

# SCOT Approximation, Training and Asymptotic Inference

**Mikhail Malyutov**

*Mathematics Department*

*Northeastern University*

*360 Huntington Avenue, Boston, MA 02115, USA*

M.MALIOUTOV@NEU.EDU

**Paul Grosu**

*College of Computer and Information Science*

*Northeastern University*

*360 Huntington Avenue, Boston, MA 02115, USA*

PGROSU@GMAIL.COM

**Editor:** Alex Gammerman, Vladimir Vovk, Zhiyuan Luo, and Harris Papadopoulos

## Abstract

Approximation of stationary strongly mixing processes by SCOT models and the Le Cam-Hajek-Ibragimov-Khasminsky locally minimax theory of statistical inference for them is outlined. SCOT is an  $m$ -Markov model with sparse memory structure. In our previous papers we proved SCOT equivalence to 1-MC with state space—alphabet consisting of the SCOT contexts. For the fixed alphabet size and growing sample size, the Local Asymptotic Normality is proved and applied for establishing asymptotically optimal inference. We outline what obstacles arise for a large SCOT alphabet size and not necessarily vast sample size. Training SCOT on a large string using clusters of computers and statistical applications are described.

**Keywords:** Strong mixing, strongly stationary sequences, Local Asymptotic Normality, Local Asymptotic Minimavity, SCOT models, Edgeworth expansion.

## Appendix. Simulation of an Asymmetric Cyclic Random Walk (RW)

The simulation program — written using the R statistical language — determines the number of visits of  $k$  for large numbers of sample size  $N$ , with an alphabet  $A$ . The algorithms have been augmented with pseudo-code in order to make the syntax easier to read. For the first algorithm we will determine the visits of  $k$  for a specific value of  $A$  and  $N$ , assuming the probability of moving left and remaining on the node in the ring is the same ( $P(\text{move\_left}) = P(\text{remain}) = 0.5$ ):

---

**Algorithm 1:** GET\_K\_VISITS\_FOR\_N\_JUMPS: Get visits of  $k$  after  $N$  Jumps given  $A$

---

**Inputs:**

$A$ : Number of nodes, where  $k = A/2$ . ( $A = \{10, 20, \dots, 100\}$ )

$N$ : Number of jumps. ( $N = \{20, 40, 60, 80\}$ )

$P(\text{move\_left}) = P(\text{remain}) = 0.5$

**Outputs:**

Returns the number visits at  $k$  ( $A/2$ ) after  $N$  jumps.

```

1 begin
2    $k \leftarrow A/2$ 
3    $visits \leftarrow 0$ 
4    $A \leftarrow 1$ 
5   // Randomize the random number generator with
6   // a seed based on the timer
7    $t = t \leftarrow as.numeric(Sys.time())$ 
8    $seed \leftarrow 1e8 * (t - floor(t))$ 
9    $set.seed(seed)$ 
10  while  $N > 0$  do
11     $move\_or\_stay \leftarrow sample(c(-1,0), prob =$ 
12       $c(P\_move\_left, P\_remain), size = 1, replace = TRUE)$ 
13     $A\_previous \leftarrow A$ 
14     $N \leftarrow N - 1$ 
15    if  $move\_or\_stay == -1$  then
16       $A \leftarrow A - 1$ 
17    if  $A < 1$  then
18       $A \leftarrow k * 2$ 
19    if  $A\_previous \neq A$  AND  $A == k$  then
20       $visits \leftarrow visits + 1$ 
21  return ( $visits$ )

```

---

Next we had to build the table of  $N \times A$  for one round, which is used in performing multi-round simulations for a table of  $N \times A$  containing distributions of individual rounds. Both algorithms are presented in subsequent pages.

---

**Algorithm 2:** BUILD\_N\_BY\_A: Build the  $N \times A$  table for one round.

---

**Inputs:***A\_vals*: List of number of nodes. ( $A = \{10, 20, \dots, 100\}$ )*N\_vals*: List of number of jumps. ( $N = \{20, 40, 60, 80\}$ ) $P(\text{move\_left}) = P(\text{remain}) = 0.5$ **Outputs:**Returns the table of the number visits at  $k$  ( $A/2$ ) after  $N$  jumps for different values of  $A$  and  $N$ .

```

1 begin
2    $N\_by\_A\_Table \leftarrow matrix(|N\_vals|, |A\_vals|)$ 
3   for  $A \in A\_vals$  do
4     for  $N \in N\_vals * A$  do
5        $N\_by\_A\_Table[N/A, A] \leftarrow$ 
6          $c(GET\_K\_VISITS\_FOR\_N\_JUMPS(A, N, P\_move\_left,$ 
7            $P\_remain))$ 
7   return ( $N\_by\_A\_Table$ )

```

---

Finally we run multi-round simulations for table of  $N \times A$  using the following algorithm, which gives us a distribution for each element of the  $N \times A$  table:

---

**Algorithm 3:** MULTI\_ROUND\_N\_BY\_A: Populate the  $N \times A$  table for multi-round simulations .

---

**Inputs:**

$A\_vals$ : List of number of nodes. ( $A = \{10, 20, \dots, 100\}$ )

$N\_vals$ : List of number of jumps. ( $N = \{20, 40, 60, 80\}$ )

$Rounds = 100$  is the rounds of repetition.

$P(move\_left) = P(remain) = 0.5$

**Outputs:**

Returns the table of the number visits at  $k$  ( $A/2$ ) after  $N$  jumps for different values of  $A$  and  $N$ .

```

1 begin
2   # Construct the initial table to which more values of the distributions
   # will be appended  $NA\_Table \leftarrow$ 
    $BUILD\_N\_BY\_A(A\_vals, N\_vals, P\_move\_left, P\_remain)$ 
3   for  $r \in \{1, \dots, Rounds - 1\}$  do
4     # Get the values of one round to append to each cell's distribution
      $N\_by\_A\_Table\_Sample \leftarrow$ 
      $BUILD\_N\_BY\_A(A\_vals, N\_vals, P\_move\_left, P\_remain)$ 
5     for  $A \in A\_vals$  do
6       for  $N \in N\_vals$  do
7          $value\_list \leftarrow N\_by\_A\_Table\_Sample[N, A]$ 
8          $cell\_list \leftarrow N\_by\_A\_Table[N, A]$ 
9          $cell\_list \leftarrow c(cell\_list, value\_list)$ 
10         $N\_by\_A\_Table[N, A] \leftarrow cell\_list$ 
11  return ( $N\_by\_A\_Table$ )

```

---