

# Program Synthesis

Беляев Станислав

Руководитель: Алексей Шпильман

# Prophet

<http://groups.csail.mit.edu/pac/patchgen/papers/prophet-popl16.pdf>

Идея:

- ищем потенциально плохие места в программе
- изменяем, конечным числом преобразований
- ранжируем патчи, обучая вероятностную модель

Фундаментальная проблема:

- корректного патча с большой вероятностью может не оказаться в генерируемом пространстве патчей
  - преобразований мало
  - сложно добавить новые преобразования (нужно писать код)

## Insert Control Flow

$S \Rightarrow \text{if } (E) \text{ return } K; S$

$S \Rightarrow \text{if } (E) \text{ goto } L; S$

$S \Rightarrow \text{if } (E) \text{ break; } S$

## Insert Guard

$S \Rightarrow \text{if } (E) \{ S \}$

## Change Condition

$\text{if } (C) \{ \dots \} \text{ else } \{ \dots \} \Rightarrow \text{if } (C \ \&\& \ E) \{ \dots \} \text{ else } \{ \dots \}$

$\text{if } (C) \{ \dots \} \text{ else } \{ \dots \} \Rightarrow \text{if } (C \ || \ E) \{ \dots \} \text{ else } \{ \dots \}$

## Replace

$S \Rightarrow S[e1/e2]$

$S \Rightarrow S[c1/c2]$

$S \Rightarrow S[f(e1, \dots, en)/g(e1, \dots, en)]$

## Copy

$S \Rightarrow Q[e1/e2]; S$

$S \Rightarrow Q[c1/c2]; S$

$S \Rightarrow Q[f(e1, \dots, en)/g(e1, \dots, en)]; S$

## Initialize

$S \Rightarrow \text{memset}(\&e, 0, \text{sizeof}(e)); S$

# Bayesian encoder & decoder

<https://arxiv.org/pdf/1503.00075.pdf>

- задача - определить распределение  $P(X|\Theta)$ 
  - $X$  - программа или ее часть
  - $\Theta$  - evidence, например, последовательность методов
- делим задачу на две -  $P(\Psi|\Theta)$  (encoder) и  $P(X|\Psi)$  (decoder)
  - $\Psi$  - латентный вектор, в котором заключен “смысл программы”
- encoder собирает смысл по программе
  - идея: перейти от программе как последовательности токенов к программе как дереву ast
- decoder итеративно генерирует новую программу
  - как-бы разворачивает собранный смысл в новую программу
  - начинаем с пустого ast, генерируем ребенка, а дальше его братьев
- проблемы
  - что делать с литералами при генерации?
  - можно ли абстрагироваться от языка?

$$P(X|\Theta) = \int P(X|\Psi)P(\Psi|\Theta)d\Psi$$

# Encoder

<https://arxiv.org/pdf/1503.00075.pdf>

- программа = слово из КС грамматики = дерево ast
  - парсинг реализован как форк русparser
- векторизация word2vec со связью родитель-ребенок
  - реализована на tensorflow
- child sum tree LSTM
  - проходимся от листьев к корню, пересчитывая вектора внутренних вершин
  - вектор корня и есть вектор, в котором собран весь “смысл” дерева

$$\tilde{h}_j = \sum_{k \in C(j)} h_k,$$

$$i_j = \sigma \left( W^{(i)} x_j + U^{(i)} \tilde{h}_j + b^{(i)} \right),$$

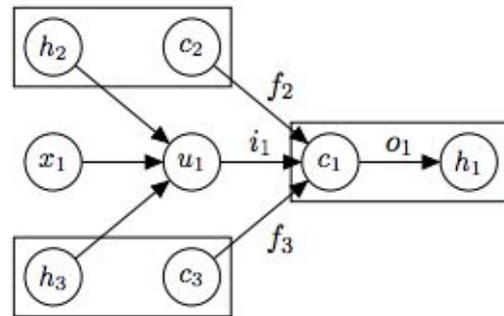
$$f_{jk} = \sigma \left( W^{(f)} x_j + U^{(f)} h_k + b^{(f)} \right),$$

$$o_j = \sigma \left( W^{(o)} x_j + U^{(o)} \tilde{h}_j + b^{(o)} \right),$$

$$u_j = \tanh \left( W^{(u)} x_j + U^{(u)} \tilde{h}_j + b^{(u)} \right)$$

$$c_j = i_j \odot u_j + \sum_{k \in C(j)} f_{jk} \odot c_k,$$

$$h_j = o_j \odot \tanh(c_j),$$



# Semantic relatedness

- хорошо работает для естественных языков и для синтаксических деревьев
- есть интуиция, что также будет хорошо работать и для ast
- К степеней похожести, по возрастанию
- KL - мера удаленности распределений
- обучающая выборка - посылки задачек с codeforces.ru

$$h_x = h_L \odot h_R,$$

$$h_+ = |h_L - h_R|,$$

$$h_s = \sigma \left( W^{(x)} h_x + W^{(+)} h_+ + b^{(h)} \right),$$

$$\hat{p}_\theta = \text{softmax} \left( W^{(p)} h_s + b^{(p)} \right),$$

$$\hat{y} = r^T \hat{p}_\theta,$$

$$r^T = [1 \ 2 \ \dots \ K]$$

$$J(\theta) = \frac{1}{m} \sum_{k=1}^m \text{KL} \left( p^{(k)} \parallel \hat{p}_\theta^{(k)} \right) + \frac{\lambda}{2} \|\theta\|_2^2.$$

accuracy  $\approx$  0.82

# Результаты

- изучил предметную область
- реализована encoder-часть
  - парсинг
  - векторизация
  - tree LSTM
- протестирована на задаче semantic-relatedness для программ
  - 0.82
- ИСПОЛЬЗОВАЛ
  - python (theano, tensorflow, русparser, pickle)

Спасибо за внимание

<https://github.com/StasBel/program-synthesis>