## **Quiz-03 Results for Ziyu Han**

① Correct answers will be available on Feb 6 at 12am.

Score for this attempt: **7.5** out of 10

Submitted Feb 5 at 3:55pm

This attempt took 2,386 minutes.

## **Question 1**

1 / 1 pts

For this question, please read the paper: Rumelhart, Hinton and Williams (1986 (http://www.cs.toronto.edu/~hinton/absps/naturebp.pdf)
) (http://www.cs.toronto.edu/~hinton/absps/naturebp.pdf)

[Can be found at: http://www.cs.toronto.edu/~hinton/absps/naturebp.pdf]

One version of gradient descent changes each weight by an amount proportional to the accumulated  $\delta E/\delta w$ .

$$\Delta w = -\epsilon rac{\delta E}{\delta w}$$

Select all that are true about this method:

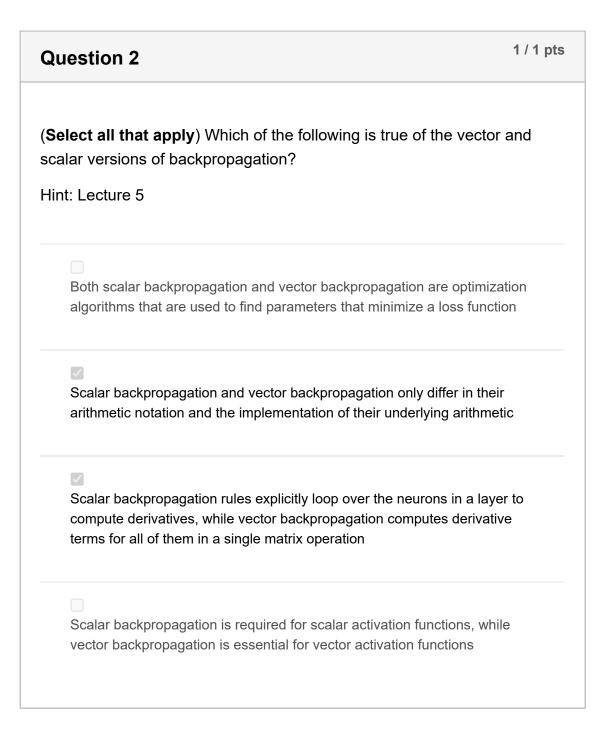
This method converges as rapidly as methods that make use of second derivatives.

- It cannot be implemented by local computations in parallel hardware.
- It's simpler than methods that use second derivatives.

"This method does not converge as rapidly as methods which make use of the second derivatives, but it is much simpler [...]" p535

It can be improved without sacrificing simplicity and locality.

"It can be significantly improved, without sacrificing the simplicity and locality, [...]" p535



Question 3 1/1 pts

Consider a perceptron in a network that has the following vector activation:

$$y_j = \prod_{j 
eq i} z_i$$

Where  $y_i$  is the j-th component of column vector y, and  $z_i$  is the i-th component of column vector z. Using the notation from lecture, which of the following is true of the derivative of y w.r.t. z? (select all that are true)

Hint: Vector calculus notes 1 (and beyond)

- It will be a matrix whose diagonal entries are all 0.
- It is a row vector whose i-th component is given by  $\prod_{i \neq i} z_i$
- It is a matrix whose (i,j)th component is given by  $z_i z_j$
- It is a matrix whose (i, j)th component where  $i \neq j$  is given by  $\prod_{k \neq i, k \neq j} z_k$
- $oxed{\hspace{0.5cm}}$  It is a column vector whose i-th component is given by  $\prod_{j \neq i} z_j$

Question 4	
Backpropagation can be applied to any differentiable activation function	

- True
- False

1 / 1 pts

You can propagate derivatives backward through any differentiable activation function, or even activations that only have a finite subgradient.

Partial Question 5 0.5 / 1 pts

(Select all that apply) At any point, the gradient of a scalar function with multivariate inputs...

Hint: Lecture 4, "Gradient of a scalar function of a vector" and "properties of a gradient".

- Is parallel to equal-value contours of the function
- Is in the direction of steepest ascent
- Is the vector of local partial derivatives w.r.t. all the inputs
- Is in the direction of steepest descent

Question 6

Which of the following are valid subgradients of a RELU, given by Relu(x), at x = 0? We will represent the subgradient as  $\nabla_{subgrad}RELU(x)$ 

Hint: Lecture 5, slides 112-114.

- $\square \ 
  abla_{subgrad}RELU\left(0
  ight)=0$
- $lacksquare 
  abla_{subgrad}RELU\left( 0
  ight) =1$
- $lacksquare 
  abla_{subgrad}RELU\left(0
  ight) = -0.5$
- $\square \ 
  abla_{subgrad}RELU(0) = 0.5$

## Incorrect Question 7 0 / 1 pts

The KL divergence between the output of a multi-class network with softmax output  $y=[y_1\dots y_K]$  and desired output  $d=[d_1\dots d_K]$  is defined as  $KL=\sum_i d_i \log d_i - \sum_i d_i \log y_i$ . The first term on the right hand side is the entropy of d, and the second term is the Crossentropy between d and g, which we will represent as Xent(y,d). Minimizing the KL divergence is strictly equivalent to minimizing the crossentropy, since  $\sum_i d_i \log d_i$  is not a parameter of network parameters. When we do this, we refer to Xent(y,d) as the cross-entropy loss.

Defined in this manner, which of the following is true of the cross-entropy loss Xent(y, d)? Recall that in this setting both y and d may be viewed as probabilities (i.e. they satisfy the properties of a probability distribution).

- ✓ It is always non-negative
- It's derivative with respect to  $m{y}$  goes to zero at the minimum (when  $m{y}$  is exactly equal to  $m{d}$  )
- $\ lue{}$  It goes to 0 when  $oldsymbol{y}$  equals  $oldsymbol{d}$
- ☐ It only depends on the output value of the network for the correct class

If d is not one hot (e.g. when we use label smoothing), the cross entropy may not be 0 when d = y.

For one-hot d, we saw in class that the KL divergence is equal to the cross entropy. Also, in this case, at d=y, the gradient of the DL divergence (and therefore Xent(y,d)) is not 0.

Question 8

In order to maximize the possibility of escaping local minima and finding the global minimum of a generic function, the best strategy to manage step sizes during gradient descent is:

Hint: Lecture 6, "Issues 2"

To start with a large, non-divergent step size (e.g. less than twice the optimal step size for a quadratic approximation at the initial location) and gradually decrease it over iterations

To maintain a step size consistently close to the optimal step size (e.g. close to the inverse second derivative at the current estimate)

To keep the step size low throughout to prevent divergence into a local minima

To start with a large, divergent step size (e.g. greater than twice the optimal step size for a quadratic approximation at the initial location) and gradually decrease it over iterations

See lecture for explanation.

Gradient descent with a fixed step size \_\_\_\_\_\_ for all convex functions (Fill in the blank)

Hint: Lecture 6

Always converges to a local minimum

Does not always converge

Always converges to a global minimum

Always converges to some point

## Question 10 1/1 pts

Let f(.) be a scalar-valued function with multivariate input and  $x=[x_1,x_2]$  be a two-component vector such that y=f(x). y is being minimized using RProp from lecture. In the k-th iteration, the derivative of y with respect to  $x_1$  is  $\frac{dy}{dx_1}=2$ , the derivative of y with respect to  $x_2$  is  $\frac{dy}{dx_2}=-1$ . As a result,  $x_1$  has a step size of  $\Delta x_1^{(k)}=1$  and  $x_2$  has a step size of  $\Delta x_2^{(k)}=1$ . At the (k+1)-th iteration, the derivative of y with respect to  $x_1$  is  $\frac{dy}{dx_1}=0.5$  and the derivative of y with respect to  $x_2$  is  $\frac{dy}{dx_2}=1$ . Which of the following is true about the step size at the (k+1)-th iteration?

Hint: Lecture 6, RProp

$$\bigcirc$$
  $\Delta x_1^{(k+1)} > 1$  and  $\Delta x_2^{(k+1)} > 1$ 

$$\bigcirc$$
  $\Delta x_1^{(k+1)} < 1$  and  $\Delta x_2^{(k+1)} > 1$ 

$$\bigcirc$$
  $\Delta x_1^{(k+1)} < 1$  and  $\Delta x_2^{(k+1)} < 1$ 

$$igotimes \Delta x_1^{(k+1)} > 1$$
 and  $\Delta x_2^{(k+1)} < 1$ 

Quiz Score: 7.5 out of 10