

Introduction to Neural Networks U. Minn. Psy 5038

Lecture 9 Summed vector memories, sampling, Intro. to non-linear models

Initialization

```
In[3]:= Off[SetDelayed::write]
Off[General::spell1]
```

Summed vector memories

(see SVDMEMO in Anderson, chapter 7)

Generalized Hebb rule

Taylor series expansion

Recall that a smooth function $h(x)$ can be expanded in a Taylor's series:

```
In[33]:= Series[h@xD, 8x, 0, 3<D
```

```
Out[33]= h@0D + h^@0D x + 1/2 h^@0D x^2 + 1/6 h^@0D x^3 + O@xD^4
```

where we've used Series[] to write out terms to order 3, and $O@xD^4$ means there are more terms (potentially infinitely more), but whose values fall off as the fourth power of x , so are small for $x < 1$. $h^{nnL}@0D$ means the n th derivative of h evaluated at $x=0$.

If we only include terms up to first order, this corresponds to approximating $h[x]$ near $x=0$ by a straight line. What if we have a surface $h[x,y]$? We can approximate it near $(0,0)$ by a plane:

```
In[153]:= Series[h@x, yD, 8x, 0, 1<, 8y, 0, 1<D
```

```
Out[153]= Hh@0, 0D + h^@0,1L@0, 0D y + O@yD^2 L + Hh^@0,0L@0, 0D + h^@0,1L@0, 0D y + O@yD^2 L x + O@xD^2
```

The "generalized Hebb rule":

In general, we might model the change in synaptic weights between neuron i and j by $DW@f_i, g_jD$, where as before f_i, g_j are the pre- and post-synaptic neural activities. Then with $DW@f_i, g_jD$ playing the role of $h[x,y]$ in the above expansion, we have that $DW@f_i, g_jD$ is approximately equal to:

```
In[51]:= Expand@Normal@Series@DW@f_i, g_jD, 8f_i, 0, 1<, 8g_j, 0, 1<DDD
```

```
Out[51]= DW@0, 0D + g_j DW^@0,1L@0, 0D + f_i DW^@1,0L@0, 0D + f_i g_j DW^@1,1L@0, 0D
```

(where we've used Normal[] to cut off the $O[]$ terms, and Expand[] to expand out the products). By generalizing the learning rule to any smooth function, we see that the Hebbian rule used in the linear associator models (both auto and heteroassociation) corresponds to using only the last term.

What if we used just the first term? In other words, suppose learning depended only on the input strength f_i ? Is there any useful function for such "strengthening-by-use" synapses? Suppose we have a set of k normalized input vectors $\{f^k\}$. The learning rule says that the synaptic weights s , for a single neuron would be

$$s = \sum_k f^k \quad (1)$$

Learning is easy, but it seems that the information about the set of input vectors is pretty messed up due to superposition. There are two cases where a template matching operation (i.e. dot product) could pull out useful information. Consider,

$$s \cdot f^1 = \sum_k f^1 \cdot f^k = f^1 \cdot f^1 + \sum_{1 \neq k} f^1 \cdot f^k \quad (2)$$

We know that $f^1 \cdot f^1 > f^1 \cdot f^k$ for $1 \neq k$, but we potentially have lots of terms in the sum that could swamp out $f^1 \cdot f^1$. However, if their directions are randomly distributed, we could have many cancellations.

Further, suppose one of the inputs appears more frequently than the others, say f^1 , the this term would dominate the sum s , and we might expect that a template matching operation ($s \cdot f^1$) could provide information that a high output neuron in effect is saying "yes, this input pattern looks like something I've seen frequently". Conversely, an unusually low value of the dot product would mean that "...mm this is novel, maybe I should pay more attention to this one". The latter function has at least one potential disadvantage in that a high rate of firing would be the norm.

How familiar is X, compared to what has been seen before?

Learning: input vector sums

Let's simulate the case where $DW @ \mathbf{f}_i, \mathbf{g}_j D \mathbf{f}_i$.

I

```
In[56]:= Imatrix = {
  {0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
  {0, 0, 0, 0, 0, 0, 1, 0, 0, 0},
  {0, 0, 0, 0, 0, 0, 1, 0, 0, 0},
  {0, 0, 0, 0, 0, 0, 1, 0, 0, 0},
  {0, 0, 0, 0, 0, 0, 1, 0, 0, 0},
  {0, 0, 0, 0, 0, 0, 1, 0, 0, 0},
  {0, 0, 0, 0, 0, 0, 1, 0, 0, 0},
  {0, 0, 0, 0, 0, 0, 1, 0, 0, 0},
  {0, 0, 0, 0, 0, 0, 1, 0, 0, 0},
  {0, 0, 0, 0, 0, 0, 1, 0, 0, 0},
  {0, 0, 0, 0, 0, 0, 1, 0, 0, 0},
  {0, 0, 0, 0, 0, 0, 0, 0, 0, 0}};
```

T

```
In[57]:= Tmatrix = {
  {0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
  {0, 0, 0, 0, 1, 0, 0, 0, 0, 0},
  {0, 0, 0, 0, 1, 0, 0, 0, 0, 0},
  {0, 0, 0, 0, 1, 0, 0, 0, 0, 0},
  {0, 0, 0, 0, 1, 0, 0, 0, 0, 0},
  {0, 0, 0, 0, 1, 0, 0, 0, 0, 0},
  {0, 0, 0, 0, 1, 0, 0, 0, 0, 0},
  {0, 0, 0, 0, 1, 0, 0, 0, 0, 0},
  {0, 0, 0, 0, 1, 0, 0, 0, 0, 0},
  {0, 1, 1, 1, 1, 1, 1, 1, 0, 0},
  {0, 0, 0, 0, 0, 0, 0, 0, 0, 0}};
```

P

```
General::spell1:
Possible spelling error: new symbol name "Tmatrix"
is similar to existing symbol "Imatrix".
```

```
General::spell1:
Possible spelling error: new symbol name "Tmatrix"
is similar to existing symbol "Imatrix".
```

```
General::spell1:
Possible spelling error: new symbol name "Tmatrix"
is similar to existing symbol "Imatrix".
```

```
In[58]:= Pmatrix = {
  {0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
  {0, 1, 0, 0, 0, 0, 0, 0, 0, 0},
  {0, 1, 0, 0, 0, 0, 0, 0, 0, 0},
  {0, 1, 0, 0, 0, 0, 0, 0, 0, 0},
  {0, 1, 1, 1, 1, 1, 0, 0, 0, 0},
  {0, 1, 0, 0, 0, 0, 1, 0, 0, 0},
  {0, 1, 0, 0, 0, 0, 1, 0, 0, 0},
  {0, 1, 0, 0, 0, 0, 1, 0, 0, 0},
  {0, 1, 1, 1, 1, 1, 0, 0, 0, 0},
  {0, 0, 0, 0, 0, 0, 0, 0, 0, 0}};
```

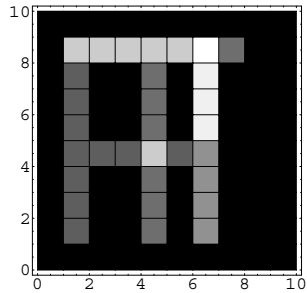
X

```
In[59]:= Xmatrix = {
  {0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
  {0, 1, 0, 0, 0, 0, 0, 0, 1, 0},
  {0, 0, 1, 0, 0, 0, 0, 1, 0, 0},
  {0, 0, 0, 1, 0, 0, 1, 0, 0, 0},
  {0, 0, 0, 0, 0, 1, 0, 0, 0, 0},
  {0, 0, 0, 0, 1, 0, 0, 0, 0, 0},
  {0, 0, 0, 1, 0, 0, 1, 0, 0, 0},
  {0, 0, 1, 0, 0, 0, 0, 1, 0, 0},
  {0, 1, 0, 0, 0, 0, 0, 0, 1, 0},
  {0, 0, 0, 0, 0, 0, 0, 0, 0, 0}};
```

```
In[103]:= normalize[x_] := N[x/Sqrt[x.x]];
Tv = normalize[Flatten[Tmatrix]];
Iv = normalize[Flatten[Imatrix]];
Pv = normalize[Flatten[Pmatrix]];
Xv = normalize[Flatten[Xmatrix]];
```

```
In[108]:= sv = Tv+Iv+Pv;
```

```
In[109]:= ListDensityPlot[Partition[sv,10]];
```

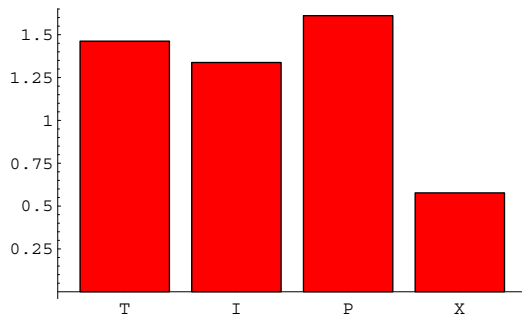


Recall: Matched filter (cross-correlator)

Let's look at the outputs of the summed vector memory to the three inputs it has seen before (T,I,P) and to a new input X. We will use a *Mathematica* graphics package that has some extra plot styles in it--in particular, the `BarChart[]`.

```
In[73]:= <<Graphics`Graphics`
```

```
In[74]:= matchedfilterout = {sv.Tv,sv.Iv,sv.Pv,sv.Xv};
BarChart[matchedfilterout,BarLabels->{"T","I","P","X"}];
```



Signal-to-noise ratio

How well does the output separate the familiar from the unfamiliar (X)? We'd like to compare the output of the model neuron when the input is the novel stimulus X, vs. the output we might expect for familiar inputs. There are several ways of summarizing performance, but one simple formula calculates the ratio of the squared output to the average squared input.

We first read in an add-on statistics package that provides us with extra functions, including `Mean[]`.

```
In[78]:= <<Statistics`DescriptiveStatistics`
```

```
In[141]:= Hsv.XvL^2 HMean 88Hsv.TvL^2, Hsv.IvL^2, Hsv.PvL^2<<L
```

```
Out[141]= 80.0245508<
```

Center vectors about zero , i.e. each has zero mean

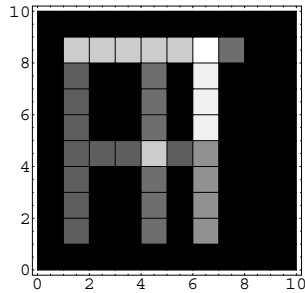
In the above representation of the letters, all the input vectors lived in the positive "quadrant", so their dot products are all positive. What if we center the vectors about zero?

```
In[142]:= normalize[x_] := N[x/Sqrt[x.x]];
Tv = normalize[Flatten[Tmatrix]];
Iv = normalize[Flatten[Imatrix]];
Pv = normalize[Flatten[Pmatrix]];
Xv = normalize[Flatten[Xmatrix]];
```

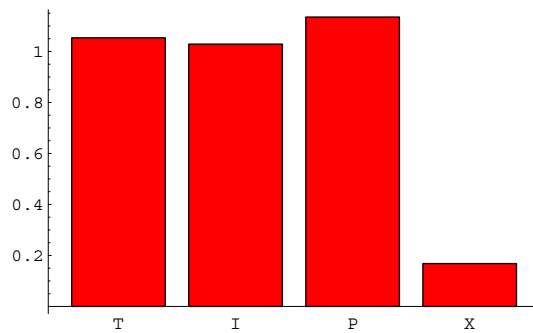
```
In[147]:= Tv = Tv - Mean@TvD;
Iv = Iv - Mean@IvD;
Pv = Pv - Mean@PvD;
Xv = Xv - Mean@XvD;
```

```
In[137]:= sv = Tv+Iv+Pv;
```

```
In[121]:= ListDensityPlot[Partition[sv,10]];
```



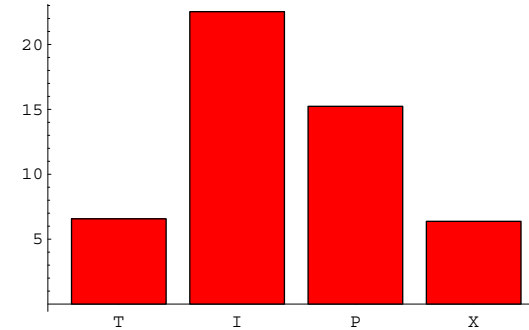
```
In[122]:= matchedfilterout = {sv.Tv,sv.Iv,sv.Pv,sv.Xv};
BarChart[matchedfilterout,BarLabels->{"T","I","P","X"}];
```



```
In[141]:= Hsv.XvL^2 HMean 88Hsv.TvL^2, Hsv.IvL^2, Hsv.PvL^2<<L
```

What happens if the summed vector memory has seen many I's and P's, but only one T?

```
In[23]:= sv = Tv + 20 Iv+ 10 Pv;
matchedfilterout = {sv.Tv,sv.Iv,sv.Pv,sv.Xv};
BarChart[matchedfilterout,BarLabels->{"T","I","P","X"}];
```



Side-note: Optimality of matched filter

The field of signal detection theory has shown that if one is given a vector input x , and required to detect whether it is due to a signal in noise ($s+n$), or just noise (n), then under certain conditions, one cannot do any better than to base one's decision on the dot product $x \cdot s$. The conditions are simple: the elements of the noise vector are assumed to be identical and independently distributed gaussian random variables, and s is assumed to be known exactly.

Learning information about the relative frequencies

We've seen how a very simple form of the generalized Hebbian learning rule can capture useful information about the relative frequencies of stimulus occurrence. This is a crude form of self-organization. We know from statistics that there are standard devices for estimating frequency of occurrence--namely, histograms. The vector sum has accumulated a kind of histogram, in the sense that it counts the number of times a particular synapse has been activated. But it is sub-optimal for our function, because what we'd like to have ideally is a device that told us how often Ts, Is, Ps occur, in a way that doesn't muddle up their representational elements.

Later we will look at the statistical framework for self-organization and the problem of measuring histograms and using these data to model probability densities, or "density estimation" as it is called.

Overview of Statistical learning theory

Statistical learning theory

A common distinction in neural networks is between supervised and unsupervised learning. The heteroassociative network was supervised, in the sense that a "teacher" supplies the proper output to associate with the input. Learning in autoassociative networks is unsupervised in the sense that they just take in inputs, and try to organize an internal representation based on the inputs.

Over the past decade, there has been considerable progress in establishing the theoretical foundations of neural networks in the larger domain of statistical learning theory. In particular, neural networks can be seen to be solving several standard problems in statistics: regression, classification, and probability density estimation. Here is a summary:

Supervised learning:

Supervised learning: Training set $\{f_i, g_i\}$

Regression: Find a function $f: \mathbf{f} \rightarrow \mathbf{g}$, i.e. where \mathbf{g} takes on continuous values.

Classification: Find a function $f: \mathbf{f} \rightarrow \{0, 1, 2, \dots, n\}$, i.e. where \mathbf{g}_i takes on discrete values or labels.

Many problems require discrete decisions. A problem with linear regression networks that we've studied so far is that they don't. Below there is a simple exercise to illustrate the how the linear associator deals with inputs that it hasn't seen before.

Next time we will take a look at the binary classification problem

$f: \mathbf{f} \rightarrow \{0, 1\}$

and see how the Perceptron solves it:

Unsupervised learning:

Unsupervised learning: Training set $\{f_i\}$

Estimate probability density: $p(\mathbf{f})$, e.g. so that the statistics of $p(\mathbf{f})$ match those of $\{f_i\}$, but generalizes well beyond the data.

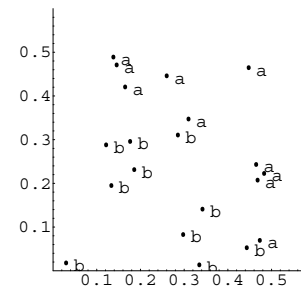
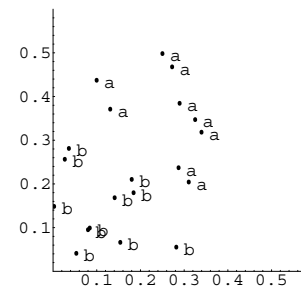
In general, \mathbf{f} is a vector with many elements that may depend on each other, so density estimation is a hard problem, and involves much more than compiling histograms.

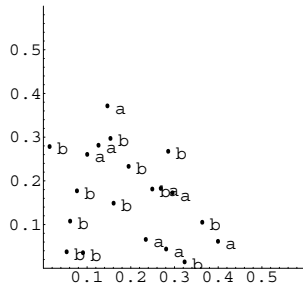
Generative modeling and Statistical sampling

Generative modeling

Whether and how well a particular learning method works depends on how the data is generated. We spend most of our time thinking about how to model learning and inference, i.e. estimation and classification. But it is also important to understand how to model incoming data. A powerful way to do this is to develop "generative models" that when implemented produce artificial data that resembles what the real data looks like. This corresponds to the statistics problem of "filling a hat with the appropriate slips of paper" and then "drawing samples from the hat".

We will begin doing "Monte Carlo" simulations of neural network behavior. This means that rather than using real data, we will use the computer to generate random samples for our inputs. Monte Carlo simulations help to see how the structure of the data determines network performance. It is useful to know how to generate random variables (or vectors) with the desired characteristics. For example, suppose you had a one unit network whose job was to take two scalar inputs, and from that decide whether the input belonged to group "a" or group "b". The complexity of the problem, and thus of the network computation depends on the data structure. The next three plots illustrate how the data determines the complexity of the decision boundary that separates the a's from the b's.





Later we'll see how the simplest Perceptron can always solve problems of the first category, but that we'll need more complex models to classify patterns whose separating boundaries are not straight.

Inner product of random vectors

In another application of Monte Carlo techniques, in the problem set you will see how the inner product of random vectors is distributed as a function of the dimensionality of the vectors.

The assumption of orthogonality for the input patterns for the linear associator would seem to make it useless as a memory advice for arbitrary patterns. However, if the dimensionality of the input space is large, the odds are pretty good that the cosine of the angle between any two random vectors is close to zero. In the exercise, you will calculate the histograms for the distributions of the cosines of random vectors for dimensions 10, 50, and 250 to show that they get progressively narrower (see Anderson, p. 187).

Probability densities and discrete distributions

As we noted earlier, most standard programming languages come with standard subroutines for doing pseudo-random number generation. Unlike the Poisson or Gaussian distribution, these numbers are **uniformly distributed**--that is, the probability of the random variable taking on a certain value is the same over the sampling range. *Mathematica* comes with a standard function, **Random[]** that enables us to generate (pseudo) random numbers that are uniform, Poisson, Normal, and others. (Why are they "pseudo" random numbers?)

Later in the course, we'll see that there is a close connection between Gaussian random numbers and linear estimators.

There are two packages **DiscreteDistributions.m**, and **ContinuousDistributions.m** which contain the definitions of distributions, cumulative distributions, and provide the means to draw samples.

The alternative package function to the built-in function **Random[]**, is **UniformDistribution[]** that generates uniformly distributed random numbers.

```
<<Statistics`DiscreteDistributions`
<<Statistics`ContinuousDistributions`
```

```
udist = UniformDistribution[0,1];
```

We can define a function, **sample[]**, to generate **ntimes** samples, and then make a list of a 1000 values like this:

```
sample[ntimes_] :=
  Table[Random[Real], {ntimes}];
```

Or like this:

```
sample[ntimes_] :=
  Table[Random[udist], {ntimes}];
```

The second way is more general, because we can use other distributions in our simulations later.

Now let us do a sampling experiment to get the list.

```
z = sample[1000];
```

Count up how many times the result was 20 or less. To do this, we will use two built-in functions: **Count[]**, and **Thread[]**. You can obtain their definitions using the ?? query.

```
Count[Thread[z<=.5], True]
```

```
517
```

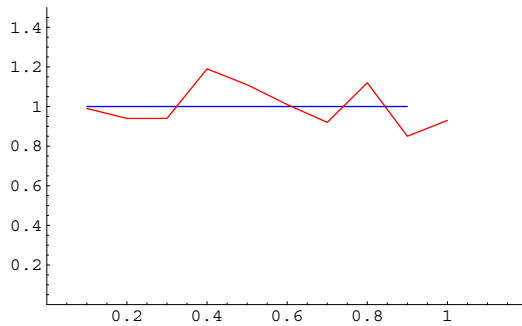
So far, we have good agreement with what we expect--about half (500/1000) of the samples should be less than 0.5. We can make a better comparison by comparing the plots of the histogram from the sampling experiment with the theoretical prediction. Let's make a table that summarizes the frequency. We do this by testing each sample to see if it lies within the bin range between x and $x + 0.1$. We count up how many times this is true to make a histogram.

```
bin = 0.1;
Freq = Table[Count[Thread[x<z<=x+bin], True], {x, 0, 1-bin, bin}];
```

Now we will plot up the results. Note that we normalize the **Freq** values by the number values in **z** using **Length[]**.

```
i=1;
theoreticalz = Table[{x, PDF[udist,x]}, {x, bin, .99, bin}];
simulatedz = Table[{x, (1/bin) N[Freq/Length[z]][[i++]]},
  {x, bin, 1, bin}];
theoreticalg = ListPlot[theoreticalz,
  PlotJoined->True, PlotStyle->{RGBColor[0,0,1]},
  DisplayFunction->Identity, PlotRange->{{0,1.2},{0,1.5}}];
simulatedg = ListPlot[simulatedz,
  PlotJoined->True, PlotStyle->{RGBColor[1,0,0]},
  DisplayFunction->Identity, PlotRange->{{0,1.2},{0,1.5}}];
```

```
Show[theoreticalg,simulatedg,
      DisplayFunction->$DisplayFunction];
```



As you can see, the computer simulation matches fairly closely what theory predicts.

Central Limit theorem Demonstration

Now we'd like to see what happens when we make new random numbers by adding up the squares of uniformly distributed ones. Why squares? It turns out that our observation below doesn't matter what we do the uniformly distributed numbers, and squaring is easy to do.

Let's define a function, `rv`, that generates random 20-dimensional vectors whose elements are uniformly distributed between 0.5 and -0.5.

```
rv := Table[Random[Real]-0.5,{i,1,20}];
ipsample = Table[rv,rv,{5000}];
```

`ipsample` is a list of 5000 elements, each of which is the dot product of a random vector.

```
bin = 0.1;
Freq = Table[Count[Thread[x<=x+bin],True],{x,-1,1-bin,bin}];
```

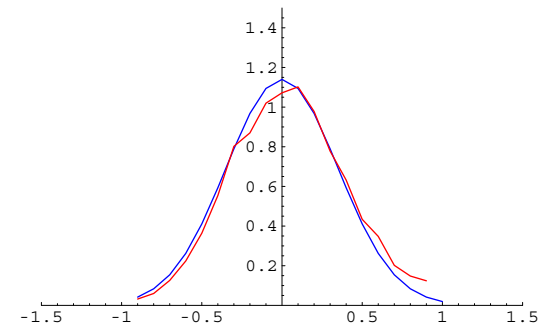
```
i=1;
simulatedz = Table[{x,(1/bin) N[Freq/Length[ipsample]][[i++]]},
                  {x,-1+bin,1-bin,bin}];
simulatedg = ListPlot[simulatedz,
                      PlotJoined->True, PlotStyle->{RGBColor[1,0,0]},
                      DisplayFunction->Identity, PlotRange->{{-1.5,1.5},{0,1.5}}];
```

Now what is the theoretical distribution? The **Central Limit Theorem** states that the sum of n independent random variables approaches the Gaussian distribution as n gets large. The n independent random variables can come from any "reasonable" distribution-- the uniform distribution is reasonable, so is the distribution of the random variable $z = x.y$, where x and y are uniform random variables.

We don't know (although we could do some theory to find out) what the standard deviation of the theoretical distribution is, but it should be normal by the Central Limit Theorem. And we know the mean has to be zero, by symmetry. So we can try out various theoretical standard deviations to see what fits the simulation best:

```
standdev = 0.35;
ndist = NormalDistribution[0,standdev];
theoreticalz = Table[{x,PDF[ndist,x]}, {x,-1+bin,1,bin}];
theoreticalg = ListPlot[theoreticalz,
                       PlotJoined->True, PlotStyle->{RGBColor[0,0,1]},
                       DisplayFunction->Identity, PlotRange->{{-1.5,1.5},{0,1.5}}];
```

```
Show[theoreticalg,simulatedg,
      DisplayFunction->$DisplayFunction];
```



Exercises

Exercise

Instead of generating samples of the inner product of random vectors, add up the elements:

```
rv := Table[Random[Real]-0.5,{i,1,20}];
ipsample = Table[Apply[Plus,rv],{5000}];
```

Calculate what the theoretical mean and standard deviation should be using the following rule:

1. The mean of a sum of independent random variables equals the sum of their means
2. The variance of a sum of independent random variables equals the sum of the variances

(And remember that the standard deviation equals the square root of the variance).

Plot up the simulated and theoretical distributions.

Linear interpolation interpretation of linear heterassociative learning and recall

```
f1 = 80, 1, 0<;
f2 = 81, 0, 0<;
g1 = 80, 1, 3<;
g2 = 81, 0, 5<;
```

```
W = Outer@Times, g1, f1D;
```

```
W.f1
```

```
80, 1, 3<
```

```
W = Outer@Times, g2, f2D;
W.f2
```

```
81, 0, 5<
```

```
Wt = Outer@Times, g1, f1D + Outer@Times, g2, f2D;
Wt.f1
```

```
80, 1, 3<
```

```
a = 0.45;
fi = a * f1 + H1 - aL * f2;
gt = a * g1 + H1 - aL * g2;
Wt.fi
gt
```

```
80.55, 0.45, 4.1<
```

```
80.55, 0.45, 4.1<
```

Introduction to non-linear models

Perceptron (Rosenblatt, 1958)

The original perceptron was fairly sophisticated--input layer ("retina" of sensory units), associator units, and response units. There was feedback between associator and response units.

The neuron models were threshold logic units (TLU)--i.e. the generic connectionist unit with a step threshold function.

These networks were difficult to analyse theoretically, but a simplified single-layer perceptron can be analyzed. Next lecture we will look at linear separability, the perceptron learning rule, and the work of Minsky and Papert (1969).

References

Bishop, C. M. (1995). [Neural Networks for Pattern Recognition](#). Oxford: Oxford University Press.

Duda, R. O., & Hart, P. E. (1973). [Pattern classification and scene analysis](#). New York.: John Wiley & Sons.

Vapnik, V. N. (1995). [The nature of statistical learning](#). New York: Springer-Verlag.

© 1998, 2001 Daniel Kersten, Computational Vision Lab, Department of Psychology, University of Minnesota.