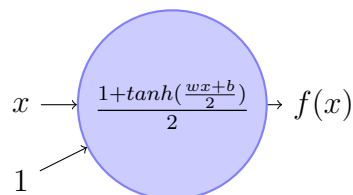


Instructions:

- This assignment is meant to help you grok certain concepts we will use in the course. Please don't copy solutions from any sources.
- Avoid verbosity.
- Questions marked with * are relatively difficult. Don't be discouraged if you cannot solve them right away!
- The assignment needs to be written in latex using the attached tex file. The solution for each question should be written in the solution block in space already provided in the tex file. **Handwritten assignments will not be accepted.**

1. Partial Derivatives

(a) Consider the following computation ,



where $f(x) = \frac{1 + \tanh(\frac{wx+b}{2})}{2}$ and by definition : $\tanh(z) = \frac{e^z - e^{-z}}{e^z + e^{-z}}$

The value L is given by,

$$L = \frac{1}{2}(y - f(x))^2$$

Here, x and y are constants and w and b are parameters that can be modified. In other words, L is a function of w and b .

Derive the partial derivatives, $\frac{\partial L}{\partial w}$ and $\frac{\partial L}{\partial b}$.

Solution:

$$\begin{aligned} h(t) &= \tanh(t) = \frac{e^t - e^{-t}}{e^t + e^{-t}} \\ h'(t) &= \frac{(e^t + e^{-t})(e^t + e^{-t}) - ((e^t - e^{-t})(e^t - e^{-t}))}{(e^t + e^{-t})^2} \\ h'(t) &= 1 - \frac{(e^t - e^{-t})^2}{(e^t + e^{-t})^2} \\ h'(t) &= 1 - \left(\frac{e^t - e^{-t}}{e^t + e^{-t}} \right)^2 \end{aligned}$$

$$h'(t) = 1 - \tanh^2(t) \tag{1}$$

to calculate $\frac{\partial L}{\partial w}$

Apply chain rule $\frac{dg(u)}{dw} = \frac{dg}{du} \cdot \frac{du}{dw}$

$$g = u^2, \quad u = \left(y - \frac{1 + \tanh\left(\frac{wx+b}{2}\right)}{2} \right)$$

$$= \frac{1}{2} \frac{\partial}{\partial u} (u^2) \frac{\partial}{\partial w} \left(y - \frac{1 + \tanh\left(\frac{wx+b}{2}\right)}{2} \right)$$

$$= \frac{1}{2} * 2 * u \left(\frac{\partial}{\partial w} (y) - \frac{\partial}{\partial w} \left(\frac{1 + \tanh\left(\frac{wx+b}{2}\right)}{2} \right) \right)$$

Using eq(1) $\frac{\partial}{\partial w} \left(\frac{1 + \tanh\left(\frac{wx+b}{2}\right)}{2} \right) = \frac{x(1 - \tanh^2\left(\frac{b+wx}{2}\right))}{4}$

$$= \frac{1}{2} * 2 * u \left(0 - \frac{x(1 - \tanh^2\left(\frac{b+wx}{2}\right))}{4} \right)$$

Substitute $u = \left(y - \frac{1 + \tanh\left(\frac{wx+b}{2}\right)}{2} \right)$

$$= \frac{1}{2} * 2 * \left(y - \frac{1 + \tanh\left(\frac{wx+b}{2}\right)}{2} \right) * \left(-\frac{x(1 - \tanh^2\left(\frac{b+wx}{2}\right))}{4} \right)$$

by simplifying

$$= -\frac{x(1 - \tanh^2\left(\frac{b+wx}{2}\right))(2y - \tanh\left(\frac{b+wx}{2}\right) - 1)}{8}$$

to calculate $\frac{\partial L}{\partial b}$ applying similar approach as above

$$\frac{\partial L}{\partial b} = -\frac{1}{4} (1 - \tanh^2\left(\frac{wx+b}{2}\right)) * \left(y - \frac{1 + \tanh\left(\frac{wx+b}{2}\right)}{2} \right)$$

(b) Consider the evaluation of E as given below,

$$E = g(x, y, z) = \sigma(c(ax + by) + dz)$$

Here x, y, z are inputs (constants) and a, b, c, d are parameters (variables). σ is the logistic sigmoid function defined as:

$$\sigma(x) = \frac{1}{1 + e^{-x}}$$

Note that here E is a function of a, b, c, d .

Compute the partial derivatives of E with respect to the parameters a, b and d i.e. $\frac{\partial E}{\partial a}$, $\frac{\partial E}{\partial b}$ and $\frac{\partial E}{\partial d}$.

Solution:

$$\begin{aligned}
 \frac{d}{dx}\sigma(x) &= \frac{d}{dx} \left[\frac{1}{1+e^{-x}} \right] \\
 &= \frac{d}{dx} (1+e^{-x})^{-1} \\
 &= -(1+e^{-x})^{-2}(-e^{-x}) \\
 &= \frac{e^{-x}}{(1+e^{-x})^2} \\
 &= \frac{1}{1+e^{-x}} \cdot \frac{e^{-x}}{1+e^{-x}} \\
 &= \frac{1}{1+e^{-x}} \cdot \frac{(1+e^{-x})-1}{1+e^{-x}} \\
 &= \frac{1}{1+e^{-x}} \cdot \left(\frac{1+e^{-x}}{1+e^{-x}} - \frac{1}{1+e^{-x}} \right) \\
 &= \frac{1}{1+e^{-x}} \cdot \left(1 - \frac{1}{1+e^{-x}} \right) \\
 &= \sigma(x) \cdot (1 - \sigma(x)) \tag{1}
 \end{aligned}$$

to calculate $\frac{\partial E}{\partial a}$
 $= \frac{\partial(\sigma(c(ax+by)+dz))}{\partial a}$

$= \frac{\partial\sigma(p)}{\partial a}$ where $p = (c(ax+by)+dz)$

Apply chain rule

$= \frac{\partial\sigma(p)}{\partial a} = \frac{\partial\sigma(p)}{\partial p} \cdot \frac{\partial p}{\partial a}$

using eq (1)

$= \sigma(p) * (1 - \sigma(p)) * \frac{\partial p}{\partial a}$

$= \sigma(p) * (1 - \sigma(p)) * c * x$

putting p value in the equation and simplifying

$= \sigma((c(ax+by)+dz)) * (1 - \sigma((c(ax+by)+dz))) * c * x$

similarly to calculate $\frac{\partial E}{\partial b}$

$= \sigma((c(ax+by)+dz)) * (1 - \sigma((c(ax+by)+dz))) * c * y$

similarly to calculate $\frac{\partial E}{\partial d}$

$= \sigma((c(ax+by)+dz)) * (1 - \sigma((c(ax+by)+dz))) * z$

2. Erroneous Estimates

The first order derivative of a real valued function f is defined by the following limit (if it exists),

$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \tag{2}$$

On observing the above definition we see that the derivative of a function is the ratio of change in the function value to the change in the function input, when we change the input by a small quantity (infinitesimally small).

Consider the function $f(x) = x^2 - 2x + 1$.

- (a) Using the limit definition of derivative, show that the derivative of $f(x)$ is $\frac{df(x)}{dx} = 2x - 2$.

Solution:

$$\begin{aligned}\frac{df(x)}{dx} &= \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h} \\ &= \lim_{h \rightarrow 0} \frac{(x+h)^2 - 2(x+h) + 1 - x^2 + 2x - 1}{h} \\ &= \lim_{h \rightarrow 0} \frac{x^2 + 2xh + h^2 - 2(x+h) + 1 - x^2 + 2x - 1}{h} \\ &= \lim_{h \rightarrow 0} \frac{h(2x - 2 + h)}{h} \\ &= \lim_{h \rightarrow 0} 2 * x - 2 + h \\ &= 2 * x - 2\end{aligned}$$

- (b) The function evaluates to 0 at 1 i.e. $f(1) = 0$.

Say we wanted to estimate the value of $f(1.01)$ and $f(1.5)$ without using the definition of $f(x)$. We could think of using the definition of derivative to “extrapolate” the value of $f(1)$ to obtain $f(1.01)$ and $f(1.5)$.

A first degree approximation based on 2 would be the following.

$$f(x+h) \approx f(x) + h \frac{df(x)}{dx} \quad (3)$$

Estimate $f(1.01)$ and $f(1.5)$ using the above formula.

Solution: using above formula to calculate

$$\begin{aligned}f(1 + .01) &\approx f(1) + 0.01 * \left. \frac{dy}{dx} \right|_{=1} \\ &= 0 + 0.01 * 0 \\ &= 0\end{aligned}$$

using above formula to calculate

$$\begin{aligned} f(1 + .50) &\approx f(1) + 0.50 * \frac{dy}{dx}|_{=1} \\ &= 0 + 0.05 * 0 \\ &= 0 \end{aligned}$$

- (c) Compare it to the actual value of $f(1.01) = 0.0001$, and $f(1.5) = 0.25$.

Solution: original $f(1.01) = 0.0001$
 calculated by above formula
 $f(1.01) = 0$
 and $f(1.5) = 0.25$
 calculated by above formula
 $f(1.50) = 0$

- (d) Explain the discrepancy from the actual value. Why does it increase/decrease when we move further away from 1?

Solution: Discrepancy from the actual value as h by definition is very small
 $\lim_{h \rightarrow 0} h \rightarrow 0$ which is not for $h = .01$ or $h = .5$ so we need Taylor Series expansion

As we further away from 1 the value of the function either increase or decrease,
 and as $\frac{d^2f}{dx^2}|_{=1} = +2$ so the point $x = 1$ is local minima

- (e) Can we get a better estimate of $f(1.01)$ and $f(1.5)$ by “correcting” our estimate from part (a)? Can you suggest a way of doing this?

Solution: according to Taylor Series expansion
 $f(x) \approx f(x_0) + f'(x_0)(x - x_0) + \frac{f''(x_0)}{2!}(x - x_0)^2 + \frac{f'''(x_0)}{3!}(x - x_0)^3 + \dots$
 if we put $x + h$ in place of x we get
 $f(x+h) = f(x) + f'(x)(x - (x+h)) + \frac{f''(x)}{2!}(x - (x+h))^2 + \frac{f'''(x)}{3!}(x - (x+h))^3 + \dots$
 $= f(x) + hf'(x) + \frac{h^2}{2!}f''(x) + \frac{h^3}{3!}f'''(x) + \dots$
 now putting $x = 1$ and $h = .01$ we get
 $f(1.01) = f(1) + .01 * f'(1) + \frac{(.01)^2}{2!}f''(1) + \frac{(.01)^3}{3!}f'''(1) + \dots$
 Simplifying
 $= 0 + .01 * 0 + \frac{(.01)^2}{2!} * 2 + \frac{(.01)^3}{3!} * 0 = .0001$
 now putting $x = 1$ and $h = .5$ we get

$$\begin{aligned}
f(1.5) &= f(1) + .5 * f'(1) + \frac{(.5)^2}{2!} f''(1) + \frac{(.5)^3}{3!} f'''(1) + \dots \\
\text{Simplifying} \\
&= 0 + .5 * 0 + \frac{(.5)^2}{2!} * 2 + \frac{(.5)^3}{3!} * 0 = .25
\end{aligned}$$

3. Differentiation w.r.t. Vectors and matrices

Consider vectors $\mathbf{u}, \mathbf{x} \in \mathbb{R}^d$, and matrix $\mathbf{A} \in \mathbb{R}^{n \times n}$.

The gradient of a scalar function f w.r.t. a vector \mathbf{x} is a vector by itself, given by

$$\nabla_{\mathbf{x}} f = \left(\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n} \right)$$

Gradient of a scalar function w.r.t a matrix is a matrix.

$$\nabla_{\mathbf{A}} f = \begin{bmatrix} \frac{\partial f}{\partial A_{11}} & \frac{\partial f}{\partial A_{12}} & \dots & \frac{\partial f}{\partial A_{1n}} \\ \frac{\partial f}{\partial A_{21}} & \frac{\partial f}{\partial A_{22}} & \dots & \frac{\partial f}{\partial A_{2n}} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f}{\partial A_{n1}} & \frac{\partial f}{\partial A_{n2}} & \dots & \frac{\partial f}{\partial A_{nn}} \end{bmatrix}$$

Gradient of the gradient of a function w.r.t. a vector is a matrix. It is referred to as Hessian.

$$\mathbf{H}_{\mathbf{x}} f = \nabla_{\mathbf{x}}^2 f = \begin{bmatrix} \frac{\partial^2 f}{\partial x_1^2} & \frac{\partial^2 f}{\partial x_1 \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f}{\partial x_2 \partial x_1} & \frac{\partial^2 f}{\partial x_2^2} & \dots & \frac{\partial^2 f}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f}{\partial x_n \partial x_1} & \frac{\partial^2 f}{\partial x_n \partial x_2} & \dots & \frac{\partial^2 f}{\partial x_n^2} \end{bmatrix}$$

(a) Derive the expressions for the following gradients.

1. $\nabla_{\mathbf{x}} \mathbf{u}^T \mathbf{x}$
2. $\nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{x}$
3. $\nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x}$
4. $\nabla_{\mathbf{A}} \mathbf{x}^T \mathbf{A} \mathbf{x}$
5. $\nabla_{\mathbf{x}}^2 \mathbf{x}^T \mathbf{A} \mathbf{x}$

(Aside: Compare your results with derivatives for the scalar equivalents of the above expressions ax and x^2 .)

The gradient of a scalar f w.r.t. a matrix \mathbf{X} , is a matrix whose (i, j) component is $\frac{\partial f}{\partial X_{ij}}$, where X_{ij} is the (i, j) component of the matrix \mathbf{X} .)

Solution: 1.

$$\begin{aligned}\nabla_{\mathbf{x}} \mathbf{u}^T \mathbf{x} &= \frac{\partial}{\partial x} (u^T x) \\ \text{as } \frac{\partial \sum_{i=1}^n u_i x_i}{\partial x_i} &= u_i \\ &= (u_1, u_2, u_3, u_4 \dots, u_n) \\ &= u^T\end{aligned}$$

2.

$$\begin{aligned}\nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{x} &= \frac{\partial}{\partial x} (x^T x) \\ \text{as } \frac{\partial \sum_{i=1}^n x_i \cdot x_i}{\partial x_i} &= \frac{\partial \sum_{i=1}^n x_i^2}{\partial x_i} = 2 * x_i \\ &= 2 * (x_1, x_2, x_3, x_4 \dots, x_n) \\ &= 2 * x^T\end{aligned}$$

3.

$$\begin{aligned}\nabla_{\mathbf{x}} \mathbf{x}^T \mathbf{A} \mathbf{x} &= \frac{\partial}{\partial x} (x^T A x) \\ &= \frac{\partial \bar{x}^T A x}{\partial x} + \frac{\partial x^T A \bar{x}}{\partial x} \\ &\text{where } \bar{x} \text{ and } \bar{x}^T \text{ as constant so applying chain rule} \\ &\text{to compute the above derivative as } v_1 = A \bar{x} \quad v_2^T = \bar{x}^T A \\ &= \frac{\partial x^T v_1}{\partial x} + \frac{\partial v_2^T x}{\partial x} \\ \text{applying eq(1) result} &= v_1^T + v_2^T \\ &= x^T A^T + x^T A \\ &= x^T (A + A^T)\end{aligned}$$

Solution: 4.

$$\begin{aligned}
\nabla_{\mathbf{A}} \mathbf{x}^T \mathbf{A} \mathbf{x} &= \frac{\partial}{\partial a_{ij}} \left(\begin{bmatrix} x_1 & x_2 & \dots & x_n \end{bmatrix} \begin{bmatrix} a_{11} & a_{12} & a_{13} & \dots & a_{1n-2} & a_{1n-1} & a_{1n} \\ \vdots & \vdots & \vdots & \vdots & \ddots & \vdots & \vdots \\ a_{n1} & a_{n2} & a_{n3} & \dots & a_{nn-2} & a_{nn-1} & a_{nn} \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \right) \\
&= \frac{\partial}{\partial a_{ij}} \left(\begin{bmatrix} a_{11}x_1 + a_{21}x_2 + \dots + a_{n-11}x_{n-1} + a_{n1}x_n \\ a_{12}x_1 + a_{22}x_2 + \dots + a_{n-12}x_{n-1} + a_{n2}x_n \\ \vdots \\ a_{1n}x_1 + a_{2n}x_2 + \dots + a_{n-1n}x_{n-1} + a_{nn}x_n \end{bmatrix}^T \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_n \end{bmatrix} \right) \\
&= \begin{bmatrix} x_1^2 & x_1x_2 & \dots & x_1x_{n-1} & x_1x_n \\ x_1x_2 & x_2^2 & \dots & x_2x_{n-1} & x_2x_n \\ \vdots & \vdots & \dots & \vdots & \vdots \\ x_1x_n & x_2x_n & \dots & x_{n-1}x_n & x_n^2 \end{bmatrix}
\end{aligned}$$

5.

$$\nabla_{\mathbf{x}}^2 \mathbf{x}^T \mathbf{A} \mathbf{x} = \frac{\partial^2}{\partial x^2} x^T \mathbf{A} x$$

using result from eq(3)

$$\begin{aligned}
\nabla_{\mathbf{x}}^2 \mathbf{x}^T \mathbf{A} \mathbf{x} &= \frac{\partial}{\partial x} x^T (\mathbf{A} + \mathbf{A}^T) \\
\frac{\partial^2}{\partial x^2} x^T \mathbf{A} x &= \mathbf{A} + \mathbf{A}^T
\end{aligned}$$

- (b) Use the equations obtained in the previous part to get the Linear regression solution that you studied in ML or PR. Suppose X as input example-feature matrix, Y as given outputs and \mathbf{w} as weight vector.

Solution: $Y = XW + \epsilon$

The vector of residuals ϵ is given by:

$$\epsilon = Y - X\hat{w}$$

mean square error

$$\epsilon^T \epsilon = (Y - X\hat{w})^T (Y - X\hat{w})$$

$$= Y^T Y - \hat{w} X^T Y - Y^T X \hat{w} + \hat{w} X^T X \hat{w}$$

$$= Y^T Y - 2\hat{w} X^T Y + \hat{w} X^T X \hat{w}$$

as the transpose of a scalar is scalar $(\hat{w} X^T Y)^T = Y^T X \hat{w}$

$$\frac{\partial \epsilon^T \epsilon}{\partial \hat{w}} = -2X^T Y + 2X^T X \hat{w} = 0 \text{ using eq(1) and eq(2)}$$

$$\begin{aligned} \text{so } X^T X w &= X^T Y \\ (X^T X)^{-1} X^T X w &= (X^T X)^{-1} X^T Y \\ w &= (X^T X)^{-1} X^T Y \end{aligned}$$

- (c) By now you must have the intuition. Gradient w.r.t. a 1 dimensional array was 1 dimensional. Gradient w.r.t. a 2-dimensional array was 2 dimensional. Higher order arrays are referred to as tensors. Let \mathbf{T} be a 3 dimensional tensor. Write the expression of $\nabla_{\mathbf{T}} f$. You can use gradients w.r.t. a vector or a matrix in the expression.

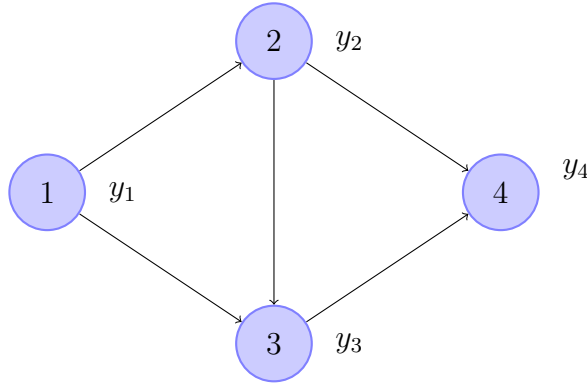
Solution: $\nabla_{\mathbf{T}} f = [\nabla_{\mathbf{T}} A_1 \quad \nabla_{\mathbf{T}} A_2 \quad \cdots \quad \nabla_{\mathbf{T}} A_n]$
 where $A_1, A_2, \cdots A_n$ are 2 dimensional array

4. Ordered Derivatives

An ordered network is a network where the state variables can be computed one at a time in a specified order.

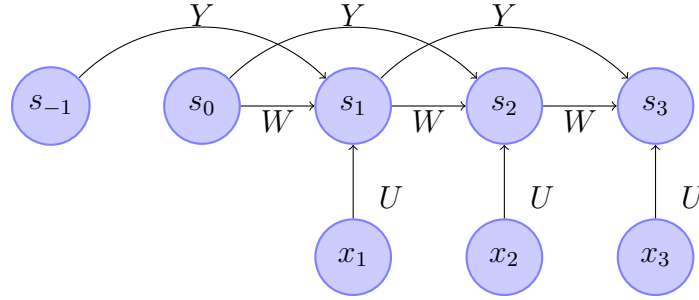
Answer the following questions regarding such a network.

- (a) Given the ordered network below, give a formula for calculating the ordered derivative $\frac{\partial y_4}{\partial y_1}$ in terms of partial derivatives w.r.t. y_1 and y_2 where y_1, y_2 and y_3 are the outputs of nodes 1, 2 and 3 respectively.



Solution:
$$\begin{aligned} \frac{\partial y_4(y_3, y_2)}{\partial y_1} &= \frac{\partial y_4(y_3, y_2)}{\partial y_3} * \frac{\partial y_3(y_2, y_1)}{\partial y_1} + \frac{\partial y_4(y_3, y_2)}{\partial y_2} * \frac{\partial y_2(y_1)}{\partial y_1} \\ &= \frac{\partial y_4(y_3, y_2)}{\partial y_3} * \frac{\partial y_3(y_2, y_1)}{\partial y_2} * \frac{\partial y_2(y_1)}{\partial y_1} + \frac{\partial y_4(y_3, y_2)}{\partial y_3} * \frac{\partial y_3(y_2, y_1)}{\partial y_1} + \frac{\partial y_4(y_3, y_2)}{\partial y_2} * \frac{\partial y_2(y_1)}{\partial y_1} \end{aligned}$$

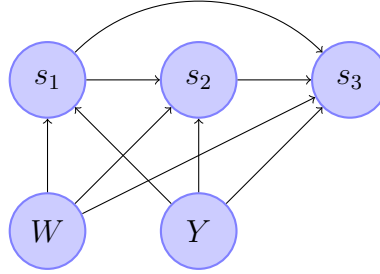
- (b) The figure above can be viewed as a dependency graph as it tells us which variables in the system depend on which other variables. For example, we see that y_3 depends on y_1 and y_2 which in turn also depends on y_1 . Now consider the network given below,



Here, $s_i = \sigma(Ws_{i-1} + Ys_{i-2} + Ux_i + b) \quad (\forall i \geq 1)$.

Can you draw a dependency graph involving the variables s_3, s_2, s_1, W, Y ?

Solution:



(c) Give a formula for computing $\frac{\partial s_3}{\partial W}$, $\frac{\partial s_3}{\partial Y}$ and $\frac{\partial s_3}{\partial U}$ for the network shown in part (b)

Solution: applying chain rule for multiplication we get

$$\frac{\partial s_3}{\partial W} = (s_3(1 - s_3))(s_2 + s_2(1 - s_2)(s_1 + s_0 + Ws_{-1}) + Y(s_0 + Ws_{-1}))$$

$$\frac{\partial s_3}{\partial Y} = (s_3(1 - s_3))(Ws_2(1 - s_2)(Ws_1(1 - s_1)s_{-1} + s_0) + Y(s_1(1 - s_1)s_{-1} + s_0)) + s_1$$

$$\frac{\partial s_3}{\partial U} = s_3(1 - s_3)x_3$$

5. Baby Steps

From basic calculus, we know that we can find the minima (local and global) of a function by finding the first and second order derivatives. We set the first derivative to zero and verify if the second derivative at the same point is positive. The reasoning behind the following procedure is based on the interpretation of the derivative of a function as the slope of the function at any given point.

The above procedure, even though correct can be intractable in practice while trying to minimize functions. And this is not just a problem for the multivariable case, but even

for single variable functions. Consider minimizing the function $f(x) = x^5 + 5\sin(x) + 10\tan(x)$. Although the function f is a contrived example, the point is that the standard derivative approach, might not always be a feasible way to find minima of functions.

In this course, we will be routinely dealing with minimizing functions of multiple variables (in fact millions of variables). Of course we will not be solving them by hand, but we need a more efficient way of minimizing functions. For the sake of this problem, consider we are trying to minimize a convex function of one variable $f(x)$,¹ which is guaranteed to have a single minima. We will now build an iterative approach to finding the minima of functions.

The high level idea is the following:

Start at a (random) point x_0 . Verify if we are at the minima. If not, change the value so that we are moving closer to the minima. Keep repeating until we hit the minima.

- (a) Use the intuition built from Q.3 to find a way to change the current value of x while still ensuring that we are improving (i.e. minimizing) the function.

Solution:

repeat until convergence: $\{ \quad x_j := x_j - \alpha \frac{d}{dx}(f(x)) \text{ for } j := 0 \dots n \}$
where α is learning rate

- (b) How would you use the same idea, if you had to minimize a function of multiple variables ?

Solution:

repeat until convergence: $\{ \quad x_j := x_j - \alpha \nabla_{x_j}(f(x_0, x_1, x_2 \dots x_n)) \text{ for } j := 0 \dots n \}$
where α is learning rate

- (c) Does your method always lead to the global minima (smallest value) for non convex functions (which may have multiple local minima)? If yes, can you explain (prove or argue) why? If not, can you give a concrete example of a case where it fails?

Solution: No My method does not guarantee to find global minima. It can stuck in local minima as we choose α the rate at which the function converges to local minima. As it is large it can overshoot and jump out of the local minima and fall in other part of the graph or if choose very low value of α it may take very large time to converge. Or it can stuck in steep local minima rather than gradual global minima

¹https://en.wikipedia.org/wiki/Convex_function

- (d) Do you think this procedure always works for convex functions ? (*i.e.*, are we always guaranteed to reach the minima)

Solution: Yes, I think this procedure works for convex function as there is only global minima is present so no other option to stuck in local minima.

- (e) (Extra) Can you think of the number of steps needed to reach the minima ?

Solution: It depends on the learning rate α as it is optimally chosen the faster the convergence and less step to reach minima.

- (f) (Extra) Can you think of ways to improve the number of steps needed to reach the minima ?

Solution: There are other method as Nesterov momentum ,Optimal gradient method etc used for faster convergence so the steps need less.

6. **Constrained Optimization** Let $f(x, y)$ and $g(x, y)$ be smooth (continuous, differentiable etc.) real valued functions of two variables. We want to minimize f , which is a convex function of x and y .

- (a) Argue that at the minima of f , the partial derivative $\frac{\partial f}{\partial x}$ and $\frac{\partial f}{\partial y}$ will be zero. Thus setting the partial derivatives to zero is a possible method for finding the minima.

Solution: To find a stationary point we need to set the partial derivative $\frac{\partial f}{\partial x}$ and $\frac{\partial f}{\partial y}$ to zero , this gives two equation for two unknowns x and y . after finding the points there are 3 possibilities

1. the point is **saddle point**
2. the point is either a **maximum point**
3. the point is either a **minimum point**

so to find minimum we need to further check value of $f''_{xx}f''_{yy} - f_{xy}^2$ either > 0 or < 0

Intuitively, in terms of graphs, minima of multivariable functions are down bottom, just as they are with single variable functions. The gradient of a multivariable function at a minimum point will be the zero vector, which corresponds to the graph having a flat tangent plane.

- (b) Suppose we are only interested in minimizing f in the region where $g(x, y) = c$, where c is some constant. Suppose this region is a curve in the x - y plane. Call this the feasible curve. Will our previous technique still work in this case? Why or why not?

Solution: our previous technique will work if we can represent $g(x, y) = c$ in terms of $x = t(y)$ or $y = s(z)$ and replace x or y respectively in f and apply previous method.

And not possible if we fail to do so then we need Lagrange method for constraint optimization.

- (c) What is the component of ∇g along the feasible curve, when computed at points lying on the curve?

Solution: The component of ∇g along the feasible curve is perpendicular.

- (d) * At the point on the feasible curve, which achieves minimum value of f , what will be the component of ∇f along the curve?

Solution: The component of ∇f along the feasible curve is perpendicular.

- (e) Using the previous answers, show that at the point on the feasible curve, achieving minimum value of f , $\nabla f = \lambda \nabla g$ for some real number λ . Thus, this equation, combined with the constraint $\nabla g = 0$ should enable us to find the minima.

Solution: gradient of a function is perpendicular to the contour lines, the contour lines of f and g are parallel if and only if the gradients of functions parallel ($\nabla f = \lambda \nabla g$). Thus we want points (x,y) where $g(x, y) = 0$

- (f) * Using the insights from discussion so far, solve the the following optimization problem:

$$\max_{x,y,z} x^a y^b z^c$$

where

$$x + y + z = 1$$

and given $a, b, c > 0$.

Solution: We define the Lagrangian as $\mathcal{L}(x, y, \lambda) = (x^a y^b z^c - \lambda * (x + y + z - 1))$ gradient of a function is perpendicular to the contour lines, the contour lines of $x^a y^b z^c$ and $(x + y + z - 1)$ are parallel if and only if the gradients of functions parallel. Thus we want points (x,y,z)

$$\nabla_{x,y,z,\lambda}(\mathcal{L}(x, y, \lambda)) = 0$$

$$\text{so } \nabla_{x,y,z}(x^a y^b z^c) = \lambda \nabla_{x,y,z}(x + y + z - 1) \quad \text{and}$$

$$\nabla_{\lambda}(\mathcal{L}(x, y, \lambda)) = 0 \quad \text{means} \quad (x + y + z - 1) = 0$$

so we are getting equation for $\nabla_{x,y,z}(\mathcal{L}(x, y, \lambda)) = 0$ are

$$a * x^{a-1} y^b z^c = \lambda * 1$$

$$b * x^a y^{b-1} z^c = \lambda * 1$$

$$c * x^a y^b z^{c-1} = \lambda * 1$$

and solving with bellow equation

$$(x + y + z - 1) = 0$$

we are getting

$$x = \frac{a}{a + b + c}$$

$$y = \frac{b}{a + b + c}$$

$$z = \frac{c}{a + b + c}$$

$$\max_{x,y,z} x^a y^b z^c = \frac{a^a b^b c^c}{(a + b + c)^{(a+b+c)}}$$

7. Billions of Balloons

Consider a large playground filled with 1 billion balloons. Of these there are k_1 blue, k_2 green and k_3 red balloons. The values of k_1 , k_2 and k_3 are not known to you but you are interested in estimating them. Of course, you cannot go over all the 1 billion balloons and count the number of blue, green and red balloons. So you decide to randomly sample 1000 balloons and note down the number of blue, green and red balloons. Let these counts be \hat{k}_1 , \hat{k}_2 and \hat{k}_3 respectively. You then estimate the total number of blue, green and red balloons as $1000000 * \hat{k}_1$, $1000000 * \hat{k}_2$ and $1000000 * \hat{k}_3$.

- (a) Your friend knows the values of k_1 , k_2 and k_3 and wants to see how bad your estimates are compared to the true values. Can you suggest some ways of calculating this difference? [Hint: Think about probability!]

Solution:

- (b) * Consider two ways of converting \hat{k}_1 , \hat{k}_2 and \hat{k}_3 to a probability distribution:

$$p_i = \frac{\hat{k}_i}{\sum_i \hat{k}_i}$$

$$q_i = \frac{e^{\hat{k}_i}}{\sum_i e^{\hat{k}_i}}$$

Would you prefer the distribution $\mathbf{q} = [q_1, q_2, \dots, q_n]$ over $\mathbf{p} = [p_1, p_2, \dots, p_n]$ for the above task? Give reasons and provide an example to support your choice.

Solution: q_i in terms of differentiability it is easy to use for optimization, its good fit for 1-N classification. the function q_i is called softmax classifier which can represent any N-class probability function over the feature space as Maximum Likelihood probability has the Universal Approximation Property

8. ** Let X be a real-valued random variable with p as its probability density function (PDF). We define the cumulative density function (CDF) of X as

$$F(x) = \Pr(X \leq x) = \int_{y=-\infty}^{y=x} p(y)dy$$

What is the value of $\mathbb{E}_X[F(X)]$ (the expected value of the CDF of X)? **The answer is a real number** (Hint: The expectation can be formulated as a double integral. Try to plot the area over which you need to integrate in the x-y plane. Now look at the area over which you are not integrating. Do you notice any symmetries?)

Solution: by definition of expectation for any function

$$\begin{aligned}\mathbb{E}_X[F(X)] &= \int_{-\infty}^{\infty} F(x)p(x)dx \\ \text{by definition } p(x) &= \frac{dF_X(x)}{dx} \text{ replacing } p(x)dx \text{ with } dF_X(x) \\ &= \int_{-\infty}^{\infty} F(x) dF_X(x) \\ \text{as } F_X(x) &\text{ is CDF so maximum value 1 and minimum value 0} \\ &= \int_0^1 F(x) dF_X(x) \\ &= \left[\frac{1}{2} F_X(x)^2 \right]_{F_X(x)=0}^{F_X(x)=1} \\ &= \frac{1}{2}\end{aligned}$$

9. * **Intuitive Urns**

An urn initially contains 3 red balls and 3 blue balls. One of the balls is removed without being observed. To find out the color of the removed ball, Alice and Bob independently perform the same experiment: they randomly draw a ball, record the color, and put it back. This is repeated several times and the number of red and blue balls observed by each of them is recorded.

Alice draws 6 times and observes 6 red balls and 0 blue balls.

Bob draws 600 times and observes 303 red balls and 297 blue balls.

Obviously, both of them will predict that the removed ball was blue.

- (a) Intuitively, who do you think has stronger evidence for claiming that the removed ball was blue, and why? (**Don't cheat by computing the answer. This subquestion has no marks, but is compulsory!**)

Solution: As of intuitively I think **bob** has stronger evidence as he done experiment more so his error probability is less compared to alice

- (b) What is the exact probability that the removed ball was blue, given Alice's observations? (Hint: Think Bayesian Probability)

Solution:

$$\begin{aligned}
 P\left(\frac{Blue}{alice}\right) &= \frac{P\left(\frac{alice}{Blue}\right) * P(Blue)}{P\left(\frac{alice}{Blue}\right) * P(Blue) + P\left(\frac{alice}{Red}\right) * P(Red)} \\
 &= \frac{\frac{3^6}{5} * \frac{1}{2}}{\frac{3^6}{5} * \frac{1}{2} + \frac{2^6}{5} * \frac{1}{2}} \\
 &= \frac{3^6}{3^6 + 2^6}
 \end{aligned}$$

- (c) What is the exact probability that the removed ball was blue, given Bob's observations? (Hint: Think Bayesian Probability)

Solution:

$$\begin{aligned}
 P\left(\frac{Blue}{bob}\right) &= \frac{P\left(\frac{bob}{Blue}\right) * P(Blue)}{P\left(\frac{bob}{Blue}\right) * P(Blue) + P\left(\frac{bob}{Red}\right) * P(Red)} \\
 &= \frac{\binom{600}{297} \frac{3^{303}}{5} * \frac{2^{297}}{5} * \frac{1}{2}}{\binom{600}{297} \frac{3^{303}}{5} * \frac{2^{297}}{5} * \frac{1}{2} + \binom{600}{303} \frac{2^{303}}{5} * \frac{3^{297}}{5} * \frac{1}{2}} \\
 &= \frac{3^6}{3^6 + 2^6}
 \end{aligned}$$

- (d) Computationally, who do you think has stronger evidence for claiming that the removed ball was blue?

Solution: Both have stronger evidence to claim that the removed ball was blue by calculation of probability.

Did your intuition match up with the computations? If yes, awesome! If not, remember that probability can often be seen deceptively straightforward. Try to avoid intuition when dealing with probability by grounding it in formalism.

10. Plotting Functions for Great Good

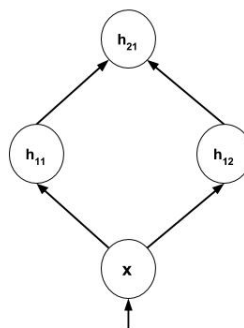
- (a) Consider the variable x and functions $h_{11}(x)$, $h_{12}(x)$ and $h_{21}(x)$ such that

$$h_{11}(x) = \frac{1}{1 + e^{-(400x+24)}}$$

$$h_{12}(x) = \frac{1}{1 + e^{-(400x-24)}}$$

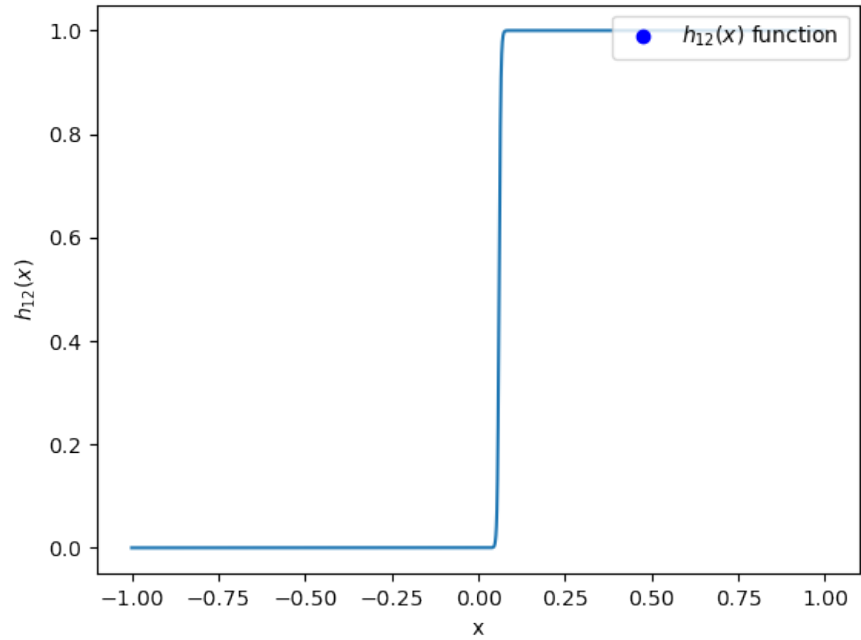
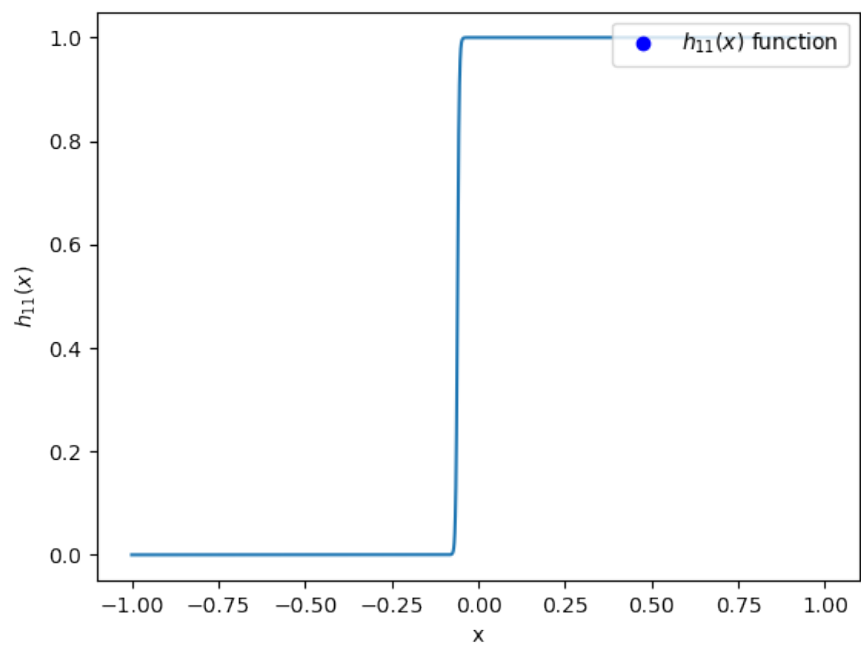
$$h_{21} = h_{11}(x) - h_{12}(x)$$

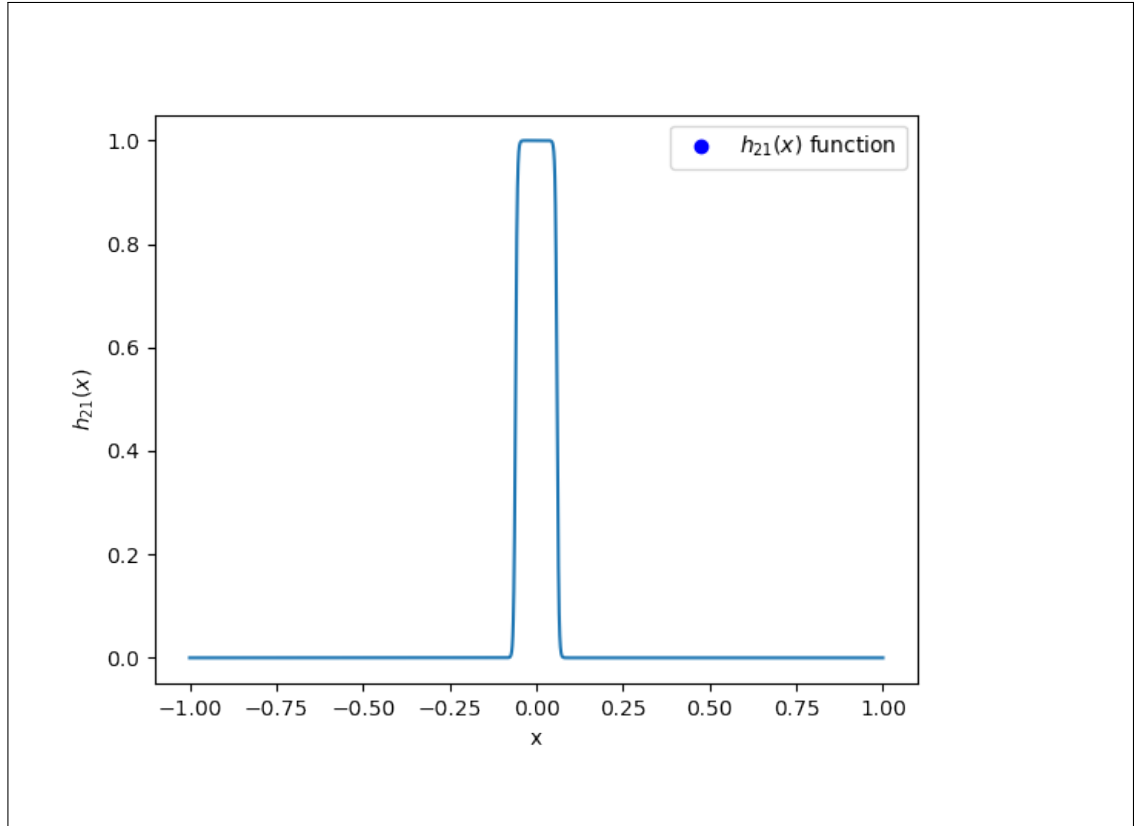
The above set of functions are summarized in the graph below.



Plot the following functions: $h_{11}(x)$, $h_{12}(x)$ and $h_{21}(x)$ for $x \in (-1, 1)$

Solution:





- (b) Now consider the variables x_1, x_2 and the functions $h_{11}(x_1, x_2), h_{12}(x_1, x_2), h_{13}(x_1, x_2), h_{14}(x_1, x_2), h_{21}(x_1, x_2), h_{22}(x_1, x_2), h_{31}(x_1, x_2)$ and $f(x_1, x_2)$ such that

$$h_{11}(x_1, x_2) = \frac{1}{1 + e^{-(x_1 + 100x_2 + 200)}}$$

$$h_{12}(x_1, x_2) = \frac{1}{1 + e^{-(x_1 + 100x_2 - 200)}}$$

$$h_{13}(x_1, x_2) = \frac{1}{1 + e^{-(100x_1 + x_2 + 200)}}$$

$$h_{14}(x_1, x_2) = \frac{1}{1 + e^{-(100x_1 + x_2 - 200)}}$$

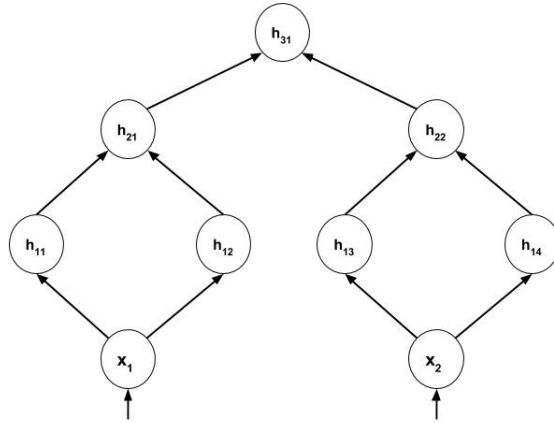
$$h_{21}(x_1, x_2) = h_{11}(x_1, x_2) - h_{12}(x_1, x_2)$$

$$h_{22}(x_1, x_2) = h_{13}(x_1, x_2) - h_{14}(x_1, x_2)$$

$$h_{31}(x_1, x_2) = h_{21}(x_1, x_2) + h_{22}(x_1, x_2)$$

$$f(x_1, x_2) = \frac{1}{1 + e^{-(50h_{31}(x) - 100)}}$$

The above set of functions are summarized in the graph below.



Plot the following functions: $h_{11}(x_1, x_2)$, $h_{12}(x_1, x_2)$, $h_{13}(x_1, x_2)$, $h_{14}(x_1, x_2)$, $h_{21}(x_1, x_2)$, $h_{22}(x_1, x_2)$, $h_{31}(x_1, x_2)$ and $f(x_1, x_2)$ for $x_1 \in (-5, 5)$ and $x_2 \in (-5, 5)$

Solution:

