# Quiz-03 Results for Jiayi Huang

⚠ Correct answers are hidden.

Score for this attempt: 8 out of 10
Submitted Sep 14 at 11:48pm
This attempt took 23 minutes.

⁞

PartialQuestion 1

0.5 / 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 $\delta E/\delta w$.

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

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.
"[...] and can easily be implemented by local computations in parallel hardware [...]" p535

☑ 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 2

1 / 1 pts

(**Select all that apply**) Which of the following is true of the vector and scalar versions of backpropagation?

Hint: Lecture 5

☐
Both scalar backpropagation and vector backpropagation are optimization algorithms that are used to find parameters that minimize a loss function

☐

Scalar backpropagation is required for scalar activation functions, while vector backpropagation is essential for vector activation functions

☑

Scalar backpropagation rules explicitly loop over the neurons in a layer to compute derivatives, while vector backpropagation computes derivative terms for all of them in a single matrix operation

☑

Scalar backpropagation and vector backpropagation only differ in their arithmetic notation and the implementation of their underlying arithmetic

⠿

## Question 3

1 / 1 pts

Which of the following is true given the Hessian of a scalar function with multivariate inputs?

Hint: Lec 4 "Unconstrained minimization of a function". Also note that an eigen value of 0 indicates that the function is flat (to within the second derivative) along the direction of the corresponding Hessian Eigenvector.

☑   The eigenvalues are all strictly negative at a local maximum.

☑   The eigenvalues are all strictly positive at a local minimum.

☐   The eigenvalues are all strictly positive at global minima, but not at local minima.

☐   The eigenvalues are all non-negative at local minima.

⠿

## Question 4

1 / 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 in the direction of steepest ascent

☐   Is in the direction of steepest descent

☐   Is parallel to equal-value contours of the function

☑   Is the vector of local partial derivatives w.r.t. all the inputs

⠿

## Question 5

1 / 1 pts

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

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

Where $y_j$ 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 (lecture 5, slide 135 and beyond)

☐ It is a matrix whose (i,j)th component is given by $z_i z_j$

☐ It is a column vector whose i-th component is given by $\prod_{j \neq i} z_j$

☑ It is a matrix whose (i, j)th component where $i \neq j$ is given by $\prod_{k \neq i, k \neq j} z_k$

☐ It is a row vector whose i-th component is given by $\prod_{j \neq i} z_j$

☑ It will be a matrix whose diagonal entries are all 0.

⋮⋮

IncorrectQuestion 6

0 / 1 pts

Tom decides to construct a new vector activation function based on the RELU or softplus to output probabilities.  He considers the following two options:

1)  $SR(z_i) = \dfrac{RELU(z_i)}{\sum_j RELU(z_j)}$

2)  $SR(z_i) = \dfrac{Softplus(z_i)}{\sum_j Softplus(z_j)}$

where $z_i$ are the affine terms just before the activation, where he defines the softplus as  Softplus(z) = log(1+exp(z).

Which of the following statements is true.

Hint: Lecture 5, Slides 99-115.

○ The RELU-based activation is a valid vector activation that produces probabilities, whereas the softplus-based one is not.

○ The Softplus-based activation is a valid vector activation that produces probabilities, whereas the RELU-based one is not.

◉ Both are valid vector activations to produce probabilities.

⋮⋮

Question 7

1 / 1 pts

The KL divergence between the output of a multi-class network with softmax output $y = [y_1 \ldots y_K]$ and *desired* output $d = [d_1 \ldots d_K]$ 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 *Cross-entropy* between $d$ and $y$ , which we will represent as $Xent(y, d)$ . Minimizing the KL divergence is strictly equivalent to minimizing the cross entropy, 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's derivative with respect to $y$ goes to zero at the minimum (when $y$ is exactly equal to $d$ )

☐ It only depends on the output value of the network for the correct class

☐ It goes to 0 when $y$ equals $d$

☑ It is always non-negative

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
1 / 1 pts

Gradient descent yields a solution that is not sensitive to how a network's weights are initialized.

Hint: Basic gradient descent from lecture 5 - slide 5

○ True

◉ False

⠿

PartialQuestion 9
0.5 / 1 pts

Which of the following update rules explicitly computes second-order derivatives or their approximations? (select all that apply)

Hint: second half of Lecture 6

☐ Quickprop

☐ RProp

☐ Gradient descent

☑ Newton's Method

⁞

## Question 10

1 / 1 pts

Let f be a quadratic function such that at $x = 1$, $f(x) = 10$, $f'(x) = -4$, and $f''(x) = 1$. The

minimum has a value of $x =$ | 5 | and a value of $f(x) =$ | 2 | . (Truncate

your answer to 1 digit after the decimal point i.e. enter your answer in the format x.x,  e.g. 4.5)

Hint: Lecture 6 "Convergence for quadratic surfaces"

**Answer 1:**
5
**Answer 2:**
2

Quiz Score: 8 out of 10