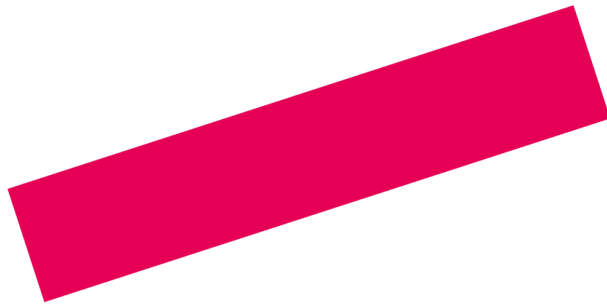


Master Engineering Systems

Big data & Small data



wk3 – linear regression assignment



linear regression assignment gradient descent

- Update thetas, loop over thetas

```
1 function [theta, cost_h] = mvgd_doubleloop2(X, y, theta, alpha, num_iters)
2 %MVGD Performs multi-variable gradient descent to learn theta
3 % theta = MVGD(x, y, theta, alpha, num_iters) updates theta by
4 % taking num_iters gradient steps with learning rate alpha
5
6 %% Initialize some useful values
7 m = length(y); %Number of training examples
8 n = length(theta);
9 cost_h = zeros(num_iters, 1); %Preallocate for speed
10 %% Perform gradient descent
11 for iter = 1:num_iters %loop for iterations
12     hypothesis = X*theta;
13     gradient = (1/m)*X'*(hypothesis-y); % compute gradient
14     for index_theta = 1:n %loop over thetas
15         %update theta
16         theta(index_theta) = theta(index_theta) - alpha*gradient(index_theta);
17     end
18     % Save the cost J in every iteration
19     cost_h(iter) = 1/(2*m)*sum((X*theta-y).^2);
20 end
21 end
```

linear regression assignment

gradient descent

- Update thetas, loop over thetas
- Vectorized, update all thetas in one line

```
1 function [theta, cost_h] = mvgd(X, y, theta, alpha, num_iters)
2 %MVGD Performs multi-variable gradient descent to learn theta
3 % theta = MVGD(x, y, theta, alpha, num_iters) updates theta by
4 % taking num_iters gradient steps with learning rate alpha
5
6 %% Initialize some useful values
7 m = length(y); %Number of training examples
8 cost_h = zeros(num_iters, 1); %Preallocate for speed
9
10 %% Perform gradient descent
11 for iter = 1:num_iters
12     gradient = (1/m)*X'*(X*theta-y);
13     % Perform a single gradient step on the parameter vector
14     theta = theta - alpha*gradient;
15     % Save the cost J in every iteration
16     cost_h(iter) = 1/(2*m)*sum((X*theta-y).^2);
17 end
end
```

linear regression assignment

gradient descent

- Update thetas, loop over thetas

- Vectorized, update

- Stopping criteria

- Fixed number of iterations
- Check for convergence

```
1 function [theta, cost_h] = mvgd2(X, y, theta, alpha, num_iters)
2 %MVG2 Performs multi-variable gradient descent to learn theta
3 % theta = MVGD(x, y, theta, alpha, num_iters) updates theta by
4 % taking at most num_iters gradient steps with learning rate alpha
5 % Stops when the improvement in the cost function is below eps = 10e-10
6
7 %% Initialize some useful values
8 m = length(y); %Number of training examples
9 cost_h = zeros(num_iters, 1); %Preallocate for speed
10 %% Perform gradient descent
11 iter=1;
12 delta_cost=1; % just an initial value to make sure it passes the first check
13 eps = 10e-12; % computer precision (~10^-16)
14 while iter <= num_iters && delta_cost>eps
15     gradient = 1/m*X'*(X*theta-y);
16     % Perform a single gradient step on the parameter vector
17     theta = theta - alpha*gradient;
18     % Save the cost J in every iteration
19     cost_h(iter) = 1/(2*m)*sum((X*theta-y).^2);
20     % Check for improvement in the cost
21     if iter>1
22         delta_cost = cost_h(iter-1)-cost_h(iter);
23     end
24     iter = iter +1;
25 end
26 %remove trailing zeros from cost history
27 cost_h = cost_h(1:iter-1);
28 end
29
30
```

gradient descent weights and bias

```
1 function [weights, bias, cost_h] = mvgd_weights_bias(X_nobias, y, weights, bias, alpha, num_iters)
2 %MVGD Performs multi-variable gradient descent to learn theta
3 % theta = MVGD(x, y, theta, alpha, num_iters) updates theta by
4 % taking num_iters gradient steps with learning rate alpha
5
6 %% Initialize some useful values
7
8 m = length(y); %Number of training examples
9 cost_h = zeros(num_iters, 1); %Preallocate for speed
10 %% Perform gradient descent
11 for iter = 1:num_iters
12     hypothesis = X_nobias*weights + bias*ones(m,1);
13
14     % compute the gradient
15     gradient_w = (1/m)*X_nobias'*(hypothesis-y);
16     gradient_b = (1/m)*sum(hypothesis-y);
17
18     % Perform a single gradient step on the parameter vector and the
19     % bias
20     weights = weights - alpha*gradient_w;
21     bias = bias - alpha*gradient_b;
22
23     % Save the cost J in every iteration
24     hypothesis = X_nobias*weights + bias*ones(m,1);
25     cost_h(iter) = 1/(2*m)*sum((hypothesis-y).^2);
26 end
27 end
```

- Separate gradient and update for weights and bias

linear regression assignment normalization

- Shift with mean
- Scale with standard deviation or range
- A vectorized version
- What can go wrong?
 - Constant feature leads to division by zero
 - Do not include the bias vector in X

```
function [X,mu,sigma] = normalizeFeatures(X)

mu=zeros(1,size(X,2)); % creating a row of zeros
sigma=zeros(1,size(X,2)); % creating a row of zeros

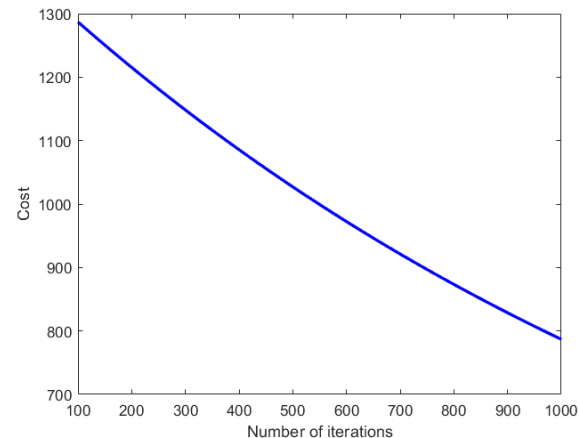
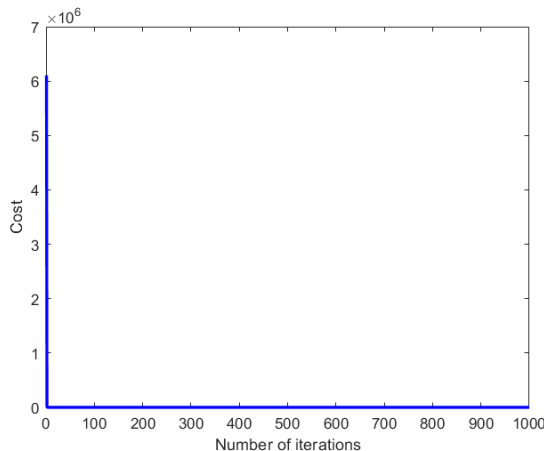
tic
for i=1:size(mu,2)
    mu(1,i)=mean(X(:,i)); % iterative calculation of mu
    sigma(1,i)=std(X(:,i)); % iterative calculation of sigma
    X(:,i)=(X(:,i)-mu(1,i))/sigma(1,i); % iterative normalization
end
toc
```

```
function [X, mu, sigma] = normalizeFeatures(X)
%Preparing parameters
mu = mean(X); %
sigma=max(X) std(X) %
%Data normalization
X = (X-mu)./sigma; %
end
```

linear regression assignment

finding learning rate and number of iterations

- Non-normalized data
 - Features with ranges in different orders of magnitude
 - $X^T X$ eigenvalues $\sim 0.5 - 3.8 \cdot 10^9$
 - Slope very high in one direction, very low in other
 - Computational problem
 - Very small learning rate, many iterations
 - How to check for convergence?
 - Plot may be misleading



compare normalized and non-normalized

- normalized:
 - $\alpha = 0.3$
 - iterations = ~ 1000
 - $\hat{\theta} = (23.45, -0.8415, 2.082, -0.6524, -5.499, 0.2223, 2.766, 1.149)'$
 - $J(\hat{\theta}) = 5.424$
- non-normalized
 - $\alpha \leq 2 \cdot 10^{-7}$
 - iterations $\geq 10^7$
 - $\hat{\theta} = (0.9280, 0.2528, 0.0059, -0.0369, -0.0062, -0.0927, 0.5552, 1.0695)'$
 - $J(\hat{\theta}) = 5.774$

direct vs gradient descent

- normalized data:
 - direct: 0.0005 s
 - gradient descent: 0.03 s
 - theta values ($\hat{\theta}$) are nearly the same:
 - $\hat{\theta}_d = (23.446, -0.842, 2.082, -0.653, -5.499, 0.222, 2.766, 1.149)'$
 - $\hat{\theta}_g = (23.446, -0.841, 2.081, -0.652, -5.499, 0.222, 2.766, 1.149)'$
- non-normalized data:
 - direct: 0.0005 s
 - gradient descent: ~ 100 s (10^7 iterations)
 - theta values ($\hat{\theta}$) are completely different:
 - $\hat{\theta}_d = (-17.2, -0.49, 0.020, -0.017, -0.0065, 0.081, 0.75, 1.43)'$
 - $\hat{\theta}_g = (0.93, 0.25, 0.0059, -0.037, -0.0062, -0.093, 0.56, 1.07)'$