

The Modal Logic of Pure Provability

SAMUEL R. BUSS*

Abstract We introduce a propositional modal logic PP of “pure” provability in arbitrary theories (propositional or first-order) where the \Box operator means “*provable in all extensions*”. This modal logic has been considered in another guise by Kripke. An axiomatization and a decision procedure are given and the $\Box\Diamond$ subtheory is characterized.

1 Introduction This paper discusses a modal logic PP of *pure* provability; that is to say, of provability in arbitrary theories (propositional or first-order). The modal formula $\Box\phi$ is intended to mean “ ϕ is provable in all possible extensions of the present theory”; the subtleties arise in the interpretation of iterated modalities. The modal theory we are studying is very different from the provability interpretation of Solovay [6]; we allow theories which are much weaker than Peano Arithmetic and which may be unable to formalize metamathematics. Our theory PP was briefly mentioned by Kripke [3]; however, Kripke’s aims were somewhat different from ours and he did not explore PP in depth. The aim and motivation of this paper is to give a modal theory of provability; Kripke’s program (fulfilled by Solovay) was to give a provability interpretation of modal logic.

It turns out that the modal theory PP of pure provability is somewhat pathological; most notably, PP is not closed under the rule of substitution of arbitrary wff’s for propositional variables. On the other hand, PP does satisfy all consequences of S4 and McKinsey’s axiom (of S4.1). Furthermore, we provide

*Supported in part by NSF postdoctoral fellowship DMS-8511465 and NSF Grant DMS-8701828.

I wish to thank Michael Kremer for pointing out some errors in an earlier version of this paper and suggesting some corrections.

Received April 7, 1988; revised July 8, 1988 and September 12, 1988

an axiomatization for PP and prove the relevant completeness theorem. We also give a simple characterization of the $\Box\Diamond$ -fragment of PP.

A restricted form of the Grzegorzcyk axiom Grz is valid for PP although the usual full Grz axiom is not. The Grzegorzcyk axiom has already been studied in connection with modal logics with provability interpretations by Boolos [1] and Goldblatt [2]; however, we know of no direct connections of this with our work.

McDermott [4] and McDermott and Doyle [5] defined a modal logic intended for nonmonotonic logic with an operator M meaning ‘is consistent’. In spite of the linguistic similarity between our definition of \Diamond and their definition of M, the resulting mathematical systems appear to be radically different. It would be interesting to investigate whether our modal logic of pure provability can be utilized to describe nonmonotonic logics.

2 Semantics of pure provability logic We define the modal theory PP by defining the truth of a modal formula ϕ with respect to a theory T . The theory T is a propositional theory with, say, binary connectives \wedge , \vee , and \rightarrow , unary connective \neg , and propositional variables p, q, r, \dots . The modal formula ϕ may involve these connectives plus the modal necessitation operator \Box . The symbol \Diamond is an abbreviation for $\neg\Box\neg$.

Definition We define ϕ is true for T under the pure provability interpretation, denoted $T \vDash_{\text{PP}} \phi$, inductively on the number of occurrences of the modal operator \Box in ϕ :

- (1) If ϕ is \Box -free then $T \vDash_{\text{PP}} \phi$ if and only if $T \vdash \phi$ (i.e., ϕ is a consequence of T);
- (2) If ϕ is $\Box\psi$ then $T \vDash_{\text{PP}} \phi$ if and only if for every consistent extension S of T , $S \vDash_{\text{PP}} \psi$;
- (3) Otherwise, let $\Box\psi_1, \dots, \Box\psi_k$ be the maximal subformulas of ϕ which have outermost connective \Box . Let ϕ^* be obtained from ϕ by replacing each $\Box\psi_j$ by the tautology $p \vee \neg p$ if $T \vDash_{\text{PP}} \Box\psi_j$ and by $p \wedge \neg p$ if $T \not\vDash_{\text{PP}} \Box\psi_j$. Now ϕ^* is \Box -free and we define $T \vDash_{\text{PP}} \phi$ to hold if and only if $T \vDash_{\text{PP}} \phi^*$ (i.e., if and only if $T \vdash \phi^*$).

Definition The *pure provability theory* PP is the theory consisting of the modal formulas which are true for every consistent propositional theory under the pure provability interpretation. A formula is said to be *PP-valid* if it is in PP.

It is obvious from the definition that PP contains every tautology. PP is also closed under modus ponens; to prove this suppose that ϕ and $\phi \rightarrow \psi$ are PP-valid. Then for an arbitrary propositional theory T , define ϕ^* and ψ^* as in case (3) above. By assumption, $T \vdash \phi^*$ and $T \vdash \phi^* \rightarrow \psi^*$ so, by modus ponens, $T \vdash \psi^*$. Hence $T \vDash_{\text{PP}} \psi$. Since T was arbitrary, ψ is PP-valid.

However, PP is not closed under substitution. For example, it is easy to see that $\Box\Diamond p \rightarrow \Box p$ is in PP and yet $\Box\Diamond(q \rightarrow \Box q) \rightarrow \Box(q \rightarrow \Box q)$ is not. To prove the latter assertion note that $T \vDash_{\text{PP}} (q \rightarrow \Box q)$ if and only if either $T \vdash q$ or $T \vdash \neg q$; hence $\Box\Diamond(q \rightarrow \Box q)$ is PP-valid. On the other hand, $\Box(q \rightarrow \Box q)$ is not PP-valid.

Next we give a second characterization of PP by giving a Kripke modal structure M in which the valid modal formulas comprise PP. The modal structure M consists of all “worlds” (T, τ) where T is a consistent propositional theory and τ is a truth assignment to the propositional variables such that $\tau \vDash T$, or in words, τ is a model of T . The reachability relation R between worlds is defined thus: (T, τ) is reachable from (S, σ) if and only if T is an extension of S (denoted by $(S, \sigma)R(T, \tau)$). We write $(T, \tau) \vDash_M \phi$, ϕ a modal formula, to denote ϕ being true in the world (T, τ) in the modal structure M . The definition is by induction on the complexity of ϕ ; for ϕ atomic, $(T, \tau) \vDash_M \phi$ means that $\tau(\phi) = \text{True}$, and the usual definitions for the propositional and modal connectives define $(T, \tau) \vDash_M \phi$ for more complicated ϕ . In particular, $(T, \tau) \vDash_M \Box \phi$ holds if and only if $(S, \sigma) \vDash_M \phi$ for all consistent $S \supseteq T$ and all truth assignments σ satisfying S . We say that ϕ is M -valid if $(T, \tau) \vDash_M \phi$ for all worlds (T, τ) of M .

Theorem 1 *A modal formula is PP-valid if and only if it is M -valid.*

Lemma 1 *For any wff ϕ and any consistent T , ϕ is true for T if and only if $(T, \tau) \vDash_M \phi$ for all $\tau \vDash T$.*

Proof: Theorem 1 follows immediately from Lemma 1 so we need only prove the lemma. We proceed by induction on the complexity of ϕ .

Case (1): If ϕ is \Box -free then the claim is an immediate consequence of the completeness theorem for propositional logic.

Case (2): If ϕ is $\Box\psi$ then ϕ is true for T if and only if ψ is true for every $S \supseteq T$. By the induction hypothesis this is equivalent to $(S, \sigma) \vDash_M \psi$ for all $S \supseteq T$ and $\sigma \vDash S$. And that is equivalent to $(T, \tau) \vDash_M \phi$ for all $\tau \vDash T$.

Case (3): If ϕ is $\psi \wedge \chi$ then ϕ is true for T if and only if both ψ and χ are. By the induction hypothesis, the latter holds if and only if $(T, \tau) \vDash_M \psi$ and $(T, \tau) \vDash_M \chi$ for all $\tau \vDash T$. And this is equivalent to $(T, \tau) \vDash_M \phi$ for all $\tau \vDash T$.

Case (4): The case for disjunction and negation is slightly more complicated: we first consider the case where ϕ is of the form

$$\Box\psi_1 \vee \Box\psi_2 \vee \dots \vee \Box\psi_k \vee \neg\Box\chi_1 \vee \neg\Box\chi_2 \vee \dots \vee \neg\Box\chi_m \vee \gamma$$

where γ is \Box -free. By the induction hypothesis, we have that $\Box\psi_i$ being true for T is equivalent to $(T, \tau) \vDash_M \Box\psi_i$ for all $\tau \vDash T$, and similarly for $\neg\Box\chi_i$ and for γ . Thus if ϕ is true for T then one of $\Box\psi_i$, $\neg\Box\chi_j$, and γ is true for T and hence $(T, \tau) \vDash_M \phi$ for all $\tau \vDash T$. For the converse, suppose that $(T, \tau) \vDash_M \phi$ for all $\tau \vDash T$. Because of the definition of the reachability relation R , if there is *some* $\tau \vDash T$ such that $(T, \tau) \vDash_M \Box\psi_i$ or $(T, \tau) \vDash_M \neg\Box\chi_i$ then for *every* $\tau \vDash T$, $(T, \tau) \vDash_M \Box\psi_i$ or $(T, \tau) \vDash_M \neg\Box\chi_i$, respectively. Otherwise, $(T, \tau) \vDash_M \gamma$ for every $\tau \vDash T$. In any case the induction hypothesis applied to $\Box\psi_i$, $\neg\Box\chi_i$, or γ shows that ϕ is true for T .

Case (5): When ϕ is a Boolean combination of formulas we can put ϕ in conjunctive normal form, use Case (4) to show that each conjunct satisfies Theorem 1, and then use Case (3) to show that ϕ satisfies Theorem 1.

Theorem 1 helps to justify our definition of the pure provability logic PP: although some objections could be raised to our definition of PP, the fact that PP is the set of M -valid formulas shows that PP has a natural semantics and is a reasonable theory. Kripke [3] defined a structure similar to M except that his worlds are pairs (T, \mathcal{A}) , where T is a first-order theory extending Peano arithmetic and \mathcal{A} is a structure such that $\mathcal{A} \models T$. The modal theory of provability in Kripke's structure is easily seen to be equivalent to the modal theory of our structure M .

3 Axioms for pure provability logic We have already noted that PP contains all tautologies and is closed under modus ponens. It is also easy to see that PP includes all theorems of S4. This is because of Theorem 1 and the fact that the reachability relation R is reflexive and transitive. In addition, as Kripke [3] noted, PP includes the McKinsey S4.1 axiom $\diamond(\Box\phi \vee \Box\neg\phi)$ for all wff's ϕ and hence PP includes S4.1. To prove this, since every theory can be extended to a complete theory, it suffices to show that if T is a complete theory and ϕ is a wff then $T \vDash_{\text{PP}} \Box\phi$ or $T \vDash_{\text{PP}} \Box\neg\phi$. Since T is complete, there is a unique truth assignment τ such that $\tau \models T$ and hence in the modal structure described above the only world reachable from (T, τ) is (T, τ) itself. Thus if $(T, \tau) \vDash_M \phi$ then $(T, \tau) \vDash_M \Box\phi$ and, by Lemma 1, $T \vDash_{\text{PP}} \Box\phi$. Similarly, if $(T, \tau) \vDash_M \neg\phi$ then $T \vDash_{\text{PP}} \Box\neg\phi$. Thus either $T \vDash_{\text{PP}} \Box\phi$ or $T \vDash_{\text{PP}} \Box\neg\phi$.

We have established that PP contains the theory S4.1; however, PP is not closed under the rule of substitution of arbitrary wff's for propositional variables and hence cannot be equal to S4.1. Indeed, PP cannot be axiomatized by an axiom scheme closed under substitution. Nonetheless, we can axiomatize PP. We let Greek letters ϕ, ψ, χ, \dots range over arbitrary modal wff's, capital Roman letters A, B, C, \dots range over \Box -free wff's and lower case Roman letters p, q, r, \dots range over propositional variables.

Definition Let S4 + X be the modal theory made up of the S4 axioms, the S4 rules of inference of modus ponens and necessitation, and the additional axioms (denoted by X):

$$\Box(\Box A \rightarrow \Box B_1 \vee \dots \vee \Box B_k \vee C) \leftrightarrow \Box(A \rightarrow B_1) \vee \dots \vee \Box(A \rightarrow B_k) \vee \Box(A \rightarrow C)$$

where A, B_1, \dots, B_k, C are \Box -free formulas and $k \geq 0$.

A useful special case of the X axioms is when A is a tautology and the axioms reduce to

$$\Box(\Box B_1 \vee \dots \vee \Box B_k \vee C) \leftrightarrow \Box B_1 \vee \dots \vee \Box B_k \vee \Box C.$$

We now claim that the X axioms are PP-valid. Actually, the right-to-left implication of the X axioms is a consequence of S4; thus the X axioms could equivalently be stated as an implication instead of an equivalence. To show that an X axiom is PP-valid it will suffice to show that it is M -valid. We assume for simplicity that $k = 1$; now suppose that $\tau \models T$ and $(T, \tau) \vDash_M \Box(\Box A \rightarrow \Box B \vee C)$. By the definition of the reachability relation in the modal structure M and by Lemma 1 this is equivalent to $T \vDash_{\text{PP}} \Box(\Box A \rightarrow \Box B \vee C)$. This is equivalent to the condition that if $S \supseteq T$ and $S \vdash A$ then either $S \vdash B$ or $S \vdash C$. This is fur-

ther equivalent to $T \vdash A \rightarrow B$ or $T \vdash A \rightarrow C$; which is the same as $\Box(A \rightarrow B) \vee \Box(A \rightarrow C)$ being true for T . Again by Lemma 1, this is equivalent to $(T, \tau) \vDash_M \Box(A \rightarrow B) \vee \Box(A \rightarrow C)$ for all (or some) $\tau \vDash T$.

We have established that $S4 + X \subseteq PP$; we show below that equality holds. First we prove a lemma which is interesting in its own right – it states that the \Box -operator does not need to be iterated.

Lemma 2 *Every formula is equivalent in $S4 + X$ to a formula with no nested modal operators. Every formula of the form $\Box\phi$ is equivalent to a positive Boolean combination of formulas of the form $\Box A$ with A \Box -free.*

Proof: The second assertion of the lemma implies the first; we prove the second assertion by induction on the complexity of ϕ . The base case where ϕ is \Box -free is trivial since ϕ is already in the desired form. For the induction step we can, without loss of generality, assume that ϕ is in conjunctive normal form, $\phi = \phi_1 \wedge \dots \wedge \phi_n$. By $S4$, $\Box\phi$ is equivalent to $\Box\phi_1 \wedge \dots \wedge \Box\phi_n$. By the induction hypothesis and the fact that ϕ is in conjunctive normal form, each ϕ_i is equivalent to a formula of the form

$$\Box A_1 \wedge \dots \wedge \Box A_r \rightarrow \Box B_1 \vee \dots \vee \Box B_s \vee C$$

where each A_j, B_j, C is \Box -free. This is $S4$ -equivalent to

$$\Box A \rightarrow \Box B_1 \vee \dots \vee \Box B_s \vee C$$

where A is $A_1 \wedge \dots \wedge A_r$. So by axiom X , $\Box\phi_i$ is equivalent to $\Box(A \rightarrow B_1) \vee \dots \vee \Box(A \rightarrow B_s) \vee \Box(A \rightarrow C)$. Each $\Box(A \rightarrow B_i)$ occurs positively and the lemma is proved.

Theorem 2 $PP = S4 + X$.

Proof: Since $PP \supseteq S4 + X$ it suffices to show that every PP -valid formula ϕ is a consequence of $S4 + X$. By Lemma 2, $\Box\phi$ is $(S4 + X)$ -equivalent and hence PP -equivalent to a formula of the form

$$\psi = \bigwedge_i \bigvee_j \Box B_{i,j}$$

with each $B_{i,j}$ \Box -free. (Note that we are using the fact that $\Box\phi$ is equivalent to a *positive* Boolean combination of the formulas $\Box B_{i,j}$.) If any $\Box B_{i,j}$ is PP -valid then $B_{i,j}$ is a tautology and hence $S4 \vdash \Box B_{i,j}$. Thus, since $\Box\phi$ is PP -valid, ψ is a consequence of $S4$ and thus ϕ is a consequence of $S4 + X$.

Theorem 2 gives an axiomatization for PP . Note also that the proofs of Lemma 2 and Theorem 2 give a decision procedure for PP .

Our axiomatization for PP is not a finite scheme since the description of the X axioms included a variable k . The author does not know if there is a finite set of axiom schemes for PP .

4 The Grzegorzcyk axiom The Grzegorzcyk axiom Grz is $\Box(\Box(\phi \rightarrow \Box\phi) \rightarrow \phi) \rightarrow \phi$. To understand what this means, take the contrapositive and let ψ be $\neg\phi$. The Grzegorzcyk axiom becomes

$$\psi \rightarrow \diamond(\psi \wedge \square(\neg\psi \rightarrow \square\neg\psi)). \quad (1)$$

In essence this says: “If ψ is true in some world then there is some reachable world W where ψ is true such that in any world reachable from W if ψ is false then ψ is necessarily false.” One way to interpret this is as a discreteness property of the reachability relation on worlds, in that it says that there is a point where ψ last changes from being true to being false.

In the modal logic of pure provability the Grzegorzcyk axiom is not valid for arbitrary wff’s; indeed, we show below that Grz is not PP-valid when ϕ is $p \rightarrow \square p$. However, when ϕ is restricted to be \square -free then the Grzegorzcyk axioms *are* valid. This follows immediately from the X axioms as follows: We say that two wff’s ψ_1 and ψ_2 are PP-equivalent, written as $\psi_1 \equiv \psi_2$, if and only if $\psi_1 \leftrightarrow \psi_2$ is PP-valid. Now if A is \square -free then $\square(A \rightarrow \square A)$ is PP-equivalent to $\square(A \rightarrow A)$, which is PP-valid. Thus the Grzegorzcyk axiom $\square(\square(A \rightarrow \square A) \rightarrow A) \rightarrow A$ is PP-equivalent to $\square A \rightarrow A$. Since the latter formula is an axiom of S4, the Grzegorzcyk axiom with $\phi = A$ a \square -free formula is valid.

Next we show that the Grzegorzcyk axiom is not valid when ϕ is the wff $p \rightarrow \square p$. Indeed in this case the Grzegorzcyk axiom is equivalent to ϕ itself as the following chain of PP-equivalences shows:

$$\begin{aligned} & \square(\square((p \rightarrow \square p) \rightarrow \square(p \rightarrow \square p)) \rightarrow (p \rightarrow \square p)) \rightarrow (p \rightarrow \square p) \\ & \equiv \square(\square((p \rightarrow \square p) \rightarrow \square(\square p \vee \neg p)) \rightarrow (p \rightarrow \square p)) \rightarrow (p \rightarrow \square p) \\ & \equiv \square(\square((p \rightarrow \square p) \rightarrow \square p \vee \square \neg p) \rightarrow (p \rightarrow \square p)) \rightarrow (p \rightarrow \square p) \\ & \equiv \square(\square(p \vee \square p \vee \square \neg p) \rightarrow (p \rightarrow \square p)) \rightarrow (p \rightarrow \square p) \\ & \equiv \square((\square p \vee \square p \vee \square \neg p) \rightarrow (p \rightarrow \square p)) \rightarrow (p \rightarrow \square p) \\ & \equiv \square(\square \neg p \rightarrow \neg p \vee \square p) \rightarrow (p \rightarrow \square p) \\ & \equiv \square(\neg p \rightarrow \neg p) \vee \square(\neg p \rightarrow p) \rightarrow (p \rightarrow \square p) \\ & \equiv \square(\neg p \rightarrow p) \rightarrow (p \rightarrow \square p) \\ & \equiv \square p \rightarrow (p \rightarrow \square p) \\ & \equiv p \rightarrow \square p. \end{aligned}$$

The above equivalences are obtained by alternately using the X axioms and replacing a subformula by a tautologically equivalent subformula; the third from last equivalence follows because $\square(\neg p \rightarrow \neg p)$ is S4-valid. However, $p \rightarrow \square p$ is valid for a theory T if and only if either T proves p or T disproves p ; hence $p \rightarrow \square p$ and the Grzegorzcyk axiom for $p \rightarrow \square p$ are not PP-valid.

5 The $\square\diamond$ -fragment of PP We next show that the $\square\diamond$ -fragment of PP has a very simple characterization. Let $\phi^{-\square}$ be the formula obtained by removing all modal operators from ϕ while leaving the propositional connectives intact.

Theorem 3 $\square\diamond\phi$ is PP-valid if and only if $\phi^{-\square}$ is a tautology.

Proof: Note that if T is complete then $T \vDash_{\text{PP}} A \leftrightarrow \square A$ for all \square -free A . (Actually $T \vDash_{\text{PP}} \phi \leftrightarrow \square\phi$ for all ϕ .) Hence the modal operators in ϕ can be eliminated one at a time from ϕ to get that $T \vDash_{\text{PP}} \phi^{-\square} \leftrightarrow \phi$ for all complete T .

So if $\phi^{-\square}$ is a tautology ϕ is true for every complete theory, and thus $\square\diamond\phi$ is PP-valid since every consistent theory can be extended to a complete theory. Conversely, if $\square\diamond\phi$ is PP-valid then ϕ is true for every complete theory and hence every complete theory $T \vdash \phi^{-\square}$, i.e., $\phi^{-\square}$ is a tautology.

REFERENCES

- [1] Boolos, G., "On systems of modal logic with provability interpretations," *Theoria*, vol. 46 (1980), pp. 7–18.
- [2] Goldblatt, R., "Arithmetical necessity, provability and intuitionistic logic," *Theoria*, vol. 44 (1978), pp. 38–46.
- [3] Kripke, S., "Semantical considerations on modal logic," *Acta Philosophica Fennica*, vol. 16 (1963), pp. 83–94.
- [4] McDermott, D., "Nonmonotonic logic II: Nonmonotonic modal theories," *Journal of the Association for Computing Machinery*, vol. 29 (1982), pp. 33–57.
- [5] McDermott, D. and J. Doyle, "Nonmonotonic logic I," *Artificial Intelligence*, vol. 13 (1980), pp. 41–72.
- [6] Solovay, R., "Provability interpretations of modal logic," *Israel Journal of Mathematics*, vol. 25 (1976), pp. 287–304.

*Department of Mathematics
University of California, San Diego
La Jolla, California 92093*