



EURECOM  
Department of Multimedia Communications  
2229, route des Crêtes  
B.P. 193  
06560 Sophia-Antipolis  
FRANCE

Research Report RR-11-261

## **Online Non-Negative Convolutional Pattern Learning for Speech Signals**

Dec 7<sup>th</sup>, 2011  
Last update June 12<sup>th</sup>, 2012

Dong Wang, Ravichander Vipera, Nicholas Evans and Thomas Fang Zheng

Tel : (+33) 4 93 00 81 00  
Fax : (+33) 4 93 00 82 00  
Email : dong.wang@ed.ac.uk, {vipera, evans}@eurecom.fr, fzheng@tsinghua.edu.cn

<sup>1</sup>EURECOM's research is partially supported by its industrial members: BMW Group, Cisco, Monaco Telecom, Orange, SAP, SFR, Sharp, STEricsson, Swisscom, Symantec, Thales.



# Online Non-Negative Convolutional Pattern Learning for Speech Signals

Dong Wang, Ravichander Vipplera, Nicholas Evans and Thomas Fang Zheng

## Abstract

The unsupervised learning of spectro-temporal patterns within speech signals is of interest in a broad range of applications. Where patterns are non-negative and convolutional in nature, relevant learning algorithms include convolutional non-negative matrix factorization (CNMF) and its sparse alternative, convolutional non-negative sparse coding (CNSC). Both algorithms, however, place unrealistic demands on computing power and memory which prohibit their application in large scale tasks. This technical report presents a new online implementation of CNMF and CNSC which processes input data piece-by-piece and updates learned patterns gradually with accumulated statistics. The proposed approach facilitates pattern learning with huge volumes of training data that are beyond the capability of existing alternatives. We show that, with unlimited data and computing resources, the new online learning algorithm almost surely converges to a local minimum of the objective cost function. In more realistic situations, where the amount of data is large and computing power is limited, online learning tends to obtain lower empirical cost than conventional batch learning.

## Index Terms

Non-negative matrix factorization, convolutional NMF, online pattern learning, sparse coding, speech processing, speech recognition



## CONTENTS

<b>I</b>	<b>Introduction</b>	1
<b>II</b>	<b>Online convolutive pattern learning</b>	2
II-A	Problem formulation . . . . .	2
II-B	Online CNSC . . . . .	3
<b>III</b>	<b>Complexity and convergence analysis</b>	5
III-A	Computational complexity . . . . .	5
III-B	Convergence of perfect learning . . . . .	6
III-C	Convergence of imperfect learning . . . . .	12
III-D	Batch learning and online learning . . . . .	14
<b>IV</b>	<b>Experiments</b>	15
IV-A	Speech separation . . . . .	15
IV-A1	Convergence and computing resource . . . . .	15
IV-A2	Convergence and piece length . . . . .	16
IV-A3	Convergence and data volume . . . . .	16
IV-A4	Speech separation . . . . .	17
IV-B	Denoising for speech recognition . . . . .	18
IV-B1	Experimental setup . . . . .	19
IV-B2	Incomplete noise pattern learning . . . . .	20
IV-B3	Complete noise pattern learning . . . . .	21
<b>V</b>	<b>Conclusion</b>	22
	<b>Appendix</b>	22
	<b>References</b>	23

## LIST OF FIGURES

1	An example of patterns learned with active learning. . . . .	16
2	Value of the cost function for the first 100 iterations with online and batch learning. . . . .	16
3	Average run-time for the first 100 iterations with online and batch learning. . . . .	17
4	Value of the cost function for $U = 1$ to 10 pieces and after 10 iterations for active and inertial online learning. . . . .	17
5	Average run-time for between $U = 1$ to 10 pieces and after 10 iterations for active and inertial online learning. . . . .	18
6	Value of the cost function with the multiplicative update fixed to 100 iterations and the speech data duplicated up to 50 times. . . . .	18
7	Value of the cost function with the multiplicative update fixed to 10 iterations and the speech data duplicated 50 times. . . . .	19
8	SDR of speech separation. . . . .	19
9	Average run-time for 10 iterations with online and batch learning. . . . .	20
10	ASR accuracies with CNSC-based noise cancelation where the background noise patterns are learnt on a set of randomly sampled background audio segments with batch, active online and inertial online learning. . . . .	20
11	Accuracies with CNSC-based noise cancelation where the background noise patterns are learnt using active and inertial online CNSC. . . . .	21
12	Accuracy with online CNSC-based noise cancelation where the 100 background noise patterns are learnt based on random or entire data respectively. Results using both active and inertial learning are presented. . . . .	22

## I. INTRODUCTION

Many signals exhibit clear spectro-temporal patterns; the discovery and learning of such patterns with automatic approaches is often needed for signal interpretation and for the design of suitable algorithms in practical applications. In speech signals, for instance, patterns of interest might be related to the speaker identity or the phonetic content. Whilst some of these patterns might be readily defined and learned with supervised approaches, e.g. neural networks, other, more complex patterns are difficult to pre-define and annotate particularly when they involve large datasets, hence the need for unsupervised approaches.

Various unsupervised learning techniques have been developed for automatic pattern discovery. The general idea behind such learning approaches involves the search for a number of patterns which can be used to reconstruct a set of training signals according to a certain cost function, e.g. minimum reconstruction loss, and an appropriate set of constraints. This can be written formally as:

$$\tilde{W} = \arg \min_W \left\{ \min_H \ell(X, \tilde{X}(W, H)) \right\} \quad s.t. \quad \{g_i(W, H)\} \quad (1)$$

where  $X$  represents a set of training signals and  $\tilde{X}$  their reconstructed approximations.  $\ell(\cdot, \cdot)$  represents the objective function and  $\{g_i(W, H)\}$  represents the set of constraints. The reconstruction usually takes a linear form:

$$\tilde{X}(W, H) = W \times H$$

where  $H$  represents the projection of  $\tilde{X}$  onto a set of patterns  $W$ . Pattern learning is thus closely related to matrix factorization, a field that has been studied extensively in mathematics and statistics. In signal processing and pattern learning,  $W$  is referred to as a *dictionary* whereas in statistics,  $W$  is referred to as a *basis*. The coefficient matrix  $H$  is known as a *factor matrix* or a *code matrix* in some literature. In this paper we refer to  $W$  and  $H$  as ‘patterns’ and ‘coefficients’ respectively.

Different cost functions and constraints lead to different learning techniques. An  $l_2$  reconstruction loss or Kullback-Leibler divergence cost function and a non-negative constraint applied to both patterns and coefficients leads to non-negative matrix factorization (NMF) [1]–[6]. In contrast to other pattern learning approaches NMF is capable of learning partial patterns and has thus proved to be popular in applications such as data analysis, speech processing, image processing and pattern recognition [3], [7]–[19].

A number of extensions have been introduced to improve the basic NMF approach, e.g. [20]–[31]. Convolutional NMF (CNMF) [32], [33] and sparse NMF [34]–[36] are among the most significant. Patterns learned with convolutional NMF span a number of consecutive frames and thus capture spectro-temporal features. With sparse NMF, sparsity constraints imposed on both patterns and coefficients generally lead to improved representation and noise robustness. The two extensions can be combined, resulting in a more powerful learning approach referred to as convolutional non-negative sparse coding (CNSC) [37]–[40].

While promising results have been demonstrated in some tasks such as speech enhancement [41] and source separation [42], NMF and its variants such as CNMF and CNSC place high demands on both computing resources and memory when the training database is large. The original form of the multiplicative update procedure [4] requires all the signals to be read into memory and processed in each iteration; this is prohibitive in large scale applications such as large vocabulary speech recognition which usually involves thousands of gigabytes of training data. This problem is more pronounced for both CNMF, where patterns cover a greater number of signal frames and so are usually large, and CNSC, which not only involves larger patterns but also a greater number of patterns. In some cases their number might even be larger than the signal dimension (sometimes referred as ‘over-complete patterns’). Most related publications in speech processing accordingly focus on small databases, e.g. TIDIGITS or TIMIT and learning is often based on even smaller subsets or random samples [43], [44]. Such ad-hoc learning schemes are clearly unacceptable for complex tasks.

To address this problem, we propose in this article a novel on-line learning approach for CNMF and CNSC, which processes input signals piece-by-piece and updates learned patterns gradually using accumulated statistics. With this approach, only a limited segment of the input signal is processed at a time. This approach resolves the problem of memory usage and computing cost from which conventional NMF suffers, thereby facilitating learning from large databases. As is the case for batch learning, we prove that the proposed online approach *almost surely* converges to a local minimum of the objective cost function when the amount of data and computational resources

are unlimited. We furthermore demonstrate that the new online approach tends to obtain lower empirical cost than batch learning in practical applications.

In the following section, we first formulate the learning task and present the online CNSC algorithm (CNMF can be regarded as a special case of CNSC with zero sparsity). Section III presents a complexity analysis and convergence study. Experimental results are reported in Section IV. Our conclusions are presented in Section V with ideas for further work.

## II. ONLINE CONVOLUTIVE PATTERN LEARNING

### A. Problem formulation

CNSC can be formulated according to different cost functions [39], [45]. We adopt the formulation in [39] which defines the learning problem as the minimization of the following cost function:

$$f(W, H) = \ell(X, \tilde{X}(W, H)) \quad \text{s.t.} \quad W_{i,j,k}, H_{i,j} \geq 0 \quad (2)$$

where  $X \in \mathbb{R}_{0,+}^{M \times N}$  represents the original signal of length  $N$  in  $M$ -dimensional space<sup>1</sup> and  $\tilde{X}$  is its reconstructed approximation. It is obtained from a pattern matrix  $W \in \mathbb{R}_{0,+}^{M \times R \times P}$  with  $R$  patterns of convolution range  $P$  and a coefficient matrix  $H \in \mathbb{R}_{0,+}^{R \times N}$  according to:

$$\tilde{X}(W, H) = \sum_{p=0}^{P-1} W(p) \overset{p \rightarrow}{H} \quad \text{s.t.} \quad H_{i,j} \geq 0 \quad (3)$$

where  $\overset{p \rightarrow}{H}$  shifts  $H$  by  $p$  columns to the right and where  $W(p) \in \mathbb{R}_{0,+}^{M \times R}$  is the pattern matrix corresponding to  $\overset{p \rightarrow}{H}$ . Finally, the cost function is the sparse-regularized least square distance given by:

$$\ell(X, \tilde{X}) = \|X - \tilde{X}\|_2^2 + \lambda \|H\|_l \quad (4)$$

where  $\|\cdot\|_l$  denotes the element-wise  $l$ -norm, which is equivalent to the sum of squares of the matrix elements when  $l = 2$  or the sum of their absolute values when  $l = 1$ . The factor  $\lambda$  is introduced to control the sparsity of  $H$ . To avoid a nullified  $H$ , the pattern matrix in  $W$  are forced to be unity, i.e.,  $\|W(p)\|_2 = 1$  for any  $p$ .

Note that in the optimization problem (2), both patterns  $W$  and coefficients  $H$  are free variables and need to be optimized simultaneously, even if patterns  $W$  are the primary target. This co-optimization problem is not convex and it is impossible to find a globally optimal solution. A multiplicative update approach is presented in [39] to search for a local minimum solution by extending the procedure presented in the seminal NMF paper [4]. This is formulated as follows:

$$H \leftarrow H \odot \frac{[W(p)]^T \overset{\leftarrow p}{\tilde{X}}}{[W(p)]^T \overset{\leftarrow p}{\tilde{X}} + \lambda \Xi} \quad (5)$$

$$W(p) \leftarrow W(p) \odot \frac{X \overset{p \rightarrow T}{H}}{\tilde{X} \overset{p \rightarrow T}{H}} \quad (6)$$

where  $\odot$  is the element-wise product and where the division is also element-wise.  $\Xi$  is a matrix whose elements are all equal to 1. Note that the update of  $H$  is different for different  $p$  and so, in practice,  $H$  is averaged over all  $p$  to obtain the updated coefficients.

The above equations show that most of the computation is involved in calculating the reconstruction  $\tilde{X}$ , which has a complexity of  $O(M \times N \times R \times P)$ . This is highly demanding for large pattern sets (large  $R$ ) and large databases (large  $N$ ). More importantly, since all signals must be loaded into memory and processed together, both the memory and computational requirements become prohibitive when the training corpus is large.

<sup>1</sup>The term ‘‘signal’’ here denotes any sequential data which may be non-negative in its natural form. In general, however, they are alternative or transformed non-negative representations of the original signal, such as power spectra.

## B. Online CNSC

In order to extend the application of CNSC to large scale tasks which involve large volumes of training data and complex patterns, we present an online learning approach which reads in and processes only a part of the training data at a time and updates patterns gradually until the whole training corpus is processed. The online approach has been presented previously to train probabilistic models in machine learning (e.g., [46], [47]), however it has seldom been studied in a multiplicative update setting such as in CNSC. A recent contribution presented by Mairal *et al.* is online dictionary learning (ODL) [48]. ODL reads in and decomposes signals frame-by-frame and updates patterns as each frame is processed. The authors show that such ‘partial learning’ almost surely converges to a stationary point of the objective function, given unlimited training data and a few reasonable assumptions. Similar research can be found in [49], [50].

In this paper we present an alternative online learning approach which is based on the simple NMF-style multiplicative update rule, and employ this approach to learn temporal patterns based on CNSC. Note that, while ODL supports NMF or sparse NMF by enforcing the non-negativity constraint on both patterns and coefficients, our work is the first to couple online and convolutive learning.

We start by designing a partial learning formulation which retains the temporal information within training data. We define a signal *piece* as a number of neighboring frames within which the signal is correlated, while different pieces are assumed to be independent. Temporal patterns can be learned by processing pieces sequentially and separately. Through a simple re-arrangement, the pattern update rule (6) can be re-written as follows:

$$W(p) \leftarrow W(p) \odot \frac{\sum_u \dot{B}(p, u)}{\sum_q W(q) \sum_u \dot{A}(q, p, u)} \quad (7)$$

where  $u$  is the piece index and

$$\begin{aligned} \dot{A}(q, p, u) &= \overset{q \rightarrow p \rightarrow T}{H_u} H_u \\ \dot{B}(p, u) &= X_u \overset{p \rightarrow T}{H_u} \end{aligned}$$

are the statistics contributed by piece  $u$ . The contribution of the first  $u$  pieces can then be ‘memorized’ in two auxiliary variables defined as follows:

$$A(q, p; u) = \sum_{t=1}^u \dot{A}(q, p, t)$$

and

$$B(p; u) = \sum_{t=1}^u \dot{B}(p, t).$$

The most significant difference between rules (7) and (6) is that the training data are broken into small pieces and processed sequentially. The application of rule (7) is thus more suitable in applications involving live, streamed data and the adaptive learning of new patterns in time-variant data. Second, through rule (7) the contribution of processed signals is stored in two auxiliary variables. Their size is independent of training data quantities and thus reduces memory and computational demands and hence enables the learning of complex patterns from large corpora. Finally, piece-wise learning allows the updating of patterns with each new signal piece. This leads to ‘early learning’ which significantly increases convergence speed as presented in Sections III and IV.

This leads to online CNSC which is presented in Algorithm 1. The flag variable *activeW* defines two different learning schemes: if *activeW* = *True*, both patterns and coefficients are updated  $K$  times iteratively when processing each piece; if *activeW* = *False*, only coefficients are iteratively updated. The former approach is referred to as *active learning* whereas the second approach is referred to as *inertial learning*. Note that, in both cases, patterns are nonetheless learned actively with the first piece to ensure reasonable initialization of the pattern matrix. In general, active learning converges with fewer iterations than inertial learning but places a greater demand on computing resources. We address this point further in Sections III and IV. Algorithm 2 illustrates the pattern update process (7). Matlab code for these algorithms is available online<sup>2</sup>.

<sup>2</sup><http://audio.eurecom.fr/software>

**Algorithm 1** Online CNSC learning

---

```

1: U: number of pieces
2: K: iteration
3:  $A(i, j) \in \mathbb{R}^{R \times R}$ ,  $0 < i, j < P$ 
4:  $B(i) \in \mathbb{R}^{M \times R}$ ,  $0 < i < P$ 
5:  $A(i, j) \leftarrow 0$ ,  $\forall i, j$ 
6:  $B(i) \leftarrow 0$ ,  $\forall i$ 
7: for  $u := 0$  to  $U-1$  do
8:   randomize(H)
9:   for  $k := 0$  to  $K-1$  do
10:    if (activeW=true) or ( $k = 0$ ) then
11:       $W = \text{updateW}(A, B, X_u, W, H)$ 
12:    end if
13:     $H = \text{updateH}(X, W, H)$  (Eq.5)
14:  end for
15:   $[\dot{A}, \dot{B}, W] = \text{updateW}(A, B, X_u, W, H)$ 
16:   $A(i, j) \leftarrow A(i, j) + \dot{A}(i, j)$ 
17:   $B(i) \leftarrow B(i) + \dot{B}(i)$ 
18: end for

```

---

**Algorithm 2** CNSC pattern update**Require:**  $A, B, X, W, H$ 


---

```

1:  $\dot{A} \in \mathbb{R}^{R \times R}$ ,  $0 < i, j < P$ 
2:  $\dot{B} \in \mathbb{R}^{M \times R}$ ,  $0 < i < P$ 
3:  $\dot{A}(i, j) = H \overset{i \rightarrow j}{\rightarrow} H^T$   $\forall i, j$ 
4:  $\dot{B}(i) = X \overset{i \rightarrow}{\rightarrow} H^T$   $\forall i$ 
5:  $A = A + \dot{A}$ 
6:  $B = B + \dot{B}$ 
7: for  $p := 0$  to  $P-1$  do
8:    $F \leftarrow 0$ 
9:   for  $q := 0$  to  $P-1$  do
10:     $F = F + W_q A(q, p)$ 
11:  end for
12:   $\dot{W}_p = W_p \odot \frac{B(p)}{F}$ 
13: end for
14:  $W_p = \frac{\dot{W}_p}{\|\dot{W}_p\|_2^2}$   $\forall p$  s.t.  $W_p \in \mathbb{R}_{0,+}^{M \times R}$ 
15: return  $[A, B, W]$ 

```

---

We note that the online CNSC algorithm emphasizes pattern learning and thus coefficients obtained for each signal piece might be sub-optimal as a result of early learning. This is in contrast to batch learning where patterns and coefficients are optimized simultaneously with respect to the objective cost function. With the learned patterns, however, the coefficients can be optimized easily either by iteratively applying (5) or by more efficient techniques such as quadratic optimization; both are amenable to parallel computation. Finally we note that the choice of segmentation of signals into pieces is a trade-off between intra-piece correlation and inter-piece independence. For speech signals, a segmentation according to sentence boundaries avoids the splitting of voiced patterns and is thus a natural choice.

### III. COMPLEXITY AND CONVERGENCE ANALYSIS

In this section we analyze the computational complexity and convergence of the online CNSC algorithm. We show that online learning requires comparable (for active learning) or less (for inertial learning) computations and *almost surely* converges to a local minimum of the same cost as in the batch learning. Here ‘almost surely’ means that the convergence is guaranteed with probability one. In practice, with the same computational load, online learning tends to learn better patterns than batch learning.

#### A. Computational complexity

The computing demand for both conventional batch learning and online learning consists of updating the coefficient matrix  $H$  and the pattern matrix  $W$ . We start the analysis with the computation required for one iteration.

Firstly for the coefficient update rule (5) which is shared by both batch and online learning, one update requires

$$(2M + PM + 2)RN$$

multiplications and  $RN$  divisions. Considering updates for various shifts  $p$ , the computational complexity is in the order of

$$O(M \times N \times R \times P)$$

for one iteration and is identical for both batch and online learning.

In addition, rule (6) shows that one iteration of pattern update in the batch learning requires

$$(3N + 1)MRP$$

multiplications and  $MRP$  divisions for one iteration, while online pattern update (7) requires

$$(N + M)R^2P^2 + (N + 1)MRP$$

multiplications and  $MRP$  divisions. If the signals are segmented into  $U$  pieces, the computing demand to process the entire signals amounts to

$$\sum_{u=0}^{U-1} (N_u + M)R^2P^2 + (N_u + 1)MRP$$

or

$$NR^2P^2 + NMRP + UMR^2P^2 + UMRP$$

multiplications and  $UMRP$  divisions.

In the case where  $N$  is dominant, batch learning and online learning require approximately

$$3MRPN$$

and

$$(RP + M)RPN$$

multiplications respectively.

A simple calculation shows that online learning is more efficient if  $2M > RP$ . This implies that, with high dimensional features and when learning a small number of patterns with a small convolution range, online learning update rule (7) is more efficient than batch learning update rule (6). For example, in speech processing we usually choose  $M = 128$  for power spectra, and the convolution range is often chosen to be small, e.g.,  $P = 4$ . If we learn a modest number of patterns, i.e.,  $R < 64$ , then the online algorithm is more efficient than the batch algorithm. For

sparse coding where the pattern matrix is over-complete e.g.,  $R > M$ , then online learning is slower than batch learning. The compensation, however, is that a greater volume of training data can be handled.

Now consider computing complexity for multiple iterations. To simplify the comparison, we assume that both online and batch learning use update rule (7) and therefore have the same complexity for a single update. For batch learning,  $K$  iterations require  $K$  times the computation required for one iteration. For online learning, without considering the special treatment for the first piece, the active learning invokes  $K$  iterations for coefficient update and  $K + 1$  iterations for pattern update, while inertial learning invokes  $K$  iterations for coefficient update and 1 iteration for pattern update. Thus active learning always requires more computation than batch learning, while inertial learning is always more efficient than batch learning. As we discuss in the next section, both active and inertial learning converge almost surely and the greater computational demand for active learning is compensated for by faster convergence.

### B. Convergence of perfect learning

We here study the convergence behavior of the online CNSC algorithm. As in ODL, we define the learning task as an optimization problem which aims to minimize an *objective cost function*  $f_u(W)$  with respect to the pattern matrix  $W$ , where  $f_u(W)$  is defined as follows:

$$f_u(W) \equiv \frac{1}{u} \sum_{t=1}^u \ell_{X_t}(W)$$

where

$$\ell_X(W) = \min_H \frac{1}{|X|} \ell(X, \tilde{X}(W, H))$$

is the cost of signal  $X$  and  $|X|$  denotes the number of frames of  $X$ . The limit of the objective cost function is defined as the *expected cost function*, denoted by  $f(W)$  and given as follows:

$$\begin{aligned} f(W) &\equiv \mathbb{E}_X[\ell_X(W)] \\ &= \lim_{u \rightarrow \infty} f_u(W) \end{aligned}$$

where  $\mathbb{E}_X$  represents expectation over  $X$ .

Note the definitions of  $f_u(W)$  and  $f(W)$  are independent of specific learning process. In order to study the convergence of the proposed CNSC online learning algorithm, we define  $W_t$  as the pattern matrix learned at the  $t^{\text{th}}$  step, and  $H_t$  as the coefficient matrix of the  $t^{\text{th}}$  signal obtained in learning. An *empirical cost function* is defined as follows to evaluate quality of the learning process:

$$\hat{f}_u(W) = \frac{1}{u} \sum_{t=1}^u \ell_{X_t}(W, H_t) \quad (8)$$

where

$$\ell_X(W, H) = \frac{1}{|X|} \ell(X, \tilde{X}(W, H))$$

is the empirical cost of signal  $X$ . Note that the coefficients  $H_t$  are ‘imperfect’ in general, which means the multiplicative update does not converge for each piece of signal, usually due to limitations on computing resource. We therefore call  $\hat{f}_u(W)$  *imperfect empirical cost function*, and  $\ell_{X_t}(W, H_t)$  *imperfect cost* of  $X_t$ .

If the computing resource is unlimited and the learning is ‘perfect’, the coefficient matrix  $H_t$  is optimized and explicitly denoted by:

$$\hat{H}_t = \arg \min_H \ell(X_t, \tilde{X}(W_t, H))$$

or

$$\hat{H}_t = \arg \min_H \ell(X_t, \tilde{X}(W_{t-1}, H))$$

with active and inertial learning respectively. The corresponding empirical cost function is referred to as *perfect empirical cost function* and is explicitly denoted by  $\hat{f}$ :

$$\hat{f}_u(W) = \frac{1}{u} \sum_{t=1}^u \ell_{X_t}(W, \hat{H}_t) \quad (9)$$

where  $\ell_{X_t}(W, \hat{H}_t)$  is referred to as *perfect cost* of  $X_t$ .

We now seek to prove that the empirical cost  $\hat{f}_u(W_u)$  of online learning almost surely converges to the objective cost  $f_u(W_u)$  when  $u$  approaches infinity. Before presenting the proof, we first state a few necessary assumptions.

- (A) **The training signals follow a time-invariant distribution supported by a compact set.** This assumption is reasonable for many signals such as audio and video due to the acquisition process.
- (B) **The update rules (5) and (6) are well-defined so that the multiplicative update converges to local minima with unlimited iterations.** Theoretically the convergence with this simple rule is not guaranteed, but it has little impact in practice. In addition, with a simple modification, the convergence can be enforced [6], [15]. We do not consider such complexity in this work, and just assume convergence. Reasonable initialization for the update process and a bounded denominator matrix in the update rules help to respect this assumption.
- (C) **The empirical cost function  $\hat{f}_u$  is strictly convex with lower-bounded Hessians.** This assumption can be guaranteed with a threshold on the smallest eigenvalue of the accumulated statistics  $\frac{1}{u}A(u)$  [48]<sup>3</sup>.
- (D) **The existence of a unique solution for the coefficient matrix is satisfied for signal pieces.** For CNMF, this means the rank of the pattern matrix is not larger than the feature dimension; for CNSC, this means a unique sparse code exists and can be found by  $l_1$  optimization [51].

First notice that assumption (B) implies batch learning converges to a local minimum of the cost function (4). Convergence proof for online learning is not straightforward and involves several lemmas regarding the convergence of variation and the Lipschitz property of the objective and empirical cost functions.

We first prove several lemmas regarding convergence of variation and the Lipschitz property of the objective and empirical cost functions. These lemmas are required to prove convergence of online learning.

*Lemma 31:* Given assumptions (A)-(D), the following convergence properties hold for the objective cost:

- (1)  $|f_{u+1}(W) - f_u(W)| = O(\frac{1}{u})$
- (2)  $|f_{u+1}(W_{u+1}) - f_u(W_u)| = O(\frac{1}{u})$

*Proof:*

$$\begin{aligned} f_{u+1}(W) - f_u(W) &= \frac{uf_u(W) + \ell_{X_{u+1}}(W)}{u+1} - f_u(W) \\ &= \frac{\ell_{X_{u+1}}(W) - f_u(W)}{u+1} \end{aligned}$$

Note that both  $\ell_{X_{u+1}}(W)$  and  $f_u(W)$  are bounded due to assumption (A) and the non-negative constraint. This leads to result (1).

Moreover, if  $f_{u+1}(W_{u+1}) \geq f_u(W_u)$ , we have:

$$\begin{aligned} |f_{u+1}(W_{u+1}) - f_u(W_u)| &= f_{u+1}(W_{u+1}) - f_u(W_u) \\ &\leq f_{u+1}(W_u) - f_u(W_u) \\ &= O(\frac{1}{u}) \end{aligned}$$

Since  $f_u(W_u)$  is lower bounded by 0 and the positive sequence  $f_{u+1}(W_{u+1}) - f_u(W_u)$  is  $O(\frac{1}{u})$ , the negative sequence must be at most  $O(\frac{1}{u})$ , i.e.,  $|f_{u+1}(W_{u+1}) - f_u(W_u)| = O(\frac{1}{u})$  if  $f_{u+1}(W_{u+1}) < f_u(W_u)$ . Therefore result (2) is proved. ■

*Lemma 32:* Given assumptions (A)-(D), if the update process converges for each signal piece, then  $\hat{f}_{u+1} - \hat{f}_u$  is Lipschitz with a factor  $k = O(\frac{1}{u})$

*Proof:*

First notice for any matrix  $Q$ ,

$$\|Q\|_2^2 = \mathbf{Tr}(QQ^T)$$

and for any square matrix  $Q$

<sup>3</sup>Precisely,  $A(u) = \sum_{t=1}^u \frac{\hat{A}(t)}{|X_t|}$  following the utterance-averaged cost function (9). This is different from the statistics in Algorithm 1 where the cost function is frame-averaged. The use of utterance-averaged costs in the proof simplifies notation; the proof presented here can be applied similarly to the frame-averaged cost function.

$$\mathbf{Tr}\left(\frac{Q}{\|Q\|_2}\right) \leq \mathbf{dim}(Q).$$

where  $\mathbf{Tr}(\cdot)$  denotes trace, and  $\mathbf{dim}(\cdot)$  denotes dimension.

We then can re-write the empirical cost function in a trace form as follows:

$$\begin{aligned} \hat{f}_u(W) &= \frac{1}{u} \sum_{t=1}^u \frac{1}{|X_t|} (\|X_t - WH_t\|_2^2 + \lambda \|H_t\|_1) \\ &= \frac{1}{u} \{ \mathbf{Tr}(W^T W \tilde{A}_u) - 2 \mathbf{Tr}(W^T \tilde{B}_u) \} \\ &\quad + \frac{1}{u} \sum_{t=1}^u \frac{1}{|X_t|} (\mathbf{Tr}(X_t X_t^T) + \lambda \|H_t\|_1) \end{aligned}$$

where  $\tilde{A}_u$  and  $\tilde{B}_u$  are the accumulated statistics at the  $u^{\text{th}}$  step. This results in variation of the empirical cost as follows:

$$\begin{aligned} \delta_{\hat{f}}(W) &= \hat{f}_{u+1}(W) - \hat{f}_u(W) \\ &= \frac{1}{u} \mathbf{Tr}(W^T W (\tilde{A}_{u+1} - \tilde{A}_u)) \\ &\quad + \frac{2}{u} \mathbf{Tr}(W^T (\tilde{B}_u - \tilde{B}_{u+1})) \\ &\quad + \epsilon \end{aligned}$$

where  $\epsilon$  is a factor independent of  $W$ . Therefore we have:

$$\begin{aligned} \delta_{\hat{f}}(W_1) - \delta_{\hat{f}}(W_2) &= \frac{1}{u} \|\Delta W\|_2^2 \mathbf{Tr}\left(\frac{\Delta(W^T W) \Delta \tilde{A}}{\|\Delta W\|_2^2}\right) \\ &\quad + \frac{2}{u} \|\Delta W\|_2 \mathbf{Tr}\left(\frac{\Delta W^T \Delta \tilde{B}}{\|\Delta W\|_2}\right) \end{aligned}$$

where  $\Delta W = W_1 - W_2$  and  $\Delta(W^T W) = W_1^T W_1 - W_2^T W_2$ .

According to our assumptions,  $W$  is in a compact set and so  $\|\Delta W\|_2$  is bounded; furthermore,  $\Delta \tilde{A}$  and  $\Delta \tilde{B}$  are bounded since the signals and coefficients are bounded. Then there exists a value  $k$  with which  $\delta_{\hat{f}}(W)$  is Lipschitz, i.e.,

$$|\delta_{\hat{f}}(W_1) - \delta_{\hat{f}}(W_2)| \leq \frac{k}{u} \|W_1 - W_2\|_2. \quad \blacksquare$$

With these results, the following property holds for online learning:

*Lemma 33:* Given assumptions (A)-(D), if the update process converges for each signal piece, then inertial online learning ensures that:

$$\hat{f}_u(W_{u+1}) - \hat{f}_u(W_u) = O\left(\frac{1}{u}\right)$$

*Proof:*

$$\begin{aligned}
\hat{f}_u(W_{u+1}) - \hat{f}_u(W_u) &= \hat{f}_u(W_{u+1}) - \hat{f}_{u+1}(W_{u+1}) \\
&\quad + \hat{f}_{u+1}(W_{u+1}) - \hat{f}_{u+1}(W_u) \\
&\quad + \hat{f}_{u+1}(W_u) - \hat{f}_u(W_u) \\
&\leq [\hat{f}_u(W_{u+1}) - \hat{f}_{u+1}(W_{u+1})] \\
&\quad - [\hat{f}_u(W_u) - \hat{f}_{u+1}(W_u)] \\
&= -\delta_{\hat{f}}(W_{u+1}) + \delta_{\hat{f}}(W_u) \\
&\leq \frac{k}{u} \|W_{u+1} - W_u\|_2
\end{aligned}$$

where  $k$  is a constant. Note that we have applied Lemma 32 and used the fact that  $\hat{f}_{u+1}(W_{u+1}) \leq \hat{f}_{u+1}(W_u)$ . ■

Compared to inertial learning where  $H_u$  is first optimized with respect to  $W_{t-1}$  and then  $W_u$  is updated once but not optimized, active learning jointly optimizes the patterns  $W_u$  and coefficients  $H_u$  at the  $u^{\text{th}}$  step. This joint optimization usually leads to faster pattern learning with active learning, although it comes at the risk of overfit to initial data. The following lemma states that active learning converges faster than inertial learning.

*Lemma 34:* Given assumptions (A)-(D), if the update process converges for each signal piece, active online learning ensures that:

$$\hat{f}_u(W_{u+1}) - \hat{f}_u(W_u) = O\left(\frac{1}{u^2}\right)$$

*Proof:*

Following the same procedure as in inertial learning, it can be proved that

$$\hat{f}_u(W_{u+1}) - \hat{f}_u(W_u) \leq \frac{k}{u} \|W_{u+1} - W_u\|_2.$$

On the other hand, since the empirical cost  $\hat{f}_u$  is strictly convex according to Assumption (C), the second-order growth condition can be verified as:

$$\hat{f}_u(W_{u+1}) - \hat{f}_u(W_u) \geq k' \|W_{u+1} - W_u\|_2^2$$

where  $k'$  is a constant.

This results in

$$\|W_{u+1} - W_u\|_2 \leq \frac{kk'}{u}$$

and therefore:

$$\hat{f}_u(W_{u+1}) - \hat{f}_u(W_u) \leq \frac{k^2k'}{u^2}.$$

We can now prove the convergence of online learning using Lemma 33 and Lemma 34. For that, we first prove that the empirical cost variation is  $O(\frac{1}{u})$  with both active and inertial learning. ■

*Proposition 35:* Given assumptions (A)-(D), if the multiplicative update converges for each signal piece, the following property holds with both active and inertial online learning:

$$\hat{f}_{u+1}(W_{u+1}) - \hat{f}_u(W_u) = O\left(\frac{1}{u}\right).$$

*Proof:*

For active learning,

$$\hat{f}_{u+1}(W_{u+1}) = \frac{u\hat{f}_u(W_{u+1}) + \ell_{X_{u+1}}(W_{u+1})}{u+1}.$$

Therefore

$$\begin{aligned}\hat{f}_{u+1}(W_{u+1}) - \hat{f}_u(W_u) &= \frac{u(\hat{f}_u(W_{u+1}) - \hat{f}_u(W_u))}{u+1} \\ &\quad + \frac{\ell_{X_{u+1}}(W_{u+1}) - \hat{f}_u(W_u)}{u+1}\end{aligned}$$

It is easy to verify that  $\hat{f}_{u+1}(W_{u+1}) - \hat{f}_u(W_u)$  is  $O(\frac{1}{u})$  according to Lemma 34 and the fact that  $\ell(W)$  and  $\hat{f}(W)$  are bounded.

For inertial learning,

$$\hat{f}_{u+1}(W_{u+1}) = \frac{u\hat{f}_u(W_{u+1}) + \ell_{X_{u+1}}(W_u)}{u+1}.$$

We therefore have

$$\begin{aligned}\hat{f}_{u+1}(W_{u+1}) - \hat{f}_u(W_u) &= \frac{u(\hat{f}_u(W_{u+1}) - \hat{f}_u(W_u))}{u+1} \\ &\quad + \frac{\ell_{X_{u+1}}(W_u) - \hat{f}_u(W_u)}{u+1}\end{aligned}$$

Again, this is  $O(\frac{1}{u})$  according to Lemma 33 and the bounded  $\ell(W)$  and  $\hat{f}(W)$ . ■

The convergence of variation  $\hat{f}_{u+1} - \hat{f}_u$  with  $O(\frac{1}{u})$  suggests that online learning tends to be stable with a large amount of data, however it does not ensure convergence of  $\hat{f}_u$ . A faster convergence of the variation is required to prove convergence of the learning process.

We show that this convergence is almost sure with both active and inertial learning, conditioned on sufficient training data. As Proposition (3) in [48], our proof applies Theorem A1 from [52], which states that if the sum of the positive variations of a sequence  $v_u$  is bounded, then  $v_u$  is quasi-martingale, which converges with probability one (see Theorem A1 in Appendix). We first prove the convergence of active learning, and then apply the same approach to inertial learning with a minor change.

*Proposition 36:* Given assumptions (A)-(D), if the multiplicative update converges for each signal piece, the empirical cost converges to the object cost almost surely with active online learning, i.e.,

$$\lim_{u \rightarrow \infty} \hat{f}_u(W_u) = \lim_{u \rightarrow \infty} f_u(W_u) \quad a.s.$$

*Proof:*

Defining  $v_u \equiv \hat{f}_u(W_u)$ , we have:

$$\begin{aligned}v_{u+1} - v_u &= \hat{f}_{u+1}(W_{u+1}) - \hat{f}_u(W_u) \\ &\leq \frac{u(\hat{f}_u(W_{u+1}) - \hat{f}_u(W_u))}{u+1} \\ &\quad + \frac{\ell_{X_{u+1}}(W_u) - \hat{f}_u(W_u)}{u+1} \\ &\quad + \frac{f_u(W_u) - \hat{f}_u(W_u)}{u+1}\end{aligned} \tag{10}$$

where we have applied

$$\ell_{X_{u+1}}(W_{u+1}) < \ell_{X_{u+1}}(W_u).$$

To use Theorem A1, define the filter of the past information as  $\mathcal{F}_u$ . Considering the causal nature of  $\mathcal{F}_u$  and applying the fact that  $f_u - \hat{f}_u \leq 0$ , we have

$$\begin{aligned} \mathbb{E}[v_{u+1} - v_u | \mathcal{F}_u] &\leq \frac{u(\hat{f}_u(W_{u+1}) - \hat{f}_u(W_u))}{u+1} \\ &\quad + \frac{\mathbb{E}[\ell_{X_{u+1}}(W_u)] - f_u(u)}{u+1} \\ &= \frac{u(\hat{f}_u(W_{u+1}) - \hat{f}_u(W_u))}{u+1} \\ &\quad + \frac{f(W_u) - f_u(W_u)}{u+1} \\ &\leq \frac{u(\hat{f}_u(W_u) - \hat{f}_u(W_u))}{u+1} \\ &\quad + \frac{\|f - f_u\|_\infty}{u+1} \end{aligned}$$

where

$$\|f - f_u\|_\infty = \sup_W |f(W) - f_u(W)|.$$

According to Lemma 34,

$$\hat{f}_u(W_{u+1}) - \hat{f}_u(W_u) = O\left(\frac{1}{u^2}\right)$$

and according to Lemma A2 from [53] (see in Appendix),

$$\mathbb{E}[\|f - f_u\|_\infty] = O\left(\frac{1}{\sqrt{u}}\right).$$

We therefore have:

$$\mathbb{E}[\mathbb{E}[v_{u+1} - v_u | \mathcal{F}_u]^+] = O\left(\frac{1}{u^{\frac{3}{2}}}\right).$$

Thus the expected positive variance of the process  $v$  is bounded, so Theorem A1 can be applied to prove  $v = \hat{f}_u(W_u)$  converges with probability one, and that

$$\sum_{u=1}^{\infty} |\mathbb{E}[v_{u+1} - v_u | \mathcal{F}_u]| < +\infty \quad a.s.$$

This further implies that:

$$\sum_{u=1}^{\infty} \mathbb{E}[v_{u+1} - v_u | \mathcal{F}_u] > -\infty \quad a.s.$$

Returning to (10) and by summing over all variations on the two sides, we can prove that:

$$\sum_{u=1}^{\infty} \frac{\hat{f}_u(W_u) - f_u(W_u)}{u+1} < +\infty. \quad a.s.$$

Let  $a_u = \hat{f}_u(W_u) - f_u(W_u)$  and  $b_u = \frac{1}{u+1}$ , we have:

$$\begin{aligned} |a_{u+1} - a_u| &\leq (\hat{f}_{u+1}(W_{u+1}) - \hat{f}_u(W_u)) + \\ &\quad |f_{u+1}(W_{u+1}) - f_u(W_u)|. \end{aligned}$$

According to Proposition 35 and Lemma 31, the two items on the right side are  $O(\frac{1}{u})$ , and so  $|a_{u+1} - a_u| = O(\frac{1}{u})$ . Applying Lemma A3, we obtain

$$\lim_{u \rightarrow \infty} a_u = \lim_{u \rightarrow \infty} \hat{f}_u(W_u) - \lim_{u \rightarrow \infty} f_u(W_u) = 0 \quad a.s.$$

This proves that, with unlimited data, the empirical cost with active online learning converges to the objective cost almost surely.

*Proposition 37:* Given assumptions (A)-(D), if the multiplicative update converges for each signal piece, the empirical cost converges to the object cost almost surely with inertial online learning, i.e.,

$$\lim_{u \rightarrow \infty} \hat{f}_u(W_u) = \lim_{u \rightarrow \infty} f_u(W_u) \quad a.s.$$

*Proof:*

A similar approach to that applied for active learning can be applied to the proof for inertial learning. As has been shown in Lemma 35, we have:

$$\begin{aligned} v_{u+1} - v_u &= \hat{f}_{u+1}(W_{u+1}) - \hat{f}_u(W_u) \\ &\leq \hat{f}_{u+1}(W_u) - \hat{f}_u(W_u) \\ &= \frac{\ell_{X_{u+1}}(W_u) - f_u(u)}{u+1} \\ &\quad + \frac{f_u(W_u) - \hat{f}_u(W_u)}{u+1} \end{aligned}$$

We then have the same form of variation sequence as in active learning except an  $O(\frac{1}{u^2})$  term has been omitted. Following the same process, we can show that:

$$\mathbb{E}[\mathbb{E}[v_{u+1} - v_u | \mathcal{F}_u]^+] = O\left(\frac{1}{u^2}\right)$$

and  $\hat{f}_u(W_u)$  converges to the objective cost  $f_u(W_u)$  almost surely.

It can be further proved that the objective cost function  $f_u$  converges to the expected objective cost function, i.e.,

$$\|f_u - f\|_\infty \rightarrow_{u \rightarrow \infty} 0$$

and therefore

$$\lim_{u \rightarrow \infty} \hat{f}_u(W_u) - \lim_{u \rightarrow \infty} f(W_u) = 0 \quad a.s.$$

This result shows that with unlimited data, the empirical cost  $\hat{f}(W_u)$  converges to the expected cost  $f(W_u)$  almost surely. Since  $W_u$  is a stationary point of  $\hat{f}$ , it can be proved that  $W_u$  converges to a stationary point of the expected cost function  $f$ . The proof for this stronger result is similar to that in ODL [48].

### C. Convergence of imperfect learning

Here we prove the convergence of the imperfect empirical cost function  $\hat{f}_u(W_u)$ . This corresponds to ‘imperfect learning’ where the multiplicative update does not converge for each signal piece. This is generally the case in practice due to limited computing resources.

First define  $\delta_u(W)$  as the bias of the imperfect empirical cost shifted away from the perfect empirical cost of the  $u^{th}$  piece of signal, i.e.,

$$\delta_u(W) = \ell_{X_u}(W, H_u) - \ell_{X_u}(W, \hat{H}_u).$$

The following proposition states that online learning converges with imperfect coefficients.

*Proposition 38:* Suppose online learning does not converge for each signal piece due to limited computational resources, and that the bias  $\delta_u(W)$  has the same expectation in spite of  $u$ , i.e.,

$$\mathbb{E}_{X,H}[\delta_u(W)] = \epsilon(W) \quad \forall u \quad (11)$$

where the expectation is taken on both signals  $X$  and coefficients  $H$ . Given assumptions (A)-(D), online learning converges almost surely and that:

$$\lim_{u \rightarrow \infty} \hat{f}_u(W_u) = \lim_{u \rightarrow \infty} f_u(W_u) + \epsilon(W_u) \quad a.s.$$

*Proof:*

This proof is similar to the convergence proof in Proposition 36 and 37. For simplicity, we prove the condition with active learning; it is straight forward to verify the application of the same proof to inertial learning with minor changes.

As in the proof of Proposition 36, let us define  $v_u \equiv \dot{f}_u(W_u)$ ; we have:

$$\begin{aligned}
v_{u+1} - v_u &= \dot{f}_{u+1}(W_{u+1}) - \dot{f}_u(W_u) \\
&\leq \dot{f}_{u+1}(W_u) - \dot{f}_u(W_u) \\
&= \frac{\ell_{X_{u+1}}(W_u) + \delta_{u+1}(W_u) - f_u(W_u)}{u+1} \\
&\quad + \frac{f_u(W_u) - (\hat{f}_u(W_u) + \frac{1}{u} \sum_{i=1}^u \delta_i(W_u))}{u+1} \\
&= \frac{\ell_{X_{u+1}}(W_u) - f_u(W_u)}{u+1} \\
&\quad + \frac{f_u(W_u) - \hat{f}_u(W_u)}{u+1} \\
&\quad + \frac{\delta_{u+1}(W_u) - \frac{1}{u} \sum_{i=1}^u \delta_i(W_u)}{u+1}
\end{aligned}$$

This leads to:

$$\begin{aligned}
\mathbb{E}_{X,H}[v_{u+1} - v_u | \mathcal{F}_u] &\leq \frac{\mathbb{E}[\ell_{X_{u+1}} | \mathcal{F}_u] - f_u(W_u)}{u+1} \\
&\quad + \frac{\epsilon(W_u) - \frac{1}{u} \sum_{i=1}^u \delta_i(W_u)}{u+1} \\
&\leq \frac{\|f - f_u\|_\infty}{u+1} \\
&\quad + \frac{\epsilon(W_u) - \frac{1}{u} \sum_{i=1}^u \delta_i(W_u)}{u+1}
\end{aligned}$$

Then we have:

$$\mathbb{E}_{X,H}[\mathbb{E}_{X,H}[v_{u+1} - v_u | \mathcal{F}_u]] = \frac{\mathbb{E}_{X,H}[\|f - f_u\|_\infty]}{u+1}$$

where we have applied:

$$\epsilon(W_u) - \frac{1}{u} \sum_{i=1}^u \mathbb{E}_{X,H}[\delta_i(W_u)] = 0.$$

This results in the same convergence speed as in Proposition 36, i.e.,

$$\mathbb{E}_{X,H}[\mathbb{E}_{X,H}[v_{u+1} - v_u | \mathcal{F}_u]^+] = O\left(\frac{1}{u^{\frac{3}{2}}}\right). \tag{12}$$

The convergence of  $\dot{f}_u(W_u)$  can be proved following the same process in the proof of Proposition 36. Furthermore, noting that:

$$\dot{f}_u(W_u) = \hat{f}_u(W_u) + \frac{\sum_{i=1}^u \delta_i(W_u)}{u}$$

and  $\hat{f}_u(W_u)$  converges to  $f_u(W_u)$ , we have:

$$\begin{aligned}\lim_{u \rightarrow \infty} \dot{f}_u(W_u) &= \lim_{u \rightarrow \infty} f_u(W_u) + \lim_{u \rightarrow \infty} \frac{\sum_{i=1}^u \delta_i(W_u)}{u} \quad a.s. \\ &= \lim_{u \rightarrow \infty} f_u(W_u) + \epsilon(W_u) \quad a.s.\end{aligned}$$

Theoretically, the perfect empirical cost  $\dot{f}_u(W_u)$  does not necessarily converge, but if we assume that the update on  $W_u$  in each step improves  $\dot{f}_u(W_u)$ , then it can be verified that  $\dot{f}_u(W_u)$  converges indeed and

$$\lim_{u \rightarrow \infty} \dot{f}_u(W_u) = \lim_{u \rightarrow \infty} f(W_u) + \epsilon(W_u) \quad a.s. \quad (13)$$

This means that online learning converges to a stationary point of a new function  $g(W) = f(W) + \epsilon(W)$ , which is obviously suboptimal for our task which intends to minimize the expected cost  $f(W)$ . If the multiplicative update approaches to convergence,  $\epsilon(W)$  approaches to 0 and the learning process converges to  $f(W)$ , which is just the form of a perfect learning.

Note that the above analysis relies on strong assumptions. First the identical expected bias assumption (11) is not always respected, although some support can be found from the random coefficient initialization and identical number of multiplicative iterations. Second, the assumption of improvement on  $\dot{f}_u(W_u)$  can be simply false if the coefficients of each piece of signal are highly imperfect. This suggests that the convergence property with imperfect learning is not guaranteed in practice, and may highly task-dependent.

Some interesting results can be obtained from the empirical convergence proposition. First notice that under the proposed assumptions, both active learning and inertial learning converge to local minima of cost functions in the same form  $g(w) = f(W) + \epsilon(W)$ , which means the two learning may obtain similar empirical cost. Second, since  $\ell$  is an utterance-based average cost, the convergence behavior is not impacted by the length of signal pieces, and involving pieces of variable lengths does not impact the convergence property.

#### D. Batch learning and online learning

The convergence analysis in the previous section shows that both batch and online learning converge to a stationary point of the expected objective function  $f(W)$  with unlimited data and unlimited computing resources. This situation is only valid in theory. For small scale tasks where data are limited but computing resource is abundant, batch learning converges to a stationary point of the objective cost function  $f_u(W)$  while online learning fails to converge, resulting in suboptimal patterns. For large scale tasks, the more common situation is that the training data is in a large amount while the computing resource affords only imperfect learning. In this situation, the online learning tends to obtain lower empirical cost than the batch learning. This is stated as the following proposition:

*Proposition 39:* If the learning is imperfect, for any batch learning, there is an online learning that attains lower empirical cost if the data are unlimited.

*Proof:*

First notice that the statistics  $\frac{A(q,p;u)}{u}$  and  $\frac{B(p;u)}{u}$  converge with  $u$  being large, due to the stationary distribution assumption according to the central limit theorem. This means that, batch learning with limited iterations converges to an empirical cost  $\dot{f}_B$  when the data volume approaches to infinity. This can be written as

$$|\dot{f}_B(u) - \dot{f}_B| < \sigma \quad \forall u > N \quad (14)$$

where  $\dot{f}_B(u)$  is the empirical cost with batch learning and  $u$  pieces of signals, and  $\sigma$  is a small value.

We then split the training data into blocks where each block consists of  $N$  pieces of data. Active online learning can be started by treating each block as a learning piece. Proposition 38 states that the empirical cost converges to a level  $\dot{f}$ . This leads to a relation as follows:

$$\dot{f}_B(N) = \dot{f}_N(W_N) > \dot{f}$$

where  $\dot{f}$  is the converged empirical cost with online learning. Combined with (14), we have:

$$\dot{f} < \dot{f}_B + \sigma.$$

This means this online learning always obtains lower cost than the batch learning if  $\sigma$  is small enough, which is guaranteed by a large  $N$ .

This result indicates that online learning not only enables pattern learning with large databases; moreover, it learns patterns of better quality compared to batch learning, even with less computation (in the case of inertial learning). However, we must emphasize that the assumption behind this statement is that computing resources are limited and training data are unlimited. If computing resources are unlimited and data are limited, batch learning is ensured to converge to a local optimum, while online learning usually arrives at suboptimal solutions. If both computing resources and data are unlimited, then the theoretical convergence analysis presented Section III-B applies, i.e., online and batch learning converge to the same empirical cost. Finally, in small scale tasks where both the computing resources and data are limited, there is no guarantee which learning approach is superior; it will depend on which factor (computing resource or data) is more dominant. ■

#### IV. EXPERIMENTS

We present two experiments which demonstrate the characteristics and benefits of the proposed online CNSC approach. The first experiment is a small-scale speech separation task which aims to compare the behavior of the two online learning approaches with batch learning; the second experiment involves a noise cancellation task for large-scale speech recognition and was defined by the CHiME challenge<sup>4</sup>. It aims to demonstrate the power of online learning in real applications.

##### A. Speech separation

In this experiment, we study the behavior of the proposed online pattern learning algorithm using a toy experiment proposed by P. Smaragdis<sup>5</sup> in a study of CNMF [32]. The task is to learn two sets of patterns from individual speech signals of a male and female speaker respectively, and then to use the corresponding patterns to separate the two voices from a segment of mixed speech. Smaragdis showed that speaker specific patterns can be learned using CNMF and then employed to separate the speech signal according to the constituent speakers. It has also been shown in [42] that sparse coding can deliver improved performance in a similar signal separation task.

The two individual speech segments used for pattern learning are in the order of 30 seconds in length, are sampled at 16kHz and are mixed together by simple addition with appropriate zero-padding being applied to the shorter speech recording. Signals are windowed into frames of 32ms with a frame shift of 16ms, thereby resulting in a frame rate of 62.5 frames per second. The discrete Fourier transform is applied to each frame and the magnitude spectrum is used as a non-negative representation which is suitable for processing with NMF and CNSC. All experiments reported here are based on fixed parameters of  $R = 20$ ,  $P = 4$  and  $\lambda = 0.01$  which are all chosen heuristically. Finally, all experiments were conducted on a desktop machine with two dual-core 2.60GHz CPUs and memory of 4GB.

Before presenting the separation task, we investigate several factors which impact on the convergence behavior of online learning. The recording of male speech is used to conduct pattern learning and signal reconstruction; learning quality is measured using the value of the cost function (4). Fig. 1 illustrates example patterns learned for the male speaker.

1) *Convergence and computing resource*: The first factor which impacts on convergence is the number of multiplicative update iterations, which is directly related to computing resources. The male speech signals are divided into 10 pieces to conduct online learning, and batch learning is implemented as an active online learning with the number of pieces set to 1.

Fig. 2 presents the cost values obtained with various learning approaches for the first 100 iterations. We observe that active online learning converges in the first 5 iterations while inertial online learning requires 10 iterations to converge. An interesting observation is that the cost obtained with the two online learning approaches increases with a higher number of iterations, indicating some over-fitting to the first few signal pieces. Upon comparison of the two online learning approaches, we see that active learning leads to lower cost, which is expected considering the more aggressive early learning.

Batch learning converges much more slowly than online learning: it requires 15 and 30 iterations to reach the same cost obtained with inertial and active learning respectively, and requires more than 80 iterations to converge

<sup>4</sup><http://spandh.dcs.shef.ac.uk/projects/chime/challenge.html>

<sup>5</sup><http://www.cs.illinois.edu/~paris/demos/>

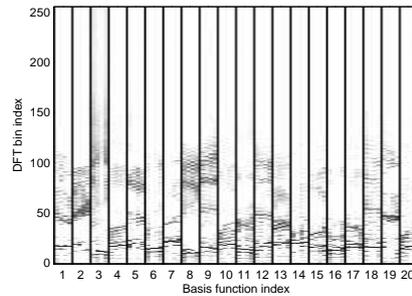


Fig. 1: An example of patterns learned with active learning.

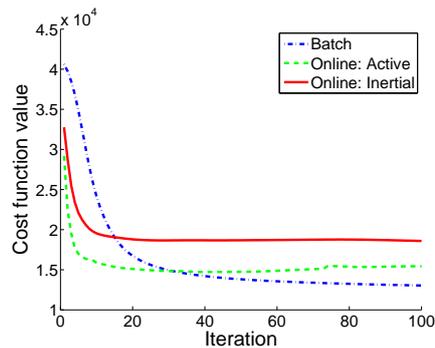


Fig. 2: Value of the cost function for the first 100 iterations with online and batch learning.

itself. In spite of slow convergence, batch learning delivers lower cost if the number of iterations is sufficiently large, thereby demonstrating its advantage in small scale tasks. These observations are consistent with the convergence analysis presented in Section III.

Fig. 3 shows the corresponding average run-time for the first 100 iterations of the three learning approaches. We can see that the inertial online learning is the most efficient while the active learning is the most expensive. This observation is consistent with the computational complexity analysis presented in Section III-A.

2) *Convergence and piece length*: The second factor which impacts on the convergence of online learning is the manner in which signals are split into pieces. Generally speaking, the splitting of signals into more pieces leads to more frequent pattern updates and hence more aggressive early learning; we therefore expect lower cost with smaller pieces for online learning, subject to the assumption of independence among pieces being held.

To test this conjecture, the signals of the male speaker are split into  $U$  pieces, and then the patterns are learned by setting the number of multiplicative iterations to 10. The cost of the two online learning approaches is shown in Fig. 4 for  $U = 1$  to 10. Note that active learning with  $U = 1$  is equivalent to batch learning. As expected, we observe that the two online learning approaches obtain substantially lower cost than batch learning ( $U = 1$ ) and that smaller pieces lead to lower cost. Note that, with increasing number of pieces, the cost function exhibits some variation. This can be attributed to the boundary effect stemming from signal segmentation.

The corresponding average run-time is shown in Fig. 5. We first observe that active learning requires more computational resources as the training data are split into more pieces, due to the increased number of pattern updates. Inertial learning exhibits different behavior: the computational demand first decreases when the data are split into a small number of pieces; with increasing number of pieces, the computational requirements increase steadily by a small factor. This is because the initial pattern update for the first piece (ref. Algorithm 1) is less costly when the data are split into smaller pieces. As the number of pieces increases, the computational saving with the initial pattern update becomes marginal while the cost associated with pattern update for each piece increases.

3) *Convergence and data volume*: In the third experiment, we study the impact of the amount of training data, for which the male speech signals are duplicated and concatenated to simulate increasing data volume. This simulation

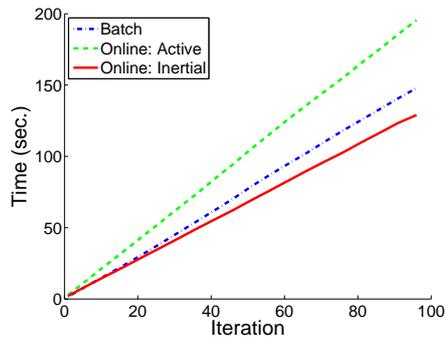


Fig. 3: Average run-time for the first 100 iterations with online and batch learning.

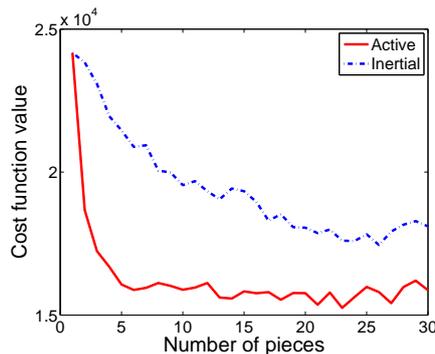


Fig. 4: Value of the cost function for  $U = 1$  to 10 pieces and after 10 iterations for active and inertial online learning.

certainly can not fully represent practical scenarios with large amount of data, however it does approximate a stationary compact distribution.

We first study the case of perfect learning. From Fig. 2, we see that with 100 iterations both online and batch learning can be regarded as converged. We therefore set up an experiment where the iteration number is fixed to 100 and the amount of training data is increased by data duplication. Results are shown in Fig. 6 where the x-axis denotes the number of duplications and the y-axis is the cost. We see that, with limited data, batch learning obtains significantly lower cost than the two online learning algorithms; with more and more data, however, the cost obtained with online learning approaches that obtained with batch learning. Although not reported in the figure, when the duplication number increases to over 500, the three learning approaches obtain very similar cost, thus demonstrating the convergence theory presented in Section III. Note that for batch learning, the cost profile exhibits some fluctuation which can be attributed to the boundary effect between duplications.

In another experiment we study the case of imperfect learning. To simulate this situation, the number of multiplicative updates is set to 10 and the training data are again increasingly duplicated. Results are shown in Fig. 7. We observe that the online learning obtains significantly lower cost than batch learning and that the two online learning approaches converge to the same cost. When compared to the results in the case of perfect learning (Fig. 6), we find that cost obtained with perfect and imperfect learning is comparable when empirical convergence is reached. This suggests that a few iterations might be sufficient for online learning on large scale tasks, as we will see in the noise cancelation experiment presented in Section IV-B.

4) *Speech separation*: In the speech separation task, both the male and female speech utterances are split into 30 pieces which has been shown effective for online learning. The male and female patterns are learned using corresponding training speech by applying either the online and or the batch learning approach. The spectrum of the mixed speech signal is then projected independently onto the two sets of patterns and the reconstruction cost of the resulting magnitude spectrum is computed for the individual male and female speech signals respectively<sup>6</sup>.

<sup>6</sup>The original and reconstructed speech waveforms in this experiment are available at <http://audio.eurecom.fr/software/ol>

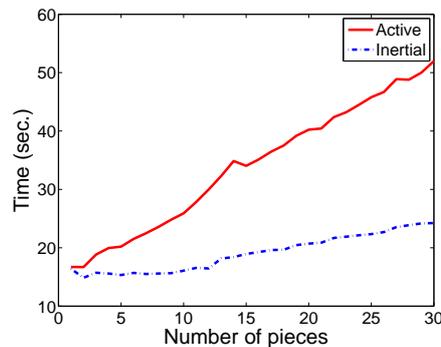


Fig. 5: Average run-time for between  $U = 1$  to 10 pieces and after 10 iterations for active and inertial online learning.

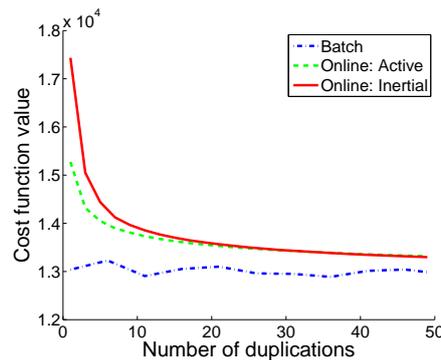


Fig. 6: Value of the cost function with the multiplicative update fixed to 100 iterations and the speech data duplicated up to 50 times.

The separation performance is evaluated in terms of the signal to distortion ratio (SDR), defined as follows:

$$SDR = \frac{1}{2} \sum_{i \in \{male, female\}} 10 \log_{10} \frac{\|s_{dist}^i\|^2}{\|e_{interf}^i + e_{noise}^i + e_{artif}^i\|^2}$$

where  $i$  denotes channels (female or male),  $s_{dist}$  is the original speech signal,  $e_{interf}, e_{noise}, e_{artif}$  denote interference among channels, noise and artifacts introduced by separation [54]. The BSS Eval tool was used to conduct the evaluation<sup>7</sup>. Results are shown in Fig. 8 where the x-axis represents the number of multiplicative iterations and the y-axis represents SDR. It can be observed that, with a small number of iterations, the two online learning algorithms result in better separation than batch learning. With an increasing number of iterations, inertial learning converges to an approximate SDR of 4.5 while batch and active learning give an approximate SDR of 6.0, though batch learning ultimately outperforms active learning. Note that this does not mean online learning is less valuable: with a large amount of data, batch learning simply runs out of memory whereas online learning requires far less resources. In this case, online learning is thus the only viable choice.

### B. Denoising for speech recognition

In the second experiment we apply the online learning approach to suppress multi-source noise from speech signals for improved automatic speech recognition (ASR). The basic idea is to learn the patterns of clean speech and background noise. Clean speech representations are then obtained by distributing the signal energy among the speech and noise patterns and by discarding that attributed to noise. The procedure is described in [32].

<sup>7</sup>[http://bass-db.gforge.inria.fr/bss\\_eval](http://bass-db.gforge.inria.fr/bss_eval).

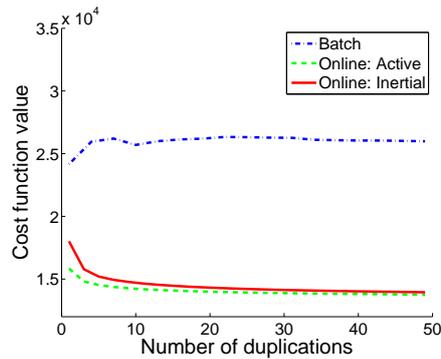


Fig. 7: Value of the cost function with the multiplicative update fixed to 10 iterations and the speech data duplicated 50 times.

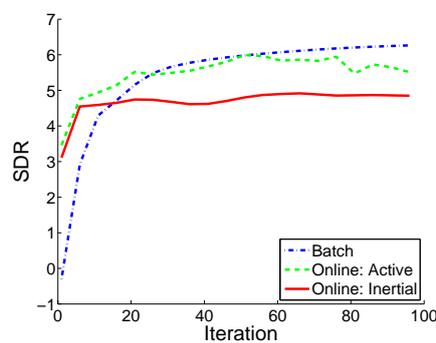


Fig. 8: SDR of speech separation.

1) *Experimental setup*: Our experiments are set up within the framework of the CHiME challenge [55], where the task is to recognize the speech utterances in a home environment with various kinds of background noise under six different signal-to-noise ratio (SNR) conditions. The background noise comprises voices, television sounds, music, noise from home appliances and a host of other ambient noises typically observed in a home environment. All audio signal were recorded with a binaural microphone array. The location of the target speaker is specified to be 2 meters directly in-front of the microphone array, while the type of noise sources and their locations are unknown and variable.

The database contains recordings from 34 speakers and a set of 84 recordings of ambient noise, each of which is 5 minutes in duration. The training set comprises 500 utterances per speaker amounting to approximately 15.3 minutes of audio per speaker. The test set comprises 600 utterances under each of the following SNR conditions:  $-6\text{dB}$ ,  $-3\text{dB}$ ,  $0\text{dB}$ ,  $3\text{dB}$ ,  $6\text{dB}$  and  $9\text{dB}$ . All utterances follow a simple grammar that involves digits and letters; only the hypotheses for the letter and digit are scored to evaluate recognition performance.

As the first step, the two channels were mixed with zero delay to obtain a mono-channel audio signal for further processing. We used the standard ASR setup provided for the CHiME challenge. The setup uses speaker dependent acoustic models trained on Mel frequency cepstral coefficients (MFCC) with energy plus the first and second order derivatives. Cepstral mean normalization (CMN) is applied to improve robustness to additive noise. The language model is a simple lattice that covers all possible sequences in the grammar mentioned above and the utterances are decoded using HTK [56].

Spectral representations are extracted using a window size of 25ms and an overlap of 10ms. Speaker patterns were learnt from the training set available for each speaker using a convolutional span of 4 frames. This is equivalent to capturing prominent spectro-temporal patterns that span about 70ms, i.e. subphone patterns. For each speaker, a pattern matrix is learnt using batch learning [57]. Its dimension was empirically set to 100.

The learning of noise patterns is more complex. Since the noise is highly diverse and variable, we would ideally

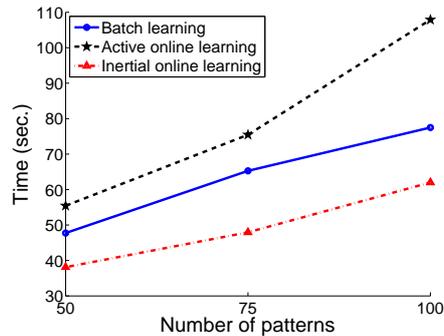


Fig. 9: Average run-time for 10 iterations with online and batch learning.

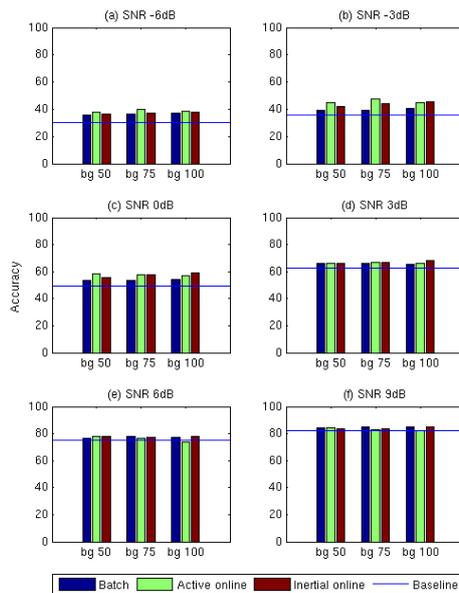


Fig. 10: ASR accuracies with CNSC-based noise cancellation where the background noise patterns are learnt on a set of randomly sampled background audio segments with batch, active online and inertial online learning.

like to learn as many patterns as possible from the 7 hours of noise recordings provided in the development set. However, the memory and computational requirements to store and process such large amounts of data are very demanding and the classical batch learning approach simply fails. This problem can be avoided through incomplete training [43], [44], or through online learning as proposed in this work.

2) *Incomplete noise pattern learning*: In this experiment, we choose a subset of the background training data to learn the noise patterns. On the one hand this avoids the prohibitive computing and memory demands with batch learning and, on the other hand, provides an opportunity to compare batch learning and online learning in real applications. In our experiment, 10 seconds of audio segments from each of the 5 minute waveforms were randomly sampled to obtain 840 seconds of background noise.

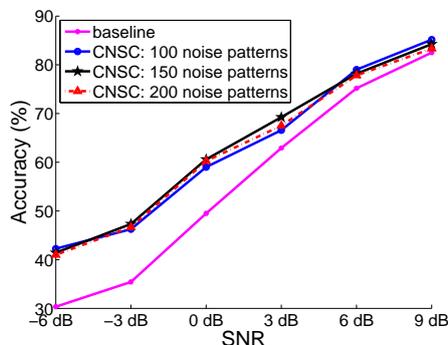
With this data, patterns of size 50, 75 and 100 were learnt using the batch, active online and inertial online learning approaches. Each of the 10 seconds of audio segments acts as a piece in online learning. Similar to the experiments in Section IV-A, a single dual core 2.6GHz processor with 4GB memory was used.

The time taken for pattern learning with 10 iterations using the three learning methods is shown in Fig. 9. As in the experiments of Section IV-A, the values presented in the figure are averaged over 100 runs to avoid computational fluctuations. Results confirm that active online learning takes longer than batch learning while inertial online learning outperforms batch mode in terms of computational time, as discussed in Section III-A. ASR performance on the evaluation data is shown in Fig. 10. We observe considerable improvement in accuracies for low SNR conditions

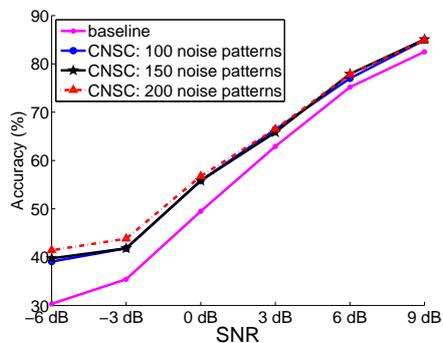
and a marginal improvement over the baseline for high SNR conditions with all the three learning approaches. An interesting trend that can be observed in Fig. 10 is that, at lower SNR conditions, active online learning outperforms the other two approaches, while at high SNR conditions and in particular, with higher number of noise patterns, active learning tends to give the worst performance. This might be explained by the quick convergence to local (partial) patterns with active learning as a result of which the speech energy may be incorrectly attributed to the noise patterns. This problem is more severe in high SNR conditions where the noise is low while the number of noise patterns is large. This leads to the incorrect attribution of speech energy.

3) *Complete noise pattern learning*: Methods employing random sampling techniques are not desirable in practice since a large proportion of the training data remains unused. In order to learn noise patterns from the entire training data, we apply the proposed online pattern learning algorithm. The background training data are provided in segments of 5 minutes each. We use the same partition structure as pieces in the online training algorithm. We learn background bases with 100, 150 and 200 dimensions respectively, all with a convolutional span of 4 frames.

ASR accuracies on the evaluation data with noise cancellation using background patterns learnt with active and inertial online learning algorithms are presented in Fig. 11. We observe that pattern-based denoising significantly improves ASR performance with both active and inertial online learning. Again, active learning is more effective in low SNR conditions than inertial learning; it however does not show much advantage at high SNRs. Simply increasing the number of patterns does not result in significant gains.



(a) Active Online learning



(b) Inertial Online learning

Fig. 11: Accuracies with CNSC-based noise cancellation where the background noise patterns are learnt using active and inertial online CNSC.

Fig. 12 presents a comparison between incomplete and complete learning. We see that the patterns learned from the entire background noise data provide improved accuracies over those learned with random sampling in the case of active learning; for inertial learning, the advantage of using the entire data is not evident, indicating that the training data does not fit a stationary distribution and thus slow inertial learning cannot reach empirical convergence with the use of additional data. Nevertheless, these results clearly demonstrate the capability of online learning in real, large-scale applications.

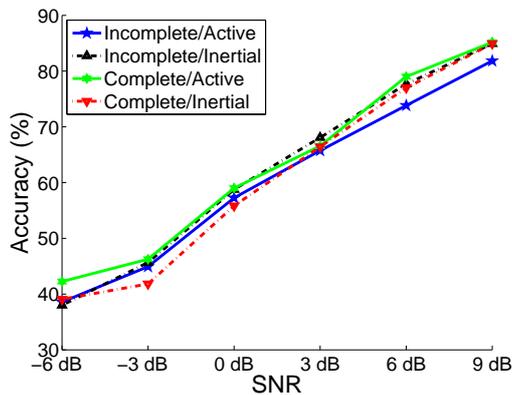


Fig. 12: Accuracy with online CNSC-based noise cancellation where the 100 background noise patterns are learnt based on random or entire data respectively. Results using both active and inertial learning are presented.

## V. CONCLUSION

This paper presents a new online CNSC algorithm to learn convolutive non-negative patterns with sparse coding. Compared to conventional batch learning, the proposed approach is able to learn complex patterns from large volumes of training data and is thus suited to large-scale applications. The theoretical analysis shows that the online algorithm almost surely converges to a stationary point of the objective cost function with unlimited computational resources and training data. In real applications where the computational resources are limited and the training data volume is large, the online approach tends to gain lower empirical cost than the batch learning. This analysis is confirmed by the results we obtained with a toy experiment on speech separation. A noise cancellation task for a large-scale speech recognition system we constructed within the CHiME challenge framework demonstrates that the online learning approach is efficient in learning complex patterns from large corpora in real applications.

Future work includes the study of incremental patterns with online learning; another direction is to extend the online CNSC approach to other unsupervised learning techniques such as sparse PCA.

## APPENDIX

In this section, we cite some theorems that are used for the convergence proof in this paper. These theorem are mostly reproduced from [48] for convenience of readers.

*Theorem A1:* [Sufficient condition of convergence for a stochastic process. See [52], [58], [59]]

Let  $(Q, F, P)$  be a measurable probability space,  $u_t$ , for  $t \geq 0$ , be the realization of a stochastic process and  $F_t$  be the filtration determined by the past information at time  $t$ . Let

$$\delta_t = \begin{cases} 1 & \text{if } \mathbb{E}[u_{t+1} - u_t | F_t] > 0, \\ 0 & \text{otherwise.} \end{cases}$$

If for all  $t$ ,  $u_t \geq 0$  and  $\sum_{t=1}^{\infty} \mathbb{E}[\delta_t(u_{t+1} - u_t)] < \infty$ , then  $u_t$  is a quasi-martingale and converges almost surely. Moreover,

$$\sum_{t=1}^{\infty} |\mathbb{E}[u_{t+1} - u_t | F_t]| < +\infty \quad a.s.$$

*Lemma A2:* [A corollary of Donsker theorem for  $O(\frac{1}{\sqrt{n}})$  of  $|f_n - f|$ . See [53], chap. 19.]

Let  $F = \{f_\theta : \chi \rightarrow \mathbb{R}, \theta \in \Theta\}$  be a set of measurable functions indexed by a bounded subset  $\Theta$  of  $\mathbb{R}^d$ . Suppose that there exists a constant  $K$  such that

$$|f_{\theta_1} - f_{\theta_2}| \leq K \|\theta_1 - \theta_2\|_2$$

for every  $\theta_1$  and  $\theta_2$  in  $\Theta$  and  $x$  in  $\chi$ . Then  $F$  is P-Donsker. For any  $f$  in  $F$ , define  $\mathbb{P}_n f$ ,  $\mathbb{P}f$  and  $\mathbb{G}_n f$  as

$$\mathbb{P}_n f = \frac{1}{n} \sum_{i=1}^n f(X_i), \mathbb{P}_f = \mathbb{E}_X[f(X)], \mathbb{G}_n f = \sqrt{n}(\mathbb{P}_n f - \mathbb{P}_f).$$

Further suppose for all  $f$ ,  $\mathbb{P}f^2 < \delta^2$  and  $\|f\|_\infty < M$  and that the random elements  $X_i$  are Borel-measurable. Then we have

$$\mathbb{E}_P \|\mathbb{G}_n\|_F = O(1),$$

where  $\|\mathbb{G}_n\|_F = \sup_{f \in F} |\mathbb{G}_n f|$ .

**Lemma A3:** [A lemma on positive converging sums. See [60], prop 1.2.4.]

Let  $a_n, b_n$  be two real sequences such that for all  $n$ ,  $a_n \geq 0$  and  $b_n \geq 0$ ,  $\sum_{n=1}^\infty a_n = \infty$ ,  $\sum_{n=1}^\infty a_n b_n < \infty$ ,  $\exists K > 0$  s.t.  $|b_{n+1} - b_n| < K a_n$ . Then  $\lim_{n \rightarrow +\infty} b_n = 0$ .

#### ACKNOWLEDGMENT

This work was conducted when Dong Wang was at EURECOM as a post-doctoral research fellow and was completed when he was a visiting researcher in Tsinghua University and a senior research engineer in Nuance. It was partially supported by the French Ministry of Industry (Innovative Web call) under contract 09.2.93.0966, collaborative Annotation for Video Accessibility (ACAV) and by the Adaptable Ambient Living Assistant (ALIAS) project funded through the joint national Ambient Assisted Living (AAL) programme.

#### REFERENCES

- [1] P. Paatero and U. Tapper, "Positive matrix factorization: A non-negative factor model with optimal utilization of error estimates of data values," *Environmetrics*, vol. 5, no. 12, pp. 111–126, 1994.
- [2] P. Paatero, "Least squares formulation of robust non-negative factor analysis," *Chemometrics and Intelligent Laboratory Systems*, vol. 37, pp. 23–35, 1997.
- [3] D. D. Lee and H. Seung, "Learning the parts of objects by non-negative matrix factorization," *Nature*, vol. 401, no. 6755, pp. 788–791, October 1999.
- [4] —, "Algorithms for non-negative matrix factorization," in *Proc. NIPS*, 2000, pp. 556–562.
- [5] D. L. Donoho and V. Stodden, "When does non-negative matrix factorization give a correct decomposition into parts?" in *Proc. NIPS 2003*, 2003.
- [6] C. Lin, "On the convergence of multiplicative update algorithms for nonnegative matrix factorization," *IEEE Transactions on Neural Networks*, vol. 18, no. 6, pp. 1589–1596, Nov. 2007.
- [7] A. Shashua and T. Hazan, "Non-negative tensor factorization with applications to statistics and computer vision," in *Proc. International Conference on Machine Learning (ICML)*, 2005, pp. 792–799.
- [8] S. Z. Li, X. Hou, H. Zhang, and Q. Cheng, "Learning spatially localized, parts-based representation." in *Proc. CVPR'01*, 2001, pp. 207–212.
- [9] P. Smaragdis and J. C. Brown, "Non-negative matrix factorization for polyphonic music transcription," in *Proc. IEEE Workshop on Applications of Signal Processing to Audio and Acoustics*, 2003, pp. 177–180.
- [10] D. Guillaumet, J. Vitrià, and B. Schiele, "Introducing a weighted non-negative matrix factorization for image classification," *Pattern Recognition Letters*, vol. 24, no. 14, pp. 2447–2454, 2003.
- [11] M. Rajapakse, J. Tan, and J. C. Rajapakse, "Color channel encoding with NMF for face recognition." in *Proc. ICIP'04*, 2004, pp. 2007–2010.
- [12] J.-P. Brunet, P. Tamayo, T. R. Golub, and J. P. Mesirov, "Metagenes and molecular pattern discovery using matrix factorization," *PNAS*, vol. 101, no. 12, pp. 4164–4169, 2004.
- [13] B. Wang and M. D. Plumbley, "Investigating single-channel audio source separation methods based on non-negative matrix factorization," in *Proc. ICA Research Network International Workshop*, 2006, pp. 17–20.
- [14] G. Wang, A. V. Kossenkov, and M. F. Ochs, "LS-NMF: A modified non-negative matrix factorization algorithm utilizing uncertainty estimates," *BMC Bioinformatics*, vol. 7, p. 175, 2006.
- [15] M. W. Berry, M. Browne, A. N. Langville, V. P. Pauca, and R. J. Plemmons, "Algorithms and applications for approximate non-negative matrix factorization," *Computational Statistics and Data Analysis*, vol. 52, no. 1, pp. 155–173, 2006.
- [16] V. P. Pauca, J. Piper, and R. J. Plemmons, "Nonnegative matrix factorization for spectral data analysis," *Linear Algebra and its Applications*, vol. 416, no. 1, pp. 29–47, 2006.
- [17] W. Xu, X. Liu, and Y. Gong, "Document clustering based on non-negative matrix factorization." in *Proc. SIGIR*, 2003, pp. 267–273.
- [18] K. Devarajan and B. Bryant, "Nonnegative matrix factorization: An analytical and interpretive tool in computational biology," *PLoS Computational Biology*, vol. 4, 2008.
- [19] N. Bertin, R. Badeau, and E. Vincent, "Enforcing harmonicity and smoothness in Bayesian non-negative matrix factorization applied to polyphonic music transcription," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 18, pp. 538–549, March 2010.
- [20] M. Welling and M. Weber, "Positive tensor factorization," *Pattern Recognition Letters*, vol. 22, no. 12, pp. 1255–1261, 2001.

- [21] I. S. Dhillon and S. Sra, "Generalized non-negative matrix approximations with Bregman divergences," in *Proc. Neural Information Proc. Systems*, 2005, pp. 283–290.
- [22] R. Zdunek and A. Cichocki, "Non-negative matrix factorization with quasi-newton optimization," in *Proc. Eighth International Conference on Artificial Intelligence and Soft Computing, ICAISC*. Springer, 2006, pp. 870–879.
- [23] A. Cichocki, R. Zdunek, and S. ichi Amari, "Csiszár's divergences for non-negative matrix factorization: Family of new algorithms," in *Lecture Notes in Computer Science (LNCS)*. Springer, 2006, pp. 32–39.
- [24] C. Lin, "Projected gradient methods for non-negative matrix factorization," *Neural Computation*, vol. 19, no. 10, pp. 2756–2779, 2004.
- [25] R. Zdunek and A. Cichocki, "Nonnegative matrix factorization with constrained second-order optimization," *Signal Processing*, vol. 87, no. 8, pp. 1904–1916, 2007.
- [26] C. Févotte, N. Bertin, and J.-L. Durrieu, "Nonnegative matrix factorization with the Itakura-Saito divergence: With application to music analysis," *Neural Computation*, pp. 793–830, 2009.
- [27] R. Kompass, "A generalized divergence measure for non-negative matrix factorization," *Neural Computation*, vol. 19, no. 3, pp. 780–791, 2007.
- [28] H. Kim and H. Park, "Nonnegative matrix factorization based on alternating nonnegativity constrained least squares and active set method," *SIAM J. Matrix Anal. Appl.*, vol. 30, pp. 713–730, July 2008.
- [29] A. T. Cemgil, "Bayesian inference for non-negative matrix factorisation models," *Intell. Neuroscience*, vol. 2009, pp. 4:1–4:17, January 2009.
- [30] M. N. Schmidt, O. Winther, and L. K. Hansen, "Bayesian non-negative matrix factorization," in *Independent Component Analysis and Signal Separation, International Conference on*, ser. Lecture Notes in Computer Science (LNCS), vol. 5441. Springer, 2009, pp. 540–547.
- [31] Z. Yang and E. Oja, "Linear and non-linear projective non-negative matrix factorization," *IEEE Transactions on Neural Network*, vol. 21, no. 5, pp. 734–749, 2010.
- [32] P. Smaragdis, "Convolutional speech bases and their application to supervised speech separation," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 1, pp. 1–12, January 2007.
- [33] D. FitzGerald and E. Coyle, "Shifted non-negative matrix factorisation for sound source separation," in *Proc. IEEE conference on Statistics in Signal Processing*, 2005.
- [34] P. Hoyer, "Non-negative matrix factorization with sparseness constraints," *Journal of Machine Learning Research*, vol. 5, pp. 1457–1469, 2004.
- [35] M. Heiler and C. Schnorr, "Learning sparse representations by non-negative matrix factorization and sequential cone programming," *Journal of Machine Learning Research*, vol. 7, pp. 1385–1407, 2006.
- [36] H. Kim and H. Park, "Sparse non-negative matrix factorizations via alternating non-negativity-constrained least squares for microarray data analysis," *Bioinformatics*, vol. 23, no. 12, pp. 1495–1502, 2007.
- [37] P. D. O'Grady and B. A. Pearlmutter, "Convolutional non-negative matrix factorisation with a sparseness constraint," in *Proc. IEEE International Workshop on Machine Learning for Signal Processing (MLSP 2006)*, Maynooth, Ireland, Sep. 2006, pp. 427–432.
- [38] T. Virtanen, "Monaural sound source separation by nonnegative matrix factorization with temporal continuity and sparseness criteria," *IEEE Transactions on Audio, Speech, and Language Processing*, vol. 15, no. 3, pp. 1066–1074, 2007.
- [39] W. Wang, "Convolutional non-negative sparse coding," in *Proc. IJCNN'08*, 2008, pp. 3681–3684.
- [40] W. Wang, A. Cichocki, and J. A. Chamber, "A multiplicative algorithm for convolutional non-negative matrix factorization based on squared Euclidean distance," *IEEE Transactions on Signal Processing*, vol. 57, no. 5, pp. 2858–2864, 2009.
- [41] C. D. Sigg, T. Dikk, and J. M. Buhmann, "Speech enhancement with sparse coding in learned dictionaries," in *Proc. ICASSP'10*, 2010.
- [42] T. Virtanen, "Separation of sound sources by convolutional sparse coding," in *Proc. SAPA'2004*, 2004.
- [43] W. Smit and E. Barnard, "Continuous speech recognition with sparse coding," *Computer Speech and Language*, vol. 23, no. 2, 2009.
- [44] G. Sivaram, S. Nemala, M. Elhilali, T. Tran, and H. Hermansky, "Sparse coding for speech recognition," in *Proc. ICASSP'10*, 2010.
- [45] P. O'Grady and B. Pearlmutter, "Convolutional non-negative matrix factorisation with a sparseness constraint," in *Proc. IEEE workshop on Machine Learning for Signal Processing*, September 2006, pp. 427–432.
- [46] H. Robbins and S. Monro, "A stochastic approximation method," *Annals of Mathematical Statistics*, vol. 22, pp. 400–407, 1951.
- [47] J. MacQueen, "Some methods for classification and analysis of multivariate observations," in *Proc. Fifth Berkeley Symposium on Mathematical Statistics and Probability*, vol. 1, 1967, pp. 281–297.
- [48] J. Mairal, F. Bach, J. Ponce, and G. Sapiro, "Online learning for matrix factorization and sparse coding," *Journal of Machine Learning Research*, vol. 2010, no. 11, pp. 19–60, January 2010.
- [49] S. S. Bucak and B. Günsel, "Incremental subspace learning via non-negative matrix factorization," *Pattern Recognition*, vol. 42, no. 5, pp. 788–797, 2009.
- [50] A. Lefevre, F. Bach, and C. Févotte, "Online algorithms for nonnegative matrix factorization with the itakura-saito divergence," in *Proc. IEEE Workshop on Applications of Signal Processing to Audio and Acoustics (WASPAA)*, Mohonk, NY, October 2011.
- [51] J. Jacques Fuchs, "On sparse representations in arbitrary redundant bases," *IEEE Transactions on Information Theory*, vol. 50, no. 6, pp. 134–1344, 2004.
- [52] D. Fisk, "Quasi-martingales," *Transactions of the American Mathematical Society*, vol. 120, no. 3, pp. 359–388, 1965.
- [53] A. W. Van der Vaart, *Asymptotic Statistics*. Cambridge University Press, 1998.
- [54] E. Vincent, C. Févotte, and R. Gribonval, "Performance measurement in blind audio source separation," *IEEE Trans. Audio, Speech and Language Processing*, vol. 14, no. 4, pp. 1462–1469, 2006.
- [55] H. Christensen, J. Barker, N. Ma, and P. Green, "The CHiME corpus: A resource and a challenge for computational hearing in multisource environments," in *Interspeech*, 2010.
- [56] S. Young, G. Evermann, M. Gales, T. Hain, D. Kershaw, X. Liu, G. Moore, J. Odell, D. Ollason, D. Povey, V. Valtchev, and P. Woodland, *The HTK Book (for Hidden Markov Model Toolkit Version 3.4)*, 2006.
- [57] R. Vipperla, S. Bozonnet, D. Wang, and N. Evans, "Robust speech recognition in multi-source noise environments using convolutional non-negative matrix factorization," in *CHiME: Workshop on Machine Listening in Multisource Environments*, 2011, pp. 74–79.

- [58] L. Bottou, "Online algorithms and stochastic approximations," in *Online Learning and Neural Networks*, D. Saad, Ed., 1998.
- [59] M. Métivier, *Semi-martingales*. Walter de Gruyter, 1983.
- [60] D. Bertsekas, *Nonlinear Programming*. Athena Scientific Belmont, 1999.