

Abstract which is the only  
one.

## Efficiently Ranking Hypotheses in Machine Learning

Steve Chien, Andre Stechert, and Darren Mutz

Jet Propulsion Laboratory, California Institute of Technology  
4800 (irk Grove Drive, M/S 525-3660, Pasadena, CA 91 109-8099

Hypothesis ranking problems are an abstract class of learning problems where an algorithm is given a set of hypotheses to rank according to *expected utility* over some unknown distribution, where the expected utility must be estimated from training data.

Hypothesis ranking problems are an extension of hypothesis selection problems (Chien95), in which a learning system attempts to select the best alternative hypothesis from a set of hypotheses. The distinction between hypothesis ranking and hypothesis selection is that in selection the learning algorithm is interested in a *single best hypothesis*, while in ranking the learning algorithm must determine the relative order of all of the hypotheses<sup>1</sup>.

Hypothesis evaluation is an important aspect of many machine learning problems. For example, the utility problem in speedup learning can be viewed as a selection problem where a single problem-solving heuristic or strategy is chosen from a larger set of candidates. In this case, the expected utility is typically defined as the average time to solve a problem. The attribute selection problem in machine learning can also be viewed as a hypothesis selection problem in which one must select the best attribute split from a set of possible attribute splits and utility is often measured by information gain. In reinforcement learning, a system must learn the appropriate action for each context, where utility is interpreted as expected reward (with immediate feedback).

In many of these applications, a system chooses a single alternative and never revisits the decision. In contrast, if the system is able to investigate several options (either serially or in parallel), such as in beam search or iterative broadening, the ranking formulation is most appropriate. Also, as is the case with evolutionary approaches, a system may need to populate future alternative hypotheses on the basis of the ranking of the current population (Goldberg89).

In any hypothesis evaluation problem, always achieving a correct ranking is impossible in practice, because the actual underlying probability distributions are unavailable and

there is always a (perhaps vanishingly) small chance that the algorithms will be unlucky because only a finite number of samples can be taken. Consequently, rather than always requiring an algorithm to output a correct ranking, we impose probabilistic criteria on the rankings to be produced. While several families of such requirements exist, in this paper we examine two, the *probably approximately correct* (PAC) requirement from the computational learning theory community (Valiant84) and the *expected loss* (EL) requirement frequently used in decision theory and gaming problems (?). With the PAC requirement, an algorithm produces a ranking that with high probability is close to correct (e.g., incorrect orderings are between hypotheses with similar expected utilities). The EL requirement bounds the expected loss, where loss represents the difference in utilities between two incorrectly ordered hypotheses.

The principal contributions of this paper are:

- We define two families of hypothesis ranking algorithms based on recursive selection and adjacency. We provide specific details on how to apply them to a probably approximately correct (PAC) and expected loss (EL) decision criteria.
- We provide empirical results demonstrating the effectiveness of these algorithms at achieving requested decision criteria on synthetic data.
- We provide empirical results showing how these algorithms significantly outperform existing statistical methods on real-world data from a spacecraft design optimization application.

*Ranking as Recursive Selection:* One obvious way to determine a ranking  $H_1, \dots, H_k$  is to view ranking as recursive selection from the set of remaining candidate hypotheses. In this view, the overall ranking error, as specified by the desired confidence in PAC algorithms and the loss threshold in EL algorithms, is first distributed among  $k$  — *selection errors* which are then further subdivided into *pairwise comparison errors*. Data is then sampled until the estimates of the pairwise comparison error (as dictated by equation ?? or ??) satisfy the bounds set by the algorithm.

*Ranking by Adjacency Comparison:* Another interpretation of ranking confidence (or loss) is that only adjacent elements in the ranking need be compared. In this case, the

<sup>1</sup>This paper describes research conducted by the Jet Propulsion Laboratory, California Institute of Technology, under contract with the National Aeronautics and Space Administration.

<sup>2</sup>The algorithms and results described in this paper trivially extend to hybrid ranking-selection problems in which the system must select and rank the top  $M$  out of  $N$  hypotheses.

overall ranking error is divided directly into  $k - 1$  pairwise comparison errors. This leads to the following confidence equation for the PAC criteria:

### References

S. A. Chien, J. M. Gratch and M. C. Burl, "On the Efficient Allocation of Resources for 1 hypothesis Evaluation: A Statistical Approach," *IEEE Trans. Pattern Analysis and Machine Intelligence* **17 (7)**, July 1995, pp. 652-665.

D. Goldberg, *Genetic Algorithms in Search, Optimization and Machine Learning*, Addison Wesley, 1989.

Turnbull and Weiss, "A class of sequential procedures for k-sample problems concerning normal means with unknown unequal variances," in *Design of Experiments: ranking and selection*, T. J. Santner and A. C. Tamhane (eds.), Marcel Dekker, 1984.

L. G. Valiant, "A Theory of the Learnable," *Communications of the ACM* **27**, (1984), pp. 1134-1142.

*Full paper, will be made available  
later, submitted for publication  
only. Not to be used with IEEE*

## Efficiently Ranking Hypotheses in Machine Learning

**Abstract ID: A827**

Content Areas: Machine Learning

### Abstract

This paper considers the problem of learning the ranking of a set of alternatives based upon incomplete information (e.g., a limited number of observations). At each decision cycle, the system can output a complete ordering on the hypotheses or decide to gather additional information (e.g., observations) at some cost. Balancing the expected utility of the additional information against the cost of acquiring the information is the central problem we address.

The hypothesis ranking problem is a generalization of the previously studied hypothesis selection problem - in selection, an algorithm must select the single best hypothesis, while in ranking, an algorithm must order all the hypotheses. We describe two algorithms for hypothesis ranking and their application for probably approximately correct (PAC) and expected loss (EL) learning criteria. Empirical results are provided to demonstrate the effectiveness of these ranking procedures on both synthetic datasets and real-world data from a spacecraft design optimization application.

### Introduction

In many learning applications, the cost of information can be quite high, imposing a requirement that the learning algorithms glean as much usable information as possible with a minimum of data. For example:

- in speedup learning, the expense of processing each training example can be significant (Tadepalli92).
- In decision tree learning, the cost of using all available training examples when evaluating potential attributes for partitioning can be computationally expensive (Musick93).
- In evaluating medical treatment policies, additional training examples imply suboptimal treatment of human subjects.
- In data-poor applications, training data may be very scarce and learning as well as possible from limited data may be key.

When one wishes some sort of guarantee on the quality of a solution, a statistical decision theoretic framework is useful. The framework answers the questions: How much

information is enough? At what point do we have adequate information to rank the alternatives with some requested confidence?

This paper focuses on parametric ranking problems, a general class of statistical machine learning problems in which the goal is to rank a set of alternative hypotheses where the goodness of a hypothesis is a function of a set of unknown parameters (e.g., (Gratch92; Greiner92; Kaelbling93; Moore94; Musick93)). The learning system determines and refines estimates of these parameters by using training examples, with a secondary goal of minimizing learning cost.

The principal contributions of this paper are:

- We define two families of hypothesis ranking algorithms based on recursive selection and adjacency. We provide specific details on how to apply them to a probably approximately correct (PAC) and expected loss (EL) decision criteria.
- We provide empirical results demonstrating the effectiveness of these algorithms at achieving requested decision criteria on synthetic data.
- We provide empirical results showing how these algorithms significantly outperform existing statistical methods on real-world data from a spacecraft design optimization application.

The remainder of this paper is structured as follows. First, we describe the hypothesis ranking problem more formally, including definitions for the probably approximately correct (PAC) and expected loss (EL) decision criteria. We then define two algorithms for establishing these criteria for the hypothesis ranking problem - a recursive hypothesis selection algorithm and an adjacency based algorithm. Next, we describe empirical tests demonstrating the effectiveness of these algorithms as well as documenting their improved performance over a standard algorithm from the statistical ranking literature. Finally, we describe related work and future extensions to the algorithms.

### Hypothesis Ranking Problems

Hypothesis ranking problems are an abstract class of learning problems where an algorithm is given a set of hypotheses

to rank according to *expected utility* over some unknown distribution, where the expected utility must be estimated from training data.

Hypothesis ranking problems are an extension of hypothesis selection problems (Chien95), in which a learning system attempts to select the best alternative hypothesis from a set of hypotheses. The distinction between hypothesis ranking and hypothesis selection is that in selection the learning algorithm is interested in a single best hypothesis, while in ranking the learning algorithm must determine the relative order of all of the hypotheses<sup>1</sup>.

Hypothesis evaluation is an important aspect of many machine learning problems. For example, the utility problem in speedup learning can be viewed as a selection problem where a single problem-solving heuristic or strategy is chosen from a larger set of candidates. In this case, the expected utility is typically defined as the average time to solve a problem (Gratch92; Greiner92; Minton88). The attribute selection problem in machine learning can also be viewed as a hypothesis selection problem in which one must select the best attribute split from a set of possible attribute splits and utility is often measured by information gain (Musick93). In reinforcement learning, a system must learn the appropriate action for each context, where utility is interpreted as expected reward (Kaelbling93)<sup>2</sup>

in many of these applications, a system chooses a single alternative and never revisits the decision. In contrast, if the system is able to investigate several options (either serially or in parallel), such as in beam search or iterative broadening, the ranking formulation is most appropriate. Also, as is the case with evolutionary approaches, a system may need to populate future alternative hypotheses on the basis of the ranking of the current population (Goldberg89).

In any hypothesis evaluation problem, always achieving a correct ranking is impossible in practice, because the actual underlying probability distributions are unavailable and there is always a (perhaps vanishingly) small chance that the algorithms will be unlucky because only a finite number of samples can be taken. Consequently, rather than always requiring an algorithm to output a correct ranking, we impose probabilistic criteria on the rankings to be produced. While several families of such requirements exist, in this paper we examine two, the *probably approximately correct* (PAC) requirement from the computational learning theory community (Valiant84) and the *expected loss* (EL) requirement frequently used in decision theory and gaming problems (Russell92). With the PAC requirement, an algorithm produces a ranking that with high probability is close to correct (e.g., incorrect orderings are between hypotheses with similar expected utilities). The EL requirement bounds the expected loss, where loss represents the difference in utilities

<sup>1</sup>The algorithms and results described in this paper trivially extend to hybrid ranking-selection problem in which the system must select and rank the top  $M$  out of  $N$  hypotheses.

<sup>2</sup>Note that the analogous reinforcement learning problem is the one in which we are learning the appropriate action with immediate feedback rather than delayed feedback.

between two incorrectly ordered hypotheses.

The expected utility of a hypothesis can be estimated by observing its values over a finite set of training examples. However, to satisfy the PAC and EL requirements, an algorithm must also be able to reason about the potential difference between the estimated and true utilities of each hypothesis. Let  $U_i$  be the true expected utility of hypothesis  $i$  and let  $\hat{U}_i$  be the estimated expected utility of hypothesis  $i$ . Without loss of generality, let us presume that the proposed ranking of hypotheses is  $U_1 > U_2 > \dots > U_{k-1} > U_k$ . The PAC requirement states that for some user-specified  $\epsilon$  with probability  $1 - \delta$ :

$$\bigwedge_{i=1}^{k-1} [(U_i + \epsilon) > \text{MAX}(U_{i+1}, \dots, U_k)] \quad (1)$$

Correspondingly, let the loss  $L$  of selecting a hypothesis  $H_1$  to be the best from a set of  $k$  hypotheses  $H_1, \dots, H_k$  be as follows.

$$L(H_1, \{H_1, \dots, H_k\}) = \text{MAX}(0, \text{MAX}(U_2, \dots, U_k) - U_1) \quad (2)$$

and let the loss  $RL$  of a ranking  $H_1, \dots, H_k$  be as follows.

$$RL(H_1, \dots, H_k) = \sum_{i=1}^{k-1} L(H_i, \{H_{i+1}, \dots, H_k\}) \quad (3)$$

A hypothesis ranking algorithm which obeys the expected loss requirement must produce rankings that on average have less than the requested expected loss bound. Consider ranking the hypotheses with expected utilities:  $U_1 = 1.0$ ,  $U_2 = 0.95$ ,  $U_3 = 0.86$ . The ranking  $U_2 > U_1 > U_3$  is a valid PAC ranking for  $\epsilon = 0.06$  but not for  $\epsilon = 0.01$  and has an observed loss of  $0.05 + O = 0.05$ .

However, while the confidence in a pairwise comparison between two hypotheses is well understood, it is less clear how to ensure that desired confidence is met in the set of comparisons required for a selection or the more complex set of comparisons required for a ranking. Equation 4 defines the confidence that  $U_i + \epsilon > U_j$ , **when** the distribution underlying the utilities is normally distributed with unknown and unequal variances.

$$\gamma = \phi \left( \frac{(\hat{U}_{i-j} + \epsilon) \sqrt{n}}{\hat{S}_{i-j}} \right) \quad (4)$$

where  $\phi$  represents the cumulative standard normal distribution function,  $n$ ,  $\hat{U}_{i-j}$ , and  $\hat{S}_{i-j}$  are the size, sample mean, and sample standard deviation of the blocked differential distribution, respectively<sup>3</sup>.

Likewise, computation of the expected loss for asserting an ordering between a pair of hypotheses is well understood, but the estimation of expected loss for an entire ranking is less clear. Equation 5 defines the expected loss for drawing

<sup>3</sup>Note that in our approach we *block* examples to further reduce sampling complexity. Blocking forms estimates by using the difference in utility between competing hypotheses on each observed example. Blocking can significantly reduce the variance in the data when the hypotheses are not independent. It is trivial to modify the formulas to address the cases in which it is not possible to block data (see (Moore94; Chien95) for farther details).

the conclusion  $U_i > U_j$ , again under the assumption of normality (see (Chien95) for further details).

$$EL\{U_i > U_j\} = \frac{\hat{s}_{i-j} e^{-0.5n(\frac{\hat{U}_{i-j}}{\hat{s}_{i-j}})^2}}{\sqrt{2\pi n}} + \frac{\hat{U}_{i-j}}{\sqrt{2\pi}} \int_{-\frac{\hat{U}_{i-j}}{\hat{s}_{i-j}\sqrt{n}}}^{\infty} e^{-0.5z^2} dz \quad (5)$$

In the next two subsections, we describe two interpretations for estimating the likelihood that an overall ranking satisfies the PAC or EL requirements by estimating and combining pairwise PAC errors or EL estimates. Each of these interpretations lends itself directly to an algorithmic implementation as described below.

### Ranking as Recursive Selection

One obvious way to determine a ranking  $H_1, \dots, H_k$  is to view ranking as recursive selection from the set of remaining candidate hypotheses. In this view, the overall ranking error, as specified by the desired confidence in PAC algorithms and the loss threshold in EL algorithms, is first distributed among  $k-1$  selection errors which are then further subdivided into pairwise comparison errors. Data is then sampled until the estimates of the pairwise comparison error (as dictated by equation 4 or 5) satisfy the bounds set by the algorithm.

Thus, another degree of freedom in the design of recursive ranking algorithms is the method by which the overall ranking error is ultimately distributed among individual pairwise comparisons between hypotheses. Two factors influence the way in which we compute error distribution. First, our model of error combination determines how the error allocated for individual comparisons or selections combines into overall ranking error and thus how many candidates are available as targets for the distribution. Using Bonferroni's inequality, one combine errors additively, but a more conservative approach might be to assert that because the predicted "best" hypothesis may change during sampling in the worst case the conclusion might depend on all possible pairwise comparisons and thus the error should be distributed among all  $\binom{n}{2}$  pairs of hypotheses<sup>4</sup>.

Second, our policy with respect to allocation of error among the candidate comparisons or selections determines how samples will be distributed. For example, in some contexts, the consequences of early selections far outweigh those of later selections. For these scenarios, we have implemented ranking algorithms that divide overall ranking error unequally in favor of earlier selections<sup>5</sup>. Also, it is possible to divide selection error into pairwise error unequally based on estimates of hypothesis parameters in order to reduce sampling cost (for example, (Gratch94) allocates error rationally).

Within the scope of this paper, we only consider algorithms that: (1) combine pairwise error into selection error additively, (2) combine selection error into overall ranking error additively and (3) allocate error equally at each level.

<sup>4</sup>For a discussion of this issue, see pp. 18-20 of (Gratch93)

<sup>5</sup>Space constraints preclude their description here.

One disadvantage of recursive selection is that once a hypothesis has been selected, it is removed from the pool of candidate hypotheses. This causes problems in rare instances when, while sampling to increase the confidence of some later selection, the estimate for a hypothesis' mean changes enough that some previously selected hypothesis no longer dominates it. In this case, the algorithm is restarted taking into account the data sampled so far.

These assumptions result in the following formulations (where  $\delta(U_1 \triangleright_{\epsilon} \{U_2, \dots, U_k\})$  is used to denote the error due to the action of selecting hypothesis 1 under Equation 1 from the set  $\{H_1, \dots, H_k\}$  and  $\delta(U_1 \triangleright_{\epsilon} \{U_2, \dots, U_k\})$  denotes the error due to selection loss in situations where Equation 2 applies):

$$\delta_{rec}(U_1 > U_2 > \dots > U_k) = \delta_{rec}(U_2 > U_3 > \dots > U_k) + \delta(U_1 \triangleright_{\epsilon} \{U_2, \dots, U_k\}) \quad (6)$$

where  $\delta_{rec}(U_k) = 0$  (the base case for the recursion) and the selection error is as defined in (Chien95):

$$\delta(U_1 \triangleright_{\epsilon} \{U_2, \dots, U_k\}) = \sum_{i=2}^k \delta_{1,i} \quad (7)$$

using Equation 4 to compute pairwise confidence.

Algorithmically, we implement this by:

1. sampling a default number of times to seed the estimates for each hypothesis mean and variance,
2. allocating the error to selection and pairwise comparisons as indicated above,
3. sampling until the desired confidences for successive selections is met, and
4. restarting the algorithm if any of the hypotheses means changed significantly enough to change the overall ranking.

An analogous recursive selection algorithm based on expected loss is defined as follows.

$$EL_{rec}(U_1 > U_2 > \dots > U_k) = EL_{rec}(U_2 > U_3 > \dots > U_k) + EL(U_1 \triangleright \{U_2, \dots, U_k\}) \quad (8)$$

where  $EL_{rec}(U_k) = 0$  and the selection EL is as defined in (Chien95):

$$EL(U_1 \triangleright \{U_2, \dots, U_k\}) = \sum_{i=2}^k EL(U_1, U_i) \quad (9)$$

### Ranking by Adjacency Comparison

Another interpretation of ranking confidence (or loss) is that only adjacent elements in the ranking need be compared. In this case, the overall ranking error is divided directly into  $k-1$  pairwise comparison errors. This leads to the following confidence equation for the PAC criteria:

$$\delta_{adj}(U_1 > U_2 > \dots > U_k) = \sum_{i=1}^{k-1} \delta_{i,i+1} \quad (10)$$

And the following equation for the EL criteria.

$$EL_{adj}(U_1 > U_2 > \dots > U_k) = \sum_{i=1}^{k-1} EL(U_i, U_{i+1} | \mathbf{1}) \quad (11)$$

Because ranking by comparison of adjacent hypotheses does not establish the dominance between non-adjacent hypotheses (where the hypotheses are ordered by observed mean utility), it has the advantage of requiring fewer comparisons than recursive selection (and thus may require fewer samples than recursive selection). However, for the same reason, adjacency algorithms may be less likely to correctly bound probability of correct selection (or average loss) than the recursive selection algorithms. In the case of the PAC algorithms, this is because  $\epsilon$ -dominance is not necessarily transitive. In the case of the EL algorithms, it is because expected loss is not additive when considering two hypothesis relations sharing a common hypothesis. For instance, the size of the blocked differential distribution may be different for each of the pairs of hypotheses being compared.

### Other Relevant Approaches

Most standard statistical ranking/selection approaches make strong assumptions about the form of the problem (e.g., the variances associated with underlying utility distribution of the hypotheses might be assumed known and equal). Among these, Turnbull and Weiss (Turnbull84) is most comparable to our PAC-based approach<sup>6</sup>. Turnbull and Weiss treat hypotheses as normal random variables with unknown mean and unknown and unequal variance. However, they make the additional stipulation that hypotheses are independent. So, while it is still reasonable to use this approach when the candidate hypotheses are not independent, excessive statistical error or unnecessarily large training set sizes may result. In the case that the hypotheses are truly independent, Turnbull and Weiss' technique should be able to exploit this knowledge and outperform our methods which do not adopt this assumption.

### Empirical Performance Evaluation

We now turn to empirical evaluation of the hypothesis ranking techniques on both synthetic and real-world datasets. This evaluation serves three purposes. First, it demonstrates that the techniques perform as predicted (in terms of bounding the probability of incorrect selection or expected loss). Second, it validates the performance of the techniques as compared to standard algorithms from the statistical literature. Third, the evaluation demonstrates the robustness of the new approaches to real-world hypothesis ranking problems.

---

<sup>6</sup>“PAC-based approaches have been investigated extensively in the statistical ranking and selection literature under the topic of *confidence interval based* algorithms (see El asceb85) for a review of the recent literature).

### Methodology

An experimental trial consists of solving a hypothesis ranking problem with a given technique and a given set of problem and control parameters. We measure performance by (1) how well the algorithms satisfy their respective criteria; and (2) the number of samples taken. Since the performance of these statistical algorithms on any single trial provides little information about its overall behavior, each trial is repeated multiple times and the results are averaged across trials. Synthetic experimental trials were repeated 500 times, while trials on the real-world data were repeated 100 times. Because the PAC and expected loss criteria are not directly comparable, the approaches are analyzed separately.

### Evaluation on Synthetic Datasets

Evaluation on synthetic data is used to show that: (1) the techniques correctly bound probability of incorrect ranking and expected loss as predicted when the underlying assumptions are valid even when the underlying utility distributions are inherently hard to rank, and (2) that the PAC techniques compare favorably to the algorithm of Turnbull and Weiss in a wide variety of circumstances.

For the synthetic datasets, the utility distributions of the hypotheses were modeled as random variables defined on some underlying parameterized distribution. Thus, characterizing a ranking problem consists of choosing some number of hypotheses to rank and then assigning a distribution and values for its parameters to the random variables representing the utility distributions for these hypotheses. In our case, we model the utilities as independent normal random variables with some mean and standard deviation. Thus, if we let  $k$  be the number of hypotheses, then each hypothesis ranking problem is described by the  $2k$  parameters specifying the expected utility and utility standard deviation for each hypothesis. In general, while several more parameters may be required to characterize a ranking problem fully<sup>7</sup>, the number of hypotheses and the choices for the parameters of the utility distributions underlying these hypotheses characterize the overall difficulty of the ranking problem.

The statistical ranking and selection community uses a standard family of selection problems with known difficulty to analyze the performance of hypothesis selection strategies. The method, called the least favorable configuration (LFC) of the population means is that assignment of the parameters to distributions which is most likely to cause a technique to choose a wrong hypothesis and thus provides the most severe test of the technique's abilities. Under this configuration, all utilities are independent normally distributed variables of equal variance.  $k-1$  of the hypotheses have utilities with equal expectation,  $\mu$ , and the remaining hypothesis has expected utility  $\mu + \epsilon$ .

---

<sup>7</sup>For instance, when samples are allocated rationally in (Chien95), it becomes necessary to assign parameters to a cost distribution as well, or if only a few of the candidate hypotheses were to be ranked, the number of hypotheses to rank would be another problem parameter.

Table 1: Estimated expected total number of observations by PAC algorithms in the stepped means configuration. Achieved probability of correct ranking is shown in parenthesis.

k	$\gamma^*$	$\frac{\sigma}{\epsilon}$	TURNBULL	$PAC_{rec}$	$PAC_{adj}$
3	0.75	2	62 (0.88)	55 (0.95)	38 (0.78)
3	0.75	3	117 (0.89)	101 (0.86)	49 (0.80)
3	0.90	2	97 (0.96)	86 (0.94)	58 (0.92)
3	0.90	3	183 (0.97)	152 (0.96)	96 (0.89)
3	0.95	2	130 (0.97)	122 (0.97)	89 (0.97)
3	0.95	3	231 (0.96)	204 (0.95)	146 (0.94)
5	0.75	2	177 (0.87)	165 (0.95)	105 (0.87)
5	0.75	3	321 (0.95)	314 (0.93)	161 (0.75)
5	0.90	2	245 (0.98)	245 (0.97)	163 (0.91)
5	0.90	3	445 (0.98)	409 (0.91)	290 (0.92)
5	0.95	2	299 (0.98)	294 (0.98)	216 (1.00)
5	0.95	3	541 (0.98)	538 (0.98)	377 (0.92)
10	0.75	2	558 (0.92)	624 (0.91)	345 (0.85)
10	0.75	3	1,015 (0.94)	1,042 (0.95)	635 (0.83)
10	0.90	2	700 (0.97)	742 (0.96)	523 (0.91)
10	0.90	3	1,254 (0.97)	1,359 (0.97)	883 (0.90)
10	0.95	2	821 (1.00)	877 (0.97)	661 (0.94)
10	0.95	3	1,462 (0.99)	1,569 (0.98)	1,164 (0.93)

Because we are interested in hypothesis ranking problems rather than selection problems, we use a generalization of the LFC that we call stepped means. In this configuration, one of the hypotheses is assigned expected utility  $\mu$  and successive hypotheses are assigned expected utility  $\mu - i\epsilon$  for  $i$  from  $1, \dots, k - 1$ .

In general, problems based on the least favorable configuration become more difficult (i.e., require more samples) when the number of hypotheses  $k$  increases, the common utility variance  $\sigma^2$  increases, or the difference in the means of the utility distributions decreases. In the standard methodology, a technique is evaluated by its ability to achieve a confidence of correct selection  $\gamma^*$  using several settings for  $k$  and  $\frac{\sigma}{\epsilon}$ . This last ratio combines  $\sigma$  and  $\epsilon$  into a single quantity which, as it increases, makes the problem more difficult. This methodology extends to stepped means directly.

The hypothesis ranking strategies themselves have *algorithm control parameters* that govern how they attack a problem. The PAC techniques have time control parameters: an initial sample size  $n_0$ , a desired confidence of correct ranking  $\gamma^*$  and an indifference setting  $\epsilon^8$ . The expected loss techniques have two control parameters: an initial sample size  $n_0$  and a loss threshold  $11^*$ .

For our experiments,  $n_0 = 7$ ,  $\mu = 50$ ,  $\sigma = 64$ , and all other parameters are varied as indicated.

The observed number of samples required and achieved accuracy of the PAC techniques on the stepped means configuration are shown in Table 1. The results indicate that all systems are roughly comparable in the number of examples required to choose a hypothesis. As expected, the number of examples increases with  $k$ ,  $\gamma^*$ , and  $\frac{\sigma}{\epsilon}$ . The  $PAC_{adj}$

<sup>8</sup>Note that in our formulation of the stepped means test for the PAC approaches,  $\epsilon$  is both the difference in the expected mean of successive hypotheses and the indifference interval of the algorithm. Thus,  $\epsilon$  plays the roles of both problem parameter and control parameter here.

Table 2: Estimated expected total number of observations of EL algorithms in stepped means configuration. Observed average loss of produced rankings.

k	$\epsilon$	$H^*$	$EL_{rec}$		$EL_{adj}$	
			Samples	Loss	Samples	Loss
3	2	1.0	96	0.6	43	1.2
3	2	0.75	102	0.5	56	1.0
3	2	0.5	139	0.2	73	0.6
3	2	0.25	235	0.1	139	0.4
5	2	1.0	320	0.7	140	1.3
5	2	0.75	343	0.4	169	1.2
5	2	0.5	464	0.4	247	0.7
5	2	0.25	575	0.2	350	0.5
10	2	1.0	1,136	0.5	572	1.4
10	2	0.75	1,325	0.5	668	1.1
10	2	0.5	1,533	0.3	872	0.7
10	2	0.25	1,856	0.1	1,153	0.4

algorithm required the least number of samples but was inconsistent in meeting the desired accuracy bound. It is interesting that the Turnbull and Weiss method did not significantly outperform the PAC techniques despite the fact that the algorithm assumes that the hypotheses are independent (as is the case in the stepped means configuration), while the PAC approaches do not make this assumption.

In the expected loss experiments, we ran the expected loss hypothesis ranking algorithms on the same stepped means configurations described above with a range of expected loss bounds. Table 2 shows the results of this experiment, displaying the number of samples required to produce a ranking and the average observed loss for each configuration. These results show that the  $EL_{rec}$  algorithm correctly bounded the loss and that the  $EL_{adj}$  algorithm required less samples than the  $EL_{rec}$  algorithm, but did not correctly bound the expected loss.

## Evaluation on Real Datasets

The test of real-world applicability is based on data drawn from an actual NASA spacecraft design optimization application. This data provides a strong test of the applicability of the techniques in that all of the statistical techniques make some form of normality assumption - yet the data in this application is highly non-normal.

The goal of the spacecraft design problem is to determine a good set of physical dimensions for a penetrator - a small, robust probe designed to impact a surface at extremely high velocity with the goal of performing deep soil sample analysis. Specifically, we use design data from the New Millennium Deep Space Two mission penetrator design.

For our casting of the design problem, we hold the shape of the penetrator constant and rank designs based on the variables of penetrator diameter and length. For a specific design a sample is taken by choosing impact orientation, impact velocity, and soil density from a parameterized multivariate distribution and then calling a complex physical simulation to determine if and to what depth the penetrator bored into the Martian surface. The goal of the penetrator design problem is to determine the dimensions that maximize the probability of penetration, and in cases of

Table 3: Estimated expected total number of observations to rank DS-2 spacecraft designs. Achieved probability of correct ranking is shown in parenthesis.

k	$\gamma^*$	$\frac{\sigma}{\epsilon}$	TURNBULL	$PAC_{ecc}$	$PAC_{adj}$
10	0.75	2	534 (0.96)	144 (1.00)	92 (0.98)
10	0.90	2	667 (0.98)	160 (1.00)	98 (1.00)
10	0.95	2	793 (0.99)	177 (1.00)	103 (0.99)

Table 4: Estimated expected total number of observations and expected loss of an incorrect ranking of DS-2 penetrator designs.

Parameters		$EL_{ecc}$		$EL_{adj}$	
k	$H^*$	Samples	Loss	Samples	Loss
10	0.10	152	0.005	77	0.014
10	0.05	200	0.003	90	0.006
10	0.02	378	0.003	139	0.003

penetration, maximize penetration depth.

Tables 3 and 4 show the results of applying the PAC-based, Turnbull, and expected loss algorithms to a ranking problem in which the system is requested to rank 10 penetrator designs<sup>9</sup>. In this problem the utility function is the depth of penetration of the penetrator, with those cases in which the penetrator dots not penetrate being assigned zero utility. As shown in Table 3, both PAC algorithms significantly outperformed the Turnbull algorithm, which is to be expected because the hypotheses are somewhat correlated (via impact orientations and soil densities). Table 4 shows that the  $EL_{ecc}$  expected loss algorithm effectively bounded actual loss but the  $EL_{adj}$  algorithm was inconsistent,

## Discussion and Conclusions

There are a number of areas of related work. First, there has been considerable analysis of hypothesis selection problems. Selection problems have been formalized using a Bayesian framework (Moore94; Rivest88) that does not require an initial sample, but uses a rigorous encoding of prior knowledge. Howard (Howard70) also details a Bayesian framework for analyzing learning cost for selection problems. If one uses a hypothesis selection framework for ranking, allocation of pairwise errors can be performed rationally (Gratch94). Reinforcement learning work (Kaelbling93) with immediate feedback can also be viewed as a hypothesis selection problem.

In summary, this paper has described the hypothesis ranking problem, an extension to the hypothesis selection problem. We defined the application of two decision criteria, *probably approximately correct* and *expected loss*, to this problem. We then defined two families of algorithms, recursive selection and adjacency, for solution of hypothesis ranking problems. Finally, we demonstrated the effectiveness of these algorithms on both synthetic and real-world

<sup>9</sup>True expected utility values are computed by perturbing a deep sample of 20,000 samples. These expected utilities can then be used to compute  $PAC_{\epsilon}$  - validity of rankings and actual loss,

datasets, documenting improved performance over existing statistical approaches.

## References

- R.E. Bechhofer, "A Single-sample Multiple Decision Procedure for Ranking Means of Normal Populations with Known Variances," *Annals of Math. Statistics* (25) 1, 1954 pp. 16-39.
- S. A. Chien, J. M. Gratch and M.C. Burl, "On the Efficient Allocation of Resources for Hypothesis Evaluation: A Statistical Approach," *IEEE Trans. Pattern Analysis and Machine Intelligence* 17 (7), July 1995, pp. 652-665.
- D. Goldberg, *Genetic Algorithms in Search, Optimization and Machine Learning*, Addison Wesley, 1989.
- Z. Govindarajulu, "The Sequential Statistical Analysis," American Sciences Press, Columbus, OH, 1981.
- J. Gratch and G. DeJong, "COMPOSER: A Probabilistic Solution to the Utility Problem in Speed-up Learning," Proc. AAAI92, San Jose, CA, July 1992, pp. 235-240.
- J. Gratch, "COMPOSER: A Decision-theoretic Approach to Adaptive Problem Solving," Tech. Rep. UIUCDCS-R-93-1 806, Dept. Comp. Sci., Univ. Illinois, May 1993.
- J. Gratch, S. Chien, and G. DeJong, "Improving Learning Performance Through Rational Resource Allocation," Proc. AA AI94, Seattle, WA, August 1994, pp. 576-582.
- R. Greiner and I. Jurisica, "A Statistical Approach to Solving the EBL Utility Problem," Proc. AAAI92, San Jose, CA, July 1992, pp. 241-245.
- R. M. Haseeb, *Modern Statistical Selection*, Columbus, OH: Am. Sciences Press, 1985.
- R. V. Hogg and A. T. Craig, *Introduction to Mathematical Statistics*, Macmillan inc., London, 1978.
- R. A. Howard, *Decision Analysis: Perspectives on Inference, Decision, and Experimentation*, Proceedings of the IEEE 58, 5 (1970), pp. 823-834.
- L. I. Kaelbling, *Learning in Embedded Systems*, MIT Press, Cambridge, MA, 1993.
- S. Minton, *Learning Search Control Knowledge: An Explanation-Based Approach*, Kluwer Academic Publishers, Norwell, MA, 1988.
- A. W. Moore and M. S. Lee, "Efficient Algorithms for Minimizing Cross Validation Error," Proc. MI.94, New Brunswick, MA, July 1994.
- R. Musick, J. Catlett and S. Russell, "Decision Theoretic Subsampling for Induction on Large Databases," Proc. MI.93, Amherst, MA, June 1993, pp. 212-219.
- R. I. Rivest and R. Sloan, "A New Model for Inductive Inference," Proc. 2nd Conference on Theoretical Aspects of Reasoning about Knowledge, 1988.
- S. Russell and E. Wefald, *Do the Right Thing: Studies in Limited Rationality*, MIT Press, Cambridge, MA.
- P. Tadepalli, "A theory of unsupervised speedup learning," Proc. AAAI92, S. Jose, CA, 1992, pp. 229-234.
- Turnbull and Weiss, "A class of sequential procedures for k-sample problems concerning normal means with unknown unequal variances," in *Design of Experiments: ranking and selection*, T. J. Santner and A. C. Tamhane (eds.), Marcel Dekker, 1984.
- I. G. Valiant, "A Theory of the Learnable," *Communications of the ACM* 27, (1984), pp. 1134-1142.