

# Introduction to Neural Networks

## Self-organization and efficient neural coding

### Initialization

```
Off[SetDelayed::write];  
Off[General::spell1];
```

```
SetOptions[ArrayPlot, ColorFunction → "GrayTones",  
  DataReversed → False, Frame → False, AspectRatio → Automatic,  
  Mesh → False, PixelConstrained → True, ImageSize → Small];
```

### Reminders

Next week Monday, Dec 5th. Final Homework due

Next week Friday, Dec 9th. Complete Draft of final project due by 5 pm.

Monday, Dec 12th. Peer comments due.

Monday, Dec 12th. Class presentations

10 minutes each.

Wednesday, Dec 14th. Class presentations.

10 minutes each.

---

## Introduction

### Lateral organization

Efficiency principles:

Self-organizing topographical maps that minimize “wiring length”

**Codes that minimize the number of dimensions required for reconstruction**

**Codes that minimize the number of active neurons required to represent a high-dimensional signal--“sparse coding”**

Efficient codes for learning new associations

### Machine learning view of “lateral organization”

Unsupervised learning--what does it mean to learn without a teacher?

Statistical view:

Learning as probability density estimation and exploitation. From  $N$  samples (e.g. images), can one improve the representation and transmission of information? What does "improve" mean?

Smaller number of dimensions? Noise resistance? Fewer spikes per bit (less energy), fewer neurons?

Straightforward to estimate probability with a single random variable, e.g. histogram, and to reduce dimensionality, e.g. to summary or sufficient statistics: mean, standard deviation

Less straightforward to estimate density as dimensions go up, i.e. a vector random variable with unknown parametric form. Less straightforward because of the problem of getting enough samples to fill all the bins in the histogram. Can still try to parametrically model distribution constrained by vector means, autocovariance, higher-order statistics,

For natural patterns, it is typically the case that even tho' the data maybe  $n$ -dimensional vectors, the ensemble lives in a much smaller space. How can we find that space? Once density is found, what is it good for? We'll first look at examples from vision and learn some basic principles of coding. Principles of efficient coding are general, and are applicable to other perceptual, and cognitive domains as well.

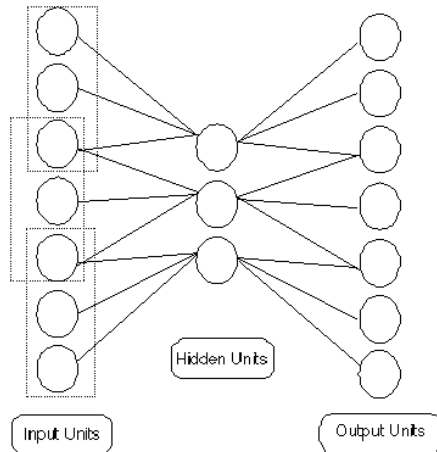
While this lecture covers "flat" or lateral organization, there are unsupervised learning methods can also be applied to learning hierarchical structure, cf. Zhu, L., Chen, Y., & Yuille, A. (2011).

---

## Unsupervised learning, self-organization, and data compression

When we introduced backpropagation, we saw the increased computational power of a multi-layer network. The reason for the increased power is that the hidden units discover effective ways of representing contingencies in the training data set. For example, a solution to the XOR problem in effect discovers how to do AND as well as OR relations, and piece these together.

One of the problems with multi-layer nets is understanding exactly what they have discovered and are representing in the hidden layers. It is perhaps easiest to begin tackling this problem by setting up a network to do autoassociation. So let's turn the supervised backprop network with three layers of nodes (two layers of weights) into an unsupervised learning algorithm simply by setting the output equal to the input. Then we are seeking weights that achieve the goal that the outputs come as close as possible to matching the inputs, in a least squares sense. If we have a smaller number of hidden units than input or output units, we can ask: What has the network discovered about the input ensemble that is captured with the smaller dimensionality of the hidden unit layer?

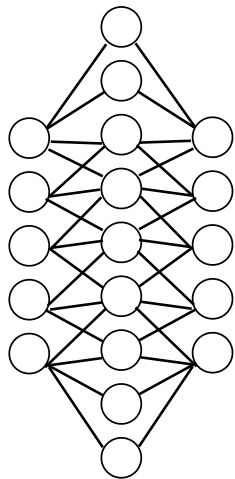


A theoretical result (Baldi and Hornik, 1988) showed that for the linear case this kind of network is closely related to a standard statistical technique called "Principal Component Analysis", (PCA) that dates back to 1901 (Karl Pearson) and independently discovered and named in 1933 (Hotelling, H. (1933). Analysis of a complex of statistical variables into principal components. J. Educ. Psych., 24, 417-441,498-520). The idea is that the variability in a data set consisting of  $n$ -dimensional vectors may concentrate along certain axes or subspaces of dimension  $m < n$ . Principal component analysis is one standard technique for finding the dimensions that capture the most variation.

In the autoencoder case, suppose there are  $n$  input units, and  $m$  hidden units. Baldi and Hornik showed that the hidden units were finding the  $m$ -dimensional sub-space that capture the most variance. This is not the same as principal component analysis, but it is closely related.

In this notebook you will learn about PCA, and then see how it can be done by neural-like systems that use local hebbian learning. In order to understand PCA, let's start with the following simple two neuron system. This is the simplest system for which correlational structure can be analyzed. The rationale is that the two neuron system will give us insight into the problem of finding structure in data sets with really high dimensionality, such as images or speech.

Note that we can also ask: What has the network discovered about the input ensemble that is captured with the larger dimensionality of the hidden unit layer? We will return to this in the context of sparse coding later.



## A simple generative model: Statistical correlations of a 2D input ensemble

### The generative model

Recall that the idea behind generative modeling is to understand the structure of the data coming into the network. This gives us a better idea of where the network should have problems. And as we'll see later, the generative model is needed to specify an optimal network. Our goal is first create data that has a particular correlational structure, and then to see how a neural network can discover the structure.

Consider a "two-neuron" system whose inputs are correlated. The random variable,  $\mathbf{rv}$ , is a 2D vector representing the input. The scatter plot for this vector has a slope of  $\text{Tan}[\theta] = 0.41$ . The variances along the axes are:

$$\text{bigstd}^2 = 4^2 = 16 \text{ and } \text{smallstd}^2 = .25^2 = .0625.$$

`gprincipalaxes` is a graph of the principal axes which we will use for later comparison with simulations. *Mathematica* provides functions for sampling from Gaussian and other distributions, rather than the standard uniform distribution that `RandomReal[]` provides.

We define `ndist` so that `RandomReal[ndist]` returns numbers whose average is zero, and whose standard deviation is 1. Variance is equal to standard deviation squared. If we multiple `RandomReal[ndist]` by `bigstd`, we get random numbers with variance `bigstd^2`. And similarly for `smallstd`.  $\mathbf{rv}$  is the projection of these numbers onto the x and y axes.

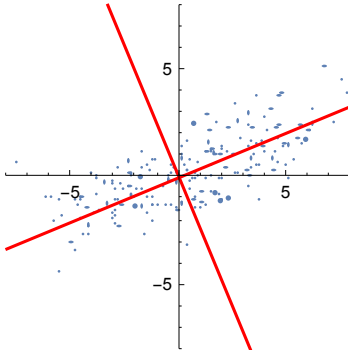
```
ndist = NormalDistribution[0,1];theta =Pi/8;
bigstd = 4.0; smallstd = 0.25;
smallstd = 1.0;
alpha = N[Cos[theta]]; beta = N[Sin[theta]];
rv :=
{bigstd x1 alpha + smallstd y1 beta, bigstd x1 beta -
smallstd y1 alpha} /.
{x1-> RandomReal[ndist],y1-> RandomReal[ndist]};

gprincipalaxes = Plot[{x (beta/alpha),
x (-1/(beta/alpha))}, {x,-8,8},
PlotRange->{{-8,8},{-8,8}},
PlotStyle->{RGBColor[1,0,0]},
AspectRatio->1];
```

$x_1$  and  $y_1$  are said to be *correlated*. Let's view a scatterplot of samples from these two correlated Gaussian random variables.

```
npoints = 200;
rvsamples = Table[rv,{n,1,npoints}];
```

```
g1 = ListPlot[rvsamples, PlotRange -> {{-8, 8}, {-8, 8}},
  AspectRatio -> Automatic]; Show[g1, gprincipalaxes]
```



The slope of the dominant axis is:

```
{Tan[theta], N[Tan[theta]]}
```

```
{Tan[π/8], 0.414214}
```

---

## Principal Component Analysis (PCA) using standard statistical methods

In the previous section, we developed a model for synthetic data--so we know the statistical structure. Usually, we don't have a good model of the statistical structure of an ensemble, but want to discover something about it. Where there are lots of dimensions, it can be hard to find a complete model of the underlying distribution; however, we can still analyze the data to find means, variances and other statistics. And when dealing with a high-dimensional ensemble, sometimes most of the variation occurs in some subspace. In other words, there is a particular direction with the largest variance.

In the example above, most of the variation is occurring in the 1D subspace defined by a line of slope determined by the angle **theta**. How could we have discovered that only from the data--i.e. without having the model before us? We could graph it and look. But what if the dimensionality of the space is really big?

PCA provides a method. PCA seeks out a new coordinate system that is just a rotation of the original that does three very interesting things:

- 1) The data when projected onto the new rotated coordinates are no longer correlated--in fact the autocovariance matrix (see below) is diagonal;
- 2) The new coordinates can be ordered so that the main or "principal component" has the most variance, the next coordinate has the second most, and so forth. How is this done? It turns out that the *eigenvectors of the autocovariance matrix are the principal component vectors that point in the direction of the new coordinates, and the eigenvalues are the variances of the data when projected onto the new coordinates (or axes).*
- 3) The error when reconstructing the data by projection along the new coordinates is minimized.

What good is PCA for a data set? If the variance of some of the projections is near zero, one can in fact dispense with these components and achieve a good approximate coding of the data with just the remaining coordinates. We are going to see that in the 2D example, we can get an economical coding

of the data with just one number, rather than two.

## Calculate the autocovariance matrix

We need to quantify the idea of a correlation or covariation between two random variables, and then between random vectors.

A quick review of covariance. First consider two (scalar) random variables,  $x$  and  $y$ . Let  $E[\cdot]$  stand for the expected or average of a random variable,  $\bullet$ . Suppose we want a measure of the covariation of two random variables,  $x$  and  $y$ . One way is to first subtract off the mean from both giving:  $x_1 = x - E[x]$ , and  $y_1 = y - E[y]$ . Then we find the average of the product  $E[x_1 y_1]$ . This is called **covariance**. If  $x_1$  and  $y_1$  tend to go up together (a high value of  $x_1$  predicts a high value of  $y_1$ , etc..), then the product will tend to be big on average--they are **correlated**. (If  $x=y$ , then this is just the standard definition of variance.) (The correlation coefficient in statistics is a normalized version of covariance.)

Suppose we have a vector (like  $\mathbf{rv}$ , of dimension 2 or perhaps higher) whose elements are random variables. These vectors make up an ensemble, from which we may have a particular subset of random samples, as shown in the above scatterplot. We can define a **cross-covariance** analogous to covariance. But for our purposes in this lecture, all we need is a measure of how the various elements of a vector covary with each other. The **autocovariance** matrix of a vector random variable,  $\mathbf{x}$ , can be defined by first subtracting off the mean vector  $E[\mathbf{x}]$ , and then averaging over all the *outer-products* of each vector,  $\mathbf{x} - E[\mathbf{x}]$ , with itself:

$$E[ (\mathbf{x} - E[\mathbf{x}])(\mathbf{x} - E[\mathbf{x}])^T ]$$

(Note this could be estimated using the auto-associative Hebbian memory rule we studied earlier--i.e. add up all the outer products, and then normalize by the number of learning pairs.) The diagonal values of the autocovariance matrix correspond to the variances of each element. The off-diagonals are the covariances for each pair of elements. Because of the symmetry of the outer-product (multiplication commutes), the autocovariance matrix is symmetric.

Let's compute the autocovariance matrix for  $\mathbf{rv}$ . The calculations are simpler because the average value of  $\mathbf{rv}$  is zero. As we would expect, the matrix is symmetric:

The variances of the two inputs (the diagonal elements) are due to the projections onto the horizontal and vertical axis of the generating random variable.

```
rvsamples2 = Map[# - Mean[rvsamples] &, rvsamples];

autolist = Table[
  Outer[Times, rvsamples2[[i]], rvsamples2[[i]],
    {i, Length[rvsamples2]}];
MatrixForm[auto =
  Sum[autolist[[i]],
    {i, Length[autolist]}] / (Length[autolist] - 1)]
Clear[autolist];
( 13.9458  5.15545 )
( 5.15545  3.13867 )
```

Note that the built-in function **Covariance[ ]** takes care of the above details in one step:

```
Covariance[rvsamples, rvsamples] // MatrixForm
( 13.9458  5.15545 )
( 5.15545  3.13867 )
```

- ▶ 1. Check to see if the empirical autocovariance matrix gives what you might predict from the standard deviations in the generative model. E.g. Try changing `npoints=1000`. Set `theta` to a small number (0.0001).

## Calculate the principal components (eigenvectors of the autocovariance matrix)

Now calculate the eigenvectors of the autocovariance matrix

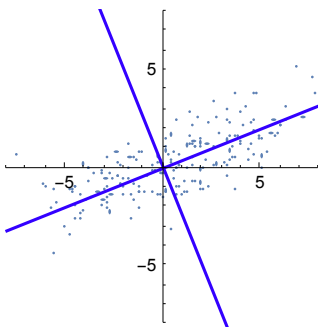
```
eigauto = Eigenvectors[auto]
{{-0.928311, -0.371805}, {0.371805, -0.928311}}
```

- ▶ 2. Recall that the eigenvectors of a symmetric matrix are orthogonal. You can check that these are.

Let's graph the principal axes corresponding to the eigenvectors of the autocovariance matrix together with the scatterplot we plotted earlier. One axis is defined by:  $y = \text{eigauto}[[1,2]]/\text{eigauto}[[1,1]] x$ , and the other axis by:  $y = \text{eigauto}[[2,2]]/\text{eigauto}[[2,1]] x$ .

```
gPCA = Plot[{eigauto[[1,2]]/eigauto[[1,1]] x,
            eigauto[[2,2]]/eigauto[[2,1]] x},
            {x,-8,8}, AspectRatio->Automatic,
            PlotStyle->{RGBColor[.2,0,1]}];
```

```
Show[g1, gPCA]
```



You can see that PCA has *discovered* important structure in the input ensemble as shown by the blue lines.

- ▶ 3. Compare above graph by including `gprincipalaxes` in `Show[]`
- ▶ 4. Try changing `npoints=5`. Then plot the true principal axes with `gPCA`, using `show[g1,gPCA,gprincipalaxes]`. Do `gPCA` and `gprincipalaxes` match? Why not?

## Calculating the variances (eigenvalues)

The eigenvalues give the ratio of the variances of the projections of the random variables  $\mathbf{rv}[[1]]$ , and  $\mathbf{rv}[[2]]$  along the principal axes:

```
eigenvalues = Eigenvalues[auto]
{16.0106, 1.07381}
```

Let's project the data onto the principal axes and calculate the autocovariance matrix in the new coordinate system. The projected values are given by: `eigauto.rvsamples`.

The new coordinates for this transformation of the data have zero correlation.

```
autolist = Table[
  Outer[Times,eigauto.rvsamples2[[i]],
    eigauto.rvsamples2[[i]]],{i,Length[rvsamples2]};
MatrixForm[Chop[
  Sum[autolist[[i]],
    {i,Length[autolist]}/Length[autolist]]]
Clear[autolist];

$$\begin{pmatrix} 15.9306 & 0 \\ 0 & 1.06844 \end{pmatrix}$$

```

Note that the off-diagonal elements (the terms that measure the covariation of the transformed random variables) are zero. Further, the diagonal elements are estimates of the population variances along the principal axes. They should be approximately equal to the population variances specified by the generative model: i.e. the **bigstd**<sup>2</sup>, and **smallstd**<sup>2</sup> from our synthetic data process.

We've demonstrated the fundamental properties of PCA:

- 1) PCA decorrelates previously correlated input data;
- 2) PCA discovers principal axes that capture the dimensions with the dominant variance.

A property that we haven't shown is that PCA minimizes the reconstruction error when projecting the data onto rotated coordinates.

- ▶ 5. Suppose you had 2-dimensional data where the first element represented age, and the second height for a population of people. Does PCA as defined above provide a unique set of principal components? HINT: What happens to the scatter plot if you change the measurement units?

*Shortly we will address the question: How can we do PCA using "brain-style" computation? But first, SVD.*

## Singular Value Decomposition (SVD)

For computational reasons, it is usually better to calculate the principal components using Singular Value Decomposition (SVD). The homework assignment illustrates the advantages of SVD, and shows the relationship between the PCA and SVD. For a formal explanation of the connection, see [wiki entry](#).

```
Clear[u, v, w]
{u, w, v} = SingularValueDecomposition[rvsamples2];
Dimensions /@ {u, w, v}
{{200, 200}, {200, 2}, {2, 2}}
```

The columns of u are orthogonal unit vectors called the left singular vectors.

### Singular values in w. Relation to eigenvalues and variance

The diagonal elements of w are related to the variances of the projections along the principal components that we calculated using PCA, i.e. the eigenvector of covariance method:



```
w[[1 ;; 7, 1 ;; 2]]
{{56.4457, 0.}, {0., 14.6181}, {0., 0.}, {0., 0.}, {0., 0.}, {0., 0.}, {0., 0.}}
```

Specifically, the variances are given by squaring the singular values and dividing by the number of samples -1:

```
Diagonal[w]^2 / (Length[rvsamples] - 1)
{16.0106, 1.07381}
```

Compare with PCA from the previous section:

```
eigvalues
{16.0106, 1.07381}
```

The columns of  $v$  are the “right” singular vectors, which point in the same direction as the principal components

$v[[All,1]]$  and  $v[[All,2]]$  are the column vectors of  $v$ , which represent the two principal component directions. In the figure below we represent those two vectors into the original data space.

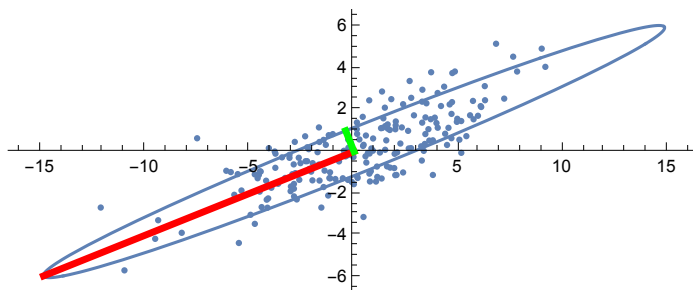
```
v[[All, 1]].eigauto[[1]]
```

1.

```
v // MatrixForm
```

```
(-0.928311 -0.371805 )
(-0.371805  0.928311 )
```

```
Show[{ParametricPlot[auto.{Cos[t], Sin[t]}, {t, 0, 2 Pi}, AspectRatio -> Automatic,
  Epilog -> {Thickness[0.01], {Red, Line[{{0, 0}, auto.v[[All, 1]]]}},
  {Green, Line[{{0, 0}, auto.v[[All, 2]]]}]}], g1 ]]
```



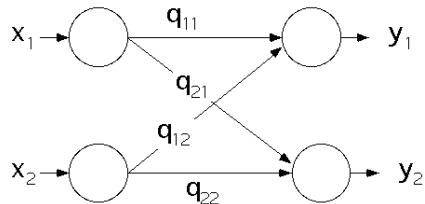
For an excellent introduction to SVD and a “psychological interpretation” in a application to identifying user preferences, see: “Singular Value Decomposition | Stanford University”, <https://youtu.be/P5ml-g91as1c>

## Principal Component Analysis (PCA) with neural networks

Neural network model using Hebb together with Oja's rule for extracting the dominant principal component

## Introduction to Oja's network: weight normalization

Consider the following linear neural network. The input and output values are represented by vectors  $\mathbf{x}$ , and  $\mathbf{y}$  respectively. The connection weights are represented by matrix  $\mathbf{Q}$ .



We will combine the outer product form of Hebb's rule, together with Oja's modification. Without Oja's rule, the Hebb rule does not place a limit on the size of the weights.

$$\Delta q_{ij} = \alpha (x_j y_i - q_{ij} y_i^2)$$

Oja's rule automatically tends to normalize the weights so that their squared sum is finite. In fact:

- 6. Show, when the weights are no longer changing, that:

$$\sum_{i,j} q_{ij}^2 = 1$$

## Implementing Oja's network

Oja's rule constrains the sum of the squares of the weights to approach 1. We will set the initial values of the weight matrix to random values between 0 and 1.

```
npoints = 400; p1 = {};  $\alpha$  = 0.01;
size = 2;
Q = Table[RandomReal[], {size}, {size}];
```

Note that a space or \* in *Mathematica* between two expressions does an element by element multiplication. We use this notation as economical way of writing Oja's rule. An example is:

```
Clear[a,b,c,d,x,y]
{{a,b},{c,d}} {x,y}
{{a x, b x}, {c y, d y}}
```

Note that this different from standard matrix/vector multiplication.

For each random input  $\mathbf{rv}$ , we compute the output  $\mathbf{y} = \mathbf{Q} \cdot \mathbf{x}$ , and the weights  $\mathbf{Q}$  get adjusted on-line by Oja's rule. Unlike computing eigenvectors, we don't have to store all npoints of the random vectors to compute an autocovariance matrix.

```

For[i=1,i<=npoints,i++,
  x = rv; y = Q.x;
  Q = Q +  $\alpha$  (Outer[Times,y,x] - Q y y);

(*p1 keeps track of the evolution of the weights
  in terms of the slopes of the rotated coordinates*)
  If[Mod[i,5]==0,
    p1 = Join[p1,{{Q[[1,2]]/Q[[1,1]],
      Q[[2,2]]/Q[[2,1]] }}]];
];

```

The slopes of the rotated coordinates can be calculated from the elements of p1. The last computed values are:

```

Q
Q[[1, 2]] / Q[[1, 1]]
Q[[2, 2]] / Q[[2, 1]]
{{0.915031, 0.407817}, {0.915031, 0.407817}}
0.445687
0.445687

```

### Graph the evolution of the network: slopes of the projection axes

Let's plot the slopes of projection axes as a function of iterations. We've sampled every 5th value, using **Mod[i,5]**, and stored it in **p1**. These values should approach the slope of the scatter plot, **Tan[theta]**.

```

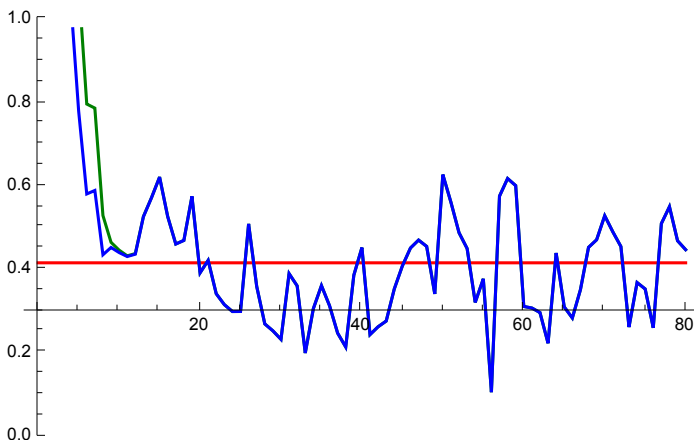
gg1=Plot[ Tan[theta],{x,0,Length[p1]}, AxesOrigin->{0,.3}, PlotRange->{0,1},
  PlotStyle->{RGBColor[1,0,0]};

gg2=ListPlot[Map#[[2]]&,p1], AxesOrigin->{0,0},
  PlotJoined->True, PlotStyle->{RGBColor[0,.5,0]};

gg3=ListPlot[Map#[[1]]&,p1], AxesOrigin->{0,0},
  PlotJoined->True, PlotStyle->{RGBColor[0,0,1]};

```

```
Show[gg1,gg2,gg3]
```

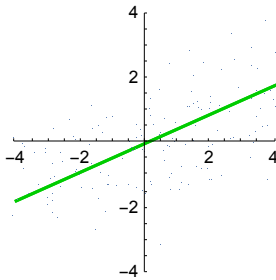


There is some random fluctuation in the weights. We can obtain more stability by having a time constant over which the Hebbian term and the variance of  $y$  are averaged.

## Graph the slope of the slope of the projection axes (ratio of network weights) together with the data

We can see how well the coordinate transformation fits the principal axes of a sample scatter plot:

```
gnetwork = Plot[
  {s p1[[Length[p1]]][[1]], s p1[[Length[p1]]][[2]]},
  {s, -4, 4}, PlotRange->{{-4,4},{-4,4}},
  AspectRatio->1,PlotStyle->{RGBColor[0,.8,0]};
  Show[gnetwork,g1]
```



You can verify that the network does a good job of extracting the principal component. Recall that the slope for the population distribution is **Tan[theta]**:

```
N[Tan[theta]]
```

```
0.414214
```

It doesn't take long for the green and blue lines to converge to identical values.

The only problem with this network is that having two output neurons is redundant--they both pull out the same principal component--the dominant axis. The slopes for both are:

```
p1[[Length[p1]]]
```

```
{0.445687, 0.445687}
```

This isn't surprising because the network was symmetrical--both output neurons saw the same inputs and updated their weights using the same rule. How can this be fixed to pick out the other principal components? Some kind of asymmetry has to be introduced.

The problem with Oja's network is that it extracts just the main principal component. Our example had two outputs, but the network is symmetric, and both outputs were the same--the projection of the input onto the main principal axes. Sanger (1989) proposed a modification to the Oja network that can extract *all* of the principal components.

The generalization of Oja's term to update weights  $q_i$  is given by:

$$\Delta q_{ij} = \alpha \left( x_j y_i - y_i \sum_{k=1}^i q_{kj} y_k \right)$$

See the Appendix for a simulation.

# Non-orthogonal decorrelation: contingent adaptation and the McCollough effect

Celeste McCollough (1965).

## Make stimulus

```
width = 128; freq=8;
grating[x_,y_,xfreq_,yfreq_] := Sign[Cos[(2. Pi)*(xfreq*x + yfreq*y)]];
```

## Vertical red adapting grating

```
xfreq = freq; theta = Pi / 2;
yfreq = xfreq / Tan[theta];
```

```
gvertred = ArrayPlot[Table[grating[y, x, xfreq, yfreq], {x, 0, 1, .01}, {y, 0, 1, .01}],
  Mesh → False, Frame → False, ColorFunction → (RGBColor[#, 0, 0] &)];
```

## Horizontal green adapting grating

```
xfreq = 0; theta = 0;
yfreq = freq;
```

```
ghorizgreen =
  ArrayPlot[Table[grating[y, x, xfreq, yfreq], {x, 0, 1, .01}, {y, 0, 1, .01}],
    Mesh → False, Frame → False, ColorFunction → (RGBColor[0, #, 0] &)];
```

## Horizontal gray test grating

```
xfreq = 0;
yfreq = freq;
```

```
ghorizgray =
  ArrayPlot[Table[grating[y, x, xfreq, yfreq], {x, 0, 1, .01}, {y, 0, 1, .01}],
    Mesh → False, Frame → False, ColorFunction → (RGBColor[#, #, #] &)];
```

## Vertical gray test grating

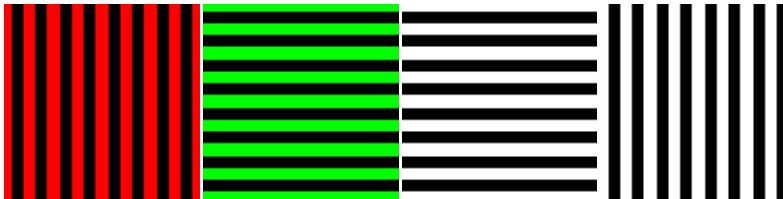
```
xfreq = freq;
yfreq = 0;

gvertgray =
  ArrayPlot[Table[grating[y, x, xfreq, yfreq], {x, 0, 1, .01}, {y, 0, 1, .01}],
    Mesh → False, Frame → False, ColorFunction → (RGBColor[#, #, #] &)];
```

## Test: Try it

1. First look at the two black and white gratings below on the right. They should look black and white, and the white should have little or no tinges of color.
2. Now look at the left vertical red grating, then the right horizontal green, then back at the left and so forth for 2-4 minutes. (You can use **Pause[]** for timing in seconds).
3. Then when adapted, look at the black and white test gratings again. The white bars of horizontal B&W should look pinkish, and the vertical bars greenish.

```
Show[GraphicsGrid[{{gvertred, ghorizgreen, ghorizgray, gvertgray}},
  Spacings → {Scaled[0.01`], Scaled[0.01`] }]]
```



## Is there a functional explanation for what is going on? A neural network explanation?

Barlow suggested that adaptation is the result of the visual system adjusting to new statistical dependencies between features. It isn't just that neurons are "getting tired". Specifically, he suggested that perceptual systems adjust their representations of features to be independent or uncorrelated with each other (technically, independence is a stronger condition). This is a problem of neural self-organization, perceptual learning based on experience. Next we'll look at several neural networks that re-organize to decrease the correlations between their firing. One network may provide an explanation for McCulloch's color-contingent after-effect.

## Non-orthogonal decorrelation & contingent adaptation

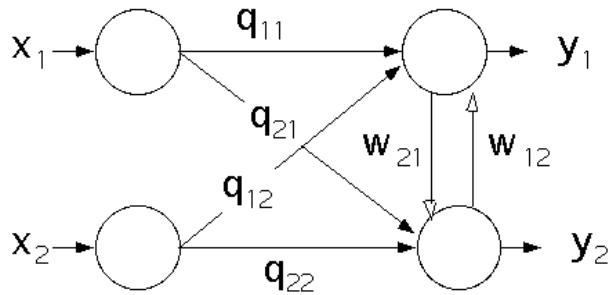
### Foldiak's scheme combining Hebbian and anti-Hebbian learning

Rigid rotations aren't the only possible transformations that decorrelate the inputs. Further, one might want a new coordinate system that shares the variance equally--after all, we do not have strong evidence that neurons vary greatly in their ability to code the range of variation. This section looks at a

network due to Peter Foldiak.

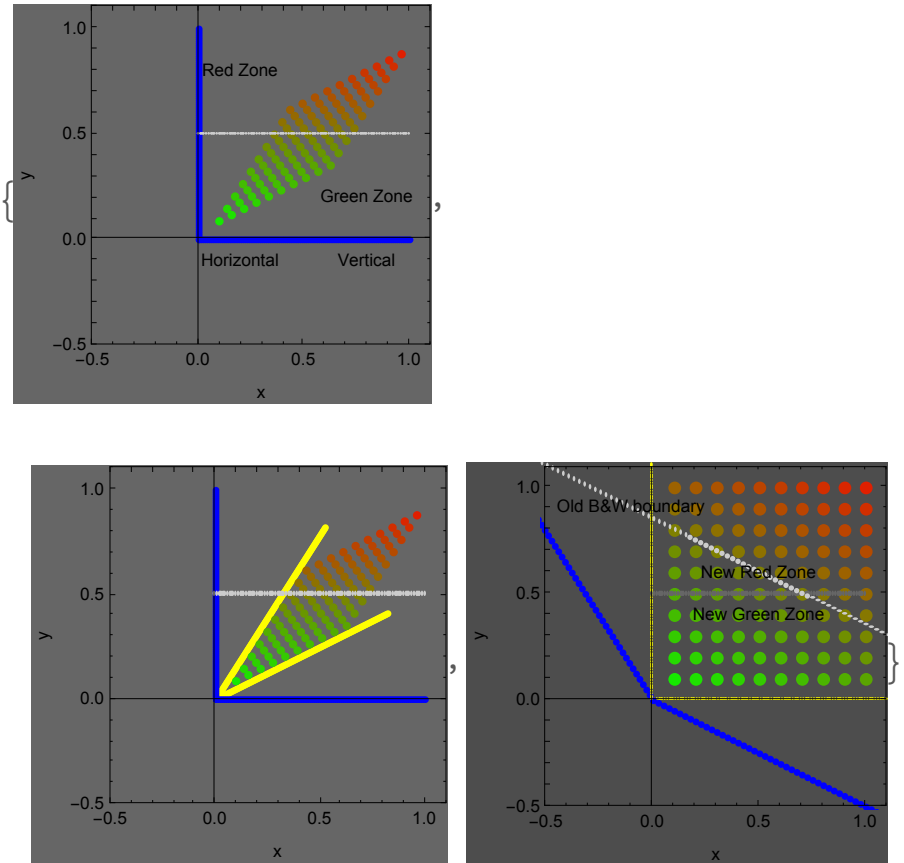
Foldiak and Barlow devised a neural network that combined a Hebbian learning rule on the forward connections with anti-Hebbian learning on the inhibitory connections between the output units. Oja's rule was used to normalize the weights. It can be shown that decorrelated output values are steady-state solutions for the weight changes.

One of the reasons for interest in this kind of model are the potential relations with the physiology. Inhibitory links are well-known, and evidence for anti-Hebbian learning is something to be looked for empirically.



The figures below illustrate how features that covary in color and orientation space (first panel) can be projected onto non-orthogonal axes (yellow in middle panel). When the data are re-plotted in the new coordinate system, they are no longer correlated (third panel). The white dotted line and blue lines in the third panel show the re-mapped color boundary, and original axes in the new space, respectively.

See Non-orthogonal Decorrelation Code in the Appendix for the code to illustrate decorrelation.



## Sparse codes

Some of this material was briefly covered in Lecture 5 and Lecture 22.

In PCA, the number of basis functions or vectors is less than or equal to the dimensionality of the input. But what if “efficiency” has another meaning, e.g. represent a high-dimensional input (an image) with as few features as possible? ...and we allow for over-complete representations where the number of feature detectors could be more than the dimensionality of the input. In Lecture 5, we introduced the general idea of a linear decomposition of an image into combinations of basis elements, where the  $A_i$ 's represent “features”, and the  $s_i$ 's (neuron activity) how much of each feature is needed to reconstruct the image  $I(x,y)$ :

$$I(x,y) = \sum_{i=1}^n A_i(x,y)s_i \quad s_i = \sum_{x,y} W_i(x,y)I(x,y)$$

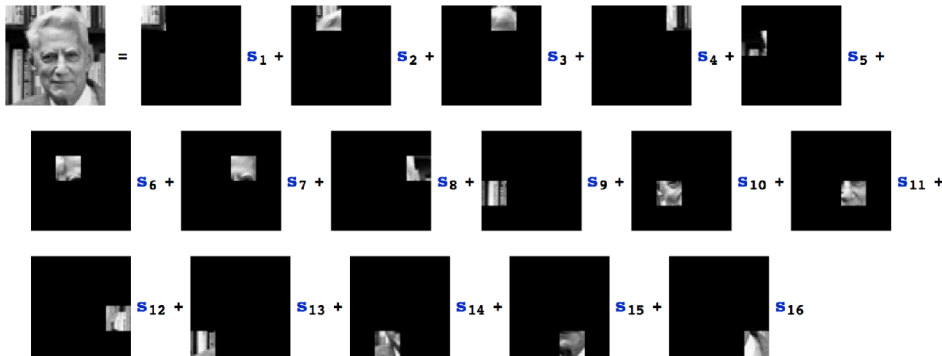
(While A isn't necessarily the same as W, it is often the case in neural models of V1 cells.)

The principal components provided by PCA are just one possibility to represent the features. But there is an infinite range of other possibilities.



$$\text{Image} = s_1 \cdot \text{Patch}_1 + s_2 \cdot \text{Patch}_2 + \dots + s_k \cdot \text{Patch}_k$$

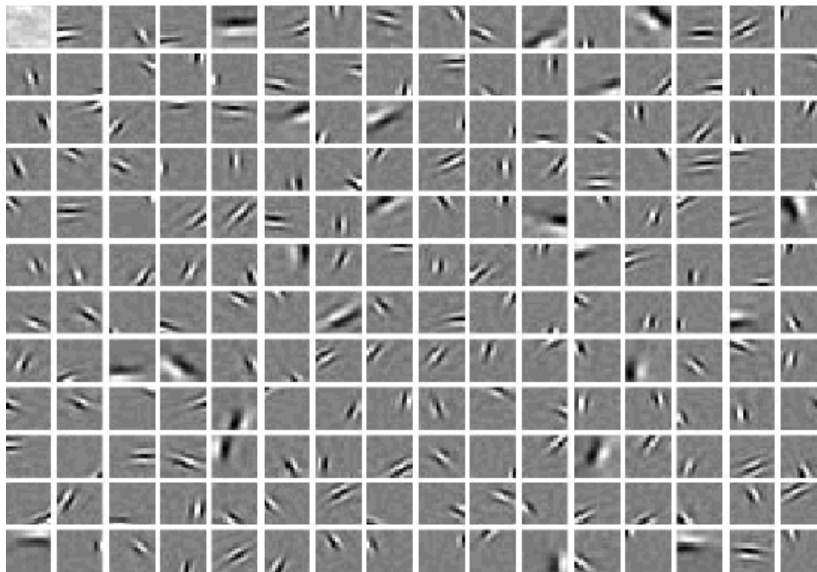
We saw that a single image could be represented not only by the usual basis, i.e. pixels but also by “patches”, such as:



Here only a few features are required for one image...but what if we wanted to have a set of features, or “dictionary” that was in “good” for *all* natural images? Good, efficient representation is interpreted as finding the receptive field weights that minimize the sum of squared errors AND # active neurons. So given  $L(x,y)$  over a set of images find the  $A_i(x,y)$ 's that minimize:

$$\left[ L(x, y) - \sum_i s_i A_i(x, y) \right]^2 + \sum_i B(s_i)$$

averaged over all the images. The first term penalizes loss of information about the image, and the second penalizes too many active neurons.



Olshausen & Field's model of V1 receptive fields captures localized sensitivities to orientation and spatial frequency.

Higher-order statistical dependencies between these filters also seem to be discovered and removed by the neural circuitry in V1, through a mechanism called “divisive normalization”. Schwartz, O., & Simoncelli, E. P. (2001). Natural signal statistics and sensory gain control. *Nature Neuroscience*, 4(8), 819–825.

## Appendix

### A generalization of Oja’s rule for extracting all of the principal components (Sanger, 1989)

#### Introduction to Sanger's network for PCA

The problem with Oja's network is that it extracts just the main principal component. Our example had two outputs, but the network is symmetric, and both outputs were the same—the projection of the input onto the main principal axes. Sanger proposed a modification to the Oja network that can extract *all* of the principal components.

We will use the same network as in the above example. However, the normalization part of learning rule will be asymmetric. The generalization of Oja's term to update weights  $\mathbf{q}$ , is given by:

$$\Delta q_{ij} = \alpha \left( x_j y_i - y_i \sum_{k=1}^i q_{kj} y_k \right)$$

#### Implementing the Sanger network

The above learning rule can be evaluated in *Mathematica* as: `LT Outer[Times,y,y].Q`, where `LT` is a lower triangular matrix. The entries above the diagonal are all zero, and the entries below and including the diagonal are one. You can verify the Sanger weight update formula with the following expressions:

```
Clear[Q,q,y,yy];
n = 2;
LT = Table[If[i>=j,1,0],{i,n},{j,n}];
Q = Array[q,{n,n}];
yy = Array[y,{n}];
(LT Outer[Times,yy,yy]).Q//MatrixForm
```

$$\begin{pmatrix} q[1, 1] y[1]^2 & q[1, 2] y[1]^2 \\ q[1, 1] y[1] y[2] + q[2, 1] y[2]^2 & q[1, 2] y[1] y[2] + q[2, 2] y[2]^2 \end{pmatrix}$$

OK, let's try Sanger's network out on our synthetic data.

```
npoints = 10000;
p1 = {};  $\alpha$  = 0.01;
size = 2;
LT = Table[If[i>=j,1,0],{i,size},{j,size}];
Q = Table[.3 RandomReal[], {size}, {size}];
```

```

For[i=1,i<=npoints,i++,
  x = rv; y = Q.x;
  deltaQ = (Outer[Times,y,x] - (LT Outer[Times,y,y]).Q);
  Q = Q +  $\alpha$  deltaQ;
  If[Mod[i,10]==0,
    p1 = Join[p1,{{Q[[1,2]]/Q[[1,1]], Q[[2,2]]/Q[[2,1]] }}]];
];

```

- 7. You may have to adjust the learning constant. It can take 1000's of iterations to converge, so don't give up easily. So run the above cell, graph the result using the input cell below. Then go back up and evaluate the cell above again. Check graphically below, and so forth until you see convergence.

## Graph the evolution of the network weights in terms of the slope

From the model of our synthetic data, the two slopes should be:

**Tan[theta]** and **-1/Tan[theta]**. Let's take a look at the slopes of the transformed coordinates in the list **p1**:

```

gh1=Plot[-1/Tan[theta],{x,0,Length[p1]},PlotRange->{-3,1},
  PlotStyle->{RGBColor[1,.5,0]};
gh2=Plot[Tan[theta],{x,0,Length[p1]},
  PlotStyle->{RGBColor[1,0,1]};

```

```

gh3=ListPlot[Map[#[[2]]&,p1],
  PlotJoined->True,
  PlotStyle->{RGBColor[0,.5,0]};

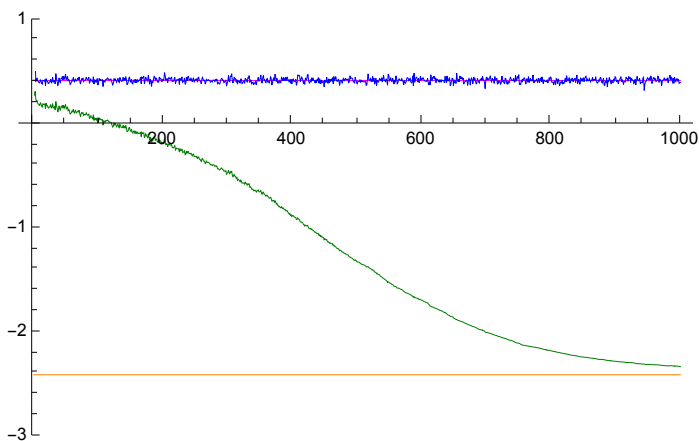
```

```

gh4=ListPlot[Map[#[[1]]&,p1],
  PlotJoined->True,
  PlotStyle->{RGBColor[0,0,1]};

```

```
Show[gh1,gh2,gh3,gh4]
```



```
p1[[Length[p1]]]
```

```
{0.395296, -2.33564}
```

Compare your results with the slopes of gprincipalaxes:

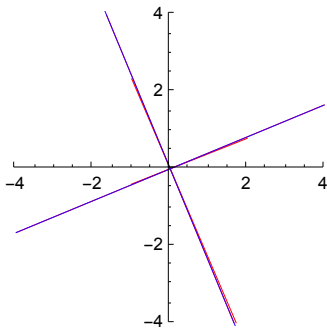
```
slope = (beta / alpha)
(-1 / slope)
0.414214
-2.41421
```

Note that the number of iterations is plotted in multiples determined by the Mod[] function above.

## Graph the transformed axes of the Sanger network and compare them to those from the underlying distribution

Let's plot up the transformation axes of the Sanger network (**gnetwork2**), and compare them with the axes from the population distribution (**gprincipalaxes**), and the calculated principal component axes (**gPCA**):

```
gnetwork2 = Plot[{s p1[[Length[p1]]][[1]], s p1[[Length[p1]]][[2]]}, {s, -1, 2},
  PlotRange -> {{-4, 4}, {-4, 4}}, PlotStyle -> {RGBColor[1, 0, 0]}, AspectRatio -> 1];
Show[gnetwork2, gprincipalaxes, gPCA]
```



## Initialization

```
Off[SetDelayed::write]
Off[General::spell1]
```

## Non-orthogonal Decorrelation Code

Linear neural transform, weight matrix A

Setup coordinate frame graphics

```
lineg=Plot[{}, {x, -1, 1}, PlotRange -> {{0, 1.1}, {0, 1.1}}, AspectRatio -> 1, Frame -> True, FrameLab
```

Define the neural net transform as a function called: **transform[]**:

```
A = 2.5*{{1, -7/11}, {-1/2, 1}};
transform[x_] := A.x;
invtransform[x_] := Inverse[A].x;
```

## Pre-adaptation feature space (x-y): (xp,yp)=Identity.(x,y)

Define adaptation points, and coordinate lines in x-y space

Generate adaptation points

```
t1 = Flatten[Table[{x, y}, {x, 0.1, 1, 0.1}, {y, 0.1, 1, 0.1}], 1];
color = Table[Min[1, invtransform[t1[[i]][[2]]], {i, 1, Length[t1]}];
adaptp = invtransform /@ t1;
preadapaptg =
  Graphics[({RGBColor[color[[#1]], 1 - color[[#1]], 0], Disk[adaptp[[#1]], 0.02]}) & /@
    Range[1, Length[adaptp]], AspectRatio -> Automatic];
```

Generate GRAY y=0.5 coordinate lines in x-y space: Separates RED from GREEN

```
grayyl = Table[{x, 0.5}, {x, 0, 1, 0.01}];
grayylg = Graphics[{PointSize[0.01], RGBColor[0.8, 0.8, 0.8], Point /@ grayyl}];
```

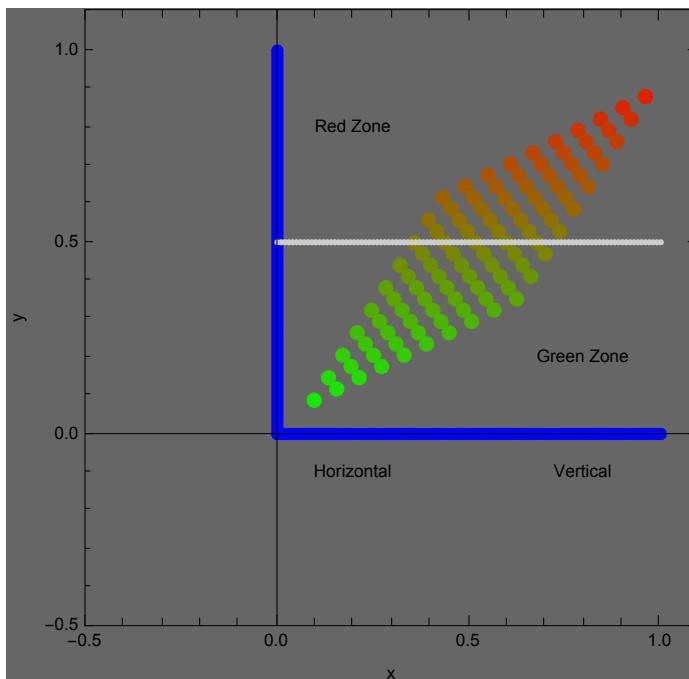
Generate BLUE x-y coordinate lines in x-y space

```
xl=Table[{x,0},{x,0,1,.01}];
yl=Table[{0,y},{y,0,1,.01}];

xlg = Graphics[{PointSize[0.02], RGBColor[0, 0, 1], Point /@ xl}];
ylg = Graphics[{PointSize[0.02], RGBColor[0, 0, 1], Point /@ yl}];
```

Show plot of correlated data in (BLUE) x-y space

```
Show[Lineg, xlg, ylg, preadapaptg, grayylg,
  Graphics[{Text["Red Zone", {0.2, 0.8}], Text["Green Zone", {0.8, 0.2}],
    Text["Horizontal", {0.2, -0.1}], Text["Vertical", {0.8, -0.1}]}],
  PlotRange -> {{-0.5, 1.1}, {-0.5, 1.1}}, Background -> GrayLevel[0.4]]
```



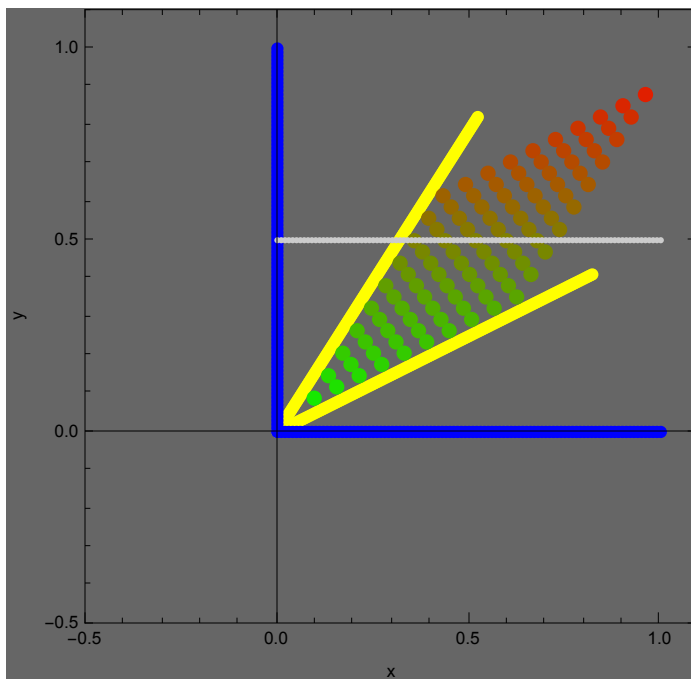
## Plot of pre-adaptation feature space with desired transformed coordinates (YELLOW) in x-y space

Generate YELLOW xp-yp coordinate lines in x-y space

```
xpl=Map[invtransform,Table[{x,0},{x,0,1.4,.01}]];
ypl=Map[invtransform,Table[{0,y},{y,0,1.4,.01}]];

xplg = Graphics[{PointSize[0.02`], RGBColor[1, 1, 0], Point /@xpl}];
yplg = Graphics[{PointSize[0.02`], RGBColor[1, 1, 0], Point /@ypl}];

Show[lineg, xplg, yplg, xlg, ylg, preadapptg, grayylg,
  PlotRange -> {{-0.5`, 1.1`}, {-0.5`, 1.1`}}, Background -> GrayLevel[0.4]]
```



Post-adaptation remapped feature space (xp,yp):  $(xp,yp) = A.(x,y)$

Generate new DARK GRAY y=0.5 coordinate lines in x-y space

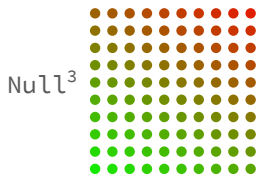
```
newgrayyl = Table[{x, 0.5`}, {x, 0, 1, 0.01`}];
newgrayylg =
  Graphics[{PointSize[0.01`], RGBColor[0.4`, 0.4`, 0.4`], Point /@newgrayyl}];
```

Generate old LIGHT GRAY y=0.5 coordinate lines in x-y space

```
oldgrayyl = Table[{x, 0.5`}, {x, 0, 1, 0.01`}];
oldgrayyl = transform /@oldgrayyl;
oldgrayylg =
  Graphics[{PointSize[0.01`], RGBColor[0.8`, 0.8`, 0.8`], Point /@oldgrayyl}];
```

## Remap RED/GREEN adaptation points

```
(t1 = Flatten[Table[{x, y}, {x, 0.1, 1, 0.1}, {y, 0.1, 1, 0.1}], 1];)
(color = Table[Min[1, invtransform[t1[[i]]][2]], {i, 1, Length[t1]}];)
(adaptp = t1;) (postadaptg =
Graphics[({RGBColor[color[[#1]], 1 - color[[#1]], 0], Disk[adaptp[[#1]], 0.03]} &) /@
Range[1, Length[adaptp]], AspectRatio -> Automatic])
```



## Remap x-y coordinate lines (BLUE) in xp-yp space

```
remapxl=Map[transform,xl];
remapyl=Map[transform,yl];

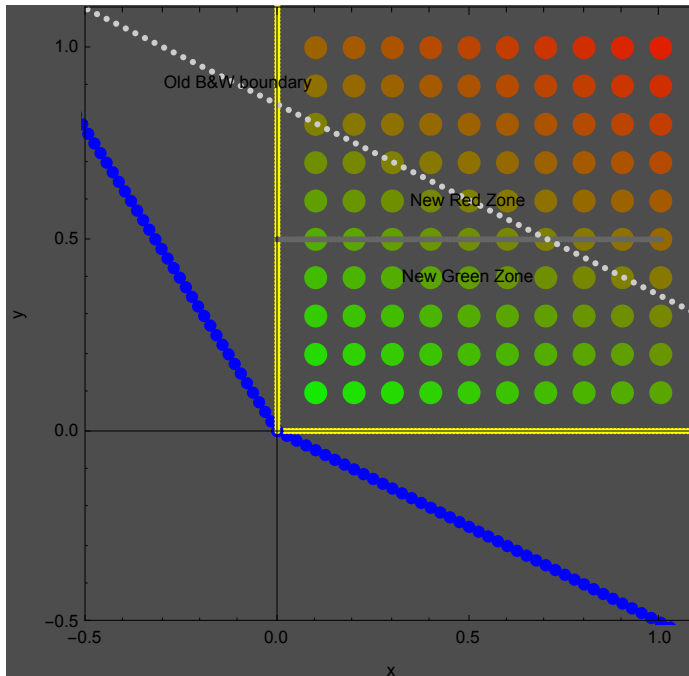
remapxlg = Graphics[{PointSize[0.02], RGBColor[0, 0, 1], Point /@ remapxl}];
remapylg = Graphics[{PointSize[0.02], RGBColor[0, 0, 1], Point /@ remapyl}];
```

## Remap xp-yp coordinate lines (YELLOW) in xp-yp space

```
remapxpl=Map[transform,xpl];
remapypl=Map[transform,ypl];

remapxplg = Graphics[{PointSize[0.01], RGBColor[1, 1, 0], Point /@ remapxpl}];
remapyplg = Graphics[{PointSize[0.01], RGBColor[1, 1, 0], Point /@ remapypl}];
```

```
Show[lineg, postadaptg, remapxlg, remapxplg, remapylg, remapyplg,
newgrayylg, oldgrayylg, Graphics[{Text["New Red Zone", {0.5`, 0.6`}],
Text["New Green Zone", {0.5`, 0.4`}], Text["Old B&W boundary", {-0.1`, 0.9`}]}],
PlotRange -> {{-0.5`, 1.1`}, {-0.5`, 1.1`}}, Background -> GrayLevel[0.3`]]
```



## References

- Atick, J. J., & Redlich, A. N. (1992). What does the retina know about natural scenes? *Neural Computation*, 4(2), 196–210.
- Baldi, P., & Hornik, K. (1989). Neural networks and principal components analysis: Learning from examples without local minima. 2, 53-58.
- Barlow, H. (1990). Conditions for versatile learning, Helmholtz's unconscious inference, and the task of perception. *Vision Research*, 30(11), 1561-1572.
- Barlow, H. B., & Foldiak, P. (1989). Adaptation and decorrelation in the cortex. In C. Miall, R. M. Durban, & G. J. Mitchison (Ed.), *The Computing Neuron* Addison-Wesley.
- Bell A. J. and Sejnowski T. J. . An information-maximization approach to blind separation and blind deconvolution. *Neural Computation*, 7(6):1129-1159, 1995.
- Roweis, S., & Ghahramani, Z. (1999). A unifying review of linear gaussian models. *Neural Comput*, 11(2), 305-45.
- Barlow, H. B., & Foldiak, P. (1989). Adaptation and decorrelation in the cortex. In C. Miall, R. M. Durban, & G. J. Mitchison (Ed.), *The Computing Neuron* Addison-Wesley.
- Clifford, C. W., Wenderoth, P., & Spehar, B. (2000). A functional angle on some after-effects in cortical vision. *Proc R Soc Lond B Biol Sci*, 267(1454), 1705-1710.
- Clifford, C. W. (2002). Perceptual adaptation: motion parallels orientation. *Trends Cogn Sci*, 6(3), 136-143.
- Hotelling, H. (1933). Analysis of a complex of statistical variables into principal components. *J. Educ. Psych.*, 24, 417-441,498-520



- McCollough, C. (1965, 3 September 1965). Color Adaptation of Edge-Detectors in the Human Visual System. *Science*, 149, 1115-1116.
- Oja, E. (1982). A simplified neuron model as a principal component analyzer. *Journal of Mathematical Biology*, 15, 267-273
- Olshausen, B. A., & Field, D. J. (1996). Emergence of simple-cell receptive field properties by learning a sparse code for natural images. *Nature*, 381, 607-609.
- Ratliff, C. P., Borghuis, B. G., Kao, Y. H., Sterling, P., & Balasubramanian, V. (2010). Retina is structured to process an excess of darkness in natural scenes. *Proceedings of the National Academy of Sciences of the United States of America*, 107(40), 17368–17373.
- Rhodes, G., Jeffery, L., Watson, T. L., Clifford, C. W., & Nakayama, K. (2003). Fitting the mind to the world: face adaptation and attractiveness aftereffects. *Psychol Sci*, 14(6), 558-566.
- Sanger, T. (1989). Optimal unsupervised learning in a single-layer linear feedforward neural network. *Neural Networks*, 2, 459-473.
- Schwartz, O., & Simoncelli, E. P. (2001). Natural signal statistics and sensory gain control. *Nat Neurosci*, 4(8), 819-825.
- Stilp, C. E., & Kluender, K. R. (2012). Efficient Coding and Statistically Optimal Weighting of Covariance among Acoustic Attributes in Novel Sounds. *PLoS ONE*, 7(1), e30845. doi:10.1371/journal.pone.0030845.t001
- Zhu, L., Lin, C., Huang, H., Chen, Y., & Yuille, A. (2008). Unsupervised structure learning: Hierarchical recursive composition, suspicious coincidence and competitive exclusion (pp. 1–14). Presented at the European Conference on Computer Vision, 2008.
- Zhu, L., Chen, Y., & Yuille, A. (2011). Recursive Compositional Models for Vision: Description and Review of Recent Work. *Journal of Mathematical Imaging and Vision*, 41(1-2), 122–146. <http://doi.org/10.1007/s10851-011-0282-2>