

Introduction to Neural Networks

U. Minn. Psy 5038

Daniel Kersten

Lecture 4

Introduction

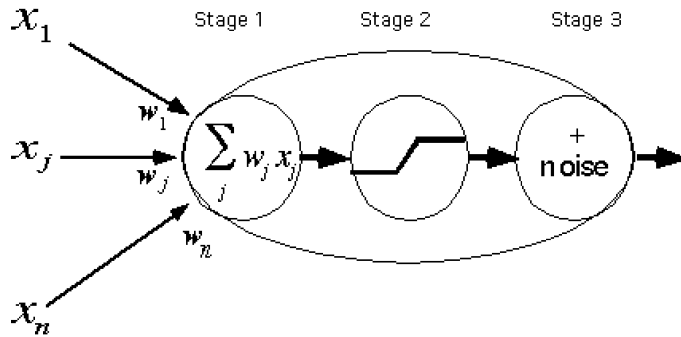
Last time

- Slow-potential qualitative neuron model.
- Various types of neuron models: Levels of abstraction
- McCulloch-Pitts threshold logic: Discrete time, discrete state, no spatial structure

Today

- We develop a "structure-less, continuous signal, and discrete time" generic neuron model and from there build a network. This "connectionist" model is one of several abstractions that we saw in the previous lecture.
- Review basic linear algebra. Motivate linear algebra concepts from neural networks.

The generic neuron model



The generic neuron model abstracts the basic properties of the integrate and fire neuron, and makes provision for saturation as well.

Stage 1:

Linear weighted sum of inputs

fixed bias term ()

The weights correspond to the synaptic efficiency of the inputs to the neuron which model the net effect on the input current. The bias term models a threshold.

Stage 2:

non-linearity

Popular forms are: logistic function, arctan(), limit function

A point non-linearity models both the effects of small signal compression (e.g. threshold) and large signal saturation on the output frequency of firing.

(Stage 3:)

noise

The number of spikes in a neuron's discharge is not strictly determined by the input, but varies statistically. This can be modeled assuming some form of additive (or other) stochastic component to the neural discharge frequency.

Much of our intuition and theory about neural networks is based on studies of *linear neural networks*. Linear neural networks don't include the non-linearity of stage 2, and their behavior is determined by our knowledge of linear algebra.

Now we will develop some *Mathematica* tools to model Stage 1 and 2 of the generic connectionist model of the neuron.

- **Defining functions.** Let's define a function to model the non-linearity (in this case, the "logistic function" mentioned earlier):

```
squash[x_] := N[1/(1 + Exp[-x])];
```

Recall, that the underscore, **x_** is **important** because it tells Mathematica that x represents a slot, not an expression. Note that we've used a colon followed by equals (**:=**) instead of just an equals sign (=). When you use an equals sign, the value is calculated immediately. When there is a colon in front of the equals, the value is calculated only when called on later.

So here we use `:=` because we need to define the function for later use. Also note that our squashing function was defined with `N[]`. *Mathematica* tries to keep everything exact as long as possible and thus will try to do symbol manipulation if we don't explicitly tell it that we want numbers.

Graphics. Adjust the input scale of `squash[]` to plot a very steep squashing function for $-5 < x < 5$. I.e. it should look like the step function we used when modeling the McCulloch-Pitts neuron.

- **Lists.** We already introduced lists when we studied the McCulloch-Pitts model. In this course, we are going to do a lot of work with lists, in particular with vectors (a vector is a list of scalar elements) and matrices (a matrix is a list of vector elements). Here is a four-dimensional vector which we'll call \mathbf{x} . \mathbf{x} could represent the input signals to a neuron.

```
x = {2, 3, 0, 1};
```

As we discovered earlier, by ending a line with a semi-colon, you suppress the output after hitting the return key. Now let's make another vector, this one will be a list of "weights", say, representing the efficiency with which the inputs at the synapses are transmitted to the neuron hillock (we'll allow negative weights for the time being):

```
w = {2, 1, -2, 3};
```

■ Model linear neuron

Now the output of a model neuron that simply takes a weighted sum of the inputs is just the *dot product* of the input with the weights:

```
y = w . x
```

```
10
```

This kind of operation is sometimes referred to as a "cross-correlator". It takes a signal \mathbf{x} , and cross-correlates it with a template, \mathbf{w} . Later on in an exercise you will show that for signals of fixed vector length (or "norm"), the cross-correlator gives the biggest response to the signal that exactly matches the template.

We can see what the dot product does algebraically by defining the input and weights algebraically:

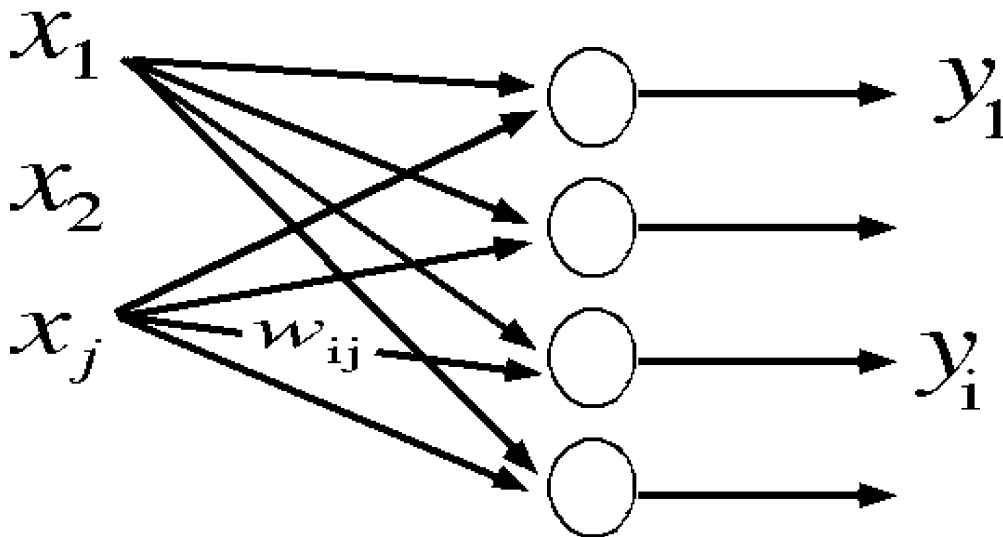
```
Clear[w1, w2, w3, w4, x1, x2, x3, x4];  
y = {w1, w2, w3, w4} . {x1, x2, x3, x4}
```

```
w1 x1 + w2 x2 + w3 x3 + w4 x4
```

Now let's include the non-linear squashing function to complete our model of the first two stages of the generic neuron:

```
y = squash[w.x];
```

Modeling a simple neural network



What if the input is applied to four neurons, each with a different set of weights? We can represent the weights by a "weight matrix", which is just a list of the four weight lists or vectors. Here is a 4x4 matrix W :

```
W = {{2, 1, -2, 3}, {3, 1, -2, 2}, {4, 6, 5, -3}, {1, -2, 2, 1}}
```

$$\begin{pmatrix} 2 & 1 & -2 & 3 \\ 3 & 1 & -2 & 2 \\ 4 & 6 & 5 & -3 \\ 1 & -2 & 2 & 1 \end{pmatrix}$$

Each successive element of the list W is a *row* of matrix W . Verify this by displaying W in `MatrixForm`

Now what are the outputs of the four neurons? It is just the product of the matrix W times the input vector x :

```
y = W.x
```

```
{10, 11, 23, -3}
```

In traditional form, this matrix multiplication is written as:

$$y_i = \sum w_{i,j} x_j$$

```
Sum[x_j w_{i,j}]
```

So to multiply an input vector by a matrix, we take the dot product of the input vector with each successive *row* of the matrix.

Note that in *Mathematica*, a dot is used for multiplying vectors by themselves, vectors by a matrix, or to multiply two matrices together. If you want to multiply a vector or matrix by a scalar, *c*, you don't use a dot. For example, to normalize *x* by its vector length:

```
c = 1/Sqrt[x.x];  
x2 = N[c x]
```

```
{0.534522, 0.801784, 0., 0.267261}
```

Take the dot product of x2 above with itself to confirm normalization

Now let's apply our squashing function to the output *y*. Note how the big positive values are set close to one, and the negative value is set close to zero.

```
y  
squash[y]
```

```
{10, 11, 23, -3}
```

```
{0.999955, 0.999983, 1., 0.0474259}
```

By default, our function `squash[]` is a **listable** function. This means that even though it was first defined to operate on a scalar, when applied to a list, it automatically gets applied to each element of the list in turn.

We can do everything at once in our four-neuron network, producing the four outputs of four generic neurons to an input *x*:

```
y = squash[W.x]
```

```
{0.999955, 0.999983, 1., 0.0474259}
```

There we have it--a model for a simple neural network. It is also "feedforward" in that it maps inputs to outputs, in contrast to a network in which outputs get fed back to inputs.

This equation will occur many times in the rest of the course, so it is worth taking some time to understand it.

Our example has four inputs, and four outputs. Try making a graphical sketch of the net to illustrate what is connected to what, label the inputs x_j ; the weights w_{ij} , and the outputs y_i .

You can access the components of vectors. For example here is the second element of y , and the element in the second row, third column of \mathbf{W} :

```
y[[2]]
```

```
0.999983
```

```
w[[2,3]]
```

```
-2
```

```
w[[2]][[3]]
```

```
-2
```

```
w[[1]]
```

```
{2, 1, -2, 3}
```

Modeling noise (Stage 3): Generic neuron plus noise

We'd like to add a Stage 3 to our model of the neuron in which we take account of the noisiness of neural transmission. For this, we need the notion of a *probability distribution*. We could develop the routines we need using basic *Mathematica* functions. However, much of the work is built into *Mathematica*. As a first approximation the maintained action potential discharge can be modeled as a Poisson distribution.

Statistics and stochastic processes

Statistical routines are useful for both theoretical aspects of modeling as well as for Monte Carlo simulations where one simulates drawing random samples to use as inputs to an algorithm. So it is worth a little effort to get acquainted with some fundamental tools and definitions. Let's define Poisson distribution with a mean of λ . And then specify $\lambda=50$ (e.g. 50 spikes per second of a neuron).

Discrete distributions

```
Clear[a];
PDF[PoissonDistribution[λ], a]
```

$$\frac{e^{-\lambda} \lambda^a}{a!}$$

Let's specify a Poisson distribution with mean $\lambda = 50$:

```
pdist = PoissonDistribution[50];
```

The probability distribution function is given by:

```
PDF[pdist, a]
```

$$\frac{50^a}{e^{50} a!}$$

The output shows *Mathematica's* definition of the function. You can obtain the mean, variance and standard deviation (which is the square root of the variance) of the distribution we've defined. Try it:

```
Mean[pdist]
Variance[pdist]
StandardDeviation[pdist]
```

```
50
```

```
50
```

```
 $5\sqrt{2}$ 
```

What is your guess of the general relationship between the mean and variance for the Poisson distribution?

We are going to approximate the noisiness of neural discharge with a Normal or Gaussian distribution. The Gaussian distribution is continuous, rather than discrete. It is a fairly good approximation of a Poisson distribution for large values of the mean, and the theory is usually easier. We do have to be careful tho', because the Gaussian is defined over negative numbers too.

Continuous distributions, probability densities

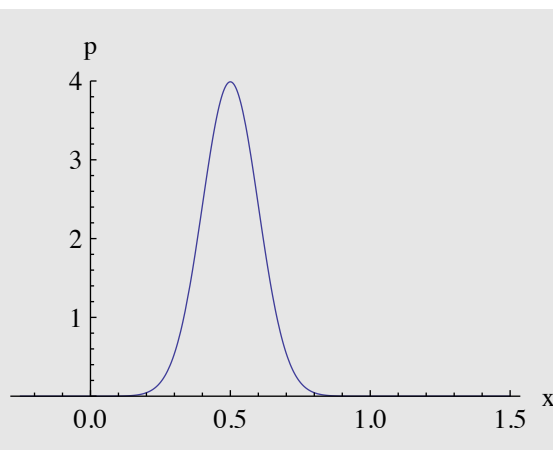
```
ndist = NormalDistribution[0.5,.1];
```

```
Print[Mean[ndist],", ",Variance[ndist],", ",  
StandardDeviation[ndist]]
```

0.5,0.01,0.1

A plot of the probability distribution function for this normal distribution looks like:

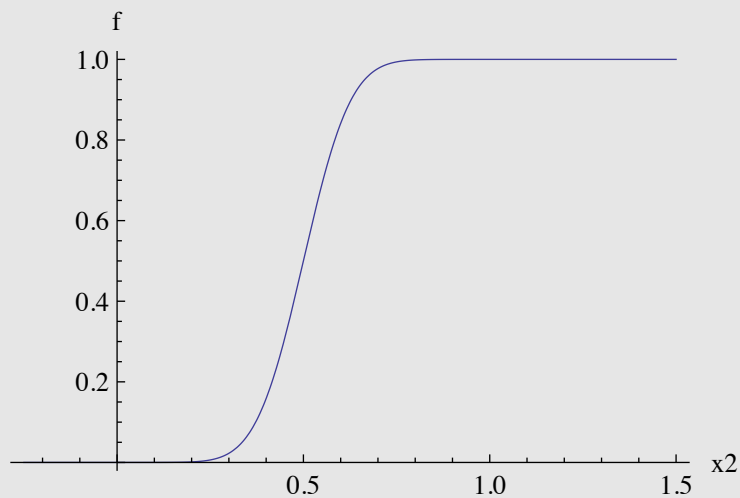
```
Plot[PDF[ndist, x], {x, -0.25`, 1.5`}, PlotRange -> {0, 4},  
AxesLabel -> {"x", "p"}]
```



A given ordinate value is a "*density*", rather than probability. But we can talk about the probability that x takes on some value in an interval, say dx . For a small interval, dx , the probability $\approx p(x)dx$. What is the probability that x takes on some value between $+\infty$ and $-\infty$? What is area under this curve?

The *cumulative distribution*, $f(x) = P(x < x_2)$, tells us the probability of x being less than a particular value of x_2 :


```
Plot[CDF[ndist, x2], {x2, -0.25, 1.5}, AxesLabel -> {"x2", "f"}]
```



You can see from the graph that for this distribution, once x_2 is greater than 0.7 or so, the probability of x being less than that is virtually certain, i.e. is essentially 1. If we set the mean = 0, and the standard deviation, we'd have a graph of the "cumulative normal".

■ Statistical Sampling

Having defined the normal distribution, how can we draw samples from it? In other words, can we simulate a process in which we fill a hat with slips of paper in such a way that the proportions for each value mimic what we obtain from a theoretical distribution?

Most standard programming languages come with subroutines for doing pseudo-random number generation. Unlike the Poisson or Gaussian distribution, these numbers are usually **uniformly distributed**--that is, the probability of being a certain value (or within a tiny range) is constant over the entire sampling range of the random variable.

This is like filling the hat with slips of paper where the number of slips is the same for each value.

As we noted earlier, *Mathematica* comes with standard functions, **RandomReal[]** and **RandomInteger[]**, that enables us to generate random numbers that are uniformly distributed. With the appropriate argument, we can also define Poisson, Normal, and other kinds of random numbers.

```
RandomReal[ndist]
```

```
0.517909
```

Simulate drawing Poisson distributed integers with a mean value of 50 (See RandomInteger[])

Putting together stages 1, 2 and 3 together

We can do everything at once, producing the output of a generic neuron, with synaptic weights w , neural noise with a mean of 0.0 and std. dev. 0.1 to an input x :

```
Clear[x1,w,y];  
w = {2,1,-2,3};  
ndist2 = NormalDistribution[0.0,.1];  
y[x1_] := N[squash[w.x1] + RandomReal[ndist2]];
```

```
y[{2, 3, 0, 1}]
```

```
0.870993
```

Note that this is a model with what is called *additive noise*. The noise level doesn't depend on the activity value of the deterministic part. Although often a good first approximation and starting point for a model, additive noise isn't necessarily true, and real neurons often depart from this model.

If we invoke the `y[]` function again, we get a different response:

```
y[{2,3,0,1}]
```

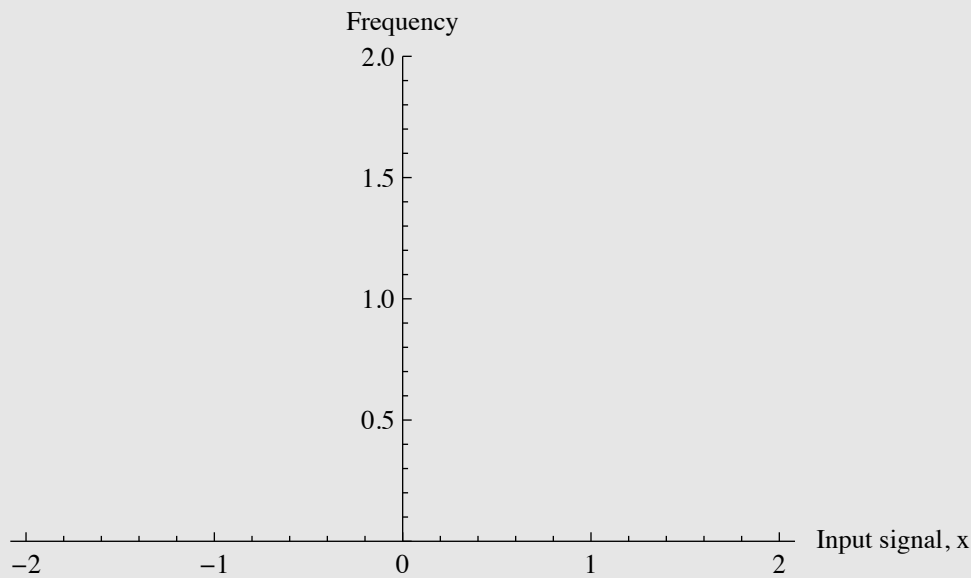
```
0.860516
```

To sum up, the model you should have in mind is that at any given time interval (which is implicit in this continuous-response, discrete-time model), the neuron computes the sum of its weighted inputs, and the output signal, y , is a spike rate over this interval. With a sigmoidal non-linearity, there is small-signal suppression, and large-signal saturation.

Exercise

Suppose all the inputs except the first are clamped at zero. What does the response, y look like as a function of x for various levels of input signal? Fill in the argument for `Plot[]`:

```
Plot[Null, {x, -2, 2}, PlotRange -> {0, 2},
  AxesLabel -> {"Input signal, x", "Frequency"}]
```

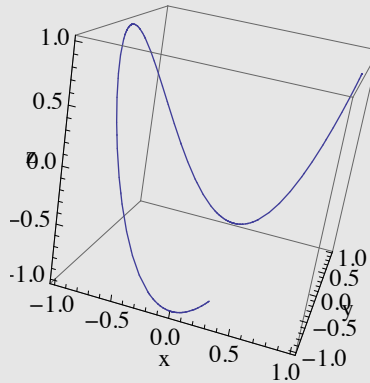


Vector operations and patterns of neural activity

■ State space and state vectors.

In neural networks, we are often concerned with a vector whose components represent the activities of neurons which are changing in time. So sometimes we will talk about state vectors. There isn't anything profound about this terminology—it just reflects that we are interested in the value of the vector when the system is in a particular state at time t . It is often very useful to think of an n -dimensional vector as a point in an n -dimensional space. This space is often referred to as state space. Suppose, we have a 3 neuron system. We can describe the state of this system as a 3-dimensional vector where each component represents the activity of the neuron. Further, suppose just for the sake of an example to visualize, the activities of the first, second, and third neurons (i.e. components) of a 3-dimensional vector are given by: $y = \{\text{Cos}[t], \text{Sin}[t], t\}$. We can use the Mathematica function, `ParametricPlot3D[]` to get a picture of how this state vector evolves in time through state space:

```
Clear[y];  
y[t_] := {Cos[t], Sin[t], Cos[2 t]};  
ParametricPlot3D[y[t], {t,0,5}, AxesLabel->{"x", "y", "z"}]
```



■ Dimension of a vector.

```
v = {2.1, 3, -0.45, 4.9};
```

Dimensions[], will give you the dimensions of a matrix, while **Length[]** tells you the number of elements in the list. For example,

```
M = {{2,4,2}, {1,6,4}};
```

```
Length[M]
```

```
2
```

Use `Dimensions[]` to print the dimensions of `v`

Compare `Length[M]` with `Dimensions[M]`. Compare `Length[v]` with `Dimensions[v]`.

■ Transpose of a vector.

The transpose of a column vector is just the same vector arranged in a row. However, because of the way Mathematica uses lists to represent vectors you don't have to distinguish between row and column vectors. In standard math notation, transpose of a vector \mathbf{x} , is often written \mathbf{x}^T . You can see a vector in column form by typing `v//MatrixForm`, or:

```
MatrixForm[v]
```

```

$$\begin{pmatrix} 2.1 \\ 3 \\ -0.45 \\ 4.9 \end{pmatrix}$$

```

■ Vector addition

is accomplished by simply adding the components of each vector to make a new vector.

Note that the vectors all have the same dimension.

```
a = {3, 1, 2};  
b = {2, 4, 8};  
c = a + b
```

```
{5, 5, 10}
```

Vectors can be multiplied by a constant. We saw an example of this earlier.

```
2 a
```

```
{6, 2, 4}
```

■ Euclidean "length" of a vector

It is unfortunate terminology, but **Length[]** does NOT give you the metrical or Euclidean length of the vector, which is the Euclidean distance from the origin to the end of the vector. But **Norm[]** does. To get the length of a vector, you calculate the Euclidean distance from the origin to the end-point of the vector. squaring each component, adding up the squares, and taking the square root.

```
Remove[x1, x2, x3]
Norm[{x1, x2, x3}]
```

$$\sqrt{|x1|^2 + |x2|^2 + |x3|^2}$$

To get a little more practice with *Mathematica*, you can also do this with the **Apply[]** function, where the **Plus** operation is applied to all the elements of the list. Note that the operation of exponentiation (raising to the power of 2) is "listable", that is it is applied to each element of the vector:

```
{x1, x2, x3}^2
Sqrt[Apply[Plus, {x1, x2, x3}^2]]
```

$$\{x1^2, x2^2, x3^2\}$$

$$\sqrt{x1^2 + x2^2 + x3^2}$$

Norm[expr] generalizes to **Norm[expr, p]**. The second argument says that you want the vector 2-norm (ie. Euclidean length). **Norm[x,1]** would return the "city-block" norm.

```
Norm[u, 1]
```

$$\|u\|_1$$

Compare: $\{x1, x2, x3\} \{x1, x2, x3\}$, $\{x1, x2, x3\} * \{x1, x2, x3\}$, $\{x1, x2, x3\}^2$, $\{x1, x2, x3\} . \{x1, x2, x3\}$

The norm or vector length of a vector **a** is often written as **||a||** in standard math notation. In the next section, we use the inner or dot product to calculate the Euclidean length of a vector.

- **Dot or Inner product.** To calculate the inner product of two vectors, you multiply the corresponding components and add them up:

```
Clear[u1, u2, u3, u4, v1, v2, v3, v4];
u = {u1, u2, u3, u4};
v = {v1, v2, v3, v4};
u.v
```

```
u1 v1 + u2 v2 + u3 v3 + u4 v4
```

The **inner product** is also called the **dot product**. Later we will see what is meant by **outer product**. The inner product between two vectors **a** and **b** is traditionally written either as:

$$\mathbf{a} \cdot \mathbf{b} \text{ or } [\mathbf{a}, \mathbf{b}], \text{ or } \mathbf{a}^T \mathbf{b}$$

Mathematica uses the dot notation.

One use of the inner product is to calculate the length of a vector. $\mathbf{a} \cdot \mathbf{a}$ is just the sum of the squares of the elements of **a**, so gives us another way of calculating the length of a vector.

```
N[Sqrt[a.a]]
```

```
3.74166
```

```
Sqrt[u.u]
```

$$\sqrt{u_1^2 + u_2^2 + u_3^2 + u_4^2}$$

Define your own function that will return the L2-norm (regular Euclidean length of a vector) **x: Vector-length[x_] := N[Sqrt[x.x]]**

■ Projection

Projection is an important concept in linear neural networks.

When a pattern of activity, **x**, is input to a linear neural network, the weight matrix **W** transforms the input pattern to a new output pattern **y** of activities. This linear transformation works by "projecting" the input onto a new set of dimensions by taking the dot product of the input with each *row* of the weight matrix.

Suppose we have 3 inputs and 2 outputs to our network. Inputs to the network live in a 3-dimensional space. Outputs live in 2 dimensions.

With a 3-dimensional vector that inputs into 2 neurons, one can visualize the output as a 2-dimensional vector whose length and direction is determined by the 2-dimensional vector sum of three column weight vectors:

$\begin{pmatrix} w_{11} \\ w_{21} \end{pmatrix}$, $\begin{pmatrix} w_{12} \\ w_{22} \end{pmatrix}$, $\begin{pmatrix} w_{13} \\ w_{23} \end{pmatrix}$, where the amplitude of each vector is scaled by the input activity levels x_1 , x_2 , x_3 , to give: $\begin{pmatrix} w_{11} \\ w_{21} \end{pmatrix} x_1 + \begin{pmatrix} w_{12} \\ w_{22} \end{pmatrix} x_2 + \begin{pmatrix} w_{13} \\ w_{23} \end{pmatrix} x_3$. Models of image representation in the primary visual cortex have been analyzed in terms of whether the high-dimensional images received (at the retina) are projected into lower or higher dimensional spaces in cortex, and what the consequences might be for biological image processing.

Prove :
$$\begin{pmatrix} w_{11} & w_{12} & w_{13} \\ w_{21} & w_{22} & w_{23} \end{pmatrix} \cdot \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} = \begin{pmatrix} w_{11} \\ w_{21} \end{pmatrix} x_1 + \begin{pmatrix} w_{12} \\ w_{22} \end{pmatrix} x_2 + \begin{pmatrix} w_{13} \\ w_{23} \end{pmatrix} x_3$$

In problem set 1, you calculate the output of a linear neuron model as the dot product between an input vector and a weight vector. Both the weight and input lists can be thought of as vectors in an n-dimensional space. Suppose the weight vector has unit length. Recall that you can normalize any vector to unit length by dividing by its length:

$$v_n = v / \text{Sqrt}[v \cdot v];$$

or using the built-in function `Normalize[]`:

$$v_n = \text{Normalize}[v]$$

$$\left\{ \frac{v_1}{\sqrt{|v_1|^2 + |v_2|^2 + |v_3|^2 + |v_4|^2}}, \frac{v_2}{\sqrt{|v_1|^2 + |v_2|^2 + |v_3|^2 + |v_4|^2}}, \frac{v_3}{\sqrt{|v_1|^2 + |v_2|^2 + |v_3|^2 + |v_4|^2}}, \frac{v_4}{\sqrt{|v_1|^2 + |v_2|^2 + |v_3|^2 + |v_4|^2}} \right\}$$

The dot product, $\mathbf{a} \cdot \mathbf{b}$, is equal to:

$$|\mathbf{a}| |\mathbf{b}| \cos(\text{angle between } \mathbf{a} \text{ and } \mathbf{b})$$

Geometrically, we can think of the output of a neuron as the projection of the activity of the neuron input activity vector onto the weight vector direction. Suppose the input vector is already perpendicular to the weight vector, then the output of the neuron is zero, because the cosine of 90 degrees is zero. As you found or will find with the cross-correlator of Problem Set 1, the further the input pattern is away from the weight vector, as measured by the cosine between them, the poorer the "match" between input and weight vectors, and the lower the response.

Consider the simple case of a 2-dimensional input onto one neuron. The output lives in a 1-dimensional space, i.e. a line pointing in the direction of the weight vector \mathbf{w} . Here are three lines of code that calculate the two-dimensional vector \mathbf{z} in the direction of \mathbf{w} , with a length determined by "how much of \mathbf{x} projects in the \mathbf{w} direction". The vector \mathbf{z} , represents the output activity level of the neuron.


```
x = .85*{1,2};
w = N[{2/Sqrt[5],1/Sqrt[5]}];
z = (x.w) w
```

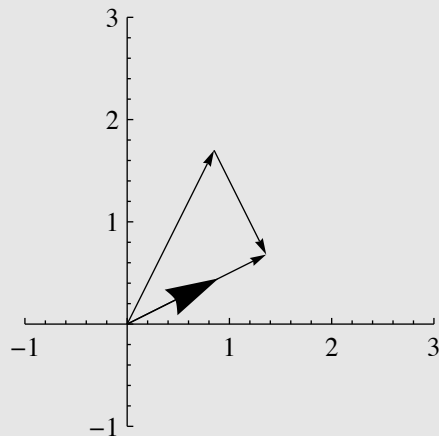
```
{1.36, 0.68}
```

Projection is a fundamental concept, and later versions of *Mathematica* provide a function for it:

```
z = Projection[x, w]
```

```
{1.36, 0.68}
```

```
Show[
Graphics[{Tooltip[Arrow[{{0, 0}, x}], "x"],
  Tooltip[Arrow[{{0, 0}, z}], "z"],
  Tooltip[Arrowheads[Large], Arrow[{{0, 0}, w}], "w"],
  Tooltip[Arrow[{x, z}], "z-x"]}], PlotRange -> {{-1, 3}, {-1, 3}},
Axes -> True, AspectRatio -> 1]
```



Try moving your mouse over the arrows above. **Tooltips[]** provides rollover labels for x, w, z, and z-x on the graph.

■ Angle between two vectors and orthogonality: Similarity measure between patterns

Often we will want some measure of the similarity between two patterns of neural firings. As we have just seen, one measure of comparison is the degree to which the two state vectors point in the same direction. The cosine of the angle between two vectors is one possible measure:

```
cosine[x_,y_] := x.y/(Norm[x] Norm[y])
```

Note that if two vectors point in the same direction, the angle between them is zero, and the cosine of the angle is 1:

```

a = {2,1,3,6};
b = {6, 3, 9, 18};
VectorAngle[a, b]
cosine[a,b]

```

0

1

Try verifying that w and z from the previous section point in the same direction.

If two vectors point in the opposite directions, the cosine of the angle between them is -1:

```

a = {-2,-1,-3,-6};
b = {6, 3, 9, 18};
VectorAngle[a, b]
cosine[a,b]

```

 π

-1

And if the vectors **a2**, **b2**, are orthogonal, then:

```

a = {-2, -1, -3, -6};
b = {6, 3, 9, 12};
{a2, b2} = Orthogonalize[{a, b}];
cosine[a2, b2]
VectorAngle[a2, b2]

```

0

 $\frac{\pi}{2}$

■ Euclidean distance between two vectors

Two vectors may point in the same direction, but could be quite different because they have different lengths. Another measure of similarity is the Euclidean length of the difference between two vectors, or the "distance between the tips of their vectors":

```
N[Norm[a - b]]
```

```
23.4094
```

Or you can use:

```
EuclideanDistance[a, b]
```

```
2  $\sqrt{137}$ 
```

By thinking about the geometry, what is the Norm of $\{3,0\}-\{0,4\}$?

- **Orthogonality.** The case where vectors are at right angles to each other is an important special case that is worth spending a little time on. Consider an 8-dimensional space. One very familiar set of orthogonal vectors is the following:

```
u1 = {1,0,0,0,0,0,0,0};
u2 = {0,1,0,0,0,0,0,0};
u3 = {0,0,1,0,0,0,0,0};
u4 = {0,0,0,1,0,0,0,0};
u5 = {0,0,0,0,1,0,0,0};
u6 = {0,0,0,0,0,1,0,0};
u7 = {0,0,0,0,0,0,1,0};
u8 = {0,0,0,0,0,0,0,1};
```

Each vector has unit length, and it is easy to see just by inspection that the inner product between any two is zero. On the other hand, here is another set of 8 vectors in 8-space for which it is not immediately obvious that they are all orthogonal. Here is a set of what are called Walsh functions:

```
v1 = {1, 1, 1, 1, 1, 1, 1, 1};
v2 = {1,-1,-1, 1, 1,-1,-1, 1};
v3 = {1, 1,-1,-1,-1,-1, 1, 1};
v4 = {1,-1, 1,-1,-1, 1,-1, 1};
v5 = {1, 1, 1, 1,-1,-1,-1,-1};
v6 = {1,-1,-1, 1,-1, 1, 1,-1};
v7 = {1, 1,-1,-1, 1, 1,-1,-1};
v8 = {1,-1, 1,-1, 1,-1, 1,-1};
```

The above example is just one of an infinite number of possible orthogonal sets.

You can calculate the inner products between any two Walsh vectors, and you will find out that they are all zero. Note that with the first set of vectors, $\{u_i\}$, you can tell which vector it is just by looking for where the 1 is. For the second set, $\{v_i\}$, you can't tell by looking at just one component. For example, the first component of all of the Walsh functions has a 1. You have to look at the pattern to tell which Walsh function you are looking at.

See example section in the `UnitStep[]` *Mathematica* reference for general code for Walsh functions defined on a continuum:

```
Walsh[{n_, k_}, y_] :=
Module[
  {li =
    Split[
      Extract[
        Nest[
          Sequence@@@ {Join[Riffle[#, #] & /@#,
            Join[Riffle[#, -#] & /@Reverse[#]]} &, {{1}}, n, {k}}},
-1 -
    2 Total[MapIndexed[(-1) ^ #2 UnitStep[Mod[y, 2^n] - #1] &,
      Most[FoldList[Plus, 0, Length /@ li]], Infinity]]
```

■ Grandmother cell coding

Let's pursue this further. I have 8 neurons that can fire in response to me viewing one of my relatives. Let's assign specific meanings to each of the patterns--each pattern is a code for some thing, like "grandma Tompkins", "grandma Wilke", "uncle Heine", "aunt Mabel", and so forth. If we use the **u**'s, then to decide who I'm looking at (to "read my mind"), one could look for the one neuron that lights up to find out which relative it is representing--then the neuron activity represented, for example, by the third element of the pattern could mean "grandma Wilke". The third unit $u_{3[3]}$ is 1 if and only if it is grandma Wilke.

This strategy wouldn't work if we encoded a collection of relatives using the **v**'s. Suppose relatives are represented by coding in the Walsh set, and that grandma Wilke is represented by the third pattern of activity v_3 , then although grandma Wilke implies that $v_{3[3]}$ is -1 (or in general that $v_{3[i]}$ is a particular value for any of the *i*s), the reverse isn't true--the activity level of any single neuron does not uniquely specify which relative I'm looking at.

The **v**'s give us a simple example of what is sometimes referred to as a **distributed code**. The **u**'s are examples of a **grandmother cell code**. The reason for this obscure terminology can be traced to debates on whether there may be single cells in the brain whose firing uniquely determines the recognition of one's grandmother.

- **Orthonormality.** The Walsh set is orthogonal, but they are not of unit length. We have already seen some of the advantages of working with unit length vectors. The general issue of normalization comes up all the time in neural networks both in terms of limiting overall neural activity, and limiting synaptic weights. So it is sometimes convenient to normalize an orthogonal set, producing what is known as an *orthonormal* set of vectors:

```
w1 = v1/Norm[v1];
w2 = v2/Norm[v2];
w3 = v3/Norm[v3];
w4 = v4/Norm[v4];
w5 = v5/Norm[v5];
w6 = v6/Norm[v6];
w7 = v7/Norm[v7];
w8 = v8/Norm[v8];
```

Vector representations, linear algebra

The issue of how information is to be represented is fundamental in the information sciences generally, as well as for neural network theory. A pattern of activity over a set of neurons is presumed to mean something, and there are different ways of coding the same meaning. But different codes have different properties. A code may not be sufficient to uniquely code all the possible things we need to represent. A code could be redundant and have more than one way of representing the same thing. This section continues with our review of the basics of vector and linear algebra by going a little more deeply into the subject. The pay-off will be some mathematics that provides intuition about issues of neural representation. You can think of this as a first lesson in the "psychology of linear algebra".

■ Basis sets

It is pretty clear that given any vector whatsoever in 8-space, you can specify how much of it gets projected in each of the eight directions specified by the unit vectors v_1, v_2, \dots, v_8 . But you can also build back up an arbitrary vector by adding up all the contributions from each of the component vectors. This is a consequence of vector addition and can be easily seen to be true in 2 dimensions. We can verify it ourselves. Pick an arbitrary vector \mathbf{g} , project it onto each of the basis vectors, and then add them back up again:

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

```
 $(\mathbf{g} \cdot \mathbf{u}_1) \mathbf{u}_1 + (\mathbf{g} \cdot \mathbf{u}_2) \mathbf{u}_2 + (\mathbf{g} \cdot \mathbf{u}_3) \mathbf{u}_3 + (\mathbf{g} \cdot \mathbf{u}_4) \mathbf{u}_4 +$   
 $(\mathbf{g} \cdot \mathbf{u}_5) \mathbf{u}_5 + (\mathbf{g} \cdot \mathbf{u}_6) \mathbf{u}_6 + (\mathbf{g} \cdot \mathbf{u}_7) \mathbf{u}_7 + (\mathbf{g} \cdot \mathbf{u}_8) \mathbf{u}_8$ 
```

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

Exercise

What happens if you project \mathbf{g} onto the normalized Walsh basis set defined by $\{\mathbf{w}_1, \mathbf{w}_2, \dots\}$ above, and then add up all 8 components?


```
Simplify[%]
```

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

The projections, $\mathbf{g} \cdot \mathbf{w}_i$ are sometimes called the **spectrum** of \mathbf{g} . This terminology comes from the Fourier basis set used in Fourier analysis. A discrete version of a Fourier basis set is similar to the Walsh set, except that the elements fit a sine wave pattern, and so are not binary-valued.

The orthonormal set of vectors we've defined above is said to be **complete**, because *any* vector in 8-space can be expressed as a linear weighted sum of these **basis vectors**. The weights are just the projections. If we had only 7 vectors in our set, then we would not be able to express all 8-dimensional vectors in terms of this basis set. The seven vector set would be said to be **incomplete**. A basis set which is orthonormal and complete is very nice from a mathematical point of view. Another bit of terminology is that these seven vectors would not **span** the 8-dimensional space. But they would span some sub-space, that is of smaller dimension, of the 8-space.

As mentioned above, there has been much interest in describing the effective weighting properties of visual neurons in primary visual cortex of higher level mammals (cats, monkeys) in terms of basis vectors. One issue is if the input (e.g. an image) is projected (via a collection of receptive fields) onto a set of neurons, is information lost? If the set of weights representing the receptive fields of the collection of neurons is complete, then no information is lost.

■ Linear dependence

What if we had 9 vectors in our basis set used to represent vectors in 8-space? For the u 's, it is easy to see that in a sense we have too many, because we could express the 9th in terms of a sum of the others. This set of nine vectors would be said to be linearly dependent. A set of vectors is linearly dependent if one or more of them can be expressed as a linear combination of some of the others. Sometimes there is an advantage to having an "over-complete" basis set (e.g. more than 8 vectors for 8-space; cf. Simoncelli et al., 1992).

Theorem: A set of mutually orthogonal vectors is linearly independent.

However, note it is quite possible to have a linearly independent set of vectors which are not orthogonal to each other. Imagine 3-space and 3 vectors which do not jointly lie on a plane. This set is linearly independent.

If we have a linearly independent set, say of 8 vectors for our 8-space, then no member can be dropped without a loss in the dimensionality of the space spanned.

It is useful to think about the meaning of linear independence in terms of geometry. A set of three linearly independent vectors can completely span 3-space. So any vector in 3-space can be represented as a weighted sum of these 3. If one of the members in our set of three can be expressed in terms of the other two, the set is not linearly independent and the set only spans a 2-dimensional subspace. That is, the set can only represent vectors which lay on a plane in 3-space. This can be easily seen to be true for the set of u 's, but is also true for the set of v 's.

Thought exercise

Suppose there are three inputs feeding into three neurons in the simple linear network such as defined at the beginning of this lecture. If the weight vectors of the three neurons are not linearly independent, do we lose information?

Linearity, real neural networks, and what's up next time?

From a computational standpoint, the squashing function has both advantages and disadvantages. It is what makes our neural network model *non-linear*, and as we will see later, this non-linearity enables networks to compute functions that can't be computed with a linear network. On the other hand, non-linearities make the analysis complicated because we leave the well-understood domain of linear algebra. In fact, there are cases for which most of the neural activities are in the mid-range of the squashing function, and here one can approximate the network as a purely linear one--just matrix operations on vector inputs, and the analysis becomes relatively simple.

Compared to the complexity of real neurons and networks, assuming linearity might seem to be just too simple. But we will see in the next lecture, that a linear model can be quite good model for some biological subsystems. We will apply the techniques of linear vector algebra to model a network discovered in the visual system of the horseshoe crab. Later we'll see how some aspects of associative memory can be modeled using linear systems.

References

A mathematica-based linear algebra course. <http://library.wolfram.com/infocenter/MathSource/4611/>

Olver, Peter J. and Shakiban, Chehrzad (2005) Applied Linear Algebra, Prentice-Hall, Upper Saddle River, N.J., 2005
<http://www.math.umn.edu/~olver/ala.html>

Simoncelli, E. P., Freeman, W. T., Adelson, E. H., & Heeger, D. J. (1992). Shiftable Multi-scale Transforms. IEEE Trans. Information Theory, 38(2), 587--607.

Strang, G. (1988). Linear Algebra and Its Applications (3rd ed.). Saunders College Publishing Harcourt Brace Jovanovich College Publishers.

Also, take a look at Strang's lecture notes on video:[http : web.mit.edu/18.06/www/Video/video – fall – 99. html](http://web.mit.edu/18.06/www/Video/video-fall-99.html)

(Also available through iTunesU --- Start iTunes)

© 1998, 2001, 2003, 2005, 2007, 2009 Daniel Kersten, Computational Vision Lab, Department of Psychology, University of Minnesota.