

Machine Learning

Derek Li

Contents

1	Octave	4
1.1	Plotting data	4
1.2	For, while, if statements, and functions	4
1.2.1	For statements	4
1.2.2	While statements	4
1.2.3	If statements	4
1.2.4	Functions	4
1.3	Vectorization	4
2	Univariate Linear Regression	5
2.1	Model representation	5
2.2	Cost function	5
2.3	Gradient descent algorithm	5
2.4	Gradient descent for linear regression	5
3	Multivariate Linear Regression	7
3.1	Model representation	7
3.2	Gradient descent for multivariate linear regression	7
3.2.1	Feature scaling	7
3.2.2	Learning rate	7
3.3	Normal equation	7
4	Logistic Regression	9
4.1	Hypothesis representation	9
4.2	Decision boundary	9
4.3	Cost function	9
4.4	Gradient descent for logistic regression	9
4.5	Advanced optimization	10
4.6	Multi-class classification	10
5	Regularization	11
5.1	Over-fitting	11
5.1.1	Addressing over-fitting	11
5.2	Cost function	11
5.3	Regularized linear regression	11
5.3.1	Gradient descent	11
5.3.2	Normal equation	11
5.4	Regularized logistic regression	12
5.4.1	Cost function	12

5.4.2	Gradient descent	12
5.4.3	Advanced optimization	12
6	Neural Network	13
6.1	Model representation	13
6.2	Cost function	14
6.3	Back-Propagation algorithm	14
6.4	Unrolling parameters	14
6.5	Gradient checking	15
6.6	Random initialization	15
6.7	Implementation of neural network	16
7	Evaluation of Learning Algorithm	17
7.1	Evaluation of hypothesis	17
7.1.1	Training procedure for linear regression	17
7.1.2	Training procedure for logistic regression	17
7.2	Model selection and training/validation/test sets	17
7.3	Diagnosing Bias and Variance	17
7.3.1	Regularization	18
7.4	Learning curves	18
7.4.1	High bias	18
7.4.2	High variance	18
7.5	Debugging a learning algorithm	18
8	Machine Learning System Design	19
8.1	Error analysis	19
8.2	Error metrics for skewed classes	19
8.3	Data for machine learning	19
9	Support Vector Machines	20
9.1	Optimization objective	20
9.2	Kernels	20
9.2.1	Kernels and similarity	20
9.2.2	SVM with kernels	20
9.2.3	SVM parameters	20
9.3	SVM in practice	21
9.3.1	Gaussian kernel	21
9.3.2	Other choices of kernel	21
9.3.3	Multi-class classification	21
9.3.4	Logistic regression vs. SVM	21
10	Clustering	22
10.1	K -Means algorithm	22
10.2	Optimization objective	22
10.3	Random initialization	22
11	Principal Component Analysis	23
11.1	Principal component analysis algorithm	23
11.2	Reconstruction from compressed representation	23
11.3	Choosing the number of principal components	23

11.4 Supervised learning speedup	24
12 Anomaly Detection	25
12.1 Gaussian distribution	25
12.2 Anomaly detection algorithm	25
12.3 Developing and evaluating an anomaly detection system	25
12.4 Anomaly detection vs. supervised learning	25
12.5 Multivariate Gaussian distribution	26
12.6 Anomaly detection using the multivariate Gaussian distribution	26
13 Recommender Systems	27
13.1 Content-Based recommendations	27
13.2 Collaborative filtering	27
13.3 Vectorization: Low rank matrix factorization	28
14 Large Scale Machine Learning	29
14.1 Stochastic gradient descent	29
14.2 Mini-batch gradient descent	29
14.3 Stochastic gradient descent convergence	29
14.4 Online learning	29
14.5 Map reduce and data parallelism	29

1 Octave

1.1 Plotting data

- Save the plot: `print -dpng 'filename'`
- Divide plot a 1×2 grid, access first element: `subplot(1, 2, 1)`
- Clears a figure: `clf`
- Visualizes the matrix A : `imagesc(A)`

1.2 For, while, if statements, and functions

1.2.1 For statements

```
for condition,  
    execution;  
end;
```

1.2.2 While statements

```
while condition,  
    execution;  
end;
```

1.2.3 If statements

```
if condition,  
    execution;  
elseif condition,  
    execution;  
else  
    execution;  
end;
```

1.2.4 Functions

```
function (output(s)) = FunctionName(variable(s))  
    execution;
```

1.3 Vectorization

Example 1.1 ($h_{\theta}(x) = \sum_{j=0}^n \theta_j x_j = \theta^T \mathbf{x}$).

```
prediction = theta' * x;
```

2 Univariate Linear Regression

2.1 Model representation

The hypothesis function is

$$h_{\theta}(x) = \theta_0 + \theta_1 x,$$

where θ_i 's are parameters.

2.2 Cost function

For univariate linear regression, we define a cost function

$$J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2,$$

where m is the size of the training set. Our objective is to solve

$$\min_{\theta_0, \theta_1} J(\theta_0, \theta_1).$$

The cost function $J(\theta_0, \theta_1)$ is also called squared error cost function.

2.3 Gradient descent algorithm

Algorithm 2.1 (Gradient Descent Algorithm). Repeat

$$\left\{ \theta_j := \theta_j - \alpha \frac{\partial}{\partial \theta_j} J(\theta_0, \theta_1) \right\}$$

until convergence.

Note 1. α is called the learning rate, which controls the velocity of gradient descent. If α is too small, gradient descent can be slow; if α is too large, it may fail to converge or even diverge.

Note 2. The algorithm requires to update θ_j 's simultaneously.

Note 3. Generally, depending on where it is initializing, it can end up with different local optima.

2.4 Gradient descent for linear regression

Since for linear regression, $J(\theta_0, \theta_1) = \frac{1}{2m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)})^2$, then

$$\begin{aligned} \frac{\partial}{\partial \theta_0} J(\theta_0, \theta_1) &= \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}), \\ \frac{\partial}{\partial \theta_1} J(\theta_0, \theta_1) &= \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)}. \end{aligned}$$

Algorithm 2.2 (Gradient Descent Algorithm for Linear Regression). Repeat

$$\left\{ \begin{aligned} \theta_0 &:= \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) \\ \theta_1 &:= \theta_1 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(x^{(i)}) - y^{(i)}) x^{(i)} \end{aligned} \right\}$$

until convergence.

Note 1. Cost function for linear regression is convex and thus, it does not have any local optima except for one global optimum. Hence, gradient descent always makes it converge to the global optimum.

Note 2. The algorithm is also called Batch Gradient Descent, and "Batch" means each step of gradient descent uses all the training examples.

3 Multivariate Linear Regression

3.1 Model representation

The hypothesis function is

$$h_{\theta}(x) = \theta_0 x_0 + \theta_1 x_1 + \cdots + \theta_n x_n = \theta^T \mathbf{x},$$

where $x_0 = 1$.

3.2 Gradient descent for multivariate linear regression

Algorithm 3.1 (Gradient Descent Algorithm for Multivariate Linear Regression). Repeat

$$\left\{ \theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)}) x_j^{(i)} \right\}$$

until convergence.

3.2.1 Feature scaling

Gradient descent can converge more quickly if the features are on a similar scale. Hence, we want every feature into approximately a $-1 \leq x_i \leq 1$ range.

Method 1. Dividing by the maximum value.

Method 2. (Mean normalization)

$$x_i := \frac{x_i - \mu_i}{s_i}.$$

3.2.2 Learning rate

For sufficiently small α , $J(\theta)$ decreases on every iteration. To choose α , try

$$..., 0.001, 0.03, 0.01, 0.03, 0.1, 0.3, 1, ...$$

and plot $J(\theta)$ as a function of the number of iterations. Try to pick the largest possible value or just slightly smaller than the largest reasonable value.

3.3 Normal equation

Normal equation is a method to solve for θ analytically, compared to gradient descent.

Suppose there are m training examples and n features, then

$$\mathbf{x}^{(i)} = \begin{pmatrix} x_0^{(i)} \\ x_1^{(i)} \\ \vdots \\ x_n^{(i)} \end{pmatrix} \in \mathbb{R}^{n+1}.$$

Then the design matrix is

$$X = \begin{pmatrix} (\mathbf{x}^{(1)})^T \\ (\mathbf{x}^{(2)})^T \\ \vdots \\ (\mathbf{x}^{(m)})^T \end{pmatrix}.$$

We have

$$\theta = (X^T X)^{-1} X^T \mathbf{y},$$

and such θ minimize the $J(\theta)$. The code in Octave is:

```
pinv(X' * X) * X' * y
```

Even $X^T X$ is not invertible (singular/degenerate), **pinv** can still compute $X^T X$. But in this case, we can delete redundant features (linearly dependent). If there are too many features ($m \geq n$), then delete some features or use regularization.

Comparison between gradient descent and normal equation:

Gradient Descent	Normal Equation
Needs to choose α	Does not need to choose α
Needs many iterations	Does not need to iterate
Works well even when n is large	Need to compute $(X^T X)^{-1}$, $O(n^3)$
	Slow if n is very large

4 Logistic Regression

4.1 Hypothesis representation

The hypothesis function is

$$h_{\theta}(\mathbf{x}) = g(\theta^T \mathbf{x}),$$

where

$$g(z) = \frac{1}{1 + e^{-z}},$$

and $g(z)$ is called sigmoid/logistic function. Therefore,

$$h_{\theta}(\mathbf{x}) = \frac{1}{1 + e^{-\theta^T \mathbf{x}}}.$$

We have $h_{\theta}(\mathbf{x}) = P(y = 1 | \mathbf{x}; \theta)$.

4.2 Decision boundary

Suppose $h_{\theta}(\mathbf{x}) \geq 0.5 \Rightarrow y = 1$, $h_{\theta}(\mathbf{x}) < 0.5 \Rightarrow y = 0$, then

$$\theta^T \mathbf{x} \geq 0 \Rightarrow y = 1, \theta^T \mathbf{x} < 0 \Rightarrow y = 0.$$

The decision boundary is the line that separates the area where $y = 0$ and $y = 1$.

4.3 Cost function

For logistic regression, we define

$$\text{Cost}(h_{\theta}(\mathbf{x}), y) = \begin{cases} -\ln(h_{\theta}(\mathbf{x})) & \text{if } y = 1 \\ -\ln(1 - h_{\theta}(\mathbf{x})) & \text{if } y = 0 \end{cases} = -y \ln(h_{\theta}(\mathbf{x})) - (1 - y) \ln(1 - h_{\theta}(\mathbf{x})),$$

then the cost function is

$$\begin{aligned} J(\theta) &= \frac{1}{m} \sum_{i=1}^m \text{Cost}(h_{\theta}(\mathbf{x}^{(i)}), y^{(i)}) \\ &= -\frac{1}{m} \sum_{i=1}^m [y^{(i)} \ln(h_{\theta}(\mathbf{x}^{(i)})) + (1 - y^{(i)}) \ln(1 - h_{\theta}(\mathbf{x}^{(i)}))]. \end{aligned}$$

4.4 Gradient descent for logistic regression

Algorithm 4.1 (Gradient Descent Algorithm for Logistic Regression). Repeat

$$\left\{ \theta_j := \theta_j - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)}) x_j^{(i)} \right\}$$

until convergence.

Note. A vectorized implementation is $\theta := \theta - \frac{\alpha}{m} X^T (g(X\theta) - \mathbf{y})$.

4.5 Advanced optimization

In addition to gradient descent, there are some optimization algorithms: conjugate gradient, BFGS, L-BFGS. They do not need to manually pick α and are often faster than gradient descent, but they are more complex. In Octave,

```
function [jVal, gradient] = costFunction(theta)
    jVal = ____; % Code to compute J(theta)
    gradient = [____]; % Code to compute paritial derivative of J(theta) w.r.t. theta_0,
        ..., theta_n
end
options = optimset('GradObj', 'on', 'MaxIter', 100);
initialTheta = zeros(2, 1);
[optTheta, functionVal, exitFlag] = fminunc(@costFunction, initialTheta, options);
```

4.6 Multi-class classification

Train a logistic regression classifier $h_{\theta}^{(i)}(\mathbf{x})$ for each class i to predict the probability that $y = i$. On a new input \mathbf{x} , to make a prediction, pick the class i that maximizes $h_{\theta}^{(i)}(\mathbf{x})$.

5 Regularization

5.1 Over-fitting

If we have too many features, the learned hypothesis may fit the training set very well, $J(\theta) \approx 0$, but fail to generalize to new examples.

5.1.1 Addressing over-fitting

Method 1. Reduce number of features

- (1) Manually select which features to keep.
- (2) Model selection algorithm.

Method 2. Regularization

- (1) Keep all the features, but reduce magnitude/values of parameters θ_j .
- (2) Works well when we have a lot of features, each of which contributes a bit to predicting y .

5.2 Cost function

The cost function is

$$J(\theta) = \frac{1}{2m} \left[\sum_{i=1}^m (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)})^2 + \lambda \sum_{j=1}^n \theta_j^2 \right],$$

where $\lambda \sum_{j=1}^n \theta_j^2$ is the regularization term and λ is the regularization parameter.

5.3 Regularized linear regression

5.3.1 Gradient descent

Algorithm 5.1 (Gradient Descent Algorithm for Regularized Linear Regression). Repeat

$$\left\{ \begin{array}{l} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)}) x_0^{(i)} \\ \theta_j := \theta_j (1 - \alpha \frac{\lambda}{m}) - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)}) x_j^{(i)}, j = 1, \dots, n \end{array} \right\}$$

until convergence.

5.3.2 Normal equation

The θ minimizes $J(\theta)$ is

$$\theta = \left(X^T X + \lambda \begin{pmatrix} 0 & & \\ & 1 & \\ & & \ddots \\ & & & 1 \end{pmatrix} \right)^{-1} X^T \mathbf{y},$$

where $\begin{pmatrix} 0 & & \\ & 1 & \\ & & \ddots \\ & & & 1 \end{pmatrix}$ is a $(n+1) \times (n+1)$ matrix.

5.4 Regularized logistic regression

5.4.1 Cost function

The cost function is

$$J(\theta) = - \left[\frac{1}{m} \sum_{i=1}^m y^{(i)} \ln(h_{\theta}(\mathbf{x}^{(i)})) + (1 - y^{(i)}) \ln(1 - h_{\theta}(\mathbf{x}^{(i)})) \right] + \frac{\lambda}{2m} \sum_{j=1}^n \theta_j^2.$$

5.4.2 Gradient descent

Algorithm 5.2 (Gradient Descent Algorithm for Regularized Logistic Regression). Repeat

$$\left\{ \begin{array}{l} \theta_0 := \theta_0 - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)}) x_0^{(i)} \\ \theta_j := \theta_j (1 - \alpha \frac{\lambda}{m}) - \alpha \frac{1}{m} \sum_{i=1}^m (h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)}) x_j^{(i)}, j = 1, \dots, n \end{array} \right\}$$

until convergence.

5.4.3 Advanced optimization

In Octave,

```
function [jVal, gradient] = costFunction(theta)
jVal = ____; % Code to compute J(theta)
gradient = [____]; % Code to compute paritial derivative of J(theta) w.r.t. theta_0,
    ..., theta_n
end
options = optimset('GradObj', 'on', 'MaxIter', 100);
initialTheta = zeros(2, 1);
[optTheta, functionVal, exitFlag] = fminunc(@costFunction, initialTheta, options);
```

6 Neural Network

6.1 Model representation

In neural network, the first layer is called the input layer and the final layer is called the output layer.

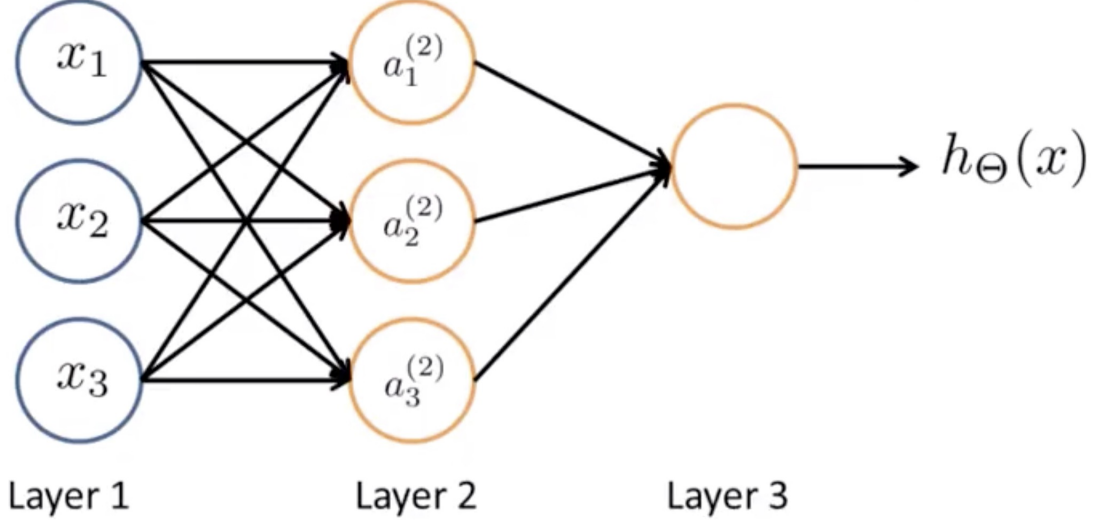


Figure 6.1: Model Representation for Neural Network

$a_i^{(j)}$ is the activation of unit i in layer j , and $\Theta^{(j)}$ is the matrix of weights controlling function mapping from layer j to layer $j + 1$. For this figure, we have

$$\begin{aligned} a_1^{(2)} &= g(\Theta_{10}^{(1)}x_0 + \Theta_{11}^{(1)}x_1 + \Theta_{12}^{(1)}x_2 + \Theta_{13}^{(1)}x_3) = g(z_1^{(2)}), \\ a_2^{(2)} &= g(\Theta_{20}^{(1)}x_0 + \Theta_{21}^{(1)}x_1 + \Theta_{22}^{(1)}x_2 + \Theta_{23}^{(1)}x_3) = g(z_2^{(2)}), \\ a_3^{(2)} &= g(\Theta_{30}^{(1)}x_0 + \Theta_{31}^{(1)}x_1 + \Theta_{32}^{(1)}x_2 + \Theta_{33}^{(1)}x_3) = g(z_3^{(2)}), \\ h_{\Theta}(\mathbf{x}) &= a_1^{(3)} = g(\Theta_{10}^{(2)}a_0^{(2)} + \Theta_{11}^{(2)}a_1^{(2)} + \Theta_{12}^{(2)}a_2^{(2)} + \Theta_{13}^{(2)}a_3^{(2)}) = g(z^{(3)}). \end{aligned}$$

Generally, if network has s_j units in layer j , s_{j+1} units in layer $j + 1$, then $\Theta^{(j)}$ will be of dimension $s_{j+1} \times (s_j + 1)$.

For this figure, a vectorized implementation is

$$\begin{aligned} \mathbf{z}^{(2)} &= \Theta^{(1)}\mathbf{x} = \Theta^{(1)}\mathbf{a}^{(1)}, \\ \mathbf{a}^{(2)} &= g(\mathbf{z}^{(2)}), \end{aligned}$$

$$\text{where } \mathbf{x} = \begin{pmatrix} x_0 \\ x_1 \\ x_2 \\ x_3 \end{pmatrix}, \mathbf{z}^{(2)} = \begin{pmatrix} z_1^{(2)} \\ z_2^{(2)} \\ z_3^{(2)} \end{pmatrix}.$$

Add $a_0^{(2)} = 1$, then

$$\begin{aligned} z^{(3)} &= \Theta^{(2)}\mathbf{a}^{(2)}, \\ h_{\Theta}(\mathbf{x}) &= a^{(3)} = g(z^{(3)}). \end{aligned}$$

The process is called forward propagation.

In general,

$$\begin{aligned}
\mathbf{a}^{(1)} &= \mathbf{x} \\
\mathbf{z}^{(2)} &= \Theta^{(1)} \mathbf{a}^{(1)} \\
\mathbf{a}^{(2)} &= g(\mathbf{z}^{(2)}) \quad (\text{add } a_0^{(2)}) \\
\mathbf{z}^{(3)} &= \Theta^{(2)} \mathbf{a}^{(2)} \\
\mathbf{a}^{(3)} &= g(\mathbf{z}^{(3)}) \quad (\text{add } a_0^{(3)}) \\
&\vdots \\
\mathbf{z}^{(L-1)} &= \Theta^{(L-2)} \mathbf{a}^{(L-2)} \\
\mathbf{a}^{(L-1)} &= g(\mathbf{z}^{(L-1)}) \quad (\text{add } a_0^{(L-1)}) \\
\mathbf{z}^{(L)} &= \Theta^{(L-1)} \mathbf{a}^{(L-1)} \\
\mathbf{a}^{(L)} &= g(\mathbf{z}^{(L)}) = h_{\Theta}(\mathbf{x})
\end{aligned}$$

6.2 Cost function

The cost function is

$$J(\Theta) = -\frac{1}{m} \left[\sum_{i=1}^m \sum_{k=1}^K \mathbf{y}_k^{(i)} \ln(h_{\Theta}(\mathbf{x}^{(i)}))_k + (1 - \mathbf{y}_k^{(i)}) \ln(1 - (h_{\Theta}(\mathbf{x}^{(i)}))_k) \right] + \frac{\lambda}{2m} \sum_{l=1}^{L-1} \sum_{i=1}^{s_l} \sum_{j=1}^{s_{l+1}} (\Theta_{ji}^{(l)})^2,$$

where K is the number of output units, L is the total number of layers in network, and s_l is the number of units (not counting bias unit) in layer l .

6.3 Back-Propagation algorithm

Algorithm 6.1 (Back-Propagation Algorithm). Suppose training set is $\{(\mathbf{x}^{(1)}, \mathbf{y}^{(1)}), \dots, (\mathbf{x}^{(m)}, \mathbf{y}^{(m)})\}$, and set $\Delta_{ij}^{(l)} = 0$.

For $i = 1$ to m :

Set $\mathbf{a}^{(1)} = \mathbf{x}^{(i)}$.

Perform forward propagation to compute $\mathbf{a}^{(l)}$ for $l = 2, \dots, L$.

Use $\mathbf{y}^{(i)}$ to compute $\delta^{(L)} = \mathbf{a}^{(L)} - \mathbf{y}^{(i)}$.

Compute $\delta^{(L-1)}, \dots, \delta^{(2)}$, where $\delta^{(l)} = (\Theta^{(l)})^T \delta^{(l+1)} \cdot g'(\mathbf{z}^{(l)})$, and $g'(\mathbf{z}^{(l)}) = \mathbf{a}^{(l)} \cdot (1 - \mathbf{a}^{(l)})$.

$\Delta_{ij}^{(l)} := \Delta_{ij}^{(l)} + a_j^{(l)} \delta_i^{(l+1)}$. In vector form, $\Delta^{(l)} := \Delta^{(l)} + \delta^{(l+1)} (\mathbf{a}^{(l)})^T$.

$$\frac{\partial}{\partial \Theta_{ij}^{(l)}} D_{ij}^{(l)} = D_{ij}^{(l)} := \begin{cases} \frac{1}{m} \Delta_{ij}^{(l)} & \text{if } j = 0 \\ \frac{1}{m} (\Delta_{ij}^{(l)} + \lambda \Theta_{ij}^{(l)}) & \text{if } j \neq 0 \end{cases}.$$

6.4 Unrolling parameters

Example 6.1. Suppose $s_1 = 10, s_2 = 10, s_3 = 1$, then we have

$$\begin{aligned}
\Theta^{(1)} &\in \mathbb{R}^{10 \times 11}, \Theta^{(2)} \in \mathbb{R}^{10 \times 11}, \Theta^{(3)} \in \mathbb{R}^{1 \times 11}, \\
D^{(1)} &\in \mathbb{R}^{10 \times 11}, D^{(2)} \in \mathbb{R}^{10 \times 11}, D^{(3)} \in \mathbb{R}^{1 \times 11}.
\end{aligned}$$

Hence, in Octave,

```
thetaVec = [Theta1(:); Theta2(:); Theta3(:)];
DVec = [D1(:); D2(:); D3(:)];
```

We can get back the original matrices from the unrolled versions:

```
Theta1 = reshape(thetaVec(1:110), 10, 11)
Theta2 = reshape(thetaVec(111:220), 10, 11)
Theta3 = reshape(thetaVec(221:231), 1, 11)
```

To summarize, we have initial parameters $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ and unroll to get initialTheta to pass fminunc(@costFunction, initialTheta, options).

For the cost function, we get $\Theta^{(1)}, \Theta^{(2)}, \Theta^{(3)}$ from thetaVec and use FP/BP to compute $D^{(1)}, D^{(2)}, D^{(3)}$ and $J(\Theta)$. Then unroll $D^{(1)}, D^{(2)}, D^{(3)}$ to get gradientVec.

6.5 Gradient checking

For unrolled $\theta = (\theta_1, \dots, \theta_n)$, we have

$$\begin{aligned} \frac{\partial}{\partial \theta_1} J(\theta) &\approx \frac{J(\theta_1 + \varepsilon, \theta_2, \dots, \theta_n) - J(\theta_1 - \varepsilon, \theta_2, \dots, \theta_n)}{2\varepsilon}, \\ &\vdots \\ \frac{\partial}{\partial \theta_n} J(\theta) &\approx \frac{J(\theta_1, \theta_2, \dots, \theta_n + \varepsilon) - J(\theta_1, \theta_2, \dots, \theta_n - \varepsilon)}{2\varepsilon}. \end{aligned}$$

In Octave,

```
for i = 1:n,
    thetaPlus = theta;
    thetaPlus(i) = thetaPlus(i) + EPSILON;
    thetaMinus = theta;
    thetaMinus(i) = thetaMinus(i) - EPSILON;
    gradApprox(i) = (J(thetaPlus) - J(thetaMinus)) / (2 * EPSILON);
end;
```

Check if gradApprox \approx DVec.

To summarize:

- Implement BP to compute DVec (unrolled).
- Implement numerical gradient check to compute gradApprox.
- Make sure they give similar values.
- Turn off gradient checking. Using BP code for learning.

6.6 Random initialization

Initializing all theta weights to zero does not work with neural networks. Instead we can randomly initialize weights for Θ in Octave:

```
Theta_num = rand(a, b) * (2 * INIT_EPSILON) - INIT_EPSILON;
```

6.7 Implementation of neural network

Step 1. Pick a network architecture (connectivity pattern between neurons).

Note 1. Number of input units is the dimension of features $\mathbf{x}^{(i)}$.

Note 2. Number of output units is the number of classes.

Note 3. The reasonable default for hidden layer is 1, or if the number is greater than 1, the hidden layers should have same number of hidden units in every layer. Usually, the more units the better.

Step 2. Train a neural network.

Step 2.1 Randomly initialize weights.

Step 2.2 Implement forward propagation to get $h_{\Theta}(\mathbf{x}^{(i)})$ for any $\mathbf{x}^{(i)}$.

Step 2.3 Implement code to compute cost function $J(\Theta)$.

Step 2.4 Implement backward propagation to compute $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$.

Step 2.5 Use gradient checking to compare $\frac{\partial}{\partial \Theta_{jk}^{(l)}} J(\Theta)$ computed using BP with using numerical estimate of gradient of $J(\Theta)$, then disable gradient checking code.

Step 2.6 Use gradient descent of advanced optimization method with BP to try to minimize $J(\Theta)$ as a function of parameters Θ .

7 Evaluation of Learning Algorithm

7.1 Evaluation of hypothesis

7.1.1 Training procedure for linear regression

- Learn parameter θ from training data (minimizing training error $J(\theta)$ with 70% of all data).
- Compute test set error:

$$J_{\text{test}}(\theta) = \frac{1}{2m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} (h_{\theta}(\mathbf{x}_{\text{test}}^{(i)}) - \mathbf{y}_{\text{test}}^{(i)})^2.$$

7.1.2 Training procedure for logistic regression

- Learn parameter θ from training data.
- Compute test set error:

$$J_{\text{test}}(\theta) = -\frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} [\mathbf{y}_{\text{test}}^{(i)} \ln h_{\theta}(\mathbf{x}_{\text{test}}^{(i)}) + (1 - \mathbf{y}_{\text{test}}^{(i)}) \ln(1 - h_{\theta}(\mathbf{x}_{\text{test}}^{(i)}))].$$

- Or compute misclassification error(0/1 misclassification error): Define

$$\text{Err}(h_{\theta}(\mathbf{x}), \mathbf{y}) = \begin{cases} 1 & \text{if } h_{\theta}(\mathbf{x}) \geq 0.5, y = 0 \text{ or if } h_{\theta}(\mathbf{x}) < 0.5, y = 1 \\ 0 & \text{otherwise} \end{cases}.$$

The test error is

$$\frac{1}{m_{\text{test}}} \sum_{i=1}^{m_{\text{test}}} \text{Err}(h_{\theta}(\mathbf{x}_{\text{test}}^{(i)}), \mathbf{y}_{\text{test}}^{(i)}).$$

7.2 Model selection and training/validation/test sets

- Break down the dataset into three sets: Training set(60%), cross validation set(20%), and test set(20%).
- Optimize the parameters in Θ using the training set for each polynomial degree.
- Find the polynomial degree d with the least error using the cross validation set.
- Estimate the generalization error using the test set with $J_{\text{test}}(\Theta^{(d)})$.

7.3 Diagnosing Bias and Variance

- Bias(under-fit): $J_{\text{train}}(\theta)$ will be high and $J_{\text{cv}}(\theta) \approx J_{\text{train}}(\theta)$.
- Variance(over-fit): $J_{\text{train}}(\theta)$ will be high and $J_{\text{cv}}(\theta) \gg J_{\text{train}}(\theta)$.

7.3.1 Regularization

- Create a list of λ , i.e.,

$$\lambda \in \{0, 0.01, 0.02, 0.04, 0.08, 0.16, 0.32, 0.64, 1.28, 2.56, 5.12, 10.24\}.$$

- Create a set of models with different degrees or any other variants.
- Iterate through the λ 's and for each λ go through all the models to learn some Θ .
- Compute the cross validation error using the learned Θ computed with λ on the $J_{cv}(\Theta)$ without regularization.
- Select the best combo that produces the lowest error on the cross validation set.
- Using the best combo Θ and λ , apply it on $J_{test}(\Theta)$ to see if it has a good generalization of the problem.

7.4 Learning curves

7.4.1 High bias

- Low training set size causes $J_{train}(\Theta)$ to be low and $J_{CV}(\Theta)$ to be high.
- Large training set size causes both $J_{train}(\Theta)$ and $J_{CV}(\Theta)$ to be high with $J_{train}(\Theta) \approx J_{CV}(\Theta)$.
- If a learning algorithm is suffering from high bias, getting more training data will not help much.

7.4.2 High variance

- Low training set size causes J_{train} to be low and $J_{CV}(\Theta)$ to be high.
- $J_{train}(\Theta)$ increases with training set size and $J_{CV}(\Theta)$ continues to decrease without leveling off. Also $J_{train}(\Theta) < J_{CV}(\Theta)$ but the difference between them remains significant.
- If a learning algorithm is suffering from high variance, getting more training data is likely to help.

7.5 Debugging a learning algorithm

- Get more training examples - Fixes high variance.
- Try smaller sets of features - Fixes high variance.
- Try getting additional features - Fixes high bias.
- Try adding polynomial features - Fixes high bias.
- Try decreasing λ - Fixes high bias.
- Try increasing λ - Fixes high variance.

8 Machine Learning System Design

8.1 Error analysis

The recommended approach to solving machine learning problems is to:

- Start with a simple algorithm, implement it quickly, and test it on cross validation data.
- Plot learning curves to decide if more data, more features, etc. are likely to help.
- Error analysis: Manually examine the errors on examples in the cross validation set and try to spot a trend where most of the errors were made.

8.2 Error metrics for skewed classes

We use precision/recall to measure the error for skewed classes.

- Precision = $\frac{\text{True positive}}{\text{Predicted positive}} = \frac{\text{True positive}}{\text{True positive} + \text{False positive}}$.
- Recall = $\frac{\text{True positive}}{\text{Actual positive}} = \frac{\text{True positive}}{\text{True positive} + \text{False negative}}$.

One way to trade off between precision and recall is using F score:

$$F = 2 \frac{PR}{P + R}.$$

8.3 Data for machine learning

Large data rationale: Assume feature $\mathbf{x} \in \mathbf{R}^{n+1}$ has sufficient information to predict y accurately. Using a learning algorithm with many parameters causes low bias ($J_{\text{train}}(\theta)$ will be small). Using a very large training set (unlikely to over-fit) causes low variance ($J_{\text{train}}(\theta) \approx J_{\text{test}}(\theta) \Rightarrow J_{\text{test}}(\theta)$ will be small).

9 Support Vector Machines

9.1 Optimization objective

We want to solve

$$\min_{\theta} C \sum_{i=1}^m [y^{(i)} \text{cost}_1(\theta^T \mathbf{x}^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T \mathbf{x}^{(i)})] + \frac{1}{2} \sum_{j=1}^n \theta_j^2.$$

Unlike logistic regression, the support vector machine does not output the probability but makes a prediction of y .

9.2 Kernels

A Gaussian kernel function is

$$f_i = \text{similarity}(\mathbf{x}, \mathbf{l}^{(i)}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{l}^{(i)}\|^2}{2\sigma^2}\right).$$

9.2.1 Kernels and similarity

For example, take the first landmark

$$f_1 = \text{similarity}(\mathbf{x}, \mathbf{l}^{(1)}) = \exp\left(-\frac{\|\mathbf{x} - \mathbf{l}^{(1)}\|^2}{2\sigma^2}\right) = \exp\left(-\frac{\sum_{j=1}^n (x_j - l_j^{(1)})^2}{2\sigma^2}\right).$$

If $\mathbf{x} \approx \mathbf{l}^{(1)}$, $f_1 \approx 1$. If \mathbf{x} is far from $\mathbf{l}^{(1)}$, $f_1 \approx 0$.

9.2.2 SVM with kernels

Given $(\mathbf{x}^{(1)}, y^{(1)}), \dots, (\mathbf{x}^{(m)}, y^{(m)})$, choose $\mathbf{l}^{(1)} = \mathbf{x}^{(1)}, \dots, \mathbf{l}^{(m)} = \mathbf{x}^{(m)}$. Compute $\mathbf{f}^{(i)} = \begin{pmatrix} f_0^{(i)} \\ \vdots \\ f_m^{(i)} \end{pmatrix}$, where $f_0^{(i)} = 1$, and solve

$$\min_{\theta} C \sum_{i=1}^m y^{(i)} \text{cost}_1(\theta^T \mathbf{f}^{(i)}) + (1 - y^{(i)}) \text{cost}_0(\theta^T \mathbf{f}^{(i)}) + \frac{1}{2} \sum_{j=1}^n \theta_j^2,$$

where $n = m$.

9.2.3 SVM parameters

- Large C causes lower bias and higher variance.
- Small C causes higher bias and lower variance.
- Large σ^2 causes features f_i to vary more smoothly, higher bias and lower variance.
- Small σ^2 causes features f_i to vary less smoothly, lower bias and higher variance.

9.3 SVM in practice

9.3.1 Gaussian kernel

If we use Gaussian kernel, we should perform feature scaling before using the Gaussian kernel.

9.3.2 Other choices of kernel

Not all similarity functions make valid kernels, and it need to satisfy Mercer's Theorem to make sure SVM packages' optimization run correctly and do not diverge. Here is some kernels available:

- Polynomial kernel
- String kernel
- chi-square kernel
- Histogram intersection kernel, etc.

9.3.3 Multi-class classification

Many SVM packages already have built-in multi-class classification functionality. Otherwise, use one-vs.-all method: Train K SVMs, one to distinguish $y = i$ from the rest, for $i = 1, \dots, K$, get $\theta^{(1)}, \dots, \theta^{(K)}$ and pick class i with largest $(\theta^{(i)})^T \mathbf{x}$.

9.3.4 Logistic regression vs. SVM

- If n is large (relative to m ,) then we use logistic regression, or SVM without a kernel (linear kernel).
- If n is small, m is intermediate, then we use SVM with Gaussian kernel.
- If n is small, m is large, then we add more features, then use logistic regression.
- Neural network is likely to work well for most of these settings but may be slower to train.

10 Clustering

10.1 K -Means algorithm

Algorithm 10.1 (K -Means Algorithm). Take number of clusters K and training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$ as inputs, where $\mathbf{x}^{(i)} \in \mathbb{R}^n$. Randomly initialize K cluster centroids $\mu_1, \dots, \mu_K \in \mathbb{R}^n$. Repeat

$$\left\{ \begin{array}{l} \text{for } i = 1 \text{ to } m, c^{(i)} := \text{index (from 1 to } K) \text{ of cluster centroid closest to } x^{(i)} \\ \text{for } k = 1 \text{ to } K, \mu_k := \text{average (mean) of points assigned to cluster } k \end{array} \right\}$$

until convergence.

10.2 Optimization objective

Let $c^{(i)}$ be the index of cluster $(1, \dots, K)$ to which example $\mathbf{x}^{(i)}$ is currently assigned, μ_k be the cluster centroid k , and $\mu_{c^{(i)}}$ be the cluster centroid of cluster to which example $\mathbf{x}^{(i)}$ has been assigned.

The cost function is

$$J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K) = \frac{1}{m} \sum_{i=1}^m \|x^{(i)} - \mu_{c^{(i)}}\|^2.$$

Then we solve

$$\min_{c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K} J(c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K).$$

10.3 Random initialization

We can randomly initialize K -means and run, then get $c^{(1)}, \dots, c^{(m)}, \mu_1, \dots, \mu_K$ and compute cost function (distortion). We can take random initialization many times and pick clustering that gives lowest cost.

11 Principal Component Analysis

11.1 Principal component analysis algorithm

Before principal component analysis, we should do data preprocessing. Given training set $x^{(1)}, \dots, x^{(m)}$, let $\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$. Replace each $x_j^{(i)}$ with $x_j - \mu_j$. If different features on different scales, scale features to have comparable range of values.

Algorithm 11.1 (Principal Component Analysis Algorithm). Reduce data from n dimensions to k dimensions: Compute covariance matrix:

$$\Sigma = \frac{1}{m} \sum_{i=1}^n (\mathbf{x}^{(i)})(\mathbf{x}^{(i)})^T.$$

Compute eigenvectors of matrix Σ . In Octave:

```
Sigma = (1 / m) * X' * X;  
[U, S, V] = svd(Sigma);
```

We get $U \in \mathbb{R}^{n \times n}$, and choose the first k columns to get a $n \times k$ matrix U_{reduce} , then $\mathbf{z} = U_{\text{reduce}}^T \mathbf{x}$. In Octave:

```
Ureduce = U(:, 1:k);  
z = Ureduce' * x;
```

11.2 Reconstruction from compressed representation

$$\mathbf{x}_{\text{approx}} = U_{\text{reduce}} \mathbf{z}.$$

11.3 Choosing the number of principal components

Define average squared projection error $\frac{1}{m} \sum_{i=1}^m \|\mathbf{x}^{(i)} - \mathbf{x}_{\text{approx}}^{(i)}\|^2$ and total variation in the data $\frac{1}{m} \sum_{i=1}^m \|\mathbf{x}^{(i)}\|^2$.

Typically, we choose k to be smallest value so that

$$\frac{\frac{1}{m} \sum_{i=1}^m \|\mathbf{x}^{(i)} - \mathbf{x}_{\text{approx}}^{(i)}\|^2}{\frac{1}{m} \sum_{i=1}^m \|\mathbf{x}^{(i)}\|^2} \leq 0.01.$$

In Octave:

```
[U, S, V] = svd(Sigma);
```

We get $S \in \mathbb{R}^{n \times n}$ and for given k , we just check if

$$\frac{\sum_{i=1}^k S_{ii}}{\sum_{i=1}^n S_{ii}} \geq 0.99.$$

11.4 Supervised learning speedup

We extract inputs, i.e., unlabeled dataset: $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)} \xrightarrow{PCA} \mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}$. Then we get the new training set: $((\mathbf{z}^{(1)}, y^{(1)}), \dots, (\mathbf{z}^{(m)}, y^{(m)}))$.

12 Anomaly Detection

12.1 Gaussian distribution

Say $x \in \mathbb{R}$. If x is a distributed Gaussian with mean μ , variance σ^2 , then $x \sim \mathcal{N}(\mu, \sigma^2)$. $p(x; \mu, \sigma^2) = \frac{1}{\sqrt{2\pi}\sigma} \exp\left(-\frac{(x-\mu)^2}{2\sigma^2}\right)$.

12.2 Anomaly detection algorithm

Algorithm 12.1 (Anomaly Detection Algorithm). First, choose features x_i that might be indicative of anomalous examples. Then, fit parameters $\mu_1, \dots, \mu_n, \sigma_1^2, \dots, \sigma_n^2$, and $\mu_j = \frac{1}{m} \sum_{i=1}^m x_j^{(i)}$, $\sigma_j^2 = \frac{1}{m} \sum_{i=1}^m (x_j^{(i)} - \mu_j)^2$. Given new example \mathbf{x} , compute $p(\mathbf{x})$:

$$p(\mathbf{x}) = \prod_{j=1}^n p(x_j; \mu_j, \sigma_j^2) = \prod_{j=1}^n \frac{1}{\sqrt{2\pi}\sigma_j} \exp\left(-\frac{(x_j - \mu_j)^2}{2\sigma_j^2}\right).$$

Anomaly if $p(\mathbf{x}) < \varepsilon$.

12.3 Developing and evaluating an anomaly detection system

Fit model $p(\mathbf{x})$ on training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$. On a cross validation/test example \mathbf{x} , predict

$$y = \begin{cases} 1 & \text{if } p(x) < \varepsilon \text{ (anomaly)} \\ 0 & \text{if } p(x) \geq \varepsilon \text{ (normal)} \end{cases}.$$

Possible evaluation metrics:

- True positive, false positive, false negative, true negative
- Precision/Recall
- F score

Use cross validation set to choose parameter ε .

12.4 Anomaly detection vs. supervised learning

- Anomaly detection
 - * Very small number of positive examples ($y = 1$) (0-20 is common).
 - * Large number of negative ($y = 0$) example.
 - * Application: Fraud detection, manufacturing, monitoring machines in a data center, and etc.
- Supervised learning
 - * Large number of positive and negative examples.
 - * Application: Email spam classification, weather prediction, cancer classification, and etc.

12.5 Multivariate Gaussian distribution

Say $\mathbf{x} \in \mathbb{R}^n, \mu \in \mathbb{R}^n, \Sigma \in \mathbb{R}^{n \times n}$ (covariance matrix).

$$p(\mathbf{x}; \mu, \Sigma) = \frac{1}{(2\pi)^{\frac{n}{2}} |\Sigma|^{\frac{1}{2}}} \exp \left(-\frac{1}{2} (\mathbf{x} - \mu)^T \Sigma^{-1} (\mathbf{x} - \mu) \right).$$

12.6 Anomaly detection using the multivariate Gaussian distribution

Given training set $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\}$. $\mu = \frac{1}{m} \sum_{i=1}^m \mathbf{x}^{(i)}, \Sigma = \frac{1}{m} \sum_{i=1}^m (\mathbf{x}^{(i)} - \mu)(\mathbf{x}^{(i)} - \mu)^T$.

Algorithm 12.2 (Anomaly Detection Algorithm). First, fit model $p(\mathbf{x})$ by setting μ, Σ . Then given a new example \mathbf{x} , compute $p(\mathbf{x})$ and flag an anomaly if $p(\mathbf{x}) < \varepsilon$.

The original model $p(\mathbf{x}) = p(x_1; \mu_1, \sigma_1^2) \times \dots \times p(x_n; \mu_n, \sigma_n^2)$ corresponds to multivariate Gaussian

where $\Sigma = \begin{pmatrix} \sigma_1^2 & 0 & \dots & 0 \\ 0 & \sigma_2^2 & \dots & 0 \\ 0 & 0 & \ddots & 0 \\ 0 & 0 & \dots & \sigma_n^2 \end{pmatrix}$.

The differences are:

- Original model:
 - * Manually create features to capture anomalies where x_1, \dots, x_n take unusual combinations of values.
 - * Computationally cheaper (alternatively, scales better to large n).
 - * Works even if m is small.
- Multivariate Gaussian:
 - * Automatically captures correlations between features.
 - * Computationally more expensive.
 - * Must have $m > n$ ($m \geq 10n$) or else Σ is non-invertible.

13 Recommender Systems

13.1 Content-Based recommendations

Let $r(i, j) = 1$ if user j has rated movie i (0 otherwise), $y^{(i,j)}$ be the rating by user j on movie i (if defined), $\theta^{(j)}$ be the parameter vector for user j , and $\mathbf{x}^{(i)}$ be the feature vector for movie i . Hence, for user j and movie i , predicted rating is $(\theta^{(j)})^T \mathbf{x}^{(i)}$. Let $m^{(j)}$ be the number of movies rated by user j . To learn $\theta^{(j)}$, we solve

$$\min_{\theta^{(j)}} \frac{1}{2} \sum_{i:r(i,j)=1} ((\theta^{(j)})^T \mathbf{x}^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (\theta_k^{(j)})^2.$$

To learn $\theta^{(1)}, \dots, \theta^{(n_u)}$, we solve

$$\min_{\theta^{(1)}, \dots, \theta^{(n_u)}} \frac{1}{2} \sum_{j=1}^{n_u} \sum_{i:r(i,j)=1} ((\theta^{(j)})^T \mathbf{x}^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2.$$

The optimization algorithm is gradient descent:

$$\begin{aligned} \theta_k^{(j)} &:= \theta_k^{(j)} - \alpha \sum_{i:r(i,j)=1} ((\theta^{(j)})^T \mathbf{x}^{(i)} - y^{(i,j)}) x_k^{(i)}, k = 0 \\ \theta_k^{(j)} &:= \theta_k^{(j)} - \alpha \left(\sum_{i:r(i,j)=1} ((\theta^{(j)})^T \mathbf{x}^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right), k \neq 0 \end{aligned}$$

13.2 Collaborative filtering

Given $\theta^{(1)}, \dots, \theta^{(n_u)}$, to learn $\mathbf{x}^{(i)}$, we need to solve

$$\min_{\mathbf{x}^{(i)}} \frac{1}{2} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T \mathbf{x}^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{k=1}^n (x_k^{(i)})^2.$$

Given $\theta^{(1)}, \dots, \theta^{(n_u)}$, to learn $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n_m)}$, we need to solve

$$\min_{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n_m)}} \frac{1}{2} \sum_{i=1}^{n_m} \sum_{j:r(i,j)=1} ((\theta^{(j)})^T \mathbf{x}^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2.$$

To improve the procedure, we can minimize $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n_m)}$ and $\theta^{(1)}, \dots, \theta^{(n_u)}$ simultaneously, and we have

$$J(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}) = \frac{1}{2} \sum_{(i,j):r(i,j)=1} ((\theta^{(j)})^T \mathbf{x}^{(i)} - y^{(i,j)})^2 + \frac{\lambda}{2} \sum_{i=1}^{n_m} \sum_{k=1}^n (x_k^{(i)})^2 + \frac{\lambda}{2} \sum_{j=1}^{n_u} \sum_{k=1}^n (\theta_k^{(j)})^2,$$

and we solve

$$\min_{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}} J(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}).$$

Algorithm 13.1 (Collaborative Filtering Algorithm). First, initialize $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)}$ to small random values. Then minimize $J(\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(n_m)}, \theta^{(1)}, \dots, \theta^{(n_u)})$ using gradient descent (or an

advanced optimization algorithm). If we use gradient descent, for every $j = 1, \dots, n_u, i = 1, \dots, n_m$:

$$\begin{aligned} x_k^{(i)} &:= x_k^{(i)} - \alpha \left(\sum_{j:r(i,j)=1} ((\theta^{(j)})^T \mathbf{x}^{(i)} - y^{(i,j)}) \theta_k^{(j)} + \lambda x_k^{(i)} \right) \\ \theta_k^{(j)} &:= \theta_k^{(j)} - \alpha \left(\sum_{i:r(i,j)=1} ((\theta^{(j)})^T \mathbf{x}^{(i)} - y^{(i,j)}) x_k^{(i)} + \lambda \theta_k^{(j)} \right), \end{aligned}$$

where $\mathbf{x}, \theta \in \mathbb{R}^n$. Finally, for a user with parameters θ and a movie with learned features \mathbf{x} , predict a star rating of $\theta^T \mathbf{x}$.

13.3 Vectorization: Low rank matrix factorization

Suppose the predicted ratings matrix is

$$Y = \begin{pmatrix} (\theta^{(1)})^T(\mathbf{x}^{(1)}) & (\theta^{(2)})^T(\mathbf{x}^{(1)}) & \dots & (\theta^{(n_u)})^T(\mathbf{x}^{(1)}) \\ (\theta^{(1)})^T(\mathbf{x}^{(2)}) & (\theta^{(2)})^T(\mathbf{x}^{(2)}) & \dots & (\theta^{(n_u)})^T(\mathbf{x}^{(2)}) \\ \vdots & \vdots & \ddots & \vdots \\ (\theta^{(1)})^T(\mathbf{x}^{(n_m)}) & (\theta^{(2)})^T(\mathbf{x}^{(n_m)}) & \dots & (\theta^{(n_u)})^T(\mathbf{x}^{(n_m)}) \end{pmatrix}.$$

We can let

$$X = \begin{pmatrix} (\mathbf{x}^{(1)})^T \\ (\mathbf{x}^{(2)})^T \\ \vdots \\ (\mathbf{x}^{(n_m)})^T \end{pmatrix}, \Theta = \begin{pmatrix} (\theta^{(1)})^T \\ (\theta^{(2)})^T \\ \vdots \\ (\theta^{(n_u)})^T \end{pmatrix},$$

then $Y = X\Theta^T$.

14 Large Scale Machine Learning

14.1 Stochastic gradient descent

Let

$$\text{cost}(\theta, (\mathbf{x}^{(i)}, y^{(i)})) = \frac{1}{2}(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)})^2, J_{\text{train}}(\theta) = \frac{1}{m} \sum_{i=1}^m \text{cost}(\theta, (\mathbf{x}^{(i)}, y^{(i)})).$$

First, randomly reorder dataset. Then, repeat

$$\left\{ \text{for } i := 1, \dots, m \left\{ \theta_j := \theta_j - \alpha(h_{\theta}(\mathbf{x}^{(i)}) - y^{(i)})x_j^{(i)} \text{ for every } j = 0, \dots, n \right\} \right\}.$$

14.2 Mini-batch gradient descent

Say $b = 10, m = 1000$. Repeat

$$\left\{ \text{for } i := 1, 11, 21, 31, \dots, 991 \left\{ \theta_j := \theta_j - \alpha \frac{1}{10} \sum_{k=i}^{i+9} (h_{\theta}(\mathbf{x}^{(k)}) - y^{(k)})x_j^{(k)} \text{ for every } j = 0, \dots, n \right\} \right\}.$$

14.3 Stochastic gradient descent convergence

During learning, compute $\text{cost}(\theta, (\mathbf{x}^{(i)}, y^{(i)}))$ before updating θ using $(\mathbf{x}^{(i)}, y^{(i)})$. Every 1000 (or other numbers) iterations, plot $\text{cost}(\theta, (\mathbf{x}^{(i)}, y^{(i)}))$ averaged over the last 1000 examples processed by algorithm.

14.4 Online learning

Repeat

$$\begin{aligned} &\{\text{Get } (x, y) \text{ corresponding to user. Update } \theta \text{ using } (x, y) : \\ &\quad \theta_j := \theta_j - \alpha(h_{\theta}(\mathbf{x}) - y)x_j \text{ for } j = 0, \dots, n\}. \end{aligned}$$

14.5 Map reduce and data parallelism

Many learning algorithms can be expressed as computing sums of functions over the training set. Therefore we can speed up via more machines or cores.