

C477: Computing for Optimal Decisions Background in Linear Algebra & Calculus

Ruth Misener & Panos Parpas

October 7, 2018

Note: These notes summarise the background material in linear algebra and calculus necessary for C447. If this material is unfamiliar, please revise independently. It is in your best interest to revise before the first week of classes ends; C477 will assume that you understand the following notes. To assist revision, we offer the following:

- Post questions to Piazza: piazza.com/imperial.ac.uk/fall2018/477/home;
- Quiz 0 posted on CATE will help you assess your level of preparation for the module. It does not affect your final mark. The answers to the quiz will be posted on the morning of 16 October.
- Revise using notes from C145 (Mathematical Methods) and C233 (Computational Techniques); we have posted these notes on Piazza: <https://piazza.com/imperial.ac.uk/fall2018/477/resources>.

We will also assume familiarity with the material in C496 (Mathematics for Inference and Machine Learning).

Contents

1	Linear Algebra	3
1.1	Vectors & Matrices	3
	Vector Definitions (column vector, row vector, transpose)	3
	Matrix Definitions (matrix, symmetric matrix)	3
	Vector Operations (equality test, summation, scalar multiplication) .	4
	Matrix Multiplication	4
1.2	Linear Systems of Equations	4
	Linear Independence & Matrix Rank	5
1.3	Eigenvalues & Eigenvectors	6
	Positive definite & positive semidefinite matrices	6
1.4	Inner product	6
	Orthogonality	6
1.5	Norms	7
	Examples: 1-norm, 2-norm, ∞ -norm	7
1.6	Angles Between Vectors	8
1.7	Cauchy-Schwarz Inequality	8
2	Analysis	8
2.1	Continuous Functions	8
3	Multivariable Calculus	9
3.1	Gradients & Jacobian	10
	Example	10
3.2	Hessian	11
3.3	The Chain Rule	11
3.4	Mean Value Theorem	12
3.5	Taylor Series	13

1 Linear Algebra

Alternate revision sources: C145 notes *Matrices and Vectors* by Dr Mahdi Cheraghchi & Dr Marc Deisenroth (provided on Piazza: <https://piazza.com/imperial.ac.uk/fall2018/477/resources>); Chapters 2-3, *An Introduction to Optimization*, Chong & Żak; C233 *General Lecture Notes* by Prof István Maros gives a complete development.

1.1 Vectors & Matrices

Vector Definitions (column vector, row vector, transpose)

- We define a **column vector** in \mathbb{R}^n as an array of n numbers,

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

- A **row vector** in \mathbb{R}^n is

$$\mathbf{x} = (x_1, x_2, \dots, x_n)$$

- The **transpose** of a column vector is a row vector. If x is a column vector in \mathbb{R}^n then,

$$\mathbf{x}^\top = (x_1, x_2, \dots, x_n)$$

- In C477 all vectors are column vectors; we use lower-case bold font to represent vectors \mathbf{x} rather than arrow notation \vec{x} . Note that this is opposite from C496. In C496, all vectors are row vectors. The courses are following two different books; this is the reason for the discrepancy.

Matrix Definitions (matrix, symmetric matrix)

- A **matrix** is a rectangular array of entries with n rows and m columns.
- A matrix B with n rows and m columns belongs to $\mathbb{R}^{n \times m}$
- A **symmetric matrix** is one that is equal to its transpose, i.e., $B = B^\top$.

Vector Operations (equality test, summation, scalar multiplication)

- Two vectors \mathbf{x} and \mathbf{y} are **equal**, $\mathbf{x} = \mathbf{y}$ if $x_i = y_i \forall i = 1, \dots, n$
- The **sum of two vectors** $\mathbf{x} + \mathbf{y}$ is the vector

$$(x_1 + y_1, x_2 + y_2, \dots, x_n + y_n)^\top$$

- **Multiplication** of a vector \mathbf{x} by a scalar α is defined as,

$$\alpha \mathbf{x} = (\alpha x_1, \alpha x_2, \dots, \alpha x_n)^\top$$

Matrix Multiplication

- **Multiplication:** If $A \in \mathbb{R}^{n \times m}$ and $B \in \mathbb{R}^{m \times k}$, then $C = A B$ is a matrix in $\mathbb{R}^{n \times k}$. The $(i, j)^{th}$ entry of C is,

$$C_{ij} = \sum_{l=1}^m A_{il} \cdot B_{lj}$$

- For two matrices A and B , the **transpose of their product** is $(AB)^\top = B^\top A^\top$.

1.2 Linear Systems of Equations

Definitions Suppose that we are given m equations in n unknowns of the form,

$$\begin{array}{ccccccc} a_{11} x_1 & + a_{12} x_2 & + \dots & + a_{1n} x_n & = & b_1 \\ a_{21} x_1 & + a_{22} x_2 & + \dots & + a_{2n} x_n & = & b_2 \\ \vdots & & & & & \vdots \\ a_{m1} x_1 & + a_{m2} x_2 & + \dots & + a_{mn} x_n & = & b_m. \end{array}$$

We can also represent the set of equations above in vector notation,

$$\mathbf{a}_1 x_1 + \mathbf{a}_2 x_2 + \dots + \mathbf{a}_n x_n = \mathbf{b}$$

where

$$\mathbf{a}_j = \begin{pmatrix} a_{1j} \\ a_{2j} \\ \vdots \\ a_{mj} \end{pmatrix} \quad \text{and} \quad \mathbf{b} = \begin{pmatrix} b_1 \\ b_2 \\ \vdots \\ b_m \end{pmatrix}$$

We associate the matrix,

$$A = [\mathbf{a}_1, \mathbf{a}_2, \dots, \mathbf{a}_n]$$

with the system of equations, and represent the system as follows,

$$A \mathbf{x} = \mathbf{b}$$

where

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

Linear Independence & Matrix Rank

- A vector \mathbf{y} is said to be a **linear combination of vectors** $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$ if there are scalars $\alpha_1, \alpha_2, \dots, \alpha_m$ such that,

$$\mathbf{y} = \alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2 + \dots + \alpha_m \mathbf{x}_m = \sum_{i=1}^m \alpha_i \mathbf{x}_i$$

- A set of vectors $\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_m$ is said to be **linearly independent** if the equality,

$$\alpha_1 \mathbf{x}_1 + \alpha_2 \mathbf{x}_2 + \dots + \alpha_m \mathbf{x}_m = \sum_{i=1}^m \alpha_i \mathbf{x}_i = \mathbf{0},$$

implies that all coefficients $\alpha_i, i = 1, \dots, m$ are equal to zero.

- **Exercise:** A set of vectors contains the zero vector, is the set linearly independent?

Definition 1.1 (Rank of a Matrix). *The maximum number of linearly independent columns of A is called the rank of the matrix, denoted by $\text{rank}(A)$*

Theorem 1.1 (Solution of linear systems of equations). *The system of equations*

$$A \mathbf{x} = \mathbf{b},$$

where $A \in \mathbb{R}^{m \times n}$, $\mathbf{x} \in \mathbb{R}^n$, $\mathbf{b} \in \mathbb{R}^m$ has a solution **if and only if**,

$$\text{rank}(A) = \text{rank}([A, \mathbf{b}]) = n.$$

If $m = n$ and the condition above holds then the system has a unique solution

1.3 Eigenvalues & Eigenvectors

Definitions

- An **eigenvector** of square matrix $A \in \mathbb{R}^{n \times n}$ is a vector $\mathbf{v} \in \mathbb{R}^n$ such that the product $A \mathbf{v}$ is equal to a scalar multiple ($\lambda \in \mathbb{R}$) of \mathbf{v} :

$$A \mathbf{v} = \lambda \mathbf{v}$$

- The scalars λ are called **eigenvalues**; a matrix is positive definite if all eigenvalues are positive.

Positive definite & positive semidefinite matrices

- A matrix $A \in \mathbb{R}^{n \times n}$ is called **positive semidefinite** if for all $\mathbf{d} \in \mathbb{R}^n$,

$$\mathbf{d}^\top A \mathbf{d} \geq 0.$$

We use the notation $A \succeq 0$. All eigenvalues of $(A + A^\top)/2$ are non-negative.

- If the above inequality is satisfied strictly, i.e. if

$$\mathbf{d}^\top A \mathbf{d} > 0 \quad \forall \mathbf{d} \in \mathbb{R}^n \setminus \{0\},$$

then A is called **positive definite**. We use the notation $A \succ 0$. All eigenvalues of $(A + A^\top)/2$ are positive.

1.4 Inner product

Definition: $\mathbf{x}^\top \mathbf{y}$ is called the **inner product** between two vectors and it is sometimes written as $\langle \mathbf{x}, \mathbf{y} \rangle$, where

$$\langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y} = \sum_{i=1}^n x_i y_i$$

Orthogonality

Two vectors \mathbf{x} and \mathbf{y} are **orthogonal** if $\langle \mathbf{x}, \mathbf{y} \rangle = 0$. By this definition, the zero vector is orthogonal to every other vector. An example of two nonzero orthogonal vectors:

$$\langle \mathbf{x}, \mathbf{y} \rangle = [1, 2] \begin{bmatrix} -3 \\ 3/2 \end{bmatrix}$$

1.5 Norms

Definition: A vector norm is a function $\|\cdot\| : \mathbb{R}^n \mapsto \mathbb{R}$ satisfying the following properties:

1. Positivity: $\|\mathbf{x}\| \geq 0$ and $\|\mathbf{x}\| = 0$ if and only if $\mathbf{x} = \mathbf{0}$;
2. Homogeneity: $\|r \mathbf{x}\| = |r| \|\mathbf{x}\|$ for all $r \in \mathbb{R}$;
3. Triangle Inequality: $\|\mathbf{x} + \mathbf{y}\| \leq \|\mathbf{x}\| + \|\mathbf{y}\|$

Examples: 1-norm, 2-norm, ∞ -norm

1-norm $\|\mathbf{x}\|_1 = |x_1| + |x_2| + \dots + |x_n|$

2-norm $\|\mathbf{x}\|_2 = \sqrt{x_1^2 + x_2^2 + \dots + x_n^2}$

∞ -norm $\|\mathbf{x}\|_\infty = \max_i |x_i|$

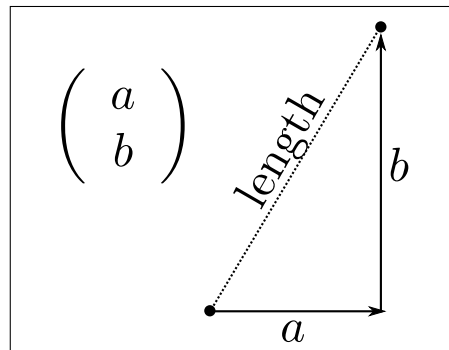


Figure 1: The length of vector $\mathbf{x} = [a, b]^\top$ is the Euclidean norm $\|\mathbf{x}\|_2 = \sqrt{a^2 + b^2}$. The 1-norm of vector \mathbf{x} is $a + b$ and the ∞ -norm is b .

- The **2-norm** is also called the **Euclidean norm**. When no index is specified on a norm, e.g., $\|\cdot\|$, this is considered to be the Euclidean norm.
- $\|\mathbf{x}\|_1 \geq \|\mathbf{x}\|_2 \geq \|\mathbf{x}\|_\infty$

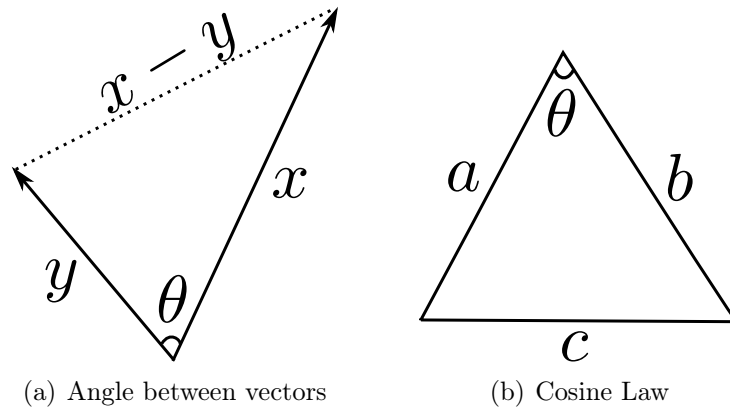


Figure 2: Illustration of angles between vectors and the cosine law in Section 1.6

1.6 Angles Between Vectors

The **angle between two vectors** $\mathbf{x}, \mathbf{y} \in \mathbb{R}^n$ is the angle $\theta \in [0, \pi]$,

$$\cos(\theta) = \frac{\mathbf{x}^\top \mathbf{y}}{\|\mathbf{x}\|_2 \|\mathbf{y}\|_2}.$$

Cosine Law: For any triangle with sides of length a , b and c the angle opposite c satisfies (see Fig. 2),

$$c^2 = a^2 + b^2 - 2 a b \cos(\theta)$$

1.7 Cauchy-Schwarz Inequality

The **Cauchy-Schwarz Inequality** states: $|\langle \mathbf{x}, \mathbf{y} \rangle| \leq \|\mathbf{x}\|_2 \|\mathbf{y}\|_2$

2 Analysis

2.1 Continuous Functions

Definitions

- We write $f : X \rightarrow Y$ to mean that a function takes points from the set X (domain) to the set Y (range).

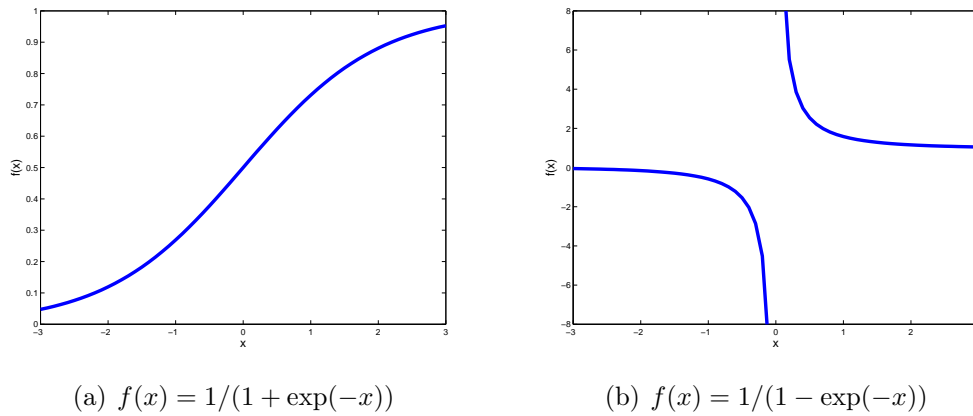


Figure 3: Example of a continuous and a discontinuous function at $x = 0$

- Let $X \subset \mathbb{R}^n$, then the function $f : X \rightarrow \mathbb{R}^m$ is a **continuous function** at \mathbf{x}_0 if and only if for every $\epsilon > 0$ there exists a $\delta > 0$ such that when $\|\mathbf{x} - \mathbf{x}_0\|_2 < \delta$ then $\|f(\mathbf{x}) - f(\mathbf{x}_0)\|_2 < \epsilon$.

Example: The function,

$$f(\mathbf{x}) = f(x_1, x_2) = \begin{bmatrix} x_1^2 + x_2 \\ \exp(x_1) + x_2 \\ x_2 \end{bmatrix},$$

evaluated at the point $(2, -1)$ is,

$$f(\mathbf{x}) = f(2, -1) = \begin{bmatrix} 3 \\ \exp(2) - 1 \\ -1 \end{bmatrix}$$

3 Multivariable Calculus

Alternate revision sources: C145 notes *Calculus* by Dr Mahdi Cheraghchi & Dr Marc Deisenroth (provided on Piazza: <https://piazza.com/imperial.ac.uk/fall2018/477/resources>); Chapter 5, *An Introduction to Optimization*, Chong & Āak; C233 *General Lecture Notes* by Prof István Maros gives a complete development.

Definitions

- The derivative of a function in one dimension is defined below.

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

When the above holds, in the sense that the limit exists, then we say that the **function is differentiable at the point x** .

- When a function is differentiable and its derivative is also continuous we say that the function is **continuously differentiable**.
- If a function is defined over \mathbb{R}^n , then its **partial derivative** with respect to dimension i is defined as,

$$\frac{\partial f(\mathbf{x})}{\partial x_i} = \lim_{h \rightarrow 0} \frac{f(x_1, \dots, x_i + h, \dots, x_n) - f(x_1, \dots, x_i, \dots, x_n)}{h}$$

3.1 Gradients & Jacobian

Definitions

- The vector of the partial derivatives is the **gradient**: $\nabla f(\mathbf{x}) = \begin{bmatrix} \frac{\partial f(\mathbf{x})}{\partial x_1} \\ \vdots \\ \frac{\partial f(\mathbf{x})}{\partial x_n} \end{bmatrix}$;
- If $f : \mathbb{R}^n \rightarrow \mathbb{R}^m$ we evaluate the gradient of each function, and call the matrix of derivatives $\nabla f(\mathbf{x})$ the gradient of f :

$$= \begin{bmatrix} \frac{\partial f_1(\mathbf{x})}{\partial x_1} & \dots & \frac{\partial f_m(\mathbf{x})}{\partial x_1} \\ \vdots & & \vdots \\ \frac{\partial f_1(\mathbf{x})}{\partial x_n} & \dots & \frac{\partial f_m(\mathbf{x})}{\partial x_n} \end{bmatrix};$$

- The transpose of the gradient is called the **Jacobian** matrix.

Example

Find the gradient matrix of,

$$f(x_1, x_2) = \begin{bmatrix} x_1 x_2 \\ \sin(x_1 + x_2) \\ x_1^2 - x_2^2 \end{bmatrix}$$

Then,

$$\nabla f_1(\mathbf{x}) = \begin{bmatrix} x_2 \\ x_1 \end{bmatrix}, \quad \nabla f_2(\mathbf{x}) = \begin{bmatrix} \cos(x_1 + x_2) \\ \cos(x_1 + x_2) \end{bmatrix}, \quad \nabla f_3(\mathbf{x}) = \begin{bmatrix} 2x_1 \\ -2x_2 \end{bmatrix}$$

We therefore have the gradient matrix.

$$\nabla f(\mathbf{x}) = \begin{bmatrix} \nabla f_1(\mathbf{x}) & \nabla f_2(\mathbf{x}) & \nabla f_3(\mathbf{x}) \end{bmatrix} = \begin{bmatrix} x_2 & \cos(x_1 + x_2) & 2x_1 \\ x_1 & \cos(x_1 + x_2) & -2x_2 \end{bmatrix}$$

3.2 Hessian

Definitions

- Given a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$,

$$\frac{\partial^2 f}{\partial x_i \partial x_j},$$

to denote the i^{th} partial derivative of $\frac{\partial f}{\partial x_j}$.

- We define the matrix $\nabla^2 f(\mathbf{x})$ to denote the matrix whose $(i, j)^{th}$ entry is given by $\frac{\partial^2 f}{\partial x_i \partial x_j}$ as the **Hessian** matrix of f , i.e.

$$Hf(\mathbf{x}) = \nabla^2 f(\mathbf{x}) = \begin{bmatrix} \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_1} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_1 \partial x_n} \\ \vdots & & \vdots \\ \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_1} & \cdots & \frac{\partial^2 f(\mathbf{x})}{\partial x_n \partial x_n} \end{bmatrix}$$

Example: Find the Hessian matrix of the function $f : \mathbb{R}^3 \rightarrow \mathbb{R}$, given by $f(x_1, x_2, x_3) = x_1^2 + x_1 x_3 + x_1 x_2 + x_3^2 + \exp(x_1 x_3)$

$$\nabla^2 f(\mathbf{x}) = \begin{bmatrix} 2 + x_3^2 \exp(x_1 x_3) & 1 & 1 + (1 + x_1 x_3) \exp(x_1 x_3) \\ 1 & 0 & 0 \\ 1 + (1 + x_1 x_3) \exp(x_1 x_3) & 0 & 2 + x_1^2 \exp(x_1 x_3) \end{bmatrix}$$

3.3 The Chain Rule

The *chain rule* is used for differentiating the composition $g(f(t))$ of a function $g : \mathbb{R}^n \rightarrow \mathbb{R}$, and $f : \mathbb{R} \rightarrow \mathbb{R}^n$. In particular, suppose that $g : \mathbb{R}^n \rightarrow \mathbb{R}$ is differentiable,

and suppose that $f : \mathbb{R} \rightarrow \mathbb{R}^n$ is differentiable in \mathbb{R} . Then the derivative of the composite function $h(t) = g(f(t))$ is given by,

$$\frac{dh(t)}{dt} = h'(t) = \nabla g(f(t))^\top \begin{pmatrix} \frac{df_1(t)}{dt} \\ \frac{df_2(t)}{dt} \\ \vdots \\ \frac{df_n(t)}{dt} \end{pmatrix}$$

Let $x(t) = [e^t + t^3, t^2, t + 1]^\top$, $t \in \mathbb{R}$, and $f(\mathbf{x}) = x_1^3 x_2 x_3^2 + x_1 x_2 + x_3$, $x = [x_1, x_2, x_3] \in \mathbb{R}^3$. Find $\frac{df(x(t))}{dt}$ in terms of t .

$$\begin{aligned} \frac{df(x(t))}{dt} &= \nabla f(x(t))^\top \frac{dx(t)}{dt} \\ &= [3x_1(t)^2 x_2(t) x_3(t)^2 + x_2(t), x_1(t)^3 x_3(t)^2 + x_1(t), 2x_1(t)^3 x_2(t) x_3(t) + 1] \begin{bmatrix} e^t + 3t^2 \\ 2t \\ 1 \end{bmatrix} \\ &= 12t(e^t + 3t^2)^3 + 2te^t + 6t^2 + 2t + 1 \end{aligned}$$

3.4 Mean Value Theorem

1-D: For function defined over the closed interval $[a, b]$ there must exist a point $a < c < b$ such that:

$$f(b) = f(a) + \frac{df(x)}{dx} \Big|_{x=c} (b - a)$$

n-D If a function is differentiable in U and suppose that the interval $[a, b]$ is contained in U , then there exists a point $c \in [a, b]$ such that,

$$f(b) = f(a) + \nabla f(\mathbf{x})^T \Big|_{\mathbf{x}=c} (b - a).$$

Consider the following function, illustrated in Figure 4, $f : \mathbb{R} \rightarrow \mathbb{R}$,

$$f(x) = (x - 2)(x - 3) + (x - 1)^3.$$

On the interval $[-2, 3]$. Show that there exist two points c_1 and c_2 such that,

$$f(3) = f(-2) + 5 \frac{df}{dx} \Big|_{x=c_i}.$$

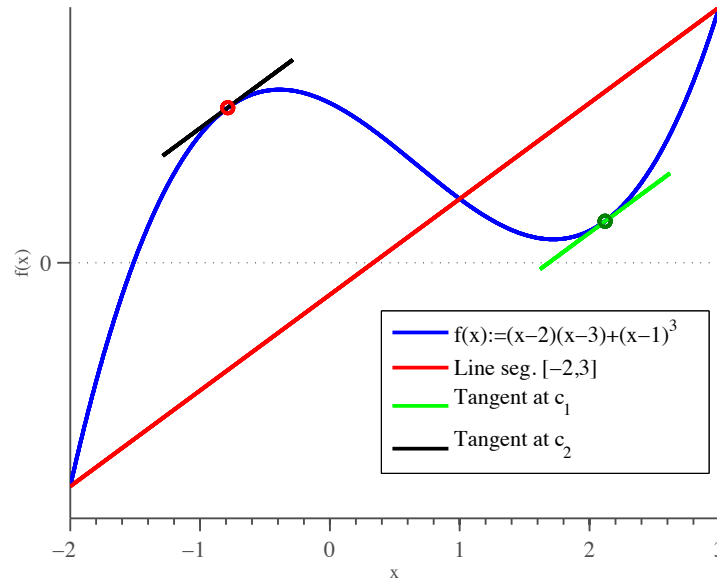


Figure 4: Illustration of the example from Section 3.4

Since $\frac{df}{dx} = 2x - 5 + 3(x-1)^2$ we obtain the following quadratic equation for c ,

$$3c^2 - 4c - 5 = 0$$

It has two roots $c_1 = (4 - \sqrt{76})/6$ and $c_2 = (4 + \sqrt{76})/6$. It is easy to verify that these two points satisfy,

$$f(b) = f(a) + \nabla f(\mathbf{x}) \Big|_{\mathbf{x}=c} (b - a),$$

with $a = -2$ and $b = 3$.

3.5 Taylor Series

In many dimensions a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ has the first order expansion about \mathbf{x}_0

$$f(\mathbf{x}) = f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0) + r(\mathbf{x} - \mathbf{x}_0)$$

Where the residual $r(\mathbf{x} - \mathbf{x}_0)$ satisfies $\lim_{\mathbf{x} - \mathbf{x}_0 \rightarrow 0} r(\mathbf{x} - \mathbf{x}_0) = 0$.

The second order expansion is given by,

$$f(\mathbf{x}) = f(\mathbf{x}_0) + \nabla f(\mathbf{x}_0)^T (\mathbf{x} - \mathbf{x}_0) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_0)^T \nabla^2 f(\mathbf{x}_0) (\mathbf{x} - \mathbf{x}_0) + r(\mathbf{x} - \mathbf{x}_0)$$