

Introduction to Neural Networks

U. Minn. Psy 5038

Unifying neural network computations using Bayesian decision theory

Initialize

- Read in Statistical Add-in packages:

```
Off[General::spell1];  
<< Statistics`DescriptiveStatistics`  
<< Statistics`DataManipulation`  
<< Statistics`NormalDistribution`  
<< Statistics`MultiDescriptiveStatistics`  
<< Statistics`MultinormalDistribution`  
<< Statistics`LinearRegression`  
<< Graphics`MultipleListPlot`
```

Previously

Fisher's Linear Discriminant

Statistical learning, model selection & the bias/variance dilemma

Above we summarized optimal rules for minimizing risk, assuming that we know the distributions of the generative model.

But what if we don't? This is the topic of statistical learning theory.

Use the following link for the notes:

http://gandalf.psych.umn.edu/~kersten/kersten-lab/courses/Psy5038WF2003/MathematicaNotebooks/Lect_27_BiasVariance/biasvarianceNotes.pdf

Consider the regression problem, fitting data that may be a complex function of the input.

The problem in general is how to choose the function that both remembers the relationship between \mathbf{x} and \mathbf{y} , and generalizes with new values of \mathbf{x} . At first one might think that it should be as general as possible to allow all kinds of maps.

For example, if one is fitting a curve, you might wish to use a very high-order polynomial, or a back-prop network with lots of hidden units. There is a draw back to the flexibility afforded by extra degrees of freedom in fitting the data. We can get drastically different fits for different sets of data that are randomly drawn from the same underlying process. The fact that we get different fit parameters (e.g. slope of a regression line) each time means that although we may exactly fit the data each time, we introduce variation between the average fit (over all data sets) and the fit for a single data set. We could get around this problem with a huge amount of data, but the problem is that the amount of required data can grow exponentially with the order of the fit--an example of the so-called "curse of dimensionality".

On the other hand, if the function is restrictive, (e.g. straight lines through the origin), then we will get similar fits for different data sets, because all we have to adjust is one parameter--the slope. The problem here, is that the fit is only good if the underlying process is in fact a straight line through the origin. If it isn't a straight line for instance, there will be a fixed error or **bias** that will never go away, no matter how much data we collect. Statisticians refer to this problem as the *bias/variance* dilemma.

To sum up, lots of parameter flexibility (or lots of hidden units) has the benefit of fitting anything, but at the cost of sensitivity to variability in the data set--there is *variance* introduced by the fits found over multiple training sets (e.g. of a small fixed size).

A fit with very few parameters is not as sensitive to the inevitable variability in the training set, but can give large constant errors or *bias* if the data do not match the underlying model.

There is no general way of getting around this problem, and neural networks are no exception. We generalized linear regression to non-linear fits using error back-propagation. Because back-propagation models can have lots of hidden layers with many units and weights, they form a class of very flexible approximators and can fit almost any function. But these models can show high variability in their fits from one data set to the next, even when the data comes from the same underlying process. Lots of hidden units can mean low bias, but at a high cost in variance.

Summary: Three types of Graphical models and their relationship to neural network problems

Three types of nodes in a graphical model: known, unknown to be estimated, unknown to be integrated out (marginalized)

We have three basic states for nodes in a graphical model: known, unknown to be estimated, unknown to be integrated out (marginalization). We have causal states of the world S , that get mapped to some image data I , through some intermediate parameters M : $S \rightarrow M \rightarrow I$.

Consider faces described in terms of geometrical shape S . M could be a set of parameters that link face identity S to the resulting images I .

Generative model (forward, data synthesis): We want to model I through $p(I|S)$. In our example, we want to specify the shape S of "Bill's face", and then $p(I|S="Bill")$ can be implemented as an algorithm to spit out images of Bill. M gets integrated out.

Hypothesis inference: we want to model samples for S : $p(S|I)$. Given an image, we want a routine to output likely object identities, so that we can minimize risk (see below), or do MAP classification for accurate object identification. M gets integrated out again.

Classification: discrete labels. E.g. face identity ("Bill", "George",...) given an image

Regression: estimating continuous quantities, like the shape of Bill's face.

Learning: we want to model M : $p(M|I,S)$, to learn how the intermediate variables are distributed. Given lots of samples of objects and their images, we want to learn the mapping parameters between them. These parameters could have a simple interpretation (e.g. A face label S determines a shape M which determines an image I).

Alternatively, do a mental switch and consider a neural network in which an input S gets mapped to an output I through intermediate variables M . We can think of M as representing synaptic weights to be learned as in back-prop, but now we want to know the probability distribution of the weights, not just a particular set of values.

Density estimation: Can also be thought of as regression on a probability density functions, with the added condition that the area under the fitted curve must sum to one.

Bayesian Decision Theory

Bayes Decision theory, loss, and risk

Minimizing error isn't always the best thing to do.

The costs of certain kinds of errors (e.g. a high cost to false alarms) could affect the decision criterion. Consider light discrimination. Even though the sensitivity of the observer is essentially unchanged (e.g. the separation between the signal and noise means, and their standard deviations' for the two Gaussian distributions (bright vs. dim light) remains unchanged), increasing the criterion for deciding bright or dim can increase the overall error rate. This isn't necessarily bad. (d' is the signal-to-noise ratio, given by the difference between the two Gaussian means divided by the standard deviation is assumed to be the same.)

A doctor might say that since stress EKG's have about a 30% false alarm rate, it isn't worth doing. The cost of a false alarm is high--at least for the HMO, with the resulting follow-ups, angiograms, etc.. And some increased risk to the patient of extra unnecessary tests. But, of course, false alarm rate isn't the whole story, and one should ask what the hit rate (or alternatively the miss rate) is? Miss rate is about 10%. (Thus, d' is actually pretty high--what is it?). From the patient's point of view, the cost of a miss is very high, one's life. So a patient's goal would not be to minimize errors (i.e. probability of a mis-diagnosis), but rather to minimize a measure of subjective cost that puts a very high cost on a miss, and low cost on a false alarm.

Although decision theory has traditionally been applied to analogous trade-offs that are more cognitive than perceptual, recent work has shown that perception has implicit, unconscious trade-offs in the kinds of errors that are made.



One example is in shape from shading. An image provides the "test measurements" that can be used to estimate an object's shape and/or estimate the direction of illumination. Accurate object identification often depends crucially on an object's shape, and the illumination is a confounding (secondary) variable. This suggests that visual recognition should put a high cost to errors in shape perception, and lower costs on errors in illumination direction estimation. So the process of perceptual inference depends on a task. The effect of marginalization in the fruit example illustrated task-dependence. Now we show how marginalization can be generalized through decision theory to model other kinds of goals than error minimization (MAP) in task-dependence.

Bayes Decision theory provides the means to model perceptual/cognitive performance as a function of utility.

■ Decision theory

Some terminology. We clearly distinguish the decision space from the state or hypothesis space, and introduce the idea of a loss $L(d,s)$, which is the cost for making the decision d , when the actual state is s .

Often we can't directly measure s , and we can only infer it from observations. Thus, given an observation (image measurement) x , we define a risk function that represents the *average loss* over signal states s :

$$R(d; x) = \sum_s L(d, s) p(s | x) \quad (1)$$

This suggests a decision rule: $\alpha(x) = \underset{d}{\operatorname{argmin}} R(d;x)$. But not all x are equally likely. This decision rule minimizes the expected risk average over all observations:

$$R(\alpha) = \sum_x R(\alpha; x) p(x) \quad (2)$$

We won't show them all here, but with suitable choices of likelihood, prior, and loss functions, we can derive standard estimation procedures (maximum likelihood, MAP, estimation of the mean) as special cases of risk minimization.

For the MAP estimator (e.g. energy minimization in the Hopfield net),

$$R(d; x) = \sum_s L(d, s) p(s | x) = \sum_s (1 - \delta_{d,s}) p(s | x) = 1 - p(d | x) \quad (3)$$

where $\delta_{d,s}$ is the discrete analog to the Dirac delta function--it is zero if $d \neq s$, and one if $d = s$. Basically you score big if you get it right, but are penalized uniformly for anything less than perfect.

Thus minimizing risk with the loss function $L = (1 - \delta_{d,s})$ is equivalent to maximizing the posterior, $p(d|x)$. And we saw that for a large class of distributions (exponential), maximizing the posterior is equivalent to minimizing an energy function, as with the Hopfield and Boltzmann networks.

What about marginalization? You can see from the definition of the risk function, that this corresponds to a uniform loss:

$L = -1$.

$$R(s_1; \mathbf{x}) = \sum_{s_2} L(d_2, s_2) p(s_1, s_2 | \mathbf{x}) \quad (4)$$

So for our face recognition example, a really huge error in illumination direction has the same cost as getting it right.

For the fruit example, optimal classification of the fruit identity required marginalizing over fruit color--i.e. effectively treating fruit color identification errors as equally costly

...even tho', doing MAP after marginalization effectively means we are not explicitly identifying color.

Other kinds of loss functions do different kinds of estimation. So for example, if $L(d,s) = -(s - d)^2$, then minimizing risk is equivalent to finding the mean (rather than the mode as in MAP estimation).

For a recent discussion of Bayesian decision theory in the context of object perception, see:

Kersten, D., Mamassian P & Yuille A (2004) Object perception as Bayesian inference. Annual Review of Psychology.

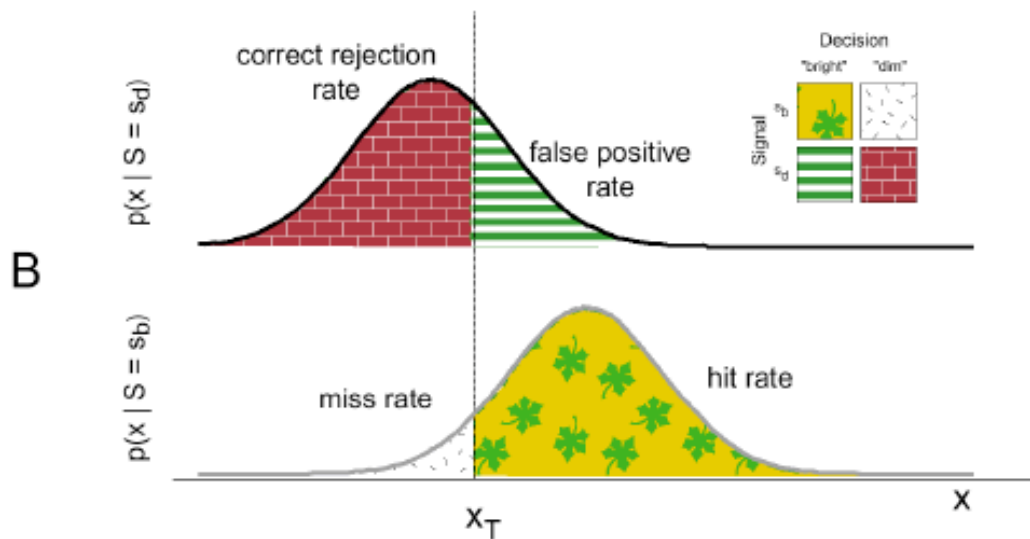
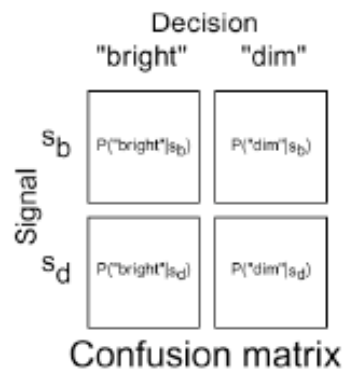
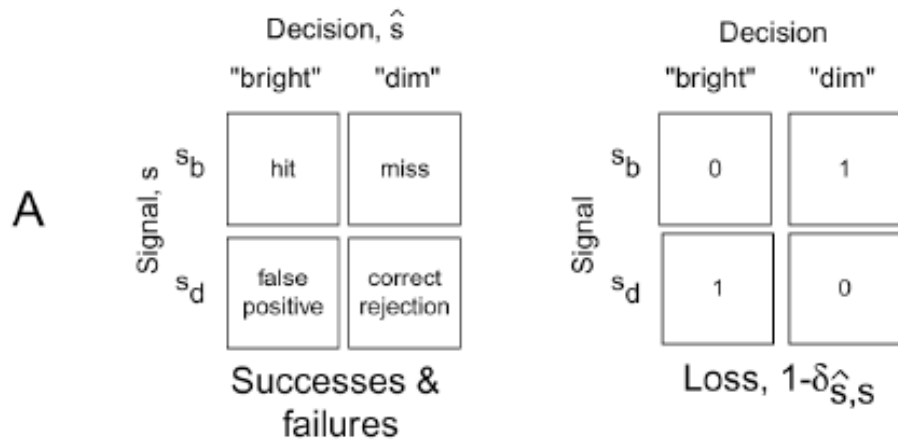
<http://arjournals.annualreviews.org/doi/pdf/10.1146/annurev.psych.55.090902.142005>

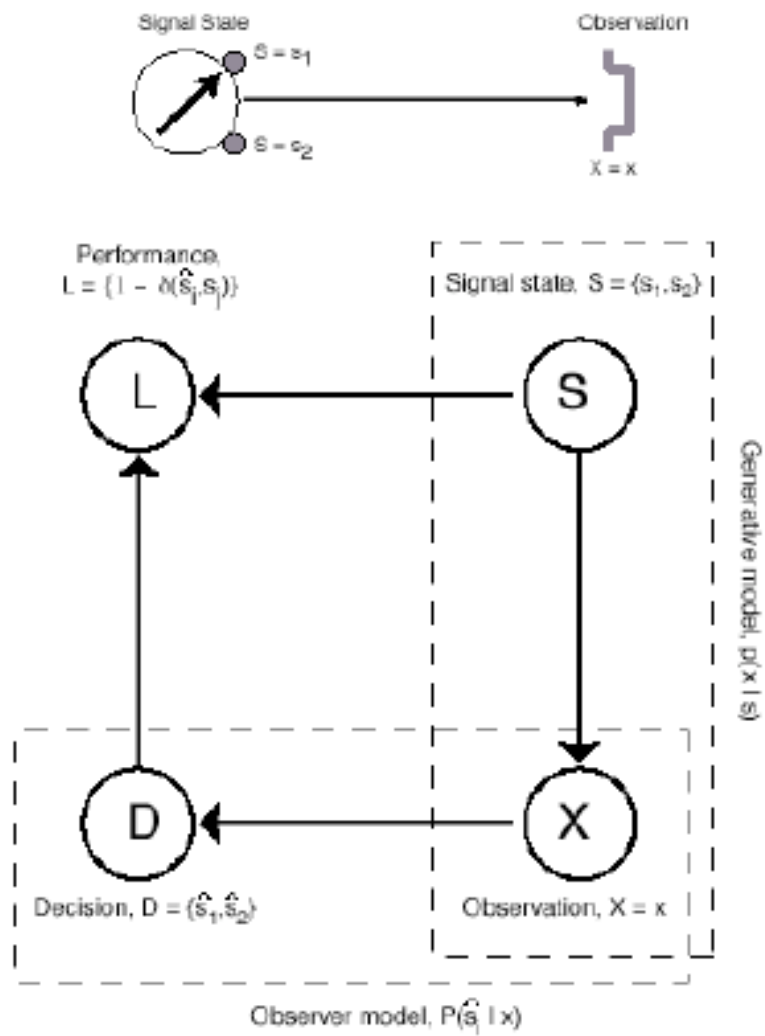
Appendix

Hypothesis inference: Three types

■ Detection

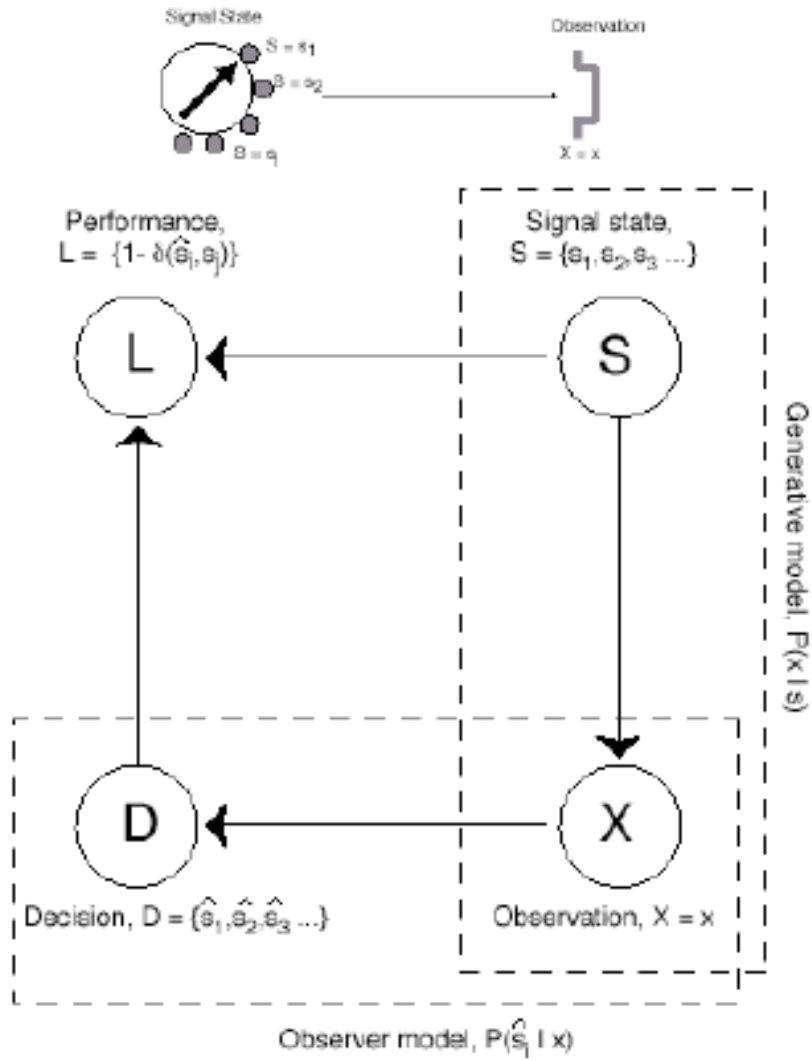
Let the decision variable d , be represented by \hat{s} .





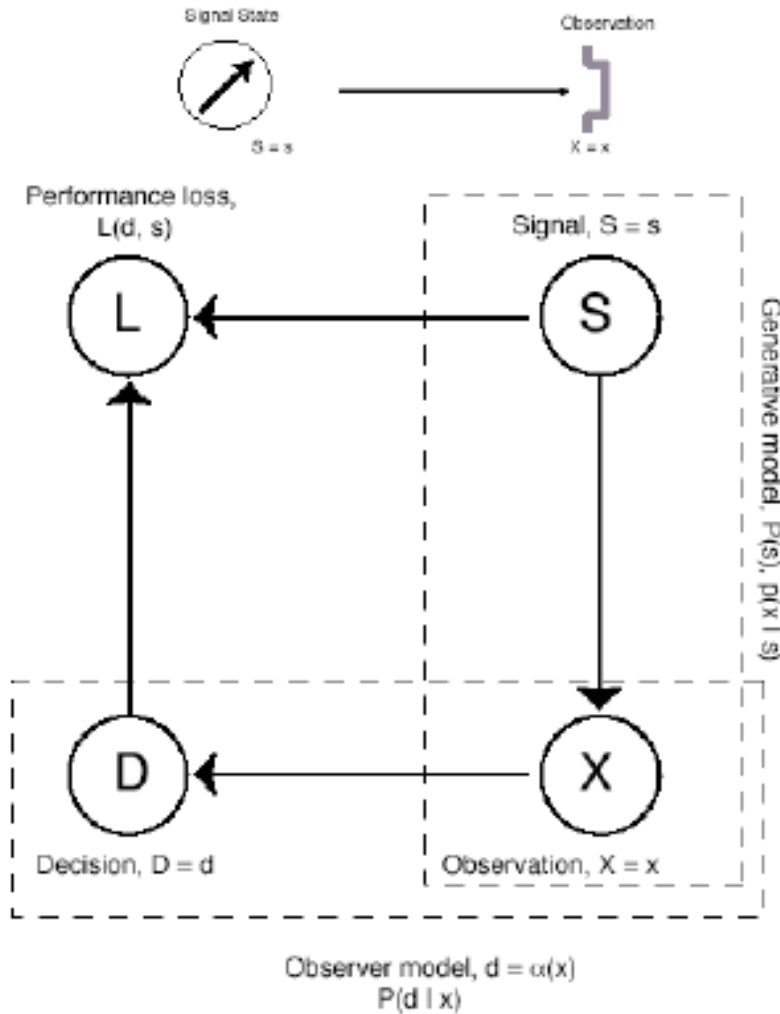
■ Classification

MAP rule: $\underset{i}{\operatorname{argmax}} \{p(S_i | x)\}$.



■ Continuous estimation

$$\operatorname{argmax}_S \{p(S | x)\}$$



As described above, one can show that $L(d,s) = -(d-s)^2$ produces an estimator that finds the mean, $L(d,s) = -\delta(d-s)$, does MAP (i.e. finds the mode), and $L(d,s) = 1$ is equivalent to marginalization (integrating out s).

MAP minimizes probability of error: Proof for detection

Here is why MAP minimizes average error. Suppose that x is fixed at a value for which $P(S = s_b | x) > P(S = s_d | x)$. This is exactly like the problem of guessing “heads” or “tails” for a biased coin, say with a probability of heads $P(S = s_b | x)$. Imagine the light discrimination experiment repeated many times and you have to decide whether the switch was set to bright or not –but only on those trials for which you measured exactly x . The optimal strategy is to always say “bright”. Let’s see why. First note that:

$$p(\text{error} | x) = p(\text{say "bright", actually dim} | x) + p(\text{say "dim", actually bright} | x) = p(\hat{s}_1, s_2 | x) + p(s_1, \hat{s}_2 | x)$$

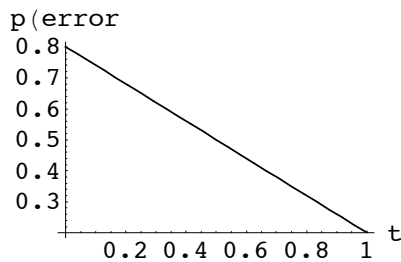
Given x , the response is independent of the actual signal state (see graphical model for detection above--"response is conditionally independent of signal state, given observation x "), so the joint probabilities factor:

$$p(\text{error}|x) = p(\text{say "bright" } |x)p(\text{actually dim } |x) + p(\text{say "dim" } |x)p(\text{actually bright } |x)$$

Let $t = p(\text{say "bright" } |x)$, then

$$p(\text{error}|t,x) = t * p(\text{actually dim } |x) + (1-t) * p(\text{actually bright } |x).$$

$p(\text{error}|t)$, as a function of t , defines a straight line with slope $p(\text{actually dim } |x) - p(\text{actually bright } |x)$. (Just take the partial derivative with respect to t .) We've assumed $P(S = \text{sb} |x) > P(S = \text{sd} |x)$, so $p(\text{error}|t)$ has a negative slope, with the smallest non-negative value of t being one. So, error is minimized when $t = p(\text{say "bright" } |x) = 1$. I.e. Always say "bright".



Always saying "bright" results in a probability of error $P(\text{error} |x) = P(S = \text{sd} |x)$. That's the best that can be done on average. On the other hand, if the observation is in a region for which $P(S = \text{sd} |x) > P(S = \text{sb} |x)$, the minimum error strategy is to always pick "dim" with a resulting $P(\text{error} |x) = P(S = \text{sb} |x)$. Of course, x isn't fixed from trial to trial, so we calculate the total probability of error which is determined by the specific values where signal states and decisions don't agree:

$$\begin{aligned} p(\text{error}) &= \sum_{i \neq j} p(\hat{s}_i, s_j) \\ &= \sum_{i \neq j} \int p(\hat{s}_i, s_j | x) p(x) dx = \sum_{i \neq j} \int p(\hat{s}_i | x) p(s_j | x) p(x) dx \end{aligned}$$

Because the MAP rule ensures that $p(\hat{s}_i, s_j | x)$ is the minimum for each x , the integral over all x minimizes the total probability of error.

Nearest neighbor classifier

<http://cgm.cs.mcgill.ca/~soss/cs644/projects/simard/>

References

Duda, R. O., & Hart, P. E. (1973). *Pattern classification and scene analysis*. New York.: John Wiley & Sons.

Duda, R. O., Hart, P. E., & Stork, D. G. (2001). *Pattern classification* (2nd ed.). New York: Wiley.

(Amazon.com)

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

Kersten, D., & Yuille, A. (2003). Bayesian models of object perception. *Current Opinion in Neurobiology*, 13(2), 1-9.

Kersten, D., Mamassian, P., & Yuille, A. (2004). Object perception as Bayesian Inference. *Annual Review of Psychology*, 55.

MacKay, D. J. C. (1992). Bayesian interpolation. *Neural Computation*, 4(3), 415-447.

Ripley, B. D. (1996). *Pattern Recognition and Neural Networks*. Cambridge, UK: Cambridge University Press.

Vapnik, V. N. (1995). *The nature of statistical learning*. New York: Springer-Verlag.

<http://neuron.eng.wayne.edu/software.html>

© 1998, 2001, 2003, 2005 Daniel Kersten, Computational Vision Lab, Department of Psychology, University of Minnesota.