

Introduction to Neural Networks

U. Minn. Psy 5038

Spring, 1998

Linear systems

Introduction to Learning and Memory

Linear systems, matrices & neural networks

Introduction

Consider the generic 2-layer network. It consists of a weighted average of the inputs (stage 1), followed by a point-nonlinearity (the squash function of stage 2), and added noise (stage 3). Although later we will see how the non-linearity enables computations that are not possible without it, useful functions can be realized with just the linear or stage 1 part of the network. We've seen one application already with the model of the limulus eye. In the next lecture, we will see how linear networks can be used to model associative memory. But first, let us take what we've learned so far about modeling linear networks and look at in the general context of linear systems theory. Our 2-layer net is a matrix of weights that operates on a vector of input activities by computing a weighted sum. One property of such a system is that it satisfies the fundamental definition of a "linear system", which we will define shortly.

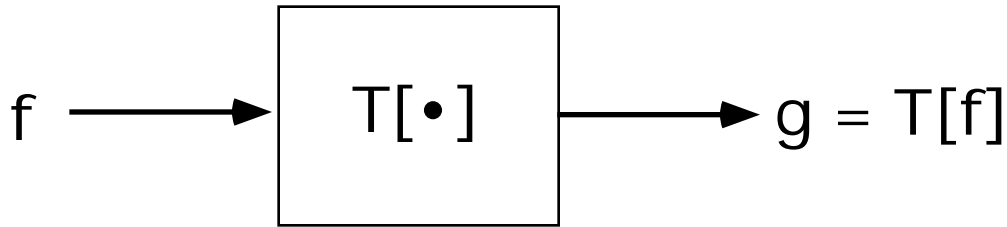
The world of input/output systems can be divided up into linear and non-linear systems. Linear systems are nice because the mathematics that describes them is not only well-known, but also has a mature elegance. On the other hand, it is a fair statement to say that most real-world systems are not linear, and thus hard to analyze...but fascinating if for that reason alone. Scientists were lucky with the limulus eye. That nature is usually non-linear doesn't mean one shouldn't familiarize oneself with the basics of linear system theory. Many times non-linear systems can be approximated by a linear one over some restricted range of parameter values.

So exactly what is a "linear system"?

The notion of a "linear system" is a generalization of the input/output properties of a straight line passing through zero. The matrix equation $\mathbf{W} \cdot \mathbf{x} = \mathbf{y}$ is a linear system. This means that if \mathbf{W} is a matrix, $\mathbf{x1}$ and $\mathbf{x2}$ are vectors, and a and b are scalars:

$$\mathbf{W} \cdot (a \mathbf{x1} + b \mathbf{x2}) = a \mathbf{W} \cdot \mathbf{x1} + b \mathbf{W} \cdot \mathbf{x2}$$

This is a consequence of the laws of matrix algebra. The idea of a linear system has been generalized beyond matrix algebra. Imagine we have a box that takes inputs such as f , and outputs $g = T[f]$.



The abstract definition of a linear system is that it satisfies:

$$T[a f + b g] = a T[f] + b T[g]$$

where T is the transformation that takes the sum of scaled inputs f, g (which can be functions or vectors) to the sum of the scaled transformation of f and g . The property, that the output of a sum is the sum of the outputs, is sometimes known as the *superposition principle* for linear systems. The fact that linear systems show superposition is good for doing theory, but as we will see later, it limits the kind of computations that can be done with linear systems, and thus with linear neural network models.

Characterizing a linear system by its response to an orthonormal basis set

Suppose we have an unknown physical system, which we model as a linear system T :

```
T = Table[Random[], {i, 1, 8}, {j, 1, 8}];
```

We would like to make a simple set of measurements that could characterize T in such a way that we could predict the output of T to any input. This is the sort of task that engineers face when wanting to characterize, say a stereo amplifier (as a model linear system), so that the output sound can be predicted for any input sound. What kind of measurements would tell us what T is? Well, we could just "stimulate" the system with cartesian vectors $\{1, 0, 0, 0, 0, 0, 0, 0\}$, $\{0, 1, 0, 0, 0, 0, 0, 0\}$, and so forth and collect the responses which would be the columns of T . This has two practical problems: 1) for a real physical system, such as your stereo, or a neuron in the limulus eye, this would require stimulating it with a high-intensity audio or light intensity spike, which could damage what you are trying to study; 2) Characterizing the linear system by a matrix T , requires n^2 numbers, where n is the input signal vector length--and n can be pretty big for both audio and visual systems. Problem 2) has a nice solution when T is symmetric, and is addressed later. Problem 1) can be addressed by showing that we can characterize T with any basis set--so we can pick one that won't blow out the physical system being tested.

The set of Walsh functions we looked at earlier is just one possible set that has the advantage that the elements that contribute to the "energy", i.e. (the square of the length) are distributed across the vector.

```
Vectorlength[x_] := N[Sqrt[x.x]]
```

```

v1 = {1, 1, 1, 1, 1, 1, 1, 1}; w1 = v1/Vectorlength[v1];
v2 = {1,-1,-1, 1, 1,-1,-1, 1}; w2 = v2/Vectorlength[v2];
v3 = {1, 1,-1,-1,-1,-1, 1, 1}; w3 = v3/Vectorlength[v3];
v4 = {1,-1, 1,-1,-1, 1,-1, 1}; w4 = v4/Vectorlength[v4];
v5 = {1, 1, 1, 1,-1,-1,-1,-1}; w5 = v5/Vectorlength[v5];
v6 = {1,-1,-1, 1,-1, 1, 1,-1}; w6 = v6/Vectorlength[v6];
v7 = {1, 1,-1,-1, 1, 1,-1,-1}; w7 = v7/Vectorlength[v7];
v8 = {1,-1, 1,-1, 1,-1, 1,-1}; w8 = v8/Vectorlength[v8];

```

We have already seen that the set $\{w_i\}$ spans 8-space in such a way that we can easily express any vector as a linear sum of these basis vectors.

$$\mathbf{g} = \sum (\mathbf{g} \cdot \mathbf{w}_i) \mathbf{w}_i \quad (1)$$

So as we saw before, an arbitrary vector, \mathbf{g}

```
g = {2,6,1,7,11,4,13, 29};
```

is the sum of its own projections onto the basis set:

```

(g.w1) w1 + (g.w2) w2 + (g.w3) w3 + (g.w4) w4 +
(g.w5) w5 + (g.w6) w6 + (g.w7) w7 + (g.w8) w8

```

```
{2., 6., 1., 7., 11., 4., 13., 29.}
```

Suppose we now do an "experiment" to find out how \mathbf{T} transforms the vectors of our basis set; and we put all of these transformed basis elements into a new set of vectors **newW**[[i]]. **newW** is a matrix for which each row is the response of \mathbf{T} to a basis vector.

```
newW = {T.w1,T.w2,T.w3,T.w4,T.w5,T.w6,T.w7,T.w8};
```

Note that **newW** is an 8x8 matrix. So how can we calculate the output of \mathbf{T} , given \mathbf{g} without actually running the input through \mathbf{T} ? If we do run the input through \mathbf{T} we get:

```
T.g
```

```
{26.4304, 36.2209, 23.9967, 32.6211, 24.4136, 44.8792, 41.2066, 46.1586}
```

But by the principle of linearity, we can also calculate the output by finding the "spectrum" of \mathbf{g} as in Problem Set 2, and then scaling each of the transformed basis elements by the spectrum and adding them up:

$$\mathbf{T} \cdot \mathbf{g} = \mathbf{T} \cdot \left\{ \sum (\mathbf{g} \cdot \mathbf{w}_i) \mathbf{w}_i \right\} = \sum (\mathbf{g} \cdot \mathbf{w}_i) \mathbf{T} \cdot \mathbf{w}_i \quad (2)$$

```
(g.w1) T.w1 + (g.w2) T.w2 + (g.w3) T.w3 + (g.w4) T.w4 +
(g.w5) T.w5 + (g.w6) T.w6 + (g.w7) T.w7 + (g.w8) T.w8
```

```
{26.4304, 36.2209, 23.9967, 32.6211, 24.4136, 44.8792, 41.2066, 46.1586}
```

Of course, we have already done our "experiment", so we know what the transformed basis vectors are, we stored them as rows of the matrix **newW**. We can calculate what the spectrum (**g.wi**) is, so the output of **T** is:

```
(g.w1) newW[[1]] + (g.w2) newW[[2]] + (g.w3) newW[[3]] +
(g.w4) newW[[4]] + (g.w5) newW[[5]] + (g.w6) newW[[6]] +
(g.w7) newW[[7]] + (g.w8) newW[[8]]
```

```
{26.4304, 36.2209, 23.9967, 32.6211, 24.4136, 44.8792, 41.2066, 46.1586}
```

■ Same thing in more concise notation

Let the basis vectors be the rows of a matrix **W**:

```
W = {w1, w2, w3, w4, w5, w6, w7, w8};
```

So again, we can project **g** onto the rows of **W**, and then reconstitute it in terms of **W** to get **g** back again:

```
(W.g).W
```

```
{2., 6., 1., 7., 11., 4., 13., 29.}
```

```
g.Transpose[W].newW
```

```
{26.4304, 36.2209, 23.9967, 32.6211, 24.4136, 44.8792, 41.2066, 46.1586}
```

Show that $T == \text{Transpose}[\text{newW}].W$, and that the system output is thus: $\text{Transpose}[\text{newW}].W.g$

These new basis vectors do span 8-space, but they are not necessarily orthonormal. Under what conditions would they be orthogonal? What if the matrix **T** was symmetric, and we deliberately chose the basis set to describe our input to be the eigenvectors of **T**?

What if the choice of basis set is the set of eigenvectors of T ?

We've seen how linearity provides us with a method for characterizing a linear system in terms of the responses of the system to the basis vectors. The problem is that if the input signals are long vectors, say with dimension 40,000, then this set of basis vector responses is really big-- 1.6×10^9 .

Construct a symmetric matrix transformation, T . Show that if the elements of the basis set are the eigenvectors of T , then the transformation of any arbitrary input vector \mathbf{x} is given by:

$$T[\mathbf{x}] = \sum_i \lambda_i \mathbf{e}_i$$

Having the eigenvectors of T enables us to express the input and output of T in terms of the same basis set--the eigenvectors. All T does to the input is to scale its projection onto each eigenvector by the eigenvalue for that eigenvector. The set of these eigenvalues is sometimes called the *modulation transfer function* because it describes how the amplitude of the eigenvectors change as they pass through T .

Linear systems analysis is the foundation of Fourier analysis, and is why it makes sense to characterize your stereo amplifier in terms of frequency response. Sinewave inputs are the eigenvectors of your stereo system. The dimensionality is much higher--if you are interested in frequencies up to 20,000 Hz, your eigenvector for this highest frequency would have at least 40,000 elements--not just 8!

This kind of analysis has been applied not only to physical systems, but to a wide range of neural sensory systems. For the visual system alone, linear systems analysis has been applied to the cat retina (Enroth-Cugell and Robson, 1964), the monkey visual cortex, and the human contrast sensitivity system as a whole (Campbell and Robson, 1968).

Much empirical analysis has been done using linear systems theory to characterize neural sensory systems, and other neural systems such as those for eye movements. It works wonderfully as long as the linear system approximation holds. And it does do quite well for the lateral eye of the limulus, X-cells and P-cells of the mammalian visual system, over restricted ranges for so-called "simple" cells in the visual cortex, among others. The optics of the simple eye is another example of an approximately linear system. Many non-linear systems can be approximated as linear systems over smooth subdomains.

■ Learning and memory: Brief overview

Definition of learning and memory

It is curious to note that historically, the simple linear model that we will discuss came after much research had been devoted to non-linear learning models, based in particular on a special case of McCulloch-Pitts neurons, called the perceptron. Linear models have obvious limitations, (an obvious one is a consequence of superposition principle discussed above). Nevertheless, many of the interesting properties of non-linear systems can be understood in terms of small signal linearity properties. And linear systems are easy to study.

Memory is the adaptive change in the behavior of an organism as a consequence of past experience. More precisely one can distinguish learning and memory:

Learning has to do with acquisition, and memory with storage and retrieval of information.

Psychology and biology of learning and memory

■ Associative and non-associative memory

We are going to be talking about *associative memory*, so it is useful to bear in mind that not all memory is considered associative.

Associative memory: animal learns about the relationship of one stimulus to another or about the relationship between a stimulus and the organism's response.

Nonassociative memory: exposure to a single stimulus either once, or repeated offers opportunity to learning about it

■ Examples of nonassociative learning

Habituation (Ivan Pavlov)

decrease in behavioral reflex response to repeated non-noxious stimulus

Sensitization (pseudo-conditioning)

exposure to noxious stimuli increases sensitivity

More complex examples of nonassociative learning

memory for sensory record...although what one calls a stimulus vs. response is problematic

imitation learning

■ Associative learning

Classical conditioning (Ivan Pavlov)

CS -- US -- R

(tone) -- food -- salivation

tone -- air puff - eye blink

Temporal contiguity of CS-US, and frequency important, but not the only factors

Operant conditioning (Thorndike of Columbia U.)

also called instrumental or trial-and-error learning

e.g. hungry rat in a cage with a lever

occasionally the rat may press the lever out of curiosity, accident or whatever, but if it does and promptly receives some food (reward), the lever pressing (called the "operant") will increase above a spontaneous rate.

■ Animal learning

Conditioning is dependent on the degree to which the US (e.g. food) is expected (e.g. Kamin's 2 part experiment)

1. light - shock

measure strength of conditioning by how the light affects on-going behavior

2. (light-tone pair) - shock

3. test with tone alone, no suppression of ongoing behavior

"blocking phenomenon"

(Rescorla-Wagner model can account for this, and Widrow-Hoff).

Ecological constraints on learning

■ Human memory

■ Stages of memory

Main results for studies of cognitive psychology -- stages of memory

Iconic -- e.g. visual afterimages (1 sec)

Working memory (short-term memory, minutes to hours)

small capacity, disrupted by being knocked unconscious encoding?

short-term neural plastic events

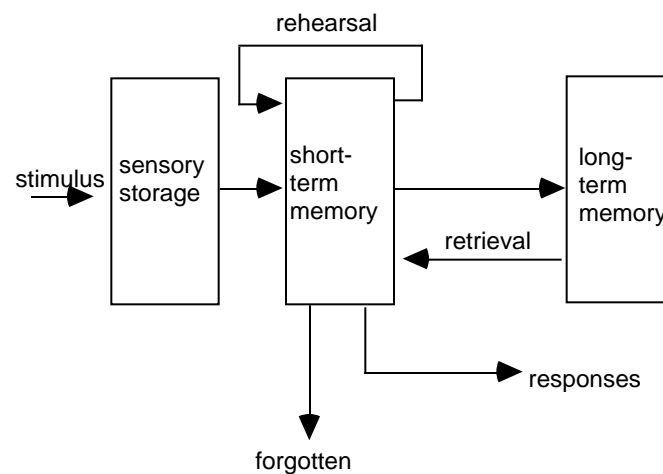
possible mechanisms?

rapid synaptic modification?

reverberating circuits?

Long-term memory

large capacity, relatively permanent (years)



■ Implicit (reflexive) vs. explicit (declarative) memory

reflexive -- automatic readout not dependent on awareness, or cognitive processes of comparison, evaluation

declarative -- "I saw a yellow canary yesterday"

relies on inference, comparison, and evaluation

Learning to driving a car -- memory moves from declarative towards reflexive with experience.

In our consideration of models, we will be primarily concerned with simple reflexive associative memory in which an neural input event leads to the reconstruction of an predictive response based on experience.

■ Hebbian rule for synaptic modification

■ Introduction

The brain has developed mechanisms that allow the animal to distinguish events that reliably and predictably occur together from those that do not. What kinds of neural models could be used to capture and retrieve associations? How can one pattern of neural activity come to be associated with another?

Assumptions

- physical basis of memory is in synaptic modification which alters the mapping of inputs to outputs
- the strength of the modification is determined by how often input and output activity occurs together, and by the strengths of the input and output activities.

The idea of learning as association goes back to William James (1890) (See Anderson text).

"When two elementary brain processes have been active together or in immediate succession, one of them on recurring, tends to propagate its excitement into the other (Psychology: Briefer Course)."

James also has an "activation equation" that sounds a lot like our stage 1 of the 2-layer feedforward network:

"The amount of activity at any given point in the brain-cortex is the sum of the tendencies of all other points to discharge into it"

Replace "point" by "neuron", and we have our limulus equation.

There are similar statements elsewhere (e.g. Kenneth Craik of Cambridge in the 1940's). But the clearest and most explicit statement of an associative learning rule is credited to Canadian Psychologist Donald Hebb:

"When an axon of cell A is near enough to excite a cell B and repeatedly or persistently 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" (Hebb in the "Organization of behavior", 1949).

Until recently, there was little evidence to support this notion of synaptic modification. But direct tests of Hebbian conjecture have now been made in rat hippocampus. See Anderson, chapter 6 for a review.

■ Hebbian synaptic modification: Modeling use the outer product

The fundamental Hebbian assumption is that synaptic strength grows with co-activity of pre- and post-synaptic activity. How can we quantify this? We will use a simple model of synaptic modification that assumes that the connection strength between neuron i and j is proportional to the product of the input and output activities.

$$W_{ij} = f_i g_j$$

If you remember the definition of the outer product of two vectors, you can see that the above rule is just that--an outer product of the the input and output activities.

■ Linear model of associative memory

■ Heteroassociation vs. autoassociation

We will look at two types of associative memory models. In heteroassociation, an input \mathbf{f} is associated with \mathbf{g} . So at some later time, when the system is stimulated with \mathbf{f} , it should produce \mathbf{g} . In autoassociation, an input \mathbf{f} is associated with itself. In the next lecture, we will see the use of autoassociation.

■ Summary of linear association model

1. Let $\{\mathbf{f}_n, \mathbf{g}_n\}$ be a set of input/output activity pairs. Memories are stored by superimposing new weight changes on old ones. Information from many associations is present in *each* connection strength.

$$\mathbf{W}_{n+1} = \mathbf{W}_n + \mathbf{g}_n \mathbf{f}_n^T$$

2. Let \mathbf{f} be an input possibly associated with output pattern \mathbf{g} . For recall, the neuron acts as a linear summer:

$$\mathbf{g} = \mathbf{W}\mathbf{f}$$

$$g_i = \sum_j w_{ij} f_j$$

3. If $\{\mathbf{f}_n\}$ are orthonormal, the system shows perfect recall:

$$\begin{aligned} \mathbf{W}_n \mathbf{f}_m &= (\mathbf{g}_1 \mathbf{f}_1^T + \mathbf{g}_2 \mathbf{f}_2^T + \dots + \mathbf{g}_n \mathbf{f}_n^T) \mathbf{f}_m \\ &= \mathbf{g}_1 \mathbf{f}_1^T \mathbf{f}_m + \mathbf{g}_2 \mathbf{f}_2^T \mathbf{f}_m + \dots + \mathbf{g}_m \mathbf{f}_m^T \mathbf{f}_m + \dots + \mathbf{g}_n \mathbf{f}_n^T \mathbf{f}_m \\ &= \mathbf{g}_m \\ \mathbf{f}_n^T \mathbf{f}_m &= \begin{cases} 1, & n = m \\ 0, & n \neq m \end{cases} \end{aligned}$$

So an $n \times n$ matrix has a memory capacity of n for orthogonal inputs. For random vectors, it is about 10-20%.

But linear association is also useful as a mapping, in which one wants to generalize in a smooth (actually linear) way to novel inputs. For example, linear regression in 2D can be done with a 2×2 matrix. It can be "trained" with a few input/output pairs of points $\{x, y\}$. Once trained, if the data are well-modeled by a straight line, the network will do a nice job of "predicting" the output value (y) given a novel input value (x).

Optional: eigenvectors, matrices and solving algebraic equations

Eigenvectors, eigenvalues: algebraic manipulation

Last time we used the *Mathematica* function **Eigenvectors[]** and **Eigenvalues[]** to produce the eigenvectors and eigenvalues for the matrix equation: $\mathbf{Ax} = \lambda \mathbf{x}$. It is worth spending a little time to understand what is being done algebraically. Consider the following pair of equations specified by a 2x2 matrix \mathbf{W} acting on \mathbf{xv} to produce a scaled version of \mathbf{xv} :

```
W = {{1,2},{2,1}};
xv := {x,y};
lambda xv == W.xv
```

```
{lambda x, lambda y} == {x + 2 y, 2 x + y}
```

Finding the eigenvalues is a problem in solving this set of equations. If we eliminate x and y from the pair of equations, we end up with a quadratic equation in λ :

■ Eliminate[] & Solve[]

```
Eliminate[{lambda xv == W.xv, lambda != 0}, {x, y}]
```

```
lambda - 3 != 0 & lambda != 0 & lambda + 1 != 0 & lambda^2 - 2 lambda == 3
```

```
Solve[-2 lambda + lambda^2 == 3, lambda]
```

```
{{lambda -> -1}, {lambda -> 3}}
```

So our eigenvalues are -1 and 3. We can plug these values of λ into our equations to solve for the eigenvectors:

```
Solve[{-x == x + 2 y, - y == 2 x + y}, {x,y}]
Solve[{3 x == x + 2 y, 3 y == 2 x + y}, {x,y}]
```

Solve::svars: Equations may not give solutions for all "solve" variables.

```
{{x -> y}}
```

Solve::svars: Equations may not give solutions for all "solve" variables.

```
{{x -> y}}
```

■ Reduce

Mathematica is smart enough that we can use **Reduce[]** to do it all in one line:

```
Reduce[{lambda x v == W.xv}, {x,y,lambda}]
```

```
x == -y & lambda == -1 & x == y & lambda == 3 & lambda - 3 != 0 & lambda + 1 != 0 & x == 0 & y == 0
```

The eigenvectors are unique only up to a scale factor, so one can choose how to normalize them. For example, we could arbitrarily set x to 1, and then the eigenvectors are: {1,1}, and {1,-1}. Alternatively, we could normalize them to {1/Sqrt[2], 1/Sqrt[2]} and {1/Sqrt[2], -1/Sqrt[2]}.

```
Eigenvectors[W]
```

```
 $\begin{pmatrix} -1 & 1 \\ 1 & 1 \end{pmatrix}$ 
```

■ Side note: Solve[] vs. Reduce[]

Solve[] makes assumptions about constraints left unspecified, so the following returns the solution true for any lambda:

```
Solve[{lambda x == x + 2 y, lambda y == 2 x + y}, {x,y,lambda}]
```

Solve::svars: Equations may not give solutions for all "solve" variables.

```
{{x -> 0, y -> 0}}
```

```
Solve[{lambda x == x + 2 y, lambda y == 2 x + y,
lambda != 0, x != 0, y != 0}, {x,y,lambda}]
```

Solve::svars: Equations may not give solutions for all "solve" variables.

```
{{x -> -y, lambda -> 1}, {x -> 0, lambda -> 3}}
```

Reduce[] gives all the possibilities without making specific assumptions about the parameters:

Either of the following forms will work too:

```
Reduce[lambda xv == W.xv, {xv[[1]],xv[[2]],lambda}]
Reduce[{{lambda x, lambda y} == {x + 2 y, 2 x + y}},
{x,y,lambda}]
```

```
x == -y & lambda == -1 || x == y & lambda == 3 || lambda - 3 != 0 & lambda + 1 != 0 & x == 0 & y == 0
```

```
x == -y & lambda == -1 || x == y & lambda == 3 || lambda - 3 != 0 & lambda + 1 != 0 & x == 0 & y == 0
```

References

Campbell, F. W., & Robson, J. R. (1968). Application of Fourier Analysis to the Visibility of Gratings. *Journal of Physiology*, 197, 551-566.

Enroth-Cugell, C., & Robson, J. G. (1966). The contrast sensitivity of retinal ganglion cells of the cat, *Journal of Physiology, London*, 187, 517-552.

© 1998 Daniel Kersten, Computational Vision Lab, Department of Psychology, University of Minnesota.