

Lab 4: Time Series Prediction with GP¹

Released: March 9, 2018

Deadline: March 23, 2018

Weight: 10 % for 06-27819, or 11.25 % for 06-27818 %

You need to implement one program that solves Exercises 1-3 using any programming language. In Exercise 5, you will run a set of experiments and describe the result using plots and a short discussion.

(In the following, replace `abc123` with your username.) You need to submit one zip file with the name `niso4-abc123.zip`. The zip file should contain one directory named `niso4-abc123` containing the following files:

- the source code for your program
- a Dockerfile (instructions will be provided later)
- a PDF file for Exercises 4 and 5



In this last lab, we will do a simple form of time series prediction. We assume that we are given some historical data, (e.g. bitcoin prices for each day over a year), and need to predict the next value in the time series (e.g., tomorrow's bitcoin value).

We formulate the problem as a regression problem. The training data consists of a set of m input vectors $\mathcal{X} = (x^{(0)}, \dots, x^{(m-1)})$ representing historical data, and a set of m output values

¹Revised: Sunday 11th March, 2018 at 18:41.

$\mathcal{Y} = (x^{(0)}, \dots, x^{(m-1)})$, where for each $0 \leq j \leq m-1$, $x^{(j)} \in \mathbb{R}^n$ and $y^{(j)} \in \mathbb{R}$. We will use genetic programming to evolve a prediction model $f : \mathbb{R}^n \rightarrow \mathbb{R}$, such that $f(x^{(j)}) \approx y^{(j)}$.

Candidate solutions, i.e. programs, will be represented as *expressions*, where each expression *evaluates* to a value, which is considered the output of the program. When evaluating an expression, we assume that we are given a current input vector $x = (x_0, \dots, x_{n-1}) \in \mathbb{R}^n$. Expressions and evaluations are defined recursively. Any floating number is an expression which evaluates to the value of the number. If e_1, e_2, e_3 , and e_4 are expressions which evaluate to v_1, v_2, v_3 and v_4 respectively, then the following are also expressions

- `(add e_1 e_2)` is addition which evaluates to $v_1 + v_2$, e.g. `(add 1 2)` \equiv 3
- `(sub e_1 e_2)` is subtraction which evaluates to $v_1 - v_2$, e.g. `(sub 2 1)` \equiv 1
- `(mul e_1 e_2)` is multiplication which evaluates to $v_1 v_2$, e.g. `(mul 2 1)` \equiv 2
- `(div e_1 e_2)` is division which evaluates to v_1/v_2 if $v_2 \neq 0$ and 0 otherwise, e.g., `(div 4 2)` \equiv 2, and `(div 4 0)` \equiv 0,
- `(pow e_1 e_2)` is power which evaluates to $v_1^{v_2}$, e.g., `(pow 2 3)` \equiv 8
- `(sqrt e_1)` is the square root which evaluates to $\sqrt{v_1}$, e.g. `(sqrt 4)` \equiv 2
- `(log e_1)` is the logarithm base 2 which evaluates to $\log(v_1)$, e.g. `(log 8)` \equiv 3
- `(exp e_1)` is the exponential function which evaluates to e^{v_1} , e.g. `(exp 2)` \equiv $e^2 \approx 7.39$
- `(max e_1 e_2)` is the maximum which evaluates to $\max(v_1, v_2)$, e.g., `(max 1 2)` \equiv 2
- `(ifleq e_1 e_2 e_3 e_4)` is a branching statement which evaluates to v_3 if $v_1 \leq v_2$, otherwise the expression evaluates to v_4 e.g. `(ifleq 1 2 3 4)` \equiv 3 and `(ifleq 2 1 3 4)` \equiv 4
- `(data e_1)` is the j -th element x_j of the input, where $j \equiv \lfloor v_1 \rfloor \bmod n$.
- `(diff e_1 e_2)` is the difference $x_k - x_\ell$ where $k \equiv \lfloor v_1 \rfloor \bmod n$ and $\ell \equiv \lfloor v_2 \rfloor \bmod n$
- `(avg e_1 e_2)` is the average $\frac{1}{|k-\ell|} \sum_{t=\min(k,\ell)}^{\max(k,\ell)-1} x_t$ where $k \equiv \lfloor v_1 \rfloor \bmod n$ and $\ell \equiv \lfloor v_2 \rfloor \bmod n$

We can build large expressions from the recursive definitions. For example, the expression

$$\text{code}(\text{add}(\text{mul } 2 \ 3) \ (\text{log } 4))$$

evaluates to

$$2 \cdot 3 + \log(4) = 6 + 2 = 8.$$

To evaluate the fitness of an expression e on a training data $(\mathcal{X}, \mathcal{Y})$ of size m , we use the mean square error

$$f(e) = \frac{1}{m} \sum_{j=0}^{m-1} \left(y^{(j)} - e(x^{(j)}) \right)^2,$$

where $e(x^{(j)})$ is the value of the expression e when evaluated on the input vector $x^{(j)}$.

Exercise 1. (30 % of the marks)

Implement a routine to evaluate expressions. You can assume that the input describes a syntactically correct expression. Hint: Use a library for parsing s-expressions.

Input arguments:

- `-expr` an expression
- `-n` the dimension of the input vector n
- `-x` the input vector

Output:

- the value of the expression

Example:

```
[pk1@phi ocamlc]$ app_niso_lab4 -question 1 -n 1 -x "1.0" \
    -expr "(mul (add 1 2) (log 8))"
9.0
[pk1@phi ocamlc]$ app_niso_lab4 -question 1 -n 2 -x "1.0 2.0" \
    -expr "(max (data 0) (data 1))"
2.0
```

Exercise 2. (10 % of the marks) Implement a routine which computes the fitness of an expression given a training data set.

Input arguments:

- `-expr` an expression
- `-n` the dimension of the input vector
- `-m` the size of the training data $(\mathcal{X}, \mathcal{Y})$
- `-data` the name of a file containing the training data in the form of m lines, where each line contains $n + 1$ values separated by tab characters. The first n elements in a line represents an input vector x , and the last element in a line represents the output value y .

Output:

- The fitness of the expression, given the data.

Exercise 3. (30 % of the marks)

Design a genetic programming algorithm to do time series forecasting. You can use any genetic operators and selection mechanism you find suitable.

Input arguments:

- `-lambda` population size
- `-n` the dimension of the input vector
- `-m` the size of the training data $(\mathcal{X}, \mathcal{Y})$
- `-data` the name of a file containing training data in the form of m lines, where each line contains $n + 1$ values separated by tab characters. The first n elements in a line represents an input vector x , and the last element in a line represents the output value y .
- `-time.budget` the number of seconds to run the algorithm

Output:

- The fittest expression found within the time budget.

Exercise 4. (10 % of the marks)

Describe your algorithm from Exercise 3 in the form of pseudo-code (see Lab 1 for an example). The pseudo-code should be sufficiently detailed to allow an exact re-implementation.

Exercise 5. (20 % of the marks)

In this final task, you should try to determine parameter settings for your algorithm which lead to as fit expressions as possible.

Your algorithm is likely to have several parameters, such as the population size, mutation rates, selection mechanism, and other mechanisms components, such as diversity mechanisms.

Choose parameters which you think are essential for the behaviour of your algorithm. Run a set of experiments to determine the impact of these parameters on the solution quality. For each parameter setting, run 100 repetitions, and plot box plots of the fittest solution found within the time budget.