



HAL
open science

Sparse imputation for large vocabulary noise robust ASR

Jort Florent Gemmeke, Bert Cranen, Ulpu Remes

► **To cite this version:**

Jort Florent Gemmeke, Bert Cranen, Ulpu Remes. Sparse imputation for large vocabulary noise robust ASR. Computer Speech and Language, 2010, <10.1016/j.csl.2010.06.004>. <hal-00692187>

HAL Id: hal-00692187

<https://hal.science/hal-00692187v1>

Submitted on 29 Apr 2012

HAL is a multi-disciplinary open access archive for the deposit and dissemination of scientific research documents, whether they are published or not. The documents may come from teaching and research institutions in France or abroad, or from public or private research centers.

L'archive ouverte pluridisciplinaire **HAL**, est destinée au dépôt et à la diffusion de documents scientifiques de niveau recherche, publiés ou non, émanant des établissements d'enseignement et de recherche français ou étrangers, des laboratoires publics ou privés.



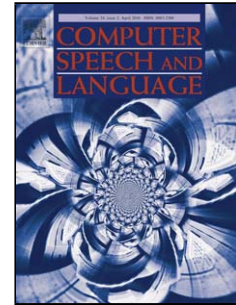
HAL Authorization

Accepted Manuscript

Title: Sparse imputation for large vocabulary noise robust ASR

Authors: Jort Florent Gemmeke, Bert Cranen, Ulpu Remes

PII: S0885-2308(10)00065-3
DOI: doi:10.1016/j.csl.2010.06.004
Reference: YCSLA 473



To appear in:

Received date: 29-1-2010
Revised date: 29-6-2010
Accepted date: 30-6-2010

Please cite this article as: Gemmeke, J.F., Cranen, B., Remes, U., Sparse imputation for large vocabulary noise robust ASR, *Computer Speech & Language* (2010), doi:10.1016/j.csl.2010.06.004

This is a PDF file of an unedited manuscript that has been accepted for publication. As a service to our customers we are providing this early version of the manuscript. The manuscript will undergo copyediting, typesetting, and review of the resulting proof before it is published in its final form. Please note that during the production process errors may be discovered which could affect the content, and all legal disclaimers that apply to the journal pertain.

Sparse imputation for large vocabulary noise robust ASR

Jort Florent Gemmeke^a, Bert Cranen^a, Ulpu Remes^b

^a*Centre for Language and Speech Technology, Radboud University Nijmegen, P.O. Box 9103, NL-6500 HD Nijmegen, The Netherlands*

^b*Adaptive Informatics Research Centre, Aalto University School of Science and Technology, P.O. Box 15400, FI-00076 Aalto, Finland*

Abstract

An effective way to increase noise robustness in automatic speech recognition is to label the noisy speech features as either reliable or unreliable (‘missing’), and replace (‘impute’) the missing ones by clean speech estimates. Conventional imputation techniques employ parametric models and impute the missing features on a frame-by-frame basis. At low SNR’s, frame-based imputation techniques fail because many time frames contain few, if any, reliable features. In previous work, we introduced an exemplar-based method, dubbed *sparse imputation*, which can impute missing features using reliable features from neighbouring frames. We achieved substantial gains in performance at low SNR’s for a connected digit recognition task. In this work, we investigate whether the exemplar-based approach can be generalised to a large vocabulary task.

Experiments on artificially corrupted speech show that sparse imputation substantially outperforms a conventional imputation technique when the ideal ‘oracle’ reliability of features is used. With error-prone estimates of feature reliability, sparse imputation performance is comparable to our baseline imputation technique in the cleanest conditions, and substantially better at lower SNR’s. With noisy speech recorded in realistic noise conditions, sparse imputation performs slightly worse than our baseline imputation technique in the cleanest conditions, but substantially better in the noisier conditions.

Key words: missing data techniques, noise robustness, automatic speech

Email addresses: J.Gemmeke@let.ru.nl (Jort Florent Gemmeke),
B.Cranen@let.ru.nl (Bert Cranen), ulpu.remes@tkk.fi (Ulpu Remes)

Preprint submitted to Computer Speech and Language

June 29, 2010

1
2
3
4
5 recognition, sparse imputation
6

7 8 **1. Introduction** 9

10 Automatic speech recognition (ASR) performance drops rapidly when speech
11 is corrupted with increasing levels of unfamiliar background noise (i.e., noise
12 not seen during training) since the observed acoustic features no longer match
13 the acoustic models. Although in real-world environments, speech is often
14 corrupted by several unknown and time-varying noise sources, few techniques
15 other than multi-condition training have been proposed to enhance robust-
16 ness towards non-stationary noise. Missing Data Techniques (MDT) [1] are
17 among the most promising alternative proposals.

18 MDT, first proposed in [2], build on the assumption that one can estimate—
19 prior to decoding—which spectro-temporal elements in the acoustic represen-
20 tation of noisy speech are reliable (i.e., dominated by speech) and which are
21 unreliable (i.e., dominated by background noise). In the unreliable elements,
22 the clean speech information is considered *missing*, and the challenge is then
23 to do speech recognition with partially observed data. In this work, we focus
24 on the so-called *imputation* approach [3] which handles the missing elements
25 by replacing them with clean speech estimates. Classic imputation methods
26 include e.g. correlation and cluster-based reconstruction [4, 1] and methods
27 for reconstruction in the cepstral and PROSPECT domains [5], while the
28 state-based imputation method proposed in [6] combines front-end imputa-
29 tion and classifier modification.
30

31 Imputation has been proven an effective technique in both small and large
32 vocabulary tasks [7], performing better than conventional feature-enhancement
33 techniques (cf. [8, 9]). However, a number of issues concerning the applicabil-
34 ity of imputation methods in different ASR tasks remain under-investigated.
35 First, since most work on MDT has been done on artificially constructed
36 databases (see e.g. [10, 4, 7]), the potentials and limitations of the missing
37 data approach in real-world environments are not well known. Using artifi-
38 cially corrupted data is attractive as it allows estimating which features are
39 reliable based on exact knowledge of the speech and noise power in each time-
40 frequency cell. This facilitates comparison of different MDT approaches and
41 allows for analysis of the influence of errors in reliability estimation. The re-
42 sults from such experiments are not, however, truly indicative for real-world
43 conditions where the observed signal is rarely a simple addition of clean
44
45
46
47
48
49
50
51
52
53

1
2
3
4
5 speech and noise: in many cases, channel effects, the Lombard effect, and
6 room reverberation also affect the observations.

7 Another issue is the imputation performance when the signal-to-noise ra-
8 tio (SNR) is low and a substantial number of frames contains few, if any,
9 reliable features. In previous work [11], we suggested that the observed per-
10 formance loss when using conventional imputation methods at low SNR's is
11 at least partly due to the fact that these methods work on a frame-by-frame
12 basis. We argued that taking into account the time-context and utilising
13 reliable features from neighbouring frames could reduce the number of impu-
14 tation errors significantly. The use of time context for imputation has been
15 explored in various other studies [12, 13]. Our approach to harnessing the
16 information in neighbouring frames is by using a novel, non-parametric impu-
17 tation method, *sparse imputation* (SI). We showed that the use of SI results
18 in large performance gains and allows for successful missing data imputation
19 at lower SNR's provided that the locations of the reliable time-frequency cells
20 are estimated accurately [14].

21 The key concept in sparse imputation is that any speech fragment can be
22 represented as a linear combination of a small number of example speech to-
23 kens. First, a dictionary of *exemplars* is constructed using fixed-length clean
24 speech tokens. Then, a sparse linear combination of exemplars is sought us-
25 ing only the reliable speech features. Imputation of the unreliable features is
26 accomplished by replacing them with the corresponding features of the linear
27 combination of clean speech dictionary exemplars. Initially, we illustrated SI
28 on an isolated digit recognition task where each fixed-length exemplar in
29 the dictionary corresponded to a complete word [11]. In [15], we success-
30 fully adapted the technique for continuous digits by using a sliding window
31 approach and a dictionary that consists of randomly selected, fixed-length
32 segments of clean speech. During imputation, the reliable features of each
33 window of the speech signal are treated as a sparse linear combination of
34 clean speech windows in the dictionary. At every instant in time, the final
35 estimates for each spectral feature vector are then calculated as an average
36 over the reconstructions in overlapping windows.

37 In [15], the sparse imputation approach for continuous digits was evalu-
38 ated on the AURORA-2 digit recognition task [16] which is frequently used
39 in noise robust ASR experiments. However, it is well known that results for
40 small vocabulary tasks are difficult to generalise to large vocabulary contin-
41 uous speech recognition. One reason for this is that speech sounds occur in
42 a much larger number of different contexts in large vocabulary tasks, which
43
44

1
2
3
4
5 might make it more difficult to model speech as a sparse combination of a
6 small number of examples. This problem will only become more serious if
7 the number of context frames in the SI approach is increased. In this paper,
8 we will investigate to what extent the increased number of reliable features
9 that comes from using multiple time-frames, in combination with the natural
10 coherence of speech signals, will result in performance gains at low SNR's,
11 despite the potential loss in accuracy due to increased variation.
12
13

14 In this work, we apply the sliding window approach for sparse imputa-
15 tion proposed in [15] on large vocabulary continuous speech data from the
16 Finnish SPEECON database [17]. The data used in the experiments are either
17 the original SPEECON data recorded in real-world noisy environments or arti-
18 ficially constructed from mixing clean speech SPEECON data and noise from
19 the NOISEX-92 database [18]. By experimenting on different window sizes
20 and noise types, we will investigate to what extent using more time-context
21 can improve recognition accuracy. We will compare the results obtained with
22 sparse imputation to results obtained with a standard frame-based paramet-
23 ric method, cluster-based imputation [4, 1], which has been shown to work
24 well for the SPEECON database [19].
25
26

27 The rest of the paper is organised as follows. In Section 2, we discuss
28 Missing Data Techniques for ASR and introduce the two types of reliabil-
29 ity estimates and missing data masks used in this work. In Section 3, we
30 briefly describe the baseline cluster-based imputation method. In Section 4,
31 we describe the sparse imputation approach and discuss the generalisation
32 to imputing large vocabulary speech by using a sliding time window. In
33 Section 5, we present the experimental setup, while the results appear in
34 Section 6 and are discussed in Section 7. Conclusions and suggestions for
35 future research are given in Section 8.
36
37
38
39

40 2. Missing Data Techniques in ASR

41 2.1. Motivation

42 In this section, we briefly discuss the MDT framework as used for noise
43 robust ASR [20, 21]. In ASR, the basic representation of speech is a spectro-
44 temporal distribution of acoustic power, a *spectrogram*. In noise-free condi-
45 tions, the value of each time-frequency cell in this two-dimensional matrix is
46 determined only by the speech signal. In noisy conditions, the value in each
47 cell represents a combination of speech and background noise power.
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

1
2
3
4
5 Assuming noise is additive, the power spectrogram of noisy speech can
6 be approximately described as the sum of the individual power spectrograms
7 of clean speech and noise. To mimic human hearing, often a Mel-frequency
8 scale and logarithmic compression of the power scale are employed. Since
9 the logarithmic compression of a sum can be approximated by the logarithm
10 of the largest of the two terms [22], it approximately holds for noisy speech
11 features that:
12
13

$$14 \quad \mathbf{Y} \approx \max(\mathbf{S}, \mathbf{N}) \quad (1)$$

15
16 with the (Mel-frequency) log-power spectrograms \mathbf{Y} denoting noisy speech,
17 \mathbf{S} denoting clean speech and \mathbf{N} representing the background noise. From (1)
18 we can infer that the noisy speech features dominated by clean speech energy
19 remain approximately uncorrupted and can be used directly as estimates of
20 the clean speech features. The noise dominated features, on the other hand,
21 provide only an upper bound for the clean speech, which means the clean
22 speech features cannot be observed and are effectively missing.
23
24
25

26 *2.2. Missing data masks*

27
28 Elements of \mathbf{Y} that predominantly contain speech or noise energy are
29 distinguished by introducing a spectrographic mask \mathbf{M} . The elements of
30 a mask \mathbf{M} are either 1, meaning that the corresponding element of \mathbf{Y} is
31 dominated by speech ('reliable') or 0, meaning that it is dominated by noise
32 ('unreliable' c.q. 'missing'). Thus, we write:
33
34

$$35 \quad \mathbf{M}(k, t) = \begin{cases} 1 \stackrel{def}{=} \text{reliable} & \text{if } \mathbf{S}(k, t) - \mathbf{N}(k, t) > \theta \\ 0 \stackrel{def}{=} \text{unreliable} & \text{otherwise} \end{cases} \quad (2)$$

36
37 with \mathbf{M} , \mathbf{Y} , \mathbf{S} , and \mathbf{N} two-dimensional matrices of $K \times T$, with frequency-
38 band index k , $1 \leq k \leq K$ and time-frame index t , $1 \leq t \leq T$. θ denotes a
39 constant SNR-threshold.
40
41

42
43 Smaller values of θ will result in more elements being considered as reliable
44 in the mask but the proportion of errors implied in the assumption that
45 $\mathbf{S}(k, t) = \mathbf{Y}(k, t)$ will be larger. Larger values of θ lead to a safer model but
46 also to fewer reliable elements for estimating the missing values.
47
48

49 *2.3. Estimating missing data masks*

50
51 In experiments with artificially added noise, an *oracle mask* can be com-
52 puted directly with (2) using knowledge of the corrupting noise and the clean
53
54

1
2
3
4
5 speech signal. The oracle mask is useful to assess the potential of missing
6 data imputation techniques and to compare the performances of different
7 techniques in ideal conditions.

8
9 In realistic situations, however, the masks must be estimated from the
10 noisy speech. Many different estimation techniques have been proposed,
11 such as SNR based estimators [23, 19], machine learning approaches to mask
12 estimation [24, 25, 26], methods that focus on speech characteristics, e.g.
13 harmonicity based SNR estimation [8, 27], and mask estimation exploiting
14 binaural cues [28] or correlogram structure [29] (cf. [30] and the references
15 therein for a more complete overview of mask estimation techniques).

16
17 In this work, we use the mask estimation approach described in [19].
18 Local SNR's are obtained from comparing the noisy speech to a noise estimate
19 which is calculated based on frames identified as non-speech by a speech/non-
20 speech classifier. Implementation details are given in Section 5.4. Since
21 imputation accuracy is dependent on the quality of the missing data mask,
22 we investigate the influence of mask estimation errors by additionally using
23 oracle masks when the test set contains speech artificially corrupted with
24 background noise.

25 26 27 28 29 *2.4. Use of MDT in ASR*

30 Techniques for speech recognition with missing data can be divided in two
31 categories: marginalisation and imputation. In the marginalisation approach
32 [10, 31], acoustic likelihoods are calculated by integrating over the range of
33 possible values of the missing features and recognition is carried out primarily
34 based on the reliable features. In the imputation approach [12, 4], the missing
35 features are replaced by clean speech estimates, after which recognition can
36 proceed without modification of the recognition system.

37
38 The marginalisation approach has been shown to be more robust against
39 data sparsity at low SNR's than the traditional imputation methods [10].
40 Imputation methods are, however, attractive for two reasons. First, after
41 the missing features have been replaced with clean speech estimates, any
42 recogniser developed for clean speech can be deployed without further modi-
43 fications. Another benefit is that the reconstructed features can be converted
44 to an arbitrary domain, like the cepstral domain. This is advantageous since
45 cepstral features are known to be less correlated and better suited for process-
46 ing with the state-of-the-art HMM-based ASR techniques [32]. Therefore, in
47 this work, only imputation techniques are investigated.

3. Cluster-based imputation

As mentioned in the introduction, we use the cluster-based imputation method proposed in [4] as the baseline approach for missing data speech recognition. It is a frame-based method where the unreliable feature values are estimated based on information in the observed features and a parametric clean speech model.

3.1. Modelling assumptions

The clean speech distribution model used in cluster-based imputation [4, 1] assumes the clean speech vectors $\mathbf{s}(t)$ are independent and identically distributed (*i.i.d.*). Thus, the model will capture the statistical dependencies between spectral channels but not between time-frames. It is also assumed the clean speech data can be clustered so that the features in each cluster are approximately normally distributed, and the clean speech can be modelled using a Gaussian mixture model (GMM):

$$P(\mathbf{s}(t)) = \sum_{\nu} P(z(t) = \nu) N[\mathbf{s}(t); \boldsymbol{\mu}(\nu), \boldsymbol{\Sigma}(\nu)] \forall t, \quad (3)$$

where ν are the cluster indices, $z(t)$ indicates the current cluster, $P(z(t) = \nu)$ are the cluster weights i.e. prior probabilities for $z(t)$, and $\boldsymbol{\mu}(\nu)$ are the cluster means and $\boldsymbol{\Sigma}(\nu)$ the covariance matrices. Here, the cluster identities $z(t)$ underlying $\mathbf{s}(t)$ are assumed unknown and modelled as a latent variable. In this work, the clusters and the distribution parameters $\boldsymbol{\mu}(\nu)$ and $\boldsymbol{\Sigma}(\nu)$ were jointly estimated from a clean speech training corpus using the expectation-maximisation (EM) algorithm.

3.2. Missing data imputation

The noisy observations $\mathbf{y}(t)$ corresponding to individual frames of the spectrogram \mathbf{Y} are divided into mutually exclusive reliable and unreliable regions $\mathbf{y}_r(t)$ and $\mathbf{y}_u(t)$ as indicated by the missing data mask (2). In cluster-based imputation [4, 1], the clean speech estimates or *reconstructions* for the unreliable features are chosen so that 1) the reconstructed vectors $\hat{\mathbf{s}}(t) = \mathbf{s}_r(t) \cup \hat{\mathbf{s}}_u(t)$ are similar to clean speech i.e. provide the best possible fit with the clean speech distribution model while 2) the reconstructed values $\hat{\mathbf{s}}_u(t)$ are constrained not to exceed the observed values $\mathbf{y}_u(t)$. Finding such reconstruction can be written as a bounded maximum a posteriori (BMAP)

estimation task, where the BMAP estimator for the unreliable features is given as

$$\hat{\mathbf{s}}_u = \operatorname{argmax}_{\mathbf{s}_u \in \mathbb{R}^U} \{P(\mathbf{s}_u | \mathbf{s}_r, \mathbf{s}_u \leq \mathbf{y}_u, \Lambda)\}, \quad (4)$$

where Λ are the parameters estimated for the GMM (3) and U is the number of unreliable features in $\mathbf{y}(t)$. We dropped the explicit notation to indicate the dependency on t , i.e., $\mathbf{s} = \mathbf{s}(t)$ and $\mathbf{y} = \mathbf{y}(t)$. Note that the model contains two unknown variables: in addition to the unreliable feature values, the cluster identities $z = z(t)$ are unknown. In (4), the latent variable z has been marginalised, but the dependency can be made explicit and (4) written as

$$\hat{\mathbf{s}}_u = \operatorname{argmax}_{\mathbf{s}_u \in \mathbb{R}^U} \left\{ \sum_{\nu} P(z = \nu | \mathbf{s}_r, \mathbf{s}_u \leq \mathbf{y}_u, \Lambda) P(\mathbf{s}_u | \mathbf{s}_r, \mathbf{s}_u \leq \mathbf{y}_u, \Lambda, \nu) \right\}, \quad (5)$$

where the first probability term is the posterior probability for the ν -th GMM cluster given the reliable features $\mathbf{y}_r(t)$ and the upper bound given by the unreliable features $\mathbf{y}_u(t)$, and the second is the cluster-conditional posterior probability for the unreliable features.

In practice, finding maximum a posteriori estimates for GMM-distributed variables is difficult. Therefore, in cluster-based imputation [4, 1], (5) is approximated as

$$\hat{\mathbf{s}}_u = \sum_{\nu} P(z = \nu | \mathbf{s}_r, \mathbf{s}_u \leq \mathbf{y}_u, \Lambda) \operatorname{argmax}_{\mathbf{s}_u \in \mathbb{R}^U} \{P(\mathbf{s}_u | \mathbf{s}_r, \mathbf{s}_u \leq \mathbf{y}_u, \Lambda, \nu)\}, \quad (6)$$

where the latter term is the cluster-conditional BMAP estimate for the unreliable features \mathbf{s}_u . The cluster-conditional estimate for \mathbf{s}_u is weighted with the posterior probability for cluster ν which is calculated based on the prior probability $P(z(t) = \nu)$ and cluster-conditional observation probability. In this work, we use full covariance matrices $\Sigma(\nu)$ and calculate the cluster-conditional BMAP estimates iteratively over the frequency channels k as proposed in [4, 1]. The covariance matrices are only assumed diagonal when evaluating the posterior probabilities for $z(t) = \nu$.

4. Sparse Imputation

In sparse imputation, speech tokens are represented as a linear combination of tokens from an overcomplete dictionary of noise-free exemplars

1
2
3
4
5 represented as fixed-length vectors. For an unknown speech token, a sparse
6 linear combination is sought in the dictionary using all reliable features in the
7 entire token. Imputation of the unreliable features is then accomplished by
8 replacing them with the corresponding values from this linear combination
9 of the clean speech dictionary exemplars. In [11], the tokens were chosen
10 to constitute time-normalised complete words, but for the continuous large
11 vocabulary speech used in this work, we must apply the sliding window ap-
12 proach proposed in [15].
13
14

15 4.1. Sparse representation of speech

16 The log-power spectrogram of clean speech, \mathbf{S} , is reshaped to a single vec-
17 tor \mathbf{s} of dimension $D = K \cdot T$ by concatenating T subsequent K -dimensional
18 time frames. For now, we assume that T is fixed. Inspired by a similar
19 approach in the field of face recognition [33], we assume that \mathbf{s} can be repre-
20 sented exactly (or at least approximated with sufficient accuracy) by a linear,
21 non-negative, combination of exemplar spectrograms \mathbf{a}_n , where n denotes a
22 specific exemplar ($1 \leq n \leq N$) in the dictionary of N available exemplars:
23
24

$$25 \mathbf{s} = \sum_{n=1}^N x_n \mathbf{a}_n = \mathbf{A} \mathbf{x} \quad \text{subject to} \quad \mathbf{x} \geq 0 \quad (7)$$

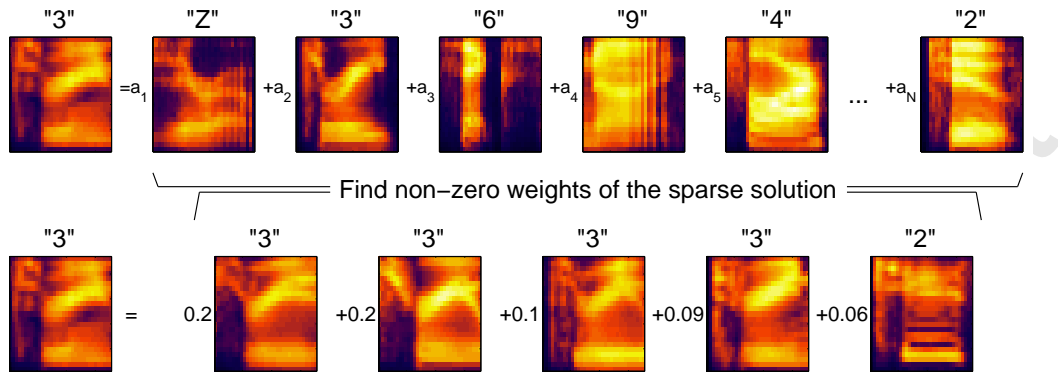
26 with \mathbf{x} an N -dimensional activation vector. The matrix \mathbf{A} denotes an over-
27 complete dictionary: $\mathbf{A} = [\mathbf{a}_1 \ \mathbf{a}_2 \ \dots \ \mathbf{a}_N]$, with dimensions $D \times N$ with
28 $N \gg D$. A schematic representation of this process for a set of non-noisy
29 spoken digits (1 through 9, “zero”, and “oh”) is displayed in Fig. 1 A.
30

31 Although the system of linear equations in (7) has no unique solution,
32 research in the field of Compressive Sensing [34, 35, 36] has shown that under
33 mild conditions on \mathbf{A} , the activation vector \mathbf{x} can be *uniquely* determined if
34 \mathbf{x} is sufficiently *sparse*. This is accomplished by solving:
35
36

$$37 \mathbf{x} = \underset{\tilde{\mathbf{x}} \in \mathbb{R}^N}{\operatorname{argmin}} \{ \|\mathbf{A} \tilde{\mathbf{x}} - \mathbf{s}\|_2 + \lambda \|\tilde{\mathbf{x}}\|_1 \} \quad (8)$$

38 with a regularisation parameter λ . The requirement that the linear combina-
39 tion must be sparse means that it must be possible to represent speech tokens
40 with a small number of exemplars, resulting in a small number of nonzero
41 values in \mathbf{x} . For spoken digits, it was shown in [14] that the representation
42 is indeed sparse.
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

A: Clean speech



B: Masked noisy speech

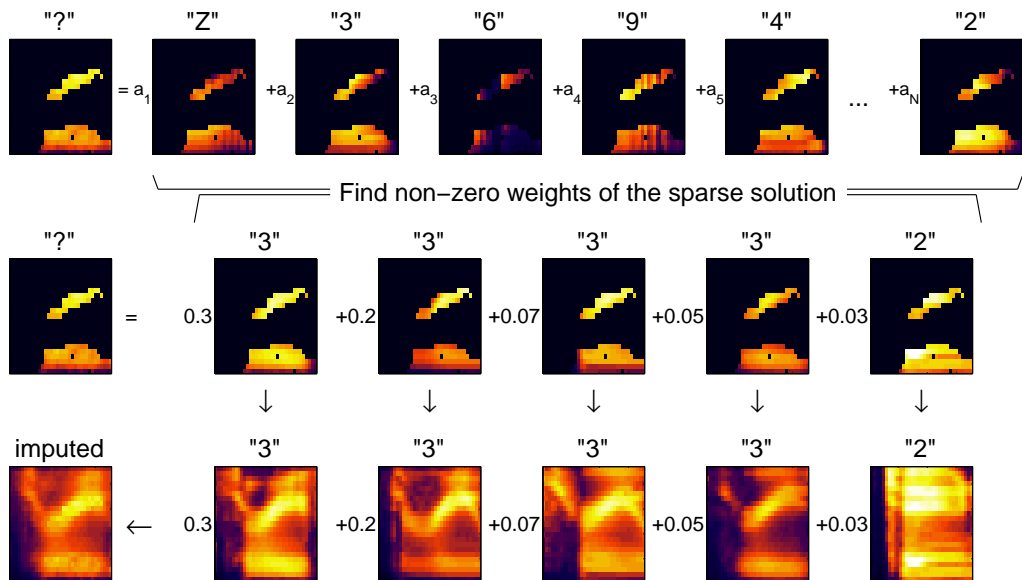


Figure 1: Schematic representation of sparse imputation, using isolated, pre-segmented digits as an example. Digit labels are at the top of the log-power spectrograms, where "Z" denotes "zero" and "O" denotes "oh". Panel A shows the sparse representation of the digit "three" in the case of clean, unmasked speech. Panel B shows the same digit with background noise at -5 dB SNR. The missing data (in black) is replaced by the corresponding features of the linear combination of clean speech dictionary exemplars found. In both panels, only the five largest nonzero weights of the linear combination are shown.

4.2. Missing data imputation

If the data contains missing values, we begin by concatenating subsequent time frames of the spectrographic mask \mathbf{M} discussed in Section 2.2 to form a mask vector \mathbf{m} ; this is done similarly as described for \mathbf{s} in the previous section. Using the same approach for the noisy speech spectrogram \mathbf{Y} we construct a noisy observation vector \mathbf{y} . The elements of \mathbf{y} corresponding to elements of the mask vector \mathbf{m} that are equal to 1 are the reliable coefficients \mathbf{y}_r . We use the reliable elements \mathbf{y}_r as an approximation for the corresponding elements of \mathbf{s} , so problem (8) becomes:

$$\mathbf{x} = \underset{\tilde{\mathbf{x}} \in \mathbb{R}^N}{\operatorname{argmin}} \{ \|\mathbf{A}_r \tilde{\mathbf{x}} - \mathbf{y}_r\|_2 + \lambda \|\tilde{\mathbf{x}}\|_1 \} \quad (9)$$

with \mathbf{A}_r pertaining to the rows of \mathbf{A} for which $\mathbf{m} = 1$. The sparse representation \mathbf{x} obtained by solving problem (9) could be directly used to estimate the clean observation vector as $\hat{\mathbf{s}} = \mathbf{A}\mathbf{x}$. A schematic representation of this process for a set of pre-segmented spoken digits (1 through 9, “zero”, and “oh”) is displayed in Fig. 1 B.

In practice, the sparse representation is not directly used as a clean speech estimate since the reconstruction error for the reliable coefficients will generally be non-zero if we solve problem (9), so it is better to only impute the unreliable elements. Furthermore, under the assumption that noise and speech are additive in the power domain, the observed noisy speech \mathbf{y} is an upper limit for $\hat{\mathbf{s}}$. Incorporating these two modifications we obtain:

$$\hat{\mathbf{s}} = \begin{cases} \hat{\mathbf{s}}_r = \mathbf{y}_r \\ \hat{\mathbf{s}}_u = \min(\mathbf{A}_u \mathbf{x}, \mathbf{y}_u) \end{cases} \quad (10)$$

with \mathbf{A}_u and $\hat{\mathbf{s}}_u$ pertaining to the rows of \mathbf{A} and $\hat{\mathbf{s}}$ for which $\mathbf{m} = 0$ and with the *min*-operator taking the element-wise minimum of two values. A version of $\hat{\mathbf{s}}$ that is reshaped into a $K \times T$ matrix $\hat{\mathbf{S}}$ can be considered a denoised spectrogram representing the underlying speech signal, and as illustrated in Fig. 1 B, it can be directly used in speech recognition.

4.3. Sliding window approach

The approach described above is suitable for imputation of noisy speech tokens that can be adequately represented by a fixed number of time frames T [11]. Since arbitrary length utterances clearly do not satisfy this constraint, we adopt a sliding window approach introduced in [15]. In this approach,

1
2
3
4
5 each window is imputed separately using sparse imputation as described in
6 Section 4.2. Subsequently, at every time frame, the different clean speech
7 estimates resulting from any overlapping windows are combined.

8
9 Consider a noisy speech utterance \mathbf{Y}_{tot} represented as a spectrogram with
10 K frequency bands and T_{tot} time-frames. The goal of the missing data impu-
11 tation process is to provide an estimate $\hat{\mathbf{S}}_{tot}$ of the underlying clean speech
12 \mathbf{S}_{tot} using a missing data mask \mathbf{M}_{tot} .

13
14 We slide a window of length T_w through \mathbf{Y}_{tot} , with shifts of $\Delta, 1 \leq \Delta \leq T_w$
15 frames (cf. Fig. 2). \mathbf{Y}_w and \mathbf{M}_w describe the noisy speech and associated
16 missing data mask for each window $w, 1 \leq w \leq W$. The ratio of Δ and T_w
17 determines the degree with which subsequent windows overlap. Larger step
18 sizes Δ reduce computational effort, but can decrease imputation accuracy
19 [15]. Throughout this paper, we keep the window shift constant at $\Delta = 1$
20 frame. The total number of windows we process is $W = T_{tot} - T_w + 1$.

21
22 We then use, for each window, the sparse imputation approach described
23 in Section 4.2 to provide a clean speech estimate $\hat{\mathbf{S}}_w$ of the underlying clean
24 speech \mathbf{S}_w . Since windows overlap, each frame in \mathbf{Y} is associated with multi-
25 ple clean speech estimate candidates, with the number of candidates ranging
26 from 1 (at the beginning and end of an utterance) to T_w . For each frame, the
27 feature values of the final clean speech estimate $\hat{\mathbf{S}}_{tot}$ are created by averaging
28 over the multiple clean speech estimate candidates pertaining to that frame
29 (cf. Fig. 2). The clean speech estimate is calculated using only clean speech
30 estimates derived from windows with a nonzero number of reliable elements.

31
32 However, in very noisy conditions and particularly at the start and end of
33 an utterance, it may happen that many adjacent windows do not contain re-
34 liable features, leaving the method unable to provide a clean speech estimate
35 from averaging. Yet, despite the lack of information about the underlying
36 signal, input must be provided to the ASR engine. Thus, we opted to impute
37 silence (the average feature values per frequency band for silence states as
38 observed in a training database) for such frames.

39 40 41 42 43 44 45 **5. Experimental setup**

46 47 *5.1. Speech recognition system and performance evaluation*

48 The speech recognition system used in this work is the large vocabulary
49 continuous speech recognition system developed in the Adaptive Informatics
50 Research Centre at the Aalto University School of Science and Technology.
51 The acoustic models are trained with 30-hours of clean speech recorded with a
52

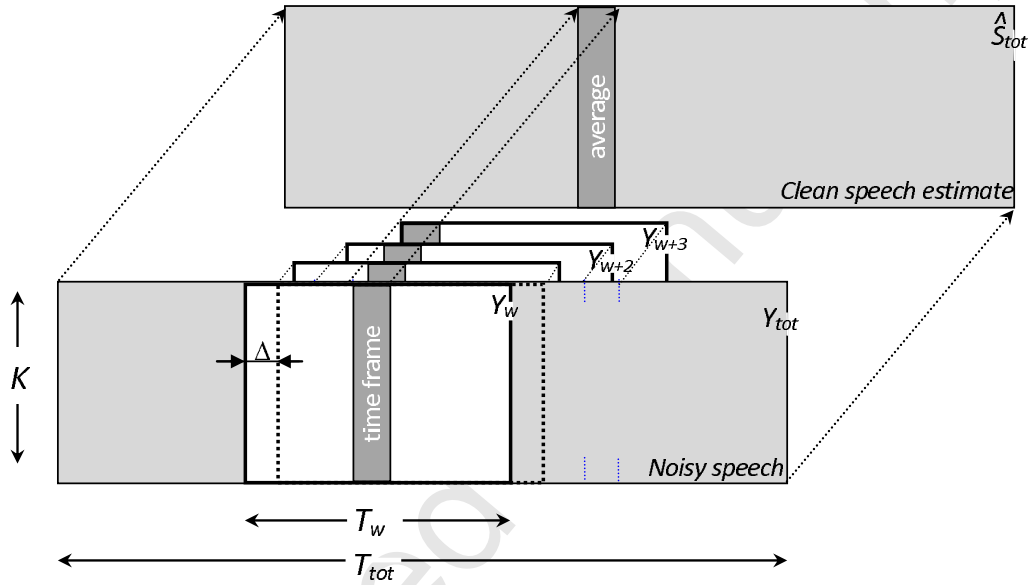


Figure 2: Schematic diagram of the sliding window approach for imputation. The dark shaded time-frame in \mathbf{Y}_{tot} is processed in several fixed-length imputation windows, of which we have shown \mathbf{Y}_w through \mathbf{Y}_{w+3} . Within each window, the given frame takes a different position due to the window shift Δ . The corresponding time-frame in the clean speech estimate $\hat{\mathbf{S}}_{tot}$ is the average over these individual window-based imputations.

1
2
3
4
5 headset in quiet conditions and selected from the Finnish SPEECON database
6 (recorded with a 16 kHz sampling rate) [17]. The training set comprises
7 293 speakers (142 female and 151 male). The utterances used for training
8 contain words, read sentences and spontaneous speech in order to have a
9 general acoustic model valid for multiple tasks.

10
11 The decoder used in the system is a time-synchronous beam-pruned Viterbi
12 token-pass system described in [37] and the acoustic models are state-clustered,
13 hidden Markov triphone models, constructed with a decision-tree method [38].
14 There are acoustic models for 13 250 triphones and two silences. The tri-
15 phones are modelled as left-right HMM with three states and silences with
16 one state each. In total we used 1564 individual states modelled with approx-
17 imately 28 Gaussians per state. Each state is also associated with gamma
18 probability functions to model the state durations [39].

19
20 The language model employs morpheme-like subword units, called *statistical*
21 *morphs*, discovered in an unsupervised, data-driven manner [40]. These
22 are used because word-based modelling is not feasible for highly inflected lan-
23 guages such as Finnish, Estonian, or Turkish. The statistical morph lexicon
24 with 25k morpheme-like units was learned from the 160k most common words
25 extracted from 145 million words of Finnish book and newspaper data [41].
26 The variable-length, growing n-gram language model [42] used in this work
27 was trained on the same text corpus and contains 52 million n-grams. The
28 decoding vocabulary is in practice unlimited since all words and word forms
29 can be represented using the statistical morphs [43].

30
31 Finally, in this work, the speech recognition performance is measured
32 primarily in letter error rates (LER). This is because the words in Finnish
33 are often long and consist of several morphemes so that measuring the word
34 error rate (WER) would correspond better to measuring sentence or phrase
35 error rates in languages such as English. Using the word error rate is also
36 considered to over-penalise misrecognised word breaks.

37 38 39 40 41 42 43 44 45 46 47 48 49 50 51 52 53 54 55 56 57 58 59 60 61 62 63 64 65

5.2. Recognition task

The imputation methods are evaluated with clean speech recordings arti-
ficially corrupted with noise at different SNR's as well as with speech recorded
in real-world noisy environments. In both conditions, the speech data consists
of read sentences selected from the Finnish SPEECON database [17]. The arti-
ficially corrupted clean speech was constructed by mixing headset-recorded
clean speech utterances with a randomly selected sample of the babble noise
from the NOISEX-92 database [18] at SNR's 15 dB, 10 dB, 5 dB and 0 dB.

The real-world noisy speech data is recorded in two types of environments: in a car and in public places both indoors and outdoors. These recordings have been made with three microphones: 1) with a headset, 2) with a lavalier microphone, and 3) with a microphone from 0.5 m–1 m distance (in public environments) or with a microphone mounted on the rear-view mirror (in car environments).

In the SPEECON documentation, the average SNR’s in the public environments are estimated to be 24 dB for the headset microphone data, 14 dB for the lavalier microphone data, and 9 dB for the far-field microphone data. For the car recordings, the estimated average SNR’s are 13 dB for the headset microphone data, 5 dB for the lavalier microphone data, and 8 dB for the rear-view mirror (RVM) microphone data. The RVM microphone data has a higher SNR than the lavalier microphone data because the RVM microphone (AKG Q400 Mk3T) has a limited frequency response, specifically designed for in-car use and suppressing low frequency noise.

The speech material in each of the three scenarios (public environments, car environments, and artificially added babble noise) is divided in development and evaluation sets. The composition of the sets in terms of number of utterances (#u), duration (d), number of speakers (#s), number of female speakers (#f) and number of male speakers (#m) is shown in Table 1. None of the sets share speakers with each other or with the speech data used for training the speech recognition system.

Table 1: Composition of development and test set.

	public					car					babble				
	#u	d (min)	#s	#f	#m	#u	d (min)	#s	#f	#m	#u	d (min)	#s	#f	#m
Development	587	60	20	7	13	288	29	10	2	8	1093	115	40	22	18
Evaluation	878	94	30	13	17	575	57	20	12	8	1118	113	40	21	19

5.3. Feature extraction

Feature extraction was carried out using a 16 ms Hamming window with 8 ms overlap between subsequent windows. First-order pre-emphasis was applied to the signal using a coefficient of 0.97. After Fourier transformation, the log-power was computed in 21 triangular-shaped Mel-frequency bands. Imputation was carried out on these log-power Mel-frequency spectra to obtain clean speech estimates.

1
2
3
4
5 After imputation, the resulting spectra were transformed to 12 Mel-
6 frequency cepstral coefficients (MFCC) and a log-energy feature, augmented
7 with first and second-order time derivatives for a total of 39 features per time
8 frame. Channel normalisation was applied using cepstral mean subtraction,
9 and as a final step, a maximum likelihood linear transformation (MLLT). The
10 MLLT, optimised during training of the acoustic models, is applied on the
11 normalised features to improve the modelling of any remaining correlation in
12 the normalised MFCCs as proposed in [44].
13
14
15

16 5.4. Missing data mask estimation

17 In this work, the missing data masks are constructed based on local SNR
18 estimates obtained from comparing the noisy speech to a static noise estimate
19 calculated during speech pauses. These speech pauses were detected using an
20 HMM-based speech/non-speech classifier described in [19]. Additionally, we
21 used the MATLAB command `bwareaopen` to discard small, isolated regions
22 of reliable features from the estimated mask since it was suggested in [45]
23 that such *glimpses* are not detectable to human listeners and are therefore
24 unlikely to contain usable information. Experiments on the SPEECON de-
25 velopment data also confirmed that removing glimpses comprising less than
26 five spectro-temporal components improves speech recognition results. In ex-
27 periments with artificially corrupted speech, we also computed oracle masks
28 (Section 2.3) from which glimpses were not removed.
29
30
31
32

33 The SNR threshold θ for deciding whether a time-frequency component
34 is treated as reliable or unreliable was determined by maximising recognition
35 accuracy on the development sets described in Section 5.2. For both real-
36 world noisy speech sets, the optimum value for estimated masks was at $\theta =$
37 3 dB for both imputation methods. For artificially corrupted speech, we
38 determined an SNR-independent threshold using the development data sets
39 containing noise at 10 and 5 dB SNR. The optimum mask threshold value
40 for the estimated masks was at $\theta = 4$ dB for both imputation methods and
41 for the oracle masks at $\theta = -2$ dB for sparse imputation and $\theta = -1$ dB for
42 cluster-based imputation.
43
44
45

46 5.5. Cluster-based imputation

47 The clean speech model used in this work is a 5-component GMM trained
48 using a 52-minute dataset of 500 read sentences randomly selected from the
49 SPEECON training data described in Section 5.1. The clusters and distri-
50 bution parameters are jointly estimated using the expectation-maximisation
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

(EM) algorithm implemented in the `GMMBayes` toolbox¹. The cluster-conditional bounded maximum a posteriori (BMAP) estimates are calculated in an iterative manner as described in [4] using a MATLAB implementation. The estimates are calculated in 10 iterations over the mel-frequency bands.

The number of 5 clusters was selected as a reasonable balance between imputation performance and computational complexity. Non-exhaustive tests on the development data showed that while more clusters do improve recognition accuracy, the increase is slight compared to the extra computational effort required.

5.6. Sparse imputation

The sparse imputation was implemented in MATLAB. The l^1 minimisation (9) was carried out using the `SolveLasso` solver.² Due to instability issues, the slower, but more robust `l1_ls_nonneg` solver³ [46] was used whenever the (fast) solver `SolveLasso` appeared to crash. When using the `l1_ls_nonneg` solver, the regularisation parameter λ was determined using the utility function `find_lambdamax_l1_ls_nonneg`. The stopping criterion of the `l1_ls_nonneg` solver was a duality gap of 0.01. The `SolveLasso` solver was run for 30 iterations.

For each window size $T_w \in \{1, 5, 10, 15, 20, 25, 30, 35\}$ being considered, an initial dictionary is created by randomly extracting 4 spectrogram segments of the desired duration T_w from each of the 8139 read sentences (containing 14 hours of speech) in the SPEECON training data described in Section 5.1. From this initial dictionary spanning 32 556 exemplars, we then randomly extract 8000 exemplars to form the final dictionary used for imputation. The dictionary size of 8000 exemplars was chosen because pilot tests showed that while using larger dictionaries improves recognition accuracy, the increase is slight enough to consider it a reasonable balance between recognition performance and computational complexity. In exemplar selection, no effort was made to balance a possible over-representation of spectra containing silence.

After reshaping the spectrograms to one-dimensional vectors as described in Section 4.1, the feature vectors form the columns of the dictionary matrix. The zero-dB level in the spectra is set such that the lowest occurring feature

¹This toolbox is publicly available from www2.it.lut.fi/project/gmmbayes/

²This solver is implemented as part of the `SparseLab` toolbox which is publicly available from <http://www.sparselab.stanford.edu>

³This solver is publicly available from http://www.stanford.edu/~boyd/l1_ls/

value in the dictionary is zero. Finally, the columns of the dictionary are normalised to Euclidean unit norm.

6. Results

6.1. Public and car environment data

The speech recognition results from our experiments with the data recorded in the public and car real-world environments are displayed in Fig. 3. The results depict the letter error rate (LER) of the sparse imputation (SI) method as a function of window size ($T_w \in \{1, 5, 10, 15, 20, 25, 30, 35\}$). The figure also shows the performance of the cluster-based imputation (CI) method and the baseline (B) recogniser. The latter has no noise compensation other than what is implicit in the feature extraction (see Section 5.3).

One point of interest is the fact that the LERs for SI with window length $T_w = 1$ are much higher than for any other window length. In fact, in many cases the performance at $T_w = 1$ is even worse than the baseline performance. Possible causes for this effect will be discussed in Section 7. In the rest of this section, we will largely ignore the data points at $T_w = 1$, and focus on a comparison of CI and B with SI at window lengths $T_w \geq 5$.

Another remark that holds for virtually all testing conditions is that with SI the performance seems to have an optimum in the range $T_w = [5, 20]$. Generally speaking, however, the differences obtained with various window sizes are quite small. Taking window size into account when comparing the SI results with those of the other methods would unnecessarily complicate matters. In the description below we will therefore focus on gross effects that can be observed in the range $T_w = [5, 20]$.

The first row in Fig. 3 illustrates results on the headset recorded data, the condition which resembles clean speech the most. Although some of the observed differences are statistically significant, the differences in performance of CI, SI, and B are quite small.

In the case of the lavalier microphone data, we can observe more differences between the car and public environments. In the car environment, both the baseline and CI method achieve 7 to 11 % absolute lower accuracies than on data recorded in the public environment. In the public environment, SI performs comparable with CI for window lengths in the range $T_w = [5, 20]$, while in the car environment SI outperforms CI by some 4 % (absolute).

In the case of the far-field microphone (public environment) or rear-view mirror (RVM) microphone (car environment), the results are similar as for

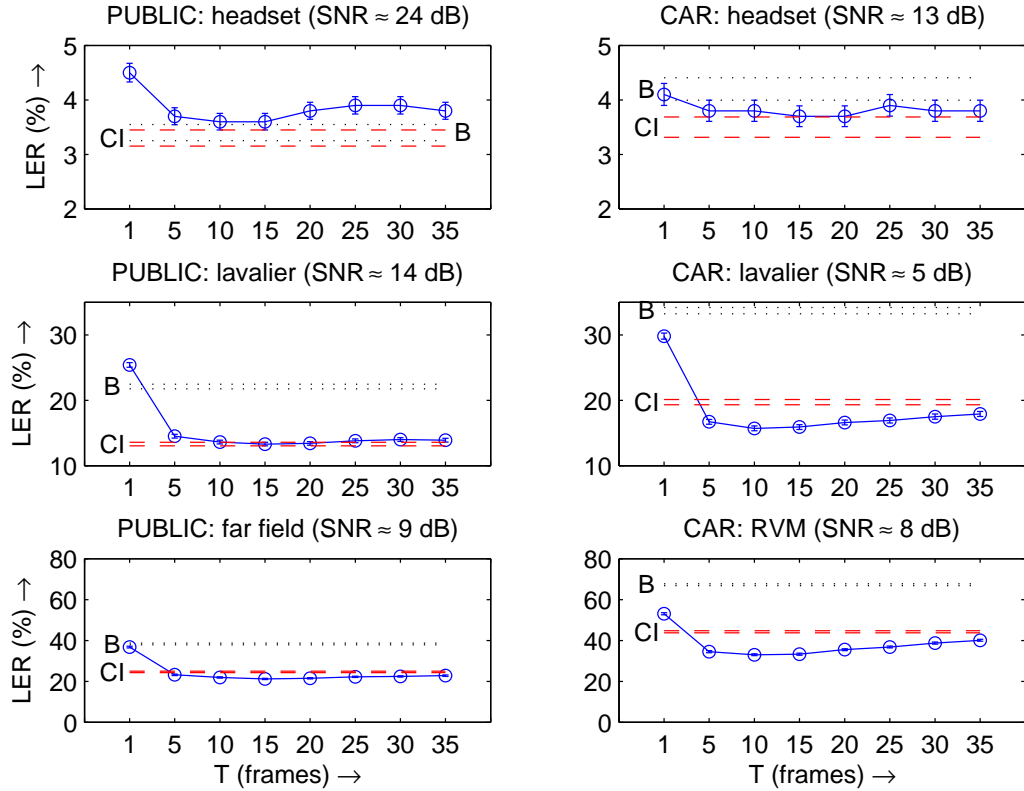


Figure 3: Recognition accuracy expressed as letter error rates (LER) for the public (left pane) and the car environment (right pane). From top to bottom, the rows correspond to headset, lavalier, and far-field microphone (public environment) or rear-view mirror (RVM) microphone (car environment). In each panel, the LER is shown as a function of window size T_w (in frames) for the sparse imputation (SI) method (solid line) with vertical bars around the data points indicating the 95% confidence intervals. The 95% confidence intervals for the cluster-based imputation (CI) method are indicated by dashed lines and that of the baseline recogniser (B) by dotted lines.

1
2
3
4
5 the lavalier microphone. Again, the car environment proves to be the more
6 difficult recognition environment with a baseline LER score of 67.3 %, com-
7 pared to 38.3 % in the public environment. As with the lavalier microphone
8 data, the SI method performs substantially better than CI in the car envi-
9 ronment, doing up to 11 % (absolute) better for window lengths in the range
10 $T_w = [5, 15]$. In the public environment, SI also performs much better than
11 CI, although by a smaller margin.

12
13 All in all, both CI and SI have a positive effect on recognition perfor-
14 mance in comparison to the baseline, but clearly SI performs better in the
15 more difficult conditions (lavalier and RVM microphone data from the car en-
16 vironment and the far-field microphone data from the public environments)
17 at the cost of a small performance loss in relatively clean environments.
18
19

20 21 *6.2. Artificially corrupted speech: babble noise scenario*

22 Recognition performance using the clean speech data artificially corrupted
23 with babble noise is displayed in Fig. 4. The baseline result for the clean
24 speech signal is $LER = 3.3$ %.

25
26 In the $SNR = 15$ dB condition, SI and CI achieve comparable accuracies
27 with $LER \approx 7.5$ %, again with the exception of $T_w = 1$ when using an es-
28 timated missing data masks. We can also observe that using an error-free
29 oracle masks leads to a much lower number of recognition errors: SI now
30 comes closer to clean speech baseline recognition scores ($LER = 4.1$ %) while
31 CI achieves 4.9 % LER. While in the case of the estimated mask, there is an
32 indication of an optimum window length in the range $T_w = [10, 15]$, no such
33 optimum can be seen when using the oracle mask.

34
35 In the $SNR=10$ dB condition, SI does significantly better than CI both
36 for estimated and oracle missing data mask. The same can be observed in the
37 $SNR=5$ dB condition, although the gap between SI and CI performance be-
38 comes much larger at lower SNR's, particularly when using the oracle mask.
39 This indicates that SI gains more performance from the extra, error-free
40 information contained in the oracle mask. In both SNR conditions, it is dif-
41 ficult to see a clear relation between the window length and SI performance,
42 although LER seems to show a shallow minimum around $T_w = [10, 20]$. Re-
43 markably, even in the $T_w = 1$ condition, SI does better than CI when using
44 the oracle mask.

45
46 In the $SNR=0$ dB condition when using an estimated mask, neither SI
47 nor CI can reconstruct the clean speech signal to a sufficient degree to achieve
48 a usable performance, since both methods have LER around 75 %. When
49
50
51
52
53

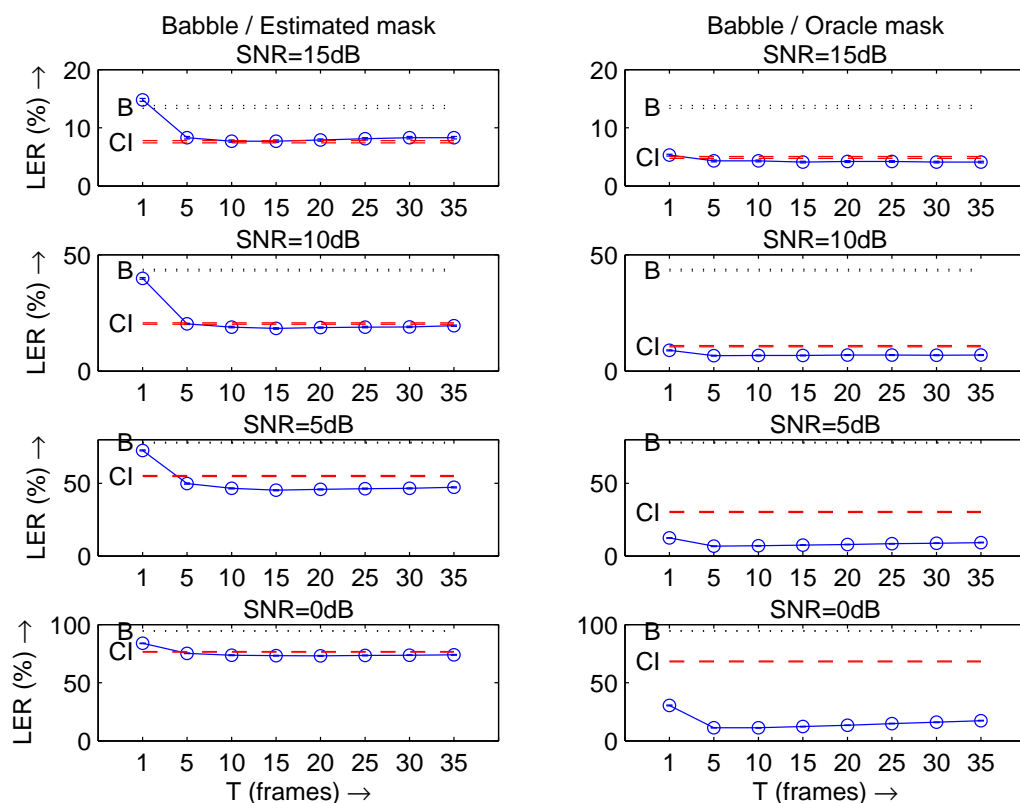


Figure 4: Recognition accuracy expressed as letter error rates (LER) for the dataset containing clean speech artificially corrupted with babble noise. The results are shown for the estimated missing data mask (left pane) and oracle missing data mask (right pane). Row one through four apply to different signal-to-noise ratio's (SNR). In each panel, the LER is shown as a function of window size T_w (in frames) for the sparse imputation (SI) method (solid line) with vertical bars around the data points indicating the 95% confidence intervals. The 95% confidence intervals for the cluster-based imputation (CI) method are indicated by dashed lines and that of the baseline recogniser (B) by dotted lines.

1
2
3
4
5 using an oracle mask, the situation is quite different. While CI achieves a
6 LER of 68.3 %, gaining less than 10 % from using an oracle mask, SI achieves
7 error rates of only 11.2 % at SNR = 0 dB.
8
9

10 7. Discussion

11 7.1. Sparse imputation for large vocabulary continuous speech recognition

12
13 Research on noise robust MDT started out with experiments on small
14 vocabulary tasks artificially corrupted by noise [10, 5]. While vocabulary
15 sizes have since increased, in the majority of cases still artificially corrupted
16 speech has been employed [4, 47, 27].
17
18

19 It is not until recently that research has turned to MDT on large vocab-
20 ulary speech recorded in realistic conditions [19, 48]. In this work, we in-
21 vestigated whether the improvements in recognition accuracy obtained with
22 SI on the AURORA-2 digit recognition task [15] would generalise to a large
23 vocabulary task.
24

25 Experiments on the SPEECON data indicate that SI is 1) indeed capa-
26 ble of significantly improving large vocabulary continuous speech recognition
27 (LVCSR) performance on noisy speech data and 2) performs equally well on
28 artificially corrupted data and noisy speech recorded in real-world environ-
29 ments where different types of microphones and different microphone-speaker
30 distances were used.
31
32

33 Compared to the cluster-based imputation (CI) method, SI improved
34 the speech recognition performance especially at low SNR's. When using
35 an estimated missing data mask, SI performance on the cleanest conditions
36 (headset-recorded data or artificially corrupted data at SNR=15 dB) was
37 comparable or slightly lower than the CI performance, but SI performed
38 better than CI in all the noisier conditions. The difference between SI and
39 CI was most notable when using SI on the car data recorded with a rear-view-
40 mirror (RVM) microphone, which resulted in 26 % relative error reduction in
41 LER. In the experiments using oracle masks, SI outperformed CI at all SNR
42 levels, and as with the estimated masks, the differences in performance grew
43 larger at low SNR's. On the artificially corrupted data at SNR=0 dB, the
44 relative error reduction in LER was 84 % compared to the CI performance.
45
46
47
48

49 7.1.1. Sparse representation of speech

50 Initially, we expressed some concern regarding the performance limits of
51 an exemplar-based method in modelling large vocabulary speech, especially
52
53

1
2
3
4
5 when using windows spanning multiple time frames. The oracle mask results
6 reported in Section 6, however, serve as empirical evidence that such concerns
7 are unfounded, since the use of 8000 exemplars suffices for SI to find a proper
8 reconstruction of the underlying clean speech, even when many features are
9 missing.

10 To study this issue in more detail, we did a small additional experiment in
11 which we investigated for three window lengths the sparsity of clean (uncor-
12 rupted) speech of a random subset of 10 utterances of the SPEECON test
13 database. At window lengths $T_w \in \{10, 20, 30\}$, the utterances contain
14 6794, 6694, 6594 windows, respectively.

15 For each of the windows, the observation vector was first normalised to
16 a Euclidean unit norm after which we recovered its sparse representation \mathbf{x}
17 using the `l1_ls_nonneg` solver solver (cf. Section 5.6) and sorted the elements
18 of this vector with respect to weight. Finally, we averaged the sorted weight
19 vectors over all windows.

20 The result is an average weight vector ordered with respect to weight
21 which indicates how many exemplars (on average) are needed to represent
22 the windows of the selected test utterances. For each window length, the 45
23 largest weights are shown in Figure 5. From this figure it can be deduced that
24 the fixed length spectrogram segments of large vocabulary SPEECON speech in
25 the test set can indeed be sparsely represented. The results show that, within
26 the accuracy of the solver, on average windowed spectra can be sparsely
27 represented using no more than approximately 30 exemplars.

28 7.1.2. Low-noise conditions

29 With CI performance being consistently lower than the SI performance
30 in all the noisier conditions, it is interesting to explore in more detail why CI
31 performs comparable or better than SI when using an estimated mask in the
32 cleanest conditions, i.e., the headset-recorded data or artificially corrupted
33 data at SNR=15 dB. A factor that may contribute to the difference is that
34 CI uses the unreliable features as an upper bound during the missing fea-
35 ture reconstruction. This has been shown to improve the MDT performance
36 in various noise conditions [10]. In SI, the upper bound is applied on the
37 reconstructed features *after* selection of the exemplars as indicated in Equa-
38 tion (10). As a consequence, it is not taken into account that if some of the
39 estimates are considered incorrect because they are larger than the observed
40 upper bound, there is no guarantee the other features are correct since they
41 stem from the same linear combination of exemplars.

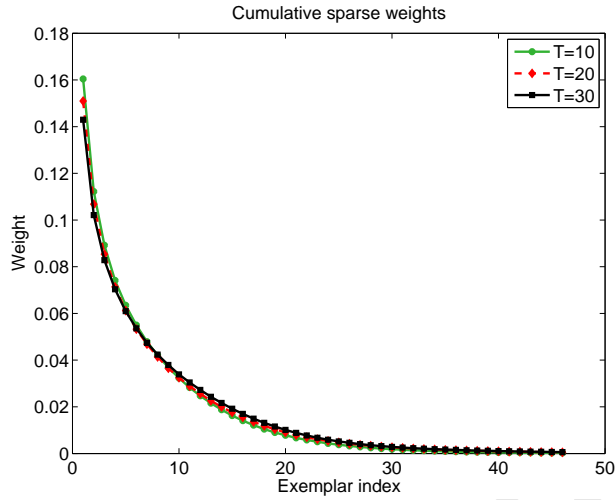


Figure 5: The sparsity of clean speech in a subset of the SPEECON test database. The graph shows the average weight of the 45 largest nonzero elements of \mathbf{x} of each sparsely represented window in a random subset of 10 utterances. The results displayed here pertain to the window lengths $T_w \in \{10, 20, 30\}$.

The upper bounds are likely to have the largest impact in low-noise conditions, which is where CI performs best, but the treatment of upper bounds does not explain why SI does better than CI at SNR=15 dB if the estimated masks are replaced with oracle masks. Analysing the recognition errors in detail (not shown) revealed that when CI outperforms SI, the difference is mainly due to insertion errors. This is in line with the observations in [15]. The SI method is prone to insertion errors because even a single, isolated, reliable feature results in the entire T_w -frame window being represented as a linear combination of clean speech exemplars. Depending on the decoder, such a single reliable feature could lead to the insertion of a segment/letter/word. These insertion errors are most likely when using an estimated mask that erroneously marks spurious features reliable.

It is noteworthy that in the noisier conditions, an isolated reliable feature may be all that is left of the underlying clean speech, and indeed, the same property that makes SI prone to insertion errors in relatively clean conditions is what causes the performance gain in the more difficult noise conditions where SI performs better than CI. A solution for this issue may

1
2
3
4
5 be found in techniques such as weighted Viterbi decoding [13] or uncertainty
6 decoding [49], in which the decoder takes the uncertainty of the accuracy with
7 which clean speech was estimated into account. Preliminary experiments on
8 artificially corrupted speech showed that SI performance indeed increases
9 both at high and low SNR's when using a proper measure of uncertainty (cf.
10 [50]).
11

12 7.2. Optimal time context

13 7.2.1. Using multiple frames of time-context

14
15
16 Figures 3 and 4 in Section 6 clearly indicate that increasing the time
17 context beyond $T_w = 1$ in the SI method improves the speech recognition
18 performance. In most conditions, the optimal recognition performance is
19 achieved when the time context in SI is approximately 5–20 frames. A de-
20 tailed analysis of the recognition errors (not shown) revealed that increasing
21 the time context past the optimum range systematically leads to an increased
22 number of deletion errors and that in general, when SI outperformed CI, it
23 was due to a reduced number of deletion errors.
24

25
26 Analysis of the clean estimates produced by SI at window sizes larger than
27 5–20 frames (not shown) revealed that the increased number of deletions is
28 due to difficulties in finding a sparse linear combination of exemplars that
29 would describe the high-dimensional observation sufficiently accurately. Con-
30 sequently, the resulting sparse representations only capture the high energy
31 regions of the spectrogram window and do not describe the details accu-
32 rately. In addition, the number of spectrograms over which we average to
33 calculate the final estimate increases with increasing window length, results
34 in smoothed approximations for the clean speech spectrograms and tends to
35 decrease the dynamic range of the feature values. Since the decoder employs
36 mean normalisation, the smoothed, mean-normalised features result in a re-
37 duced contrast between states and may even start to resemble silence. This,
38 in turn, can lead to recognition errors such as deletions.
39

40
41
42 Finally, as the results reported in Section 6 and the previous results in [14]
43 indicate that missing data imputation using SI benefits from increasing the
44 time context, the question could come to mind whether using more time-
45 context might also be advantageous for CI. While we do not address this
46 question in this work, it should be noted that the computational complexity
47 of CI is $\mathcal{O}(K^4)$, where K denotes the feature dimension, which shows that
48 adding time context (effectively increasing the number of features per frame)
49 would quickly become infeasible. Moreover, there is a quadratic increase
50
51
52
53

1
2
3
4
5 in the amount of data required to accurately estimate the full covariances
6 needed for CI, which can lead to training data scarcity.
7

8 9 7.2.2. *SI performance for single time frames*

10 Although the SI performance is not better than CI performance for every
11 noise scenario and window length T_w , the performance is never much worse
12 than CI, except for $T_w = 1$ in combination with an estimated mask. In this
13 case, the SI performance is occasionally even worse than the uncompensated
14 baseline system performance. However, if oracle masks are used instead, the
15 SI performance is substantially better than CI, even with $T_w = 1$.
16

17 There are two issues that could explain the difference. First, the way
18 frames are treated when all features are unreliable is different in SI and CI.
19 In SI, frames without any reliable features are imputed as silence. If several
20 of such frames occur during a speech segment, the decoder is more or less
21 forced to recognise these segments as silence. Since it becomes less likely that
22 a speech segment of T_w frames does not contain any reliable frames when
23 T_w is large, SI achieves better results at longer window sizes. In contrast,
24 CI imputes frames without any reliable values by making use of the prior
25 probability given the noisy data. In this approach, the imputed frame is less
26 likely to be interpreted as silence all the time.
27

28 To test this hypothesis, we ran an additional experiment on the artificial
29 babble noise data. Here, CI was modified to also impute silence for all the
30 frames that contain no reliable features. The results (not shown) confirmed
31 that this decreases the system performance when an estimated missing data
32 mask is used, although the performance was still better than the SI perfor-
33 mance at $T_w = 1$. Interestingly, the modification *improved* the CI perfor-
34 mance when an oracle mask was used, although the obtained accuracies were
35 still substantially lower than the SI accuracies at $T_w = 1$. In summary, this
36 experiment shows not only that imputing silence when all features are unre-
37 liable is only a good approach in the absence of mask estimation errors, but
38 also that the difference in treating frames without any reliable values does
39 not fully explain the differences between SI using $T_w = 1$ and CI.
40

41 A second factor that may contribute to the sub-optimal performance of
42 SI when $T_w = 1$ and an estimated mask is used, is that at $T_w = 1$, the
43 minimisation problem (9) gets extremely under-determined as the number
44 of reliable features decreases. It is a property of the applied technique that
45 the quality of the reconstruction does not deteriorate gradually, but suffers
46 from a sudden break-down once the number of reliable features gets below
47
48
49
50
51
52
53

1
2
3
4
5 a certain threshold. This happens because with very few reliable values, the
6 sparsest representation can suddenly be constituted of rather arbitrary ex-
7 emplar vectors. As discussed in detail in [14], it is very difficult to give an
8 estimate on the minimum number of (reliable) features needed for imputa-
9 tion, but the results in the paper can be seen as an empirical indication that
10 in low SNR conditions, a sizable proportion of single time frames does not
11 contain enough reliable features for the SI method to work properly.
12

13
14 In light of these results, we may need to reconsider the SI approach to
15 impute silence for frames that contain no reliable features, at least when a
16 short window is used. Possibly, a better approach would be to interpolate
17 between neighbouring frames, to use the output of the CI method for frames
18 without any reliable values, or to use an approach related to the uncertainty
19 decoding approach mentioned above by setting all state likelihoods to an
20 equal value for frames which do not contain reliable features.
21
22

23 *7.3. Influence of mask quality*

24

25 It is well known that recognition accuracy improvements obtainable with
26 MDT are highly dependent on the quality of the applied missing data mask.
27 When comparing the CI results for estimated and oracle mask, it is obvious
28 that the absence of mask estimation errors substantially improves recogni-
29 tion accuracy, with relative LER improvements ranging from 10% to 45%.
30 The same holds for sparse imputation, although the differences are larger,
31 particularly at low SNR's. As can be seen in Figure 4, at an SNR of 0 dB,
32 SI achieves an accuracy of 11.2 % LER as opposed to the 73.1% with the
33 estimated mask at $T_w = 20$. As in our previous experiments on a digit recog-
34 nition task [15], the large performance gap between the two types of missing
35 data masks indicates that SI can potentially perform much better than CI,
36 provided mask estimation errors are reduced.
37

38 An equally valid conclusion is that SI is more sensitive to mask estima-
39 tion errors than CI. The reason for this is that mask estimation errors which
40 incorrectly label features reliable will mislead the search for the 'true' sparse
41 representation associated with the underlying clean speech. Depending on
42 the location of these features, imputed features can become very different
43 from the underlying clean speech. Using a more conservative estimation
44 method is no solution, since mask estimation errors which incorrectly label
45 features as unreliable reduce the overall number of reliable features avail-
46 able for imputation and cause the search to miss out on features useful for
47 distinguishing between exemplars [14].
48
49
50
51
52
53

1
2
3
4
5 The recognition performance using error-free oracle masks suggest an even
6 larger real-world potential for SI, provided substantially better mask estima-
7 tion methods can be found. Alternatively, the SI method could be made
8 more robust against mask estimation errors by a number of algorithmical
9 improvements. An attractive approach would be to change the search for
10 exemplars, i.e., the minimisation problem (9), to include the constraint that
11 the reconstructed speech should not exceed the noisy observation. Other
12 possibilities for improvements are the additional use of derivative features in
13 the exemplars rather using only static features, or the use of *soft* missing
14 data masks [51]. In soft masks the binary reliability score is replaced by the
15 probability that a spectral component is reliable, providing more robustness
16 against mask estimation errors. For a more extensive discussion on these
17 potential improvements, we refer the reader to [14, 52].
18
19
20
21

22 7.4. Computational effort

23 To roughly characterise the computational effort needed, we did a small
24 test of the running time of the CI algorithm on a machine with a Core 2 Duo
25 E6550 2.33 GHz processor. The running time for an utterance of 756 frames
26 (6 seconds of speech) containing 11 734 unreliable values was 80 seconds. For
27 the SI algorithm with a window length of $T_w = 10$ frames, the running time
28 on this utterance was 61 seconds.
29
30

31 While SI is faster for this particular utterance, SI still performs at about
32 10 times real-time. It is therefore interesting to study its computational
33 complexity. SI using the `SolveLasso` solver has a computational complexity
34 of $\mathcal{O}(W((RT_w)^3 + NRT_w))$, where W denotes the number of windows to be
35 processed, N the number of clean speech exemplars in the dictionary, R the
36 average number of reliable features per frame, and T_w the window length in
37 frames. In practice, the computational complexity is completely dominated
38 by the term NRT_w .
39
40
41

42 There are three ways to reduce the computational effort, each of which
43 scale approximately linearly. The first is to increase the window shift Δ which
44 decreases the number of windows W . In [15] it was shown that increasing
45 Δ too much reduces imputation performance, but also that small increases
46 do not decrease accuracy. Moreover, Δ does not have to be constant over
47 the utterance and could for example be made dependent on the number of
48 reliable features. The second approach is to reduce the number of features in
49 each window. Aside from reducing the window length, which in Section 7.2
50 was shown to decrease performance when reduced to much, it is also possible
51
52
53

1
2
3
4
5 to apply dimensionality reduction to the features in a window, as used in [33].
6 Finally, the dictionary size might be reduced by methods such as clustering
7 and an algorithmic way to handle shift-invariance rather than by including
8 time-shifted variants of the same phenomena.
9

10 11 **8. Conclusions and future work**

12
13 In this work, we investigated the performance of the sparse imputation
14 missing data technique on read sentences of the Finnish SPEECON corpus us-
15 ing real-world noisy speech recordings as well as clean speech recordings arti-
16 ficially corrupted with babble noise. In previous research sparse imputation
17 was shown to be an effective method for improving the noise robustness of a
18 connected digit recognizer using the artificially noisified AURORA-2 data [15].
19 The current results show that the method can be readily extended to a large
20 vocabulary recognition task in which the speech suffers from corruptions
21 found in real-world environments. We found that also in the SPEECON cor-
22 pus, fixed-length spectrogram windows can be adequately represented by a
23 sparse, linear combination of exemplars: On average less than 30 are needed.
24 As with the AURORA-2 task, sparse imputation on SPEECON greatly ben-
25 efits from using additional time-context for imputation. With a dictionary
26 size of 8000 randomly selected exemplars used in this study, we typically
27 found a context of 5-20 frames (i.e., 50-200 ms) to yield the best recognition
28 accuracies.
29

30
31 Experiments on artificially corrupted speech indicated that sparse im-
32 putation outperforms a conventional imputation technique by a significant
33 margin when the ideal ‘oracle’ reliability of noisy speech features are used.
34 With error-prone reliability estimates, sparse imputation performs slightly
35 worse than our baseline imputation technique in the cleanest conditions, but
36 significantly better at lower SNR’s.
37

38
39 Future work in the sparse imputation framework will focus on improving
40 the robustness toward mask estimation errors and better handling the frames
41 which do not contain any reliable values.
42

43 44 **Acknowledgement**

45
46 The research by Jort Florent Gemmeke was carried out in the MIDAS
47 project, granted under the Dutch-Flemish STEVIN program. The research
48 by Ulpu Remes was supported by the Helsinki Graduate School in Computer
49

1
2
3
4
5 Science and Engineering and by the Academy of Finland in the projects *Au-*
6 *ditory approaches to automatic speech recognition* and *Adaptive Informatics*
7 *Research Centre*. We acknowledge Lou Boves and Kalle Palomäki for useful
8 discussions.
9

10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

Accepted Manuscript

References

- [1] B. Raj, R. M. Stern, Missing-feature approaches in speech recognition, *IEEE Signal Processing Magazine* 22 (5) (2005) 101–116.
- [2] M. Cooke, P. Green, M. Crawford, Handling missing data in speech recognition, in: *Proc. ICSLP, Yokohama, Japan, 1994*, pp. 1555–1558.
- [3] B. Raj, R. Singh, R. Stern, Inference of missing spectrographic features for robust automatic speech recognition, in: *Proc. ICSLP, Sydney, Australia, 1998*, pp. 1491–1494.
- [4] B. Raj, M. Seltzer, R. Stern, Reconstruction of missing features for robust speech recognition, *Speech Communication* 43 (4) (2004) 275–296.
- [5] H. Van hamme, PROSPECT features and their application to missing data techniques for robust speech recognition, in: *Proc. INTER-SPEECH, Jeju Island, Korea, 2004*, pp. 101–104.
- [6] L. Josifovski, M. Cooke, P. Green, A. Vizinho, State based imputation of missing data for robust speech recognition and speech enhancement, in: *Proc. EUROSPEECH, Budapest, Hungary, 1999*, pp. 2837–2840.
- [7] M. V. Segbroeck, Robust large vocabulary continuous speech recognition using missing data techniques, Ph.D. thesis, K.U.Leuven (2010).
- [8] H. Van hamme, Robust speech recognition using cepstral domain missing data techniques and noisy masks, in: *Proc. ICASSP, Montreal, Quebec, Canada, 2004*, pp. 213–216.
- [9] V. Stouten, Robust automatic speech recognition in time-varying environments, Ph.D. thesis, K.U.Leuven (2006).
- [10] M. Cooke, P. Green, L. Josifovski, A. Vizinho, Robust automatic speech recognition with missing and unreliable acoustic data, *Speech Communication* 34 (3) (2001) 267–285.
- [11] J. Gemmeke, B. Cranen, Using sparse representations for missing data imputation in noise robust speech recognition, in: *Proc. EUSIPCO, Lausanne, Switzerland, 2008*.

- 1
2
3
4
5 [12] B. Raj, Reconstruction of incomplete spectrograms for robust speech
6 recognition, Ph.D. thesis, Carnegie Mellon University (2000).
7
- 8 [13] Z.-H. Tan, P. Dalsgaard, B. Lindberg, Exploiting temporal correlation of
9 speech for error-robust and bandwidth-flexible distributed speech recog-
10 nition, *IEEE Transactions on Audio, Speech and Language Processing*
11 15 (4) (2007) 1391–1403.
12
13
- 14 [14] J. F. Gemmeke, H. V. hamme, B. Cranen, L. Boves, Compressive sensing
15 for missing data imputation in noise robust speech recognition, *IEEE*
16 *Journal of Selected Topics in Signal Processing* 4 (2) (2010) 272–287.
17
18
- 19 [15] J. Gemmeke, B. Cranen, Missing data imputation using compressive
20 sensing techniques for connected digit recognition, in: *Proc. DSP, San-*
21 *torini, Greece, 2009*, pp. 1–8.
22
23
- 24 [16] H. Hirsch, D. Pearce, The Aurora experimental framework for the per-
25 formance evaluation of speech recognition systems under noisy condi-
26 tions, in: *Proc. ISCA Tutorial and Research Workshop ASR2000, Paris,*
27 *France, 2000*, pp. 181–188.
28
29
- 30 [17] D. Iskra, B. Grosskopf, K. Marasek, H. V. D. Heuvel, F. Diehl,
31 A. Kiessling, SPEECON - speech databases for consumer devices:
32 Database specification and validation, in: *Proc. LREC, Las Palmas,*
33 *Canary Islands, Spain, 2002*, pp. 329–333.
34
35
- 36 [18] A. Varga, H. Steeneken, Assessment for automatic speech recognition:
37 II. NOISEX-92: a database and an experiment to study the effect of
38 additive noise on speech recognition systems, *Speech Communication*
39 12 (3) (1993) 247–51.
40
41
- 42 [19] U. Remes, K. J. Palomäki, M. Kurimo, Missing feature reconstruction
43 and acoustic model adaptation combined for large vocabulary contin-
44 uous speech recognition, in: *Proc. EUSIPCO, Lausanne, Switzerland,*
45 *2008*.
46
47
- 48 [20] J.-C. Junqua, J.-P. Haton, *Robustness in Automatic Speech Recog-*
49 *nition: Fundamentals and Applications*, Kluwer Academic Publishers,
50 1996.
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4
5 [21] J. Barker, M. Cooke, D. Ellis, Decoding speech in the presence of other
6 sources, *Speech Communication* 45 (1) (2005) 5–25.
7
8 [22] A. Nadas, D. Nahamoo, M. Picheny, Speech recognition using noise-
9 adaptive prototypes, *IEEE Transactions on Acoustics, Speech and Sig-
10 nal Processing* 37 (10) (1989) 1495–1503.
11
12 [23] A. Vizinho, P. Green, M. Cooke, L. Josifovski, Missing data theory,
13 spectral subtraction and signal-to-noise estimation for robust ASR: An
14 integrated study, in: *Proc. EUROSPEECH*, Budapest, Hungary, 1999,
15 pp. 2407–2410.
16
17 [24] M. Seltzer, B. Raj, R. Stern, A Bayesian classifier for spectrographic
18 mask estimation for missing feature speech recognition, *Speech Com-
19 munication* 43 (4) (2004) 379–393.
20
21 [25] W. Kim, R. M. Stern, Band-independent mask estimation for missing-
22 feature reconstruction in the presence of unknown background noise, in:
23 *Proc. ICASSP*, Toulouse, France, 2006, pp. 305–308.
24
25 [26] R. Weiss, D. Ellis, Estimating single-channel source separation masks:
26 Relevance Vector Machine classifiers vs. pitch-based masking, in: *Proc.
27 Workshop on Statistical and Perceptual Audition SAPA-06*, Pittsburgh,
28 Pennsylvania, USA, 2006, pp. 31–36.
29
30 [27] M. Van Segbroeck, H. Van hamme, Vector-Quantization based mask
31 estimation for missing data automatic speech recognition, in: *Proc. IN-
32 TERSPEECH*, Antwerp, Belgium, 2007, pp. 910–913.
33
34 [28] S. Harding, J. Barker, G. J. Brown, Mask estimation for missing data
35 speech recognition based on statistics of binaural interaction, *IEEE
36 Transactions on Audio, Speech and Language Processing* 14 (1) (2006)
37 58–67.
38
39 [29] N. Ma, P. Green, J. Barker, A. Coy, Exploiting correlogram structure
40 for robust speech recognition with multiple speech sources, *Speech Com-
41 munication* 49 (12) (2007) 874–891.
42
43 [30] C. Cerisara, S. Demange, J.-P. Haton, On noise masking for automatic
44 missing data speech recognition: A survey and discussion, *Computer
45 Speech and Language* 21 (3) (2007) 443–457.
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4
5 [31] J. Barker, M. Cooke, P. Green, Robust ASR based on clean speech mod-
6 els: An evaluation of missing data techniques for connected digit recog-
7 nition in noise, in: Proc. EUROSPEECH, Aalborg, Denmark, 2001, pp.
8 213–216.
9
- 10 [32] S. B. Davis, P. Mermelstein, Comparison of parametric representations
11 for monosyllabic word recognition in continuously spoken sentences,
12 IEEE Transactions on Acoustics, Speech, and Signal Processing 28 (4)
13 (1980) 357–366.
14
- 15 [33] J. Wright, A. Y. Yang, A. Ganesh, S. S. Sastry, Y. Ma, Robust face
16 recognition via sparse representation, IEEE Transactions on Pattern
17 Analysis and Machine Intelligence 31 (2) (2009) 210–227.
18
- 19 [34] D. L. Donoho, Compressed sensing, IEEE Transactions on Information
20 Theory 52 (4) (2006) 1289–1306.
21
- 22 [35] D. L. Donoho, For most large underdetermined systems of linear equa-
23 tions the minimal L1-norm solution is also the sparsest solution, Com-
24 munications on Pure and Applied Mathematics 59 (6) (2006) 797–829.
25
- 26 [36] E. J. Candès, J. Romberg, T. Tao, Stable signal recovery from incom-
27 plete and inaccurate measurements, Communications On Pure and Ap-
28 plied Mathematics 59 (8) (2006) 1207–1223.
29
- 30 [37] J. Pyllkkönen, An efficient one-pass decoder for Finnish large vocabu-
31 lary continuous speech recognition, in: Proc. 2nd Baltic Conference on
32 Human Language Technologies, Tallinn, Estonia, 2005, pp. 167–172.
33
- 34 [38] J. J. Odell, The use of context in large vocabulary speech recognition,
35 Ph.D. thesis, University of Cambridge (1995).
36
- 37 [39] J. Pyllkkönen, M. Kurimo, Duration modeling techniques for continuous
38 speech recognition, in: Proc. INTERSPEECH, Jeju Island, Korea, 2004,
39 pp. 385–388.
40
- 41 [40] M. Creutz, K. Lagus, Unsupervised discovery of morphemes, in:
42 Proc. ACL-02 Workshop on Morphological and Phonological Learning,
43 Philadelphia, Pennsylvania, USA, 2002, pp. 21–30.
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4
5 [41] CSC Tieteellinen laskenta Oy, The language bank of Finland (2001).
6 URL www.csc.fi/languagebank/
7
- 8 [42] V. Siivola, B. Pellom, Growing an n-gram language model, in: Proc.
9 INTERSPEECH, Lisbon, Portugal, 2005, pp. 1309–1312.
10
- 11 [43] T. Hirsimäki, M. Creutz, V. Siivola, M. Kurimo, S. Virpioja,
12 J. Pylkkönen, Unlimited vocabulary speech recognition with morph lan-
13 guage models applied to Finnish, *Computer Speech and Language* 20 (4)
14 (2006) 515–541.
15
16
- 17 [44] M. Gales, Semi-tied covariance matrices for hidden Markov models,
18 *IEEE Transactions on Speech and Audio Processing* 7 (3) (1999) 272–
19 281.
20
21
- 22 [45] M. Cooke, A glimpsing model of speech perception in noise, *J. Acoust.*
23 *Soc. Am.* 119 (3) (2006) 1562–1573.
24
25
- 26 [46] S. Kim, K. Koh, M. Lustig, S. Boyd, D. Gorinevsky, An interior-point
27 method for large-scale l_1 -regularized least squares, *IEEE Journal on*
28 *Selected Topics in Signal Processing* 1 (4) (2007) 606–617.
29
30
- 31 [47] S. Srinivasan, N. Roman, D. Wang, Binary and ratio time-frequency
32 masks for robust speech recognition, *Speech Communication* 48 (11)
33 (2006) 1486–1501.
34
- 35 [48] J. F. Gemmeke, Y. Wang, M. V. Segbroeck, B. Cranen, H. V. hamme,
36 Application of noise robust MDT speech recognition on the SPEECON
37 and SpeechDat-Car databases, in: Proc. INTERSPEECH, Brighton,
38 UK, 2009.
39
40
- 41 [49] H. Liao, M. J. F. Gales, Issues with uncertainty decoding for noise robust
42 automatic speech recognition, *Speech Communication* 50 (4) (2008) 265–
43 277.
44
45
- 46 [50] J. F. Gemmeke, U. Remes, K. J. Palomäki, Observation uncertainty
47 measures for sparse imputation, in: Submitted to INTERSPEECH 2010.
48
- 49 [51] J. Barker, L. Josifovski, M. Cooke, P. Green, Soft decisions in miss-
50 ing data techniques for robust automatic speech recognition, in: Proc.
51 ICSLP, Beijing, China, 2000, pp. 373–376.
52
53
54
55
56
57
58
59
60
61
62
63
64
65

- 1
2
3
4
5 [52] J. Gemmeke, B. Cranen, Sparse imputation for noise robust speech
6 recognition using soft masks, in: Proc. ICASSP, Taipei, Taiwan, 2009,
7 pp. 4645–4648.
8
9
10
11
12
13
14
15
16
17
18
19
20
21
22
23
24
25
26
27
28
29
30
31
32
33
34
35
36
37
38
39
40
41
42
43
44
45
46
47
48
49
50
51
52
53
54
55
56
57
58
59
60
61
62
63
64
65