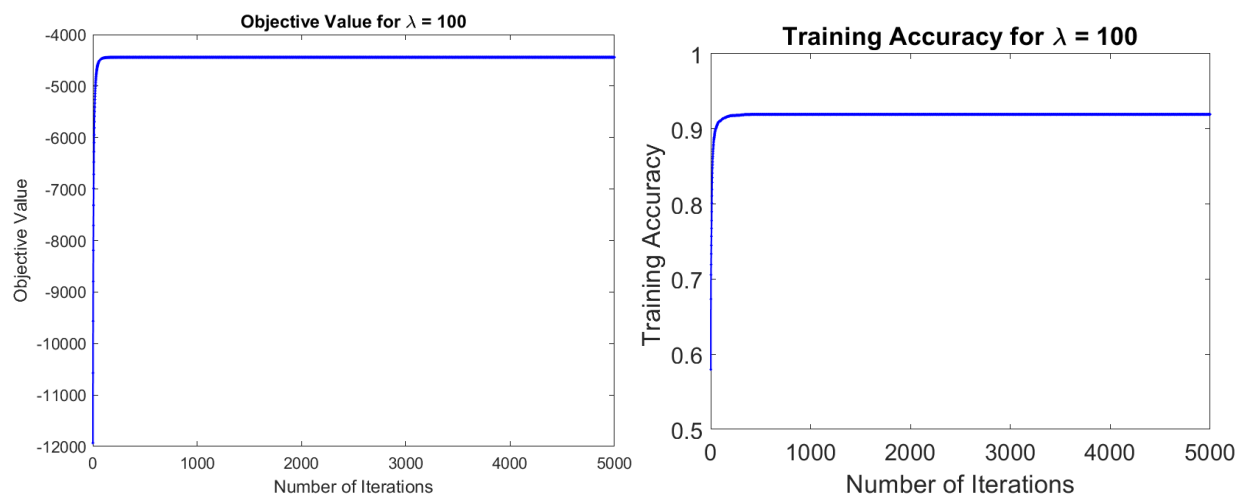


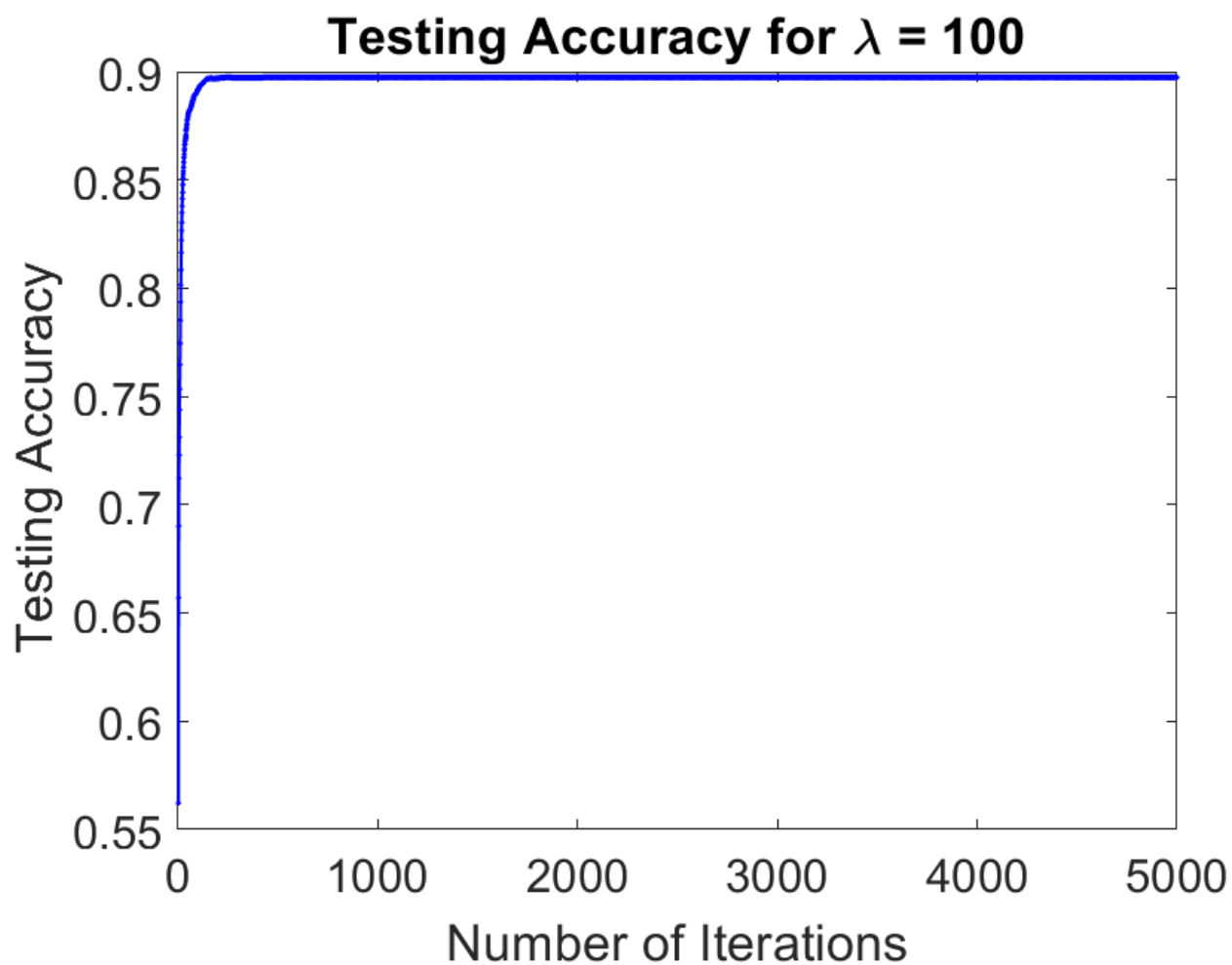
When $\lambda = 10$



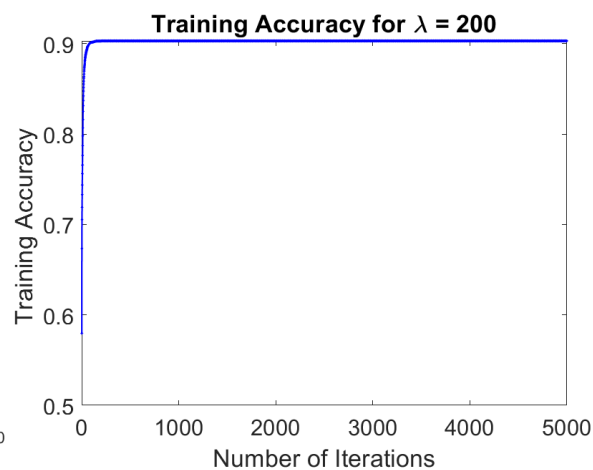
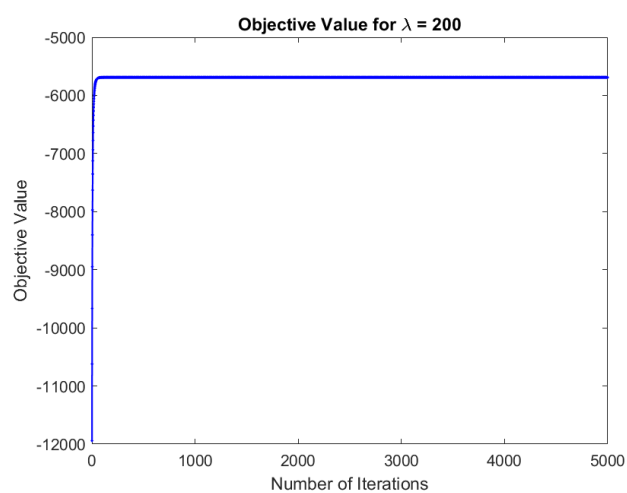


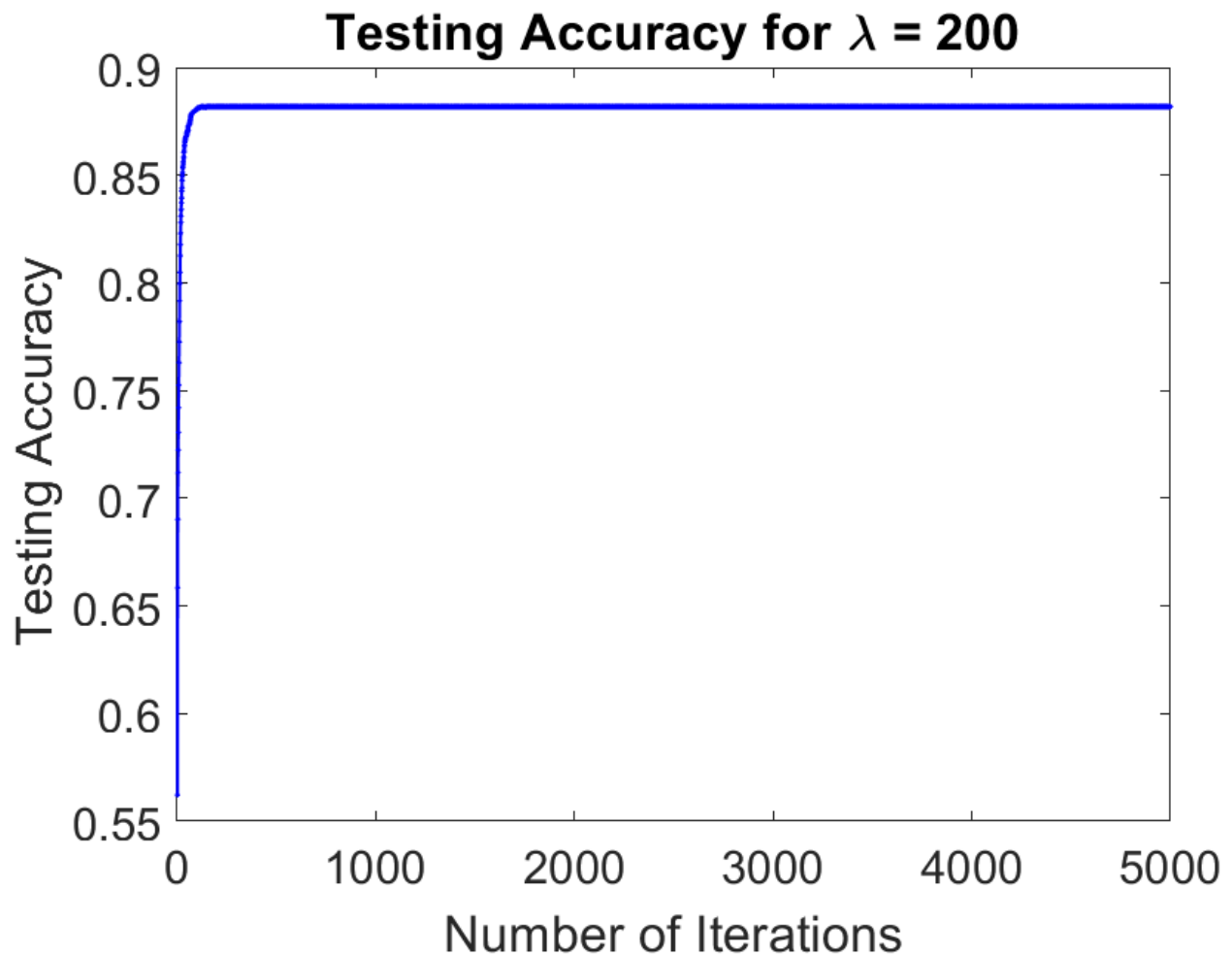
When lambda =100





When lambda =200





From the above charts we recognize that:

- Objective value:
 - In general the obj value decreases over iterations, showing convergence for all lambda values.
 - While for smaller lambda values ($\lambda = 0$), the convergence is faster and reach a lower final value.
 - And for larger values, the convergence is slower and the final obj is higher.
 - To wrap up: a moderate lambda balances fast convergence and regularization.
- Training Accuracy:
 - For smaller lambda values, training accuracy is higher as the model overfits to the training data. While for larger values, training accuracy decreases because the model becomes over regularized.
 - To wrap up, moderate lambda has in between better training accuracy.
- Testing Accuracy:
 - This improves with moderate values as the model gernalization better to unseen data.

- While for small values, testing accuracy is lower since the model overfits to the training sets.
- And for larger values, it declines again due to underfitting since the model becomes too simplistic.