

# Introduction to Neural Networks

## Daniel Kersten

### Belief Propagation

#### Initialize

```
Off[General::spell1];  
Needs["ErrorBarPlots`"]
```

---

## Last time

Formulating an energy or cost function, and then use gradient descent to derive update rules.

Application to interpolation in preception

Many low-level vision problems can be cast as minimizing a cost that depends on the data and a smoothness constraint on the quantities to be estimated.

In Bayesian terms, this translations to maximizing a posterior that depends on a likelihood and prior.

Bias variance and over-fitting

---

## Bayesian learning of univariate Gaussian mean: MAP

We'll spend more time on unsupervised learning later, but it is useful to see how to apply what we've learned about to one of the simplest applications of Bayes to unsupervised learning.

From a statistical point of view, one form of unsupervised learning is "density estimation" which can be done from histogram measurements. In high dimensions this is hard, but is easier if we have a low-dimensional parametric model for the density--i.e. the density is modeled in terms of a few parameters. So for example, the 1D Gaussian could be approximated by a huge list of numbers ("statistics")--one for each bin, each number is an estimate of the probability of the value of the random variable falling in that bin. But because it is Gaussian, we can be more efficient by representing the density in terms of just two numbers (also "statistics", but just the mean and variance), and a formula. In this context, learning becomes *parameter estimation*.

**A Bayesian learning example: Suppose we believe the data comes from a Gaussian generative process, but we don't know the mean.**

For simplicity, suppose we have a set of samples that come from a Gaussian distribution with known variance  $\sigma^2$ , but unknown mean  $\mu$ .

$x_i = \text{noise}$ , where  $\text{noise} \sim N[\mu, \sigma]$ , or equivalently  
 $x_i = \mu + \text{noise}$ , where  $\text{noise} \sim N[0, \sigma]$

In[1]:= **ndist0 = NormalDistribution**[ $\mu$ ,  $\sigma$ ];

This is the forward model. We now want to do “inverse probability”, i.e. estimate the mean and its distribution based on prior assumptions and data.

So lets treat the unknown mean,  $\mu$  as a random variable, and assume a Gaussian prior on it:

$p(\mu)$ , with  $\mu \sim N[\mu_0, \sigma_0]$

In[4]:= **ndist $\mu$  = NormalDistribution**[ $\mu_0$ ,  $\sigma_0$ ];  
**PDF**[**ndist $\mu$** ,  $\mu$ ]

Out[5]= 
$$\frac{e^{-\frac{(\mu-\mu_0)^2}{2\sigma_0^2}}}{\sqrt{2\pi}\sigma_0}$$

In[6]:= 
$$\frac{e^{-\frac{(\mu-\mu_0)^2}{2\sigma_0^2}}}{\sqrt{2\pi}\sigma_0}$$

Out[6]= 
$$\frac{e^{-\frac{(\mu-\mu_0)^2}{2\sigma_0^2}}}{\sqrt{2\pi}\sigma_0}$$

We can think of  $\mu_0$  as an initial guess of the mean's mean, with standard deviation ( $\sigma_0$ ). But we are willing to change our estimate of the mean given new data. If we are really uncertain at the beginning,, we can start of assuming a large standard deviation  $\sigma_0$ , and as we gather data, the uncertainty about the value of the mean will decrease.

Suppose the generative model  $N[\mu, \sigma]$  produces three i.i.d. (independent, identically distributed) samples  $x_1, x_2, x_3$ . What is the MAP estimate of  $\mu$ ? Which value of  $\mu$  makes the posterior biggest? We use Bayes rule:

$$p(\mu | x_1, x_2, x_3) = \frac{p(x_1, x_2, x_3 | \mu) p(\mu)}{p(x_1, x_2, x_3)}$$

$p(x_1 | \mu)$  is given by :

In[7]:= **PDF**[**ndist0**,  $x_1$ ]

Out[7]= 
$$\frac{e^{-\frac{(-\mu+x_1)^2}{2\sigma^2}}}{\sqrt{2\pi}\sigma}$$

Because the samples are drawn independently, the  $p(x_1, x_2, x_3 | \mu)$  is the product of three terms, so the numerator is  $p(x_1 | \mu) p(x_2 | \mu) p(x_3 | \mu)$  times the prior  $p(\mu)$ :

In[9]:= **PDF**[**ndist0**,  $x_1$ ] \* **PDF**[**ndist0**,  $x_2$ ] \* **PDF**[**ndist0**,  $x_3$ ] \* **PDF**[**ndist $\mu$** ,  $\mu$ ]

Out[9]= 
$$\frac{e^{-\frac{(-\mu+x_1)^2}{2\sigma^2} - \frac{(-\mu+x_2)^2}{2\sigma^2} - \frac{(-\mu+x_3)^2}{2\sigma^2} - \frac{(\mu-\mu_0)^2}{2\sigma_0^2}}}{4\pi^2\sigma^3\sigma_0}$$

## Calculating the MAP estimate of mean

To find the value of the mean that is largest given our three samples, and our prior assumption, we need to find  $\mu$  where  $p(x_1, x_2, x_3 | \mu) p(\mu)$  is biggest. This is the same as finding the value where the log of  $p(x_1, x_2, x_3 | \mu) p(\mu)$  is biggest—a simpler expression. And with calculus, we can do that by finding where the slope of the derivative of the log with respect to  $\mu$  is zero.

```
g = PDF[ndist0, x1] * PDF[ndist0, x2] * PDF[ndist0, x3] * PDF[ndistμ, μ];
t = Log[g];
t1 = PowerExpand[t]
t2 = D[t1, μ]
Solve[-t2 == 0, μ]
```

$$-2 \log[2] - 2 \log[\pi] - 3 \log[\sigma] - \log[\sigma_0] - \frac{(-\mu + x_1)^2}{2\sigma^2} - \frac{(-\mu + x_2)^2}{2\sigma^2} - \frac{(-\mu + x_3)^2}{2\sigma^2} - \frac{(\mu - \mu_0)^2}{2\sigma_0^2}$$

$$\frac{-\mu + x_1}{\sigma^2} + \frac{-\mu + x_2}{\sigma^2} + \frac{-\mu + x_3}{\sigma^2} - \frac{\mu - \mu_0}{\sigma_0^2}$$

$$\left\{ \left\{ \mu \rightarrow \frac{\mu_0 \sigma^2 + x_1 \sigma_0^2 + x_2 \sigma_0^2 + x_3 \sigma_0^2}{\sigma^2 + 3 \sigma_0^2} \right\} \right\}$$

Divide the above expression by  $\sigma_0^2 \sigma^2$ .

$$\left\{ \left\{ \mu \rightarrow \frac{\frac{\mu_0}{\sigma_0^2} + \frac{1}{\sigma^2} (x_1 + x_2 + x_3)}{\frac{3}{\sigma^2} + \frac{1}{\sigma_0^2}} \right\} \right\}$$

We've done the calculation with just three data points, but one can see the pattern. So in general,  $\mu$  can be estimated from  $n$  samples in batch mode by the following rule:

$$\left\{ \left\{ \mu \rightarrow \frac{\frac{\mu_0}{\sigma_0^2} + \frac{1}{\sigma^2} \sum_{i=1}^n x_i}{\frac{n}{\sigma^2} + \frac{1}{\sigma_0^2}} \right\} \right\}$$

With a bit more effort, one can also derive the solution to estimate the standard deviation too, when it is unknown.

So to sum up, we have a method to estimate the probability distribution of the mean given the data. This is because the mean and standard deviation are sufficient statistics for the gaussian distribution. (For the multi-variate case, cf. Duda and Hart.)

What is the influence of the initial estimate of the mean as learning goes on? What is the estimate of the mean as  $n$  gets large?

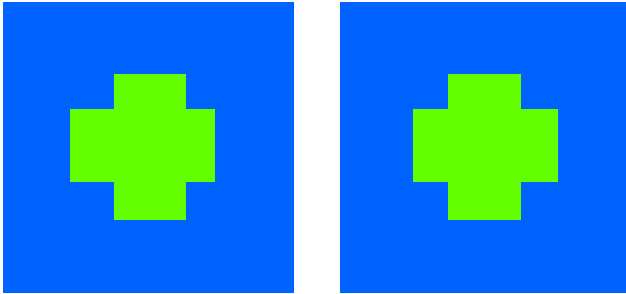
---

Formulate the estimation process as a sequential rule where the posterior gets iteratively updated as each new data point arrives.

---

## Interpolation and perceptual surface completion revisited

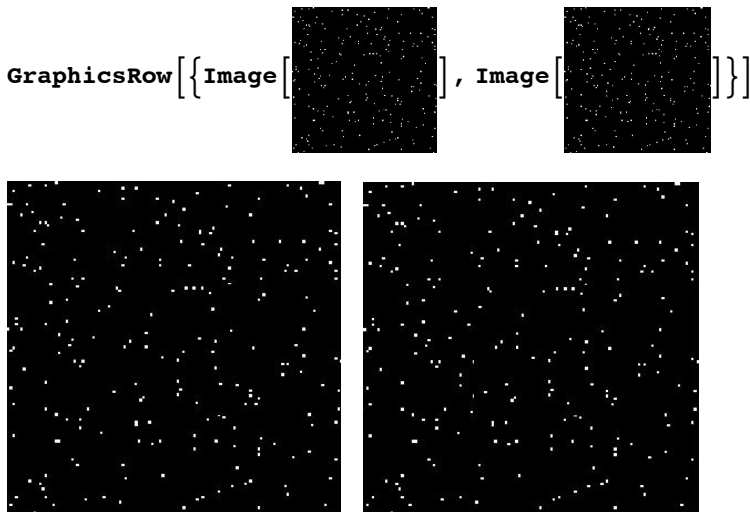
Recall this demonstration. When you view the left and right images below in stereo, you may see a horizontal rectangle floating out in front of the vertical rectangle in the back:



The figure below gives some idea of what it looks like, in case you can't cross-fuse the above (I've changed the color of the horizontal bar for the purposes of illustration).



This is an example of your visual system interpolating a smooth surface from sparse data (cf. Nakayama and Shimojo, 1992). The data is sparse because the information about depth comes from the disparities in the left and right vertical edges of the horizontal rectangle. The random dot stereogram that we saw in an earlier lecture was “dense” rather than sparse. The data is dense because there were lots of potential features to match throughout the background and the square that floated out in depth. Here's an example of a sparse random dot stereogram:

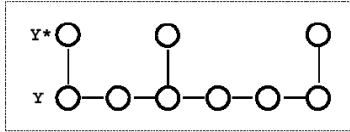


One sees the white points on dark, central surface floating in front of a background surface, also with a white point texture. How can one model processes of interpolation? Earlier we approached the problem by constructing a cost function to be minimized.

Here we re-formulate the problem from a Bayesian perspective in order to introduce Belief Propagation.

## Review of interpolation using smoothness: Gradient descent

For simplicity, we'll assume 1-D as in the lecture on sculpting the energy function. In anticipation of formulating the problem in terms of a graph that represents conditional probability dependence, we represent *observable* depth cues by  $y^*$ , and the true ("hidden") depth estimates by  $y$ .



(Figure from Weiss (1999).)

## First-order smoothness

Earlier we saw that under specific assumptions, biologically plausible neural updating can be seen to decrease the value of an energy or cost function.

One can also start off with an assumed cost function, determined by a set of constraints, and use gradient descent to derive an update rule that minimizes the cost. (See supplementary material). However, such an update rule will not necessarily resemble a biologically plausible neural mechanism.

For an interpolation problem, we can write the energy or cost function by:

$$J(Y) = \sum_k w_k (y_k - y_k^*)^2 + \lambda \sum_i (y_i - y_{i+1})^2$$

where  $w_k (= xs[[k]]$  below) is an "indicator function", and  $y_k^* = d$ , are the data values. The indicator function is 1 if there is data available, and zero otherwise. The second sum represents the sum of the squared differences between neighboring  $y$ -values. Minimizing the first sum encourages the estimates of  $y$  to be close to the measured values. Minimizing the second sum encourages nearby  $y$ -values to be the same. Thus minimizing  $J(Y)$  encourages fidelity to the data where present, and similarity to nearby values where there is no data.

Gradient descent gives the following local update rule:

$$y_k \leftarrow y_k + \eta_k \left( \lambda \left( \frac{y_{k-1} + y_{k+1}}{2} - y_k \right) + w_k (y_k^* - y_k) \right)$$

$\lambda$  is a free parameter that controls the degree of smoothness, i.e. smoothness at the expense of fidelity to the data.

At each step, the rule encourages the estimates to get closer to the data (where it exists), and also closer to the values of neighbors.

There are various choices for how to change the smoothness as a function of iterations.

E.g. standard methods in numerical analysis include:

Gauss-Seidel:  $\eta[k] := 1 / (\lambda + xs[[k]])$ . And successive over-relaxation (SOR):

```
 $\eta_2[k\_]:=1.9/(\lambda+xs[[k]]);$ 
```

## A simulation: Straight line with random missing data points

### Make the data

Consider the problem of interpolating a set of points with missing data, marked by an indicator function with the following notation:

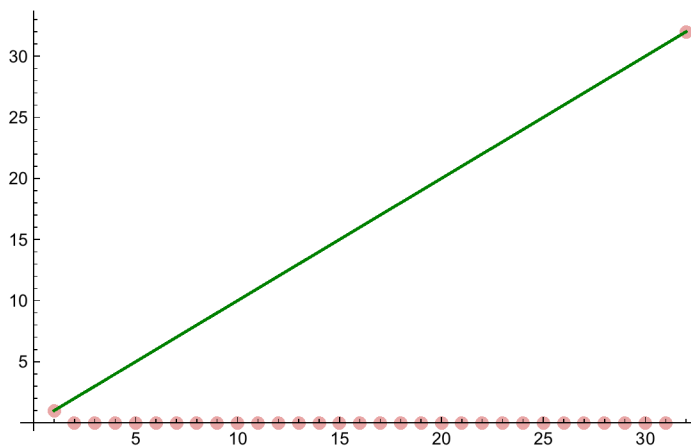
$w_k = xs[[k]]$ ,  $y^* = \text{data}$ ,  $y=f$ .

We'll assume the true model is that  $f = y = j$ , where  $j=1$  to  $size$ . `data` is a function of the sampling process on  $f = j$

```
size = 32; xs = Table[0, {i, 1, size}]; xs[[1]] = 1;
xs[[size]] = 1; data = Table[N[j] xs[[j]], {j, 1, size}];
g3 = ListPlot[Table[N[j], {j, 1, size}], Joined -> True,
  PlotStyle -> {RGBColor[0, 0.5, 0]}; g2 = ListPlot[data, Joined -> False,
  PlotStyle -> {Opacity[0.35], RGBColor[0.75, 0., 0], PointSize[Large]}];
```

The green line shows the a straight line connecting the data points. The red dots on the abscissa mark the points where data are missing.

```
Show[{g2, g3}, ImageSize -> Medium]
```



The update rule is linear, so we can represent in terms of two matrices, **Tm** and **Sm** such that the gradient of the energy is equal to:

$$\mathbf{Tm} \cdot \mathbf{f} - \mathbf{Sm} \cdot \mathbf{f}.$$

**Sm** will be our filter to exclude non-data points specified by the indicator function  $w_k = xs[[k]]$ .

**Tm** will express the "smoothness" constraint.

```
Sm = DiagonalMatrix[xs];
Tm = Table[0, {i, 1, size}, {j, 1, size}];
For[i=1, i<=size, i++, Tm[[i, i]] = 2];
Tm[[1, 1]] = 1; Tm[[size, size]] = 1; (*Adjust for the boundaries*)
For[i=1, i<size, i++, Tm[[i+1, i]] = -1];
For[i=1, i<size, i++, Tm[[i, i+1]] = -1];
```

## Run gradient descent

```
Clear[f, d, λ]
(λ * Tm.Array[f, size] - Sm.((Array[d, size]) - Array[f, size])) // MatrixForm ;

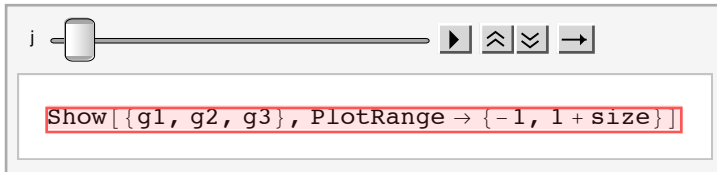
Clear[Tf, f1, j];
dt = 1; λ=2;
Tf[f1_] := f1 - dt*(1/(λ+xs))*(Tm.f1 - λ*Sm.(data-f1));
```

We will initialize the state vector to zero, and then run the network for **iter** iterations:

```
f0 = Table[0, {i, 1, size}];
f0 = Table[RandomReal[{0, 30}], {i, 1, size}];
result=f0;
f=f0;
iter=25;
```

Now plot the interpolated function.

```
Animate[result = Nest[Tf, f, iter];
  f = result;
  g1 = ListPlot[result, Joined → False, AspectRatio → Automatic,
    PlotRange → {{0, size}, {-1, size + 1}}, ImageSize → Medium];
  Show[{g1, g2, g3}, PlotRange → {-1, size + 1}], {j, 1, 10, 1},
  AnimationRunning → False]
```



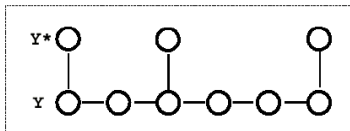
Try starting with  $f$  = random values. Try various numbers of iterations.

Try different sampling functions  $xs[[i]]$ .

## Interpolation using Belief Propagation

Same interpolation problem, but now using belief propagation

Example is taken from Yair Weiss.(Weiss, 1999)



Probabilistic generative model

$$\text{data}[[i]] = y^*[i] = xs[[i]] y[[i]] + \text{dnoise}, \text{dnoise} \sim N[0, \sigma_D]$$

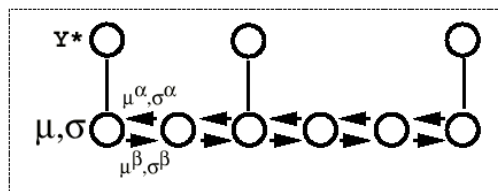
$$y[[i+1]] = y[[i]] + \text{znoise}, \quad \text{znoise} \sim N[0, \sigma_R]$$

The first term is the "data formation" model, i.e. how the data is directly influenced by the interaction of the underlying influences or causes:

$y^*$  is determined by an underlying hidden "y" which can't be directly measured. But we assume we can measure  $y^*$ , which is determined by sampling some values of  $y$  and adding noise.

The second term reflects our prior assumptions about the smoothness of  $y$ , i.e. nearby  $y$ 's are correlated, and in fact assumed identical except for some added noise. So with no noise the prior reflects the assumption that lines are horizontal--all  $y$ 's are the same. This is sometimes called a "soft constraint", because it is a tendency--we don't insist that it be satisfied exactly.

## Summary of the message passing rules



We'd like to know the distribution of the random variables at each node  $i$ , conditioned on all the data: i.e. we want the posterior

$$p(Y_i = u \mid \text{all the data})$$

If we could find this, we'd be able to: 1) say what the most probable value of  $y$  is, and 2) give a measure of confidence.

The appendix derives a solution, based on belief propagation, which sequentially updates local conditional distributions for the estimates of the mean of  $Y$  at each point, as well as estimating its standard deviation.

Let's go over the rules of message passing. We follow Weiss, and make a (hopefully not too confusing) notation change to avoid the square superscripts, using the notation substitution:  $\sigma_D^2 \rightarrow \sigma_D$ ,  $\sigma_R^2 \rightarrow \sigma_R$ .

So  $\sigma_D$  represents the variability or uncertainty in the data. If  $\sigma_D$  is small, we trust the data  $Y^*$ .  $\sigma_R$  represents our prior belief

in the "roughness" of the curve--the  $Y$ 's--to be estimated. If  $\sigma_R$  is big, then we tolerate more deviations from straightness, while small values would bias the estimates towards smooth, straight lines.

Let's look at the messages passed for our example:



$$\mu_i \leftarrow \frac{\frac{w_i}{\sigma_D} Y_i^* + \frac{1}{\sigma_i^\alpha} \mu_i^\alpha + \frac{1}{\sigma_i^\beta} \mu_i^\beta}{\frac{w_i}{\sigma_D} + \frac{1}{\sigma_i^\alpha} + \frac{1}{\sigma_i^\beta}}$$

$$\sigma_i \leftarrow \frac{1}{\frac{w_i}{\sigma_D} + \frac{1}{\sigma_i^\alpha} + \frac{1}{\sigma_i^\beta}}$$

$$\mu_i^\alpha \leftarrow \frac{\frac{1}{\sigma_{i-1}^\alpha} \mu_{i-1}^\alpha + \frac{w_{i-1}}{\sigma_D} Y_{i-1}^*}{\frac{1}{\sigma_{i-1}^\alpha} + \frac{w_{i-1}}{\sigma_D}}$$

$$\sigma_i^\alpha \leftarrow \sigma_R + \left( \frac{1}{\sigma_{i-1}^\alpha} + \frac{w_{i-1}}{\sigma_D} \right)^{-1}$$

Side note: Notice the underlying operation in which estimates get combined weighted by uncertainty, similar to how we updated the mean in the Bayesian learning example. This is a recurring theme. In your next assignment, you show that the maximum a posteriori estimate of for cue combination is:

$$\mu = \frac{1/\sigma_1^2}{1/\sigma_1^2 + 1/\sigma_2^2} \mu_1 + \frac{1/\sigma_2^2}{1/\sigma_1^2 + 1/\sigma_2^2} \mu_2 = \frac{\mu_1/\sigma_1^2 + \mu_2/\sigma_2^2}{1/\sigma_1^2 + 1/\sigma_2^2}$$

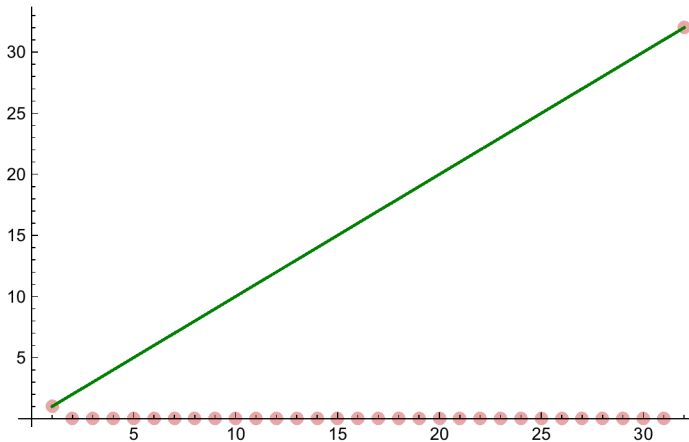
where  $r_i = 1/\sigma_i^2$ , and  $\mu_1$  and  $\mu_2$  are the estimates of the means obtained separately. How this could be done neurally has inspired recent work in neural population codes, mentioned below. Combining information weighted by reliability underlies a number of models of visual and multimodal cue integration in psychophysics over the past decade.

## A simulation: Belief propagation for interpolation with missing data

```
size = 32; xs = Table[0, {i, 1, size}]; xs[[1]] = 1;
xs[[size]] = 1; data = Table[N[j] xs[[j]], {j, 1, size}];
g3bp = ListPlot[Table[N[j], {j, 1, size}], Joined -> True,
  PlotStyle -> {RGBColor[0, 0.5, 0]}]; g2bp = ListPlot[data, Joined -> False,
  PlotStyle -> {Opacity[0.35], RGBColor[0.75, 0., 0], PointSize[Large]}];
```

The green line shows the a straight line connecting the data points. The red dots on the abscissa mark the points where data are missing.

```
Show[{g2bp, g3bp}, ImageSize -> Medium]
```



## Initialization

```
size = 32;
μ0 = 1;
μα = 1; σα = 100 000; (*large uncertainty *)
μβ = 1; σβ = 100 000; (*large*)
σR = 4.0; σD = 1.0;
μ = Table[μ0, {i, 1, size}];
σ = Table[σα, {i, 1, size}];
μα = Table[μ0, {i, 1, size}];
σα = Table[σα, {i, 1, size}];
μβ = Table[μ0, {i, 1, size}];
σβ = Table[σβ, {i, 1, size}];
iter = 0;
i = 1;
j = size;
```

The code below implements the above iterative equations, taking care near the boundaries. The plot shows the estimates of  $y_i = \mu$ , and the error bars show  $\pm \sigma_i$ .

## Belief Propagation Routine

```
yfit = Table[{0, 0}, {i1, 1, size}];
g1b = ErrorListPlot[{yfit}];
Dynamic[Show[
  {g1b, g2bp, g3bp, Graphics[{Text["Iteration=" <> ToString[iter], {
     $\frac{\text{size}}{2}$ , size}]}]}],
  PlotRange -> {-50, 50}, Axes -> {False, True}]]
```

Execute the next cell to run 31 iterations. The display is slowed down so that you can see the progression of the updates in the above graph.

```

Do[
  Pause[.5];
  
$$\mu[[i]] = \frac{\frac{xs[[i]] data[[i]]}{\sigma D} + \frac{\mu\alpha[[i]]}{\sigma\alpha[[i]]} + \frac{1. \mu\beta[[i]]}{\sigma\beta[[i]]}}{\frac{xs[[i]]}{\sigma D} + \frac{1}{\sigma\alpha[[i]]} + \frac{1}{\sigma\beta[[i]]}};$$

  
$$\sigma[[i]] = \frac{1.}{\frac{xs[[i]]}{\sigma D} + \frac{1}{\sigma\alpha[[i]]} + \frac{1}{\sigma\beta[[i]]}};$$

  
$$\mu[[j]] = \frac{\frac{xs[[j]] data[[j]]}{\sigma D} + \frac{\mu\alpha[[j]]}{\sigma\alpha[[j]]} + \frac{1. \mu\beta[[j]]}{\sigma\beta[[j]]}}{\frac{xs[[j]]}{\sigma D} + \frac{1}{\sigma\alpha[[j]]} + \frac{1}{\sigma\beta[[j]]}};$$

  
$$\sigma[[j]] = \frac{1.}{\frac{xs[[j]]}{\sigma D} + \frac{1}{\sigma\alpha[[j]]} + \frac{1}{\sigma\beta[[j]]}};$$

  nextj = j - 1;
  
$$\mu\alpha[[nextj]] = \frac{\frac{xs[[j]] data[[j]]}{\sigma D} + \frac{1. \mu\alpha[[j]]}{\sigma\alpha[[j]]}}{\frac{xs[[j]]}{\sigma D} + \frac{1}{\sigma\alpha[[j]]}};$$

  
$$\sigma\alpha[[nextj]] = \sigma R + \frac{1.}{\frac{xs[[j]]}{\sigma D} + \frac{1}{\sigma\alpha[[j]]}};$$

  nexti = i + 1;
  
$$\mu\beta[[nexti]] = \frac{\frac{xs[[i]] data[[i]]}{\sigma D} + \frac{1. \mu\beta[[i]]}{\sigma\beta[[i]]}}{\frac{xs[[i]]}{\sigma D} + \frac{1}{\sigma\beta[[i]]}};$$

  
$$\sigma\beta[[nexti]] = \sigma R + \frac{1.}{\frac{xs[[i]]}{\sigma D} + \frac{1}{\sigma\beta[[i]]}};$$

  j--;
  i++;
  iter++;
  yfit = Table[{μ[[i1]], σ[[i1]]}, {i1, 1, size}];
  glb = ErrorListPlot[{yfit};
    , {size - 1}];

```

---

## Relation to neural networks

We've seen two very different ways of estimating states through iterative updating. It is easy to see how the first, derived from gradient descent, is related to updating in a traditional neural network. The second, belief propagation, gives us a very different view of how information might be represented and updated in a network. For example, in the second case applied to interpolation, the nodes represent probability distributions over depths (summarized by the mean and variance) not just estimates of depth. A critical new feature is the explicit representation of uncertainty, i.e. in the standard deviations of

the node values. Later in the course, we will discuss whether the brain might have explicit representations of such uncertainties, how these might be represented in populations of neurons, and how the brain might do computations on these. Does the brain do belief propagation, and if so how?

For a preview of the general question of whether the brain represents and computes on distributions, see: Ma et al. (2006, 2012), Knill & Pouget (2004), Zemel & Pouget (1998).

---

## Discussion: relation to behavior?

Yair Weiss and colleagues showed how the probabilistic combination of local motion cues and prior motion assumptions could explain a number of visual illusions: Weiss, Y., Simoncelli, E. P., & Adelson, E. H. (2002). Motion illusions as optimal percepts. *Nature Neuroscience*.

But their conclusion (and probably all other similar applications to perceptual behavior) didn't depend on the algorithm.

Can you think of an application and test of belief propagation to perception?

---

## Exercises

Run the gradient descent algorithm using successive over-relaxation (SOR):  $\eta_2[k_] := 1.9 / (\lambda + xs[[k]])$ .

How does convergence compare with Gauss-Seidel?

---

Run Belief Propagation using:  $\sigma_R = 1.0$ ;  $\sigma_D = 4.0$ ; How does fidelity to the data compare with the original case ( $\sigma_R = 4.0$ ;  $\sigma_D = 1.0$ ).

---

BP with missing sine wave data

---

### Generate sine wave with missing data

```
size = 64; xs = Table[RandomInteger[1], {i, 1, size}];
data = Table[N[Sin[ $\frac{2 \pi j}{20}$ ] xs[[j]]], {j, 1, size}];
g3b = ListPlot[Table[N[Sin[ $\frac{2 \pi j}{20}$ ]]], {j, 1, size}],
  Joined → True, PlotStyle → {RGBColor[0, 0.5, 0]};
g2b = ListPlot[data, Joined → False, PlotStyle → {RGBColor[0.75, 0., 0]}];
```

## Initialize

```

μ0 = 1;
μα = 1; σα = 100 000; (*large uncertainty *)
μβ = 1; σβ = 100 000; (*large*)
σR = .5; σD = .1;
μ = Table[μ0, {i, 1, size}];
σ = Table[σα, {i, 1, size}];
μα = Table[μ0, {i, 1, size}];
σα = Table[σα, {i, 1, size}];
μβ = Table[μ0, {i, 1, size}];
σβ = Table[σβ, {i, 1, size}];
iter = 0;
i = 1;
j = size;

yfit = Table[{0, 0}, {i1, 1, size}];
g1bb = ErrorListPlot[{yfit}];
Dynamic[Show[{g1bb, g2b, g3b}, PlotRange → {-2, 2}, Axes → {False, True}]]
Show[{g1bb, g2b, g3b}, PlotRange → {-2, 2}, Axes → {False, True}]

```



## SINE WAVE DEMO: Belief Propagation Routine

```

Do[
  Pause[0.2];

  
$$\mu[i] = \frac{\frac{x_s[i] \text{data}[i]}{\sigma_D} + \frac{\mu\alpha[i]}{\sigma\alpha[i]} + \frac{1 \cdot \mu\beta[i]}{\sigma\beta[i]}}{\frac{x_s[i]}{\sigma_D} + \frac{1}{\sigma\alpha[i]} + \frac{1}{\sigma\beta[i]}};$$


  
$$\sigma[i] = \frac{1.}{\frac{x_s[i]}{\sigma_D} + \frac{1}{\sigma\alpha[i]} + \frac{1}{\sigma\beta[i]}};$$


  
$$\mu[j] = \frac{\frac{x_s[j] \text{data}[j]}{\sigma_D} + \frac{\mu\alpha[j]}{\sigma\alpha[j]} + \frac{1 \cdot \mu\beta[j]}{\sigma\beta[j]}}{\frac{x_s[j]}{\sigma_D} + \frac{1}{\sigma\alpha[j]} + \frac{1}{\sigma\beta[j]}};$$


  
$$\sigma[j] = \frac{1.}{\frac{x_s[j]}{\sigma_D} + \frac{1}{\sigma\alpha[j]} + \frac{1}{\sigma\beta[j]}};$$


  nextj = j - 1;

  
$$\mu\alpha[\text{nextj}] = \frac{\frac{x_s[j] \text{data}[j]}{\sigma_D} + \frac{1 \cdot \mu\alpha[j]}{\sigma\alpha[j]}}{\frac{x_s[j]}{\sigma_D} + \frac{1}{\sigma\alpha[j]}};$$


  
$$\sigma\alpha[\text{nextj}] = \sigma_R + \frac{1.}{\frac{x_s[j]}{\sigma_D} + \frac{1}{\sigma\alpha[j]}};$$


  nexti = i + 1;

  
$$\mu\beta[\text{nexti}] = \frac{\frac{x_s[i] \text{data}[i]}{\sigma_D} + \frac{1 \cdot \mu\beta[i]}{\sigma\beta[i]}}{\frac{x_s[i]}{\sigma_D} + \frac{1}{\sigma\beta[i]}};$$


  
$$\sigma\beta[\text{nexti}] = \sigma_R + \frac{1.}{\frac{x_s[i]}{\sigma_D} + \frac{1}{\sigma\beta[i]}};$$


  j--;
  i++;
  iter++;
  yfit = Table[{μ[[i1]], σ[[i1]]}, {i1, 1, size}];
  glbb = ErrorListPlot[{yfit}];
  , {size - 1}]

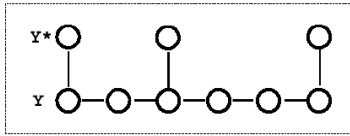
```

---

## Appendix

Derivation of the update rules, called “message passing”

## Updating the mean and variance given the data at point i, and current beliefs about mean and variance before and after i



Let  $p(Y_i=u \mid \text{all the data})$  be normally distributed:  $\text{NormalDistribution}[\mu_i, \sigma_i]$ .

Consider the  $i$ th unit. The posterior  $p(Y_i=u \mid \text{all the data}) =$

$$p(Y_i=u \mid \text{all the data}) \propto p(Y_i=u \mid \text{data before } i) p(\text{data at } i \mid Y_i=u) p(Y_i=u \mid \text{data after } i)$$

(“before” and “after” means to the left and right of  $i$ , respectively.)

Suppose that  $p(Y_i=u \mid \text{data before } i)$  is also gaussian:

$$p(Y_i=u \mid \text{data before } i) = \alpha_i[u] \sim \text{NormalDistribution}[\mu\alpha, \sigma\alpha]$$

and so is probability conditioned on the data after  $i$ :

$$p(Y_i=u \mid \text{data after } i) = \beta_i[u] \sim \text{NormalDistribution}[\mu\beta, \sigma\beta]$$

And the noise model for the data:

$$p(\text{data at } i \mid Y_i=u) = L_i[u] \sim$$

$\text{NormalDistribution}[y_p, \sigma_D]$

$$y_p = \text{data}[[i]]$$

So in terms of these functions, the posterior probability of the  $i$ th unit taking on the value  $u$  can be expressed as proportional to a product of the three factors:

$$p(Y_i=u \mid \text{all the data}) \propto \alpha_i[u] * L_i[u] * \beta_i[u]$$

```
ClearAll[\mu\alpha, \sigma\alpha, \mu\beta, \sigma\beta, y_p, \sigma_D];
```

```
\alphaudist = NormalDistribution[\mu\alpha, \sigma\alpha];
```

```
\alpha[u] = PDF[\alphaudist, u];
```

```
\Ddist = NormalDistribution[y_p, \sigma_D];
```

```
L[u] = PDF[Ddist, u];
```

```
\budist = NormalDistribution[\mu\beta, \sigma\beta];
```

```
\beta[u] = PDF[budist, u];
```

```
\alpha[u] * L[u] * \beta[u]
```

$$\frac{e^{-\frac{(u-y_p)^2}{2\sigma_D^2} - \frac{(u-\mu\alpha)^2}{2\sigma\alpha^2} - \frac{(u-\mu\beta)^2}{2\sigma\beta^2}}}{2\sqrt{2}\pi^{3/2}\sigma_D\sigma\alpha\sigma\beta}$$

This just another gaussian distribution on  $Y_i=u$ . What is its mean and variance? Finding the root enables us to complete the square to see what the numerator looks like. In particular, what the mode (=mean for gaussian) is.

$$\text{Solve}\left[-\mathbf{D}\left[-\frac{(\mathbf{u}-\mu\alpha)^2}{2\sigma\alpha^2}-\frac{(\mathbf{u}-\mu\beta)^2}{2\sigma\beta^2}-\frac{(\mathbf{u}-\mathbf{y}_p)^2}{2\sigma_p^2}, \mathbf{u}\right] == \mathbf{0}, \mathbf{u}\right]$$

$$\left\{\left\{\mathbf{u} \rightarrow \frac{\frac{\mu\alpha}{\sigma\alpha^2} + \frac{\mu\beta}{\sigma\beta^2} + \frac{\mathbf{y}_p}{\sigma_p^2}}{\frac{1}{\sigma\alpha^2} + \frac{1}{\sigma\beta^2} + \frac{1}{\sigma_p^2}}\right\}\right\}$$

This suggests that if we had estimates of  $\mu\alpha, \mu\beta, \sigma\alpha^2, \sigma\beta^2$  and the data, we could update the mean of node  $i$  using:

$$\mathbf{u} \leftarrow \frac{\frac{\mu\alpha}{\sigma\alpha^2} + \frac{\mu\beta}{\sigma\beta^2} + \frac{\mathbf{y}_p}{\sigma_p^2}}{\frac{1}{\sigma\alpha^2} + \frac{1}{\sigma\beta^2} + \frac{1}{\sigma_p^2}}$$

Similarly, the update rule for the variance is:

$$\sigma^2 \leftarrow \frac{1}{\sigma\alpha^2} + \frac{1}{\sigma\beta^2} + \frac{1}{\sigma_p^2}$$

### How do we get $\mu\alpha, \mu\beta, \sigma\alpha, \sigma\beta$ ?

We express the probability of the  $i$ th unit taking on the value  $\mathbf{u}$  in terms of the values of the neighbor before, conditioning on what is known (the observed measurements), and marginalizing over what isn't (the previous "hidden" node value,  $\mathbf{v}$ , at the  $i-1$ th location).

We have three terms to worry about that depend on nodes in the neighborhood preceding  $i$ :

$$\alpha[\mathbf{u}] = \int_{-\infty}^{\infty} \alpha_p[\mathbf{v}] * \mathbf{S}[\mathbf{u}] * \mathbf{L}[\mathbf{v}] \, d\mathbf{v} \propto \int_{-\infty}^{\infty} e^{-\frac{(\mathbf{v}-\mathbf{y}_p)^2}{2\sigma_p^2} - \frac{(\mathbf{u}-\mathbf{v})^2}{2\sigma_R^2} - \frac{(\mathbf{v}-\mu\alpha_p)^2}{2\sigma\alpha_p^2}} \, d\mathbf{v}$$

$\alpha_p = \alpha_{i-1} \cdot \mathbf{S}[\mathbf{u}]$  is our smoothing term, or transition probability:  $\mathbf{S}[\mathbf{u}] = p(\mathbf{u} | \mathbf{v})$ .

$\mathbf{L}[\mathbf{v}]$  is the likelihood of  $\mathbf{v}$  given the data previous at the previous node.

**Rdist = NormalDistribution[v,  $\sigma_R$ ];**

**S[u] = PDF[Rdist, u];**

**avdist = NormalDistribution[ $\mu\alpha_p, \sigma\alpha_p$ ];**

**$\alpha_p[\mathbf{v}] = \text{PDF}[\text{avdist}, \mathbf{v}]$ ;**

**$\mathbf{L}_p[\mathbf{v}] = \text{PDF}[\text{Ddist}, \mathbf{v}]$ ;**

**Integrate[ $\alpha_p[\mathbf{v}] * \mathbf{S}[\mathbf{u}] * \mathbf{L}_p[\mathbf{v}]$ , {v, -Infinity, Infinity}]**

$$\text{ConditionalExpression}\left[\frac{e^{-\frac{u^2 \sigma_D^2 + y_p^2 \sigma_R^2 + \mu\alpha_p^2 (\sigma_D^2 + \sigma_R^2) - 2 \mu\alpha_p (u \sigma_D^2 + y_p \sigma_R^2) + u^2 \sigma\alpha_p^2 - 2 u y_p \sigma\alpha_p^2 + y_p^2 \sigma\alpha_p^2}{2 (\sigma_D^2 \sigma\alpha_p^2 + \sigma_R^2 (\sigma_D^2 + \sigma\alpha_p^2))}}}{2 \pi \sigma_D \sigma_R \sqrt{\frac{1}{\sigma_D^2} + \frac{1}{\sigma_R^2} + \frac{1}{\sigma\alpha_p^2}}} \sigma\alpha_p, \text{Re}\left[\frac{1}{\sigma_D^2} + \frac{1}{\sigma_R^2} + \frac{1}{\sigma\alpha_p^2}\right] \geq 0\right]$$

### Some uninspired Mathematica manipulations

Let's find an expression for the mode of the above calculated expression for  $\alpha[\mathbf{u}]$



$$\mathbf{D} \left[ - \left( (\mathbf{u} - \mu\alpha_p)^2 \sigma_D^2 + \mu\alpha_p^2 \sigma_R^2 + \mathbf{u}^2 \sigma\alpha_p^2 + \mathbf{y}_p^2 (\sigma_R^2 + \sigma\alpha_p^2) - 2 \mathbf{y}_p (\mu\alpha_p \sigma_R^2 + \mathbf{u} \sigma\alpha_p^2) \right) / \right. \\ \left. \left( 2 (\sigma_R^2 \sigma\alpha_p^2 + \sigma_D^2 (\sigma_R^2 + \sigma\alpha_p^2)) \right), \mathbf{u} \right] \\ = \frac{2 (\mathbf{u} - \mu\alpha_p) \sigma_D^2 + 2 \mathbf{u} \sigma\alpha_p^2 - 2 \mathbf{y}_p \sigma\alpha_p^2}{2 (\sigma_R^2 \sigma\alpha_p^2 + \sigma_D^2 (\sigma_R^2 + \sigma\alpha_p^2))}$$

**Solve**[-% == 0, u]

$$\left\{ \left\{ \mathbf{u} \rightarrow \frac{\frac{\mu\alpha_p \sigma_D^2}{\sigma_R^2 \sigma\alpha_p^2 + \sigma_D^2 (\sigma_R^2 + \sigma\alpha_p^2)} + \frac{\mathbf{y}_p \sigma\alpha_p^2}{\sigma_R^2 \sigma\alpha_p^2 + \sigma_D^2 (\sigma_R^2 + \sigma\alpha_p^2)}}{\frac{\sigma_D^2}{\sigma_R^2 \sigma\alpha_p^2 + \sigma_D^2 (\sigma_R^2 + \sigma\alpha_p^2)} + \frac{\sigma\alpha_p^2}{\sigma_R^2 \sigma\alpha_p^2 + \sigma_D^2 (\sigma_R^2 + \sigma\alpha_p^2)}} \right\} \right\}$$

$$\text{Simplify} \left[ \left( \frac{\mu\alpha_p \sigma_D^2}{\sigma_R^2 \sigma\alpha_p^2 + \sigma_D^2 (\sigma_R^2 + \sigma\alpha_p^2)} + \frac{\mathbf{y}_p \sigma\alpha_p^2}{\sigma_R^2 \sigma\alpha_p^2 + \sigma_D^2 (\sigma_R^2 + \sigma\alpha_p^2)} \right) / (\sigma_D^2 * \sigma\alpha_p^2) \right] \\ = \frac{\mu\alpha_p \sigma_D^2 + \mathbf{y}_p \sigma\alpha_p^2}{\sigma_D^2 \sigma_R^2 \sigma\alpha_p^4 + \sigma_D^4 \sigma\alpha_p^2 (\sigma_R^2 + \sigma\alpha_p^2)}$$

$$\text{Simplify} \left[ \left( \frac{\sigma_D^2}{\sigma_R^2 \sigma\alpha_p^2 + \sigma_D^2 (\sigma_R^2 + \sigma\alpha_p^2)} + \frac{\sigma\alpha_p^2}{\sigma_R^2 \sigma\alpha_p^2 + \sigma_D^2 (\sigma_R^2 + \sigma\alpha_p^2)} \right) / (\sigma_D^2 * \sigma\alpha_p^2) \right] \\ = \frac{\sigma_D^2 + \sigma\alpha_p^2}{\sigma_D^2 \sigma_R^2 \sigma\alpha_p^4 + \sigma_D^4 \sigma\alpha_p^2 (\sigma_R^2 + \sigma\alpha_p^2)}$$

$$\left( \frac{\mu\alpha_p \sigma_D^2 + \mathbf{y}_p \sigma\alpha_p^2}{\sigma_D^2 \sigma_R^2 \sigma\alpha_p^4 + \sigma_D^4 \sigma\alpha_p^2 (\sigma_R^2 + \sigma\alpha_p^2)} \right) / \left( \frac{\sigma_D^2 + \sigma\alpha_p^2}{\sigma_D^2 \sigma_R^2 \sigma\alpha_p^4 + \sigma_D^4 \sigma\alpha_p^2 (\sigma_R^2 + \sigma\alpha_p^2)} \right) \\ = \frac{\mu\alpha_p \sigma_D^2 + \mathbf{y}_p \sigma\alpha_p^2}{\sigma_D^2 + \sigma\alpha_p^2}$$

We now have a rule that tells us how to update the  $\alpha(u)=p(y_i=u|\text{data before } i)$ , in terms of the mean and variance parameters of the previous node:

$$\mu\alpha \leftarrow \frac{\mu\alpha_p \sigma_D^2 + \mathbf{y}_p \sigma\alpha_p^2}{\sigma_D^2 + \sigma\alpha_p^2} = \frac{\frac{\mu\alpha_p \sigma_D^2}{\sigma\alpha_p^2 \sigma_D^2} + \frac{\mathbf{y}_p \sigma\alpha_p^2}{\sigma\alpha_p^2 \sigma_D^2}}{\frac{\sigma_D^2}{\sigma\alpha_p^2 \sigma_D^2} + \frac{\sigma\alpha_p^2}{\sigma\alpha_p^2 \sigma_D^2}} = \frac{\frac{\mu\alpha_p}{\sigma\alpha_p^2} + \frac{\mathbf{y}_p}{\sigma_D^2}}{\frac{1}{\sigma\alpha_p^2} + \frac{1}{\sigma_D^2}}$$

The update rule for the variance is:

$$\sigma\alpha^2 \leftarrow \sigma_R^2 + \frac{1}{\frac{1}{\sigma_D^2} + \frac{1}{\sigma\alpha_p^2}}$$

A similar derivation gives us the rules for  $\mu\beta, \sigma\beta^2$

$$\mu\beta \leftarrow \frac{\frac{\mu\beta_a}{\sigma\beta_a^2} + \frac{\mathbf{y}_a}{\sigma_D^2}}{\frac{1}{\sigma\beta_a^2} + \frac{1}{\sigma\alpha_a^2}}$$

$$\sigma\beta^2 \leftarrow \sigma_R^2 + \frac{1}{\frac{1}{\sigma_b^2} + \frac{1}{\sigma\beta_a^2}}$$

Where the subscript index  $p$  (for "previous", i.e. unit  $i-1$ ) is replaced by  $a$  (for "after", i.e. unit  $i+1$ ).

Recall that sometimes we have data and sometimes we don't. So replace:

$$y_p \rightarrow \text{xs}[i-1] \text{ data}[i-1] = w_{i-1} y_{i-1}^*$$

And similarly for  $y_a$ .

The ratio,  $\left(\frac{\sigma_D}{\sigma_R}\right)^2$  plays the role of  $\lambda$  above. If  $\sigma_D^2 \gg \sigma_R^2$ , there is greater smoothing. If  $\sigma_D^2 \ll \sigma_R^2$ , there is more fidelity to the data. (Recall  $y^* \rightarrow \text{data}.w_k \rightarrow \text{xs}[[k]]$ ). But now we have a principled way of assigning the relative amount of smoothing.

## References

- Applebaum, D. (1996). *Probability and Information*. Cambridge, UK: Cambridge University Press.
- Frey, B. J. (1998). *Graphical Models for Machine Learning and Digital Communication*. Cambridge, Massachusetts: MIT Press.
- Jepson, A., & Black, M. J. (1993). Mixture models for optical flow computation. Paper presented at the Proc. IEEE Conf. Comput. Vision Pattern Recog., New York.
- Kersten, D. and P.W. Schrater (2000), *Pattern Inference Theory: A Probabilistic Approach to Vision, in Perception and the Physical World*, R. Mausfeld and D. Heyer, Editors. , John Wiley & Sons, Ltd.: Chichester. (pdf)
- Kersten, D., & Madarasmí, S. (1995). The Visual Perception of Surfaces, their Properties, and Relationships. *DIMACS Series in Discrete Mathematics and Theoretical Computer Science*, 19, 373-389.
- Knill, D. C., & Pouget, A. (2004). The Bayesian brain: the role of uncertainty in neural coding and computation. *Trends in Neurosciences*, 27(12), 712–719. doi:10.1016/j.tins.2004.10.007
- Ma, W. J., Beck, J. M., Latham, P. E., & Pouget, A. (2006). Bayesian inference with probabilistic population codes. *Nature Neuroscience*, 9(11), 1432–1438. doi:10.1038/nn1790
- Ma, W. J. (2012). Organizing probabilistic models of perception. *Trends in Cognitive Sciences*, 16(10), 511–518. doi:10.1016/j.tics.2012.08.010
- Madarasmí, S., Kersten, D., & Pong, T.-C. (1993). The computation of stereo disparity for transparent and for opaque surfaces. In C. L. Giles & S. J. Hanson & J. D. Cowan (Eds.), *Advances in Neural Information Processing Systems 5*. San Mateo, CA: Morgan Kaufmann Publishers.
- Nakayama, K., & Shimojo, S. (1992). Experiencing and perceiving visual surfaces. *Science*, 257(5075), 1357–1363. doi:10.1126/science.1529336
- Pearl, Judea. (1997) *Probabilistic Reasoning in Intelligent Systems : Networks of Plausible Inference*. (amazon.com link)
- Ripley, B. D. (1996). *Pattern Recognition and Neural Networks*. Cambridge, UK: Cambridge University Press.
- Weiss Y. (1999) *Bayesian Belief Propagation for Image Understanding* submitted to SCTV 1999. (zipped postscript 297K)
- Weiss, Y. (1997). Smoothness in Layers: Motion segmentation using nonparametric mixture estimation. Paper presented at the Proceedings of IEEE conference on Computer Vision and Pattern Recognition.
- Weiss, Y., Simoncelli, E. P., & Adelson, E. H. (2002). Motion illusions as optimal percepts. *Nature Neuroscience*, 5(6), 598–604. doi:10.1038/nn858

Yuille, A., Coughlan J., Kersten D.(1998) (pdf)

Zemel, R. S. & Pouget, A. (1998). Probabilistic interpretation of population codes. *Neural Computation*, 10(2), 403–430.

For notes on Graphical Models, see:<http://www.cs.berkeley.edu/~murphyk/Bayes/bayes.html>

© 2000, 2001, 2003, 2005, 2007, 2009, 2011, 2012, 2014 Daniel Kersten, Computational Vision Lab, Department of Psychology, University of Minnesota.  
(<http://vision.psych.umn.edu/www/kersten-lab/kersten-lab.html>)