

Mathematics for Machine Learning

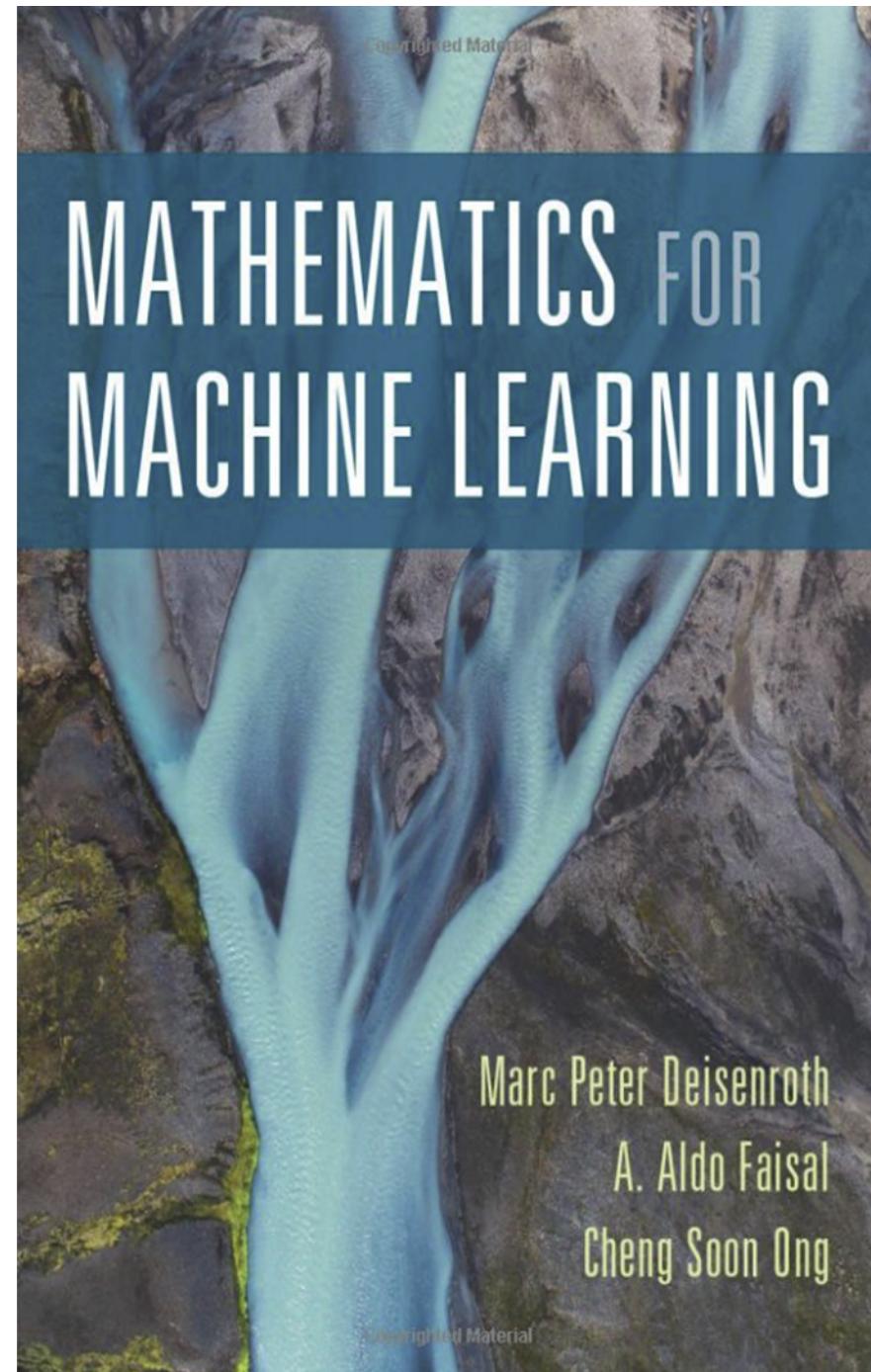
In this session ...

- Calculus
- Linear Algebra
- Probability
- Other methods

There is no ML in this session!

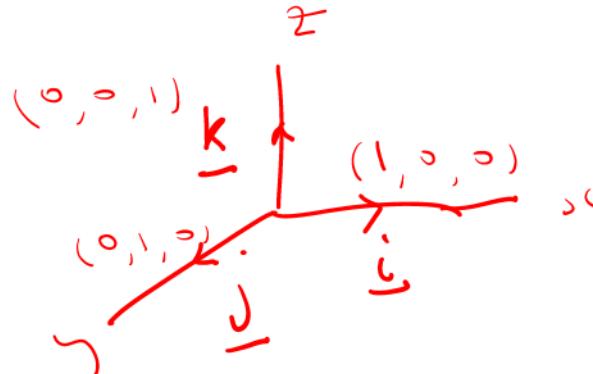
This module presents a number of themes associated with Machine Learning (ML). The lectures in ~~this~~ modules cover several advanced topics in ~~this~~ field.
~~these~~

This lecture only covers the advanced mathematical methods that are/to be employed by other tutors in their Machine Learning lectures. We will not be doing ML or applications of the maths to ML.



Extrema

Vector Calculus



1D

$$\frac{df}{dx} = \left(\frac{d}{dx} \right) f$$

If $f = f(x, y, z)$ is a function of three variables, then the *gradient of f* is

$$\nabla f(x, y, z) = f_x(x, y, z) \mathbf{i} + f_y(x, y, z) \mathbf{j} + f_z(x, y, z) \mathbf{k}.$$

The symbol ∇ , called the **del** operator, is a vector differential operator symbolised by

$$\begin{aligned} \mathbf{i} &= (1, 0, 0) & \mathbf{k} &= (0, 0, 1) & (\nabla = \mathbf{i} \frac{\partial}{\partial x} + \mathbf{j} \frac{\partial}{\partial y} + \mathbf{k} \frac{\partial}{\partial z}) & (f_x, f_y, f_z) \\ \mathbf{j} &= (0, 1, 0) \end{aligned}$$

It has properties similar to the operator d/dx . Standing alone it is meaningless; however, if it operates on $f(x, y, z)$ it produces the three-dimensional vector function given above.

In applications, the gradient $\nabla f(x, y, z)$ is sometimes denoted by grad f(x, y, z).

∇ : to find derivs in higher dim^s. Gives rate of increase in each direction

Also called Hamilton operator

Example 1: Calculate ∇f for $f(x, y, z) = x^2 + yz$.

$$f_x = 2x; \quad f_y = z; \quad f_z = y,$$

to give

partial
differentiation

$$\nabla f(x, y, z) = 2x\mathbf{i} + z\mathbf{j} + y\mathbf{k}$$

$$\nabla f = i \frac{\partial f}{\partial x} + j \frac{\partial f}{\partial y}$$

Example 2: If $f(x, y, z) = yz^3 - 2x^2$; find the gradient of f at the point $P(2, -3, 1)$.

$$f_x = -4x; f_y = z^3; f_z = 3yz^2,$$

to give

$$\nabla f = (f_x, f_y, f_z)$$

$$\nabla f(x, y, z) = -4x\mathbf{i} + z^3\mathbf{j} + 3yz^2\mathbf{k}$$

Hence, at $P(2, -3, 1)$

$$\begin{matrix} \uparrow & \uparrow & \leftarrow \\ x & y & z \end{matrix}$$

$$\nabla f(x, y, z) = -8\mathbf{i} + \mathbf{j} - 9\mathbf{k}.$$

$$\nabla = \left(\frac{\partial}{\partial x} + \frac{\partial}{\partial y} + \frac{\partial}{\partial z} \right)$$

scalar

$$\nabla \cdot \underline{f}$$

divergence (vector fields source at each point)

vector product

curl (infinitesimal circulation) =

$$\underline{f} = f_1 \underline{i} + f_2 \underline{j} + f_3 \underline{k}$$

$$\begin{vmatrix} \underline{i} & \underline{j} & \underline{k} \\ \frac{\partial}{\partial x} & \frac{\partial}{\partial y} & \frac{\partial}{\partial z} \\ f_1 & f_2 & f_3 \end{vmatrix}$$

Lagrange Multipliers

Let $f(x, y, z)$ and $g(x, y, z)$ have continuous first order partial derivatives, and suppose f has an extremum $f(x_0, y_0, z_0)$ when (x, y, z) is subject to the constraint $g(x, y, z) = 0$. If $\nabla g(x_0, y_0, z_0) \neq 0$ then there is a real number λ such that

$$\nabla f(x_0, y_0, z_0) = \lambda \nabla g(x_0, y_0, z_0)$$
*
*
Theorem

λ is called a *Lagrange multiplier*.



α

β

g is $\propto f$

$$\underbrace{R_{\underline{\underline{\cdot}}} \underline{\downarrow} - ?}_{g} = 0 \quad R_{\overline{\overline{\cdot}}} \underline{\downarrow} = ?$$

$\Delta . 1$

Example 3: Find the extrema of $f(x, y) = xy$; if (x, y) is restricted to the ellipse $4x^2 + y^2 = 4$. *constraint*

Solution: In this example the constraint is $g(x, y) = 4x^2 + y^2 - 4 = 0$.
Setting $\nabla f(x, y) = \lambda \nabla g(x, y)$, we obtain

$$\nabla f$$

$$y\mathbf{i} + x\mathbf{j} = \lambda(8x\mathbf{i} + 2y\mathbf{j})$$

$$\lambda \nabla g$$

Equating coefficients

$$y = 8\lambda x; \quad x = 2\lambda y,$$

together with $4x^2 + y^2 - 4 = 0$. A number of ways to solve, here we eliminate y .

$$x - 2\lambda y = 0$$

$$x - 2\lambda(8\lambda x) = 0.$$

So solve

$$x(1 - 16\lambda^2) = 0$$

$$x = 0; \quad 16\lambda^2 = 1$$

to get the values $x = 0$ or $\lambda = \pm 1/4$.

$x = 0$: using $4x^2 + y^2 - 4 = 0 \rightarrow y = \pm 2$.

Possible choices for extrema: $(0, 2)$; $(0, -2)$.

If $\lambda = \pm 1/4$, then $y = 8\lambda x = \pm 2x$. Substitute in $4x^2 + y^2 - 4 = 0$ to get $x = \pm \sqrt{2}/2$.

The corresponding y values are $y = \pm \sqrt{2}$.

This gives the following points

Exercise : Verify for
yourselves.

$$\begin{aligned} & (\sqrt{2}/2, \pm \sqrt{2}) \\ & (-\sqrt{2}/2, \pm \sqrt{2}) \end{aligned}$$



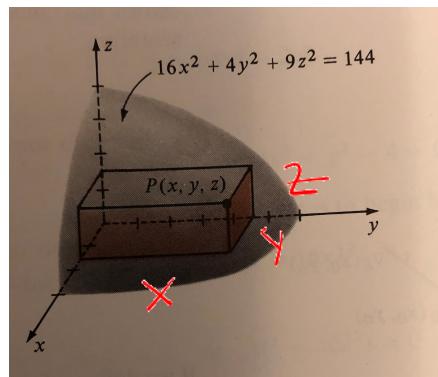
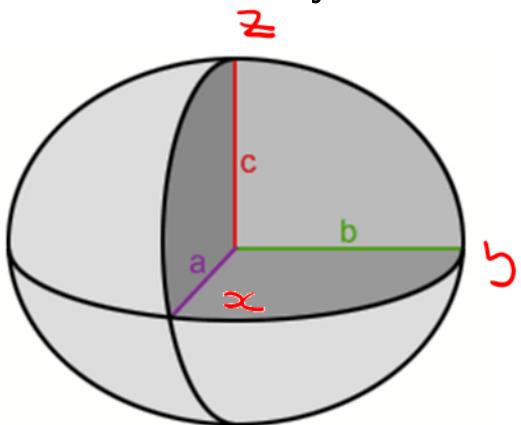
(x, y)	$f(x, y)$
$(0, 2)$	0
$(0, -2)$	0
$(\sqrt{2}/2, \sqrt{2})$	1 ✓
$(\sqrt{2}/2, -\sqrt{2})$	-1
$(-\sqrt{2}/2, \sqrt{2})$	-1
$(-\sqrt{2}/2, -\sqrt{2})$	1 ✓

So we see $f(x, y) = 1$ is the maximum value at 1 at either $(\sqrt{2}/2, \sqrt{2})$ or $(-\sqrt{2}/2, -\sqrt{2})$.

A minimum value of $f(x, y) = -1$ is attained at $(\sqrt{2}/2, -\sqrt{2})$ or $(-\sqrt{2}/2, \sqrt{2})$.

Example 4 (harder): Find the volume of the largest rectangular box with faces parallel to the coordinate planes that can be inscribed in the ellipsoid $16x^2 + 4y^2 + 9z^2 = 144$. *Constraint*

Start by considering a sketch of the problem.



f not given

$$V = 2x \times 2y \times 2z$$

$$\left(\frac{x}{a}\right)^2 + \left(\frac{y}{b}\right)^2 + \left(\frac{z}{c}\right)^2 = 1$$

We wish to maximise

$$V = f(x, y, z) = 8xyz.$$

subject to the constraint

$$g(x, y, z) = 16x^2 + 4y^2 + 9z^2 - 144 = 0$$

$$\left(\frac{x}{12a}\right)^2 + \left(\frac{y}{12b}\right)^2 + \left(\frac{z}{12c}\right)^2 = 1$$

$$\nabla f = 2\nabla g$$

- Consider

$$\nabla f(x, y, z) = \lambda \nabla g(x, y, z) \quad \text{←}$$

or

$$8yz\mathbf{i} + 8xz\mathbf{j} + 8xy\mathbf{k} = \lambda(32x\mathbf{i} + 8y\mathbf{j} + 18z\mathbf{k}) \quad \text{←}$$

- Together with $g(x, y, z) = 0$ gives the following system of 4 equations:

$$8yz = 32\lambda x, \quad 8xz = 8\lambda y, \quad 8xy = 18\lambda z, \quad 16x^2 + 4y^2 + 9z^2 - 144 = 0.$$

$$yz = 4\lambda x^2$$

$$xz = 2\lambda y^2$$

$$4xy = 9\lambda z^2$$

$$16x^2 + 4y^2 + 9z^2 = 144$$

- After some algebra we have $xyz = 12\lambda$

$$xyz = 4\lambda x^2 = 2\lambda y^2 = 9\lambda z^2 = 12\lambda$$

$$\frac{16x^2}{4\lambda} + \frac{4y^2}{2\lambda} + \frac{9z^2}{2\lambda} = 144$$

- The first equation gives

$$\lambda^2 = 12\lambda$$

$$9\lambda z^2 = 4(12\lambda)$$

so either $\lambda = 0$ or $x = \sqrt{3}$.

$$32\lambda(3 - x^2) = 0$$

$$12xyz = 144\lambda$$

$$xyz = 12\lambda$$

- Similarly, second gives

$$8\lambda(12 - y^2) = 0$$

hence $y = 2\sqrt{3}$.

- Finally multiplying the third equation $8xy = 18\lambda z$ by z

$$8xyz = 18\lambda z^2$$

and using the values found earlier

$$8(12\lambda) = 18\lambda z^2$$

which gives $z = 4/\sqrt{3}$.

Hence we obtain the desired volume

~~$$V = 8(\sqrt{3})(2\sqrt{3})(4/\sqrt{3}) = 64\sqrt{3}\text{unit}^3.$$~~



Constrained Extrema

Finding relative max. / min. of $f(x, y, z)$ s.t. constraints condition

$g(x, y, z) = 0$ consists of the formation of the Auxiliary Function,

A.F.
Lagrangian

$$G(x, y, z; \lambda) = f(x, y, z) + \lambda g(x, y, z)$$

λ - Lagrange multiplier

Subject to conditions

$$\underline{G_x = 0} ; \underline{G_y = 0} ; \underline{G_z = 0} ; \underline{G_\lambda = 0}$$

Ex: Optimise the function $f(x, y) = 4x^2 + 3xy + 6y^2$ s.t constraint
 $x+y=56$.

Lagrangian $G(x, y) = f + \lambda g = 4x^2 + 3xy + 6y^2 + \lambda(56 - x - y) \quad (*)$

$$\left. \begin{array}{l} G_x = 8x + 3y - \lambda = 0 \\ G_y = 3x + 12y - \lambda = 0 \\ G_\lambda = 56 - x - y = 0 \end{array} \right\} \text{3 eq's in 3 unknowns.}$$

$$x_0 = 36 ; y_0 = 20 ; \lambda = 348$$

Subst in $G(x, y)$. Exercise: ① Re-do this problem using

$$\nabla f = \lambda \nabla g$$

② Redo earlier examples by expressing as a Lagrangian.

$$\nabla f = 2\nabla g$$

$$\nabla(4x^2 + 3xy + 6y^2) = 2\nabla(x+y-5b)$$

$$(8x+3y, 3x+12y) = 2(1,1)$$

$$Eq^{\hat{1}}: 8x+3y=2$$

$$Eq^{\hat{2}}: 3x+12y=2$$

$$Eq^{\hat{3}}: x+y-5b=0$$

Essential Linear Algebra

Given $\mathbf{x} \in \mathbb{R}^N; \mathbf{y} \in \mathbb{R}^N$. Define the inner (scalar) product:

$\langle \cdot, \cdot \rangle$

$$\rightarrow \langle \mathbf{x}, \mathbf{y} \rangle = \mathbf{x}^\top \mathbf{y}$$

generalising

$$\langle \mathbf{x}, \mathbf{x} \rangle = \sum_{i=1}^N x_i y_i.$$

Inner products and norms are closely related

$$\text{Recall } \|\mathbf{x}\|_2 = \sqrt{\sum_{i=1}^N x_i^2}$$

ℓ_2 norm

$$\|\mathbf{x}\| = \sqrt{\mathbf{x}^\top \mathbf{x}} = (\underbrace{x_1, \dots, x_n}_{\mathbf{x}}) \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}$$

This gives

$\langle \mathbf{x}, \mathbf{x} \rangle$

$$\cos \theta = \frac{\mathbf{x}^\top \mathbf{y}}{\sqrt{\mathbf{x}^\top \mathbf{x} \mathbf{y}^\top \mathbf{y}}}$$

$$x_1^2 + \dots + x_n^2$$

$\langle \cdot, \cdot \rangle$

represents a mathematical op.

$$\langle f, g \rangle = \int_0^1 f g$$

Field

\mathbb{R}^n

$$\sum_{i=1}^m \lambda_i \underline{v}_i \text{ l.c}$$

Linear Dependence/Independence

Set of vectors $\underline{v}_1, \underline{v}_2, \underline{v}_3, \dots, \underline{v}_m \in V$

λ_i

1. The set $(\underline{v}_1, \underline{v}_2, \underline{v}_3, \dots, \underline{v}_m)$ is said to be linearly dependent (l.d)
= if \exists scalars $\{\lambda_1, \dots, \lambda_m\} \in \mathbb{F}$ not all zero s.t

$$\lambda_1 \underline{v}_1 + \lambda_2 \underline{v}_2 + \dots + \lambda_m \underline{v}_m = \underline{0}$$

The left hand side is called a LINEAR COMBINATION of $\underline{v}_1, \dots, \underline{v}_m$.

2. If $\{\underline{v}_1, \dots, \underline{v}_m\}$ are not l.d then they are linearly independent (l.i.). In this case an equation such as

If

$$\lambda_1 \underline{v}_1 + \lambda_2 \underline{v}_2 + \dots + \lambda_m \underline{v}_m = \underline{0}$$

$$\Rightarrow \lambda_1 = \lambda_2 = \dots = \lambda_m = 0$$

\mathbb{Z}^3

\mathbb{Z}

\mathbb{R}^n

Space of
vectors

Vector Space

\mathbb{R}

field.

Examples $\underline{v}_i, \omega_i \in \mathbb{R}^3$; $\lambda_i \in \mathbb{R}$

field

$$\textcircled{1} \quad \underline{v}_1 = \begin{pmatrix} 0 \\ -1 \\ 1 \end{pmatrix}; \quad \underline{v}_2 = \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix}; \quad \underline{v}_3 = \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix}$$

linear system

$$\sum_{i=1}^3 \lambda_i \underline{v}_i = \underline{0} = \lambda_1 \underline{v}_1 + \lambda_2 \underline{v}_2 + \lambda_3 \underline{v}_3 \quad \leftarrow$$

$$= \lambda_1 \begin{pmatrix} 0 \\ -1 \\ 1 \end{pmatrix} + \lambda_2 \begin{pmatrix} 1 \\ 0 \\ -1 \end{pmatrix} + \lambda_3 \begin{pmatrix} -1 \\ 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$$

augmented matrix

$$\left(\begin{array}{ccc|c} 0 & 1 & -1 & 0 \\ -1 & 0 & 1 & 0 \\ 1 & -1 & 0 & 0 \end{array} \right)$$

$$\lambda_2 = \lambda_3 = \alpha \quad (\text{say}) \text{ free var.}$$

$$\lambda_1 = \lambda_3 \rightarrow \underline{\alpha} \begin{pmatrix} 1 \\ 1 \\ 1 \end{pmatrix}$$

w.l.o.g put $\alpha = 1$

$$\therefore \boxed{\lambda_i = 1 \quad \forall i}$$

$$(\underline{v}_1 + \underline{v}_2 + \underline{v}_3 = \underline{0}) \uparrow$$

$$\text{e.g. } \underline{v}_1 = -\underline{v}_2 - \underline{v}_3 \leftarrow$$

$$\textcircled{2} \quad \underline{\omega}_1 = \begin{pmatrix} 2 \\ 2 \\ 1 \end{pmatrix}; \quad \underline{\omega}_2 = \begin{pmatrix} -1 \\ 6 \\ 5 \end{pmatrix}; \quad \underline{\omega}_3 = \begin{pmatrix} 1 \\ 0 \\ 0 \end{pmatrix} \quad \boxed{\sum_{i=1}^3 \lambda_i \underline{\omega}_i = 0}$$

$$\left(\begin{array}{ccc|c} 2 & -1 & 1 & 0 \\ 2 & 6 & 0 & 0 \\ 1 & 5 & 0 & 0 \end{array} \right) \sim \left(\begin{array}{ccc|c} 2 & -1 & 1 & 0 \\ 0 & 7 & -1 & 0 \\ 0 & 11 & -1 & 0 \end{array} \right) \begin{matrix} R_2 - R_1 \\ 2R_3 - R_2 \end{matrix} \sim \left(\begin{array}{ccc|c} 2 & -1 & 1 & 0 \\ 0 & 7 & -1 & 0 \\ 0 & 4 & 0 & 0 \end{array} \right)$$

$$4\lambda_1 = 0 \rightarrow \lambda_1 = 0; \quad 7\lambda_2 - \lambda_3 = 0 \rightarrow \lambda_3 = 0 \quad \therefore \lambda_1 = 0$$

$\underline{\omega}_1, \underline{\omega}_2, \underline{\omega}_3$ are linearly indep \therefore all λ' s = 0

$$\textcircled{3} \quad \text{A special set } \underline{e}_1 = (1, 0, 0)^T; \quad \underline{e}_2 = (0, 1, 0)^T; \quad \underline{e}_3 = (0, 0, 1)^T$$

$$\left(\begin{array}{ccc|c} 1 & 0 & 0 & 0 \\ 0 & 1 & 0 & 0 \\ 0 & 0 & 1 & 0 \end{array} \right) \Rightarrow \lambda_1 = 0 = \lambda_2 = \lambda_3 \quad \therefore \text{Lin. Indep}$$

$\underline{e}_i \in \mathbb{R}^3$ standard basis for \mathbb{R}^3

$$④ \quad \underline{u} = (1, 6, 9) ; \quad \underline{v} = (2, 4, 8) ; \quad \underline{w} = (0, 0, 0) \in \mathbb{R}^3$$

Write $\lambda_1 \underline{u} + \lambda_2 \underline{v} + \lambda_3 \underline{w} = \underline{0}$

choose $\lambda_1 = 0 = \lambda_2 \quad \lambda_3 = 6$

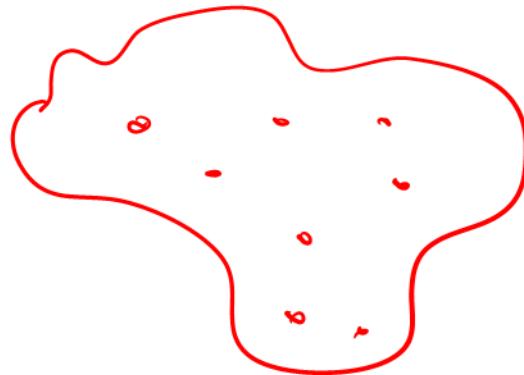
Then $0 \times \underline{u} + 0 \times \underline{v} + 6 \times \underline{w} = \underline{0} \quad \text{i.e. } \lambda_i \text{'s not all zero}$

$\therefore \underline{u}, \underline{v}, \underline{w}$ lin. dep.

Moral: If n vectors in \mathbb{R}^n and contain zero vector
then set is linearly dependent.

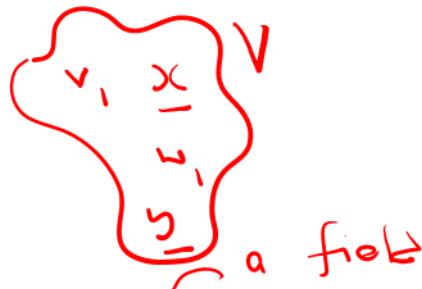
$\underline{v}_1, \underline{v}_2, \dots, \underline{v}_n \in \mathbb{R}^n$ one of the vectors is $\underline{0}$

Vector Spaces



We are interested in the non-abstract treatment of this subject. Throughout we will use the term *field* denoted F to refer to a set of scalars.

Definition:



A *vector space* V over F is a set with a binary operation called *Vector Addition*, denoted $V \times V \rightarrow V$, $(x, y) \mapsto x + y$, and a function $F \times V \rightarrow V$, $(c, x) \mapsto cx$, ($c \in F$), called *Scalar Multiplication* such that the following eight rules hold:

Don't worry.

+1) + is associative $\forall x, y, z \in V \quad (x + y) + z = x + (y + z)$

+2) + is commutative $\forall x, y \in V \quad x + y = y + x$

+3) + has a neutral $\exists 0 \in V \quad \forall x \in V \quad x + 0 = x$

+4) + has inverse $\forall x \in V \quad \exists y \in V \quad x + y = 0 \quad y \text{ is denoted } (-x)$

.1) · is associative $\forall c, d \in F, \quad \forall x \in V \quad c(dx) = (cd)x$

·2) · is commutative

$$\forall x, y \in V \quad xy = yx$$

·3) · has a neutral

$$\forall x \in V \quad 1 \cdot x = x \quad (1 \neq 0)$$

·4) · has an inverse

$$\forall x \in V \quad (x \neq 0) \Rightarrow \quad (\exists y \in V \quad xy = 1) \quad y$$

is denoted (x^{-1})

+·1) Right distributive $\forall c \in F \quad \forall x, y \in V \quad c(x + y) = cx + cy$ scalar multiplication is distributive over vector addition

+·2) Left distributive

$$(c + d)x = cx + dx$$

Remarks:

field

\mathbb{R} vector space

\mathbb{R}^n

1. Elements of F are called SCALARS and elements of V are called VECTORS.

2. If $F = \mathbb{R}$ we say V is a real vector space

$= \mathbb{C}$ complex

$= \mathbb{Q}$ rational

3. At this stage we have

2	+	's	addition
2	·	's	multiplication
2	0	's	neutrals
2	-	's	inverses

(and things usually get a lot worse) even in ^{1st} year
university algebra

4. The axioms can be used to deduce various rules.

Examples:

1. Let $m, n \in \mathbb{N}^+$ then ${}^m F^n$ is a vector space over F with respect to operations of matrix addition and scalar multiplication.

$$\begin{pmatrix} a & b & c \\ d & e & f \end{pmatrix} + \begin{pmatrix} A & B & C \\ D & E & F \end{pmatrix} = \begin{pmatrix} a+A & b+B & c+C \\ d+D & e+E & f+F \end{pmatrix}$$

$$M_1, M_2 \in {}^2 F^3; \lambda \in F$$

$$\lambda M_1 = \begin{pmatrix} \lambda a & \lambda b & \lambda c \\ \lambda d & \lambda e & \lambda f \end{pmatrix} \in {}^2 F^3$$

2. Let $V = \mathbb{R}[x]$ denote the set of all polynomials

$$\sum_{n=0}^N a_n x^n \quad n \in \mathbb{N}, \quad a_i (i = 1, \dots, N) \in \mathbb{R}$$

Then V is a vector space over \mathbb{R} w.r.t. addition of polynomials and multiplication by a constant.

$$\sum_{n=0}^N a_n x^n + \sum_{n=0}^N b_n x^n = \sum_{n=0}^N (a_n + b_n) x^n = \sum_{n=0}^N c_n x^n$$

$$\lambda \sum_{n=0}^N a_n x^n$$

3. Let F be an arbitrary field and V the set of all n dimensional vectors with vector addition

$$(a_1, \dots, a_n) + (b_1, \dots, b_n) = (a_1 + b_1, \dots, a_n + b_n)$$

and scalar multiplication

$$k(a_1, a_2, \dots, a_n) = (ka_1, ka_2, \dots, ka_n)$$

where $a_i, b_i, k \in F$. Then V is a vector space over F .

3.1 Subspaces

Special subset of a
vector space

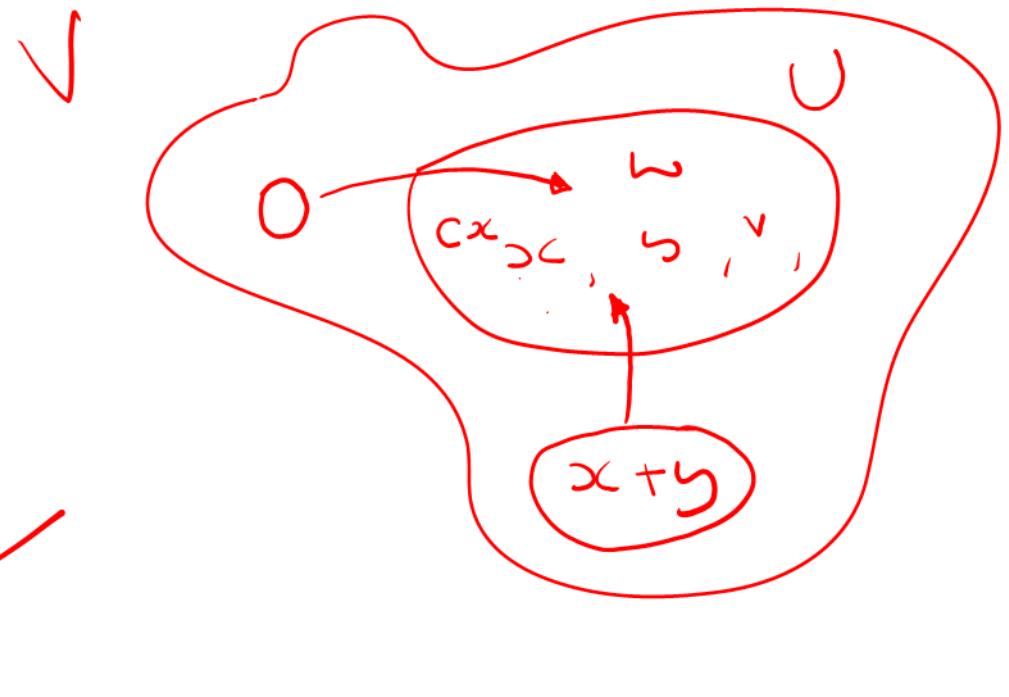
Definition:

A subspace U of a vector space V (over F) is a subset , i.e. $U \subset V$ such that

(i) $0 \in U$ ✓

(ii) $\forall x, y \in U \quad x + y \in U$ ✓

(iii) $\forall c \in F \quad \forall x \in U \quad cx \in U$ ✓



Note:

U is then a vector space with vector addition and scalar multiplication calculated in V .

✓ $U \times U \rightarrow U$ is a binary operation on U

$$(x, y) \mapsto x + y$$

✓ $F \times U \rightarrow U$, is a function

$$(c, x) \mapsto cx$$

The eight rules obviously hold because they are satisfied in the larger space and will automatically be satisfied in every subspace.

Examples:

$$\mathbb{R}^3$$

(1) Let $V = \text{set of all } 3 \times 3 \text{ matrices}$ and $U, W \subset V$ such that

$U = \text{set of lower triangular matrices}$

$$\begin{pmatrix} a & 0 & 0 \\ d & b & 0 \\ e & f & c \end{pmatrix}$$

$W = \text{set of symmetric matrices.}$

$$\begin{pmatrix} A & \alpha & \beta \\ \alpha & B & \gamma \\ \beta & \gamma & C \end{pmatrix}$$

Suppose $A, B \in U$; $C, D \in W$, then $A + B \in U$ and $cA \in U$ where $c \in \mathbb{R}$. The sums $A + B$ and cA inherit properties of A and B ; and

similarly $C + D \in W$ and $kC \in W$ where $k \in \mathbb{R}$. 0 is in both spaces.
Hence U and W are subspaces of V .



(2) Consider the vector space $V = \mathbb{R}^2$ over $F = \mathbb{R}$. $U = \{(x, y) / x, y \in [0, \infty)\}$ is the subset consisting of vectors whose components are ≥ 0 , i.e. the first quadrant, all co-ordinates $(x, y) \geq 0$. Now $\underline{0} \in U$ and U is closed under vector addition $(x + y) \in U$. What about closure under scalar multiplication? Suppose $\underline{u} = (1, 1) \in U$ and $c = \underline{-1} \in \mathbb{R}$, then $cu = (-1, -1) \notin U$ ∵ scalar multiplication fails. Hence \underline{U} is not a subspace of \mathbb{R}^2 .

Now suppose $W = \{(-a, -b) / a, b \in [0, \infty)\} \subset V$, i.e. the third quadrant. Define a new subset of V such that $S = \underline{U + W} \subset V$. We see that for any vector \underline{w} in S , $k\underline{w} \in S$ where $k \in \mathbb{R}$, so closure under scalar multiplication. However now addition fails because $(2, 1) + (-1, -3) = (1, -2)$



which is in neither quadrant.

So the smallest subspace containing the 1st quadrant is the whole space \mathbb{R}^2 .

Vector Spaces

Summary

$\underline{a}, \underline{b}, \underline{c}, \dots, \underline{x}, \underline{y}, \underline{z}$ vectors

$\alpha, \beta, \gamma, \lambda, \dots$ scalars

Vector Space: set V on which are defined ops

① scalar mult

② vector add

scalar multiplication: $\underline{x} \in V$ $\lambda \in F$ $\underline{\underline{w}} = \lambda \underline{x} \in V$

field

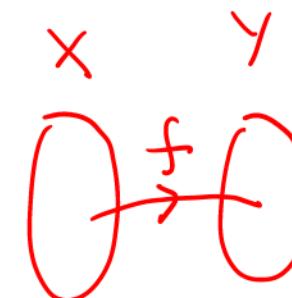
vector addition: $\underline{x}, \underline{y} \in V$
 $\underline{\underline{v}} = \underline{x} + \underline{y} \in V$

8 properties satisfied:

comm.; associative; \exists neutral element (e.g. $0(+)$ | $1(\times)$);
inverse, distributive

Linear Mappings:

V, W vector spaces over some field F .



Linear mapping f such that

$$f: \underline{V} \rightarrow \underline{W}$$

$$f: X \rightarrow Y$$

is a mapping satisfying

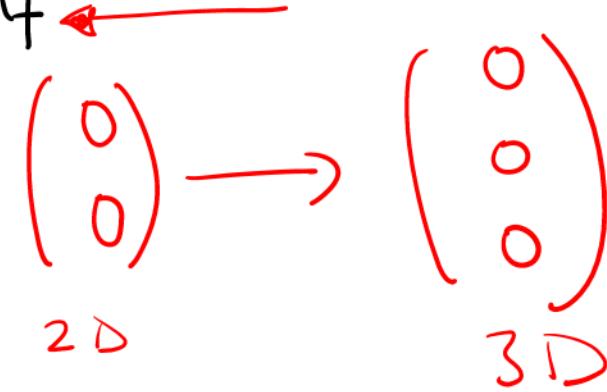
② $\forall \underline{v}, \underline{w} \in V$
$$f(\underline{v} + \underline{w}) = f(\underline{v}) + f(\underline{w})$$

① any scalar $\lambda \in F$ and any $\underline{v} \in V$

$$f(\lambda \underline{v}) = \lambda f(\underline{v})$$

Linear map also preserves $\underline{0}$ onto itself

f is called a linear transf[^]



V — departure space of f
 W — arrival space of f

Example 1: $f: \mathbb{R}^2 \rightarrow \mathbb{R}^3$

i)

$$f(2\mathbf{v}) = f\left(2\begin{pmatrix} x \\ y \end{pmatrix}\right) = f\left(\begin{pmatrix} 2x \\ 2y \end{pmatrix}\right) = \begin{pmatrix} 2 \cdot 2x + 2y \\ 0 \\ 2y - 2x \end{pmatrix} = 2 \begin{pmatrix} x + y \\ 0 \\ y - x \end{pmatrix} = 2f\begin{pmatrix} x \\ y \end{pmatrix}$$

ii)

$$\begin{aligned} f\left(\underline{v} + \underline{w}\right) &= f\left(\underbrace{\begin{pmatrix} x_1 \\ y_1 \end{pmatrix}}_{\underline{v}} + \underbrace{\begin{pmatrix} x_2 \\ y_2 \end{pmatrix}}_{\underline{w}}\right) = \begin{pmatrix} 2x_1 + y_1 + 2x_2 + y_2 \\ 0 + 0 \\ y_1 - x_1 + y_2 - x_2 \end{pmatrix} \\ &= \begin{pmatrix} 2x_1 + y_1 \\ 0 \\ y_1 - x_1 \end{pmatrix} + \begin{pmatrix} 2x_2 + y_2 \\ 0 \\ y_2 - x_2 \end{pmatrix} = f\begin{pmatrix} x_1 \\ y_1 \end{pmatrix} + f\begin{pmatrix} x_2 \\ y_2 \end{pmatrix} = f(\underline{v}) + f(\underline{w}) \end{aligned}$$

iii)

$$f\begin{pmatrix} 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 2 \cdot 0 + 0 \\ 0 - 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 0 \end{pmatrix} \quad \therefore \quad 0 \text{ preserved}$$

$\underline{0} \in \mathbb{R}^2$

$\underline{0} \in \mathbb{R}^3$

$$\text{i.e. } f\left(\underline{0}_{\mathbb{R}^2}\right) := \underline{0}_{\mathbb{R}^3}$$

Example 2: $g: \mathbb{R}^2 \rightarrow \mathbb{R}^3$ $g\begin{pmatrix} x \\ y \end{pmatrix} = \begin{pmatrix} x+3y \\ 1 \\ y \end{pmatrix}$!!

$g\begin{pmatrix} 0 \\ 0 \end{pmatrix} = \begin{pmatrix} 0 \\ 1 \\ 0 \end{pmatrix} \neq \begin{pmatrix} 0 \\ 0 \\ 0 \end{pmatrix}$ $\underline{0}$ is not preserved $\therefore f$ not a linear mapping.

Linear maps can be represented by matrices.

Consider standard basis for \mathbb{R}^2 $\left\{ e_1 = \begin{pmatrix} 1 \\ 0 \end{pmatrix}, e_2 = \begin{pmatrix} 0 \\ 1 \end{pmatrix} \right\}$

and earlier $f(x, y) = (2x+y, 0, y-x)$

$$f(e_1) = f\begin{pmatrix} 1 \\ 0 \end{pmatrix} = \begin{pmatrix} 2 \\ 0 \\ -1 \end{pmatrix} \quad f(e_2) = f\begin{pmatrix} 0 \\ 1 \end{pmatrix} = \begin{pmatrix} 1 \\ 0 \\ 1 \end{pmatrix}$$

$$\boxed{\begin{pmatrix} 2 & 1 \\ 0 & 0 \\ -1 & 1 \end{pmatrix}}$$

This is the matrix rep. of f wrt the standard basis of \mathbb{R}^2 .

e_i are lin. indep and orthogonal

All vectors in \mathbb{R}^n can be written uniquely as a lin. comb.

Matrix Diagonalisation

$$f = \lambda^2 \quad (x_1 + y_1)^2 + x_1^2 + y_1^2$$

Given $A \in \mathbb{R}^n \times \mathbb{R}^n$. $|A - \lambda I| = 0$ characteristic eqn

Form the matrix \underline{P} from the eigenvectors of A . Then A can be factorised as $P^{-1} = P^T$

$$A = P D P^{-1} \quad P D P^T$$

where $D = \begin{pmatrix} \lambda_1 & & 0 \\ & \ddots & \\ 0 & \cdots & \lambda_n \end{pmatrix}$ check: $D = P^{-1} A P$

We wish to calculate A^n .

$$A^n = (\cancel{P} D \cancel{P}^{-1}) (\cancel{P} D \cancel{P}^{-1}) (\cancel{P} D \cancel{P}^{-1}) \dots \dots \dots (\cancel{P} D \cancel{P}^{-1})$$

n lots of D

$$= P D \dots D P^{-1} = (P D^n P^{-1})$$

Example :

Given $M = \begin{pmatrix} 0 & -1 \\ -5 & 4 \end{pmatrix}$ with $\lambda_1 = 5$ $\underline{v}_1 = \begin{pmatrix} -1 \\ 5 \end{pmatrix}$
 $\lambda_2 = -1$ $\underline{v}_2 = \begin{pmatrix} 1 \\ 1 \end{pmatrix}$

$$P = \begin{pmatrix} -1 & 1 \\ 5 & 1 \end{pmatrix} \quad P^{-1} = -\frac{1}{6} \begin{pmatrix} 1 & -1 \\ -5 & -1 \end{pmatrix} \quad D = \begin{pmatrix} 5^{\lambda_1} & 0 \\ 0 & (-1)^{\lambda_2} \end{pmatrix}$$

$$\hat{M} = (P D P^{-1})^n = (P D^n P^{-1})$$

$$= \left(-\frac{1}{6} \begin{pmatrix} -1 & 1 \\ 5 & 1 \end{pmatrix} \begin{pmatrix} 5^n & 0 \\ 0 & (-1)^n \end{pmatrix} \begin{pmatrix} 1 & -1 \\ -5 & -1 \end{pmatrix} \right) P^{-1}$$

$$= -\frac{1}{6} \begin{pmatrix} -1 & 1 \\ 5 & 1 \end{pmatrix} \begin{pmatrix} 5^n - (-5)^n \\ -5(-1)^n - (-1)^n \end{pmatrix} = -\frac{1}{6} \begin{pmatrix} -5^n - 5(-1)^n & 5^n - (-1)^n \\ 5^{n+1} - 5(-1)^n & -5(5^n) - (-1)^n \end{pmatrix}$$

$$\hat{M} = \frac{1}{6} \begin{pmatrix} 5^n + 5(-1)^n & (-1)^n - 5^n \\ 5(-1)^n - 5^{n+1} & 5(5^n) + (-1)^n \end{pmatrix}$$

Try this by
induction

Principal Component Analysis

PCA is a technique that takes a dataset \mathbf{X} (consisting of a set of N tuples in p -dimensional space) and finds the directions along which the tuples line up best.

$$\begin{matrix} \text{N rows of } p \text{ col vectors} \\ (\mathbf{x}_1, \mathbf{x}_2, \dots, \mathbf{x}_p) \end{matrix} \left(\begin{matrix} \alpha_1 \\ \vdots \\ \alpha_p \end{matrix} \right)$$

In short, a dimension reduction technique:

1. Taking linear combinations of the p original variables. Say $\mathbf{Z} = \mathbf{X}\alpha_j$.

Note: α_j is a vector.

$$j=1, \dots$$

2. Each of these linear combinations explains as much as possible of the variance in the data:

$$\max_{\alpha_j} V_j = \mathbb{V}(\mathbf{X}\alpha_j) \quad \text{with } V_j \geq V_k \text{ if } j < k.$$

3. Each linear combination is uncorrelated with all the others: $\alpha_j^T \alpha_k = 0 \forall j \neq k.$
also think in
4. Each linear combination has unit length: $\alpha_j^T \alpha_j = 1 \forall j.$ vector norm
for
- $\rightarrow \|\underline{x}\| = \sqrt{\underline{x}^T \underline{x}}$

PCA Derivation

Similar to least squares.

Denote the variance-covariance matrix of our dataset \mathbf{X} as Σ .

Let's focus on the first linear combination with maximum variance. In order to find this, we need to solve the following constrained optimization problem:

$$G = f + \lambda \mathcal{S} = 0$$

Take the first derivative and set to zero:

Note that we want to do maximisation

$$\max_{\alpha_1} \mathbb{V}(\mathbf{X}\alpha_1) \text{ s.t. } \alpha_1^\top \alpha_1 = 1$$

Constrained constraint

The Lagrangian is

$$\begin{aligned} L &= (\mathbf{X}\alpha_1)^\top \mathbf{X}\alpha_1 - \lambda_1 (\alpha_1^\top \alpha_1 - 1) & \mathbf{X}^\top \mathbf{X} = \Sigma \\ &= \alpha_1^\top \Sigma \alpha_1 - \lambda_1 (\alpha_1^\top \alpha_1 - 1) \\ &\quad \text{Treat as } \alpha_1^2 = \alpha_1^\top \alpha_1 = \sum_i \alpha_i^2 - 1 \end{aligned}$$

done with
Ses.

"Power Method"

Take first derivative and set to zero:

$$\frac{\partial L}{\partial \alpha_1} = \Sigma \alpha_1 - 2 \alpha_1 = 0 \Rightarrow (\Sigma - 2 \alpha_1) \alpha_1 = 0 \Rightarrow \Sigma = 2 \alpha_1$$

Note we want to maximise $\alpha_1^T \Sigma \alpha_1 = \lambda_1 \alpha_1^T \alpha_1 = \lambda_1$.

Solution: λ_1 = largest eigenvalue, α_1 = corresponding eigenvector.

To find the second linear combination, we need to take a further constraint into account

$$\max_{\alpha_2} \mathbb{V}(X\alpha_2) \quad \text{s.t.} \quad \begin{cases} \alpha_2^T \alpha_2 = 1 \\ \alpha_2^T \alpha_1 = 0 \end{cases}$$

earlier def⁻s

The Lagrangian is

$$L = \alpha_2^T \Sigma \alpha_2 - \lambda_2 (\alpha_2^T \alpha_2 - 1) - \theta_2 \alpha_2^T \alpha_1$$

behaves like $\sum \alpha_2^2 - \lambda_2 (\alpha_2^2 - 1) - \theta_2 \alpha_2^2$

Or use "product rule".

Optimising yields

$$\frac{\partial L}{\partial \alpha_2} = 2 \sum \alpha_2 - 2 \lambda_2 \alpha_2 - 2 \theta_2 \alpha_2$$

$$\frac{\partial L}{\partial \alpha_2} = \sum \alpha_2 - \lambda_2 \alpha_2 - \theta_2 \alpha_1 = 0$$

$$\text{Left-multiplying } \alpha_1^T : \underbrace{\alpha_1^T \sum \alpha_2}_{=0} - \underbrace{\lambda_2 \alpha_1^T \alpha_2}_{=0} - \underbrace{\theta_2 \alpha_1^T \alpha_1}_{=1} = 0 \rightarrow \theta_2 = 0$$

Therefore,

$$(\Sigma - \lambda_2 I_p) \alpha_2 = 0 \quad \xrightarrow{\text{to factorise}} \quad \Sigma = \lambda_2 I_p$$

$\Rightarrow \lambda_2$ = second largest eigenvalue, α_2 = corresponding eigenvector.

$$\left\{ L = \lambda_2^T \bar{\Sigma} \alpha_2 = \alpha_2^T \lambda_2 I_p \alpha_2 = \lambda_2 \right\}$$

Differentiating wrt vectors \Rightarrow
in module 2 More later this
module

Some background calculus *for earlier work on extremes.*

If f is a function of two variables x, y such that $f(x, y) \begin{cases} \leq \\ \geq \end{cases} f(x_0, y_0)$ at all points (x, y) inside some sufficiently small circle, centre (x_0, y_0) , then $f(x, y)$ has a local $\begin{cases} \text{maximum} \\ \text{minimum} \end{cases}$ at (x_0, y_0) .

We define a critical point (x_0, y_0) at which either

$$\underbrace{f_x(x_0, y_0) = 0}_{\text{or}} \text{ or } \underbrace{f_y(x_0, y_0) = 0}$$

or

$f_x(x_0, y_0)$ or $f_y(x_0, y_0)$ do not exist.

For example consider $f(x, y) = 1 - x^2 + y^2$. Then

$$f_x = -2x \text{ and } f_y = 2y.$$

So if $f_x = f_y = 0$ then $x = 0, y = 0$. But the point $(0, 0)$ corresponds neither to a local maxima nor a minima. Such a point is called a *saddle point*.

Testing for Local Extrema of a function of Two Variables

Suppose $f(x, y)$ has continuous partial derivatives of order one and two within some disc with centre (x_0, y_0) . Let $f_x(x_0, y_0) = 0$ and $f_y(x_0, y_0) = 0$ and not all the partial derivatives of f are zero at (x_0, y_0) . Define

$$\begin{aligned}\Delta &= \begin{vmatrix} f_{xx}(x_0, y_0) & f_{xy}(x_0, y_0) \\ f_{xy}(x_0, y_0) & f_{yy}(x_0, y_0) \end{vmatrix} \quad \text{C.P.S.} \\ &= f_{xx}f_{yy}|_{(x_0, y_0)} - f_{xy}^2|_{(x_0, y_0)}\end{aligned}$$

Then we have the following

- $\Delta > 0$ { Local maximum at (x_0, y_0) if $f_{xx}(x_0, y_0) < 0$
 Local minimum at (x_0, y_0) if $f_{xx}(x_0, y_0) > 0$
- $\Delta < 0$ The function has a saddle point at (x_0, y_0)
- $\Delta = 0$ More work required to determine the nature of (x_0, y_0)

+ test fails.

T.J.E.

Probability

Maximum Likelihood Estimators (MLE)

- Suppose we draw x_1, x_2, \dots, x_n iid with density $p(\cdot; \theta)$.
 \downarrow parameter
 $=$ vars. x_i
- The probability of drawing a number in $[x_1, x_1 + dx]$ is $p(x_1; \theta) dx$
 $\underbrace{}$
- Similarly the probability of drawing a number from x_2 to $x_2 + dx$ is $p(x_2; \theta) dx$ and so forth.
 $x_i \rightarrow x_i + dx \quad p(x_i; \theta)$
- The probability of all these events occurring is

$$\rightarrow p(x_1; \theta) \cdots p(x_n; \theta) dx^n = L(\theta; x_1, x_2, \dots, x_n) dx^n$$

multⁿ rule

$$\prod_{i=1}^n p(x_i; \theta)$$

$$dx_i \in \mathbb{R}$$

parameter θ

vector \underline{x}

- The likelihood function is

$$L(\theta; \underbrace{x_1, x_2, \dots, x_n}_{\text{vector}}) = L(\theta, \mathbf{x}),$$

which we think of as a function of θ .

- We ignore dx^n as this term is 'constant'.
- Maximum likelihood estimation works on the principle that the 'best' estimate of the parameters θ is given by maximising $L(\theta; \mathbf{x})$. →

- As likelihoods are positive, and given independence are of the form

$$L(\theta, \mathbf{x}) = \prod_{i=1}^n p(x_i; \theta) = p(x_1; \theta) \cdots p(x_n; \theta)$$

it is often easier working with the log-likelihood maths is easier.

$$\rightarrow \ell(\theta, \mathbf{x}) = \log L(\theta, \mathbf{x}) = \sum_{i=1}^n \log p(x_i; \theta)$$

if L then $\ell = \log L$

- SCORE Function $\frac{\partial \ell}{\partial \theta} = S(\theta, x)$

- If ℓ is differentiable in θ , then the maximum is when the Score function (derivative of ℓ) is zero,

$$S(\theta, x) = \frac{d\ell}{d\theta} = 0.$$

- Usually, ℓ cannot have local minima or saddle points, so there is a maximum.

Theorem from Calculus:

Suppose we pick n i.i.d random variables x_i such that $X_i \sim N(\mu, \sigma^2)$. Then the probability of drawing a number in the interval

$$\underline{x_i \text{ to } x_i + dx} = p(x_i; \underbrace{\mu, \sigma}_2) \quad \text{for } i = 1, \dots, n$$

where $p(x_i; \mu, \sigma^2)$ is the pdf for the normal distribution.

$$\begin{aligned} & L(x_1, x_2, \dots, x_n; \mu, \sigma) \\ &= \underbrace{p(x_1; \mu, \sigma)}_{\text{can omit}} \dots \times \underbrace{p(x_n; \mu, \sigma)}_{\text{can omit}} (dx)^n. \end{aligned}$$

Recall, we ignore the term $(dx)^n$ as it is a constant. Hence

$$L(x_1, x_2, \dots, x_n; \mu, \sigma) := \prod_{i=1}^n p(x_i; \mu, \sigma)$$

We now choose the parameters μ & σ in a way that maximises the likelihood of observing this. This is done by setting the **Score function** equal to zero

2 eq's ; 2 unknowns $S(\mu, \mathbf{x}) = \frac{\partial \ell}{\partial \mu} = 0, \quad S(\sigma, \mathbf{x}) = \frac{\partial \ell}{\partial \sigma} = 0.$

We have already seen (in ~~the~~ calculus course), in principle that these values can be local maxima/minima or saddle points.

However due to the nature of this problem, i.e. L is a likelihood function ensures that L can only have a maximum.

The problem can be simplified by noting that L attains a maximum iff $\ell = \log L$ has a maximum. In this case ℓ is called the *log-Likelihood function*.

So we have

$$\log L = \ell(x_1, x_2, \dots, x_n; \mu, \sigma) = \sum_{i=1}^n \log p(x_i; \mu, \sigma)$$

and our problem becomes solving

$$\frac{\partial \ell}{\partial \mu} = 0, \quad \frac{\partial \ell}{\partial \sigma} = 0.$$

Example 1:

Suppose we draw n positive numbers X_i ($i = 1, \dots, n$) from i.i.d samples with distribution

$$f(x) = \begin{cases} 2\lambda^3 \sqrt{\frac{x}{\pi}} e^{-\lambda^2 x} & \text{for } x \geq 0 \\ 0 & \text{for } x < 0 \end{cases}$$

\wedge estimate

where the unknown parameter $\lambda > 0$.

How can we obtain a maximum Likelihood estimator $\hat{\lambda}$ for λ ?

$$\left(\frac{2\lambda^3}{\sqrt{\pi}} \right)^n \prod_{i=1}^n e^{-\lambda^2 x_i} = \left(\frac{2\lambda^3}{\sqrt{\pi}} \right)^n (x_1 x_2 \cdots x_n) e^{-\lambda^2 (x_1 + x_2 + \cdots + x_n)}$$

~~MATX~~
~~STATX~~
G. Hoefer

The Likelihood function is

$$L(x_1, x_2, \dots, x_n; \lambda) = \prod_{i=1}^n f(x_i; \lambda)$$
$$= \frac{2^n \lambda^{3n}}{\pi^{n/2}} \sqrt{(x_1 x_2 \dots x_n)} e^{-\lambda^2(x_1 + x_2 + \dots + x_n)},$$

Text - f
basic
algebra

and for $\lambda \rightarrow 0$ & ∞ , $L \rightarrow 0$, so there must be at least one maximum.

We use $\ell = \log L$ (much easier), as ℓ has maximum iff L has a maximum.

$$\ell = \log \frac{2^n}{\pi^{n/2}} + 3n \log \lambda + \frac{1}{2} \log (x_1 x_2 \dots x_n) - \lambda^2 (x_1 + x_2 + \dots + x_n).$$

$$\log \left(\frac{2^n}{\pi^{n/2}} \cdot \lambda^{3n} \cdot (x_1 x_2 \dots x_n)^{1/2} \cdot e^{-\lambda^2(x_1 + \dots + x_n)} \right) := \text{sum of logs}$$

$$\text{Score } f \sim \frac{\partial \ell}{\partial \lambda} = 0$$

So

$$\begin{aligned}\frac{\partial \ell}{\partial \lambda} &= \frac{3n}{\lambda} - 2\lambda(x_1 + x_2 + \dots + x_n) = 0 \\ \Rightarrow & 2\hat{\lambda}(x_1 + x_2 + \dots + x_n) \\ &= \frac{3n}{\hat{\lambda}}\end{aligned}$$

Hence

$$\hat{\lambda}^2 = \frac{3n}{2(x_1 + x_2 + \dots + x_n)}$$

and

$$\hat{\lambda} = \sqrt{\frac{3n}{2(x_1 + x_2 + \dots + x_n)}}.$$

The Maximum Likelihood estimator method is a popular method due to its simplicity when the form of the probability density function is known.

Example 2: Suppose $X_1, \dots, X_n \sim N(\mu, \sigma^2)$ iid. What is the MLE of (μ, σ^2) ?

Repeat for $\frac{\partial l}{\partial \sigma} = 0$

In this case,

pdf

$$p(x_i; \mu, \sigma^2) = (2\pi\sigma^2)^{-1/2} \exp\left(-\frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

so

Likelihood

$$L(x_i; \mu, \sigma^2) = (2\pi\sigma^2)^{-n/2} \exp\left(-\sum_i \frac{(x_i - \mu)^2}{2\sigma^2}\right)$$

log Likelihood

$$l(x_i; \mu, \sigma^2) = -\frac{n}{2} \log(2\pi\sigma^2) - \sum_i \frac{(x_i - \mu)^2}{2\sigma^2}$$

Note that for μ , this is just least squares.

$$\frac{\partial l}{\partial \mu} = \cancel{-\frac{1}{2\sigma^2}} \sum_i (x_i - \mu) \times \cancel{-2} = 0 = \sum_{i=1}^n x_i - \hat{\mu} \sum_i = 0$$

$$n\mu = \sum x_i; \quad \hat{\mu} = \underbrace{\frac{1}{n} \sum x_i}_{\text{ave exp. dist}}$$

Example 3: Suppose $X_1, \dots, X_n \sim \exp(-\lambda x_i)$ iid. What is the MLE of λ ?

In this case,

$$p(x_i; \lambda) = \lambda \exp(-\lambda x_i)$$

so

$$L(x_i; \lambda) = \lambda^n \exp\left(-\lambda \sum_i x_i\right)$$

$$l(x_i; \lambda) = n \log(\lambda) - \lambda \sum_i x_i$$

hence the score function

$$S(x_i; \lambda) = \frac{n}{\lambda} - \sum_i x_i \quad \text{diff} = 0$$

which gives the MLE

$$\hat{\lambda} = \frac{n}{\sum_i x_i} = \bar{x}^{-1}$$

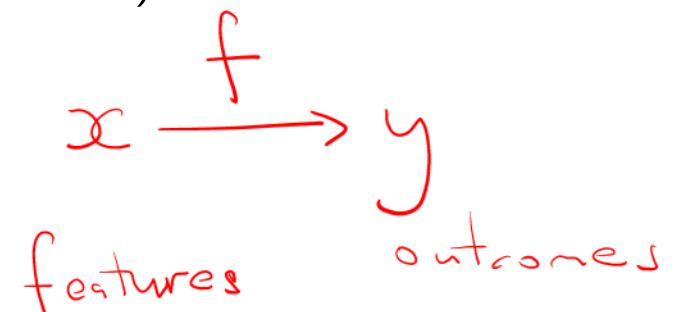
$$\hat{\lambda} = \frac{1}{\bar{x}}$$

Loss Functions

- We have an outcome measurement, which can be quantitative (such as stock price). Let's denote these **outcomes** as y .
- We wish to predict the outcome based on a set of **features**. Let's denote these features as x .
- We have a **training dataset**, where we observe the features and the outcome for a set of objects.
- Using the training data we build a **prediction model**, or **learner**, which we denote as $f(x)$. The model is a (possibly highly non-linear) transformation of features into outcomes:

$$y = f(x) + \varepsilon$$

f(x) - true fn. to approximate



Other types of notation are

$$\begin{aligned}y &= f_{\beta}(\mathbf{x}) = f(\mathbf{x}; \boldsymbol{\beta}) = \boldsymbol{\beta}^T \mathbf{x} \\&= (\beta_1, \dots, \beta_N) \begin{pmatrix} x_1 \\ \vdots \\ x_N \end{pmatrix} = \sum_{i=1}^N \beta_i x_i.\end{aligned}$$

input— Calculus $x \in \mathbb{R}$. Vector Calculus: $\mathbf{x} \in \mathbb{R}^D$. ML: can also work for discrete inputs, strings, tree, graphs, ...

output— Classification: $y \in \{0, 1\}$. Regression: $y \in \mathbb{R}$

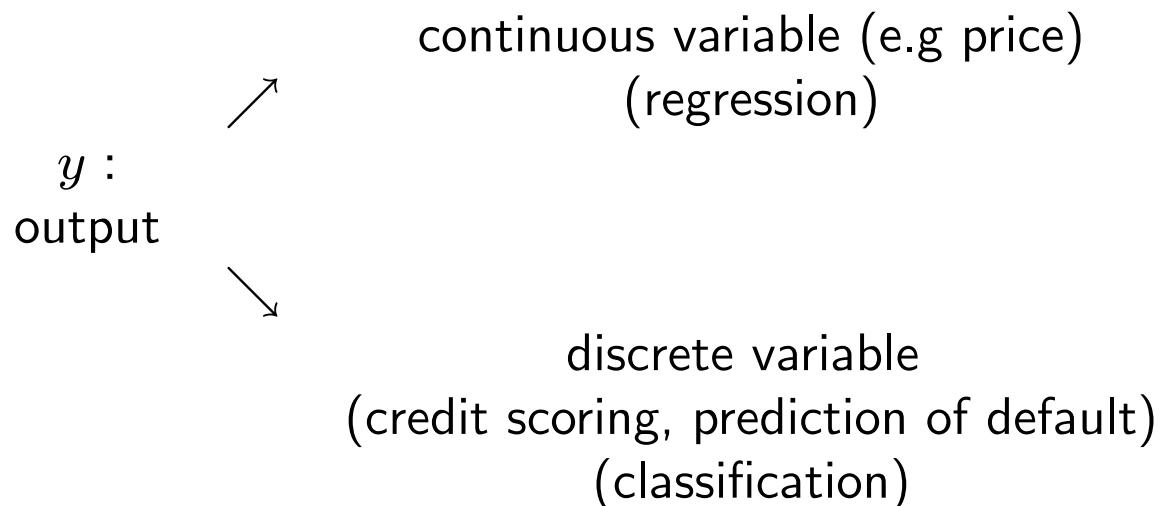
methods— Linear classifiers, neural networks, decision trees, ensemble models,...

probabilistic classifiers, ...

β — parameters/weights: $\beta \in \mathbb{R}$. $\boldsymbol{\beta} \in \mathbb{R}^D$

To set the scene: Here's a typical scenario of interest

Outcome y (this is the output) which is what we wish to predict



\mathbf{x} (set of features): input. So \mathbf{x} can be a scalar or vector.

p : number of features - size of vector, i.e. $\mathbf{x} \in \mathbb{R}^p$.

Set of input and output pairs $\mathcal{S} = \left\{ \left(x^{(i)}, y^{(i)} \right) : i = 1, \dots, N \right\}$;
data set for training

$x^{(i)} \in \mathbb{R}$, $y^{(i)} \in \mathbb{R}$ we want to estimate the parameter (weights) vector $\beta = (\beta_1, \dots, \beta_n)$ of unknown parameters.

$$(x^{(i)}, y^{(i)}) : i = 1, \dots, N \text{ pairs/2-tuples}$$

GOAL: Predict the relationship between the input and output.

So how is $x \rightarrow y$ achieved?



So output is some (unknown) function of the input, together with some noise.

Let's write

$$y = f(x) + \varepsilon$$

The estimate for f , which is sometimes denoted as \hat{f} .

$$\rightarrow \mathbb{E}[\varepsilon] = 0 \quad \mathbb{V}[\varepsilon] = \sigma^2 \leftarrow$$

$$\mathbb{E}[y] = \mathbb{E}[f + \varepsilon] = \mathbb{E}[f] + \mathbb{E}[\varepsilon] = f$$

The Loss Function L

Here we define the *loss function*

$$L(y, f(x)) = y - f(x)$$

as a measure of how far our prediction is from the outcome. A popular definition is

$$L(y, f(x)) = (y - f(x))^2 \quad \text{P} \leftarrow \text{L}$$

y — actual outcome. $f(x)$ — predicted outcome.

We want to minimise L so that the distance between the two outcomes is the smallest value possible.

Loss functions are used in ML to estimate how badly models are performing
- a measure of how wrong the model is in terms of its ability to estimate the relationship between input and output.

Look for $f(x)$ which minimises this 'error function' L , which can be written as

$$\arg \min_{f(x)} \{ \mathbb{E} [L(y, f(x))] \}$$

or in component form

$$\arg \min_{f(x)} \left\{ \frac{1}{N} \sum_{i=1}^N (y^{(i)} - f(x^{(i)}))^2 \right\}$$

MSE for the training set is

$$\frac{1}{N} \sum_{i=1}^N (y^{(i)} - f(x^{(i)}))^2$$

Loss fn. ML

In OR a classic optimisation problem is to maximise/minimise an objective function subject to a system of constraints. The function we want to minimize or maximize is called the objective function. However when we are minimizing it, we may also call it the **cost function**, **loss function**, or **error function**.

- suppose that X is a p -dimensional real-valued random input vector and Y is a real valued random output variable.

- The (unknown) joint distribution of X and Y is $P(X, Y)$.

- We seek to find a function $f(x)$ for predicting Y given X .

- To find the “optimal” $f(x)$, we need an objective (or loss) function, which we assume to be quadratic here

$$L(Y, f(X)) = (Y - f(X))^2$$

Econometrics

The expected loss function

math def. of an Expectation

$$\mathbb{E}[L] = \int_{\mathbb{R}} \int_{\mathbb{R}} (y - f(x))^2 P(x, y) dx dy$$

$$\begin{aligned}
 (y - f(x))^2 &= \underbrace{\{y - \mathbb{E}[y|x] + \mathbb{E}[y|x] - f(x)\}^2}_{\textcircled{1}} \\
 &= \underbrace{\{y - \mathbb{E}[y|x]\}^2}_{\textcircled{2}} + 2 \{y - \mathbb{E}[y|x]\} \{\mathbb{E}[y|x] - f\} + \underbrace{\{y - f\}^2}_{\textcircled{3}} \\
 &\quad \mathbb{E} \quad = 0
 \end{aligned}$$

$$\mathbb{E}[\textcircled{1}] = \mathbb{E}[(y - \mu_y)^2] = \textcircled{V}[y] \leftarrow \text{variance}$$

$$\mathbb{E}[\textcircled{2}] = 0 \quad \therefore \mathbb{E}[y - \mathbb{E}[y]] = \mathbb{E}[y] - \mathbb{E}[\mu_y] = \mu_y - \mu_y = 0$$

$$\mathbb{E}[\textcircled{3}] = \mathbb{E}\left[\underbrace{(\mathbb{E}[y] - \hat{f})^2}_{B(y, \hat{f})}\right] = \textcircled{bias}^2$$

$$(\textcircled{B}(\theta) = \mathbb{E}[\hat{\theta}] - \theta)$$

MSE of estimator =
 var. of est. + bias²

Bias-Variance decomposition

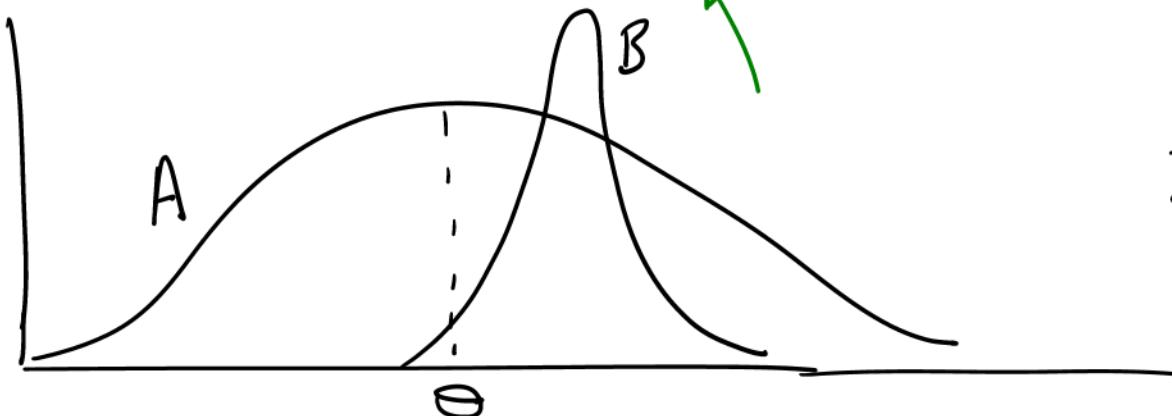
There is a trade-off between bias and variance, with very flexible models having low bias and high variance, and relatively rigid models having high bias and low variance. The model with the optimal predictive capability is the one that leads to the best balance between bias and variance.

$$y = f + \varepsilon$$

Writing $y - \hat{f} = f + \varepsilon - \hat{f}$

$$\mathbb{E}[(f + \varepsilon - \hat{f})] \rightarrow \mathbb{E}[\varepsilon^2] = \sigma^2 \quad \text{i.e. } \mathbb{V}[\text{noise}]$$

$$\mathbb{E}[L] = \mathbb{V}(\hat{f}) + \sigma^2 + \text{Bias}^2(\hat{f})$$



A: unbiased but larger var.

B: smaller var