

Sol: 3

(a) Mini batch gradient descent seeks to find a balance between the robustness of stochastic gradient descent and efficiency of batch gradient descent. It is the most common.

The Advantages :-

- The model update frequency is higher than batch gradient descent with which allows for a more robust convergence, avoiding local minima.
- Also, the batched update provides a computationally more efficient process.
- Batching, allows both the efficiency of not having all training data in memory (case when batch size = 1 or without mini batch) and algorithm implementations.

(b) A partial solution to the vanishing / Exploding gradients in NN is better or more careful choice of the random initialization of weights.

Suppose we have symmetrically initialized our weights.

$w[l] = \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix} > \text{Identity}$ l : hidden layer number.

Then if we ignore biases for a minute,

$y' = w[l] w[l-1] \dots w[2] w[1] x$

$\Rightarrow y' = w[l] \begin{bmatrix} 1.5 & 0 \\ 0 & 1.5 \end{bmatrix}^{(L-1)} x = (1.5)^L$ (which will be very large)
(exploding gradients)

also, if $w[l] = \begin{bmatrix} 0.9 & 0 \\ 0 & 0.9 \end{bmatrix} < \text{Identity}$.

$y' = w[l] \begin{bmatrix} 0.9 & 0 \\ 0 & 0.9 \end{bmatrix}^{(L-1)} x = (0.9)^L$ (which will be very small)
(vanishing gradients)

(c) Most often used regularization technique is L2 regularization

$J(w, b) = (1/m) * \text{Sum}(L(y(i), y'(i))) + \text{lambda}/2m * \text{Sum}(|w[i]|^2)$

lambda: regularization parameter.

- (i) If lambda is too large - a lot of w's will be close to zeroes which will make the NN simpler.
- (ii) If lambda is good enough, it will just reduce some weights that make the neural network overfit. Also, it will just make some of tanh activations roughly linear which will prevent overfitting.

Other methods to deal with overfitting :- ^{Collect} More data & try a different model

(d) Batch Normalization

We generally ~~are~~

Before training the model, we normalize the input by subtracting the mean and dividing by variance. This helped a lot for the shape of the cost function and for reaching the minimum point faster.

So, for any hidden layers, ~~can be~~ ^{we can} normalize $A[l]$ to train $w[l+1]$ and $b[l]$ faster.

The batch normalization reduces the problem of input values changing (shifting).

Sol: 4

The number of parameters here (without including bias terms)
will be = ~~$(3 \times 3) \times (3)$~~

$$(f_h \times f_w) \times (\text{no. of input channels}) \times (\text{no. of output channels})$$

$$= 3 \times 3 \times 3 \times 8 = 216$$

$$\text{Input image} = N_1 \times N_2 \times 3 \quad \cdot \quad N_1: \text{height} \quad N_2: \text{width}.$$

$$\text{So, output image dimensions : } O_H = (N_1 - f_h + 2P)/s + 1$$

$$\text{So, } O_H = (N_1 - 3 + 2)/2 + 1$$

$$O_H = (N_1 - 1)/2 + 1 = \left(\frac{N_1 + 1}{2}\right)$$

$$O_W = (N_2 - f_w + 2P)/s + 1$$

$$= (N_2 - 3 + 2)/2 + 1$$

$$O_W = (N_2 - 1)/2 + 1$$

$$\text{Output image} = O_H \times O_W \times (\text{no. of output channels})$$

$$= O_H \times O_W \times 8$$

Classification algorithms such as logistic regression learn a probability distribution over classes conditioned on the input instead of directly producing an output y because

~~As After producing~~

After using different activations functions to be scale down the output values of our algorithm to $(0,1)$ is more efficient and more useful.

Introducing loss function which would directly classify the output to either 0 & 1 (in binary classification) will introduce large scale errors into the model, and so accuracy of the model will be low.

After Using loss functions such and activations functions such that the output is a probability, we can easily classify them on the basis of what they are more likely to be..

Q-2

Stochastic Gradient Descent with Momentum.

SGD with momentum is a method which helps accelerate gradient vectors in the right directions, thus leading to fast fast convergences.

For fast convergence of our model, what we want to do with the data is, we want some kind of 'moving' average which would 'denoise' the data and bring it closer to the original function. This can be achieved by using an exponential exponentially weighted averages, which are defined as

$$V_t = \beta V_{t-1} + (1-\beta) S_t$$

(V : new sequence; S : old sequence)

For We can use this exponential weighted averages to average over our gradients. Momentum is the moving average of our gradients. Our update equation can be written as

$V_t = \beta V_t + (1-\beta) \nabla_w L(W, x, y)$ $W = W - \alpha V_t$	<p>L: Loss function ∇_w: gradient w.r.t weight α: learning rate.</p>
--	---

Momentum in Physics refers to a push or the effect something is possessed by the body due to its motion and its mass. The momentum here accelerates the Gradient descent algorithm at so, the convergence is achieved faster and more smoothly.