

3/27/01

IDSA:

Title, Abstract, of first talk April 4, (Wed) 11 AM

System  
+ Progressive Machine Learning.  
High Level

130 - 230P  
Gates

## The Scientist's Assistant: a System for Progressive Machine Learning.

The "Scientist's Assistant" has been a long term subject in our research. ~~and a part in~~ <sup>Recent progress in</sup> machine learning. We will describe "The Scientist's Assistant" — a program for machine learning that we have been working on for many years. Very little of the program ~~has~~ <sup>it has been</sup> considered in software. Our discussions will be largely at ~~theoretical~~ <sup>learning to solve</sup> operation of Data System ~~WISLIP~~ — how it ~~solves~~ <sup>problems, and some of its</sup> apparent difficulties, and our motivation for a system of this kind.

The system begins with some initially simple ~~techniques for solving problems~~ <sup>we want to</sup> ~~problems, which it solves easily.~~ Having solved these easy problems, we move to somewhat more difficult problems — which it solves using some of the techniques required in solving ~~initial~~ <sup>the</sup> set. We continue to find that the system problem sets of progressive difficulty, until it ~~reaches~~ <sup>reaches</sup> and at this time it's able ~~has arrived~~ <sup>to</sup> the desired level of intelligence ~~by using~~ <sup>to solve them</sup> ~~techniques~~ <sup>it acquires</sup> the desired level of intelligence by using ~~the~~ techniques acquired in solving previous problems. We continue this ~~until~~ <sup>training</sup> until it acquires the desired level of intelligence.

It sov

We begin training the system by giving it an initial set of simple problems. Using a kind of Levin-Universal Search, it is able to solve those ~~easy~~ <sup>nearly simple</sup> problems, guided by an probability distribution over possible solutions. Having solved these problems, the guiding probability distribution is modified in view of the problems solved. We then give it more difficult problems than it is able to solve ~~using~~ <sup>using</sup> ~~the~~ modified probability distribution. This process of problem ~~the~~ probability distribution is then again modified to accommodate the newly solved problems. This iteration continues ~~until~~ <sup>the</sup> attainment of ~~gradual~~ <sup>gradual</sup> modification of the probability distribution due to ~~new~~ <sup>new</sup> problems of progressively greater difficulty. ~~difficulty~~ <sup>difficulty</sup> is alternated with modification of the probability distribution, until ~~can be acquired in this manner.~~ <sup>can be acquired in this manner.</sup>

At very high levels of problem solving experience ~~the~~ <sup>the</sup> system ~~is~~ <sup>is</sup> acquired.

Finally we will discuss some of the details of program construction, the kinds of ~~problems~~ <sup>as well as the</sup> problems solved, the method of modification of the probability distribution ~~the~~ <sup>capabilities and</sup> limitations of the system.

autonella@idmex.ch.

D

## Abstract

We will describe a System for Machine Learning based on Lenn's Universal Search Algorithm to solve problems of progressively increasing difficulty.

First a discussion of the search algorithm — its necessary limitations and how they can be overcome — partly by modifying the algorithm itself and partly by progressive modification of the probability distribution that guides it.

That uses

(Lenn's)

~~discusses~~ Title: A System for Machine Learning based on Lenn's Search Algorithm.

We will discuss Lenn's search algorithm in various forms with WM and MAP for two kinds of problems: ~~concrete~~ <sup>Universal</sup>.

The main necessary limitation is:

One limitation that normally limits use to simple problems <sup>we</sup> will show how to overcome this's limitation by using a probability distribution suitable conditions for probability distribution to guide the search.

by updating this conditional probability distribution, the system <sup>is able to integrate knowledge</sup> about base of solutions to problems it has solved. The updating system is based on ~~experience with problem solutions~~. <sup>algorithm</sup> global credit assignment function

We will discuss the updating algorithm and the Global credit assignment function ~~and the~~

function ~~used to update~~ to the update of that governs the update algorithm.

In previous papers, we claimed that ~~the universal search system~~ was very plausible that within a factor of 4 of optimum, it was restricted the behavior of competing systems in certain ways. If we relax those restrictions and modify our system accordingly, the factor of 4 becomes closer to "a factor of one" but the probability arguments are more uncertain.

Remember ~~universal search~~

If we relax those ~~restrictions~~ on competing systems and suitably augment our own system, it appears that we may be very close to optimum, but the "plausibility arguments" are not as strong as they were formerly.

Within a certain restricted class of problem solving systems, the system proposed, we have probably to show it was optimum within a factor of 4. Using a more efficient <sup>for Lenn's</sup> algorithm, we can obtain a factor of 2. ~~It is~~ It was too compete our system with unregistered sets of problem solving algorithms — if we augment our system about, it appears to be within a factor of 2 of unregistered learning algorithms solving the same kinds of problems — but the arguments justifying this claim are much less <sup>heuristic</sup> better.

3/22/01

3

Fdsz

so  $\lambda^{226}$  now follows the criterion in general as well as more recent claims of ~~as~~ near optimality under less restricted conditions.

3/27/01

IDSA:

Title, Abstract, of first talk April 4, (Wed) 11 AM

System  
+ Progressive Machine Learning.  
High Level

130 - 230P  
Gates

## The Scientist's Assistant: a System for Progressive Machine Learning.

The "Scientist's Assistant" has been a long term subject in our research. ~~and a part in~~ <sup>Recent progress in</sup> machine learning. We will describe "The Scientist's Assistant" — a program for machine learning that we have been working on for many years. Very little of the program ~~has~~ <sup>it has been</sup> considered in software. Our discussions will be largely at ~~theoretical~~ <sup>learning to solve</sup> operation of Data System ~~WISLIP~~ — how it ~~solves~~ <sup>problems, and some of its</sup> apparent difficulties, and our motivation for a system of this kind.

The system begins with some initially simple ~~techniques for solving problems~~ <sup>we want to</sup> ~~problems, which it solves easily.~~ Having solved these easy problems, we move to somewhat more difficult problems — which it solves using some of the techniques required in solving ~~initial~~ <sup>the</sup> set. We continue to find that the system problem sets of progressive difficulty, until it ~~reaches~~ <sup>reaches</sup> and at this time it's able ~~has arrived~~ <sup>to</sup> the desired level of intelligence ~~by using~~ <sup>to solve them</sup> ~~techniques~~ <sup>it acquires</sup> the desired level of intelligence by using ~~the~~ techniques acquired in solving previous problems. We continue this ~~until~~ <sup>training</sup> until it acquires the desired level of intelligence.

It sov

We begin training the system by giving it an initial set of simple problems. Using a kind of Levin-Universal Search, it is able to solve those ~~easy~~ <sup>nearly simple</sup> problems, guided by an probability distribution over possible solutions. Having solved these problems, the guiding probability distribution is modified in view of the problems solved. We then give it more difficult problems than it is able to solve ~~using~~ <sup>using</sup> ~~the~~ modified probability distribution. This process of problem ~~the~~ probability distribution is then again modified to accommodate the newly solved problems. This iteration continues ~~until~~ <sup>the</sup> attainment of ~~gradual~~ <sup>gradual</sup> modification of the probability distribution due to ~~new~~ <sup>new</sup> problems of progressively greater difficulty. ~~difficulty~~ <sup>difficulty</sup> is alternated with modification of the probability distribution, until ~~can be acquired in this manner.~~ <sup>can be acquired in this manner.</sup>

At very high levels of problem solving experience ~~the~~ <sup>the</sup> system ~~is~~ <sup>is</sup> acquired.

Finally we will discuss some of the details of program construction, the kinds of ~~problems~~ <sup>as well as the</sup> problems solved, the method of modification of the probability distribution ~~the~~ <sup>capabilities and</sup> limitations of the system.

autonella@idmex.ch.

D

## Abstract

We will describe a System for Machine Learning based on Lenn's Universal Search Algorithm to solve problems of progressively increasing difficulty.

First a discussion of the search algorithm — its necessary limitations and how they can be overcome — partly by modifying the algorithm itself and partly by progressive modification of the probability distribution that guides it.

That uses

(Lenn's)

~~discusses~~ Title: A System for Machine Learning based on Lenn's Search Algorithm.

We will discuss Lenn's search algorithm in various forms with WM and MAP for two kinds of problems: ~~concrete~~ <sup>Universal</sup>.

The main necessary limitation is:

One limitation that normally limits use to simple problems <sup>we</sup> will show how to overcome this's limitation by using a probability distribution suitable conditions!

Probability distribution to guide the search.

by / updating This conditional probability distribution, the system <sup>is able to integrate knowledge</sup> about base of solutions to problems it has solved. The updating system is based on ~~experience with problem solutions~~. <sup>algorithm</sup> global credit assignment function

We will discuss the updating algorithm and the Global Credit assignment function ~~and the~~

function ~~used~~ <sup>describes</sup> ~~in the update~~ to the update of that governs the update algorithm.

In previous papers, we claimed that ~~the universal search system~~ was very plausible that was very close to the optimum, if we restricted the behavior of competing systems within a factor of 4 of optimum, it was restricted the behavior of competing systems in certain ways. If we relax those restrictions and modify our system accordingly, the factor of 4 becomes closer to "a factor of one" but the probability arguments are more uncertain.

Remember ~~universal search~~

If we relax those ~~restrictions~~ restrictions on competing systems and suitably augment our own system, it appears that we may be very close to optimum, but the "probability arguments" are not as strong as they were formerly.

~~demonstrates~~ Within a certain restricted class of problem solving systems, the system proposed, we have probably to show it was optimum within a factor of 4. Using a more efficient <sup>for Lenn's</sup> algorithm, we can obtain a factor of 2. ~~It is~~ It was too compete our system with an unregistered set of problem solving algorithms —

If we augment our system about, it appears to be within a factor of 2 of unregistered learning algorithms solving the same kinds of problems — but the arguments justifying this claim are much less <sup>heuristic</sup> better.

3/22/01

3

Fdsz

so  $\lambda^{226}$  now follows the criterion in general as well as more recent claims of ~~as~~ near optimality under less restricted conditions.

Talk  $\Leftarrow \sum_{-\infty}^{\infty} \sum_{-\infty}^{\infty}$

1.00: 94.40

Dealt w/ it more fully, after we ~~will~~ discuss ~~the estimation of probabilities~~

The particular method of probability estimation we use is called Algorithmic Probability.

Binary

Suppose we ~~are~~ doing sequential prediction! We have a ~~long~~ / Sequential

~~string~~,  $S$  and we want to know the relative probability of the next symbol

being a value from  $\Omega$  — ~~S is a sequence of symbols in an alphabet of which~~

~~and to our memory~~

Let us assume first we know an a priori distribution on all possible strings ~~of length~~,

~~R, binary alphabet~~.  $P(R)$  assigns a probability to each symbol string.

Then the relative probabilities of two configurations  $\Omega \rightarrow \Omega$  will be  $\frac{P(S_1)}{P(S_0)}$ .

This would solve the sequential prediction problem, if we knew  $P(R)$ .

(How to find  $P(R)$ )

For a heuristic understanding of what follows, consider Ockham's ~~Idea~~ idea:

That simpler hypotheses are more likely than complex ones. We will quantify this idea ~~and~~ describe some results on ~~the~~

~~the accuracy of the resultant system.~~

To quantify Ockham: Suppose  $B$  is a finite binary sequence.

then  $D_B$  is "description" of  $B$  if  $M(D_B) = B$ .

Here  $M$  is an algorithm or ~~program~~ computing which no part is able to take the

binary string  $D_B$  as input, for which it gives  $B$  as output.  $D_B$  can be

regarded as ~~the~~ <sup>Description</sup> of  $B$  in terms of Machin function  $M_r$ .

I use the subscript  $r$ , because  $M_r$  is to be regarded as a reference machine.

(Usually, by changing reference machines, ~~the~~ the correct descriptions for  $B$  will change.)

~~The~~ Shortest description will be done with fewest bits —

containing latest information. We could approximate the probability of  $B$  by  $2^{-D_B}$ .

For more accurate estimation, we use  $\approx 2^{-D_{B,r}}$ . <sup>(Also known as 'incisional' theory)</sup> Here we are making the

total probability due to all possible descriptions of  $B$ . We will call  $\approx 2^{-D_{B,r}}$

(The Algorithmic probability assigned to  $B$  by  $M_r$ ) =  $p^{M_r}(B)$

X

OP: 95, fo; There are 2 kinds of losses here: First, what reference machine does for our probabilistic calcn? There is a group of machines called "Universal Machines", that are particularly good for describing things. They have the property that if we compare 2 of them, there will always be a constant factor that tells how much they differ from one another.

Two machines  $M_1, M_2$  then  $P^{M_1}(B)$  and  $P^{M_2}(B)$  are always w.p.  $\geq$  constant factor  $C_{1,2}$  of one another. This constant factor depends on  $M_1$  and  $M_2$  but is independent of  $B$ . The constant factor can be quite large — so it usually makes a lot of difference as to what reference machine we use.  $\square$

In our system for Machine Learning ~~we~~ will be periodically changing the reference machine as a way to accommodate new things that the system ~~keeps~~ hasn't learned.

Second: The method of computing probabilities of strings may seem a bit arbitrary, — and after several years after I thought of this method, I wasn't sure it would work — and

surprised myself by working. However, I ~~initially did work out~~  $E \sum_{i=1}^m (p_i - p'_i)^2 \leq \frac{1}{2} m^2 \ln 2$ : ~~but~~ ~~it's~~ ~~not~~ ~~the~~ ~~best~~ ~~estimate~~.

Proof that the system was ~~very~~ accurate in estimating probabilities from empirical data.

~~If we use universal reference machine~~

Third: It turns out that  $\sum 2^{-D_B}$  is not computable in a finite amount of time, ~~but~~  $\approx 97.0\%$ . We can't even do ~~shortest~~ ~~well~~ ~~no~~ found two shortest code for  $B$ ! That algorithmic probability is incomputable may seem like  $\rightarrow 97.00$

string representation appears — but it is not — this incomputability is an essential part of probability — of science itself!

The short codes that contribute most to the  $\sum 2^{-D_B}$  correspond to strings

repetitives in the data. Consider the string  $(01)^{2048}$ . We can easily program a machine

to write that sequence by simply telling it to write 01, 2048 times —

.. Intuitively

a OZ problem.

which takes 2 bits to say "01" and 10 bits to say 1002 + plus a few overheads —

$\Rightarrow$  base compression of 2048 bits in the original string,

~~but~~ ~~it's~~ ~~not~~ ~~the~~ ~~best~~ ~~estimate~~,

96.23 ! The reason is not something we can try all possible codes! —  
 Very long codes don't contribute much to the sum. The main problem is that  
 Programs contain short strings that one puts into the machine, and the machine  
 runs and runs, and after a long time we can't know what will output Board stop —  
 and there is no survey to tell. → **Back to 96.23**

**REMARKABLE** In general in science, when you have a batch of data  
 and you are looking for regularities in the data rather than to be used for  
 prediction — you can never be sure that if you spending so many minutes  
 hunting for regularities would not give you much better than the best  
 you have had yet.

The corresponding prob occurs with random sequences that seem  
 to be "RANDOM" — i.e. no discernible regularities. It is balanced by  
 many, ~~that~~ stock market prices are random — yet there is no finite amount  
 of investigation that could be convincing on this question. On the other  
 hand, once a regularity is found, it becomes  
 very unlikely that the sequence is random.

The problem of finding  $\sum_{i=1}^{m+1} 2^{-l_i}$  good approximations to  
 this done by  $\left[ \sum_{i=1}^{m+1} 2^{-l_i} \text{ summing over not } > 10^6 \text{ of codes.} \right] \leftarrow P'$

$P'$  is an approximation to  $P$ . To make codes one finds the closer  $P' \rightarrow P$ .  
 To maximize  $P'$  by finding as many codes as possible in the available  
 time, is a time limited optimization problem. It is equivalent to  
 having an minimum error ( $\|P - P'\|$ ) in ones estimates  $P'$ .

Induction inference The prediction problem and its solution as a time limited optimization  
 problem are very important in the Machine Learning System that I'm describing.  
 One reason is that the problem of updating the probability distribution that  
 guides search for solutions to problems — is an induction problem.

How is this so?

Consider all of the [problem description, solution program, time] triplets  
 that have occurred thus far.

[Problem Descr., Solution, Prgm, Time]

4/3/01  
IDSIA

.00 : 97.40 ! What we want is a bunch of relatively short codes to describe all of this data.

From a set of codes of this sort, we can then extrapolate the data to say now,

G-PD<sub>1</sub>(problem, solution) → Probability distribution on Time to Solution.

From G-PD<sub>1</sub>, it is possible by a process of induction, to create

G-PD<sub>2</sub>(problem, soln.) → Probability that this is Reacheable fastest solution to the problem.

This last distribution is the Griding distribution based search for solutions to new problems.

I will now discuss Levin's universal search procedure!

Suppose I'm in a Gambling house and there is a kind of lottery with a ~~test single~~ single ticket price. ~~I am the only customer~~  
~~A ticket costs \$1 and I am allowed to choose any book~~

Each lottery ticket has a certain probability of winning — which is printed on the ticket.

In the first ~~test~~ lottery, all tickets cost \$1. The best ticket to get would be

one with zero win probability. If it doesn't win a does next largest → ad so on.

In the next case, each ticket costs a different amount of money,

( $p_1, m_1$ ) are the probability, cost associated with the 1st ticket type.

The best choice to make is the ticket with maximum  $\frac{p_1}{m_1}$ . You get the maximum probability of winning per dollar spent. If you continue to buy tickets in  $\frac{p_1}{m_1}$  order you are ~~continuing to play~~ ~~more~~ ~~worse~~ ~~expected~~ money spent before winning.

The Gambling house ~~can~~ can be described into a problem Solving Environment.

Instead of money for trial, we ~~will~~ spend TIME.

The best trial to choose one with maximum  $\frac{p_i}{t_i}$ .

Normally in trying to solve problems, one does not know  $p_i$ , but one doesn't know  $t_i$ . In this ~~random~~ case, every good strategy is Levin's time share strategy → (which is practically likely to be the best possible.)

The way it works! You work on all trials simultaneously, but you work harder on trials with larger  $p_i$  values. The rate at which you work on trials, is proportional to  $p_i$ . Suppose we use this strategy and after  $j$  trials the  $j$ th trial gives a solution after spending a total time  $t_j$  on it. Since other trials have amounts of time spent on them proportional to their  $p_i$ 's. The  $i$ th trial will have  $\frac{p_i}{p_j} \cdot t_j$  time spent on it.

0.0019840: The total time spent on trials will be  $\sum \frac{p_i}{p_j} t_j = \frac{T_1}{p_1} \leq p_2$   
 $\sum p_i$  will be around 1. So the total time needed to solve the problem  
 this way is  $\frac{T_1}{p_1}$ .

So this is Levin's time shared search and it finds its best as possible.

For PRO leads of problems in which it is used i.e. INVERSION problems

An application for PRO search

This technique can also be used for time limited optimization problems.

In this case our trials are not ~~done~~ programs that attempt

to solve inversion problems! They are ~~optimization techniques~~

each one takes as input, the problem description,  $f(x)$  the function

to be minimized, and the time available for solution.

As before we have a conditional probability distribution that looks at the problem

description and assigns probabilities to various optimization techniques to

be applied to the problem. However since we know the time limit for

each trial (say  $p_i = T$ ), the best strategy is to try the optimization

techniques in ~~ps~~ order.

The way we do the ~~search~~: starting with some time limit  $T$

for some small  $i$  we spend time  $T \cdot p_i$  in the ~~if~~ O.T.

and slowly increase  $T$ . We keep track of the O.T. first has best G value

of all the O.T.'s. When  $T \cdot p_i = T$  for any O.T., we stop ~~spend~~ spend time on

that O.T. When  $T \cdot p_i = T$  for the ~~last~~ O.T., we stop. That is our best value.

Total time  $\leq T = T \cdot p_i = T \leq p_2 \leq T : T = \frac{T}{p_2}$  so time to obtain

best soln. in time  $T$ , is  $\leq \frac{T}{p_2}$  so we have efficiency factor of  $\frac{1}{p_2}$  — as we did for ENO problems

fraction  
 = 71.13

61.55  
 All into  
 into P.D.

59.33  
 63.20  
 on OZ

.00 : : Neck narrow System : Factors of epidemiology in system.

① ON OPTIMALITY: L-shaped ~~group~~ LIR problems is E

- $\Rightarrow$  optimum if All info used in such as in P.D. ① direct  
② learning is recording & reflecting back.  
③ Human P.D.'s vs. Synthetic Data.

④ No learning allowed during ~~complex structures~~ complex structures

So we modify System so it uses ~~new~~<sup>same</sup> learning factor. Main problem is problem of proving GPD

The new system would seem to be much better than a system that allows no  
learning before trials — I have no idea on how closely could be to  
optimum

② On ■ inductions → place ■ resources & limited or problem! (85.00ff)

- While it may be true for easy (like P/B ratio system), it is not true in general.  
If you have a very large mass of data, say information about the stock market, and you want to do prediction, ~~in general~~ for next day's price, and you have one hour to make the prediction, you will not be able to process all of the data — i.e. try to find short codes for it. Typically you will consider only recent data that scales more relevant and trying to choose best short codes using "Shallow Model".  
~~the best data using Shallow Model~~ The trade-off between the desirability of

~~Use as much as possible, since the more complex, the better the quality.~~  
The ROM API can't be used, so I find many good codes  
in very long data sequences, in a short time.

~~There are two competing goals in induction =~~  
~~It is easier to find regularities in small data sets.~~

Using large data sets gives more accurate attack probabilities. — but takes more time.

③ How system can transcend LSRCH if necessary ! That cards can be static points down learning binary search.

(4) How to eat food and at it's a concessionary set.

Things left over in talk:

" needed for notes / lecture.

- 1. Details on how induction works for BAG, Operator, Sequential! What are parameters?  
What are Gray Codes etc., etc. ↗ GPD<sub>1</sub> is off by 1 per
- 2. Some ways in which induction is actually done — particularly for Bag & Operator prob.
- 3. Details of update process → for OZ for INV, for  
induction (Bag, Operator, Sequential).

- 4. How new (non-local) items are stored (mention Marcus/Jay's ideas on Proj), maybe in Appendix, (how systems transcend Larch)  
(Q - Are they really analogous? — do they really work?)
- 5. ↑ but I made clear how induction was a OZ problem

- 6. Discussion of optimality of Larch in what sense is it optimal?

• for INV prob

Q's Jungen was concerned  
w. execution in which  
GPD<sub>1</sub> got stuck was obtained  
method short code, but long  
execution time.  
Main Q is how long does GPD take  
to update w/ [many PCs]

- for OZ prob ← (Motivation (little work)) see 102.04

- 7. Discussion of optimality of System as a whole: The "Power of 2" question, i.e. "how many days will it take"

- 8. Clear statement of Delft: induction is a "pure OZ problem" (how soon of course enters (do I have a soln?)) see, e.g. 100.14 (Best way to use Gray code)

- 9. How to start w. Goodset of O.I.'s — Break into "factor" down into a stock larch, first extrapolatory to see: This can be done (call it partly) by humans.

- It can be easily  
It can do, but not quite so easily



— needs a lot of work.

- 10. perhaps discuss details of Larch for both INV & OZ problems (Along w/ optimality) ← see 100.14

- 11. That the incomputability of ALP is not a valid criticism of it.

Riess uses Computable version of ALP & if has serious difficulties,

- 12. Discuss "random" Larch. Have any papers been written using actual Larch?

Also, other kinds of random Larch. **N.B.** My Arg is about Random Larch only correct if there is 1 or no soln.

- 13. Q of Optimality if learning during Larch (modelling GPD during Larch) is allowed.

- 14. Math & Logical reasoning: Emptiness in prob Solving, yet not explicitly considered in Proj System

- 15. Training Sequences & how to Write them "CJS" (my favorite CJS — its very nice teaching).

- 16. Serious diff'ren! in converting from GPD<sub>1</sub> to GPD<sub>2</sub>: we write  $\sum_{i=1}^{2^n} p_i(G) T_i P_i(C) \rightarrow$  112.28  
Hm. Pro  $p_i(G)$  are usually highly correlated (say 0.9) so  $P_i(C)$  is useless!

- 17. OSL! how most induction schemes (≈ MSL) have to be modified to accommodate it.