

# Markov substitute models and statistical inference in linguistics

Olivier Catoni, joint work with Thomas Mainguy

CNRS, INRIA – CLASSIC

Département de Mathématiques et Applications,

ENS, 45 rue d'Ulm, 75 230 Paris Cedex 05,

`Olivier.Catoni@ens.fr`

*Séminaire Parisien de Statistique,*

*Institut Henri Poincaré,*

*April 7, 2014*

# Toric grammars in action

A training sample:

- 1 [0 He is a clever guy .
- 1 [0 He is doing some shopping .
- 1 [0 He is laughing .
- 1 [0 He is not interested in sports .
- 1 [0 He is walking .
- 1 [0 He likes to walk in the streets .
- 1 [0 I am driving a car .
- 1 [0 I am riding a horse too .
- 1 [0 I am running .
- 1 [0 Paul is crossing the street .
- 1 [0 Paul is driving a car .
- 1 [0 Paul is riding a horse .
- 1 [0 Paul is walking .
- 1 [0 Peter is walking .
- 1 [0 While I was walking , I saw Paul crossing the street .

The estimated toric grammar:

10 [0 He likes to walk ]6 ]3 streets .  
2 [0 ]1 ]8 clever guy .  
2 [0 ]1 doing some shopping .  
2 [0 ]1 laughing .  
2 [0 ]1 not interested ]6 sports .  
2 [0 ]1 riding ]8 horse .  
2 [0 ]1 riding ]8 horse ]2 .  
2 [0 ]1 running .  
24 [0 ]7 am ]5 .  
28 [0 Paul is ]5 .  
40 [0 He is ]5 .  
4 [0 ]1 crossing ]3 street .  
4 [0 ]1 driving ]8 car .  
5 [0 ]4 is ]5 .  
6 [0 ]1 walking .  
7 [0 Peter is ]5 .  
8 [0 While ]7 was ]5 , ]7 saw ]4 ]5 .  
10 [1 He is  
2 [1 Peter is

2 [1 While ]7 was ]5 , ]7 saw ]4  
6 [1 ]7 am  
8 [1 Paul is  
2 [2 too  
30 [3 the  
14 [4 Paul  
1 [4 Peter  
16 [5 crossing ]3 street  
16 [5 driving ]8 car  
16 [5 riding ]8 horse  
34 [5 walking  
8 [5 ]5 too  
8 [5 ]8 clever guy  
8 [5 doing some shopping  
8 [5 laughing  
8 [5 not interested ]6 sports  
8 [5 running  
20 [6 in  
50 [7 I  
50 [8 a

## New sentences discovered:

- 1 [0 Paul is driving a car too .
- 1 [0 Paul is doing some shopping .
- 1 [0 Paul is laughing .
- 1 [0 Paul is riding a horse too .
- 1 [0 Paul is running too .
- 1 [0 Paul is running .
- 1 [0 Paul is not interested in sports too .
- 1 [0 Paul is not interested in sports .
- 1 [0 Paul is a clever guy too .
- 1 [0 Paul is a clever guy .
- 1 [0 Paul is walking too .
- 1 [0 Peter is driving a car too .
- 1 [0 Peter is driving a car .
- 1 [0 Peter is doing some shopping .
- 1 [0 Peter is laughing .
- 1 [0 Peter is riding a horse too .
- 1 [0 Peter is riding a horse .
- 1 [0 Peter is running too .
- 1 [0 Peter is running .
- 1 [0 Peter is not interested in sports .

1 [0 Peter is a clever guy .  
1 [0 Peter is crossing the street .  
1 [0 He is driving a car too .  
1 [0 He is driving a car .  
1 [0 He is riding a horse too .  
1 [0 He is riding a horse .  
1 [0 He is running too .  
1 [0 He is running .  
1 [0 He is not interested in sports too .  
1 [0 He is crossing the street too .  
1 [0 He is crossing the street .  
1 [0 He is walking too .  
1 [0 I am driving a car too .  
1 [0 I am doing some shopping .  
1 [0 I am laughing too .  
1 [0 I am laughing .  
1 [0 I am riding a horse .  
1 [0 I am not interested in sports .  
1 [0 I am a clever guy .  
1 [0 I am crossing the street too .  
1 [0 I am crossing the street .  
1 [0 I am walking too .  
1 [0 I am walking .

1 [0 While I was driving a car , I saw Paul doing some shopping too .  
1 [0 While I was driving a car , I saw Paul doing some shopping .  
1 [0 While I was driving a car , I saw Paul riding a horse .  
1 [0 While I was driving a car , I saw Paul crossing the street .  
1 [0 While I was driving a car , I saw Paul walking .  
1 [0 While I was driving a car , I saw Peter riding a horse .  
1 [0 While I was doing some shopping , I saw Paul riding a horse .  
1 [0 While I was doing some shopping , I saw Paul walking .  
1 [0 While I was laughing too , I saw Peter crossing the street .  
1 [0 While I was laughing , I saw Peter riding a horse .  
1 [0 While I was riding a horse , I saw Paul driving a car too .  
1 [0 While I was riding a horse , I saw Paul driving a car .  
1 [0 While I was riding a horse , I saw Paul laughing .

1 [0 While I was riding a horse , I saw Paul running .  
1 [0 While I was riding a horse , I saw Paul walking .  
1 [0 While I was riding a horse , I saw Peter not interested in sports .  
1 [0 While I was running , I saw Paul laughing .  
1 [0 While I was running , I saw Paul not interested in sports .  
1 [0 While I was running , I saw Paul a clever guy .  
1 [0 While I was running , I saw Paul walking .  
1 [0 While I was not interested in sports , I saw Paul driving a car .  
1 [0 While I was not interested in sports , I saw Paul riding a horse .  
1 [0 While I was a clever guy , I saw Paul running .  
1 [0 While I was a clever guy , I saw Paul crossing the street .  
1 [0 While I was a clever guy , I saw Paul walking .  
1 [0 While I was crossing the street , I saw Paul riding a horse .  
1 [0 While I was crossing the street , I saw Paul running .  
1 [0 While I was crossing the street , I saw Paul crossing the street .  
1 [0 While I was crossing the street , I saw Paul walking .  
1 [0 While I was crossing the street , I saw Peter walking .  
1 [0 While I was walking , I saw Paul driving a car .  
1 [0 While I was walking , I saw Paul laughing .  
1 [0 While I was walking , I saw Paul riding a horse .  
1 [0 While I was walking , I saw Paul running .  
1 [0 While I was walking , I saw Paul not interested in sports .  
1 [0 While I was walking , I saw Paul crossing the street too .  
1 [0 While I was walking , I saw Paul walking .  
1 [0 While I was walking , I saw Peter not interested in sports .  
1 [0 While I was walking , I saw Peter walking .



## Definition of Markov substitute sets

Let  $D$  be a dictionary of words,  $D^+ = \bigcup_{j=1}^{\infty} D^j$  the set of finite sequences of words and  $D^* = \{\epsilon\} \cup D^+$  the set of possibly empty finite sequences of words.

Let  $S \in D^+$  be a random sentence, and  $(S_i, 1 \leq i \leq n)$  a sample of  $n$  independent copies of  $S$ .

Given a context  $x = (x_1, x_2) \in (D^*)^2$ , and an expression  $y \in D^+$ , we define the insertion operator

$$\alpha(x, y) = x_1 y x_2 \in D^+,$$

that inserts the expression  $y$  in the context  $x$ .

## Definition

A subset  $B \subset D^+$  is a Markov substitute set for  $S$  when there is a probability measure  $q_B \in \mathcal{M}_+^1(B)$  on  $B$  (called the substitute measure) such that for any context  $x \in (D^*)^2$  and any  $y \in B$ ,

$$\mathbb{P}[S = \alpha(x, y)] = \mathbb{P}[S \in \alpha(x, B)] q_B(y),$$

where  $\alpha(x, B) = \{\alpha(x, y), y \in B\}$ .

In simple words, the conditional distribution of  $y$  in context  $x$  is independent of the context  $x$ .

Equivalently, for any  $x, x' \in (D^*)^2$ , any  $y, y' \in B$ ,

$$\mathbb{P}_S(x_1 y x_2) \mathbb{P}_S(x'_1 y' x'_2) = \mathbb{P}_S(x_1 y' x_2) \mathbb{P}_S(x'_1 y x'_2).$$

*(The model could be broadened further by imposing restrictive conditions on the context  $x$ .)*

## Markov chains are a special case of Markov substitute models

If  $S = (Z_1, \dots, Z_L)$ , where  $(Z_t, t \in \mathbb{N})$  is a Markov chain, then for any  $x = (w_1, w_2) \in D^2$ ,  $\alpha(x, D)$  is a Markov substitute set.

If  $S = (Z_1, \dots, Z_\tau)$ , where  $\tau$  is the first hitting time of  $C \subset D$ , then for any  $x = (x_1, x_2) \in (D^+)^2$ , any  $B \subset (D \setminus C)^+$ ,  $\alpha(x, B)$  is a Markov substitute set.

## Basic properties of Markov substitute sets

Any one point set  $\{y\}$ ,  $y \in D^+$ , is a Markov substitute set.

A subset of a Markov substitute set is itself a Markov substitute set.

If  $B$  and  $C$  are Markov sets such that  $B \cap C \neq \emptyset$ ,  $B \cup C$  is also a Markov substitute set.

The relation

$$y \sim y' \iff \{y, y'\} \text{ is a Markov substitute pair}$$

is an equivalence relation and  $D^+ / \sim$  forms a partition of  $D^+$  into maximal Markov substitute sets.

## Basic properties of Markov substitute sets

The set  $B$  is Markov if and only if there is a connected undirected spanning graph  $\mathcal{G} \subset B^2$  such that for any  $(y, y') \in \mathcal{G}$ ,  $\{y, y'\}$  is a Markov substitute pair.

If  $B_j, 1 \leq j \leq \ell$  are Markov substitute sets (including possibly some one point sets), then

$$\gamma(B_1 \dots B_\ell) = \{s = y_1 \dots y_\ell : y_j \in B_j, 1 \leq j \leq \ell\}$$

is also a Markov substitute set, and

$q_B(y_1 \dots y_\ell) = C_B \prod_{j=1}^{\ell} q_{B_j}(y_j)$ , where  $C_B$  is a suitable normalizing constant. (The map  $(y_1, y_\ell) \mapsto y_1 \dots y_\ell$  may not be one to one, in which case  $C_B$  may be different from one!)

# Characterization of Markov substitute sets in terms of random parsing

Let  $B \subset D^+$  be some subset.

Let us define the set of splits of any sentence  $s \in D^+$  as

$$\mathcal{S}(s, B) = \{(x, y), x \in (D^*)^2, y \in B, \alpha(x, y) = s\}.$$

Let us consider some conditional probability kernel  $(\pi(s; x, y), s \in D^+, x \in (D^*)^2, y \in B \cup \{\epsilon\})$  such that

$$\mathcal{S}(s, B) \subset \text{supp}(\pi(s; \cdot)) \subset \mathcal{S}(s, B) \cup \{((s, \epsilon), \epsilon)\}.$$

We can for instance take  $\pi(s; x, y) = |\mathcal{S}(s, B)|^{-1}$ ,  $(x, y) \in \mathcal{S}(s, B)$ .

Let us define the random  $B$ -parse  $X, Y$  of the random sentence  $S$  on the same probability space by its conditional distribution

$$\begin{aligned} & \mathbb{P}_{X, Y|S = s}(x, y) \\ &= \begin{cases} \min_{y' \in B} \pi(\alpha(x, y'), x, y'), & (x, y) \in \mathcal{S}(s, B), \\ 1 - \sum_{(x', y') \in \mathcal{S}(s, B)} \mathbb{P}_{X, Y|S = s}(x', y'), & x = (s, \epsilon), y = \epsilon. \end{cases} \end{aligned}$$

### Lemma

*The set  $B$  is a Markov substitute set for  $S$  if and only if one of the following equations is true*

$$\begin{aligned} & \mathbb{P}_{X, Y|Y \in B} = \mathbb{P}_{X|Y \in B} \otimes \mathbb{P}_{Y|Y \in B}, \\ & \mathbb{P}_{X|Y = y} = \mathbb{P}_{X|Y = y'}, \quad y, y' \in B, \\ & \mathbb{P}_{Y|X = x, Y \in B} = \mathbb{P}_{Y|Y \in B} \end{aligned}$$

# Invariant dynamics

Let  $(B_j, 1 \leq j \leq t)$  be a family of Markov substitute sets.

Let us consider a conditional probability kernel  $(\pi(s; x, y, j) : s \in D^+, x \in (D^*)^2, 1 \leq j \leq t, y \in B_j \cup \{\epsilon\}, \alpha(x, y) = s)$ .

The kernel  $(k(s, s') : s, s' \in D^+)$ , defined as

$$k(s, s') = \begin{cases} \sum_{\substack{x \in (D^*)^2, \\ (y, y') \in (D^*)^2, j}} \pi(s; x, y, j) q_{B_j}(y') \left( \frac{\pi(s', x, y', j)}{\pi(s, x, y, j)} \wedge 1 \right), & s \neq s', \\ 1 - \sum_{s'' \in D^+ \setminus \{s\}} k(s, s''), & s' = s, \end{cases}$$

is reversible with respect to  $\mathbb{P}_S$ .



# Proof

$$\begin{aligned}\mathbb{P}_S(s)k(s, s') &= \sum_{\substack{x \in (D^*)^2, \\ (y, y') \in (D^*)^2, j}} \mathbb{P}_S(\alpha(x, B_j)) \\ &\quad \times q_{B_j}(y)q_{B_j}(y')[\pi(s; x, y, j) \wedge \pi(s'; x, y', j)] \\ &= \mathbb{P}_S(s')k(s', s).\end{aligned}$$

## Basic properties of reversible dynamics

If  $k(y, y') > 0$ ,  $\{y, y'\}$  is a Markov substitute pair and  $q_{\{y, y'\}}(y)k(y, y') = q_{\{y, y'\}}(y')k(y', y)$ .

For any Markov substitute set  $B$ , (including one point sets),  $\text{supp}\left(\sum_{t=0}^{\infty} q_B k^t\right)$  is a Markov substitute set.

The communicating classes  $\left\{ \text{supp}\left(\sum_{t=0}^{\infty} \delta_s k^t\right), s \in D^+ \right\}$  forms a partition of  $D^+$  into Markov substitute sets.

## Basic properties of reversible dynamics

For any domain  $\mathcal{D} \subset D^+$ , the reflected dynamics

$$k_{\mathcal{D}}(s, s') = \begin{cases} k(s, s'), & s, s' \in \mathcal{D}, s \neq s', \\ 0, & s' \notin \mathcal{D} \cup \{s\}, \\ 1 - \sum_{s'' \neq s} k(s, s''), & \text{otherwise.} \end{cases}$$

is reversible with respect to  $\mathbb{P}_{\mathcal{S}}$ .

## Test functions

Let us consider a family  $\Theta$  of subsets  $\theta \subset (D^*)^2$ , containing all the one point sets  $\{x\}$ ,  $x \in (D^*)^2$ .

For any pair  $B_1, B_2$  of Markov substitute sets, such that  $B_1 \cap B_2 = \emptyset$ . let us put  $B = B_1 \cup B_2$  and let us consider some  $B$ -parse process  $(X_B, Y_B)$  and the random variables

$$\begin{aligned} F_{B_1, B_2, \theta}(X_B, Y_B, p) \\ = \mathbf{1}(X_B \in \theta) [\mathbf{1}(Y_B \in B_1) - p\mathbf{1}(Y_B \in B)], \quad \theta \in \Theta, p \in [0, 1]. \end{aligned}$$

The set  $B$  is a Markov substitute set if and only if there is  $p \in [0, 1]$  such that for any  $\theta \in \Theta$ ,  $\mathbb{E}[F_{B_1, B_2, \theta}(X_B, Y_B, p)] = 0$ . In this case  $q_B(B_1) = p$ .

## Test functions

To estimate  $\mathbb{E}[F_{B_1, B_2, \theta}(X_B, Y_B, p)]$ , we can simulate from the sample  $(S_i, 1 \leq i \leq n)$  an i.i.d. sample  $(S_i, X_{B,i}, Y_{B,i})$  such that  $\mathbb{P} X_{B,i}, Y_{B,i} | S_i = \mathbb{P} X_B, Y_B | S$ , or we can compute

$$\begin{aligned} F_{B_1, B_2, \theta}(s, p) &\stackrel{\text{def}}{=} \mathbb{E}[F_{B_1, B_2, \theta}(X_B, Y_B, \theta) | S = s] \\ &= \sum_{x \in (D^*)^2} \sum_{y \in B} \mathbb{1}(x \in \theta) \min_{y' \in B} \pi(\alpha(x, y'), x, y') \\ &\quad \times [\mathbb{1}(y \in B_1) - p \mathbb{1}(y \in B)] \mathbb{1}(s = \alpha(x, y)), \end{aligned}$$

and consider the i.i.d. samples  $F_{B_1, B_2, \theta}(S_i, p)$ .

## Alternative test functions

Another choice of test functions is

$$\begin{aligned} F_{B_1, B_2, \theta}(S, p) &= \sum_{x \in (D^*)^2} \sum_{(y_1, y_2) \in B_1 \times B_2} \mathbf{1}(x \in \theta) \\ &\quad \times [\pi(\alpha(x, y_1), x, y_1) \wedge \pi(\alpha(x, y_2), x, y_2)] \\ &\quad \times [\mathbf{1}(S = \alpha(x, y_1)) q_{B_2}(y_2)(1-p) - \mathbf{1}(S = \alpha(x, y_2)) q_{B_1}(y_1)p] \\ &= \sum_{x \in (D^*)^2} \sum_{y, y' \in (D^*)^2} \pi(S, x, y) \mathbf{1}(x \in \theta) \\ &\quad \times [\mathbf{1}(y \in B_1)(1-p)q_{B_2}(y') - \mathbf{1}(y \in B_2)pq_{B_1}(y')] \\ &\quad \times \left( \frac{\pi(\alpha(x, y'), x, y')}{\pi(\alpha(x, y), x, y)} \wedge 1 \right). \end{aligned}$$

## Simulations

Consider  $\mathbb{P}_{X', Y' | S} = \pi$ ,

$\mathbb{P}_{Y'' | X', Y'} = \mathbb{1}(Y' \in B_1)q_{B_2} + \mathbb{1}(Y' \in B_2)q_{B_1}$ ,

$$w(X', Y') = \mathbb{E} \left( \left( \frac{\pi(\alpha(X', Y''), X', Y'')}{\pi(\alpha(x', Y'), X', Y')} \wedge 1 \right) \mid X', Y' \right),$$

$\mathbb{P}_{X, Y | X', Y'} =$

$w(X', Y')\delta_{X', Y'} + (1 - w(X', Y'))\delta_{(\alpha(X', Y'), \epsilon), \epsilon}$ .

$$F_{B_1, B_2, \theta}(S, p) = \mathbb{E} \left[ \mathbb{1}(X \in \theta) [\mathbb{1}(Y \in B_1) - p\mathbb{1}(Y \in B)] \mid S \right].$$

# Statistical tests

Let  $\mathcal{B}$  be a set of known Markov substitute sets (to start with, we can take  $\mathcal{B} = \{\{y, \cdot\}, y \in D\}$ ).

$$\text{Let } C_1 = \max_{y \in D^+} \sum_{B \in \mathcal{B}} \mathbf{1}(y \in B) < \infty,$$

$$C_2 = \max_{x \in (D^*)^2} \sum_{\theta \in \Theta} \mathbf{1}(x \in \theta) < \infty,$$

$$h(B, s) = \mathbf{1}(\mathcal{S}(s, B) \neq \emptyset),$$

$$g(B, \theta, s) = \mathbf{1}(\exists (x, y) \in \mathcal{S}(s, B) : x \in \theta),$$

$$\nu(B) = \frac{1}{n} \sum_{i=1}^n \left( \sum_{B' \in \mathcal{B}} h(B', S_i) \right)^{-1} h(B, S_i),$$



# Statistical tests

$$\xi(\theta|B_1, B_2) = \frac{1}{n} \sum_{i=1}^n \left( \sum_{\theta' \in \Theta} g(B_1 \cup B_2, \theta', S_i) \right)^{-1} g(B_1 \cup B_2, \theta, S_i),$$

$$\mu(B_1, B_2, \theta) = \nu(B_1)\nu(B_2)\mu(\theta|B_1, B_2),$$

$$\bar{h}(s) = \sum_{B \in \mathcal{B}} h(B, s) \leq C_1 \ell(s)(\ell(s) + 1)/2,$$

$$\bar{g}(s) = \sup_{B_1, B_2 \in \mathcal{B}} \sum_{\theta \in \Theta} g(B_1 \cup B_2, \theta, s) \leq C_2 \ell(s)(\ell(s) + 1)/2,$$

$$\begin{aligned} \mu(B_1, B_2, \theta) &\geq \underbrace{n^{-3} \left( \max_{1 \leq i \leq n} \bar{h}(S_i) \right)^{-2} \left( \max_{1 \leq i \leq n} \bar{g}(S_i) \right)^{-1}}_{\geq 8n^{-3} C_1^{-2} C_2^{-1} L^{-3} (L+1)^{-3},} \mathbf{1}(\mu(B_1, B_2, \theta) > 0). \\ &\quad \text{where } L = \max_{1 \leq i \leq n} \ell(S_i) \end{aligned}$$

# Statistical tests

## Proposition

Consider some finite set  $\Lambda \subset ]0, 1[$ . With probability at least  $1 - 2\varepsilon$ , for any  $\lambda \in \Lambda \cup \{-\Lambda\}$ , any  $p \in \mathcal{P} \subset [0, 1]$ , any  $\rho \in \mathcal{M}_+^1(\mathcal{B}^2 \times \Theta)$ ,

$$\begin{aligned} & \sum_{i=1}^n \frac{(k-1)\lambda}{k} \int_{\theta \in \mathcal{B}^2 \times \Theta} F_{\theta}(S_i, p) \, d\rho(\theta) - \frac{\lambda^2}{2(1-|\lambda|)^2} \int F_{\theta}(S_i, p)^2 \, d\rho(\theta) \\ & \leq \int \sum_{i=1}^n \left[ \log(1 + \lambda F_{\theta}(S_i, p)) - \frac{\lambda}{k} F_{\theta}(S_i, p) \right] \, d\rho(\theta) \\ & \leq \frac{(k-1)n\lambda}{k} \int \mathbb{E}[F_{\theta}(S, p)] \, d\rho(\theta) + \mathcal{K}(\rho, \mu) + 3\log(k) + \log(|\Lambda| |\mathcal{P}| / \varepsilon). \end{aligned}$$

## Statistical tests

Let  $p(B_1, B_2, \theta) = \mathbb{P}(Y_{B_1 \cup B_2} \in B_1 \mid X \in \theta, Y_{B_1 \cup B_2} \in B_1 \cup B_2)$ ,  
so that  $\mathbb{E}(F_{B_1, B_2, \theta}(S, p(B_1, B_2, \theta))) = 0$ ,

$$p_+(B_1, B_2) = \sup \left\{ p(B_1, B_2, \theta) : \theta \in \Theta, \right.$$

$$\left. \mathbb{P}(X_{B_1 \cup B_2} \in \theta, Y_{B_1 \cup B_2} \in B_1 \cup B_2) > 0 \right\},$$

$$p_-(B_1, B_2) = \inf \left\{ p(B_1, B_2, \theta) : \theta \in \Theta, \right.$$

$$\left. \mathbb{P}(X_{B_1 \cup B_2} \in \theta, Y_{B_1 \cup B_2} \in B_1 \cup B_2) > 0 \right\},$$

$$\psi(z) = \log(1+z) - z/k,$$

## Statistical tests

We will say that  $(B_1, B_2)$  is an  $\eta$ -Markov substitute pair of sets when  $B = B_1 \cup B_2$  is a Markov substitute set such that  $q_B(B_1) \in [\eta, 1 - \eta]$ . We will say that  $(B_1, B_2)$  is a  $\gamma$ -weak  $\eta$ -Markov substitute pair of sets when

$$\eta \leq p_-(B_1, B_2) \leq p_+(B_1, B_2) \leq 1 - \eta, \text{ and } p_+(B_1, B_2) - p_-(B_1, B_2) \leq \gamma.$$

# Statistical tests

## Proposition

Let  $\Lambda$  be a finite subset of  $]0, 1[$ . With probability at least  $1 - 2\varepsilon$ , for any pair  $(B_1, B_2) \in \mathcal{B}^2$ ,

$$\begin{aligned} B_-(p_+(B_1, B_2)) &\stackrel{\text{def}}{=} \sup_{\rho \in \mathcal{M}_+^1(\Theta), \lambda \in \Lambda} \int \sum_{i=1}^n \psi\left(\lambda F_{B_1, B_2, \theta}(S_i, p_+(B_1, B_2))\right) d\rho(\theta) \\ &\quad - \mathcal{K}(\rho, \mu_1) - 3\log(k) - \log\left(\frac{|\Lambda|}{\varepsilon \nu_1(B_1) \nu_1(B_2)}\right) \leq 0 \\ B_+(p_-(B_1, B_2)) &\stackrel{\text{def}}{=} \sup_{\rho \in \mathcal{M}_+^1(\Theta), \lambda \in \Lambda} \int \sum_{i=1}^n \psi\left(-\lambda F_{B_1, B_2, \theta}(S_i, p_-(B_1, B_2))\right) d\rho(\theta) \\ &\quad - \mathcal{K}(\rho, \mu_1) - 3\log(k) - \log\left(\frac{|\Lambda|}{\varepsilon \nu_1(B_1) \nu_1(B_2)}\right) \leq 0 \end{aligned}$$

## Statistical tests

Therefore, if we reject the hypothesis that  $B_1 \cup B_2$  is a Markov substitute set when

$$\inf_{p \in [0,1]} \max\{B_-(p), B_+(p)\} > 0,$$

the probability of making a false rejection (after testing all pairs in  $\mathcal{B}^2$ ) is at most  $2\varepsilon$ .

In the same way we can reject the hypothesis that  $(B_1, B_2) \in \mathcal{B}^2$  is an  $\eta$ -Markov substitute pair of sets when

$$\inf_{p \in [\eta, 1-\eta]} \max\{B_-(p), B_+(p)\} > 0,$$

with a probability of rejecting one of the true  $\eta$ -Markov pairs (after testing all pairs in  $\mathcal{B}^2$ ), not greater than  $2\varepsilon$ .

## Statistical tests

With probability at least  $1 - 2\epsilon$ , for any  $\gamma$ -weak  $\eta$ -Markov substitute pair,

$$\inf_{p \in [\eta, 1 - \eta - \gamma]} \max\{B_-(p + \gamma), B_+(p)\} \leq 0.$$

For this test, the probability of false rejection is not greater than  $2\epsilon$ .

# Probability of false acceptance

## Lemma

For any  $p \in [0, 1]$ , any  $\lambda \in ]-1, 1[$ , any  $B_1, B_2 \in \mathcal{B}$ , any  $\theta \in \Theta$ , let  $r(B_1, B_2, \theta) = \mathbb{E}\left(\mathbb{1}(X_{B_1 \cup B_2} \in \theta, Y_{B_1 \cup B_2} \in B_1 \cup B_2)\right)$ .

With probability at least  $1 - 2\varepsilon$ ,

$$\sum_{i=1}^n \psi(\lambda F_{B_1, B_2, \theta}(S, p)) \geq \log(\varepsilon) - nr(\theta) \left[ \lambda \frac{k-1}{k} (p - p(\theta)) + \frac{\lambda^2}{1 - |\lambda|} \left( \frac{k-1}{k} + \frac{\varphi(k^{-1})}{2k^2} \right) (p(\theta)(1 - p(\theta)) + (p - p(\theta))^2) \right],$$

where  $\varphi(z) = 2z^{-2}(\exp(z) - 1 - z)$ .



## Probability of false acceptance

Let

$$\delta = \frac{1}{n} \log \left[ k^3 n^3 \left( \max_{1 \leq i \leq n} \bar{h}(S_i) \right) \left( \max_{1 \leq i \leq n} \bar{g}(S_i) \right) |\Lambda| \varepsilon^{-2} \right],$$

$$\chi = \sup_{x \in [(2n)^{-1/2}, (2n)^{1/2}]} \inf_{\lambda \in \Lambda} \cosh \left[ \log \left( \frac{\lambda x}{1 - \lambda} \right) \right],$$

$$a = \frac{4\chi^2 k}{k-1} \left( 1 + \frac{\varphi(k^{-1})}{2k(k-1)} \right) \leq 4.47\chi^2 \text{ when } k = 10,$$

$$b = \frac{(2 + \sqrt{2})k}{k-1} \leq 3.8 \text{ when } k = 10.$$

## Probability of false acceptance

Let us assume that there are  $B_1, B_2 \in \mathcal{B}$ ,  $\theta_+, \theta_- \in \Theta$  such that  $\bar{p}_+ = p(B_1, B_2, \theta_+)$ ,  $\bar{p}_- = p(B_1, B_2, \theta_-)$ ,  $r_+ = r(B_1, B_2, \theta_+)$ , and  $r_- = r(B_1, B_2, \theta_-)$  are such that

$$r_- \wedge r_+ \geq \frac{16 k \chi^2 \delta}{k-1},$$
$$\bar{p}_+ - \bar{p}_- \geq \sqrt{\frac{a \bar{p}_+ (1 - \bar{p}_+) \delta}{r_+}} \left(1 + \frac{a \delta}{r_+}\right) + \frac{b \delta}{r_+}$$
$$+ \sqrt{\frac{a \bar{p}_- (1 - \bar{p}_-) \delta}{r_-}} \left(1 + \frac{a \delta}{r_-}\right) + \frac{b \delta}{r_-}.$$

## Probability of false acceptance

With probability at least  $1 - 2\varepsilon$ ,  $\inf_{p \in [0,1]} \max\{B_-(p), B_+(p)\} > 0$ ,  
so that the probability of false acceptance of  $B_1 \cup B_2$  as a  
Markov substitute set is at most equal to  $2\varepsilon$  in this case.

More precisely, with probability at least  $1 - 2\varepsilon$ ,

$$B_- \left( \bar{p}_+ - \sqrt{\frac{a\bar{p}_+(1-\bar{p}_+)\delta}{r_+}} \left( 1 + \frac{a\delta}{r_+} \right) - \frac{b\delta}{r_+} \right) > 0,$$

$$B_+ \left( \bar{p}_- + \sqrt{\frac{a\bar{p}_-(1-\bar{p}_-)\delta}{r_-}} \left( 1 + \frac{a\delta}{r_-} \right) + \frac{b\delta}{r_-} \right) > 0.$$

## Probability of false acceptance

If we assume now that

$$\bar{p}_+ - 1 + \eta \geq \sqrt{\frac{a\bar{p}_+(1-\bar{p}_+)\delta}{r_+}} \left(1 + \frac{a\delta}{r_+}\right) + \frac{b\delta}{r_+},$$

$$\text{or that } \eta - \bar{p}_- \geq \sqrt{\frac{a\bar{p}_-(1-\bar{p}_-)\delta}{r_-}} \left(1 + \frac{a\delta}{r_-}\right) + \frac{b\delta}{r_-},$$

$$\begin{aligned} \text{or that } \bar{p}_+ - \bar{p}_- \geq & \gamma + \sqrt{\frac{a\bar{p}_+(1-\bar{p}_+)\delta}{r_+}} \left(1 + \frac{a\delta}{r_+}\right) + \frac{b\delta}{r_+} \\ & + \sqrt{\frac{a\bar{p}_-(1-\bar{p}_-)\delta}{r_-}} \left(1 + \frac{a\delta}{r_-}\right) + \frac{b\delta}{r_-}. \end{aligned}$$

## Probability of false acceptance

the false acceptance probability of the test that  $(B_1, B_2) \in \mathcal{B}^2$  is an  $\gamma$ -weak  $\eta$ -Markov substitute pair of sets is not greater than  $2\epsilon$ .

## Building syntax trees

Starting from the obvious family of Markov substitute sets  $\mathcal{A}_0 = \{\{w\}, w \in D\}$ , and assuming that  $\mathcal{A}_k \subset 2^{D^+}$  is already constructed, consider the family of Markov sets  $\mathcal{B} = \{\gamma(e), e \in \mathcal{A}_k^+\}$ .

We can use the above tests to find out new Markov substitute sets of the form  $\gamma(e_1) \cup \gamma(e_2)$ , where  $e_1, e_2 \in \mathcal{A}_k^+$ , and add them to  $\mathcal{A}_k$  to form  $\mathcal{A}_{k+1}$ .

To compute the tests, we need to define a kernel  $(\pi(s; x, y), s \in D^+, x \in (D^*)^2, y \in \gamma(e_1) \cup \gamma(e_2))$ .

## Building syntax trees

To do this, we can use two kernels

$(t(e, e'), e \in \mathcal{A}_{j-1}^+, e' \in \mathcal{A}_j^+, 1 \leq j \leq k, \gamma(e) \subset \gamma(e'))$ , and  
 $(\bar{\pi}(s, e; x, y), s \in D^+, e \in \mathcal{A}_k^+, s \in \gamma(e), e = \alpha(\bar{x}, \bar{y}), \bar{x} \in (\mathcal{A}_k^*)^2, \bar{y} \in \{e_1, e_2\}, x \in \gamma(\bar{x}_1) \times \gamma(\bar{x}_2), y \in \gamma(\bar{y}))$ .

The  $k$  th iterate of  $t$ ,  $t^k$ , builds a random syntax tree, and we can put  $\pi(s; x, y) = (t^k \bar{\pi})(s; x, y)$ .

The incremental construction of  $\mathcal{A}_k$  can be described by rewriting rules  $B \rightarrow e_1, B \rightarrow e_2$ , where  $B \in \mathcal{A}_j$ , and  $e_1, e_2 \in \mathcal{A}_{j-1}^+$ , forming a context free grammar.

## Estimating the language distribution

If  $B$  is a Markov substitute set such that  $B \cap \text{supp}(\mathbb{P}_S) \neq \emptyset$ , then  $B \subset \text{supp}(\mathbb{P}_S)$  and  $\mathbb{P}_{S|S \in B} = q_B$ .

Given a collection of Markov substitute sets  $B_j$ ,  $1 \leq j \leq t$  and the above defined reversible dynamics  $k$ , we may define the random Markov substitute sets

$$C_i = \text{supp} \left( \delta_{S_i} \sum_{j=0}^{\infty} k^j \right),$$

and estimate  $\mathbb{P}_S$  by

$$\hat{\mathbb{P}} = \frac{1}{n} \sum_{i=1}^n q_{C_i},$$



## Estimating the language distribution

and consequently  $\text{supp}(\mathbb{P}_S)$  by  $\bigcup_{i=1}^n C_i$ .

To compute  $\mathbb{1}[s \in \text{supp}(\hat{\mathbb{P}})]$  for a given  $s$  and answer the question : is  $s$  a sentence of the language ? we need to compute  $\mathbb{1}(s \in C_i)$ . The syntax tree can help here, since

$$\mathbb{1}(s \in C_i) = \mathbb{1} \left[ \text{supp} \left( \delta_s \sum_{j=0}^{\infty} k^j t^k \right) = \text{supp} \left( \delta_{S_i} \sum_{j=0}^{\infty} k^j t^k \right) \right].$$