

SHANMUGA PERUMAL
HARIKUMAR

641831911

Dec 15 '16

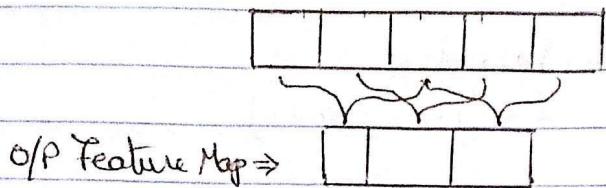
Deep Learning Final Exam

i) Backpropagation:

- a) Nb, the algorithm may get stuck at a local minimum depending upon initialization. Since we are using entire dataset to update weights, the algorithm has no way of escaping local minima.
- b) Depending upon the learning rate, the loss might get better or worse. If the learning rate is small enough, the loss must decrease.
- c) Nb, we might get stuck at a local minimum with no way of escaping from it.
- d) No, we are not trying to optimize that directly, hence during training, there is no guarantee the classification error never increases.

ConvNets : Basics:

- 2) a) Let us consider an feature of size $n \times 5$ and a kernel of size $K = 3$



The Size of the output feature map m is,

$$m = n - K + 1$$

\Rightarrow The Size of f output feature maps is
 $f * (n - K + 1)$

- b) Not Counting the biases, the total number of parameters is equal to the number of weights in the layer.

Given kernel size K and f such kernels, total number of parameters is $K * f$

- c) The process of convolving the input with kernels consists of multiplying and accumulation of the values in the input with kernel values.

For input dimension n and Kernel dimension K ,
Each layer would have $n - K + 1$ multiply-accumulate operations.
Considering f such kernels, the total number of multiply accumulate operations as $f * (n - K + 1)$

d)

With a stride of 's', the size of each feature map is

$$m = (n-k)/s + 1$$

For f such kernels

$$m = f * ((n-k)/s + 1)$$

3) Convnet Object Detection:

c)

Since the minimum height of pedestrians is 30, width must always be less than 15. We will be using a sliding window at different scales to find candidate regions. The minimum size of windows will be ~~30x60~~ 60×30 , in order to detect even the 15 pixel wide pedestrians.

b)

Since a sliding window at different scales is being used to detect candidate regions, scale normalization and cropping is taken care of.

To select negative examples, the following procedure can be followed:

1) Starting from training dataset, randomly sample N regions without pedestrians

2) Compute color histogram of regions and normalize it

3) Compute euclidean distance of cumulative color probability of each pair of negative regions

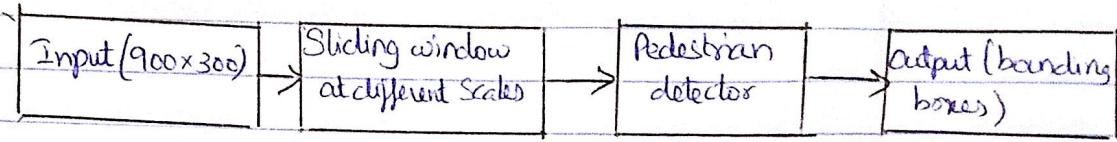
4) Select negative region with highest distance and remove from pool

5) Repeat 1 to required number of regions selected

c)

In order to reduce false negatives the sliding window is very accurate as it passes through all regions of image at different scales. Also, falsely negative classifications of images in validation set can be added to training set.

d)



The Sliding window slices the image into smaller patches at different scales. This is sent to the pedestrian detector, since a sliding window is used the size range of pedestrians depends on the limit on the scale at which the window searches. Also each patch is independently being processed, so multiple pedestrians will not be a problem -

The size of the bounding box would depend on the size of the patch pedestrian was found at -

4) ConvNets: Weak Supervision

a) The architecture consists of consecutive convolution and pooling layers followed by deconvolution layer consisting of deconvolution and unpooling layers. The reverse operations depend on the masks, weights in the forward operation. Based on the featturemap of highest abstraction level details of features are stored using deconvolutional layers. The entire feature map is then expanded and concatenated. This is followed by a convlution to generate class-specific activation map. These are finally softmaxed and aggregated into single vector, which is compared with image level vector label.

5) Word, Text and Image Embedding:

a) Distance Function, $D_w(\vec{x}_1, \vec{x}_2) = \|G_w(\vec{x}_1) - G_w(\vec{x}_2)\|_2$

where \vec{x}_1, \vec{x}_2 are pair of input vectors

then Loss function is

$$L(w) = \sum_{i=1}^P L(w, y, \vec{x}_1, \vec{x}_2)^i$$

where, $L(w, y, \vec{x}_1, \vec{x}_2)^i = (1-y)L_s(D_w) + yL_d(D_w)$

$\Leftrightarrow L_s \rightarrow$ partial loss function for similar pair
of points

$L_d \rightarrow$ Partial loss function for dissimilar
pair of points

The exact loss function being

$$L(w, y, \vec{x}_1, \vec{x}_2) = (1-y) \frac{1}{2} (D_w)^2 + (y) \frac{1}{2} \{ \max(0, m - D_w) \}^2$$

- b) A weighted least square regression model defined below can be used as an objective function in the given case of word embedding by GloVe

$$J = \sum_{i,j=1}^V f(x_{ij}) (w_i^T \tilde{w}_j + b_i + \tilde{b}_j - \log x_{ij})^2$$

where $V \rightarrow$ size of vocabulary

The weighting functions $f(n)$ should obey following properties

1) $f(0)=0$, if f is seen as continuous, it must vanish as $n \rightarrow 0$ fast enough that $\lim_{n \rightarrow 0} f(n) \log^2 n$ is finite

2) $f(n)$ should be non-decreasing, so that rare co-occurrences are not over-weighted.

3) $f(x)$ should be relatively small for large values of n , so that frequent co-occurrences are not over-weighted

c) Ord Grankng error function

$$\text{err}(\text{for}, y) = L(\text{rank}_y(f(x)))$$

where $\text{rank}_y(f(x))$ is the rank of true label y given by $f(x)$

$$\text{rank}_y(f(x)) = \sum_{i \neq y} I(f_i(x) \geq f_y(x))$$

where I is indicator function and $L(\cdot)$ transforms this into loss

$$L(K) = \sum_{j=1}^K \alpha_j \text{, with } \alpha_1 \geq \alpha_2 \geq \dots \geq 0$$

This lets us use different $L(\cdot)$ with different minimizers.

6) Recurrent Nets:

a) The gradients for Recurrent Neural Networks (RNNs) are calculated by an algorithm called Back Propagation Through Time (BPTT). This is similar to normal Back Propagation but the recurrent network is expanded to how many time steps are required and then backpropagations are done.

The gradient Jacobian Matrix elements are pointwise derivatives. Both tanh and Sigmoid functions have 0 derivatives at one point. This zero gradient drives gradients in previous layers towards zero. As those small valued matrices get multiplied, they shrink and vanish in few time steps. Gradients from far away time steps become zero and do not contribute in learning. This is vanishing gradient problem.

For other activation functions the matrix values can become large leading to exploding gradient problem.

b) Initializing the hidden weight matrices to identity, so when the network BPTT is carried out without inputs the matrices stay the same, which is not the case in non-identity initializations

c) LSTM [Long Short Term Memory] Networks are designed to remember information for a long time. Unlike normal RNNs, LSTMs have a different recurring module. Instead of single neural network layer, there are four interacting in special way. The cell state in LSTMs can pass information along uninterrupted. The LSTM modules have the power to change cell states, regulated by structures called gates.

Gates are a way to optionally let information through. LSTM has three such gates

Gated Recurrent Units (GRU) are similar to LSTMs but have fewer parameters and hence generally preferred. GRUs have two gates

i) Energy-Based Learning: Weakly Supervised Object Localization

a) The entire image may be separated into 32×32 or 50×50 slices. Each slice is given as input to a deep neural network of convolutional and subsampling layers which produce a final vector output.

This network is replicated for all slices, producing a vector for each slice. This process is repeated across multiple scales.

The network →

b) The loss function

$$L_0(w, o, x^i) = K \exp[-E(o, z, x^i)]$$

where K is positive constant

$$L(w, o, z, x) = E_w(o, z, x)$$

The loss function combine is

$$L(w, s) = \frac{1}{|S_1|} \sum_{i \in S_1} L(w, z^i, x^i) + \frac{1}{|S_0|} \sum_{i \in S_0} L_0(w, x^i)$$

S_1 - Subset of training sample with pedestrian

S_0 - Subset of training sample without pedestrian.

The network can be trained by optimizing a loss function of three labels, pedestrian presence, location and image. When three variables match the energy function is trained to have a smaller value and larger value if they don't match.

8) Unsupervised Learning and Auto-Encoders:

a) The shape of the energy function $J(y, w)$

in the space of y should be a surface with local points along the data points. The energy should increase as one moves away from the function that the learning algorithm is trying to approximate.

b)

objective function

$$J_{\text{sparseness}}(w, b) = J(w, b) + \beta \sum_{j=1}^{S_2} KL(p_j || \hat{p}_j)$$

where $w \rightarrow$ weights

$b \rightarrow$ bias

$\beta \rightarrow$ weighting of penalty term

$\ell \rightarrow$ sparsity parameter

$$KL(p_j || \hat{p}_j) = \ell \log \frac{\ell}{\hat{p}_j} + (1-\ell) \log \frac{1-\ell}{1-\hat{p}_j}$$

is the penalty term

Average activation of hidden layer $a_j^{(2)}$ is given by

$$\hat{a}_j^{(2)} = \frac{1}{m} \sum_{i=1}^m (a_j^{(2)})(n_i^{(1)})$$

Actual loss without penalty is

$$J(w, b) = \left[\frac{1}{m} \sum_{i=1}^m J(w, b; n_i^{(1)}, y_i^{(1)}) \right] + \frac{\lambda}{2} \sum_{l=1}^{L-1} \sum_{f=1}^{S_l} \sum_{j=1}^{S_{l+1}} (w_{j,l}^{(1)})^2$$

where

$$f(w, b) =$$

$$J(w, b; x, y) = \frac{1}{2} \|h_w(b(x)) - y\|^2$$
 is individual loss term

Sparsity Constraint:

Enforcing a constraint that

$$\hat{e}_j = e$$

where $\hat{e}_j \rightarrow$ average activation of hidden units j

$e \rightarrow$ sparsity parameter

sparsity parameter is generally a small value
To satisfy this constraint, hidden units activation
must be near 0.

- c) A denoising autoencoder is similar to a normal autoencoder, but the input is corrupted. A denoising autoencoder is trained to retrieve original input data from the corrupted data.

9) Optimization:

a) The standard deviation of inputs is inversely proportional to the learning rates

$$\sigma \propto \frac{1}{\gamma}, \quad \gamma \propto \frac{1}{\sigma}$$

$$\frac{\gamma_1}{\gamma_2} = \frac{1/\sigma_1}{1/\sigma_2} = \frac{\gamma_1}{\gamma_3} = 3:1$$

$$\gamma_1 : \gamma_2 = 3:1$$

where, σ - standard deviations γ - learning rates

b)

$$\text{Hinge loss} = \ell(y) = \max_{t \neq y} (0, 1 + \max(w_t^T x - w_t n))$$

the degree is 1

c)

This suggests that there exists numerous landscapes that have their own local minima, which the algorithm will reach. This means there are multiple approximations of the same function that the neural network will reach. This symmetry also suggests that the gradient descent algorithms will do similar actions in all such landscapes. Additionally since in a fully connected network every neuron affects every following neuron, the symmetry is understandable.

10) General Questions:

- a) Even though it is possible to approximate any well behaved function with a two layered network, the number of parameters required to do so is ~~not~~ feasible. So we need deep neural networks to better approximate. Also having more layers makes it easier to detect and learn hierarchical features.
- b) Since most real world functions are non-linear, deep networks use non-linear layers to better approximate such functions. Linear layers could also be replaced by the sum of such non-linearities as a single layer, which in turn loses the freedom to better approximate functions.
- c) The commonly used modules include
 - Convolutional layers
 - max pooling layers
 - Fully connected layers
 - flatten module
 - dropout layers
 - Non-linear activation module.
- d) Siamese networks contain two or more sub-networks that are similar in every sense, the weights, parameters, even parameter updating is mirrored in both networks. They are used to find similarity or connections between comparable things. ~~or e.g.) Signature Verification~~