

# CPSC 540: Assignment 2

Yiwei Hou(84435156)  
Xiaomeng Ju(86475150)  
Tingting Yu(74439118)

## 1 Convergence Rates

### 1.1 Gradient Descent

1. If  $f$  is a  $C^1$  function, strong convexity means  $f(y) \geq f(x) + \langle f'(x), y - x \rangle + \frac{u}{2} \|x - y\|^2$  for  $u > 0$  and  $y, x$  in its domain. Applying this inequality at  $x^t$  and  $x^*$ , we have  $f(x^t) \geq f(x^*) + \langle f'(x^*), x^t - x^* \rangle + \frac{u}{2} \|x^t - x^*\|^2$ . Since  $f'(x^*) = 0$ ,  $\frac{u}{2} \|x^t - x^*\|^2 \leq f(x^t) - f(x^*) = O(\rho^t)$ . This means the iteration  $\|x^t - x^*\|$  has a convergence rate of  $O(\rho^{t/2})$ .

If  $f$  is a  $C^2$  function, by Taylor's theorem, there is a  $z$  such that  $f(x^t) = f(x^*) + \nabla f(x^*)(x^t - x^*) + \frac{1}{2}(x^t - x^*)^T \nabla^2 f(z)(x^t - x^*)$ . By strong convexity,  $f(x^*) \geq f(x^t) + \nabla f(x^*)(x^t - x^*) + \frac{u}{2} \|x^t - x^*\|^2$ . Again, as  $\nabla f(x^*) = 0$ , we have  $\frac{u}{2} \|x^t - x^*\|^2 \leq f(x^t) - f(x^*) = O(\rho^t)$ , which means the iteration has a convergence rate of  $O(\rho^{t/2})$ .

2. From the descent lemma, we know  $f(x^{t+1}) \leq f(x^t) + \nabla f(x^t)^T(x^{t+1} - x^t) + \frac{L}{2} \|x^{t+1} - x^t\|^2$ . With the constant step size,  $x^{t+1} - x^t = \alpha \nabla f(x^t)$ . Plugging this into the descent lemma, we get

$$\begin{aligned} f(x^{t+1}) &\leq f(x^t) - \alpha \nabla f(x^t)^T \nabla f(x^t) + \frac{L\alpha^2}{2} \|\nabla f(x^t)\|^2 \\ &= f(x^t) - (\alpha - \frac{L\alpha^2}{2}) \|\nabla f(x^t)\|^2 \end{aligned}$$

Comparing the upper bound with the proof in class, we define  $\frac{1}{2\tilde{L}} = \alpha - \frac{L\alpha^2}{2}$  and only need to show  $\tilde{L} \geq L$  in order to preserve the linear convergence rate. Since  $\alpha < \frac{L}{2}$ , the minimum value of  $\tilde{L}$  is  $L$ . Thus, we show that the convergence rate is linear when using a constant step size.

3. By descent lemma, we have

$$\begin{aligned} f(x^{t+1}) &\leq f(x^t) - \frac{1}{L^t} \nabla f(x^t)^T \nabla f(x^t) + \frac{L}{2L^{t2}} \|\nabla f(x^t)\|^2 \\ &= f(x^t) - (\frac{1}{L^t} - \frac{L}{2L^{t2}}) \|\nabla f(x^t)\|^2 \end{aligned}$$

To satisfy the inequality, we need  $\frac{1}{L^t} - \frac{L}{2L^{t2}} \geq \frac{1}{2L^t}$  which gives us  $L^t \geq L$ . Since we start with some  $L^0$  that is smaller than  $L$  and double the value when the inequality is not

satisfied,  $L^t \leq 2L$  always holds. Following the steps in class, we have  $f(x^t) - f(x^*) \leq (1 - \frac{u}{L^0})(1 - \frac{u}{L^1}) \cdots (1 - \frac{1}{L^{t-1}})(f(x^0) - f(x^*)) \leq (1 - \frac{u}{2L})^t (f(x^0) - f(x^*))$ .

4. When  $L^t > L$  for any  $t$ , the convergence rate would be faster than  $\rho = 1 - \frac{u}{L}$ .

## 1.2 Sign-Based Gradient Descent

1. By the Lipschitz continuity in  $\infty$ -norm, we have

$$\begin{aligned}
f(x^{t+1}) &\leq f(x^t) + \nabla f(x^t)^T (x^{t+1} - x^t) + \frac{L_\infty}{2} \|x^{t+1} - x^t\|_\infty^2 \\
&= f(x^t) + \nabla f(x^t)^T \left( -\frac{\|\nabla f(x^t)\|_1}{L_\infty} \text{sign}(\nabla f(x^t)) \right) + \frac{L_\infty}{2} \left\| -\frac{\|\nabla f(x^t)\|_1}{L_\infty} \text{sign}(\nabla f(x^t)) \right\|_\infty^2 \\
&= f(x^t) - \frac{1}{L_\infty} \|\nabla f(x^t)\|_1^2 + \frac{1}{2L_\infty} \|\nabla f(x^t)\|_1^2 \\
&= f(x^t) - \frac{1}{2L_\infty} \|\nabla f(x^t)\|_1^2 \\
&\leq f(x^t) - \frac{1}{2L_\infty} \|\nabla f(x^t)\|_2^2
\end{aligned}$$

By the strong convexity, we have

$$-\|\nabla f(x^t)\|_2^2 \leq -2\mu(f(x^t) - f(x^*)).$$

Therefore, we get  $f(x^{t+1}) - f(x^*) \leq (1 - \frac{u}{L_\infty})(f(x^t) - f(x^*))$ .

2. If  $f$  is  $L_\infty$  Lipschitz continuous,

$$\begin{aligned}
\|f'(x) - f'(y)\|_2 &\leq \|f'(x) - f'(y)\|_1 \\
&\leq L_\infty \|x - y\|_\infty \\
&\leq L_\infty \|x - y\|_2
\end{aligned}$$

As  $L$  is the minimum value of the Lipschitz  $L_2$ -norm constant, we have  $L \leq L_\infty$ . Similarly, if  $f$  is  $L_2$  Lipschitz continuous, we have

$$\begin{aligned}
\|f'(x) - f'(y)\|_1 &\leq \sqrt{d} \|f'(x) - f'(y)\|_2 \\
&\leq \sqrt{d} L \|x - y\|_2 \\
&\leq \sqrt{d} L \sqrt{d} \|x - y\|_\infty \\
&\leq dL \|x - y\|_\infty
\end{aligned}$$

As  $L_\infty$  is the minimum value of the Lipschitz  $L_\infty$ -norm constant, we have  $L_\infty \leq dL$ . In conclusion,  $L \leq L_\infty \leq dL$ .

### 1.3 Block Coordinate Descent

1. By block-wise strong-smoothness and  $p(b_t = b) = \frac{1}{k}$ , we have

$$\begin{aligned}
f(x^{t+1}) &\leq f(x^t) - \frac{1}{L} \nabla_{b_t} f(x^t)^T (\nabla f(x^t) \circ e_{b_t}) + \frac{1}{2L} \|\nabla f(x^t) \circ e_{b_t}\|^2 \\
&= f(x^t) - \frac{1}{2L} \|\nabla f(x^t) \circ e_{b_t}\|^2 \\
\mathbf{E}[f(x^{t+1})] &\leq \mathbf{E}[f(x^t) - \frac{1}{2L} \|\nabla f(x^t) \circ e_{b_t}\|^2] \\
&= \frac{1}{k} \sum_{b=1}^k (f(x^t) - \frac{1}{2L} \|\nabla f(x^t) \circ e_{b_t}\|^2) \\
&= f(x^t) - \frac{1}{2kL} \sum_{b=1}^k \|\nabla f(x^t) \circ e_{b_t}\|^2 \\
&= f(x^t) - \frac{1}{2kL} \|\nabla f(x^t)\|^2
\end{aligned}$$

From strong convexity, we can re-write the inequality into:

$$\begin{aligned}
\mathbf{E}(f(x^{t+1})) - f(x^*) &\leq f(x^t) - f(x^*) - \frac{u}{kL} [f(x^t) - f(x^*)] \\
&= (1 - \frac{u}{kL}) [f(x^t) - f(x^*)]
\end{aligned}$$

2. Assuming that each block  $b$  has its own strong-smoothness constant  $L_b$  and the blocks are sampled proportional to  $L_{b_t}$ , we have

$$\begin{aligned}
f(x^{t+1}) &\leq f(x^t) - \frac{1}{L_{b_t}} \nabla_{b_t} f(x^t)^T (\nabla f(x^t) \circ e_{b_t}) + \frac{1}{2L_{b_t}} \|\nabla f(x^t) \circ e_{b_t}\|^2 \\
&= f(x^t) - \frac{1}{2L_{b_t}} \|\nabla f(x^t) \circ e_{b_t}\|^2 \\
\mathbf{E}[f(x^{t+1})] &\leq \mathbf{E}[f(x^t) - \frac{1}{2L_{b_t}} \|\nabla f(x^t) \circ e_{b_t}\|^2] \\
&= \sum_{b=1}^k \frac{L_{b_t}}{\sum L_{b_t}} (f(x^t) - \frac{1}{2L_{b_t}} \|\nabla f(x^t) \circ e_{b_t}\|^2) \\
&= f(x^t) - \sum_{b=1}^k \frac{1}{2 \sum L_{b_t}} \|\nabla f(x^t) \circ e_{b_t}\|^2 \\
&\leq f(x^t) - \frac{u}{\sum L_{b_t}} [f(x^t) - f(x^*)] \\
&= (1 - \frac{u}{\sum L_{b_t}}) [f(x^t) - f(x^*)]
\end{aligned}$$

The second last line is from applying the strong convexity block-wise. Since  $L = \max_{b_t} \{L_{b_t}\}$ ,  $L_{b_t} \leq L$ ,  $\sum L_{b_t} \leq kL$ . If there exists some  $L_{b_t} \neq L$ ,  $1 - \frac{u}{\sum L_{b_t}} \leq 1 - \frac{u}{kL}$ , meaning we have a faster convergence rate.

## 2 Large-Scale Algorithms

### 2.1 Coordinate Optimization

1. The final function value is 1.4724e+02 and the total time is 2.328713 seconds.

```

1  function [model] = logisticL2(X,y,lambda)
2
3  % Add bias variable
4  [n,d] = size(X);
5  X = [ones(n,1) X];
6  d = d+1;
7
8  % Initial values of regression parameters
9  w = zeros(d,1);
10
11 % Optimizaion parameters
12 maxPasses = 500;
13 progTol = 1e-4;
14 L = .25*max(sum(X.^2, 1))+ lambda;
15
16 w_old = w;
17
18 Xw = X*w;
19
20 for t = 1:maxPasses*d
21
22     % Choose variable to update 'j'
23     j = randi(d);
24
25     % Compute partial derivative 'g-j'
26
27     yXw = y.*(Xw);
28     sigmoid = 1./(1+exp(-yXw));
29     g-j = -(y.*(1-sigmoid))*X(:,j) + lambda*w(j);
30
31     % Update variable
32     w(j) = w(j) - (1/L)*g-j;
33     Xw = Xw - (1/L)*X(:,j)*g-j;
34
35     % Check for lack of progress after each "pass"
36     if mod(t,d) == 0
37         change = norm(w-w_old,inf);
38         fprintf('Passes = %d, function = %.4e, change = %.4f\n',t/d,
39             logisticL2_loss(w,X,y,lambda),change);
39         if change < progTol
40             fprintf('Parameters changed by less than progTol on pass\n');
41             break;

```

```

42         end
43         w_old = w;
44     end
45 end
46
47
48 model.w = w;
49 model.predict = @predict;
50 end

```

2. The final function value is 1.4418e+02 and the total number of passes is 153.

```

1  function [model] = logisticL2(X,y,lambda)
2
3  % Add bias variable
4  [n,d] = size(X);
5  X = [ones(n,1) X];
6  d = d+1;
7
8  % Initial values of regression parameters
9  w = zeros(d,1);
10
11 % Optimizaion parameters
12 maxPasses = 500;
13 progTol = 1e-4;
14 L_candidate = .25*(sum(X.^2, 1))+ lambda;
15 w_old = w;
16
17 Xw = X*w;
18
19 for t = 1:maxPasses*d
20
21     % Choose variable to update 'j'
22     %j = randi(d);
23
24     p = L_candidate/(sum(L_candidate));
25     j = sampleDiscrete(p);
26     L = L_candidate(j);
27
28     % Compute partial derivative 'g-j'
29
30     yXw = y.*(Xw);
31     sigmoid = 1./(1+exp(-yXw));
32     g_j = -(y.*(1-sigmoid))'*X(:,j) + lambda*w(j);
33
34     %g-j = g(j);
35
36     % Update variable
37     w(j) = w(j) - (1/L)*g_j;
38     Xw = Xw - (1/L)*X(:,j)*g_j;
39
40     % Check for lack of progress after each "pass"
41     if mod(t,d) == 0
42         change = norm(w-w_old, inf);
43         fprintf('Passes = %d, function = %.4e, change = %.4f\n',t/d,

```

```

        logisticL2_loss(w,X,y,lambda),change);
44     if change < progTol
45         fprintf('Parameters changed by less than progTol on pass\n');
46         break;
47     end
48     w_old = w;
49 end
50 end
51
52
53 model.w = w;
54 model.predict = @predict;
55 end

```

3. The final function value is 1.4091e+02 and the total number of passes is 149.

```

1  function [model] = logisticL2(X,y,lambda)
2
3  % Add bias variable
4  [n,d] = size(X);
5  X = [ones(n,1) X];
6  d = d+1;
7
8  % Initial values of regression parameters
9  w = zeros(d,1);
10
11 % Optimizaion parameters
12 maxPasses = 500;
13 progTol = 1e-4;
14 L_candidate = .25*(sum(X.^2, 1))+ lambda;
15 w_old = w;
16
17 Xw = X*w;
18
19 for t = 1:maxPasses*d
20
21     % Choose variable to update 'j'
22     j = randi(d);
23     L = L_candidate(j);
24
25     % Compute partial derivative 'g-j'
26
27     yXw = y.*(Xw);
28     sigmoid = 1./(1+exp(-yXw));
29     g-j = -(y.*(1-sigmoid))*X(:,j) + lambda*w(j);
30
31     % Update variable
32     w(j) = w(j) - (1/L)*g-j;
33     Xw = Xw - (1/L)*X(:,j)*g-j;
34
35     % Check for lack of progress after each "pass"
36     if mod(t,d) == 0
37         change = norm(w-w_old,inf);
38         fprintf('Passes = %d, function = %.4e, change = %.4f\n',t/d,
            logisticL2_loss(w,X,y,lambda),change);

```

```

39         if change < progTol
40             fprintf('Parameters changed by less than progTol on pass\n');
41             break;
42         end
43         w_old = w;
44     end
45 end
46
47
48 model.w = w;
49 model.predict = @predict;
50 end

```

4. We stated in class that  $L_{j_t}$  is used as the step size and sampling  $j_t$  proportional to  $L_{j_t}$  gives

$$E[f(x^t)] - f(x^*) \leq (1 - \frac{\mu}{d\bar{L}})[f(x^0) - f(x^*)].$$

By using the uniform sampling,  $E[\bar{L}]_{Uniform} = \frac{1}{d} \sum_{i=1}^d L_i$

By using the Lipschitz sampling,  $E[\bar{L}]_{Lipschitz} = \sum_{i=1}^d p_i L_i = \frac{1}{\sum_i L_i} \sum_{i=1}^d L_i^2$

To prove that uniform sampling has a tighter bound than Lipschitz sampling, we need to show that

$$E[\bar{L}]_{Uniform} \leq E[\bar{L}]_{Lipschitz},$$

which is equivalent to showing

$$\frac{1}{d} \left( \sum_{i=1}^d L_i \right)^2 \leq \sum_{i=1}^d L_i^2.$$

According to the CauchySchwarz inequality, the above inequality holds and thus uniform sampling outperforms Lipschitz sampling.

## 2.2 Proximal-Gradient

1. errTrain = 0.0060; errTest = 0.2740

The number of non-zero parameters: 500

The number of original features that the model uses: 100

```

1 function [model] = softmaxClassifierL2(X,y,lambda)
2
3 % Compute sizes
4 [n,d] = size(X);
5 k = max(y);
6
7 W = zeros(d,k); % Each column is a classifier
8 W(:) = findMin(@softmaxLoss,W(:),500,1,X,y,k,lambda);
9
10 model.W = W;
11 model.predict = @predict;

```

```

12 end
13
14 function [yhat] = predict(model,X)
15 W = model.W;
16 [~,yhat] = max(X*W,[],2);
17 end
18
19 function [nll,g,H] = softmaxLoss(w,X,y,k,lambda)
20
21 [n,p] = size(X);
22 W = reshape(w,[p k]);
23
24 XW = X*W;
25 Z = sum(exp(XW),2);
26
27 ind = sub2ind([n k],[1:n]',y);
28 nll = -sum(XW(ind)-log(Z))+ sum(sum((lambda/2)*W.^2));
29
30
31 g = zeros(p,k);
32 for c = 1:k
33     g(:,c) = X'*(exp(XW(:,c))./Z-(y==c)) + lambda*W(:,c);
34 end
35
36 g = reshape(g,[p*k 1]);
37 end

```

2. The number of non-zero parameters: 35

The number of original features that the model uses: 19

```

1 function [model] = softmaxClassifierL1(X,y,lambda,maxIter)
2
3 % Compute sizes
4 [n,d] = size(X);
5 k = max(y);
6
7 W = zeros(d,k); % Each column is a classifier
8 W(:) = proxGradL1(@softmaxLoss,W(:),X,y,k,lambda,maxIter);
9
10
11 model.W = W;
12 model.predict = @predict;
13 end
14
15 function [yhat] = predict(model,X)
16 W = model.W;
17 [~,yhat] = max(X*W,[],2);
18 end
19
20
21 function [nll,g] = softmaxLoss(w,X,y,k,lambda)
22
23 [n,p] = size(X);
24 W = reshape(w,[p k]);
25

```



```

26 XW = X*W;
27 Z = sum(exp(XW),2);
28
29 ind = sub2ind([n k],[1:n]',y);
30 nll = -sum(XW(ind)-log(Z));
31
32 g = zeros(p,k);
33 for c = 1:k
34     g(:,c) = X'*(exp(XW(:,c))./Z-(y == c));
35 end
36 g = reshape(g,[p*k 1]);
37 end

```

3. errTrain = 0.0220; errTest = 0.0540 The number of non-zero parameters: 115  
The number of original features that the model uses: 23

```

1 function [model] = softmaxClassifierGL1(X,y,lambda,maxIter)
2
3 % Compute sizes
4 [n,d] = size(X);
5 k = max(y);
6
7 W = zeros(d,k); % Each column is a classifier
8 W(:) = proxGradGroupL1(@softmaxLoss,W(:),X,y,k,lambda,maxIter);
9
10
11 model.W = W;
12 model.predict = @predict;
13 end
14
15 function [yhat] = predict(model,X)
16 W = model.W;
17 [~,yhat] = max(X*W,[],2);
18 end
19
20
21 function [nll,g] = softmaxLoss(w,X,y,k,lambda)
22
23 [n,p] = size(X);
24 W = reshape(w,[p k]);
25
26 XW = X*W;
27 Z = sum(exp(XW),2);
28
29 ind = sub2ind([n k],[1:n]',y);
30
31 tmp = sqrt(sum(W.^2,2));
32
33 nll = -sum(XW(ind)-log(Z));
34
35 g = zeros(p,k);
36 for c = 1:k
37     g(:,c) = X'*(exp(XW(:,c))./Z-(y == c));
38 end
39 g = reshape(g,[p*k 1]);

```

```

40 end
41
42 function [w,f] = proxGradGroupL1(funObj,w,X,y,k,lambda,maxIter)
43 % Minimize funOb(w) + lambda*sum(abs(w))
44
45 % Evaluate initial objective and gradient of smooth part
46 [f,g] = funObj(w,X,y,k,lambda);
47 funEvals = 1;
48
49 L = 1;
50 [n,p] = size(X);
51 k = max(y);
52 while funEvals < maxIter
53
54     % proximal-gradient step
55     alpha = 1/L;
56     W = reshape(w,[p k]);
57     W_new = W;
58     W_update = reshape(w - alpha*g,[p k]);
59     for i = 1: p
60         W_new(i,:) = softThreshold(W_update(i,:),alpha*lambda);
61     end
62     w_new = W_new(:);
63     %w_new = softThreshold(w - alpha*g,alpha*lambda);
64     [f_new,g_new] = funObj(w_new,X,y,k,lambda);
65     funEvals = funEvals + 1;
66
67     % adaptive step-size
68     while f_new > f + g'*(w_new - w) + (L/2)*norm(w_new-w)^2
69         L = L*2;
70         alpha = 1/L;
71         w_new = softThreshold(w - alpha*g,alpha*lambda);
72         [f_new,g_new] = funObj(w_new,X,y,k,lambda);
73         funEvals = funEvals + 1;
74     end
75
76     w = w_new;
77     f = f_new;
78     g = g_new;
79
80     % Print out how we are doing
81     optCond = norm(w-softThreshold(w-g,lambda),'inf');
82     fprintf('%6d %15.5e %15.5e %15.5e\n',funEvals,alpha,f + lambda*sum(abs(w)
83         ),optCond);
84
85     if optCond < 1e-1
86         break;
87     end
88 end
89
90 function [w] = softThreshold(w,threshold)
91     w = (w./(norm(w))).*max(0,norm(w)-threshold);
92 end
93

```

```

94 function [nll,g] = softmaxLoss(w,X,y,k,lambda)
95
96 [n,p] = size(X);
97 W = reshape(w,[p k]);
98
99 XW = X*W;
100 Z = sum(exp(XW),2);
101
102 ind = sub2ind([n k],[1:n]',y);
103
104 tmp = sqrt(sum(W.^2,2));
105
106 nll = -sum(XW(ind)-log(Z));
107
108 g = zeros(p,k);
109 for c = 1:k
110     g(:,c) = X'*(exp(XW(:,c))./Z-(y == c));
111 end
112 g = reshape(g,[p*k 1]);
113 end

```

## 2.3 Stochastic Gradient

1. We tried the step-sizes defined by  $1/(0.1t), 1/(0.2t), \dots, 1/(2t)$ , and  $1/(0.1\sqrt{t}), 1/(0.2\sqrt{t}), \dots, 1/(2\sqrt{t})$ . The maximum number of passes is set to 10. The step-size with the best performance is  $\frac{1}{1.8\sqrt{t}}$  and its corresponding function value is 2.7145e+04.
2. We choose to weight the  $w$  with equal weights. The step-size with the best performance is  $1/\sqrt{t}$  and its corresponding function value is 2.7076e+04.
3. We tried the step-sizes defined by  $1/(0.1\sqrt{t}), 1/(0.2\sqrt{t}), \dots, 1/(2\sqrt{t})$ , and  $\delta$  defined by 1, 101, 201, 301, 401. The sequence with the best performance is step-size =  $1/(0.1\sqrt{t})$  and  $\delta = 301$ , the function value is 2.7083e+04. There are other combinations that gives close results: e.g. step-size =  $1/(0.2\sqrt{t})$  and  $\delta = 301$ ; step-size =  $1/(0.3\sqrt{t})$  and  $\delta = 301$ .

```

1 load quantum.mat
2 [n,d] = size(X);
3 lambdaFull = 1;
4
5 % Initialize
6 maxPasses = 10;
7 progTol = 1e-4;
8 w = zeros(d,1);
9 lambda = lambdaFull/n; % The regularization parameter on one example
10 dt = zeros(d,d);
11 tmp = zeros(d,1); % save the summation of the squared gradient in D
12 % Stochastic gradient
13 w_old = w;
14 c_candidate = 0.1:0.1:2;
15 for j = 1: length(c_candidate) %search for the constant value in the step
    sie
16     for k = 1:100:500 %search for the optimal delta

```

```

17     w = zeros(d,1);
18     w_old = w;
19     w_pre_iter = w;
20     tmp = zeros(d,1);
21     delta = k;
22     c_candidate(j);
23     for t = 1:maxPasses*n
24
25         % Choose variable to update
26         i = randi(n);
27
28         % Evaluate the gradient for example i
29
30         [f,g] = logisticL2_loss(w,X(i,:),y(i),lambda);
31         tmp = tmp + g.^2;
32         D = diag(1./ (sqrt(delta + tmp)));
33
34         % Choose the step-size
35         alpha = 1./ (c_candidate(j)*t^(0.5));
36
37         w = w - alpha*D*g;
38
39         if mod(t,n) == 0
40             change = norm(w-w_old,inf);
41             %fprintf('Passes = %d, function = %.4e, change = %.4f\n',
42                     t/n,logisticL2_loss(w,X,y,lambdaFull),change);
43             if change < progTol
44                 % fprintf('Parameters changed by less than progTol on
45                     pass\n');
46                 break;
47             end
48             w_old = w;
49         end
50     end
51
52 end

```

4. The function value is  $2.7068e+04$  after 10 passes. It appears to converge faster than the algorithms in 1-3.

```

1  load quantum.mat
2  [n,d] = size(X);
3  lambdaFull = 1;
4
5  % Initialize
6  maxPasses = 10;
7  progTol = 1e-6;
8  w = zeros(d,1);
9  lambda = lambdaFull/n; % The regularization parameter on one example
10
11 % Stochastic gradient
12 w_old = w;

```

```

13 L = .25*max(sum(X.^2,2)) + lambda;
14
15 D = zeros(d,1);    %represent the sum of the gradients calculated from n
    samples
16 Y_i = zeros(n,d);
17 g_old = zeros(d,1);
18 for t = 1:maxPasses*n
19
20     % Choose variable to update
21     i = randi(n);
22
23     % Evaluate the gradient for example i
24
25     [f,g] = logisticL2_loss(w,X(i,:),y(i),lambda);
26
27     D = D - (Y_i(i,:))' + g;
28     Y_i(i,:) = g';
29
30     % Choose the step-size
31     alpha = 1/L;
32     w = w - (alpha*D./n);
33
34     if mod(t,n) == 0
35         change = norm(w-w_old,inf);
36         fprintf('Passes = %d, function = %.4e, change = %.4f\n',t/n,
            logisticL2_loss(w,X,y,lambdaFull),change);
37         if change < progTol
38             fprintf('Parameters changed by less than progTol on pass\n');
39             break;
40         end
41         w_old = w;
42     end
43 end

```

### 3 Kernels and Duality

#### 3.1 Fenchel Duality

1. Let  $X = X, z = X\omega$ .

$$\begin{aligned}
f^*(u) &= \sup_z \left\{ u^T z - \frac{1}{2} \|z - y\|^2 \right\} \\
&= \sup_z \left\{ -\frac{1}{2} \|z - u - y\|^2 + \frac{1}{2} \|u\|^2 + u^T y \right\} \\
&= \frac{1}{2} \|u\|^2 + u^T y. \\
g^*(u) &= \sup_\omega \left\{ u^T \omega - \frac{\lambda}{2} \|\omega\|^2 \right\} \\
&= \sup_\omega \left\{ -\frac{\lambda}{2} \|\omega - \frac{1}{\lambda} u\|^2 + \frac{1}{2\lambda} \|u\|^2 \right\} \\
&= \frac{1}{2\lambda} \|u\|^2.
\end{aligned}$$

Therefore, we have

$$D(z) = -f^*(-z) - g^*(X^T z) = -\left(\frac{1}{2} \|z\|^2 - z^T y\right) - \frac{1}{2\lambda} \|X^T z\|^2.$$

2. Let  $X = X, z = X\omega$ .

$$\begin{aligned}
f^*(u) &= \sup_z \left\{ u^T z - \|z - y\|_1 \right\} \\
&= \sup_z \left\{ \sum_{i=1}^n (u_i z_i - |z_i - y_i|) \right\} \\
&= \begin{cases} u^T y & \text{if } |u_i| \leq 1 \text{ for all } i \\ \infty & \text{else} \end{cases} \\
g^*(u) &= \sup_\omega \left\{ u^T \omega - \lambda \|\omega\|_1 \right\} \\
&= \begin{cases} 0 & \text{if } |u_i| \leq \lambda \text{ for all } i \\ \infty & \text{else} \end{cases}
\end{aligned}$$

Therefore, we have

$$D(z) = -f^*(-z) - g^*(X^T z) = \begin{cases} z^T y & \text{if } |z_i| \leq 1 \text{ and } |(X^T z)_i| \leq \lambda \text{ for all } i \\ -\infty & \text{else} \end{cases}.$$

3. Let  $X = - \begin{pmatrix} y^1 x^1 \\ y^2 x^2 \\ \vdots \\ y^n x^n \end{pmatrix}$  and  $z = X\omega$ .

$$\begin{aligned} f^*(u) &= \sup_z \left\{ u^T z - \sum_{i=1}^n \log(1 + \exp(z_i)) \right\} \\ &= \begin{cases} \sum_{i=1}^n \left\{ u_i \log(u_i) + (1 - u_i) \log(1 - u_i) \right\} 1_{(0 < u_i < 1)} & \text{if } 0 \leq u_i < 1 \text{ for all } i \\ \infty & \text{else} \end{cases} \end{aligned}$$

$$\begin{aligned} g^*(u) &= \sup_{\omega} \left\{ u^T \omega - \frac{\lambda}{2} \|\omega\|^2 \right\} \\ &= \sup_{\omega} \left\{ -\frac{\lambda}{2} \|\omega - \frac{1}{\lambda} u\|^2 + \frac{1}{2\lambda} \|u\|^2 \right\} \\ &= \frac{1}{2\lambda} \|u\|^2. \end{aligned}$$

Therefore, we have

$$\begin{aligned} D(z) &= -f^*(-z) - g^*(X^T z), \\ &= \begin{cases} -\sum_{i=1}^n \left\{ -z_i \log(-z_i) + (1 + z_i) \log(1 + z_i) \right\} 1_{(-1 < z_i < 0)} - \frac{1}{2\lambda} \|X^T z\|^2, & \text{if } -1 < z_i \leq 0 \text{ for all } i \\ -\infty & \text{else} \end{cases} \end{aligned}$$

## 3.2 Stochastic Dual Coordinate Ascent

The code is given below. The optimal values of the primal and dual objectives with  $\lambda = 1$  are 90.6829 and 90.6827 respectively. The number of support vectors is 99.

```
%%%%%%%%%%%%%% Matlab code
X = load('statlog.heart.data');
y = X(:,end);
y(y==2) = -1;
X = X(:,1:end-1);
n = size(X,1);

% Add bias and standardize
X = [ones(n,1) standardizeCols(X)];
d = size(X,2);

% Set regularization parameter
lambda = 1;

% Some values used by the dual
```

```

YX = diag(y)*X;
G = YX*YX';

% Find min for -D(z)
z = quadprog(G/lambda, -ones(n,1), [],[],[],[], zeros(n,1), ones(n,1));

% Convert from dual to primal variables
w = (1/lambda)*(YX'*z);

% Evaluate primal objective:
P = sum(max(1-y.*(X*w),0)) + (lambda/2)*(w'*w)

% Evaluate dual objective:
D = sum(z) - (z'*G*z)/(2*lambda)

% count number of support vectors
sum(z> 0.01)
%%%%%%%%%%%% end-Matlab code

```

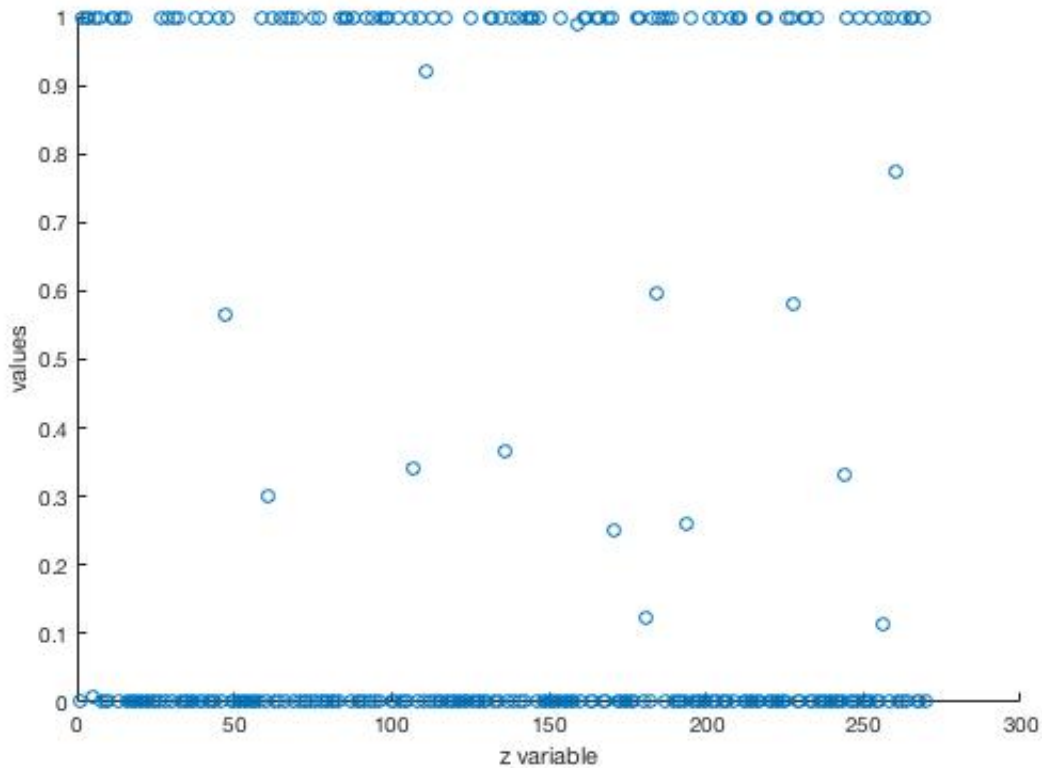


Figure 1: Values of  $z$  variables.



### 3.3 Large-Scale Kernel Methods

1. The best value for  $\sigma$  is 0.5 and the best value for  $\lambda$  is  $\lambda_1 = 4.8828e - 04$ . The  $\sigma$  with Gaussian kernels is the same as the one with Gaussian RBFs. The  $\lambda$  for Gaussian RBFs was  $\lambda_2 = 1.47e - 05$ . The approximate relationship between  $\lambda_1$  and  $\lambda_2$  is that  $\lambda_2 \approx \lambda_1 * \text{mean}(K(:))$ , where  $K = \text{rbfBasis}(X, X, \text{sigma})$  is the matrix of pair-wise evaluations of the kernel for all training patterns.

2. The code is provided below and the performances with different  $m$  values ( $m = 10, 30, 40, 60, 200$ ) are visually examined by Figures 2-6. From the figures we see that the performance gets better when the  $m$  value increases. The performance of  $m=60$  is as good as using the whole dataset (i.e.  $m=200$ ). The squaredTestError values are  $2.0787e+04$  with  $m=10$ ,  $3.6422e+03$  with  $m=30$ ,  $429.0070$  with  $m=40$ ,  $82.8489$  with  $m=60$ , and  $83.4175$  with  $m=200$ .

```
%%%%%%%%%%%%%% Matlab code
lambda=4.8828e-04;
sigma = 0.5;
m=10; % or 30, 40, 60, 200
model = kernelRegression332(X,y,lambda,sigma,m);
yhat = model.predict(model, Xtest,m);
squaredTestError = sum((yhat-ytest).^2)/t

figure(1);
plot(X,y,'b. ');
hold on
plot(Xtest,ytest,'g. ');
Xhat = [min(X):.1:max(X)]'; % Choose points to evaluate the function
yhat = model.predict(model, Xhat, m); % Gaussian kernel
plot(Xhat,yhat,'r ');
ylim([-300 400]);

function [model] = kernelRegression332(X,y,lambda,sigma, m)

% Compute sizes
[n,d] = size(X);

% new
K11 = rbfBasis(X(1:m,:),X(1:m,:),sigma);
K21 = rbfBasis(X((m+1):n,:),X(1:m,:),sigma);
K1 = [K11;K21];
z=(K1'*K1+lambda*K11)\K1'*y;

model.X = X;
model.z = z;
model.sigma = sigma;
```

```

    model.predict = @predict;
end

function [yhat] = predict(model,Xhat,m)
    Khat = rbfBasis(Xhat,model.X(1:m,:),model.sigma);
    yhat = Khat*model.z;
end

function [Xrbf] = rbfBasis(X1,X2,sigma)
    n1 = size(X1,1);
    n2 = size(X2,1);
    d = size(X1,2);
    Z = 1/sqrt(2*pi*sigma^2);
    D = X1.^2*ones(d,n2) + ones(n1,d)*(X2')).^2 - 2*X1*X2';
    Xrbf = Z*exp(-D/(2*sigma^2));
end
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%% end- Matlab code

```

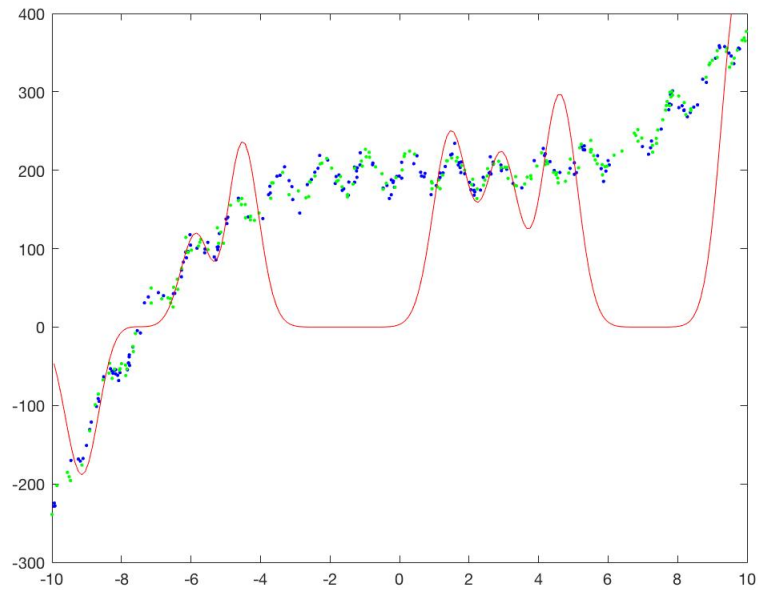


Figure 2: Performance with  $m = 10$ .

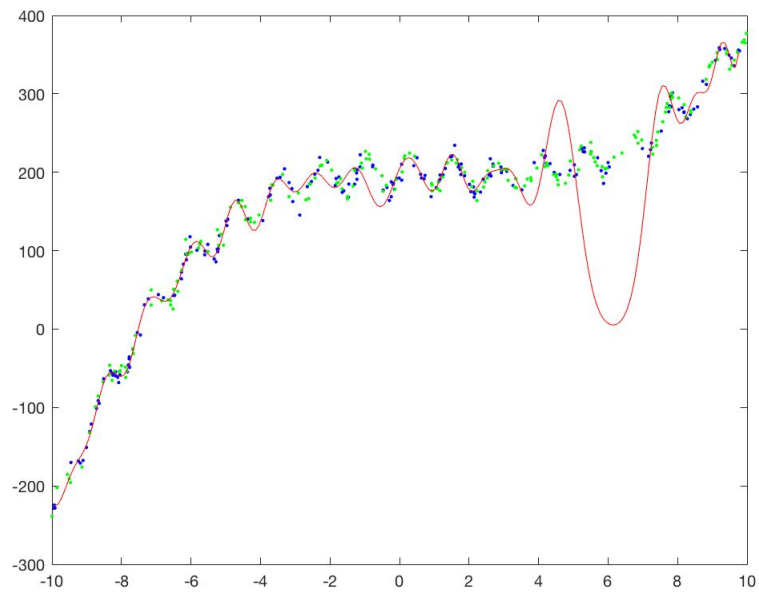


Figure 3: Performance with  $m = 30$ .

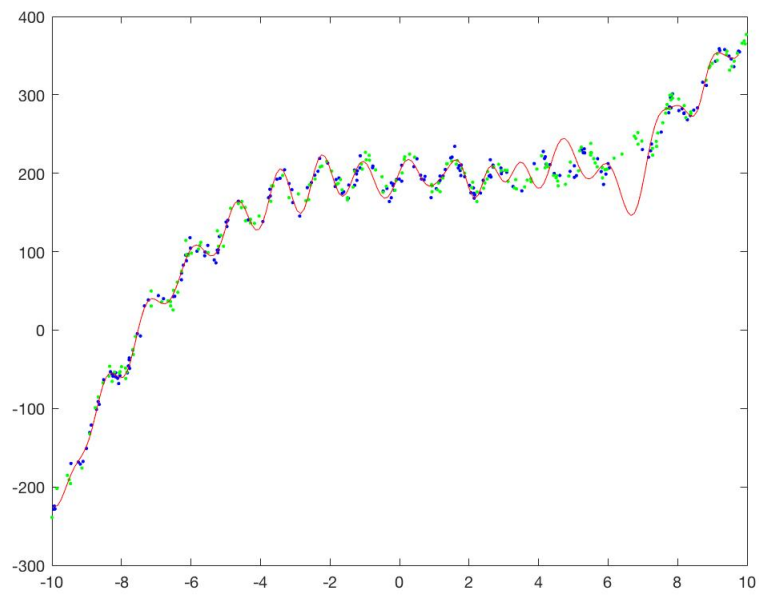


Figure 4: Performance with  $m = 40$ .

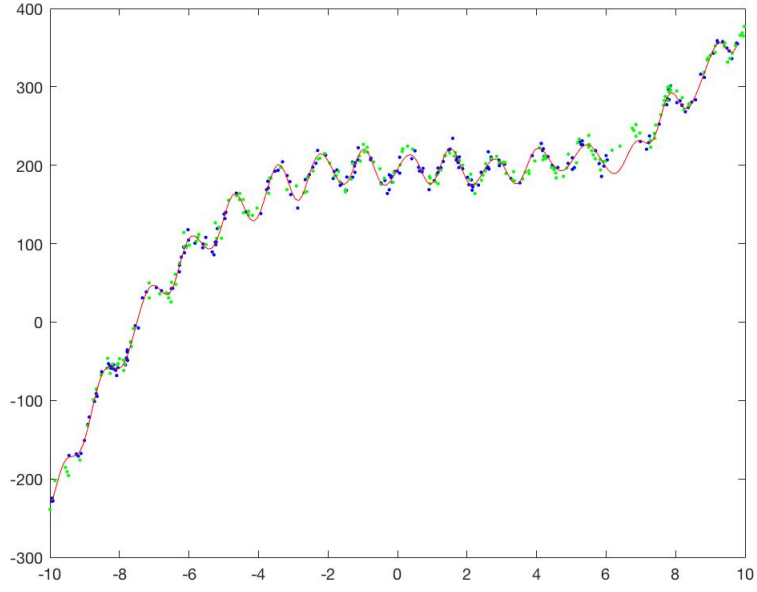


Figure 5: Performance with  $m = 60$ .

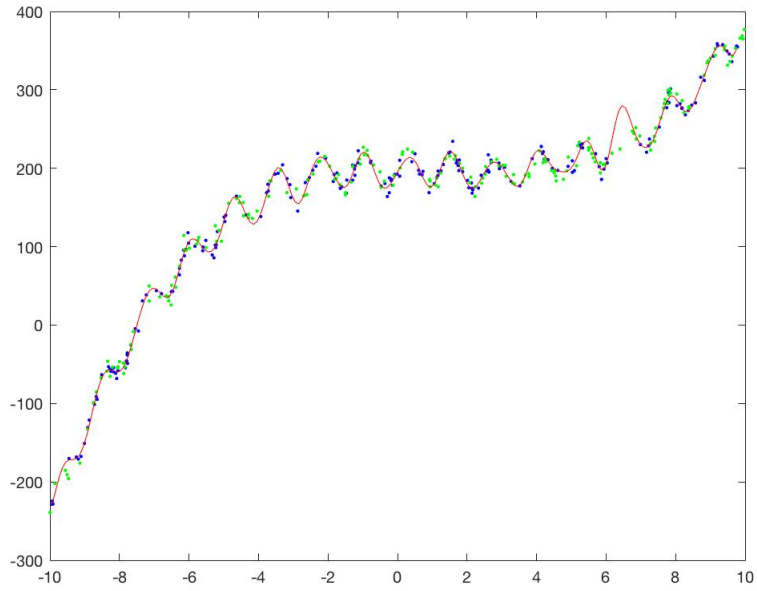


Figure 6: Performance with  $m = 200$ .

3. The code is given below and the performances with different  $m$  values ( $m = 3, 5, 8, 10, 20$ ) are visually examined by Figures 7-11. From the figures we see that the performance gets better when the  $m$  value increases. This method can achieve a fitting curve that captures the overall pattern well with a small value of  $m$  ( $m=10$ ), which is much smaller than the  $m$  value of the

subset of regressors model ( $m=60$ ). However, this method does not capture the wiggly pattern in the data as the subset of regressors method does.

```
%%%%%%%%%%%%%% Matlab code
clear all
close all

% Load data
load nonLinear.mat % Loads {X,y,Xtest,ytest}
[n,d] = size(X);
[t,~] = size(Xtest);

lambda=4.8828e-04;
sigma = 0.5;

rng(1);
m=3; % or 5, 8, 10, 20
model = kernelRegression333(X,y,lambda,sigma,m);
yhat = model.predict(model, Xtest);
squaredTestError = sum((yhat-ytest).^2)/t
figure(1);
plot(X,y,'b. ');
hold on
plot(Xtest,ytest,'g. ');
Xhat = [min(X):.1:max(X)]'; % Choose points to evaluate the function
yhat = model.predict(model, Xhat);
plot(Xhat,yhat,'r ');
ylim([-300 400]);

function [model] = kernelRegression333(X,y,lambda,sigma, m)

% Compute sizes
[n,d] = size(X);

R = randn(d, m)*sigma; % dxm
Z = exp(sqrt(-1)*X*R); % n,m

% Solve least squares problem(
z = (Z'*Z + lambda*eye(m))\Z'*y; % m,1

model.R = R; % dxm
model.z = z; % mx1
model.sigma = sigma;
model.predict = @predict;
end
```

```

function [yhat] = predict(model,Xhat)
    Zhat = exp(sqrt(-1)*Xhat*model.R); % n2xm
    yhat = Zhat*model.z;
end
%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%%% end- Matlab code

```

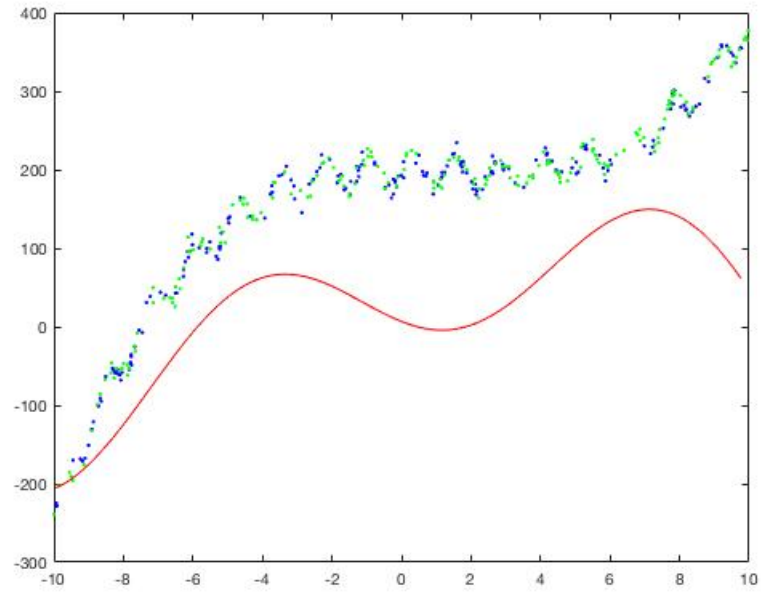


Figure 7: Performance with  $m = 3$ .

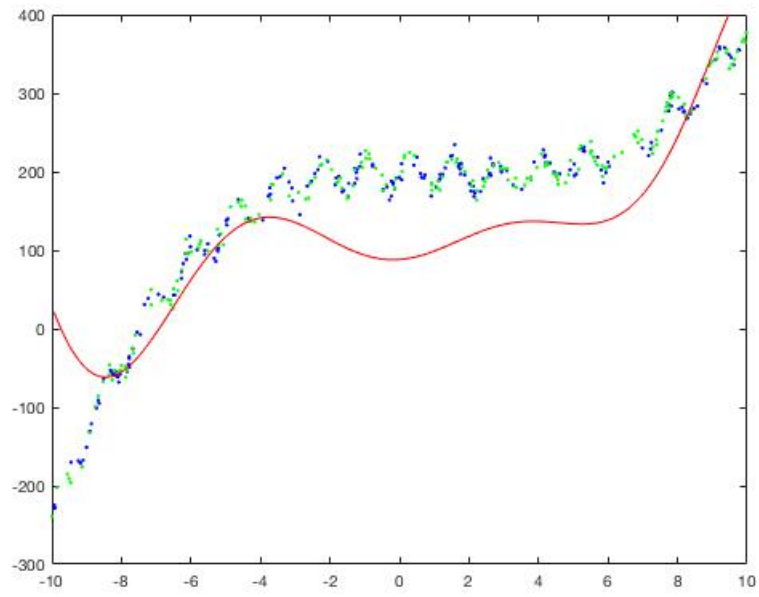


Figure 8: Performance with  $m = 5$ .

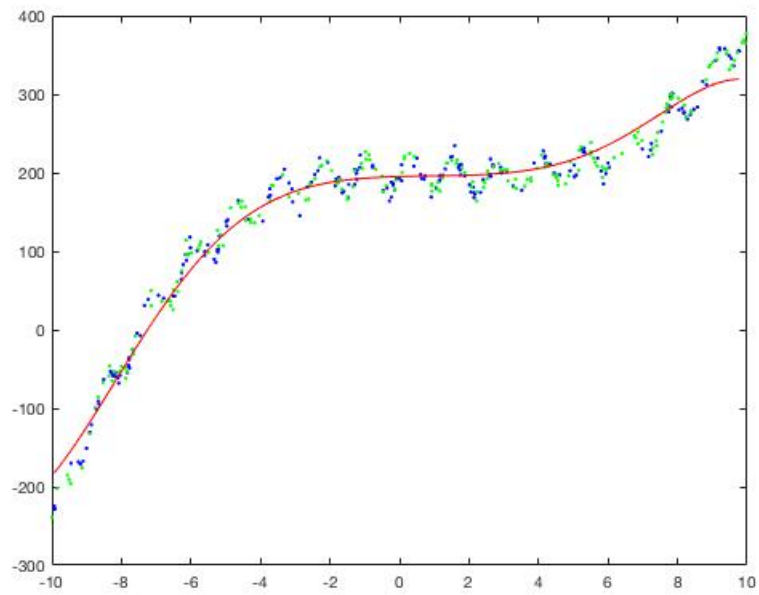


Figure 9: Performance with  $m = 8$ .

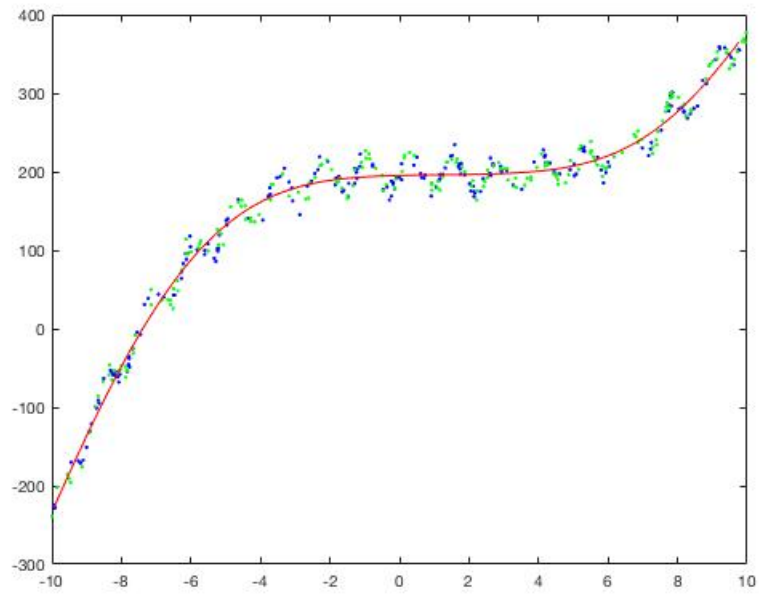


Figure 10: Performance with  $m = 10$ .

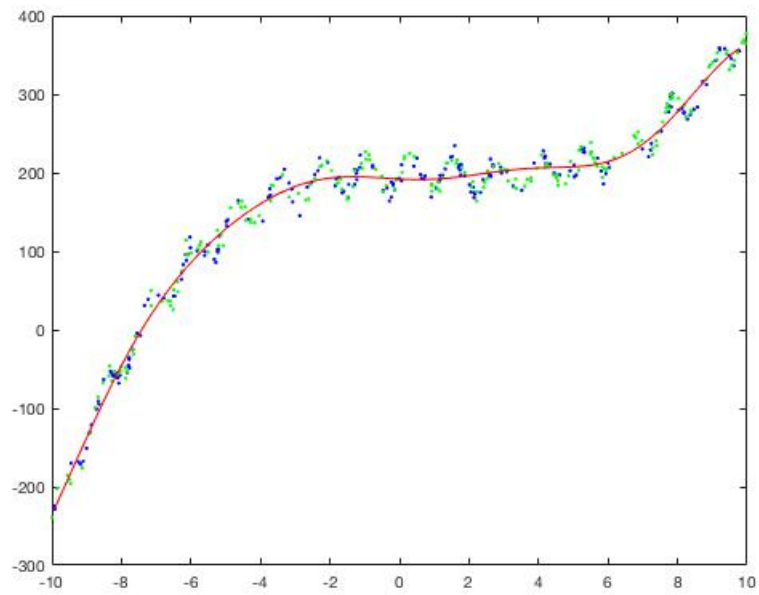


Figure 11: Performance with  $m = 20$ .