

Computability in ergodic theory

Jeremy Avigad

Department of Philosophy and Department of Mathematical Sciences
Carnegie Mellon University

(joint work with Philipp Gerhardy, Ksenija Simic, and Henry Towsner)

November, 2007

Ergodic theory

A *discrete dynamical system* consists of a structure, \mathcal{X} , and an map T from \mathcal{X} to \mathcal{X} :

- Think of the underlying set of \mathcal{X} as the set of states of a system.
- If x is a state, Tx gives the state after one unit of time.

In ergodic theory, \mathcal{X} is assumed to be a finite measure space (X, \mathcal{B}, μ) :

- \mathcal{B} is a σ -algebra (the “measurable subsets”).
- μ is a σ -additive measure, with $\mu(X) = 1$.

T is assumed to be a *measure preserving transformation*, i.e. $\mu(T^{-1}A) = \mu(A)$ for every $A \in \mathcal{B}$.

Ergodic theory

Call (X, \mathcal{B}, μ, T) a *measure preserving system*.

- These can model physical systems (e.g. Hamilton's equations preserve Lebesgue measure).
- They can model probabilistic processes.
- They have applications to number theory and combinatorics.

The metamathematics of ergodic theory

Ergodic theory emerged from seventeenth century dynamics and nineteenth century statistical mechanics.

Since Poincaré, the emphasis has been on characterizing structural properties of dynamical systems, especially with respect to long term behavior (stability, recurrence).

Today, the field uses structural, infinitary, and nonconstructive methods that are characteristic of modern mathematics.

These are often at odds with computational concerns.

The metamathematics of ergodic theory

Central questions:

- To what extent can the methods and objects of ergodic theory be given a direct computational interpretation?
- How can we locate the “constructive content” of the nonconstructive methods?

I'll start with an overview of some results:

- the von Neumann and Birkhoff ergodic theorems
- negative results
- positive results

Then, as time allows, I'll present some of the details.

The ergodic theorems

Consider the orbit x, Tx, T^2x, \dots , and let $f : \mathcal{X} \rightarrow \mathbb{R}$ be some measurement. Consider the averages

$$\frac{1}{n}(f(x) + f(Tx) + \dots + f(T^{n-1}x)).$$

For each $n \geq 1$, define $A_n f$ to be the function $\frac{1}{n} \sum_{i < n} f \circ T^i$.

Theorem (Birkhoff). For every f in $L^1(\mathcal{X})$, $(A_n f)$ converges pointwise almost everywhere, and in the L^1 norm.

A space is *ergodic* if for every A , $T^{-1}(A) = A$ implies $\mu(A) = 0$ or $\mu(A) = 1$.

If \mathcal{X} is *ergodic*, then $(A_n f)$ converges to the constant function $\int f d\mu$.

The ergodic theorems

Recall that $L^2(\mathcal{X})$ is the Hilbert space of square-integrable functions on \mathcal{X} modulo a.e. equivalence, with inner product

$$(f, g) = \int fg \, d\mu$$

Theorem (von Neumann). For every f in $L^2(\mathcal{X})$, $(A_n f)$ converges in the L^2 norm.

A measure-preserving transformation T gives rise to an isometry \hat{T} on $L^2(\mathcal{X})$,

$$\hat{T} f = f \circ T.$$

Riesz showed that the von Neumann ergodic theorem holds, more generally, for any nonexpansive operator \hat{T} on a Hilbert space (i.e. satisfying $\|\hat{T} f\| \leq \|f\|$ for every f in \mathcal{H} .)

Bounding the rate of convergence

Can we compute a bound on the rate of convergence of $(A_n f)$ from the initial data $(T$ and $f)$?

In other words: can we compute a function $r : \mathbb{Q} \rightarrow \mathbb{N}$ such that for every rational $\varepsilon > 0$,

$$\|A_m f - A_{r(\varepsilon)} f\| < \varepsilon$$

whenever $m \geq r(\varepsilon)$?

Krengel (et al.): convergence can be arbitrarily slow.

But computability is a different question.

Note that the question depends on suitable notions of computability in analysis (I'll come back to this).

Observations

If $(a_n)_{n \in \mathbb{N}}$ is a sequence of reals that decreases to 0, no matter how slowly, one can compute a bound on the rate of convergence from (a_n) .

But there are bounded, computable, decreasing sequences (b_n) of rationals that do not have a computable limit.

There are also computable sequences (c_n) of rationals that converge to 0, with no computable bound on the rate of convergence.

Conclusion: at issue is not the *rate* of convergence, but its *predictability*.

A negative result

Theorem (A-S). There are a computable measure-preserving transformation of $[0, 1]$ under Lebesgue measure and a computable characteristic function $f = \chi_A$, such that if $f^* = \lim_n A_n f$, then $\|f^*\|_2$ is not a computable real number.

In particular, f^* is not a computable element of $L^2(\mathcal{X})$, and there is no computable bound on the rate of convergence of $(A_n f)$ in either the L^2 or L^1 norm.

A positive result

Theorem (A-G-T). Let \hat{T} be a nonexpansive operator on a separable Hilbert space and let f be an element of that space. Let $f^* = \lim_n A_n f$. Then f^* , and a bound on the rate of convergence of $(A_n f)$ in the Hilbert space norm, can be computed from f , \hat{T} , and $\|f^*\|$.

In particular, if \hat{T} arises from an ergodic transformation T , then f^* is computable from T and f .

A constructive mean ergodic theorem

When there is no computable bound on the rate of convergence, is there anything more we can say?

The assertion that the sequence $(A_n f)$ converges can be represented as follows:

$$\forall \varepsilon > 0 \exists n \forall m \geq n (\|A_m f - A_n f\| < \varepsilon).$$

This is classically equivalent to the assertion that for any function K ,

$$\forall \varepsilon > 0 \exists n \forall m \in [n, K(n)] (\|A_m f - A_n f\| < \varepsilon).$$

A constructive mean ergodic theorem

Theorem (A-G-T). Let \hat{T} be any nonexpansive operator on a Hilbert space, let f be any element of that space, and let $\varepsilon > 0$, and let K be any function. Then there is an $n \geq 1$ such that for every m in $[n, K(n)]$, $\|A_m f - A_n f\| < \varepsilon$.

In fact, we provide a bound on n expressed solely in terms of K and $\rho = \|f\|/\varepsilon$ (and independent of \hat{T}).

As special cases, we have the following:

- If $K = n^{O(1)}$, then $n(f, \varepsilon) = 2^{2^{O(\rho^2 \log \log \rho)}}$.
- If $K = 2^{O(n)}$, then $n(f, \varepsilon) = 2^1_{O(\rho^2)}$.
- If $K = O(n)$ and \hat{T} is an isometry, then $n(f, \varepsilon) = 2^{O(\rho^2 \log \rho)}$.

A constructive pointwise ergodic theorem

The following is classically equivalent to the pointwise ergodic theorem:

Theorem (A-G-T). For every f in $L^2(\mathcal{X})$, $\lambda_1 > 0$, $\lambda_2 > 0$, and K there is an $n \geq 1$ satisfying

$$\mu(\{x \mid \max_{n \leq m \leq K(n)} |A_n f(x) - A_m f(x)| > \lambda_1\}) \leq \lambda_2.$$

We provide explicit bounds on n in terms of f , λ_1 , λ_2 , and K .

Hard and soft analysis

On his blog, Terence Tao recently emphasized the distinction between “hard” and “soft” analysis.

“Hard” (or “quantitative,” or “finitary”) analysis deals with the cardinality of finite sets, the measure of bounded sets, the value of convergent integrals, the norm of finite-dimensional vectors, etc.

“Soft” analysis deals with infinitary objects, like sequences, measurable sets and functions, σ -algebras, Banach spaces, etc.

“To put it more symbolically, hard analysis is the mathematics of ε , N , $O()$, and \leq ; soft analysis is the mathematics of 0 , ∞ , \in , and \rightarrow .”

Tao independently observed that the methods described here provide “hard” analogues of “soft” results.

Hard and soft analysis

Theorem (Tao). Let T_1, \dots, T_l be commuting measure preserving transformations of \mathcal{X} , and $f_1, \dots, f_l \in L^\infty(\mathcal{X})$. Then the sequence of “diagonal averages”

$$\frac{1}{N} \sum_{n=0}^{N-1} f_1(T_1^n x) \cdots f_l(T_l^n x)$$

converges in the L^2 norm.

When $l = 1$, this is essentially the mean ergodic theorem.

Tao’s method: run the “Furstenberg correspondence” in reverse, and prove a finitary combinatorial statement by induction.

When $l = 1$, this statement is an instance of our constructive MET.

Details

Thus ends the overview. Now for some of the details:

- Notions of computability in analysis.
- A proof of the mean ergodic theorem.
- Noncomputability of the rate of convergence.
- Computability of the rate of convergence from $\|f^*\|$.
- Our constructive mean ergodic theorem.

Computability and analysis

Definition. A real number $r \in \mathbb{R}$ is *computable* if there is a computable function $\alpha : \mathbb{N} \rightarrow \mathbb{Q}$ such that $\lim_{n \rightarrow \infty} \alpha(n) = r$ and

$$\forall n \forall m \geq n (|\alpha(m) - \alpha(n)| < 1/2^n).$$

In other words, α is a computable Cauchy sequence, with an explicit rate of convergence, representing r .

Definition. A function $f : \mathbb{R} \rightarrow \mathbb{R}$ is computable if there is a computable function $F(\alpha, n)$, such that whenever α represents a real number x , $\lambda n.F(\alpha, n)$ represents the real number $f(x)$.

Fact. With the obvious extension to binary functions, addition and multiplication are computable.

Note: computable implies continuous.

Computability in analysis

What is special about \mathbb{R} ?

- There is a countable dense subset, \mathbb{Q} .
- One can construct \mathbb{R} as the Cauchy completion of \mathbb{Q} .
- With a natural encoding of \mathbb{Q} , operations we care about are computable.

The idea generalizes to separable metric spaces, and structures built on these.

For example: a computable Hilbert space is given by a countable set S , operations $+$ and $x \mapsto q \cdot x$ for $q \in \mathbb{Q}$, and an inner product (x, y) , such that S is an inner product space in the usual sense; the corresponding Hilbert space is the Cauchy completion.

Fact. A bounded linear operator T can be defined, equivalently, by its operation on S .

Computability in analysis

How to handle measure spaces? Think of $[0, 1]$ under Lebesgue measure, or $\{0, 1\}^\omega$ under coin-flipping measure.

Define a countable algebra of “simple” sets \mathcal{C} , and a σ -additive measure μ on those.

Then define the σ -algebra of measurable functions to be the completion of \mathcal{C} under the metric $d(A, B) = \mu(A \Delta B)$, modulo the relation $C \approx D$ given by $\mu(C \Delta D) = 0$.

Alternative approach: define a countable set \mathcal{C} of *simple functions*, with an “integration” operation f . Define the L^1 space as a completion.

In the usual cases, these turn out to be equivalent. Note that a measurable set of function is only defined up to points of measure 0.

Recap

Let (X, \mathcal{B}, μ, T) be a measure preserving system, and let $f : X \rightarrow \mathbb{R}$ be a measurable function. For every $n \geq 1$, let

$$(A_n f)(x) = \frac{1}{n} \sum_{i < n} f(T^i x).$$

The pointwise ergodic theorem says that for f in L^1 , $(A_n f)$ converges pointwise a.e.

The mean ergodic theorem says that for f in L^2 , $(A_n f)$ converges in the L^2 norm.

In general, the rate of convergence cannot be computed from T and f . But it *can* be computed from T , f , and $\|f^*\|$.

Recap

Of course, in *particular* cases, one can compute rates of convergence.

For example, the law of large numbers is a special case of the ergodic theorem, and there one has explicit bounds.

General computability results are only useful in fixing the outer limits of what can be done.

On the other hand, our constructive ergodic theorems give explicit bounds on how far one has to look to find pockets of “local” stability.

The mean ergodic theorem

Theorem. If T is any nonexpansive linear operator on a Hilbert space and f is any element, then the sequence $(A_n f)$ converges.

Proof. Let $M = \{h \mid Th = h\}$ be the subspace consisting of fixed-points of T . Clearly $A_n h = h$ for every $h \in M$.

Let N be subspace generated by vectors of the form $u - Tu$.

For any g of the form $u - Tu$ we have

$$\|A_n g\| = \frac{1}{n} \|u - T^n u\| \leq 2\|u\|/n, \text{ which converges to } 0.$$

Passing to limits (using the fact that A_n satisfies $\|A_n v\| \leq \|v\|$ for any v), we have that $A_n g$ converges to 0 for every $g \in N$.

The mean ergodic theorem

For arbitrary f , write $f = g + h$, where g is the projection of f on N , and $h = f - g$. It suffices to show that h is in M .

We have

$$\begin{aligned}\|Th - h\|^2 &= \|Th\|^2 - 2\langle Th, h \rangle + \|h\|^2 \\ &\leq \|h\|^2 - 2\langle Th, h \rangle + \|h\|^2 \\ &= 2\langle h, h \rangle - 2\langle Th, h \rangle \\ &= 2\langle h - Th, h \rangle,\end{aligned}$$

and the right-hand side is equal to 0, since h is orthogonal to N . So $Th = h$.

What can go wrong

What is so nonconstructive about that?

Answer: assuming the existence of the projection of f on N .

One can show (constructively) that if the projection exists, it is equal to $f - \lim_{n \rightarrow \infty} A_n f = \lim_{n \rightarrow \infty} (y_n - T y_n)$, where

$$y_n = \frac{n-1}{n} f + \frac{n-2}{n} T f + \dots + \frac{1}{n} T^{n-1} f.$$

(Simic and I learned this from Bas Spitters.) But, in general, $(y_n - T y_n)$ will not have a computable rate of convergence.

Theorem (A-S). There are a computable measure-preserving transformation of $[0, 1]$ under Lebesgue measure and a computable characteristic function $f = \chi_A$, such that if $f^* = \lim_n A_n f$, then $\|f^*\|_2$ is not a computable real number.

What can go wrong

Example (Bishop): imagine a vat of clear liquid partitioned into two parts.

Question: how can you determine whether there is a leak?

Answer: Put red dye on one side. If there is a leak then, “in the limit,” all the liquid will turn pink.

Formally: given a Turing machine, M , define a computable real number r_M such that $r_M \neq 0$ if and only if M halts on input 0.

Then design a transformation of $[0, 1]$ that is either the identity or a rotation by r_M . Apply it to $\chi_{[0,1/2]}$.

Dividing $[0, 1]$ into countably many pieces and doing this for each Turing machine yields a solution to the halting problem from $\|f^*\|$.

What can go wrong

The system in the example just described is not ergodic, and the limit of $(A_n f)$ has a noncomputable norm.

Can one find an example where the system is ergodic, $(A_n f)$ converges to 0, but at a noncomputable rate?

Recall that there are bounded, computable, decreasing sequences (b_n) of rationals that do not have a computable limit.

There are also computable sequences (c_n) of rationals that converge to 0, with no computable bound on the rate of convergence.

But if (a_n) is monotone and has a computable limit, then it has a computable bound on the rate of convergence.

A positive result

Theorem (A-G-T). Let \hat{T} be a nonexpansive operator on a separable Hilbert space and let f be an element of that space. Let $f^* = \lim_n A_n f$. Then f^* , and a bound on the rate of convergence of $(A_n f)$ in the Hilbert space norm, can be computed from f , \hat{T} , and $\|f^*\|$.

In particular, if \hat{T} arises from an ergodic transformation T , then f^* is computable from T and f .

The idea: because the mean ergodic theorem is proved by taking a projection, it is similar to the first example on the previous slide.

Formal axiomatic frameworks for analysis

The *finite types* are defined as follows:

- N is a finite type
- If σ and τ are finite types, so are $\sigma \times \tau$ and $\sigma \rightarrow \tau$

For example, the reals can be represented as type $N \rightarrow N$ functionals. Functions from \mathbb{R} to \mathbb{R} can be represented by type $(N \rightarrow N) \rightarrow (N \rightarrow N)$.

The *primitive recursive functionals of finite type* allow:

- λ abstraction, application, pairing, projection
- Higher-type primitive recursion:

$$F(0) = G, \quad F(n + 1) = H(F(n), n)$$

The theory PRA^ω (i.e. Gödel's theory T) axiomatizes these.

Formal axiomatic frameworks for analysis

Adding induction on the natural numbers provide higher-order variants of classical (Peano) and constructive (Heyting) first-order arithmetic.

- $PA^\omega = PRA^\omega + \text{induction}$
- $HA^\omega = PRA_i^\omega + \text{induction}$

These are conservative extensions of PA and HA respectively.

In fact, one can add quantifier-free choice axioms ($QF-AC$) to PA^ω , and full choice (AC) to HA^ω . One can also add Markov's principle (MP) and an “independence of premise” principle (IP_\forall) to HA^ω .

With appropriate coding, the language of PA^ω provides natural means of representing common structures in analysis, like complete separable metric spaces, Hilbert spaces, Banach spaces, and so on.

Formal axiomatic frameworks for analysis

The Gödel-Gentzen double-negation translation interprets classical logic in intuitionistic logic:

- $A^N \equiv \neg\neg A$ for atomic A
- $(\varphi \vee \psi)^N \equiv \neg(\neg\varphi^N \wedge \neg\psi^N)$
- $(\exists x \varphi)^N \equiv \neg\forall x \neg\varphi^N$.

The translation commutes with \forall , \wedge , \rightarrow .

Theorem. If $\Gamma \vdash \varphi$ classically, $\Gamma^N \vdash \varphi^N$ in intuitionistic logic.

Corollary. If $PA^\omega + (QF-AC) \vdash \varphi$, then $HA^\omega + (MP) + (IP_\forall) + (AC) \vdash \varphi^N$

Corollary. If $PA^\omega + (QF-AC) \vdash \forall x \exists y R(x, y)$ for a primitive recursive R , then $HA^\omega + (MP) + (IP_\forall) + (AC) \vdash \forall x \exists y R(x, y)$.

The Dialectica interpretation

Assigns to every formula φ in the language of PRA^ω a formula

$$\varphi^D \equiv \exists x \forall y \varphi_D(x, y)$$

where x and y are sequences of variables and φ_D is quantifier-free.

Idea: $\forall y \varphi_D(x, y)$ asserts that x is an explicit witness to the truth of φ , carrying extra information.

Inductively one shows:

Theorem (Gödel). If $HA^\omega + (MP) + (IP_\forall) + (AC)$ proves φ , there is a sequence of terms t such that quantifier-free PRA^ω proves $\varphi_D(t, y)$.

The Dialectica interpretation

Define the translation inductively, assuming

$$\varphi^D = \exists x \forall y \varphi_D \quad \text{and} \quad \psi^D = \exists u \forall v \psi_D.$$

1. For θ an atomic formula, $\theta^D = \theta_D = \theta$.
2. $(\varphi \wedge \psi)^D = \exists x, u \forall y, v (\varphi_D \wedge \psi_D)$.
3. $(\varphi \vee \psi)^D = \exists z, x, u \forall y, v ((z = 0 \wedge \varphi_D) \vee (z = 1 \wedge \psi_D))$.
4. $(\forall z \varphi(z))^D = \exists X \forall z, y \varphi_D(X(z), y, z)$.
5. $(\exists z \varphi(z))^D = \exists z, x \forall y \varphi_D(x, y, z)$.
6. $(\varphi \rightarrow \psi)^D = \exists U, Y \forall x, v (\varphi_D(x, Y(x, v)) \rightarrow \psi_D(U(x), v))$.

The last clause is a Skolemization of the formula

$$\forall x \exists u \forall v \exists y (\varphi_D(x, y) \rightarrow \psi_D(u, v)).$$

Applying the Dialectica interpretation

Recipe:

- Start with a nonconstructive proof.
- Formalize it in $PA^\omega + (QF-AC)$.
- Apply a double-negation translation.
- Get a proof in $HA^\omega + (MP) + (IP_\forall) + (AC)$.
- Apply the Dialectica interpretation

There are ways of handling certain nonconstructive set existence principles, like weak König's lemma and arithmetic comprehension.

There is also a modification of the D-interpretation, due to Kohlenbach, that makes it easier to extract bounds instead of witnesses.

The no-counterexample interpretation

Consider a sentence of the following form:

$$\forall x \exists y \forall z A(x, y, z)$$

where A is quantifier-free. The ND-interpretation yields:

$$\forall x, Z \exists y A(x, y, Z(y))$$

If the first statement is true, one can compute a y from x and Z in the second statement. Such a y foils the putative counterexample function, Z .

The Skolemization of this formula is Kreisel's *no-counterexample* interpretation:

$$\exists Y \forall x, Z A(x, Y(x, Z), z(Y(x, Z))),$$

This works for any number of quantifiers.

The no-counterexample interpretation

Recall our constructive version of the mean ergodic theorem:

Theorem (A-G-T). Let \hat{T} be any nonexpansive operator on a Hilbert space, let f be any element of that space, and let $\varepsilon > 0$, and let K be any function. Then there is an $n \geq 1$ such that for every m in $[n, K(n)]$, $\|A_m f - A_n f\| < \varepsilon$.

This is just a nicer formulation of the no-counterexample interpretation of the MET.

Analyzing the proof of the mean ergodic theorem

The bad news: the proof of the MET cannot be carried out in PA^ω . It uses a nonconstructive set existence principle (“arithmetic comprehension”).

Specifically: recall the subspace N generated by vectors of the form $u - Tu$. In fact, it suffices to consider u of the form $T^i f$. We need the projection of f onto this space.

For each n , let g_n be the projection of f onto the space spanned by $\{f - Tf, Tf - T^2 f, \dots, T^{n-1} f - T^n f\}$. Then $g = \lim_n g_n$ is the projection we want.

For each n , let $a_n = \|g_n\|$. Then the sequence (a_n) is increasing, and bounded by $\|f\|$.

The projection on N can be computed from a rate of convergence for the sequence (a_n) . But this is not always computable.

Analyzing the proof of the mean ergodic theorem

The fact that every bounded increasing sequence of real numbers converges can be expressed as follows:

$$\forall a : \mathbb{N} \rightarrow \mathbb{R}, c \in \mathbb{R} (\forall i (a_i \leq a_{i+1} \leq c) \rightarrow \\ \forall \varepsilon > 0 \exists n \forall m \geq n (|a_m - a_n| \leq \varepsilon)).$$

“Arithmetic comprehension” implies that there is a function, r , bounding the rate of convergence:

$$\forall a : \mathbb{N} \rightarrow \mathbb{R}, c \in \mathbb{R} (\forall i (a_i \leq a_{i+1} \leq c) \rightarrow \\ \exists r \forall \varepsilon > 0 \forall m \geq r(\varepsilon) (|a_m - a_{r(\varepsilon)}| \leq \varepsilon)).$$

In general, r cannot be computed from the sequence (a_i) .

The good news: a trick due to Kohlenbach can be applied here.

Analyzing the proof of the mean ergodic theorem

Note that we *can* prove

(a_n) converges with rate $r \rightarrow (A_n f)$ converges

constructively.

The ND-interpretation gives us explicit witnesses to the translation of the conclusion from a weakening of the hypothesis:

$$\forall a : \mathbb{N} \rightarrow \mathbb{R}, c \in \mathbb{R} (\forall i (a_i \leq a_{i+1} \leq c) \rightarrow \\ \forall \varepsilon > 0, M \exists n (M(n) \geq n \rightarrow (|a_{M(n)} - a_n| \leq \varepsilon))).$$

This last principle can be given a clear computational interpretation: iteratively compute $0, M(0), M(M(0)), \dots$ until one finds a value of n such that $|a_{M(n)} - a_n| \leq \varepsilon$.

Analyzing the proof of the mean ergodic theorem

In practice, one never formalizes a source proof completely.

But the metamathematical result provides a powerful heuristic: formalize intermediate lemmas and statements, and then fill in the gaps.

In our case, the proof involves working backwards to obtain the relevant “counterexample” functions on the sequence (g_n) .

For example, here is one of our lemmas:

Lemma. Let $\varepsilon > 0$, let $m \geq 1$, let $d''' = \lceil 2^9 m^4 \|f\|^4 / \varepsilon^4 \rceil$. Further suppose $\|g_i - g_{i+d'''}\| \leq \varepsilon/8$. Then for any $n \leq m$,
 $\|A_m(f - g_i) - A_n(f - g_i)\| \leq \varepsilon$.

About five or six such lemmas yield the desired theorem.

Conclusions

General remarks:

- Looking at general results in analysis through the lens of computability raises interesting issues.
- Ergodic theory is a natural market for such analyses, since it combines “hard” computational concerns with “soft” structural characterizations.
- Mathematical logic and methods of “proof mining” provide general tools and insights.

Future plans:

- Towsner and I are analyzing applications of ergodic methods in combinatorics.
- There seems to be no shortage of areas where “soft” results can be mined for quantitative bounds and dependences.