

### 3.5 Herbrand's theorem (see e.g. U. Schöning: Logic for computer science)

Definition 3.5.1: Let  $\varphi$  be a sentence in  $L$ .

The "Herbrand universe  $D(\varphi)$ " of  $\varphi$  is the set of all closed terms which can be built up out of the symbols occurring in  $\varphi$  (plus some distinguished constant symbol  $c$  in case  $\varphi$  does not contain any constant symbol):

(i) every constant in  $\varphi$  is in  $D(\varphi)$  ( $c \in D\varphi$ ) if no constant in  $\varphi$ ).

(ii) if  $t_1, \dots, t_n \in D(\varphi)$  and the many function symbol  $f$  occurs in  $\varphi$ , then  $f(t_1, \dots, t_n) \in D(\varphi)$ .

Example:  $\varphi = \forall x \exists y P(x, f(x))$ ,  $P$  binary predicate symbol. Then  $D(\varphi) = \{c, f(c), f(f(c)), \dots\}$ .

Definition 3.5.2: Let  $\varphi$  be a sentence in  $L$ .

A structure  $A = \langle A, \dots \rangle$  is called a "Herbrand structure" for  $L$  if it

(i)  $A = D(\varphi)$

(ii) for all many function symbols  $f$  in  $\varphi$  and all terms  $t_1, \dots, t_n \in D(\varphi)$

$$f^A(t_1, \dots, t_n) = f(t_1, \dots, t_n)$$

For all constant symbols  $c$  in  $\varphi$ :  $c^A = c$ .

A Herbrand structure for  $\varphi$  which is a model of  $\varphi$  is called a "Herbrand model" of  $\varphi$ . Analogous for sets  $P$  of sentences.

- Remark: (i) Herbrand structures  $\mathcal{A}$  do not need to interpret closed terms by themselves:  $t^A = t$ . That is why Herbrand models are also called "term models".  
(ii) There is no special requirement to the interpretation of predicate symbols  $P$  in Herbrand structures.  
(iii) We consider Herbrand structures only in connection with logic without equality.

Proposition 3.5.3 (Existence of Herbrand models):

Let  $\varphi$  be a sentence (without  $=$ ) that is purely universal, i.e.  $\varphi = \forall \exists \varphi_{\text{qt}}$  ( $\exists 1$ ), where  $\varphi_{\text{qt}}$  is quantifier-free and  $\exists = x_1, \dots, x_n$ .

$\varphi$  has a model iff  $\varphi$  has a Herbrand model.

Proof: " $\Leftarrow$ " is trivial.

" $\Rightarrow$ ": Let  $\mathcal{A} = \langle A, \dots \rangle$  be an arbitrary model of  $\varphi$ . The interpretation of the function and constant symbols in the Herbrand model  $\mathcal{A}^H = \langle D(\varphi), \dots \rangle$  we are going to construct is fixed by the definition of Herbrand structures. Hence it suffices to interpret the predicate symbols  $P$  in  $\varphi$ :

$$(t_1, \dots, t_n) \in P^{\mathcal{A}^H} \Leftrightarrow (t_1^H, \dots, t_n^H) \in P^{\mathcal{A}},$$

where  $t_1, \dots, t_n \in D(\varphi)$ .

One verifies that  $\mathcal{A}^H \models \forall \exists \varphi_{\text{qt}} (\exists 1)$  by induction on the length  $k$  of  $\exists = x_1, \dots, x_k$ .

Proposition 3.5.3 implies the corresponding dual statement for purely existential sentences:

Proposition 3.5.4: Let  $\varphi = \exists x_1 \dots x_n \varphi_{xt} (\pm)$ , be a purely existential sentence (without  $=$ ). Then

$\models \varphi$  iff (for all Herbrand structures  $A^H$  for  $\varphi$ :  $A^H \models \varphi$ ).

Proof: Obvious using that the negation of  $\varphi$  is (logically equivalent) to a purely universal sentence so that proposition 3.5.3 applies.

Definition 3.5.5 (Herbrand expansion):

Let  $\varphi = \forall x_1 \dots x_n \varphi_{xt} (\pm)$  be a purely universal sentence. Then the "Herbrand expansion"  $E(\varphi)$  of  $\varphi$  is defined as  $E(\varphi) := \{\varphi_{xt}(t_1, \dots, t_n) : t_1, \dots, t_n \in D(\varphi)\}$ .

Proposition 3.5.6: Let  $\varphi$  be a purely universal sentence (without  $=$ ). Then

$\varphi$  has a model iff  $E(\varphi)$  is satisfiable in the sense of propositional logic.

Proof: By the previous results it suffices to show that  $\varphi$  has a Herbrand model iff  $E(\varphi)$  is satisfiable in the sense of propositional logic: let  $A^H$  be a Herbrand model of  $\varphi$ . Then  $\nu(R(t_1, \dots, t_n)) = \begin{cases} 1, & \text{if } A^H \models R(t) \\ 0, & \text{otherwise} \end{cases}$  is a satisfying assignment for  $E(\varphi)$ .

Conversely, if such an assignment  $v$  with  $v \models E(\varphi)$  is given, then (as many predicate symbols  $P$  in  $\varphi$ )  
 $P^{\neq^H} := \{(t_1, \dots, t_n) \in D(\varphi)^n : v(P(t_1, \dots, t_n)) = 1\}$   
defined a Herbrand model  $H^H$  of  $\varphi$ .  $\square$

Theorem 3.5.7 (Herbrand's Theorem, J. Herbrand 1930)

Let  $\varphi = \exists x \varphi_{\exists x} (\leq)$  be a purely existential sentence (without  $=$ ). Let  $\vdash_{\neg}$  denote deriving in ND without the  $=$ -rules.

Then the following holds:

$$\vdash_{\neg} \varphi \text{ iff } \exists \underbrace{t_{1,1}, \dots, t_{1,m}, \dots, t_{k,1}, \dots, t_{k,n}}_m \in D(\varphi) : \\ \bigvee_{i=1}^m \varphi_{\exists x} (t_{1,i}, \dots, t_{k,i}) \in \text{TAUT}.$$

Proof:

" $\Leftarrow$ " If  $\bigvee_{i=1}^m \varphi_{\exists x} (t_{1,i}, \dots, t_{k,i}) \in \text{TAUT}$ , then  
 $\vdash_{\neg} \bigvee_{i=1}^m \varphi_{\exists x} (t_{1,i}, \dots, t_{k,i})$  by the completeness theorem for propositional logic.

$\vdash_{\neg} \varphi_{\exists x} (t_{1,i}, \dots, t_{k,i}) \rightarrow \exists x_1, \dots, x_k \varphi_{\exists x} (x_1, \dots, x_k)$   
by  $\exists$ -introduction. Hence

$\vdash_{\neg} \bigvee_{i=1}^m \varphi_{\exists x} (t_{1,i}, \dots, t_{k,i}) \rightarrow \varphi \vee \underbrace{\dots \vee \varphi}_{m-\text{times}}$  and  
so by contraction  $\vdash_{\neg} \varphi$ .

" $\Rightarrow$ " We give an ineffective model-theoretic proof for the contrapositive formulation: if  $\bigvee_{i=1}^m q_{4i}(t_i) \notin \text{TAUT}$  for all  $t_1, \dots, t_m \in D(\varphi)$ , then

$\{\neg q_{4i}(t_1, \dots, t_n) : t_1, \dots, t_n \in D_f(\varphi)\}$  is satisfiable (in the sense of propositional logic) for every finite subset  $D_f(\varphi)$  of  $D(\varphi)$ . By propositional compactness this implies that  $E(\forall x_1 \dots x_n \neg q_{4i})$  is satisfiable (in the sense of prop. logic). Hence by proposition 3.56  $\varphi$  has a model  $\mathcal{A}$ . So  $\mathcal{A} \not\models \varphi$ . The completeness theorem now yields that  $\mathcal{A} \models \varphi$ .  $\square$

Remark: Herbrand's very complicated proof was purely proof-theoretic (with some corrections due to Gold 1946, resp. Andreka, Nandris, Dabben 1966) and provides an algorithm for the extraction of a tautological Herbrand disjunction from a given proof of  $\varphi$  in  $ND_-=$ . Subsequently, other syntactic proofs were given by the Hilbert  $\varepsilon$ -substitution method

(D. Hilbert / P. Bernays: Grundlagen der Mathematik I, Springer 1938) and by Gentzen as a consequence of his cut-elimination theorem (G. Gentzen 1936).

See e.g. "J. Shoenfield: Mathematical Logic" (1967) or "S. Buss (ed.): Handbook of Proof Theory" (1998) for syntactic proofs.

Using the Herbrand normal form  $\varphi^H$  of a prenex sentence  $\varphi$ , Herbrand's Theorem immediately extends to arbitrary sentences in the following form:

Herbrand's Theorem (general form) 3.5.8:

Let  $\varphi$  be a sentence in prenex normal form without equality and  $\varphi^H = \exists \bar{x} \varphi_{\bar{x}}^H (\pm)$  its Herbrand normal form. Then

$$\vdash_{\sim} \varphi \text{ iff } \exists \bar{t}_1, \dots, \bar{t}_n \in D(\varphi^H) \quad \bigvee_{i=1}^n \varphi_{\bar{x}}^H(\bar{t}_i) \in \text{TACT.}$$

(note that now the Herbrand terms  $\bar{t}_i$  are built up also by using the Herbrand index functions  $f$  used to form  $\varphi^H$  from  $\varphi$ ).

Proof: " $\Rightarrow$ "  $\vdash_{\sim} \varphi$  implies trivially that  $\vdash_{\sim} \varphi^H$ .

Now apply theorem 3.5.7 to the purely existential sentence  $\varphi^H$ .

" $\Leftarrow$ " By " $\Leftarrow$ " in thm. 3.5.7  $\bigvee_{i=1}^n \varphi_{\bar{x}}^H(\bar{t}_i) \in \text{TACT}$  implies that  $\vdash \varphi^H$  also  $\vdash \varphi$  (3.46). To get from the tautological H-disjunction  $\bigvee_{i=1}^n$  a proof of  $\varphi$  in  $\vdash_{\sim}$  one needs a more involved argument which we only sketch in the example below:

Example: Clearly:  $\vdash_{\sim} \exists x \forall y (P(x) \vee P(y))$ .

$$(\exists x \forall y (P(x) \vee P(y)))^H = \exists x (P(x) \vee P(f(x))).$$

$$(P(c) \vee P(f(c))) \vee (P(f(c)) \vee P(f(f(c)))) \in \text{TACT}.$$

So  $t_1 := c$ ,  $t_2 := f(c)$  provide a valid Herbrand disjunction.

There is also a version of the general form of Herbrand's theorem that is formulated without the use of Herbrand index functions:

Let  $\varphi^H$  be the Herbrand normal form of a prenex sentence  $\varphi$  (without equality) and

$$\varphi^D = \bigvee_{i=1}^m \varphi^H(\underline{t}_i) \in \text{Taut} \text{ a Herbrand disjunction.}$$

Now replace in  $\underline{t}_i$  all terms starting with a Herbrand function symbol (e.g.  $f(s_1, \dots, s_n)$ ) by a new variable starting from terms with largest size. Let the resulting  $\underline{f}$ -free disjunction be denoted by  $\varphi^D$ . With  $\varphi^{H,D}$  also  $\varphi^D$  is a tautology and there is a direct proof from  $\varphi^D$  to  $\varphi$ , i.e. a proof that essentially only uses quantifier-introduction and contraction. The general procedure is somewhat complicated to describe and so we just treat an example:

Consider again  $\varphi := \exists x \forall y (P(x) \vee P(y))$  and

$$\varphi^H = \exists x (P(x) \vee \forall y P(y)).$$

$$\varphi^{H,D} = (P(c) \vee \neg P(f(c))) \vee (P(f(c)) \vee \neg P(f(f(c)))) \in \text{Taut}.$$

Now replace  $f(f(c))$  by the variable  $z$  and  $f(c)$  by  $y$ .

Then

$$\varphi^D := (P(c) \vee \neg P(y)) \vee (P(y) \vee \neg P(z)) \in \text{Taut}.$$

91

In ND<sub>→</sub> one easily shows that the following reasoning can be carried out

$$\begin{array}{c}
 \varphi^D \\
 \hline
 \frac{}{\exists x \forall y (P(x) \vee P(y)) \vee \exists y \forall z (P(y) \vee \neg P(z))} \text{ " } \forall\text{-Intro" } \\
 \hline
 \frac{\exists x \forall y (P(x) \vee \neg P(y)) \vee \exists y \forall z (P(y) \vee \neg P(z))}{\exists x \forall y (P(x) \vee \neg P(y)) \vee \exists y \forall z (P(y) \vee \neg P(z))} \text{ " } \exists\text{-Intro" } \\
 \hline
 \frac{\exists x \forall y (P(x) \vee \neg P(y)) \vee \exists y \forall z (P(y) \vee \neg P(z))}{\exists x \forall y (P(x) \vee \neg P(y))} \text{ " } \exists\text{-Intro" } \\
 \end{array}$$

"Extracting modals  
from d-variables"

Herbrand's theorem for sentences with equality 3.59.

Let  $\varphi$  be a prenex sentence and  $\varphi^H = \exists \underline{x} \varphi_{\text{at}}^H(\underline{x})$  its Herbrand normal form. Let  $\forall \underline{u} E_{\text{at}}(\underline{u})$  be the prenex normal form of the conjunction of the axioms  $I_1 - I_q$  for all the function and predicate symbols occurring in  $\varphi$ . Then  $\vdash \varphi$  implies that  $\vdash \forall \underline{u} E_{\text{at}}(\underline{u}) \rightarrow \varphi$  and also  $\vdash \forall \underline{u} E_{\text{at}}(\underline{u}) \rightarrow \exists \underline{x} \varphi_{\text{at}}^H(\underline{x})$  and so

$$\vdash \underbrace{\exists \underline{x}, \underline{u} (E_{\text{at}}(\underline{u}) \rightarrow \varphi_{\text{at}}^H(\underline{x}))}_{\varphi :=}.$$

Now apply theorem 3.5.7 to  $\varphi$  to obtain closed terms  $s_1, \dots, s_n, t_1, \dots, t_m \in D(\varphi^H)$  s.t.

$$\left( \bigwedge_{i=1}^k E_{\text{st}}(S_i) \rightarrow \bigvee_{i=1}^m q_{\text{st}}^{(i)}(t_i) \right) \in \text{TAUT}.$$

So we no longer get a tautological Herbrand disjunction for  $\exists z q_{\text{st}}^{(z)}$  but a disjunction that is a tautological consequence of finitely many closed instances of =-axioms, i.e. a disjunction that is a so-called quasi-tautology.

The converse direction " $(\Lambda \rightarrow V) \xrightarrow{\lambda} q$ " follows as before and even holds if  $\bigwedge_{i=1}^k E_{\text{st}}(S_i)$  contains instances of equality axioms.

$$r_1 = r'_1 \wedge \dots \wedge r_l = r'_l \rightarrow f(r_1, \dots, r_l) = f(r'_1, \dots, r'_l)$$

be function symbols used in form  $q^{(i)}$  above.

In that latter case, though, the aforementioned process of replacing f.-terms by new variables no longer results in a valid quasi-tautology.

Herbrand's theorem for open theories: Let us recall that an open theory T is a theory axiomatized by purely universal sentences. Let  $\alpha, \alpha'$  as above and  $T \models \alpha \rightarrow \beta$  in the case of the purely universal =-axioms, one can shift all the purely universal axioms of T used in proving  $\alpha$  as an implicative assumption. Applying Herbrand's theorem yields a disjunction that is

a tautological consequence of finitely many closed instances of these axioms (and equality axioms). In particular:

Theorem 3.5.10: Let  $T$  be an open theory,  $\varphi$  a sentence in  $L(T)$  in prenex normal form and  $\varphi^H$  its Herbrand normal form. Then (for the extension  $T[f]$  of  $T$  by the  $f$ -functions)

$$T \vdash \varphi \text{ iff } \exists \underline{t}_1, \dots, \underline{t}_n : T[\varphi] \vdash \bigvee_{i=1}^m \varphi_{\text{at}}^H(\underline{t}_i),$$

where the closed terms  $\underline{t}_i$  are built up out of  $\varphi^H$ -material and the constant and function symbols occurring in the non-logical axioms of  $T$ .

Remark: Every theory  $T$  can be extended to its Skolem extension  $T^{\text{Sk}}$ , that is an open theory to which theorem 3.5.10 extends. Then, however, the Herbrand terms in general also involve the Skolem function symbols used in forming  $T^{\text{Sk}}$ !

Herbrand's Theorem does not apply directly to non-open theories such as PA. Consider PA augmented by a new unary function symbol  $f$ :

$$\text{PA}[f] \vdash \exists x \forall y (f(x) \leq f(y)) =: \varphi$$

$\varphi^H = \exists x (f(x) \leq f(g(x)))$ . One only has a

"H-disjunction" of variable length  $f(0)$

$$\checkmark^{(6)} f(t_i) \leq f(g(t_i)), \text{ where } t_i := g^{(i)}(0).$$

But:  $\overset{i=0}{\text{no H-disjunction of fixed length }} n!$

Corollary to the theorem 3.5.9: Let  $\varphi$  be a sentence without " $=$ ". If  $\vdash \varphi$  then  $\vdash_{=} \varphi$ .

Proof: W.l.o.g. we may assume that  $\varphi$  is in prenex normal form. By Thm. 3.5.9  $\vdash \varphi$  ad so  $\vdash \varphi^k$  yield a quasi-tautology

$$\bigvee_{i=1}^n \varphi_{qf}^k(t_i) \text{ whl is a tautology since " $=$ "}$$

does not occur ad so we may ~~and~~ interpret any formula " $s = t$ " in  $\bigwedge_{i=1}^k E_\varphi(s_i)$  as "true" if  $s$  and  $t$  are identical as terms and "false" otherwise.

" $\in$ " from Thm. 3.5.8 then yields  $\vdash_{=} \varphi$ .  $\square$

A lower bound for the length of Herbrand disjunctions (as well as the complexity of normalizing proofs or elimination of cuts):

A growth example (Statman, Orevkov, Yang, Paddd 79-91).

Consider the following open first order theory  $\mathcal{T}$ :

(i) Language  $\mathcal{L}(\mathcal{T})$  of  $\mathcal{T}$ :  $=, +, \cdot, 2^{(\cdot)}, I(x), 0, 1$

(ii) Non-logical axioms:

$$x + (y + z) = (x + y) + z, y + 0 = y,$$

$$2^0 = 1, 2^x + 2^y = 2^{x+y}$$

$$I(0), I(x) \rightarrow I(1+x).$$

Intended meaning of " $I(x)$ ": " $x$  is a (readable) natural number".

Let  $\vdash_{\mathcal{Q}_0}(x)$  be the prenex part of the conclusion of the universal closure of the non-logical axioms.

Prop.:  $\mathcal{T} \vdash I(2^{\varphi_0(x)})$  for any fixed  $k$ , i.e.

$$\vdash \vdash_{\mathcal{Q}_0}(x) \rightarrow I(2^{\varphi_0(x)}), \text{ i.e.}$$

$$\vdash \exists x (\varphi_0(x) \rightarrow I(2^{\varphi_0(x)}))$$

Proof: We define inductivity relations  $R_i$ :

$R_0 := I$ . Let  $R_i$  be defined by I.H. then

$$R_{i+1}(x) := \forall y (R_i(y) \rightarrow R_i(2^x + y)).$$

By meta-induction on  $i$  we prove that

$$\mathcal{J} \vdash R_i(0) \wedge \forall x (R_i(x) \rightarrow R_i(1+x)).$$

For  $i=0$  this is clear. Assume

$$(+) \quad R_i(0) \wedge (R_i(x) \rightarrow R_i(1+x)).$$

Using  $2^0 = 1$  we get

$$R_{i+1}(0) \equiv \forall y (R_i(y) \rightarrow R_i(2^0 + y)) \leftarrow \\ \forall y (R_i(y) \rightarrow R_i(1+y)).$$

Also

$$R_{i+1}(x) \equiv \forall y (R_i(y) \rightarrow R_i(2^x + y)) \\ \rightarrow \forall y ((R_i(y) \rightarrow R_i(2^x + y)) \wedge \\ (R_i(2^x + y) \rightarrow R_i(\underbrace{2^x + (2^x + y)}_{= (2^x + 2^x) + y}))) \\ = (2^x + 2^x) + y = 2^{x+x} + y$$

$$\rightarrow \forall y (R_i(y) \rightarrow R_i(2^{x+x} + y)) \equiv R_{i+1}(2x).$$

Now  $R_{i+1}(x) \xrightarrow{y=0} (R_i(0) \rightarrow R_i(\underbrace{2^x + 0}_{= 2^x}))$ , i.e.

$$(**) \quad R_{i+1}(x) \rightarrow R_i(2^x). \text{ To}$$

$$R_k(0) \xrightarrow{(**)} R_{k-1}(?)$$

$$\vdots$$

$$\xrightarrow{(**)} R_0(2^{\cdot\frac{?}{?}}) = I(2^{\cdot\frac{?}{?}}).$$

Corollary: There exist a constant  $c \in \mathbb{N}$  s.t.

for any  $k \in \mathbb{N}$  there is a <sup>Hebrew</sup> ND-proof  
 $d_k \vdash_{ND} \exists x (A_0(x) \rightarrow I(2^{\cdot\frac{?}{?}}))$  s.t. the depth  
of  $d_k$  is  $\leq c \cdot k$ .

Theorem (Thm 92, Pudlák 92): Any valid  
H-deduction for " $\exists x (A_0(x) \rightarrow I(2^{\cdot\frac{?}{?}}))$ "  
has length  $\geq 2^{\cdot\frac{?}{?}}$ . The same applies to the  
length of a normal or cut-free proof of " $\exists x (A_0(x) \rightarrow I(2^{\cdot\frac{?}{?}}))$ ".

Prove (of the first claim): We prove that the  
Hebrew d-prf has to "contain" all instances  
 $I(0) \rightarrow I(1), I(1) \rightarrow I(\bar{n}), \dots, I(\bar{n}) \rightarrow I(\bar{n}+1)$   
for all  $n < 2^{\cdot\frac{?}{?}}$ . Suppose it's not, i.e.  
 $\vdash \bigvee_t (A_0(x) \rightarrow I(2^{\cdot\frac{?}{?}}))$  s.t.

but for some  $\bar{n} < 2^{\cdot\frac{?}{?}}$  we don't have  $I(\bar{n}) \rightarrow I(\bar{n}+1)$ .  
We interpret  $I(x)$  as " $x \leq \bar{n}$ ". □