

Epistemic Arithmetic

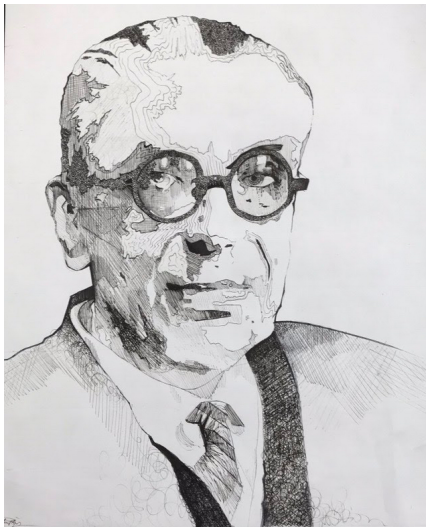
Eric Pacuit

July 28, 2025

Plan

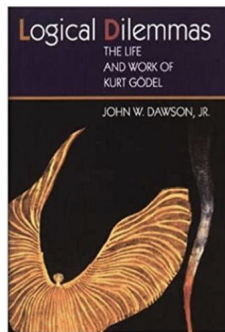
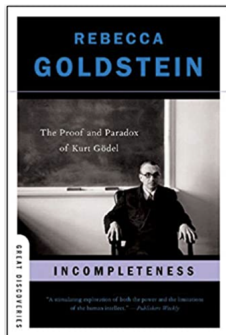
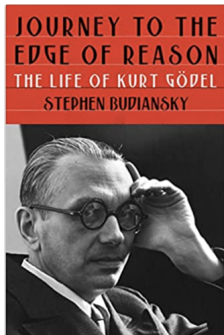
- ▶ Introduction: Smullyan's Machine
- ▶ Background
 - ▶ Formal Arithmetic
 - ▶ Gödel's Incompleteness Theorems
 - ▶ Names and Gödel numbering
 - ▶ Fixed Point Theorem
- ▶ Provability predicate and Löb's Theorem
- ▶ Provability logic
- ▶ Truth predicate and Tarski's Theorem
- ▶ A Primer on Epistemic and Doxastic Logic
- ▶ Anti-Expert Paradoxes
- ▶ Predicate approach to modality
- ▶ The Knower Paradox and variants
- ▶ Epistemic Arithmetic
- ▶ Gödel's Disjunction

Introduction



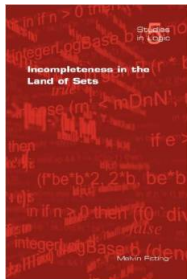
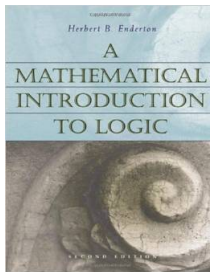
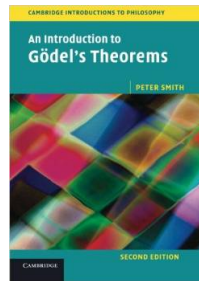
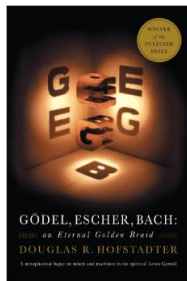
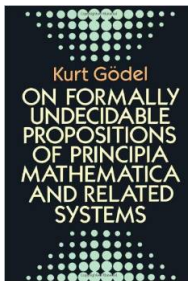
Kurt Gödel (1906 - 1978)

plato.stanford.edu/entries/goedel/



- 1929 Completeness of First-Order Logic
- 1931 First and Second Incompleteness Theorems
- 1933 Translation of classical logic in intuitionistic logic
- 1936 Speed-up Theorems
- 1938 Consistency of the Continuum Hypothesis
- 1949 Work on General Relativity
- 1958 The “Dialectica interpretation”

- 1929 Completeness of First-Order Logic
- 1931 First and Second Incompleteness Theorems
- 1933 Translation of classical logic in intuitionistic logic
- 1936 Speed-up Theorems
- 1938 Consistency of the Continuum Hypothesis
- 1949 Work on General Relativity
- 1958 The “Dialectica interpretation”

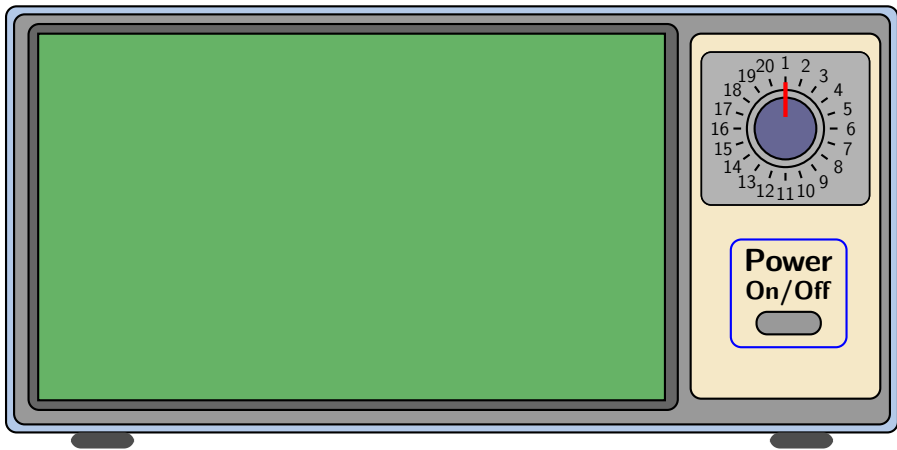


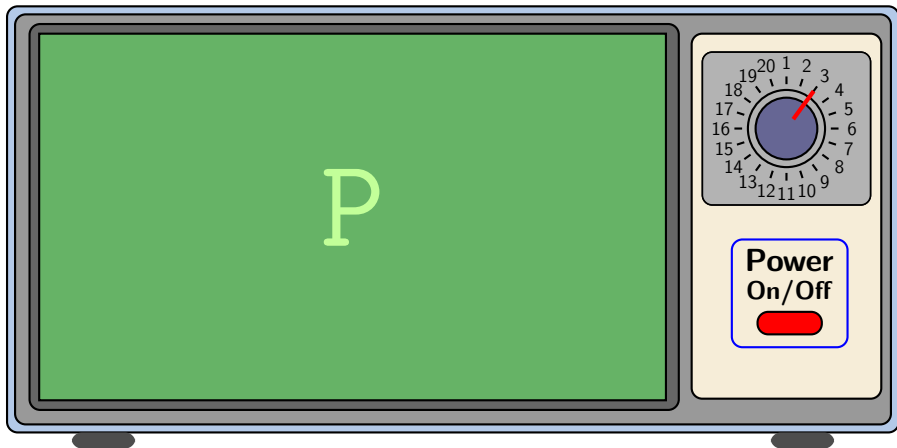
Smullyan's machine

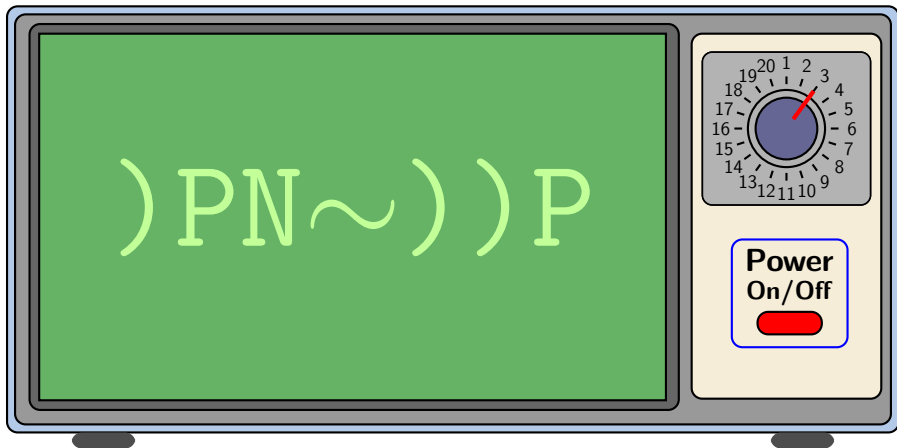
R. Smullyan. *Chapter 1: The General Idea Behind Gödel's Proof, In Gödel's Incompleteness Theorems*. Oxford University Press, 1992.

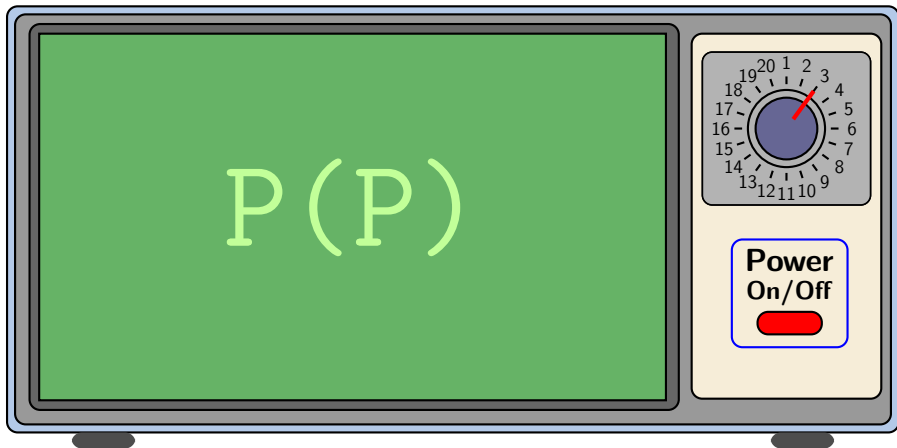
Consider a machine that displays strings of the following symbols:

) (P N ~









An **expression** is any finite string of $)$, $($, P , N or \sim .

An **expression** is any finite string of $)$, $($, P , N or \sim .

Given an expression X , the **norm** of X is $X(X)$.

An **expression** is any finite string of $\neg, (, P, N$ or \sim .

Given an expression X , the **norm** of X is $X(X)$.

Question: What are the norms of $\neg P$, $N)P$, $P(P)$, PN and $\sim PN$?

An **expression** is any finite string of $)$, $($, P , N or \sim .

Given an expression X , the **norm** of X is $X(X)$.

Question: What are the norms of $\sim P$, $N)P$, $P(P)$, PN and $\sim PN$?

Answer:

1. The norm of $\sim P$ is $\sim P(\sim P)$
2. The norm of $N)P$ is $N)P(N)P)$
3. The norm of $P(P)$ is $P(P)(P(P))$

An **expression** is any finite string of $)$, $($, P , N or \sim .

Given an expression X , the **norm** of X is $X(X)$.

Question: What are the norms of $\sim P$, $N)P$, $P(P)$, PN and $\sim PN$?

Answer:

1. The norm of $\sim P$ is $\sim P(\sim P)$
2. The norm of $N)P$ is $N)P(N)P)$
3. The norm of $P(P)$ is $P(P)(P(P))$
4. The norm of PN is $PN(PN)$
5. The norm of $\sim PN$ is $\sim PN(\sim PN)$

A **statement** is any expression of the following form:

$$P(X)$$

$$\sim P(X)$$

$$PN(X)$$

$$\sim PN(X)$$

Statement is true if...

$P(X)$

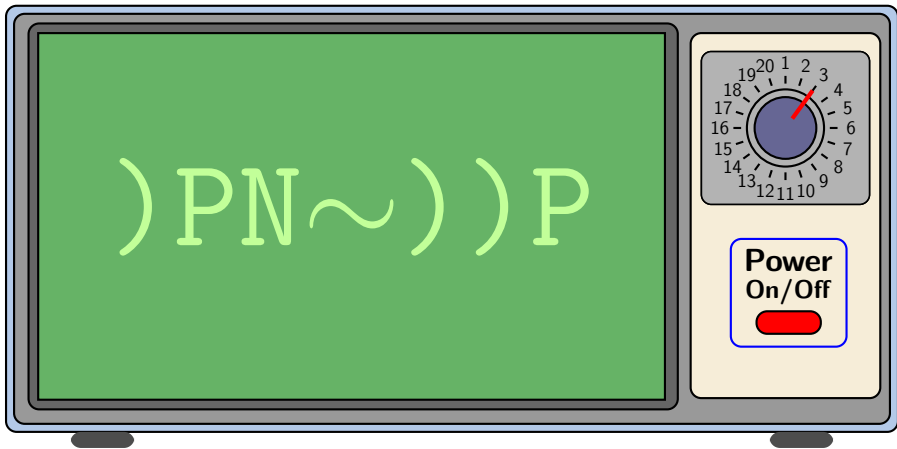
$\sim P(X)$

$PN(X)$

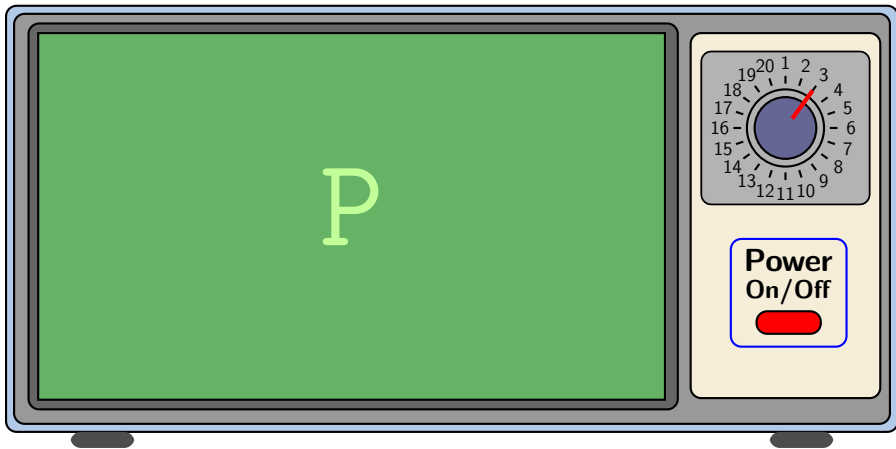
$\sim PN(X)$

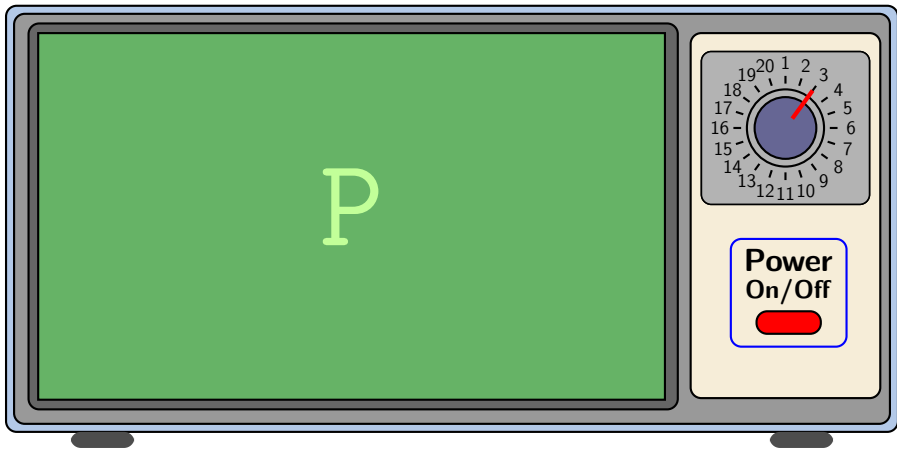
Statement	is true if...
$P(X)$	the expression X is printable.
$\sim P(X)$	the expression X is not printable.
$PN(X)$	
$\sim PN(X)$	

Statement	is true if...
$P(X)$	the expression X is printable.
$\sim P(X)$	the expression X is not printable.
$PN(X)$	the norm of X is printable.
$\sim PN(X)$	the norm of X is not printable.

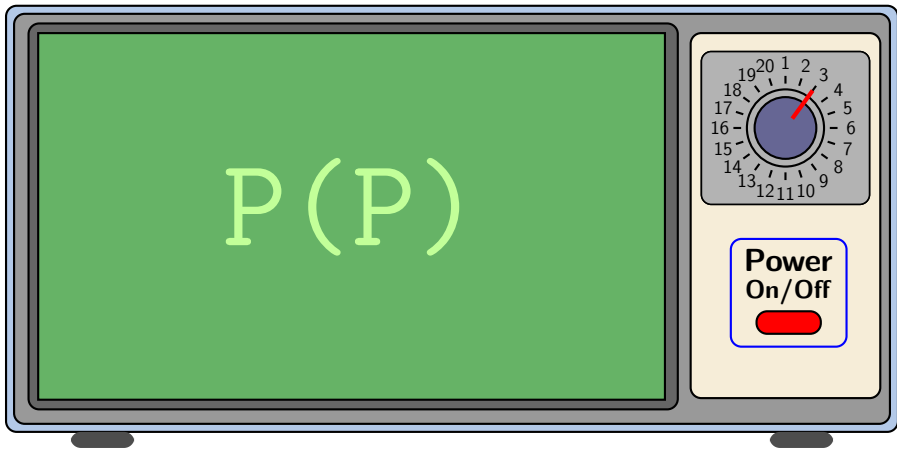


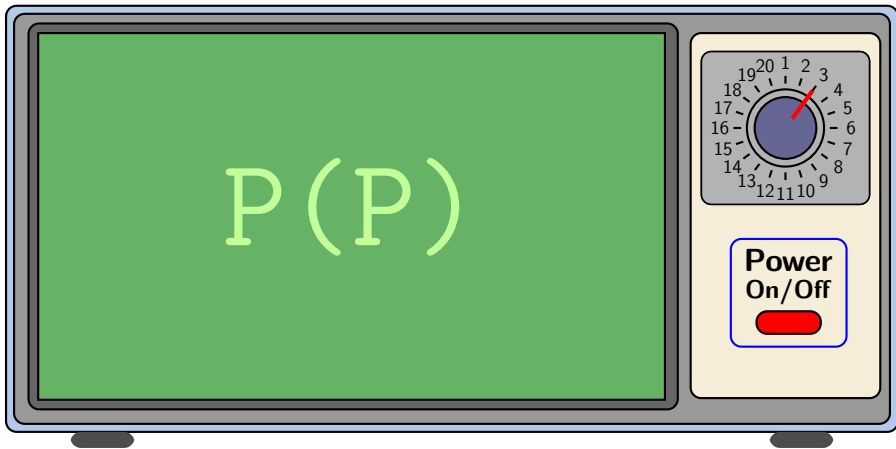
Not a statement, so neither true nor false.



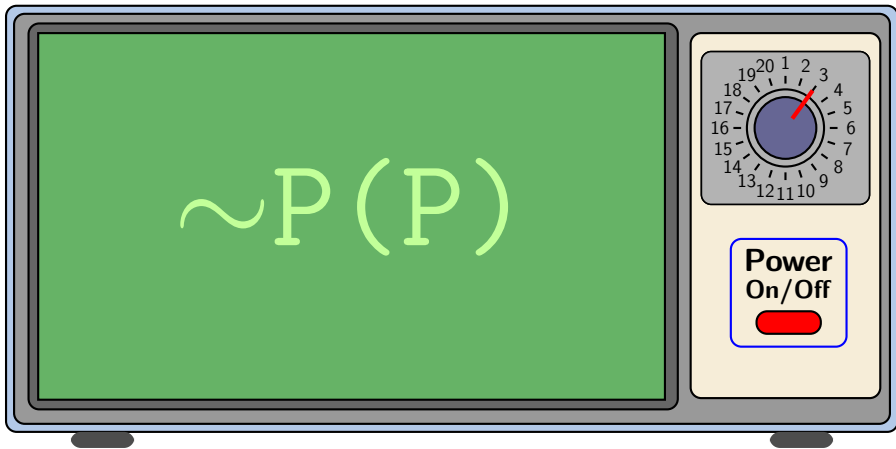


Not a statement, so neither true nor false.

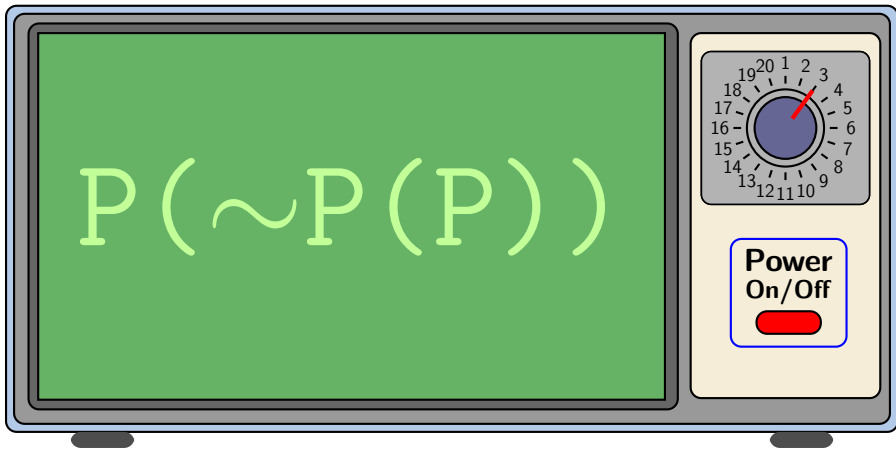




This is true.



This is false.

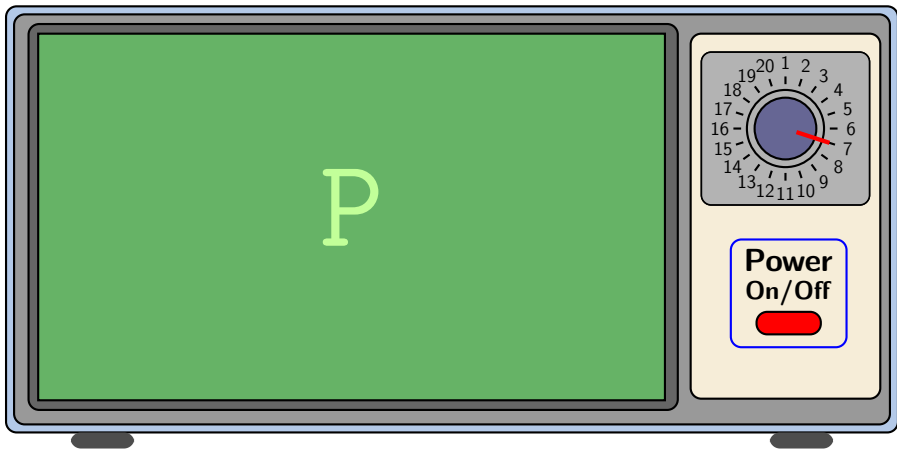


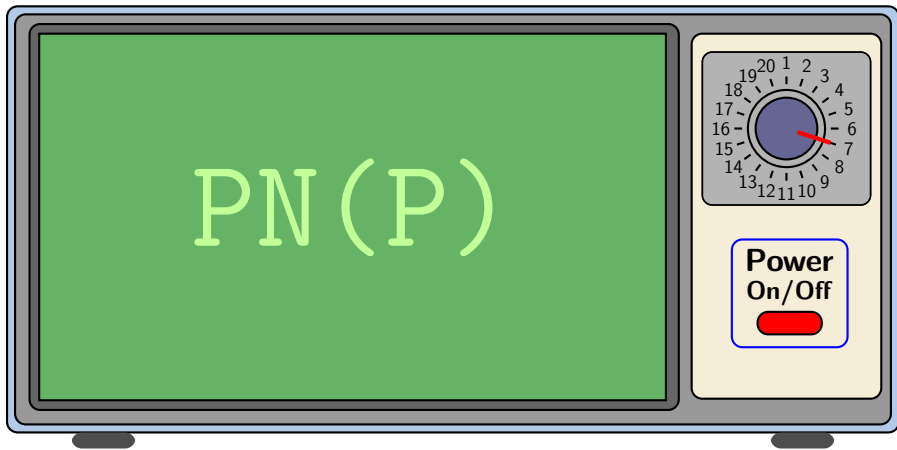
This is true.

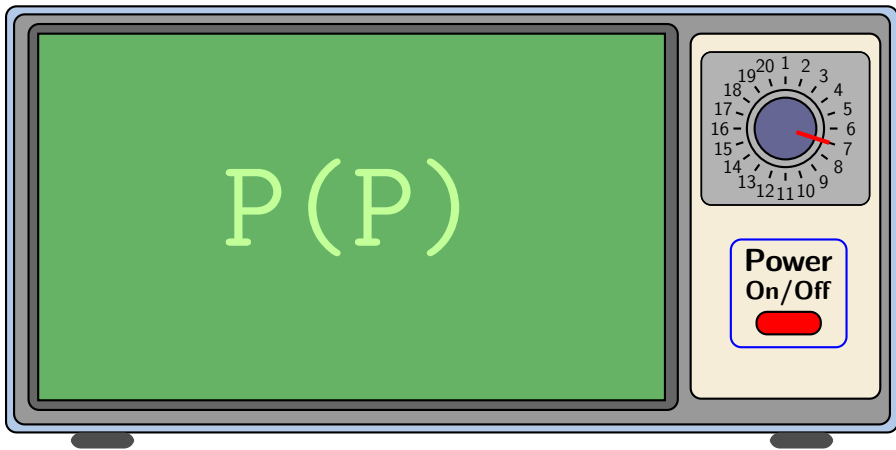
Assumption: The machine only prints true statements
(if the machine prints a statement, then it is true).

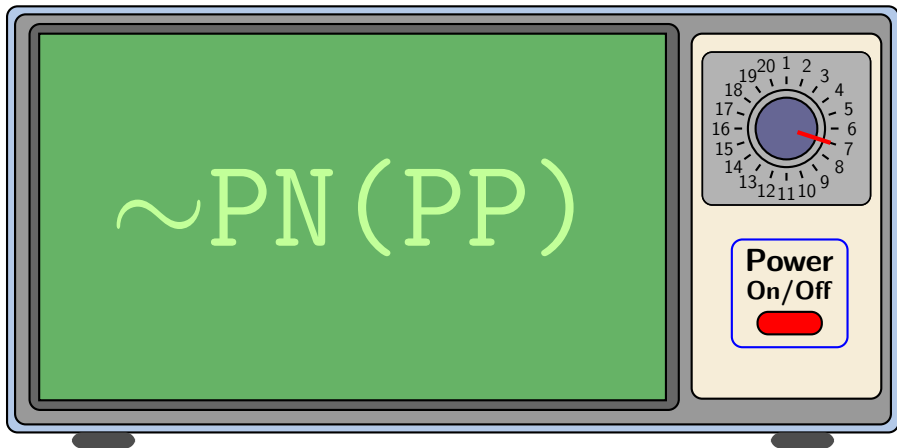
Assumption: The machine only prints true statements
(if the machine prints a statement, then it is true).

Is it possible to construct a machine that print *all* true statements?

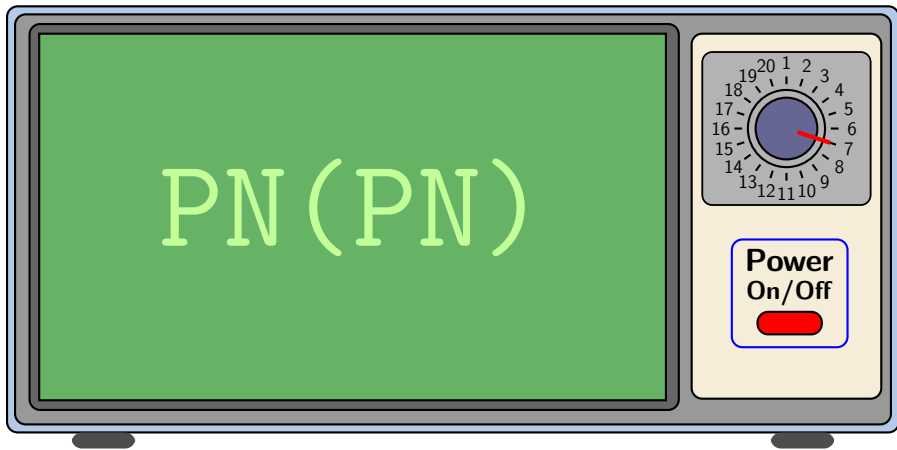




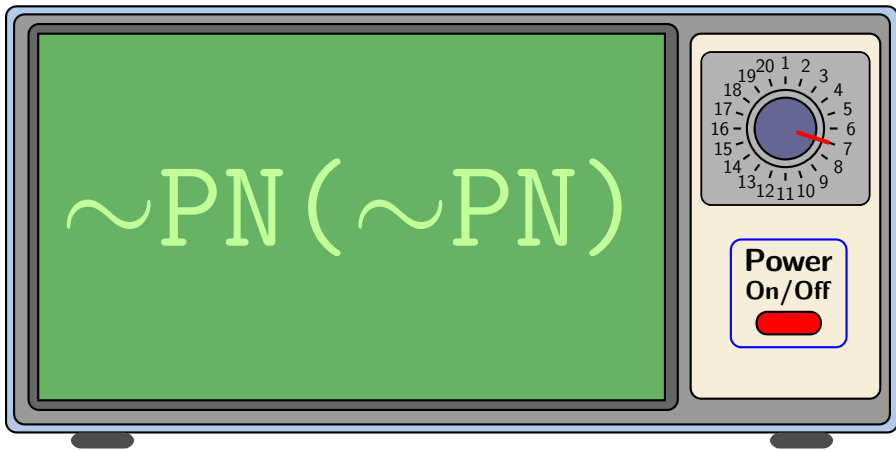




The machine is designed so that $PP(PP)$ will not be printed



$PN(PN)$ is true
if, and only if,
the norm of PN is printable
if, and only if,
 $PN(PN)$ is printable.



$\sim\text{PN}(\sim\text{PN})$ is true

if, and only if,

the norm of $\sim\text{PN}$ is not printable

if, and only if,

$\sim\text{PN}(\sim\text{PN})$ is not printable.

$\sim\text{PN}(\sim\text{PN})$ is true if, and only if, $\sim\text{PN}(\sim\text{PN})$ is not printable.

$\sim\text{PN}(\sim\text{PN})$ is true if, and only if, $\sim\text{PN}(\sim\text{PN})$ is not printable.

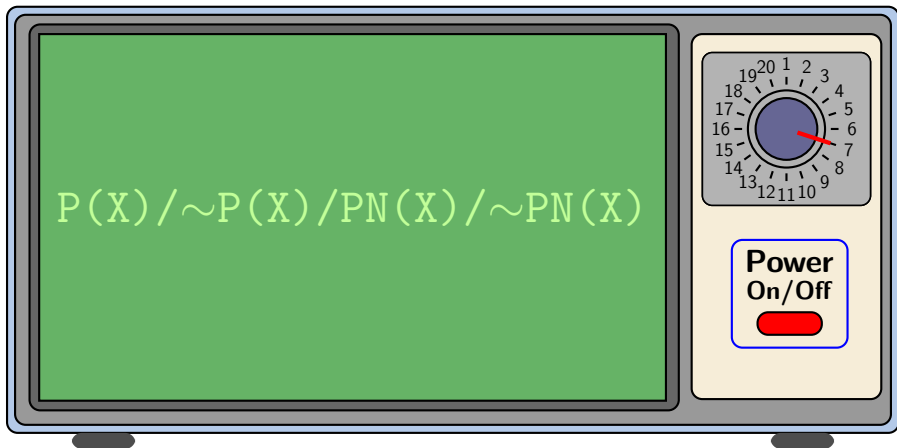
Two possibilities:

1. The machine is designed to print $\sim\text{PN}(\sim\text{PN})$
2. The machine is designed to not print $\sim\text{PN}(\sim\text{PN})$

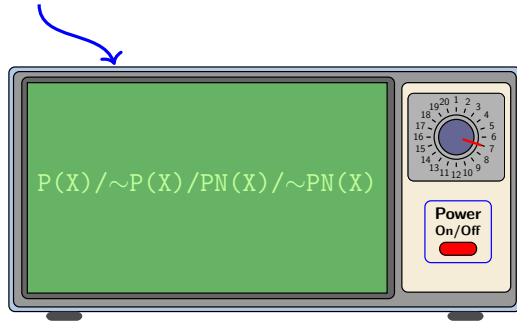
$\sim\text{PN}(\sim\text{PN})$ is true if, and only if, $\sim\text{PN}(\sim\text{PN})$ is not printable.

Two possibilities:

1. The machine is designed to print $\sim\text{PN}(\sim\text{PN})$: There is a statement that is printable, but not true. (Contradicts the assumption.)
2. The machine is designed to not print $\sim\text{PN}(\sim\text{PN})$: There is a statement that is true, but is not printable.

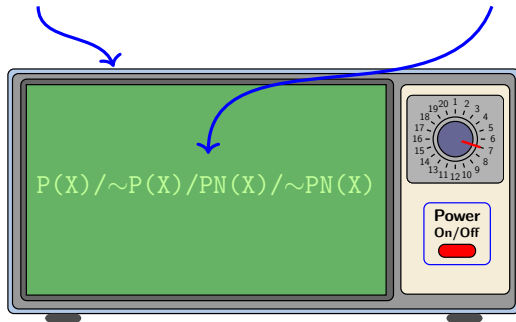


Formal system producing statements
about the natural numbers



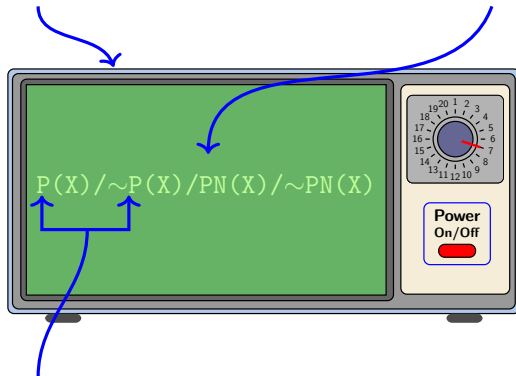
Formal system producing statements
about the natural numbers

Statements about the
natural numbers



Formal system producing statements
about the natural numbers

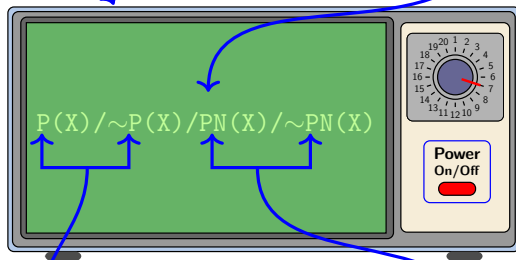
Statements about the
natural numbers



It is/is not provable in the
formal system that...

Formal system producing statements
about the natural numbers

Statements about the
natural numbers



It is/is not provable in the
formal system that...

The *diagonal* of the expression
is/is not provable...

Background

“...It would seem reasonable, therefore, to surmise that these **axioms and rules of inference are sufficient to decide all mathematical questions which can be formulated in the system concerned.**

“...It would seem reasonable, therefore, to surmise that these **axioms and rules of inference are sufficient to decide all mathematical questions which can be formulated in the system concerned**. In what follows it will be shown that this is not the case, but rather that, in both cited systems, there exists relatively simple problems of the theory of ordinary whole numbers which cannot be decided on the basis of the axioms.” (Gödel)

K. Gödel. *Über formal unentscheidbare Sätze der Principia Mathematica und verwandter Systeme, I*. Monatshefte für Mathematik und Physik, v. 38 n. 1, pp. 173 - 198, 1931.

Hilbert's Program

Hilbert's Program had two goals:

1. A complete axiomatization of mathematics, one which will settle every question in mathematics.
2. A proof using strictly finitary means to analyze the formal aspects of the above theory that the axiomatization is *reliable* (i.e., consistent).

R. Zach. *Hilbert's Program*. Stanford Encyclopedia of Philosophy, 2019, <https://plato.stanford.edu/entries/hilbert-program/>.

► $0 \neq 1$

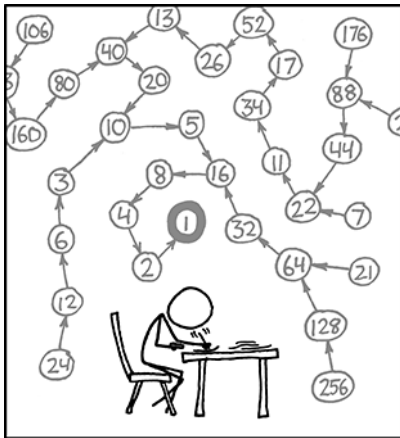
► $2 \times (1 + 4) = 5 + 3 \times 1 + 2 \times 1 + 0$

► $3 \times 2 = 2 \times 3$

- ▶ for all n , if $n \neq 0$, then there is a m such that $m + 1 = n$
- ▶ for all n, m , $n \times m = m \times n$
- ▶ There is no number smaller than 0

- ▶ There is no biggest prime number
- ▶ there are no a, b, c such that $a^n + b^n = c^n$ for $n > 2$

- ▶ every even number is the sum of two prime numbers
- ▶ there are infinitely many primes that differ by 2
- ▶ for every number n there is a sequence of numbers k_0, k_1, \dots, k_m such that $k_0 = n$, for each $0 < i \leq m$, $k_m = k_{m-1}/2$ if k_{m-1} is even and $k_m = 3k_{m-1} + 1$ if k_{m-1} is odd, and $k_m = 1$



THE COLLATZ CONJECTURE STATES THAT IF YOU PICK A NUMBER, AND IF IT'S EVEN DIVIDE IT BY TWO AND IF IT'S ODD MULTIPLY IT BY THREE AND ADD ONE, AND YOU REPEAT THIS PROCEDURE LONG ENOUGH, EVENTUALLY YOUR FRIENDS WILL STOP CALLING TO SEE IF YOU WANT TO HANG OUT.

Language of Arithmetic \mathcal{L}_A

Each of these statements can be expressed in the **language of arithmetic**.

Language of Arithmetic \mathcal{L}_A

Each of these statements can be expressed in the **language of arithmetic**.

Terms	$0 \mid x \mid S(x) \mid (x + y) \mid x \times y$
Formulas of \mathcal{L}_A	$(t = s) \mid (t < s) \mid \neg\varphi \mid (\varphi \wedge \psi) \mid (\forall x)\varphi$

Language of Arithmetic \mathcal{L}_A

Each of these statements can be expressed in the **language of arithmetic**.

Terms	$0 \mid x \mid S(x) \mid (x + y) \mid x \times y$
Formulas of \mathcal{L}_A	$(t = s) \mid (t < s) \mid \neg\varphi \mid (\varphi \wedge \psi) \mid (\forall x)\varphi$

If we could specify some axioms and inference rules that pin down the number sequence and characterize S , $+$ and \times , then we should be able to *decide* any statement about the natural numbers.

The Standard Model

$$\mathcal{N} = (\mathbb{N}, 0, S, +, *, <)$$

- ▶ $0^{\mathcal{N}} = 0$
- ▶ $S^{\mathcal{N}} : \mathbb{N} \rightarrow \mathbb{N}$ is the successor function: for all $n \in \mathbb{N}$, $S^{\mathcal{N}}(n) = n + 1$
- ▶ $+^{\mathcal{N}} : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$ is addition: for all $n, m \in \mathbb{N}$, $+^{\mathcal{N}}(n, m) = n + m$
- ▶ $\times^{\mathcal{N}} : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$ is multiplication: for all $n, m \in \mathbb{N}$, $\times^{\mathcal{N}}(n, m) = n * m$
- ▶ $<^{\mathcal{N}} \subseteq \mathbb{N} \times \mathbb{N}$ is less-than: for all $n, m \in \mathbb{N}$, $(n, m) \in <^{\mathcal{N}}$ provided that $n < m$.

Numerals

For each $n \in \mathbb{N}$ we write \bar{n} for the term representing n :

$$\bar{n} \text{ is } \underbrace{S(\cdots(S(0))\cdots)}_{n \text{ times}}$$

For instance, $\bar{3}$ is $S(S(S(0)))$

To simplify the notation, we often drop the parentheses in the terms \bar{n} . For instance, we write $SSS(0)$ instead of $S(S(S(0)))$.

Robinson's Q

$$S1. \quad \forall x(0 \neq S(x))$$

$$S2. \quad \forall x \forall y (S(x) = S(y) \rightarrow x = y)$$

$$S3. \quad \forall x (x \neq 0 \rightarrow \exists y (x = S(y)))$$

Robinson's Q

$$S1. \quad \forall x(0 \neq S(x))$$

$$S2. \quad \forall x \forall y (S(x) = S(y) \rightarrow x = y)$$

$$S3. \quad \forall x (x \neq 0 \rightarrow \exists y (x = S(y)))$$

$$A1. \quad \forall x (x + 0 = x)$$

$$A2. \quad \forall x \forall y (x + S(y) = S(x + y))$$

Robinson's \mathbf{Q}

$$S1. \quad \forall x(0 \neq S(x))$$

$$S2. \quad \forall x \forall y (S(x) = S(y) \rightarrow x = y)$$

$$S3. \quad \forall x (x \neq 0 \rightarrow \exists y (x = S(y)))$$

$$A1. \quad \forall x (x + 0 = x)$$

$$A2. \quad \forall x \forall y (x + S(y) = S(x + y))$$

$$M1. \quad \forall x (x \times 0 = 0)$$

$$M2. \quad \forall x \forall y (x \times S(y) = x \times y + x)$$

We write $\mathbf{Q} \vdash A$ when there is a derivation of A in which the only open assumptions are the axioms of \mathbf{Q} .

Defining $<$

$$x < y \leftrightarrow \exists z(x + S(z) = y)$$

Exercises

- ▶ $\mathbf{Q} \vdash \bar{1} \neq \bar{2}$
- ▶ $\mathbf{Q} \vdash \bar{1} + \bar{1} = \bar{2}$
- ▶ $\mathbf{Q} \vdash 0 + \bar{3} = \bar{3} + 0$

Exercises

- ▶ $\mathbf{Q} \vdash \bar{1} \neq \bar{2}$
- ▶ $\mathbf{Q} \vdash \bar{1} + \bar{1} = \bar{2}$
- ▶ $\mathbf{Q} \vdash 0 + \bar{3} = \bar{3} + 0$
- ▶ For all closed terms s, t ,
 - ▶ if $\mathcal{N} \models s = t$, then $\mathbf{Q} \vdash s = t$
 - ▶ if $\mathcal{N} \models s \neq t$, then $\mathbf{Q} \vdash s \neq t$

Exercises

- ▶ $\mathbf{Q} \vdash \bar{1} \neq \bar{2}$
- ▶ $\mathbf{Q} \vdash \bar{1} + \bar{1} = \bar{2}$
- ▶ $\mathbf{Q} \vdash 0 + \bar{3} = \bar{3} + 0$
- ▶ For all closed terms s, t ,
 - ▶ if $\mathcal{N} \models s = t$, then $\mathbf{Q} \vdash s = t$
 - ▶ if $\mathcal{N} \models s \neq t$, then $\mathbf{Q} \vdash s \neq t$
- ▶ $\mathbf{Q} \vdash \forall x(x + 0 = x)$

Exercises

- ▶ $\mathbf{Q} \vdash \bar{1} \neq \bar{2}$
- ▶ $\mathbf{Q} \vdash \bar{1} + \bar{1} = \bar{2}$
- ▶ $\mathbf{Q} \vdash 0 + \bar{3} = \bar{3} + 0$

- ▶ For all closed terms s, t ,
 - ▶ if $\mathcal{N} \models s = t$, then $\mathbf{Q} \vdash s = t$
 - ▶ if $\mathcal{N} \models s \neq t$, then $\mathbf{Q} \vdash s \neq t$

- ▶ $\mathbf{Q} \vdash \forall x (x + 0 = x)$

- ▶ For all $n \in \mathbb{N}$, $\mathbf{Q} \vdash 0 + \bar{n} = \bar{n}$

Exercises

- ▶ $\mathbf{Q} \vdash \bar{1} \neq \bar{2}$
- ▶ $\mathbf{Q} \vdash \bar{1} + \bar{1} = \bar{2}$
- ▶ $\mathbf{Q} \vdash 0 + \bar{3} = \bar{3} + 0$

- ▶ For all closed terms s, t ,
 - ▶ if $\mathcal{N} \models s = t$, then $\mathbf{Q} \vdash s = t$
 - ▶ if $\mathcal{N} \models s \neq t$, then $\mathbf{Q} \vdash s \neq t$

- ▶ $\mathbf{Q} \vdash \forall x(x + 0 = x)$

- ▶ For all $n \in \mathbb{N}$, $\mathbf{Q} \vdash 0 + \bar{n} = \bar{n}$

- ▶ $\mathbf{Q} \not\vdash \forall x(0 + x = x)$

Peano Arithmetic (**PA**)

The axioms of **PA** (Peano Arithmetic) are all the axioms of **Q** with every instance of the following axiom schema:

Induction Scheme: For all formulas φ of \mathcal{L}_A ,

$$(\varphi(0) \wedge \forall x(\varphi(x) \rightarrow \varphi(S(x)))) \rightarrow \forall x \varphi(x)$$

We write **PA** $\vdash \varphi$ when there is a derivation of φ in which the only open assumptions are the axioms of **PA**.

Exercises, continued

► **PA** $\vdash \forall x(0 + x = x)$

The Theory of True Arithmetic

True arithmetic: $Th(\mathcal{N}) = \{\varphi \mid \mathcal{N} \models \varphi\}$

1. Is there a computational procedure we can use to test if a sentence is in $Th(\mathcal{N})$?

The Theory of True Arithmetic

True arithmetic: $Th(\mathcal{N}) = \{\varphi \mid \mathcal{N} \models \varphi\}$

1. Is there a computational procedure we can use to test if a sentence is in $Th(\mathcal{N})$?

A **theory** is a set of sentences that is closed under entailment, i.e., \mathbf{T} is a theory if $\mathbf{T} = \{\varphi \mid \mathbf{T} \models \varphi\}$

A theory is **axiomatizable** if there is a *decidable* set of sentences \mathbf{T}_0 such that $\mathbf{T} = \{\varphi \mid \mathbf{T}_0 \models \varphi\}$

The Theory of True Arithmetic

True arithmetic: $Th(\mathcal{N}) = \{\varphi \mid \mathcal{N} \models \varphi\}$

1. Is there a computational procedure we can use to test if a sentence is in $Th(\mathcal{N})$?

A **theory** is a set of sentences that is closed under entailment, i.e., \mathbf{T} is a theory if $\mathbf{T} = \{\varphi \mid \mathbf{T} \models \varphi\}$

A theory is **axiomatizable** if there is a *decidable* set of sentences \mathbf{T}_0 such that $\mathbf{T} = \{\varphi \mid \mathbf{T}_0 \models \varphi\}$

2. Is there an axiomatizable theory \mathbf{T} such that $\mathbf{T} = Th(\mathcal{N})$? This is equivalent to asking whether \mathbf{T} is **complete**: For every sentence φ , either $\mathbf{T} \models \varphi$ or $\mathbf{T} \models \neg\varphi$.

The answer to both questions is no.

The answer to both questions is no.

Gödel's first incompleteness theorem (informal statement):

Any consistent formal theory within which a certain amount of elementary arithmetic can be carried out is **incomplete**.

Arithmetic Hierarchy

- ▶ A quantifier is **bounded** if it is the form ' $\forall x \leq t$ ' or ' $\exists x \leq t$ ', where t is a term not involving x .
- ▶ A formula is a **bounded formula** (denoted Δ_0^0) if all of its quantifiers are bounded.

Arithmetic Hierarchy

- ▶ A quantifier is **bounded** if it is the form ' $\forall x \leq t$ ' or ' $\exists x \leq t$ ', where t is a term not involving x .
- ▶ A formula is a **bounded formula** (denoted Δ_0^0) if all of its quantifiers are bounded.
- ▶ For $n \geq 0$, the classes of formulas Σ_n^0 and Π_n^0 are defined as follows:
 - ▶ $\Sigma_0^0 = \Pi_0^0 = \Delta_0^0$.
 - ▶ Σ_{n+1}^0 is the set of formulas of the form $\exists \vec{x} \varphi$ where φ is a Π_n^0 formula and \vec{x} is a (possibly empty) list of variables.
 - ▶ Π_{n+1}^0 is the set of formulas of the form $\forall \vec{x} \varphi$ where φ is a Σ_n^0 formula and \vec{x} is a (possibly empty) list of variables.

Definition

Σ_1^0 -sound A theory **T** is Σ_1^0 -**sound** iff for every Σ_1^0 -formula φ , if $\mathbf{T} \vdash \varphi$, then φ is true (in the standard model).

Definition

Σ_1^0 -complete A theory **T** is Σ_1^0 -**complete** iff for every Σ_1^0 -formula φ , if φ is true (in the standard model), then $\mathbf{T} \vdash \varphi$.

Proposition

PA (in fact, even **Q**) is Σ_1^0 -complete.

Theorem (Gödel's First Incompleteness Theorem)

Assume that **PA** is Σ_1^0 -sound. Then there is a Π_1^0 -sentence φ such that **PA** neither proves φ nor $\neg\varphi$.

Theorem (Gödel's Second Incompleteness Theorem)

Assume that **PA** is consistent. Then **PA** cannot prove Con_{PA} .

Con_{PA} is a Π_1^0 -statement that informally asserts “for all x , x does not code a proof of a contradiction from the axioms of **PA**”

- ▶ Gödel numbering
- ▶ Gödel-Carnap Fixed Point Theorem
- ▶ (Naming systems)
- ▶ Representing functions/relations

Gödel Numbering

Gödel-numbering assigns numbers to the syntactic objects of a logic (i.e., the terms, the formulas, and the derivations).

Suppose that χ is a syntactic object (i.e., a term, formula or a derivation). We use the following notation:

$gn(\chi)$: The Gödel number of χ (an integer)

$\ulcorner \chi \urcorner$: The numeral of the Gödel number of χ (a numeral). That is:

$$\ulcorner \chi \urcorner \equiv \overline{gn(\chi)}$$

Fixed-Point Theorem

Theorem (Gödel-Carnap Fixed-Point Theorem)

Let $A(x)$ be any formula of \mathcal{L}_A with one free variable x . Then there is a sentence B such that

$$\mathbf{Q} \vdash B \leftrightarrow A(\ulcorner B \urcorner).$$

Substitution

Suppose that A is a formula where x is a free variable. We write $A(x)$ to when the formula A has at most one free variable x .

If t is a term, then $A(x)[x/t]$ is A with every instance of x replaced with t . We sometimes abuse notation and write $A(t)$ instead of $A(x)[x/t]$.

Substitution

Suppose that A is a formula where x is a free variable. We write $A(x)$ to when the formula A has at most one free variable x .

If t is a term, then $A(x)[x/t]$ is A with every instance of x replaced with t . We sometimes abuse notation and write $A(t)$ instead of $A(x)[x/t]$.

Let $Sub : \mathbb{N} \times \mathbb{N} \rightarrow \mathbb{N}$ be a function, where for each $n, m \in \mathbb{N}$, $Sub(n, m)$ is the code of $\alpha(x)[x/\overline{m}]$ where n is the code of $\alpha(x)$. So, for any formula $\alpha(x)$ and $m \in \mathbb{N}$:

$$Sub(gn(\alpha(x)), m) = gn(\alpha(\overline{m}))$$

We sketch a proof under the assumption that sub is a function symbol in the language \mathcal{L}_A and the theory \mathbf{Q} “represents” Sub in the following sense:

For any formula $A(x)$ and $n \in \mathbb{N}$,

$$\mathbf{Q} \vdash \text{sub}(\ulcorner A(x) \urcorner, \bar{n}) = \ulcorner A(\bar{n}) \urcorner$$

- ▶ Let $A^*(x)$ be $A(\text{sub}(x, x))$
Let B be $A^*(\ulcorner A^*(x) \urcorner)$
- ▶ $A^*(\ulcorner A^*(x) \urcorner)$ is the formula $A(\text{sub}(\ulcorner A^*(x) \urcorner, \ulcorner A^*(x) \urcorner))$

- ▶ Let $A^*(x)$ be $A(\text{sub}(x, x))$
Let B be $A^*(\ulcorner A^*(x) \urcorner)$
- ▶ $A^*(\ulcorner A^*(x) \urcorner)$ is the formula $A(\text{sub}(\ulcorner A^*(x) \urcorner, \ulcorner A^*(x) \urcorner))$

$$\mathbf{Q} \vdash \quad \text{sub}(\ulcorner A^*(x) \urcorner, \ulcorner A^*(x) \urcorner) \quad = \quad \ulcorner A^*(\ulcorner A^*(x) \urcorner) \urcorner$$

- ▶ Let $A^*(x)$ be $A(\text{sub}(x, x))$
Let B be $A^*(\ulcorner A^*(x) \urcorner)$
- ▶ $A^*(\ulcorner A^*(x) \urcorner)$ is the formula $A(\text{sub}(\ulcorner A^*(x) \urcorner, \ulcorner A^*(x) \urcorner))$

$$\mathbf{Q} \vdash \text{sub}(\ulcorner A^*(x) \urcorner, \ulcorner A^*(x) \urcorner) = \ulcorner A^*(\ulcorner A^*(x) \urcorner) \urcorner$$

$$\mathbf{Q} \vdash A(\text{sub}(\ulcorner A^*(x) \urcorner, \ulcorner A^*(x) \urcorner)) \leftrightarrow A(\ulcorner A^*(\ulcorner A^*(x) \urcorner) \urcorner)$$

- ▶ Let $A^*(x)$ be $A(\text{sub}(x, x))$
Let B be $A^*(\ulcorner A^*(x) \urcorner)$

- ▶ $A^*(\ulcorner A^*(x) \urcorner)$ is the formula $A(\text{sub}(\ulcorner A^*(x) \urcorner, \ulcorner A^*(x) \urcorner))$

$$\mathbf{Q} \vdash \text{sub}(\ulcorner A^*(x) \urcorner, \ulcorner A^*(x) \urcorner) = \ulcorner A^*(\ulcorner A^*(x) \urcorner) \urcorner$$

$$\mathbf{Q} \vdash A(\text{sub}(\ulcorner A^*(x) \urcorner, \ulcorner A^*(x) \urcorner)) \leftrightarrow A(\ulcorner A^*(\ulcorner A^*(x) \urcorner) \urcorner)$$

$$\mathbf{Q} \vdash B \leftrightarrow A(\ulcorner B \urcorner)$$

From Cantor to Gödel...

H. Gaifman (2006). *Naming and Diagonalization, From Cantor to Gödel to Kleene*. Logic Journal of the IGPL, pp. 709 - 728.

What's in a name?

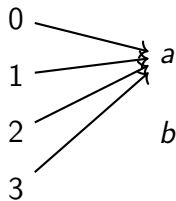
Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

	0	1	2	3

What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

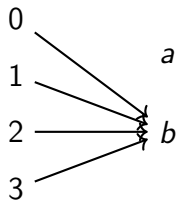
	0	1	2	3
	a	a	a	a



What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

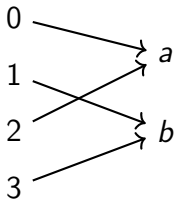
	0	1	2	3
a	a	a	a	a
b	b	b	b	b



What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

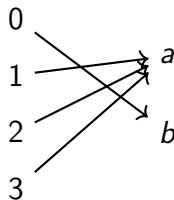
	0	1	2	3
f	a	a	a	a
g	b	b	b	b
h	a	b	a	b



What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

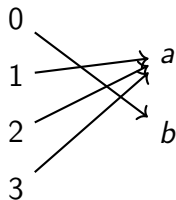
	0	1	2	3
f	a	a	a	a
g	b	b	b	b
h	a	b	a	b
i	b	a	a	a



What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

	0	1	2	3
α	a	a	a	a
β	b	b	b	b
γ	a	b	a	b
δ	b	a	a	a



What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

	0	1	2	3
α	a	a	a	a
β	b	b	b	b
γ	a	b	a	b
δ	b	a	a	a

$$g(n) = \begin{cases} b & \text{if } \gamma(n) = a \\ a & \text{if } \gamma(n) = b \end{cases}$$

0	
1	a
2	b
3	

What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

	0	1	2	3
α	a	a	a	a
β	b	b	b	b
γ	a	b	a	b
δ	b	a	a	a

$$g(n) = \begin{cases} b & \text{if } \gamma(n) = a \\ a & \text{if } \gamma(n) = b \end{cases}$$

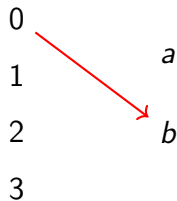
0	
1	a
2	b
3	

What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

	0	1	2	3
α	a	a	a	a
β	b	b	b	b
γ	a	b	a	b
δ	b	a	a	a

$$g(n) = \begin{cases} b & \text{if } \gamma(n) = a \\ a & \text{if } \gamma(n) = b \end{cases}$$

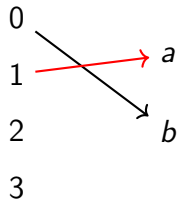


What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

	0	1	2	3
α	a	a	a	a
β	b	b	b	b
γ	a	b	a	b
δ	b	a	a	a

$$g(n) = \begin{cases} b & \text{if } \gamma(n) = a \\ a & \text{if } \gamma(n) = b \end{cases}$$

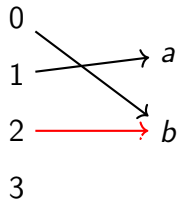


What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

	0	1	2	3
α	a	a	a	a
β	b	b	b	b
γ	a	b	a	b
δ	b	a	a	a

$$g(n) = \begin{cases} b & \text{if } \gamma(n) = a \\ a & \text{if } \gamma(n) = b \end{cases}$$

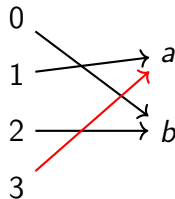


What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

	0	1	2	3
α	a	a	a	a
β	b	b	b	b
γ	a	b	a	b
δ	b	a	a	a

$$g(n) = \begin{cases} b & \text{if } \gamma(n) = a \\ a & \text{if } \gamma(n) = b \end{cases}$$



What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

	0	1	2	3
[0]	<i>a</i>	<i>a</i>	<i>a</i>	<i>a</i>
[1]	<i>b</i>	<i>b</i>	<i>b</i>	<i>b</i>
[2]	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>
[3]	<i>b</i>	<i>a</i>	<i>a</i>	<i>a</i>

0

a

1

2

b

3

What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

	0	1	2	3
[0]	a	a	a	a
[1]	b	b	b	b
[2]	a	b	a	b
[3]	b	a	a	a

$$diag(n) = \begin{cases} b & \text{if } n = a \\ a & \text{if } n = b \end{cases}$$

0

1

a

2

b

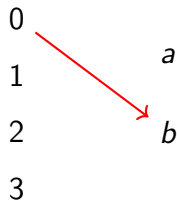
3

What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

	0	1	2	3
[0]	<i>a</i>	<i>a</i>	<i>a</i>	<i>a</i>
[1]	<i>b</i>	<i>b</i>	<i>b</i>	<i>b</i>
[2]	<i>a</i>	<i>b</i>	<i>a</i>	<i>b</i>
[3]	<i>b</i>	<i>a</i>	<i>a</i>	<i>a</i>

$$diag(n) = \begin{cases} b & \text{if } n = a \\ a & \text{if } n = b \end{cases}$$

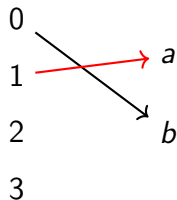


What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

	0	1	2	3
[0]	a	a	a	a
[1]	b	b	b	b
[2]	a	b	a	b
[3]	b	a	a	a

$$diag(n) = \begin{cases} b & \text{if } n = a \\ a & \text{if } n = b \end{cases}$$

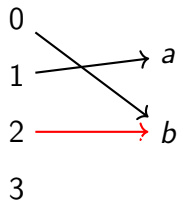


What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

	0	1	2	3
[0]	a	a	a	a
[1]	b	b	b	b
[2]	a	b	a	b
[3]	b	a	a	a

$$diag(n) = \begin{cases} b & \text{if } n = a \\ a & \text{if } n = b \end{cases}$$

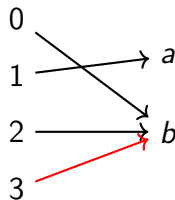


What's in a name?

Functions from $\{0, 1, 2, 3\}$ to $\{a, b\}$

	0	1	2	3
[0]	a	a	a	a
[1]	b	b	b	b
[2]	a	b	a	b
[3]	b	a	a	a

$$\text{diag}(n) = \begin{cases} b & \text{if } n = a \\ a & \text{if } n = b \end{cases}$$



Cantor's Diagonalization Proof

Functions from \mathbb{N} to $\{0, 1\}$

	0	1	2	3	\dots	n	\dots
	0	0	0	0	\dots	0	\dots
	0	1	0	1	\dots	1	\dots
	0	1	1	0	\dots	0	\dots
	1	0	1	0	\dots	1	\dots
	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
	0	0	1	0	\dots	1	\dots
	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

Cantor's Diagonalization Proof

Functions from \mathbb{N} to $\{0, 1\}$

	0	1	2	3	...	n	...
[0]	0	0	0	0	...	0	...
[1]	0	1	0	1	...	1	...
[2]	0	1	1	0	...	0	...
[3]	1	0	1	0	...	1	...
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
[n]	0	0	1	0	...	1	...
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

$$d : \mathbb{N} \rightarrow \{0, 1\}$$

$$d(n) = 1 - n$$

Cantor's Diagonalization Proof

Functions from \mathbb{N} to $\{0, 1\}$

	0	1	2	3	...	n	...
[0]	0	0	0	0	...	0	...
[1]	0	1	0	1	...	1	...
[2]	0	1	1	0	...	0	...
[3]	1	0	1	0	...	1	...
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
[n]	0	0	1	0	...	1	...
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

$$d : \mathbb{N} \rightarrow \{0, 1\} \quad d(n) = 1 - n$$

Then, $d \neq [n]$ for any $n \in \mathbb{N}$.

Cantor's original statement is phrased as a non-existence claim: there is no function mapping all the members of a set S onto the set of all 0, 1-valued functions over S . But the proof establishes a positive result: given any way of correlating functions with members of S , one can construct a function not correlated with any member of S .

(Gaiffman, p. 711)

Richard's Paradox (1905)

Consider all the definitions (in English) of real numbers.

Richard's Paradox (1905)

Consider all the definitions (in English) of real numbers. Since any such definition is a finite sequence of letters, the definitions can be listed in order.

Richard's Paradox (1905)

Consider all the definitions (in English) of real numbers. Since any such definition is a finite sequence of letters, the definitions can be listed in order.

Let u_i be the real number defined by the i th definition and $f_i(n)$ be the n th member of the decimal expansion of u_i .

Richard's Paradox (1905)

Consider all the definitions (in English) of real numbers. Since any such definition is a finite sequence of letters, the definitions can be listed in order.

Let u_i be the real number defined by the i th definition and $f_i(n)$ be the n th member of the decimal expansion of u_i .

Let u^* be the number whose decimal expansion is $0.g(1)g(2)\cdots g(n)\cdots$ where g is defined by $g(n) = f_n(n) + 1$ if $f_n(n) < 8$, $g(n) = 1$ otherwise.

Richard's Paradox (1905)

Consider all the definitions (in English) of real numbers. Since any such definition is a finite sequence of letters, the definitions can be listed in order.

Let u_i be the real number defined by the i th definition and $f_i(n)$ be the n th member of the decimal expansion of u_i .

Let u^* be the number whose decimal expansion is $0.g(1)g(2)\cdots g(n)\cdots$ where g is defined by $g(n) = f_n(n) + 1$ if $f_n(n) < 8$, $g(n) = 1$ otherwise.

But the previous description defines a number, so $u^* = u_i$ for some i . But, this is impossible.

Richard's Paradox (1905)

1. Let A be the set of all positive integers that can be defined in under 100 words. Since there are only finitely many of these, there must be a smallest positive integer n that does not belong to A .

Richard's Paradox (1905)

1. Let A be the set of all positive integers that can be defined in under 100 words. Since there are only finitely many of these, there must be a smallest positive integer n that does not belong to A .

But haven't I just defined n in under 100 words?

Richard's Paradox (1905)

1. Let A be the set of all positive integers that can be defined in under 100 words. Since there are only finitely many of these, there must be a smallest positive integer n that does not belong to A .

But haven't I just defined n in under 100 words?

2. Let B be the set of all reasonably interesting positive integers. Let n be the smallest integer not belonging to B .

Richard's Paradox (1905)

1. Let A be the set of all positive integers that can be defined in under 100 words. Since there are only finitely many of these, there must be a smallest positive integer n that does not belong to A .

But haven't I just defined n in under 100 words?

2. Let B be the set of all reasonably interesting positive integers. Let n be the smallest integer not belonging to B .

But surely this defining property of n makes it reasonably interesting.

Let f be a function that associates each number $x \in \mathbb{N}$ with a subset of \mathbb{N} , i.e., for all $x \in \mathbb{N}$, $f(x) \subseteq \mathbb{N}$.

Let f be a function that associates each number $x \in \mathbb{N}$ with a subset of \mathbb{N} , i.e., for all $x \in \mathbb{N}$, $f(x) \subseteq \mathbb{N}$.

Define S^* by:

$$x \in S^* \Leftrightarrow x \notin f(x)$$

Let f be a function that associates each number $x \in \mathbb{N}$ with a subset of \mathbb{N} , i.e., for all $x \in \mathbb{N}$, $f(x) \subseteq \mathbb{N}$.

Define S^* by:

$$x \in S^* \Leftrightarrow x \notin f(x)$$

The assumption that there is some z such that $f(z) = S^*$ leads to a contradiction.

	0	1	2	3	\dots	n	\dots	$S \subseteq \mathbb{N}$
$f(0)$	0	0	0	0	\dots	0	\dots	
$f(1)$	0	1	0	1	\dots	1	\dots	
$f(2)$	0	1	1	0	\dots	0	\dots	
$f(3)$	1	0	1	0	\dots	1	\dots	
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	
$f(n)$	0	0	1	0	\dots	1	\dots	
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	

	0	1	2	3	...	n	...	$S \subseteq \mathbb{N}$
$f(0)$	0	0	0	0	...	0	...	\emptyset
$f(1)$	0	1	0	1	...	1	...	
$f(2)$	0	1	1	0	...	0	...	
$f(3)$	1	0	1	0	...	1	...	
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	
$f(n)$	0	0	1	0	...	1	...	
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	

	0	1	2	3	...	n	...	$S \subseteq \mathbb{N}$
$f(0)$	0	0	0	0	...	0	...	\emptyset
$f(1)$	0	1	0	1	...	1	...	$\{1, 3, \dots, n, \dots\}$
$f(2)$	0	1	1	0	...	0	...	
$f(3)$	1	0	1	0	...	1	...	
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	
$f(n)$	0	0	1	0	...	1	...	
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	

	0	1	2	3	...	n	...	$S \subseteq \mathbb{N}$
$f(0)$	0	0	0	0	...	0	...	\emptyset
$f(1)$	0	1	0	1	...	1	...	$\{1, 3, \dots, n, \dots\}$
$f(2)$	0	1	1	0	...	0	...	$\{1, 2\}$
$f(3)$	1	0	1	0	...	1	...	
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	
$f(n)$	0	0	1	0	...	1	...	
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	

	0	1	2	3	...	n	...	$S \subseteq \mathbb{N}$
$f(0)$	0	0	0	0	...	0	...	\emptyset
$f(1)$	0	1	0	1	...	1	...	$\{1, 3, \dots, n, \dots\}$
$f(2)$	0	1	1	0	...	0	...	$\{1, 2\}$
$f(3)$	1	0	1	0	...	1	...	$\{0, 2, \dots, n, \dots\}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	
$f(n)$	0	0	1	0	...	1	...	
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	

	0	1	2	3	...	n	...	$S \subseteq \mathbb{N}$
$f(0)$	0	0	0	0	...	0	...	\emptyset
$f(1)$	0	1	0	1	...	1	...	$\{1, 3, \dots, n, \dots\}$
$f(2)$	0	1	1	0	...	0	...	$\{1, 2\}$
$f(3)$	1	0	1	0	...	1	...	$\{0, 2, \dots, n, \dots\}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
$f(n)$	0	0	1	0	...	1	...	$\{2, \dots, n, \dots\}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

$$n \in S^* \text{ iff } n \notin f(n)$$

	0	1	2	3	...	n	...	$S \subseteq \mathbb{N}$
$\varphi_0(x)$	0	0	0	0	...	0	...	\emptyset
$\varphi_1(x)$	0	1	0	1	...	1	...	$\{1, 3, \dots, n, \dots\}$
$\varphi_2(x)$	0	1	1	0	...	0	...	$\{1, 2\}$
$\varphi_3(x)$	1	0	1	0	...	1	...	$\{0, 2, \dots, n, \dots\}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
$\varphi_n(x)$	0	0	1	0	...	1	...	$\{2, \dots, n, \dots\}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

$n \in S^*$ iff $n \notin$ set defined by $\varphi_n(x)$

	0	1	2	3	...	n	...	$S \subseteq \mathbb{N}$
$\varphi_0(x)$	0	0	0	0	...	0	...	\emptyset
$\varphi_1(x)$	0	1	0	1	...	1	...	$\{1, 3, \dots, n, \dots\}$
$\varphi_2(x)$	0	1	1	0	...	0	...	$\{1, 2\}$
$\varphi_3(x)$	1	0	1	0	...	1	...	$\{0, 2, \dots, n, \dots\}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
$\varphi_n(x)$	0	0	1	0	...	1	...	$\{2, \dots, n, \dots\}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

$n \in S^*$ iff $n \notin$ set defined by $\varphi_n(x)$

Suppose that S^* is definable in our language (say by $\varphi_m(x)$)

	0	1	2	3	...	n	...	$S \subseteq \mathbb{N}$
$\varphi_0(x)$	0	0	0	0	...	0	...	\emptyset
$\varphi_1(x)$	0	1	0	1	...	1	...	$\{1, 3, \dots, n, \dots\}$
$\varphi_2(x)$	0	1	1	0	...	0	...	$\{1, 2\}$
$\varphi_3(x)$	1	0	1	0	...	1	...	$\{0, 2, \dots, n, \dots\}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
$\varphi_n(x)$	0	0	1	0	...	1	...	$\{2, \dots, n, \dots\}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

$n \in S^*$ iff $n \notin$ set defined by $\varphi_n(x)$

Write $\varphi_m(\bar{n})$ for “ $\varphi_m(x)$ is true of n ”

	0	1	2	3	...	n	...	$S \subseteq \mathbb{N}$
$\varphi_0(x)$	0	0	0	0	...	0	...	\emptyset
$\varphi_1(x)$	0	1	0	1	...	1	...	$\{1, 3, \dots, n, \dots\}$
$\varphi_2(x)$	0	1	1	0	...	0	...	$\{1, 2\}$
$\varphi_3(x)$	1	0	1	0	...	1	...	$\{0, 2, \dots, n, \dots\}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots
$\varphi_n(x)$	0	0	1	0	...	1	...	$\{2, \dots, n, \dots\}$
\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots	\vdots

$n \in S^*$ iff $n \notin$ set defined by $\varphi_n(x)$

$$\varphi_m(\bar{n}) \leftrightarrow \neg \text{True}(\ulcorner \varphi_n(\bar{n}) \urcorner)$$

where $\ulcorner \varphi_n(\bar{n}) \urcorner$ is the term in the language representing the code of $\varphi_n(\bar{n})$

D-Liar

$$\varphi_m(\overline{m}) \leftrightarrow \neg \text{True}(\ulcorner \varphi_m(\overline{m}) \urcorner)$$

“ m is true of $\varphi_m(x)$ iff it is not true that m is true of $\varphi_m(x)$ ”

Gödel's Idea

$$\varphi_m(\overline{m}) \leftrightarrow \neg \text{True}(\ulcorner \varphi_m(\overline{m}) \urcorner)$$

Gödel's Idea

$$\varphi_m(\overline{m}) \leftrightarrow \neg \text{True}(\ulcorner \varphi_m(\overline{m}) \urcorner)$$

$$\varphi_m(\overline{m}) \leftrightarrow \neg \text{Prov}(\ulcorner \varphi_m(\overline{m}) \urcorner)$$

" $\varphi_m(\overline{m})$ is true iff $\varphi_m(\overline{m})$ is not provable."

Gödel's Idea

$$\varphi_m(\overline{m}) \leftrightarrow \neg \text{True}(\ulcorner \varphi_m(\overline{m}) \urcorner)$$

$$\varphi_m(\overline{m}) \leftrightarrow \neg \text{Prov}(\ulcorner \varphi_m(\overline{m}) \urcorner)$$

" $\varphi_m(\overline{m})$ is true iff $\varphi_m(\overline{m})$ is not provable."

$$\varphi_m(\overline{m}) \leftrightarrow \neg \text{Prov}(\ulcorner \varphi_m(\overline{m}) \urcorner)$$

$$\varphi_m(\overline{m}) \leftrightarrow \neg \text{Prov}(\ulcorner \varphi_m(\overline{m}) \urcorner)$$

$\varphi_m(\overline{m})$ is not provable: Suppose $\varphi_m(\overline{m})$ is provable. Then, since we can only prove true statements, $\varphi_m(\overline{m})$ is true. This means that $\neg \text{Prov}(\ulcorner \varphi_m(\overline{m}) \urcorner)$ is true. So, $\varphi_m(\overline{m})$ is not provable. Contradiction.

$$\varphi_m(\overline{m}) \leftrightarrow \neg \text{Prov}(\ulcorner \varphi_m(\overline{m}) \urcorner)$$

$\varphi_m(\overline{m})$ is not provable: Suppose $\varphi_m(\overline{m})$ is provable. Then, since we can only prove true statements, $\varphi_m(\overline{m})$ is true. This means that $\neg \text{Prov}(\ulcorner \varphi_m(\overline{m}) \urcorner)$ is true. So, $\varphi_m(\overline{m})$ is not provable. Contradiction.

$\neg \varphi_m(\overline{m})$ is not provable: Suppose that $\neg \varphi_m(\overline{m})$ is provable. Since our system only proves true statements, $\neg \varphi_m(\overline{m})$ is true. Then $\neg \neg \text{Prov}(\ulcorner \varphi_m(\overline{m}) \urcorner)$ is true. So, $\varphi_m(\overline{m})$ is provable. This contradicts the assumption that the system is consistent.

$$\varphi_m(\overline{m}) \leftrightarrow \neg \text{Prov}(\ulcorner \varphi_m(\overline{m}) \urcorner)$$

$\varphi_m(\overline{m})$ is not provable: Suppose $\varphi_m(\overline{m})$ is provable. Then, since we can only prove true statements, $\varphi_m(\overline{m})$ is true. This means that $\neg \text{Prov}(\ulcorner \varphi_m(\overline{m}) \urcorner)$ is true. So, $\varphi_m(\overline{m})$ is not provable. Contradiction.

$\neg \varphi_m(\overline{m})$ is not provable: Suppose that $\neg \varphi_m(\overline{m})$ is provable. Since our system only proves true statements, $\neg \varphi_m(\overline{m})$ is true. Then $\neg \neg \text{Prov}(\ulcorner \varphi_m(\overline{m}) \urcorner)$ is true. So, $\varphi_m(\overline{m})$ is provable. This contradicts the assumption that the system is consistent.

Conclusion: Neither $\varphi_m(\overline{m})$ nor $\neg \varphi_m(\overline{m})$ is provable.

$$\varphi_m(\overline{m}) \leftrightarrow \neg \text{Prov}(\ulcorner \varphi_m(\overline{m}) \urcorner)$$

1. Apply Richard's move to Cantor's construction to get the D-Liar
2. Replace 'true' with 'provable' on the right-hand side of the sentence
3. Proceed with the difficult task of *arithmetizing syntax* to construct the right-side of the sentence ($\text{Prov}(v)$).
4. Show that the above sentence is provable within the formal system eliminating any appeal to the concept of "truth". The assumption that provable implies truth is replaced with (ω -)consistency.

H. Gaifman (2006). *Naming and Diagonalization, From Cantor to Gödel to Kleene*. Logic Journal of the IGPL, pp. 709 - 728.

Naming systems

Naming systems are intended as a basic framework for studying situations in which functions can be applied to their names....In a naming system we do not specify how the names are attached to functions, we assume only that there is such a correlation and that it satisfies certain minimal requirements.

H. Gaifman (2006). *Naming and Diagonalization, From Cantor to Gödel to Kleene*. Logic Journal of the IGPL, pp. 709 - 728.

Naming systems I

$$\mathcal{D} = (D, \text{type}, \{ \})$$

such that:

- ▶ D is a non-empty set.
- ▶ type assigns to each $a \in D$ its type: $\text{type}(a)$ tells us if a is a name (of a function) and, if it is, the function's arity.

A name of arity n , or n -ary name, is one that names an n -ary function.

Types can be construed as tuples: (0) —if a is not a name, $(1, n)$ —if it is an n -ary name.

- ▶ $\{ \}$ is a mapping that assigns to every n -ary name, a , a function:

$$\{a\} : D^n \rightarrow D$$

Naming systems II

- ▶ There is at least one named function of arity greater than 0

Naming systems II

- ▶ There is at least one named function of arity greater than 0
- ▶ Substitution of names (SN): If f is an n -ary named function, where $n > 0$, then, for every name a :

$$\lambda x_2, \dots, x_n f(a, x_2, \dots, x_n) \text{ is named}$$

Naming systems II

- ▶ There is at least one named function of arity greater than 0
- ▶ Substitution of names (SN): If f is an n -ary named function, where $n > 0$, then, for every name a :

$$\lambda x_2, \dots, x_n f(a, x_2, \dots, x_n) \text{ is named}$$

- ▶ Variable permutation (VP): If f is an n -ary named function, where $n > 0$, and π is a permutation of $\{1, \dots, n\}$, then

$$\lambda x_1, \dots, x_n f(x_{\pi(1)}, x_{\pi(2)}, \dots, x_{\pi(n)}) \text{ is named}$$

n -Diagonal Function

For $n > 0$, an n -diagonal function, denoted dl_n , is a function that maps each n -ary name a to a name of the function:

$$\lambda x_2, \dots, x_n \{a\}(a, x_2, \dots, x_n)$$

Thus, $dl_n(a)$ is the name of the above function.

n -Diagonal Function

For $n > 0$, an n -diagonal function, denoted dl_n , is a function that maps each n -ary name a to a name of the function:

$$\lambda x_2, \dots, x_n \{a\}(a, x_2, \dots, x_n)$$

Thus, $dl_n(a)$ is the name of the above function.

For all n -ary names a ,

$$\{dl_n(a)\}(x_2, \dots, x_n) = \{a\}(a, x_2, \dots, x_n)$$

General Fixed-Point Theorem

GFP Theorem. If F is an $(n + 1)$ -ary named function, $n \geq 0$, and the composition $F(dl_{n+1}(x_0), x_1, \dots, x_n)$ is named, then there is an n -ary name, e , such that:

$$\{e\}(x_1, \dots, x_n) = F(e, x_1, \dots, x_n)$$

General Fixed-Point Theorem

GFP Theorem. If F is an $(n + 1)$ -ary named function, $n \geq 0$, and the composition $F(dl_{n+1}(x_0), x_1, \dots, x_n)$ is named, then there is an n -ary name, e , such that:

$$\{e\}(x_1, \dots, x_n) = F(e, x_1, \dots, x_n)$$

Proof. Let c be the name of $F(dl_{n+1}(x_0), x_1, \dots, x_n)$ and let $e = dl_{n+1}(c)$. Then for and $\vec{x} = (x_1, \dots, x_n)$,

General Fixed-Point Theorem

GFP Theorem. If F is an $(n + 1)$ -ary named function, $n \geq 0$, and the composition $F(dl_{n+1}(x_0), x_1, \dots, x_n)$ is named, then there is an n -ary name, e , such that:

$$\{e\}(x_1, \dots, x_n) = F(e, x_1, \dots, x_n)$$

Proof. Let c be the name of $F(dl_{n+1}(x_0), x_1, \dots, x_n)$ and let $e = dl_{n+1}(c)$. Then for and $\vec{x} = (x_1, \dots, x_n)$,

$$\{e\}(\vec{x}) = \{dl_{n+1}(c)\}(\vec{x}) \quad (\text{definition of } e)$$

General Fixed-Point Theorem

GFP Theorem. If F is an $(n + 1)$ -ary named function, $n \geq 0$, and the composition $F(dl_{n+1}(x_0), x_1, \dots, x_n)$ is named, then there is an n -ary name, e , such that:

$$\{e\}(x_1, \dots, x_n) = F(e, x_1, \dots, x_n)$$

Proof. Let c be the name of $F(dl_{n+1}(x_0), x_1, \dots, x_n)$ and let $e = dl_{n+1}(c)$. Then for and $\vec{x} = (x_1, \dots, x_n)$,

$$\begin{aligned}\{e\}(\vec{x}) &= \{dl_{n+1}(c)\}(\vec{x}) && \text{(definition of } e) \\ &= \{c\}(c, \vec{x}) && \text{(definition of } dl_{n+1}(c))\end{aligned}$$

General Fixed-Point Theorem

GFP Theorem. If F is an $(n + 1)$ -ary named function, $n \geq 0$, and the composition $F(dl_{n+1}(x_0), x_1, \dots, x_n)$ is named, then there is an n -ary name, e , such that:

$$\{e\}(x_1, \dots, x_n) = F(e, x_1, \dots, x_n)$$

Proof. Let c be the name of $F(dl_{n+1}(x_0), x_1, \dots, x_n)$ and let $e = dl_{n+1}(c)$. Then for and $\vec{x} = (x_1, \dots, x_n)$,

$$\begin{aligned}\{e\}(\vec{x}) &= \{dl_{n+1}(c)\}(\vec{x}) && \text{(definition of } e) \\ &= \{c\}(c, \vec{x}) && \text{(definition of } dl_{n+1}(c)) \\ &= F(dl_{n+1}(c), \vec{x}) && \text{(definition of } c)\end{aligned}$$

General Fixed-Point Theorem

GFP Theorem. If F is an $(n + 1)$ -ary named function, $n \geq 0$, and the composition $F(dl_{n+1}(x_0), x_1, \dots, x_n)$ is named, then there is an n -ary name, e , such that:

$$\{e\}(x_1, \dots, x_n) = F(e, x_1, \dots, x_n)$$

Proof. Let c be the name of $F(dl_{n+1}(x_0), x_1, \dots, x_n)$ and let $e = dl_{n+1}(c)$. Then for and $\vec{x} = (x_1, \dots, x_n)$,

$$\begin{aligned}\{e\}(\vec{x}) &= \{dl_{n+1}(c)\}(\vec{x}) && \text{(definition of } e\text{)} \\ &= \{c\}(c, \vec{x}) && \text{(definition of } dl_{n+1}(c)\text{)} \\ &= F(dl_{n+1}(c), \vec{x}) && \text{(definition of } c\text{)} \\ &= F(e, \vec{x}) && \text{(definition of } e\text{)}\end{aligned}$$

- ✓ Gödel numbering
- ✓ Gödel-Carnap Fixed Point Theorem
- ✓ (Naming systems)
- ▶ Representing functions/relations

Representability

Definition

Suppose that $f : \mathbb{N}^k \rightarrow \mathbb{N}$. We say that f is **representable** in \mathbf{Q} when there is a formula $A_f(x_0, \dots, x_{k-1}, y)$ such that for all $n_0, \dots, n_{k-1} \in \mathbb{N}$: if $f(n_0, \dots, n_{k-1}) = m$ then

1. $\mathbf{Q} \vdash A_f(\overline{n_0}, \dots, \overline{n_{k-1}}, \overline{m})$
2. $\mathbf{Q} \vdash \forall y (A_f(\overline{n_0}, \dots, \overline{n_{k-1}}, y) \rightarrow y = \overline{m})$

Equivalent definitions of representability

- ▶ f is representable in \mathbf{Q} iff there is a formula $A_f(x_0, \dots, x_{k-1}, y)$ such that for all $n_0, \dots, n_{k-1} \in \mathbb{N}$, if $f(n_0, \dots, n_{k-1}) = m$ then:

$$\mathbf{Q} \vdash \forall y (A_f(\overline{n_0}, \dots, \overline{n_{k-1}}, y) \leftrightarrow y = \overline{m})$$

Equivalent definitions of representability

- ▶ f is representable in \mathbf{Q} iff there is a formula $A_f(x_0, \dots, x_{k-1}, y)$ such that for all $n_0, \dots, n_{k-1} \in \mathbb{N}$, if $f(n_0, \dots, n_{k-1}) = m$ then:

$$\mathbf{Q} \vdash \forall y (A_f(\overline{n_0}, \dots, \overline{n_{k-1}}, y) \leftrightarrow y = \overline{m})$$

- ▶ f is representable in \mathbf{Q} iff there is a formula $A_f(x_0, \dots, x_{k-1}, y)$ such that for all $n_0, \dots, n_{k-1} \in \mathbb{N}$:
 1. If $f(n_0, \dots, n_{k-1}) = m$, then $\mathbf{Q} \vdash A_f(\overline{n_0}, \dots, \overline{n_{k-1}}, \overline{m})$
 2. If $f(n_0, \dots, n_{k-1}) \neq m$, then $\mathbf{Q} \vdash \neg A_f(\overline{n_0}, \dots, \overline{n_{k-1}}, \overline{m})$

Equivalent definitions of representability

- ▶ f is representable in \mathbf{Q} iff there is a formula $A_f(x_0, \dots, x_{k-1}, y)$ such that for all $n_0, \dots, n_{k-1} \in \mathbb{N}$, if $f(n_0, \dots, n_{k-1}) = m$ then:

$$\mathbf{Q} \vdash \forall y (A_f(\overline{n_0}, \dots, \overline{n_{k-1}}, y) \leftrightarrow y = \overline{m})$$

- ▶ f is representable in \mathbf{Q} iff there is a formula $A_f(x_0, \dots, x_{k-1}, y)$ such that for all $n_0, \dots, n_{k-1} \in \mathbb{N}$:
 1. If $f(n_0, \dots, n_{k-1}) = m$, then $\mathbf{Q} \vdash A_f(\overline{n_0}, \dots, \overline{n_{k-1}}, \overline{m})$
 2. If $f(n_0, \dots, n_{k-1}) \neq m$, then $\mathbf{Q} \vdash \neg A_f(\overline{n_0}, \dots, \overline{n_{k-1}}, \overline{m})$
- ▶ f is representable in \mathbf{Q} iff there is a formula $A_f(x_0, \dots, x_{k-1}, y)$ such that for all $n_0, \dots, n_{k-1} \in \mathbb{N}$:
 1. if $f(n_0, \dots, n_{k-1}) = m$ then $\mathbf{Q} \vdash A_f(\overline{n_0}, \dots, \overline{n_{k-1}}, \overline{m})$
 2. $\mathbf{Q} \vdash \exists! y A_f(\overline{n_0}, \dots, \overline{n_{k-1}}, y)$

Exercise

Prove that all of the definitions of representability are equivalent.

Representing Relations

A relation $R \subseteq \mathbb{N}^k$ is **representable** in \mathbf{Q} provided that the characteristic function χ_R is representable in \mathbf{Q} . It is not hard to see that this is equivalent to saying that $R \subseteq \mathbb{N}^k$ is representable in \mathbf{Q} provided that there is a formula A_R such that for all $n_0, \dots, n_{k-1} \in \mathbb{N}$:

1. if $(n_0, \dots, n_{k-1}) \in R$, then $\mathbf{Q} \vdash A_R(\overline{n_0}, \dots, \overline{n_{k-1}})$
2. if $(n_0, \dots, n_{k-1}) \notin R$, then $\mathbf{Q} \vdash \neg A_R(\overline{n_0}, \dots, \overline{n_{k-1}})$

All of the following relations are representable in **Q**:

- ▶ $Sent(x)$: x is the Gödel number of a sentence of \mathcal{L}_A
- ▶ $Form(x)$: x is the Gödel number of a formula of \mathcal{L}_A
- ▶ $Term(x)$: x is the Gödel number of a term of \mathcal{L}_A
- ▶ $Axiom(x)$: x is the Gödel number of an axiom of **Q**
- ▶ $Prf_{\mathbf{PA}}(x, y)$: x is the Gödel number of a derivation in **PA** of a formula with Gödel number y .
- ▶ ...

Plan

- ✓ Introduction: Smullyan's Machine
- ✓ Background
 - ✓ Formal Arithmetic
 - ✓ Gödel's Incompleteness Theorems
 - ✓ Names and Gödel numbering
 - ✓ Fixed Point Theorem
- ▶ Provability predicate and Löb's Theorem
- ▶ Provability logic
- ▶ Truth predicate and Tarski's Theorem
- ▶ A Primer on Epistemic and Doxastic Logic
- ▶ Anti-Expert Paradoxes
- ▶ Predicate approach to modality
- ▶ The Knower Paradox and variants
- ▶ Epistemic Arithmetic
- ▶ Gödel's Disjunction

Proof Predicate

The proof relation $Prf_{\mathbf{PA}}(x, y)$ is represented by a formula $\text{Prf}_{\mathbf{PA}}$.

Proof Predicate

The proof relation $Prf_{\mathbf{PA}}(x, y)$ is represented by a formula $\text{Prf}_{\mathbf{PA}}$.

The *proof predicate*, denoted $\text{Prov}_{\mathbf{PA}}(y)$, is defined as follows:

$$\exists x \text{Prf}_{\mathbf{PA}}(x, y)$$

Derivability Conditions

It can be shown that the provability predicate $\text{Prov}_{\mathbf{PA}}$ satisfies the following:

D1. If $\mathbf{PA} \vdash A$, then $\mathbf{PA} \vdash \text{Prov}_{\mathbf{PA}}(\ulcorner A \urcorner)$

D2. $\mathbf{PA} \vdash \text{Prov}_{\mathbf{PA}}(\ulcorner A \rightarrow B \urcorner) \rightarrow (\text{Prov}_{\mathbf{PA}}(\ulcorner A \urcorner) \rightarrow \text{Prov}_{\mathbf{PA}}(\ulcorner B \urcorner))$

D3. $\mathbf{PA} \vdash \text{Prov}_{\mathbf{PA}}(\ulcorner A \urcorner) \rightarrow \text{Prov}_{\mathbf{PA}}(\ulcorner \text{Prov}_{\mathbf{PA}}(\ulcorner A \urcorner) \urcorner)$

Derivability Conditions

A provability predicate for \mathbf{T} , denoted $\text{Prov}_{\mathbf{T}}$, satisfies the following:

D1. If $\mathbf{T} \vdash A$, then $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner)$

D2. $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \rightarrow B \urcorner) \rightarrow (\text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \rightarrow \text{Prov}_{\mathbf{T}}(\ulcorner B \urcorner))$

D3. $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \rightarrow \text{Prov}_{\mathbf{T}}(\ulcorner \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \urcorner)$

Reflection Principle

The reflection principle for **T** is the schema

$$\text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \rightarrow A$$

Monotonicity Inference for the Provability Predicate

Lemma

For any theory \mathbf{T} , if $\text{Prov}_{\mathbf{T}}$ satisfies $D1$ and $D2$, then:

From $\mathbf{T} \vdash A \rightarrow B$, infer $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \rightarrow \text{Prov}(\ulcorner B \urcorner)$.

Löb's Theorem

Theorem (Löb's Theorem)

Let \mathbf{T} be an axiomatizable theory extending \mathbf{Q} , and suppose $\text{Prov}_{\mathbf{T}}(y)$ is a formula satisfying conditions $D1$ - $D3$.

If $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \rightarrow A$, then $\mathbf{T} \vdash A$.

Suppose A is a sentence such that $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \rightarrow A$. Let $B(y)$ be the formula

$$\text{Prov}_{\mathbf{T}}(y) \rightarrow A$$

Suppose A is a sentence such that $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \rightarrow A$. Let $B(y)$ be the formula

$$\text{Prov}_{\mathbf{T}}(y) \rightarrow A$$

By the Fixed-Point Theorem, there is a sentence D such that

$$\mathbf{T} \vdash D \leftrightarrow B(\ulcorner D \urcorner)$$

Suppose that $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \rightarrow A$.

Suppose A is a sentence such that $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \rightarrow A$. Let $B(y)$ be the formula

$$\text{Prov}_{\mathbf{T}}(y) \rightarrow A$$

By the Fixed-Point Theorem, there is a sentence D such that

$$\mathbf{T} \vdash D \leftrightarrow B(\ulcorner D \urcorner)$$

Suppose that $\mathbf{T} \vdash \text{Prov}_{\mathbf{T}}(\ulcorner A \urcorner) \rightarrow A$.

To simplify the notation, we write $\text{Prov}(\cdot)$ instead of $\text{Prov}_{\mathbf{T}}$

1. $D \leftrightarrow (\text{Prov}(\ulcorner D \urcorner) \rightarrow A)$ FPT
2. $\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \rightarrow A \urcorner)$ Lemma: 1
3. $\text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \rightarrow A \urcorner) \rightarrow (\text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner))$ D2
4. $\text{Prov}(\ulcorner D \urcorner) \rightarrow (\text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner))$ PC: 2, 3

- | | | |
|----------|---|----------|
| 1. | $D \leftrightarrow (\text{Prov}(\ulcorner D \urcorner) \rightarrow A)$ | FPT |
| \vdots | \vdots | \vdots |
| 4. | $\text{Prov}(\ulcorner D \urcorner) \rightarrow (\text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner))$ | PC: 2, 3 |
| 5. | $\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner)$ | D3 |

- | | | |
|----------|---|----------|
| 1. | $D \leftrightarrow (\text{Prov}(\ulcorner D \urcorner) \rightarrow A)$ | FPT |
| \vdots | \vdots | \vdots |
| 4. | $\text{Prov}(\ulcorner D \urcorner) \rightarrow (\text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner))$ | PC: 2, 3 |
| 5. | $\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner)$ | D3 |
| 6. | $\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner)$ | PC: 4, 5 |

- | | | |
|----------|---|------------|
| 1. | $D \leftrightarrow (\text{Prov}(\ulcorner D \urcorner) \rightarrow A)$ | FPT |
| \vdots | \vdots | \vdots |
| 4. | $\text{Prov}(\ulcorner D \urcorner) \rightarrow (\text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner))$ | PC: 2, 3 |
| 5. | $\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner)$ | D3 |
| 6. | $\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner)$ | PC: 4, 5 |
| 7. | $\text{Prov}(\ulcorner A \urcorner) \rightarrow A$ | Assumption |

1.	$D \leftrightarrow (\text{Prov}(\ulcorner D \urcorner) \rightarrow A)$	FPT
\vdots	\vdots	\vdots
4.	$\text{Prov}(\ulcorner D \urcorner) \rightarrow (\text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner))$	PC: 2, 3
5.	$\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner)$	D3
6.	$\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner)$	PC: 4, 5
7.	$\text{Prov}(\ulcorner A \urcorner) \rightarrow A$	Assumption
8.	$\text{Prov}(\ulcorner D \urcorner) \rightarrow A$	PC: 6, 7

- | | | |
|----------|---|------------|
| 1. | $D \leftrightarrow (\text{Prov}(\ulcorner D \urcorner) \rightarrow A)$ | FPT |
| \vdots | \vdots | \vdots |
| 4. | $\text{Prov}(\ulcorner D \urcorner) \rightarrow (\text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner))$ | PC: 2, 3 |
| 5. | $\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner)$ | D3 |
| 6. | $\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner)$ | PC: 4, 5 |
| 7. | $\text{Prov}(\ulcorner A \urcorner) \rightarrow A$ | Assumption |
| 8. | $\text{Prov}(\ulcorner D \urcorner) \rightarrow A$ | PC: 6, 7 |
| 9. | D | PC: 1, 8 |

1.	$D \leftrightarrow (\text{Prov}(\ulcorner D \urcorner) \rightarrow A)$	FPT
\vdots	\vdots	\vdots
4.	$\text{Prov}(\ulcorner D \urcorner) \rightarrow (\text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner))$	PC: 2, 3
5.	$\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner)$	D3
6.	$\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner)$	PC: 4, 5
7.	$\text{Prov}(\ulcorner A \urcorner) \rightarrow A$	Assumption
8.	$\text{Prov}(\ulcorner D \urcorner) \rightarrow A$	PC: 6, 7
9.	D	PC: 1, 8
10.	$\text{Prov}(\ulcorner D \urcorner)$	D1 from 9

1.	$D \leftrightarrow (\text{Prov}(\ulcorner D \urcorner) \rightarrow A)$	FPT
\vdots	\vdots	\vdots
4.	$\text{Prov}(\ulcorner D \urcorner) \rightarrow (\text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner))$	PC: 2, 3
5.	$\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner)$	D3
6.	$\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner)$	PC: 4, 5
7.	$\text{Prov}(\ulcorner A \urcorner) \rightarrow A$	Assumption
8.	$\text{Prov}(\ulcorner D \urcorner) \rightarrow A$	PC: 6, 7
9.	D	PC: 1, 8
10.	$\text{Prov}(\ulcorner D \urcorner)$	D1 from 9
11.	A	PC: 8, 10

1.	$D \leftrightarrow (\text{Prov}(\ulcorner D \urcorner) \rightarrow A)$	FPT
\vdots	\vdots	\vdots
4.	$\text{Prov}(\ulcorner D \urcorner) \rightarrow (\text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner))$	PC: 2, 3
5.	$\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner \text{Prov}(\ulcorner D \urcorner) \urcorner)$	D3
6.	$\text{Prov}(\ulcorner D \urcorner) \rightarrow \text{Prov}(\ulcorner A \urcorner)$	PC: 4, 5
7.	$\text{Prov}(\ulcorner A \urcorner) \rightarrow A$	Assumption
8.	$\text{Prov}(\ulcorner D \urcorner) \rightarrow A$	PC: 6, 7
9.	D	PC: 1, 8
10.	$\text{Prov}(\ulcorner D \urcorner)$	D1 from 9
11.	A	PC: 8, 10

'PA couldn't be more modest about its own veracity'

By Löb's Theorem, it is not true that for all sentences φ ,

$$\mathbf{PA} \vdash \text{Prov}(\ulcorner \text{Prov}(\ulcorner \varphi \urcorner) \rightarrow \varphi \urcorner)$$

'PA couldn't be more modest about its own veracity'

By Löb's Theorem, it is not true that for all sentences φ ,

$$\mathbf{PA} \vdash \text{Prov}(\ulcorner \text{Prov}(\ulcorner \varphi \urcorner) \rightarrow \varphi \urcorner)$$

Statement

$\mathbf{PA} \vdash \text{Prov}(\ulcorner \varphi \urcorner)$
implies $\mathbf{PA} \vdash \varphi$

It is not true that...

$\mathbf{PA} \vdash \text{Prov}(\ulcorner \varphi \urcorner) \rightarrow \varphi$

'PA couldn't be more modest about its own veracity'

By Löb's Theorem, it is not true that for all sentences φ ,

$$\mathbf{PA} \vdash \text{Prov}(\ulcorner \text{Prov}(\ulcorner \varphi \urcorner) \rightarrow \varphi \urcorner)$$

Statement

$\mathbf{PA} \vdash \text{Prov}(\ulcorner \varphi \urcorner)$
implies $\mathbf{PA} \vdash \varphi$

$\mathbf{PA} \vdash \text{Prov}(\ulcorner \neg \varphi \urcorner)$
implies $\mathbf{PA} \not\vdash \text{Prov}(\ulcorner \varphi \urcorner)$

It is not true that...

$\mathbf{PA} \vdash \text{Prov}(\ulcorner \varphi \urcorner) \rightarrow \varphi$

$\mathbf{PA} \vdash \text{Prov}(\ulcorner \neg \varphi \urcorner) \rightarrow \neg \text{Prov}(\ulcorner \varphi \urcorner)$

'PA couldn't be more modest about its own veracity'

By Löb's Theorem, it is not true that for all sentences φ ,

$$\mathbf{PA} \vdash \text{Prov}(\ulcorner \text{Prov}(\ulcorner \varphi \urcorner) \rightarrow \varphi \urcorner)$$

Statement

$\mathbf{PA} \vdash \text{Prov}(\ulcorner \varphi \urcorner)$
implies $\mathbf{PA} \vdash \varphi$

$\mathbf{PA} \vdash \text{Prov}(\ulcorner \neg \varphi \urcorner)$
implies $\mathbf{PA} \not\vdash \text{Prov}(\ulcorner \varphi \urcorner)$

$\mathbf{PA} \vdash \text{Prov}(\ulcorner \neg \text{Prov}(\ulcorner \varphi \urcorner) \urcorner)$
implies $\mathbf{PA} \vdash \neg \text{Prov}(\varphi)$

It is not true that...

$\mathbf{PA} \vdash \text{Prov}(\ulcorner \varphi \urcorner) \rightarrow \varphi$

$\mathbf{PA} \vdash \text{Prov}(\ulcorner \neg \varphi \urcorner) \rightarrow \neg \text{Prov}(\ulcorner \varphi \urcorner)$

$\mathbf{PA} \vdash \text{Prov}(\ulcorner \neg \text{Prov}(\ulcorner \varphi \urcorner) \urcorner) \rightarrow \neg \text{Prov}(\ulcorner \varphi \urcorner)$

Plan

- ✓ Introduction: Smullyan's Machine
- ✓ Background
 - ✓ Formal Arithmetic
 - ✓ Gödel's Incompleteness Theorems
 - ✓ Names and Gödel numbering
 - ✓ Fixed Point Theorem
- ✓ Provability predicate and Löb's Theorem
 - ▶ Provability logic
 - ▶ Truth predicate and Tarski's Theorem
 - ▶ A Primer on Epistemic and Doxastic Logic
 - ▶ Anti-Expert Paradoxes
 - ▶ Predicate approach to modality
 - ▶ The Knower Paradox and variants
 - ▶ Epistemic Arithmetic
 - ▶ Gödel's Disjunction

Rineke Verbrugge (2024). *Provability Logic*. The Stanford Encyclopedia of Philosophy (Summer 2024 Edition), Edward N. Zalta & Uri Nodelman (eds.), <https://plato.stanford.edu/archives/sum2024/entries/logic-provability/>.

Propositional Modal Logic

$$p \mid \neg\varphi \wedge \varphi \wedge \psi \mid \Box\varphi$$

where $p \in AT$ (at set of atomic propositions).

The intended interpretation of $\Box\varphi$ is “there is a proof (in **PA**) of φ ”.

Propositional Modal Logic

$$p \mid \neg\varphi \wedge \varphi \wedge \psi \mid \Box\varphi$$

where $p \in AT$ (at set of atomic propositions).

The intended interpretation of $\Box\varphi$ is “there is a proof (in **PA**) of φ ”.

A **frame** is a tuple (W, R) such that $W \neq \emptyset$ and $R \subseteq W \times W$.

Propositional Modal Logic

$$p \mid \neg\varphi \wedge \varphi \wedge \psi \mid \Box\varphi$$

where $p \in \text{AT}$ (at set of atomic propositions).

The intended interpretation of $\Box\varphi$ is “there is a proof (in **PA**) of φ ”.

A **frame** is a tuple (W, R) such that $W \neq \emptyset$ and $R \subseteq W \times W$.

A **model** is a tuple (W, R, V) where (W, R) is a frame and $V : \text{AT} \rightarrow \wp(W)$.

Truth/Validity

For a model $\mathcal{M} = (W, R, V)$ and $w \in W$, we write $\mathcal{M} \models \varphi$ when φ is true at w in \mathcal{M} .

- ▶ $\mathcal{M}, w \models p$ iff $w \in V(p)$
- ▶ $\mathcal{M}, w \models \neg\varphi$ iff $\mathcal{M}, w \not\models \varphi$
- ▶ $\mathcal{M}, w \models \varphi \wedge \psi$ iff $\mathcal{M}, w \models \varphi$ and $\mathcal{M}, w \models \psi$
- ▶ $\mathcal{M}, w \models \Box\varphi$ iff for all $v \in W$, if $w R v$, then $\mathcal{M}, v \models \varphi$

For a frame $\mathcal{F} = (W, R)$, φ is **valid on \mathcal{F}** , denoted $\mathcal{F} \models \varphi$, when $\mathcal{M}, w \models \varphi$ for all models \mathcal{M} based on \mathcal{F} and $w \in W$.

Provability Logic: **GL**

$$\text{K} \quad \Box(\varphi \rightarrow \psi) \rightarrow (\Box\varphi \rightarrow \psi)$$

$$\text{L} \quad \Box(\Box\varphi \rightarrow \varphi) \rightarrow \Box\varphi$$

$$\text{MP} \quad \varphi, \varphi \rightarrow \psi \therefore \psi$$

$$\text{NEC} \quad \varphi \therefore \Box\varphi$$

Some Results

► **GL** $\vdash \Box\varphi \rightarrow \Box\Box\varphi$.

Some Results

- ▶ **GL** $\vdash \Box\varphi \rightarrow \Box\Box\varphi$.
- ▶ $\Box(\Box\varphi \rightarrow \varphi) \rightarrow \Box\varphi$ is valid on a frame (W, R) if, and only if, R is transitive and converse well-founded (there are no infinite ascending sequences, that is sequences of the form $w_1 R w_2 R w_3 \cdots$).

Some Results

- ▶ **GL** $\vdash \Box\varphi \rightarrow \Box\Box\varphi$.
- ▶ $\Box(\Box\varphi \rightarrow \varphi) \rightarrow \Box\varphi$ is valid on a frame (W, R) if, and only if, R is transitive and converse well-founded (there are no infinite ascending sequences, that is sequences of the form $w_1 R w_2 R w_3 \dots$).
- ▶ The logic **GL** is not compact:

$$\Gamma = \{\Diamond p_0, \Box(p_0 \rightarrow \Diamond p_1), \Box(p_1 \rightarrow \Diamond p_2), \dots, \Box(p_n \rightarrow \Diamond p_{n+1}), \dots\}.$$

is finitely satisfiable, but not satisfiable.

Some Results

- ▶ **GL** $\vdash \Box\varphi \rightarrow \Box\Box\varphi$.
- ▶ $\Box(\Box\varphi \rightarrow \varphi) \rightarrow \Box\varphi$ is valid on a frame (W, R) if, and only if, R is transitive and converse well-founded (there are no infinite ascending sequences, that is sequences of the form $w_1 R w_2 R w_3 \dots$).
- ▶ The logic **GL** is not compact:

$$\Gamma = \{\Diamond p_0, \Box(p_0 \rightarrow \Diamond p_1), \Box(p_1 \rightarrow \Diamond p_2), \dots, \Box(p_n \rightarrow \Diamond p_{n+1}), \dots\}.$$

is finitely satisfiable, but not satisfiable.

- ▶ The logic **GL** is sound and weakly complete with respect to the class of frames that are transitive and converse well-founded.

Arithmetic Completeness

An **arithmetic translation** is a function t such that

1. For all $p \in \text{At}$, $t(p)$ is a sentence of \mathcal{L}_A
2. t commutes with the boolean connectives: $t(\neg\varphi) = \neg t(\varphi)$, $t(\varphi \wedge \psi) = t(\varphi) \wedge t(\psi)$, etc.
3. $t(\Box\varphi) = \text{Prov}_{\mathbf{PA}}(\ulcorner t(\varphi) \urcorner)$

Theorem (Solovay 1976).

GL $\vdash \varphi$ iff for every arithmetic translation t , **PA** $\vdash t(\varphi)$.