

Part I

Background

Chapter 3

Learning

Compare the following two somewhat distinct definitions of learning:

- (From The Random House Dictionary, 1980:) 1. to acquire knowledge (of) or skill (in) by study, instruction or experience. 2. To become informed (of). 3. To memorize (something).
- (From An Introduction to Theories of Learning, B. Hergenhahn, M. Olson, 1997) Learning is a relatively permanent change in behavior or in behavioral potentiality that results from experience and [the change in behavior] cannot be attributed to temporary body states such as those induced by illness, fatigue or drugs.

There is something unsatisfying about both of these definitions. Both would state that simple memorization can be defined as learning. Consider the following examples:

Example 1: In calculus, we learn that the derivative of x^n is nx^{n-1} . It is very common for beginning students to make the following error:

$$\frac{d}{dx} 3^x = x3^{x-1}$$

Example 2: In a study of learning in chimpanzees, a game was played. A chimpanzee was presented with two bowls of chocolate candy, one bowl had much more candy than the other. The chimp pointed at the bowl with more candy, and it was promptly given to another chimp. A while later, the same chimp was presented with a similar set of bowls, and again he chose the one with more candy, and it was given to another chimp. The chimp became very frustrated with the game, voicing great displeasure at having the bowl with more candy removed; he never pointed at the bowl with less candy.

Example 3: You put together a lookup table for your computer. That is, you list a series of possible keyboard inputs and the corresponding response for the computer to make. For example, “When I press the letter L, you will list the contents of the current directory”.

In the first example, something was learned, just not the right thing- but the goal is well defined (computing the derivative). In this case, the student has memorized a pattern, but has not realized what the salient items in the pattern referred to (e.g., x^n refers to a variable being raised to a constant power, not simply that we have one thing being raised to a power). In this case, there was an error in distinguishing between inputs (x^3 and 3^x).

The second example illustrates the basic question we must ask in learning problems- Does a solution actually exist? Are we using an algorithm that will pick up that solution? We might hypothesize that the chimp cannot learn to overcome its natural instinct. In machine learning, we need to be sure that we are presenting an association that *can* be learned.

In the third example, has the computer “learned” something? We have changed its behavior, and this change was due to instruction by the programmer. If we told the system to wipe out these commands after 30 days, then the change in behavior would be only relatively permanent. Of course, the authors of the definition had human beings in mind so it may be unfair to apply this to an inanimate object. However, it gives us a starting position in what we should mean when we talk about “Machine Learning”, or “Artificial Intelligence”. Suppose I were to change the third example:

Example 3: You put together a lookup table for your computer. That is, you list a series of possible keyboard inputs and the corresponding response for the computer to make. For example, “When I press the letter L, list the contents of the current directory”. Sometime later, I press the lower case “l”, and the computer listed the contents of the directory, *even though it had not been explicitly told to do that*. Now has the computer learned something?

We’re in a position now to talk about the learning experience, as it might relate to a machine:

- The learning experience must be initialized with enough samples so that differences can be distinguished, and similarities established. This means that we must use enough samples so that we will adequately represent all the desirable nuances. For example, if I present to you the sequence:

$$1, 2, \dots$$

have I given you enough information to infer the remaining elements? Probably, if I meant that the sequence is $1, 2, 3, 4, 5, \dots$. Probably not, if I meant the sequence $1, 2, 4, 8, 16, \dots$, or possibly I meant the sequence $1, 2, 3, 5, 8, 13, \dots$? Or possibly $1, 2, -1, -2, 3, 4, -3, -4, \dots$? In this example, we mean that if we have a more exotic sequence, then we will need more samples. This translates to a rule of thumb: If your desired input-output associations are more complex, then you will need a lot of samples.

At this point, it would be useful to see how this translates into a mathematical condition. Here is one way to say that a function (what we have been calling a relationship) should not vary wildly-

Let f be a function over a domain D . The function f is said to be *Lipschitz* if there exists an $M > 0$ so that, for all $a, b \in D$,

$$\text{dist}(f(b), f(a)) \leq M \text{dist}(b, a)$$

Exercise: Let $f : \mathbb{R} \rightarrow \mathbb{R}$, so that the definition becomes:

$$|f(b) - f(a)| \leq M|b - a|$$

- Show that a Lipschitz function is continuous at any point in its domain. (Therefore, all Lipschitz functions are continuous, but not vice-versa)
- Show that if f is differentiable on its domain and the derivative is bounded, then it is Lipschitz.

So, it is in this sense that having a Lipschitz function means that there is a restriction in how fast it can grow (or vary).

- Input-Output associations should not be so rigorous as to extinguish any possibility of extrapolating to novel stimulus. By this we mean that there is usually a trade off between how accurate a machine can learn, versus how well it can extrapolate. If we force the machine to have zero error during training, it might do very poorly on extrapolation. Furthermore, if the machine always predicts the mean of the outputs, it might be doing pretty well in extrapolation, but is very far away from any particular input-output relationship. An analogy in learning might be this: If a student is punished severely for all errors, it might reduce the possibility of the student answering novel questions in the future; novel stimulus may elicit no response (or a random response, since all incorrect answers get punished equally). On the other hand, if a student is rewarded for any answer that even remotely resembles the correct one, the student may never get an answer that is correct¹.

This balance between accuracy and extrapolation goes by other names as well: In learning, we might assign it to memorization versus understanding. In statistics, it is known as the trade off between bias and variance.

In any case, we can say that learning is indexed by appropriate responses to novel input- By appropriate, we will mean that small changes in input characteristics should correspond to small changes in response. Mathematically, again we are talking about a bound on the rate of change.

There are other aspects to learning. For example, we might think of learning as the process of constructing an internalized representation of the outside world, and this construction allows us to draw inferences or expectations about the

¹We have used the terms punishment and reward in a colloquial context- a behavioral psychologist would object to our example, but this is only meant to be a rough analogy to illustrate a point.

environment. The representation, implicitly referred to in the first item above, begs the question: If input from the environment can be viewed as data streams coming from sensation, how is the information processed?

Patients that have undergone cochlear implants report that initially, sounds are not recognized as hearing- rather, the patients say that they feel strange sensations in their brain. In time, they learn to translate these into recognizable sounds. The point is that we do not hear with our ears; we do not see with our eyes- we hear and see with our brain, which receives the environment as it appears after it has been processed and filtered through the sensory organs. We again ask the question: how is the information processed? If the brain were to receive and try to process every piece of data coming into our senses, we might hypothesize that it would quickly be overwhelmed (think about it- every square millimeter of skin would be constantly sending information at the same time that our ears are picking up sounds, our eyes are picking up images, our nose is picking up smells, we might be thinking about what to have for dinner, etc.). In conclusion, it seems reasonable to assume that each sensation is doing some processing at the signal's point of origin.

When psychologists study learning in babies, they assume that stimulus that has not been internalized elicits attention. For example, studies have shown that even very small babies have developed fairly sophisticated expectations about the world. One study did the following:

The baby is placed in front of a stage. On the stage stands a small screen. A puppet appears from the right side, and is shown to go behind the screen. Another puppet appears from the left side, and again goes behind the screen. In one set of trials, the screen is raised and the two puppets appear. In another set of trials, the screen is raised, and only one puppet appears. In the former case, babies looked briefly at the puppets, but soon began looking elsewhere. In the latter case, the babies stared substantially longer at the single puppet. The babies have begun to understand the concept of addition...?

The conclusion here is that novelty elicits attention, which implies that there is an expectation of what is “normal” based on prior experience.

3.1 Putting it Together: The 3 Phases of Learning

In a broad sense, *learning is the process of building a desirable association between stimulus and response (domain and range), and is measured through resulting behavior on stimulus that has not been previously seen.*

It is important to give a context in which to study machine learning. We will think about learning as a three step process:

1. Building Experience. In hearing this is the processing done at the ear; in sight, this is the processing done at the eye. In mathematics, this might be characterized as getting an understanding of numbers (as a precursor

to Algebra), or of functions (as a precursor to Calculus). In our discussion topic, this is the stage at which we convert alphabetic characters to numerical values. This is where we build the initial tools to distinguish between input patterns that are important, and remove patterns that are not.

2. Building an Internalization.

3. Evaluation of the Internalization. Here we check the internalization that has been built to see if it is robust to changes in the stimulus. Errors here can be coming from either of the first two stages; perhaps the experience base was not broad enough, or the ability to distinguish has not been established.

There could probably be much debate over this kind of categorization, but it is useful as a starting point. In fact, many might argue that points (1) and (2) are essentially the same. We want to distinguish between these for reasons we'll discuss later.

In a machine learning sense at least, the process of learning can be cast in one of two ways:

- Build a function (when the outputs are unique), $y = f(x)$.
- Estimation of a probability distribution (when each input has a whole range of possible outputs), $y = p(x)$, where $p(x) > 0$ for all x , and $\int_{-\infty}^{\infty} p(x) dx = 1$

The first model is probably somewhat easier than the second- both in terms of the theory needed, and in the algorithms used for training. We will focus mostly on this first paradigm. The second is studied more in depth in a field known as Reinforcement Learning (in computer science terms, which is different than the psychology term). These two models also distinguish between two different types of learning paradigms: Supervised versus Unsupervised Learning.

In *supervised learning*, we are given examples of proper behavior, and we want the computer to emulate that behavior. In *unsupervised learning*, no desired outputs are given. This might seem to be an impossible task: How can we learn something if we don't know what it is we're learning? Before continuing, let's give some examples of each type:

- **Discovery Learning:** You might have some experience with this. This is where you give a group of students a problem to solve, but do not specify appropriate tools or give examples on how to solve it. This is unsupervised learning.
- Input-output pairs are given, and we wish to find the association. This is supervised learning, since for each input, we're told what the output should be.

- A data set is given, and we'd like to determine *clusters* for the data. That is, we seek a membership function, $m(x)$ where we input a given data point, and output an integer $1, 2, 3, \dots, k$ where the integer corresponds to the cluster the data is in. In this paradigm, no value of k is given beforehand. This is unsupervised learning, since we don't have any specific goal in mind (number of clusters).
- We can change the last problem to become supervised learning if, for each data point, we also supply a label- for example, each data point belongs to either a "red group", a "blue group", or a "green group".
- **The Mountain Car Problem** Suppose a car is sitting in the valley of a hill. If we apply full throttle, there is not enough power to get us to the top of the hill. If we apply full throttle for a few minutes, then allow gravity to take us down the hill and up the other side, we can start to build momentum by stepping on the throttle when we start to come back up the hill again. Therefore, we must go back and forth across the valley, building up more and more momentum until we have enough power to get us over the top. The problem: For each hillside position, give the sequence of throttles and slides that get us to the top. In this example, we have a goal in mind, but we don't tell the machine how to reach that goal. This again is unsupervised learning.

In general, supervised learning is easier than unsupervised learning in the sense that in unsupervised learning, there will probably be a lot of time "wasted" in trial-and-error where we have to explore the entire input space. In terms of the three processes of learning that was given before:

- In supervised learning, the step of building experience can be separated from the internalization process. We are told explicitly how to differentiate between inputs- we are given some measure of what are the distinguishing characteristics. In this setting, the "expert" sets up the examples of proper behavior.
- In unsupervised learning, these two steps are occurring simultaneously- as we build experience, the internalization is being updated, then we go back to gain more experience to change the internalization again, and the process repeats itself until some goal has been attained. In this setting, we will sometimes get feedback from some "expert". For example, rather than being told that an association is the right association, we might be told that "we're getting warmer or getting colder".

In supervised learning, we can be more specific. The three parts to learning can be recast:

1. The preprocessing phase. In this phase, we are interested in analyzing the desired input-output relationships to see if there are undesirable elements to be removed. That is, we want to find some kind of essential representation of the data. This phase could be doing things like noise removal,

and in general, can be cast as dimensionality reduction. Sometimes this phase is not necessary, but removing redundant information can make the next phase easier- that is, we don't want our learning to become confused. Ideally:

- The compression results in an optimally small set (we will need to construct a measure of what we will mean by “optimal”) - small in the sense that the representation carries no redundant information. On the other hand, the resulting compression should also be maximal in the sense that all training data is well represented; we want to be able to maintain an ability to distinguish between distinct stimuli.
 - The compression results in a true representation of the original set. If the function f performs the compression, then this will imply that f should be 1-1 on its range, so that it is invertible (on the appropriately restricted domain and range).
2. The learning phase. In this phase, we select an appropriate mathematical model and an appropriate metric so we can measure how well the model is associating. Normally, the model will depend on its *parameters*, and learning is performed via maximization or minimization of the metric. Here we will study some general techniques of optimization theory.
 3. The postprocessing phase. Once the associations have been established, it is critical to test the model on new data- is the model generalizing well? It could be that the model is too simple or too complex.

3.2 Discussion:

1. Consider the concept of *superstition*: This is a belief that one must engage in certain behaviors in order to receive a reinforcement, where in reality, the reinforcement did not depend on those behaviors. Is it possible for a computer to engage in superstitious activity? Discuss in terms of the supervised versus unsupervised learning paradigms.
2. In psychology, an *unconditioned response* is something that is automatic and not something that needs to be learned. For example, salivation with the presence of food, and constriction of the pupil of the eye in the presence of light are unconditioned responses. An *unconditioned stimulus* is the stimulus that elicits an unconditioned response. In classical conditioning, we can train animals to perform a given task by first associating a conditioned stimulus to an unconditioned stimulus. Soon, the conditioned stimulus will elicit the unconditioned response (therefore, the conditioned stimulus becomes a conditioned response). Can such a paradigm be useful in machine learning? Discuss classical conditioning in terms of what we called “learning” in this chapter.

3. Probabilities

A signal light comes on and is followed by one of two other lights. The goal is to predict which of the lights comes on given that the signal light comes on. The experimenter is free to arrange the pattern of the two response lights in any way- for example, one might come on 75% of the time.

Let E_1, E_2 denote the event that the first (second) light comes on, and let A_1, A_2 denote the prediction that the first (second) light comes on (respectively). Let π be the probability that E_1 occurs.

- (a) If the goal is to maximize your reward through accurate predictions, what should you do in this experiment? Just give a heuristic answer- you do not have to formally justify it.
- (b) How would you program a machine to maximize it's prediction accuracy? Can you state this in mathematical terms?
- (c) What do you think happens with actual subject (human) trials?

4. Harlow's Monkeys (Harry F. Harlow)

In a test similar to the previous experiment, Harlow had monkeys learn a discrimination task. On each problem, they were to pick one of two objects that had a reinforcer placed under it. There were a total of 344 tasks, each with a different set of objects. In light of the previous problem, how do you think the monkeys performed early on in the trials versus later in the trials?

5. In Jean Piaget's theory of learning, *assimilation* and *accommodation* play a key part. Assimilation refers to a matching between a given cognitive structure and the environment. For example, a cognitive structure defines the boundaries on what a child can learn- if only sucking, looking, reaching and grasping are available to a baby, then experiences are assimilated in those terms. Accommodation is the process by which the cognitive structure is modified. In Calculus I, we learn what a derivative is. Describe this process in terms of assimilation and accommodation.

3.3 A Case Study in Reinforcement Learning

In this section, we will summarize much of what we've been talking about in learning theory in terms of a special case: The problem of the n -armed bandit (Barto and Sutton). The name comes from a slang term for a slot machine, also known as a one-armed bandit.

In the n -armed bandit problem, we have n slot machines (or equivalently, one machine with n -arms), each having a different probability of winning. The goal is to play the slots and to maximize our returns. Before continuing, let us set up some notation: Define

$$Q(a) = \text{The expected return for playing slot machine } a$$

You can also think of $Q(a)$ as the *mean* of the payoffs for slot machine a . If we knew $Q(a)$ for each machine a , our strategy would be very simple: Play the machine with the largest payoff, and ignore the others. The interesting part of the problem is that we don't know what the returns are! So, we have to estimate the returns. Let

$$Q_t(a) = \text{Our estimation of } Q(a) \text{ at time } t$$

Our general algorithm should be set up so that our estimates get better over time:

$$\lim_{t \rightarrow \infty} Q_t(a) = Q(a) \quad (3.1)$$

Let's start things off with a straightforward estimate. Suppose we play slot machine a a total of n_a times, with payoffs r_1, \dots, r_{n_a} (note that these values could be negative!). Then we might estimate $Q(a)$ as the mean of these values:

$$Q_t(a) = \frac{r_1 + r_2 + \dots + r_{n_a}}{n_a}$$

In statistical terms, we are using the sample mean to estimate the actual mean. In statistics, the Law of Large Numbers states that as the number of samples increase, the sample mean converges to the population mean. We should also agree on a starting value: By definition, we'll take $Q_0(a) = 0$.

3.3.1 Strategies for the n -armed bandit

In this section, we'll describe different strategies one might take to *learn* $Q(a)$.

The Greedy Algorithm

This strategy is straightforward: Always play the slot machine with the largest (estimated) payoff. If a_{t+1} is the machine we'll play at time $t + 1$, then:

$$a_{t+1} = \arg \max \{Q_t(1), Q_t(2), \dots, Q_t(n)\}$$

where "arg" refers to the argument of the maximum (which is an integer from 1 to n corresponding to the max). If there is a tie, then choose one of them at random. Let's take a moment to see how we might do this in Matlab. The `find` command will be used to find the maximum values:

```
idx=find(x==max(x))
```

will return all indices of the vector x that are equal to the max. We'll talk about the double equality in class. You should try this command with something like:

```
x=[1 2 3 0 3];
idx=find(x==max(x));
```

and the result will be a vector, `idx`, will contain the values 3 and 5.

Going back to the original strategy, I think you'll see a problem- What if the estimations are wrong? Then its very possible that you'll get stuck on a suboptimal machine. This problem can be dealt with in the following way: Every once in a while, try out the other machines to see what you get. This is what we'll do in the next section.

The ϵ -Greedy Algorithm

In this algorithm we will choose, with probability ϵ , a machine totally at random. Notice that in this case, the number of trials gets larger and larger, $n_a \rightarrow \infty$ for *all* machines a , and so we will be guaranteed convergence to the proper estimates of $Q(a)$ for all a machines. On the flip side, because we're always investigating other machines every once in a while, we'll never maximize our returns.

You might note the relationship of what we've just said to Estes' probability matching experiments!

```
for j=1 to Number of Trials

    • Select an action:

        - Sometimes choose one at random

        - Otherwise, select the action with greatest return. Check
          for ties, and if there is a tie, pick on of them at random.

    • Get your payoff

    • Update the estimates  $Q$ 
```

The first programming problem will be to decide on how to sometimes select a machine at random, and sometimes not. If $\epsilon = E$ is the probability of this event, and N is the number of trials, then one way of selection is to set up an N -vector, *greedy*, that will "flag" the events- that is, on trial j , if *greedy*(j) = 1, choose a machine at random. Otherwise, choose using the greedy method. The following code will do just that (N is the number of trials)

```
greedy=zeros(1,N);
if E>0
    m=round(E*N); %Total number of times we should choose at random
    greedy(1:m)=ones(1,m);
    m=randperm(N); %Randomly permute the vector indices
    greedy=greedy(m);
    clear m
end
```

And here's the full function. We assume that the actual rewards for each of the bandits is given in the vector A_q , and that when machine a is played, the sample reward will be chosen from a normal distribution with unit variance and mean $A_q(a)$.

```

function [As,Q,R]=banditE(N,Aq,E)
%FUNCTION [As,Q,R]=banditE(N,Aq,E)
% Performs the N-armed bandit example using epsilon-greedy
% strategy.
% Inputs:
%     N=number of trials total
%     Aq=Actual rewards for each bandit (these are the mean rewards)
%     E=epsilon for epsilon-greedy algorithm
% Outputs:
%     As=Action selected on trial j, j=1:N
%     Q are the reward estimates
%     R is N x 1, reward at step j, j=1:N

numbandits=length(Aq);      %Number of Bandits
ActNum=zeros(numbandits,1); %Keep a running sum of the number of times
                             % each action is selected.
ActVal=zeros(numbandits,1); %Keep a running sum of the total reward
                             % obtained for each action.
Q=zeros(1,numbandits);      %Current reward estimates
As=zeros(N,1);              %Storage for action
R=zeros(N,1);               %Storage for averaging reward

%*****
% Set up a flag so we know when to choose at random (using epsilon)
%*****
greedy=zeros(1,N);
if E>0
    m=round(E*N); %Total number of times we should choose at random
    greedy(1:m)=ones(1,m);
    m=randperm(N);
    greedy=greedy(m);
    clear m
end
if E>=1
    error('The epsilon should be between 0 and 1\n');
end
%*****
%
% Now we're ready for the main loop
%*****
for j=1:N
    %STEP ONE: SELECT AN ACTION (cQ) , GET THE REWARD (cR) !
    if greedy(j)>0
        cQ=ceil(rand*numbandits);
        cR=randn+Aq(cQ);
    else
        [val,idx]=find(Q==max(Q));
        m=ceil(rand*length(idx)); %Choose a max at random
        cQ=idx(m);
        cR=randn+Aq(cQ);
    end
    As(j)=cQ;
    ActNum(j)=ActNum(j)+1;
    ActVal(j)=ActVal(j)+cR;
    Q(cQ)=Q(cQ)+cR/ActNum(cQ);
    R(j)=cR;
end

```

```

    end
    R(j)=cR;
    %UPDATE FOR NEXT GO AROUND!
    As(j)=cQ;
    ActNum(cQ)=ActNum(cQ)+1;
    ActVal(cQ)=ActVal(cQ)+cR;
    Q(cQ)=ActVal(cQ)/ActNum(cQ);
end

```

Next we'll create a test bed for the routine. We will call the program 2,000 times, and each call will consist of 1,000 plays. We will set the number of bandits to 10, and change the value of ϵ from 0 to 0.01 to 0.1, and see what the average reward per play is over the 1000 plays.

Here's a script file that we'll use to call the `banditE` routine:

```

Ravg=zeros(1000,1);
E=0.1;
for j=1:2000
    m=randn(10,1);
    [As,Q,R]=banditE(1000,m,E);
    Ravg=Ravg+R;
    if mod(j,10)==0
        fprintf('On iterate %d\n',j);
    end
end
Ravg=Ravg./2000;
plot(Ravg);

```

The output of the algorithms are shown in Figure 3.1.

The Softmax Action Selection

This technique is fairly common across many algorithms. In this particular instance, we can modify the ϵ -greedy algorithm so that the current estimates of rewards are interpreted as probabilities. That is, if one of the machines is giving a poor return, then we won't visit that machine as often as the rest. In this case, we will use a Gibbs (or Boltzmann) distribution, and create a probability measure in the following way: Let

$$S_t = \sum_{k=1}^{\text{numbandits}} \exp\left(\frac{Q_t(k)}{\tau}\right)$$

and the probability of choosing bandit a :

$$\pi_t(a) = \frac{\exp\left(\frac{Q_t(a)}{\tau}\right)}{S_t}$$

Remarks about this method:

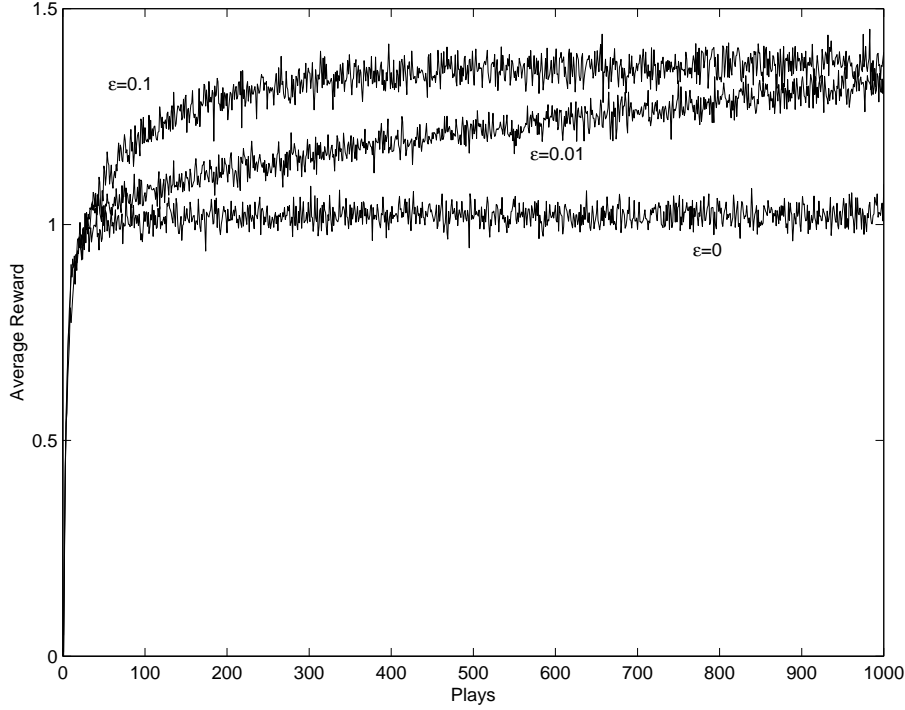


Figure 3.1: Results of the testbed on the 10-armed bandit. Shown are the rewards given per play, averaged over 2000 trials.

1. Using the exponential function has the effect of making all the results positive, so that we indeed have a probability measure. We are also making them all sum to 1.
2. τ plays the role of a temperature:

- The effect of hot temperatures, $\tau \rightarrow \infty$, would be to set all the probabilities approximately equal, since $\lim_{\tau \rightarrow \infty} Q_t(a)/\tau = 0$
- The effect of cold temperatures, $\tau \rightarrow 0$ would be effectively the same as a greedy selection. To see this, we show it for the case of two bandits, and let $Q_t(1)$ be the maximum reward. Then:

$$\lim_{\tau \rightarrow \infty} \left(\frac{\exp\left(\frac{Q_t(1)}{\tau}\right)}{\exp\left(\frac{Q_t(1)}{\tau}\right) + \exp\left(\frac{Q_t(2)}{\tau}\right)} \right) = \lim_{\tau \rightarrow \infty} \frac{1}{1 + \exp\left(\frac{(Q_t(2) - Q_t(1))}{\tau}\right)} = 1$$

and similarly,

$$\lim_{\tau \rightarrow \infty} \left(\frac{\exp\left(\frac{Q_t(2)}{\tau}\right)}{\exp\left(\frac{Q_t(1)}{\tau}\right) + \exp\left(\frac{Q_t(2)}{\tau}\right)} \right) = \lim_{\tau \rightarrow \infty} \frac{\exp\left(\frac{(Q_t(2) - Q_t(1))}{\tau}\right)}{1 + \exp\left(\frac{(Q_t(2) - Q_t(1))}{\tau}\right)} = 0$$

3. In Matlab, these probabilities are easy to program. Let Q be a vector holding the current estimates of the returns, as before, and let $t = \tau$, the temperature. Then we construct a vector of probabilities using the softmax algorithm:

```
P=exp(Q./t);
P=P./sum(P);
```

Exercise in Programming

How to select action a with probability $p(a)$? We could do what we did before, and create a vector of choices with those probabilities fixed, but our probabilities change. We can also use the uniform distribution, so that if $x = \text{rand}$, and $x \leq p(1)$, use action 1. If $p(1) < x \leq p(1) + p(2)$, choose action 2. If $p(1) + p(2) < x \leq p(1) + p(2) + p(3)$, choose action 3, and so on. There is an easy way to do this, but it is not optimal (in terms of speed). We introduce two new Matlab functions, `cumsum` and `histc`.

The function `cumsum`, which means *cumulative sum*, takes a vector x as input, and outputs a vector y so that `y=cumsum(x)` creates:

$$y_k = \sum_{n=1}^k x_n = x_1 + x_2 + \dots + x_k$$

The function `histc` (for *histogram count*) has the form: `n=histc(x,y)`, where the vector y is monotonically increasing. The elements of y form “bins” so that $n(k)$ counts the number of values in x that fall between the elements $y(k)$ (inclusive) and $y(k+1)$ (exclusive) in the vector y . In our particular case, x is a scalar, so that $n(k)$ will have exactly one value of 1, and the rest will be zero. Here is a script that will test these functions:

```
p=[.30,.10,.20,.40]
P=cumsum([0 p]);
for j=1:1000
    x=rand;
    n=histc(x,P);
    v(j)=find(n==1); %v holds the action choice
end
```

In the script, we put the action choice in a vector so that we could inspect it later. Normally, we wouldn't need to do this. In one sample run of the script, the algorithm output:

Choice	Percent Chosen	Actual
1	0.297	0.30
2	0.104	0.10
3	0.197	0.20
4	0.402	0.40

We can see from this table that the algorithm is performing quite well. Your exercise is to modify the program `banditE` to create a program `banditS` that will implement the softmax strategy. You may assume that the temperature `tau` is an additional input parameter. Modify this program so that `tau` begins large and becomes small during the training.

“Win-Stay, Lose-Shift” Strategy

The “Win-Stay, Lose-Shift” strategy discussed in terms of Harlow’s monkeys and perhaps the probability matching experiments of Estes might be re-formulated here for the n -armed bandit experiment.

In this experiment, we interpret the strategy as: If I’m winning, make the probability of choosing that action stronger. If I’m losing, make the probability of choosing that action weaker. This brings us to the class of *pursuit* methods. That is, if we take $a_{t+1}^* = \max_a Q_t(a)$, then we can update the probability of choosing the best action by incrementing its probability towards 1. If we define $\pi_t(a)$ as the probability of choosing action a at time t , then (we leave off the subscript on the action, assume the time is given):

$$\pi_{t+1}(a^*) = \pi_t(a^*) + \beta [1 - \pi_t(a^*)]$$

and the rest of the probabilities decrease towards zero:

$$\pi_{t+1}(a) = \pi_t(a) + \beta [0 - \pi_t(a)]$$

Exercises with the Pursuit Strategy

1. Suppose we have three probabilities, P_1, P_2, P_3 , and P_1 is the unique maximum. Show that, for any $\beta > 0$, the updated values still sum to 1.
2. Using the same values as before, show that, for any $\beta > 0$, the updated values will stay between 0 and 1- that is, If $0 \leq P_i \leq 1$ for all i before the update, then after the update, $0 \leq P_i \leq 1$.
3. Using your ideas from the previous two problems, how should we update if we have a tie for the maximum?
4. Suppose that for some fixed j , P_j is always a loser (never a max). Show that the update rule guarantees that $P_j \rightarrow 0$ as $t \rightarrow \infty$. HINT: Show that $P_j(t) = (1 - \beta)^t P_j(0)$
5. Suppose that for some fixed j , P_j is always a winner (with no ties). Show that the update rule guarantees that $P_j \rightarrow 1$ as $t \rightarrow \infty$.
6. Modify the previously written program, `banditS`, to create a new program, `banditP` that implements the pursuit strategy.

7. (*) Compare the three reinforcement learning algorithms using the test script. We'll have one group test the softmax strategy with different values of temperature τ , and another group test the pursuit strategy with different values of β . Finally, we'll compare the outputs of these algorithms with the ϵ -greedy method, with $\epsilon = 0.1$.

3.4 Donald Olding Hebb and Hebbian Learning

D.O. Hebb (1904-1985) was a physiological psychologist at McGill University. We will discuss his influence on machine learning, but he also did some very interesting experiments to determine how the brain functions. Among the research²:

- Evidence for the plasticity of the brain:

Much to Hebb's amazement, he found that even after substantial loss of tissue from the frontal lobes of the brain, there was no loss in intelligence, and in some cases, he even detected a gain in intelligence³.

- "Monkey Horror" in chimpanzees.
- Sensory Deprivation experiments: Sensory input is vital for brain functioning. Over time without sensory stimulation, subjects hallucinate and their personalities begin to break down.

In Hebb's view, learning could be described physiologically. That is, there is a physical change in the nervous system to accommodate learning, and that change is summarized by what we now call Hebb's postulate (from his 1949 book):

When an axon of cell A is near enough to excite a cell B and repeatedly takes part in firing it, some growth process or metabolic change takes place in one or both cells such that A's efficiency, as one of the cells firing B, is increased.

As with many named theorems and postulates, this was not an idea that was completely new, but he does give the postulate in a form that can be used as a basis for machine learning.

3.4.1 Linear Neural Nets and Hebbian Learning

We will go into the formal details later for defining *neural nets*, but this is a good place to get a feel for what they're all about.

²As usual, these footnote sketches are somewhat oversimplified. Read about Hebb's work from his writings, especially "The Organization of Behavior" (1949), and his autobiography in 1980 "A History of Psychology through Autobiography", vol VII

³An Introduction to the Theories of Learning, p. 394

Let us first build a simple model for a neuron. A neuron has three body parts- The Dendrites, which carry information to the Cell Body, the Cell Body, and the Axon, which carries information away from the Cell Body.

Multiple signals come in to the cell body from the dendrites. Mathematically, we will assume they all arrive at the same time, and the action of the dendrites (or the arrival site of the cell body) is that each signal is changed by the physiology of the cell. That is, if x_i is information along dendrite i , arrival at the cell body changes it to $w_i x_i$, where w_i is some real value. Next, the cell body collates this information by summing these signals together. So far, then, this action is simply a dot product of the vector of w 's (the *weights*) to the signal x . For the purposes of this section, we will assume no further processing.

If the signal is passed to a *cell assembly*, or group of neurons, then the overall action of the group is simply a linear mapping:

$$\mathbf{x} \rightarrow W\mathbf{x} = \mathbf{y}$$

If we have k neurons and \mathbf{x} is a vector in \mathbb{R}^n , then W is a $k \times n$ matrix, and each row corresponds to a signal neuron's weights (that is, W_{ij} refers to the weight taking x_j to neuron i).

A further assumption will be to specify the mathematical realization of Hebb's postulate. Suppose that we present the network with a pattern, $\mathbf{x} \in \mathbb{R}^n$, and it outputs a pattern, $\mathbf{y} \in \mathbb{R}^k$. In terms of each weight, we see that

W_{ij} connects the j^{th} value of the input to the i^{th} value of the output.

Thus we might take the following as Hebb's Rule:

$$\Delta W_{ij} = \alpha y_j x_i$$

where α is called **the learning rate**. If both x_i and y_j match in sign, then W_{ij} becomes larger. If there is a mismatch in sign, W_{ij} gets smaller. This is the *unsupervised Hebbian rule*. What should we do if we want the network to learn a specific association?

Let t be the target (or desired) value for the input \mathbf{x} . In the supervised Hebbian rule, we need to update the weights based on what we want the network to do, rather than what the network is already doing:

$$\Delta W_{ij} = \alpha t_j x_i$$

There is still something unsatisfying here- When should we stop the training? It seems like the weights could diverge as training progresses, and furthermore, we're not taking the actual outputs of the network into account. Heuristically, we would like for the update to go to zero as the target values approach the network outputs. This leads us to our final modification of our basic Hebb rule, and is called the Widrow-Hoff learning rule⁴:

$$\Delta W_{ij} = \alpha(t_j - y_j)x_i$$

⁴Also goes by the names Least Mean Squares rule, and the delta rule.

If we put this in matrix form, the learning rule becomes:

$$W^{\text{new}} = W^{\text{old}} + \alpha (\mathbf{t} - \mathbf{y}) \mathbf{x}^T$$

where (\mathbf{x}, \mathbf{t}) is a desired input-output relation, and $\mathbf{y} = W\mathbf{x}$.

Additionally, sometimes it is appropriate to add a bias term so that the network has the output:

$$\mathbf{y} = W\mathbf{x} + \mathbf{b}$$

Later we'll describe the details, but for now recall that a pure linear mapping must take the zero vector to the zero vector, which may not be what we want. The bias update is similar to the previous update:

$$\mathbf{b}^{\text{new}} = \mathbf{b}^{\text{old}} + \alpha (\mathbf{t} - \mathbf{y})$$

3.4.2 Exercise

Let's put all of this together to solve a pattern classification problem. Suppose we are given the following associations:

Point	Class
(1, 1)	1
(1, 2)	1
(2, -1)	2
(2, 0)	2
(-1, 2)	3
(-2, 1)	3
(-1, -1)	4
(-2, -2)	4

Graphically, we can see the classes in the plane in Figure 3.2. There are several ways of performing the desired mapping- for example, the outputs could be 1, 2, 3, 4. But this may have unintended consequences. In this case, the metric between outputs would imply that Class 4 is much farther away from Class 1 than Class 3. A better method may be to take Class 1 to be the vector $[-1, -1]^T$, Class 2 is the vector $[-1, 1]^T$, Class 3 is $[1, -1]^T$, and Class 4 is $[1, 1]^T$. Now the 4 classes are on the vertices of a square.

Now for the details of the program. First write the inputs as an 2×8 matrix, with a corresponding output matrix that is also 2×8 . Parameters that can be placed first will be the maximum number of times through the data N and the learning rate, a , which we will set to 0.04. We can also set an error bound so that we might stop early. Set the initial weights to the 2×2 identity, and the bias vector b to $[1, 1]^T$.

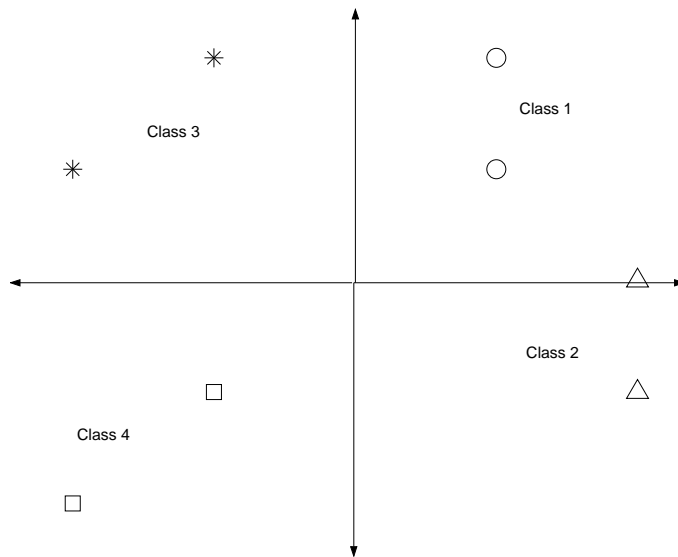


Figure 3.2: Pattern Classification Problem. Each point is a sample of one of the four classes.

Check your algorithm using the following data: After using the first data point, W and b should have the values:

$$W = \begin{bmatrix} 0.88 & -0.12 \\ -0.12 & 0.88 \end{bmatrix}, \quad b = \begin{bmatrix} 0.88 \\ 0.88 \end{bmatrix}$$

Let the program run until you think it has converged (this is your choice). Then, plot the points $\mathbf{x} = [x, y]^T$ so that $W\mathbf{x} + \mathbf{b} = \mathbf{0}$, which will be the two lines:

$$W_{11}x + W_{12}y + b_1 = 0, \quad W_{21}x + W_{22}y + b_2 = 0$$

These lines form what is called the *decision boundary*.

