

Problem Statement

- Most source identification algorithms are based on stationarity of data.
- Consider a non-stationary data stream in which data statistics may change abruptly from sample to sample.
- Identify sources (the models and parameters) from multiple observations.

Data Generation Model

- A single stream (block) of length T , x_1^T , has been generated by multiple sources.
- Neither the sources nor the switching times are known.
- Let y_1^T be the (unknown) indexes of the sources that generated the sequence x_1^T
- Probability of x_t (sample at time t) depends on
 - the past samples x_1^{t-1}
 - the source that generated it, y_t ,

$$P(x_1^T | y_1^T, \Delta) = \prod_{t=1}^T P_{y_t}(x_t | x_1^{t-1})$$

Δ is the set of mixture parameters.

Switching between Sources

- At time t , y_t , the index of the active source, depends on the previous ones, y_1^{t-1} .
- For simplicity, independent from x_1^{t-1} .
- Assume it is governed by a hidden Markov(1) source

$$P(y_1^T | \Delta) = w_{y_1} \prod_{t=2}^T P_h(y_t | y_{t-1})$$

w_k : the initial probability of source k

$P_h(k|l)$: the probability of switching from source l to k

$$P(x_1^T | \Delta) = \sum_{y_1^T} P(x_1^T | y_1^T, \Delta) P(y_1^T | \Delta)$$

Tree Source Modeling

Each source, S_i , can be described by

- A set \mathcal{M}_i consisting of all sequences, *contexts*, such that
 - No sequence in \mathcal{M}_i is a suffix of another one
 - For any arbitrary sequences $x_{-\infty}^{-1}$ there is a unique context $\mathbf{c} \in \mathcal{M}_i$ such that \mathbf{c} is a suffix of $x_{-\infty}^{-1}$ and $P(X_0 = a | X_{-\infty}^{-1} = x_{-\infty}^{-1}) = P(X_0 = a | X_{-l}^{-1}(\mathbf{c}) = \mathbf{c}) =: \theta(\mathbf{c}, a)$
- Conditional probability distributions on contexts

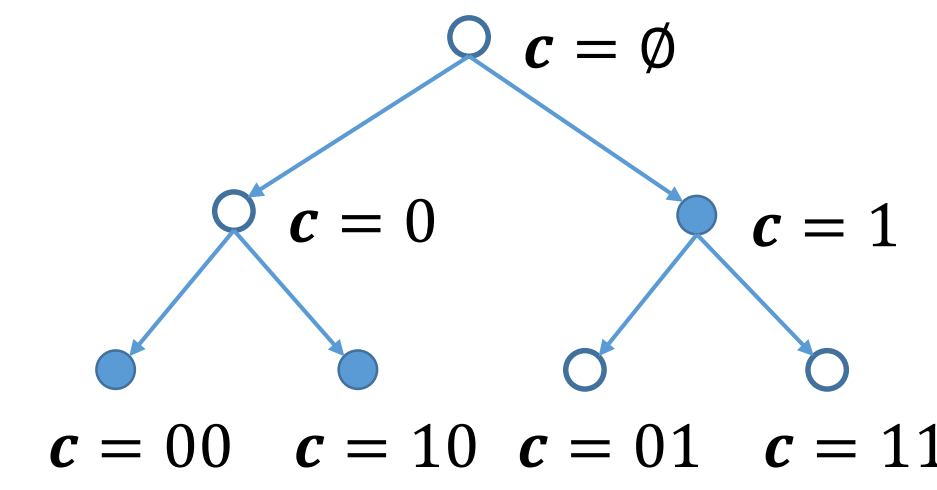
$$\Theta_i = \{\theta(\mathbf{c}, a) : a \in \mathcal{A}, \mathbf{c} \in \mathcal{M}_i\}$$

Set \mathcal{M}_i can be represented as a tree.

For example:

$$\mathcal{M} = \{00, 10, 1\}$$

$$P(011000) = P(01)P(1|1)P(0|1)P(0|10)P(0|00)$$



BIC Context Tree Estimation for Single Source

For a hypothetical tree, \mathcal{M} , of maximum depth D

- Number of free parameters $d = (|\mathcal{A}| - 1) |\mathcal{M}|$
- Histogram over tree model, $\forall \mathbf{c} \in \mathcal{M}$ and $a \in \mathcal{A}$,

$$n_x(\mathbf{c}, a) = |\{i: D < i \leq n, x_{i-l}^{i-1} = \mathbf{c}, x_i = a\}|$$

$$n_x(\mathbf{c}) = \sum_{a \in \mathcal{A}} n_x(\mathbf{c}, a)$$

- Maximum Log-Likelihood:

$$\mathcal{L}_{\mathcal{M}}(\mathbf{x}) \approx \sum_{\mathbf{c} \in \mathcal{M}, a \in \mathcal{A}} n_x(\mathbf{c}, a) \log \left(\frac{n_x(\mathbf{c}, a)}{n_x(\mathbf{c})} \right)$$

- Bayesian Information Criteria (BIC) of model w.r.t. \mathbf{x} of length T :

$$BIC_{\mathcal{M}}(\mathbf{x}) = -\mathcal{L}_{\mathcal{M}}(\mathbf{x}) + \frac{d}{2} \log T$$

- BIC model estimator

$$\hat{\mathcal{M}}_{BIC}(\mathbf{x}) = \operatorname{argmin}_{\mathcal{M}} BIC_{\mathcal{M}}(\mathbf{x})$$

Theorem [Talata and Duncan, 2011]

For a stationary ergodic source, S , with context tree \mathcal{M}_S , for any constant D , $\hat{\mathcal{M}}_{BIC}(\mathbf{x})|_D \rightarrow \mathcal{M}_S|_D$ almost surely as length of observed signal increases. Also, the maximum likelihood estimates

$$\hat{\theta}(\mathbf{c}, a) = \frac{n_x(\mathbf{c}, a)}{n_x(\mathbf{c})} \text{ converges to } P_S(a|\mathbf{c}).$$

Model Estimation for Multiple Sources

- Observed a non-stationary sequence \mathbf{x} of length T
- Effective histogram for the k^{th} source

$$\bar{n}_k(\mathbf{c}, a; \mathbf{x}) = \sum_{t=1}^T P(Y_t = k | \mathbf{x}) \mathbb{I}_{\mathbf{c}, a}(t; \mathbf{x})$$

where

$$\mathbb{I}_{\mathbf{c}, a}(t; \mathbf{x}) = \begin{cases} 1 & x_t = a, x_{t-l}^{t-1} = \mathbf{c} \\ 0 & \text{otherwise} \end{cases}$$

- For multiple sequences, $\mathcal{X} = \{\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \dots, \mathbf{x}^{(N)}\}$

$$\bar{n}_k(\mathbf{c}, a) = \sum_{n=1}^N \bar{n}_k(\mathbf{c}, a; \mathbf{x}^{(n)})$$

$$\bar{n}_k(\mathbf{c}) = \sum_{a \in \mathcal{A}} \bar{n}_k(\mathbf{c}, a)$$

- BIC model estimator for the k -th source

$$\mathcal{L}_{\mathcal{M}}(\mathcal{X}; k) \approx \sum_{\mathbf{c} \in \mathcal{M}, a \in \mathcal{A}} \bar{n}_k(\mathbf{c}, a) \log \frac{\bar{n}_k(\mathbf{c}, a)}{\bar{n}_k(\mathbf{c})}$$

$$BIC_{\mathcal{M}}(\mathcal{X}; k) = -\mathcal{L}_{\mathcal{M}}(\mathcal{X}; k) + \frac{d}{2} \log \bar{n}_k$$

$$\hat{\mathcal{M}}_k(\mathcal{X}) = \operatorname{argmin}_{\mathcal{M}} BIC_{\mathcal{M}}(\mathcal{X}; k)$$

Lemma

Assuming ergodicity and stationarity of sources, $\frac{\bar{n}_k(\mathbf{c}, a)}{\bar{n}_k(\mathbf{c})} \rightarrow P_k(a|\mathbf{c})$ almost surely as $N \rightarrow \infty$ provided that $P(Y_t = k | \mathbf{x}^{(n)}) \neq 0$ and $l(\mathbf{x}^{(n)}) > l(\mathbf{c})$ infinitely often.

Theorem

For a constant D , assume that $l(\mathbf{x}^{(n)}) > D$. Then, $\hat{\mathcal{M}}_k(\mathcal{X})|_D \rightarrow \mathcal{M}_k|_D$ almost surely as $N \rightarrow \infty$.

EM-BIC Model Estimator

- Exact values of $P(Y_t = k | \mathbf{x}^{(n)})$ are unknown.
- Iteratively estimate the probabilities

Outline of the Algorithm

- Start from an initial estimate for sources.
- Repeat until convergence:
 - Use Baum-Welch algorithm and current estimates of sources to compute $P(Y_{t-1} = l, Y_t = k | \mathbf{x}^{(n)})$ and $P(Y_t = k | \mathbf{x}^{(n)})$
 - Using updated posteriors, for the k^{th} source, estimate the model

$$\hat{\mathcal{M}}_k(\mathcal{X}) = \operatorname{argmin}_{\mathcal{M}} BIC_{\mathcal{M}}(\mathcal{X}; k)$$

- Update sources' parameters and the switching probabilities

Simulation Results

- 3 random sources over alphabet $\mathcal{A} = \{0,1,2,3\}$
- Different trees with depths 2, 2 and 3 for each source
- Entropy of each source is I
- Transition probability between source

$$A = \begin{bmatrix} 0.857 & 0.052 & 0.091 \\ 0.041 & 0.879 & 0.08 \\ 0.142 & 0.038 & 0.82 \end{bmatrix}$$

- Optimum compression rate: 1.4648 bits per symbol
- 100 sequence of length 1000 created for training and another set for test.
- KL divergence between true sources and estimated ones: less than 0.001
- Compression performance compared to PAQ8, Zip, using different number of estimated sources in the mixture model.

Algorithm	Zip	PAQ8	$\bar{R} = 1$	$\bar{R} = 2$	$\bar{R} = 3$	$\bar{R} = 4$	$\bar{R} = 5$
Redundancy (bits/1K symbols)	277.6	19.2	68.8	20.0	2.4	3.2	4.8