

Introduction to Neural Networks

Heteroassociation and autoassociation

Initialization

```
In[120]:= Off[SetDelayed::write]
Off[General::spell1]
SetOptions[ListPlot, Joined -> True];
SetOptions[ArrayPlot, ColorFunction -> "GrayTones", ImageSize->Tiny, Frame->False];
```

Introduction

Last time

Linear systems overview

Introduction to learning and memory

Outer product

Outer product, $\mathbf{g}\mathbf{f}^T$ models the increase in weight strength when input \mathbf{f} is associated with an output \mathbf{g} :

```
Clear[f, g];
Outer[Times, Array[g, 5], Array[f, 4]] // MatrixForm
```

$$\begin{pmatrix} f[1] g[1] & f[2] g[1] & f[3] g[1] & f[4] g[1] \\ f[1] g[2] & f[2] g[2] & f[3] g[2] & f[4] g[2] \\ f[1] g[3] & f[2] g[3] & f[3] g[3] & f[4] g[3] \\ f[1] g[4] & f[2] g[4] & f[3] g[4] & f[4] g[4] \\ f[1] g[5] & f[2] g[5] & f[3] g[5] & f[4] g[5] \end{pmatrix}$$

Note the argument order in `Outer[Times, g, f]`, i.e. that the input \mathbf{f} comes after \mathbf{g} . So read left to right.

Learning & recall

1. Learning

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. Recall

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. Condition for perfect recall

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

$$\begin{aligned} \mathbf{W}_n \mathbf{f}_m &= (\mathbf{g}_1 \mathbf{f}_1^T + \mathbf{g}_2 \mathbf{f}_2^T + \dots + \mathbf{g}_n \mathbf{f}_n^T) \mathbf{f}_m \\ &= \mathbf{g}_1 \mathbf{f}_1^T \mathbf{f}_m + \mathbf{g}_2 \mathbf{f}_2^T \mathbf{f}_m + \dots + \mathbf{g}_m \mathbf{f}_m^T \mathbf{f}_m + \dots + \mathbf{g}_n \mathbf{f}_n^T \mathbf{f}_m \\ &= \mathbf{g}_m \end{aligned}$$

since the dot products are:

$$\mathbf{f}_n^T \mathbf{f}_m = \begin{cases} 1, & n = m \\ 0, & n \neq m \end{cases}$$

Today

A basic distinction in neural networks and machine learning is between supervised and unsupervised learning ("self-organization"). The heteroassociative network is 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 a useful internal representation based on the inputs. What "useful" means depends on the application. We explore the idea that if memories are stored with autoassociative weights, it is possible to later "recall" the whole pattern after seeing only part of the whole. Later we'll see how heteroassociative learning can be treated as a special case of autoassociative learning.

Simulation examples

Heteroassociation

Autoassociation

Superposition and interference

Heteroassociation

In this and the next section, we will use *Mathematica* to simulate a process of association between image representations of the letters, T, I, and P. You will learn more about how to manipulate lists in *Mathematica*. Critically, you will learn some of the limitations of linear recall. There are several simple exercises/questions you should try to answer.

Simulation of heteroassociative learning - Learning "IT"

Stimuli

If after seeing I, the letter T follows, you might expect that T would become associated with I. The letter I might later act as a stimulus that should elicit T as a response.

```
In[76]:= 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, 0, 0, 0, 0}};
```

```
In[77]:= Tmatrix = {
  {0, 0, 0, 0, 0, 0, 0, 0, 0, 0},
  {0, 1, 1, 1, 1, 1, 1, 1, 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, 0, 0}};
```

```
In[78]:= Pmatrix = {
  {0, 0, 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, 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, 0, 0, 0, 0, 0}};
```

We now turn our 2D image stimuli into 1D vectors. (We compute the maximum values, maxv and maxi, just to determine plot ranges for use later.)

```
In[79]:= Tv = N[Normalize[Flatten[Tmatrix]]];
Iv = N[Normalize[Flatten[Imatrix]]];
Pv = N[Normalize[Flatten[Pmatrix]]];
size = Dimensions[Imatrix][[1]];
maxv = Max[Flatten[Tmatrix]];
maxi = Max[Tv];
```

Sidenote: Making images into vectors: Flatten[] and Partition[]

You've already had practice with Flatten[]. And you may have had experience with Matlab and Python's

reshape(), and numpy.reshape() functions. **Flatten[]** takes a list of lists and turns it into a list of elements, that is, it removes all of the inner braces:

```
In[86]:= Clear[a,b,c,d]
         Flatten[{{a,b},{c,d}}]
```

```
Out[87]= {a, b, c, d}
```

Partition[] does the reverse of **Flatten[]** and takes a list of elements and structures it back into a list of lists:

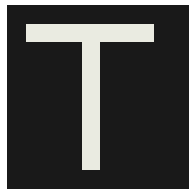
```
In[88]:= Partition[%,2]
```

```
Out[88]= {{a, b}, {c, d}}
```

For our purposes, **Flatten[]** turns a matrix representing a 2D picture (e.g. the letter **I**, or **T**) into a vector that we can store in a weight matrix memory. Later, we use **Partition[]** to turn whatever the matrix remembers into a 2D picture for comparison with the input picture originally learned. You can visualize an $n \times m$ matrix using **MatrixPlot**, **ArrayPlot** or **Image**. We'll use **ArrayPlot**:

```
In[124]:= ArrayPlot[Tmatrix, PlotRange -> {0, maxv}]
```

```
Out[124]=
```



```
ArrayPlot[Partition[Tv, size], PlotRange -> {0, maxi}, Mesh -> False]
```



Learning an association $I \rightarrow T$

Learning

Let's use the outer product to represent the change in the synaptic weights caused by the simultaneous activity of **T** and **I** which assuming a Hebb-type rule, is proportional to the product of the activities:

```
In[127]:= weights = Outer[Times,Tv,Iv];
```

```
In[128]:= Dimensions[weights]
```

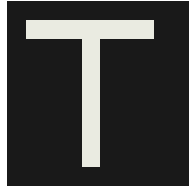
```
Out[128]= {100, 100}
```

Now if sometime later, the weights matrix is "stimulated" with the letter **I**, it produces as a response the letter **T**:

Recall: Remembering T from I

```
In[129]:= response = weights.Iv;
ArrayPlot[Partition[response, size], Mesh → False, Frame → False]
```

Out[130]=



Note that we expect this result from the rules of algebra: $(\mathbf{T}_v \cdot \mathbf{I}_v^T) \cdot \mathbf{I}_v = \mathbf{T}_v \cdot (\mathbf{I}_v^T \cdot \mathbf{I}_v) = \mathbf{T}_v$, because we normalized the input vectors.

But...

- ▶ 1. ...what if \mathbf{T}_v was the input?

What if a random vector was the input? What response would you expect?

- ▶ 2. What if you used stimulated the Transpose of weights with \mathbf{T}_v ?

What if you right-multiply \mathbf{T}_v by **weights**? (right-multiply means you put the matrix to the right of the vector, i.e. $\mathbf{T}_v \cdot (\mathbf{T}_v \cdot \mathbf{I}_v^T)$)

Add a second association $\mathbf{T} \rightarrow \mathbf{P}$ to the mix

Learning \mathbf{P} from \mathbf{T} , storing it with the association between \mathbf{I} and \mathbf{T}

```
In[131]:= weights = weights + Outer[Times,Pv,Tv];
```

Recall: Stimulate with \mathbf{T} : Response? \mathbf{P} , \mathbf{I} , or a mixture?

```
In[132]:= response = weights.Tv;
ArrayPlot[Partition[response, size]]
```

Out[133]=



- ▶ 3. Exercise

You should see evidence for *interference*. Why might you expect this based on the two inputs \mathbf{I}_v and \mathbf{T}_v ? Try comparing the dot products of the various inputs.

Can you think of a network modification for recall that might help to reduce the interference?

- ▶ 4. Exercise

The **ArrayPlot** functions don't always give the best way of seeing variations in a function or list. The ArrayPlot function automatically scales the plot range so that white corresponds to the maximum value

in the list. Try `ListPlot[response]`. What do you notice?

You can also use `Histogram[]` to see how the responses are distributed, or `Max[Tv]` to find the peak value in a list.

Add a third association $P \rightarrow I$ to the mix

Store it with the associations between $I \& T$, $P \& T$

```
In[141]:= weights = weights + Outer[Times,Iv,Pv];
```

Recall: Stimulate with P : Response? T , I , or a mixture?

```
In[143]:= response = weights.Pv;
ArrayPlot[Partition[response, size]]
```

Out[144]=



We see increased interference during recall. But we might have expected this given that the inputs are not orthogonal.

Autoassociation

If $\mathbf{f} = \mathbf{g}$, then we have an autoassociative system. There is only one set of units, and each element potentially connects to each other element. Later we will see how this architecture is used in non-linear networks. Autoassociation stores information about the relationships between the elements or features of a stimulus pattern (vector).

Let's see how this kind of knowledge representation can be used to predict or reconstruct missing information. Neural networks of this sort build internal models of the statistical structure of the ensemble they are exposed to. They can discover "2nd order statistics". But they are limited because only pairwise product information is encoded.

Reconstructive property

Autoassociation can reconstruct missing parts of a stimulus.

Suppose a vector \mathbf{x} has two sets of elements shared by parts of vectors \mathbf{f} and \mathbf{g} .

$$\mathbf{x} = \begin{pmatrix} f_1 \\ f_2 \\ \cdot \\ \cdot \\ f_m \\ 0 \\ \cdot \\ \cdot \\ \cdot \\ 0 \end{pmatrix} + \begin{pmatrix} 0 \\ 0 \\ \cdot \\ \cdot \\ 0 \\ g_1 \\ \cdot \\ \cdot \\ \cdot \\ g_n \end{pmatrix} = \begin{pmatrix} f_1 \\ f_2 \\ \cdot \\ \cdot \\ f_m \\ g_1 \\ \cdot \\ \cdot \\ \cdot \\ g_n \end{pmatrix}$$

We calculate $\mathbf{x}\mathbf{x}^T$, and then later compute: $(\mathbf{x}\mathbf{x}^T).\mathbf{f}$

What do we get?

Let's try this with \mathbf{x} , consisting of two parts $\{f_1, f_2, f_3\}$, and $\{g_1, g_2, g_3, g_4\}$, and the association of vector \mathbf{x} with itself is represented by the outer product in matrix \mathbf{W} :

```
In[165]:= Clear[f1, f2, f3, g1, g2, g3, g4];
          f={f1, f2, f3, 0, 0, 0, 0};
          g={0, 0, 0, g1, g2, g3, g4};
          x=f+g;
          W=Outer[Times, x, x];
```

```
In[170]:= x = f + g
```

```
Out[170]:= {f1, f2, f3, g1, g2, g3, g4}
```

Sometime later, we input a "version" of \mathbf{x} , but with "missing" elements--i.e. \mathbf{x} with some elements set to zero--namely, vector \mathbf{f} . What do we get in response?

```
In[180]:= Simplify[W.f]
```

```
Out[180]:= {f1 (f1^2 + f2^2 + f3^2), f2 (f1^2 + f2^2 + f3^2), f3 (f1^2 + f2^2 + f3^2),
           (f1^2 + f2^2 + f3^2) g1, (f1^2 + f2^2 + f3^2) g2, (f1^2 + f2^2 + f3^2) g3, (f1^2 + f2^2 + f3^2) g4}
```

If we replace $(f_1^2 + f_2^2 + f_3^2)$ by α , we see that $\mathbf{W}\mathbf{f}$ is proportional to \mathbf{x} . In general, the matrix \mathbf{W} restores \mathbf{f} to the pattern \mathbf{x} up to a scale factor α :

```
In[181]:= % /. (f1^2 + f2^2 + f3^2) -> alpha
```

```
Out[181]:= {f1 alpha, f2 alpha, f3 alpha, g1 alpha, g2 alpha, g3 alpha, g4 alpha}
```

If \mathbf{f} was initially normalized to 1, then $\mathbf{W}\mathbf{f}$ would be equal to \mathbf{x} .

► 5. Exercise

What are examples of a "missing" part of a pattern?

Autoassociation includes heteroassociation

At first it may seem that an autoassociative system is a more restrictive type of association than heteroassociation. But suppose we have an input/output pair $\{\mathbf{f}, \mathbf{g}\}$. If we form a new vector \mathbf{f}' in which we stack \mathbf{f} on top of \mathbf{g} , then autoassociation outer product matrix contains within it, the heteroassociation between \mathbf{f} and \mathbf{g} .

```
In[192]:= Clear[f, g];
          fprime = Join[Array[f, 3], Array[g, 3]];
```

```
Outer[Times, fprime, fprime] // MatrixForm
```

```
Out[194]//MatrixForm=
```

$$\begin{pmatrix} f[1]^2 & f[1] f[2] & f[1] f[3] & f[1] g[1] & f[1] g[2] & f[1] g[3] \\ f[1] f[2] & f[2]^2 & f[2] f[3] & f[2] g[1] & f[2] g[2] & f[2] g[3] \\ f[1] f[3] & f[2] f[3] & f[3]^2 & f[3] g[1] & f[3] g[2] & f[3] g[3] \\ f[1] g[1] & f[2] g[1] & f[3] g[1] & g[1]^2 & g[1] g[2] & g[1] g[3] \\ f[1] g[2] & f[2] g[2] & f[3] g[2] & g[1] g[2] & g[2]^2 & g[2] g[3] \\ f[1] g[3] & f[2] g[3] & f[3] g[3] & g[1] g[3] & g[2] g[3] & g[3]^2 \end{pmatrix}$$

Note that we are only representing pair-wise relationships (through products) between elements in the vectors. Later we'll see the connection to "correlational structure" and "2nd order statistics".

```
In[197]:= Outer[Times, Array[f, 3], Array[g, 3]] // MatrixForm
```

```
Out[197]//MatrixForm=
```

$$\begin{pmatrix} f[1] g[1] & f[1] g[2] & f[1] g[3] \\ f[2] g[1] & f[2] g[2] & f[2] g[3] \\ f[3] g[1] & f[3] g[2] & f[3] g[3] \end{pmatrix}$$

► 6. Question:

What can you say about the eigenvectors of the weight matrix from autoassociative learning?

Take a look at : `Outer[Times,Array[ff,4],Array[ff,4]]//MatrixForm`

Hint: Is the autoassociative matrix symmetric?

In general, is the heteroassociative matrix symmetric?

Autoassociative example with TIP pictures

Learn about T, learn about I, and store the associations together by superimposing their weight matrices

```
In[201]:= Clear[weights];
          weights = Outer[Times,Tv,Tv] + Outer[Times,Iv,Iv];
```

Sometime later, stimulate the network with an impoverished T, missing some bits

First, let's delete the 2nd row of T:


```
In[207]:= forgettingTmatrix = Partition[Tv, size];
forgettingTmatrix[[2]] = Table[0, {10}];
forgettingT = Flatten[forgettingTmatrix];
ArrayPlot[Partition[forgettingT, size]]
```

Out[210]=



Recall of the original T, given the remaining bits in forgettingT:

```
In[211]:= rememberingT = weights.forgettingT;
ArrayPlot[Partition[rememberingT, size]]
```

Out[211]=



Interference: Corrupt T again, this time with some other random bits missing

Let's do something a little more drastic to T. We'll randomly "delete" pixels of the picture:

```
In[244]:= pepper = RandomInteger[1, Dimensions[Tv][[1]]];
peppermatrix = DiagonalMatrix[pepper];
forgettingT = peppermatrix.Tv;
ArrayPlot[Partition[forgettingT, size]]
```

Out[244]=

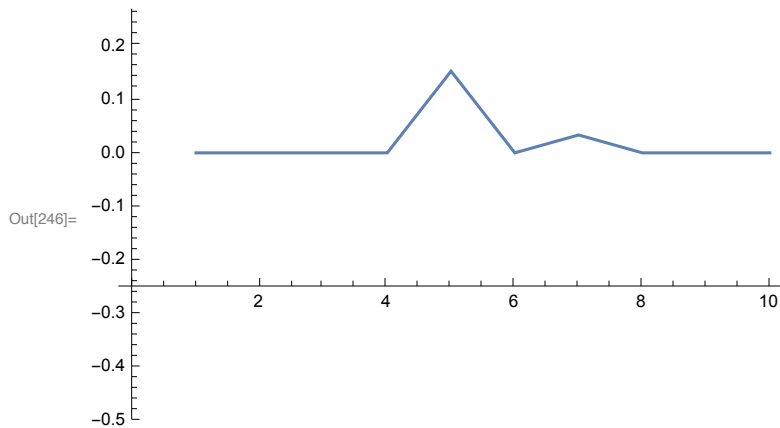


```
In[245]:= rememberingT = weights.forgettingT;
ArrayPlot[Partition[rememberingT, size]]
```

Out[245]=



```
In[246]:= ListPlot[Partition[rememberingT, size][[5]],
  AxesOrigin -> {0, -0.25}, PlotRange -> {- .5, maxi}, Joined -> True]
```



Mathematica note: Recall that you can get information about the options as well as the functions with a ?? query:

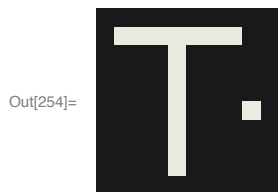
```
In[247]:= ??AxesOrigin
```

AxesOrigin is an option for graphics functions that specifies where any axes drawn should cross. >>

```
Attributes[AxesOrigin] = {Protected}
```

Interference: Corrupt T, with added noise

```
In[254]:= forgettingT = Tv;
forgettingT[[59]] = 0.27;
ArrayPlot[Partition[forgettingT, size]]
```



```
In[255]:= rememberingT = weights.forgettingT;
ArrayPlot[Partition[rememberingT, size]]
```



Although the memory looks pretty good, it is not perfect because although **Tv** and **lv** were almost orthogonal, with a cosine of about .09, they were not perfectly orthogonal. In fact, we can get a measure of how close **rememberingT** is to **Tv** in terms of the cosine of the angle between them:

```
In[256]:= Tv.Normalize[rememberingT]
```

```
Out[256]= 0.995682
```

Interference with more autoassociations: I, T, and now P too

If we have the connection matrix, weights store another letter, **P**, then we will begin to get even more interference when we try to recall **T** from a fragment of **T**:

```
In[257]:= weights = weights +
           Outer[Times,Pv,Pv];
```

```
In[258]:= rememberingT = weights.forgettingT;
           ArrayPlot[Partition[rememberingT, size]]
```

```
Out[258]=
```



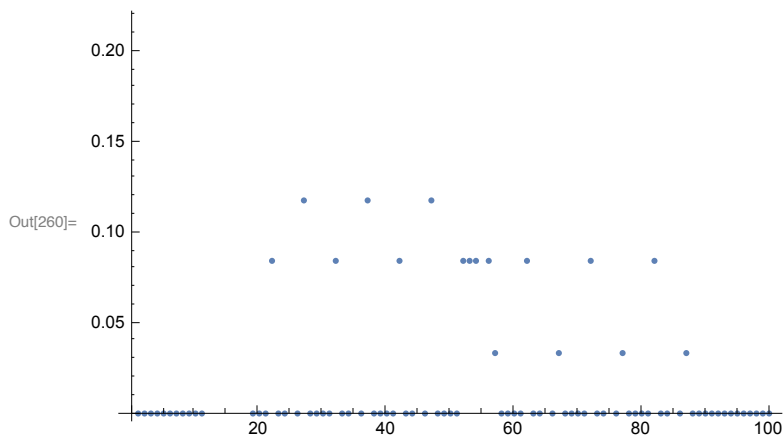
This is because the patterns we've stored are not mutually orthogonal, and in particular, **P** is too close to **I** and **T**:

```
In[259]:= {Iv.Pv, Tv.Pv, Tv.Iv}
```

```
Out[259]= {0.243332, 0.367884, 0.0944911}
```

We can picture the range of values that rememberingT takes on:

```
In[260]:= ListPlot[rememberingT, Joined -> False]
```



Include a threshold. Applying a function over a list

Define a non-linear threshold, **step[x_]**, which when applied to **rememberingT** removes the interference. A critical parameter is the threshold.

Note: When you define a new function, it is not necessarily "**Listable**". If not, here are two solutions.

```
In[275]:= step[x_, t_] := If[x > t, 1, 0.0];
```

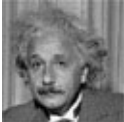

This is the same sort of exercise as above, but uses graylevel patterns, and shows a few tools for importing image data.

Reading in images

Note that you can import data from any accessible URL. E.g.

```
In[290]:= Import["http://gandalf.psych.umn.edu/users/kersten/kersten-lab/courses/
  NeuralNetworksKoreaUF2011/MathematicaNotebooks/Lect_8_HeterAuto/einstein64x64.
  jpg"]
```

Out[290]=



```
Import["http://gandalf.psych.umn.edu/users/kersten/kersten-lab/courses/
  NeuralNetworksKoreaUF2011/MathematicaNotebooks/Lect_8_HeterAuto/
  shannon64x64.jpg"];
```

You can import a file, such as `einstein.jpg` with a dialog box:

```
ieinstein = Import[SystemDialogInput["FileOpen"], "Data"];
with the size extracted with:
size = Dimensions[ieinstein][[1]];
```

Or you can just drag an image into an appropriate argument slot of a function. This is what we've done below.

Autoassociation with some other patterns: Einstein or Shannon

Turn `einstein64x64.jpg` into a 64x64 matrix of intensities:

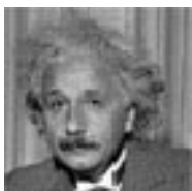
```
In[295]:= ieinstein = ImageData[];
```

```
In[296]:= size2 = Dimensions[ieinstein][[1]]
```

Out[296]= 64

```
In[297]:= geinstein = ArrayPlot[ieinstein, ColorFunction -> GrayLevel]
einstein = N[Normalize[Flatten[ieinstein]]];
```


Out[297]=



Turn shannon64x64.jpg into a 64x64 matrix of intensities:

```
In[299]:= ishannon = ImageData[];
```

```
In[300]:= gshannon = ArrayPlot[ishannon, ColorFunction -> GrayLevel]
shannon = N[Normalize[Flatten[ishannon]]];
```

Out[300]= 


Autoassociatively store einstein

```
In[302]:= weights = Outer[Times, einstein, einstein];
```

► 8. What are the dimensions of "weights"?

Make an einstein picture with some "salt and pepper noise"--random intensities at random locations, **einstein64x64missing.jpg**:


```
In[303]:= forgettingeinstein = einstein;
Table[forgettingeinstein[[RandomInteger[{1, 4096}]]] =
  RandomReal[{0.0, Max[einstein]}], {1000}];
gforgettingeinstein = ArrayPlot[Partition[forgettingeinstein, size2],
  ColorFunction -> GrayLevel]
```

Out[305]= 

Recall einstein, with only a part of the original pixels in einstein, i.e. with forgettingeinstein as input:

```
In[306]:= rememberingeinstein = weights.forgettingeinstein;
grememberingeinstein = ArrayPlot[Partition[rememberingeinstein, size2], ColorFunction -> Gra

In[308]:= Show[GraphicsRow[{geinstein, gforgettingeinstein, grememberingeinstein}]]
```

Out[308]= 

Now superimpose the outer product weight matrices for einstein and shannon:

```

In[309]:= deltaweights = Outer[Times,shannon,shannon];
In[310]:= weights = weights + deltaweights;
In[311]:= rememberingeinstein = weights.forgettingeinstein;
grememberingeinstein = ArrayPlot[Partition[rememberingeinstein,size2], ColorFunction -> Gra
In[313]:= Show[GraphicsRow[{geinstein, gshannon, gforgettingeinstein, grememberingeinstein}]]

```

Out[313]=



We see there is considerable interference with just two pictures. And we only have two images so far!

How could the images be recoded so as to reduce the interference?

- ▶ 9. Try replacing einstein and shannon with a sparser representation, e.g. in terms of edges, using `EdgeDetect[]`.
- ▶ 10. What if you computed the spectra, with respect to some orthonormal basis, of each image and used these as inputs to your outer product learning rule?

Where are we headed?

The above simulations helped to show how far we could with a hebbian rule for storage, and a linear rule for recall. The answer is “not very far”.

Improve learning or improve recall?

We've seen that the linear associator has problems of interference. How do we improve the network? We have two general strategies:

- 1) improve the learning and storage, so that linear recall will do better;
- 2) improve the recall mechanism, so that a simple Hebbian outer product rule can still be used.

Later we will derive a new learning rule that builds better association matrices for certain problems like interpolation, regression, and generalization. And we will look at non-linear recall mechanisms that make cleaner classifications for memory problems.

Neural networks as statistical pattern processing

We are going to step back, and ask ourselves what these networks are doing from the perspective of statistical pattern processing. This will lead naturally to a framework for understanding both supervised and unsupervised learning networks from a probabilistic perspective, and machine learning.

Heterassociation will lead to regression.

Autoassociation will lead to second-order statistical learning.

A deeper understanding of the information processing roles of 1) our learning rule; and 2) our recall mechanism, helps in two ways. First, we will have a better idea of what we need to do to get the net-

work to solve a given problem. Second, we may discover that there is a different problem for which it is better suited.

So linear heterassociation recall may be a poor model on which to build classification, but it might be a good interpolation model. Linear autoassociation learning may be a poor way to store information about specific memories, but may be a good way of learning about the statistics of ensembles.

As a preview of the latter, consider another application of autoassociation learning, where the goal is not to remember a particular previous input, but rather to learn something about an ensemble of inputs that all belong to the same class. Practical examples are collections of networks that can learn about their own special environments. For example, you want to build an automated driving system. The problem is complex, in part, due to different kinds of demands placed by the driving environment. So you decide you need three "experts": one for single-lane country roads, another for two-lane highways, and another for four-lane divided highways. Now you put them all in one vehicle and let them all monitor their sensory inputs as the vehicle is being driven (by one of them). When given a new environment, these experts "compare notes" to see how well this new environment fits their internal models or domain of expertise. Which ever expert "knows" the new environment the best, gets to drive the car. This is related to the idea of mental modules. The part of the driving expert that validates the environment could be realized by an autoassociative network. It can do this by testing how well its prediction of the environment's input to its sensors fits the actual sensor measurements.

Next time

Introduction to non-linear networks and classifiers. Demonstrate learning and classification with a single-layer perceptron.

Then in the lecture after next, we'll revisit today's linear model but with an improved learning rule, called Widrow-Hoff. This will allow us to make direct ties to traditional methods of regression in statistics. Further, an extension of this rule to the "error back-propagation" rule will enable us to learn weights in non-linear neural networks with more than one layer of weights that can solve non-linear problems of both regression and classification.