
Label Ranking under Ambiguous Supervision for Learning Semantic Correspondences

Abstract

This paper studies the problem of learning from ambiguous supervision, focusing on the task of learning semantic correspondences. A learning problem is said to be ambiguously supervised when, for a given training input, a set of output candidates is provided with no prior of which one is correct. We propose to tackle this problem by solving a related unambiguous task with a label ranking approach and show how and why this performs well on the original task, via the method of task-transfer. We apply it for learning to match natural language sentences to structured representation of their meaning and empirically demonstrate that this competes with the state-of-the-art on two benchmarks.

1. Introduction

Annotating training data for supervised learning algorithms is often costly and time-consuming, and depending on the task can even require highly-advanced expertise on the part of the labeler. One opportunity to bypass this requirement is that for many tasks an automatic use of multimodal environments can provide training corpora with little or no human processing. For instance, the time synchronisation of several media can generate annotated corpora: matching movies with subtitles (Cour et al., 2008) can be used for speech recognition or information retrieval in videos, matching vision sensors and other sensors can be used to improve robotic vision (as in (Angelova et al., 2007)), matching natural language and perceptive events (such as audio commentaries and soccer actions in RoboCup (Chen & Mooney, 2008)) can be used to learn semantics. Indeed, the Internet is abundant with such sources, for example one could think to use the text surrounding pictures in a webpage as image labeling candidates.

Such automatic procedures can build large corpora of ambiguously supervised examples. Indeed, every resulting input instance (picture, video, speech, ...) is paired with a set of candidate output labels (text caption, subtitle, ...). The automation of the data collection makes it impossible to directly know which one is correct among them, or even if there exists a correct label. To conceive systems able to efficiently learn out of such noisy and ambiguous supervision would be a huge leap forward in machine learning. These methods could then benefit from large training sets obtained with drastically reduced costs.

A domain for which data collection is particularly expensive is semantic parsing (Mooney, 2004). The goal of semantic parsing is to build systems able to understand questions or instructions in natural language in order to bring about a major improvement in human-computer interfacing. Formally, this consists in mapping natural language sentences into structured representations of their meaning which are domain-specific and directly interpretable by a computer. Recent machine learning work (Zettlemoyer & Collins, 2009; Branavan et al., 2009; Ge & Mooney, 2009) exhibits promising progress on this task. Building training data for semantic parsing requires the precise alignment of sentences and formal representations with a costly process that forbids the creation of large-scale corpora. However, for many topics such as finance, music or sports, huge databases paired with corresponding texts are readily available and can be automatically aligned to provide large quantities of ambiguously annotated examples. Unfortunately, this data cannot be used by most semantic parsing methods.

In this paper, we tackle the problem of learning semantic correspondences for natural language (Snyder & Barzilay, 2007; Liang et al., 2009). More precisely, this task consists in aligning texts with corresponding database entries in order to provide disambiguated training examples for semantic parsing. To learn it under ambiguous supervision, we propose to solve an associated task and make use of task-transfer. We derive a label ranking approach to a related unambiguous task and demonstrate that a solution to this problem

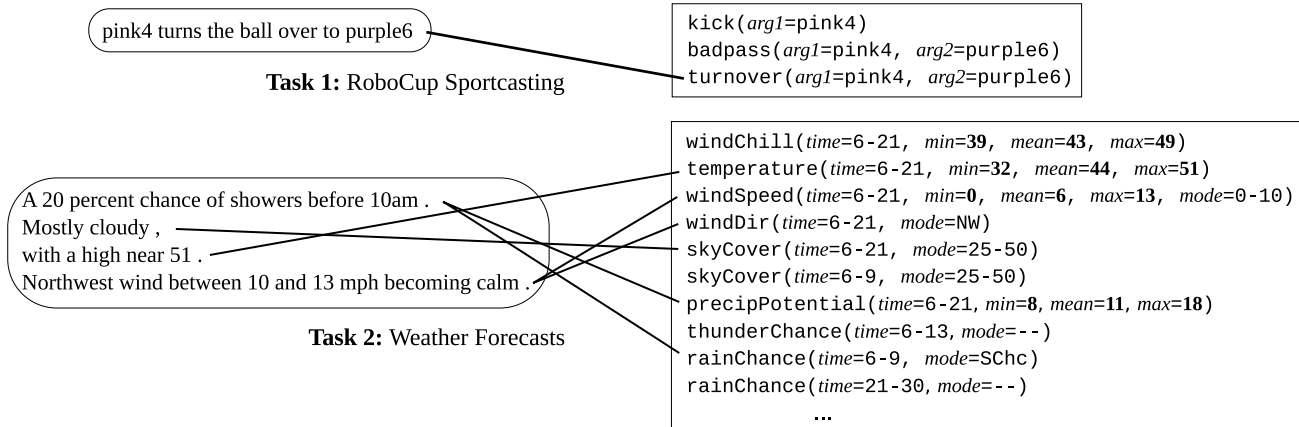


Figure 1. **Examples of scenarios** for the two tasks studied in this paper. A scenario is composed of a text (left) and a set of records (right). Some of those correspond to interpretations of either, all, or segments of the text (black lines). However these gold alignments are unknown in training. All record structures are identical: a type (**kick**, **windDir**, ...) and several fields (*arg1*, *time*, *mean*, ...) whose values can be integer (in **bold**) or categorical (in **typewriter**).

performs well on the original task. In other words, we show that one can bypass the difficulty of ambiguous supervision and still reach good performances on its desired target. We then propose an intuitive framework to directly apply standard ranking algorithms to successfully and quickly learn semantic correspondences even with a poor supervision level. On two concrete problems, RoboCup sportscasting and Weather forecasting, we empirically demonstrate that this new simple approach is competitive with the generative method proposed by (Liang et al., 2009), which is the best current approach to the best of our knowledge.

The rest of the paper is organized as follows. Section 2 details the task as well as the datasets we used. In Section 3 we formalize the problem of ranking under ambiguous supervision while in Section 4 we establish the transfer of performance. We explain in Section 5 how to use it for learning semantic matching and describe experimental results in Section 6.

2. Learning Semantic Correspondences

Our interest consists in learning to match a natural language text with a structured representation of its meaning, which is composed of one or several domain-specific database records. All records share the same pattern: a record type followed by a set of fields, which can take either categorical or integer values. However, the kinds of types and fields are task-dependent. We detail them for the specific datasets we consider in this paper in the following subsections.

The training algorithm is given pairs composed by a

text and a set of records. Following (Liang et al., 2009), we use the term *scenario* to refer to such a pair. The set of records, which we call the *candidate set*, is gathered via a cheap automatic process which introduces noise and ambiguity in the supervision. Hence, the candidate set typically groups together pertinent records regarding the associated text as well as many irrelevant ones. Learning semantic correspondences aims at detecting these relevant records, i.e. the *gold alignments*, among all the records of the candidate set. Therefore, at test time, one is provided scenarios, i.e. texts paired with candidate sets, as well as gold alignments, which are not given to training algorithm but are essential for evaluation purposes.

2.1. RoboCup Sportscasting

Our first specific task is to learn to match commentaries with records describing actions of RoboCup soccer games. We used the data collected by (Chen & Mooney, 2008) corresponding to four RoboCup finals and composed of text commentaries automatically paired with records representing the actions that occurred within 5 seconds of them. The whole dataset groups 1,872 scenarios, and each of them is composed by one sentence and a candidate set containing between 1 and 12 records (with a mean of 2.4). One of those is supposed to correspond to the commentary but it is worth noting that, for more than 15% of the scenarios, the correct record is not in the candidate set: in that case, any prediction is automatically wrong.

In total, there are 9 record types (e.g. **pass**, **kick**, **ballstopped**, ...) and each record can have at most

2 fields (e.g. $arg1=$ purple1, $arg2=$ pink4, ...) indicating the player(s) involved in the action. As illustration, Figure 1 (top) depicts a RoboCup scenario, composed by a commentary (left) and a candidate set of 3 records (right), and containing a gold alignment. In this dataset, all record fields are categorical so that, for example, no a priori association between the field value pink4 and the word “pink4” is possible. Category names chosen are purely for ease of explanation.

2.2. Weather Forecasts

The second task concerns learning correspondences between local forecast reports and records representing weather events. These records actually consist in measurements of meteorological indicators such as temperature, wind speed/direction or chance of sleet automatically extracted from the database of www.weather.gov. The dataset, created by (Liang et al., 2009), groups 22,146 scenarios collected, each day and night over 3 days, from the local forecasts of 3,753 US cities. Each candidate set contains exactly 36 records.

An example of scenario is given at the bottom of Figure 1. There are a total of 12 record types (e.g. temperature, windDir, thunderChance, ...) and each of them can have up to 5 fields. Two of them take categorical values: *time* which indicates the time range of the event and *mode* which can precise some of its characteristics (e.g. the direction of the wind). The other three take integer values. Denoted *min*, *mean* and *max*, they provide exact values of some quantifiable indicators like the temperature or the wind speed.

Learning for this task is harder than for RoboCup because the candidate sets group more records and each text refers to more than one record (5.8 on average). Indeed, the reports have been split by punctuation into lines, giving an average of 4.6 text lines per scenario, all sharing the same candidate set. During evaluation, gold alignments must be performed at this level and there is approximately 1.2 alignments per line.

3. Ranking and Ambiguous Supervision

The task of label ranking commonly considers a measurable space of observations \mathcal{X} , a finite set of labels $\mathcal{Y} = \{1, \dots, |\mathcal{Y}|\}$ and a function ϕ that maps any (input, label) pair to a measurable space \mathcal{M} . For the particular case of semantic matching, \mathcal{X} is the set of possible sentences, \mathcal{Y} is the set of all possible records (identified with $\{1, \dots, |\mathcal{Y}|\}$) and ϕ is the joint representation of a sentence and a record. A *label scoring function* (LSF in brief) is a real-valued function of the form $h = g \circ \phi$ with $g: \mathcal{M} \rightarrow \mathbb{R}$. The score $h(x, i)$ of the

label i for the input x is denoted $h_i(x)$. From now on, we only consider measurable LSFs.

3.1. Ambiguous Supervision

We now describe our setting for label ranking with ambiguous supervision. The data is modeled by a triplet of jointly distributed random variables (X, Z, Y) taking values in $\mathcal{X} \times \{0, 1\}^{|\mathcal{Y}|} \times \{0, 1\}^{|\mathcal{Y}|}$. For a given realization (x, z, y) of (X, Z, Y) , x is the observation, z is a subset of \mathcal{Y} called the *target set* for x ($z_i = 1$ when i is in the target set), and y is the *candidate set* for x .

A *target task* is defined by a risk functional \mathcal{R} which uses the full knowledge of (X, Z, Y) . In the label ranking framework, \mathcal{R} measures the ability of any LSF to give higher scores to the labels of the target set. Learning with ambiguous supervision means that the target sets are unknown at training time: the training data only consists of n realizations $(x^k, y^k)_{k=1}^n$ of (X, Y) . For example, on our datasets x is a sentence, z represents the gold alignment which is unknown during training and y is the candidate set for x . The target task is to rank the gold alignments above the other candidates.

We consider learning with ambiguous supervision as a kind of transfer learning: we intend to find an LSF of low risk as measured by the target task. But since we cannot measure this risk in general, our approach is based on defining a proxy risk functional which only depends on (X, Y) , so that the training data helps us to find a performing function for this proxy. We then address the issue of how to define the proxy risk so that we can transfer what we have learnt on the training data to the target task, under some assumptions that are (almost) satisfied by real-life datasets.

3.2. Categories of Supervision

Before describing the specific target tasks of ranking we consider in this paper, we may distinguish several characteristics of the data, which influence the difficulty of their learning. We say that (X, Z, Y) is:

noisy if $\mathbb{P}(Z^*(X) \neq Z) > 0$, where:

$$\forall x, Z^*(x) = \arg \max_{z \in \{0, 1\}^{|\mathcal{Y}|}} \mathbb{P}(Z = z | X = x),$$

ambiguous if $\mathbb{P}(Z \neq Y) > 0$,

incomplete if $\mathbb{P}(Z \cap Y \neq Z) > 0$.

The data exhibits some *noise* when the target set is not deterministic. The supervision is *ambiguous* when the candidate set can be different from the target set. It may be bigger, or even completely uncorrelated (although learning is probably impossible in the latter case). We also qualify an ambiguous supervision as *incomplete* when some target labels are not

candidates or when the target set is empty while the candidate set is not. This occurs frequently in the RoboCup dataset as many training sentences have no gold alignment, but rarely in the Weather dataset. The latter is however very noisy: many lines (e.g. “Mostly cloudy”) appear several times with different target sets (e.g. `skycover(time=6-21, mode=25-50)` and `skycover(time=17-30, mode=25-50)`).

3.3. Target Tasks

In this paper, we consider the following target tasks:

Full Ranking This task is defined by the *Full Ranking Risk* $\mathcal{R}^{\text{Full}}$, which measures the ability of h to rank the target labels above all others:

$$\mathcal{R}^{\text{Full}}(h) = \mathbb{E}[\ell^{\text{Full}}(h, X, Z)], \text{ with}$$

$$\ell^{\text{Full}}(h, x, z) = \frac{1}{P} \sum_{i \neq j} \mathbf{I}_{\{(z_i - z_j) \tilde{h}_{ij}(x) < 0\}} \quad (1)$$

for all (x, z) , where $P = |\mathcal{Y}|(|\mathcal{Y}| - 1)$ is a normalization factor, $\mathbf{I}_{\{\cdot\}}$ is the indicator function and $\tilde{h}_{ij}(x) = \text{sign}(h_i(x) - h_j(x))$ with $\text{sign}(t) = 2\mathbf{I}_{\{t \geq 0\}} - 1$.

Candidate Set Ranking This task ignores the labels that are not in the candidate set, independently of whether they are target labels or not. The corresponding risk is defined as:

$$\mathcal{R}^{\text{CSet}}(h) = \mathbb{E}[\ell^{\text{CSet}}(h, X, Z, Y)], \text{ with}$$

$$\ell^{\text{CSet}}(h, x, z, y) = \frac{1}{P} \sum_{i \neq j} \mathbf{I}_{\{y_i = y_j = 1\}} \mathbf{I}_{\{(z_i - z_j) \tilde{h}_{ij}(x) < 0\}} \quad (2)$$

for any (x, z, y) .

These risks correspond to standard pairwise ranking risks of label ranking (see e.g. Har-Peled et al., 2002), up to the normalization factor P .¹ The Full Ranking Risk increases linearly with the number of pairs of (non-target, target) labels for which the relative ordering is incorrectly predicted. The Candidate Set Ranking Risk behaves the same way, but restricted to the labels of the candidate set. The difference between these two tasks is clarified by the following lemma:

Lemma 1 Define, for all x, i ,

- $\eta_i^{\text{Full}}(x) = \mathbb{P}(Z_i = 1 | X = x)$,
- $\eta_i^{\text{CSet}}(x) = \mathbb{P}(Z_i = 1 | Y_i = 1, X = x)$,

¹In practice, the risks are normalized in order to be equal to 1 in the worst case. We use the same normalization factor P in all our definitions as it simplifies the notations and does not essentially change the results.

and denote:

$$\mathcal{R}_{\text{bayes}}^{\text{Full}} = \mathcal{R}^{\text{Full}}(\eta^{\text{Full}}) \text{ and } \mathcal{R}_{\text{bayes}}^{\text{CSet}} = \mathcal{R}^{\text{CSet}}(\eta^{\text{CSet}}).$$

Then, for any LSF h , we have: $\mathcal{R}_{\text{bayes}}^{\text{Full}} \leq \mathcal{R}^{\text{Full}}(h)$.

Besides, if for any labels i and j with $i \neq j$, $Z_i = 1$ is conditionally independent of Y_j given $Y_i = 1$ and X then:

$$\mathcal{R}_{\text{bayes}}^{\text{CSet}} \leq \mathcal{R}^{\text{CSet}}(h) \text{ for any LSF } h.$$

The proof of the lemma, as well as all other proofs, are given in supplementary material. The two target tasks are inherently different, as they have different optimal LSF in general. In the candidate set task, a label can be very rarely a target, whilst being top-ranked any time it appears. On the other hand, the full ranking setup only depends on the frequency of a label in the target set given x .

The conditional independence assumption used for Candidate Set Ranking is rather natural because it only requires that a candidate label is a target label independently of the appearance of the other labels. Notice that no assumption is made on how the candidates are selected, so they can be quite correlated.

4. Transfer of Performance

We now describe our approach to transfer label ranking with ambiguous supervision to our target tasks. In the absence of prior knowledge on the target task, and more precisely on the distribution of Z , we cannot expect more than learning to rank the labels i according to $\mathbb{P}(Y_i = 1 | X = x)$. By Lemma 1 (exchanging Z with Y) we could attempt to minimize the risk defined by $\mathbb{E}[\ell^{\text{Full}}(h, X, Y)]$. However, minimizing the corresponding empirical risk on the whole training set can be inefficient when \mathcal{Y} is large (as illustration, for Weather, $|\mathcal{X}| \approx 10^5$ lines and $|\mathcal{Y}| \approx 10^6$ records).

4.1. Ranking Risk for Training

In this work, we thus consider the following alternative, which measures the ability of an LSF to rank the candidate labels higher than a *random sample* of \mathcal{Y} :

Proposition 2 For any LSF h , the Ambiguous Label Ranking Risk of h , denoted $\mathcal{R}^{\text{Amb}}(h)$, is defined by:

$$\mathcal{R}^{\text{Amb}}(h) = \mathbb{E}[\ell^{\text{Amb}}(h, X, Y^+, Y^-)]$$

where Y^+ and Y^- are $\{0, 1\}$ -valued r.v. defined by:

$$Y^+ = Y \text{ and } \mathbb{P}(Y^- = 1 | Y = y, X = x) = \mathbb{P}(Y^- = 1) = \frac{s}{|\mathcal{Y}|}$$

and, for any x, y^+, y^- :

$$\ell^{\text{Amb}}(h, x, y^+, y^-) = \frac{1}{P} \sum_{i, j: i \neq j} \mathbf{I}_{\{y_i^+ = 1\}} \mathbf{I}_{\{y_j^- = 1\}} \mathbf{I}_{\{\tilde{h}_{ij}(x) < 0\}}$$

Then, denoting $\eta_i^{\text{Amb}}(x) = \mathbb{P}(Y_i = 1 | X = x)$ we have, for any LSF h :

$$\mathcal{R}_{\text{bayes}}^{\text{Amb}} \leq \mathcal{R}^{\text{Amb}}(h) \text{ where } \mathcal{R}_{\text{bayes}}^{\text{Amb}} = \mathcal{R}^{\text{Amb}}(\eta^{\text{Amb}}).$$

In the pointwise loss $\ell^{\text{Amb}}(h, x, y^+, y^-)$, y^+ corresponds to the candidate set for x , and y^- , called the *negative set*, is a random subsample of \mathcal{Y} . One way to create it is simply to randomly sample s labels of \mathcal{Y} without replacement, where s is called the *size parameter*. In practice, given the training data $(x^k, y^k)_{k=1}^n$, we create a random subsample $y^{k,-}$ of size s for each k , and apply an existing algorithm for ranking. Note that nothing forbids a label to appear in both the candidate and the negative sets.

4.2. Coherent Supervision and Task-Transfer

We now address the following issue: to what extent the performance of the function learnt by minimizing the (empirical) Ambiguous Label Ranking Risk is transferred to the target tasks? Such a task-transfer is bound to the following notion of *coherence* between the ambiguous supervision and an arbitrary LSF:

Definition 1 (Coherence) Denote $\lfloor t \rfloor_+$ the positive part of t . The ambiguous supervision is coherent with the LSF ρ if there is $\alpha > 0$ such that, for any x, i, j :

$$\lfloor \eta_i^{\text{Amb}}(x) - \eta_j^{\text{Amb}}(x) \rfloor_+ \geq \alpha \lfloor \rho_i(x) - \rho_j(x) \rfloor_+.$$

Thus, η^{Amb} is coherent with a LSF when it preserves the relative ordering of labels as well as the relative differences of scores. Our main result is that coherence with one of the Bayes-optimal LSF defined in Lemma 1 implies that the ranking performance on the ambiguous task defined by Proposition 2 is transferred to the corresponding target task:

Theorem 3 (Coherence implies transfer)

If η^{Amb} is coherent with η^{Full} , then there is a constant $\beta^{\text{Full}} > 0$ such that, for any LSF h , we have:

$$\mathcal{R}^{\text{Amb}}(h) - \mathcal{R}_{\text{bayes}}^{\text{Amb}} \geq \beta^{\text{Full}} (\mathcal{R}^{\text{Full}}(h) - \mathcal{R}_{\text{bayes}}^{\text{Full}}).$$

Moreover, under the conditional independence assumption of Lemma 1, if η^{Amb} is coherent with η^{CSet} , then there is $\beta^{\text{CSet}} > 0$ such that for any LSF h :

$$\mathcal{R}^{\text{Amb}}(h) - \mathcal{R}_{\text{bayes}}^{\text{Amb}} \geq \beta^{\text{CSet}} (\mathcal{R}^{\text{CSet}}(h) - \mathcal{R}_{\text{bayes}}^{\text{CSet}}).$$

The constants in the theorem increase with the value of α of Definition 1 and decrease with the size parameter s of Proposition 2. Roughly speaking, η^{Amb} is coherent with η^{Full} when the most frequent candidate labels

are also the most frequent target labels and, η^{Amb} is coherent with η^{CSet} when the most frequent candidate labels are supposed to be top-ranked in the candidate sets they appear. In any case, when the supervision is coherent with a Bayes-optimal LSF of Lemma 1, there is a strong transfer of performance: an approximately optimal LSF for the ambiguous ranking risk is also approximately optimal for the target task. Let us now discuss our results on our applications.

RoboCup Sportscasting In this dataset, most sentences have a single possible target label. Thus, for a sentence x , there is a single $i^*(x)$ such that $\eta_{i^*(x)}^{\text{Full}}(x) > 0$. It happens sometimes that $\eta_{i^*(x)}^{\text{Full}}(x) \neq 1$ because the target set is empty rather often (the supervision is incomplete), but if $i^*(x)$ is in the candidate set, we are sure it is also in the target set, so we have $\eta_{i^*(x)}^{\text{CSet}}(x) = 1$, and for any $j \neq i^*(x)$, $\eta_j^{\text{Full}}(x) = \eta_j^{\text{CSet}}(x) = 0$. Then, if there exists $\epsilon > 0$ such that $\eta_{i^*(x)}^{\text{Amb}}(x) > \eta_j^{\text{Amb}}(x) + \epsilon$ for all x and $j \neq i^*(x)$, by Definition 1, η^{Amb} is coherent with both η^{Full} and η^{CSet} . This assumption only requires the correct label to be the most frequent in the candidate sets, without any other constraint on the other labels. It is very likely to be verified in the dataset, so we should be able to learn a function performing for both the Full Ranking and the Candidates Set Ranking tasks.

Weather Forecasting As discussed in Section 3.2, the Weather dataset is rather noisy, meaning that $\eta_j^{\text{Full}}(x) > 0$ for many labels j . However, given a sentence x and a candidate set y , a candidate label i is either correct or incorrect independent of the other labels, so $\eta_i^{\text{CSet}}(x) = 0$ or $\eta_i^{\text{CSet}}(x) = 1$. Following the same reasoning as for RoboCup, to be coherent with η^{CSet} , it is sufficient for η^{Amb} that the labels i with $\eta_i^{\text{CSet}}(x) = 1$ appear more frequently in the candidate sets of x than those with $\eta_i^{\text{CSet}}(x) = 0$. This is likely to be true because the records in the candidate sets are not totally correlated.

Related Work Cour et al. (2009) propose a proxy risk for multiclass classification under ambiguous supervision, and prove a result similar in essence to Theorem 3 for that case. We can notice two fundamental differences with our approach. First, the formulation of their result is strictly weaker, since they do not prove that the Bayes-optimal point of their proxy risk is also Bayes-optimal for their target risk, even under strong assumptions similar to our notion of coherence. Secondly, their framework for ambiguous supervision does not treat the case of incomplete supervision.

5. Practical Ranking Setup

To be more concrete, we now describe how we employed ranking to learn semantic correspondences.

5.1. Learning Model

The training data is a set of pairs $(x^k, y^k)_{k=1}^n$. On RoboCup, x^k is a sentence and y^k the candidate set. On the Weather dataset, x^k is a line, and y^k is the candidate set associated to the line’s scenario. On the latter dataset, the candidate sets have many uninformative records² which are constantly expressed while never correct. As they artificially introduce ambiguity, we discard them from both the training and test sets.

To apply Proposition 2, we define, for each k , $y^{k,+}=y^k$, and create the negative set $y^{k,-}$ by sampling a number $s.n_k$ of random records (without replacement) among all the records present in the data, with $n_k = |y^k|$. The hyperparameter s is actually employed as a multiplicative factor of the size of the candidate set. Given a joint feature function ϕ (discussed below) of (text, record) pairs, we learn a linear LSF with a regularized convex relaxation of the empirical risk corresponding to \mathcal{R}^{Amb} similar to SVM^{rank} (Joachims, 2006):³

$$\min_w \frac{1}{2} \|w\|^2 + \sum_{k=1}^n \frac{C}{sn_k^2} \sum_{\substack{i:y_i^{k,+}=1 \\ j:y_j^{k,-}=1}} [1 - \langle w, \phi(x^k, i) - \phi(x^k, j) \rangle]_+$$

In the experimental section, we refer to this learning model as the ARank algorithm.

5.2. Feature Representations

We use slightly different feature systems to encode the texts and records of each task in the function ϕ .

RoboCup Sportscasting For this dataset, each commentary $x \in \mathcal{X}$ is encoded using a binary vector based on a bag of its words, bigrams and trigrams. Similarly, each record is characterized by a binary vector indicating its type and its different categorical field values. The joint representation $\phi(x, y)$ of a sentence x and a record y is then obtained by performing an outer product of their respective encoding vectors.

Weather Forecasting We also employ bag of word representations (including bigrams and trigrams) for the sentences, and binary ones for the records. However, some extra-characteristics must be added. First,

²The records with only *time* and *mode=--* arguments.

³One may note that in Proposition 2, the normalization factor and the negative set size are constant, while they may vary. It is a matter of implementation, with no real influence as these values are close to their means. Likewise, the mean of sn_k^2 is not P , but it only changes the C scale.

Table 1. Datasets and parameters used in the experiments.

	ROBOCUP	WEATHER
Scenarios	1,872	2,2146
Records (per scenario)	2.4	36.0
Gold align. (per scenario)	0.8	5.8
Negative set size: s	$\times 50$	$\times 1$
SVM regularizer: C	10^{-2}	10^{-5}

special care must be taken with the integer valued record fields. Following (Liang et al., 2009), we incorporate to ϕ , features that express the crucial information of whether a word m matches the value of a record field (e.g. in the example of Figure 1, the number 51 corresponds to a word of the 3rd line and a field of the **temperature** record). We also consider cases where an approximation of the record value might be used in the text, in place of the exact one. So we check for m , m rounded to 5 (up/low), $m+/-1$ and $m+/-2$. Then, we add an extra-feature indicating whether the *time* field value of a record corresponds to the majority value of the candidate records of $y^{k,+}$.

5.3. Prediction Strategies

In the next section, we evaluate our model with a ranking measure similar to Candidate Set Ranking loss. However, we also have to set up a prediction process in order to compare its performances with previous works using a precision/recall measure. For a given scenario, a prediction consists in returning the target set. As our algorithm cannot directly determine the number of records to pick out, we used post-processing heuristics. On the RoboCup dataset, we have to detect when the target set is empty, and for the Weather one, we need to predict its size (1 or 2).

In Section 4.2 it is explained that, for a RoboCup scenario, a single record is correct and if it is not a candidate, then the target set is empty. Since we can learn the Full Ranking for this task, we use the following decision rule: for a given scenario, we rank a random subset of \mathcal{Y} (mixing the n candidates and $s.n$ negative records). We return the top-ranked record if it belongs to the candidate set, and the empty set otherwise.

For the Weather dataset, a sentence can have different target records depending on its candidate set so we must use another decision rule. We propose these two simple heuristics. For every scenario, (1) we rank among the candidate set only and always return the top-ranked element, (2) if this record is wind-related (i.e. with the type **windSpeed**, **windDir** or **windChill**), we also return the second element.

Table 2. **Top-k ranking accuracies on both datasets.** Results (in %) obtained by 4-fold cross-validations.

MODEL	ROBOCUP	WEATHER
Random Baseline	59.0	9.6
ARank	91.9	75.7

6. Experiments

In this section, we evaluate our ranking formulation, ARank, and compare it to two reference models for the tasks presented in Section 2.

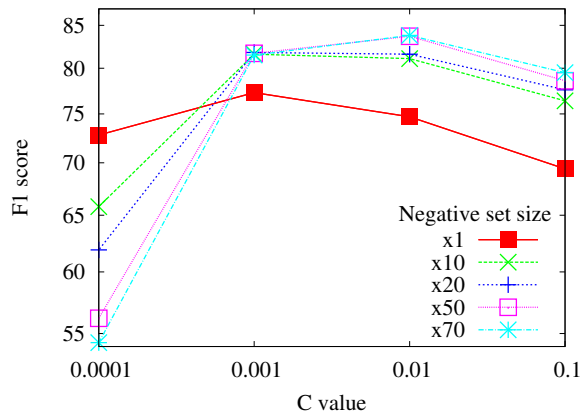


Figure 2. **Hyperparameter influence.** The y-axis displays the F₁ score of ARank on RoboCup. The x-axis (in log-scale) covers several values for C . Each curve represents a different size of the negative set, ranging from a size equal to the positive one ($\times 1$) to 70 times larger ($\times 70$).

Figure 2 displays the performances of ARank on RoboCup for different values of C and s that clearly indicates that hyperparameters have a non-trivial influence. However, the way to set them for ambiguously supervised systems is an open issue at the moment, because one cannot conduct any direct evaluation without using gold alignments. As this paper mainly targets to introduce our ranking approach, we leave this question for future work. We did not perform any exhaustive search but tried several reasonable values: 10^{-5} , 10^{-4} , 10^{-3} , 10^{-2} , 10^{-1} for C and 1, 10, 50 for s .⁴ In the following we present the results we obtained with the values listed in Table 1 (we kept the same value for all experiments concerning a dataset).

6.1. Ranking Evaluation

We first assess the ability of ARank to discover the hidden ranking within the candidate set. Splitting

⁴We also used $s = 20$ and $s = 70$ to draw Figure 2.

Table 3. **F₁ scores on RoboCup** obtained by 4-fold cross-validation as defined in (Chen & Mooney, 2008).

MODEL	F ₁
Random Baseline	48.0
Krisper	67.0
(Liang et al., 2009)	75.7
ARank	83.0

both datasets in four, our model is evaluated by cross-validation⁵ using a ranking metric: the *top-k accuracy*. For a given scenario, if k alignments must be predicted, we measure the proportion of these belonging to the top k elements of the list returned by the algorithm, and average that for all testing scenarios. The results displayed in Table 2, clearly indicate that ARank actually learns to correctly rank the candidates, even for a complex task like weather forecasting for which the random baseline is very low.

6.2. Alignment Comparison

In the literature, the standard evaluation metric for semantic matching is a F₁ score based on the number of actual gold alignments detected among the predictions. In this section, we use it to compare the predictions performed by ARank using the strategies defined in Section 5.3 to those of two state-of-the-art methods.

Baselines The first one is Krisper (Kate & Mooney, 2007) which obtained the best results on RoboCup in (Chen & Mooney, 2008). This algorithm works by repeatedly building noisy, unambiguous datasets from the ambiguous one, and training a parser designed for unambiguous supervision only. Recently, Liang et al. (2009) proposed a hierarchical hidden semi-Markov model for learning under ambiguous supervision directly. Their generative approach models the correspondences between text and records using latent variables and is trained with a sophisticated 3-stages process based on EM. They achieve the best current performances on both RoboCup and Weather.

Results In Table 3, we provide cross-validation scores on RoboCup. They express that ARank, despite its simple prediction strategy, attains strong performance behavior and outperforms both baselines, thanks to the good quality of the learnt ranking function. On the Weather data, ARank reaches a F₁ score of 76.4 in cross-validation but there is no cross-

⁵We used the split of (Chen & Mooney, 2008) for RoboCup and a random one for Weather.

Table 4. **Alignment results on both datasets.** Following (Liang et al., 2009), these results were obtained by training and testing on all scenarios. The table displays F_1 scores as well as [precision/ recall] values.

MODEL	ROBOCUP	WEATHER
Liang et al.	80.5 [77.3/ 84.0]	75.0 [76.3/ 73.8]
ARank	83.7 [76.6/ 92.3]	76.6 [78.0/ 75.3]

validated comparison available in the literature.

Indeed, in (Liang et al., 2009), in addition to cross-validation, another evaluation scheme is proposed for which models are both trained and tested on all scenarios. For both tasks, we report results in this setting in Table 4. They demonstrate that ARank remains very competitive. We can notice that its performances for both evaluation schemes are somewhat similar, unlike the method of (Liang et al., 2009) which lose almost 5% when being cross-validated on RoboCup.

Candidate Set vs Full Ranking In Section 5.3, we explain that for RoboCup we use a solution to the Full Ranking task to predict whereas for Weather we employ a solution to the Candidate Set Ranking one. That is the reason why the prediction strategy for RoboCup uses a ranking on mixed candidate and negative records and the one for weather involves only candidates. Now, if we reverse the strategies and use only candidate records for predicting on RoboCup and mix of candidates and negatives for Weather (and with no other parameter change), we observe that F_1 score respectively drops from 83.7 to 78.0 and from 76.6 to 72.4. This confirms the intuitions developed in Section 4.2 which tend to indicate that employing the Full Ranking is actually appropriate to predict on RoboCup but useless on Weather.

7. Conclusion

This paper casts a new light on the task of learning under ambiguous supervision: we demonstrated that solving a derived label ranking problem allows to perform a transfer of performance to the original task. As illustration, we empirically validated the efficiency of this approach by proposing a concrete application for learning semantic correspondences which happens to be very competitive with state-of-the-art methods.

References

Angelova, A., Matthies, L., Helmick, D. M., and Perona, P. Dimensionality Reduction Using Automatic

Supervision for Vision-Based Terrain Learning. In *Robotics: Science and Systems*. The MIT Press, 2007.

Branavan, S.R.K., Chen, H., Zettlemoyer, L., and Barzilay, R. Reinforcement learning for mapping instructions to actions. In *Proceedings of the 47th Annual Meeting of the ACL*, 2009.

Chen, D.L. and Mooney, R.J. Learning to Sportscast: A Test of Grounded Language Acquisition. In *Proceedings of ICML '08*, 2008.

Cour, T., Jordan, C., Miltsakaki, E., and Taskar, B. Movie/Script: Alignment and Parsing of Video and Text Transcription. In *Proceedings of ECCV '08*, 2008.

Cour, Timothee, Sapp, Ben, Jordan, Chris, and Taskar, Ben. Learning from ambiguously labeled images. In *Proceedings of CVPR*, 2009.

Ge, R. and Mooney, R. J. Learning a compositional semantic parser using an existing syntactic parser. In *Proceedings of the 47th Annual Meeting of the ACL*, 2009.

Har-Peled, S., Roth, D., and Zimak, D. Constraint Classification for Multiclass Classification and Ranking. In *Proceedings of Adv. in Neural Inf. Processing Syst.*, 2002.

Joachims, Thorsten. Training linear svms in linear time. In *Proceedings of the 12th ACM SIGKDD conference*, 2006.

Kate, R.J. and Mooney, R.J. Learning Language Semantics from Ambiguous Supervision. In *Proceedings of AAAI '07*, 2007.

Liang, P., Jordan, M. I., and Klein, D. Learning semantic correspondences with less supervision. In *Proceedings of the 47th Annual Meeting of the ACL*, 2009.

Mooney, R.J. Learning Semantic Parsers: An Important But Under-Studied Problem. In *Proceedings of AAAI '04*, 2004.

Snyder, B. and Barzilay, R. Database-text alignment via structured multilabel classification. In *In Proc. of the International Joint Conference on Artificial Intelligence*, pp. 1713–1718, 2007.

Zettlemoyer, L. and Collins, M. Learning context-dependent mappings from sentences to logical form. In *Proceedings of the 47th Annual Meeting of the ACL*, 2009.