

# Prerequisites for Foundations of Machine Learning and Data Science

Luke Dickens

September 24, 2018

## 1 A note on these prerequisites

This module is designed to give you the **mathematical** intuition behind a number of core/common machine learning algorithms. There is a significant practical aspect too, in which you will learn to use the methods and evaluate their performance on real world data. However, the module is not purely about learning to use machine learning libraries. It is about learning how and why they work too. Incoming students are expected to have **some background in mathematics** (see the next section) and to have **the rudiments of programming in Python**. If you have never programmed in Python before, then you should have all you need by simply working your way through the following online tutorial: DataCamp – Intro to Python for Data Science <sup>1</sup> (it should take you about 4 hours).

## 2 Mathematical Prerequisites

Below is a brief list of the mathematical notation, results and techniques that you would be expected to be familiar with before starting the module. Much of this material is taught in a UK mathematics A-level (school leaving certificate). However, some of the notation may be unfamiliar, and some ideas might go a little beyond high school mathematics. These represent many of the concepts that we will be building on for the content of the module. You should therefore expect to encounter them early in the module, and for the module material to extend and generalise from these concepts. There are also some questions within the table. These represent ideas that are less essential you know in advance, but still useful.

If there are a small number of items that are hazy or unfamiliar, then there is no reason you cannot familiarise/refamiliarise yourself with these in advance to taking the module. However, if a significant proportion is new to you, then you may find the mathematics in the module too challenging. Please contact me ([1.dickens@ucl.ac.uk](mailto:1.dickens@ucl.ac.uk)) if you would like some advice/links to supporting material.

Prerequisite	Symbols	Examples & Notes
<b>Sets/Discrete</b>		
Set membership	$\in, \notin$	$x \in X$ ( $x$ is member of $X$ ) $x \notin X$ ( $x$ is not a member of $X$ )
Set operations	$\cup$ (union)	If $x \in A \cap B$ then $x \in A$ and $x \in B$

Table 1: You should be relatively familiar with the majority of the mathematical notation, results and techniques outlined in this table.

<sup>1</sup><https://www.datacamp.com/courses/intro-to-python-for-data-science>

Prerequisite	Symbols	Examples & Notes
Set relations	$\cap$ (intersection)	If $x \in A \cup B$ then $x \in A$ or $x \in B$
	$\setminus$ (set-minus)	If $x \in A \setminus B$ then $x \in A$ and $x \notin B$
	$\subseteq$ (subset of) $\subset$ (strict subset of)	How would you formalise these concepts using the above notation?
Types of numbers	$\mathbb{N}$ (natural numbers), $\mathbb{Z}$ (integers) $\mathbb{R}$ (real numbers)	The integers are part of the reals, $\mathbb{Z} \subset \mathbb{R}$ . Often we denote a range of real numbers as $[a, b] \subset \mathbb{R}$ meaning all numbers between $a$ and $b$ (inclusive).
Partition (of a set)		A set's partition is a grouping of that set's elements into non-empty subsets, such that every element is included in one and only one of the subsets. Example: if $\{A_1, A_2, \dots, A_n\}$ is a partition of set $A$ , then $A_i \cap A_j = \emptyset$ (the emptyset) for any $i \neq j$ and $A_1 \cup A_2 \cup \dots \cup A_n = A$ .
Summation	$\sum$	$\sum_{i=1}^n a_i = a_1 + \dots + a_n$
Product	$\prod$	$\prod_{i=1}^n a_i = a_1 \times \dots \times a_n$
<b>Functions</b>		
Integer Powers	$a^n$ ( $a \in \mathbb{R}, n \in \mathbb{Z}$ )	$a^n = a \cdot a \cdot \dots \cdot a$ for $n > 0$ $b^{-1} = \frac{1}{b}$ for $b \neq 0$ (what if $b = 0$ ?) we normally say $a^0 = 1$ (why?)
General Powers	$a^b$ ( $a, b \in \mathbb{R}$ )	$b^{m+n} = b^m \cdot b^n$ $(b^m)^n = b^{m \cdot n}$ $(b \cdot c)^n = b^n \cdot c^n$ How do we define $a^b$ for $a, b \in \mathbb{R}$ ?
Logarithms	$\ln$ (natural log) $\log_b$ (log base $b$ )	$\ln a = x$ means $a = e^x$ , $\ln_b a = c$ means $a = b^c$ , $\log_{10} 100 = 2$ , $\log_2 64 = 8$ , $\ln xy = \ln x + \ln y$ , $\ln \frac{x}{y} = \ln x - \ln y$ , $\ln x^n = n \ln x$
Factorials	$!$	$n! = n \cdot (n-1) \cdot \dots \cdot 1$ for $n \in \mathbb{N}$ We normally say $0! = 1$ (why?)

Table 1: You should be relatively familiar with the majority of the mathematical notation, results and techniques outlined in this table.

Prerequisite	Symbols	Examples & Notes
Functions	$f : \mathbb{R} \rightarrow \mathbb{R}$  $f : \mathbb{R}^2 \rightarrow \mathbb{R}$  $f : \mathbb{R}^m \rightarrow \mathbb{R}^n$	<p>A function <math>f</math> which takes a single real number input and outputs a single real number, e.g. <math>f(x) = y</math> where <math>x, y \in \mathbb{R}</math>, can be written <math>f : \mathbb{R} \rightarrow \mathbb{R}</math>.</p> <p>A function can take 2 inputs, and output a scalar. Here, <math>f</math> takes a two real number inputs and outputs a single real number, e.g. <math>f(x, y) = z</math> where <math>x, y, z \in \mathbb{R}</math>.</p> <p>More generally, a function can take <math>m</math> inputs (or an <math>m</math>-vector of inputs) and output an <math>n</math>-vector.</p>
<b>Vector/Matrix</b>		
Vectors	$\mathbf{v} \in \mathbb{R}^n$	<p><math>\mathbf{v} \in \mathbb{R}^n</math> means <math>\mathbf{v}</math> is a <i>stack</i> of <math>n</math> real numbers.</p> <p><math>v_i</math> used to refer to the <math>i</math>th element of <math>\mathbf{v}</math></p>
Matrices	$A \in \mathbb{R}^{m \times n}$	<p><math>A \in \mathbb{R}^{m \times n}</math> is a matrix with <math>m</math> rows and <math>n</math> columns we can visualise it as:</p> $A = \begin{pmatrix} a_{11} & a_{12} & \cdots & a_{1n} \\ a_{21} & a_{22} & \cdots & a_{2n} \\ \vdots & \vdots & \ddots & \vdots \\ a_{m1} & a_{m2} & \cdots & a_{mn} \end{pmatrix} \in \mathbb{R}^{m \times n}$ <p>We can refer to the <math>(i, j)</math>th element as <math>[A]_{ij}</math></p> <p><math>A^T</math> is the transpose of <math>A</math> (switching rows and columns) so <math>[A^T]_{ij} = [A]_{ji}</math></p>
The dot product	$\mathbf{a} \cdot \mathbf{b}$ or $\mathbf{a}^T \mathbf{b}$	<p><math>\mathbf{a} \cdot \mathbf{b} = \sum_{i=1}^n a_i b_i</math></p> <p>The dot product is sometimes written <math>\mathbf{a}^T \mathbf{b}</math> assuming column vectors (why?)</p> <p>The (Euclidean) length of a vector is:</p> $\ \mathbf{v}\  = \sqrt{\mathbf{v}^T \mathbf{v}} = \sqrt{\mathbf{v}^2}$ <p>The (Euclidean) distance between two vectors is:</p> $\ \mathbf{a} - \mathbf{b}\  = \sqrt{(\mathbf{a} - \mathbf{b})^T (\mathbf{a} - \mathbf{b})} = \sqrt{(\mathbf{a} - \mathbf{b})^2}$
Cosine Rule		<p>If <math>\theta</math> is the angle between <math>\mathbf{a}</math> and <math>\mathbf{b}</math>:</p> $\mathbf{a}^T \mathbf{b} = \ \mathbf{a}\  \ \mathbf{b}\  \cos \theta$

Table 1: You should be relatively familiar with the majority of the mathematical notation, results and techniques outlined in this table.

Prerequisite	Symbols	Examples & Notes
Matrix multiplication		$AB$ for $A \in \mathbb{R}^{m \times n}$ and $B \in \mathbb{R}^{n \times p}$ $[AB]_{ij} = \sum_k a_{ik} b_{kj}$
Matrix Determinant	$\det(A)$ or $ A $	<p>The determinant of a square matrix <math>A</math> is a positive number that can be computed from <math>A</math> elements. For <math>A \in \mathbb{R}^{2 \times 2}</math>,</p> $ A  = \begin{vmatrix} a & b \\ c & d \end{vmatrix} = ad - bc$ <p>Geometrically, <math> A </math> is the scaling factor of the linear transformation described by matrix <math>A</math>.</p>
Matrix Rank	$\text{rank}(A)$	<p>The rank of a matrix <math>A</math> is the number of linearly independent columns (or rows) of <math>A</math>. Equally, if we treat every column of <math>A</math> as a vector, then <math>\text{rank}(A)</math> is the dimension of the space spanned by these vectors.</p>
Identity Matrix	$I, I_n$	<p>The square matrix <math>I \in \mathbb{R}^{n \times n}</math> has 1s on the diagonal and zeros elsewhere, so <math>[I]_{ii} = 1</math> and <math>[I]_{ij} = 0</math> when <math>i \neq j</math>. <math>I</math> is the matrix analog of unity (1). If the dimension of <math>I</math> is unclear, we can denote it with a subscript, e.g. <math>I_m \in \mathbb{R}^{m \times m}</math>.  Multiplying by <math>I</math> leaves your matrix (or vector) unchanged, e.g. for <math>A \in \mathbb{R}^{m \times n}</math>, <math>I_m A = A = A I_n</math></p>
Matrix Inversion	$A^{-1}$	<p>If a square matrix <math>A \in \mathbb{R}^{n \times n}</math> is invertible, then its inverse <math>A^{-1}</math> multiplied (left or right) by <math>A</math> gives the identity matrix, i.e.</p> $A^{-1}A = I = AA^{-1}$ <p>Not all (square) matrices are invertible. How would you check if a matrix is invertible?  The inverse is the matrix analog of the reciprocal <math>r^{-1} = \frac{1}{r}</math> of a scalar <math>r \in \mathbb{R}</math>. Do you know how to find the inverse of a <math>2 \times 2</math> matrix?</p>

Table 1: You should be relatively familiar with the majority of the mathematical notation, results and techniques outlined in this table.

Prerequisite	Symbols	Examples & Notes
Eigenvalues Eigenvectors	&	<p>An eigenvector <math>\mathbf{x} \in \mathbb{R}^n</math> of a square-matrix <math>A \in \mathbb{R}^{n \times n}</math> satisfies the following equation:</p> $A\mathbf{x} = \lambda\mathbf{x}$ <p>for some corresponding (scalar) eigenvalue <math>\lambda \in \mathbb{R}</math>. Geometrically an eigenvector, corresponding to a real nonzero eigenvalue, points in a direction that is stretched by the transformation and the eigenvalue is the factor by which it is stretched.</p>
<b>Calculus</b>		
Differentiation	$f'(x), \frac{df}{dx}$	<p>For a function <math>f : \mathbb{R} \rightarrow \mathbb{R}</math>, then the slope at <math>x</math> (if it exists) is</p> $f'(x) = \frac{df}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$
Simple Derivatives		<p>e.g. <math>\frac{d}{dx}c = 0</math> for constant <math>c</math></p> $\frac{d}{dx}x^a = a \cdot x^{a-1}$ $\frac{d}{dx}e^x = e^x$ $\frac{d}{dx}\ln x = \frac{1}{x}$ $\frac{d}{dx}\sin(x) = \cos(x)$ $\frac{d}{dx}\cos(x) = -\sin(x)$ <p>What are the product, quotient and chain rules?</p>
Integration	$\int_a^b f(x)dx,$ $\int f(x)dx$	<p>To know what is meant by a definite and indefinite integral (Riemann integration), and to have some idea of the requirements on <math>f(x)</math>, e.g. finite etc.</p>
Common Integrals		<p>e.g. <math>\int kdx = kx + C</math></p> $\int x^a dx = \frac{x^{a+1}}{(a+1)} + C$ $\int e^{ax} dx = \frac{1}{a}e^{ax} + C \text{ etc}$
Taylor Series		<p>A representation of a function as an infinite sum of terms that are calculated from the values of the function's derivatives at a single point.</p> <p>For function <math>f</math> (with appropriate properties)</p> $f(x) = f(a) + f'(a)(x-a) + \frac{f''(a)}{2!}(x-a)^2 + \dots$ <p>Do you know the Taylor series for <math>e^x</math>?</p>

Table 1: You should be relatively familiar with the majority of the mathematical notation, results and techniques outlined in this table.

Prerequisite	Symbols	Examples & Notes
Partial derivative	$\frac{\partial f}{\partial x}(x, y),$ $\frac{\partial}{\partial x}(f(x, y))$	The partial derivative of a function of several variables is its derivative with respect to one of those variables, with the others held constant.
Gradient	$\nabla f, \nabla_x f$	The gradient is a multi-variable generalization of the derivative. Like the derivative, the gradient represents the slope of the tangent of the graph of the function. For function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ , the gradient <div style="text-align: center;"> <math display="block">\nabla f(\mathbf{x}) = \left( \frac{\partial f}{\partial x_1}, \dots, \frac{\partial f}{\partial x_n} \right)^T</math> </div> points in the direction of the greatest rate of increase of the function $f$ at point $\mathbf{x}$ .
Using Derivatives		How do you find the maxima, minima and points of inflection of a function? E.g. set <div style="text-align: center;"> <math display="block">f'(x) = 0</math> </div> How do you determine what type of point you have found? What if your function takes 2 or more inputs e.g. $f : \mathbb{R}^m \rightarrow \mathbb{R}$ ? How do you find the maxima, minima and other stationary points of this function?
<b>Statistics and Probability</b>		
Mean		The mean of some set of values $x_1, x_2 \dots x_N$ is <div style="text-align: center;"> <math display="block">\bar{x} = \frac{1}{N} \sum_{i=1}^N x_i</math> </div> If these values are the outcomes (samples) of $N$ repetitions of some experiment, this might be called the <i>sample mean</i> .
Probability distribution		A random variable $X$ with a finite number of outcomes $x_1, x_2, \dots, x_k$ can be represented by finite number of weights $p_1, p_2, \dots, p_k \in [0, 1]$ such that $\sum_{i=1}^k p_i = 1$ , where $p_i$ is the probability of $X$ taking value $x_i$ .

Table 1: You should be relatively familiar with the majority of the mathematical notation, results and techniques outlined in this table.

Prerequisite	Symbols	Examples & Notes
		<p>For a 1-dimensional random variable, <math>X</math>, which can take any value, <math>x</math>, in a continuous range, its distribution is typically described by a <i>probability density function</i>, <math>p(x)</math>. The probability that <math>a &lt; X &lt; b</math> for some <math>a</math> and <math>b</math> is given by:</p> $p(X \in [a, b]) = \int_a^b p(x) dx$
Cumulative Distribution Function (CDF)	$F(x), F_X(x)$	<p>For a random variable <math>X</math> the CDF specifies the probability that the random variable takes a value less than or equal to <math>x</math>, i.e.</p> $F_X(x) = p(X \leq x)$
Expectations	$\mathbf{E}$	<p>For a random variable <math>X</math> with a finite number of outcomes <math>x_1, x_2, \dots, x_k \in \mathbb{R}</math> occurring with respective probabilities <math>p_1, p_2, \dots, p_k</math>, the expectation of <math>X</math> is</p> $\mathbf{E}[X] = \sum_{i=1}^k p_i x_i \quad (1)$ <p>For a continuous random variable <math>a &lt; X &lt; b</math> with probability density function <math>p(x)</math> the expectation of <math>X</math> is</p> $\mathbf{E}[X] = \int_a^b x p(x) dx$ <p>We sometimes call <math>\mathbf{E}[X]</math> the mean of <math>X</math>.</p>
Variance		<p>The variance of a finite random variable <math>X</math>, with outcomes <math>x_1, x_2, \dots, x_k \in \mathbb{R}</math> and respective probabilities <math>p_1, p_2, \dots, p_k</math> is</p> $\text{Var}[X] = \sum_{i=1}^k p_i (x_i - \mu)^2$ <p>where <math>\mu</math> is the mean of <math>X</math> (i.e. <math>\mathbf{E}[X]</math>). What is the variance of a random variable <math>X</math> with probability density function <math>p(x)</math>?</p>

Table 1: You should be relatively familiar with the majority of the mathematical notation, results and techniques outlined in this table.

Prerequisite	Symbols	Examples & Notes
Normal Distribution		<p>The normal (Gaussian) distribution is a very commonly encountered continuous probability distribution, with probability density function:</p> $f(x \mu, \sigma^2) = \frac{1}{\sqrt{2\pi\sigma^2}} e^{-\frac{(x-\mu)^2}{2\sigma^2}} \quad (2)$ <p>where <math>\mu</math> is the mean, and <math>\sigma^2</math> is the variance.</p>
Linear Regression		<p>In the 1-dimensional case, we assume that we observe <math>n</math> pairs of values <math>(x_i, y_i) \in \mathbb{R}^2</math> for <math>i = 1, \dots, n</math>, and we want to find the straight line</p> $y = f(x) = \beta_0 + \beta_1 x$ <p>that <i>best fits</i> these values, i.e. <math>y_i \approx f(x_i)</math> for each <math>i</math>. The meaning of this can vary, but in the simplest case we want to <i>minimise the sum of squared errors</i>, i.e. find the parameters <math>(\beta_0, \beta_1)</math> that minimises</p> $\sum_{i=1}^n (y_i - f(x_i))^2 = \sum_{i=1}^n (y_i - \beta_0 - \beta_1 x_i)^2$ <p>How do we find the <i>best values</i> for <math>(\beta_0, \beta_1)</math>?  We might call the <math>x_i</math>s the input variables or independent variables, and the <math>y_i</math>s the output variables or dependent variables.  What changes if we want to find the best fit for this function</p> $y = \beta_0 + \beta_1 x + \beta_2 x^2 \quad ?$ <p>What if we allow each input variable to be a <math>k</math>-dimensional vector, i.e. <math>\mathbf{x}_i = (x_{i1}, x_{i2}, \dots, x_{ik})^T \in \mathbb{R}^k</math>, and want to ensure that</p> $y_i \approx \beta_0 + \sum_{j=1}^k \beta_j x_{ij} \quad ?$ <p>What assumptions are we typically making when we use these techniques or interpret their results?</p>

Table 1: You should be relatively familiar with the majority of the mathematical notation, results and techniques outlined in this table.