

COMPSCI 671D Fall 2019
 Homework 4
 Vinayak Gupta
 vg101

1 EM Algorithm for Topic Modeling (35 points)

In this question we will try to design an algorithm for discovering the abstract topics that occur in a set of documents. In the problem, we have a set of M abstract documents in the universe, \mathcal{D} , which are mixtures of latent topics. We also have a dictionary of words, \mathcal{W} , with size N . Intuitively, \mathcal{W} is the list of all unique words that occur at least once within \mathcal{D} . Using these two sets we can denote the number of times word w_j occurs in document d_i as $n(w_j, d_i)$. It follows that we can represent a document d_i as a vector of size N , each entry of which corresponds to how many times word w_j appears in it. The topics z are latent variables in a topic set \mathcal{Z} sized K . Intuitively, if a document is about a certain topic, it will include some particular words with higher probability. For example, “music”, “show” and “film” appear frequently in documents about Arts while “school”, “student” and “teachers” usually occur if the document is about Education. Meanwhile, a document can be a mixture of several topics, e.g., a document can be about arts education for high school students. Inspired by the intuition, a reasonable approach to model the generative process is as follows.

$$\Pr(d_i, z_k, w_j) = \Pr(d_i) \Pr(z_k|d_i) \Pr(w_j|z_k)$$

which means the topics only depend on the documents and the words only depend on topics. For computational convenience, let $\Pr(d_i)$ be a uniform distribution and $\Pr(z_k|d_i)$ and $\Pr(w_j|z_k)$ be Multinomial distributions. i.e.

$$\begin{aligned} \Pr(d_i) &= \frac{1}{M} \\ \Pr(z_k|d_i) &= \alpha_{ik}, \quad \sum_{k=1}^K \alpha_{ik} = 1 \\ \Pr(w_j|z_k) &= \beta_{kj}, \quad \sum_{j=1}^N \beta_{kj} = 1. \end{aligned}$$

We are interested in learning α_{ik} and β_{kj} because they are substantively interesting: α_{ik} represents the proportion of topic k that makes up document i , and β_{kj} the probability of seeing word j under topic k . Therefore, a single word w_j in a document d_i is generated in this way: (1) Choose a topic z_k from the topic distribution $\Pr(z_k|d_i)$, the probability of choosing topic k is α_{ik} ; (2) Choose a word from the vocabulary distribution $p(w_j|z_k)$ in this topic, the probability is β_{kj} . We want to learn these parameters by maximizing the likelihood of the data using the EM algorithm.

- a) For our set of N documents and W word of choices, write down the log-likelihood of the model above, i.e. $\log(p(D = d_i, W = w_j | \alpha, \beta))$.

Remember that in the EM (Expectation-Maximization) algorithm, we can figure out the parameters and the hidden variables by iteratively computing (1) the distribution of latent variables using the old parameters and (2) find new parameters that maximize the likelihood using the old latent variable distribution. So, you can do the following:

- b) E-step: Derive the distribution of latent variable $p(z_k | w_j, d_i, \alpha^{old}, \beta^{old})$ by fixing the old parameters.
- c) M-step: Find the α^{new} and β^{new} that optimize the log likelihood using $p(z_k | w_j, d_i, \alpha^{old}, \beta^{old})$ as the distribution of z .

You would iterate these until convergence to get the final parameters in practice.

Just so you know, the model you have just been working with in this problem is a very famous and very popular model. :)

1. EM Algorithm for Topic Modeling

Universe : $D \Rightarrow M$ documents

Dictionary : $W \Rightarrow N$ words

Number of times w_j in document $d_i \Rightarrow n(w_j, d_i)$

Topic Set : $Z \Rightarrow K$ topics

$$P_{\pi}(d_i, z_k, w_j) = P_{\pi}(d_i) P_{\pi}(z_k | d_i) P_{\pi}(w_j | z_k)$$

$$P_{\pi}(d_i) = 1/M$$

$$P_{\pi}(z_k | d_i) = \alpha_{ik}, \sum_{k=1}^K \alpha_{ik} = 1$$

$$P_{\pi}(w_j | z_k) = \beta_{kj}, \sum_{j=1}^N \beta_{kj} = 1$$

$$(a) \text{ likelihood} = P(D_1, W_1, \dots, D_M, W_M | \Theta) \\ \text{where } \Theta = (\alpha, \beta)$$

$$= \prod_{i=1}^M \prod_{j=1}^N P_{\pi}(D=d_i, W=w_j | \Theta)^{n(w_j, d_i)}$$

where $n(w_j, d_i)$ term comes from the fact that there are $n(w_j, d_i)$ occurrences of word w_j in document d_i ; so we can group them (multiply them) together while calculating likelihood

$$\text{log likelihood} = \sum_{i=1}^M \sum_{j=1}^N n(w_j, d_i) \log P_{\pi}(D=d_i, W=w_j | \Theta) \\ = \sum_{i=1}^M \sum_{j=1}^N n(w_j, d_i) \log P_{\pi}(W=w_j | D=d_i, \Theta) \quad (1)$$

By law of total probability (Definition of joint probability)

$$\log P_{\pi}(D=d_i, W=w_j | \Theta) = \log P_{\pi}(W=w_j | D=d_i, \Theta) \cdot P_{\pi}(D=d_i | \Theta)$$

$$P_{\pi}(D=d_i | \Theta) = P_{\pi}(d_i), \text{ since } \Theta \text{ and } d_i \text{ are independent}$$

$$\text{Also, } P_{\pi}(d_i) = 1/M$$

$$\text{so } \log P_{\pi}(D=d_i, W=w_j | \Theta) = \log \frac{P_{\pi}(W=w_j | D=d_i, \Theta)}{M}$$



Scanned with
CamScanner

We also have latent variables (topics):

$$z_1, \dots, z_K$$

By law of Total Probability:

$$P_{\theta} (W=w_j | D=d_i, \Theta) = \sum_{k=1}^K P_{\theta} (W=w_j | z=z_k, D=d_i, \Theta) \cdot P_{\theta} (z=z_k | D=d_i, \Theta)$$

$$= \sum_{k=1}^K \beta_{kj} \alpha_{ik} \quad \text{as } P_{\theta} (W=w_j | z=z_k, D=d_i, \Theta) = \beta_{kj} \\ P_{\theta} (z=z_k | D=d_i, \Theta) = \alpha_{ik}$$

So ① becomes:

Ans
log likelihood = $\sum_{i=1}^M \sum_{j=1}^N n(w_j, d_i) \log \frac{\sum_{k=1}^K \beta_{kj} \alpha_{ik}}{M}$

(b) Since there is a log of sum in the likelihood, we can use Jensen's Inequality and EM method to maximize the log likelihood.

$$\log \sum_{k=1}^K \frac{\beta_{kj} \alpha_{ik}}{M} \Rightarrow \log \sum_{k=1}^K P_{\theta} (z=z_k | W=w_j, D=d_i, \Theta) \cdot \beta_{kj} \alpha_{ik} \\ M \cdot P_{\theta} (z=z_k | W=w_j, D=d_i, \Theta)$$

By Jensen's Inequality:

$$\geq \sum_{k=1}^K P_{\theta} (z=z_k | W=w_j, D=d_i, \Theta) \log \frac{\beta_{kj} \alpha_{ik}}{M \cdot P_{\theta} (z=z_k | W=w_j, D=d_i, \Theta)}$$

$$\text{so } A(\Theta, \Theta_t) = \sum_{i=1}^M \sum_{j=1}^N n(w_j, d_i) \sum_{k=1}^K \gamma_{ijk} \log \frac{\beta_{kj} \alpha_{ik}}{M \cdot \gamma_{ijk}}$$

where $\gamma_{ijk}^t = P_{\theta} (z=z_k | W=w_j, D=d_i, \Theta_t)$



Scanned with
CamScanner

E-step: compute γ_{ijk}^{t+1}

$$\gamma_{ijk}^{t+1} = P_{\theta_t}(z=z_k | w=w_i, D=d_i, \theta_t) = \frac{P_{\theta_t}(w=w_i | z=z_k, D=d_i, \theta_t) \cdot P_{\theta_t}(z=z_k | D=d_i, \theta_t)}{\sum_{k=1}^K P_{\theta_t}(w=w_i | z=z_k, D=d_i, \theta_t) \cdot P_{\theta_t}(z=z_k | D=d_i, \theta_t)}$$

$$\boxed{\gamma_{ijk}^{t+1} = \frac{\beta_{kj}^t \alpha_{ik}^t}{\sum_{k=1}^K \beta_{kj}^t \alpha_{ik}^t}} \quad \text{Ans} \Rightarrow \gamma_{ijk}^{t+1} = \frac{\beta_{kj}^{\text{old}} \alpha_{ik}^{\text{old}}}{\sum_{k=1}^K \beta_{kj}^{\text{old}} \alpha_{ik}^{\text{old}}}$$

(c)

M-step: Update α , β by setting derivatives of A to 0

Since there are constraints for α and β , we need to use Lagrangian:

$$\text{for } \alpha \quad L(\theta, \theta_t) = A(\theta, \theta_t) + \lambda \left(1 - \sum_{k=1}^K \alpha_{ik} \right)$$

$$\frac{dL(\theta, \theta_t)}{d\alpha_{(ik)^t}} = \frac{dA}{d\alpha_{(ik)^t}} - \lambda = 0$$

$$= \frac{d}{d\alpha_{(ik)^t}} \left(\sum_{i=1}^m \sum_{j=1}^n \sum_{k=1}^K n(w_i, d_i) \gamma_{ijk} \log \frac{\beta_{kj}^{\text{new}} \alpha_{ik}^{\text{new}}}{M \gamma_{ijk}} \right) = \lambda$$

Since β_{kj}^{new} , M , γ_{ijk} do not depend on $\alpha_{(ik)^t}$, remove them from above eqn.

$$= \frac{d}{d\alpha_{(ik)^t}} \left(\sum_{i=1}^m \sum_{j=1}^n \sum_{k=1}^K n(w_i, d_i) \gamma_{ijk} \log \alpha_{ik}^{\text{new}} \right) = \lambda$$

$$= \sum_{i=1}^m \sum_{j=1}^n n(w_i, d_i) \gamma_{ijk} \frac{\partial}{\partial \alpha_{ik}^{\text{new}}} = \lambda$$

$$\therefore \alpha_{ik}^{\text{new}} = \frac{\lambda}{\sum_{i=1}^m \sum_{j=1}^n n(w_i, d_i) \gamma_{ijk}}$$

$$\text{Also, } \sum_{k=1}^K \alpha_{ik}^{\text{new}} = 1 = \frac{1}{\lambda} \sum_{k=1}^K \sum_{j=1}^n n(w_i, d_i) \gamma_{ijk}$$



Scanned with
CamScanner

$$\alpha_{ik}^{\text{new}} = \frac{\sum_{i=1}^N n(w_i, d_i) \gamma_{ik}}{\sum_{k=1}^K \sum_{i=1}^N n(w_i, d_i) \gamma_{ik}}$$

for β

$$L(\theta, \theta_t) = A(\theta, \theta_t) + \lambda \left(1 - \sum_{j=1}^N \beta_{kj} \right)$$

$$\frac{dL(\theta, \theta_t)}{d\beta_{kj}^t} = \frac{dA}{d\beta_{kj}^t} - \lambda = 0$$

$$= \frac{d}{d\beta_{kj}^t} \left(\sum_{i=1}^M \sum_{j=1}^N \sum_{k=1}^K n(w_i, d_i) \gamma_{ik} \log \beta_{kj}^{\text{new}} \right) = \lambda$$

because $M, \alpha_{ik}^{\text{new}}, \gamma_{ik}$ do not depend on β_{kj}^t , so we removed them

$$\Rightarrow \sum_{i=1}^M n(w_i, d_i) \gamma_{ik} = \lambda$$

$$\beta_{kj}^{\text{new}} = \frac{\sum_{i=1}^M n(w_i, d_i) \gamma_{ik}}{\lambda}$$

$$\sum_{j=1}^N \beta_{kj}^{\text{new}} = \frac{\sum_{j=1}^N \sum_{i=1}^M n(w_i, d_i) \gamma_{ik}}{\lambda} = 1$$

$$\text{so } \lambda = \sum_{j=1}^N \sum_{i=1}^M n(w_i, d_i) \gamma_{ik}$$

$$\beta_{kj}^{\text{new}} = \frac{\sum_{i=1}^M n(w_i, d_i) \gamma_{ik}}{\sum_{j=1}^N \sum_{i=1}^M n(w_i, d_i) \gamma_{ik}}$$



2 Neural Networks (65 points)

In this problem you will get the chance to construct a neural network with architecture $784 - H - H - 1$, where H is the number of hidden nodes you choose. This neural network can be used for binary classification, such 0, 1 digits classification for MNIST. The model can be represented as

$$f(\mathbf{x}; \theta) : \mathbb{R}^{784} \rightarrow [0, 1]$$

where $\theta = \{\mathbf{W}_1 \in \mathbb{R}^{784 \times H}, \mathbf{b}_1 \in \mathbb{R}^H, \mathbf{W}_2 \in \mathbb{R}^{H \times H}, \mathbf{b}_2 \in \mathbb{R}^H, \mathbf{w}_3 \in \mathbb{R}^H, b_3 \in \mathbb{R}\}$. Explicitly, $f(\mathbf{x})$ is

$$\mathbf{h}_1 = \sigma(\mathbf{W}_1^\top \mathbf{x} + \mathbf{b}_1)$$

$$\mathbf{h}_2 = \sigma(\mathbf{W}_2^\top \mathbf{h}_1 + \mathbf{b}_2)$$

$$f(\mathbf{x}) = \sigma(\mathbf{w}_3^\top \mathbf{h}_2 + b_3)$$

where $\sigma(z) = \frac{1}{1+e^{-z}}$, and it is an element-wise operator, which means if $\mathbf{x} = (x^1, x^2, \dots, x^d) \in \mathbb{R}^d$, we have $\sigma(\mathbf{x}) = (\sigma(x^1), \sigma(x^2), \dots, \sigma(x^d))$.

The loss function for this model (or the negative log-likelihood) is the usual one:

$$\mathcal{L}(\theta) = - \sum_{i=1}^N (y_i \log f(\mathbf{x}_i) + (1 - y_i) \log(1 - f(\mathbf{x}_i)))$$

a) Backpropagation

Derive the gradients $\nabla_{w_3} \mathcal{L}(\theta)$, $\nabla_{b_3} \mathcal{L}(\theta)$, $\nabla_{w_2} \mathcal{L}(\theta)$, $\nabla_{b_2} \mathcal{L}(\theta)$, $\nabla_{w_1} \mathcal{L}(\theta)$, $\nabla_{b_1} \mathcal{L}(\theta)$ by backpropagation.

$$\text{2(a)} \quad z_1 = w_1^T x + b_1; \quad z_2 = w_2^T h_1 + b_2; \quad z_3 = w_3^T h_2 + b_3 \\ \text{where } h_1 = \sigma(z_1); \quad h_2 = \sigma(z_2) \quad \text{and } \sigma \text{ is a sigmoid function}$$

$$\mathcal{L}(\theta) = - \sum_{i=1}^N (y_i \log f(x_i) + (1-y_i) \log (1-f(x_i)))$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial f(x)} = \frac{f(x) - y}{f(x)(i_N - f(x))^T}; \quad \text{where } i_N \text{ is } (1 \dots 1)$$

$$\frac{\partial \mathcal{L}(\theta)}{\partial w_3} = \nabla_{w_3} \mathcal{L}(\theta) = \frac{\partial \mathcal{L}(\theta)}{\partial f(x)} \cdot \frac{\partial f(x)}{\partial z_3} \cdot \frac{\partial z_3}{\partial w_3}$$

$$\frac{\partial z_3}{\partial w_3} = h_2; \quad \cancel{\frac{\partial f(x)}{\partial z_3}} = f(x)(i_N - f(x))^T$$

$$\Rightarrow \nabla_{w_3} \mathcal{L}(\theta) = f(x) - y \cdot f(x)(i_N - f(x))^T \cdot h_2 \\ = [(f(x) - y) h_2]$$

$$\nabla_{b_3} \mathcal{L}(\theta) = \frac{\partial \mathcal{L}(\theta)}{\partial f(x)} \cdot \frac{\partial f(x)}{\partial z_3} \cdot \frac{\partial z_3}{\partial b_3} \\ = (f(x) - y) \cdot 1 = [f(x) - y]$$

$$\nabla_{w_2} \mathcal{L}(\theta) = \frac{\partial \mathcal{L}(\theta)}{\partial f(x)} \cdot \frac{\partial f(x)}{\partial z_3} \cdot \frac{\partial z_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial w_2} \cdot \frac{\partial z_2}{\partial h_2} \cdot \frac{\partial h_2}{\partial w_3} \\ = (f(x) - y) w_3 h_2 (i_N - h_2)^T h_1$$

$$\nabla_{b_2} \mathcal{L}(\theta) = \frac{\partial \mathcal{L}(\theta)}{\partial f(x)} \cdot \frac{\partial f(x)}{\partial z_3} \cdot \frac{\partial z_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial b_2} \cdot \frac{\partial z_2}{\partial h_2} \cdot \frac{\partial h_2}{\partial b_2} \\ = (f(x) - y) w_3 h_2 (i_N - h_2)^T$$

$$\nabla_{w_1} \mathcal{L}(\theta) = \frac{\partial \mathcal{L}(\theta)}{\partial f(x)} \cdot \frac{\partial f(x)}{\partial z_3} \cdot \frac{\partial z_3}{\partial h_2} \cdot \frac{\partial h_2}{\partial w_2} \cdot \frac{\partial z_2}{\partial h_2} \cdot \frac{\partial h_2}{\partial w_1} \\ = (f(x) - y) w_3 h_2 (i_N - h_2)^T w_2 h_1 (i - h_1)^T x$$



Scanned with
CamScanner

$$\nabla_{b_1} L(\theta) = \frac{\partial L(\theta)}{\partial b_1} = \frac{\partial f(x)}{\partial z_3} \frac{\partial z_3}{\partial h_2} \frac{\partial h_2}{\partial z_2} \frac{\partial z_2}{\partial h_1} \frac{\partial h_1}{\partial z_1} \frac{\partial z_1}{\partial b_1}$$

$$= (f(x) - y) w_3 h_2 (i_H - h_2)^T w_2 h_1 (i - h_1)^T$$

As we can see, the formulas recursively are

$$\nabla_{z_{l-1}} L(\theta) = \nabla_{z_l} L(\theta) w_l h_{l-1} (i - h_{l-1})^T$$

$$\nabla_{w_l} L(\theta) = \nabla_{z_l} L(\theta) h_{l-1}$$

$$\nabla_{b_l} L(\theta) = \nabla_{z_l} L(\theta)$$



Scanned with
CamScanner

b) Weight Initialization

Neural Networks are normally trained by gradient descent based optimization algorithms, for which we have to give initial values for our parameters (weight and bias). The bias parameters are usually initialized as 0's, while the weight parameters are usually initialized by independently sampling from $N(0, \sigma^2)$.

(i) What are the potential problems if we choose σ to be either too small or too large?

(ii) One nice property we would like our weights to have is that by proper initialization, we want the values before activation to be similar in distribution for each layer. It turns out we can achieve this by choosing σ properly. Now suppose we use ReLU activation in our neural network and the inputs are normalized. We want to choose the proper σ_l for the l -th hidden layer such that

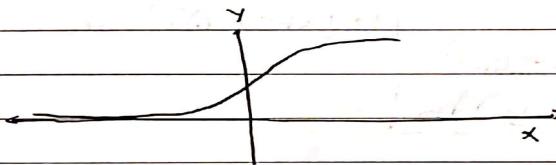
$$Var(z_l) = Var(z_{l-1}).$$

We assume the neurons in each layer are mutually independent and share the same distribution and we use z_l to represent the random variable that generates \mathbf{z}_l , where $\mathbf{z}_l = \mathbf{W}_l^T \mathbf{h}_{l-1} + \mathbf{b}_l$ and $\mathbf{h}_{l-1} = \max(0, \mathbf{z}_{l-1})$ and $\mathbf{W}_l \in \mathbb{R}^{n_{l-1} \times n_l}$.

Find σ_l such that $Var(z_l) = Var(z_{l-1})$ and prove your result.

(b) (ii) Weight parameters are initialized from $N(0, \sigma^2)$
 If σ is too small, all parameters are initialized close to 0. This causes the derivative w.r.t loss function nearly same for all weights. This cascades into subsequent iterations and all hidden neural nodes act similarly for all iterations thus defeating the purpose of neural networks.

If σ is very large, since the mean is 0, the individual weights are initialized with very high values. If the weight are very high or very low, $z = W^T h + b$ is also very high or very low. Activation fns. like sigmoid have a vanishing gradient at extremely high/low values as can be seen by their graph. This causes slow learning and network might get stuck at a local optima.



Sigmoid graph.

ii) $V_{an}(\text{ReLU}) = \frac{1}{2} V_{ar}(x)$ where $\text{ReLU} = \max(0, x)$
 because $f(x) = x$'s right part (+ve part) is same as RELU's +ve part and the negative part is just 0. This can also be shown as:

$$\begin{aligned} V_{an}(\text{ReLU}) &= E(\text{ReLU}^2) = \int_{-\infty}^{\infty} \text{ReLU}^2 p(\text{ReLU}) dx \\ &\quad (\text{Inputs normalized}) = \int_{-\infty}^{\infty} \max(0, x)^2 p(x) dx \\ &= \int_0^{\infty} \max(0, x)^2 p(x) dx = \int_0^{\infty} x^2 p(x) dx \\ &= \frac{1}{2} \int_{-\infty}^{\infty} x^2 p(x) dx \quad \text{since } x^2 \text{ is symmetric} \end{aligned}$$

Also, $E(x) = 0$ as for $f(x) = x$, $E(x) = 0$

$$\text{so } = \frac{1}{2} \int_{-\infty}^{\infty} (x - E(x))^2 p(x) dx = \frac{1}{2} V_{ar}(x)$$



✓ Hidden layers

$$\text{Var}(z_k^d) = \text{Var}\left(\sum_j w_{kj}^d a_j^{d-1}\right)$$

kth observation

$$= \sum_{j=1}^H (\mathbb{E}[w_{kj}^d]^2 \text{Var}(a_j^{d-1}) + \text{Var}(w_{kj}^d) \mathbb{E}[a_j^{d-1}]^2 + \text{Var}(w_{kj}^d) \text{Var}(a_j^{d-1}))$$

$$= \sum_{j=1}^H \text{Var}(w_{kj}^d) \text{Var}(a_j^{d-1}) \quad \text{because } \mathbb{E}[w_{kj}^d] = 0 \text{ (mean = 0)}$$

$$= h \text{Var}(w^d) \text{Var}(a^{d-1})$$

as variance of all individual weights same and

variance of all a 's also same

$$\text{so } \text{Var}(z^d) = h \text{Var}(w^d) \frac{\text{Var}(z^{d-1})}{2} \quad \left(\text{as } \text{Var}(a^{d-1}) = \frac{\text{Var}(z^{d-1})}{2} \right)$$

$$\text{so } \text{Var}(w^d) = \frac{2}{h} \text{ if } \text{Var}(z^d) = \text{Var}(z^{d-1})$$

$$\text{so } \sigma_w^2 = \frac{2}{h}$$

	$\sigma_w^2 = \sqrt{\frac{2}{h}}$	Ans
--	-----------------------------------	-----



c) Implement your Neural Network

Use the preprocessed MNIST data provided in the attachment. MNIST is a dataset of images of handwritten digits that also contain labels for the number they correspond to. We want to train a neural network to learn to recognize these handwritten digits. Use gradient descent with the gradients you have derived in part a) to train a neural network on the data. Report the accuracy for the network on the train and test data. Since the sample size is large, we can approximate the true gradient with stochastic gradient for computational convenience:

$$\nabla \mathcal{L}(\theta) \approx \frac{1}{B} \sum_{j=1}^B \nabla \mathcal{L}(\theta)_j$$

which means at each step, instead of calculating the gradients with all samples in the dataset, we randomly pick a mini-batch of B samples and calculate the gradient using these samples. In implementing the neural network, it would be beneficial to build and train the model in a general way. That is, build a NN with d hidden layers such that we can easily modify d .

The code was developed using the framework from the coursera deeplearning.ai course I did https://www.coursera.org/learn/neural-networks-deep-learning?fbclid=IwAR0bRdW3V3cUGbR1vi1v7-SCgJBGpg4ag9_2mmbOXOn-FEJtu8FDsHtUtLQ

```

1 import numpy as np
2 import h5py
3 import matplotlib.pyplot as plt
4 import time
5 import scipy.io as sio
6 from PIL import Image
7 from scipy import ndimage
8 import pylab
9 np.random.seed(1)

10
11 def sigmoid(Z):
12     A = 1/(1+np.exp(-Z))
13     cache = Z
14     return A, cache

15
16
17 def sigmoid_diff(Z):
18
19     s = 1/(1+np.exp(-Z))
20     dZ = s * (1-s)
21     return dZ

22
23 def initialize_parameters_deep(layer_dims):
24     parameters = {}
25     L = len(layer_dims)
26     np.random.seed(1)
27     for l in range(1, L):
28
29         parameters['W' + str(l)] = np.random.randn(layer_dims[l], layer_dims[l-1])
30         * 0.01
31         parameters['b' + str(l)] = np.zeros((layer_dims[l], 1))
32
33     return parameters

34 def linear_forward(A, W, b):

```

```

35     Z = np.dot(W,A)+b
36     cache = (A, W, b)
37     return Z, cache
38
39 def linear_activation_forward(A_prev, W, b):
40
41     Z, linear_cache = linear_forward(A_prev, W, b)
42     A, activation_cache = sigmoid(Z)
43
44     cache = (linear_cache, activation_cache)
45
46     return A, cache
47
48 def L_model_forward(X, parameters):
49     caches = []
50     A = X
51     L = len(parameters) // 2
52
53
54     for l in range(1, L):
55         A_prev = A
56         A, cache = linear_activation_forward(A_prev,parameters['W' + str(l)] ,
57 parameters['b' + str(l)])
58         caches.append(cache)
59
60         AL, cache = linear_activation_forward(A, parameters['W' + str(L)] , parameters
61 ['b' + str(L)])
62         caches.append(cache)
63
64     return AL, caches
65
66 def compute_cost(AL,Y):
67     m=Y.shape[1]
68     cost=-(1/m)*(np.dot(Y,np.transpose(np.log(AL)))+np.dot(1-Y,np.transpose(np.log
69 (1-AL))))
70     cost=np.squeeze(cost)
71     return cost
72
73 def linear_backward(dZ, cache):
74
75     A_prev, W, b = cache
76     m=A_prev.shape[1]
77     dW = (1/m)*(np.dot(dZ,np.transpose(A_prev)))
78     db = (1/m)*(np.dot(dZ,np.ones((dZ.shape[1],1))))
79     dA_prev = np.dot(np.transpose(W),dZ)
80
81     return dA_prev, dW, db
82
83 def linear_activation_backward(dA, cache):
84
85     linear_cache, activation_cache = cache
86
87     dZ = dA*sigmoid_diff(activation_cache)
88     dA_prev, dW, db = linear_backward(dZ, linear_cache)
89
90     return dA_prev, dW, db
91
92 def L_model_backward(AL, Y, caches):

```

```

91
92     grads = {}
93     L = len(caches) # the number of layers
94     m = AL.shape[1]
95     Y = Y.reshape(AL.shape) # after this line, Y is the same shape as AL
96
97     dAL = - (np.divide(Y, AL) - np.divide(1 - Y, 1 - AL))
98
99     current_cache = caches[L-1]
100    grads["dA" + str(L-1)], grads["dW" + str(L)], grads["db" + str(L)] =
101        linear_activation_backward(dAL, current_cache)
102
103    for l in reversed(range(L-1)):
104
105        current_cache = caches[l]
106        dA_prev_temp, dW_temp, db_temp = linear_activation_backward(grads["dA" +
107            str(l + 1)], current_cache)
108        grads["dA" + str(l)] = dA_prev_temp
109        grads["dW" + str(l + 1)] = dW_temp
110        grads["db" + str(l + 1)] = db_temp
111
112    return grads
113
114 def update_parameters(parameters, grads, learning_rate):
115
116     L = len(parameters) // 2
117
118     for l in range(L):
119         parameters["W" + str(l+1)] = parameters["W" + str(l+1)]-(learning_rate*
120             grads["dW" + str(l + 1)])
121         parameters["b" + str(l+1)] = parameters["b" + str(l+1)]-(learning_rate*
122             grads["db" + str(l + 1)])
123
124     return parameters
125
126 def L_layer_model(X, Y, layers_dims, check, learning_rate = 0.1, num_iterations =
127     3000):
128
129     costs = []
130     parameters = initialize_parameters_deep(layers_dims)
131
132     for i in range(0, num_iterations):
133
134         AL, caches = L_model_forward(X, parameters)
135
136         cost = compute_cost(AL, Y)
137
138         grads = L_model_backward(AL, Y, caches)
139
140         if(check==1):
141             return np.linalg.norm(grads['dW1'], ord='fro')
142
143         parameters = update_parameters(parameters, grads, learning_rate)

```

```

145     return parameters
146
147 def predict(X, y, parameters):
148
149
150     m = X.shape[1]
151     n = len(parameters)
152     p = np.zeros((1,m))
153
154
155     probas, caches = L_model_forward(X, parameters)
156
157
158     for i in range(0, probas.shape[1]):
159         if probas[0,i] > 0.5:
160             p[0,i] = 1
161         else:
162             p[0,i] = 0
163
164
165     print("Accuracy: " + str(np.sum((p == y)/m)))
166
167     return p
168
169 tmp = sio.loadmat("mnist_all.mat")
170 x1=tmp['train0']
171 x2=tmp['train1']
172 x3=tmp['test0']
173 x4=tmp['test1']
174
175 y1=[[0] for i in range(x1.shape[0])]
176 y2=[[1] for i in range(x2.shape[0])]
177 y3=[[0] for i in range(x3.shape[0])]
178 y4=[[1] for i in range(x4.shape[0])]
179
180 train_x_orig=np.concatenate((x1,x2))
181 test_x_orig=np.concatenate((x3,x4))
182 train_y=np.concatenate((y1,y2))
183
184 test_y=np.concatenate((y3,y4))
185
186
187 layers_dims = [784]
188 d=2
189 for j in range(d):
190     layers_dims.append(20)
191 layers_dims.append(1)
192
193 train_y=np.transpose(train_y)
194 test_y=np.transpose(test_y)
195
196 m_train = train_x_orig.shape[0]
197 num_px = train_x_orig.shape[1]
198 m_test = test_x_orig.shape[0]
199
200 print ("Number of training examples: " + str(m_train))
201 print ("Number of testing examples: " + str(m_test))
202 print ("Each image is of size: (" + str(num_px) + ", " + str(num_px) + ", 3)")
203 print ("train_x_orig shape: " + str(train_x_orig.shape))

```

```

204 print ("train_y shape: " + str(train_y.shape))
205 print ("test_x_orig shape: " + str(test_x_orig.shape))
206 print ("test_y shape: " + str(test_y.shape))
207
208
209
210 train_x_flatten = train_x_orig.reshape(train_x_orig.shape[0], -1).T
211 test_x_flatten = test_x_orig.reshape(test_x_orig.shape[0], -1).T
212
213
214 train_x = train_x_flatten/255.
215 test_x = test_x_flatten/255.
216
217 print ("train_x's shape: " + str(train_x.shape))
218 print ("test_x's shape: " + str(test_x.shape))
219
220 parameters = L_layer_model(train_x, train_y, layers_dims,0, num_iterations = 3000)
221
222
223
224 pred_train = predict(train_x, train_y, parameters)
225
226 pred_test = predict(test_x, test_y, parameters)
227
228 layers_dims = [784]
229 for i in range(1,11):
230
231     layers_dims.append(20)
232     layers_dims.append(1)
233     if i==1:
234         normval=[L_layer_model(train_x, train_y, layers_dims, 1,num_iterations =
235         1)]
236     else:
237         normval.append(L_layer_model(train_x, train_y, layers_dims,1,
238         num_iterations = 1))
239     layers_dims.pop()
240
241 print(normval)
242
243 fig = plt.figure(figsize=(20,10))
244 ax = fig.add_subplot(2, 2, 1)
245 line, = ax.plot(normval, color='blue', lw=2)
246
247 ax.set_yscale('log')
248 plt.xlabel('Number of hidden layers')
249 plt.ylabel('Log of frobenius norm of dW1')
250 pylab.show()

```

Accuracy: 0.9978681405448085

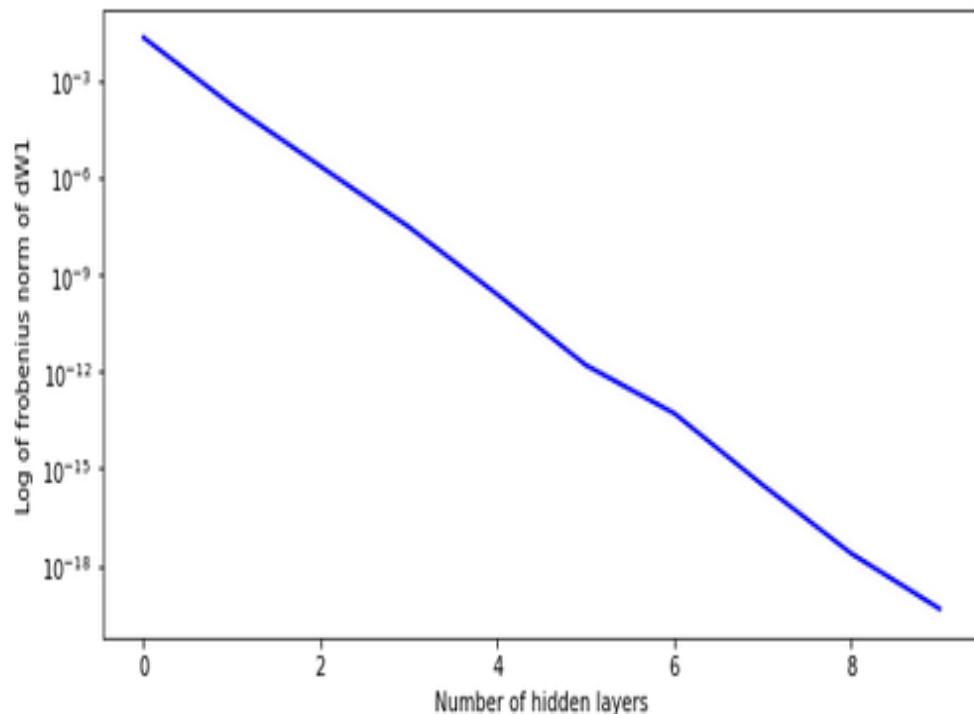
Accuracy: 0.9990543735224583

d) Vanishing Gradient

(i) In our current model, we have 2 hidden layers. What is the magnitude of $\nabla_{\mathbf{W}_1} \mathcal{L}(\theta)$ (you can use Frobenius norm or induced matrix norm)? Now try to modify the the number of hidden layers d , from 1 to 10. What happens to the magnitude of $\nabla_{\mathbf{W}_1} \mathcal{L}(\theta)$ when we increase d ? Plot the magnitude against d on log scale. In this question, you don't need to train the whole neural network, just initialize it with random weights, e.g. $\mathcal{N}(0, 1)$, and calculate the gradient.

The above code in c also contains the code for d(i). The graph of log of Frobenius Norm with number of hidden layers are shown below along with their actual values. We can see that the Norm decreases with increasing hidden layers.

```
[0.020929657601947354,  
 0.00016853966838273834,  
 2.206016544453013e-06,  
 2.9407276490753038e-08,  
 2.408563213511275e-10,  
 1.5851636785460938e-12,  
 5.1170501426214545e-14,  
 3.127652475763494e-16,  
 2.3980611774373226e-18,  
 4.6875823257695946e-20]
```



(ii) Suppose that all weights are finite and the weight matrices $\mathbf{W}_2, \mathbf{W}_3, \dots, \mathbf{W}_d$ are diagonalizable matrices in which the absolute value of the eigenvalues are upper-bounded by $4 - \epsilon$. Prove that

$$\lim_{d \rightarrow \infty} \nabla_{\mathbf{W}_1} \mathcal{L}(\theta) = \mathbf{0}$$

d (ii) In a), we got the following relation:

$$\nabla_{z_{e-1}} L(\theta) = \nabla_{z_e} L(\theta) w_e h_{e-1} (i - h_{e-1})^T$$

where $i = [1 \dots 1]$

We can use Cauchy-Schwartz inequality on the following eqn:

$$\frac{\nabla_{z_{e-1}} L(\theta)}{\nabla_{z_e} L(\theta)} = w_e h_{e-1} (i - h_{e-1})^T$$

Using inequality:

$$\left\| \frac{\nabla_{z_{e-1}} L(\theta)}{\nabla_{z_e} L(\theta)} \right\|_2 = \|w_e h_{e-1} (i - h_{e-1})^T\|_2 \leq \|w_e\|_2 \|h_{e-1}\|_2 \|i - h_{e-1}\|_2$$

$\|h_{e-1} (i - h_{e-1})^T\|_2$ can have maximum value $1/4$ as $h(1-h)$ if $h \in (0,1)$ is $1/4$.

This means that $\left\| \frac{\nabla_{z_{e-1}} L(\theta)}{\nabla_{z_e} L(\theta)} \right\|_2$ is bounded by $1 - \frac{(4-\epsilon)}{4}$,

as $\|w_e\|_2$ is bounded by $4 - \epsilon$. (Because eigenvalues are bounded by $4 - \epsilon$)

Now, suppose there are d hidden layers, then

$$\left\| \frac{\nabla_{z_{d+1}} L(\theta)}{\nabla_{z_d} L(\theta)} \right\|_2 = \left\| \frac{\nabla_{z_1} L(\theta)}{\nabla_{z_2} L(\theta)} \cdot \frac{\nabla_{z_2} L(\theta)}{\nabla_{z_3} L(\theta)} \cdots \frac{\nabla_{z_d} L(\theta)}{\nabla_{z_{d+1}} L(\theta)} \right\|_2$$

By using Cauchy-Schwartz inequality, we get

$$\begin{aligned} &\leq \left\| \frac{\nabla_{z_1} L(\theta)}{\nabla_{z_2} L(\theta)} \right\|_2 \left\| \frac{\nabla_{z_2} L(\theta)}{\nabla_{z_3} L(\theta)} \right\|_2 \cdots \left\| \frac{\nabla_{z_d} L(\theta)}{\nabla_{z_{d+1}} L(\theta)} \right\|_2 \\ &= \left(\frac{4-\epsilon}{4} \right)^d \leq 1^d \end{aligned}$$

so this $\frac{4-\epsilon}{4} < 1$ and as $d \rightarrow \infty$ $\left(\frac{4-\epsilon}{4} \right)^d \rightarrow 0$



Hence $\lim_{d \rightarrow \infty} \left\| \frac{\nabla_{z_1} L(\theta)}{\nabla_{z_{d+1}} L(\theta)} \right\|_2 = 0$

$$\text{Also, } \|\nabla_{z_d} L(\theta)\| = \|f(x) - y\|$$

This value is finite as it depends only on observations which are constant, so $\|\nabla_{z_d} L(\theta)\|_2 \rightarrow 0$

$$\text{Also, } \nabla_{w_d} L(\theta) = \nabla_{z_d} L(\theta) h_{d-1}$$

Again by Cauchy-Schwarz:

$$\|\nabla_{w_d} L(\theta)\| = \|\nabla_{z_d} L(\theta) h_{d-1}\|_2 \leq \|\nabla_{z_d} L(\theta)\|_2 \|h_{d-1}\|_2$$

Any Sigmoid fn. can have max. 2 norm. 1 so $\|h_{d-1}\|_2 \leq 1$

$$\therefore \|\nabla_{w_d} L(\theta)\|_2 \leq \|\nabla_{z_d} L(\theta)\|_2 \|h_{d-1}\|_2 < \|\nabla_{z_d} L(\theta)\|_2$$

$$\|\nabla_{w_d} L(\theta)\|_2 \rightarrow 0 \text{ as } d \rightarrow \infty \text{ because } \|\nabla_{z_d} L(\theta)\|_2 \rightarrow 0$$

$$\text{so } \nabla_{w_d} L(\theta) \rightarrow 0$$

Hence, $\lim_{d \rightarrow \infty} \nabla_{w_d} L(\theta) = 0$



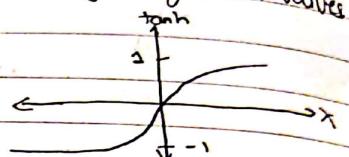
Scanned with
CamScanner

(iii) The result above is known as the vanishing gradient problem in the feedforward neural network. Will there still be such a problem if we use a tanh activation? What if we use a ReLU activation?

Give an intuitive explanation of why vanishing gradient happens in light of the shapes of activation functions and provide 2 ways of dealing with this problem.

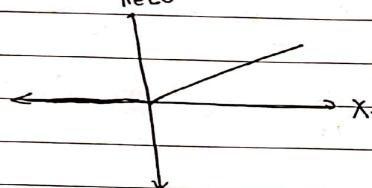
d(iii) Yes, tanh behaves in a similar fashion to sigmoid function. While sigmoid function maps $R \rightarrow [0, 1]$, while tanh maps $R \rightarrow (-1, 1)$. Other than that the shapes are very similar with flat gradients for very high/low values. Also, the formula for tanh is:

$$\tanh(x) = \frac{e^x - e^{-x}}{e^x + e^{-x}}$$



However, in case of ReLU, the gradient is non-zero for activated part ($x > 0$). There is no vanishing gradient problem in case of ReLU.

ReLU

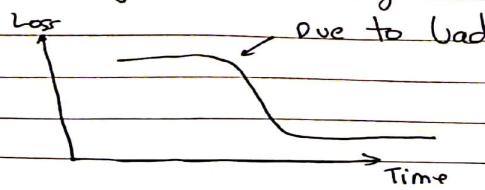


Vanishing Gradient can be ~~dealt~~ dealt with in the following ways:

1) Batch Normalization: Outputs of several nodes (mini-batch) in same layers are normalized. The mean and standard deviation are treated as separate parameters which are required to be learned. This process is done before non-linear activation function and adds regularization and helps gradient flow in the network.

2) Leaky ReLU: It also removes vanishing gradients as it works similar to ReLU ($\max(0.1x, x)$).

3. Initialization of network weights is also important.



Scanned with
CamScanner

e) Sigmoid Activation

One potential problem of using the sigmoid activation function for hidden layers is that it maps from \mathbb{R} to $[0, 1]$, which means the output of the sigmoid function is always positive. What happens to the signs of the gradients of hidden layer weights in this situation if we only have one example? What if we are adding gradients up across a batch of data?

2(e)

If we are using only one example, the gradients are given by the following formula:

$$\frac{dE}{dW_{a,b}} = \frac{dE}{dO_b} \cdot dO_b \cdot \frac{dnet_b}{dW_{a,b}}$$

$$= \left(\sum_{e \in L} \delta_e w_{b,e} \right) O_b(1-O_b) O_a$$

where l are the output nodes from b , a is the node from where weight originates, b is the node to which it comes.

Suppose, there are 2 nodes (a_1, a_2 from which weight originates) and one to which converges b , then

$$\frac{dE}{dW_{a_1,b}} = \left(\sum_{e \in L} \delta_e w_{b,e} \right) O_b(1-O_b) O_{a_1}$$

$$\frac{dE}{dW_{a_2,b}} = \left(\sum_{e \in L} \delta_e w_{b,e} \right) O_b(1-O_b) O_{a_2}$$

The formula shows that sign of $\frac{dE}{dW_{a,b}}$ will remain same for all incoming weights $w_{b,e}$ to a node. This is a limitation as either all weights increase/decrease simultaneously.

If we are adding gradients up across a batch of data, this issue is AVOIDED. Since, we sum gradients, while the sign of gradients for individual observation for all incoming weights are same, their sum may have different signs. This can be shown by the following example.

Obs 1 : $\frac{dE}{dW_{a_1,b}} = -2 \quad \frac{dE}{dW_{a_2,b}} = -2$

Obs 2 : $\frac{dE}{dW_{a_1,b}} = +1 \quad \frac{dE}{dW} = +3$



Scanned with
CamScanner

The signs of gradients are same but while summing across observations:

$$\frac{dE}{dW_{a_1,b}} = -2+1 = -1 \quad ; \quad \frac{dE}{dW_{a_2,b}} = -2+3 = +1$$

We are able to get different signs. (because of different magnitudes of σ_{a_1} and σ_{a_2})



Scanned with
CamScanner