

利用卷积神经网络对MNIST分类

- 一. cost function
- 二. 结构
- 三. error backpropagation
- 四. 求梯度
- 五. 代码

## 利用卷积神经网络对MNIST分类

本文是参照UFDLF中的教程写的。

### 一. cost function

一般有两种cost function: 均方误差(MSE)和交叉熵(cross entropy), 本文实现的卷积神经网络采用交叉熵。

均方误差

$$C(W, B) = \frac{1}{2m} \sum_x \|y(x) - a\|^2 \quad (1)$$

其中 $W$ 为网络中的权值,  $B$ 为神经元的偏置,  $m$ 为样本数,  $x$ 为样本,  $y(x)$ 为样本对应的标签,  $a$ 为样本对应的输出。

交叉熵

$$C(W, B) = - \sum_x \log P(y(x) = k) \quad (2)$$
$$k = 0, 1, \dots, C - 1$$

其中 $C$ 为目标类的个数。本文的卷积神经网络采用的cost function是交叉熵。

### 二. 结构

包含3层(不包括输入层):

- 第一层: 卷积层(convolutional layer)。该层直接与输入图像相连, 我用的数据集是MNIST, 所以输入图像大小为28x28, 用来卷积的filter size为9x9, 步长为1, 共20个filter, 因此卷积层的大小为20(height 28 - 9 + 1) x 20(width) x 20(feature map)。
- 第二层: 池化层(pooling layer)。该层的输入为卷积层的输出, 我用的是average pooling, 而不是max pooling。采用的filter大小为2x2, 每次subsample不重叠, 因此步长为2。池化层可以大幅减少数据规模, 减少参数的数量并减少计算量。一个2x2的filter可以减少上一层75%的输入。池化层的大小为20个10x10。
- 第三层: 输出层。该层与池化层之间为全连接, 共10个神经元对应10个类, 输出为该类的概率。这两层相连的作用实际上就是一次softmax分类。在一般的softmax分类中, 输入为整个图像, 而在这里的输入为池化层的输出。

### 三. error backpropagation

feedforward过程没什么说的, 接下来重点说error backpropagation过程。

#### 1. 输出层的 $\delta$

根据定义, 第 $l$ 层的 $\delta$ 为

$$\delta^L = \nabla_{z^L} C \quad (3)$$

其中 $z^L$ 为最后一层的输入,  $C$ 为cost function。由(3)可知,  $\delta$ 与每个神经元一一对应。接下来进行推导, 为了简化, 只考虑一个样本, 且该样本属于类 $k$ 。

$$\begin{aligned}
\delta_i^L &= \frac{\partial C}{\partial z_i^L} \\
&= \frac{\partial}{\partial z_i^L} [-\log(P(y(x) = k))] \\
&= -\frac{1}{P(y(x) = k)} \cdot \frac{\partial}{\partial z_i^L} [P(y(x) = k)] \\
&= -\frac{1}{P(y(x) = k)} \cdot \frac{\partial}{\partial z_i^L} \left( \frac{e^{z_k^L}}{\sum_n e^{z_n^L}} \right) \\
&= -\frac{1}{P(y(x) = k)} \cdot \left[ \frac{\partial e^{z_k^L}}{\partial z_i^L} \cdot \frac{1}{\sum_n e^{z_n^L}} - \frac{e^{z_k^L}}{\sum_n e^{z_n^L}} \cdot \frac{e^{z_i^L}}{\sum_n e^{z_n^L}} \right] \\
&= -\frac{1}{P(y(x) = k)} [1\{k = i\} \cdot P(y(x) = k) - P(y(x) = k) \cdot P(y(x) = i)] \\
&= -(1\{k = i\} - P(y(x) = i))
\end{aligned}$$

其中，n为该层神经元的个数，这里为10。写成向量形式为

$$\delta = -(e(k) - y(x)) \quad (4)$$

其中 $e(k)$ 为一个C维列向量，C为类的个数，只有第k个元素 ( $k = 0, 1, \dots, C - 1$ ) 为1，其余为0， $y(x)$ 为网络的输出。

## 2. 池化层的delta

这一步， $\delta$ 要从最后一层传播到池化层，这两层之间是全连接的。下面推导两个全连接层之间的误差传播公式：

$$\begin{aligned}
\delta_i^l &= \frac{\partial C}{\partial z_i^l} \\
&= \sum_n \frac{\partial C}{\partial z_n^{l+1}} \cdot \frac{\partial z_n^{l+1}}{\partial z_i^l} \quad (5) \\
&= \sum_n \delta_n^{l+1} \cdot \frac{\partial (\sum_j w_{nj}^{l+1} a_j^l + b_n)}{\partial z_i^l} \\
&= \sum_n \delta_n^{l+1} \cdot w_{ni}^{l+1} \cdot \frac{\partial a_i^l}{\partial z_i^l} \\
&= \sum_n \delta_n^{l+1} \cdot w_{ni}^{l+1} \cdot \frac{\partial f(z_i^l)}{\partial z_i^l}
\end{aligned}$$

其中 $w_{ij}^l$ 为第l-1层的第j个神经元，与第l层的第i个神经元之间的权值，f(x)为激活函数，本实验中池化层的激活函数为f(x) = x。写成向量形式为

$$\delta^l = (W^{l+1})^T \delta^{l+1} \odot f'(z^l) \quad (6)$$

其中，符号 $\odot$ 为对应元素分别相乘。

## 3. 卷积层的 $\delta$

以2x2的池化filter为例，每个池化层的神经元对应卷积层中的4个，本实验的池化层是没有重叠的，所以每个池化层神经元均对应不同的卷积层神经元，因此（5）中的求和符号就不需要了。

$$\begin{aligned}
\delta_i^l &= \frac{\partial C}{\partial z_i^l} \\
&= \frac{\partial C}{\partial z_k^{l+1}} \cdot \frac{\partial z_k^{l+1}}{\partial z_i^l} \\
&= \delta_k^{l+1} \cdot \frac{\partial \frac{1}{4} (f(z_i^l) + f(z_j^l) + f(z_m^l) + f(z_n^l))}{\partial z_i^l} \\
&= \frac{1}{4} \delta_k^{l+1} \cdot f'(z_i^l)
\end{aligned}$$

向量形式为

$$\delta^l = \frac{1}{\text{poolDim}^2} \cdot \text{upsample}(\delta^{l+1}) \odot f'(z^l)$$

本实验中激活函数f(x)为sigmoid函数,poolDim为2。

例：卷积层经sigmoid函数激活后的矩阵A为

$$\begin{bmatrix} 0.1000 & 0.2000 & 0 & 0.1000 \\ 0.9000 & 0 & 0.1000 & 0.3000 \\ 0.4000 & 0.1000 & 0.1000 & 0 \\ 0 & 0 & 0 & 0.1000 \end{bmatrix}$$

池化层的 $\delta$ 为

$$\begin{bmatrix} 3 & 1 \\ 0 & 2 \end{bmatrix}$$

则可求得卷积层的 $\delta'$ 为

$$\begin{bmatrix} 0.0750 & 0.1500 & 0 & 0.0250 \\ 0.6750 & 0 & 0.0250 & 0.0750 \\ 0 & 0 & 0.0500 & 0 \\ 0 & 0 & 0 & 0.0500 \end{bmatrix}$$

求最后结果时可以利用matlab中的kron函数计算，该函数可计算两矩阵的Kronecker product。代码为

```
1. kron(delta, ones(2,2)) .* A / (poolDim)^2
```

## 四. 求梯度

1. 对输出层到池化层的权值 $W_{ij}^l, b_i^l$

这两层是全连接的，所以

$$\begin{aligned} \frac{\partial C}{\partial W_{ij}^l} &= \frac{\partial C}{\partial z_i^l} \cdot \frac{\partial z_i^l}{\partial W_{ij}^l} \\ &= \delta_i^l \cdot a_j^{l-1} \end{aligned}$$

$$\begin{aligned} \frac{\partial C}{\partial b_i^l} &= \frac{\partial C}{\partial z_i^l} \cdot \frac{\partial z_i^l}{\partial b_i^l} \\ &= \delta_i^l \end{aligned}$$

2. 输入层与卷积层之间的权值 $W_{ij}^l, b^l$

令卷积层的 $\delta$ 为 $\Delta$ ,卷积层的输入为Z,w为某权值,则

$$\frac{\partial C}{\partial w} = \text{sum}\{\Delta \odot \frac{\partial Z}{\partial w} \odot \nabla_Z C\} \quad (7)$$

sum{·}为矩阵中所有值相加。

$$\frac{\partial C}{\partial b} = \text{sum}\{\Delta \odot \nabla_Z C\} \quad (8)$$

例：

$$\Delta = \begin{bmatrix} 1 & 2 \\ 0 & 1 \end{bmatrix}, \text{输入图像 } A = \begin{bmatrix} 2 & 0 & 1 \\ 1 & 2 & 3 \\ 1 & 1 & 1 \end{bmatrix}, \text{某filter为 } F = \begin{bmatrix} x & y \\ p & q \end{bmatrix}, \text{偏置为 } b \text{ 则}$$

$$\frac{\partial C}{\partial x} = 1 * 2 * f'(2x + p + 2q + b) + 2 * 0 * f'(y + 2p + 3q + b) + 0 * 1 * f'(x + 2y + p + q + b) + 1 * 2 * f'(2x + \dots)$$

$$\frac{\partial C}{\partial b} = 1 * f'(2x + p + 2q + b) + 2 * f'(y + 2p + 3q + b) + 0 * f'(x + 2y + p + q + b) + 1 * f'(2x + 3y + p + q + b)$$

## 五.代码

只给出需要自己写的部分的代码：

cnnCost.m

```

1. function [cost, grad, preds] = cnnCost(theta,images,labels,numClasses,...
2.                                     filterDim,numFilters,poolDim,pred)
3. % Calcualte cost and gradient for a single layer convolutional
4. % neural network followed by a softmax layer with cross entropy
5. % objective.
6. %
7. % Parameters:
8. % theta      - unrolled parameter vector
9. % images     - stores images in imageDim x imageDim x numImges
10. %            array
11. % numClasses - number of classes to predict
12. % filterDim  - dimension of convolutional filter
13. % numFilters - number of convolutional filters
14. % poolDim    - dimension of pooling area
15. % pred       - boolean only forward propagate and return
16. %              predictions
17. %
18. %
19. % Returns:
20. % cost       - cross entropy cost
21. % grad       - gradient with respect to theta (if pred==False)
22. % preds      - list of predictions for each example (if pred==True)
23.
24.
25. if ~exist('pred','var')
26.     pred = false;
27. end;
28.
29.
30. imageDim = size(images,1); % height/width of image
31. numImages = size(images,3); % number of images
32.
33. %% Reshape parameters and setup gradient matrices
34.
35. % Wc is filterDim x filterDim x numFilters parameter matrix
36. % bc is the corresponding bias
37.
38. % Wd is numClasses x hiddenSize parameter matrix where hiddenSize
39. % is the number of output units from the convolutional layer
40. % bd is corresponding bias
41. [Wc, Wd, bc, bd] = cnnParamsToStack(theta,imageDim,filterDim,numFilters,...
42.                                     poolDim,numClasses);
43.
44. % Same sizes as Wc,Wd,bc,bd. Used to hold gradient w.r.t above params.
45. Wc_grad = zeros(size(Wc));
46. Wd_grad = zeros(size(Wd));
47. bc_grad = zeros(size(bc));
48. bd_grad = zeros(size(bd));
49.
50. %=====
51. %% STEP 1a: Forward Propagation
52. % In this step you will forward propagate the input through the
53. % convolutional and subsampling (mean pooling) layers. You will then use
54. % the responses from the convolution and pooling layer as the input to a
55. % standard softmax layer.
56.
57. %% Convolutional Layer
58. % For each image and each filter, convolve the image with the filter, add
59. % the bias and apply the sigmoid nonlinearity. Then subsample the
60. % convolved activations with mean pooling. Store the results of the
61. % convolution in activations and the results of the pooling in

```

```

62. % activationsPooled. You will need to save the convolved activations for
63. % backpropagation.
64. convDim = imageDim-filterDim+1; % dimension of convolved output
65. outputDim = (convDim)/poolDim; % dimension of subsampled output
66.
67. % convDim x convDim x numFilters x numImages tensor for storing activations
68. activations = zeros(convDim,convDim,numFilters,numImages);
69.
70. % outputDim x outputDim x numFilters x numImages tensor for storing
71. % subsampled activations
72. activationsPooled = zeros(outputDim,outputDim,numFilters,numImages);
73.
74. %%% YOUR CODE HERE %%%
75. activations = cnnConvolve(filterDim, numFilters, images, Wc, bc);
76. activationsPooled = cnnPool(poolDim, activations);
77. % Reshape activations into 2-d matrix, hiddenSize x numImages,
78. % for Softmax layer
79. activationsPooled = reshape(activationsPooled,[],numImages);
80.
81. % Softmax Layer
82. % Forward propagate the pooled activations calculated above into a
83. % standard softmax layer. For your convenience we have reshaped
84. % activationPooled into a hiddenSize x numImages matrix. Store the
85. % results in probs.
86.
87. % numClasses x numImages for storing probability that each image belongs to
88. % each class.
89. % probs = zeros(numClasses,numImages);
90. output = exp(Wd * activationsPooled + repmat(bd, 1, numImages));
91. probs = bsxfun(@divide, output, sum(output));
92. %%% YOUR CODE HERE %%%
93.
94. %=====
95. % STEP 1b: Calculate Cost
96. % In this step you will use the labels given as input and the probs
97. % calculate above to evaluate the cross entropy objective. Store your
98. % results in cost.
99.
100. %%% YOUR CODE HERE %%%
101. %cost = 0; % save objective into cost
102. idx = sub2ind(size(probs), labels', 1:numImages);
103. cost = -sum(log(probs(idx))) / numImages;
104.
105. % Makes predictions given probs and returns without backproagating errors.
106. if pred
107.     [~,preds] = max(probs,[],1);
108.     preds = preds';
109.     grad = 0;
110.     return;
111. end;
112.
113. %=====
114. % STEP 1c: Backpropagation
115. % Backpropagate errors through the softmax and convolutional/subsampling
116. % layers. Store the errors for the next step to calculate the gradient.
117. % Backpropagating the error w.r.t the softmax layer is as usual. To
118. % backpropagate through the pooling layer, you will need to upsample the
119. % error with respect to the pooling layer for each filter and each image.
120. % Use the kron function and a matrix of ones to do this upsampling
121. % quickly.
122.
123. %%% YOUR CODE HERE %%%
124. %Wd: numClasses x hiddenSize
125. %Wc: filterDim x filterDim x numFilters parameter matrix
126. %delta_L: numClasses x numImages
127. %delta_pooling: hiddenSize x numImages
128. %delta_conv: convDim x convDim x numFilters x numImages
129. %images: imageDim x imageDim x numImages
130. e = zeros(numClasses, numImages);
131. idx = sub2ind(size(e), labels', 1:numImages);
132. e(idx) = 1;
133. delta_L = probs - e;
134. delta_pooling = reshape((Wd' * delta_L) .* ones(size(activationsPooled)), outputDim, outputDi

```

```

m, numFilters, numImages);
135. delta_conv = zeros(convDim, convDim, numFilters, numImages);
136.
137. for i = 1 : numFilters
138.     for j = 1 : numImages
139.         delta_conv(:, :, i, j) = kron(delta_pooling(:, :, i, j), ones(poolDim, poolDim)) / (poolDim)^2;
140.     end
141. end
142.
143. %%=====
144. %% STEP 1d: Gradient Calculation
145. % After backpropagating the errors above, we can use them to calculate the
146. % gradient with respect to all the parameters. The gradient w.r.t the
147. % softmax layer is calculated as usual. To calculate the gradient w.r.t.
148. % a filter in the convolutional layer, convolve the backpropagated error
149. % for that filter with each image and aggregate over images.
150.
151. %%% YOUR CODE HERE %%%
152. for i = 1 : numFilters
153.     for j = 1 : numImages
154.         Wc_grad(:, :, i) = Wc_grad(:, :, i) + conv2(images(:, :, j), rot90(delta_conv(:, :, i, j)...
155.             .* activations(:, :, i, j) .* (1 - activations(:, :, i, j)), 2), 'valid');
156.         bc_grad(i) = bc_grad(i) + sum(sum(delta_conv(:, :, i, j)...
157.             .* activations(:, :, i, j) .* (1 - activations(:, :, i, j))));
158.     end
159.     Wc_grad(:, :, i) = Wc_grad(:, :, i) / numImages;
160.     bc_grad(i) = bc_grad(i) / numImages;
161. end
162.
163. for i = 1 : numImages
164.     Wd_grad = Wd_grad + delta_L(:, i) * activationsPooled(:, i)';
165. end
166. Wd_grad = Wd_grad / numImages;
167.
168. bd_grad = sum(delta_L, 2) / numImages;
169.
170. %% Unroll gradient into grad vector for minFunc
171. grad = [Wc_grad(:) ; Wd_grad(:) ; bc_grad(:) ; bd_grad(:)];
172. end

```

## minFuncSGD.m

```

1. function [opttheta] = minFuncSGD(funObj, theta, data, labels, ...
2.     options)
3. % Runs stochastic gradient descent with momentum to optimize the
4. % parameters for the given objective.
5. %
6. % Parameters:
7. % funObj - function handle which accepts as input theta,
8. % data, labels and returns cost and gradient w.r.t
9. % to theta.
10. % theta - unrolled parameter vector
11. % data - stores data in m x n x numExamples tensor
12. % labels - corresponding labels in numExamples x 1 vector
13. % options - struct to store specific options for optimization
14. %
15. % Returns:
16. % opttheta - optimized parameter vector
17. %
18. % Options (* required)
19. % epochs* - number of epochs through data
20. % alpha* - initial learning rate
21. % minibatch* - size of minibatch
22. % momentum - momentum constant, defaults to 0.9
23.
24.
25. %%=====
26. %% Setup
27. assert(all(isfield(options, {'epochs', 'alpha', 'minibatch'})), ...

```

```

28.         'Some options not defined');
29.     if ~isfield(options,'momentum')
30.         options.momentum = 0.9;
31.     end;
32.     epochs = options.epochs;
33.     alpha = options.alpha;
34.     minibatch = options.minibatch;
35.     m = length(labels); % training set size
36.     % Setup for momentum
37.     mom = 0.5;
38.     momIncrease = 20;% what??
39.     velocity = zeros(size(theta));
40.
41.     %%=====
42.     %% SGD loop
43.     it = 0;
44.     for e = 1:epochs
45.
46.         % randomly permute indices of data for quick minibatch sampling
47.         rp = randperm(m);
48.
49.         for s=1:minibatch:(m-minibatch+1)
50.             it = it + 1;
51.
52.             % increase momentum after momIncrease iterations
53.             if it == momIncrease
54.                 mom = options.momentum;
55.             end;
56.
57.             % get next randomly selected minibatch
58.             mb_data = data(:, :, rp(s:s+minibatch-1));
59.             mb_labels = labels(rp(s:s+minibatch-1));
60.
61.             % evaluate the objective function on the next minibatch
62.             [cost grad] = funObj(theta,mb_data,mb_labels);
63.
64.             % Instructions: Add in the weighted velocity vector to the
65.             % gradient evaluated above scaled by the learning rate.
66.             % Then update the current weights theta according to the
67.             % sgd update rule
68.
69.             %%% YOUR CODE HERE %%%
70.             velocity = mom * velocity - alpha * grad;
71.             theta = theta + velocity;
72.             fprintf('Epoch %d: Cost on iteration %d is %f\n',e,it,cost);
73.         end;
74.
75.         % anneal learning rate by factor of two after each epoch
76.         alpha = alpha/2.0;
77.
78.     end;
79.
80.     opttheta = theta;
81.
82. end

```

参考资料：

- 1.UFLDL
- 2.Neural networks and deep learning
- 3.Deep learning：五十一(CNN的反向求导及练习)