

Diagonal-Free Proofs of the Diagonal Lemma

SAEED SALEHI

University of Tabriz & IPM

WORMSHOP 2017, Moscow

What is the Diagonal Lemma?

For any formula $\Psi(x)$ there exists some sentence η such that

$$\mathbb{N} \models \Psi(\bar{\eta}) \leftrightarrow \eta.$$

(The Semantic Diagonal Lemma)

This is usually provable in a Σ_1 -complete theory:

$$\mathbb{T} \vdash \Psi(\bar{\eta}) \leftrightarrow \eta.$$

(The Syntactic Diagonal Lemma)

What is the Diagonal Lemma good for?

(E.G.) For Proving the Following Theorems:

- ▶ Gödel's (1st & 2nd) Incompleteness Theorem(s);
- ▶ Gödel–Rosser's Incompleteness Theorem;
- ▶ Tarski's Undefinability (of Truth) Theorem;
- ▶ Löb's Theorem; $T \vdash \Box_T(\Box_T \mathcal{A} \rightarrow \mathcal{A}) \longrightarrow \Box_T \mathcal{A}$.

What is wrong with the Diagonal Lemma?

Does anybody remember its proof? What about the sketch?
Even after (so) many years of teaching the lemma?

- ▶ SAMUEL BUSS, *Handbook of Proof Theory* (Elsevier 1998, p. 119):
“[Its] proof [is] quite simple but rather tricky and difficult to conceptualize.”
- ▶ GYÖRGY SERÉNYI, *The Diagonal Lemma as the Formalized Grelling Paradox*, in: *Gödel Centenary 2006* (Eds.: M. Baaz & N. Preining), *Collegium Logicum* vol. 9, Kurt Gödel Society, Vienna, 2006, pp. 63–66. <https://arxiv.org/pdf/math/0606425.pdf>
<http://math.bme.hu/~sereny/poster.pdf>
- ▶ WAYNE URBAN WASSERMAN, *It Is “Pulling a Rabbit Out of the Hat”: Typical Diagonal Lemma “Proofs” Beg the Question*, (*Social Science Research Network*) SSRN (2008).
<http://dx.doi.org/10.2139/ssrn.1129038>

What is **really** wrong with the (proof of the) Diag. Lem.?

Vann McGee (2002) <http://web.mit.edu/24.242/www/1stincompleteness.pdf>

“The following result is a cornerstone of modern logic:

Self-referential Lemma. For any formula $\Psi(x)$, there is a sentence ϕ such that $\phi \leftrightarrow \Psi[\ulcorner \phi \urcorner]$ is a consequence of Q.

Proof: You would hope that such a deep theorem would have an insightful proof. No such luck. I am going to write down a sentence ϕ and verify that it works. What I won't do is give you a satisfactory explanation for why I write down the particular formula I do. I write down the formula because Gödel wrote down the formula, and Gödel wrote down the formula because, when he played the logic game he was able to see seven or eight moves ahead, whereas you and I are only able to see one or two moves ahead. I don't know anyone who thinks he has a fully satisfying understanding of why the Self-referential Lemma works. It has a rabbit-out-of-a-hat quality for everyone.”

The Problem of Eliminating the Diagonal Lemma!?

- ▶ HENRYK KOTLARSKI, *The Incompleteness Theorems After 70 Years*, *APAL* 126:1-3 (2004) 125–138.

The Diagonal Lemma, “being very intuitive in the natural language, is highly unintuitive in formal theories like Peano arithmetic. In fact, the usual proof of the diagonal lemma ... is short, but tricky and difficult to conceptualize. The problem was to eliminate this lemma from proofs of Gödel’s result. This was achieved only in the 1990s”.

Chaitin (1971) — Boolos (1989) — ...

Diagonal-Free Proofs ...

Some “Diagonal-Free” Proof of Tarski’s Theorem:

1. A. ROBINSON, On Languages Which Are Based On Nonstandard Arithmetic, *Nagoya Mathematical Journal* (1963).
2. H. KOTLARSKI, Other Proofs of Old Results, *MLQ* (1998).
3. G. SERÉNY, Boolos-Style Proofs of Limitative Theorems, *MLQ* (2004).
 - ▶ XAVIER CAICEDO, *Lecturas Matemáticas* (1993) (seminar 1987).
4. R. KOSSAK, Undefinability of Truth and Nonstandard Models, *APAL* (2004).

Toward a BIG SURPRISE

Tarski's Theorem (on the Undefinability of Truth) in \mathbb{N}

$$\neg \exists \Phi \forall \eta \mathbb{N} \models \Phi(\bar{\eta}) \leftrightarrow \eta$$

is equivalent with

$$\forall \Phi \exists \eta \mathbb{N} \models \neg [\Phi(\bar{\eta}) \leftrightarrow \eta]$$

or, by the propositional equivalence,

$$\neg(p \leftrightarrow q) \equiv (\neg p \leftrightarrow q)$$

with the [Semantic Diagonal Lemma](#)

$$\forall \Psi (= \neg \Phi) \exists \theta \mathbb{N} \models \Psi(\bar{\theta}) \leftrightarrow \theta.$$

A Big Surprise

So, any diagonal-free proof of Tarski's Undefinability Theorem

$$\neg \exists \Phi \forall \eta \mathbb{N} \models \Phi(\bar{\eta}) \leftrightarrow \eta$$

gives us a diagonal-free proof of the Semantic Diagonal Lemma

$$\forall \Psi \exists \theta \mathbb{N} \models \Psi(\bar{\theta}) \leftrightarrow \theta$$

by which one can prove (diagonal-freely)
the semantic version of Gödel's 1st Incompleteness Theorem

$$\forall T \exists \gamma (\mathbb{N} \models T \in \text{RE} \implies T \not\models \gamma, \neg \gamma).$$

More Surprises

H. KOTLARSKI (*APAL 2004, MLQ 1998*) proves (diagonal-freely) that

Let T be any theory in \mathcal{L}_{PA} containing PA. Assume that there exists a formula Φ such that for every sentence η ,
 $T \vdash \eta \equiv \Phi(\ulcorner \eta \urcorner)$. Then T is inconsistent.

That is to say that for any **consistent** $T \supseteq PA$,

$$\neg \exists \Phi \forall \eta T \vdash \Phi(\bar{\eta}) \leftrightarrow \eta$$

$$\forall \Phi \exists \eta T \not\vdash \Phi(\bar{\eta}) \leftrightarrow \eta$$

$\Psi = \neg \Phi : [T \not\vdash \Phi(\bar{\eta}) \leftrightarrow \eta] \iff T + [\Psi(\bar{\eta}) \leftrightarrow \eta]$ is consistent.

$\forall \Psi \exists \theta$ s.t. $T + [\Psi(\bar{\theta}) \leftrightarrow \theta]$ is **consistent**.

The Weak Diagonal Lemma

Any Diagonal-Free Proof of Tarski's Theorem for a theory T gives such a proof for the following **Weak Diagonal Lemma**.

For any consistent $T \supseteq PA$ and any formula $\Psi(x)$ there exists a sentence θ such that $T + [\Psi(\bar{\theta}) \leftrightarrow \theta]$ is consistent.

This is weak since cannot prove Gödel's 1st Incompleteness Theorem (by the way of Gödel's own proof):

Even though, for any consistent $T + [\neg \text{Pr}_T(\bar{\theta}) \leftrightarrow \theta]$ we have $T \not\vdash \theta$, we may not have $T \not\vdash \neg\theta$:

For $\theta = \perp$ we have the consistency of $[\neg \text{Pr}_T(\perp) \leftrightarrow \perp] \equiv \neg \text{Con}(T)$ with T (by Gödel's 2nd) but $T \vdash \neg\perp$ even when T is ω -consistent!

However, the Weak Diagonal Lemma **can prove Rosser's Theorem**:

The Weak Diagonal Lemma \implies Gödel–Rosser's Theorem

The following theory is consistent for some ρ :

$$T + [\forall x (\text{Proof}_T(x, \bar{\rho}) \rightarrow \exists y < x \text{Proof}_T(y, \overline{\neg\rho})) \leftrightarrow \rho].$$

Call it T' .

- ▶ If $T \vdash \rho$ then $\text{Proof}_T(k, \bar{\rho})$ for some $k \in \mathbb{N}$ and so $T' \vdash \exists y < \bar{k} \text{Proof}_T(y, \overline{\neg\rho})$ which contradicts $T' \vdash \neg \text{Proof}_T(\ell, \overline{\neg\rho})$ for all $\ell \in \mathbb{N}$ (by $T \not\vdash \neg\rho$).
- ▶ If $T \vdash \neg\rho$ then $\text{Proof}_T(k, \overline{\neg\rho})$ for some $k \in \mathbb{N}$. Also, $T' \vdash \exists a$ such that $\text{Proof}_T(a, \bar{\rho})$ and $\forall y < a \neg \text{Proof}_T(y, \overline{\neg\rho})$. Thus, $k < a$ is impossible, so $a \leq k$ whence $a \in \mathbb{N}$. This contradicts $T' \vdash \neg \text{Proof}_T(\ell, \bar{\rho})$ for all $\ell \in \mathbb{N}$ (by $T \not\vdash \rho$).

QED

The Weak Diagonal Lemma $\stackrel{?}{\implies}$ Löb's Theorem ?

Does the Weak Diagonal Lemma imply Löb's Theorem?

$$T \vdash \text{Pr}_T(\text{Pr}_T(\varphi) \rightarrow \varphi) \longrightarrow \text{Pr}_T(\varphi)$$

Only One Proof!?

Is There Any Diagonal-Free Proof For Löb's Theorem?

Is There Any Other Proof For Löb's Theorem?

Löb's Theorem \implies Gödel's 2nd Theorem

Löb's Theorem \iff Formalized Gödel's 2nd Theorem

$$\begin{aligned} & \text{Pr}_T(\text{Pr}_T(\varphi) \rightarrow \varphi) \longrightarrow \text{Pr}_T(\varphi) \\ & \neg \text{Pr}_T(\varphi) \longrightarrow \neg \text{Pr}_T(\neg \varphi \rightarrow \neg \text{Pr}_T(\varphi)) \\ & \text{Con}(T + \neg \varphi) \longrightarrow \neg \text{Pr}_{T+\neg \varphi}(\text{Con}(T + \neg \varphi)) \end{aligned}$$

for $\xi = \neg \varphi$

$$\text{Con}(T + \xi) \longrightarrow \neg \text{Pr}_{T+\xi}[\text{Con}(T + \xi)]$$

Thus Far ...

The Equivalences and The Implications:

Semantic Diagonal Lemma \iff Semantic Tarski's Theorem
 \implies Semantic Gödel's 1st Theorem

Weak Diagonal Lemma \iff Syntactic Tarski's Theorem
 \implies Gödel–Rosser's Theorem
 \implies 1st Incompleteness Theorem

Löb's Theorem \iff Formalized Gödel's 2nd Theorem
 \implies 2nd Incompleteness Theorem

Diagonal-Free Proofs for Gödel's 2nd Theorem

1. T. JECH, On Gödel's Second Incompleteness Theorem, *Proc. AMS* (1994).
2. H. KOTLARSKI, On the Incompleteness Theorems, *JSL* (1994).
3. M. KIKUCHI, A Note on Boolos' Proof of the Incompleteness Theorem, *MLQ* (1994).
4. M. KIKUCHI, Kolmogorov Complexity and the Second Incompleteness Theorem, *Arch. Math. Logic* (1997).
5. H. KOTLARSKI, Other Proofs of Old Results, *MLQ* (1998).
6. Z. ADAMOWICZ & T. BIGORAJSKA, Existentially Closed Structures and Gödel's Second Incompleteness Theorem, *JSL* (2001).
7. G. SERÉNY, Boolos-Style Proofs of Limitative Theorems, *MLQ* (2004).
8. H. KOTLARSKI, The Incompleteness Theorems After 70 Years, *APAL* (2004).

Diagonal-Free Proofs for Tarski's Theorem. I

A. Robinsion (1963):

If Φ defined truth [in \mathbb{N}](in $T \supseteq PA$) then let \mathfrak{M} be a non-standard model [$\equiv \mathbb{N}$](of T) with $\mathbb{N} < a \in \mathfrak{M}$. Put

$$\mathfrak{M}' = \{t_i(a) \mid t_i \text{ is an } \mathfrak{M}\text{-Skolem term, } i \in \mathbb{N}\} \quad (\cong \mathfrak{M}).$$

For any $n \in \mathbb{N}$ we have $\mathfrak{M}' \models \exists x \bigwedge_{i < n} x \neq t_i(a)$. So,

$$\mathfrak{M}' \models \exists x \forall y < n \Phi(\ulcorner x \neq t_y(a) \urcorner).$$

By overspill there is some $b > \mathbb{N}$ in \mathfrak{M}' such that

$$\mathfrak{M}' \models \exists x \forall y < b \Phi(\ulcorner x \neq t_y(a) \urcorner).$$

Thus for some $c \in \mathfrak{M}'$ we have $\bigwedge_{j \in \mathbb{N}} : \mathfrak{M}' \models c \neq t_j(a)$;

a contradiction!

QED

Diagonal-Free Proofs for Tarski's Theorem. II

H. Kotlarski (1998):

If Φ defined truth [in \mathbb{N}](in $T \supseteq PA$) then let

$$F(x) = \min y: \forall \bar{\varphi}, u \leq x [\exists v \Phi(\ulcorner \varphi(u, v) \urcorner) \rightarrow \exists v < y \Phi(\ulcorner \varphi(u, v) \urcorner)].$$

The unary function F is $[\mathbb{N}-](T-)$ DEFINABLE,

but **Dominates** all the unary definable functions:

If f is definable by $\varphi(u, v)$ then for any $z > \ulcorner \varphi \urcorner$ we have $F(z) > f(z)$;

a contradiction!

QED

Diagonal-Free Proofs for Tarski's Theorem. III

G. Serény (2004):

If Φ defined truth [in \mathbb{N}](in $T \supseteq PA$) then let

$$\text{Def}^{<z}(y) = \exists \varphi [\|\varphi\| < z \wedge \Phi(\ulcorner \forall \zeta [\varphi(\zeta) \leftrightarrow \zeta = y] \urcorner)].$$

where $\|\varphi\|$ is a measure of φ such that $\forall n \in \mathbb{N}$ there are finitely many ϕ with $\|\phi\| < n$

$$\text{Berry}^{<y}(u) = \neg \text{Def}^{<y}(u) \wedge \forall w < u \text{Def}^{<y}(w).$$

$$\text{Boolos}(x) = \text{Berry}^{<5\ell}(x), \text{ where } \ell = \|\text{Berry}^{<y}(x)\|.$$

$$\mathfrak{b} = \min z \neg \text{Def}^{<5\ell}(z).$$

Now we have $\|\text{Boolos}(x)\| < 5\ell$ and also

$$[\mathbb{N} \models](T \vdash) \quad \text{Boolos}(\zeta) \leftrightarrow \zeta = \mathfrak{b};$$

a contradiction!

QED

XAVIER CAICEDO, "La Paradoja de Berry, o la Indefinibilidad de la Definibilidad y las Limitaciones de los Formalismos", *Lecturas Matemáticas*, (1993) 14:37–48. Presented in "el Seminario de Lógica de la Universidad Nacional de Colombia" (1987). Revised in 2004 at <http://goo.gl/yYnstW>.

More Equivalences ...

Tarski's Semantic Theorem:

$$\{\ulcorner \varrho \urcorner \in \mathbb{N} \mid \mathbb{N} \models \varrho\} = \text{Th}(\mathbb{N}) \notin \text{Def}(\mathbb{N}) = \{X \subseteq \mathbb{N} \mid \exists \psi: n \in X \leftrightarrow \psi(n)\}$$

$$\forall T \subseteq \text{Th}(\mathbb{N}) [T \in \text{Def}(\mathbb{N}) \implies T \neq \text{Th}(\mathbb{N})]$$

$$\forall T [\mathbb{N} \models T \in \text{Def}(\mathbb{N}) \implies T \text{ is incomplete.}]$$

No Sound and Definable Theory is Complete!

Gödel–Smullyan Incompleteness Theorem

⋮

Thank You!

Thanks to

The Participants For Listening ...

and

The Organizers — For Taking Care of Everything ...