# 446 Section 3.00000001

TA: Varun Ananth

Plans for today!

1. This
2. Reminders
3. Review Vector Calculus/Approximations
    a. Will be lecture/follow along style

# Jacobians and Hessians

And how it is just a fancy way to describe gradients

Slides by Marco D.

## 1.1. Definitions

Let $f : \mathbf{R}^n \to \mathbf{R}$ and let $g : \mathbf{R}^n \to \mathbf{R}^m$. The **gradient** of $f$ (with respect to $x$) evaluated at $x$ is the vector of partial derivatives:

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \vdots \\ \frac{\partial f(x)}{\partial x_n} \end{bmatrix} \in \mathbf{R}^n$$

The **Jacobian** of $g$ (with respect to $x$) evaluated at $x$ is the matrix of partial derivatives:

$$\nabla_x g(x) = \begin{bmatrix} \frac{\partial g_1(x)}{\partial x_1} & \cdots & \frac{\partial g_1(x)}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_m(x)}{\partial x_1} & \cdots & \frac{\partial g_m(x)}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \nabla_x^T g_1(x) \\ \vdots \\ \nabla_x^T g_m(x) \end{bmatrix} \in \mathbf{R}^{m \times n}$$

WHAT IS THIS?!?!?!

Sometimes the Jacobian is denoted by $J_g(x)$, but we use $\nabla_x g(x)$ to highlight that the Jacobian is nothing more than the generalization of the gradient to functions which have a vector output.

The **Hessian** of f (with respect to x) evaluated at $x$ is the matrix of partial derivatives:

$$\nabla_x^2 f(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_n \partial x_1} & \frac{\partial^2 f(x)}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_n^2} \end{bmatrix} \in \mathbf{R}^{n \times n}$$

Sometimes the Hessian is denoted by $H_f(x)$, but we use $\nabla_x^2 f(x)$ to highlight that the Hessian is the Jacobian of the gradient of f.

Slides by Marco D.

# Gradient

# How do we calculate the gradient of a function with a vector input?

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \dots \\ x_n \end{pmatrix}$$

$$f(x) : \mathbb{R}^n \rightarrow \mathbb{R}$$

This is normal function that outputs a real number
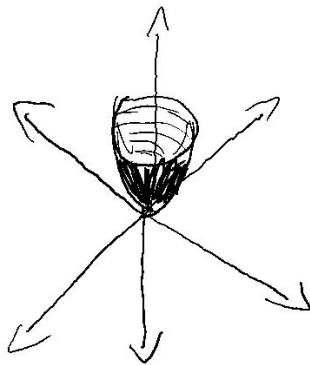
We simply do partial derivatives n times

$$\Delta_x f(x) = \begin{bmatrix} \dfrac{\partial f(x)}{\partial x_1} \\ \dfrac{\partial f(x)}{\partial x_2} \\ \vdots \\ \dfrac{\partial f(x)}{\partial x_n} \end{bmatrix}$$

$$x = \begin{pmatrix} x_1 \\ x_2 \\ \cdots \\ x_n \end{pmatrix}$$

# Visualizing the Gradient

i.e: Scalar Value

$$f: \mathbb{R}^n \longrightarrow \mathbb{R}$$
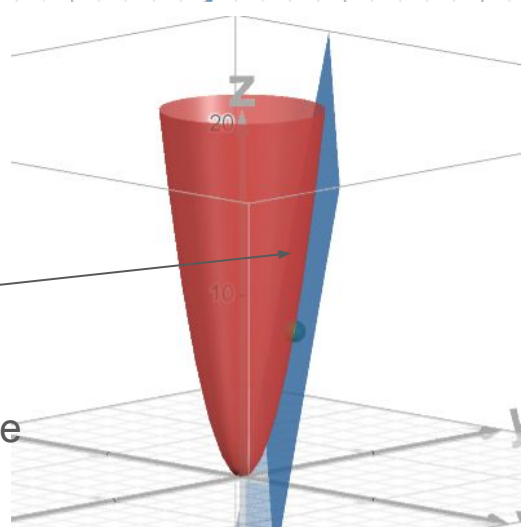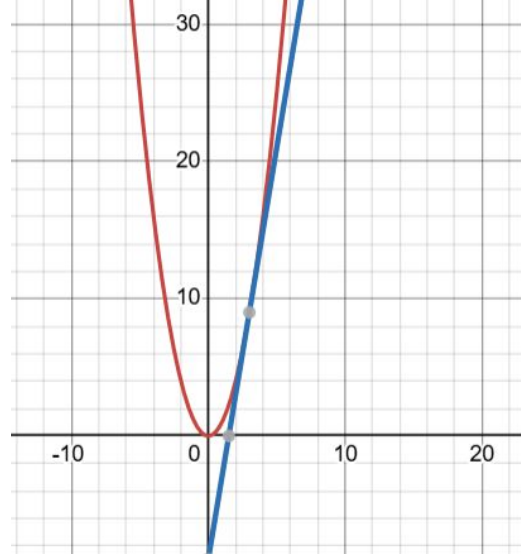
Example: $z = X_1^2 + X_2^2$

can also be thought of as $X$ and $Y$

$$\nabla_X f(X) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2X_1 \\ 2X_2 \end{bmatrix}$$

all $X$'s

Tells you the slope of the tangent plane in the $x_1$ and $x_2$ directions
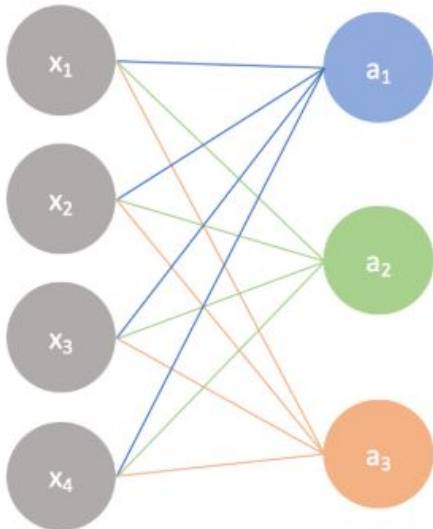
# Jacobian

In machine learning, we don't usually have the privilege of having a function that outputs a real number. Usually, the function will output a vector. For example:

Input layer          Output layer

A simple neural network



$$\begin{bmatrix} w_1 & w_2 & w_3 & w_4 \\ w_1 & w_2 & w_3 & w_4 \\ w_1 & w_2 & w_3 & w_4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \end{bmatrix} + \begin{bmatrix} b \\ b \\ b \end{bmatrix} = \begin{bmatrix} w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b \\ w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b \\ w_1x_1 + w_2x_2 + w_3x_3 + w_4x_4 + b \end{bmatrix} \xrightarrow[activation]{} \begin{bmatrix} a_1 \\ a_2 \\ a_3 \end{bmatrix}$$

What do we do now???

$$g : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

Example: $g(x) = Wx$

where: $W \in \mathbb{R}^{m \times n}$, $x \in \mathbb{R}^n$

$$m \begin{bmatrix} W \end{bmatrix} \; n \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = m \begin{bmatrix} \; \\ \; \\ \; \end{bmatrix}$$
$$\phantom{m}\underset{n}{} \qquad \underset{1}{} \qquad \underset{1}{}$$

Now what is $\nabla_x g(x)$?

dot product

$$m \begin{bmatrix} \boxed{W_1} \\ \; \\ W \\ \; \end{bmatrix} \; n \begin{bmatrix} x_1 \\ \vdots \\ x_n \end{bmatrix} = m \begin{bmatrix} W_1 \bullet x \\ W_2 \bullet x \\ \vdots \\ W_m \bullet x \end{bmatrix}$$
$$\phantom{m}\underset{n}{} \qquad \qquad \qquad \underset{1}{}$$

$$\nabla_x g(x) = m \begin{bmatrix} \boxed{\nabla_x W_1 \bullet x} \\ \nabla_x W_2 \bullet x \\ \vdots \\ \nabla_x W_m \bullet x \end{bmatrix}$$
$$\phantom{\nabla_x g(x) = m}\underset{1}{}$$

$\nabla_x (w_1^i x_1 + w_i^2 x_2 \cdots)$

Gradient!

$$\nabla_x g(x) = m \begin{bmatrix} \boxed{\nabla_x W_1 \cdot X} \\ \nabla_x W_2 \cdot X \\ \vdots \\ \nabla_x W_m \cdot X \end{bmatrix}$$
1

$\nabla_x (w_i^1 x_1 + w_i^2 x_2 \cdots)$

Gradient!

(For notational purposes, $g_i$ = computation for row *i* of the output)

Each row becomes an *n* element gradient

$$m \begin{bmatrix} \dfrac{\partial g_1(x)}{\partial x_1} & \cdots & \dfrac{\partial g_1(x)}{\partial x_n} \\ & \vdots & \\ \dfrac{\partial g_m(x)}{\partial x_1} & \cdots & \dfrac{\partial g_m(x)}{\partial x_n} \end{bmatrix}$$
n

$g_1(x) = w_1^1 x_1 + w_1^2 x_2 \cdots$

$g_2(x) = w_2^1 x_1 + w_2^2 x_2 \cdots$

$g_m(x) = w_m^1 x_1 + w_n^2 x_2 \cdots$

This is the <u>Jacobian</u> of $g(x)$

# Interpreting the Jacobian

How do we interpret the jacobian matrix?

$$\nabla_x g(x) = \begin{bmatrix} \frac{\partial g_1(x)}{\partial x_1} & \cdots & \frac{\partial g_1(x)}{\partial x_n} \\ \vdots & \ddots & \vdots \\ \frac{\partial g_m(x)}{\partial x_1} & \cdots & \frac{\partial g_m(x)}{\partial x_n} \end{bmatrix} = \begin{bmatrix} \nabla_x^T g_1(x) \\ \vdots \\ \nabla_x^T g_m(x) \end{bmatrix} \in \mathbb{R}^{m \times n}$$

This matrix gives us an idea of how the output will change if we slightly change the value of x.

For example, if we increase $x_1$, how is g(x) affected?

Slides by Marco D.

# Hessian

$$f : \mathbb{R}^n \rightarrow \mathbb{R}$$

$$\nabla_x \left( \nabla_x f(x) \right) = \nabla_x^2 f(x) = \underline{Hessian}$$

$$\underline{\nabla_x f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^n}$$

Remember the dimensionality of $\nabla_x g(x)$?

$$\downarrow$$

$$\nabla_x \left( \nabla_x f(x) \right) \in \mathbb{R}^{n \times n}$$

The Hessian is the Jacobian of the gradient of f(x)

**Important to understand!**

$$\nabla_x^2 f(x) = \begin{bmatrix} \dfrac{\partial^2 f(x)}{\partial x_1^2} & \dfrac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \dfrac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \dfrac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \dfrac{\partial^2 f(x)}{\partial x_2^2} & \cdots & \dfrac{\partial^2 f(x)}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \dfrac{\partial^2 f(x)}{\partial x_n \partial x_1} & \dfrac{\partial^2 f(x)}{\partial x_n \partial x_2} & \cdots & \dfrac{\partial^2 f(x)}{\partial x_n^2} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

Notice the combinations of variables:
- Derive by the same variable twice for the diagonal
- Derive by every combination of $x_i, x_j$ where $i \neq j$ for the off-diagonals

$$\nabla_x^2 f(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_1 \partial x_n} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2^2} & \cdots & \frac{\partial^2 f(x)}{\partial x_2 \partial x_n} \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial^2 f(x)}{\partial x_n \partial x_1} & \frac{\partial^2 f(x)}{\partial x_n \partial x_2} & \cdots & \frac{\partial^2 f(x)}{\partial x_n^2} \end{bmatrix} \in \mathbb{R}^{n \times n}$$

$$f(x) = x_1^2 + x_2^2$$

gradient: $\nabla_x f(x) = \begin{bmatrix} 2x_1 \\ 2x_2 \end{bmatrix} \Rightarrow$

$$\frac{\partial}{\partial x_1}(2x_1) = 2 \quad \bigg| \quad \frac{\partial}{\partial x_2}(2x_1) = 0$$

$$\frac{\partial}{\partial x_1}(2x_2) = 0 \quad \bigg| \quad \frac{\partial}{\partial x_2}(2x_2) = 2$$

Hessian: $\nabla_x^2 f(x) = \begin{bmatrix} 2 & 0 \\ 0 & 2 \end{bmatrix}$

## Problems 1.2 a-b

Solve them! Ask for help if you are stuck. Look at section 1.1 for help remembering how these gradients, Jacobians, and Hessians compute.

(a) Let $f(x_1, x_2) = x_1^2 + e^{x_1 x_2} + 2\log(x_2)$. What are the gradient and the Hessian of $f$?

(b) Note that $\nabla_x f : \mathbb{R}^n \to \mathbb{R}^n$. What is the Jacobian of $\nabla_x f$?

# Answers

(a) Let $f(x_1, x_2) = x_1^2 + e^{x_1 x_2} + 2\log(x_2)$. What are the gradient and the Hessian of $f$?

**Solution:**

$$\nabla_x f(x) = \begin{bmatrix} \frac{\partial f(x)}{\partial x_1} \\ \frac{\partial f(x)}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2x_1 + x_2 e^{x_1 x_2} \\ x_1 e^{x_1 x_2} + \frac{2}{x_2} \end{bmatrix} \text{ and } \nabla_x^2 f(x) = \begin{bmatrix} \frac{\partial^2 f(x)}{\partial x_1^2} & \frac{\partial^2 f(x)}{\partial x_1 \partial x_2} \\ \frac{\partial^2 f(x)}{\partial x_2 \partial x_1} & \frac{\partial^2 f(x)}{\partial x_2^2} \end{bmatrix} = \begin{bmatrix} 2 + x_2^2 e^{x_1 x_2} & e^{x_1 x_2} + x_1 x_2 e^{x_1 x_2} \\ e^{x_1 x_2} + x_1 x_2 e^{x_1 x_2} & x_1^2 e^{x_1 x_2} - \frac{2}{x_2^2} \end{bmatrix}$$

(b) Note that $\nabla_x f : \mathbb{R}^n \to \mathbb{R}^n$. What is the Jacobian of $\nabla_x f$?

**Equivalent**

**Solution:**

$$\nabla_x(\nabla_x f)(x) = \begin{bmatrix} \frac{\partial(\nabla_x f)_1(x)}{\partial x_1} & \frac{\partial(\nabla_x f)_1(x)}{\partial x_2} \\ \frac{\partial(\nabla_x f)_2(x)}{\partial x_1} & \frac{\partial(\nabla_x f)_2(x)}{\partial x_2} \end{bmatrix} = \begin{bmatrix} 2 + x_2^2 e^{x_1 x_2} & e^{x_1 x_2} + x_1 x_2 e^{x_1 x_2} \\ e^{x_1 x_2} + x_1 x_2 e^{x_1 x_2} & x_1^2 e^{x_1 x_2} - \frac{2}{x_2^2} \end{bmatrix} = \nabla_x^2 f(x)$$

# Approximations
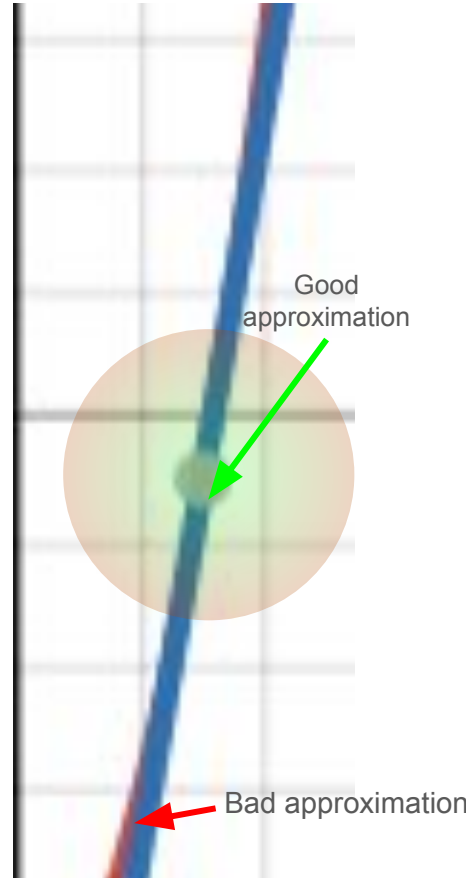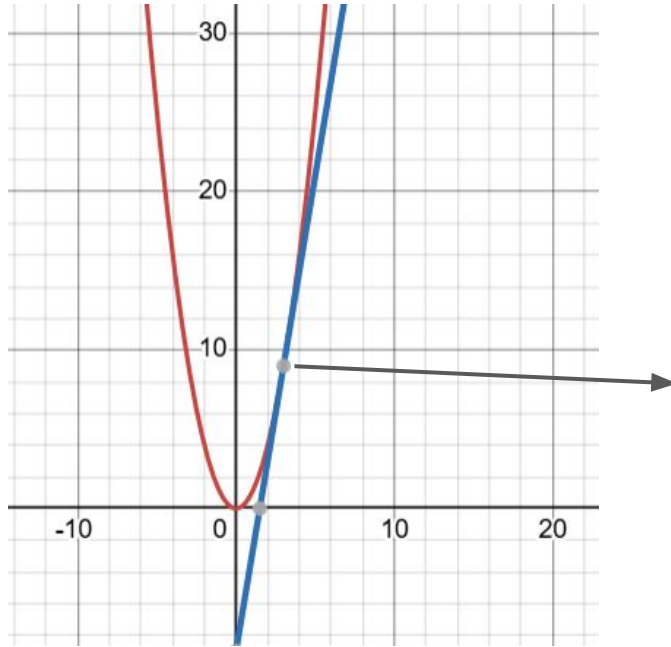
# Linear Approximation

The derivative of **f(x)** at some **(x,y)** can be used to linearly approximate **f(x ± ε)**

Where ε is very tiny!

This extends to multivariate functions… proof in your notes



Good approximation

Bad approximation

## Linear Approximation

For a "many-to-one" function, the <u>gradient</u> gives us a vector we can use to linearly approximate a small area around some **x**

What about a "many-to-many" function?

$$f : \mathbb{R}^n \rightarrow \mathbb{R}$$

Let $\epsilon = [\epsilon_1, \ldots, \epsilon_n]^T$ and $x = [x_1, \ldots, x_n]^T$

$$f(x + \epsilon) \approx f(x) + \nabla_x f(x)^T \epsilon$$

$$g : \mathbb{R}^n \rightarrow \mathbb{R}^m$$

# Problem 1.2 c

Remember that the Jacobian is just the gradient of a "many-to-many" function.

Also remember: *For a "many-to-one" function, the gradient gives us a vector we can use to linearly approximate a small area around some **x***

(c) The gradient $\nabla_x f(x)$ offers the best linear approximation of $f$ around the point $x$. What does the Jacobian of a function $g : \mathbb{R}^n \to \mathbb{R}^m$ offer?

# Answer

(c) The gradient $\nabla_x f(x)$ offers the best linear approximation of $f$ around the point $x$. What does the Jacobian of a function $g : \mathbb{R}^n \to \mathbb{R}^m$ offer?

**Solution:**

The Jacobian also offers the best linear approximation of $g$ around a point $x$, but now it approximates a vector, instead of a scalar,

$$g(x + \epsilon) \approx g(x) + \nabla_x g(x)\epsilon$$

where $\nabla_x g(x)\epsilon$ is a matrix multiplication instead of a dot product.

# Problem 1.2 d (walkthrough)

(d) If we use the gradient and the Hessian of $f : \mathbb{R}^n \to \mathbb{R}$, what type of an approximation for the function $f$ around a point $x$ can we create.

Remember Taylor expansion?

$\hookrightarrow$ To approximate a function around a point $\underline{a}$

$$f(x) \approx f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f'''(a)}{3!}(x-a)^3 \cdots$$

$\uparrow$

Exact at $\underline{a}$, close around $\underline{a}$

Better and better approximations

Remember Taylor expansion?

↳ To approximate a function around a point $\underline{a}$

$$f(x) \approx f(a) + \frac{f'(a)}{1!}(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \frac{f'''(a)}{3!}(x-a)^3 \cdots$$

Exact at $\underline{a}$, close around $\underline{a}$

Better and better approximations

Set $a = x$, we want to estimate $x + \epsilon$

$$f(x+\epsilon) \approx f(x) + \frac{f'(x)}{1!}(x+\epsilon-x) + \frac{f''(x)}{2!}(x+\epsilon-x)^2 + \cdots$$

↓

$$f(x+\epsilon) \approx f(x) + f'(x)\epsilon + \frac{1}{2}f''(x)\epsilon^2 + \cdots$$

Generalizing to vectors: $f: \mathbb{R}^n \longrightarrow \mathbb{R}$, $\begin{array}{c} x \in \mathbb{R}^n \\ \epsilon \in \mathbb{R}^n \end{array}$

$$f(x + \epsilon) \approx f(x) + \left(\nabla_x f(x)\right)^T \epsilon$$

Gradient = first order derivative of $f(x)$

So what is the second order derivative?

Second order derivative $= \nabla_x \left(\nabla_x f(x)\right) = \underline{\underline{Hessian}}$

$\hookrightarrow$ Gives us a $\underline{\text{Quadratic Approximation}}$

2nd order Taylor expansion around **x** generalized to vectors

$$f(x + \epsilon) \approx f(x) + \left(\nabla_x f(x)\right)^T \epsilon + \frac{1}{2} \epsilon^T \left(\nabla_x^2 f(x)\right)^T \epsilon$$   **Answer!**

# Problem 1.2 g (IMPORTANT!)

(g) Draw the gradient on the picture. Describe what happens to the values of the approximation of $f$ if we move from $x$ in directions $d_1, d_2, d_3$ for which $\nabla_x f(x)^T d_1 > 0, \nabla_x f(x)^T d_2 < 0, \nabla_x f(x)^T d_3 = 0$? Can the same conclusions be drawn about the function of $f$?

$$\left(\nabla_x f(x)\right)^T d_1 > 0$$

↳ Direction $d_1$ Points generally toward the gradient

$$\left(\nabla_x f(x)\right)^T d_2 < 0$$

↳ Direction $d_2$ Points generally away from the gradient

$$\left(\nabla_x f(x)\right)^T d_3 = 0$$

↳ Direction $d_3$ Points orthogonal to the gradient

# Answer

(g) Draw the gradient on the picture. Describe what happens to the values of the approximation of $f$ if we move from $x$ in directions $d_1, d_2, d_3$ for which $\nabla_x f(x)^T d_1 > 0, \nabla_x f(x)^T d_2 < 0, \nabla_x f(x)^T d_3 = 0$? Can the same conclusions be drawn about the function of $f$?

**Solution:**

- $d_1$: Value of approximation goes up.
- $d_2$: Value of approximation goes down.
- $d_3$: Value of approximation stays the same.

The same can be said for $f$, but only in the immediate vicinity of the point $x$.

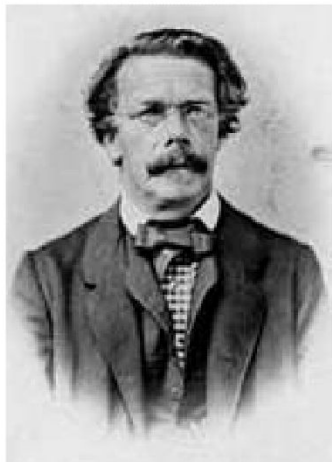Intuition used here will be useful on the exam

# Properties

# Useful rules!

Let $f : \mathbb{R}^n \to \mathbb{R}$, $g : \mathbb{R}^n \to \mathbb{R}$, . Below is a list of important gradient properties:

- **Gradient of constant:** $\nabla_x c = 0 \in \mathbb{R}^n$ for a constant $c \in \mathbb{R}^n$.

- **Linearity:** $\nabla_x(\alpha f + \beta g)(x) = \alpha \nabla_x f(x) + \beta \nabla_x g(x)$ for a scalars $\alpha, \beta \in \mathbb{R}$.

- **Product rule:** $\nabla_x(fg)(x) = \nabla_x f(x) \cdot g(x) + \nabla_x g(x) \cdot f(x)$.

Let $f : \mathbb{R}^n \to \mathbb{R}^m$, $g : \mathbb{R}^n \to \mathbb{R}^m$, $h : \mathbb{R}^m \to \mathbb{R}^k$, $l : \mathbb{R}^m \to \mathbb{R}$. Below is a list of important Jacobian properties:

- **Jacobian of constant:** $\nabla_x c = 0 \in \mathbb{R}^{n \times m}$ for a constant $c \in \mathbb{R}^n$.

- **Linearity:** $\nabla_x(\alpha f + \beta g)(x) = \alpha \nabla_x f(x) + \beta \nabla_x g(x)$ for a scalars $\alpha, \beta \in \mathbb{R}$.

- **Product rule:** $\nabla_x(f^T g)(x) = [\nabla_x f(x)]^T g(x) + [\nabla_x g(x)]^T f(x)$.

- **Chain rule:** $\nabla_x(h \circ g)(x) = \nabla_{g(x)} h(g(x)) \nabla_x g(x)$ and $\nabla_x(l \circ g)(x) = \left[[\nabla_{g(x)} l(g(x))]^T \nabla_x g(x)\right]^T$.

Ludwig Otto Hesse  VS  Carl Gustav Jacob Jacobi  VS  William Grady Hamilton  William Rowan Hamilton

Hessian  Jacobian  Gradient

# Questions/Chat Time!