M2Q1, M2Q2

## Question 1

• [3 points] We use gradient descent to find the minimum of the function $f(x) = \frac{1}{4} x^4$ with step size $\eta > 0$. If we start from the point $x_0 = 5$, how small should $\eta$ be so we make progress in the first iteration? Check all values of $\eta$ that make progress.
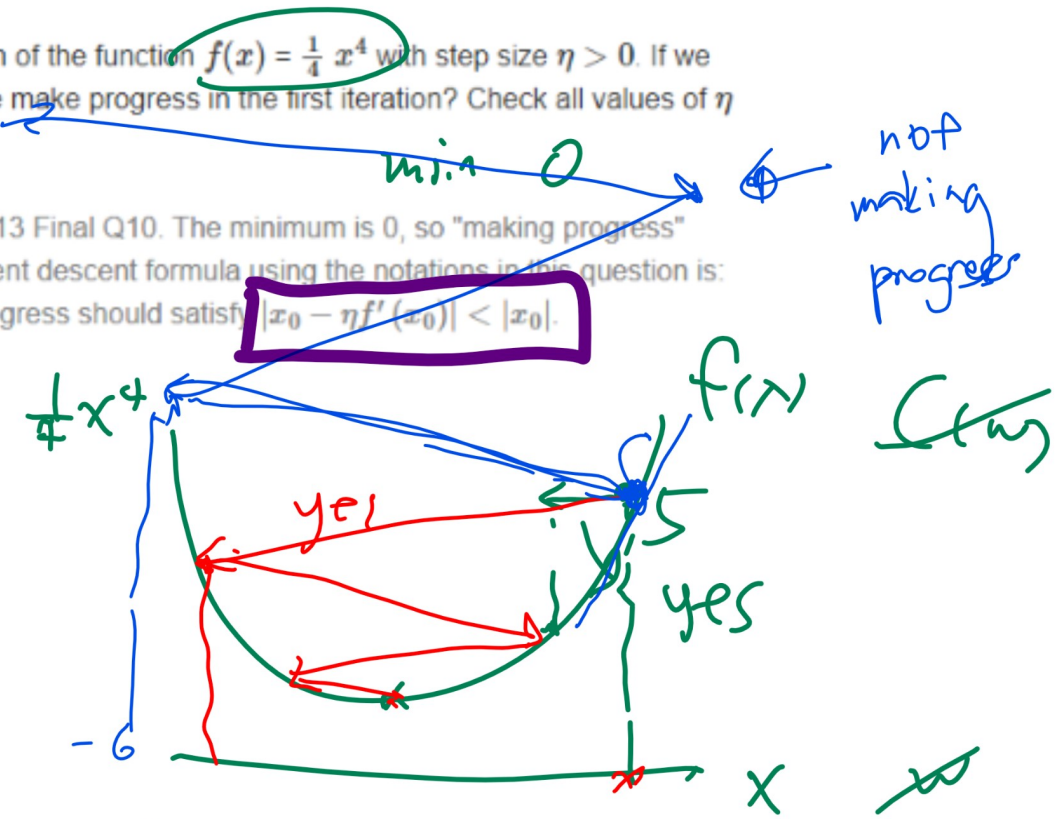
▼ Hint

See Fall 2017 Final Q7, Fall 2014 Midterm Q17, Fall 2013 Final Q10. The minimum is 0, so "making progress" means getting closer to 0 in the first iteration. The gradient descent formula using the notations in this question is:
$x_1 = x_0 - \eta f'(x_0)$. The learning rate $\eta$ that makes progress should satisfy $|x_0 - \eta f'(x_0)| < |x_0|$.

• Choices:
  - ☐ 0.1544
  - ☐ 0.0488
  - ☐ 0.0563
  - ☐ 0.0438
  - ☐ 0.0742
  - ☐ None of the above

• Calculator: [＿＿＿＿＿] [Calculate]

min 0   →   ⊕  not making progress

$\frac{1}{4} x^4$   $f(x)$   $f(x)$

yes     yes

-6

x

$|5 - \eta \cancel{x}^{\;5^3} | < |5|$

$\begin{cases} 5 - \eta 5^3 < 5 \\ \eta 5^3 - 5 < 5 \end{cases} \implies \eta < \underline{\quad}$

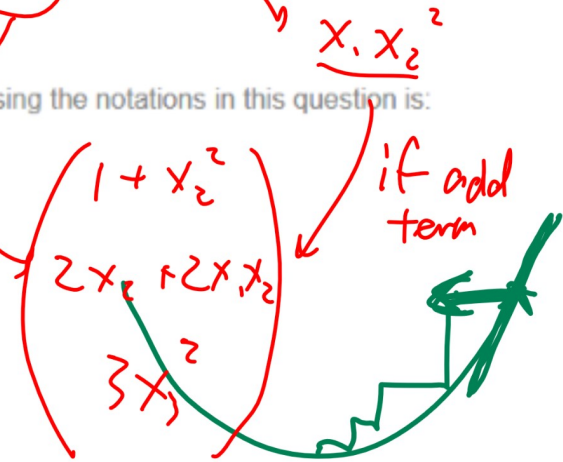converge getting closer to local min. as t → ∞

# Question 2

• [3 points] Let $x = (x_1, x_2, x_3)$. We want to minimize the objective function $f(x) = x_1 + x_2^2 + x_3^3$ using gradient descent. Let the stepsize $\eta = 0.42$. If we start at the vector $x^{(0)} = [2, 5, 1]$, what is the next vector $x^{(1)}$ produced by gradient descent?

▼ Hint

See Fall 2017 Final Q15, Fall 2010 Final Q5. The gradient descent formula using the notations in this question is:

$$x^{(1)} = x^{(0)} - \eta \nabla f(x^{(0)}) \text{ where } \nabla f(x^{(0)}) = \begin{bmatrix} \frac{\partial f}{\partial x_1^{(0)}} \\ \frac{\partial f}{\partial x_2^{(0)}} \\ \frac{\partial f}{\partial x_2^{(0)}} \end{bmatrix}$$

• Answer (comma separated vector): [          ]. [Calculate]

*Handwritten:*

M2 Q2

$x \cdot x_2^2$

$\begin{pmatrix} 1 + x_2^2 \end{pmatrix}$ if add term

$2x_2 + 2x_1 x_2$

$3x_3^2$

$\begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix}^{(0)} = \begin{pmatrix} 2 \\ 5 \\ 1 \end{pmatrix}$

$\begin{pmatrix} 2 \\ 5 \\ 1 \end{pmatrix} - 0.42 \begin{pmatrix} 1 \\ 2 \times 5 \\ 3 \times 1^2 \end{pmatrix} \rightarrow \begin{pmatrix} 1 \\ 10 \\ 3 \end{pmatrix}$

# Question 8

• [1 points] You want to design a neural network with sigmoid units to predict the academic role from his webpage. Possible roles are "professor" (label 0), "student" (label 1), "staff" (label 2). Suppose each person can take on only one of these roles at the same time. The neural network uses one-hot encoding, label 0 is encoded by $\begin{bmatrix} 1 \\ 0 \\ 0 \end{bmatrix}$, label 1 is encoded by $\begin{bmatrix} 0 \\ 1 \\ 0 \end{bmatrix}$, and label 2 is encoded by $\begin{bmatrix} 0 \\ 0 \\ 1 \end{bmatrix}$. What is the role (enter a label, not a string) if the output is

$\begin{bmatrix} 0.8 \\ 0.42 \\ 0.88 \end{bmatrix}$

▼ Hint

See Fall 2011 Midterm Q12. It is the label corresponding to the largest output value.
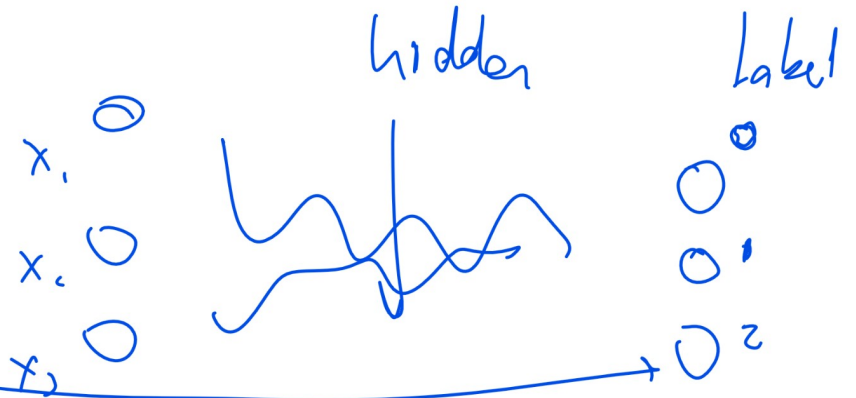
• Answer: [     2     ].

*Handwritten:*

M3 Q8

$a_i = (0,1)$

$a_i \geq 0.5 \Rightarrow \hat{y} = 1$

$a_i < 0.5 \Rightarrow \hat{y} = 0$

$a_i \quad 1 - a_i$

$P_n \, 1 \quad P_n \, 0$

$\begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$ $\begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix}$ $\begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}$

Softmax $\longrightarrow$ $\dfrac{e^{-z_1}}{e^{-z_1}+e^{-z_2}-1e^{-z_3}}$ $\quad$ $\dfrac{e^{-z_2}}{\sim}$ $\quad$ $\dfrac{e^{-z_3}}{\sim}$

$\underbrace{\qquad\qquad}$ Sum up to $\underline{\underline{1}}$

$(X_1 \quad y_1)$ $\qquad$ $(X_2, y_2)$

$\begin{pmatrix} 5 \\ 2 \\ 1 \end{pmatrix}, 2$ $\qquad$ $\begin{pmatrix} 3 \\ 3 \\ 3 \end{pmatrix}, 0$

hidden $\qquad$ Label

$x_1 \; \bigcirc$

$x_2 \; \bigcirc$ $\qquad$ $\bigcirc \; 0$

$x_3 \; \bigcirc$ $\qquad$ $\bigcirc \; 1$

$\qquad \bigcirc \; 2$

$\left[ \begin{pmatrix} 5 \\ 2 \\ 1 \end{pmatrix} \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix}, \begin{pmatrix} 3 \\ 3 \\ 3 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \right]$
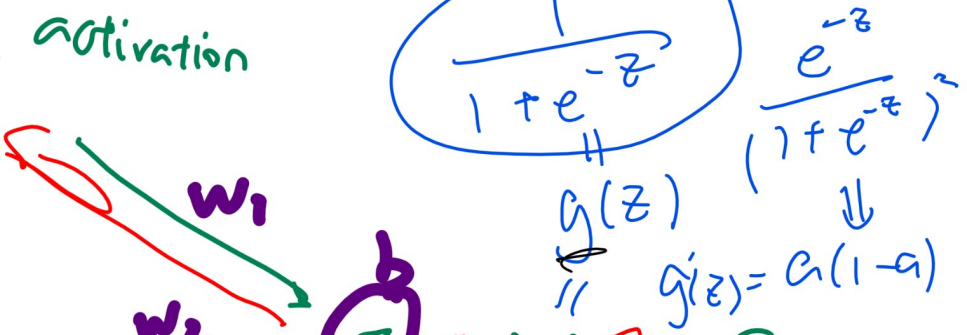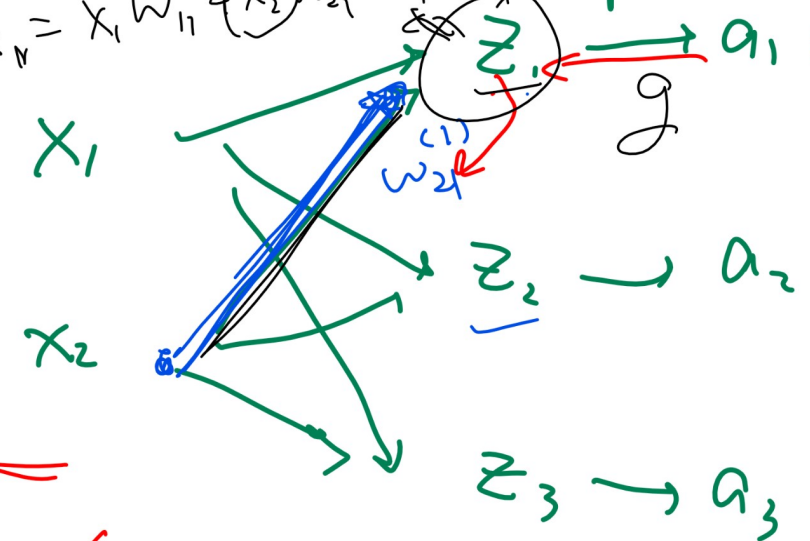
$C = (y - a)^2$

$C = \sum_{j=1}^{3} \left( \begin{pmatrix} 0 \\ 0 \\ 1 \end{pmatrix} - \begin{pmatrix} 0.5 \\ 0.2 \\ 0.3 \end{pmatrix} \right)^2$

# SGD for network

$z_1 = x_1 W_{11} + (x_2) W_{21} + b_1$    linear part    → activation

$z_1 \to a_1$    $g$

$x_1$

$w_{21}^{(1)}$

$z_2 \to a_2$

$x_2$

$z_3 \to a_3$

$w_1$

$w_2$

$w_3$

$z \leftarrow a \leftarrow C$

$$\frac{1}{1+e^{-z}}$$    $\frac{e^{-z}}{(1+e^{-z})^2}$

$g(z)$    $\Downarrow$

$g'(z) = a(1-a)$

$a_1 w_1 + a_2 w_2 + a_3 w_3 + b$    $\frac{1}{2}(a-y)^2$    ← focus on a for item i

2/21/28

$6 + 3 + 3 + 1 \Rightarrow 13$ parameters

weights   biases   ws   bias
in $L1$          in $L$          $(0,1)$    $(0.2)$    $(-1,1)$

start

$w \in [-1, 1]$ random    $w = [\text{Math.random}() * 2 - 1]$

→ back prop

$W_{21} = W_{21} - \alpha \boxed{\dfrac{\partial C}{\partial w_{21}}}$
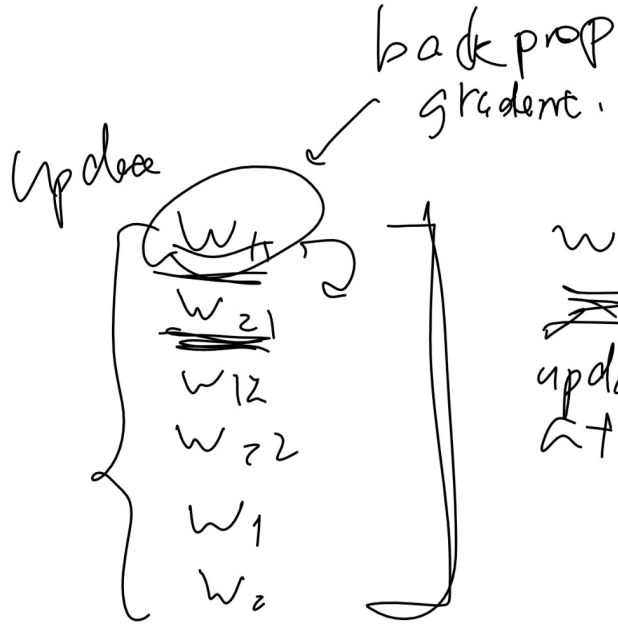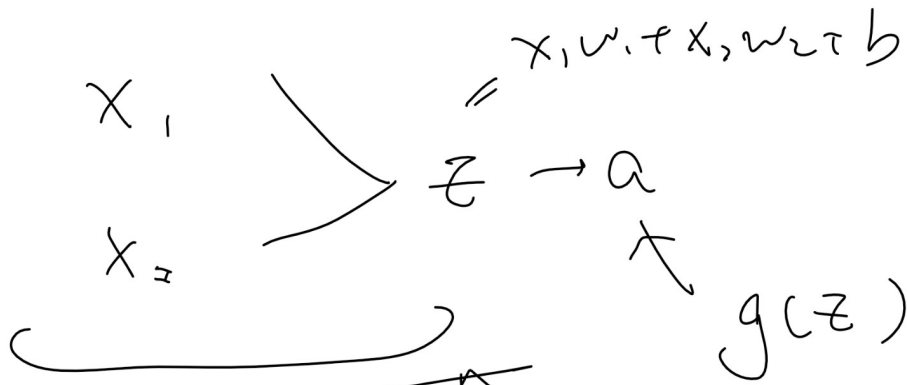
$$\frac{\partial C}{\partial a} \cdot \frac{\partial a}{\partial z} \cdot \frac{\partial z}{\partial a_1} \cdot \frac{\partial a_1}{\partial z_1} \cdot \frac{\partial z_1}{\partial w_{21}}$$

$(a-y)$    $a(1-a)$    $w_1$    $a_1(1-a_1)$    $x_2$

# Single Layer

$x_1$

$= x_1 w_1 + x_2 w_2 + b$

$z \rightarrow a$

$x_2$

$g(z)$

back prop
gradient.

feed forward

change each
iteration

Update

$\begin{bmatrix} w_{11} \\ w_{21} \\ w_{12} \\ w_{72} \\ w_1 \\ w_2 \end{bmatrix}$

$w_{11}$ from $t-1 \Longrightarrow$    all   $a, a_1, \dots$

/epoch.

update
at the same time.

$x_i$ $y_i$ stay
same