

Introduction to Neural Networks

Regression, least squares, gradient descent, Widrow-Hoff & backprop

Initialization

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

Introduction

Last time

Focused on classification:

Turning linear networks into classifiers with a hard threshold.

Perceptrons, perceptron learning rule.

Linear separability

Today

Linear regression & brain-style learning -- the Widrow-Hoff learning rule

Multiple layers of weights -- backprop learning, an example of non-linear regression

Linear regression and Widrow-Hoff learning

Introduction

We have been studying a linear matrix model of memory based on the storage of connection weights that follow a particular Hebbian rule. We have studied the "psychology" of some operations of linear algebra and have seen some interesting parallels to human memory, such as interference and pattern reconstruction.

We've learned that linear networks can be configured to be supervised or unsupervised learning devices. But linear mappings are severely limited in what they can compute. With a threshold non-linearity, the linear neuron model becomes a simple perceptron unit that makes decisions.

But we are now limited by the linearity of the decision surface, i.e. by the fact that it is a hyperplane. More complex networks can be built by adding layers making multiple layer networks ("multi-layer perceptrons"), but then we have another problem: how can we learn the weights in this more complicated setting?

Overview of our strategy

First we return to strictly linear networks (no threshold) which will allow us to introduce techniques (gradient descent learning in the context of finding the weights of a linear matrix transformation). We will learn concepts that will generalize to non-linear networks with multiple layers of weights. In particular, the technique of gradient descent will lead us to the Widrow-Hoff learning rule. By treating the linear case first, we will be able to see how the Widrow-Hoff learning rule relates to the classic problem of linear regression.

Learning, memory and generalization can be viewed from the point of view of non-linear statistical regression. The idea is to treat learning as an attempt to fit past input/output associations into a model that can be used both for recall and for future generalization.

The problem of regression in statistics is: given a set of vector inputs $\{\mathbf{x}_i\}$, and a set of corresponding vector "target" outputs $\{\mathbf{t}_i\}$, one tries to find a transformation \mathbf{W} that will map $\mathbf{x} \rightarrow \mathbf{t}$ as closely as possible over the data available. For this we need a model for \mathbf{W} , and a measure of *goodness of fit*. At first, we will assume a linear model for \mathbf{W} , so for the discrete case, \mathbf{W} is a matrix. Our measure of goodness of fit will be the *sum of the squared differences between predicted output and actual output*. In the linear case, this is linear *least squares regression*.

Then we'll see how the Widrow-Hoff rule provides a clue how to generalize learning to non-linear, multiple layer networks with smooth non-linear squashing functions, with the same measure of goodness of fit. This is the error back-propagation algorithm.

In later lectures, we will develop a more powerful notion of goodness of fit -- *probability* of parameters given experience with past inputs and outputs, for which the sum of squared differences is a special case.

More terminology regarding supervised learning

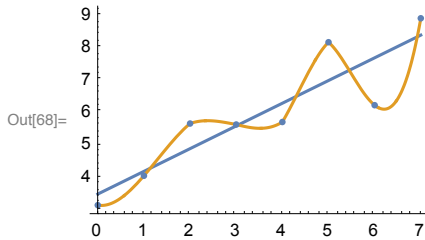
Consider **supervised learning**. We have a "training set" $\{\mathbf{x}_i, \mathbf{t}_i\}$ representing input vectors \mathbf{x}_i , and target outputs \mathbf{t}_i . The training set "samples" the larger space of possible input/output pairs $\{\mathbf{x}, \mathbf{t}\}$. We would like to learn a general mapping: $T: \mathbf{x} \rightarrow \mathbf{y}$ in such a way that $T (= \mathbf{W})$ is a good fit to the training data (i.e. \mathbf{y} is close to \mathbf{t}), and generalizes well to novel inputs. The set of target data is the feedback for the "teacher".

In general, feedback can vary in the degree of precision it provides for learning. It can specify whether the mapping is correct or not. Or the feedback can provide information as to how far off the map T 's prediction, i.e. how far is $T[\mathbf{x}_i]$ from \mathbf{t}_i ?

After training, one could require that T always maps members of the training set to exactly the target members, and generalizes appropriately for other inputs. This means that the learning should be *consistent*. **Interpolation** between data points on a graph is an example of consistent learning. A specific example of interpolation would be drawing lines connecting data points on a graph.

Or, we may require that the T maps the original members of the training set to outputs \mathbf{y} , that are close to the original targets \mathbf{t} . This is called **approximation** learning.

```
In[65]:= gdata = ListPlot[data = Table[{x, x + 4 * RandomReal[]}, {x, 0, 7}]];
approx = Fit[data, {1, x}, x];
interp = Interpolation[data];
Show[{Plot[{approx, interp[x]}, {x, 0, 7}], gdata}, ImageSize -> Small]
```



Linear regression is an example of *approximation* learning. Approximation doesn't necessarily exactly fit the data points, but the goal is that it should come close, but not at the expense of failing to generalize to new data.

In addition to remembering, regression *generalizes*, i.e. provides the basis for *prediction*. So novel input values get mapped to predicted outputs based on past "experience". We've already seen how linear heteroassociative learning does this. Later we will look at how the type of generalization affects "bias" and "variance".

Let's take a close look at linear regression, a fundamental case of approximation learning.

Least squares regression - linear models

The idea behind least squares regression is given a set of N training pairs $\{x_i, t_i\}$, where i runs from 1 to N , we would like to find a function that given an input \mathbf{x} , the function approximates well the output, i.e. $\mathbf{W} \cdot \mathbf{x} \sim \mathbf{t}$. If the function reproduces the association between input and output that it has seen before, this is like "remembering".

An example of linear regression

A fundamental assumption that affects the modeling and performance of any learning system which needs to generalize is that there is an underlying structure to the data--the relationship between associative pairs is not arbitrary.

When trying to understand how a relationship can be learned between a set of two patterns, it helps to have some understanding of the structure of the relationship. For example, one shouldn't try to fit a straight line to data when there is reason to believe that the underlying process produces data on a circle. We will study this more when we learn about the "bias/variance dilemma". So let us assume that the data have an underlying structure that we are going to try to discover or approximate using \mathbf{W} .

We will study a simple "toy" problem that has the following very specific generative structure. The inputs are randomly located points on a 2D plane, and the outputs are heights above these points. The outputs lie approximately on a planar surface that runs through the origin and whose orientation is characterized by two parameters, \mathbf{a} and \mathbf{b} .

It may seem like overkill, but we are going to estimate \mathbf{W} for the same set of data in four different ways.

The first two are drawn from standard linear algebra (a least squares solution using transpose and inverse, and the second using the pseudoinverse). The third introduces the method of "gradient descent" a common numerical technique which is used in many applications. We will use it later in several contexts. The fourth method is the most relevant to neural network theory--we estimate \mathbf{W} using a biologically plausible learning rule (the Widrow-Hoff rule).

Generative model: Synthetic training pairs

Although we will illustrate our examples with small dimensions, everything we do generalizes to higher dimensional inputs and outputs, and in fact the demonstration code below will work with higher dimensional input/output vectors.

Let x_1 and x_2 be the (scalar) inputs, and y be the (scalar) output:

$$\{x_1, x_2\} \rightarrow y,$$

and assuming the mapping is a plane through the origin and additive noise

$$y = a x_1 + b x_2 + \text{noise}$$

This is a many-to-one mapping.

We want to learn the parameters a and b from the training pairs: $\{(x_1, x_2), y\}$. (These parameters could be thought of as neural weights.)

```
In[69]:= rsurface[a_, b_] := N[Table[{x1 = 1 RandomReal[],  

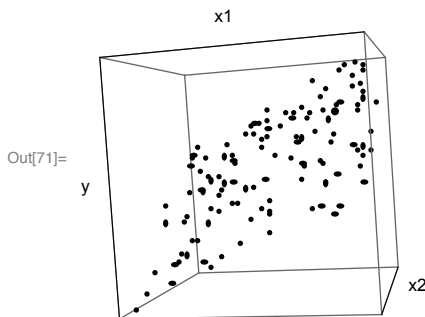
x2 = 1 RandomReal[], a x1 + b x2 + 0.5 RandomReal[] - 0.25}, {120}], 2];  

data = rsurface[2, 3];  

Graphics3D[Point /@ data, AxesLabel → {"x1", "x2", "y"},  

Axes → True, Ticks → None, BoxRatios → {1, 1, 1}, ImageSize → Small]  

Dimensions[data]
```



Out[72]= {120, 3}

You can also use: **ListPointPlot3D**[**data**]

yy represents the 1 dimensional y values, and **xx** the 2-dimensional input values $x = \{x_1, x_2\}$:

```
In[73]:= yy = data[[All, 3]];  

xx = data[[All, 1 ;; 2]];
```

I. Least squares regression to find \mathbf{W}

So let's assume we want to find a matrix \mathbf{W} that will come close to reproducing values \mathbf{y} , given inputs \mathbf{x} . Of course, because we generated the data, we know the underlying structure and what the matrix \mathbf{W} should be. The correct answer should be a 1x2 matrix = $\{\{2,3\}\}$. But let's assume we don't know the answer, and want to discover the weights from \mathbf{y} and \mathbf{x} .

Imagine you have knobs that you can twiddle that let you adjust every element of the matrix \mathbf{W} . For particular settings of these knobs, you can calculate $\mathbf{W} \cdot \mathbf{x}_i$ for any input/output pair. In general you won't be so lucky as to have $\mathbf{y}_i = \mathbf{W} \cdot \mathbf{x}_i$. So you take the opportunity to adjust the knobs to produce \mathbf{W}' so that $\mathbf{W}' \cdot \mathbf{x}_i$ is closer to \mathbf{y}_i . One measure of "how close" is

$(\text{Norm}[\mathbf{y}_i - \mathbf{W} \cdot \mathbf{x}_i])^2 = \|\mathbf{y}_i - \mathbf{W} \cdot \mathbf{x}_i\|^2 = \text{EuclideanDistance}[\mathbf{y}_i, \mathbf{W} \cdot \mathbf{x}_i]^2$. Of course, we usually have lots of data, so we could add up all the discrepancies between what \mathbf{W} predicts, and what the "teacher" has provided in the answers in the training pair (i.e. \mathbf{y}_i for all of the i 's).

In least squares regression, we try to find the values of the matrix that will minimize $e(\mathbf{W})$.

$$e(\mathbf{W}) = \sum_{i=1}^N |\mathbf{y}_i - \mathbf{W} \mathbf{x}_i|^2$$

In many cases our error function will be over a very high dimensional space determined by the number of elements in \mathbf{W} . However, it is useful to get an intuition for the problem in a low dimensional space.

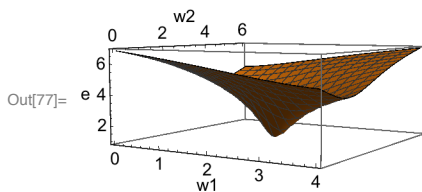
We can get a picture of the total error, $e(\mathbf{W})$, for our synthetic data as a function of the weight parameters w_1 and w_2 :

```
In[75]:= eW[w1_, w2_] := Sum[{yy[[i]] - {w1, w2}.xx[[i]]}.{yy[[i]] - {w1, w2}.xx[[i]]},
  {i, 1, Length[xx], 1}]
```

```
In[76]:= g = Plot3D[eW[w1, w2], {w1, 0, 4}, {w2, 0, 6}, ViewPoint -> {1.78, -2.861, 0.312},
  AxesLabel -> {"w1", "w2", "e"}, ImageSize -> Small];
```

If you plot the log of the error function, you can get a better view of where the minimum is:

```
In[77]:= g = Plot3D[Log[eW[w1, w2]], {w1, 0, 4}, {w2, 0, 6},
  ViewPoint -> {1.78, -2.861, 0.312}, AxesLabel -> {"w1", "w2", "e"}, ImageSize -> Small]
```

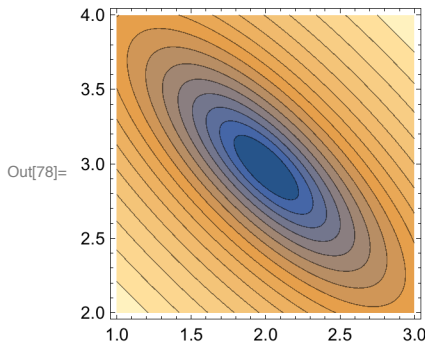


(Remember that $\log[]$ is a monotonic function, so it won't change the location of the minimum.)

Note that because our input data are 2-element vectors, and our output or target values are 1-D, the matrix \mathbf{W} we seek is not square--it has 1 row and 2 columns. So we represent \mathbf{W} above as a vector.

The minimum appears to be near $\{2,3\}$. It is easier to see the location of the minimum using `ContourPlot[]`.

```
In[78]:= gc = ContourPlot[Log[eW[w1, w2]],
  {w1, 1, 3}, {w2, 2, 4}, Contours -> 16, ImageSize -> Small]
```



Note: You can get a rough idea of the coordinates near the middle by moving your computer mouse over the above graphic after initializing the cell below:

```
In[79]:= Dynamic[MousePosition["Graphics"]]
```

Out[79]= None

Of course, in high dimensional spaces, \mathbf{W} is huge, so we can't visualize it, and we can't hope to manually "adjust knobs" or sliders. But we can find the exact location of the minimum by finding "the bottom of the bowl".

Calculus tells us that this should be where the gradient of the error function, e is zero. Although there are a lot of indices to worry about, the familiar analogy from introductory calculus is to find the point at which the derivative of a function $f(x)$ is zero, i.e. the value of x where $f'(x)=0$.

$e(\mathbf{W})$ plays the role of $f(x)$, but \mathbf{W} is a matrix of variables. The solution can be written concisely in terms of vector outerproducts, inversion, and matrix multiplication:

$$e(\mathbf{W}) = \sum_{i=1}^N |y_i - \mathbf{W}\mathbf{x}_i|^2$$

$$\frac{\partial e}{\partial \mathbf{W}} = -2 \sum_{i=1}^N (y_i - \mathbf{W}\mathbf{x}_i)\mathbf{x}_i^T = 0$$

$$\sum_{i=1}^N \mathbf{y}_i\mathbf{x}_i^T - \mathbf{W} \sum_{i=1}^N \mathbf{x}_i\mathbf{x}_i^T$$

$$\mathbf{W} = \sum_{i=1}^N \mathbf{y}_i\mathbf{x}_i^T \left(\sum_{i=1}^N \mathbf{x}_i\mathbf{x}_i^T \right)^{-1}$$

where

```
In[80]:=
```

$$\nabla e = \frac{\partial e}{\partial \mathbf{W}} = \left(\frac{\partial e}{\partial W_{11}}, \dots, \frac{\partial e}{\partial W_{ij}}, \dots \right).$$



So we see that the value of \mathbf{W} where the gradient is zero can be calculated by first adding up two sets of outer products, taking an inverse of one of them, and then multiplying. Let's translate the last line above into *Mathematica* instructions:

```
In[80]:= sumX = Sum[Outer[Times, xx[[i]], xx[[i]], {i, Length[xx]}];
sumYX = Sum[Outer[Times, {yy[[i]]}, xx[[i]], {i, Length[xx]}];
W = sumYX.Inverse[sumX]
```

```
Out[84]= {{2.00502, 2.97538}}
```

The values for \mathbf{W} come close to what we would expect from the structure of our data, a plane with parameters $\{2,3\}$.

2. Pseudoinverse to find \mathbf{W}

For linear algebra aficionados, there is another way of solving the linear least squares regression that uses the “pseudoinverse” of matrix to find the least-squares solution. For a *square* matrix \mathbf{X} , the inverse of \mathbf{X} is chosen so that $\mathbf{X}\mathbf{X}^{-1}$ is equal to the identity matrix, \mathbf{I} . The pseudoinverse, \mathbf{X}^* , of a *rectangular* matrix \mathbf{X} is chosen so that $\mathbf{X}\mathbf{X}^*$ is close to the identity matrix in the sense that the sum of the squares of all of the entries of $\mathbf{X}\mathbf{X}^* - \mathbf{I}$ is least. If interested, you can look up the derivation of \mathbf{X}^* . It can be written as:

$$\mathbf{X}^* = (\mathbf{X}^T \mathbf{X})^{-1} \mathbf{X}^T$$

or as:

$$\mathbf{X}^* = \mathbf{X}^T (\mathbf{X}\mathbf{X}^T)^{-1}$$

For our simple generative model, let's take all of the input vectors \mathbf{x} , and arrange them as columns in a matrix \mathbf{X} :

$$\mathbf{X} = \begin{pmatrix} x_1^1 & x_1^2 & \dots & x_1^N \\ x_2^1 & x_2^2 & \dots & x_2^N \\ \vdots & \vdots & \ddots & \vdots \\ x_N^1 & x_N^2 & \dots & x_N^N \end{pmatrix}$$

Now do the same for the \mathbf{y} 's:

$$\mathbf{Y} = (y_1 \quad y_2 \quad y_3 \dots y_N)$$

And what is the matrix that maps the x 's to the y 's with least squared error? Again, you can look up the derivation or try to derive it yourself. The answer is $\mathbf{X}^*\mathbf{Y}$:

```
In[85]:= Inverse[Transpose[xx].xx].Transpose[xx].yy
```

```
Out[85]= {2.00502, 2.97538}
```

\mathbf{X}^* is the *pseudoinverse* (sometimes called the *generalized inverse*) of the matrix \mathbf{X} . The **PseudoInverse[]** function is built into *Mathematica*, so now that we know about it, we can calculate the least-squares solution for our problem in one simple expression:

```
In[86]:= PseudoInverse[xx].yy
```

```
Out[86]= {2.00502, 2.97538}
```

Question: What are the dimensions of the above PseudoInverse? Verify that PseudoInverse[X].X is the identity matrix with our data.

```
In[87]:= Chop[PseudoInverse[xx].xx]//MatrixForm
```

```
Out[87]//MatrixForm=
```

$$\begin{pmatrix} 1. & 0 \\ 0 & 1. \end{pmatrix}$$

3. Gradient descent

Let's go back to the original global error function that we used above for standard linear least squares regression. There we found the weights that gave the minimum total error by setting the gradient of the error function to zero, and then solving for the weights.

Gradient descent, is a more general way of finding the minimum of an error function. It can be used when the error function is much more complicated (e.g. it isn't quadratic), and there is no linear solution. But it works best for concave error functions, where one doesn't have to worry about multiple minima. Later we will run into this last kind of situation when we try to learn weights via supervised learning for a non-linear network with more than one layer of weights.

The idea is to start off at some location, \mathbf{W}_0 in weight space (\mathbf{W}_0 is just an initial guess), and iteratively move towards the minimum by always taking a step downhill. (See **Graphical illustration of Gradient in Appendix**). Again we take the elements of \mathbf{W} and think of them as defining a point in weight space, i.e. a vector in weight space. We want to find a vector, say \mathbf{W}_{best} , that tells us which direction to move in weight space so that that the error measured at location \mathbf{W}_0 decreases by as much as possible. We don't want to take too big of step or too small, so we need to control the length of \mathbf{W}_{best} , with a scalar η . In general, if we are at the i th point of an iteration, we would like:

$$\mathbf{W}^{i+1} = \mathbf{W}^i + \eta (\mathbf{W}_{\text{best}} \text{ at } \mathbf{W}^i)$$

To derive a rule for \mathbf{W}_{best} , we assume that the error function is a smooth function of \mathbf{W} . The downhill direction is given by the negative gradient of the error function.

(This can be seen by noting that $\nabla \mathbf{e} \cdot \hat{\mathbf{n}} = |\nabla \mathbf{e}| \cos(\theta)$, which is biggest when $\theta = 0$, in other words when $\hat{\mathbf{n}}$ points in the direction of $\nabla \mathbf{e}$.)

```
In[88]:=
```

$$\nabla e = \frac{\partial e}{\partial \mathbf{W}} = \left(\frac{\partial e}{\partial W_{11}}, \dots, \frac{\partial e}{\partial W_{ij}}, \dots \right)$$

Again, the right hand side is just short-hand for

$$\mathbf{W}^{i+1} = \mathbf{W}^i - \eta \left. \frac{\partial e}{\partial \mathbf{W}} \right|_{\mathbf{W}^i}$$

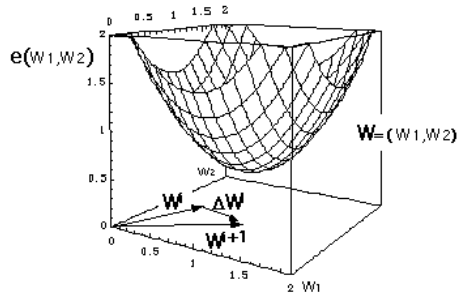
Sometimes the change in weights can be modeled as a dynamical system:

```
In[89]:=
```

$$\frac{d\mathbf{W}}{dt} = - \frac{\partial e}{\partial \mathbf{W}} = -\nabla e$$

We'll see an example of this later.

and then approximate the differential equation by a discrete update rule. $\mathbf{W}_{\text{best}} = -\nabla e$ evaluated at the i th weight space point in the iteration. The figure below gives a graphical view in the small dimensional case of our example problem.



From the expression for the gradient which we wrote earlier in terms of outer products, we can obtain an expression for $\Delta \mathbf{W}$, with η taking the place of Δt :

$$\frac{\partial e}{\partial \mathbf{W}} = -2 \left(\sum_{i=1}^N \mathbf{y}_i \mathbf{x}_i^T - \mathbf{W} \sum_{i=1}^N \mathbf{x}_i \mathbf{x}_i^T \right)$$

```
In[90]:= Clear[wg];
eWgradient[wg_] :=
  Sum[Outer[Times, {yy[[i]]}, xx[[i]]], {i, Length[xx]}] -
  wg.Sum[Outer[Times, xx[[i]], xx[[i]]],
  {i, Length[xx]}];
```

Now define a function T that specifies an updated weight matrix $\mathbf{W}=\mathbf{wg}$. You may have to adjust eta (= η) to make sure the steps are sufficiently small to get convergence, but not so small as to take long to converge.

```
In[92]:=
```

$$\mathbf{W}^{i+1} = \mathbf{W}^i - \eta \left. \frac{\partial e}{\partial \mathbf{W}} \right|_{\mathbf{W}^i}$$

```
Out[92]=
```

$$\mathbf{W}^{i+1} = \mathbf{W}^i - \eta \left. \frac{\partial e}{\partial \mathbf{W}} \right|_{\mathbf{W}^i}$$

```
In[93]:= T[wg_] := wg + eta eWgradient[wg]
```

```

In[94]:= i=0; eta = .005; wglis = {};

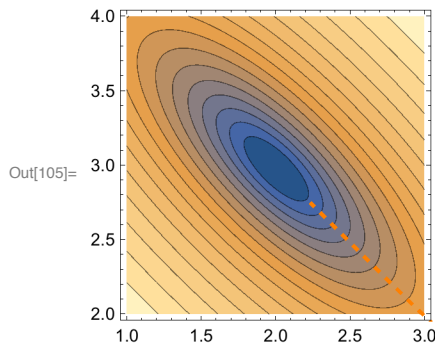
wg={RandomReal[5.0],RandomReal[5.0]};
(*wg={{0,2}};*)

T[wg_] := wg + eta eWgradient[wg];

Clear[wList]; wList = NestList[T, wg, 40];
gl = ListPlot[Flatten[wList, 1], Joined -> True,
  PlotStyle -> {Orange, Dashed, Thick}];

Show[gc, gl]
wList[[40]]

```



Out[106]= {{2.25328, 2.72744}}

We used `NestList[]` above to keep track of all the iteration values so we could graph them, but we could just output the final value after, say 50 steps, using `Nest[]`

```

In[107]:= w1 = Nest[T, wg, 50]
Out[107]= {{2.14695, 2.83364}}

```

NOTE: There are often better algorithms for finding minima depending on the problem. For example, see the conjugate gradient method.

4. "Brain-style" learning: Iterative Widrow-Hoff learning to estimate W

So far so good. But there are a several problems. First, we are interested in brain-style computation. Based on what we think we know about neurons, how could the brain compute transposes, do matrix inversion and multiplication?

Second, when we learn we don't seem to gather information on a whole set of training pairs, and then suddenly build a memory matrix. A more plausible assumption is that learning is *incremental, trial by trial*, rather than *batch*.

Another problem is purely computational. What if the dimensionality of the vectors is really big? It is computationally expensive to invert large matrices. The above gradient descent procedure avoids the problem of inverting large matrices, but it involved computing a global error term over all the training pairs. We would like a method which would learn a regression mapping without having to store all the training pairs with the accompanying computation of a global error term. Instead, we'd like *to compute an error term incrementally, trial-by-trial*.

Can we discover the mapping \mathbf{W} in such a way so as to be biologically plausible, and avoid having to invert a large matrix? The *Widrow-Hoff* rule provides an answer. The basic idea behind *Widrow-Hoff* learning is to update \mathbf{W} iteratively with each new training pair. Let's start off with an arbitrary \mathbf{W} , find out which direction we would have to go in weight space to reduce the discrepancy between what \mathbf{W} tells us \mathbf{x} should map to and what it actually is, namely \mathbf{y} . We take our clue from gradient descent, but apply it each time a new training pair comes along. We recompute the error term each time a new training pair comes along.

$$e(\mathbf{W}) = |\mathbf{y}_i - \mathbf{W} \mathbf{x}_i|^2$$

$$\frac{\partial e}{\partial \mathbf{W}} = -2(\mathbf{y}_i - \mathbf{W} \mathbf{x}_i) \mathbf{x}_i^T$$

$$\frac{\partial e}{\partial \mathbf{W}} = -\frac{d\mathbf{W}}{dt}$$

$$\mathbf{W}^{i+1} = \mathbf{W}^i + \eta_i (\mathbf{y}_i - \mathbf{W}^i \mathbf{x}_i) \mathbf{x}_i^T$$

Let's try out this update rule on our synthetic training pairs.

```
In[130]:= ww1 = {{0,0}}; ww1list = {}; ww2list = {};
```

To visualize how the learning progresses, we are using `ww1list` and `ww2list` to store the first and second weights respectively, for each learning iteration.

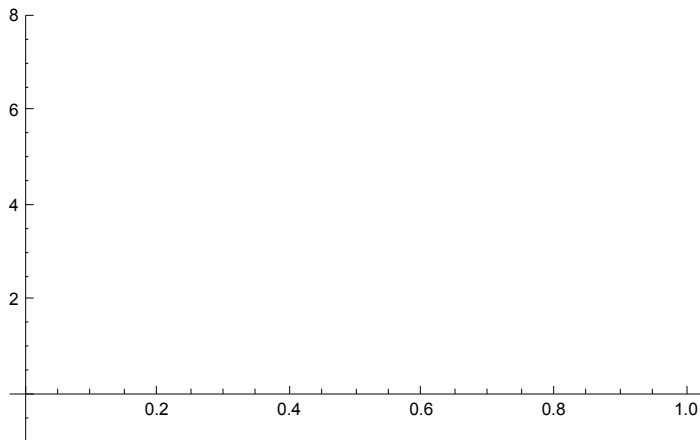
```
In[131]:= eta = 5;
```

```
In[147]:= i=0;
While[i<Length[data],
  ++i;
  in = {data[[i,1]],data[[i,2]]} ; out = {data[[i,3]]};
  ww1 = ww1 + (eta/i) Outer[Times,(out - ww1.in),in];
  ww1list = Append[ww1list,ww1[[1,1]]];
  ww2list = Append[ww2list,ww1[[1,2]]];
];
ww1
```

```
Out[149]= {{2.08527, 2.8917}}
```

Keep executing the above input cell. You may have to run through the above loop several times before reaching stable convergence. We can plot up the two weights as a function of the iteration number to see how the Widrow-Hoff rule for weight modification eventually leads to two stable weights:

```
Dynamic[ListPlot[{ww1list, ww2list},
  PlotRange -> {-1, 8}, AxesOrigin -> {0, 0}, Joined -> True]]
```



Note that we introduced a small change to the learning “constant” η , replacing it by (η/i) —it is no longer constant, but starts of large and decreases with each iteration.

Memory recall

We've seen several ways of finding the weights of a matrix that will approximately reproduce an output, given an input it has seen before. Let's try it out.

So in order to "recall" a response, from an input $\mathbf{x}[[6]]$, we run it through the "network" memory matrix $\mathbf{w1}$:

```
In[150]:= w1 . xx [ [22] ]
```

```
Out[150]= { 1.7484 }
```

And we can check to see how well it recalls:

```
In[151]:= yy [ [22] ]
```

```
Out[151]= 1.57438
```

Recall that the regression model of memory should generalize--*interpolate* and *extrapolate*. We saw a 3-dimensional example of interpolation in a graphical demonstration in the last lecture. But we can also extrapolate.

For example, $\{11, 15\}$ wasn't in the training set, but the expected output is:

```
w1 . {11, 15}
```

```
{ 66.1758 }
```

The network has "learned" a surface, (e.g. for one particular training set, $\mathbf{y} = 1.9\mathbf{x1} + 3.07\mathbf{x2}$ -- the coefficients of $\mathbf{x1}$ and $\mathbf{x2}$ will vary with each new random set of data) through the points specified in the training set.

A crucial point is that the linear associator will try to fit a plane (or hyperplane) through the data. If the data are not fitted well by that model, then the memory and generalization will not be good.

An obvious generalization is to fit hypersurfaces, rather than hyperplanes. And that is the direction we

will head. But first let us look at linear regression from a point of view that you may not have thought of before.

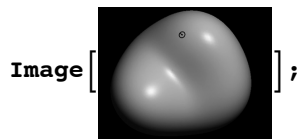
Underconstrained problems and redundancy

Learning a one-to-many mapping:
input dimensionality < output dimensionality

Motivation from vision

This example is very similar to the preceding example, except that we are going to try to learn **2D responses from 1D stimulus inputs**. At first this doesn't seem to make sense, because the mapping from 1D to 2D in general would be expected to be underconstrained by the data. But this is exactly what is needed in certain kinds of inverse problems that are said to be "ill-posed" (Poggio et al., 1985). They are made "well-posed" when there is some underlying regularity (sometimes called "smoothness") in the high-dimensional output space (e.g. Kersten et al., 1987).

In the figure below, the data reaching the eye can be represented as an intensity value at each pixel. However, you perceive a shape, and can even adjust the orientation of a "gauge figure" consisting of a circle tangent to the surface at a point, and a vector perpendicular to the tangent plane (like a "thumb tack", cf. Nefs et al., 2006). This is an example of the shape-from-shading problem in vision. The gauge figure has two degrees of freedom at each point in the image.



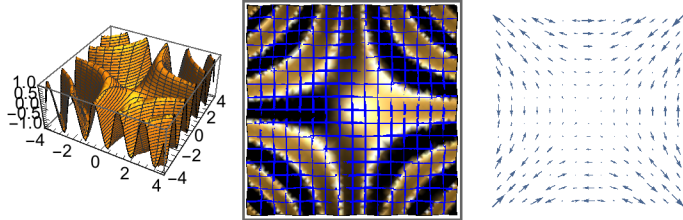
Thus "shape-from-shading" can be formulated as a problem of mapping N pixel intensities to N surface normals, where each surface normal is specified by two numbers--thus the mapping goes from N to $2N$. E.g. given a 15×15 grid of intensities, estimate the 15×15 grid of 2D vectors that are normal to each point of the underlying surface. Thus one might start off with 225 numbers as input, but outputs 550 numbers.

```
GradientFieldPlot[f_, {x_, xmin_, xmax_},
  {y_, ymin_, ymax_}, opts : OptionsPattern[]] :=
  VectorPlot[Evaluate[D[f, {{x, y}}]], {x, xmin, xmax}, {y, ymin, ymax}, opts]
```

```

ssurface[x_, y_] := Sin[x y];
ggss = Plot3D[ssurface[x, y], {x, -4, 4}, {y, -4, 4}, MeshFunctions -> {#1 &, #3 &}];
gsss = Plot3D[ssurface[x, y], {x, -4, 4}, {y, -4, 4}, ViewPoint -> {0, 0, 30},
  Axes -> False, Lighting -> {"Directional", White, {{1, 0, 1}, {1, 1, 0}}},
  Mesh -> 15, MeshStyle -> Directive[Thin, Blue]];
gf = GradientFieldPlot[ssurface[x, y], {x, -4, 4}, {y, -4, 4}, Frame -> False];
Show[GraphicsRow[{ggss, gsss, gf}]]

```



The left panel illustrates an underlying surface that generates the image intensities shown in the middle. The shape-from-shading problem tries to take the graylevel values at points (approximated by grid boxes) in the middle panel, and output a 2-D vector for that grid, as shown in the right-hand panel.

Let's try a simple version of learning a mapping from a low to high-dimensional case. We'll use low-dimensional synthetic data whose generation process we'll keep hidden for now.

Synthetic data -- try not to look closely...

Here is the generative model:

output = {x2, y} = {a x1, b x1 + noise}

We want to learn the parameters {a,b} = {w1,w2}, from training pairs of inputs and outputs:
{x1, {{x2, y}}}

```

r3Dline[a_, b_] :=
  N[Table[{x1 = 1 RandomReal[], a x1, b x1 + 0.5 RandomReal[] - 0.25}, {30}], 2];
data = r3Dline[2, 3];
xx = data[[All, 1]];
yy = data[[All, 2 ;; 3]];

```

So the input data is a list of 1D scalars, and the output data a list of 2D vectors. Let's apply the Widrow-Hoff algorithm to learn the relationship between the input stimuli, xx, and the output responses, yy.

Iterative Widrow-Hoff learning

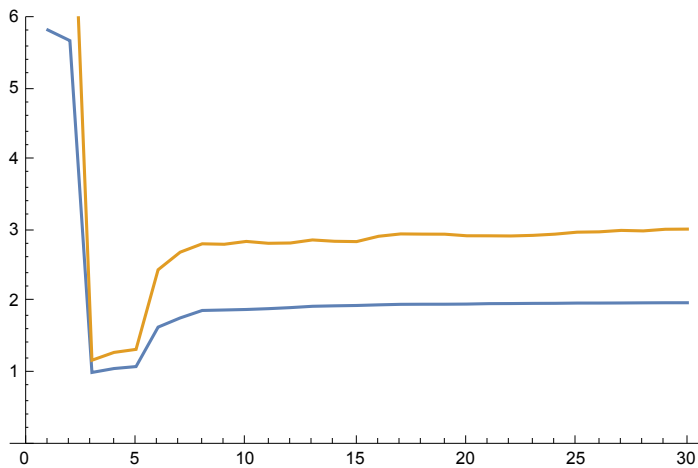
```

w1 = {{0},{0}}; w1list = {}; w2list = {}; eta = 4.0;
Dimensions[w1];

i=0;
While[i<Length[data],
  ++i;
  out = {data[[i,2]],data[[i,3]]}; in = {data[[i,1]]};
  w1 = w1 + (eta/i) Outer[Times,(out - w1.in),in];
  w1list = Append[w1list,w1[[1]][[1]]];
  w2list = Append[w2list,w1[[2]][[1]]];
];
w1
{{1.98563}, {3.023}}

```

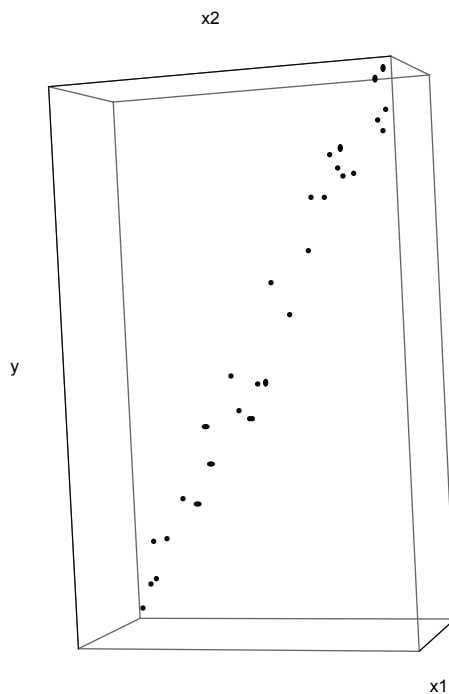
```
ListPlot[{w1list, w2list}, PlotRange -> {0, 6}, AxesOrigin -> {0, 0}, Joined -> True]
```



Visualizing the underlying structure

So why does it work to learn to predict a 2D value from a 1D input? The reason, of course, is that the underlying structure of the output data is highly constrained, and in fact lies close to a straight line in 3-space.

```
Show[Graphics3D[Point /@ data], Axes -> True,
      AxesLabel -> {"x1", "x2", "y"}, Ticks -> None, ImageSize -> Medium]
```



Here is the generative model for our synthetic data:

$x1 \rightarrow \{x2, y\}$,

and given the form is a straightline through the origin with some additive noise, the output $\{x2, y\}$ is given by:

$$\{x_2, y\} = \{a x_1, b x_1 + \text{noise}\}$$

An input, a scalar quantity x_1 specifies a plane. We've found a line that intersects the plane and for which points in our data set are close to that line. The intersection point gives us the coordinates $\{x_2, y\}$ for our regression fit, or "recall".

We learned the parameters $\{a, b\}$, corresponding to the weights, from training pairs of inputs and outputs: $\{x_1, \{x_2, y\}\}$

```
r3Dline[a_, b_] := N[Table[{x1 = 1 RandomReal[], a x1, b x1 + 0.5 RandomReal[] - 0.25}, {30}], 2];
data = r3Dline[2,3];
```

It produced the coordinates of a noisy line in 3-space.

Making estimates of high-dimensional functions from lower dimensional inputs happens a lot, and is common to many so called "early vision" problems (Poggio et al., 1985). For more recent work on learning low-level vision, see Freeman et al., 2000.

Introduction to multi-layer nets

This was an earlier exercise. For linear networks, no computational power is gained by having extra layers:

$$\begin{aligned} y_1 &:= W_0 y_0; \\ y_2 &:= W_1 y_1; \end{aligned}$$

is equivalent to:

$$y_2 := W_1(W_0 y_0) := W_1 W_0 y_0 := W_3 y_0;$$

where W_3 is just another matrix. However, if the inner product is followed by a non-linear transformation, then concatenating layers of neural elements allows one to compute more complex transformations:

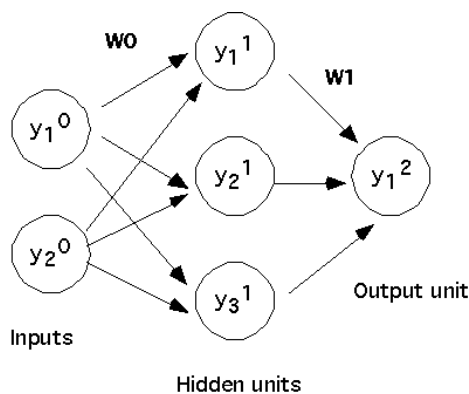
$$y_2 := f[W_1 \cdot f(W_0 y_0)];$$

where $f[\]$, for example, is a sigmoidal "squashing function" or *logistic function*:

$$f[x_] := 1 / (1 + \text{Exp}[-(x - .5)]);$$

NOTE: Notation alert--we are using f or f_2 to represent a scalar squashing function, not an input vector.

Exercise: 2 input units, 3 hidden units, 1 output unit




```
In[152]:= Clear[y2, f2];
          f2[x_] := N[1 / (1 + Exp[-(x - .5) * 100])];
          W1 = {{-13.2328, 6.06398, 6.04958}};
          W0 = {{-0.937564, -1.09841}, {-9.95589, 3.85642}, {3.79034, -10.1737}};
          y2[y0_] := Chop[f2[W1.f2[(W0.y0)]]];
```

What logical function does this network compute?

```
In[157]:= {y2[{0, 0}], y2[{0, 1}], y2[{1, 0}], y2[{1, 1}]}
```

```
Out[157]= {{0}, {1.}, {1.}, {0}}
```

What logical function will the above network compute with the following weights?

```
In[158]:= W0 = {{1.9009195689645, 1.9251997561975753}, {2.177312432892121, 2.198392078532626},
               {-2.8297064197335953, -2.8648165030618693}};
          W1 = {{1.1241018372719358, 1.7600034188676417, -10.048562036295957}};
```

```
In[160]:= Clear[y2];
          y2[y0_] := Chop[f2[W1.f2[(W0.y0)]]];
```

```
In[162]:= {y2[{0, 0}], y2[{0, 1}], y2[{1, 0}], y2[{1, 1}]}
```

```
Out[162]= {{0}, {1.}, {1.}, {1.}}
```

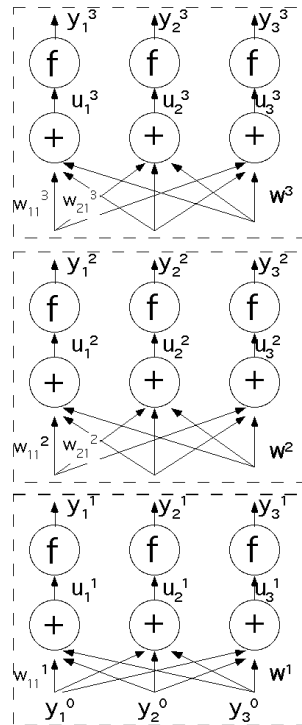
The demonstration illustrates that we can compute XOR as well as OR using a network with more than one layer of weights.

Multi-layer nets: Learning and the error back-propagation algorithm

So we've seen that the non-linear generic neuron allows increased computational power when we add an extra layer of weights.

Suppose we have a multi-layer network with 3 layers of weights. The output from the first layer is:

$\mathbf{y}^1 = f[\mathbf{u}^1] = f[\mathbf{W}^1 \cdot \mathbf{y}^0]$. The output of the second layer is: $\mathbf{y}^2 = f[\mathbf{u}^2] = f[\mathbf{W}^2 \cdot \mathbf{y}^1]$. And so forth. To simplify a bit, we'll assume each layer has the same number of units and weights:



The problem is: how to assign the weights? We'd like to find them through supervised learning.

For any complex system that is required to achieve a target goal, for the system to work, each component must contribute towards the goal. If the goal is not met, one has to figure out which component needs to be fixed. If the goal is met, each component contributed something towards the goal. How does one assign the credit for success or failure to a component?

This problem is called the *credit-assignment problem*. In particular, for the above multi-layer network, how do we adjust the weights, which intrinsics to the whole network, in a way appropriate for learning a given input/output relation?

$$\{x^p, t^p\}, y^0 = x^p, \text{ training pairs } p = 1, \dots, M$$

For a given input $\{y^0 = x^p\}$, we feed forward the information to the last layer (layer L) to produce an output $\{y = y^L\}$. We compare the output to the target value supplied by the "teacher" $\{t = t^p\}$, and compute the error as the sum of squared differences:

$$E(\{w_{ij}^\lambda\}) = \sum_{k=1}^N (y_k(y^0; w_{ij}^\lambda) - t_k)^2$$

where the sum is over all N output units. (For simplicity, we left out the superscript p in t_k^p . The subscript k in t_k^p means the component of the corresponding vector of activity t^p .) λ indexes the weight layers going from $\lambda=1$ to L. Note that the y's at any point after the input depend on the u's (the weighted sum before the non-linearity), each of which in turn depends on all the w_{ij}^λ 's before it.

The trick is to find out how to assign credit (and blame) for the error to each of the weights. Gradient descent provides the answer. Adjust the weights such that:

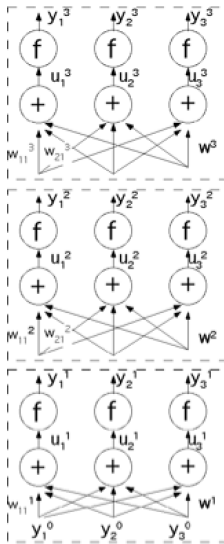
$$\text{new } \mathbf{w}_{ij}^\lambda = \text{previous } \mathbf{w}_{ij}^\lambda + \Delta \mathbf{w}_{ij}^\lambda$$

where

$$\Delta \mathbf{w}_{ij}^\lambda = -\eta \frac{\partial E}{\partial \mathbf{w}_{ij}^\lambda}$$

Again, this formula means that if we calculate the gradient $\nabla E = \partial_{\mathbf{w}_{ij}^\lambda} E$, its negative direction points in the direction of steepest descent. Thus if we update the weight vector to point in this direction, then at the next step we should in general have a lower value of E. I.e.

$$E(\text{new } \mathbf{w}_{ij}^\lambda) < E(\text{previous } \mathbf{w}_{ij}^\lambda).$$



The appendix shows how to calculate a weight adjust with just one layer of weights. This case is related to the area of Generalized Linear Regression (not to be confused with the General Linear Model!).

We won't go through the math here for the multiple layer case--it is really just a very complicated application of taking derivatives where the complication arises in keeping track of all the indices. Instead we assume we have the rule that tells us what the effective error term--the delta value-- and adjustment that needs to be for each weight, and see how it gets put together to converge on weights that tend to minimize the overall output errors.

Summary of backprop algorithm

1. Initialize the weights to small random values
2. Pick a pattern from the input/output collection, say the p th pattern: $\{x^p, t^p\}$. Run the input vector, x^p , feedforward through the network. This will generate a set of values u_i^L and y_i^L in all the nodes of the network. Keep in mind that u_i^L is the linear weighted sum of its inputs arriving from layer $L-1$: $u_i^L = \sum_{j=1}^3 w_{ij}^L y_j^{L-1}$. And $y_i^L = f(u_i^L)$.

Calculate a delta term (analogous to the Widrow-Hoff rule) for the output layer L :

$$\partial_i^L = (t_i^p - y_i^L(x^p)) f'(u_i^L)$$

Note why it is important to have an expression for the *derivative* of the squashing function $f()$. The derivative has a particularly nice form when $f(u) = 1/(1+e^{-u})$. (see Appendix, Logistic function - a smooth, differentiable non-linearity)

3. Propagate the errors back through the layers:

$$\partial_i^\lambda = f'(u_i^\lambda) \sum_{k=1}^N \partial_k^{\lambda+1} w_{ki}^{\lambda+1} \quad \lambda = L-1, \dots, 1$$

...the error back propagation or "back prop" part.

4. Calculate weight adjustments (analogous to the outerproduct part of the Widrow-Hoff) and update using:

$$\Delta w_{ij}^\lambda = \eta \partial_i^\lambda y_j^{\lambda-1}$$

5. Repeat steps 2 to 4 until convergence.

One can accumulate the weight adjustments for each training pair, and then update them all at once.

$$(\Delta w_{ij}^\lambda)_p = \eta \partial_i^\lambda y_j^{\lambda-1}$$

$$\Delta w_{ij}^\lambda = \sum_{p=1}^M (\Delta w_{ij}^\lambda)_p$$

In practice, updating the weights after each training pair often works better than accumulating a bunch of input/output pairs, and then computing the cumulative global error. The reason is that by randomly sampling a training pair, the "descent" may actually climb the global error function defined by the entire set. As we will see later with the Boltzmann machine, occasional climbing is useful to avoid local minima.

There are a number of derivations and illustrations on the web that you might find useful. See for example, http://galaxy.agh.edu.pl/~vlsi/Al/backp_t_en/backprop.html

Also, Andrew Ng has several excellent coursera videos on error backpropagation.

Backprop simulation example: XOR

With appropriate weights, 2 weight layers with 3 hidden units can solve the XOR problem. But this is still a tough problem to learn, mainly because it requires that two very different inputs map to the same output. See the supplementary material for a Mathematica demo that learns the weights for solving the XOR problem.

Next

Relationship of neural networks to machine learning

Problems of slow learning and over-fitting in MLPs.

Non-linear networks & "energy": Hopfield

Probabilistic models. From energy to probability. Boltzmann machines

Appendix

Pseudoinverse solution for the second (underconstrained redundancy) case above

```
r3Dline[a_, b_] :=
  N[Table[{x1 = 1 RandomReal[], a x1, b x1 + 0.5 RandomReal[] - 0.25}, {30}], 2];
data = r3Dline[2, 3];
xx = data[[All, 1]];
yy = data[[All, 2 ;; 3]];
```

We can calculate the memory matrix using the **PseudoInverse** in this case too:

```
Dimensions[xx];
Dimensions[yy];

matmem = Transpose[PseudoInverse[Transpose[{xx}]]].yy
{{2.}, {3.01675}}
```

Recall

Let's check out a few values to see how well the memory matrix can recall a 2 dimensional output, given a one dimensional input.

```
{"recall: ", matmem.Transpose[{xx}][[5]], matmem.Transpose[{xx}][[13]], matmem.Transpose[{"true values: ", yy[[5]], yy[[13]], yy[[22]]}]}//Grid
recall:      {0.458161, 0.691078} {1.69846, 2.56191} {1.13511, 1.71217}
true values: {0.458161, 0.93325} {1.69846, 2.58491} {1.13511, 1.54395}
```

Another function for calculating the least-squares solution

Regression and learning in linear neural networks. Last time we showed 4 different ways to find the generating parameters {2,3} for the following data:

```
rsurface[a_, b_] := N[Table[{x1 = 1 RandomReal[],
  x2 = 1 RandomReal[], a x1 + b x2 + 0.5 RandomReal[] - 0.25}, {120}], 2];
data = rsurface[2, 3];
yy = data[[All, 3]];
xx = data[[All, 1 ;; 2]];
```

Linear regression is so common that *Mathematica* has added the following function to find the least squares parameters directly:

```
LeastSquares[xx, yy]
{2.00154, 2.97608}
```

Graphical illustration of Gradient

The function $f(x,y)$ is plotted as a contour plot. The demo below calculates and plots the vector $\mathbf{g}\mathbf{v} = \nabla\mathbf{f}$. After normalized to unit length, it is plotted as red arrow. The red arrow should always point down the

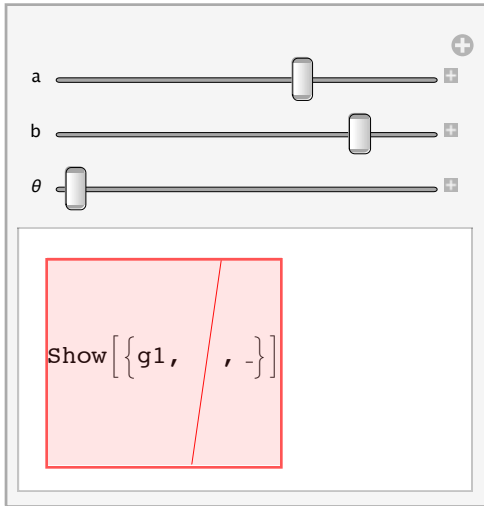
"hill" in the contour plot. You can manipulate the position of the point at which the gradient gets evaluated with the a and b sliders.

The θ slider controls the direction of a second vector **un** (in black) whose length is determined by the projection: $\nabla f \cdot \mathbf{un}$

```

Clear[a, b, gv, un, gm, g2, g3, x, y];
f[x_, y_] := 2 x^2 + 5 y^2 - Sin[x];
g1 =
  ContourPlot[f[x, y], {x, -3, 3}, {y, -3, 3}, PlotPoints -> 50, ImageSize -> Medium];
Manipulate[
  gv = N[D[2 x^2 + 6 y^2 - Sin[x], {{x, y}}] /. {x -> a, y -> b}];
  gn = -Normalize[gv];
  un = {Cos[θ], Sin[θ]};
  gm = (un.gn) un;
  g2 = Graphics[{Red, Arrow[{{a, b}, gn + {a, b}}]}];
  g3 = Graphics[Arrow[{{a, b}, gm + {a, b}}]];
  Show[{g1, g2, g3}], {{a, 1}, -3, 3}, {{b, 2}, -3, 3}, {θ, 0, 2 * Pi}]

```



Show::gcomb : Could not combine the graphics objects in Show[{g1, , -}]. >

Show::gcomb : Could not combine the graphics objects in Show[{g1, , -}]. >

Show::gcomb : Could not combine the graphics objects in Show[{g1, , -}]. >

Show::gcomb : Could not combine the graphics objects in Show[{g1, , -}]. >

Show::gcomb : Could not combine the graphics objects in Show[{g1, , -}]. >

Logistic function - a smooth, differentiable non-linearity

As we saw above, it can be useful to have a non-linearity which is smooth enough to be differentiable. **D[]** returns an expression for the derivative, so to define a function that is the derivative of another we write:

```
In[176]:= Df1[x_] := D[f[t],t] /.t->x
Df1[x_] := Evaluate[D[f[t],t]]
```

The more direct way is to use "operators" or functionals that take functions as inputs and return functions as outputs. **Derivative[]**, or **f'[x]** return functions:

```
In[178]:= (*Df[x_] := Derivative[f[x],x]*)
Df[x_] := f'[x]
```

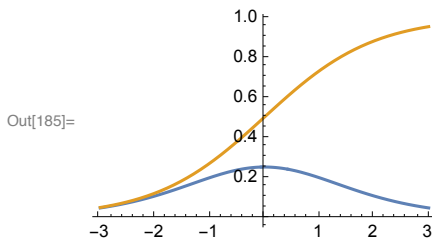
Note that the derivative has a particularly simple expression in terms of $f[x]$, which you can verify by comparing `Simplify[Df[x]]` with `Simplify[h-h[x]]`, or `Simplify[Df[x]-hh[x]]`:

```
In[183]:= f[x_] := 1 / (1 + Exp[-x]);
```

```
In[184]:= hh[x_] := f[x] (1-f[x]);
```

Here is a plot of the sigmoid and its derivative:

```
In[185]:= Plot[{Df[x], f[x]}, {x, -3, 3}]
```



Exercise: calculate the gradient for the last layer in a back-prop network

Let's see how to calculate the gradient for the last layer. Then we will generalize the result to an expression that tells us how to adjust the ij th weight for the connection in the λ th layer.

```
In[163]:= Remove["Global`*"]
```

We start with the error term for the j th training pair. Note that y_j depends on the inputs x , but these are fixed. At this point, we care about the variable weights.

```
In[164]:= (1/2) * (t_i - y_i[w_ij]) ^ 2
```

```
Out[164]= 1/2 (t_i - y_i[w_ij]) ^ 2
```

Take the derivative with respect to w_{ij}

```
In[165]:= eq1 = D[(1/2) * (t_i - y_i[w_ij]) ^ 2, w_ij]
```

```
Out[165]= -(t_i - y_i[w_ij]) y_i'[w_ij]
```

y depends on u which in turn depends on the w 's, so using the chain rule in calculus, we can write

y_i' (w_{ij}) as:

```
In[166]:= eq2 = D[y_i[u_i[w_ij]], w_ij]
```

```
Out[166]:= u_i'[w_ij] y_i'[u_i[w_ij]]
```

But $y_i'(u_i)$ is the derivative of the squashing function f : $y_i'(u_i)=f'(u_i)$ So we can substitute f' for y in eq2

```
In[167]:= eq3 = eq2 /. y_i'[u_i[w_ij]] -> f'[u_i]
```

```
Out[167]:= f'[u_i] u_i'[w_ij]
```

Now u is given by:

```
In[168]:= u_i = Sum[x_k w_i,k, {k, 1, N}]
```

```
Out[168]:= Sum[x_k w_i,k, {k, 1, N}]
```

Let's pick an arbitrary input input, say $i=5$. Suppose we take the derivative with respect to $w_{j,i}$, with $i = 5$? Then with respect to $w_{j,5}$, we have:

```
In[169]:= D[Sum[x_k w_i,k, {k, 1, 10}], w_i,5]
```

```
Out[169]:= x_5
```

So from this we can guess that in general (i.e. any i) we should have:

$$\Delta w_{ij} = (t_i - y_i(x; w_{ij})) f'(u_i) x_j$$

where x_j is the j th "input" to the last layer. This is the *delta-rule*. But this delta rule only works for the output layer. We need to know how to update all the weights.

We don't have direct access to the hidden units, and the problem is how to find the error signals for the hidden layers. It turns out that the delta error terms can be propagated back in terms of a weighted sum of the delta terms at the level above. We will see how this update rule gets applied to the weights starting from the top--i.e. those closest to the outputs, and then applied recursively down towards the inputs. Hence "back-propagation". We will relabel x_k to be y_k^λ , where λ indicates the layer.

So y_k^0 corresponds to the set of inputs to the network (bottom of figure).

To derive the recursive rules to relate deltas at an earlier layer to the next higher layer takes some work and careful book keeping with indices. We won't go through the details here. Let's look at a summary of the algorithm that results.

References

- Bishop, C. M. (1995). Neural Networks for Pattern Recognition. Oxford: Oxford Univeristy Press.
- Duda, R. O., & Hart, P. E. (1973). Pattern classification and scene analysis. New York.: John Wiley & Sons.
- Duda, R. O., Hart, P. E., & Stork, D. G. (2001). Pattern classification (2nd ed.). New York: Wiley.

William T. Freeman, Egon C. Pasztor and Owen T. Carmichael (2000). Learning Low-Level Vision. *International Journal of Computer Vision* 40, 25-47

Kersten, D., O'Toole, A. J., Sereno, M. E., Knill, D. C., & Anderson, J. A. (1987). Associative learning of scene parameters from images. *Appl. Opt.*, 26, 4999-5006.

Knill, D. C., & Kersten, D. (1990). Learning a near-optimal estimator for surface shape from shading. *CVGIP*, 50(1), 75-100.

Nefs, H.T., Koenderink, J.J. & Kappers, A.M.L. (2006). Shape-from-shading for matte and glossy objects. *Acta Psychologica*, 121 (3), 297-316

Poggio, T., Torre, V., & Koch, C. (1985). Computational vision and regularization theory. *Nature*, 317, 314-319.

Geman, S., Bienenstock, E., & Doursat, R. (1992). Neural networks and the bias/variance dilemma. *Neural Computation*, 4(1), 1-58.

© 1998, 2001, 2003, 2005, 2007, 2009, 2011, 2012, 2014 Daniel Kersten, Computational Vision Lab, Department of Psychology, University of Minnesota.