

Introduction to Neural Networks  
U. Minn. Psy 5038  
Daniel Kersten  
Belief Propagation

### Initialize

#### ■ Read in Statistical Add-in packages:

```
Off[General::spell1];
<< Statistics`NormalDistribution`
```

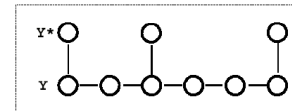
```
<< Statistics`ContinuousDistributions`
<< Statistics`MultinormalDistribution`
```

```
Off[General::spell1];
<< Graphics`MultipleListPlot`
```

### Interpolation using smoothness revisited: 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

Recall that the energy or cost function is given by:

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

where  $w_k = x_s[[k]]$  is the indicator function, and  $y_i^* = d_i$  are the data values.

Gradient descent gives the following local update rule:

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

As before,  $\lambda$  controls the degree of smoothness, i.e. smoothness at the expense of fidelity to the data.

Gauss-Seidel:  $\eta[k\_]:=1/(\lambda+x_s[[k]])$ ;

Successive over-relaxation (SOR):  $\eta2[k\_]:=1.9/(\lambda+x_s[[k]])$ ;

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

### ■ Make the data

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

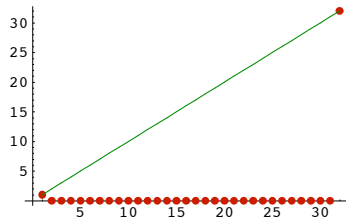
$$w_k = xs[[k]], y^* = \text{data}, y = f.$$

We'll assume the true model is that  $f=y_j$ , where  $j=1$  to  $\text{size}$ . `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; (*xs[[size/2]] = 1;*)
data = Table[N[j] xs[[j]],
{j, 1, size}];
g3 = ListPlot[Table[N[j], {j, 1, size}], PlotJoined -> True,
  DisplayFunction -> Identity, PlotStyle -> {RGBColor[0, .5, 0]}];
g2 = ListPlot[data, PlotJoined -> False,
  PlotStyle -> {RGBColor[.75, .0, 0]}, Prolog -> AbsolutePointSize[5],
  DisplayFunction -> Identity];
```

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

```
Show[g2, g3, DisplayFunction -> $DisplayFunction];
```



Let's set up two matrices,  $\mathbf{Tm}$  and  $\mathbf{Sm}$  such that the gradient of the energy is equal to:

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

$\mathbf{Sm}$  will be our filter to exclude non-data points.  $\mathbf{Tm}$  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];
```

Check the update rule code for small `size=10`:

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

### ■ Run gradient descent

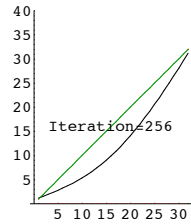
```
Clear[Tf, f1];
dt = 1; λ = 2;
Tf[f1_] := f1 - dt*(1/(λ + xs))*(Tm.f1 - λ*Sm.(data - f1));
```

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

```
iter = 256;
f = Table[0, {i, 1, size}];
result = Nest[Tf, f, iter];
```

Now plot the interpolated function.

```
g1 = ListPlot[result, PlotJoined->True,
  AspectRatio->Automatic, PlotRange->{{1, size}, {1, size}}, DisplayFunction->Identity];
Show[{g1, g2, g3, Graphics[{Text["Iteration=" <> ToString[iter], {size/2, size/2}]}]},
  DisplayFunction->$DisplayFunction, PlotRange->{0, 40}];
```



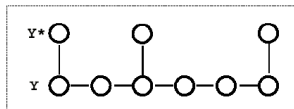
Try starting with  $f$  = random values between 0 and 40. Try various numbers of iterations.

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

## Belief Propagation

### Same interpolation problem, but now using belief propagation

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



### Probabilistic generative model

$$\text{data}[[i]] = y^*[[i]] = xs[[i]] y[[i]] + \text{dnoise}, \quad \text{dnoise} \sim N[0, \sigma_D] \quad (1)$$

$$y[[i+1]] = y[[i]] + \text{znoise}, \quad \text{znoise} \sim N[0, \sigma_R] \quad (2)$$

The first term is the "data formation" model, i.e. how the data is directly influenced by the interaction of the underlying causes,  $y$  with the sampling and noise. The second term reflects our prior assumptions about the smoothness of  $y$ , i.e. nearby  $y$ 's are correlated, and in fact 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.

### Some theory

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 the  $y$  value is, and 2) give a measure of confidence

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

$$p(y_i = u \mid \text{all the data}) \propto p(y_i = u \mid \text{data before } i) p(\text{data at } i \mid y_i = u) p(y_i = u \mid \text{data after } i) \quad (3)$$

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

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

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

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

And the noise model for the data:

$$p(\text{data at } i \mid y_i = u) = L[u] \sim \text{NormalDistribution}[y_p, \sigma_D]$$

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

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:

$$p(y_i = u \mid \text{all the data}) \propto \alpha[u] * L[u] * \beta[u] \quad (4)$$

```

audist = NormalDistribution[μa, σa];
α[u] = PDF[audist, u];

Ddist = NormalDistribution[ $y_p$ , σD];
L[u] = PDF[Ddist, u];

βudist = NormalDistribution[μβ, σβ];
β[u] = PDF[βudist, u];

α[u] * L[u] * β[u]
    
```

$$\frac{e^{-\frac{(u-\mu_a)^2}{2\sigma_a^2} - \frac{(u-\mu_\beta)^2}{2\sigma_\beta^2} - \frac{(u-y_p)^2}{2\sigma_D^2}}}{2\sqrt{2}\pi^{3/2}\sigma_a\sigma_\beta\sigma_D}$$

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.

$$\text{Solve}\left[-D\left[-\frac{(u-\mu_a)^2}{2\sigma_a^2} - \frac{(u-\mu_\beta)^2}{2\sigma_\beta^2} - \frac{(u-y_p)^2}{2\sigma_D^2}, u\right] = 0, u\right]$$

$$\left\{\left\{u \rightarrow \frac{\frac{\mu_a}{\sigma_a^2} + \frac{\mu_\beta}{\sigma_\beta^2} + \frac{y_p}{\sigma_D^2}}{\frac{1}{\sigma_\beta^2} + \frac{1}{\sigma_D^2} + \frac{1}{\sigma_a^2}}\right\}\right\}$$

The update rule for the variance is:

$$\sigma^2 \rightarrow \frac{1}{\sigma_a^2} + \frac{1}{\sigma_\beta^2} + \frac{1}{\sigma_D^2}$$

How do we get  $\mu_\alpha, \mu_\beta, \sigma_\alpha, \mu_\beta$ ?

We express the probability of the  $i$ th unit taking on the value  $u$  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,  $v$ ).

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

$$\alpha[u] = \int_{-\infty}^{\infty} \alpha_p[v] * S[u] * L[v] dv \propto \int_{-\infty}^{\infty} e^{-\frac{(v-y_p)^2}{2\sigma_\beta^2} - \frac{(u-v)^2}{2\sigma_R^2} - \frac{(v-\mu_{\alpha_p})^2}{2\sigma_{\alpha_p}^2}} dv \tag{5}$$

$\alpha_p = \alpha_{i-1}$ .  $S[u]$  is our smoothing term, or transition probability:  $S[u] = p(u | v)$ .  $L[v]$  is the likelihood of the previous data node, given its hidden node value,  $v$ .

```

Rdist = NormalDistribution[v, σR];
S[u] = PDF[Rdist, u];

avdist = NormalDistribution[μαp, σαp];
αp[v] = PDF[avdist, v];

Lp[v] = PDF[Ddist, v];
    
```

```

Integrate[αp[v] * S[u] * Lp[v], {v, -Infinity, Infinity}]
    
```

$$\frac{1}{2\sqrt{2}\pi^{3/2}\sigma_D\sigma_R\sigma_{\alpha_p}} \int_{-\infty}^{\infty} e^{-\frac{(v-y_p)^2}{2\sigma_\beta^2} - \frac{(u-v)^2}{2\sigma_R^2} - \frac{(v-\mu_{\alpha_p})^2}{2\sigma_{\alpha_p}^2}} dv$$

$$\text{If}\left[\text{Re}\left[\frac{1}{\sigma_\beta^2} + \frac{1}{\sigma_R^2} + \frac{1}{\sigma_{\alpha_p}^2}\right] > 0, \frac{e^{-\frac{(u-\mu_{\alpha_p})^2\sigma_\beta^2 + \mu_{\alpha_p}^2\sigma_R^2 + u^2\sigma_{\alpha_p}^2 + y_p^2(\sigma_R^2 + \sigma_\beta^2) - 2y_p(\mu_{\alpha_p}\sigma_R^2 + u\sigma_\beta^2)}{2(\sigma_R^2\sigma_{\alpha_p}^2 + \sigma_\beta^2(\sigma_R^2 + \sigma_{\alpha_p}^2))}}}{\sqrt{\frac{1}{\sigma_\beta^2} + \frac{1}{\sigma_R^2} + \frac{1}{\sigma_{\alpha_p}^2}}}} \sqrt{2\pi}, \right]$$

■ Some uninspired *Mathematica* manipulations

To find an expression for the mode of the above calculated expression for  $\alpha[u]$

$$D\left[-\frac{(u-\mu_{\alpha_p})^2\sigma_\beta^2 + \mu_{\alpha_p}^2\sigma_R^2 + u^2\sigma_{\alpha_p}^2 + y_p^2(\sigma_R^2 + \sigma_\beta^2) - 2y_p(\mu_{\alpha_p}\sigma_R^2 + u\sigma_\beta^2)}{2(\sigma_R^2\sigma_{\alpha_p}^2 + \sigma_\beta^2(\sigma_R^2 + \sigma_{\alpha_p}^2))}, u\right]$$

$$-\frac{2(u-\mu_{\alpha_p})\sigma_\beta^2 + 2u\sigma_{\alpha_p}^2 - 2y_p\sigma_\beta^2}{2(\sigma_R^2\sigma_{\alpha_p}^2 + \sigma_\beta^2(\sigma_R^2 + \sigma_{\alpha_p}^2))}$$

```

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

$$\left\{\left\{u \rightarrow \frac{\frac{\mu_{\alpha_p}\sigma_\beta^2}{\sigma_R^2\sigma_{\alpha_p}^2 + \sigma_\beta^2(\sigma_R^2 + \sigma_{\alpha_p}^2)} + \frac{y_p\sigma_\beta^2}{\sigma_R^2\sigma_{\alpha_p}^2 + \sigma_\beta^2(\sigma_R^2 + \sigma_{\alpha_p}^2)}}{\frac{\sigma_\beta^2}{\sigma_R^2\sigma_{\alpha_p}^2 + \sigma_\beta^2(\sigma_R^2 + \sigma_{\alpha_p}^2)} + \frac{\sigma_{\alpha_p}^2}{\sigma_R^2\sigma_{\alpha_p}^2 + \sigma_\beta^2(\sigma_R^2 + \sigma_{\alpha_p}^2)}}\right\}\right\}$$

$$\text{Simplify} \left[ \left( \frac{\mu_{\alpha_p} \sigma_b^2}{\sigma_R^2 \sigma_{\alpha_p}^2 + \sigma_b^2 (\sigma_R^2 + \sigma_{\alpha_p}^2)} + \frac{Y_p \sigma_{\alpha_p}^2}{\sigma_R^2 \sigma_{\alpha_p}^2 + \sigma_b^2 (\sigma_R^2 + \sigma_{\alpha_p}^2)} \right) / (\sigma_b^2 + \sigma_{\alpha_p}^2) \right]$$

$$\frac{\mu_{\alpha_p} \sigma_b^2 + Y_p \sigma_{\alpha_p}^2}{\sigma_b^2 \sigma_R^2 \sigma_{\alpha_p}^2 + \sigma_b^4 \sigma_{\alpha_p}^2 (\sigma_R^2 + \sigma_{\alpha_p}^2)}$$

$$\text{Simplify} \left[ \left( \frac{\sigma_b^2}{\sigma_R^2 \sigma_{\alpha_p}^2 + \sigma_b^2 (\sigma_R^2 + \sigma_{\alpha_p}^2)} + \frac{\sigma_{\alpha_p}^2}{\sigma_R^2 \sigma_{\alpha_p}^2 + \sigma_b^2 (\sigma_R^2 + \sigma_{\alpha_p}^2)} \right) / (\sigma_b^2 + \sigma_{\alpha_p}^2) \right]$$

$$\frac{\sigma_b^2 + \sigma_{\alpha_p}^2}{\sigma_b^2 \sigma_R^2 \sigma_{\alpha_p}^2 + \sigma_b^4 \sigma_{\alpha_p}^2 (\sigma_R^2 + \sigma_{\alpha_p}^2)}$$

$$\left( \frac{\mu_{\alpha_p} \sigma_b^2 + Y_p \sigma_{\alpha_p}^2}{\sigma_b^2 \sigma_R^2 \sigma_{\alpha_p}^2 + \sigma_b^4 \sigma_{\alpha_p}^2 (\sigma_R^2 + \sigma_{\alpha_p}^2)} \right) / \left( \frac{\sigma_b^2 + \sigma_{\alpha_p}^2}{\sigma_b^2 \sigma_R^2 \sigma_{\alpha_p}^2 + \sigma_b^4 \sigma_{\alpha_p}^2 (\sigma_R^2 + \sigma_{\alpha_p}^2)} \right)$$

$$\frac{\mu_{\alpha_p} \sigma_b^2 + Y_p \sigma_{\alpha_p}^2}{\sigma_b^2 + \sigma_{\alpha_p}^2}$$

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

$$\mu_{\alpha} \leftarrow \frac{\mu_{\alpha_p} \sigma_b^2 + Y_p \sigma_{\alpha_p}^2}{\sigma_b^2 + \sigma_{\alpha_p}^2} = \frac{\frac{\mu_{\alpha_p} \sigma_b^2}{\sigma_{\alpha_p}^2 \sigma_b^2} + \frac{Y_p \sigma_{\alpha_p}^2}{\sigma_{\alpha_p}^2 \sigma_b^2}}{\frac{\sigma_b^2}{\sigma_{\alpha_p}^2 \sigma_b^2} + \frac{\sigma_{\alpha_p}^2}{\sigma_{\alpha_p}^2 \sigma_b^2}} = \frac{\frac{\mu_{\alpha_p}}{\sigma_{\alpha_p}^2} + \frac{Y_p}{\sigma_b^2}}{\frac{1}{\sigma_{\alpha_p}^2} + \frac{1}{\sigma_b^2}}$$

The update rule for the variance is:

$$\sigma_{\alpha}^2 \leftarrow \sigma_R^2 + \frac{1}{\frac{1}{\sigma_b^2} + \frac{1}{\sigma_{\alpha_p}^2}}$$

A similar derivation gives us the rules for  $\mu_{\beta}, \sigma_{\beta}^2$

$$\mu_{\beta} \leftarrow \frac{\frac{\mu_{\beta_a} \sigma_b^2}{\sigma_{\beta_a}^2 \sigma_b^2} + \frac{Y_a \sigma_{\beta_a}^2}{\sigma_{\beta_a}^2 \sigma_b^2}}{\frac{\sigma_b^2}{\sigma_{\beta_a}^2 \sigma_b^2} + \frac{\sigma_{\beta_a}^2}{\sigma_{\beta_a}^2 \sigma_b^2}}$$

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

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:

$$y_p \rightarrow \text{xs}[i-1] \text{ data}[i-1] = w_{i-1} y_{i-1} \tag{6}$$

And similarly for  $y_a$ .

■ Summary of update rules

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  $y^* \rightarrow \text{data}.w_k \rightarrow \text{xs}[[k]]$ )

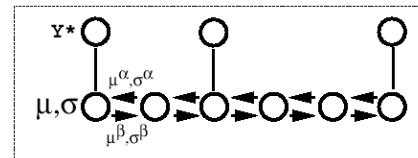
We'll follow Weiss, and also make a (hopefully not too confusing) notation change to avoid the square superscripts for  $\sigma_D^2 \rightarrow \sigma_D, \sigma_R^2 \rightarrow \sigma_R$ .

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

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

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

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



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

### ■ Initialization

```

μ0 = 1;
μα = 1; σα = 100000; (*large uncertainty *)
μβ = 1; σβ = 100000; (*large*)
σR = 4.0; σD = 1.0;
μ = Table[μ0, {i, 1, size}];
σ = Table[σα, {i, 1, size}];
μα = Table[μα, {i, 1, size}];
σα = Table[σα, {i, 1, size}];
μβ = Table[μβ, {i, 1, size}];
σβ = Table[σβ, {i, 1, size}];
iter = 0;
i = 1; j = size;

```

General::spell : Possible spelling error: new symbol name "σD" is similar to existing symbols (σ, σR). More...

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

### ■ Belief Propagation Routine: Execute this cell "manually" for each iteration

```

μ[[i]] = 
$$\frac{\frac{x_s[[i]]}{\sigma_D} * data[[i]] + \frac{1}{\sigma\alpha[[i]]} * \mu\alpha[[i]] + \frac{1.0}{\sigma\beta[[i]]} * \mu\beta[[i]]}{\frac{x_s[[i]]}{\sigma_D} + \frac{1}{\sigma\alpha[[i]]} + \frac{1}{\sigma\beta[[i]]}};$$

σ[[i]] = 
$$\frac{1.0}{\frac{x_s[[i]]}{\sigma_D} + \frac{1}{\sigma\alpha[[i]]} + \frac{1}{\sigma\beta[[i]]}};$$

μ[[j]] = 
$$\frac{\frac{x_s[[j]]}{\sigma_D} * data[[j]] + \frac{1}{\sigma\alpha[[j]]} * \mu\alpha[[j]] + \frac{1.0}{\sigma\beta[[j]]} * \mu\beta[[j]]}{\frac{x_s[[j]]}{\sigma_D} + \frac{1}{\sigma\alpha[[j]]} + \frac{1}{\sigma\beta[[j]]}};$$

σ[[j]] = 
$$\frac{1.0}{\frac{x_s[[j]]}{\sigma_D} + \frac{1}{\sigma\alpha[[j]]} + \frac{1}{\sigma\beta[[j]]}};$$

nextj = j - 1;
μα[[nextj]] = 
$$\frac{\frac{x_s[[j]]}{\sigma_D} * data[[j]] + \frac{1.0}{\sigma\alpha[[j]]} * \mu\alpha[[j]]}{\frac{x_s[[j]]}{\sigma_D} + \frac{1}{\sigma\alpha[[j]]}};$$

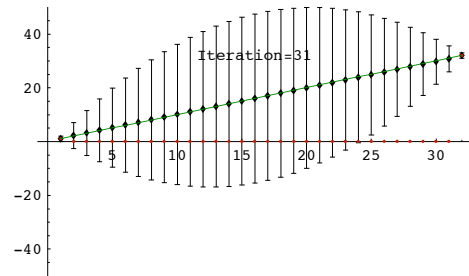
σR[[nextj]] = σR + 
$$\frac{1.0}{\frac{x_s[[j]]}{\sigma_D} + \frac{1}{\sigma\alpha[[j]]}};$$

nexti = i + 1;
μβ[[nexti]] = 
$$\frac{\frac{x_s[[i]]}{\sigma_D} * data[[i]] + \frac{1.0}{\sigma\beta[[i]]} * \mu\beta[[i]]}{\frac{x_s[[i]]}{\sigma_D} + \frac{1}{\sigma\beta[[i]]}};$$

σβ[[nexti]] = σR + 
$$\frac{1.0}{\frac{x_s[[i]]}{\sigma_D} + \frac{1}{\sigma\beta[[i]]}};$$

j--;
i++;
iter++;
yfit = Table[{μ[[i1]], ErrorBar[σ[[i1]]]}, {i1, 1, size}];
glb = MultipleListPlot[yfit, DisplayFunction -> Identity];
Show[
{glb, g2, g3, Graphics[{Text["Iteration=" <> ToString[iter], {size/2, size}]}]},
DisplayFunction -> $DisplayFunction, PlotRange -> {-50, 50}];

```



## Expectation Maximization (EM) -- A preview & simulation

We've studied the problem of interpolation given missing data. We motivated the problem by the visual phenomenon of surface completion. Another aspect of surface perception is our ability to take noisy data (e.g. depth cues), and not only interpolate the data, but also decide which of several surfaces the data belong using stereo data (Madarasmi et al., 1993; Kersten & Madarasmi, 1995) or optic flow. A number of ways have been proposed to deal with this problem, but the technique of most general application is expectation-maximization--a general statistical technique developed in the 1970's that appears in various forms in many algorithms, including belief propagation. The algorithm has been applied to the surface estimation problem, e.g. from optic flow (Jepson, 1993; Weiss, 1997). We go to Weiss again for a nice simple demo.

Consider **Generative Model 1**.

### Generative models

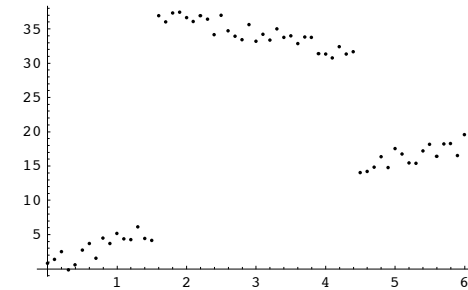
Two lines with slopes and intercepts (a1,b1) and (a2,b2).

```
ndist = NormalDistribution[0, 5];
```

#### ■ Generative model 1

```
y1[x_] := -2 * x + 40 + Random[ndist];
y2[x_] := 3 * x + 1 + Random[ndist];
data = Table[{x, If[Abs[x - 3] < 1.5, y1[x], y2[x]]}, {x, 0, 6, .1}];
{x, y} = Transpose[data];
```

```
gdata = ListPlot[data];
```



#### ■ Generative model 2

```
y1[x_] := -2 * x + 15;
y2[x_] := 3 * x + 1;
data = Table[{x, If[Random[] < .5, y1[x], y2[x]]}, {x, 0, 6, .1}];
{x, y} = Transpose[data];
```

```
gdata = ListPlot[data];
```

Which data belong together? We'll assume that we know there are two linear models, but we don't know their parameters--i.e. we don't know their slopes and intercepts.

### EM algorithm

#### ■ Initialize parameters to random values

```
σ = .1;
{a1, b1, a2, b2} = Table[10 * Random[], {4}];
```

#### ■ E-step

Compute residuals  $r_1, r_2$ , the error in the predicted and actual  $y$  values under each of the two models.

```
r1 = a1 * x + b1 - y;
r2 = a2 * x + b2 - y;
```

Using the residuals, compute weights. We'll assign these weights to data in the M-step later.

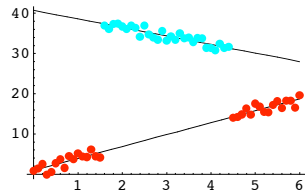
```
w1 = Exp[-r1^2 / σ] / (Exp[-r1^2 / σ] + Exp[-r2^2 / σ]);
w2 = Exp[-r2^2 / σ] / (Exp[-r1^2 / σ] + Exp[-r2^2 / σ]);
```

### ■ M-step

With above, weights compute weighted linear regression:

```
{a1, b1} = Inverse[{{w1.(x x), w1.x}, {w1.x, Apply[Plus, w1]}}].{w1.(x y), w1.y};
{a2, b2} = Inverse[{{w2.(x x), w2.x}, {w2.x, Apply[Plus, w2]}}].{w2.(x y), w2.y};
```

```
gfit = Plot[{a1 * x + b1, a2 * x + b2}, {x, 0, 6}, DisplayFunction -> Identity];
Show[
  {gfit, Graphics[{PointSize[0.03], Transpose[{Hue[#1] & /@ (w1 / 2), Point /@ data}]}]},
  PlotRange -> All, DisplayFunction -> $DisplayFunction];
```



Try repeating the E-step followed by the M-step

## Exercises

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

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

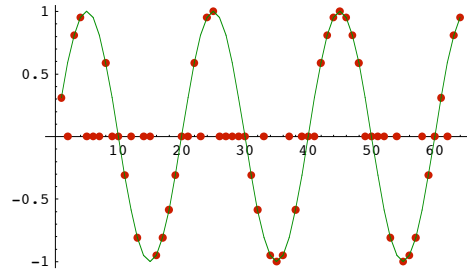
( $\sigma_R = 4.0$ ;  $\sigma_D = 1.0$ ).

BP with missing sine wave data

### ■ Generate sine wave with missing data

```
size = 64;
xs = Table[Random[Integer, 1], {i, 1, size}];
data = Table[N[Sin[2 Pi (1/20) j] xs[[j]]],
  {j, 1, size}];
g3b = ListPlot[Table[N[Sin[2 Pi (1/20) j]], {j, 1, size}],
  PlotJoined -> True, DisplayFunction -> Identity,
  PlotStyle -> {RGBColor[0, .5, 0]}];
g2b = ListPlot[data, PlotJoined -> False,
  PlotStyle -> {RGBColor[.75, .0, 0]}, Prolog -> AbsolutePointSize[5],
  DisplayFunction -> Identity];
Show[g2b, g3b, DisplayFunction -> $DisplayFunction];
```





### ■ Initialize

```

μ0 = 1;
μα = 1; σα = 100000; (*large uncertainty *)
μβ = 1; σβ = 100000; (*large*)
σR = .5; σD = .1;
μ = Table[μ0, {i, 1, size}];
σ = Table[σα, {i, 1, size}];
μα = Table[μα, {i, 1, size}];
σα = Table[σα, {i, 1, size}];
μβ = Table[μβ, {i, 1, size}];
σβ = Table[σβ, {i, 1, size}];
iter = 0;
i = 1; j = size;

```

### ■ SINE WAVE DEMO: Belief Propagation Routine: Execute this cell "manually" for each iteration

```

μ[[i]] = (xs[[i]] * data[[i]] +  $\frac{1}{\sigma\alpha}$  * μα[[i]] +  $\frac{1.0}{\sigma\beta}$  * μβ[[i]]) /
  (  $\frac{x\alpha[[i]]}{\sigma\alpha}$  +  $\frac{1}{\sigma\alpha[[i]]}$  +  $\frac{1}{\sigma\beta[[i]]}$  );
σ[[i]] =  $\frac{1.0}{\frac{x\alpha[[i]]}{\sigma\alpha} + \frac{1}{\sigma\alpha[[i]]} + \frac{1}{\sigma\beta[[i]]}}$ ;
μ[[j]] = (xs[[j]] * data[[j]] +  $\frac{1}{\sigma\alpha}$  * μα[[j]] +  $\frac{1.0}{\sigma\beta}$  * μβ[[j]]) /
  (  $\frac{x\alpha[[j]]}{\sigma\alpha}$  +  $\frac{1}{\sigma\alpha[[j]]}$  +  $\frac{1}{\sigma\beta[[j]]}$  );
σ[[j]] =  $\frac{1.0}{\frac{x\alpha[[j]]}{\sigma\alpha} + \frac{1}{\sigma\alpha[[j]]} + \frac{1}{\sigma\beta[[j]]}}$ ;

nextj = j - 1;
μα[[nextj]] =  $\frac{\frac{x\alpha[[j]]}{\sigma\alpha} * data[[j]] + \frac{1.0}{\sigma\alpha} * \mu\alpha[[j]]}{\frac{x\alpha[[j]]}{\sigma\alpha} + \frac{1}{\sigma\alpha[[j]}}$ ;
σR[[nextj]] = σR +  $\frac{1.0}{\frac{x\alpha[[j]]}{\sigma\alpha} + \frac{1}{\sigma\alpha[[j]}}$ ;

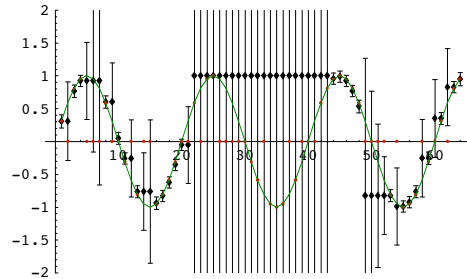
nexti = i + 1;
μβ[[nexti]] =  $\frac{\frac{x\beta[[i]]}{\sigma\beta} * data[[i]] + \frac{1.0}{\sigma\beta} * \mu\beta[[i]]}{\frac{x\beta[[i]]}{\sigma\beta} + \frac{1}{\sigma\beta[[i]}}$ ;
σβ[[nexti]] = σβ +  $\frac{1.0}{\frac{x\beta[[i]]}{\sigma\beta} + \frac{1}{\sigma\beta[[i]}}$ ;

j--;
i++;

iter++;

yfit = Table[{μ[[i1]], ErrorBar[σ[[i1]]]}, {i1, 1, size}];
g1bb = MultipleListPlot[yfit, DisplayFunction -> Identity];
Show[{g1bb, g2b, g3b}, DisplayFunction -> $DisplayFunction, PlotRange -> {-2, 2}];

```



Run EM with Generative Model 2. Increase the additive noise. How does attribution accuracy change? ("attribution" means assigning a point to its correct line)

## 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., & Madarasmi, S. (1995). The Visual Perception of Surfaces, their Properties, and Relationships. *DIMACS Series in Discrete Mathematics and Theoretical Computer Science*, 19, 373-389.
- Madarasmi, 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.
- 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. (gzipped 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.
- Yuille, A., Coughlan J., Kersten D.(1998) (pdf)

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