

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`
```

Last time

Generative modeling: Multivariate gaussian, mixtures

- Drawing samples
- Mixtures of gaussians
- Will mixture distributions in the next lecture on EM application to segmentation

Introduction to Optimal inference and task

Optimal Inference and task dependence: Fruit example

(due to James Coughlan; see Yuille, Coughlan, Kersten & Schrater).

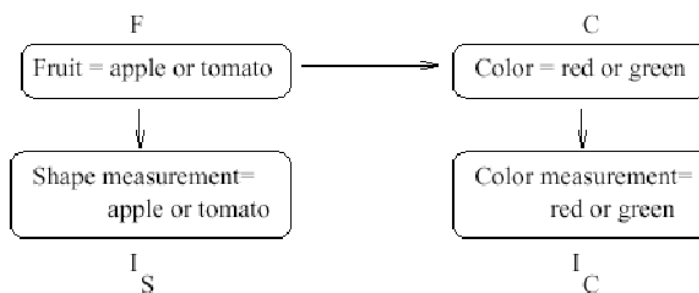


Figure from Yuille, Coughlan, Kersten & Schrater.

The graph specifies how to decompose the joint probability:

$$p[F, C, I_S, I_C] = p[I_C | C] p[C | F] p[I_S | F] p[F]$$

The prior model on hypotheses, F & C

More apples (F=1) than tomatoes (F=2), and:

```
ppF[F_] := If[F == 1, 9 / 16, 7 / 16];
TableForm[Table[ppF[F], {F, 1, 2}], TableHeadings -> {"F=a", "F=t"}]
```

F=a	$\frac{9}{16}$
F=t	$\frac{7}{16}$

The conditional probability cpCF[C|F]:

```
cpCF[F_, C_] := Which[F == 1 && C == 1, 5 / 9, F == 1 && C == 2, 4 / 9,
  F == 2 && C == 1, 6 / 7, F == 2 && C == 2, 1 / 7];
TableForm[Table[cpCF[F, C], {F, 1, 2}, {C, 1, 2}],
  TableHeadings -> {"F=a", "F=t"}, {"C=r", "C=g"}]
```

	C=r	C=g
F=a	$\frac{5}{9}$	$\frac{4}{9}$
F=t	$\frac{6}{7}$	$\frac{1}{7}$

So the joint is:

```
jpFC[F_, C_] := cpCF[F, C] ppF[F];
TableForm[Table[jpFC[F, C], {F, 1, 2}, {C, 1, 2}],
  TableHeadings -> {"F=a", "F=t"}, {"C=r", "C=g"}]
```

	C=r	C=g
F=a	$\frac{5}{16}$	$\frac{1}{4}$
F=t	$\frac{3}{8}$	$\frac{1}{16}$

We can marginalize to get the prior probability on color alone is:

$$ppC[C_] := \sum_{F=1}^2 jpFC[F, C]$$

Question: Is fruit identity independent of material color--i.e. is F independent of C?

■ Answer

No.

```
TableForm[Table[jpFC[F, C], {F, 1, 2}, {C, 1, 2}],
  TableHeadings -> {"F=a", "F=t"}, {"C=r", "C=g"}]]
TableForm[Table[ppF[F] ppC[C], {F, 1, 2}, {C, 1, 2}],
  TableHeadings -> {"F=a", "F=t"}, {"C=r", "C=g"}]]
```

	C=r	C=g
F=a	$\frac{5}{16}$	$\frac{1}{4}$
F=t	$\frac{3}{8}$	$\frac{1}{16}$

	C=r	C=g
F=a	$\frac{99}{256}$	$\frac{45}{256}$
F=t	$\frac{77}{256}$	$\frac{35}{256}$

The generative model: Imaging probabilities

Suppose that we have gathered some "image statistics" which provides us knowledge of how the image measurements for shape I_s , and for color I_c depend on the type of fruit F , and material color, C . For simplicity, our measurements are discrete and binary (a more realistic case, they would have continuous values), say $I_s = \{am, tm\}$, and $I_c = \{rm, gm\}$.

$$P(I_S=am,tm \mid F=a) = \{11/16, 5/16\}$$

$$P(I_S=am,tm \mid F=t) = \{5/8, 3/8\}$$

$$P(I_C=rm,gm \mid C=r) = \{9/16, 7/16\}$$

$$P(I_C=rm,gm \mid C=g) = \{1/2, 1/2\}$$

We use the notation am , tm , rm , gm because the measurements are already suggestive of the likely cause. So there is a correlation between apple and apple-like shapes, am ; and between red material, and "red" measurements. In general, there may not be an obvious correlation like this.

We define a function for the probability of I_c given C , $cpIc[Ic \mid C]$:

```

cpIcC[Ic_, C_] := Which[Ic == 1 && C == 1, 9 / 16, Ic == 1 && C == 2, 7 / 16,
  Ic == 2 && C == 1, 1 / 2, Ic == 2 && C == 2, 1 / 2];
TableForm[Table[cpIcC[Ic, C], {Ic, 1, 2}, {C, 1, 2}],
  TableHeadings -> {"Ic=rm", "Ic=gm"}, {"C=r", "C=g"}]]

```

	C=r	C=g
Ic=rm	$\frac{9}{16}$	$\frac{7}{16}$
Ic=gm	$\frac{1}{2}$	$\frac{1}{2}$

The probability of Is conditional on F is **cpIsF**[Is | F]:

```

cpIsF[Is_, F_] := Which[Is == 1 && F == 1, 11 / 16, Is == 1 && F == 2, 5 / 8,
  Is == 2 && F == 1, 5 / 16, Is == 2 && F == 2, 3 / 8];
TableForm[Table[cpIsF[Is, F], {Is, 1, 2}, {F, 1, 2}],
  TableHeadings -> {"Is=am", "Is=tm"}, {"F=a", "F=t"}]]

```

	F=a	F=t
Is=am	$\frac{11}{16}$	$\frac{5}{8}$
Is=tm	$\frac{5}{16}$	$\frac{3}{8}$

The total joint probability

We now have enough information to put probabilities on the 2x2x2 "universe" of possibilities, i.e. all possible combinations of fruit, color, and image measurements. Looking at the graphical model makes it easy to use the product rule to construct the total joint, which is:

$$p[F, C, Is, Ic] = p[Ic | C] p[C | F] p[Is | F] p[F]:$$

```

jpFCIsIc[F_, C_, Is_, Ic_] := cpIcC[Ic, C] cpCF[F, C] cpIsF[Is, F] ppF[F]

```

Usually, we don't need the probabilities of the image measurements (because once the measurements are made, they are fixed and we want to compare the probabilities of the hypotheses. But in our simple case here, once we have the joint, we can calculate the probabilities of the image measurements through marginalization $p(Is, Ic) = \sum_C \sum_F p(F, C, Is, Ic)$, too:

$$jpIsIc[Is_, Ic_] := \sum_{C=1}^2 \sum_{F=1}^2 jpFCIsIc[F, C, Is, Ic]$$

Three MAP tasks

We are going to show that the best guess (i.e. maximum probability) depends on the task.

■ Define argmax[] function:

```
argmax[x_] := Position[x, Max[x]];
```

■ Pick most probable fruit AND color--Answer "red tomato"

First, suppose the task is to make the best bet as to the fruit AND material color. To make it concrete, suppose that we see an "apple-like shape" with a reddish color, i.e., we measure $I_s=am=1$, and $I_c = rm=1$. The measurements suggest "red apple", but to find the most probable, we need to take into account the priors too in order to make the best guesses.

Using the total joint, $p(F,C | I_s, I_c) = \frac{p(F,C,I_s,I_c)}{p(I_s,I_c)} \propto p(F,C,I_s=1, I_c=1)$

```
TableForm[jpFCIsIcTable = Table[jpFCIsIc[F, C, 1, 1], {F, 1, 2}, {C, 1, 2}],
  TableHeadings -> {"F=a", "F=t"}, {"C=r", "C=g"}]
Max[jpFCIsIcTable]
argmax[jpFCIsIcTable]
```

	C=r	C=g
F=a	$\frac{495}{4096}$	$\frac{77}{1024}$
F=t	$\frac{135}{1024}$	$\frac{35}{2048}$

```
 $\frac{135}{1024}$ 
```

```
{{2, 1}}
```

"Red tomato" is the most probable once we take into account the difference in priors.

Calculating $p(F,C | I_s, I_c)$. We didn't actually need $p(F,C | I_s, I_c)$, but we can calculate it by conditioning the total joint on the probability of the measurements:

```
jpFCcIsIc[F_, C_, Is_, Ic_] := jpFCIsIc[F, C, Is, Ic] / jpIsIc[Is, Ic]
```

```
TableForm[
  jpFCcIsIcTable = Table[jpFCcIsIc[F, C, 1, 1], {F, 1, 2}, {C, 1, 2}],
  TableHeadings -> {"F=a", "F=t"}, {"C=r", "C=g"}]
Max[jpFCcIsIcTable]
argmax[jpFCcIsIcTable]
```

	C=r	C=g
F=a	$\frac{55}{157}$	$\frac{308}{1413}$
F=t	$\frac{60}{157}$	$\frac{70}{1413}$

$$\frac{60}{157}$$

```
{{2, 1}}
```

■ Pick most probable color--Answer "red"

Same measurements as before. But now suppose we only care about the true material color, and not the identity of the object. Then we want to integrate out or marginalize with respect to the shape or fruit-type variable, F. In this case, we want to maximize the posterior:

$$p(C | Is=1, Ic=1) = \sum_{F=1}^2 p(F, C | Is = 1, Ic = 1)$$

$$pC[C_, Is_, Ic_] := \sum_{F=1}^2 jpFCcIsIc[F, C, Is, Ic]$$

```
TableForm[pCTable = Table[pC[C, 1, 1], {C, 1, 2}],
  TableHeadings -> {"C=r", "C=g"}]
Max[pCTable]
argmax[pCTable]
```

C=r	$\frac{115}{157}$
C=g	$\frac{42}{157}$

$$\frac{115}{157}$$

```
{{1}}
```

Answer is that the most probable material color is C = r, "red".

■ Pick most probable fruit--Answer "apple"

Same measurements as before. But now, we don't care about the material color, just the identity of the fruit. Then we want to integrate out or marginalize with respect to the material variable, C. In this case, we want to maximize the posterior:

$p(F | I_s, I_c)$

$$p_F[F_, I_s_, I_c_] := \sum_{C=1}^2 j p_{FCcI_sI_c}[F, C, I_s, I_c]$$

```
TableForm[pFTable = Table[pF[F, 1, 1], {F, 1, 2}],
  TableHeadings -> {"F=a", "F=t"}]
Max[pFTable]
argmax[pFTable]
```

F=a	$\frac{803}{1413}$
F=t	$\frac{610}{1413}$

$$\frac{803}{1413}$$

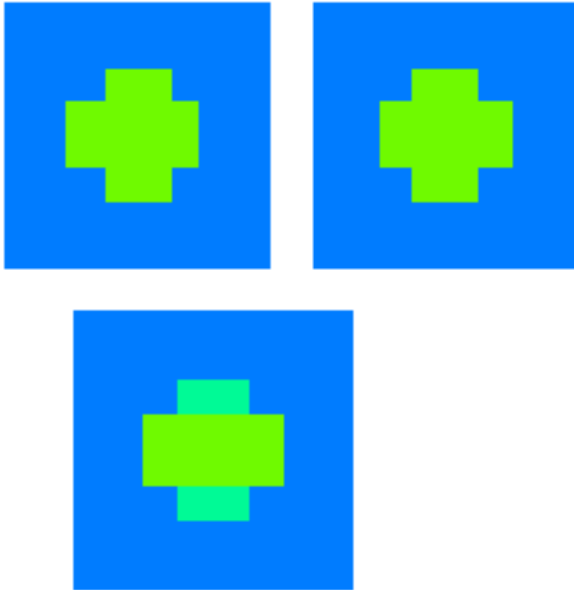
```
{{1}}
```

The answer is "apple". So to sum up, for the same data measurements, the most probable fruit AND color is "red tomato", but the most probable fruit is "apple"!

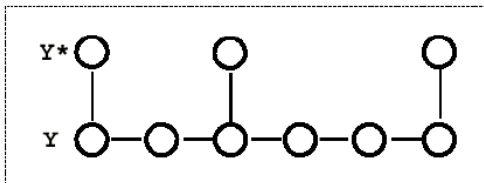
■ Important "take-home message": *Optimal inference depends on the precise definition of the task*

Try expressing the consequences using the frequency interpretation of probability.

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 $w_k = xs[[k]]$ is the indicator function, and $y_k^* = d$, 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, λ 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): $\eta_2[k_]:=1.9/(\lambda+x_s[[k]])$;

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:

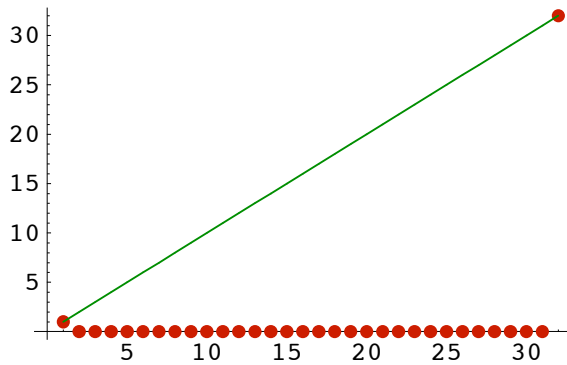
$w_k = xs[[k]]$, $y^* = \text{data}$, $y=f$.

We'll assume the true model is that $f=y=j$, where $j=1$ to 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 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, **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. **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:

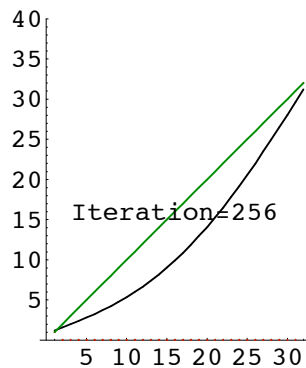
```
iter=256;
f = Table[0, {i, 1, size}];
result = Nest[Tf, f, iter];
```

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}];

```



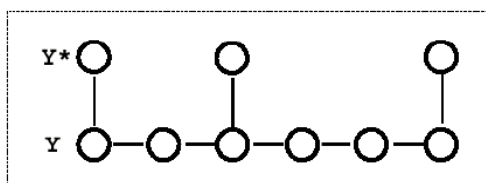
Try starting with f = random values between 0 and 40. Try various numbers of iterations.

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

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

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

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

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

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

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

$$p(Y_i = u \text{ lall the data}) \propto p(Y_i = u \text{ data before } i) p(\text{data at } i | Y_i = u) p(Y_i = u \text{ data after } i) \quad (3)$$

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

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

and so is probability conditioned on the data after i :

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

And the noise model for the data:

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

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

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:

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

```

αdist = NormalDistribution[μ $\alpha$ , σ $\alpha$ ];
α[u] = PDF[αdist, u];

Ddist = NormalDistribution[y $_p$ , σ $_D$ ];
L[u] = PDF[Ddist, u];

βdist = NormalDistribution[μ $\beta$ , σ $\beta$ ];
β[u] = PDF[βdist, u];

α[u] * L[u] * β[u]

```

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

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.

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

$$\left\{\left\{u \rightarrow \frac{\frac{\mu\alpha}{\sigma\alpha^2} + \frac{\mu\beta}{\sigma\beta^2} + \frac{y_p}{\sigma_D^2}}{\frac{1}{\sigma\beta^2} + \frac{1}{\sigma_D^2} + \frac{1}{\sigma\alpha^2}}\right\}\right\}$$

The update rule for the variance is:

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

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 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 , at the $i-1$ th location).

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

$$\alpha[u] = \int_{-\infty}^{\infty} \alpha_p[v] * S[u] * L[v] dv \propto \int_{-\infty}^{\infty} e^{-\frac{(v-y_p)^2}{2\sigma_D^2} - \frac{(u-v)^2}{2\sigma_R^2} - \frac{(v-\mu\alpha_p)^2}{2\sigma\alpha_p^2}} dv \quad (5)$$

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

$$\frac{1}{2\sqrt{2}\pi^{3/2}\sigma_D\sigma_R\sigma\alpha_p} \text{If}\left[\text{Re}\left[\frac{1}{\sigma_D^2} + \frac{1}{\sigma_R^2} + \frac{1}{\sigma\alpha_p^2}\right] > 0,\right.$$

$$\frac{e^{-\frac{(u-\mu\alpha_p)^2\sigma_D^2 + \mu\alpha_p^2\sigma_R^2 + u^2\sigma\alpha_p^2 + y_p^2(\sigma_R^2 + \sigma\alpha_p^2) - 2y_p(\mu\alpha_p\sigma_R^2 + u\sigma\alpha_p^2)}{2(\sigma_R^2\sigma\alpha_p^2 + \sigma_D^2(\sigma_R^2 + \sigma\alpha_p^2))}}}{\sqrt{\frac{1}{\sigma_D^2} + \frac{1}{\sigma_R^2} + \frac{1}{\sigma\alpha_p^2}}},$$

$$\int_{-\infty}^{\infty} e^{-\frac{(v-y_p)^2}{2\sigma_D^2} - \frac{(u-v)^2}{2\sigma_R^2} - \frac{(v-\mu\alpha_p)^2}{2\sigma\alpha_p^2}} dv]$$

■ Some uninspired *Mathematica* manipulations

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

$$D\left[-\frac{(u-\mu\alpha_p)^2\sigma_D^2 + \mu\alpha_p^2\sigma_R^2 + u^2\sigma\alpha_p^2 + y_p^2(\sigma_R^2 + \sigma\alpha_p^2) - 2y_p(\mu\alpha_p\sigma_R^2 + u\sigma\alpha_p^2)}{2(\sigma_R^2\sigma\alpha_p^2 + \sigma_D^2(\sigma_R^2 + \sigma\alpha_p^2))}, u\right]$$

$$-\frac{2(u-\mu\alpha_p)\sigma_D^2 + 2u\sigma\alpha_p^2 - 2y_p\sigma\alpha_p^2}{2(\sigma_R^2\sigma\alpha_p^2 + \sigma_D^2(\sigma_R^2 + \sigma\alpha_p^2))}$$

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

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

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

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

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

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

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

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

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:

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

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 $\mu\beta$, $\sigma\beta^2$

$$\mu\beta \leftarrow \frac{\frac{\mu\beta_a + Y_a}{\sigma\beta_a^2 \sigma_D^2}}{\frac{1}{\sigma\beta_a^2} + \frac{1}{\sigma_D^2}}$$

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

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:

$$Y_p \rightarrow \text{xs}[i-1] \text{ data}[i-1] = w_{i-1} Y_{i-1}^* \quad (6)$$

And similarly for y_a .

■ Summary of update rules

The ratio, $\left(\frac{\sigma_D}{\sigma_R}\right)^2$ plays the role of λ 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 $y^* \rightarrow \text{data}.w_k \rightarrow \text{xs}[[k]]$)

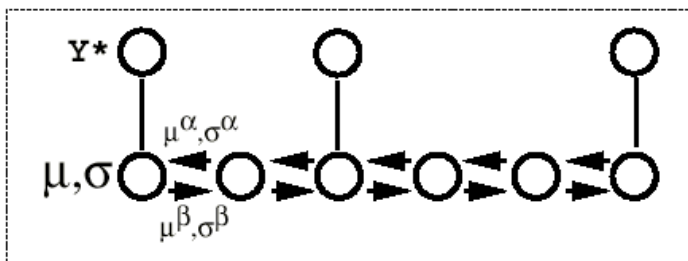
We'll follow Weiss, and also make a (hopefully not too confusing) notation change to avoid the square superscripts for $\sigma_D^2 \rightarrow \sigma_D, \sigma_R^2 \rightarrow \sigma_R$.

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

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

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

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



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[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;
```

General::spell :

Possible spelling error: new symbol name " σD " is similar to existing symbols $\{\sigma, \sigma R\}$. More...

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: Execute this cell "manually" for each iteration

```


$$\mu[[i]] = \frac{\frac{xs[[i]]}{\sigma_D} * data[[i]] + \frac{1}{\sigma_\alpha[[i]]} * \mu_\alpha[[i]] + \frac{1.0}{\sigma_\beta[[i]]} * \mu_\beta[[i]]}{\frac{xs[[i]]}{\sigma_D} + \frac{1}{\sigma_\alpha[[i]]} + \frac{1}{\sigma_\beta[[i]]}};$$


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


$$\mu[[j]] = \frac{\frac{xs[[j]]}{\sigma_D} * data[[j]] + \frac{1}{\sigma_\alpha[[j]]} * \mu_\alpha[[j]] + \frac{1.0}{\sigma_\beta[[j]]} * \mu_\beta[[j]]}{\frac{xs[[j]]}{\sigma_D} + \frac{1}{\sigma_\alpha[[j]]} + \frac{1}{\sigma_\beta[[j]]}};$$


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

nextj = j - 1;

$$\mu_\alpha[[nextj]] = \frac{\frac{xs[[j]]}{\sigma_D} * data[[j]] + \frac{1.0}{\sigma_\alpha[[j]]} * \mu_\alpha[[j]]}{\frac{xs[[j]]}{\sigma_D} + \frac{1}{\sigma_\alpha[[j]]}};$$


$$\sigma_\alpha[[nextj]] = \sigma_R + \frac{1.0}{\frac{xs[[j]]}{\sigma_D} + \frac{1}{\sigma_\alpha[[j]]}};$$

nexti = i + 1;

$$\mu_\beta[[nexti]] = \frac{\frac{xs[[i]]}{\sigma_D} * data[[i]] + \frac{1.0}{\sigma_\beta[[i]]} * \mu_\beta[[i]]}{\frac{xs[[i]]}{\sigma_D} + \frac{1}{\sigma_\beta[[i]]}};$$

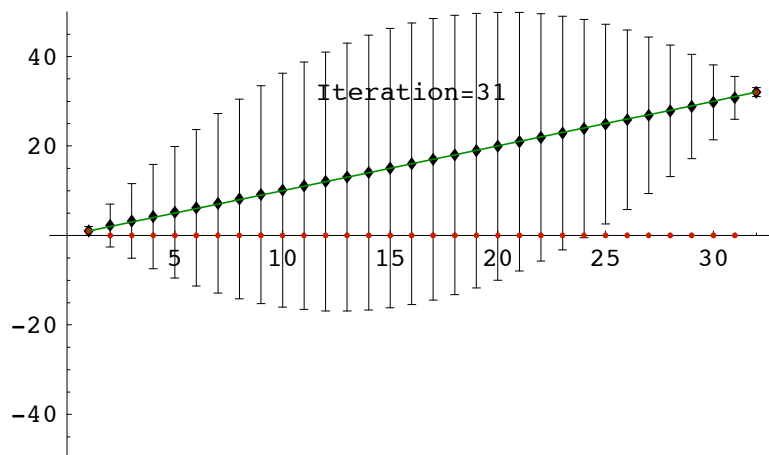

$$\sigma_\beta[[nexti]] = \sigma_R + \frac{1.0}{\frac{xs[[i]]}{\sigma_D} + \frac{1}{\sigma_\beta[[i]]}};$$

j--;
i++;

iter++;

yfit = Table[{μ[[i1]], ErrorBar[σ[[i1]]]}, {i1, 1, size}];
g1b = MultipleListPlot[yfit, DisplayFunction -> Identity];
Show[
  {g1b, g2, g3,
    Graphics[{Text["Iteration=" <> ToString[iter], {size/2, size}]}],
    DisplayFunction -> $DisplayFunction, PlotRange -> {-50, 50}];

```



Exercises

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

How does convergence compare with Gauss-Seidel?

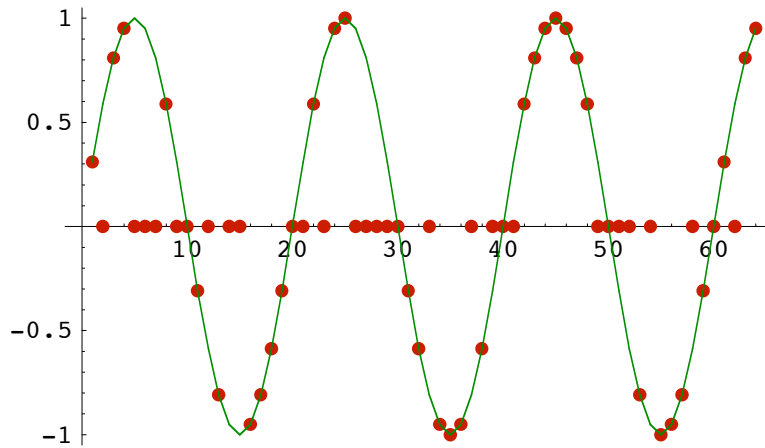
Run Belief Propagation using: $\sigma_R=1.0$; $\sigma_D=4.0$; How does fidelity to the data compare with the original case

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

BP with missing sine wave data

■ Generate sine wave with missing data

```
size = 64;
xs = Table[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[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;

```

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

```


$$\mu[[i]] = \left( \frac{\text{xs}[[i]]}{\sigma_D} * \text{data}[[i]] + \frac{1}{\sigma_\alpha[[i]]} * \mu_\alpha[[i]] + \frac{1.0}{\sigma_\beta[[i]]} * \mu_\beta[[i]] \right) /$$


$$\left( \frac{\text{xs}[[i]]}{\sigma_D} + \frac{1}{\sigma_\alpha[[i]]} + \frac{1}{\sigma_\beta[[i]]} \right);$$


$$\sigma[[i]] = \frac{1.0}{\frac{\text{xs}[[i]]}{\sigma_D} + \frac{1}{\sigma_\alpha[[i]]} + \frac{1}{\sigma_\beta[[i]]}};$$


$$\mu[[j]] = \left( \frac{\text{xs}[[j]]}{\sigma_D} * \text{data}[[j]] + \frac{1}{\sigma_\alpha[[j]]} * \mu_\alpha[[j]] + \frac{1.0}{\sigma_\beta[[j]]} * \mu_\beta[[j]] \right) /$$


$$\left( \frac{\text{xs}[[j]]}{\sigma_D} + \frac{1}{\sigma_\alpha[[j]]} + \frac{1}{\sigma_\beta[[j]]} \right);$$


$$\sigma[[j]] = \frac{1.0}{\frac{\text{xs}[[j]]}{\sigma_D} + \frac{1}{\sigma_\alpha[[j]]} + \frac{1}{\sigma_\beta[[j]]}};$$

nextj = j - 1;

$$\mu_\alpha[[\text{nextj}]] = \frac{\frac{\text{xs}[[j]]}{\sigma_D} * \text{data}[[j]] + \frac{1.0}{\sigma_\alpha[[j]]} * \mu_\alpha[[j]]}{\frac{\text{xs}[[j]]}{\sigma_D} + \frac{1}{\sigma_\alpha[[j]]}};$$


$$\sigma_\alpha[[\text{nextj}]] = \sigma_R + \frac{1.0}{\frac{\text{xs}[[j]]}{\sigma_D} + \frac{1}{\sigma_\alpha[[j]]}};$$

nexti = i + 1;

$$\mu_\beta[[\text{nexti}]] = \frac{\frac{\text{xs}[[i]]}{\sigma_D} * \text{data}[[i]] + \frac{1.0}{\sigma_\beta[[i]]} * \mu_\beta[[i]]}{\frac{\text{xs}[[i]]}{\sigma_D} + \frac{1}{\sigma_\beta[[i]]}};$$

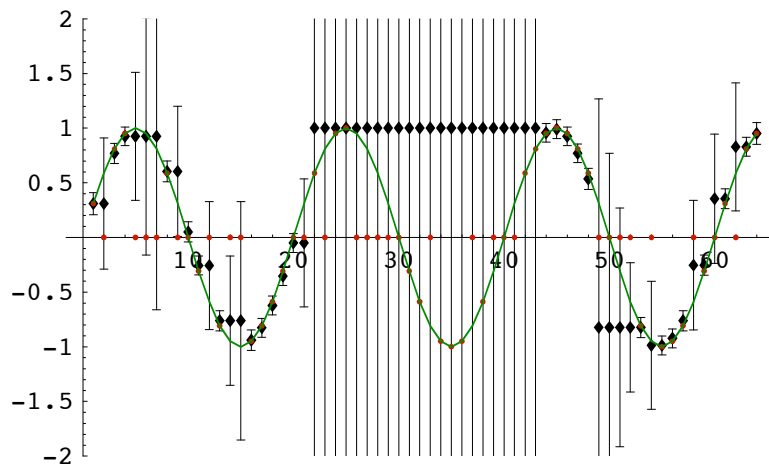

$$\sigma_\beta[[\text{nexti}]] = \sigma_R + \frac{1.0}{\frac{\text{xs}[[i]]}{\sigma_D} + \frac{1}{\sigma_\beta[[i]]}};$$

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., & Madarasmı, S. (1995). The Visual Perception of Surfaces, their Properties, and Relationships. *DIMACS Series in Discrete Mathematics and Theoretical Computer Science*, 19, 373-389.
- Madarasmı, 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>