

Automatic Optimization of Dialogue Management

Diane J. Litman, Michael S. Kearns, Satinder Singh, Marilyn A. Walker

AT&T Labs – Research

180 Park Avenue

Florham Park, NJ 07932 USA

{diane,mkearns,baveja,walker}@research.att.com

Abstract

Designing the dialogue strategy of a spoken dialogue system involves many nontrivial choices. This paper presents a reinforcement learning approach for automatically optimizing a dialogue strategy that addresses the technical challenges in applying reinforcement learning to a working dialogue system with human users. We then show that our approach measurably improves performance in an experimental system.

1 Introduction

Recent advances in spoken language understanding have made it possible to develop dialogue systems for many applications. The role of the dialogue manager in such systems is to interact in a natural way to help the user complete the tasks that the system is designed to support. Typically, an expert designs a dialogue manager by hand, and has to make many nontrivial design choices that can seriously impact system performance. This paper applies reinforcement learning (RL) to *automatically* learn design choices that optimize system performance for a chosen performance measure (Levin et al., 2000; Walker et al., 1998).

Consider the spoken dialogue system named NJFun, which we built to help users find fun places to go in New Jersey. A sample dialogue with NJFun is shown in Figure 1, with system utterances labeled S_i and user utterances labeled U_i . When NJFun greets the user, it can provide an open greeting such as utterance S_1 , expecting that the user will take the *initiative*. Alternatively, it can provide a directive greeting in which the system takes the initiative, as in “Welcome to NJFun. Please say an activity name or say ‘list activities’ for a list of activities I know about.” After the user’s response U_1 , NJFun must decide whether it should explicitly *confirm* its understanding, as in utterances S_2 and S_3 . NJFun can also simply continue on with the dialogue, as when it does not explicitly confirm that the user wants to find out about wineries. In NJFun, as shown in more detail below, decisions about initiative and confirmation strategies alone result in a

S_1 : Welcome to NJFun. How may I help you?
 U_1 : I'd like to find um winetasting in Lambertville in the morning. (*ASR output: I'd like to find out wineries the in the Lambertville in the morning.*)
 S_2 : Did you say you are interested in Lambertville?
 U_2 : Yes.
 S_3 : Did you say you want to go in the morning?
 U_3 : Yes.
 S_4 : I found a winery near Lambertville that is open in the morning. It is [...] Please give me feedback by saying ‘good’, ‘so-so’, or ‘bad’.
 U_4 : Good.

Figure 1: An example dialogue with NJFun.

search space of 2^{42} potential global dialogue strategies. Furthermore, the performance of a dialogue strategy depends on many other factors, such as the user population, the robustness of the automatic speech recognizer (ASR), and task difficulty (Kamm et al., 1998; Danieli and Gerbino, 1995).

In the main, previous research has treated the specification of the dialogue management strategy as an iterative design problem: several versions of a system are created, dialogue corpora are collected with human users interacting with different versions of the system, a number of evaluation metrics are collected for each dialogue, and the different versions are statistically compared (Danieli and Gerbino, 1995; Sanderman et al., 1998). Due to the costs of experimentation, only a few global strategies are typically explored in any one experiment.

However, recent work has suggested that dialogue strategy can be designed using the formalism of Markov decision processes (MDPs) and the algorithms of RL (Biermann and Long, 1996; Levin et al., 2000; Walker et al., 1998; Singh et al., 1999). More specifically, the MDP formalism suggests a method for optimizing dialogue strategies from sample dialogue data. The main advantage of this approach is the potential for computing an optimal dialogue strategy within a much larger search space, using a relatively small number of training dialogues.

This paper presents an application of RL to the

problem of optimizing dialogue strategy selection in the NJFun system, and experimentally demonstrates the utility of the approach. Section 2 explains how we apply RL to dialogue systems, then Section 3 describes the NJFun system in detail. Section 4 describes how NJFun optimizes its dialogue strategy from experimentally obtained dialogue data. Section 5 reports results from testing the learned strategy demonstrating that our approach improves task completion rates (our chosen measure for performance optimization). A companion paper provides only an abbreviated system and dialogue manager description, but includes additional results not presented here (Singh et al., 2000), such as analysis establishing the veracity of the MDP we learn, and comparisons of our learned strategy to strategies hand-picked by dialogue experts.

2 Reinforcement Learning for Dialogue

Due to space limitations, we present only a brief overview of how dialogue strategy optimization can be viewed as an RL problem; for more details, see Singh et al. (1999), Walker et al. (1998), Levin et al. (2000). A dialogue strategy is a mapping from a set of *states* (which summarize the entire dialogue so far) to a set of *actions* (such as the system’s utterances and database queries). There are multiple reasonable action choices in each state; typically these choices are made by the system designer. Our RL-based approach is to build a system that explores these choices in a systematic way through experiments with representative human users. A scalar performance measure, called a reward, is then calculated for each experimental dialogue. (We discuss various choices for this reward measure later, but in our experiments only terminal dialogue states have nonzero rewards, and the reward measures are quantities directly obtainable from the experimental set-up, such as user satisfaction or task completion.) This experimental data is used to construct an MDP which models the users’ interaction with the system. The problem of learning the best dialogue strategy from data is thus reduced to computing the optimal *policy* for choosing actions in an MDP — that is, the system’s goal is to take actions so as to maximize expected reward. The computation of the optimal policy *given* the MDP can be done efficiently using standard RL algorithms.

How do we build the desired MDP from sample dialogues? Following Singh et al. (1999), we can view a dialogue as a trajectory in the chosen state space determined by the system actions and user responses:

$$s_1 \rightarrow_{a_1, r_1} s_2 \rightarrow_{a_2, r_2} s_3 \rightarrow_{a_3, r_3} \dots$$

Here $s_i \rightarrow_{a_i, r_i} s_{i+1}$ indicates that at the i th ex-

change, the system was in state s_i , executed action a_i , received reward r_i , and then the state changed to s_{i+1} . Dialogue sequences obtained from training data can be used to empirically estimate the transition probabilities $P(s'|s, a)$ (denoting the probability of a transition to state s' , given that the system was in state s and took action a), and the reward function $R(s, a)$. The estimated transition probabilities and reward function constitute an MDP model of the user population’s interaction with the system.

Given this MDP, the expected cumulative reward (or *Q-value*) $Q(s, a)$ of taking action a from state s can be calculated in terms of the Q-values of successor states via the following recursive equation:

$$Q(s, a) = R(s, a) + \sum_{s'} P(s'|s, a) \max_{a'} Q(s', a').$$

These Q-values can be estimated to within a desired threshold using the standard RL *value iteration* algorithm (Sutton, 1991), which iteratively updates the estimate of $Q(s, a)$ based on the current Q-values of neighboring states. Once value iteration is completed, the optimal dialogue strategy (according to our estimated model) is obtained by selecting the action with the maximum Q-value at each dialogue state.

While this approach is theoretically appealing, the cost of obtaining sample human dialogues makes it crucial to limit the size of the state space, to minimize data sparsity problems, while retaining enough information in the state to learn an accurate model. Our approach is to work directly in a minimal but carefully designed state space (Singh et al., 1999).

The contribution of this paper is to empirically validate a practical methodology for using RL to *build* a dialogue system that optimizes its behavior from dialogue data. Our methodology involves 1) representing a dialogue strategy as a mapping from each state in the chosen state space \mathcal{S} to a set of dialogue actions, 2) deploying an initial training system that generates *exploratory* training data with respect to \mathcal{S} , 3) constructing an MDP model from the obtained training data, 4) using value iteration to learn the optimal dialogue strategy in the learned MDP, and 4) redeploying the system using the learned state/action mapping. The next section details the use of this methodology to design the NJFun system.

3 The NJFun System

NJFun is a real-time spoken dialogue system that provides users with information about things to do in New Jersey. NJFun is built using a general purpose platform for spoken dialogue systems (Levin et al., 1999), with support for modules for automatic speech recognition (ASR), spoken language

Action	Prompt
GreetS	Welcome to NJFun. Please say an activity name or say 'list activities' for a list of activities I know about.
GreetU	Welcome to NJFun. How may I help you?
ReAsk1S	I know about amusement parks, aquariums, cruises, historic sites, museums, parks, theaters, wineries, and zoos. Please say an activity name from this list.
ReAsk1M	Please tell me the activity type. You can also tell me the location and time.
Ask2S	Please say the name of the town or city that you are interested in.
Ask2U	Please give me more information.
ReAsk2S	Please tell me the name of the town or city that you are interested in.
ReAsk2M	Please tell me the location that you are interested in. You can also tell me the time.

Figure 2: Sample initiative strategy choices.

understanding, text-to-speech (TTS), database access, and dialogue management. NJFun uses a speech recognizer with stochastic language and understanding models trained from example user utterances, and a TTS system based on concatenative diphone synthesis. Its database is populated from the `nj.online` webpage to contain information about activities. NJFun indexes this database using three attributes: activity type, location, and time of day (which can assume values morning, afternoon, or evening).

Informally, the NJFun dialogue manager sequentially queries the user regarding the activity, location and time attributes, respectively. NJFun first asks the user for the current attribute (and possibly the other attributes, depending on the initiative). If the current attribute’s value is not obtained, NJFun asks for the attribute (and possibly the later attributes) again. If NJFun still does not obtain a value, NJFun moves on to the next attribute(s). Whenever NJFun successfully obtains a value, it can confirm the value, or move on to the next attribute(s). When NJFun has finished acquiring attributes, it queries the database (using a wildcard for each unobtained attribute value). The length of NJFun dialogues ranges from 1 to 12 user utterances before the database query. Although the NJFun dialogues are fairly short (since NJFun only asks for an attribute twice), the information access part of the dialogue is similar to more complex tasks.

As discussed above, our methodology for using RL to optimize dialogue strategy requires that all potential actions for each state be specified. Note that at some states it is easy for a human to make the correct action choice. We made obvious dialogue strategy choices in advance, and used learning only to optimize the difficult choices (Walker et al., 1998). In NJFun, we restricted the action choices to 1) the type of initiative to use when asking or reasking for an attribute, and 2) whether to confirm an attribute value once obtained. The optimal actions may vary with dialogue state, and are subject to active debate in the literature.

The examples in Figure 2 show that NJFun can ask the user about the first 2 attributes¹ using three types of initiative, based on the combination of the wording of the system prompt (*open* versus *directive*), and the type of grammar NJFun uses during ASR (*restrictive* versus *non-restrictive*). If NJFun uses an open question with an unrestricted grammar, it is using *user initiative* (e.g., GreetU). If NJFun instead uses a directive prompt with a restricted grammar, the system is using *system initiative* (e.g., GreetS). If NJFun uses a directive question with a non-restrictive grammar, it is using *mixed initiative*, because it allows the user to take the initiative by supplying extra information (e.g., ReAsk1M).

NJFun can also vary the strategy used to confirm each attribute. If NJFun asks the user to explicitly verify an attribute, it is using *explicit confirmation* (e.g., ExpConf2 for the location, exemplified by S2 in Figure 1). If NJFun does not generate any confirmation prompt, it is using *no confirmation* (the NoConf action).

Solely for the purposes of controlling its operation (as opposed to the learning, which we discuss in a moment), NJFun internally maintains an *operations vector* of 14 variables. 2 variables track whether the system has greeted the user, and which attribute the system is currently attempting to obtain. For each of the 3 attributes, 4 variables track whether the system has obtained the attribute’s value, the system’s confidence in the value (if obtained), the number of times the system has asked the user about the attribute, and the type of ASR grammar most recently used to ask for the attribute.

The formal state space \mathcal{S} maintained by NJFun for the purposes of learning is much simpler than the operations vector, due to the data sparsity concerns already discussed. The dialogue state space \mathcal{S} contains only 7 variables, as summarized in Figure 3. \mathcal{S} is computed from the operations vector using a hand-designed algorithm. The “greet” variable

¹“Greet” is equivalent to asking for the first attribute. NJFun always uses system initiative for the third attribute, because at that point the user can only provide the time of day.

greet	attr	conf	val	times	gram	hist
0,1	1,2,3,4	0,1,2,3,4	0,1	0,1,2	0,1	0,1

Figure 3: State features and values.

tracks whether the system has greeted the user or not (no=0, yes=1). “Attr” specifies which attribute NJFun is currently attempting to obtain or verify (activity=1, location=2, time=3, done with attributes=4). “Conf” represents the confidence that NJFun has after obtaining a value for an attribute. The values 0, 1, and 2 represent the lowest, middle and highest ASR confidence values.² The values 3 and 4 are set when ASR hears “yes” or “no” after a confirmation question. “Val” tracks whether NJFun has obtained a value for the attribute (no=0, yes=1). “Times” tracks the number of times that NJFun has asked the user about the attribute. “Gram” tracks the type of grammar most recently used to obtain the attribute (0=non-restrictive, 1=restrictive). Finally, “hist” (history) represents whether NJFun had trouble understanding the user in the earlier part of the conversation (bad=0, good=1). We omit the full definition, but as an example, when NJFun is working on the second attribute (location), the history variable is set to 0 if NJFun does not have an activity, has an activity but has no confidence in the value, or needed two queries to obtain the activity.

As mentioned above, the goal is to design a small state space that makes enough critical distinctions to support learning. The use of \mathcal{S} reduces the number of states to only 62, and supports the construction of an MDP model that is not sparse with respect to \mathcal{S} , even using limited training data.³ The state space that we utilize here, although minimal, allows us to make initiative decisions based on the success of earlier exchanges, and confirmation decisions based on ASR confidence scores and grammars.

The state/action mapping representing NJFun’s initial dialogue strategy EIC (Exploratory for Initiative and Confirmation) is in Figure 4. Only the exploratory portion of the strategy is shown, namely those states for which NJFun has an action choice. For each such state, we list the two choices of actions available. (The action choices in boldface are the ones eventually identified as optimal by the learning process, and are discussed in detail later.) The EIC strategy chooses *randomly* between these two ac-

²For each utterance, the ASR output includes not only the recognized string, but also an associated acoustic confidence score. Based on data obtained during system development, we defined a mapping from raw confidence values into 3 approximately equally populated partitions.

³62 refers to those states that can actually occur in a dialogue. For example, greet=0 is only possible in the initial dialogue state “0 1 0 0 0 0”. Thus, all other states beginning with 0 (e.g. “0 1 0 0 1 0 0”) will never occur.

State							Action Choices
g	a	c	v	t	g	h	
0	1	0	0	0	0	0	GreetS , GreetU
1	1	0	0	1	0	0	ReAsk1S , ReAsk1M
1	1	0	1	0	0	0	NoConf, ExpConf1
1	1	0	1	0	1	0	NoConf, ExpConf1
1	1	1	1	0	0	0	NoConf, ExpConf1
1	1	1	1	0	1	0	NoConf, ExpConf1
1	1	2	1	0	0	0	NoConf, ExpConf1
1	1	2	1	0	1	0	NoConf, ExpConf1
1	1	4	0	0	0	0	ReAsk1S , ReAsk1M
1	1	4	0	1	0	0	ReAsk1S , ReAsk1M
1	2	0	0	0	0	0	Ask2S , Ask2U
1	2	0	0	0	0	1	Ask2S , Ask2U
1	2	0	0	1	0	0	ReAsk2S , ReAsk2M
1	2	0	0	1	0	1	ReAsk2S , ReAsk2M
1	2	0	1	0	0	0	NoConf, ExpConf2
1	2	0	1	0	0	1	NoConf, ExpConf2
1	2	0	1	0	1	0	NoConf, ExpConf2
1	2	0	1	0	1	1	NoConf, ExpConf2
1	2	1	1	0	0	0	NoConf, ExpConf2
1	2	1	1	0	0	1	NoConf, ExpConf2
1	2	1	1	0	1	0	NoConf, ExpConf2
1	2	1	1	0	1	1	NoConf, ExpConf2
1	2	2	1	0	0	0	NoConf, ExpConf2
1	2	2	1	0	0	1	NoConf, ExpConf2
1	2	2	1	0	1	0	NoConf, ExpConf2
1	2	2	1	0	1	1	NoConf, ExpConf2
1	2	4	0	0	0	0	ReAsk2S , ReAsk2M
1	2	4	0	0	0	1	ReAsk2S , ReAsk2M
1	2	4	0	1	0	0	ReAsk2S , ReAsk2M
1	2	4	0	1	0	1	ReAsk2S , ReAsk2M
1	3	0	1	0	0	0	NoConf, ExpConf3
1	3	0	1	0	0	1	NoConf, ExpConf3
1	3	0	1	0	1	0	NoConf, ExpConf3
1	3	0	1	0	1	1	NoConf, ExpConf3
1	3	1	1	0	0	0	NoConf, ExpConf3
1	3	1	1	0	0	1	NoConf, ExpConf3
1	3	1	1	0	1	0	NoConf, ExpConf3
1	3	1	1	0	1	1	NoConf, ExpConf3
1	3	2	1	0	0	0	NoConf, ExpConf3
1	3	2	1	0	0	1	NoConf, ExpConf3
1	3	2	1	0	1	0	NoConf, ExpConf3
1	3	2	1	0	1	1	NoConf, ExpConf3

Figure 4: Exploratory portion of EIC strategy.

tions in the indicated state, to maximize exploration and minimize data sparseness when constructing our model. Since there are 42 states with 2 choices each, there is a search space of 2^{42} potential global dialogue strategies; the goal of RL is to identify an apparently optimal strategy from this large search space. Note that due to the randomization of the EIC strategy, the prompts are designed to ensure the coherence of all possible action sequences.

Figure 5 illustrates how the dialogue strategy in Figure 4 generates the dialogue in Figure 1. Each row indicates the state that NJFun is in, the ac-

State	Action	Turn	Reward
g a c v t g h			
0 1 0 0 0 0 0	GreetU	S1	0
1 1 2 1 0 0 0	NoConf	-	0
1 2 2 1 0 0 1	ExpConf2	S2	0
1 3 2 1 0 0 1	ExpConf3	S3	0
1 4 0 0 0 0 0	Tell	S4	1

Figure 5: Generating the dialogue in Figure 1.

tion executed in this state, the corresponding turn in Figure 1, and the reward received. The initial state represents that NJFun will first attempt to obtain attribute 1. NJFun executes GreetU (although as shown in Figure 4, GreetS is also possible), generating the first utterance in Figure 1. After the user’s response, the next state represents that NJFun has now greeted the user and obtained the activity value with high confidence, by using a non-restrictive grammar. NJFun then chooses the NoConf strategy, so it does not attempt to confirm the activity, which causes the state to change but no prompt to be generated. The third state represents that NJFun is now working on the second attribute (location), that it already has this value with high confidence (location was obtained with activity after the user’s first utterance), and that the dialogue history is good.⁴ This time NJFun chooses the ExpConf2 strategy, and confirms the attribute with the second NJFun utterance, and the state changes again. The processing of time is similar to that of location, which leads NJFun to the final state, where it performs the action “Tell” (corresponding to querying the database, presenting the results to the user, and asking the user to provide a reward). Note that in NJFun, the reward is always 0 except at the terminal state, as shown in the last column of Figure 5.

4 Experimentally Optimizing a Strategy

We collected experimental dialogues for both training and testing our system. To obtain training dialogues, we implemented NJFun using the EIC dialogue strategy described in Section 3. We used these dialogues to build an empirical MDP, and then computed the optimal dialogue strategy in this MDP (as described in Section 2). In this section we describe our experimental design and the learned dialogue strategy. In the next section we present results from testing our learned strategy and show that it improves task completion rates, the performance measure we chose to optimize.

Experimental subjects were employees not associated with the NJFun project. There were 54 sub-

⁴Recall that only the current attribute’s features are in the state. However, the operations vector contains information regarding previous attributes.

jects for training and 21 for testing. Subjects were distributed so the training and testing pools were balanced for gender, English as a first language, and expertise with spoken dialogue systems.

During both training and testing, subjects carried out free-form conversations with NJFun to complete six application tasks. For example, the task executed by the user in Figure 1 was: “You feel thirsty and want to do some winetasting in the morning. Are there any wineries close by your house in Lambertville?” Subjects read task descriptions on a web page, then called NJFun from their office phone. At the end of the task, NJFun asked for feedback on their experience (e.g., utterance S4 in Figure 1). Users then hung up the phone and filled out a user survey (Singh et al., 2000) on the web.

The training phase of the experiment resulted in 311 complete dialogues (not all subjects completed all tasks), for which NJFun logged the sequence of states and the corresponding executed actions. The number of samples per state for the initial ask choices are:

```
0 1 0 0 0 0 0  GreetS=155  GreetU=156
1 2 0 0 0 0 0  Ask2S=93   Ask2U=72
1 2 0 0 0 0 1  Ask2S=36   Ask2U=48
```

Such data illustrates that the random action choice strategy led to a fairly balanced action distribution per state. Similarly, the small state space, and the fact that we only allowed 2 action choices per state, prevented a data sparseness problem. The first state in Figure 4, the initial state for every dialogue, was the most frequently visited state (with 311 visits). Only 8 states that occur near the end of a dialogue were visited less than 10 times.

The logged data was then used to construct the empirical MDP. As we have mentioned, the measure we chose to optimize is a binary reward function based on the strongest possible measure of task completion, called **StrongComp**, that takes on value 1 if NJFun queries the database using exactly the attributes specified in the task description, and 0 otherwise. Then we computed the optimal dialogue strategy in this MDP using RL (cf. Section 2). The action choices constituting the learned strategy are in boldface in Figure 4. Note that no choice was fixed for several states, meaning that the Q-values were identical after value iteration. Thus, even when using the learned strategy, NJFun still sometimes chooses randomly between certain action pairs.

Intuitively, the learned strategy says that the optimal use of initiative is to begin with user initiative, then back off to either mixed or system initiative when reasking for an attribute. Note, however, that the specific backoff method differs with attribute (e.g., system initiative for attribute 1, but generally mixed initiative for attribute 2). With respect to confirmation, the optimal strategy is to

mainly confirm at lower confidence values. Again, however, the point where confirmation becomes unnecessary differs across attributes (e.g., confidence level 2 for attribute 1, but sometimes lower levels for attributes 2 and 3), and also depends on other features of the state besides confidence (e.g., grammar and history). This use of ASR confidence by the dialogue strategy is more sophisticated than previous approaches, e.g. (Niimi and Kobayashi, 1996; Litman and Pan, 2000). NJFun can learn such fine-grained distinctions because the optimal strategy is based on a comparison of 2^{42} possible exploratory strategies. Both the initiative and confirmation results suggest that the beginning of the dialogue was the most problematic for NJFun. Figure 1 is an example dialogue using the optimal strategy.

5 Experimentally Evaluating the Strategy

For the testing phase, NJFun was reimplemented to use the learned strategy. 21 test subjects then performed the same 6 tasks used during training, resulting in 124 complete test dialogues. One of our main results is that task completion as measured by StrongComp increased from 52% in training to 64% in testing ($p < .06$).⁵

There is also a significant interaction effect between strategy and task ($p < .01$) for StrongComp. Previous work has suggested that novice users perform comparably to experts after only 2 tasks (Kamm et al., 1998). Since our learned strategy was based on 6 tasks with each user, one explanation of the interaction effect is that the learned strategy is slightly optimized for expert users. To explore this hypothesis, we divided our corpus into dialogues with “novice” (tasks 1 and 2) and “expert” (tasks 3-6) users. We found that the learned strategy did in fact lead to a large and significant improvement in StrongComp for experts (EIC=.46, learned=.69, $p < .001$), and a non-significant degradation for novices (EIC=.66, learned=.55, $p < .3$).

An apparent limitation of these results is that EIC may not be the best baseline strategy for comparison to our learned strategy. A more standard alternative would be comparison to the very best hand-designed fixed strategy. However, there is no agreement in the literature, nor amongst the authors, as to what the best hand-designed strategy might have been. There is agreement, however, that the best strategy is sensitive to many unknown and unmodeled factors: the

⁵The experimental design described above consists of 2 factors: the within-in group factor *strategy* and the between-groups factor *task*. We use a two-way analysis of variance (ANOVA) to compute whether main and interaction effects of strategy are statistically significant ($p < .05$) or indicative of a statistical trend ($p < .10$). *Main effects* of strategy are task-independent, while *interaction* effects involving strategy are task-dependent.

Measure	EIC (n=311)	Learned (n=124)	p
StrongComp	0.52	0.64	.06
WeakComp	1.75	2.19	.02
ASR	2.50	2.67	.04
Feedback	0.18	0.11	.42
UserSat	13.38	13.29	.86

Table 1: Main effects of dialogue strategy.

user population, the specifics of the task, the particular ASR used, etc. Furthermore, EIC was carefully designed so that the random choices it makes never results in an unnatural dialogue. Finally, a companion paper (Singh et al., 2000) shows that the performance of the learned strategy is better than several “standard” fixed strategies (such as always use system-initiative and no-confirmation).

Although many types of measures have been used to evaluate dialogue systems (e.g., task success, dialogue quality, efficiency, usability (Danieli and Gerbino, 1995; Kamm et al., 1998)), we optimized only for one task success measure, StrongComp. However, we also examined the performance of the learned strategy using other evaluation measures (which possibly could have been used as our reward function). **WeakComp** is a relaxed version of task completion that gives partial credit: if all attribute values are either correct or wildcards, the value is the sum of the correct number of attributes. Otherwise, at least one attribute is wrong (e.g., the user says “Lambertville” but the system hears “Morristown”), and the value is -1. **ASR** is a dialogue quality measure that approximates speech recognition accuracy for the database query, and is computed by adding 1 for each correct attribute value and .5 for every wildcard. Thus, if the task is to go winetasting near Lambertville in the morning, and the system queries the database for an activity in New Jersey in the morning, StrongComp=0, WeakComp=1, and ASR=2. In addition to the objective measures discussed above, we also computed two *subjective* usability measures. **Feedback** is obtained from the dialogue (e.g. S4 in Figure 5), by mapping *good*, *so-so*, *bad* to 1, 0, and -1, respectively. User satisfaction (**UserSat**, ranging from 0-20) is obtained by summing the answers of the web-based user survey.

Table 1 summarizes the difference in performance of NJFun for our original reward function and the above alternative evaluation measures, from training (EIC) to test (learned strategy for StrongComp). For WeakComp, the average reward increased from 1.75 to 2.19 ($p < 0.02$), while for ASR the average reward increased from 2.5 to 2.67 ($p < 0.04$). Again, these improvements occur even though the learned strategy was not optimized for these measures.

The last two rows of the table show that for the

subjective measures, performance does not significantly differ for the EIC and learned strategies. Interestingly, the distributions of the subjective measures move to the middle from training to testing, i.e., test users reply to the survey using less extreme answers than training users. Explaining the subjective results is an area for future work.

6 Discussion

This paper presents a practical methodology for applying RL to optimizing dialogue strategies in spoken dialogue systems, and shows empirically that the method improves performance over the EIC strategy in NJFun. A companion paper (Singh et al., 2000) shows that the learned strategy is not only better than EIC, but also better than other fixed choices proposed in the literature. Our results demonstrate that the application of RL allows one to empirically optimize a system’s dialogue strategy by searching through a much larger search space than can be explored with more traditional methods (i.e. empirically testing several versions of a system).

RL has been applied to dialogue systems in previous work, but our approach differs from previous work in several respects. Biermann and Long (1996) did not test RL in an implemented system, and the experiments of Levin et al. (2000) utilized a simulated user model. Walker et al. (1998)’s methodology is similar to that used here, in testing RL with an implemented system with human users. However that work only explored strategy choices at 13 states in the dialogue, which conceivably could have been explored with more traditional methods (as compared to the 42 choice states explored here).

We also note that our learned strategy made dialogue decisions based on ASR confidence in conjunction with other features, and also varied initiative and confirmation decisions at a finer grain than previous work; as such, our learned strategy is not a standard strategy investigated in the dialogue system literature. For example, we would not have predicted the complex and interesting back-off strategy with respect to initiative when reasking for an attribute.

To see how our method scales, we are applying RL to dialogue systems for customer care and for travel planning, which are more complex task-oriented domains. As future work, we wish to understand the aforementioned results on the subjective reward measures, explore the potential difference between optimizing for expert users and novices, automate the choice of state space for dialogue systems, investigate the use of a learned reward function (Walker et al., 1998), and explore the use of more informative non-terminal rewards.

Acknowledgements

The authors thank Fan Jiang for his substantial effort in implementing NJFun, Wieland Eckert, Esther Levin, Roberto Pieraccini, and Mazin Rahim for their technical help, Julia Hirschberg for her comments on a draft of this paper, and David McAllester, Richard Sutton, Esther Levin and Roberto Pieraccini for helpful conversations.

References

- A. W. Biermann and P. M. Long. 1996. The composition of messages in speech-graphics interactive systems. In *Proc. of the International Symposium on Spoken Dialogue*, pages 97–100.
- M. Danieli and E. Gerbino. 1995. Metrics for evaluating dialogue strategies in a spoken language system. In *Proc. of the AAAI Spring Symposium on Empirical Methods in Discourse Interpretation and Generation*, pages 34–39.
- C. Kamm, D. Litman, and M. A. Walker. 1998. From novice to expert: The effect of tutorials on user expertise with spoken dialogue systems. In *Proc. of the International Conference on Spoken Language Processing, ICSLP98*.
- E. Levin, R. Pieraccini, W. Eckert, G. Di Fabbrizio, and S. Narayanan. 1999. Spoken language dialogue: From theory to practice. In *Proc. IEEE Workshop on Automatic Speech Recognition and Understanding, ASRU99*.
- E. Levin, R. Pieraccini, and W. Eckert. 2000. A stochastic model of human machine interaction for learning dialog strategies. *IEEE Transactions on Speech and Audio Processing*, 8(1):11–23.
- D. J. Litman and S. Pan. 2000. Predicting and adapting to poor speech recognition in a spoken dialogue system. In *Proc. of the Seventeenth National Conference on Artificial Intelligence, AAAI-2000*.
- Y. Niimi and Y. Kobayashi. 1996. A dialog control strategy based on the reliability of speech recognition. In *Proc. of the International Symposium on Spoken Dialogue*, pages 157–160.
- A. Sanderman, J. Sturm, E. den Os, L. Boves, and A. Cremers. 1998. Evaluation of the dutchtrain timetable information system developed in the arise project. In *Interactive Voice Technology for Telecommunications Applications, IVTTA*, pages 91–96.
- S. Singh, M. S. Kearns, D. J. Litman, and M. A. Walker. 1999. Reinforcement learning for spoken dialogue systems. In *Proc. NIPS99*.
- S. B. Singh, M. S. Kearns, D. J. Litman, and M. A. Walker. 2000. Empirical evaluation of a reinforcement learning spoken dialogue system. In *Proc. of the Seventeenth National Conference on Artificial Intelligence, AAAI-2000*.
- R. S. Sutton. 1991. Planning by incremental dynamic programming. In *Proc. Ninth Conference on Machine Learning*, pages 353–357.
- M. A. Walker, J. C. Fromer, and S. Narayanan. 1998. Learning optimal dialogue strategies: A case study of a spoken dialogue agent for email. In *Proc. of the 36th Annual Meeting of the Association of Computational Linguistics, COLING/ACL 98*, pages 1345–1352.