

Symbolická regrese: Jak získat kvalitní a současně kompaktní řešení?

Bc. Ondřej Končal

Abstrakt

Geometrické sémantické genetické programování (GSGP) dosahuje kvalitních výsledků při popisu dat složitými matematickými modely. Cenou za přesný popis je ale výsledná velikost řešení. Tento článek se zabývá převodem řešení získaného GSGP na instanci kartézského genetického programování (CGP) a jeho následnou optimalizací. CGP dokáže dobře redukovat velikost již vzniklých řešení a kombinace těchto dvou metod má tak potenciál vytvořit kvalitní a zároveň malý model popisující vstupní data. Úspěšným principem redukce je podstromové kartézské genetické programování (SCGP), které je představeno v tomto článku. Využívá možnosti rozdělení řešení na podstromy a následně je upravuje. Na všech testovaných úlohách z oblasti symbolické regrese se podařilo dosáhnout znatelného zmenšení řešení a pouze u jedné úlohy ze 4 došlo k přetrénování. Kombinace GSGP a SCGP tak má potenciál vytvořit přesnější model než genetické programování nebo CGP za řádově stejnou dobu. Velikost řešení je znatelně menší než samotný výstup GSGP.

Klíčová slova: geometrické sémantické genetické programování — kartézské genetické programování — symbolická regrese — evoluční algoritmus

Příložené materiály: [Naměřená data a vyvinutá řešení](#)

*xkonca00@stud.fit.vutbr.cz, *Fakulta informačních technologií, Vysoké učení technické v Brně*

1. Úvod

Evoluční algoritmy jsou dnes již zaběhlým a účinným nástrojem nejen pro optimalizaci, ale některé z nich jsou schopny automaticky vytvářet spustitelné struktury. Takovým algoritmem je genetické programování [1], které dokáže vytvářet a optimalizovat jak počítačové programy, tak i rovnice a obvody. Standardní genetické operátory v genetickém programování pracují na úrovni syntaxe a genotypu, ale cílem je úprava fenotypu. Cílenou změnu fenotypu a s tím i sémantiky přináší geometrické sémantické operátory a jejich použití v geometrickém sémantickém genetickém programování (GSGP) [2]. Cenou za kvalitní řešení a téměř žádné přetrénování je velikost výsledného řešení – počet použitých uzlů grafu. Každá aplikace geometrického sémantického operátoru jedince zvětší, ale díky speciálnímu přístupu je evaluace rychlá. Nicméně rekonstrukce výsledného řešení je časově a prostorově náročná, stejně tak i použití tohoto řešení. S možností

zmenšení již existujících řešení přichází kartézské genetické programování (CGP) [3], které v tomto směru již prokázalo svoji funkčnost.

Hlavním cílem tohoto článku je návrh algoritmu, který efektivně zmenší velikost řešení získaného pomocí GSGP. Je zde představen jak algoritmus převodu řešení získaného pomocí GSGP na instanci CGP, tak i speciální přístup využití CGP nazvaný podstromové CGP (SCGP), které pracuje vždy pouze se sémanticky spjatou částí chromozomu. Experimenty jsou provedeny na reálných farmakokinetických datech.

Článek je uspořádán následovně: Sekce 2 se zabývá stručným přehledem symbolické regrese a GSGP, sekce 3 přehledem CGP. V sekci 4 je představen algoritmus převodu výstupu GSGP na reprezentaci CGP. V sekci 5 je představeno SCGP. Výsledky experimentů a data s nastavením k nim použité jsou popsány v sekci 6. Nakonec jsou shrnuty výsledky experimentů a navrženy možnosti budoucí práce v oblasti optimalizace výstupu geometrického sémantického genetického

21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40

41 programování pomocí kartézského genetického pro-
42 gramování.

43 2. Geometrické sémantické genetické programování

Jednou z klasických úloh, ve které se genetické pro-
gramování uplatňuje, je symbolická regrese [1]. Jedná
se o nalezení výrazu

$E : \mathbb{R}^n \rightarrow \mathbb{R}$ co nejlépe popisujícího vstupní data tvo-
řená dvojicemi (v_i, y) , kde $v_i \in \mathbb{R}^n$ a $y \in \mathbb{R}$. Často je
využíván výpočet hodnoty fitness (kvality řešení) f
pomocí průměrné odchylky:

$$f = \frac{\sum_{j=1}^N |y_j^{GP} - y_j|}{N} \quad (1)$$

kde y_j^{GP} značí výstupy kandidátního programu, které
jsou porovnávány s výstupy trénovací množiny y_j o ve-
likosti N . Jedná se o úlohu minimalizace hodnoty
fitness. Problém symbolické regrese lze řešit i pomocí
GSGP.

GSGP je forma genetického programování, která
upravuje jedince přímo na úrovni sémantiky za pomoci
geometrických sémantických operátorů. Těmi jsou ge-
ometrická sémantická mutace (GSM) a geometrické
sémantické křížení (GSC):

$$\text{GSC: } T_{XO} = (T_1 \cdot T_R) + ((1 - T_R) \cdot T_2) \quad (2)$$

$$\text{GSM: } T_M = T + ms \cdot (T_{R1} - T_{R2}) \quad (3)$$

kde $T, T_1, T_2 : \mathbb{R}^n \rightarrow \mathbb{R}$ jsou rodiče, $T_R, T_{R1}, T_{R2} : \mathbb{R}^n \rightarrow \mathbb{R}$
jsou náhodné reálné funkce s oborem hodnot $\langle 0, 1 \rangle$
a ms je reálná hodnota označovaná jako *mutation step*,
která udává velikost mutace. Pro upravení náhodně
vygenerovaných funkcí tak, aby vracely hodnoty v roz-
sahu $\langle 0, 1 \rangle$, je použita speciální logistická funkce sig-
moida:

$$T_R = (1 + e^{-T_{rand}})^{-1} \quad (4)$$

44 kde T_{rand} je funkce vracející neomezenou hodnotu.
45 Tyto operace tak stále zvětšují jedince, na které jsou
46 aplikovány. Jako poslední operace je používána repli-
47 kace, která přenese jedince z minulé generace nezmě-
48 něného.

49 Implementace GSGP v C++ představená v [2]
50 využívá namísto generování stále nových náhodných
51 stromů a následného fyzického spojení stromů do jed-
52 noho metodu záznamů s odkazy. Na začátku je vy-
53 generována počáteční populace a k ní multimnožina
54 náhodně vygenerovaných stromů. Všichni jedinci jsou
55 pak tvořeni pouze odkazy na tyto stromy a záznamy
56 pak vypadají (*crossover*, $\&T_1$, $\&T_2$, $\&T_R$) nebo
57 (*mutation*, $\&T$, $\&T_{R1}$, $\&T_{R2}$, ms). Ohodnocení jedinců

Lze vytvořit vždy z ohodnocení podstromů, z kterých
se skládají, což značně urychlí evaluaci. Pro získání
výsledku je potřeba provést zpětnou rekonstrukci, která
je časově a prostorově náročná. Výsledky experi-
mentů ukázaly, že GSGP si dokáže velmi dobře poradit
s přetřénováním a řešení tak zůstává stále dostatečně
obecné (viz [4]).

Pro účely této práce je využita výše zmíněná imple-
mentace GSGP od Maura Castelli dostupná na strán-
kách <http://gsgp.sourceforge.net/>.

3. Kartézské genetické programování

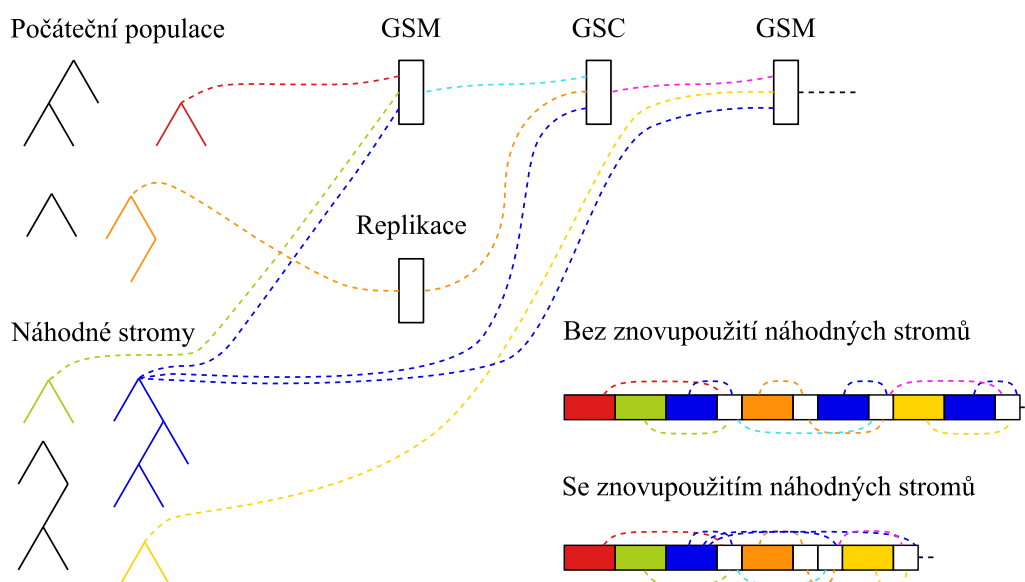
V CGP jsou reprezentovány spustitelné struktury po-
mocí acyklického grafu ve formě dvojrozměrné mřížky
dané velikosti [3]. Přirozeně tak dokáže pracovat s ob-
vody a matematickými funkcemi a zároveň zabraňuje i
vzniku *bloatu*, tzn. zvětšování stromu bez změny hod-
noty fitness, který se u standardního genetického pro-
gramování vyskytuje. Jediným genetickým operátorem
používaným v CGP je mutace, která má jak explo-
rativní, tak explozivní funkci. Je tomu tak díky
možnosti použít vstupy a vnitřní uzly i vícekrát (nebo
vůbec), což vede na několik druhů redundance. To
umožňuje znovupoužití částí grafu a možnost zakó-
dování většího fenotypu. Uzly, které jsou použity
pro tvorbu fenotypu (jsou propojeny s výstupem), se
nazývají aktivní. Evoluce používá strategií $(1 + \lambda)$,
kde počet potomků $\lambda \in \langle 4, 9 \rangle$.

Využívá se jak k návrhu nových, tak k optimali-
zaci již vzniklých řešení, které může zmenšovat nebo
zrychlovat.

4. Převod výstupu GSGP na instanci CGP

Prvním krokem je převod výstupu GSGP do CGP, aby
bylo možno jej dále zpracovat. Tato část je využitelná
i mimo SCGP a proto je zde představena samostatně.
Jak zde bylo již řečeno, řešení je uloženo jakožto
počáteční populace stromů a multimnožina stromů
náhodně vygenerovaných na začátku běhu algoritmu.
Jednotlivé generace obsahují pouze záznamy o tom,
které záznamy z minulé generace, případně počáteční
a náhodné stromy jsou spojeny. Zpětným průchodem
se všechny použité stromy a záznamy označí. Náhodné
stromy mohou být použity i vícekrát, proto je potřeba
tuto skutečnost zakomponovat do návrhu a jsou zde
představeny dva přístupy – vložení reprezentace náho-
dného stromu do chromozomu znova, nebo odkázání
se na již uloženou reprezentaci.

Algoritmus postupně prochází přes všechny gene-
race a každý záznam, který je označen jako použitý
pro tvorbu výsledku, se přidá do chromozomu. Pokud



Obrázek 1. Ilustrace převodu jedince v GSGP na jedince v CGP dvěma způsoby

107 se v záznamu objeví odkaz na náhodný strom nebo
 108 jedince z počáteční populace, je jeho první výskyt
 109 vždy přiveden na postfixovou notaci a následně na
 110 odpovídající uzly CGP. V případě opakovaného výskytu
 111 stromu je možno vytvořit strom znovu, nebo se připojit
 112 na uzel CGP se stromem již vytvořeným (obr. 1)
 113 (v druhém případě je potřeba ukládat indexy uzlů, které
 114 obsahují vrcholy stromů a vždy testovat, zda byl už
 115 strom použit). Při použití GSM, GSC a replikace je
 116 proveden odkaz vždy na uzly vytvořené v minulé ge-
 117 neraci, případně na náhodné stromy (uzly), se kterými
 118 jsou spojeny (u GSM a GSC dodány další uzly propo-
 119 jující tři podstromy).

120

5. Podstromové kartézské genetické programování

121 Hlavní částí tohoto článku je návrh algoritmu využí-
 122 vajícího CGP ke zmenšení velikosti řešení získaného
 123 pomocí GSGP. Tento přístup je zde nazván podstro-
 124 mové CGP proto, že pracuje odděleně s podstromy,
 125 které byly postupně přidávány v jednotlivých ge-
 126 neracích GSGP. Řešení získané z GSGP musí splňovat
 127 následující podmínky:

- 128 1. Nepoužívá se geometrické sémantické křížení
 129 (pouze GSM a replikace)
- 130 2. Převod do CGP je realizován bez znovupoužití
 131 náhodných stromů

132 Hlavní myšlenka je v rozdělení chromozomu na menší
 133 a sémanticky spjaté chromozomy obsahujících jed-
 134 notlivé GSM a původní strom. Díky asociativitě a
 135 komutativitě sčítání lze upravovat a vyhodnocovat jed-
 136 notlivé chromozomy nezávisle, rychleji a lokálněji.
 137 Nad každým je tak možné spustit samostatný běh CGP.
 138 Nejprve je výsledek GSGP nepoužívajícího GSC

převeden do chromozomu CGP bez znovupoužití ná- 139
 hodných stromů. Následně je chromozom rozdělen 140
 na jednotlivé GSM chromozomy tak, jak byly v gene- 141
 racích přidávány. Pro každý chromozom je vytvořen 142
 vektor odezev na vstupní vektory, což způsobí značné 143
 zrychlení výpočtu hodnoty fitness v pozdější fázi algo- 144
 ritmu. Následuje cyklus, ve kterém se postupně spouští 145
 CGP nad každým chromozomem. Pro výpočet hodnoty 146
 fitness daného chromozomu je jakožto požadovaný 147
 výstup použit výstup z trénovacích dat snížený o odezvy 148
 všech ostatních chromozomů. Dále probíhá klasické 149
 CGP s parametry (*generations, population_size,* 150
mutation_rate, selection_method). V posledním kroku 151
 jsou odstraněny uzly, které nejsou aktivní a je přepočítán 152
 vektor odezev na vstupy pro evolvovaného jedince. Po 153
 evoluci chromozomů následuje sjednocení chromo- 154
 zomů po dvou uzlem s operací součtu. Chromozomy 155
 jsou vybírány dle zadané strategie 156
subtree_select_method. V poslední fázi dojde opět 157
 k výpočtu odezev všech chromozomů na vstupní vektory. 158
 Cyklus evoluce chromozomů za pomoci CGP 159
 a následně spojování probíhá tak dlouho, dokud nezbu- 160
 de jediný chromozom, který prošel evolucí. 161

6. Experimenty

162

6.1 Nastavení

163

Parametry použité pro nastavení GSGP při těchto ex- 164
 perimentech jsou následující: 165

- Velikost populace: 2000 166
- Počet generací: 1000 (případně 300 pro *P3D*) 167
- Počet náhodných stromů: 500 168
- Metoda tvorby stromů: Ramped Half-and-Half 169
- Maximální hloubka stromu: 8 170
- Pravděpodobnost křížení: 0 171

172 Pravděpodobnost mutace: 0,9
 173 Velikost turnaje: 4
 174 Bylo provedeno 40 nezávislých běhů pro vybraná data,
 175 kdy řešení s nejmenší hodnotou fitness bylo vybráno
 176 jakožto vstup SCGP. Použité parametry SCGP jsou:
 177 Velikost populace: 8
 178 Počet generací: 50
 179 Pravděpodobnost mutace: 0,08
 180 Metoda selekce: nejkratší jedinec s fitness lepší
 181 nebo stejně dobrou jako byla počáteční
 182 Metoda selekce podstromů: náhodný výběr
 183 Pro každý vybraný výstup GSGP bylo provedeno 20
 184 nezávislých běhů.

185 6.2 Trénovací a testovací data

186 Testování proběhlo na 4 úlohách z oblasti farmako-
 187 kinetiky. Jedná se o data popisující lidskou orální
 188 biodostupnost (*bioav*), úroveň vazby plazmatických
 189 bílkovin (*PPB*), medián smrtící dávky – toxicitu (*Tox*)
 190 a 3D strukturu proteinu (*P3D*). První 3 datasety jsou
 191 důkladně rozebrány v [5]. Rozměry datasetů jsou
 192 popsány v tabulce 1 (funkce mají jeden výstup – po-
 193 slední sloupec). Každý dataset je 30-ti různými způsoby
 194 rozdělen na trénovací a testovací část v poměru 70:30 a
 195 pro experimenty bylo náhodně vybráno jedno rozdělení.
 196 Všechna rozdělení jsou nakonec použita pro otestování
 přetrénování.

	vstupní parametry (sloupce-1)	záznamy (řádky)
bioav	241	359
PPB	628	131
Tox	626	234
P3D	9	45730

197 **Tabulka 1.** Parametry použitých datasetů

197

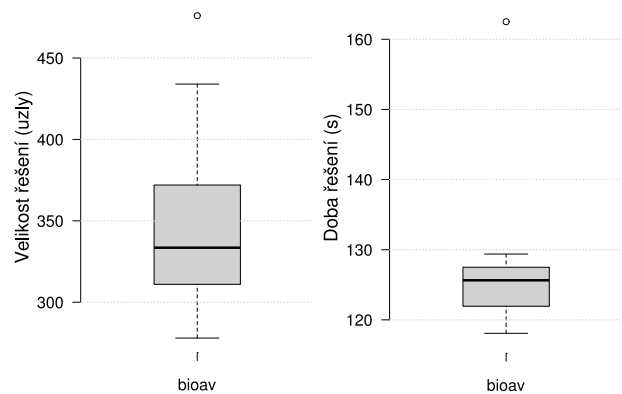
198 6.3 Výsledky

199 Základní údaje o výsledcích budou prezentovány vždy
 200 v trojici (počet aktivních uzlů; trénovací fitness; testov-
 201 vací fitness). Grafy jsou vytvořeny ze všech běhů
 202 SCGP. Nejlepší jedinec je vybrán nejen jako nejmenší,
 203 ale zároveň musí splňovat i testovací fitness, která není
 204 kontrolována programem. Čím menší (lepší) řešení
 205 bude, tím rychleji bude zpravidla evoluce ukončena.

206 Použití klasického CGP jak přímo na vstupní data
 207 nebo na výstup GSGP se ukázalo jako neefektivní.
 208 Výsledky byly vždy horší než u samotného GSGP,
 209 případně SCGP.

210 Nejlepší výstup GSGP pro *bioav* měl hodnoty
 211 (22285; 19,455; 26,938). Výsledek 20-ti běhů SCGP
 212 je zobrazen na obr. 2. Podmínku pro testovací fit-
 213 ness splňují 3 výsledky: (330; 19,412; 26,465), (332;

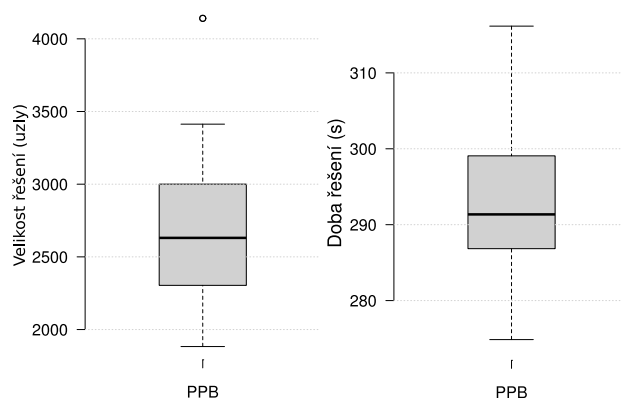
19,378; 25,839), (339; 19,443; 25,691). To značí 214
 zmenšení o 98,5% za necelých 130 sekund ve všech 215
 případech. Všechna tři řešení se ukázala stejně obecná 216
 jako výsledek získaný pomocí GSGP.



Obrázek 2. *bioav* – původní velikost 22285 uzlů

217

GSGP pro *PPB* data našlo řešení (22907; 4,312; 218
 30,768). Výsledky SCGP jsou vizualizovány na obr. 219
 3. Nalezeno bylo jediné ucházející řešení o hodnotách 220
 (2722; 4,293; 30,718). Zmenšení tak činí 88% za 295 221
 sekund. Řešení si nedokázalo poradit s přetrénováním 222
 ani v jednom případě (kromě původního trénovacího a 223
 testovacího).

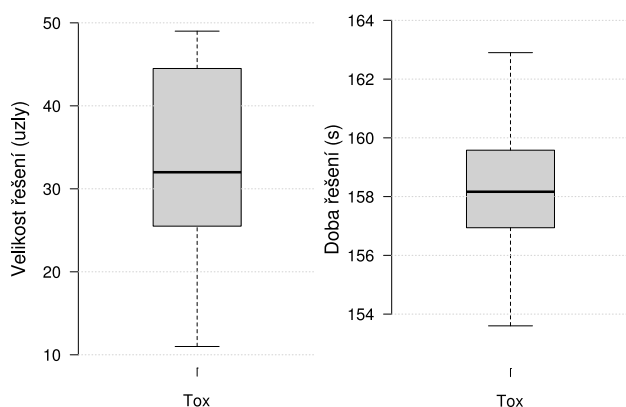


Obrázek 3. *PPB* – původní velikost 22907 uzlů

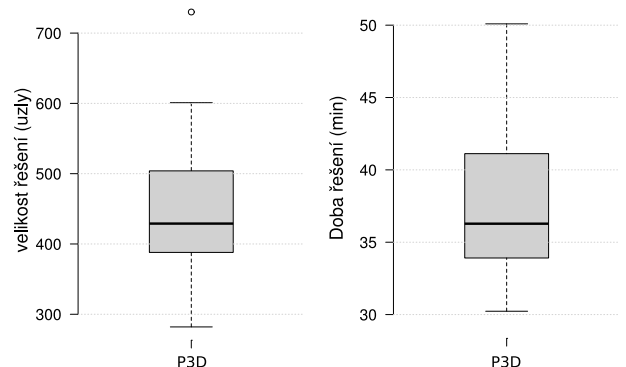
224

Výstup GSGP popisující dataset *Tox* má hodnoty 225
 (21568; 1336,956; 1495,621). Na obr. 4 je vidět 226
 naměřené hodnoty z běhů SCGP. 10 výsledků je při- 227
 jatelných s velikostí 11 až 47. Zmenšení je více jak 228
 99,5% za zhruba 160 sekund. Každé řešení si zhruba 229
 zachovalo stejnou schopnost zobecňovat jako referenční 230
 řešení poskytnuté GSGP. 231

Vstup SCGP pro úlohu *P3D* nabývá hodnot (5764; 232
 3,981; 3,999) a výstup je graficky zobrazen na obr. 5. 233
 Kvůli vysokému počtu záznamů trvá evoluce GSGP 234
 pouze 300 generací. Z 20-ti běhů je 15 dostatečně 235
 kvalitních s rozsahem velikostí 282 až 730, což dává 236
 zmenšení 87,3% až 95,1% do 3000 sekund. Tři běhy 237
 nejsou v grafu zaneseny, neboť skončily výjimkou. 238
 Všechna řešení jsou podobně obecná jako výstup GSGP. 239



Obrázek 4. Tox – původní velikost 21568 uzlů



Obrázek 5. P3D – původní velikost 5764 uzlů

vstřícnost při konzultacích a cenné rady, které mi pomohly tuto práci vyhotovit. 263
264

Literatura 265

- [1] John R. Koza. *Genetic Programming: On the Programming of Computers by Means of Natural Selection*. Cambridge : Bradford Book, London : MIT Press, 1992. 266
267
268
269
- [2] Mauro Castelli, Sara Silva, and Leonardo Vanneschi. A c++ framework for geometric semantic genetic programming. *Genetic Programming and Evolvable Machines*, 16(1):73–81, Mar 2015. 270
271
272
273
- [3] Lukáš Sekanina, Zdeněk Vašíček, Richard Růžička, Michal Bidlo, Jiří Jaroš, and Petr Švenda. *Evoluční hardware: od automatického generování patentovatelných invencí k sebedifikujícím se strojům*. Praha : Academia, 2009. ISBN 978-80-200-1729-1. 274
275
276
277
278
279
- [4] Mauro Castelli, Leonardo Trujillo, Leonardo Vanneschi, Sara Silva, Emigdio Z-Flores, and Pierrick Legrand. Geometric semantic genetic programming with local search. In *Proceedings of the 2015 Annual Conference on Genetic and Evolutionary Computation, GECCO '15*, pages 999–1006, New York, NY, USA, 2015. ACM. 280
281
282
283
284
285
286
- [5] Francesco Archetti, Stefano Lanzeni, Enza Messina, and Leonardo Vanneschi. Genetic programming for computational pharmacokinetics in drug discovery and development. *Genetic Programming and Evolvable Machines*, 8(4):413–432, Dec 2007. 287
288
289
290
291
292

240

241 7. Závěr

242 V této práci byly představeny dvě metody – Převod
243 řešení GSGP na instanci CGP a podstromové CGP.
244 GSGP pracuje se sémantickými operátory a dokáže
245 dobře zobecňovat, nicméně vytvořené řešení je velké.
246 Úlohou SCGP je tento výstup zmenšit. Experimenty
247 ukázaly, že ve třech případech (datasety *bioav*, *Tox*,
248 *P3D*) je dosažitelné zmenšení větší jak 95% a přetré-
249 nování nenastává. Pro dataset (*PPB*) zmenšení činilo
250 88%, ale přetrénování se ukázalo jako neúnosné.

251 Další práce bude probíhat na poli experimentování
252 s nastavováním parametrů evoluce, metodou výběru
253 podstromů pro sjednocení a spouštění klasického CGP
254 na již zmenšeném výsledku. Jako možnost dalšího
255 výzkumu se nabízí navrhnout období SCGP pracující
256 i s geometrickým sémantickým křížením a také vy-
257 zkoušet SCGP na nově představenou formu mutace pro
258 GSGP – geometrickou sémantickou mutaci s lokálním
259 prohledáváním.

260 Poděkování

261 Tímto bych rád poděkoval svému vedoucímu prof.
262 Ing. Lukáši Sekaninovi, Ph.D. za odborné vedení,