

Provability Logic

Benjamin Church

July 3, 2022

Contents

1	Introduction	2
1.1	First-Order Languages	2
1.2	Proof Theory	3
2	A Theory For the Natural Numbers	4
2.1	The Language of Number Theory	4
2.2	Robinson Arithmetic and Peano Arithmetic	4
2.3	Representing Functions and Relations	5
3	Computability Theory	5
3.1	μ -Recursive Functions	6
3.2	Recursive and Recursively Enumerable Sets	6
3.2.1	Church-Turing Thesis	6
3.3	The Representability Theorem	6
4	Number Theory Swallows Itself	7
4.1	Gödel Numbering	7
4.2	The Provability Predicate	8
4.3	Self-Reference	8
4.4	Gödel Incompleteness I	10
4.5	Löb's Theorem	11
4.6	Gödel Incompleteness II	12
4.7	Löb's Theorem Formalized inside Number Theory	12
4.8	Gödel Incompleteness Formalized inside Number Theory	13
5	GL Provability Logic	15
5.1	Modal Logic	15
5.2	Modal Semantics	16
5.3	Arithmetic Soundness	16
5.4	The Existence of Modal Fixed Points	16
5.5	Arithmetic Completeness	16

1 Introduction

1.1 First-Order Languages

Definition 1.1.1. A *vocabulary* or *signature* σ is a set of “non-logical” symbols which may be of three types:

- (a) Constant symbols (e.g. 0)
- (b) n-ary function symbols (e.g. +)
- (c) n-ary relation symbols (e.g. \in)

Along with the signature, a first-order language has a set of “logical” symbols:

- (a) A countable list of variable symbols: x_1, x_2, x_3, \dots
- (b) Logical connectives: $\neg, \vee, \wedge, \rightarrow$
- (c) Quantifiers: \forall (we get $\exists \iff \neg\forall\neg$ for free)
- (d) An equality relation: $=$
- (e) Punctuation: (), etc.

Definition 1.1.2. The set of *terms* of a first-order language L with vocabulary σ is defined inductively as follows:

- (a) Any variable or constant symbol is a term.
- (b) If f is an n-ary function symbol and t_1, \dots, t_n are terms then $f(t_1, \dots, t_n)$ is a term. For a binary operator (2-ary function), say \circ , we will often write $(t_1 \circ t_2)$ to mean $\circ(t_1, t_2)$.

Definition 1.1.3. The set of *formulas* of a first-order language L with vocabulary σ is defined inductively as follows:

- (a) If s, t are terms then $(s = t)$ is a formula. Furthermore if $R \in \sigma$ is an n-ary relation symbol and t_1, \dots, t_n are terms then $R(t_1, \dots, t_n)$ is a formula. For a 2-ary relation we will often write sRt to mean $R(s, t)$.
- (b) If A and B are formulas then $\neg A$, $(A \vee B)$, $(A \wedge B)$, and $(A \rightarrow B)$ are all formulas.
- (c) If x is a variable symbol and φ a formula in which x is free (φ contains x but no quantifiers over x) then $\forall x \varphi$ and $\exists x \varphi$ are formulas.

Definition 1.1.4. A *sentence* of a first-order language is a formula with no free variables.

Definition 1.1.5. A first-order theory is a first-order language L along with a set Γ of first-order L -sentences which are referred to as axioms.

1.2 Proof Theory

There are many possible first-order deduction systems each with its own unique flavor. A deduction system has logical axioms and rules of inference on formulas of L . A formal proof beginning with some assumptions is a sequence of L -formulas each of which is either a logical axiom, an assumption, or the result of a rule of inference applied to previous formulas. Here we work with an example which is a variant of Hilbert's propositional logic formal system H extended to first order logic.

Definition 1.2.1. Hilbert's system H has logical connectives $\{\neg, \rightarrow\}$ and the following axiom schemas: for any formulas A, B, C the following are axioms of H ,

$$(H1) \ A \rightarrow (B \rightarrow A)$$

$$(H2) \ (A \rightarrow (B \rightarrow C)) \rightarrow ((A \rightarrow B) \rightarrow (A \rightarrow C))$$

$$(H3) \ \neg A \rightarrow (A \rightarrow B)$$

$$(H4) \ (\neg A \rightarrow A) \rightarrow A$$

The formal system H has one rule of inference known as modus ponens (MP),

$$\frac{A \quad (A \rightarrow B)}{B}$$

We can define the formulas $A \vee B$ to stand for $\neg A \rightarrow B$ and $A \wedge B$ to stand for $\neg(A \rightarrow \neg B)$ and $A \leftrightarrow B$ to stand for $(A \rightarrow B) \wedge (B \rightarrow A)$ etc.

Definition 1.2.2. We say that a first-order theory Γ *syntactically entails* or, more simply, *proves* A if there exists a formal proof using axioms of Γ and first-order rules of inference. We write this as $\Gamma \vdash A$.

Example 1.2.3. We show that $\vdash_H A \rightarrow A$ for any formula A .

$$\begin{array}{ll} [A \rightarrow ((A \rightarrow A) \rightarrow A)] \rightarrow [(A \rightarrow (A \rightarrow A)) \rightarrow (A \rightarrow A)] & \text{axiom (L2)} \quad (1) \\ A \rightarrow ((A \rightarrow A) \rightarrow A) & \text{axiom (L1)} \quad (2) \\ (A \rightarrow (A \rightarrow A)) \rightarrow (A \rightarrow A) & \text{MP 1,2} \quad (3) \\ A \rightarrow (A \rightarrow A) & \text{axiom (L1)} \quad (4) \\ A \rightarrow A & \text{MP 3,4} \quad (5) \end{array}$$

Remark. Clearly, proofs in H are horrible. Luckily the following wonderful theorem means we will rarely need to provide explicit proofs.

Theorem 1.2.4 (Gödel). Every propositional tautology is a theorem of H .

Remark. By tautology here we mean always evaluates to true under the standard semantics for \neg and \rightarrow . In these semantics all axioms of H are tautologies and modus ponens is locally sound.

Definition 1.2.5. The formal system FO extends H by adding the additional axioms,

$$(EQ1) \ \forall x (x = x)$$

$$(EQ2) \ \forall x [(x = t) \rightarrow (A(x) \rightarrow A(t))]$$

(FO1) $\forall x A(x) \rightarrow A(t)$ where t is any term whose variables are not bound in A .

(FO2) $\forall x (A \rightarrow B(x)) \rightarrow (A \rightarrow \forall x B(x))$ where x is free in x and not in A .

and the additional rule of inference called generalization (Gen),

$$\frac{A}{\forall x A}$$

Remark. We will use the notation $A(x)$ to denote that x is a free variable in A and then $A(t)$ to denote A with t substituted for x .

Definition 1.2.6. A first-order theory Γ is *consistent* if there does not exist a statement A such that $\Gamma \vdash A$ and $\Gamma \vdash \neg A$.

Definition 1.2.7. A first-order theory Γ is *complete* if for every L -sentence A we have either $\Gamma \vdash A$ or $\Gamma \vdash \neg A$.

Lemma 1.2.8 (Categorization of Consistency). Γ is proof-theoretically consistent if and only if there exists a first-order sentence A such that $\Gamma \not\vdash A$.

Proof. If Γ is consistent and $\Gamma \vdash A$ then $\Gamma \not\vdash \neg A$. If Γ is not consistent then $\Gamma \vdash A$ and $\Gamma \vdash \neg A$ for some A . Using (H3), $\Gamma \vdash \neg A \rightarrow (A \rightarrow B)$ so applying MP twice gives $\Gamma \vdash B$ for any B . \square

2 A Theory For the Natural Numbers

2.1 The Language of Number Theory

Definition 2.1.1. The first-order language L_{NN} has signature $\sigma = \{\mathbf{0}, s, +, \cdot\}$ where $\mathbf{0}$ is a constant symbol, s is a 1-ary function symbol, and $+$ and \cdot are 2-ary function symbols.

Example 2.1.2. We be define the abbreviation $x < y$ to mean $\exists z (x + s(z) = y)$.

Definition 2.1.3. For each natural number $i \in \mathbb{N}$ we denote the term $s^i(\mathbf{0})$ by the bold-face numeral \mathbf{i} .

2.2 Robinson Arithmetic and Peano Arithmetic

Now that we have a first-order language in which to do number theory, we need an actual theory.

Definition 2.2.1. Robinson Arithmetic, denoted as \mathbf{Q} , is the first-order theory over L_{NN} with the set of axioms,

(Q1) $\forall x \neg(s(x) = \mathbf{0})$

(Q2) $\forall x \forall y [(s(x) = s(y)) \rightarrow (x = y)]$

(Q3) $\forall x (x + \mathbf{0} = x)$

(Q4) $\forall x \forall y (x + s(y) = s(x + y))$

(Q5) $\forall x (x \cdot \mathbf{0} = \mathbf{0})$

$$(Q6) \ \forall x \forall y (x \cdot s(y) = (x \cdot y) + x)$$

$$(Q7) \ \forall x [(x = \mathbf{0}) \vee \exists y (x = s(y))]$$

Remark. We see that **Q** is arithmetic without induction. You might think that we cannot do very much in **Q** since it is a very weak theory. However **Q** is sufficiently powerful to cause its own essential incompleteness. In fact, **Q** is the minimal theory necessary to prove the representability theorem. For completeness, we will now define the more familiar framework for number theory.

Definition 2.2.2. Peano Arithmetic (**PA**) is the first-order theory over L_{NN} which has axioms (Q1) - (Q6) and additionally the axiom schema of induction,

$$(\text{PA}) \ \varphi(\mathbf{0}) \rightarrow [\forall x (\varphi(x) \rightarrow \varphi(s(x))) \rightarrow \forall x \varphi(x)]$$

for each formula φ with x free. Note we have dropped (Q7) since it is a consequence of the induction axiom.

Definition 2.2.3. An extension of **Q** is a first-order theory Γ over the language L_{NN} such that $\Gamma \vdash \mathbf{Q}$, in particular if $\Gamma \supset \mathbf{Q}$.

Remark. Clearly **PA** is an extension of **Q**. In fact, the extension is proper.

2.3 Representing Functions and Relations

Definition 2.3.1. A relation $R \subset \mathbb{N}^n$ is *strongly representable* or simply *representable* in Γ , an extension of **Q** if there exists a formula $A(x_1, \dots, x_n)$ in L_{NN} with n free variables such that for all natural numbers $a_1, \dots, a_n \in \mathbb{N}$ we have,

$$\begin{aligned} R(a_1, \dots, a_n) &\implies \Gamma \vdash A(\mathbf{a}_1, \dots, \mathbf{a}_n) \\ \neg R(a_1, \dots, a_n) &\implies \Gamma \vdash \neg A(\mathbf{a}_1, \dots, \mathbf{a}_n) \end{aligned}$$

In this case we say that A *represents* R over Γ .

Definition 2.3.2. An arithmetic function $f : \mathbb{N}^n \rightarrow \mathbb{N}$ is representable over Γ iff there exists a formula $A(x_1, \dots, x_n, x_{n+1})$ of L_{NN} with $n+1$ free variables such that for all natural numbers $a_1, \dots, a_n \in \mathbb{N}$ with $b = f(a_1, \dots, a_n)$ we have,

$$\Gamma \vdash \forall x [A(\mathbf{a}_1, \dots, \mathbf{a}_n, x) \leftrightarrow (x = \mathbf{b})]$$

Remark. A function being representable is equivalent to its graph G_f being representable.

Definition 2.3.3. A relation $R \subset \mathbb{N}^n$ is *weakly representable* in Γ if there exists a formula $A(x_1, \dots, x_n)$ in L_{NN} with n free variables such that for all natural numbers $a_1, \dots, a_n \in \mathbb{N}$ we have,

$$R(a_1, \dots, a_n) \iff \Gamma \vdash A(\mathbf{a}_1, \dots, \mathbf{a}_n)$$

In this case we say that A *weakly represents* R over Γ .

Lemma 2.3.4. If Γ is consistent then weak representability implies representability.

Proof. It suffices to show that if $\Gamma \vdash A(\mathbf{a}_1, \dots, \mathbf{a}_n)$ then $R(a_1, \dots, a_n)$. Indeed, by consistency, $\Gamma \not\vdash \neg A(\mathbf{a}_1, \dots, \mathbf{a}_n)$ so therefore $R(a_1, \dots, a_n)$. \square

3 Computability Theory

We would like to construct representable functions. It turns out that there is a deep connection between computability and representability. More generally, the incompleteness theorems rely on arithmetic capturing the power of computable functions.

3.1 μ -Recursive Functions

The notion of *computability* or an *effective procedure* for computing a function is not a well-defined notion. We begin with a concrete definition for a class of clearly computable arithmetic functions. It turns out that in some sense these are *all* the computable functions.

Definition 3.1.1. An arithmetic function $F : \mathbb{N}^n \rightarrow \mathbb{N}$ is *recursive* if F is one of,

- (a) a starting function: addition $((a, b) \mapsto a + b)$, multiplication (\cdot) , projection $(U_{n,k}(a_1, \dots, a_n) = a_k)$, or less-than characteristic $(K_<(a, b) = 1 \text{ if } a < b \text{ and zero otherwise})$.
- (b) a compositions of recursive functions $F = G \circ (H_1, \dots, H_k)$
- (c) a minimalization of a regular recursive function

$$F(a_1, \dots, a_n) = \mu x[G(a_1, \dots, a_n, x) = 0]$$

where the regularity condition on G means that such a zero is always required to exist for all natural numbers $a_1, \dots, a_n \in \mathbb{N}$.

3.2 Recursive and Recursively Enumerable Sets

Definition 3.2.1. A relation $R \subset \mathbb{N}^n$ is *recursive* (R) if there exists a recursive arithmetic function $f : \mathbb{N}^n \rightarrow \mathbb{N}$ such that $R = \{(a_1, \dots, a_n) \in \mathbb{N}^n \mid f(a_1, \dots, a_n) = 0\}$.

Definition 3.2.2. A relation $R \subset \mathbb{N}^n$ is *recursively enumerable* (RE) if R can be written as $R(a_1, \dots, a_n) \iff \exists x Q(a_1, \dots, a_n, x)$ where $Q \subset \mathbb{N}^n$ is a recursive relation.

Proposition 3.2.3. A set $S \subset \mathbb{N}$ is RE iff it is enumerated by a recursive function.

Remark. This proposition explains the terminology *recursively enumerable*.

3.2.1 Church-Turing Thesis

There is no clear universally agreed upon *a priori* definition for what it means for a function to be *effectively computable*. However, logicians Alonzo Church and Alan Turing proved that a wide class of models of computation (μ -recursive functions, Turning machines, λ -calculi) are all equivalently powerful. Therefore, we define *effectively computable* functions to be exactly those computable by any of these equivalent models of computation. Often, we will invoke this thesis to show that a given function is recursive if we can find an informal effective procedure for computing it. It should be stressed that such a use of the Church-Turing thesis is never necessary for proving metalogical theorems it is simply a time-saving device for lazy logicians who don't want to explicitly construct recursive functions. It is only strictly necessary to invoke the Church-Turing thesis when computability is assumed as a hypothesis since we must develop a formal proof using some explicit model of computation.

3.3 The Representability Theorem

Theorem 3.3.1. Let $f : \mathbb{N}^n \rightarrow \mathbb{N}$ be recursive function then f is representable over \mathbf{Q} .

Proof. Very technical but conceptually easy. Show that all starting functions are representable and that given representable functions that we can construct representations of their composition and minimization. \square

Corollary 3.3.2. Let $R \subset \mathbb{N}^n$ be a recursive relation then R is representable over \mathbf{Q} .

Proof. There exists a recursive $f : \mathbb{N}^n \rightarrow \mathbb{N}$ such that f vanishes exactly on R . Then f is representable by some L_{NN} formula $A(x_1, \dots, x_{n+1})$ such that for all natural numbers $a_1, \dots, a_n \in \mathbb{N}$ and $b = f(a_1, \dots, a_n)$ then,

$$\mathbf{Q} \vdash \forall x [A(\mathbf{a}_1, \dots, \mathbf{a}_n, x) \leftrightarrow (x = \mathbf{b})]$$

Let $B(x_1, \dots, x_n) = A(x_1, \dots, x_n, \mathbf{0})$. Then I claim that,

$$\begin{aligned} R(a_1, \dots, a_n) &\implies \Gamma \vdash B(\mathbf{a}_1, \dots, \mathbf{a}_n) \\ \neg R(a_1, \dots, a_n) &\implies \Gamma \vdash \neg B(\mathbf{a}_1, \dots, \mathbf{a}_n) \end{aligned}$$

and thus B represents R . \square

4 Number Theory Swallows Itself

4.1 Gödel Numbering

We need some way of expressing the metalanguage of formulas and proofs inside of number theory such that we can use number theory to prove statements of its own meta-theory. This is accomplished by encoding formulas as natural numbers.

Theorem 4.1.1. There exists an injective function $\#_g : \text{FOR}_{L_{\text{NN}}} \rightarrow \mathbb{N}$ such that its image $S = \text{Im}(\#_g)$ is a recursive set.

Proof. Consider encoding each symbol as a unique integer and then a sequence of symbols via $p_1^{a_1} \cdots p_n^{a_n}$ where p_i is the i^{th} prime and a_i is the code of the i^{th} symbol. By uniqueness of prime factorization, this function is injective. Checking its image is recursive is highly technical so I will simply invoke the Church-Turing thesis since there exists an effective procedure to factor a number, translate it into a string of symbols, and check if this string can be produced by the rules for forming well-formed formulas. The last step is effectively computable because there are a finite number of formulas of the correct length or less (restricting to only the variables which appear in the target string) so we can simply try each. \square

Remark. The function $\#_g$ encodes each formula as a natural number such that the set of codes corresponding to well-formed formulas is computable.

Definition 4.1.2. Let A be a formula and $a = \#_g(A)$ its Gödel number. Then let $\ulcorner A \urcorner$ be the term \mathbf{a} .

Remark. This notation is intentionally suggestive of quotation in natural language. In fact, the Gödel sentence is not best described as saying “I am provable” but rather the Quine sentence,

“when preceded by its quotation is unprovable”
when preceded by its quotation is unprovable.

which is self-referential since the object of the sentence (“when preceded by its quotation is unprovable” when preceded by its quotation) is a copy of the entire sentence. This sentence accomplishes self-reference without the self-referential “machinery” of the pronoun “I” and therefore is a much better model for how such self-reference can unintentionally arise in number theory.

4.2 The Provability Predicate

Definition 4.2.1. A theory Γ with language L_{NN} is *recursively axiomatized* if $\#_g(\Gamma)$ is recursive.

Remark. Intuitively, a theory Γ is axiomatized if there exists an algorithm which can decide if a given string is an axiom of the theory.

Theorem 4.2.2. Let Γ be recursively axiomatized. We may extend $\#_g^\Gamma : \text{PRF}_\Gamma \rightarrow \mathbb{N}$ to encoding valid Γ -proofs as a sequence of formulas which, using the technique used above, we can encode in a single number. Again, we require that g_Γ be injective and have recursive image such that the codes of valid proofs comprise a computable set. Furthermore the relation, $\text{CHKPRF}_\Gamma \subset \mathbb{N}^2$ defined to contain (a, p) iff a is the code of a valid formula and p is the code of a valid proof of the formula encoded by a is a recursive relation.

Proof. We rely here on the Church-Turing thesis to show that such relations are recursive. They are effectively computable since checking a proof requires only checking each line to see if it is an axiom (which is decidable by hypothesis) or the result of applying one of finitely many rules of inference to the finitely many preceding sentences. This is clearly computable. \square

Definition 4.2.3. Since CHKPRF_Γ is recursive it is Γ -representable. Let $\mathcal{P}rf_\Gamma(x, y)$ be a formula of L_{NN} such that,

$$\begin{aligned} \text{CHKPRF}_\Gamma(a, p) &\implies \Gamma \vdash \mathcal{P}rf_\Gamma(\mathbf{a}, \mathbf{p}) \\ \neg \text{CHKPRF}_\Gamma(a, p) &\implies \Gamma \vdash \neg \mathcal{P}rf_\Gamma(\mathbf{a}, \mathbf{p}) \end{aligned}$$

Definition 4.2.4. The provability predicate $\mathcal{B}ew_\Gamma(x)$ is the formula $\exists p \mathcal{P}rf_\Gamma(x, p)$.

Remark. The notation $\mathcal{B}ew$ derives from the German word *Beweis* for proof.

Lemma 4.2.5. If $\Gamma \vdash A$ then $\Gamma \vdash \mathcal{B}ew_\Gamma(\Gamma A \neg)$.

Proof. If $\Gamma \vdash A$ then there exists a proof of A which has code p and let A have code a . Therefore, $\text{CHKPRF}_\Gamma(a, p)$ so $\Gamma \vdash \mathcal{P}rf_\Gamma(\mathbf{a}, \mathbf{p})$. Now the axiom (FO1) gives,

$$\Gamma \vdash \forall y \neg \mathcal{P}rf_\Gamma(\mathbf{a}, y) \rightarrow \neg \mathcal{P}rf_\Gamma(\mathbf{a}, \mathbf{p})$$

Thus, taking the contrapositive,

$$\Gamma \vdash \mathcal{P}rf_\Gamma(\mathbf{a}, \mathbf{p}) \rightarrow \mathcal{B}ew_\Gamma(\Gamma A \neg)$$

so by modus ponens $\Gamma \vdash \mathcal{B}ew_\Gamma(\Gamma A \neg)$. \square

Remark. We will see in the following sections that under the additionaly hypothesis of ω -consistency the provability predicate $\mathcal{B}ew_\Gamma(x)$ actually weakly represents theoremhood. However, we will also see that theoremhood is not strongly representable.

4.3 Self-Reference

Lemma 4.3.1 (Diagonalization). Let $F(x)$ be an L_{NN} formula with one free variable. Then there exists a ‘fixed-point’ sentence ψ such that,

$$\mathbf{Q} \vdash \psi \leftrightarrow F(\ulcorner \psi \urcorner)$$

Proof. There exists a recursive function $d : \mathbb{N} \rightarrow \mathbb{N}$ such that when $a = \#_g(A)$ where $A(x)$ is a formula with at least one free variable then $d(a) = \#_g(A(\mathbf{a})) = \#_g(A(\ulcorner A \urcorner))$ (for now we appeal to the Church-Turing thesis). Therefore, D is represented by some formula $D(x, y)$ such that for all $a \in \mathbb{N}$ and $b = d(a)$ we have,

$$\mathbf{Q} \vdash \forall y [D(\mathbf{a}, y) \leftrightarrow (y = \mathbf{b})]$$

Now define the formula with one free variable,

$$\varphi := \forall y [D(x, y) \rightarrow F(y)]$$

Let $a = \#_g(\varphi)$ be its Gödel number and then substitute $\mathbf{a} = \ulcorner \varphi \urcorner$ for x in φ ,

$$\psi := \varphi(\ulcorner \varphi \urcorner) := \forall y [D(\ulcorner \varphi \urcorner, y) \rightarrow F(y)]$$

The Gödel number of ψ is $q = \#_g(\varphi(\ulcorner \varphi \urcorner)) = d(a)$ so we apply the representation of d applied at $d(a) = q$,

$$\mathbf{Q} \vdash \forall y [D(\ulcorner \varphi \urcorner, y) \leftrightarrow (y = \ulcorner \varphi(\ulcorner \varphi \urcorner) \urcorner)]$$

Using the tautology,

$$\mathbf{Q} \vdash (A \leftrightarrow B) \rightarrow [(A \rightarrow C) \leftrightarrow (B \rightarrow C)]$$

we find,

$$\mathbf{Q} \vdash \forall y [D(\ulcorner \varphi \urcorner, y) \rightarrow F(y)] \leftrightarrow \forall y [(y = \ulcorner \varphi(\ulcorner \varphi \urcorner) \urcorner) \rightarrow F(y)]$$

Which we can write as,

$$\mathbf{Q} \vdash \varphi(\ulcorner \varphi \urcorner) \leftrightarrow F(\ulcorner \varphi(\ulcorner \varphi \urcorner) \urcorner)$$

and using $\psi := \varphi(\ulcorner \varphi \urcorner)$ we have,

$$\mathbf{Q} \vdash \psi \leftrightarrow F(\psi)$$

□

Remark. If we interpret $F(\ulcorner \psi \urcorner)$ to represent “the formula ψ has property F ” then the diagonal lemma proves the existence of self-referential fixed points. The sentence $\psi \leftrightarrow F(\ulcorner \psi \urcorner)$ “says” that ψ is true if and only if ψ has property F . In other words, ψ has an interpretation as the sentence: “I have property F .” As described earlier, the diagonal sentence is more accurately modeled in natural language as the Quine sentence,

“when preceded by its quotation has property F ”
when preceded by its quotation has property F .

In fact, the above proof of the diagonalization lemma closely resembles the construction of a Quine sentence: we take a sentence which refers to its object applied to (preceded by) its own quotation and apply it to (preceding it by) its own quotation. The predicate $\varphi(x)$ encodes “ x when applied to its quotation (Gödel number) has property F ” and the self-referential statement ψ is exactly φ applied to its quotation.

4.4 Gödel Incompleteness I

In this and the following sections, let \perp stand for your favorite contradiction, say $(0 = 1)$ or $(x = y) \wedge \neg(x = y)$ etc. Any choice is as good as any other as long as $\Gamma \vdash \perp$ implies that Γ is inconsistent (which the above certainly do).

Definition 4.4.1. A theory Γ is ω -consistent if for all formulas $A(x)$ with one free variable Γ cannot simultaneously prove $\exists x A(x)$ and $\neg A(\mathbf{n})$ for each natural number $n \in \mathbb{N}$.

Lemma 4.4.2. ω -consistency implies consistency.

Proof. For each formula with one free variable $A(x)$ either $\Gamma \not\vdash \exists x A(x)$ or for some $n \in \mathbb{N}$ we have $\Gamma \not\vdash \neg A(\mathbf{n})$. Therefore, there exists some formula that Γ cannot prove which implies that Γ is consistent. \square

Lemma 4.4.3. If Γ is ω -consistent and $\Gamma \not\vdash A$ then $\Gamma \not\vdash \mathcal{Bew}_\Gamma(\Gamma A^\perp)$. In particular,

$$\Gamma \vdash A \iff \Gamma \vdash \mathcal{Bew}_\Gamma(\Gamma A^\perp)$$

meaning that $\mathcal{Bew}_\Gamma(x)$ weakly represents the theoremhood relation on formulas.

Proof. Suppose that $\Gamma \not\vdash A$ and $a = g(A)$ is the Gödel number. Then for each $n \in \mathbb{N}$ we have $\neg \text{CHKPRF}_\Gamma(a, n)$ since there exist no valid proofs of A . Therefore we have $\Gamma \vdash \neg \mathcal{Prf}_\Gamma(\mathbf{a}, \mathbf{n})$ for each $n \in \mathbb{N}$ so by ω -consistency $\Gamma \not\vdash \mathcal{Bew}_\Gamma(\Gamma A^\perp)$. \square

Corollary 4.4.4. If Γ is ω -consistent then Γ is consistent so $\Gamma \not\vdash \perp$ and thus, by the previous lemma, $\Gamma \not\vdash \mathcal{Bew}_\Gamma(\Gamma \perp^\perp)$ and thus $\Gamma \not\vdash \mathcal{Bew}_\Gamma(\Gamma \mathcal{Bew}_\Gamma(\Gamma \perp^\perp)^\perp)$ etc.

Theorem 4.4.5 (Gödel). Any ω -consistent recursively axiomatized extension of \mathbf{Q} is incomplete. In particular, if Γ is a recursively axiomatized extension of \mathbf{Q} then there exists a sentence \mathcal{G}_Γ such that,

- (a) if Γ is consistent then $\Gamma \not\vdash \mathcal{G}_\Gamma$
- (b) if Γ is ω -consistent then $\Gamma \not\vdash \neg \mathcal{G}_\Gamma$.

Proof. Let Γ be a consistent recursively axiomatized extension of \mathbf{Q} . Since Γ is recursively axiomatized \mathcal{Prf}_Γ and \mathcal{Bew}_Γ exist. The fixed-point theorem proves the existence of a sentence \mathcal{G}_Γ such that,

$$\Gamma \vdash \mathcal{G}_\Gamma \leftrightarrow \neg \mathcal{Bew}_\Gamma(\Gamma \mathcal{G}_\Gamma^\perp)$$

Suppose that $\Gamma \vdash \mathcal{G}_\Gamma$ then $\Gamma \vdash \mathcal{Bew}_\Gamma(\Gamma \mathcal{G}_\Gamma^\perp)$. However, using $\Gamma \vdash \mathcal{G}_\Gamma$ and the self-reference equivalence, $\Gamma \vdash \neg \mathcal{Bew}_\Gamma(\Gamma \mathcal{G}_\Gamma^\perp)$ contradicting the consistency of Γ .

Suppose that $\Gamma \vdash \neg \mathcal{G}_\Gamma$. By the consistency of Γ we cannot have $\Gamma \vdash \mathcal{G}_\Gamma$. However, $\Gamma \vdash \neg \mathcal{G}_\Gamma$ and self-reference shows that $\Gamma \vdash \mathcal{Bew}_\Gamma(\Gamma \mathcal{G}_\Gamma^\perp)$. Because $\Gamma \not\vdash \mathcal{G}_\Gamma$, this contradicts the ω -consistency of Γ . \square

Remark. The sentence \mathcal{G}_Γ expresses “I am not provable” through Quinian self-reference. This is captured formally through $\Gamma \vdash \mathcal{G}_\Gamma \leftrightarrow \neg \mathcal{Bew}_\Gamma(\Gamma \mathcal{G}_\Gamma^\perp)$.

4.5 Löb's Theorem

Remark. In this section we assume that Γ is a recursively axiomatized extension of **PA**.

Lemma 4.5.1 (Hilbert-Bernays-Löb). The provability predicate satisfies,

- (a) $\Gamma \vdash A \implies \Gamma \vdash \mathcal{B}ew_\Gamma(\Gamma A^\neg)$
- (b) $\Gamma \vdash \mathcal{B}ew_\Gamma(\Gamma A \rightarrow B^\neg) \rightarrow (\mathcal{B}ew_\Gamma(\Gamma A^\neg) \rightarrow \mathcal{B}ew_\Gamma(\Gamma B^\neg))$
- (c) $\Gamma \vdash \mathcal{B}ew_\Gamma(\Gamma A^\neg) \rightarrow \mathcal{B}ew_\Gamma(\Gamma \mathcal{B}ew_\Gamma(\Gamma A^\neg)^\neg)$

Proof. SKETCH THIS PROOF!!!! □

Remark. The first Hilbert-Bernays derivability condition states that $\mathcal{B}ew_\Gamma(x)$ *weakly* represents theoremhood (it cannot strongly represent it however as we shall show). The second condition states that modus ponens is *provably (within Γ)* a rule of inference of Γ . Finally, the third Hilbert-Bernays derivability condition is the formalization of the first property *within the system Γ* , saying that Γ can prove that if it can prove A then it can prove that it can prove A .

Theorem 4.5.2 (Löb). If $\Gamma \vdash \mathcal{B}ew_\Gamma(\Gamma A^\neg) \rightarrow A$ then $\Gamma \vdash A$ for any sentence A .

Proof. Via the fixed point theorem applied to $\mathcal{B}ew_\Gamma(x) \rightarrow A$, there exists a sentence B such that,

$$\Gamma \vdash B \leftrightarrow (\mathcal{B}ew_\Gamma(\Gamma B^\neg) \rightarrow A)$$

Applying HB1 to one direction gives,

$$\Gamma \vdash \mathcal{B}ew_\Gamma(\Gamma B \rightarrow (\mathcal{B}ew_\Gamma(\Gamma B^\neg) \rightarrow A)^\neg)$$

and then applying HB2 twice we deduce,

$$\Gamma \vdash \mathcal{B}ew_\Gamma(\Gamma B^\neg) \rightarrow (\mathcal{B}ew_\Gamma(\Gamma \mathcal{B}ew_\Gamma(\Gamma B^\neg)^\neg) \rightarrow \mathcal{B}ew_\Gamma(\Gamma A^\neg))$$

However, HB3 gives,

$$\Gamma \vdash \mathcal{B}ew_\Gamma(\Gamma B^\neg) \rightarrow \mathcal{B}ew_\Gamma(\Gamma \mathcal{B}ew_\Gamma(\Gamma B^\neg)^\neg)$$

and thus putting the previous two together,

$$\Gamma \vdash \mathcal{B}ew_\Gamma(\Gamma B^\neg) \rightarrow \mathcal{B}ew_\Gamma(\Gamma A^\neg)$$

Now we use the hypothesis $\Gamma \vdash \mathcal{B}ew_\Gamma(\Gamma A^\neg) \rightarrow A$ to get a proof,

$$\Gamma \vdash \mathcal{B}ew_\Gamma(\Gamma B^\neg) \rightarrow A$$

but since $\mathcal{B}ew_\Gamma(\Gamma B^\neg) \rightarrow A$ is provably equivalent to B we find $\Gamma \vdash B$ so by HB1 $\Gamma \vdash \mathcal{B}ew_\Gamma(\Gamma B^\neg)$ and thus $\Gamma \vdash A$ by modus ponens. □

Remark. This theorem is truly remarkable because it says that **Q** and all extensions are “maximally modest” in the sense that they do not “believe” in their own validity (i.e. a proof of A entails A) except for statements they already know to be true. Furthermore, it answers the fascinating question posed by Henkin.

Remark. After seeing Gödel's proof of the first incompleteness theorem Henkin asked about a subtle modification. What if we apply the fixed-point lemma not to $\neg \mathcal{B}ew_{\Gamma}(x)$ but to simply $\mathcal{B}ew_{\Gamma}(x)$? Then there would exist a sentence \mathcal{H} ,

$$\Gamma \vdash \mathcal{H} \leftrightarrow \mathcal{B}ew_{\Gamma}(\Gamma \mathcal{H} \neg)$$

This sentence has the interpretation “I am provable” which seems to convey no information at all! However, clearly for such a sentence we have,

$$\Gamma \vdash \mathcal{B}ew_{\Gamma}(\Gamma \mathcal{H} \neg) \rightarrow \mathcal{H}$$

and thus by Löb's theorem we get $\Gamma \vdash \mathcal{H}$. So in fact, such a Henkin sentence which asserts its own provability must actually be provable.

4.6 Gödel Incompleteness II

Finally, Löb's theme gives us enough machinery to give an elegant proof of the second incompleteness theorem.

Definition 4.6.1. The sentence $\mathcal{C}on_{\Gamma}$ is given by $\neg \mathcal{B}ew_{\Gamma}(\Gamma \perp \neg)$ which expresses the consistency of the theory Γ .

Remark. We have shown that if $\Gamma \vdash \neg \mathcal{C}on_{\Gamma}$ then Γ is not ω -consistent. However, we are about to show a much more interesting result.

Theorem 4.6.2 (Gödel). Let Γ be a consistent recursively axiomatized extension of \mathbf{Q} then Γ cannot prove $\mathcal{C}on_{\Gamma}$.

Proof. By Löb's theorem if $\Gamma \vdash \mathcal{B}ew_{\Gamma}(\Gamma \perp \neg) \rightarrow \perp$ then $\Gamma \vdash \perp$. However, $\mathcal{B}ew_{\Gamma}(\Gamma \perp \neg) \rightarrow \perp$ is equivalent to $\mathcal{C}on_{\Gamma}$. Thus if $\Gamma \vdash \mathcal{C}on_{\Gamma}$ then $\Gamma \vdash \perp$ contradicting the consistency of Γ . Taking the contrapositive, $\Gamma \not\vdash \perp \implies \Gamma \not\vdash \mathcal{C}on_{\Gamma}$ i.e. the consistency of Γ implies that Γ cannot prove $\mathcal{C}on_{\Gamma}$. \square

Remark. Gödel's second incompleteness theorem is often stated provocatively as: a theory's proof of its own consistency establishes its inconsistency. This makes sense because an inconsistent theory can prove anything including its own consistency.

4.7 Löb's Theorem Formalized inside Number Theory

Wonderfully, we can formalize the proof of Löb's theorem inside the system Γ so that we may apply Löb *inside* formal proofs.

Theorem 4.7.1 (Löb). For any sentence A of $L_{\mathbf{NN}}$,

$$\Gamma \vdash \mathcal{B}ew_{\Gamma}(\Gamma \mathcal{B}ew_{\Gamma}(\Gamma A \neg) \rightarrow A \neg) \rightarrow \mathcal{B}ew_{\Gamma}(\Gamma A \neg)$$

Proof. Let $B := \mathcal{B}ew_{\Gamma}(\Gamma \mathcal{B}ew_{\Gamma}(\Gamma A \neg) \rightarrow A \neg)$ and $C := \mathcal{B}ew_{\Gamma}(\Gamma A \neg)$. Then HB2 gives,

$$\Gamma \vdash \mathcal{B}ew_{\Gamma}(\Gamma B \rightarrow C \neg) \rightarrow (\mathcal{B}ew_{\Gamma}(\Gamma B \neg) \rightarrow \mathcal{B}ew_{\Gamma}(\Gamma C \neg))$$

Furthermore, since $B := \mathcal{B}ew_{\Gamma}(\Gamma C \rightarrow A \neg)$,

$$\Gamma \vdash B \rightarrow (\mathcal{B}ew_{\Gamma}(\Gamma C \neg) \rightarrow \mathcal{B}ew_{\Gamma}(\Gamma A \neg))$$

and by HB3 (since B begins with $\mathcal{B}ew$),

$$\Gamma \vdash B \rightarrow \mathcal{B}ew_{\Gamma}(\Gamma B \neg)$$

Given $\mathcal{B}ew_{\Gamma}(\Gamma B \rightarrow C \neg)$ we get $\mathcal{B}ew_{\Gamma}(\Gamma B \neg) \rightarrow \mathcal{B}ew_{\Gamma}(\Gamma C \neg)$. Additionally, given B we get $\mathcal{B}ew_{\Gamma}(\Gamma B \neg)$ so we get $\mathcal{B}ew_{\Gamma}(\Gamma C \neg)$ but B also gives $\mathcal{B}ew_{\Gamma}(\Gamma C \neg) \rightarrow \mathcal{B}ew_{\Gamma}(\Gamma A \neg)$ so we get $C := \mathcal{B}ew_{\Gamma}(\Gamma A \neg)$. Thus by propositional logic,

$$\Gamma \vdash \mathcal{B}ew_{\Gamma}(\Gamma B \rightarrow C \neg) \rightarrow (B \rightarrow C)$$

Therefore, applying Löb's theorem,

$$\Gamma \vdash B \rightarrow C$$

which, expanded out is,

$$\Gamma \vdash \mathcal{B}ew_{\Gamma}(\Gamma \mathcal{B}ew_{\Gamma}(\Gamma A \neg) \rightarrow A \neg) \rightarrow \mathcal{B}ew_{\Gamma}(\Gamma A \neg)$$

□

4.8 Gödel Incompleteness Formalized inside Number Theory

Much in the way that Löb theorem can be formalized inside number theory, we can formalize the proofs of the incompleteness theorems inside the formal system itself. In fact, we can further formalize the notion that consistency implies the unprovability of the Gödel sentence and thus its truth to give an alternative proof of the second incompleteness theorem and furthermore a demonstration of the provable logical equivalence of all Gödel sentences.

Theorem 4.8.1. Let \mathcal{G}_{Γ} be a Gödel sentence for Γ then,

$$\Gamma \vdash \mathcal{C}on_{\Gamma} \leftrightarrow \mathcal{G}_{\Gamma}$$

In particular, all Gödel sentences are provably logically equivalent.

Proof. Since \mathcal{G}_{Γ} is a Gödel sentence,

$$\Gamma \vdash \mathcal{G}_{\Gamma} \leftrightarrow \neg \mathcal{B}ew_{\Gamma}(\Gamma \mathcal{G}_{\Gamma} \neg)$$

Therefore, applying HB1 and HB2,

$$\Gamma \vdash \mathcal{B}ew_{\Gamma}(\Gamma \mathcal{G}_{\Gamma} \neg) \leftrightarrow \mathcal{B}ew_{\Gamma}(\Gamma \neg \mathcal{B}ew_{\Gamma}(\Gamma \mathcal{G}_{\Gamma} \neg) \neg)$$

However, by HB3,

$$\Gamma \vdash \mathcal{B}ew_{\Gamma}(\Gamma \mathcal{G}_{\Gamma} \neg) \rightarrow \mathcal{B}ew_{\Gamma}(\Gamma \mathcal{B}ew_{\Gamma}(\Gamma \mathcal{G}_{\Gamma} \neg) \neg)$$

Furthermore, since $\Gamma \vdash \neg A \rightarrow (A \rightarrow \perp)$ by HB1 and HB2 twice we get,

$$\Gamma \vdash \mathcal{B}ew_{\Gamma}(\Gamma \neg A \neg) \rightarrow (\mathcal{B}ew_{\Gamma}(\Gamma A \neg) \rightarrow \mathcal{B}ew_{\Gamma}(\Gamma \perp \neg))$$

Applying this to $A := \mathcal{B}ew_{\Gamma}(\Gamma \mathcal{G}_{\Gamma} \neg)$ we find,

$$\Gamma \vdash \mathcal{B}ew_{\Gamma}(\Gamma \mathcal{G}_{\Gamma} \neg) \rightarrow \mathcal{B}ew_{\Gamma}(\Gamma \perp \neg)$$

However, $\Gamma \vdash \perp \rightarrow \mathcal{G}_\Gamma$ and thus applying HB1 and HB2 we find,

$$\Gamma \vdash \mathcal{Bew}_\Gamma(\Gamma \perp \neg) \rightarrow \mathcal{Bew}_\Gamma(\Gamma \mathcal{G}_\Gamma \neg)$$

In summary,

$$\Gamma \vdash \mathcal{Bew}_\Gamma(\Gamma \mathcal{G}_\Gamma \neg) \leftrightarrow \mathcal{Bew}_\Gamma(\Gamma \perp \neg)$$

However, $\Gamma \vdash \mathcal{G}_\Gamma \leftrightarrow \neg \mathcal{Bew}_\Gamma(\Gamma \mathcal{G}_\Gamma \neg)$ and $\mathcal{Con}_\Gamma := \neg \mathcal{Bew}_\Gamma(\Gamma \perp \neg)$ which implies that,

$$\Gamma \vdash \mathcal{G}_\Gamma \leftrightarrow \mathcal{Con}_\Gamma$$

□

Corollary 4.8.2. If Γ is consistent then by Gödel incompleteness I we know $\Gamma \not\vdash \mathcal{G}_\Gamma$ and thus $\Gamma \not\vdash \mathcal{Con}_\Gamma$ giving an alternative proof of incompleteness II.

Theorem 4.8.3 (Formalized Gödel I).

$$\Gamma \vdash \omega\text{-}\mathcal{Con}_\Gamma \rightarrow (\neg \mathcal{Bew}_\Gamma(\Gamma \mathcal{G}_\Gamma \neg) \wedge \neg \mathcal{Bew}_\Gamma(\Gamma \neg \mathcal{G}_\Gamma \neg))$$

Where $\omega\text{-}\mathcal{Con}_\Gamma$ is the sentence $\neg \mathcal{Bew}_\Gamma(\Gamma \mathcal{Bew}_\Gamma(\Gamma \perp \neg) \neg)$ expressing weak ω -consistency.

Proof. First, by HB3,

$$\Gamma \vdash \mathcal{Bew}_\Gamma(\Gamma \perp \neg) \rightarrow \mathcal{Bew}_\Gamma(\Gamma \mathcal{Bew}_\Gamma(\Gamma \perp \neg) \neg)$$

and therefore,

$$\Gamma \vdash \omega\text{-}\mathcal{Con}_\Gamma \rightarrow \mathcal{Con}_\Gamma$$

We have already proven above that,

$$\Gamma \vdash \mathcal{Con}_\Gamma \rightarrow \neg \mathcal{Bew}_\Gamma(\Gamma \mathcal{G}_\Gamma \neg)$$

and thus by transitivity of implication,

$$\Gamma \vdash \omega\text{-}\mathcal{Con}_\Gamma \rightarrow \neg \mathcal{Bew}_\Gamma(\Gamma \mathcal{G}_\Gamma \neg)$$

The negation of the Gödel property gives,

$$\Gamma \vdash \neg \mathcal{G}_\Gamma \leftrightarrow \mathcal{Bew}_\Gamma(\Gamma \mathcal{G}_\Gamma \neg)$$

and thus by HB1 and HB2 we have,

$$\Gamma \vdash \mathcal{Bew}_\Gamma(\Gamma \neg \mathcal{G}_\Gamma \neg) \leftrightarrow \mathcal{Bew}_\Gamma(\Gamma \mathcal{Bew}_\Gamma(\Gamma \mathcal{G}_\Gamma \neg) \neg)$$

However, by HB3,

$$\Gamma \vdash \mathcal{Bew}_\Gamma(\Gamma \neg \mathcal{G}_\Gamma \neg) \rightarrow \mathcal{Bew}_\Gamma(\Gamma \mathcal{Bew}_\Gamma(\Gamma \neg \mathcal{G}_\Gamma \neg) \neg)$$

Furthermore via $\Gamma \vdash \neg \mathcal{G}_\Gamma \rightarrow (\mathcal{G}_\Gamma \rightarrow \perp)$ and HB1 and HB2 repeatedly we find,

$$\Gamma \vdash \mathcal{Bew}_\Gamma(\Gamma \mathcal{Bew}_\Gamma(\Gamma \mathcal{G}_\Gamma \neg) \neg) \rightarrow (\mathcal{Bew}_\Gamma(\Gamma \mathcal{Bew}_\Gamma(\Gamma \neg \mathcal{G}_\Gamma \neg) \neg) \rightarrow \mathcal{Bew}_\Gamma(\Gamma \mathcal{Bew}_\Gamma(\Gamma \perp \neg) \neg))$$

and thus by transitivity of implications,

$$\Gamma \vdash \mathcal{Bew}_\Gamma(\Gamma \neg \mathcal{G}_\Gamma \neg) \rightarrow \mathcal{Bew}_\Gamma(\Gamma \mathcal{Bew}_\Gamma(\Gamma \perp \neg) \neg)$$

contradicting ω -consistency. That is, taking the contrapositive,

$$\Gamma \vdash \omega\text{-}\mathcal{Con}_\Gamma \rightarrow \neg \mathcal{Bew}_\Gamma(\Gamma \neg \mathcal{G}_\Gamma \neg)$$

giving both implications which together show that,

$$\Gamma \vdash \omega\text{-}\mathcal{Con}_\Gamma \rightarrow (\neg \mathcal{Bew}_\Gamma(\Gamma \mathcal{G}_\Gamma \neg) \wedge \neg \mathcal{Bew}_\Gamma(\Gamma \neg \mathcal{G}_\Gamma \neg))$$

□

Remark. Just as we needed ω -consistency in the standard proof of Gödel Incompleteness I, in the formalized version we require a stronger hypothesis than Con_Γ , we need $\neg \text{Bew}_\Gamma(\Gamma \text{Bew}_\Gamma(\Gamma \text{G}_\Gamma))$ which expresses the idea that ω -consistency requires Γ to be unable to prove that it can prove a contradiction. In fact this hypothesis is somewhat weaker than full ω -consistency so this is an abuse of notation.

Theorem 4.8.4 (Formalized Gödel II).

$$\Gamma \vdash \text{Con}_\Gamma \rightarrow \neg \text{Bew}_\Gamma(\Gamma \text{Con}_\Gamma)$$

Proof. Apply formalized Löb with $A = \perp$ to give,

$$\Gamma \vdash \text{Bew}_\Gamma(\Gamma \text{Bew}_\Gamma(\Gamma \perp)) \rightarrow \perp \rightarrow \text{Bew}_\Gamma(\Gamma \perp)$$

However, $\text{Bew}_\Gamma(\Gamma \perp) \rightarrow \perp$ is logically equivalent to $\neg \text{Bew}_\Gamma(\Gamma \perp)$ which is Con_Γ . Furthermore if $\Gamma \vdash A \leftrightarrow B$ then $\Gamma \vdash \text{Bew}_\Gamma(\Gamma A) \leftrightarrow \text{Bew}_\Gamma(\Gamma B)$ by HB1 and HB2 so we have,

$$\Gamma \vdash \text{Bew}_\Gamma(\Gamma \text{Con}_\Gamma) \rightarrow \text{Bew}_\Gamma(\Gamma \perp)$$

which is exactly the contrapositive of formalized Gödel incompleteness II. \square

5 GL Provability Logic

5.1 Modal Logic

Modal logics are formal systems given by standard predicate calculus with a modal predicate \square which expresses some form of “necessity.” We first define the simplest so called “normal” modal logic **K** named for Saul Kripke.

Definition 5.1.1. The formal system **K** has logical connectives $\{\rightarrow, \neg, \square\}$ and sentences are built from an infinite list of propositional variables p, q, \dots . It is defined by having as axioms,

- (a) tautologies of propositional calculus (say take Hilbert’s system *H* for concrete axiomatization)
- (b) the modal distribution axiom (K),

$$\square(A \rightarrow B) \rightarrow (\square A \rightarrow \square B)$$

and as rules of inference has,

- (a) modus ponens (MP),

$$\frac{A \quad (A \rightarrow B)}{B}$$

- (b) the necessitation rule,

$$\frac{A}{\square A}$$

Remark. **K** is the basis for most modal logics, however it is too weak to capture most modal notions. For probability logic we need a stronger extension which is named **GL** for Gödel and Löb.

Definition 5.1.2. The formal system **GL** is simply **K** with the added two axioms,

- (4) $\square A \rightarrow \square \square A$
- (L) $\square(\square A \rightarrow A) \rightarrow \square A$

5.2 Modal Semantics

5.3 Arithmetic Soundness

Definition 5.3.1. Let Γ be We define an *arithmetical realization* to be a logical map $\phi : \mathbf{GL} \rightarrow \mathbf{PA}$ (that is a map on sentences which preserves logical connectives i.e. a morphism of the Boolean algebra of sentences) which satisfies the property that for any sentence A of \mathbf{GL} ,

$$\phi(\square A) = \mathcal{Bew}_\Gamma(\Gamma \phi(A) \neg)$$

Remark. When the realization ϕ is unambiguous we will often write A^* for $\phi(A)$ the realization of A .

Theorem 5.3.2 (Arithmetical Soundness). If $\mathbf{GL} \vdash A$ then $\mathbf{PA} \vdash A^*$ for any arithmetical realization $*$.

Proof. This proceeds by induction of proofs. It suffices to show that the axioms of \mathbf{GL} are realized by provably statements of \mathbf{PA} and that rules of inference in \mathbf{GL} are sound in \mathbf{PA} . Since \mathbf{PA} contains axiom schema for all propositional tautologies and the rule of inference MP we simply need to check the necessitation deduction,

$$\mathbf{PA} \vdash A^* \implies \mathbf{PA} \vdash (\square A)^*$$

and modal axioms,

$$\begin{aligned} \mathbf{PA} \vdash (\square(A \rightarrow B) \rightarrow (\square A \rightarrow \square B))^* \\ \mathbf{PA} \vdash (\square A \rightarrow \square \square A)^* \\ \mathbf{PA} \vdash (\square(\square A \rightarrow A) \rightarrow \square A)^* \end{aligned}$$

Using the properties of $*$ we see this is equivalent to asking that,

$$\begin{aligned} \mathbf{PA} \vdash A^* \implies \mathbf{PA} \vdash \mathcal{Bew}_\Gamma(\Gamma A^* \neg) \\ \mathbf{PA} \vdash \mathcal{Bew}_\Gamma(\Gamma A^* \rightarrow B^* \neg) \rightarrow (\mathcal{Bew}_\Gamma(\Gamma A^* \neg) \rightarrow \mathcal{Bew}_\Gamma(\Gamma B^* \neg)) \\ \mathbf{PA} \vdash \mathcal{Bew}_\Gamma(\Gamma A^* \neg) \rightarrow \mathcal{Bew}_\Gamma(\Gamma \mathcal{Bew}_\Gamma(\Gamma A^* \neg) \neg) \\ \mathbf{PA} \vdash \mathcal{Bew}_\Gamma(\Gamma \mathcal{Bew}_\Gamma(\Gamma A^* \neg) \rightarrow A^* \neg) \rightarrow \mathcal{Bew}_\Gamma(\Gamma A^* \neg) \end{aligned}$$

which are exactly the three Hilbert-Bernays derivability conditions and formalized Löb's theorem which have shown to be probable in \mathbf{PA} . \square

5.4 The Existence of Modal Fixed Points

5.5 Arithmetic Completeness

Theorem 5.5.1 (Arithmetical Completeness, Solovay, 1976). If $\mathbf{PA} \vdash A^*$ for any arithmetical realization $*$ then $\mathbf{GL} \vdash A$.

Remark. This theorem is remarkable because it captures the overarching logic of \mathbf{PA} in a modal logic based of propositional calculus without quantifiers. This result is made even more remarkable by the following decidability theorem for \mathbf{GL} .

Theorem 5.5.2. The theoremhood relation for \mathbf{GL} is decidable i.e. the decision problem for \mathbf{GL} is solvable.

Remark. We know by Church and Turring that the decision problem for **PA** is unsolvable that there does not exist an algorithm which can decide theoremhood in **PA**. Therefore, it is surprising and powerful that we can capture probability logic inside **PA** with the *decidable* theory **GL**.

Remark. We will end with a application of arithmetic completeness to generating undecidable arithmetical sentences. It is not difficult to show that **GL** $\not\vdash \Box p \vee \neg \Box p$. Then Solovay's proof allows us to construct a realization such that **PA** $\not\vdash (\Box p \vee \neg \Box p)^*$. However,

$$(\Box p \vee \neg \Box p)^* = \mathcal{Bew}_\Gamma(\Gamma p^{*\neg}) \vee \neg \mathcal{Bew}_\Gamma(\Gamma p^{*\neg})$$

so if **PA** $\vdash p^*$ or **PA** $\vdash \neg p^*$ then by HB1 we would have **PA** $\vdash \mathcal{Bew}_\Gamma(\Gamma p^{*\neg})$ or **PA** $\vdash \mathcal{Bew}_\Gamma(\Gamma \neg p^{*\neg})$ contradicting Solovay's construction. Thus p^* is an undecidable arithmetic sentence giving us further examples of oddities besides Gödel sentences.