



"This sentence is wrong." Detecting errors in machine-translated sentences.

Sylvain Raybaud, David Langlois, Kamel Smaïli

► To cite this version:

Sylvain Raybaud, David Langlois, Kamel Smaïli. "This sentence is wrong." Detecting errors in machine-translated sentences.. Machine Translation, 2011, 25 (1), p. 1–34. 10.1007/s10590-011-9094-9 . hal-00606350

HAL Id: hal-00606350

<https://hal.science/hal-00606350>

Submitted on 6 Jul 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.

“This sentence is wrong.”

Detecting errors in machine-translated sentences

Sylvain Raybaud (sylvain.raybaud@loria.fr), David Langlois (david.langlois@loria.fr) and Kamel Smaili (kamel.smaili@loria.fr)
LORIA, BP 239, 54506 Nancy Cedex, France

Abstract.

Machine translation systems are not reliable enough to be used “as is”: except for the most simple tasks, they can only be used to grasp the general meaning of a text or assist human translators. The purpose of confidence measures is to detect erroneous words or sentences produced by a machine translation system. In this article after reviewing the mathematical foundations of confidence estimation we propose a comparison of several state-of-the-art confidence measures, predictive parameters and classifiers. We also propose two original confidence measures based on Mutual Information and a method for automatically generating data for training and testing classifiers. We applied these techniques to data from WMT campaign 2008 and found that the best confidence measures yielded an Equal Error Rate of 36.3% at word level and 34.2% at sentence level, but combining different measures reduced these rates to respectively 35.0% and 29.0%. We also present the results of an experiment aimed at determining how helpful confidence measures are in a post edition task. Preliminary results suggest that our system is not yet ready to efficiently help post editors, but we now have a software and protocol we can apply to further experiments, and user feedback has indicated aspects which must be improved in order to increase the level of helpfulness of confidence measures.

Keywords: machine translation, confidence measure, translation evaluation, support vector machine, mutual information, partial least squares regression, logistic regression, neural network

1. Introduction

A machine translation system generates the best translation for a given sentence according to a previously learnt or hard-coded model. However no model exists that is able to capture all the subtlety of natural language. Therefore even the best machine translation systems make mistakes, and always will — even experts make mistakes after all. Errors take a variety of forms: a word can be wrong, misplaced or missing. Whole translations can be utterly nonsensical or just slightly flawed — involving missing negation, grammatical error and so forth. Therefore, when a document is intended for publication, machine translation output cannot be used “as is”; at best, it can be used to help a human translator. A tool for detecting and pinpointing translation errors may

© 2011 Kluwer Academic Publishers. Printed in the Netherlands.

ease their work as suggested, for example, in (Ueffing and Ney, 2005). (Gandrabur and Foster, 2003) suggest the use of confidence estimation in the context of translation prediction. Confidence estimates could benefit automatic post edition systems like the one proposed in (Simard et al., 2007), by selecting which sentences are to be post edited. Even end-users using machine translation for grasping the overall meaning of a text may appreciate dubious words and sentences being highlighted, preventing them from placing too much trust in potentially wrong translations.

However, and maybe because of such high expectations, confidence estimation is a very difficult problem because if decisions are to be made based on these estimations (like modifying a translation hypothesis), they need to be very accurate in order to maintain translation quality and avoid wasting the user's time. Confidence estimation remains an active research field in numerous domains and much work remains to be done before they can be integrated into working systems.

This article is an overview of many of today's available predictive parameters for machine translation confidence estimation along with a few original predictive parameters of our own; we also evaluated different machine learning techniques — support vector machines, logistic regression, partial least squares regression and neural networks (Section 2) — to combine and optimise them. An exhaustive review would require a whole book, therefore this paper intends to give a more targeted overview of some of the most significant ideas in the domain. (Blatz et al., 2004) proposed a review of many confidence measures for machine translation. We used this work as a starting point to then carry out a thorough formalisation of the confidence estimation problem and make two contributions to the field:

- Original estimators based on Mutual Information and Part-Of-Speech tags (Section 6).
- An algorithm to automatically generate annotated training data for correct/incorrect classifiers (Section 4.3).

We implemented techniques described in (Siu and Gish, 1999) for the evaluation of the performance of the proposed confidence measures. In Sections 6.2, 7.2, we show that using a combination of all predictive parameters yields an improvement of 1.3 points absolute in terms of equal error rate over the best parameter used alone (Section 3.1).

Finally, we present the results of a post edition experiment in which we asked volunteers to correct sentences which had been automatically translated and measure their efficiency with and without confidence measures (Section 8). Unfortunately the results suggested that this

confidence estimation system was not ready to be included in a post edition software yet. However we shall provide a number of useful observations and indications about what is wrong with our system and what is really important for a user.

1.1. SENTENCE-LEVEL CONFIDENCE ESTIMATION

We intuitively recognise a wrong translation that does not have the same meaning than the source sentence, or no meaning at all, or is too disfluent. State-of-the-art natural language processing software being still unable to grasp the meaning of a sentence or to assess its grammatical correctness or fluency, we have to rely on lower level estimators. The problem is also ill-posed: sometimes one cannot decide what is the meaning of a sentence (especially without a context), let alone decide if its translation is correct or not (a translation can be correct for one of the possible meanings of the source sentence and wrong for another). In our experiments we asked human judges to assign a numerical score to machine translated sentences, ranging from one (hopelessly bad translation) to five (perfect) as described in Section 4.1. We set the confidence estimation system to automatically detect sentences with scores of three or higher (disfluencies are considered acceptable, insofar as a reader could understand the correct meaning in a reasonable time). To this end we computed simple numeric features (also called *predictive parameters*: Language Model (LM) score, length, etc. — Section 7) and combined them (Section 2).

1.2. WORD-LEVEL CONFIDENCE ESTIMATION

Defining the correctness of a word is even more tricky. Sometimes a translated word may not fit in the context of the source sentence, as may be the case when homonyms are involved (for example if the French word *vol*, speaking of a plane, is translated with the English word *theft* instead of *flight*). In this case the error is obvious but sometimes the correctness of a word could depend on how other words around it are translated. Consider the following example:

$$\text{Ces mots sont presque synonymes} \rightarrow \left\{ \begin{array}{l} 1: \text{These words are almost} \\ \text{synonyms (*correct*)} \\ 2: \text{These words have close} \\ \text{meanings (*correct*)} \\ 3: \text{These words have close} \\ \text{synonyms (*incorrect*)} \end{array} \right.$$

#3 is definitely an incorrect translation but then we have to decide which word is wrong: is it *close*? is it *synonyms*? or *have*? all of them? In the rest of the article we show how we trained classifiers to discriminate between correct and incorrect words, but this example shows that no system can ever achieve perfect classification simply because this does not exist.

1.3. MATHEMATICAL FORMULATION

Let us now state the problem in mathematically sound terms: the goal of machine translation is to generate a target sentence from a source sentence. A sentence is a finite sequence of words and punctuation marks, which are elements of a set called *vocabulary*. The sentences are represented by random variables. We use the following conventions: a random variable will be represented by a capital letter and a realisation of the variable by the corresponding lower-case letter; bold letters are non-scalar values (sentences, vectors, matrices); non-bold letters are for scalar values like words and real numbers; cursive letters are sets.

\mathcal{V}_S	:	Source language vocabulary
\mathcal{V}_T	:	Target language vocabulary
$\mathbf{S} \in \mathcal{V}_S^*$:	Sentence in the source language
$\mathbf{T} \in \mathcal{V}_T^*$:	Sentence in the target language

From these two primary random variables we then derive new variables:

$Len(\mathbf{S}) \in \mathbb{N}$:	Length of \mathbf{S} (number of words)
$Len(\mathbf{T}) \in \mathbb{N}$:	Length of \mathbf{T}
$S_i \in \mathcal{V}_S$:	i-th word of \mathbf{S}
$T_j \in \mathcal{V}_T$:	j-th word of \mathbf{T}

When estimating confidence we are given realisations of these variables and then need to guess the values of:

$C_{\mathbf{S}, \mathbf{T}} \in \{0, 1\}$:	correctness of a sentence \mathbf{T} as a translation of \mathbf{S}
$C_{\mathbf{S}, \mathbf{T}, j} \in \{0, 1\}$:	correctness of the j-th word of \mathbf{T}

To this end the following probability distribution functions (PDFs) are required and need to be estimated:

$$P(C_{\mathbf{S}, \mathbf{T}} = 1 | \mathbf{S}, \mathbf{T}) : \text{the probability that } \mathbf{T} \text{ is a correct} \quad (1) \\ \text{translation of } \mathbf{S}$$

$$P(C_{\mathbf{S}, \mathbf{T}, j} = 1 | \mathbf{S}, \mathbf{T}) : \text{the probability of correctness of the j-th} \quad (2) \\ \text{word of } \mathbf{T} \text{ given that } \mathbf{T} \text{ is a translation of } \mathbf{S}$$

As \mathbf{S} and \mathbf{T} may be any sentence, directly estimating these probabilities is impossible. We therefore opted to map the pair (\mathbf{S}, \mathbf{T}) to a vector of d_s numerical features (so-called *predictive parameters* described in Section 7.1) via the function \mathbf{x}^s . Similarly $(\mathbf{S}, \mathbf{T}, j)$ were mapped to a numerical vector of d_w features via \mathbf{x}^w (Section 6.1):

$$\mathbf{x}^s : (\mathbf{S}, \mathbf{T}) \in \mathcal{V}_S^* \times \mathcal{V}_T^* \rightarrow \mathbf{x}^s(\mathbf{S}, \mathbf{T}) \in I\!\!R^{d_s}$$

and:

$$\mathbf{x}^w : (\mathbf{S}, \mathbf{T}, j) \in \mathcal{V}_S^* \times \mathcal{V}_T^* \times \mathbb{N} \rightarrow \mathbf{x}^w(\mathbf{S}, \mathbf{T}, j) \in I\!\!R^{d_w}$$

Such parameters may include, for example, the length of source and target sentences, the score given by a translation model or a language model, etc. The following PDFs are therefore learnt (the left-hand parts are just notations) instead of Formulae (1) and (2):

$$p(C_{\mathbf{S}, \mathbf{T}}; \mathbf{S}, \mathbf{T}) \stackrel{\text{def}}{=} P(C_{\mathbf{S}, \mathbf{T}} | \mathbf{x}^s(\mathbf{S}, \mathbf{T})) \quad (3)$$

$$p(C_{\mathbf{S}, \mathbf{T}, j}; \mathbf{S}, \mathbf{T}, j) \stackrel{\text{def}}{=} P(C_{\mathbf{S}, \mathbf{T}, j} | \mathbf{x}^w(\mathbf{S}, \mathbf{T}, j)) \quad (4)$$

Note that although it does not explicitly appear in the notation, p depends on the function \mathbf{x} , which will be different in different experiments, and will also not be the same on sentence- and word-levels. These distributions were to be learnt on large data sets (described in Section 4) by standard machine learning algorithms (Section 2) such as Support Vector Machines (Cortes and Vapnik, 1995), Neural Networks (Fausett, 1994), Logistic Regression (Menard, 2002) or Partial Least Squares Regression (Tobias, 1995).

1.3.1. Classification

After this training process the probability estimates (Formulae (3) and (4)) could be used as confidence measures. It was then possible to compute a classification:

$$\hat{c} : (\mathbf{T}, \mathbf{S}) \rightarrow \hat{c}(\mathbf{T}, \mathbf{S}) \in \{0, 1\}$$

or at word-level::

$$\hat{c} : (\mathbf{T}, \mathbf{S}, j) \rightarrow \hat{c}(\mathbf{T}, \mathbf{S}, j) \in \{0, 1\}$$

In order to minimise the number of errors, classification needs to be performed according to:

$$\hat{c}(\mathbf{T}, \mathbf{S}) \stackrel{\text{def}}{=} \arg \max_{c \in \{0, 1\}} p(c; \mathbf{S}, \mathbf{T}) \quad (5)$$

$$\hat{c}(\mathbf{T}, \mathbf{S}, j) \stackrel{\text{def}}{=} \arg \max_{c \in \{0, 1\}} p(c; \mathbf{S}, \mathbf{T}, j) \quad (6)$$

However this is too strict and neither accounts for biased probability estimates nor permits the attribution of levels of importance to correct rejection or correct acceptance — that is, correct detection of good translations versus correct detection of erroneous translations (see performance estimation in Section 3). Therefore we introduced an *acceptance threshold* δ :

$$\hat{c}(\mathbf{T}, \mathbf{S}; \delta) \stackrel{\text{def}}{=} \begin{cases} 1 & \text{if } p(1; \mathbf{S}, \mathbf{T}) \geq \delta \\ 0 & \text{otherwise} \end{cases} \quad (7)$$

$$\hat{c}(\mathbf{T}, \mathbf{S}, j; \delta) \stackrel{\text{def}}{=} \begin{cases} 1 & \text{if } p(1; \mathbf{S}, \mathbf{T}, j) \geq \delta \\ 0 & \text{otherwise} \end{cases} \quad (8)$$

If $\delta = 0.5$ Formulae (7) and (8) are equivalent to (5) and (6). But setting a higher δ may compensate for a positive bias in probability estimates (3) and (4) or penalise false acceptances more heavily, while setting a lower δ compensates for a negative bias or penalise false rejections more heavily.

1.3.2. Bias

Probability estimates of Formulae (3) and (4) are often biased. This generally does not harm classification performance for two reasons. Firstly, when the bias is uniform ($p^* = \tilde{p} + b$ where b is constant), removing the bias is equivalent to setting an appropriate acceptance threshold. Secondly and most importantly, these PDFs are learnt by minimising classification cost. It is therefore not surprising that even if the probabilities are biased and even if the bias is not uniform ($p^* = \hat{p} + b(\hat{p})$), positive examples generally obtain a higher probability than negative ones.

However, biased probability estimates can harm other performance metrics and in particular will definitely harm Normalised Mutual Information (Section 3) as shown in (Siu and Gish, 1999). We therefore estimated bias on a separate corpora as explained in Siu's paper. The interval $[0, 1]$ was split into 1000 non overlapping bins \mathcal{B}_i of uniform width and bias was estimated separately on each of them.

$$\forall i \in \{1, \dots, 1000\} . \ b(\mathcal{B}_i) = \frac{\sum_{j|\hat{p}_j \in \mathcal{B}_i} (\hat{p}_j - c_j^*)}{\sum_{j|\hat{p}_j \in \mathcal{B}_i} 1} \quad (9)$$

where \hat{p}_j are the estimated probabilities of correctness of items in the training set dedicated to bias estimation, and c_j^* their true classes. Then we obtained an unbiasing function:

$$\text{if } p \in \mathcal{B}_i : \text{unbias}(p) = p - b(\mathcal{B}_i) \quad (10)$$

If \hat{p} is the probability of correctness estimated by a confidence measure, we chose to use the unbiased estimation in our applications:

$$p(1; \mathbf{S}, \mathbf{T}) = \text{unbias}(\hat{p})$$

1.3.3. Sentence quality assessment

Some applications do not require the classification of sentences as correct or incorrect, but rather the estimation of overall quality of the translation. This would ressemble BLEU score (Papineni et al., 2001) or Translation Edit Rate (Snover et al., 2006) only without using reference translations. In this case a quality metric is more suitable than a correctness probability. In Section 7 we therefore present a method for learning PDF of Formula (3) which can also perform regression against quality scores. The training set for this task was:

$$\{(\mathbf{x}^s(\mathbf{s}^n, \mathbf{t}^n); q_{\mathbf{s}^n, \mathbf{t}^n}^*)\}_{n=1..N} \subset \mathbb{R}^{d_s} \times \mathbb{R}^+$$

where $q_{\mathbf{s}^n, \mathbf{t}^n}^*$ is a score relying on expert knowledge. This can be a human evaluation, or a metric computed by comparing the sentence to expert given references, like Word Error Rate, BLEU or Translation Edit Rate. The goal is to learn the mapping $f_\Theta : \mathbb{R}^{d_s} \rightarrow \mathbb{R}^+$ minimising the mean quadratic error using regression techniques (linear regression, support vector regression, partial least squares regression...) where Θ is a set of parameters to be estimated by regression:

$$\frac{1}{N} \sum_{n=1}^N |f_\Theta(\mathbf{x}^s(\mathbf{s}^n, \mathbf{t}^n)) - q_{\mathbf{s}^n, \mathbf{t}^n}^*|^2 \quad (11)$$

1.3.4. Training sets

PDFs (Equations 3 and 4) and regression parameters Θ in Equation (11) need to be learnt using large data sets. Such data sets consist of:

- N source sentences $\mathbf{s}^1, \dots, \mathbf{s}^N$ which are realisations of \mathbf{S} .
- The corresponding N automatically translated sentences $\mathbf{t}^1, \dots, \mathbf{t}^N$ which are realisations of \mathbf{T} .
- Reference sentences classes and a quality metric $((c_{\mathbf{s}^1, \mathbf{t}^1}^*, q_{\mathbf{s}^1, \mathbf{t}^1}^*), \dots, (c_{\mathbf{s}^N, \mathbf{t}^N}^*, q_{\mathbf{s}^N, \mathbf{t}^N}^*))_{n=1..N} \in (\{0, 1\} \times \mathbb{R}^+)^N$ which are

realisations of $C_{\mathbf{S}, \mathbf{T}}$; they can be given by human experts (Section 4.1) or automatically generated from human translations (Section 4.3 and 4.2).

- Reference words classes

$\forall n \in \{1, \dots, N\}. (c_{\mathbf{s}^n, \mathbf{t}^n, 1}^*, \dots, c_{\mathbf{s}^n, \mathbf{t}^n, Len(\mathbf{t})}^*) \in \{0, 1\}^{Len(\mathbf{t})}$ which are realisations of $C_{\mathbf{S}, \mathbf{T}, j}$: also given by human experts.

2. Classification and regression techniques

The problem of confidence estimation is now reduced to standard classification and regression problems. Many well known machine learning techniques are available and we opted to experiment with the well known Support Vector Machines, Logistic Regression and Artificial Neural Network, and also with the lesser-known Partial Least Squares Regression.

2.1. LOGISTIC REGRESSION

Here we wanted to predict the correctness $C \in \{0, 1\}$ given a set of features $\mathbf{X} \in \mathbb{R}^d$; to this end we needed to estimate the distribution $P(C = 1 | \mathbf{X})$. Logistic Regression (Menard, 2002) consists of assuming that:

$$P(C = 1 | \mathbf{X}) = \frac{1}{1 + e^{\langle \Theta, \mathbf{X} \rangle + b}} \quad (12)$$

for some $\Theta \in \mathbb{R}^d$ and $b \in \mathbb{R}$ and then optimise Θ with regard to the maximum likelihood criterion on the training data. Logistic regression was used not only to combine several features but also to map the scores produced by a confidence estimator to a probability distribution.

2.2. SUPPORT VECTOR MACHINE

The well-known Support Vector Machines (SVMs) (Hsu et al., 2003) have highly desirable characteristics which made them well suited to our problem. They are able to discriminate between two non-linearly separable classes; they can also compute the probability that a given sample belongs to one class, and not only a binary decision and they can also be used to perform regression against numerical scores (Smola and Schölkopf, 2004). We used LibSVM (Chang and Lin, 2011) for feature scaling, classification and regression.

2.2.1. SVM for classification

SVM were trained to produce a probability of correctness. By doing so the acceptance threshold could be adapted (Section 1.3 and Equations (7) and (8)), making the classifier more flexible. The kernel used was a Radial Basis Function since it is simple and was reported in (Zhang and Rudnicky, 2001) as giving good results:

$$K_\gamma(\mathbf{x}(\mathbf{s}, \mathbf{t}, j), \mathbf{x}(\mathbf{s}', \mathbf{t}', j')) = e^{-\gamma \|\mathbf{x}(\mathbf{s}, \mathbf{t}, j) - \mathbf{x}(\mathbf{s}', \mathbf{t}', j')\|^2}$$

2.2.2. SVM for quality evaluation

The same kernel was used but this time to perform regression against sentence-level BLEU score (Papineni et al., 2001).

2.2.3. Meta-parameters optimisation

SVMs require two meta-parameters to be optimised: the γ parameter of the radial basis function, and the error cost C . γ and C were optimised by grid search on the development corpus with regard to equal error rate for classification, and mean quadratic error for regression.

2.3. NEURAL NETWORK

The FANN toolkit (Fast Artificial Neural Network (Nissen, 2003)) is used for building feed-forward neural networks (NN). After experimenting on a development set we decided to stick to the standard structure, namely one input layer with as many neurons as we have features, one hidden layer with half as many neurons, and an output layer made of a single neuron which returns a probability of correctness. The connection rate was 0.5 in order to keep computation time tractable. We stuck to the default sigmoid activation function. The weights were optimised by standard gradient back-propagation.

2.4. PARTIAL LEAST SQUARES REGRESSION

Partial Least Squares Regression (Wold et al., 1984; Specia et al., 2009) is a multivariate data analysis technique that finds a bilinear relation between the observable variables (our features \mathbf{X} and the response variables, namely the probability of correctness $p(1; \mathbf{X})$ or the quality score). It works by projecting both predictors and observations on a linear subspace and performs least-squares regression in this space. It has the major advantage of being robust to correlated predictors.

3. Evaluation of the classifiers

Error rate is the most obvious metric for measuring the performance of a classifier. It is however not an appropriate metric because of its sensitivity to class priors (Kononenko and Bratko, 1991; Siu and Gish, 1999). Let us exemplify the problem and consider for example a machine translation system which gives roughly 15% of wrongly translated words. Now let us consider a confidence measure such that:

$$\forall \mathbf{s}, \mathbf{t}, j \quad p^0(1; \mathbf{s}, \mathbf{t}, j) = 1$$

It makes no error on correct words (85% of total) but misclassifies all wrong words (15%). Its error rate is therefore $0 \times 0.85 + 1 \times 0.15 = 0.15$. Now let us consider a second confidence measure $p^1(1; \mathbf{s}, \mathbf{t}, j)$ which correctly detects every wrong word (if the j -th word of \mathbf{t} is wrong then $p^1(1; \mathbf{s}, \mathbf{t}, j) = 0$) but also incorrectly assigns a null probability of correctness to 20% of the words that are appropriate translations. The error rate of this measure is: $0 \times 0.15 + 0.20 \times 0.85 = 0.17$.

p^0 therefore seems to outperform p^1 . This is however not true, because p^0 does not provide the user with any useful information (or actually any information at all, strictly speaking), while if $p^0(1; \mathbf{s}, \mathbf{t}, j) > 0$ then we would be certain that the word is correct. There is a lesson here. An appropriate metric for the usefulness of a confidence measure is not the number of misclassifications it makes but *the amount of information it provides*. This is why we opted to use *Normalised Mutual Information* (Siu and Gish, 1999) to assess the performance of a measure (Section 3.2), along with Equal Error Rate (EER) and *Discrimination Error Trade-off* (DET) curves (Section 3.1). The latter is a powerful tool for the visualisation of the behaviour of a classifier with different acceptance thresholds and therefore different compromises between incorrect acceptances and incorrect rejections.

3.1. DISCRIMINATION ERROR TRADE-OFF

A classifier makes two kinds of mistakes: *False acceptance*, when an erroneous item (word or sentence) is classified as correct, also called *Type 1 error*, and *False rejection* or *Type 2 error* when a correct item is classified as incorrect. When evaluating the performance of a classifier we know the predictions \hat{c} (equations (7) and (8)) and the *actual* realisations c^* of the variables \mathbf{C} . As stated above in Section 1.3 $\hat{c}(\mathbf{t}; \mathbf{s}; \delta)$ is the *estimated* correctness of translation \mathbf{t} given source sentence \mathbf{s} with acceptance threshold δ and that $c_{\mathbf{s}, \mathbf{t}}^*$ is the *true* (expert-given) correctness (Section 1.3.4). Sentence-level false acceptance rate is:

$$e_1(\mathbf{s}, \mathbf{t}; \delta) = \begin{cases} 1 & \text{if } \hat{c}(\mathbf{t}; \mathbf{s}; \delta) = 1 \text{ and } c_{\mathbf{s}, \mathbf{t}}^* = 0 \\ 0 & \text{otherwise} \end{cases} \quad (13)$$

$$err_1(\delta) = \frac{\sum_{\mathbf{s}, \mathbf{t}} e_1(\mathbf{s}, \mathbf{t}; \delta)}{\sum_{\mathbf{s}, \mathbf{t}} (1 - c_{\mathbf{s}, \mathbf{t}}^*)} \quad (14)$$

err_1 is therefore the proportion of wrong items which are incorrectly accepted ($\sum_{\mathbf{s}, \mathbf{t}} (1 - c_{\mathbf{s}, \mathbf{t}}^*)$ is the number of wrong items). Sentence-level false rejections rate is:

$$e_2(\mathbf{s}, \mathbf{t}; \delta) = \begin{cases} 1 & \text{if } \hat{c}(\mathbf{t}; \mathbf{s}; \delta) = 0 \text{ and } c_{\mathbf{s}, \mathbf{t}}^* = 1 \\ 0 & \text{otherwise} \end{cases} \quad (15)$$

$$err_2(\delta) = \frac{\sum_{\mathbf{s}, \mathbf{t}} e_2(\mathbf{s}, \mathbf{t}; \delta)}{\sum_{\mathbf{s}, \mathbf{t}} c_{\mathbf{s}, \mathbf{t}}^*} \quad (16)$$

err_2 is the proportion of correct items which are rejected by the classifier. Adapting these formulae to word-level is straightforward.

Intuitively err_1 is the proportion of erroneous words that the classifiers wrongly accept, while err_2 is the proportion of correct words that the classifier wrongly rejects. A relaxed classifier has a low err_2 and a high err_1 while a strict one has a low err_1 and a high err_2 . Proof that err_1 and err_2 are insensitive to class priors was given in (Siu and Gish, 1999).

When δ goes from 0 to 1, more and more items are rejected. Therefore the false rejection rate (err_2) monotonically increases from 0 to 1 while the false acceptance rate (err_1) monotonically decreases from 1 to 0. The plot of $err_1(\delta)$ against $err_2(\delta)$ is called the *DET curve* (*Discrimination Error Trade-off*). See examples in Section 6.

A lower curve indicates a better classifier. All points of the DET curve should lie below the diagonal $[(0, 1), (1, 0)]$, which is the theoretical curve of a classifier using features uncorrelated with correctness (that is, inappropriate features).

Both err_1 and err_2 are generally approximations of continuous functions¹. Therefore a threshold δ_{EER} exists such that:

$$err_1(\delta_{EER}) \approx err_2(\delta_{EER}) = EER \quad (17)$$

EER is called the *equal error rate*. It can be seen as a “summary” of the DET curve when the acceptance threshold is set so that there are

¹ it actually depends on the true and estimated PDFs. When this is not the case they shall be approximated by continuous functions.

the same proportions of Type 1 and 2 errors, and can be used for direct comparisons between classifiers. However this is arbitrary, because the user may prefer to have fewer errors of one type, at the cost of more of the other type.

3.2. NORMALISED MUTUAL INFORMATION

Normalised Mutual Information (NMI) measures the level of informativeness of a predictive parameter or a set thereof in an application-independent manner (Siu and Gish, 1999). Intuitively NMI measures the reduction of entropy of the distribution of true class C over the set {“correct”, “incorrect”} when the value of the predictive parameter is known; let $\mathbf{x}(\mathbf{S}, \mathbf{T})$ be a vector of predictive parameters:

$$\begin{aligned} NMI(C, \mathbf{x}) &= \frac{I(C; \mathbf{x})}{H(C)} = \frac{H(C) - H(C|\mathbf{x})}{H(C)} \quad (18) \\ H(C) &= -p^* \log(p^*) - (1 - p^*) \log(1 - p^*) \\ H(C|\mathbf{x}) &= \int_{\mathbf{v}} \left(P(\mathbf{x}(\mathbf{S}, \mathbf{T}) = \mathbf{v}) \times \right. \\ &\quad \left. \sum_{c \in \{0,1\}} P(C = c | \mathbf{x}(\mathbf{S}, \mathbf{T}) = \mathbf{v}) \log(P(C = c | \mathbf{x}(\mathbf{S}, \mathbf{T}) = \mathbf{v})) \right) d\mathbf{v} \end{aligned}$$

where I is mutual information, H is entropy and p^* is the true prior probability of correctness. Because the true distribution $P(\mathbf{x}(\mathbf{S}, \mathbf{T}))$ is replaced with empirical frequencies observed in data, and $P(C|\mathbf{x}(\mathbf{S}, \mathbf{T}))$ is replaced with the computed estimation:

– Sentence-level NMI:

$$\begin{aligned} H(C|\mathbf{x}) &\simeq \frac{1}{N} \sum_{(\mathbf{s}, \mathbf{t}) \in \mathcal{S}} \left(p(1; \mathbf{s}, \mathbf{t}) \log(p(1; \mathbf{s}, \mathbf{t})) \right. \\ &\quad \left. + (1 - p(1; \mathbf{s}, \mathbf{t})) \log(1 - p(1; \mathbf{s}, \mathbf{t})) \right) \quad (19) \end{aligned}$$

– Word-level NMI:

$$\begin{aligned} H(C|\mathbf{x}) &\simeq \frac{1}{N_w} \sum_{(\mathbf{s}, \mathbf{t}) \in \mathcal{S}} \sum_{j=1}^{Len(\mathbf{t})} \left(p(1; \mathbf{s}, \mathbf{t}, j) \log(p(1; \mathbf{s}, \mathbf{t}, j)) \right. \\ &\quad \left. + (1 - p(1; \mathbf{s}, \mathbf{t}, j)) \log(1 - p(1; \mathbf{s}, \mathbf{t}, j)) \right) \quad (20) \end{aligned}$$

$H(C|\mathbf{x})$ can never be lower than 0 and equality is achieved when for all pair of sentences (or all words within), $p(c_{\mathbf{s},\mathbf{t}}; \mathbf{s}, \mathbf{t}) = 1$ which means that the true class is predicted with no uncertainty. On the other hand $H(C|\mathbf{x})$ can never be greater than $H(C)$ and equality is achieved when the predictive parameters are completely useless. Therefore $M(\mathbf{x})$ is theoretically a real number between 0 and 1. However the approximation of $H(C|\mathbf{x})$ can be negative in practice.

4. Training and testing data

Large data sets are needed to learn PDFs of Formulae (3) and (4). Ideally a human professional translator would read the output of a Machine Translation system and assign a label (*correct* or *incorrect*) to each item. This method would give high quality training data but would be extremely expensive. Therefore it would be preferable to use automatic or semi-automatic methods for efficiently classifying words and sentences. In the following we will discuss different methods for obtaining labelled data.

4.1. EXPERT-ANNOTATED CORPORA

This is the high-quality-high-cost whereby human experts analyse translations produced by a machine translation system and decide if each word and sentence is correct or not. The classification depends on the application, but in our setting a word is classified as erroneous if it is an incorrect translation, if it suffers from a severe agreement error or if it is completely misplaced. A sentence is considered wrong if it is not clear that it has the same meaning as the sentence of which it is supposed to be a translation, or any meaning at all, or if it contains a significant level of ambiguity that was not apparent in the source sentence. This method has two major drawbacks. The first is that it is extremely slow and therefore expensive, and the second is that it is not reproducible because a given sentence may be differently classified by different translators, or by the same translator at different times.

We needed a small corpus of real, expert-annotated machine-translated sentences for our test set. To this end we set up the statistical machine translation system described as the baseline for WMT08 evaluation campaign following the instructions on StatMT website²: it features a 5-gram language model with Kneser-Ney discounting trained with SRILM (Stolcke, 2002) on about 35 million running words, IBM-5 translation model trained on around 40 million words, and Moses

² <http://statmt.org/wmt08/baseline.html>

(Koehn et al., 2007). A hold-out set of 40,000 sentence pairs was extracted from data for the purpose of training the confidence estimation system. We annotated a small set of 150 automatically translated sentences from transcriptions of news broadcast. Because of these sentences' spontaneous style and a vocabulary which did not match that of training corpora (European Parliament) the BLEU score is not high (21.8 with only one reference). However most translations were intelligible when given some thought.

A word was annotated as “incorrect” if it was completely irrelevant, very misplaced or grammatically flawed. Sentences were given scores ranging from one (hopelessly flawed) to five (perfect). For classification purposes we considered sentences scoring three or higher (possible to get the correct meaning when given a little thought) to be correct.

Here are a few examples of expert-annotated sentences (the incorrect words are underlined):

<i>Source sentence</i>	<i>Machine translation</i>	<i>score</i>
je vous remercie monsieur le commis-saire pour votre déclaration.	thank you mr commissioner for your <u>question</u> .	2
j'ai de nombreuses questions à poser à m. le commissaire.	i have <u>some</u> questions to ask to the commissioner.	4
les objectifs de la stratégie de lisbonne ne sont pas les bons.	the lisboa strategy mistaken.	3

4.2. AUTOMATICALLY ANNOTATED CORPORA

An intuitive idea is to compare a generated translation to a reference translation, and classify as correct the candidate words that are Levenshtein-aligned to a word in the reference translation (Ueffing and Ney, 2004). However this is too strict and many correct words would be incorrectly classified, because there are often many possible translations for a given source sentence and these may have nothing in common. This problem can partly be overcome by using multiple reference translations (Blatz et al., 2004). However multiple references are not always available and are costly to produce.

4.3. ARTIFICIAL TRAINING DATA

In this section we present an algorithm aimed at getting the best of both worlds, namely automatically generating sentences (no humans involved, quickly generating huge amounts of data as with automatic annotation), without any annotation error (no errors in gold standard classes as with human annotation). Our objective was to generate enough data for training classifiers in order to combine several predictive parameters.

Starting from human-made reference translations, errors were automatically introduced in order to generate examples for training confidence measures. Given an English sentence \mathbf{t} (which is a correct translation of source sentence \mathbf{s}), we first chose where to introduce errors. As machine translation errors tend to be “bursty” (not evenly distributed but appearing in clusters) we implemented two error models whose parameters were estimated on a few human annotated sentences. These annotations were not required to be extremely precise.

Bigram error model: firstly we implemented a simple bigram model $P(C_i|C_{i-1})$; the probability that a word is correct given the correctness of the preceding word. The first word in a sentence has an a priori probability of being correct. According to this model we generated sequences of ones and zeroes corresponding to correct and incorrect words. We found that nine sentences out of ten in our human annotated test set started with a correct word, that a wrong word had approximately a 50% chance of being followed by another wrong word ($P(C_i = 0|C_{i-1} = 0) \simeq 0.5$) and that a correct word had approximately a 90% chance of being followed by another correct word ($P(C_i = 1|C_{i-1} = 1) \simeq 0.9$).

Cluster error model: the second explicitly models clusters. A sentence is a sequence of clusters of correct words and clusters of incorrect words: $\mathcal{C}_1, \dots, \mathcal{C}_n$. By definition if a cluster contains correct words, the next cluster will contain incorrect words and vice versa. Let C_i be the correctness of words in the i -th cluster. $P(\text{length}(\mathcal{C}_i)|C_i = 0)$ and $P(\text{length}(\mathcal{C}_i)|C_i = 1)$ were estimated on a hold-out set of 50 machine translations annotated by a human. Sequences of zeroes and ones were generated accordingly. This model’s parameters cannot theoretically be represented by a finite set of real number (they are distributions over \mathbb{N}). In practice, cluster lengths are bounded therefore these distributions are actually over $\{0, \dots, \text{max}(\text{length}(\mathcal{C}_i))\}$. Just to give an idea, we found that the average length of a cluster of wrong words was 1.9 ($(\sum_{k \geq 1} k \times P(\text{length}(\mathcal{C}_i) = k|C_i = 0) = 1.9$), and that of a cluster of correct words is 12.2.

Once the exact location of errors was known, we randomly introduced errors of five types: **move**, **deletion**, **substitution**, **insertion** and **grammatical error**. “Deletion” is straightforward: a word is chosen randomly according to a uniform distribution and deleted. “Move” is not much more complicated: a word is chosen at random according to the uniform distribution, and the distance it will be moved (jump length) is chosen according to a probability which is uniform

within a given range (4 in our experiments) and null beyond. “Grammatical” errors are generated by modifying the ending of randomly selected words (“preserving” may become “preserved”, “environment” may become “environmental”). “Substitution” and “insertion” are a little more subtle. Given the position i of the word to be replaced or inserted, the probability of every word in the vocabulary was computed using an IBM-1 translation model and a 5-gram language model:

$$\forall t' \in \mathcal{V}_T . p(t') = p_{IBM-1}(t'|\mathbf{s}_1^I) \times p_{5-gram}(t'|t_{i-4}, \dots, t_{i-1})$$

The new word t' was then picked among all w at random according to the distribution p . This way the generated errors were not too “silly”. WordNet (Miller, 1995) was used to check that t' was not a synonym of t (otherwise it would not be an incorrect word): t' could not belong to any synset of which t is an element. The algorithm was controlled by several parameters, which were empirically chosen:

- probability distribution P_m of the proportion of move errors in a sentence and probability distribution P_j of jump length
- probability distribution P_d of the proportion of deletions
- probability distribution P_s of the proportion of substitutions
- probability distribution P_i of the proportion of insertions
- probability distribution P_g of the proportion of grammatical errors

We chose triangle shaped distributions with $mode = 0.2$, $minimum = 0$ and $maximum = 0.5$. These may not be the real distributions but seemed reasonable. The positions of words to be moved, deleted, inserted or modified were chosen according to uniform distribution probability. For each sentence errors were inserted in the order given previously — firstly, words were moved, then some were deleted, etc. Eventually we obtained a corpus with an average 16% word error rate, which approximately matches the error rate of real machine translation output.

Below is an example of degraded translation obtained using this method, extracted from our corpus:

<i>source sentence</i>	Quant à eux, les instruments politiques doivent s’adapter à ces objectifs.
<i>reference translation</i>	Policy instruments, for their part, need to adapt to these goals.
<i>degraded translation</i>	Policy instruments, for the part, must to adapt to these goals.

We used 40,000 pairs of sentences (source: French - target: English) from WMT-2008 evaluation campaign data. We degraded the reference translations according to the above rules. We found that the *bigram error model* gave the best results in the end (classification error rates of confidence measures trained on such data are lower) so we used it for all experiments presented here. The BLEU score of the degraded corpus was 56.5 which is much higher than the score of our baseline described in Section 4.1 (21.8). The latter score was underestimated because only one reference translation was available. However this phenomenon did not affect the BLEU score of the degraded corpus as it came directly from the reference sentences and therefore there was no need for multiple references. The error rate in the degraded corpus was set to 16% to match that of real machine translation output.

Others have proposed the use of artificial corpora, for example (Blatz et al., 2004) and (Quirk, 2004). While we found that automatically generated corpora yield performances comparable to that of expert-annotated ones (Section 6.2), the later draw conclusions opposed to ours, as they found that a classifier trained on a small, human-annotated corpora, performs better than one trained on a large automatically annotated corpora. However in their experiment sentences are not automatically generated but automatically annotated. It is important to understand that automatically generated data is not the same as automatic annotation. In the latter, sentences are realistic but there is uncertainty concerning annotation. On the other hand, while automatically degraded sentences may seem less realistic, there is almost no doubt that words labelled as incorrect are actually wrong, and vice versa. Therefore automation plays a completely different role in their system and ours. Another difference is that they are evaluating sentences, while an important task for us is to evaluate words. In Section 6.2 we present an experiment showing that a classifier trained on our large artificial corpus yields better results than one trained on a small human annotated corpus (Figure 4), for a fraction of the cost.

5. Experimental framework

A single feature (for example, n -gram probability) can be used as a confidence score. It is then relatively simple to evaluate its performance because no neural network or similar machine learning machinery is necessary. Each word or sentence is attributed a score and a DET curve can be immediately computed. Computing NMI is a slightly more subtle operation because a probability is needed here, and not all predictive parameters qualify as such. In this case the score is turned

into a probability by logistic regression (Section 2.1) whose parameters are learnt from artificial data.

Combining several predictive parameters is a little more complicated. Unless otherwise specified we proceeded as follows: two artificial corpora \mathcal{T}_1 (for “training”) and \mathcal{D} (“development”) were used to find the best meta-parameters with regard to EER for SVM (γ and C — see Section 2.2) and Neural Networks (number of hidden units — Section 2.3). Once optimal meta-parameters were found (or if none was set) the classifier was trained on a larger set of automatically generated data \mathcal{T}_2 and finally tested on unseen real machine translation output \mathcal{U} . Then, if relevant, bias was estimated on a corpus of automatically generated data \mathcal{B} . $\mathcal{T}_1, \mathcal{T}_2, \mathcal{D}$ and \mathcal{B} consisted of 10,000 sentences each (around 200,000 words). \mathcal{U} consisted of 150 sentences, or approximately 3,000 words, each of them having one reference translation (Section 4.1).

6. Word-level confidence estimation

We shall now look into the details of the predictive parameters we used (the components of the vector $\mathbf{x}(\mathbf{S}, \mathbf{T}, j)$) for word-level confidence estimation. These components will be noted x_{index} where *index* is the label of the equation so that they are easier to find and refer to in the paper. Altogether these features are a numerical representation of a word in the target language (T_j), its context (the whole sentence \mathbf{T}), and the source sentence \mathbf{S} the translation of which it is supposed to be a part. Of course this representation is less expressive than the original natural words and sentences, but hopefully it is more accessible to probability estimation while still carrying enough information to enable us to determine whether a word is correct or not.

Some of these features can themselves be used as confidence measures (for example LM-based features). In this case, we provided performance evaluation. Others cannot, for example Part-Of-Speech tag, stop word indicator and rule based features.

6.1. FEATURES FOR WORD LEVEL CONFIDENCE ESTIMATION

6.1.1. *N*-gram based features

N-gram scores and backoff behaviour can provide a great deal of useful information. First, the probability of a word in a classical 3-gram language model can be used as the feature:

$$x_{21}(\mathbf{S}, \mathbf{T}, j) = P(t_j | t_{j-1}, t_{j-2}) \quad (21)$$

Intuitively, we would expect an erroneous word to have a lower n -gram probability. However this feature is generally already used in statistical translation systems, therefore even the probability levels of wrong words may not be too low.

Backward 3-gram language models, proposed for speech recognition confidence estimation in (Duchateau et al., 2002), also turned out to be useful:

$$x_{22}(\mathbf{S}, \mathbf{T}, j) = P(t_j | t_{j+1}, t_{j+2}) \quad (22)$$

This feature has the advantage of generally not being used in the decoding process.

Finally the backoff behaviour of the 3-gram and backward 3-gram models are powerful features: an n -gram not found in the language model may indicate a translation error. A score is given according to how many times the LM had to back off in order to assign a probability to the sequence, as proposed in (Uhrik and Ward, 1997) for speech recognition:

$$x_{23}(\mathbf{S}, \mathbf{T}, j) = \begin{cases} 1.0 & \text{if } t_{j-2}, t_{j-1}, t_j \text{ exists in the model} \\ 0.8 & \text{if } t_{j-2}, t_{j-1} \text{ and } t_{j-1}, t_j \text{ both exist in the model} \\ 0.6 & \text{if only } t_{j-1}, t_j \text{ exists in the model} \\ 0.4 & \text{if only } t_{j-2}, t_{j-1} \text{ and } t_j \text{ exist separately in the model} \\ 0.3 & \text{if } t_{j-1} \text{ and } t_j \text{ both exist in the model} \\ 0.2 & \text{if only } t_j \text{ exists in the model} \\ 0.1 & \text{if } t_j \text{ is completely unknown} \end{cases} \quad (23)$$

Figure 1 shows DET curves of the confidence measures based on 3-grams and backward 3-grams and scores and backoff behaviour. While 3-grams and backward 3-grams are almost indistinguishable, backoff behaviour performs better in terms of EER. Although this measure is very simple, it is less correlated with those used in the decoding or degrading process, which may explain why it achieves better discrimination results. The results are summarised in Table I.

Table I. Performances of 3-grams based confidence measures at word level.

feature	equal error rate	normalised mutual information
3-grams	42.1	4.86×10^{-3}
backward 3-grams	42.9	-3.93×10^{-3}
backoff	37.0	6.11×10^{-2}
backward backoff	38.1	1.09×10^{-2}

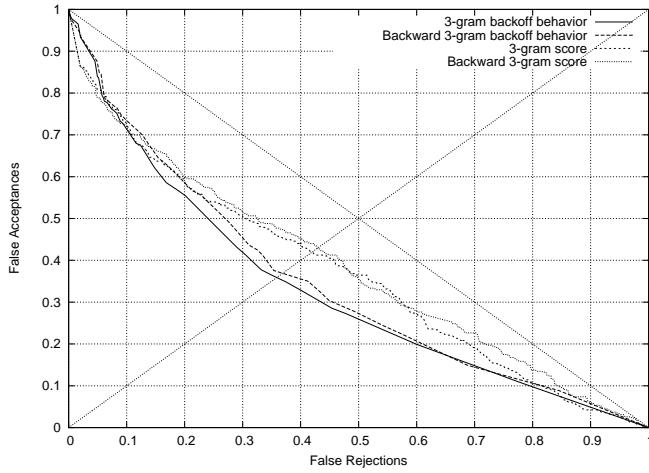


Figure 1. DET curves of 3-grams based confidence measures at word level.

NMI of backward 3-gram scores is negative. This is theoretically not possible but may be explained by a strong bias in the estimation of probabilities which our unbiassing method was unable to efficiently remove (Section 1.3.2) and because NMI was only approximated here (Section 3.2).

6.1.2. Part-Of-Speech based features

Replacing words with their POS class can help to detect grammatical errors, and also to take into account the fact that feature values do not have the same distribution for different word classes. Therefore we used syntactic POS tags as a feature, along with the score of a word in a POS 3-grams model. Tagging was performed using GPoSTTL, an open source alternative to TreeTagger (Schmid, 1994; Schmid, 1995).

$$x_{24}(\mathbf{S}, \mathbf{T}, j) = \text{POS}(t_j) \quad (24)$$

$$x_{25}(\mathbf{S}, \mathbf{T}, j) = P(\text{POS}(t_j) | \text{POS}(t_{j-2}), \text{POS}(t_{j-1})) \quad (25)$$

With our settings, POS is a non numeric features which can take 44 values, say $\{\pi_1, \dots, \pi_{44}\}$. In order to combine it with numeric features it was mapped to a vector $\pi(t_j) \in \{0, 1\}^N$ with $N = 40$, as suggested in (Hsu et al., 2003). The mapping is defined by

$$\pi(t_j)[i] = \begin{cases} 1 & \text{if } \text{POS}(t_j) = \pi_i \\ 0 & \text{otherwise} \end{cases}$$

We have chosen not to show the individual results of these confidence measures as they are only useful in combination with others.

6.1.3. Taking into account errors in the context

A common property of all n -gram based features is that a word can get a low score if it is actually correct but its neighbours are wrong. To compensate for this phenomenon we took the average score of the neighbours of the word being considered into account. More precisely, for every relevant feature x_j defined above ($x_{21}, x_{22}, x_{23}, x_{25}$) we also computed:

$$\begin{aligned}x_j^{\text{left}}(\mathbf{S}, \mathbf{T}, j) &= x_{j-2}(\mathbf{S}, \mathbf{T}) * x_{j-1}(\mathbf{S}, \mathbf{T}) * x_j(\mathbf{S}, \mathbf{T}) \\x_j^{\text{centred}}(\mathbf{S}, \mathbf{T}, j) &= x_{j-1}(\mathbf{S}, \mathbf{T}) * x_j(\mathbf{S}, \mathbf{T}) * x_{j+1}(\mathbf{S}, \mathbf{T}) \\x_j^{\text{right}}(\mathbf{S}, \mathbf{T}, j) &= x_j(\mathbf{S}, \mathbf{T}) * x_{j+1}(\mathbf{S}, \mathbf{T}) * x_{j+2}(\mathbf{S}, \mathbf{T})\end{aligned}$$

These predictive parameters were then combined using a neural network. Figure 2 and Table II show a vast improvement when using the product of 3-gram probabilities of words in the centred window.

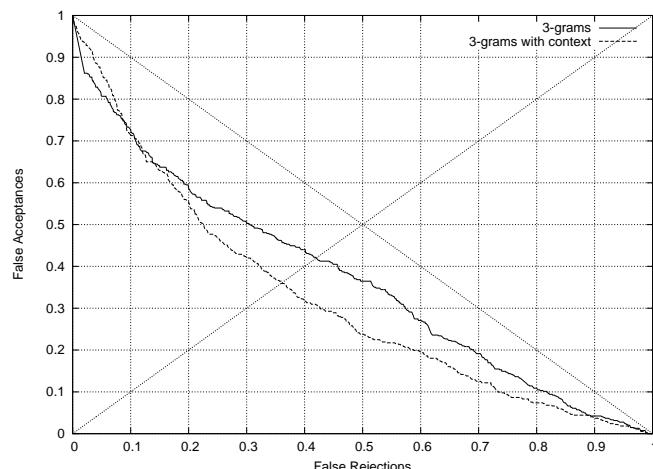


Figure 2. DET curves of 3-grams score combined with neighbours score at word level.

Table II. Influence of taking the context into account.

	equal error rate	normalised mutual information
3-grams	42.1	4.86×10^{-3}
3-grams and neighbours	36.3 (-5.7)	4.57×10^{-3}

However NMI was slightly harmed in the process. This may be because the product of 3-gram scores on the window was not a proper estimation of probability of correctness. However it is perfectly possible to have a confidence measure with good discrimination power and low NMI.

6.1.4. Intra-lingual mutual information

In (Raybaud et al., 2009a; Raybaud et al., 2009b) we introduced original predictive features based on mutual information. Mutual information is a metric for measuring how much information a random variable gives about another. Here we consider two random variables whose realisations are words, say W_1 and W_2 :

$$\begin{aligned} I(W_1, W_2) &= \sum_{w_1, w_2} P(W_1 = w_1, W_2 = w_2) \\ &\quad \times \log \left(\frac{P(W_1 = w_1, W_2 = w_2)}{P(W_1 = w_1)P(W_2 = w_2)} \right) \end{aligned}$$

We used point-wise mutual information which is the contribution of a specific pair of words to the mutual information between W_1 and W_2 (that is, a single term of the sum above).

$$MI(w_1, w_2) = P(W_1 = w_1, W_2 = w_2) \log \left(\frac{P(W_1 = w_1, W_2 = w_2)}{P(W_1 = w_1)P(W_2 = w_2)} \right)$$

The tuple $(w_1, w_2, MI(w_1, w_2))$ is called a *trigger*. Triggers are learnt on an unaltered bilingual corpus. The idea of using mutual information for confidence estimation was first expressed in (Guo et al., 2004). It has since been proved useful for computing translation tables (Lavecchia et al., 2007).

Intra-lingual mutual information (IMI) measures the level of similarity between the words in a generated sentence thus assessing its consistency. Formally W_1 and W_2 are any T_i and T_j here (words of the translation hypothesis). Let J be the length of the translation hypothesis. The feature for confidence estimation is:

$$x_{26}(\mathbf{S}, \mathbf{T}, j) = \frac{1}{J-1} \sum_{1 \leq i \neq j \leq J} MI(t_i, t_j) \quad (26)$$

6.1.5. Cross-lingual mutual information

Cross-lingual mutual information (CMI) is similar to the previous intra-lingual mutual information in that it assesses the source-translation consistency. Let I be the length of the source sentence:

$$x_{27}(\mathbf{S}, \mathbf{T}, j) = \frac{1}{I} \sum_{1 \leq i \leq I} MI(s_i, t_j) \quad (27)$$

Here W_1 and W_2 are any S_i and T_j .

Table III summarises the performances of MI-based features when used as confidence measures by themselves. Although it performs poorly, we will see that they are useful when combined with other predictive parameters (Section 6.2).

Table III. Performances of mutual information based features at word level.

feature	equal error rate	normalised mutual information
intra-lingual	45.8	9.46×10^{-4}
cross-lingual	45.7	-2.21×10^{-1}

6.1.6. IBM-1 translation model:

This feature was proposed in (Blatz et al., 2004; Ueffing and Ney, 2005):

$$x_{28}(\mathbf{S}, \mathbf{T}, j) = \frac{1}{I+1} \sum_{i=0}^I p_{IBM-1}(t_j | s_i) \quad (28)$$

where s_0 is the empty word. The performance of this predictive parameter used alone is given in Table IV. Once again the results are disappointing. The results are extremely similar to alignment probability (the sum is replaced by a max). It is surprising to note that even on a translation evaluation task, measures involving only the hypothesis yield better performances than those taking the source sentence into account.

Table IV. Performance of IBM-1 based confidence measure at word level.

feature	equal error rate	normalised mutual information
IBM-1 score	45.0	-1.84×10^{-3}

Like MI-based features, IBM-1 does not work very well when used as a confidence measure and will only be used in combination with others.

6.1.7. Stop words and rule-based features

The “stop word” predictive parameter is a simple flag indicating if the word is a stop word (*the, it, etc.*) or not. It helps a classifier to take into account the fact that the distribution of other features is not the same on stop words and on content words. This feature is less informative than Part-Of-Speech, but simpler.

$$x_{29}(\mathbf{S}, \mathbf{T}, j) = \begin{cases} 1 & \text{if } t_j \text{ is a stop word} \\ 0 & \text{otherwise} \end{cases} \quad (29)$$

The stop list was generated by picking words that are both short and frequent. Finally we implemented four binary features indicating whether the word is a punctuation symbol, numerical, a URL or a proper name (based on a list of those). These features were not of course designed to be used as standalone confidence measures.

6.2. FEATURES COMBINATION

Altogether we had 66 features for word-level confidence estimations, many of them very similar (for example 3-grams probability and average 3-grams probabilities on different windows), some very crude (for example sentence-level features like length ratio — Section 7.1.5 — used at word level). We trained four classifiers (Logistic Regression, Partial Least Squares Regression, Support Vector Machines and Neural Network) to discriminate between correct and incorrect words based on these features. Only Neural Networks brought a consistent improvement over the best feature used alone (3-gram scores on a centred window, Sections 6.1.1 and 6.1.3) for the classification task, although this was not a large improvement (-1.3 point of EER). The DET curve of neural networks is presented in Figure 3 and the results are summarised in Table V.

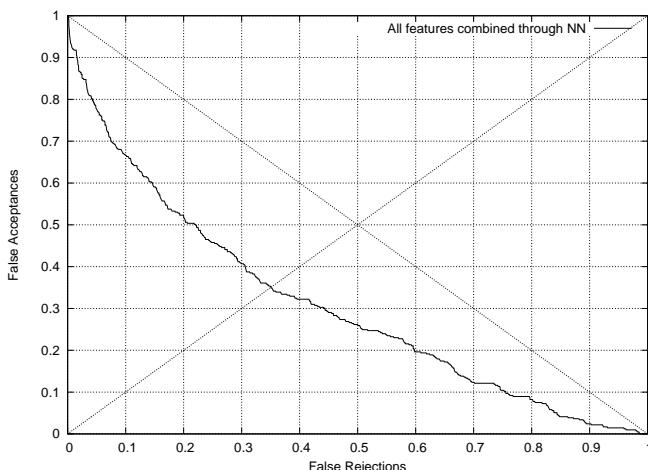


Figure 3. Combination of all features by neural network.

The network used was a fully connected three-layer perceptron with 66 input nodes, 33 hidden nodes and one output node. The activation function is sigmoid.

The NMI results were especially disappointing. As explained in Section 3.2, NMI is harmed by bias. Although we estimated bias on a dedicated set of training data and removed it from the final estimation,

Table V. Performances of all word level features combined.

classifier	equal error rate	NMI	training time	testing time
Logistic Regression	36.8	-2.61×10^{-2}	13"	5"
PLSR	37.5	-5.84×10^{-2}	15'	1"
SVM	36.7	-1.87×10^{-1}	12h	500"
Neural Network	35.0	6.06×10^{-2}	10'	2"

we believe that the poor performance may perhaps be explained by the fact that bias is very different for artificial and natural data and probably much more important on the latter.

In order to evaluate the performance gain brought by automatically generated training corpus, we also split the annotated sentences into a training set (70 sentences) a development set (30 sentences) and a testing set (50 sentences), on which we trained and evaluated the neural network. Figure 4 and Table VI show that training on annotated data do not yield better results than training on the generated corpus. The natural corpus is small, but it must be noted that the artificial corpus was generated in a few hours, while it took more than one day to annotate all the sentences. In addition, human annotations are subject to time and inter-annotator variations. Employing a trained professional may alleviate these problems but is of course more expensive.

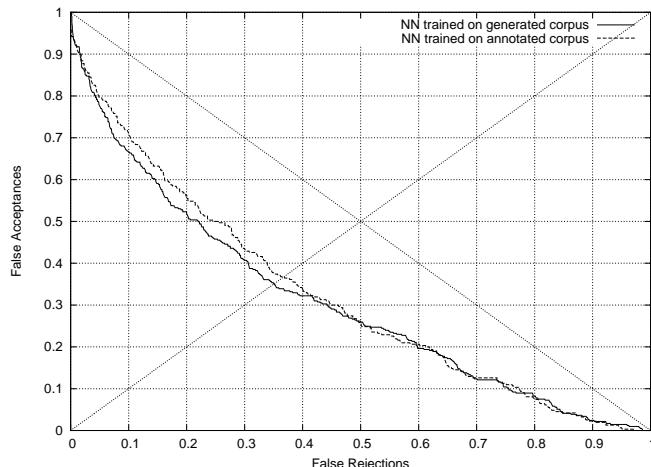


Figure 4. Training neural network on annotated or generated corpus.

In Table VII we show the modest contribution of mutual information (Sections 6.1.4 and 6.1.5) to the performance of neural network combination of the features.

Table VI. Performances of all word level features combined.

classifier	equal error rate	NMI
NN trained on generated corpus	35.0	6.06×10^{-2}
NN trained on annotated corpus	36.8	5.79×10^{-2}

Table VII. Contribution of mutual information based confidence measure to overall performance

	equal error rate	NMI
without IMI and CMI	35.6	5.32×10^{-2}
with IMI and CMI	35.0	6.06×10^{-2}
improvement	-0.60	+7.4 × 10⁻³

7. Sentence-level confidence estimation

The features described in this Section form a numerical representation of a pair made up of a source sentence and a target sentence. As in the previous section, our aim was to compute the distribution of probability of correctness on the numerical space (a subspace of $\mathbb{R}^{d_{sentence}}$). Unlike word level, the algorithm for generating degraded sentences cannot reliably tell if a degraded sentence is still correct or not. We got around the problem of creating a corpus for training classifiers (Section 7.2) but we could not automatically generate a corpus for estimating probability biases. Therefore all normalised mutual information is poor.

Many word-level features can be extended to sentence level by arithmetic or geometric averaging, for example IBM-1 translation probability, n -gram probability, etc.

7.1. FEATURES FOR SENTENCE LEVEL CONFIDENCE ESTIMATION

7.1.1. *LM based features*

The first features we propose are sentence normalised likelihood in a 3-gram model (forward and backward) and average backoff behaviour:

$$x_{30}(\mathbf{s}, \mathbf{t}) = \left(\prod_{j=1}^J P(t_j | t_{j-1}, \dots, t_{j-n+1}) \right)^{\frac{1}{J}} \quad (30)$$

$$x_{31}(\mathbf{s}, \mathbf{t}) = \left(\prod_{j=1}^J P(t_j | t_{j+1}, \dots, t_{j+n-1}) \right)^{\frac{1}{J}} \quad (31)$$

$$x_{32}(\mathbf{s}, \mathbf{t}) = \frac{1}{J} \sum_{j=1}^J x_{23}(\mathbf{S}, \mathbf{T}, j) \quad (32)$$

They can also be used as confidence measures by themselves and their performances are presented in Table VIII and Figure 5 together with intra-lingual mutual information, another kind of language model.

Table VIII. Performances of 3-gram and backoff based confidence measures at sentence-level.

feature	equal error rate	normalised mutual information
3-gram normalised likelihood	41.7	4.02×10^{-3}
backward 3-gram normalised likelihood	41.3	3.97×10^{-3}
averaged backoff behaviour	34.2	4.15×10^{-3}

The following predictive parameter is the source sentence likelihood. Its aim is to reflect how difficult to translate the source sentence is. It is obviously not designed to be used alone.

$$x_{33}(\mathbf{s}, \mathbf{t}) = \left(\prod_{i=1}^I P(s_i | s_{i-1}, \dots, s_{i-n+1}) \right)^{\frac{1}{I}} \quad (33)$$

7.1.2. Average mutual information

$$\begin{aligned} x_{34}(\mathbf{s}, \mathbf{t}) &= \frac{1}{J \times (J-1)} \sum_{i=1}^J \sum_{1 \leq j \neq i \leq J} MI(t_i, t_j) \\ &= \frac{1}{J} \sum_{j=1}^J x_{26}(\mathbf{s}, \mathbf{t}, j) \end{aligned} \quad (34)$$

$$\begin{aligned} x_{35}(\mathbf{s}, \mathbf{t}) &= \frac{1}{I \times J} \sum_{i=1}^I \sum_{j=1}^J MI(s_i, t_j) \\ &= \frac{1}{J} \sum_{j=1}^J x_{27}(\mathbf{s}, \mathbf{t}, j) \end{aligned} \quad (35)$$

We were surprised to observe that cross-lingual MI performed even worse at sentence level than at word level. We have only presented the results for intra-lingual MI in Figure 5 and Table IX, as its performance was closer to other standard confidence measures than it was at word level.

Table IX. intra-lingual mutual information CM as a sentence-level confidence measure.

feature	equal error rate	normalised mutual information
IMI	39.0	9.46×10^{-4}

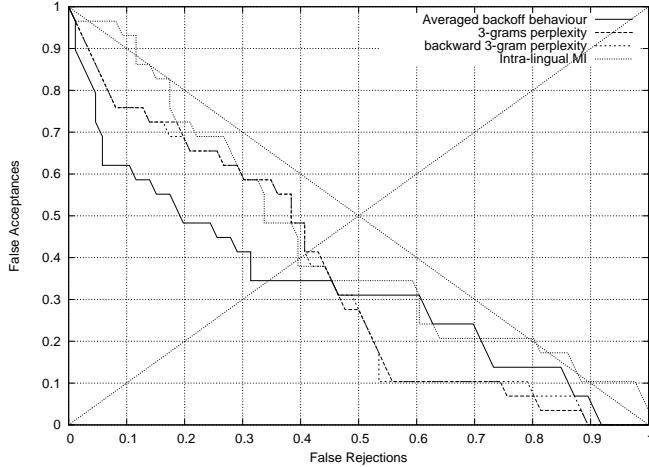


Figure 5. DET curves of 3-gram, backoff and intra-lingual mutual information based confidence measures at sentence level

7.1.3. Normalised IBM-1 translation probability

The score of a sentence is its translation probability in IBM model 1, normalised to avoid penalising longer sentences:

$$x_{36}(\mathbf{s}, \mathbf{t}) = \left(\prod_{i=1}^I \sum_{j=0}^J P(s_i | t_j) \right)^{\frac{1}{I}} \quad (36)$$

As was the case at word level it is surprising to note that although the system was tested on a translation task, confidence measures involving the source sentence do not perform better than the ones involving only the target sentence.

7.1.4. Basic syntax check

A very basic parser checks that brackets and quotation marks are matched, and that full stops, question or exclamation marks, colon or semi-colon are located at the end of the sentence (Blatz et al., 2004).

$$x_{37}(\mathbf{s}, \mathbf{t}) = \begin{cases} 1 & \text{if } \mathbf{t} \text{ is parsable} \\ 0 & \text{otherwise} \end{cases} \quad (37)$$

This feature and the following are only pieces of informations about the source and target sentences, they are not confidence measures themselves.

7.1.5. Length based features

These very basic features reflect levels of consistency between the lengths of a source sentence and of its translation (Blatz et al., 2004). The idea is that source and target sentences should be approximately of the same length, at least for language pairs such as French/English:

$$x_{38}(\mathbf{s}, \mathbf{t}) = \text{Len}(\mathbf{s}) \quad (38)$$

$$x_{39}(\mathbf{s}, \mathbf{t}) = \text{Len}(\mathbf{t}) \quad (39)$$

$$x_{40}(\mathbf{s}, \mathbf{t}) = \frac{\text{Len}(\mathbf{t})}{\text{Len}(\mathbf{s})} \quad (40)$$

7.2. COMBINATION OF SENTENCE-LEVEL FEATURES

As explained earlier in the paper, a generation algorithm cannot tell which sentences are to be considered correct and which are not. Therefore, for sentence-level confidence, it was not directly possible to train classifiers to discriminate between correct and incorrect sentences. Instead, we used SVM, Neural Networks and Partial Least Squares (PLS) to perform regression against sentence level BLEU score³. Sentences were then classified by thresholding this score.

Table X. Performances of PLSR, SVM and Neural Nets at sentence level.

feature	equal error rate	NMI
PLS	29.0	8.14×10^{-2}
SVM	38.0	-2.56×10^{-1}
Neural Net	41.3	-2.44×10^{-2}

Only PLS was found to improve (by 5.2 points, absolute) on the best stand alone confidence measure (Average backoff behavior, Section 7.1.1). Its correlation coefficient with human evaluation was 0.358.

³ It is true that BLEU is not very suited for sentence-level estimation. It has the advantage of being a well known automatic metric for which efficient toolkits are available. We also experimented with TER but too many sentences produced a null score

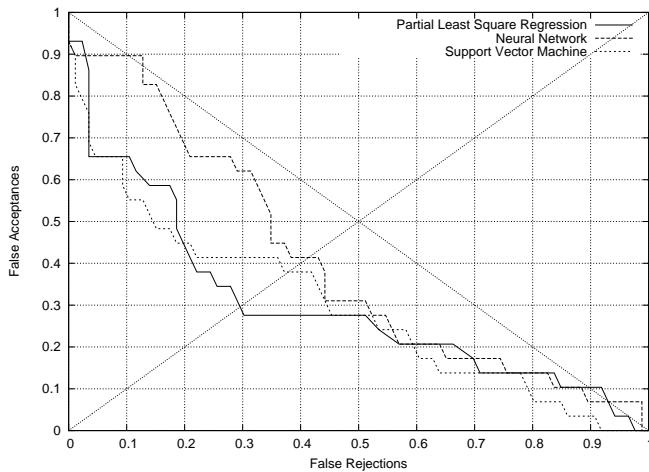


Figure 6. DET curves of PLS and Neural Network combination of sentence level features

8. Preliminary Post Edition Experiment

The previous sections have given a detailed explanation of how the proposed confidence measures work and the amount of errors they are able to detect. In this section we will describe a more subjective usability experiment. Our aim was to obtain qualitative feedback from real users of the system about the usability of confidence measures for assisted post edition. Because of the limited number of subjects, and the fact that many predictive parameters are still work-in-progress, these results are only to be interpreted as hints at what users want and find useful, at what we did right or wrong and at which direction we should follow in our research. The experimental protocol is inspired by the one described in (Plitt and Masselot, 2010). We implemented a post edition tool with confidence measures and let users correct machine translated sentences, with and without the help of confidence measures.

8.1. THE POST EDITION TOOL

The program we developed (see screenshot in Figure 7) can be seen as a simplified version of a tool for Computer Assisted Translation. It displays a source sentence (in our case, in French) and a translation generated by Moses (in English). Errors detected by the confidence measures are highlighted. The user can then opt to edit the proposed translation.

The source sentence is displayed in the top field with the candidate translation in the field below. On the left there is a slider with which the

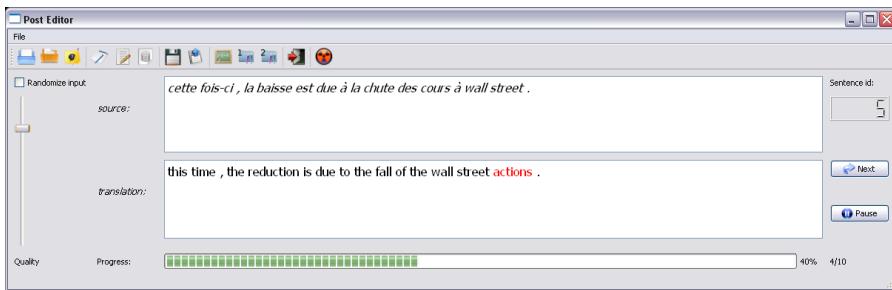


Figure 7. Screenshot of the post edition software.

user can change the acceptance threshold of the confidence estimation system (Section 1.3.1). All words with a score below this threshold are displayed in red. Simplified explanations are given to the user, who does not require a full “lecture” on confidence estimation: he or she is told that s/he may use an automatic help to detect erroneous words, and that the requested quality can be changed with this slider, if s/he so wishes. Of course, if his/her quality requirements are too high (corresponding to a threshold value of 1, that is, the point to the far right on the DET curve — Section 3.1), the system will incorrectly consider all words to be wrong. The user can edit the candidate translation if s/he thinks it is necessary. When s/he is satisfied with the translation s/he has to click on “next”. For the sake of the experiment the user may not come back to a sentence that has already been validated. If required, the user can click on “pause” to take a break thus avoiding that the program continues counting the time spent on the translation, thus making time statistics meaningless. However none of the users ever took a break. Everything else on this GUI is cosmetic (progress bar, etc.).

The total time spent on each sentence was recorded (the time between the loading of the sentence and clicking on the “next” button). This is actually the sum of three partial times, which are also recorded: time typing on the keyboard, time spent on the interface (moving the acceptance slider) and thinking time (the rest).

It should be noted that the proposed translation and confidence scores were not computed on the fly, in order to keep the program responsive and easily portable. This is quite a heavy constraint because the system cannot take the user’s editions into account to compute a new, improved translation, and cannot compute the confidence of the post edited translation (our users were of course informed of that). Also, while all users stated that the program was easy to use, an ergonomist input would be required to ensure we made the right choices with

regard to usability and that what we measure is really the influence of confidence measures and is not due to influence of the interface.

8.2. EXPERIMENTAL PROTOCOL

Because we were not expecting many volunteers, we wanted their English skills to be as homogeneous as possible (all of them are French native speakers) in order to limit the variability of the results. Seven subjects volunteered for the experiment. Six of them are English teachers and one is a master student in English. Unfortunately two of them failed to correctly follow the instructions and the corresponding data was discarded. The experiment lasted approximately two hours, divided in four stages:

First stage: introduction and training. The users were provided with some basic explanations about the domain and the task and given ten sentences to post edit along with simple instructions (see below). These sentences were just for training purposes and were not included in the final results.

Second stage: first experiment. The users were told to start the first experiment when ready. They were given 30 sentences with their corresponding machine translations and were told they could post edit these translations with the help of the confidence measures.

Third stage: second experiment. This experiment was identical to the first, except that the users did not have access to confidence measures. One volunteer out of two had the second experiment before the first, in order to compensate for the "training effect" (users complete the second experiment faster than the first one) and for fatigue (a user may be tired by the time he starts the second experiment, thus affecting post edition speed and quality).

Fourth stage: user feedback. Finally, the users were asked to complete a questionnaire, providing us with feedback on the post edition software and the confidence measures.

We gave the following instructions to the users, with the idea that translated documents must be good enough to be read without extra effort, but not necessarily in beautiful, idiomatic English:

- The goal is to obtain a correct translation, not necessarily a very fluent one. Fix mistakes, not style.

- You can use any help you want (most of them actually used paper or online dictionaries) but:
 - Don't use an online tool to re-translate the sentence
 - Don't spend too much time on details
 - Don't ask the supervisor for help

The sentences were random subsets of the test set of the WMT09 campaign, which is transcripts of news broadcast. Each user had to post edit two randomized sets of thirty sentences. This choice is questionable insofar as most "real life" applications consist of translating whole documents and not a sequence of sentences without connections to each others. However we chose randomized subsets so that the intrinsic difficulty of the task did not influence the results.

8.3. RESULTS AND ANALYSIS

Table XI summarises the most important results of the experiments. Most of these metrics are straightforward but some are worthy of more explanation.

Sentence quality: After the experiment, all the post edited translations were scored by a team member, a native French speaker also fluent in English. Each sentence received a score between 1 and 5 in the same fashion as in StatMT evaluation tasks:

1. the translation is completely unusable.
2. the translation is seriously faulty but a degree of meaning can be grasped.
3. the translation is usable although not very good.
4. the translation has minor flaws.
5. the translation is very good.

Correlation between confidence estimations and editions: our aim here was to check how the user's decisions and the machine predictions correlated. To this end every word in the machine generated hypothesis was mapped to 1 if it was Levenshtein aligned to a word in the edited hypothesis (which means it was not modified), 0 otherwise (which means it had been inserted or modified by the user). The corpus was therefore mapped to a sequence of 0 and 1 and we computed the correlation between this sequence and the estimated probabilities of correctness.

Ratio of number of editions over number of detected errors: this is the ratio of the number of edits made to the original hypothesis over the number of errors which were detected by the system. A high ratio suggests that the user could not find an appropriate trade-off between false positives and false negatives and had to lower his quality requirement (using the slider) in order to get an acceptable level of accuracy.

Table XI. Effect of confidence estimation on a post edition task.

	without CM	with CM
Average time per sentence (seconds):	77	87
Average edit rate:	30%	32%
Average sentence quality:	4.3	4.2
	1st experiment	2nd experiment
Average time per sentence:	84.22	80.12
Average edit rate:	0.29	0.33
Average sentence quality:	4.2	4.3
Ratio of corrections/detected errors	1.76	
Correlation between CM and editions	0.23	

While the results in terms of translation speed are disappointing (Table XI), this experiment was primarily designed to obtain a qualitative feedback from real users of the system. This is what the following analysis will focus on, in order to determine what must be improved and how. A finer grained analysis showed that the time difference is entirely due to “thinking” time. User feedback confirmed that they thought the help was not reliable enough to be useful and that even if it sometimes drew their attention on some mistakes, checking the systems’ recommendations wasted their time. However it must be noted that users were significantly faster during the second post edition task than the first. This suggests that more training is needed before users grow accustomed to the task and really see the program as a tool instead of a constraint. We believe that an experiment involving more users over a longer time frame is necessary. The consistently high and comparable edit rate with and without confidence measures suggests — and this is confirmed by feedback — that a lot of editing was required, but the high ratio of number of corrections over automatically detected errors suggests that confidence measures were not able to precisely discriminate between correct and incorrect words. Regardless of confidence estimation, many of our users stated that they would rather translate a sentence from scratch than edit a flawed machine translation.

As a conclusion to this experiment, we propose the following directions for further improvements and experiments:

- The users should be given a consistent task, not random sentences.

- Users need longer training time as some of them were still not sure what to do with the slider by the end of the tasks. Measurements show that their efficiency continued to increase after the training stage. We believe they need more time to get used to the tool and make the best of it.
- The program interface needs to be carefully designed with ergonomics in mind in order to really measure the influence of confidence measures and not that of the GUI.
- We need more reliable confidence measures and above all, we greatly need to focus on precision rather than recall as we observed that false alarms were very disconcerting for users

9. Conclusion

After introducing and formalising the problem, we presented a method which makes it possible to generate large amount of training data, then developed a list of predictive parameters which we consider are some of the most significant for confidence estimation, including two original measures based on mutual information. We compared different machine learning techniques combining the features we proposed. From these features, we consider Neural Networks and Partial Least Squares Regression to be the best suited, depending on the application. We have shown that combining many features improves over the best predictive parameters alone, by 1.3 points (absolute) EER at word level and 6 points at sentence level on a classification task. Finally we presented an experiment aiming at measuring how helpful confidence estimation is in a post edition task. This experiment suggested that our confidence estimation system is not mature enough to be helpful in such a setting. However the limited number of volunteers and the lack of long term observations makes the results somewhat difficult to interpret. But the knowledge we gained from this experiment and users feedback will help us improve confidence measures for the benefit of future users.

Our hope is that this paper will provide the necessary information to enable the construction of a complete confidence estimation system for machine translation from scratch and facilitate the incorporation therein of new predictive features. In addition to assisted post edition we believe there are many useful applications for confidence estimation, namely:

- Warning a user that the translation he requested may be flawed

- Automatically rejecting hypotheses generated by the decoder or combining several systems in a voting system
- Recombining good phrases from an n-best list or a word graph to generate a new hypothesis.

We have also identified important research directions in which this work could be extended to make confidence measures more helpful for users. Firstly, we would cite computing confidence estimates at phrase level which would enable users to work on semantically consistent chunks while retaining a more fine-grained analysis than with sentences. Secondly semantic features could be introduced which would make it possible to detect otherwise tricky errors like missing negations and help users to focus on meaning errors rather than grammatical errors and disfluencies which are, in some cases, arguably less important.

Acknowledgements

We would like to warmly thank all students and staff of IUFM of Lorraine (<http://www.lorraine.iufm.fr>) who volunteered for the post edition experiment and made the experiment possible. We would particularly like to thank Mr. Gilles Grateau for his careful and conscientious help throughout this work and for his patience and constant cheerfulness.

References

- Blatz, J., E. Fitzgerald, G. Foster, S. Gandrabur, C. Goutte, A. Kulesza, A. Sanchis, and N. Ueffing: 2004, ‘Confidence Estimation for Machine Translation’. In: *Proceedings of Coling 2004*. Geneva, Switzerland, pp. 315–321.
- Chang, C.-C. and C.-J. Lin: 2011, ‘LIBSVM: A library for support vector machines’. *ACM Transactions on Intelligent Systems and Technology* **2**, 27:1–27:27. Software available at <http://www.csie.ntu.edu.tw/~cjlin/libsvm>.
- Cortes, C. and V. Vapnik: 1995, ‘Support-vector networks’. *Machine learning* **20**(3), 273–297.
- Duchateau, J., K. Demuynck, and P. Wambacq: 2002, ‘Confidence scoring based on backward language models’. In: *Proceedings of IEEE International Conference on Acoustics, Speech, and Signal Processing*, Vol. 1. Orlando, Florida, pp. 221–224.
- Fausett, L. V.: 1994, *Fundamentals of neural networks*. Prentice-Hall Englewood Cliffs, NJ.

- Gandrabur, S. and G. Foster: 2003, 'Confidence estimation for translation prediction'. In: *Proceedings of the seventh Annual Conference of the North American Chapter of the Association for Computational Linguistics (HLT-NAACL)*, Vol. 4. Edmonton, Canada, pp. 95–102.
- Guo, G., C. Huang, H. Jiang, and R. Wang: 2004, 'A comparative study on various confidence measures in large vocabulary speech recognition'. In: *Proceedings of the International Symposium on Chinese Spoken Language Processing*. Hong Kong, China, pp. pp. 9–12.
- Hsu, C.-W., C.-C. Chang, and C.-J. Lin: 2003, 'A Practical Guide to Support Vector Classification'. Technical report, Department of Computer Science, National Taiwan University.
- Koehn, P., H. Hoang, A. Birch, C. Callison-Burch, M. Federico, N. Bertoldi, B. Cowan, W. Shen, C. Moran, R. Zens, et al.: 2007, 'Moses: Open source toolkit for statistical machine translation'. In: *Proceedings of the 45th Annual Meeting of the ACL on Interactive Poster and Demonstration Sessions*. pp. 177–180.
- Kononenko, I. and I. Bratko: 1991, 'Information-based evaluation criterion for classifier's performance'. *Machine Learning* **6**(1), 67–80.
- Lavecchia, C., K. Smaïli, D. Langlois, and J.-P. Haton: 2007, 'Using inter-lingual triggers for Machine translation'. In: *Proceedings of the Eighth International Conference on Speech Communication and Technology (INTERSPEECH)*. Antwerp, Belgium, pp. 2829–2832.
- Menard, S. W.: 2002, *Applied logistic regression analysis*, Sage university papers, Quantitative applications in the social sciences. Thousand Oaks, Calif. [u.a.]: Sage.
- Miller, G.: 1995, 'WordNet: a lexical database for English'. *Communications of the ACM* **38**(11), pp. 39–41.
- Nissen, S.: 2003, 'Implementation of a Fast Artificial Neural Network Library (fann)'. Technical report, Department of Computer Science University of Copenhagen (DIKU). <http://fann.sf.net>.
- Papineni, K., S. Roukos, T. Ward, and W. Zhu: 2001, 'Bleu: a method for automatic evaluation of machine translation'. In: *Proceedings of the 40th Annual of the Association for Computational linguistics*. Philadelphia, USA, pp. 311–318.
- Plitt, M. and F. Masselot: 2010, 'A Productivity Test of Statistical Machine Translation Post-Editing in a Typical Localisation Context'. *The Prague Bulletin of Mathematical Linguistics* **93**(-1), 7–16.
- Quirk, C.: 2004, 'Training a sentence-level machine translation confidence measure'. In: *Proceedings of LREC 2004*. Lisbon, Portugal, pp. 825–828.
- Raybaud, S., C. Lavecchia, D. Langlois, and K. Smaïli: 2009a, 'New Confidence Measures for Statistical Machine Translation'. In: *Proceedings of the International Conference on Agents and Artificial Intelligence*. Porto, Portugal, pp. 61–68.
- Raybaud, S., C. Lavecchia, D. Langlois, and K. Smaïli: 2009b, 'Word- and sentence-level confidence measures for machine translation'. In: *Proceedings of the 13th Annual Conference of the European Association for Machine Translation*. Barcelona, Spain, pp. 104–111.
- Schmid, H.: 1994, 'Probabilistic part-of-speech tagging using decision trees'. In: *Proceedings of International Conference on New Methods in Language Processing*, Vol. 12. pp. 44–49.
- Schmid, H.: 1995, 'Improvements In Part-of-Speech Tagging With an Application To German'. In: *In Proceedings of the ACL SIGDAT-Workshop*. pp. 47–50.

- Simard, M., N. Ueffing, P. Isabelle, and R. Kuhn: 2007, 'Rule-based translation with statistical phrase-based post-editing'. In: *Proceedings of the ACL-2007 Workshop on Statistical Machine Translation (WMT-07)*. Prague, pp. 203–206.
- Siu, M. and H. Gish: 1999, 'Evaluation of word confidence for speech recognition systems'. *Computer Speech and Language*.
- Smola, A. and B. Schölkopf: 2004, 'A tutorial on support vector regression'. *Statistics and Computing* **14**(3), 199–222.
- Snover, M., B. Dorr, R. Schwartz, L. Micciulla, and J. Makhoul: 2006, 'A study of translation edit rate with targeted human annotation'. In: *Proceedings of Association for Machine Translation in the Americas*. Cambridge, pp. 223–231.
- Specia, L., N. Cancedda, M. Dymetman, M. Turchi, and C. N.: 2009, 'Estimating the sentence-level quality of Machine Translation Systems'. In: *Proceedings of the 13th Annual Conference of the European Association for Machine Translation*. Barcelona, Spain, pp. 28–35.
- Stolcke, A.: 2002, 'SRILM – an extensible language modeling toolkit'. In: *Seventh International Conference on Spoken Language Processing*. Denver, USA, pp. 901–904.
- Tobias, R.: 1995, 'An introduction to partial least squares regression'. In: *Proceedings of the Twentieth Annual SAS Users Group International Conference, Cary, NC: SAS Institute Inc.* Orlando, Florida, pp. 1250–1257.
- Ueffing, N. and H. Ney: 2004, 'Bayes decision rules and confidence measures for statistical machine translation'. In: *Proceedings of EsTAL Espana for Natural Language Processing*. Alicante, Spain, pp. 70–81.
- Ueffing, N. and H. Ney: 2005, 'Word-level confidence estimation for machine translation using phrase-based translation models'. In: *Proceedings of the conference on Human Language Technology and Empirical Methods in Natural Language Processing (HLT '05)*. Morristown, NJ, USA, pp. 763–770.
- Uhrik, C. and W. Ward: 1997, 'Confidence Metrics Based on N-Gram Language Model Backoff Behaviors'. In: *Fifth European Conference on Speech Communication and Technology*. Rhodes, Greece, pp. 2771–2774.
- Wold, S., A. Ruhe, H. Wold, and W. Dunn III: 1984, 'The collinearity problem in linear regression. The partial least squares (PLS) approach to generalized inverses'. *SIAM Journal on Scientific and Statistical Computing* **5**(3), pp. 735–743.
- Zhang, R. and A. Rudnicky: 2001, 'Word level confidence annotation using combinations of features'. In: *Seventh European Conference on Speech Communication and Technology*. Aalborg, Denmark, pp. 2105–2108.