# Topics in Analysis

June 2, 2024

## Contents

**Lectures**

Lecture 1
Lecture 2
Lecture 3
Lecture 4
Lecture 5
Lecture 6
Lecture 7
Lecture 8
Lecture 9
Lecture 10
Lecture 11
Lecture 12
Lecture 13
Lecture 14
Lecture 15
Lecture 16
Lecture 17
Lecture 18
Lecture 19
Lecture 20

# 1 Metric spaces

## 1.1 Revision

**Definition** (Metric space). A *metric space* is a pair $(X, d)$, where $X$ is a set, and $d : X^2 \to \mathbb{R}$ is a function satisfying:

- $d(x, y) \geq 0$

- $d(x, y) = 0 \iff x = y$

- $d(x, y) = d(y, x)$

- $d(x, y) + d(y, z) \geq d(x, z)$

**Example.**

(i) $\mathbb{R}^n$ with the metric

$$d(x, y) = \|\mathbf{x} - \mathbf{y}\| = \left( \sum_{j=1}^{n} (x_j - y_j)^2 \right)^{\frac{1}{2}}$$

(ii) $\mathbb{C}$ with the metric

$$d(z, w) = |z - w|$$

**Notation.** We write $x_n \xrightarrow{d} x$ if and only if $d(x_n, x) \to 0$ as $n \to \infty$.

**Notation.** Open ball $B(x, y) = \{y : d(x, y) < r\}$ (an open set).

**Definition** (Closed and Open). A set $E \subseteq X$ is closed if for all convergent sequences $x_n \xrightarrow{d} x$, with $x_n \in E$, we have $x \in E$. $U \subseteq X$ is open if whenever $x \in U$ there exists $\delta > 0$ such that $B(x, \delta) \subseteq U$.

3

**Proposition.** If $U$ is open then $U^c$ is closed.

*Proof.* Suppose $x \in U$, then there exists $\delta > 0$ such that $B(x, \delta) \subseteq U$, so $d(x_n, x) \geq \delta$ whenever $x_n \in U^c$ so $x_n \overset{d}{\nrightarrow} x$. $\qquad\square$

**Proposition.** If $E$ is closed then $E^c$ is open.

*Proof.* Suppose $y \in E^c$. Then there does not exist $y_n \in E$ such that $y_n \overset{d}{\to} y$, so there exists $\delta$ such that $B(y, \delta) \subseteq E^c$. $\qquad\square$

**Definition** (Cauchy sequence)**.** Working in $(X, d)$, we say $(x_n)$ is *Cauchy* if

$$\forall \varepsilon > 0 \; \exists N \; \forall m, n \geq N \quad d(x_n, x_m)$$

(sometimes just say $d(x_n, x_m) \to 0$ as $n, m \to \infty$).

**Proposition.** Any convergent sequence is Cauchy.

*Proof.* If $x_n \overset{d}{\to} x$ then

$$\forall \varepsilon > 0 \; \exists N \; \forall n \geq N \quad d(x_n, x) < \frac{\varepsilon}{2}$$

and so

$$\forall \varepsilon > 0 \; \exists N \; \forall m, n \geq N \quad d(x_n, x_m) < \varepsilon \qquad\square$$

If every Cauchy sequence converges we say $(X, d)$ is complete.

The following remearks are useful.

**Proposition.** If a subsequence of a Cauchy sequence converges, then the sequence converges.

*Proof.* Suppose $(x_n)$ is Cauchy, $x_{n(j)} \to x$, with $n(j) \to \infty$. Let $\varepsilon > 0$ and choose $N$ such that $d(x_n, x_m) < \frac{\varepsilon}{2}$ for all $n, m \geq N$. Choose $n(J) \geq N$ such that $d(x_{n(J)}, x) < \frac{\varepsilon}{2}$. then if $yn \geq N$

$$d(x_n, x) \leq d(x_n, x_{n(J)}) + d(x_{n(J)}, x) < \frac{\varepsilon}{2} + \frac{\varepsilon}{2} \qquad\square$$

**Proposition.** To show $(X, d)$ complete, we need only show that for some $\varepsilon(n) \to 0$, $d(y_n, y_{n+1}) < \varepsilon(n) \implies y_n$ converges.

*Proof.* If $(x_n)$ is Cauchy, we can find $n(j)$ such that $d(x_{n(j)}, x_{n(j+1)}) < \varepsilon_j$. $\qquad\square$

**Remark.** If $(X, d)$ is a metric space and $Y \subseteq X$ then writing

$$d_Y(y_1, y_2) = d(y_1, y_2) \; \forall y_1, y_2 \in Y$$

we have $(Y, d_y)$ a metric space.

**Lemma.** If $E$ is closed in $(X, d)$ and $(X, d)$ is complete then $(Y, d_Y)$ is complete.

*Proof.* If $(y_n)$ is Cauchy in $(Y, d_Y)$ then $(y_n)$ is Cauchy in $X$, so $y_n \overset{d}{\to} x$. But $E$ is closed, so $x \in E$, so $x \in Y$, so $y_n \overset{d_Y}{\to} x$. $\qquad\square$

**Proposition.** If $(Y, d_Y)$ is complete then $Y$ is closed in $X$.

*Proof.* If $y_n \overset{d}{\to} x$, $y_n \in Y$ then $(y_n)$ is Cauchy for $d$ and so for $d_Y$, so $y_n \overset{d_Y}{\to} y$, $y \in Y$, so $y_n \overset{d}{\to} y \in Y$ so by uniqueness of limits, $x = y \in Y$. $\qquad\square$

From IA Numbers and Sets, we know that $\mathbb{R}$ is complete.

**Theorem.** $\mathbb{R}^n$ (with usual Euclidean meteric) is complete.

*Proof.* $\|\mathbf{x}(m) - \mathbf{x}(n)\| \geq |x_j(m) - x_j(n)|$, so $(\mathbf{x}(n))$ Cauchy implies $x_j(n)$ Cauchy, so $x_j(n) \to x_j$ for some $x_j$, so

$$\sum_{j=1}^n (x_{j(n)} - x_j)^2 \to 0,$$

so $\mathbf{x}_j \to \mathbf{x}$. $\qquad\square$

**Theorem** (Bolzano Weierstrass (in $\mathbb{R}$)). Let $[a, b]$ be a closed interval. If $x_n \in [a, b]$, then there exists a sequence $n(j) \to \infty$ and $x \in [a, b]$ such that $x_{n(j)} \to x$.

This theorem has an easy extension to $\mathbb{R}^m$.

**Theorem** (Bolzano Weierstrass (in $\mathbb{R}^m$)). If $\mathbf{x}(k) \in \prod_{j=1}^{m} [a_j, b_j]$, then $\mathbf{x}(k)$ has a convergent subsequence.

*Proof.* Proof by induction. True for $m = 1$ (see IA Numbers and Sets). Suppose true for $m$. Then if $\mathbf{x}(k) \in \prod_{j=1}^{m+1} [a_j, b_j]$, we write

$$\mathbf{x}(k) = (\mathbf{y}(k), z_k)$$

with $\mathbf{y}(k) \in \prod_{i=1}^{m} [a_j, b_j]$, $z_k \in [a_{m+1}, b_{m+1}]$. By inductive hypothesis, there exists $k(r) \to \infty$ such that $\mathbf{y}(k(r)) \to \mathbf{y} \in \prod_{j=1}^{m} [a_j, b_j]$. By Bolzano Weierstrass in $\mathbb{R}$, there exists $r(s) \to \infty$ such that $z_{k(r(s))} \to z$. Then

$$\mathbf{x}(k(r(s))) \to (\mathbf{y}, z). \qquad \square$$

**Theorem** (Bolzano Weierstrass (for closed and bounded sets)). If $E \subseteq \mathbb{R}^n$ is closed and bounded, then if $\mathbf{e}_n \in E$ is a sequence in $E$, then there exists a sequence $n(k) \to \infty$ such that $\mathbf{e}_{n(k)} \to \mathbf{e} \in E$. Only true if $E$ is closed and bounded.

*Proof.* Choose $R$ such that $[-R, R]^n \supseteq E$ (possible since $E$ is bounded). Then if $\mathbf{e}_r \in E$, $\mathbf{e}_r \in [-R, R]^n$ so there exists $r(j) \to \infty$ and $\mathbf{e}$ such that $\mathbf{e}_{r(j)} \to \mathbf{e}$. But $E$ is closed so $\mathbf{e} \in E$. $\qquad \square$

This is false if $E$ is not closed. Pick $x, x_n$, $x_n \in E$ with $x_n \to x$, $x \notin E$. Then any subsequence of $x_n$ also converges to $x$. If $E$ is not bounded pick $|x_n| \geq n$ with $x_n \in E$. This sequence has no convergent subsequence.

Start of

lecture 2

**Definition.** Let $(X, d)$, $(Y, e)$ be metric spaces. We say $f : X \to Y$ is continuous if given $x \in X$, $\varepsilon > 0$ there exists $\delta(\varepsilon, x)$ such that

$$d(x', x) < \delta \implies e(f(x), f(x')) < \varepsilon.$$

**Proposition.** This is equivalent to the "if $U$ open in $Y$ then $f^{-1}(U)$ open in $X$" definition.

*Proof.* We first show that if $f$ satisfies the first definition, then it satisfies the second. Suppose $U$ open in $Y$. If $x \in f^{-1}(Y)$ then $f(x) \in U$ so there existst $\varepsilon > 0$ such that $B_Y(f(x), \varepsilon) \subseteq U$. Hence there exists $\delta$ such that $f(B(x, \delta)) \subseteq B_Y(f(x), \varepsilon)$, so $B(x, \delta) \subseteq f^{-1}(U)$. So $f^{-1}(U)$ is open.

For the other direction: if $x \in X$ then given $\varepsilon > 0$, $B(f(x), \varepsilon)$ is open so $f^{-1}(B(f(x), \varepsilon))$ is open. $x \in f^{-1}(B(f(x), \varepsilon))$ so there exists $\delta$ such that $B(x, \delta) \subseteq f^{-1}(B(f(x), \varepsilon))$ and we recover the $\varepsilon, \delta$ definition. $\qquad\square$

We introduce an idea much used in the course:

**Definition** (Distance to $A$)**.** If $A \subseteq \mathbb{R}^n$ is closed and non-empty, we define

$$d(x, A) = \inf_{a \in A} \|x - a\|.$$

**Remark.** (1) $d(x, A) = 0 \iff x \in A$: the backwards direction is trivial, and for the forwards direction, if $d(x, A) = 0$ then there exists $a_n \in A$ with $d(x, a_n) \to 0$. But $A$ is closed.

(2) $x \mapsto d(x, A)$ is continuous.

*Proof.* Let $\varepsilon > 0$ then for given $x, y$ there exists $a \in A$ such that $\|x - a\| \leq d(x, A) + \frac{\varepsilon}{2}$. So $\|y - a\| \leq \|x - a\| + \|y - x\| \leq d(x, A) + \|x - y\| + \frac{\varepsilon}{2}$. So $d(y, A) \leq \|y - x\| + d(x, A) + \frac{\varepsilon}{2}$. Since $\varepsilon$ is arbitrary,

$$d(y, A) \leq d(x, A) + \|x - y\|$$

Similarly, $d(x, A) \leq d(y, A) + \|x - y\|$, so

$$|d(x, A) - d(y, A)| \leq \|x - y\|. \qquad\square$$

Now we move towards more interesting results.

**Theorem.** If $E \subseteq \mathbb{R}^n$ is compact and $f : E \to \mathbb{R}^n$ is continuous, then $f(E)$ is compact.

*Proof.* If $y_n \in f(E)$ then $y_n = f(x_n)$ for some $x_n \in E$. By Bolzano Weierstrass (in $\mathbb{R}^m$), there exists $n(j) \to \infty$ and $x \in E$ such that $x_{n(j)} \to x$. So $y_{n(j)} = f(x_{n(j)}) \to f(x) \in E$. □

> **Corollary.** If $E \subseteq \mathbb{R}^n$ is compact and non-empty and $f : E \to \mathbb{R}$ is continuous then $f$ is bounded and attains its bounds.

*Proof.* $f(E)$ is compact in $\mathbb{R}$ so closed and bounded. Further since $f(E)$ is bounded above (and non-empty), $f(E)$ has a supremum $\alpha$ so there exists $e_n$ such that $f(e_n) \to \alpha$. By compactness, there exists $n(j) \to \infty$ and $e$ such that $e_{n(j)} \to e$ so $f(e) = \alpha$. □

> **Theorem** (Fundamental Theorem of Algebra)**.** If $P$ is a non-constant polynomial then it has a root.

> **Lemma.** If $P$ is a non-constant polynomial then $|P(z)| \to \infty$ as $|z| \to \infty$.

*Proof.*

$$P(z) = \sum_{j=0}^{n} a_j z^j$$

$$= z^n \left( a_n - \sum_j a_j z^{j-n} \right)$$

$a_n \neq 0$, $n \geq 1$. Note

$$a_n - \sum_j a_j z^{j-n} \to a_n$$

as $|z| \to \infty$, hence since $|z^n| \to \infty$ as $|z| \to \infty$, we have $|P(z)| \to \infty$ as $|z| \to \infty$. □

> **Lemma.** In particular we can find an $R$ such that $|P(z)| > |P(0)|$ for all $|z| \geq R$.

*Proof.* By compactness there exists $|\alpha| \leq R$ such that

$$|P(z)| \geq |P(\alpha)| \ \forall |z| \leq R.$$

But $|P(z)| \geq |P(0)| \geq |P(\alpha)| \ \forall |z| \geq R$. So $|P(\alpha)| \leq |P(z)| \ \forall z$. □

We now show that if $\alpha$ is a minima of $|P(z)|$ then $P(\alpha) = 0$ (to some extent the rest of the proof is a matter of personal preference).

> **Remark.** By considering $P(z+\alpha)$ if necessary, we may suppose that $|P(0)| \leq |P(z)|$ for all $z$. If $P(0) = 0$ we are done (using the previous two lemmas, we know that $|P(z)|$ has a global minimum, because the minimum on $B(0, R)$, which is achieved by compactness, must be a global minimum).

*Proof of Fundamental Theorem of Algebra.* By the remark, assume $|P(0)| \leq |P(z)|$ for all $z$. If $P(0) = 0$ we are done. If not, then

$$P(z) = a_0 + \sum_{j=1}^{n} a_j z^j$$

with $a_0 \neq 0$. By considering $\frac{P(z)}{a_0}$ we may suppose $a_0 = 1$ and that

$$P(z) = 1 + a_k z^k + \sum_{j=k+1}^{n} a_j z^j$$

with $a_k \neq 0$.

By considerng $P(re^{i\theta}z)$ for suitable $r, \theta$ we may suppose $a_k = -1$. So

$$P(z) = 1 - z^k + \sum_{j=k+1}^{n} z^j$$

Since $\frac{\sum_{j=k+1}^{n} a_j z^j}{z^k} \to 0$ as $z \to 0$, there exists $\delta > 0$ such that

$$\left| \sum_{j=k+1}^{n} a_j z^j \right| < \frac{|z|^k}{2}$$

for $|z| < \delta$. If $\eta = \frac{\delta}{2}$, then

$$P(\eta) = 1 - \eta^k + \varepsilon(\eta)$$

with $|\varepsilon(\eta)| < \frac{\eta^k}{2}$. So $P(\eta) < 1 = P(0)$, contradiction. So $P(0) = 0$, and $P$ has a root as desired. $\qquad\square$

> **Corollary.** If $P$ has degree $n \geq 1$ then $P(z) = (z - a)Q(z)$ for some polynomial $Q$ of degree $n - 1$.

*Proof.* By long division (in fact proof by induction). If $P$ degree $n$, then $P(z) = (z - a)Q(z) + r$ with $\deg Q = n - 1$, $r \in \mathbb{C}$. If $P(a) = 0$ then $0 = 0 + r$ so $r = 0$, so $P(z) = (z - a)Q(z)$. $\qquad\square$

By induction if $P$ has degree $n$,

$$P(z) = A \prod_{j=1}^{n} (z - a_j)$$

for some $A \neq 0$.

> **Corollary.** If $P$ is a polynomial of degree at most $n$ which vanishes at $n + 1$ points then $P = 0$.

*Proof.* If $\deg P = k \geq 0$, then $P(z) = A \prod_{j=1}^{k} (z - a_j)$. Then $P(z) \neq 0$ for $z \neq a_1, a_2, \ldots, a_k$. $\qquad\square$

Start of

lecture 3

**\* Non-examinable material**

Laplace's equation:
$$\nabla^2 \phi = 0$$

Want to solve on $E$ with boundary conditions.

Two questions:

- Does a solution exist?

- If so, is it unique?

Dirichlet used arguments like this:

$$
\begin{aligned}
\int_E \nabla\phi \nabla phi \, dV &= \int_E \nabla \cdot (\phi \nabla \phi) - \nabla^2 \phi \, dV \\
&= \int_E \nabla \cdot (\phi \nabla \phi) \, dV && \text{if } \nabla^2 \phi = 0 \\
&= \int_{\partial E} \phi \frac{\partial \phi}{\partial n} \, dS && \text{if } \phi = 0 \text{ on surface}
\end{aligned}
$$

$\nabla\phi \cdot \nabla\phi \geq 0$, so $\nabla\phi$ constant, so $\phi$ constant so $\phi = 0$ on $E$.

It is natural to think of

$$\int \nabla\phi \cdot \nabla\phi \mathrm{d}V$$

as an energy and Dirichlet showed that the minima corresponds to a solution of Laplace's equation.

**Problems:**

- What are the conditions on $\phi$ and on the surface which actually permit us to use the Divergence Theorem.

- Does the energy $\int \nabla\phi \cdot \nabla\phi \mathrm{d}V$ have a minimum?

- If given on the boundary, does there exist a $\phi$ with $\phi = f$ on the boundary and $\int \nabla\phi \cdot \nabla\phi < \infty$.

\*\* This is the end of the non-examinable comments.


**Examinable content again**

Interested in Laplace's equation. Start by deciding what a boundary is.

> **Definition** (Closure)**.** If $E$ is a set in $(X, d)$ a metric space, then the closure of $E$ is
>
> $$\mathrm{Cl}E = \{x \in E : \exists e_n \in E, e_n \xrightarrow{d} x\}$$

Properties:

(1) Closure is closed. Suppose $z_n \in \mathrm{Cl}E$ and $z_n \xrightarrow{d} z$. Then $\exists e_n \in E$ such that $d(e_n, z_n) < \frac{1}{n}$, so by triangle inequality $e_n \to z$ so $z \in \mathrm{Cl}E$.

(2) Closure of $E$ is the smallest closed set containing $E$.

> **Definition** (Interior)**.** If $E$ is a set in $(X, d)$ a metric space, then the interior of $E$ is
> $$\mathrm{Int}E = \{x \in E : \exists \delta > 0, B(x, \delta) \subseteq E\}$$

The interior of $E$ is the largest open set contained in $E$ (proof left to reader).

$\partial E$ is the boundary of $E$ and is defined by $\partial E = \mathrm{Cl}E \setminus \mathrm{Int}E$.

**Remark.** The boundary can be a bit odd. For example,

$$E = \{x \in \mathbb{R}^n : \|x\| = 1\}$$

has the property that $E = \partial E$.

However, for simple things, the notion of the boundary of an open set is what we would expect.

**Remark.** The boundary $\partial E = \mathrm{Cl}E \cap (\mathrm{Int}E)^c$ is closed.

**Remark.** If we work in $\mathbb{R}^n$ and $E$ is bounded then $\partial E$ is bounded (and closed) hence compact.

**Theorem.** Let $\Omega \subseteq \mathbb{R}^n$ is open, non-empty and bounded. Let $f : \Omega \to \mathbb{R}$ be continuous. Let $\phi : \mathrm{Cl}\Omega \to \mathbb{R}$ be continuous, and twice differentiable on $\Omega$ with $\nabla^2 \phi = 0$ on $\Omega$.

Then there can exist at most one such $\phi$ with $\phi = f$ on $\partial\Omega$.

Suffices to prove:

**Lemma.** If $\Omega \subseteq \mathbb{R}^n$ is open, non-empty and bounded, with:

- $\phi = \mathrm{Cl}\Omega \to \mathbb{R}$ continuous

- $\phi$ twice differentiable on $\Omega$

- $\phi$ continuous on $\mathrm{Cl}\Omega$

- $\phi = 0$ on $\partial\Omega$

then $\phi = 0$.

*Proof of Theorem from Lemma (easy).* Let $\phi_1$ and $\phi_2$ satisfy the conditions of the theorem. Set $\phi = \phi_1 - \phi_2$. Then $\phi$ satisfies the conditions of the lemma, so $\phi = 0$ on $\mathrm{Cl}E$, so $\phi_1 = \phi_2$ on $\mathrm{Cl}E$. $\qquad\qquad\square$

*Proof.* **Key step:** If $\Omega$ is open, non-empty and bounded with $\phi$ continuous on $\mathrm{Cl}\Omega$ and

$\phi$ twice differentiable on $\Omega$ with $\nabla^2 \phi > 0$, then it can have no interior minimum, i.e. the minimum (know exists by compactness) must be attained on $\partial \Omega$.

Suppose $\mathbf{y} \in \Omega$. Then we can find $\delta > 0$ such that

$$\prod [y_j - \delta, y_j + \delta] \subseteq \Omega$$

Since $\nabla^2 \phi > 0$ at $\mathbf{y}$ there must exist a $k$ such that $\frac{\partial^2 \phi}{\partial x_k^2} \mathbf{y} > 0$ so $\exists 0 < \eta < \delta$ such that $\frac{\partial^2 \phi}{\partial x_k^2} > 0$ on $\prod [y_j - \eta, y_j + \eta]$. Look at

$$g(t) = \phi(y_1, y_2, \ldots, y_{k-1}, y_k + t, y_{k+1}, \ldots, y_m)$$

$g''(0) > 0$ so no minimum at $t = 0$ so $\phi$ has no minimum at $\mathbf{y}$ ※.

**Next step:** Conditions as before, but replace $\nabla^2 \phi > 0$ by $\nabla^2 \phi = 0$ on $\Omega$. Then the global minimum for $\phi$ is on the boundary.

Set $\phi_N(\mathbf{x}) = \phi(\mathbf{x}) + \frac{1}{N}(x_1^2 + x_2^2 + \cdots + x_m^2)$. Then

$$\nabla^2 \phi_N(\mathbf{x}) = \nabla^2 \phi + \frac{1}{N}(2 + 2 + 2 \cdots + 2) = \frac{2m}{N} > 0.$$

By previous result, there exists $\mathbf{x}_N \in \partial \Omega$ such that $\phi(\mathbf{x}_N) \geq \phi(\mathbf{z})$ for all $\mathbf{z} \in \mathrm{Cl}\Omega$. There exists $N(j) \to \infty$ and $\mathbf{x}^* \in \partial \Omega$ (compactness) such that $\mathbf{x}_{N(j)} \to \mathbf{x}^*$. By continuity, $\phi(\mathbf{x}_{N(j)}) \to \phi(\mathbf{x}^*)$.

$$\phi(\mathbf{x}_{N(j)}) + \frac{2m}{N(j)} \geq \phi(z) + \frac{2m}{N(j)} \geq \phi(\mathbf{z}) \qquad \forall z \in \mathrm{Cl}\Omega$$

so

$$\phi(\mathbf{x}^*) \geq \phi(\mathbf{z}) \qquad \forall \mathbf{z} \in \mathrm{Cl}\Omega.$$

We have shown $\nabla^2 = 0$ on $\Omega$ implies that we have a maximum on the boundary. But $\nabla^2 = \phi$ then $\nabla^2(-\phi) = 0$ so $-\phi$ has maximum on boundary so if $\phi = 0$ on $\partial \Omega$ we have $\phi = 0$ on $\mathrm{Cl}\Omega$. $\qquad \square$
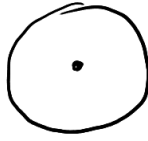
If Laplace's equation has a solution, then it is unique.

We now give an example of Zaremba which shows that (at least in the general form we have given) there need not be a solution.

We take $E = \{\mathbf{x} \in \mathbb{R}^2 : 0 < \|x\| \leq 1\}$.

Observe $\mathrm{Int}E = \{\mathbf{x} : 0 < \|x\| \leq 1\}$, $\partial E = \{x : \|x = 1\|\} \cup \{\mathbf{0}\}$.

We show that there does not exist $\phi$ such that $\phi$ is continuous on $\mathrm{Cl}E = \{\mathbf{x} : \|x\| \leq 1\}$ which is twice differentiable on $\mathrm{Int}E = \{\mathbf{x} : 0 < \|\mathbf{x}\| < 1\}$ with $\nabla^2 \phi = 0$ on $\mathrm{Int}E$ satisfying $\phi(x) = 0$ on $\|\mathbf{x}\| = 1$ with $\phi(\mathbf{0}) = 1$.

Suppose such a $\phi$ exists. Let $R_\theta$ be a rotation through $\theta$ about $\mathbf{0}$. Then $\phi_\theta(\mathbf{x}) = \phi(R_\theta \mathbf{x})$ also solve the problem, so by uniqueness we must have $\phi_\theta = \phi$ for all $\theta$. So $\phi$ is radially symmetric: $\phi(x, y) = f(r)$, $r = \sqrt{x^2 + y^2}$.

Recall from Vector Calculus the formula for $\nabla^2$ in polar coordinates. So we seek to solve

$$\frac{1}{r} \frac{\mathrm{d}}{\mathrm{d}r} (r f'(r)) = 0$$

subject to $f(0) = 1$, $f(1) = 0$. So:

$$
\begin{aligned}
& \frac{\mathrm{d}}{\mathrm{d}r} (r f'(r)) = 0 \\
& \implies r f'(r) = A \\
& \implies f'(r) = \frac{A}{r} \\
& \implies f(r) = A \log r + B
\end{aligned}
$$

But we want $f(r) \to 1$ as $r \to \infty$, so $A = 0$. So $f(1) = B$, contradiction.

This completes our discussion of Laplace's equation. One can show that the equation always has a solution in 2D, provided that the boundary is a Jordan curve, but the proof of this result is very hard. The result is similar to the Riemann mapping theorem. For higher dimensions, there are no particularly nice theorems about existence of solutions.

## 2 Brouwer's Fixed Point Theorem

> **Lemma.** If $f : [0,1] \to [0,1]$ is continuous then $\exists x \in [0,1]$ such that $f(x) = x$.

*Proof.* Set $g(t) = f(t) - t$. Then $g : [0,1] \to \mathbb{R}$ is continuous. $g(0) = f(0) \geq 0$, $g(1) = f(1) - 1 \leq 0$. So by IVT, there exists $x$ such that $g(x) = 0$. $\qquad\square$

> **Theorem** (Brouwer). If $\overline{D} = \{\mathbf{x} \in \mathbb{R}^2 : \|\mathbf{x}\| \leq 1\}$ and $f : \overline{D} \to \overline{D}$ is continuous, then there exists $\mathbf{y}$ such that $f(\mathbf{y}) = \mathbf{y}$.

> **Remark.**
>
> (a) This reusult works in $\mathbb{R}^n$. Most of our proof will work with obvious modifications in $\mathbb{R}^n$. One bit, "Sperner's Lemma" requires work (but not enromous changes).
>
> (b) If $E \subseteq \mathbb{R}^2$ and $T : E \to \overline{D}$ is a homeomorphism (i.e. $T$ bijective, $T$ continuous and $T^{-1}$ continuous) then Brouwer continues to work. Suppose $g : E \to E$ is continuous. Then $T \circ g \circ T^{-1} : \overline{D} \to \overline{D}$ is continuous, so has a fixed point $\mathbf{z}$, satisfying $TgT^{-1}(z) = z$, so $gT^{-1}(z) = T^{-1}(z)$.

The proof goes through many steps

$$A \iff B \iff C \iff D$$

so to understand the proof we need to understand the strategy (ie $A$, $B$, ... and the tactics $\iff$ ).

We start with the "no retraction" theorem.

> **Theorem** (No Retraction Theorem). There does not exist a $g : \overline{D} \to \partial D$ continuous with $g(\mathbf{x}) = \mathbf{x} \ \forall \mathbf{x} \in \partial D$.
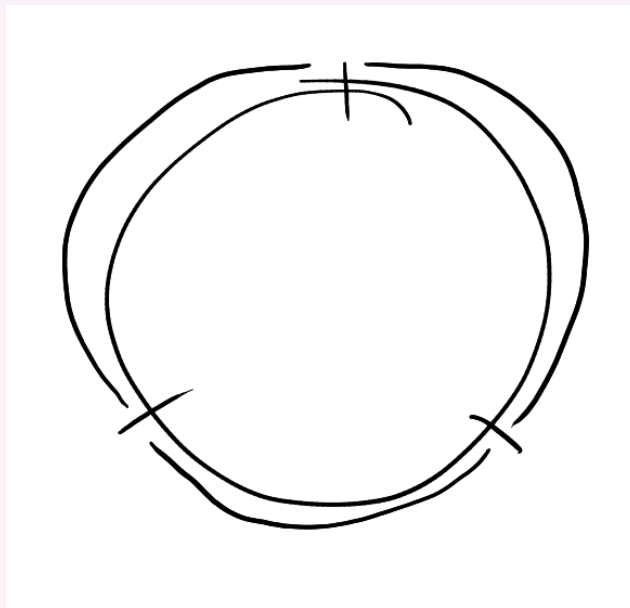
We will start by showing No Retraction Theorem $\iff$ Brouwer.

*Proof.* Suppose No Retraction Theorem is false. Then there exists $g : \overline{D} \to \overline{D}$ continuous with $g(\mathbf{x}) = \mathbf{x}$ for all $\mathbf{x} \in \partial D$. Now let $R$ be a rotation about the origin through angle $\pi$. Then $Rg : \overline{D} \to \partial D \subseteq \overline{D}$ has no fixed point. So Brouwer would have to be false.

If the No Retraction Theorem is true then we now want to show that Brouwer holds. Suppose $f : \overline{D} \to \overline{D}$ is continuous without a fixed point. Then we define $g(\mathbf{x})$ to be the point of intersection of the ray $f(\mathbf{x})$ to $\mathbf{x}$ with $\partial D$. Notice $g$ is continuous, and fixed $\partial D$, so $g$ is a retraction mapping. $\qquad\square$

Trying to prove Brouwer. We have shown that it is equivalent to No Retraction Theorem.

Next step is to establish equivalence with a weaker cousin of no No Retraction Theorem.

> **Lemma** (Cousin of No Retraction Theorem)**.** No Retraction Theorem is equivalent to the following: divide $\partial \overline{D}$ into 3 equal arcs.
>
> 
>
> In polars:
>
> $$I = \left\{ (1, \theta) : 0 \leq \theta \leq \frac{2\pi}{3} \right\}$$
>
> $$J = \left\{ (1, \theta) : \frac{2\pi}{3} \leq \theta \leq \frac{4\pi}{3} \right\}$$
>
> $$K = \left\{ (1, \theta) : \frac{4\pi}{3} \leq \theta \leq 2\pi \right\}$$
>
> Then there does not exist $g : \overline{D} \to \partial D$ continuous such that $g(I) \subseteq I$, $g(J) \subseteq J$, $g(K) \subseteq K$.

*Proof.* If this cousin is true, then No Retraction Theorem follows at once.

On the other hand if the cousin is false, and such a $g$ exists, then if $T$ is rotation about 0 through $\pi$ then $T \circ g$ has no fixed point. $\qquad\square$

The cousin has a triangle version.

---

**Theorem.** If $\overline{\triangle}$ is an equilateral triangle (closed) with sides $\tilde{I}, \tilde{J}, \tilde{K}$, then there does not exist $G : \overline{\triangle} \to \overline{\partial\triangle}$ such that $G(\tilde{I}) \subseteq \tilde{I}$, $G(\tilde{J}) = \tilde{J}$, $G(\tilde{K}) = \tilde{K}$.

---

This is equivalent to the previous result by homeomoprhism (use a homeomorphism $H : \overline{D} \to \overline{\triangle}$ with $H(I) = \tilde{I}$, $H(J) = \tilde{J}$, $H(K) = \tilde{K}$).

---

**Theorem.** The following two statements about an equilateral $\overline{\triangle}$ with sides $I, J, K$ (note sides include end point vertices) are equivalent:

(i) There does not exist $h : \overline{\triangle} \to \partial D$ continuous with $h(I) \subseteq I$, $h(J) \subseteq J$, $h(K) \subseteq K$.

(ii) There does not exist $A, B, C$ closed, $A, B, C \subseteq \overline{\triangle}$ such that $A \cup B \cup C = \overline{\triangle}$, $A \supseteq I$, $B \supseteq J$, $C \supseteq K$ and $A \cap B \cap C = \emptyset$.

---

*Proof.* If we could find $h : \overline{\triangle} \to \partial\triangle$ with $h(I) \subseteq I$, $h(J) \subseteq J$, $h(K) \subseteq K$. Then let

$$A = h^{-1}(I)$$
$$B = h^{-1}(J)$$
$$C = h^{-1}(K)$$

Then

$$A \cap B \cap C = h^{-1}(I \cap J \cap K) = h^{-1}(\emptyset).$$

For the other direction, suppose conversely that we have $A, B, C$ closed such that $A \supseteq I$, $B \supseteq J$, $C \supseteq K$, $A \cup B \cup C = \triangle$, $A \cap B \cap C = \emptyset$. We look at the triangle in $\mathbb{R}^3$ given by

$$\{(x, y, z) : x + y + z = 1, x, y, z \geq 0\}.$$

Now look at $d(\mathbf{x}, A)$, $d(\mathbf{x}, B)$, $d(\mathbf{x}, C)$. Remember that

$$d(\mathbf{x}, E) = \inf_{\mathbf{E} \in E} \|\mathbf{x} - \mathbf{e}\|,$$

and if $E$ is closed, then $d(\mathbf{x}, E) = 0$ if and only if $\mathbf{x} \in E$. We also remarked before that $\mathbf{x} \mapsto d(\mathbf{x}, E)$ is continuous. So

$$x \mapsto d(x, A), \qquad x \mapsto d(x, B), \qquad x \mapsto d(x, C)$$

are continuous. So

$$x \mapsto d(x, A) + d(x, B) + d(x, C)$$

is continuous. Since $A \cap B \cap C = \emptyset$, we have for each $x$ that at least one of the distances is positive. Thus

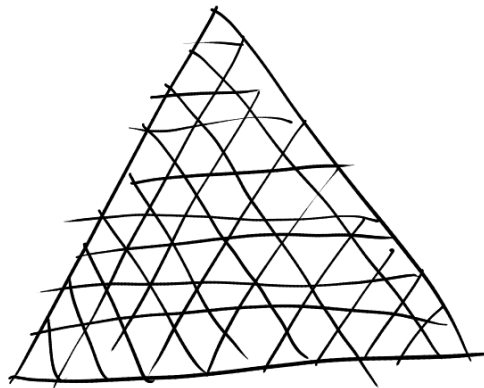$$d(\mathbf{x}, A) + d(\mathbf{x}, B) + d(\mathbf{x}, C) > 0.$$

Thus the function $F$ given by

$$\mathbf{x} \mapsto \left( \frac{d(x, A)}{d(x, A) + d(x, B) + d(x, C)}, \frac{d(x, B)}{d(x, A) + d(x, B) + d(x, C)}, \frac{d(x, C)}{d(x, A) + d(x, B) + d(x, C)} \right)$$

is well defined and continuous. Let the components of $F(x)$ be denoted by $F_i(x)$ for $i = 1, 2, 3$. Note $F_j(x) \geq 0$, and $F_1(x) + F_2(x) + F_3(x) = 1$. So $F$ maps $\triangle$ to $\overline{\triangle}$ continuously. If $\mathbf{x} \in I$ then $F_1(\mathbf{x}) = 0$, so $F(A) \subseteq I$. Similarly for the others.  $\square$

Thus we have changed out problem to a colouring problem.

We attack this by looking at a finite problem. Take an equilateral triangle and cut it up by using $n$ equally spaced lines parallel to each of the sides.



If you colour the vertices obeying the following rule by triangle $ABC$:

- All vertices on $AB$ except $B$ are red

- All vertices on $BC$ except $C$ are blue.

- All vertices on $CA$ except $A$ are green.

Remaining vertices are your choice. Then there is at least one small triangle with all 3 colours.

This is known as Sperner's Lemma.

We will prove only the 2 dimensional version of this theorem, but a more general version works in $\mathbb{R}^n$.

> **Lemma** (Sperner's Lemma)**.** Take an equilateral triangle formed from a collection of smaller triangles by using $n$ equally spaced lines parallel to each of the sides of the big triangle. Now colour all the vertices red, green or blue subject to the rule that if $XYZ$ are the vertices of the big triangle, then all the vertices on $XY$ except $X$ are red, all vertices on $YZ$ except $Z$ are blue and all vertices on $ZX$ except $X$ are green.
>
> Then at least one small triangle has vertices of all 3 colours.

*Proof.* Let $\Gamma$ be the set of small vertices. Once the colouring has been chosen, we will assign each $T \in \Gamma$ an integer $\zeta(T)$ according to a rule. For vertices of $\alpha, \beta$ of $T$, define $\zeta(\alpha\beta)$ to be:

$$\zeta(\alpha\beta) = \begin{cases} 1 & \text{RG, BR, GB} \\ 0 & \text{RR, BB, GG} \\ -1 & \text{GR, RB, BG} \end{cases}$$

Now define

$$\zeta(T) = \zeta(\alpha\beta) + \zeta(\beta\gamma) + \zeta(\gamma\alpha)$$

where $\alpha, \beta, \gamma$ are the vertices listed counter clockwise.

Observe $\zeta(T) = 3$ or $\zeta(T) = -3$ if $T$ is 3-coloured, and $\zeta(T) = 0$ otherwise (to see this, just consider the cases where we have exactly 2 of the same colour, or having 3 of the same colour). Now look at $\sum_T \zeta(T)$. Note that all inner edges cancel, so

$$\sum_T \zeta(T) = \sum_{xy} \zeta(xy)$$

where $xy$ ranges over all sides of the big triangle, pointing counter clockwise. This sum on the right equals 3 (exercise). So there must be a 3 coloured small triangle. $\qquad\square$

Now we wish to revisit our cousin of No Retraction Theorem. Suppose we have $A, B, C$ as stated. Cut up the triangle as in Sperner's Lemma. Label each vertex $\mathbf{x}$ to be R, G or B according to whether it lies in $A$, $B$ or $C$. If it lies in multiple, we pick arbitrarily, except for on the boundary, where we make sure to choose such that the boundaries are coloured in the way that Sperner's Lemma requires.

Now Sperner's Lemma tells us that there exist vertices $\mathbf{a}_n \in A$, $\mathbf{b}_n \in B$, $\mathbf{c}_n \in C$ which form a small triangle. If we cut up small enough, we can make the pairwise distances between $\mathbf{a}_n, \mathbf{b}_n, \mathbf{c}_n$ be at most $\frac{1}{n}$. By compactness, there exists $\mathbf{a} \in \overline{\triangle}$ and $n(j) \to \infty$ such that $\mathbf{a}_{n(j)} \to \mathbf{a}$. Automatically, $\mathbf{b}_{n(j)}, \mathbf{c}_{n(j)} \to \mathbf{a}$. Apply closure to deduce $\mathbf{a} \in A \cap B \cap C$.

Now, we can simply follow our chain of equivalences backwards to deduce Brouwer.

### * Non-examinable material

Brouwer came to disbelieve his theorem. Look at our 3 colour theorem (the cousin of No Retraction Theorem). It says that there is a point $\mathbf{x} \in A \cap B \cap C$, but it gives no recipe for finding such an $\mathbf{x}$.

** This is the end of the non-examinable comments.

> **Lemma.** If $A = (a_{ij})_{\substack{1 \leq i \leq 3 \\ 1 \leq j \leq 3}}$ (a $3 \times 3$ matrix). Suppose $a_{ij} \geq 0$ and $\sum_{i=1}^{3} a_{ij} = 1$ for $i = 1, 2, 3$. Then there exists $x_1, x_2, x_3 \geq 0$ not all zero such that $\sum_i a_{ij}x_j = x_i$ (i.e. $\mathbf{x}$ is an eigenvector with eigenvalue 1).

*Proof.* Consider $\overline{\triangle} = \{\mathbf{x} \in \mathbb{R}^3 \mid x_1 + x_2 + x_3 = 1, x_j \geq 0\}$. If $T\mathbf{x} = A\mathbf{x}$, then $y_i = \sum_j a_{ij}x_j$. Then $y_i \geq 0$ and

$$\sum_i y_i = \sum_i \sum_j a_{ij}x_j = \sum_j \sum_i a_{ij}x_i = \sum_j x_j = 1$$

So $\mathbf{y} \in \overline{\triangle}$, and $T$ is a continuous map $\overline{\triangle}$ to $\overline{\triangle}$, so has a fixed point. $\square$

### Nash on economics

Nash did two things. First is study of non-zero sum games.

$A$ choose strategy 1 with probability $p$, and chooses strategy 2 with probability $1 - p$. $B$ chooses strategy $1'$ with probability $q$, $2'$ with probability $1 - q$. Expected value to $A$ is

$$A(p, q) = \sum_{ij} a_{ij}p_i q_j$$

$$B(p, q) = \sum_{ij} b_{ij}p_i q_j$$

Start of

Two people two outcomes:

$$A(\mathbf{p}, \mathbf{q}) = \sum p_i a_{ij} q_j$$
$$B(\mathbf{p}, \mathbf{q}) = \sum_i p_i a_{ij} q_j$$

Von Neumann zero sum games: $a_{ij} = -b_{ij}$. Then can show there exists $\mathbf{p}^*, \mathbf{q}^*$ such that

$$A(\mathbf{p}^*, \mathbf{q}^*) = \sup_f \inf_\mathbf{q} A(\mathbf{p}^*, \mathbf{q})$$
$$B(\mathbf{p}^*, \mathbf{q}^*) = \sup_f \inf_{\mathbf{p}^*} B(\mathbf{p}^*, \mathbf{q}^*)$$

|  | you admit | you don't admit |
|---|:---:|:---:|
| he admits | 2 | 8 |
| he does not admit | $\frac{1}{2}$ | 1 |

Nash showed that there is at least one $(\mathbf{p}^*, \mathbf{q}^*)$ such that

$$A(\mathbf{p}^*, \mathbf{q}^*) \geq A(\mathbf{p}, \mathbf{q}^*) \qquad \forall \mathbf{p}$$
$$B(\mathbf{p}^*, \mathbf{q}^*) \geq B(\mathbf{p}^*, \mathbf{q}) \qquad \forall \mathbf{q}$$

That is to say there is no reason unilaterally to change your choice.

*Proof.* Use Brouwer. Let

$$\Gamma = \{(p, 1-p, q, 1-q) : 1 \geq p \geq 0, 1 \geq q \geq 0\}$$

2 dimensional square. We will define $T : \Gamma \to \Gamma$. First define $u_1$:

$$u_1 = \max(A(1, 0, q_1, q_2) - A(p_1, p_2, q_1, q_2), 0)$$

$u_1 > 0$ if worth moving towards $(1, 0)$ for $\mathbf{p}$. $u_1$ is continuous. Define

$$u_2 = \max(A(0, 1, \mathbf{q}) - A(\mathbf{p}, \mathbf{q}), 0)$$

and similarly define $v_1, v_2$ for $q$.

Then we define

$$T(\mathbf{p}, \mathbf{q}) = \left( \frac{p_1 + u_1}{1 + u_1 + u_2}, \frac{p_2 + u_2}{1 + u_1 + u_2}, \frac{q_1 + v_1}{1 + v_1 + v_2}, \frac{q_2 + v_2}{1 + v_1 + v_2} \right)$$

$u_1, u_2$ continuous, and $u_1, u_2 \geq 0$ gives that

$$\frac{p_1 + u_1}{1 + u_1 + u_2}$$

is a continuous function of $(\mathbf{p}, \mathbf{q})$. Also, $\frac{u_1 + p_1}{1 + u_1 + u_2} \geq 0$ and

$$\frac{u_1 + p_1}{1 + u_1 + u_2} + \frac{u_2 + p_2}{1 + u_1 + u_2} = \frac{u_1 + u_2 + p_1 + p_2}{1 + u_1 + u_2} = 1.$$

Similarly for $u_2, v_1, v_2$. So $T$ is a continuous map from $\Gamma$ to $\Gamma$. Thus there is at least one fixed point. Since at most one of $u_1, u_2$ is non-zero, we may suppose $u_2 = 0$. But $T(\mathbf{p}^*, \mathbf{q}^*) = (\mathbf{p}^*, \mathbf{q}^*)$, so $u_1 = u_2 = 0$ so either $1 > p_1 > 0$ and $A(1, 0, \mathbf{q}^*) = A(0, 1, \mathbf{q}^*) = A(\mathbf{p}, \mathbf{q})$ so no reasn for $A$ to change (note $A(p, q)$ is affine in $p_1$). Otherwise, $p_1 = 0$ or $p_1 = 1$. Without loss of generality, $A(1, 0, \mathbf{q}) \geq A(1 - p', p', \mathbf{q})$, but we cannot. $\qquad \square$

---

**Example** (Game of chicken). $A$ and $B$ drive cars towards each other. If both swerve, they both lose 1 prestige point. If one swerves and the other does not, then the swerver loses 5 points and the swerver gains 10 points. If neither swerves, then both lose 100 points.

$$A(p, q) = -pq - 5p(1 - q) + 10(1 - p)q - 100(1 - p)(1 - q)$$

Easiest to differentiate:

$$\frac{\partial A}{\partial p} = -q - 5(1 - q) - 10q + 100(1 - q)$$
$$= 95 - 106q$$

so if $\frac{\partial A}{\partial p} =$, then $106q = 95$, so $q = \frac{95}{106}$. This tells us that $\left(\frac{95}{106}, \frac{95}{106}\right)$ is a Nash point.

However, we have only examined interior points (because a local extrema is only found by $\frac{\partial A}{\partial p}$ if it is not on the boundary). We must check $p = 1$ for example. Upon doing this, we also find that $(1, 0)$ and $(0, 1)$ are Nash points.

---

**Remark.** Recall that Brouwer has a multidimensional extension. It is quite easy to extend from 2 people, 2 choices to 2 people, $n$ choices. It can also be extended to more than 2 players (which Von Neumann can't). However, what Nash then says is that there exist points $(p^*, q^*, v^*)$ where it is to no player's advantage to change unilaterally. This is much weaker than it sounds. Winkn, Blykn and Nod who must divide £90 by majority vote.

Start of

lecture 8

Previously we dealt with untrustworthy individuals. There are occasions when people do act trustworthy.

For example, contracts enforcable by law, or individuals whomust make repeated transactions with each other.

With high symmetry, it is not too hard. For example, if Wynken, Blykn and Noel need to divide \$90, then they will just split it evenly, if they are incentivised to behave fairly.

In our game of chicken with rules as set out before, we could do "players swerve alternately".

More different if there is no obvious symmetry. "What do we then mean by fair".

Nash has an interesting model. $n$ participants. $E \subseteq \mathbb{R}^n$. Think of $\mathbf{x} \in E$ as a choice. If $x_j > y_j$, then $j$ preferse $\mathbf{x}$ to $\mathbf{y}$.

Mathematical condition: $E$ is closed and bounded (natural if we want there to exist some kind of 'best').

More interesting:

(1) Demand $E$ be convex, i.e. $\mathbf{x}, \mathbf{y} \in E$, $0 \le \lambda \le 1$, then $\lambda \mathbf{x} + (1 - \lambda)\mathbf{y} \in E$ (for mathmos, choose $\mathbf{x}$ with probability $\lambda$, $\mathbf{y}$ with probability $1 - \lambda$; for non-mathmos, exchange sums of cash or other horse trading).

(2) There is a status quo point $\mathbf{s} \in E$. If no agreement reached then the outcome is $\mathbf{s}$. Could be "continue as before" or "strike" or "lock out". Could consider $E \cap \{\mathbf{x} : x_j \ge s_j \; \forall j\}$.

(3) Pareto optimality. It is a well kept secret amongst politicians and economists that function $f : X \to \mathbb{R}^n$ do not have maxima for $n \ge 2$. Pareto says that we should aim for a "Pareto optimum" $x^*$ such that $\nexists j$ and $\mathbf{z} \in E$ such that $x_j^* > z_j$ and $x_k^* \ge z_k$ for all $k \ne j$.

(4) Independence of irrelevant conditions. If $(E, s)$, $(E', s)$ given and $E' \supseteq E$, if $x^*$ is chosen for $(E', s)$ and $\mathbf{x}^* \in E$ then $\mathbf{x}$ is chosen for $(E, s)$.

(5) Symmetry: Suppose $(x_1, \dots, x_n) \in E$, $\sigma \in S_n$ implies $(x_{\sigma(1)}, \dots, x_{\sigma(n)} \in E$ and $\mathbf{s} = (s, \dots, s)$, then our solution $\mathbf{x}^*$ will have $x_1^* = x_2^* = \cdots = x_n^*$.

(6) You can't improve things by exaggeration. If $T(x, x_n) = (a_1 x_1 + b_1, a_2 x_2 + b_2, \dots, a_n x_n + b_n)$ then if $x^*$ is the outcome for $(E, \mathbf{s})$, then $Tx^*$ is the outcome for $(TE, T\mathbf{s})$.

> **Lemma.** If $E$ closed and convex. Suppose $(1, 1, \dots, 1) \in E$ and $\prod_{j=1}^n x_j \le 1$ for all $\mathbf{x} \in E$. Then $x_1 + x_2 + \cdots + x_n \le n$ for all $\mathbf{x} \in E$.

*Proof.* Suppose $\mathbf{x} \in E$. Then since $\mathbf{1} \in E$, convexity gives

$$(1 - \delta)\mathbf{1} + \delta\mathbf{x} \in E$$

so $\prod((1 - \delta) + \delta x_j) \leq 1$, i.e.

$$\prod(1 + \delta(x_j - 1)) \leq 1.$$

So $1 + \sum \delta(x_j - 1) + A(\delta) \leq 1$ with $\frac{A(\delta)}{\delta} \to)$ as $\delta \to 0$. Then $\sum(x_j - 1) + \frac{A(\delta)}{\delta} \leq 0$, so $\sum(x_j - 1) \leq 1$. $\qquad\square$

> **Corollary.** If $\mathbf{s} = 0$, $(1, 1, \ldots, 1) \in E$ and $\prod x_j \leq 1$ for all $x_j \in E$, then $(1, 1, \ldots, 1)$ is the choice.

*Proof.* By our previous lemma, $E \subseteq \Gamma = \{x_j \leq n\}$. Look at our problem for $(\Gamma, \mathbf{0})$. This a symmetric problem (in the sense of (5)) so solution must satisfy $x_1 = x_2 = \cdots = x_n$. By Pareto optimality, $x_1 = x_2 = \cdots = x_n = 1$. So we have a unique solution. But $E \subseteq \Gamma$ and $(1, 1, \ldots, 1) \in E$ so by (4) independence of irrelevant conditions, we have that $(1, 1, \ldots, 1)$ is optimal for $(E, \mathbf{0})$ problem. $\qquad\square$

We now use condition (6) (problem invariant under affine transformations $\mathbf{x} \mapsto (a_1 x_1 + b_1, \ldots, a_n x_n + b_n)$, with $a_j > 0$).

Suppose $(E, \mathbf{s})$ given claim solution is unique $x^*$ which maximises $\prod_{i=1}^n (x_i - s_i)$ subject to $\mathbf{x} \in E$. By compactness of $E$, such a point exists.

$$\prod_{i=1}^n (x_i^* - s_i) \geq \prod_{i=1}^n (x_i - s_i) \qquad \forall x \in X.$$

(We assume $\prod_{i=1}^n (x_i^* - s_i) > 0$). If we make the transformation $z_i = \frac{(x_i - s_i)}{(x_i^* - s_i)}$.

Start of

lecture 9

# 3 Approximation by polynomials

Taylors Theorem is not the entire answer. 2 examples:

**Lemma.** If $E : \mathbb{R} \to \mathbb{R}$ is given by

$$E(x) = \begin{cases} \exp\left(-\frac{1}{x^2}\right) & x \neq 0 \\ 0 & x = 0 \end{cases}$$

then $E$ is infinitely differentiable everywhere, with $E^{(r)}(0) = 0$ for all $r$. So the Taylor expansion about 0 is

$$\sum_r 0x^r = 0 \nrightarrow E(r)$$

for $x \neq 0$.

*Proof.* If $x \neq 0$ then standard differentiation theorems give that $E$ is infinitely differentiable with $E^{(r)}(x) = Q_r(1/x)E(x)$ where $Q_r$ is a polynomial (proved by induction). At 0, $E(0)$. Suppose $E^{(r)}$ exists with $E^{(r)}(0) = 0$. Then

$$\frac{E^{(r)}(t) - E^{(r)}(0)}{t} = \frac{1}{t}Q_r\left(\frac{1}{t}\right)E(t) \to 0.$$

(Exponential beats polynomials). So $E^{(r+1)}(0)$ exists and equals 0. $\qquad\square$

The second problem is practical. It is true that

$$\exp x = \sum_{r=0}^{\infty} \frac{x^r}{r!} \qquad \forall x$$

but computing $\exp(-20)$ by using

$$\exp(-20) \approx \sum_{r=0}^{N} \frac{(-20)^r}{r!}$$

involves large $N$ and ridiculous cancelation.

One attempt to get polynomial approximation is to use interpolation:

> **Lemma.**
>
> (i) $\mathcal{P}_n$ the collection of polynomials of degree at most $n$ is a vector space of dimension $n + 1$.
>
> (ii) If $f \in C[a, b]$, $x_0 < x_1 < \cdots < x_n$, then there is at most one polynomial $P \in \mathcal{P}_n$ such that $f(x_j) = P(x_j)$.
>
> (iii) If $e_j(x) = \prod_{x \neq j} \frac{x - x_j}{x_i - x_j}$ then these form a basis of $\mathcal{P}_n$, and writing
>
> $$P(x) = \sum_j f(x_j) e_j(x)$$
>
> we have $P \in \mathcal{P}_n$, $P(x_j) = f(x_j)$, so in fact by (ii) there is exactly one polynomial.

*Proof.*

(i) $1, t, t^2, \ldots, t^n$ is a basis for $\mathcal{P}_n$.

(ii) Suppose $f(x_j) = P(x_j) = Q(x_j)$, $[0 \leq j \leq n]$, $P, Q \in \mathcal{P}_n$. Then $P - Q \in \mathcal{P}_n$ and has $n + 1$ zeroes (given by $x_0, \ldots, x_n$), so $P - Q = 0$.

(iii) $e_j(x_i) = \delta_{ij}$, $e_j \in \mathcal{P}_n$. So if $P(x) = \sum_j f(x_j) e_j(x)$, then $P(x_i) = f(x_i)$ for $0 \leq i \leq n$.

$\square$

Practice shows interpolation to be an unreliable friend.

Chebychev introduced an interesting polynomial $T_n$ which shows how oddly polynomials can behave (borhter noticed a companion $U_n$ Chebychev polynomial of the second kind).

Chebychev says look at:

$$(\cos n\theta + i\sin n\theta) = (\cos\theta + i\sin\theta)^n$$

$$= \sum_r \binom{n}{r}(\cos\theta)^{n-r}i^r(\sin\theta)^r$$

$$= \sum_r (-1)^r \binom{n}{2r}(\cos\theta)^{n-2r}(\sin\theta)^{2r}$$

$$+ i\sum_r \binom{n}{2r+1}(\cos\theta)^{n-2r-1}(\sin\theta)^{2r+1}$$

$$= \sum_r (-1)^e \binom{n}{2r}(\cos\theta)^{n-2r}(1-\cos^2\theta)^r$$

$$+ i\sum_r (-1)^r \binom{n}{2r+1}(\cos\theta)^{n-2r-1}(\sin\theta)(1-\cos^2\theta)^r$$

so taking real and imaginary parts we get

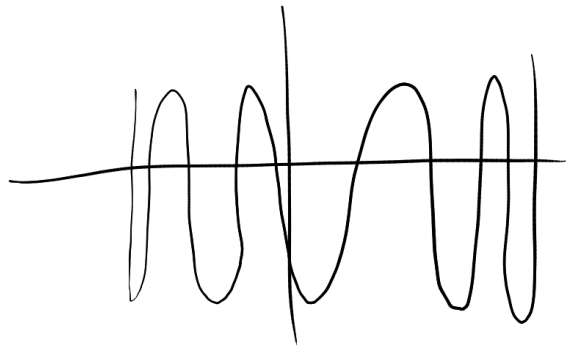$$\cos n\theta = T_n(\cos\theta), \qquad \sin n\theta = (\sin\theta)U_n(\cos\theta)$$

with

$$T_n(t) = \sum_r (-1)^r \binom{n}{2r}t^{n-2r}(1-t^2)^r$$
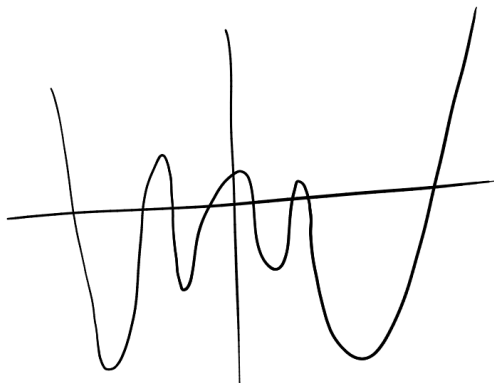
$$U_n(t) = \sum_r (-1)^r \binom{n}{2r+1}t^{2n-2r-1}(1-t^2)^r$$

$T_n \in \mathcal{P}_n$, $|T_n(t)| \le 1$ for all $t \in [-1,1]$. Leading coefficient of $T_n$ is

$$\sum_{0 \le 2r \le n} \binom{n}{2r}(-1)^r = \frac{1}{2}((1+1)^n - (1-1)^n) = 2^{n-1}$$

zeroes of $T_n$ are $\cos^{-1}(r\pi/n)$ $(T_n(\cos\theta) = \cos n\theta)$.

$\sin n\theta = U_n(\cos\theta)\sin\theta$, so $U_n(\cos\theta) = \frac{\sin n\theta}{\sin\theta}$. (if $\sin\theta = 0$, use continuity to get $U_n = \pm n$ at this point).



In spite of all this, Weierstrass showed:

**Theorem** (Weierstrass approximation). The polynomials are uniformly dense in $C[a,b]$, i.e. given $f : [a,b] \to \mathbb{R}$ continuous and $\varepsilon > 0$, there exists polynomial $P$ such that $|f(t) - P(t)| \le \varepsilon$ for all $t \in [a,b]$.

Bernstein:

**Theorem** (Bernstein). If $f : [0,1] \to \mathbb{R}$ then

$$\sum_{r=0}^{n} f\left(\frac{r}{n}\right)\binom{n}{r} t^r(t-t)^{n-r} \to f(t)$$

uniformly on $[0,1]$ as $n \to \infty$. (Can get from $[0,1]$ to $[a,b]$ by a scaling transformation).

One proof depends on reinterpreting the theorem probabilistically. Let $X_1, X_2$ be independent such that $\mathbb{P}(X_j = 1) = p$, $\mathbb{P}(X_j = 0) = 1 - p$. Then

$$\sum_{r=0}^{n} f\left(\frac{r}{n}\right)\binom{n}{r} p^r(1-p)^{n-r} = \mathbb{E}f(\overline{X})$$

where $\overline{X} = \frac{X_1 + \cdots + X_n}{n}$.

Last time we introduced Bernstein's version of the Weierstrass polynomial approximation theorem:

> **Theorem.** If $f : [0, 1] \to \mathbb{R}$ is continuous then writing
> $$P_n(p) = \sum_r \binom{n}{r} f\left(\frac{r}{n}\right) p^r (1 - p)^{n-r}$$
> then $P_n \to f$ uniformly.

We use a probabilistic interpretation and Chebyshev inequality:

> **Lemma.** If $X$ is a bounded random variable, then
> $$\mathbb{P}(|X - \mathbb{E}X| \geq a) \leq \frac{\mathrm{Var}(X)}{a^2}.$$

*Proof.* See Probability from IA. $\qquad\square$

*Proof of Bernstein's version of Weierstrass polynomial approximation theorem.* Now consider $X_1, X_2, \ldots, X_n$ independent random variables with $\mathbb{P}(X_j = 1) = p$, $\mathbb{P}(X_j = 0) = 1 - p$. Then

$$\mathrm{Var}(X_j) = \mathbb{E}((X - \mathbb{E}X)^2) = \mathbb{E}((X - p)^2) = (1 - p)p^2 + p(1 - p)^2 = p(1 - p) \leq \frac{1}{4}$$

by AM-GM. Then

$$\mathrm{Var}\left(\frac{X_1 + \cdots + X_n}{n}\right) = \frac{1}{n^2} \sum_j \mathrm{Var}(X_j) \leq \frac{1}{4n}.$$

Now consider the function $f$. By compactness there exists an $M$ such that $|f(p)| \leq M$ for all $p \in [0, 1]$. Also, $f$ is uniformly continuous – that is to say given $\varepsilon > 0$, there exists $\delta(\varepsilon) > 0$ such that
$$|s - t| \leq \delta \implies |f(s) - f(t)| < \varepsilon$$

We fix $\varepsilon > 0$, $\delta(\varepsilon) > 0$ through the proof. Then we say that we can make $\varepsilon$ as small as
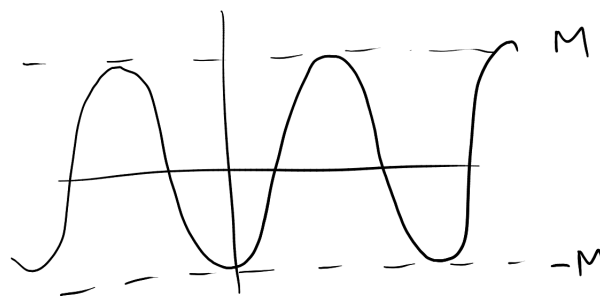
we want.

$$
\begin{aligned}
|P_n(p) - f(p)| &= |\mathbb{E}((\overline{X})) - f(p)| \\
&= |\mathbb{E}(f(\overline{X}) - f(p))| \\
&= \sum_{k=0}^{n} \mathbb{P}\left(\overline{X} = \frac{r}{n}\right) \left| f\left(\frac{r}{n}\right) - f(p) \right| \\
&= \sum_{|\frac{r}{n} - p| \leq \delta} + \sum_{|\frac{r}{n} - p| > \delta} \\
&\leq \sum_{|\frac{r}{n} - p| \leq \delta} \varepsilon \mathbb{P}\left(\left|\overline{X} - \frac{r}{n}\right| \leq p\right) + \sum_{|\frac{r}{n} - p| > \delta} M\mathbb{P}\left(\left|\overline{X} - \frac{r}{n}\right| = p\right) \\
&\leq \varepsilon + 2M\mathbb{P}\left(\left|\overline{X} - \frac{r}{n}\right| \geq \delta\right) \\
&\leq \varepsilon + 2M\frac{\mathrm{Var}(\overline{X})}{n}\frac{1}{\delta^2} \qquad\qquad\qquad \text{(Chebyshev)} \\
&\leq \varepsilon + \frac{2M}{4}\frac{1}{n} \\
&< 2\varepsilon
\end{aligned}
$$

if $n$ is large enough. $\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\qquad\quad$ $\square$

## Best uniform approximation

**Theorem** (Chebyshev equiripple criterion). If $f : [a, b] \to \mathbb{R}$ is continuous and $P$ is a polynomial of degree at most $n$ such that there exists $M \geq 0$ and $a = a_0 < a_1 < \cdots < a_n < a_{n+1} \leq b$ and $P(a_j) - f(a_j) = (-1)^j M$ (or $P(a_j) - f(a_j) = (-1)^{j+1}M$) and such that $\|P - f\|_\infty$. Then $\|P - f\|_\infty \leq \|Q - f\|_\infty$ for all polynomials $Q$ of degree $n$ or less.

*Proof.* Without loss of generality $P(a_j) - f(a_j) = (-1)^j M$. If $Q \in \mathcal{P}_n$ and $\|P - f\|_\infty > \|Q - f\|_\infty$ then $P(a_{2j} - f) \geq Q(a_{2j}) - f$, so $P(a_{2j}) > Q(a_{2j})$. Similarly for $P(a_{2j+1}) < Q(a_{2j+1})$. So there exists $c_j$ with $a_0 < c_0 < a_1 < c_1 < a_2 < \cdots < c_n < a_{n+1}$ with $(P - Q)(c_j) = 0$ (intermediate value theorem). So $P - Q$ has $n + 1$ zeroes, but $P - Q \in \mathcal{P}_n$, so $P - Q = 0$, contradiction. $\qquad\square$

Now $\|T_n\|_\infty \leq 1$ and $T_n$ is alternately $\pm 1$ at $n + 1$ points. The coefficient of $t^n$ in $T_n$ is $2^{-n+1}$ (for $n \geq 1$). So $2^{n-1}T_n(t) = t^n - Q_{n-1}(t)$ with $Q_{n-1}$ then best uniform approximation to $t^n$.

> **Corollary.** There exists an $\varepsilon_n$ such that if $P(T) = \sum_{j=0}^n a_j t^j$ and $\exists k$ such that $|a_k| \geq 1$, then $|P(t)| \geq \varepsilon_n \ \forall t \in [-1, 1]$.

*Proof.* Proof by induction. True for $n = 0, 1$ by inspection. Now suppose true for $n = N$.

$$P(t) = \sum_{j=0}^{N+1} a_j t^j = a_{N+1} t^{N+1} + Q_N(t).$$

Say $Q_N(t) = \sum_{j=0}^N a_j t^j$. If $|a_{N+1}| < \frac{\varepsilon_N}{2}$, then

$$\|P\|_\infty \geq \|Q\|_\infty - \frac{\varepsilon_N}{2} \geq \frac{\varepsilon_N}{2}.$$

If $|a_{N+1}| > \frac{\varepsilon_N}{2}$ and we know from Chebyshev that

$$|P(t)| \geq \frac{\varepsilon_N}{2} \inf_{Q \in \mathcal{P}_n} |Q(t) - t^{N+1}| = \varepsilon_N' > 0. \qquad\square$$

During the next section we shall switch between various norms on $\mathbb{R}^n$:

$$\|\mathbf{a}\|_\infty = \max |a_j|, \qquad \|\mathbf{a}\|_2 = \sqrt{\sum_j |a_j|^2}, \qquad \|\mathbf{a}\| = \sum_j |a_j|.$$

There is a general theorem that all norms on $\mathbb{R}^n$ are "equivalent":

$$\|\mathbf{a}\|_\infty \leq \|\mathbf{a}\|_1 \leq n\|\mathbf{a}\|_\infty$$
$$\|\mathbf{a}\|_\infty \leq \|\mathbf{a}\|_2 \leq n\|a\|_\infty$$

Last time we showed that there exists an $\varepsilon(n) > 0$ such that if $\mathbf{a} \in \mathbb{R}^{n+1}$ and $\|\mathbf{a}\|_\infty \geq 1$ then

$$- \sup_{t \in [a,b]} \left| \sum_{j=0}^n a_j t^j \right| \geq \varepsilon(n).$$

Thus if we write $T\mathbf{a} = P$ with $P(t) = \sum_j a_j t^j$, we have $\|T\mathbf{a}\|_\infty \to \infty$ as $\|\mathbf{a}\|_2 \to \infty$.

We now show:

> **Lemma.** If $f \in C[a, b]$, then there exists a $P \in \mathcal{P}_n$ such that
> $$\|f - P\|_\infty \leq \|f - Q\|_\infty$$
> for all $Q \in \mathcal{P}_n$, i.e. ther exists a *best* polynomial approximation (for fixed degree).

*Proof.* Without loss of generality $[a, b] = [0, 1]$.

Consider the map $S : \mathbb{R}^{n+1} \to \mathbb{R}$ given by

$$S(\mathbf{a})(t) = \left| \sum_{j=0}^n a_j t^j - f(t) \right|.$$

We claim that $S$ is continuous.

$$
\begin{aligned}
|S(\mathbf{a})(t) - S(\mathbf{b})(t)| &= \left\| \left| \sum_{j=0}^n a_j t^j - f(t) \right| - \left| \sum_{j=0}^n b_j t^j - f(t) \right| \right\| \\
&\leq \left| \sum_j (a_j t^j - f(t)) - \sum_j (b_j t^j - f(t)) \right| \\
&= \left| \sum_j (a_j - b_j) t^j \right| \\
&\leq \sum_j |a_j - b_j|
\end{aligned}
$$

so

$$\|S(\mathbf{a}) - S(\mathbf{b})\| \leq \sum_j |a_j - b_j|.$$

So $S$ is continuous as a map $(\mathbb{R}^{n+1}, \| \bullet \|_2) \to \mathbb{R}$. Now $|S(\mathbf{a})| \to \infty$ as $\|\mathbf{a}\|_2 \to \infty$ by our previous result. So we can find an $R$ such that $|S(\mathbf{a})| \geq S(\mathbf{0})$ for all $\|\mathbf{a}\|_2 \geq R$. $S$ is continuous on the compact set $\overline{B(0, R)}$ (closed ball), so has a minimum at some $\mathbf{c} \in \overline{B(0, R)}$. $S(c) \leq S(\mathbf{a}) \ \forall \|\mathbf{a}\|_2 \leq R$,

$$S(\mathbf{c}) \leq S(\mathbf{0}) \leq S(\mathbf{a}) \qquad \forall \|\mathbf{a}\|_2 \geq R. \qquad \square$$

**\* Non-examinable material**

Look at the attained minimum $(S - f)t$. Suppose that it does not satisfy the equiripple criterion. Then we can perturb it a little to improve the approximation, which is a contradiction. This shows that the equiripple is in fact a *necessary* condition (not just a sufficient condition).

\*\* This is the end of the non-examinable comments.

# 4 Gaussian Quadrature

**Problem:** Given $f(x_1), \ldots, f(x_n)$, we would like to estimate

$$\int_0^b f(t)\mathrm{d}t.$$

First idea is to use interpolation.

> **Lemma.** If $x_1, \ldots, x_n \in [a, b]$ all distinct. Then there exists unique $A_1, \ldots, A_n$ such that
> $$\int_A^b P(t)\mathrm{d}t = \sum_j A_j P(x_j)$$
> for all $P \in \mathcal{P}_{n-1}$.

So we hope that

$$\int_a^b f(t)\mathrm{d}t \stackrel{?}{\approx} \sum_j A_j f(x_j)$$

*Proof.* Recall the notation

$$e_j(x) = \prod_{i \neq j} \frac{(x - x_i)}{(x_j - x_i)},$$

so that $e_i(x_j) = \delta_{ij}$. If $P \in \mathcal{P}_{n-1}$, we know that

$$P(t) = \sum_j P(x_j) e_j(t).$$

So

$$\int_a^P (t)\mathrm{d}t \sum_j P(x_j) A_j$$

with

$$A_j = \int_a^b e_j(t) = \mathrm{d}t.$$

Conversely, if $\int_a^b P(t)\mathrm{d}t = \sum_j B_j P(x_j)$ for all $P \in \mathcal{P}_n$, then we have

$$\int_a^b e_i(t)\mathrm{d}t = \sum_j B_j e_i(x_j) = B_i.$$

So $B_i = A_i$. $\qquad\square$

Without loss of generality $[a,b] = [0,1]$. If we choose $x_r = \frac{r}{n}$ and look at $x_0, \ldots, x_n$ and $n$ is small we get quite good results. But as $n$ increases, the associated $A_j$ rapidly take very large values (allowed by the $A_j$ taking positive and negative values).

Gauss says look at Legendre polynomials.

Recall Gramm-Schmidt:

If $\mathbf{e}_1, \ldots, \mathbf{e}_k$ are orthonormal vectors in an inner product space and $\mathbf{u} \notin \mathrm{span}\{\mathbf{e}_1, \ldots, \mathbf{e}_k\}$, we can find $\mathbf{e}_{k+1}$ such that $\mathbf{e}_1, \ldots, \mathbf{e}_{k+1}$ are orthonormal and $\mathbf{u} \in \mathrm{span}\{\mathbf{e}_1, \ldots, \mathbf{e}_{k+1}\}$.

Recall also that $C([0,1])$ is an inner product space if we write

$$\langle f, g \rangle = \int_0^1 f(t)g(t)\mathrm{d}t.$$

Thus we can find $P_0, P_1, P_2, \ldots$ with $P_j \in \mathcal{P}_j$ and $P_0, P_1, \ldots$ orthonormal, and so we can find unique polynomials $p_0, p_1, \ldots$ such that $p_n$ has degree $n$ and positive leading coefficient, with $p_0, p_1, \ldots, p_n$ orthonormal. We call the $p_j$ the *Legendre polynomials*.

> **Lemma.** The Legendre polynomial $p_n$ has all its roots simple, real and lying in $[0,1]$.

*Proof.* Let $\beta_1, \ldots, \beta_m$ be the roots of odd order lying in $[0,1]$ (odd order means $p_n$ changes sign as $t$ passes through $\beta_j$). Then writing $Q(t) = \prod_j (t - \beta_j)$ we have $Q(t)p_n(t)$ is single signed and non-zero. So

$$\int_0^1 Q(t)p_n(t) \neq 0$$

so $\deg Q \geq n$, so $m \geq n$. But $p_n$ is of degree $n$, so all roots lie in $[0,1]$ and are simple. $\square$

> **Theorem.** Let $\alpha_1, \alpha_2, \ldots, \alpha_n$ be the distinct roots of $p_n$ and $A_j$ the unique constants such that
> $$\int_{-1}^1 P(t)\mathrm{d}t = \sum_j A_j P(\alpha_j)$$
> for all $P \in \mathcal{P}_{n-1}$. Then in fact
> $$\int_{-1}^1 P(t)\mathrm{d}t = \sum_{j=1}^n A_j P(\alpha_j)$$
> for all $P \in \mathcal{P}_{2n-1}$.

*Proof.* Suppose $P \in \mathcal{P}_{2n-1}$. By long division, write $P = Qp_n + R$ with $Q \in \mathcal{P}_{n-1}$ and $R \in \mathcal{P}_{n-1}$. Then

$$\begin{aligned}
\int_{-1}^{1} P \mathrm{d}t &= \int_{-1}^{1} Qp_n + R \mathrm{d}t \\
&= \int Qp_n \mathrm{d}t + \int R \mathrm{d}t \\
&= \int R \mathrm{d}t && \text{(since } p_n \perp Q) \\
&= \sum_j A_j R(\alpha_j) \\
&= \sum_j A_j (Q(\alpha_j)p(\alpha_j) + R(\alpha_j)) \\
&= \sum_j A_j P(\alpha_j) && \square
\end{aligned}$$

> **Remark.** If $\beta_1, \ldots, \beta_n, B_1, \ldots, B_n$ are such that
>
> $$\int_{-1}^{1} P(t)\mathrm{d}t = \sum_j B_j P(\beta_j)$$
>
> for all $P \in \mathcal{P}_{2n-1}$, then taking $F(t) = \prod_j (t - \beta_j)$ we have
>
> $$\begin{aligned}
> \int_{-1}^{1} F(t)Q(t)\mathrm{d}t &= \sum_j B_j F(\beta_j)Q(\beta_j) \\
> &= \sum_j B_j 0 \\
> &= 0
> \end{aligned}$$
>
> for all $Q \in \mathcal{P}_{n-1}$, using the fact that $FQ \in \mathcal{P}_{2n-1}$. So $F \perp Q$ for all $Q \in \mathcal{P}_{n-1}$. So $F$ is a scalar multiple of $p_n$, so $\{\beta_1, \ldots, \beta_n\} = \{\alpha_1, \ldots, \alpha_n\}$.

This is *not* the main point. The key observation is that for $\alpha$ the zeroes of $p_n$ and $A_j$ as before, we have:

(1) $A_j > 0$

(2) (less important) $\sum_j A_j = 2$.

*Proof.*

(1) Consider $Q_j(t) = \prod_{i \neq j}(t - \alpha_i)^2$. $Q_j$ is non-constant and positive, so $\int Q > 0$ and $Q \in \mathcal{P}_{2n-2} \subseteq \mathcal{P}_{2n-1}$. So

$$0 < \int_{-1}^{1} Q_j(t) \mathrm{d}t = \sum_i A_i Q_i(\alpha_j) = A_j Q(\alpha_j)$$

so $A_j > 0$.

(2) $2 = \int_{-1}^{1} 1 \mathrm{d}t = \sum_{j=1}^{n} A_j$. $\qquad \square$

---

**Corollary.** If $f \in C[-1,1]$ and $\|P - f\|_\infty \leq \varepsilon$, $P \in \mathcal{P}_{2n-1}$. Then

$$\left| \int_{-1}^{1} f(t) \mathrm{d}t - \sum_j A_j f(\alpha_j) \right| \leq 4\varepsilon.$$

---

*Proof.*

$$\left| \int f - \sum_j A_j f(\alpha_j) \right| = \left| \int f - \int P + \sum A_j P(\alpha_j) - \sum A_j f(\alpha_j) \right|$$

$$\leq \left| \int_{-1}^{1} (f - P) \right| + \sum |A_j||P(\alpha_j) - f(\alpha_j)|$$

$$\leq 2\|f - P\|_\infty + 2\|f - P\|_\infty$$

$$= 4\|f - P\|_\infty$$

$$\leq 4\varepsilon \qquad \square$$

---

**Corollary.** By Weierstrass's theorem, Gaussian interpolation of higher and higher degree converges to the correct answer.

---

# 5 Hausdorff Metric (and simpler things)

[ALL SETS IN THIS SECTION ARE NON-EMPTY UNLESS OTHERWISE STATED].

Recall in $\mathbb{R}^n$, if $A$ is closed and non-empty, then

$$d(\mathbf{x}, A) = \inf_{\mathbf{a} \in A} \|\mathbf{x} - \mathbf{a}\|$$

is well-defined and $\mathbf{x} \mapsto d(\mathbf{x}, A)$ is continuous.

*Proof of continuity.* If $\mathbf{x}, \mathbf{y}$ given. For any $\varepsilon > 0$, there exists $\mathbf{a} \in A$ such that $\|\mathbf{x} - \mathbf{a}\| \leq \varepsilon + d(\mathbf{x}, A)$. Then

$$\|\mathbf{y} - \mathbf{a}\| \leq \|\mathbf{x} - \mathbf{y}\| + \|\mathbf{x} - \mathbf{a}\|\| \leq \|\mathbf{x} - \mathbf{y}\| + \varepsilon + d(x, A).$$

$\varepsilon$ is arbitrary, so $\|\mathbf{y} - \mathbf{a}\| \leq \|x - y\| + d(x, A)$. So

$$d(y, A) \leq d(x, y) + d(x, A).$$

Similarly
$$d(x, A) \leq d(x, y) + d(y, A)$$
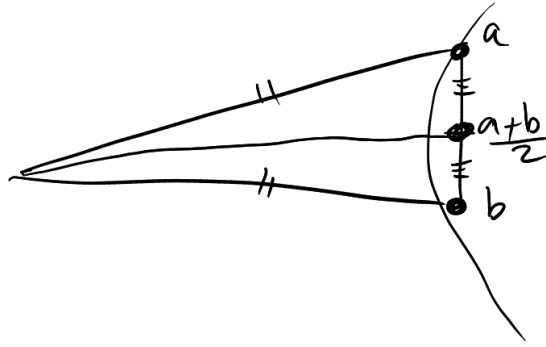
so
$$|d(x, A) - d(y, A)| \leq d(x, y). \qquad \square$$

Thus for example there exists $\mathbf{a}_x \in A$ such that $\|\mathbf{x} - \mathbf{a}_x\| = d(x, A)$.

*Proof.* Choose $R \gg 0$ such that $\overline{B(\mathbf{x}, R)} \cap A \neq \emptyset$. $A \cap \overline{B(x, R)}$ is compact, so $d(x, A)$ attains a minimum, and this must be a global minimum (if $R$ is sufficiently large). $\square$

The closest point is not necessarily unique. For example, $A = \{a : \|a\| = 1\}$, $x = 0$. If $A$ is convex, then the nearest point is unique. Suppose $a, b \in A$, $d(x, a) = d(x, b) = d(x, A)$. Then $\frac{a+b}{2} \in A$ and

$$d(x, \frac{a + b}{2}) < d(x, a)$$

unless $a = b$.

We now investigate if we can find a metric to compare compact non-empty subsets of $\mathbb{R}^n$.

Recall the required properties of a metric:

- $d(x, y) \geq 0$

- $d(x, y) = 0 \iff x = 0$

- $d(x, y) = d(y, x)$

- $d(x, y) + d(y, z) \geq d(x, z)$.

First try
$$\tau(A, B) = \inf_{a \in A} d(\mathbf{a}, B)$$

This satisfies $\tau(A, B) \geq 0$, but fails at the second hurdle:

$$\begin{aligned}
\tau(A, B) = 0 &\iff \exists a \in A \cap B \\
&\iff A \cap B \neq \emptyset
\end{aligned}$$

We will end up using
$$\sigma(A, B) = \sup_{a \in A} d(a, B)$$

Last time we tried to introduce a metric on non-empty compact sets in $\mathbb{R}^n$. We tried

$$\inf_{e \in E} d(e, F)$$

and failed. We now try

$$\sigma(E, F) = \sup_{e \in E} d(e, F).$$

Note $\sigma(E, F) \geq 0$. Now we check:

$$\sigma(E, F) = 0 \implies \sup_{e \in F} d(e, F) = 0$$
$$\implies d(e, F) = 0 \qquad \forall e \in E$$
$$\implies e \in F \qquad \forall e \in E$$
$$\implies F \supseteq E$$

so $\sigma(E, F) \iff E \subseteq F$ (the backwards direction is easy). Immediately we see that

$$\sigma(E, F) \overset{?}{=} \sigma(F, E)$$

may fail:

$$E = \{1\}, F = \{1, 2\}$$

then $\sigma(E, F) \neq \sigma(F, E)$.

However, the triangle inequality

$$\sigma(E, G) \leq \sigma(E, F) + \sigma(F, G)$$

works.

*Proof.* Let $e \in E$, $gg \in G$. Choose $f_e \in F$ such that $d(e, f_e) = d(e, F)$. Then

$$d(e, g) = d(e, f_e) + d(f_e, g)$$
$$= d(e, F) + d(f_e, g)$$

Now take a supremum over $g \in G$, and we get

$$d(e, G) \leq d(e, F) + d(f_e, G)$$
$$\leq \sigma(E, F) + \sigma(F, G)$$

so taking sup over $e \in E$ we get

$$\sigma(E, G) \leq \sigma(E, F) + \sigma(F, G). \qquad \square$$

Now we set

$$\rho(E, F) = \sigma(E, F) + \sigma(F, E)$$

Now:

- $\rho(E, F) \geq 0$

- $\rho(E, F) = \rho(F, E)$

- $\rho(E, F) = 0$ if and only if $\sigma(E, F) = 0$, $\sigma(F, E) = 0$, which happens if and only if $E \subseteq F$, $F \subseteq E$, which happens if and only if $F = E$.

- Finally,

$$\rho(E,F)+\rho(E,G) = \sigma(E,F)+\sigma(F,E)+\sigma(F,G)+\sigma(G,F) \geq \sigma(E,G)+\sigma(G,E) = \rho(E,G)$$

We call this metric the *Hausdorff metric.*

We can write it as:

$$d(E, F) = \sup_{f \in F} \inf_{e \in E} \|e - f\| + \sup_{e \in E} \inf_{f \in F} \|f - e\|.$$

Moreover, the Hausdorff metric is complete.

The proof depends on the following lemma:

> **Lemma.** If $K_1 \supseteq K_2 \supseteq \cdots$, with $K_j$ compact and non-empty, then $K = \bigcap_j K_j$ is non-empty (and compact). Furthermore, $\rho(K_j, K) \to 0$.

*Proof.* Non-empty: Choose $k_j \in K$. Since $k_j \in K_1$, $K_1$ compact, there exists $n(j)$ strictly increasing, and some $k \in K_1$ with $k_{n(j)} \to k$. But $k_{n(j)} \in K_m$ for all $m \leq n(j)$, so $k \in K_m$ for all $m \leq n(j)$, so $k \in \bigcap_m K_m = K$. So $K \neq \emptyset$.

Hausdorff metric convergence: If not then there exists a $\delta > 0$ such that $\rho(K_{n(j)}, K) > \delta$ for some $n(j) \to \infty$. But $K_1 \supseteq K_2$, so $\rho(K_n, K) \geq \delta$ for all $n$. Now choose $k_j \in K_j$ such that $d(k_j, K) \geq \delta/2$. By the previous argument, there exists $m(j) \to \infty$ such that $k_{m(j)} \to k \in K$, which gives a contradiction. $\qquad\square$

Second lemma:

> **Lemma.** If $A, B$ are compact in $\mathbb{R}^n$, then so is
> $$A + B = \{a + b : a \in A, b \in B\}.$$

*Proof.* Suppose $c_n = a_n + b_n$ with $a_n \in A$, $b_n \in B$. Then there exists $n(j) \in \infty$, $a \in A$ with $a_{n(j)} \to a \in A$ and there exists $m(k) \to \infty$, $b \in B$ such that $b_{n(m(j))} \to b$. Then $c_{n(m(j))} \to a + b$. $\qquad\square$

*Proof that Hausdorff metric is complete.* It is sufficient to show that "every sufficiently fast Cauchy sequence converges". Thus it suffices to show that if $E_n$ are compact non-empty, with

$$\rho(E_n, E_{n+1}) \leq 8^{-n}$$

then $E_n \xrightarrow{\rho} E$ for some compact non-empty set $E$.

Let $K_n = E_n + \overline{B(0, 2^{-n})}$. Then $K_n$ is compact (by the previous lemma), and $K_{n+1} \subseteq K_n$. This is because if $x \in K_{n+1}$, then $\exists y \in K_n$ such that

$$d(x, y) \leq 8^{-n},$$

and then $B(x, 2^{-n-1}) \subseteq B(y, 2^{-n})$.

We have $K_1 \supseteq K_n \supseteq \cdots$, so there exists $K$ compact non-empty such that $K_n \xrightarrow{\rho} K$ (we use $K = \bigcap_n K_n$ using the previous lemma). But $\rho(K_n, E_n) \to 0$. So $\rho(E_n, K) \to 0$. $\quad\square$

# 6 Runge's Theorem

We start with Weierstrass Theorem: real polynomials are uniformly dense in $C(I)$ for all intervals $I$.

Can we approximate continuous functions $f : K \to \mathbb{C}$ ($K$ compact in $\mathbb{C}$) by polynomials (uniformly)?

The answer is no:

Take $K = \overline{D(0,2)}$, $f(z) = \bar{z}$. Suppose $P$ is a polynomial. Then integrating around the unit disc we have:

$$\oint f(z) - P(z)\mathrm{d}z = \oint \bar{z}\mathrm{d}z = \int e^{-i\theta}e^{i\theta}i\mathrm{d}\theta = 2\pi i.$$

But then

$$2\pi \leq \oint |f(z) - P(z)|\mathrm{d}z \leq 2\pi \|f - P\_\infty,$$

so $\|f - P\|_\infty \geq 1$.

Recall (or take my word) that the uniform limit of analytic functions is analytic (see CA). Now try:

If $\Omega$ open, and $K$ compact, with $\Omega \supseteq K$, $f : \Omega \to \mathbb{C}$ analytic, then there exists $P_k$ polynomials such that $P_n \to f$ uniformly on $K$.

This is also false:

Take

$$\Omega = D(0,2) \setminus D(0,1/2)$$

with $f(z) = \frac{1}{z}$. Again,

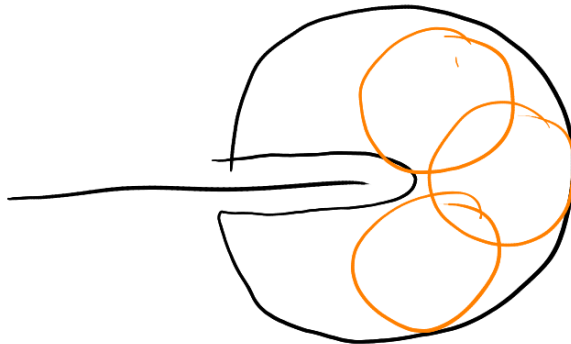$$\oint f(z)\mathrm{d}z = 2\pi i$$

so argument as before shows that we cannot approximate $f$ uniformly on $K = \{z : |z| = 1\}$ by polynomials.

Start of

lecture 14    Our last remark before looking at Runge's theorem:

Taylor's theorem is not very helpful in this case. Consider $\sqrt{z}$ with branch cut negative real axis. Take the region as sketched:

Taylor's theorem works up to the first singularity that you find.

**Definition.** $E \subseteq \mathbb{C}$ is path-connected if given $z, w \in E$, there exists $\gamma : [0, 1] \to E$ continuous such that $\gamma(0) = z$, $\gamma(1) = w$.
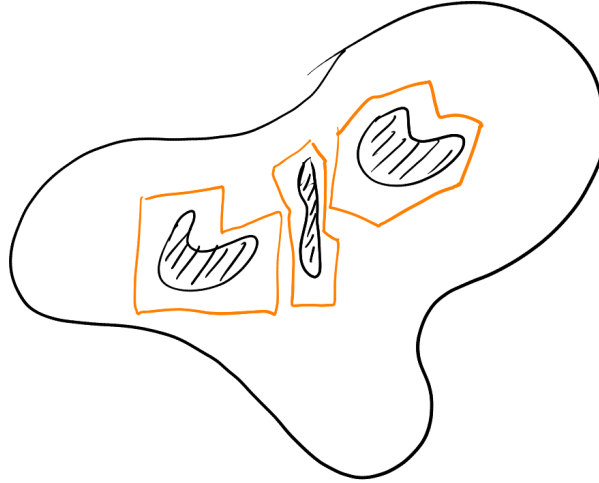
**Theorem** (Runge's Theorem). If $\Omega$ is open in $\mathbb{C}$, $f : \Omega \to \mathbb{C}$ is analytic, $K$ is a compact subset of $\Omega$ with $K^c$ path-connected, then we can find a sequence $P_n$ of polynomials such that $P_n \to f$ uniformly on $K$.

(These hypotheses stay with us for the rest of lecture 7).

**Lemma.** Let $K$, $\Omega$, $f$ as stated. We can find closed straight line segments with $l_j \subseteq \Omega \setminus K$ such that
$$f(z) = \sum_i \frac{1}{2\pi i} \int_{l_j} \frac{f(w)}{w - z} \mathrm{d}w$$

In other words, we can approximate the borders by straight lines like this:

*Proof.* Since $K$ is compact, $\Omega^c$ is closed and $K \cap \Omega^c = \emptyset$, there exists a $\delta > 0$ such that $|z - w| > \delta$ for all $z \in K$ and $w \in \Omega^c$ (proof is that $z \mapsto d(z, \Omega^c)$ is continuous $K \to \mathbb{R}$ and $K$ is compact). Choose $N \gg \delta^{-1}$ (for example $N = 100\delta^{-1} + 100$), and consider the collection of squares with vertices

$$\frac{r + si}{N}, \frac{(r+1) + si}{N}, \frac{(r+1) + (s+1)i}{N}, \frac{r + (s+1)i}{N}$$

such that

$$d\left(\frac{r + si}{N}\right) < \frac{\delta}{2}.$$

If $z \in K$ and $z$ does not lie on the boundary of a small square then

$$\frac{1}{2\pi i} \int \frac{f(w)}{w - z} \mathrm{d}w = \begin{cases} f(z) & \text{if } z \text{ lies in the small square} \\ 0 & \text{otherwise} \end{cases}$$

Summing over the set of squares (which we shall call $S$), we get

$$\sum_S \frac{1}{2\pi i} \int \frac{f(w)}{w - z} \mathrm{d}w = f(z)$$

and since interior sides cancel:

$$\sum_{l \in L} \frac{1}{2\pi i} \int_l \frac{f(w)}{w - z} \mathrm{d}w = f(z)$$

where $L$ is the set of edges that are not cancelled. Note that the $l$ in the sum lie in $\Omega \setminus K$. To extend this formula to all $z \in K$ just observe that

$$z \mapsto \frac{1}{2\pi i} \int_l \frac{f(w)}{w - z} \mathrm{d}w$$

is continuous on $\mathbb{C} \setminus l$. $\qquad\qquad \square$

Thus Runge's theorem will follow if we can show that following: If $\Omega$ open, $K$ compact subset of $\Omega$, $l \subseteq \Omega \setminus K$ a closed line segment, and $f : \Omega \to \mathbb{C}$ analytic, then we can find a sequence of polynomials $P_n$ such that $P_n(z) \to \frac{1}{2\pi i} \int_l \frac{f(w)}{w-z} \mathrm{d}w$ uniformly on $K$.

**Lemma.** If $f$ continuous on $K$, $l$ as before then given $\varepsilon >$ we can find $N$ and $w_1, w_2, \ldots, w_N \in l$ such that

$$\left| \int_l \frac{f(w)}{w - z} \mathrm{d}w - \sum_j \frac{A_j}{w_j - z} \right| < \varepsilon \qquad \forall z \in K.$$

*Proof.* $l$ is compact, $K$ is compact so $l \times K$ is compact. $G : l \times K \to \mathbb{C}$, $G(z, w) = \frac{f(w)}{w-z}$. Note $G$ is continuous, so $G$ is uniformly continuous, so if $l$ is given by $\gamma : [0, 1] \to \mathbb{C}$, $\gamma(t) = \alpha + \beta t$ then

$$\frac{f(z)}{z - w} - \sum G\left(z, \gamma\left(\frac{r}{n}\right)\right) \mathbb{1}_{\{w = \gamma(t), r/n \leq t \leq r+1/n\}}(w) \to 0$$

uniformly. $\qquad \square$

Thus Runge's theorem follows if I can show that if $w \notin K$ then there exists a sequence $P_n$ of polynoials such that $P_n(z) \to \frac{1}{w-z}$ uniformly on $K$.

**Theorem.** If $K$ is compact and $K^c$ is path-connected, $w \notin K$, then there exists polynomials $P_n$ such that $P_n(z) \to \frac{1}{w-z}$ uniformly for $z \in K$.

*Proof.* To prove this call $\Gamma$ the set of $w \notin K$ such that the result is true. Observe first that if $\overline{B(0, R)} \supseteq K$ then

$$\Gamma \supseteq \{w \in \mathbb{C} : |w| \geq 2R\}.$$

Proof:

$$\frac{1}{w - z} = \frac{1}{w\left(1 - \frac{z}{w}\right)} = \frac{1}{w} \sum_{r=1}^{\infty} \left(\frac{z}{w}\right)^r$$

and since $\left|\frac{z}{w}\right| \leq \frac{1}{2}$ the Weierstrass $M$-test tells us convergence is uniform. So

$$\frac{1}{w} \sum_{r=0}^{N} \left(\frac{z}{w}\right)^r \to \frac{1}{w - z}$$

uniformly.

Next we observe that if $w \in \Gamma$ and $d(w', K) \geq 2\eta$ $(\eta > 0)$ then if $w \in B(w', \eta)$ then we have $w \in \Gamma$. Proof:

$$\frac{1}{w-z} = \frac{1}{(w'-z)-(w'-w)} = \frac{1}{(w'-z)} \frac{1}{\left(1 - \frac{w'-w}{w'-z}\right)} = \frac{1}{w'-z} \sum_n \frac{(w'-w)^n}{w'-z}$$

convergence is uniform as $\frac{|w'-w|}{|w-z|} < \frac{1}{2}$. Thus given $\varepsilon > 0$ we can find an $N$ such that

$$\left| \frac{1}{w-z} - \sum_{n=0}^{N} \frac{(w'-w)^n}{(w-z)^{n+1}} \right| < \frac{\varepsilon}{2}$$

for $z \in K$. But we can find polynomials $P_m(z) \to \frac{1}{w'-z}$ uniformly on $K$. So for large enough $m$,

$$\left| \frac{1}{w-z} - \sum_{n=0}^{N} (w'-w)^n P_m(z)^n \right| < \varepsilon \qquad \forall z \in K.$$

So we have now shown:

(i) There exists $R$ such that $|w| \geq R$ implies $w \in \Gamma$.

(ii) If $\delta > 0$ and $w' \in \Gamma$, $\overline{B(w', 2\delta)} \cap K = \emptyset$, then $B(w', \delta) \subseteq \Gamma$.

Suppose $w_1 \notin K$. Choose $|w_0| \geq R$. Then there exists $\gamma : [0,1] \to K^c$ such that $\gamma$ is continuous and $\gamma(0) = w_1$, $\gamma(1) = w_0$. $\gamma([0,1])$ is compact, $\gamma([0,1]) \cap K = \emptyset$, so there exists $\delta > 0$ such that if $t \in [0,1]$,

$$|\gamma(t) - k| \geq 8\delta \qquad \forall k \in K.$$

$\gamma$ is continuous, hence uniformly continuous so we can find $N$ such that

$$\left| \gamma\left(\frac{r}{N}\right) - \gamma\left(\frac{r+1}{N}\right) \right| < \delta.$$

$\gamma(0) \in \Gamma$, so $\gamma\left(\frac{1}{N}\right) \in \Gamma$, so $\gamma\left(\frac{2}{N}\right) \in \Gamma$ etc, and we deduce $w_0 = \gamma(1) \in \Gamma$. $\qquad \square$

Start of

lecture 15

Thus we have proved Runge's theorem.

**Consequences**

We prove lots of examples like the following:

There exists polynomials $P_n$ such that $P_n(z) \to 1$ for $\operatorname{Im} z \geq 0$, but $P_n(z) \to 0$ for $\operatorname{Im} z < 0$. (Contrast: uniform limit of analytic functions is analytic).

Define:

$$K_n = \left\{ z - \frac{i}{n} : |z| \leq n, \operatorname{Im} z \geq 0 \right\}$$

$$K_n' = \left\{ z - \frac{4i}{n} : |z| \leq n, \operatorname{Im} z \leq 0 \right\}$$

$$\Omega_n = \left\{ z - \frac{2i}{n} : |z| < n + 1, \operatorname{Im} z > 0 \right\}$$

$$\Omega_n' = \left\{ z - \frac{3i}{n} : |z| < n + 1, \operatorname{Im} z < 0 \right\}$$

Observe that:

- $K_n \cup K_n'$ is compact.

- $\Omega_n, \Omega_n'$ are open, and disjoint (and hence their union is open).

- $\Omega_n \cup \Omega_n' \supseteq K_n \cup K_n'$.

Furthermore $(K_n \cup K_n')^c$ is path-connected. Define

$$f_n(z) = \begin{cases} 1 & z \in \Omega_n \\ 0 & z \in \Omega_n' \end{cases}$$

$f_n$ is analytic on $\Omega_n \cup \Omega_n'$ (since locally constant). So we can find a polynomial $P_n$ with

$$|P_n(z) - f_n(z)| \leq 2^{-n}$$

for all $z \in K_n \cup K_n'$. If $z^*$ is fixed and $\operatorname{Im} z^* \geq 0$ then $z^* \in \Omega_n$ for $n$ sufficiently large, so since $f_n(z^*) = 1$, $P_n(z^*) \to 1$ as $n \to \infty$. If $z^*$ is fixed and $\operatorname{Im} z^* < 0$ then $z^* \in \Omega_n'$ for $n$ sufficiently large $n$, so since $f_n(z^*) = 0$ for $n$ sufficiently large, $P_n(z^*) \to 0$ as $n \to \infty$.

# 7 Irrational and Transcendental Numbers

Proof that

$$e = \sum_{n=0}^{\infty} \frac{1}{n!}$$

is irrational.

> **Remark.** We shall use the fact that 1 is the smallest strictly positive integer.

Suppose that $e$ is rational. Then (since $e > 0$), $e = \frac{p}{q}$ for some $p, q$ coprime, with $p, q \geq 1$. Calculate:

$$e = \sum_{r=0}^{\infty} \frac{1}{r!}$$

$$\left( e - \sum_{r=0}^{q} \frac{1}{r!} \right) = \sum_{r=q+1}^{\infty} \frac{1}{r!}$$

$$q! \left( e - \sum_{r=0}^{q} \frac{1}{r!} \right) = q! \sum_{r=p+1}^{\infty} \frac{1}{r!}$$

$$q! \left( e - \sum_{n=0}^{q} \frac{1}{r!} \right) = q!e - \sum_{r=0}^{q} \frac{q!}{r!} \in \mathbb{N}$$

but also

$$q! \left( e - \sum_{r=0}^{q} \frac{1}{r!} \right) > 0.$$

But

$$\sum_{r=q+1}^{\infty} < \sum_{s=1}^{\infty} \frac{1}{(p+1)^r} = \frac{1}{p+1} \frac{1}{1 - \frac{p}{p+1}} = 1.$$

So we have found an integer between 0 and 1, which is nonsense!

A similar idea gives us $\pi$ is irrational.

*Proof.* Consider $S_n = \int_0^\pi x_n (\pi - x)^n \sin x \, dx$. We will show that $S_n$ is a polynomial in $\pi$ with coefficient of the form $A_n n!$ with $A_n \in \mathbb{Z}$. First step is to "evaluate" the integral and we shall use the fact that if $f_n(x) = x^n (\pi - x)^n$ then:

- $f_n^{(r)}(0) = 0$ if $0 \leq r \leq n - 1$

- $f_n^{(r)}(0) = B_n n! \pi^{2\pi - r}$ with $B_n$ ($n \leq r \leq 2n$.

- $f_n^{(r)}(0) = 0$ if $r > 2n$.

$$f_n(x) = \sum_k \binom{n}{k} x^{n+k} \pi^{n-k} (-1)^{n-k}.$$

$f_n^{(r)}(0) = 0$ unless $n \le r \le 2n$.

$$f_n^{(k)}(0) = \binom{n}{k}(n+k)!\pi^{n-k}(-1)^{n-k}$$
$$= B_k n! \pi^{n-k}$$

By symmetry about $\frac{\pi}{2}$, we get similar results for derivatives of $f$ when evaluated at $\pi$.

Now calculate:

$$\int_0^\pi f_n^{(r)}(x)\sin x \,\mathrm{d}x = [f_n^{(r+1)}(x)\sin x]_0^\pi + \int_0^\pi f_n^{(r+1)}(x)\cos x \,\mathrm{d}x$$
$$= \int_0^\pi f_n^{(r+1)}(x)\cos x \,\mathrm{d}x$$
$$\int_0^\pi f_n^{(r)}(x)\cos x \,\mathrm{d}x = [f_n^{(r+1)}(x)\cos x]_0^\pi - \int_0^\pi f_n^{(r+1)}(x)\sin x \,\mathrm{d}x$$
$$= -f_n^{(r+1)}(\pi) - f_n^{(r+1)}(0) - \int_0^\pi f_n^{(r+1)}(x)\cos(x)\,\mathrm{d}x$$

so by repeated integration by parts

$$\int_0^\pi x^n(\pi - x)^n \,\mathrm{d}x = \sum_{r=0}^n C_r \pi^r n!$$

with $C_r \in \mathbb{Z}$.

$$\frac{1}{n!}\int_0^\pi x^n(\pi - x)^n \,\mathrm{d}x = \sum_{r=0}^n C_r \pi^r.$$

Now

$$x^n(\pi - x)^n \le \left(\frac{\pi}{2}\right)^{-n}$$

(since $x(\pi - x) \le \frac{\pi^2}{4}$ by AM-GM), so if $\pi = \frac{p}{q}$ with $p, q > 0$, $p, q \in \mathbb{Z}$ then we have $q^n \sum_r C_r \pi^r \in \mathbb{Z}$ and

$$q^n \sum_r C_r \pi^r \le \frac{q^n}{n!}\int_0^\pi x^n(\pi - x)^n \,\mathrm{d}x \le \frac{q^n}{n!}\pi\left(\frac{\pi}{2}\right)^{-n} = u_n$$

Then $\frac{u_{n+1}}{u_n} = \frac{\pi^2 q}{(n+1)} \to 0$ as $n \to \infty$, so if $n$ is large then $u_n < 1$ so

$$0 < q^n \sum_r C_r \pi^r < 1$$

which contradicts the fact that it is an integer. $\qquad\square$

### Liouville's Transcendentals

We will always be working over the field $\mathbb{R}$ in this section.

We call a number *algebraic* if it is the root of a polynomial with integer coefficients (and transcendental otherwise). Cantor's proof of the uncountability of $\mathbb{R}$ also implies that transcendentals exist.

Let $\mathcal{E}_n$ be the collection of all real roots of polynomials $P(t) = \sum_{j=0}^{m} a_j t^j$ with $a_m \neq 0$, $m \geq 1$, $a_j \in \mathbb{Z}$, $\forall 0 \leq j \leq n$, $|a_j| \leq m$.

The collection of such polynomials is finite and such polynomial only has finitely many roots, so $\mathcal{E}_n$ is finite, so $\bigcup_n \mathcal{E}_n$ the collection of algebraic numbers is countable. The reals are not countable, so not all reals are algebraic. Cantor's proof is non-constructive.

Liouville gave another proof and his proof exhibits explicit transcendentals.

The proof depends on the following lemma:

> **Lemma** (Liouville's theorem on approximating irrational algebraic numbers). If $\xi$ is an irrational root of a polynomial $P$ of degree $n$ or less with integer coefficients then there exists an $A$ such that
>
> $$\left| \xi - \frac{p}{q} \right| \geq \frac{A}{q^n}$$
>
> whenever $p, q$ are integers with $q \geq 1$ (note that $A$ depends on $\xi$).

*Proof.* First remark is that it is sufficient to prove that there exists $N$ and $B > 0$ such that

$$\left| \xi - \frac{p}{q} \right| \geq \frac{B}{q^n} \qquad \forall p, q \in \mathbb{Z}, q \geq N.$$

Let $P$ be the polynomial in question. $P$ has only finitely many roots so we can find $\eta > 0$ such that $\eta$ is the only root of $P$ in $[\xi - \eta, \xi + \eta]$. Since $P'$ is continuous on $[\xi - n, \xi + n]$, it is bounded, i.e. $\exists M$ such that $|P'(t)| \leq M$ for all $t \in [\xi - \eta, \xi + \eta]$, so if $\frac{p}{q} \in [\xi - \eta, \xi + \eta]$ then

$$\left| P(\xi) - P\left( \frac{p}{q} \right) \right| \leq M \left| \xi - \frac{p}{q} \right|$$

(MVT). $P(\xi) = 0$, $0 \neq q^n P\left( \frac{p}{q} \right) \in \mathbb{Z}$. So $M \left| \xi - \frac{p}{q} \right| \geq \frac{1}{q^n}$, so $\left| \xi - \frac{p}{q} \right| \geq \frac{M^{-1}}{q^n}$. If

$\frac{p}{q} \in [\xi - \eta, \xi + \eta]$, then $\left| \xi - \frac{p}{q} \right| \geq \eta \geq \frac{M^{-1}}{q^n}$ for $n$ large. $\qquad\square$

Liouville's number is an example of an application of this.

$$\xi = \sum_{n=0}^{\infty} 10^{-n!}$$

If $q_m = 10^{m!}$, $p_m = 10^{m!} \sum_{r=0}^{m} 10^{-r!}$. Then

$$\left| \left( \xi - \sum_{r=0}^{\infty} 10^{-n!} \right) - \frac{p_m}{q_m} \right| = \sum_{m+1}^{\infty} 10^{-n!}$$

$$\leq 10^{-(m+1)!} \left( 1 + \frac{1}{10} + \frac{1}{10^2} + \cdots \right)$$

$$\leq 2 \cdot 10^{-(m+1)!}$$

$\xi$ irrational since it has a non-periodic decimal expansion. Then note $\frac{2 \cdot 10^{-(m+1)!}}{q_m^k} \to 0$ as $m \to \infty$. So $\xi$ is transcendental by the above lemma.

If we take $\varepsilon_r \in \{1, 2\}$, then we have that

$$\sum_r \varepsilon_r 10^{-r!}$$

is transcendental by the same argument. This gives uncountably many such examples.

# 8 Baire's Category Theorem

> **Theorem** (Baire category theorem)**.** If $(X, d)$ is a *complete* metric space and $u_1, u_2, u_3, \ldots$ are open in $X$ and are such that $u_j^c$ has empty interior (i.e. there does not exist $V$ open, $V \neq \emptyset$ such that $V \subseteq u_j^c$), then $\bigcap_{j=1}^{\infty} u_j \neq \emptyset$.

The following is the complement (so equivalent):

> **Theorem.** If $F_1, F_2$ are closed and have empty interior then $\bigcup F_j \neq X$ (i.e. there exists $x \notin \bigcup_{j=1}^{\infty} F_x$).

Another way to think of this: let $(X, d)$ be complete. If $P_j$ is a property of a point and:

(1) $P_j$ is stable, i.e. given $x$ with property $P_j$, there exists $\delta > 0$ such that $y \in B(x, \delta) \implies y$ has property $P_j$.

(2) Not $P_j$ is unstable, i.e. given $x \in X$ and $\delta > 0$, there exists $y \in B(x, \delta)$ such that $y$ has propery $P_j$.

Then there exists $x^*$ with property $P_j$ for all $j$ (set $F_j$ to be the points not satisfying $P_j$).

*Proof of Baire category theorem.* Choose $x_0 \in X$ and take $\delta_0 = 1$. We define $x_j, \delta_j$ revursively (i.e. by induction). If $x_j \in X$ and $\delta_j > 0$ constructed, then we know that there exists $x_{j+1} \notin F_{j+1}$ such that $d(x_{j+1}, x_j) < \frac{\delta_j}{4}$ (by the first hypothesis) and $0 < \delta_{j+1} < \frac{\delta_j}{4}$ such that $y \in B(x_{j+1}, \delta_{j+1}) \implies y \notin F_{j+1}$ (by the second hypothesis). Now $\delta_{j+k} \leq 4^{-k}\delta_j$ so

$$\sum_{k=0}^{\infty} d(x_{j+k}, x_{j+k+1}) \leq \sum_{r=j+1}^{\infty} 4^{-r}\delta_j \leq \frac{\delta_j}{2}$$

so $(x_j)$ is Cauchy and $x_j \to x$ with $d(x_j, x) < \frac{\delta_j}{2}$. Then $x \notin F_j$ for every $j$. $\qquad \square$

The result comes with some unsatisfactory but traditional nomenclature.

> **Definition** (First category)**.** $E$ is of first category if $E \subseteq \bigcup_{j=1}^{\infty} F_j$ with $F_j$ closed and "nowhere dense", i.e. with empty interior.

Following immediate consequences are much used:

(1) If $(X, d)$ is complete and $E$ is first category then $E \neq X$.

(2) $(X, d)$ complete. The countable union of sets of first category is of first category.

*Proof.* If $\mathcal{F}_j$ is a countable collection of nowhere dense closed sets then $\bigcup_{j=1}^{\infty} \mathcal{F}_j$ is the countable union of nowhere dense closed sets. $\square$

**Definition** (Isolated point). Let $(X, d)$ be a metric space. A point $x$ is *isolated* if there exists a $\delta > 0$ such that $B(x, \delta) = \{x\}$.

**Theorem.** A complete metric space without isolated points is uncountable.

*Proof.* Suppose $(X, d)$ is complete and countable. Let $X = \{x_1, x_2, \ldots\}$ and no $x_j$ isolated. Then $\{x_j\}$ is closed and $\{x_j\}$ is not open since $B(x_j, \delta) \neq \{x_j\}$. So $\text{Int}\{x_j\} = \emptyset$. So $X \neq \bigcup \{x_j\}$, a contradiction. $\square$

**Corollary.** $\mathbb{R}$ is uncountable.

Banach's proof that nowhere differentiable functions exist:

Work on $C([0, 1])$ with uniform norm.

**Theorem** (Baire). The set of anywhere differentiable continuous functions is of first category.

*Proof.* Banach shows that

$$E_n = \{f \in C([0, 1]) : \exists x \in [0, 1], |f(x) - f(y)| \leq n|x - y| \; \forall y \in [0, 1]\}$$

is closed and nowhere dense. If we can prove this then $\bigcup E_j$ is first category. Claim is that if $f \notin \bigcup_j E_j$ then $f$ is nowhere differentiable.

Subproof: Suppose $f$ is differentiable at $x$. Then there exists $\delta > 0$ such that

$$\left| \frac{f(x) - f(y)}{x - y} - f'(x) \right| \leq 1 \qquad \forall |y| \leq \delta$$

so $|f(x) - f(y)| \leq (|f'(x)| + 1)|x - y|$. If $|y| \geq \delta$, then

$$\left| \frac{f(x) - f(y)}{x - y} \right| \leq \frac{2\|f\|_\infty}{\delta}$$

so

$$|f(x) - f(y)| \leq \frac{2\|f\|_\infty}{\delta}|x - y|.$$

Thus $f \in E_n$ for some $n$.

Thus all we have to do is to show that $E_n$ is closed and nowhere dense (since $C([0, 1])$ is complete).

Closed: Suppose $f_m \in E_n$, $f_m \to f$ uniformly. Since $f_m \in E_n$, there exists $x_m$ such that $|f_m(x_m) - f(y)| \leq n|x_m - y|$ for all $y \in [0, 1]$. So there exists $x^*$ and $m(j) \to \infty$ such that $x_{m(j)} \to x^*$. By extracting to a subsequence we may assume $x_m \to x^*$. Then:

$$|f(x^*) - f(y)| \leq |f(x^*) - f(x_m)| + |f(x_m) - f_m(x_m)| + |f_m(x_m) - f_m(y)| + |f_m(y) - f(y)|$$
$$\leq |f(x^*) - f(x_m)| + \|f - f\|_\infty + n|x - y| + \|f_m - f\|_\infty$$
$$\to 0 + 0 + n|x^* - y| + 0$$

so $|f(x^*) - f(y)| \leq n|x^* - y|$.

Now we show that $E_n$ is nowhere dense. So we must show if $f \in C([0, 1])$ and $\delta > 0$ then there exists $g$ such that $\|f - g\|_\infty < \delta$ and $g \notin E_n$. Words of wisdom: trying to construct a "nasty function" from $g$ is hard – if $g$ happened to already be nasty in a weird way, then we might accidentally make a "nice" function out of it. So the first step is we find a nearby "nice" function, and uses this to construct a "nasty" one. Let $f \in C([0, 1])$. By Weierstrass approximation theorem, there exists a polynomial $P$ such that $\|f - P\|_\infty < \frac{\delta}{4}$. Now look at $g(x) = f(x) + \frac{\delta}{4}\sin Nx$ with $N$ to be chosen later. Observe that $\|g - f\|_\infty < \delta$ and we must show that for large $N$, $g - f \notin E_n$. Also observe that $P$ is continuously differentiable, so $|P'(t)| \leq M$ for some $M$. Suppose $x \in [0, 1]$. Without loss of generality $x \in [0, 1)$ (just reflect). Consider $N$ large, and

$$\left(r\pi + \frac{\pi}{2}\right)\frac{1}{N}, \left(r\pi + \frac{3\pi}{2}\right)\frac{1}{N}$$

with $x \leq \left(r\pi + \frac{\pi}{2}\right)\frac{1}{N}$ and $\left(\frac{r\pi}{2} + \frac{\pi}{2}\right)\frac{1}{N} < x + \frac{2\pi}{N}$. Then we have that:

$$|g(x) - g(y)| \geq |\sin Nx - \sin Ny| - |P(x) - P(y)| \geq |\sin Nx - \sin Ny| - M|x - y|.$$

So choosing $y = \left(r\pi + \frac{\pi}{2}\right)\frac{1}{N}$ or $y = \left(r\pi + \frac{3\pi}{2}\right)\frac{1}{N}$ and $N$ large enough, we get

$$|g(x) - g(y)| > (n + M)|x - y| - M|x - y| \geq n|x - y|. \qquad \square$$

Start of

lecture 18

**\* Non-examinable material**

The similarity between Baire category theorem on zero measure is immediately apparent. However they give very different kinds of genericity.

\*\* This is the end of the non-examinable comments.

Now let us prove the existence of subsets of $[0,1]$ which are closed, have empty interior and have no isolated points. (Using standard Euclidean metric).

We prove that for Hausdorff metric

$$d(E, F) = \sup_{e \in E} \inf_{f \in F} |e - f| + \sup_{f \in F} \inf_{e \in E} |e - f|,$$

except for a collection of category 1, all sets have this property.

Sufficient to prove:

(1) The collection of non-empty compact sets with no isolated points is category 1.

(2) The collection of non-empty compact sets with empty interior is of category 1.

To show collection with no isolated points is of first category, consider

$$\mathcal{E}_n = \{E \text{ compact} : \exists x \in E, E \cap B(x, 1/n) = \{x\}\}.$$

Then:

(1) $\mathcal{E}_n$ is closed in Hausdorff metric: Suppose $E_j$ such that there exists $x_j \in E$ with $B(x, 1/n) = \{x_j\}$. Then there exists $j(m) \to \infty$ such that $x_{j(m)} \to x^*$. By extracting to a subsequence, may suppose $x_j \to x^*$. If $\eta > 0$ and $n$ large enough, $|x_j - x^*| < \eta$ for all $j \leq n$, and if $m \geq n$ then $|x^* - y| \geq \frac{1}{n} - 2\eta$ for all $y \in E_j$, $y \neq x_j$. Thus $|x^* - y| \geq \frac{1}{n} - 2\eta$ for all $y \in E$, $y \neq x_n$. But $\eta$ arbitrary, so $|x^* - y| \geq \frac{1}{n} - 2\eta$ for all $\eta > 0$, $y \neq x$. So $E \in \mathcal{E}$.

On the other hand if $E \neq \emptyset$ is closed then writing

$$F_m = \left\{ \frac{r}{m} : d\left(\frac{r}{m}, E\right) \leq \frac{1}{m} \right\}$$

we know $d(F_n, E) \leq \frac{2}{n}$ and provided $n \geq 4\delta^{-1} + 1$, then there does not exist $y \in F_n$ such that $B(y, \delta) \cap E = \{y\}$.

So (1) is true.

To prove (2), let

$$\mathcal{E}_{r,n} = \left\{ E \text{ closed such that } E \supseteq \left[\frac{r}{n}, \frac{r+1}{n}\right] \right\}.$$

Then $\mathcal{E}_{r,n}$ is closed: $E_m \supseteq \left[\frac{r}{n}, \frac{r+1}{n}\right]$, $E_m \xrightarrow{d} E$, $E \supseteq \left[\frac{r}{n}, \frac{r+1}{n}\right]$ but as before, if $E$ closed and non-empty, and $\delta > 0$, then there exists $F$ finite such that $d(E, F) < \delta$ and $F \subseteq \mathcal{E}_{r,n}$.

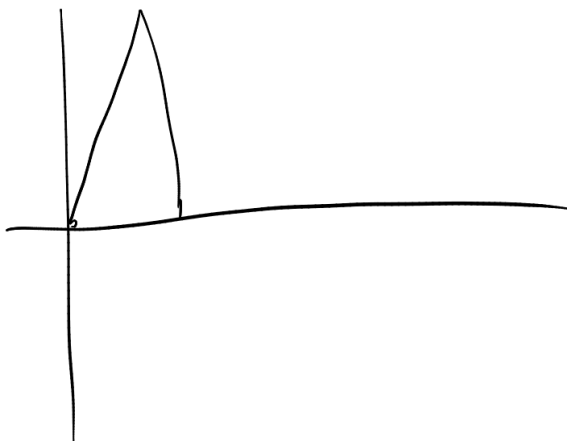If $E \notin \bigcup_{n=1}^{\infty} \bigcup_{r=0}^{n-1} \mathcal{E}_{r,n}$ then $E$ contains non non-trivial interval.

Finally we give a positive use of Baire category theorem.

Recall witch's hat:

$$f_n(x) = \begin{cases} nx & 0 \leq x \leq \frac{1}{n} \\ 2 - nx & \frac{1}{n} \leq \frac{2}{n} \\ 0 & \text{otherwise} \end{cases}$$



Standard example such that $f_n(x) \to 0$ for all $x$ (check $x = 0$ and $x \neq 0$ separately), so $f_n \to 0$ pointwise, but $\|f_n\|_\infty = 1$, so $f_n \nrightarrow 0$ uniformly.

Let $g_{m,n}(x) = f_n(m(x - [x]))$. Then $g_{m,n}(x) \to 0$ as $n \to \infty$ for all $x$ but

$$\sup_{x \in \left[\frac{r}{m}, \frac{r+1}{m}\right]} |g_{m,n}(x)| = 1.$$

Now set

$$F_n(x) = \sum_{m=1}^{\infty} 2^{-m} g_{m,n}(x).$$

This converges by the Weierstrass $M$-test. Also,

$$0 \leq F_n(x) \leq \sum_{m=1}^{M} 2^{-n} g_{m,n}(x) + \sum_{m=M+1}^{\infty} 2^{-m}$$

57

so

$$\limsup_{n \to \infty} F_n(x) \le 0 + 2^{-M}$$

so $F_n(x) \to 0$ for all $x \in [0,1]$. But $F_n(x) \ge 2^{-m} g_{n,m}(x)$ so

$$\sup_{x \in \left[\frac{r}{m}, \frac{r+1}{m}\right]} \ge 2^{-m} \qquad \forall n$$

so $F_n$ fails to converge uniformly on any non-trivial interval.

---

**Theorem** (Osgood before Baire). If $f : [0,1] \to \mathbb{R}$ is continuous and $f_n(x) \to 0$ for all $x \in [0,1]$ then given $\varepsilon > 0$ we can find a non-trivial interval $I$ and a $N$ such that

$$|f_n(x)| \le \varepsilon \qquad \forall n \ge N, \forall x \in I.$$

---

*Proof.* Let $F_n = \{x \in [0,1] : |f_n(x)| \le \varepsilon\}$. $F_n$ is closed (since $f$ is continuous). So $E_n = \bigcap_{n \ge N} F_n$ is closed. But $\bigcup E_N = [0,1]$ (because if $x \in [0,1]$, $f_n(x) \to 0$ so there exists $N$ such that $|f_n(x)| < \varepsilon$ for all $n \ge N$). So since countable union of nowhere dense sets cannot be $[0,1]$ (since it is complete), there exists $N$ such that $E_N$ is not nowhere dense, i.e. there exists an interval $I \subseteq E_N$ and this is what the theorem says. $\qquad \square$

# 9 Continued Fractions

We are sed to the representation of real numbers by decimals.

Before decimals there were fractions. What are the advantages of decimals over fractions?

(1) Easier starting from a good approximation to get a cruder approximation: we can just remove some digits from the decimal.

(2) The process of generating a decimal allows working to be continued to get a better approximation.

However there was a process with these advantages before: continued fractions.

Idea is divisiion into parts.
$$\frac{1}{3} < \frac{4}{11} < \frac{1}{2}$$
so can set
$$\frac{4}{11} = \frac{1}{2+t}$$
for some $0 < t \leq 1$. So need $2 + t = \frac{11}{4}$. So $t = \frac{3}{4}$, which can be written as $\frac{1}{1+s}$ for some $0 < s \leq 1$. In fact, $s = \frac{1}{3}$. So we can write:
$$\frac{4}{11} = \frac{1}{2 + \frac{1}{1+\frac{1}{3}}}$$

For $0 < x \leq 1$, form
$$N(x) = \left\lfloor \frac{1}{x} \right\rfloor, \qquad Tx = \frac{1}{x} - \left\lfloor \frac{1}{x} \right\rfloor.$$

Then:
$$x = \frac{1}{N(x) + Tx}$$
$$= \frac{1}{N(x) + \frac{t}{N(x)+T^2x}}$$
$$= \frac{1}{N(x) + \frac{1}{N(x)+\frac{N(x)}{T^3x}}}$$

If $x$ is irrational then the process can not terminate, and it is pretty clear that we are getting a succession of better approximations. However we will not yet consider convergence.

What happens if we start with a rational?

$$\frac{r_0}{s_0} = \frac{1}{a_1 + \frac{r_1}{s_1}} = \frac{1}{a_1 + \frac{1}{a_2 + \frac{r_2}{s_2}}}.$$

We insist that $r_j, s_j$ be coprime.

$$\frac{r_j}{s_j} = \frac{1}{a_{j+1} + \frac{r_{j+1}}{s_{j+1}}}.$$

$$\frac{r_j}{s_j} = \frac{s_{j+1}}{a_{j+1}s_{j+1} + r_{j+1}}.$$

$s_{j+1}$ is coprime to $r_{j+1}$ and hence $a_j s_{j+1} + r_{j+1}$, provided $r_{j+1} \neq 0$. So

$$s_{j+1} = r_j$$
$$r_{j+1} = a_{j+1}s_{j+1} + r_{j+1}$$

So continued fractions of rationals terminate (since these terms coincide with the definition of Euclid's algorithm, which we already know will always terminate).

**Remark.** $\frac{1}{n} = \frac{1}{(n-1)+\frac{1}{1}}$ so there can be multiple possible final forms.

If $y \in \mathbb{R}$ then $y = \lfloor y \rfloor + x$, so we often write

$$y = a_0 + \frac{1}{a_1 + \frac{1}{a_2}}.$$

**Theorem.** $\sqrt{2}$ is irrational.

*Proof.*

$$\sqrt{2} = 1 + (\sqrt{2} - 1)$$
$$= 1 + \frac{1}{\sqrt{2} + 1}$$
$$= 1 + \frac{1}{2 + (\sqrt{2} - 1)}$$
$$= 1 + \frac{1}{2 + \frac{1}{\sqrt{2}+1}}$$
$$= \frac{1}{1 + \frac{1}{2 + \frac{1}{2 + \frac{1}{2 + \cdots}}}}$$

Since this does not terminate, we deduce that $\sqrt{2}$ is irrational. $\qquad\square$

If we consider continued fraction expansion as a machine for producing digits, it is natural to compare it with decimal expansion.

Recall

$$N(x) = \left\lfloor \frac{1}{x} \right\rfloor, \qquad Tx = \frac{1}{x} - \left\lfloor \frac{1}{x} \right\rfloor$$

for continued fractions. In comparison, we have

$$Dx = \lfloor 10x \rfloor, \qquad Sx = 10x - \lfloor 10x \rfloor$$

$(0 < x \leq 1)$ for decimal expansions.

> **Remark.** If we put the uniform density on $(0, 1]$, and $X$ chosen at random from $(0, 1]$ then $DX$ has the same uniform distribution.
>
> Indeed, $X, SX, S^2X$ all have the same distribution and so $DX, DSX, DS^2X, \ldots$ are IID random variables.

Gauss observed that if take the density

$$f(x) = \frac{1}{\log 2} \frac{1}{1 + x}$$

then $X$ and $TX$ have the same distribution.

*Proof.* Let $X$ have density function $f(x) = \frac{1}{\log 2}\frac{1}{1+x}$. Then

$$\mathbb{P}(Tx \le t) = \mathbb{P}\left(\frac{1}{x} - \left\lfloor\frac{1}{x}\right\rfloor \le t\right)$$

$$= \sum_{n=1}^{\infty} \mathbb{P}\left(0 \le \frac{1}{x} - n \le t\right)$$

$$= \sum_{n=1}^{\infty} \int_{(t+n)^{-1}}^{n^{-1}} f(x)\mathrm{d}x$$

$$= \sum_{n=1}^{\infty} \int_{\frac{1}{n+t}}^{\frac{1}{n}} \frac{1}{\log 2}\frac{1}{x+1}\mathrm{d}x$$

$$= \frac{1}{\log 2}\sum_{n=1}^{\infty} [\log(1+x)]_{\frac{1}{t+n}}^{\frac{1}{n}}$$

$$= \frac{1}{\log 2}\sum_{n=1}^{\infty} \log\left(1+\frac{1}{n}\right) - \log\left(1+\frac{1}{t+n}\right)$$

$$= \frac{1}{\log 2}\sum_{n=1}^{\infty} (\log(n+1) - \log(n)) - (\log(n+t+1) - \log(n+t))$$

$$= \frac{1}{\log 2}\lim_{N\to\infty}\sum_{n=1}^{N} (\log(n+1) - \log(n) - \log(n+t+1) + \log(n+1))$$

$$= \frac{1}{\log 2}\lim_{N\to\infty} (\log(N+1) - \log(N+t+1) + \log(1+t))$$

$$= \frac{1}{\log 2}\lim_{N\to\infty} \frac{\log(N+1)}{(N+t+1)} + \log(1+t)$$

$$= \frac{1}{\log 2}\log(1+t)$$

so $TX$ has density

$$\frac{1}{\log 2}\frac{1}{t+1}. \qquad \square$$

Thus if we use the density function

$$f(t) = \frac{1}{\log 2}\frac{1}{1+t}$$

then

$$\mathbb{P}(T^n X = j) = \mathbb{P}(X = j)$$

$$= \frac{1}{\log 2} \int_{\frac{1}{j+1}}^{\frac{1}{j}} \frac{1}{1+x} dx$$

$$= \frac{1}{\log 2} [\log(1+x)]_{\frac{1}{j+1}}^{\frac{1}{j}}$$

$$= \frac{1}{\log 2} \left[ -\log\left(\frac{j+2}{j+1}\right) + \log\left(\frac{j+1}{j}\right) \right]$$

$$= \frac{1}{\log 2} \log\left(\frac{(j+1)^2}{j(j+2)}\right)$$

$$= \frac{1}{\log 2} \log\left(1 + \frac{1}{j(j+2)}\right)$$

$$\approx \frac{1}{\log 2} \frac{1}{j(j+1)}$$

$$\approx \frac{1}{\log 2} \frac{1}{j^2}$$

where the approximations are assuming $j$ is large (using the fact that $\log(1+x) \approx x$ for $x$ small.

## * Non-examinable material

Using the above observation and some more work, one can prove that if $a_j$ is the $j$-th term of the continued fraction expansion, then

$$\log(a_1 a_2 \cdots a_n)^{1/n}$$

converges as $n \to \infty$ almost everywhere.

Uses ergodic theory – see Probability and Measure.

** This is the end of the non-examinable comments.

## What about convergence

We show that if $a_0 \in \mathbb{Z}$, $a_j \in \mathbb{Z}$, $a_j \geq 1$, then writing

$$\frac{p_n}{q_n} = a_0 + \cfrac{1}{a_1 + \cfrac{1}{a_2 + \cfrac{1}{\ddots + \frac{1}{a_n}}}}$$

63

we have that $\frac{p_n}{q_n}$ converges as $n \to \infty$.

It is then easy to show that for the continued fraction expansion algorithm applied to $x$ will yield a sequence which converges to $x$.

Note we shall always take $p_n, q_n$ coprime. $\frac{p_n}{q_n}$ is called the $n$-th convergent.

Our discussion starts from the fact that we usually produce continuous fractions downwards but we evaluate them upwards.

$$\frac{r_k}{s_k} = a_k + \frac{1}{\frac{r_{k+1}}{s_{k+1}}} = a_k + \frac{s_{k+1}}{r_{k+1}} = \frac{a_k r_{k+1} + s_{k+1}}{r_{k+1}}$$

and since $r_{k+1}, s_{k+1}$ are coprime, we have

$$r_k = a_k r_{k+1} + s_{k+1}$$
$$s_k = r_{k+1}$$

We write our result in matrix form:

$$\begin{pmatrix} r_k \\ s_k \end{pmatrix} = \begin{pmatrix} a_k & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} r_{k+1} \\ s_{k+1} \end{pmatrix}$$

with

$$\begin{pmatrix} r_n \\ s_n \end{pmatrix} = \begin{pmatrix} a_n \\ 1 \end{pmatrix}.$$

We find that

$$\begin{pmatrix} p_n \\ q_n \end{pmatrix} = \begin{pmatrix} r_0 \\ s_0 \end{pmatrix} = \begin{pmatrix} a_0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} a_1 & 1 \\ 1 & 0 \end{pmatrix} \cdots \begin{pmatrix} a_{n-1} & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} a_n \\ 1 \end{pmatrix}.$$

$$\begin{pmatrix} p_{n-1} \\ q_{n-1} \end{pmatrix} = \begin{pmatrix} a_0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} a_1 & 1 \\ 1 & 0 \end{pmatrix} \cdots \begin{pmatrix} a_{n-2} & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} a_{n-1} \\ 1 \end{pmatrix}$$

$$= \begin{pmatrix} a_0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} a_1 & 1 \\ 1 & 0 \end{pmatrix} \cdots \begin{pmatrix} a_{n-2} & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} a_{n-1} & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} 1 \\ 0 \end{pmatrix}$$

$$\begin{pmatrix} p_n & p_{n-1} \\ q_n & q_{n-1} \end{pmatrix} = \begin{pmatrix} a_0 & 1 \\ 1 & 0 \end{pmatrix} \begin{pmatrix} a_1 & 1 \\ 1 & 0 \end{pmatrix} \cdots \begin{pmatrix} a_n & 1 \\ 1 & 0 \end{pmatrix}$$

Thus

$$\begin{pmatrix} p_n & p_{n-1} \\ q_n & q_{n-1} \end{pmatrix} = \begin{pmatrix} p_{n-1} & p_{n-2} \\ q_{n-1} & q_{n-2} \end{pmatrix} \begin{pmatrix} a_n & 1 \\ 1 & 0 \end{pmatrix}$$

so

$$p_n = a_n p_{n-1} + p_{n-2}$$
$$q_N = a_n q_{n-1} + q_{n-2}$$

**Remark.** $q_0 = 1$, $q_1 \geq 1$ and $q_n \geq q_{n-1} + q_{n-2}$ (since $a_n \geq 1$, $n \geq 1$). So $q_n \to \infty$.

# Index