

# Neural Computing

## Learning Guide, Lecture Summary, and Worked Examples

### Aims

The aims of this course are to investigate how biological nervous systems accomplish the goals of machine intelligence but while using radically different strategies, architectures, and hardware; and to investigate how artificial neural systems can be designed that try to emulate some of those biological principles in the hope of capturing some of their performance.

### Lectures

- **Natural versus artificial substrates of intelligence.** Comparison of the differences between biological and artificial intelligence in terms of architectures, hardware, and strategies. Levels of analysis; mechanism and explanation; philosophical issues. Basic neural network architectures compared with rule-based or symbolic approaches to learning and problem-solving.
- **Neurobiological wetware: architecture and function of the brain.** Human brain architecture. Sensation and perception; learning and memory. What we can learn from neurology of brain trauma; modular organisation and specialisation of function. Aphasias, agnosias, apraxias. How stochastic communications media, unreliable and randomly distributed hardware, slow and asynchronous clocking, and imprecise connectivity blueprints, give us unrivalled performance in real-time tasks involving perception, learning, and motor control.
- **Neural processing and signalling.** Information content of neural signals. Spike generation processes. Neural hardware for both processing and communications. Can the mechanisms for neural processing and signalling be viably separated? Biophysics of nerve cell membranes and differential ionic permeability. Excitable membranes. Logical operators.
- **Stochasticity in neural codes.** Principal Components Analysis of spike trains. Evidence for detailed temporal modulation as a neural coding and communications strategy. Is stochasticity also a fundamental neural computing strategy for searching large solution spaces, entertaining candidate hypotheses about patterns, and memory retrieval? John von Neumann's conjecture. Simulated annealing.
- **Neural operators that encode, analyse, and represent image structure.** How the mammalian visual system, from retina to brain, extracts information from optical images and sequences of them to make sense of the world. Description and modelling of neural operators in engineering terms as filters, coders, compressors, and pattern matchers.
- **Cognition and evolution. Neuropsychology of face recognition.** The sorts of tasks, primarily social, that shaped the evolution of human brains. The computational load of social cognition as the driving factor for the evolution of large brains. How the degrees-of-freedom within faces and between faces are extracted and encoded by specialised areas of the brain concerned with the detection, recognition, and interpretation of faces and facial expressions. Efforts to simulate these faculties in artificial systems.

- **Artificial neural networks for pattern recognition.** A brief history of artificial neural networks and some successful applications. Central concepts of learning from data, and foundations in probability theory. Regression and classification problems viewed as non-linear mappings. Analogy with polynomial curve fitting. General “linear” models. The curse of dimensionality, and the need for adaptive basis functions.
- **Probabilistic inference.** Bayesian and frequentist views of probability and uncertainty. Regression and classification expressed in terms of probability distributions. Density estimation. Likelihood function and maximum likelihood. Neural network output viewed as conditional mean.
- **Network models for classification and decision theory.** Probabilistic formulation of classification problems. Prior and posterior probabilities. Decision theory and minimum misclassification rate. The distinction between inference and decision. Estimation of posterior probabilities compared with the use of discriminant functions. Neural networks as estimators of posterior probabilities.

## Objectives

At the end of the course students should:

- be able to describe key aspects of brain function and neural processing in terms of computation, architecture, and communication.
- be able to analyse the viability of distinctions such as computing vs communicating, signal vs noise, and algorithm vs hardware, when these dichotomies from Computer Science are applied to the brain.
- understand the neurobiological mechanisms of vision well enough to think of ways to implement them in machine vision.
- understand basic principles of the design and function of artificial neural networks that learn from examples and solve problems in classification and pattern recognition.

## Reference books

Aleksander, I. (1989). *Neural Computing Architectures*. North Oxford Academic Press.

Bishop, C.M. (1995). *Neural Networks for Pattern Recognition*. Oxford University Press.

Haykin, S. (1994). *Neural Networks: A Comprehensive Foundation*. Macmillan.

Hecht-Nielsen, R. (1991). *Neurocomputing*. Addison-Wesley.

## Exercise 1

List five critical respects in which the operating principles that are apparent in biological nervous tissue differ from those that apply in current computers, and in each case comment upon how these might explain key differences in performance such as adaptability, speed, fault-tolerance, and ability to solve ill-conditioned problems.

## Model Answer – Exercise 1

(Five items such as the nine on this list):

1. Neural tissue is asynchronous, and there is no master clock; time is an analog (continuous) variable in neural computing. Computers are synchronous and everything happens on the edges of discrete clock ticks.
2. Neural tissue involves random connectivity, with no master blueprint; connectivity in computers is precisely specified.
3. No single element is irreplaceable in neural tissue, and function appears unaffected by the deaths of 10,000s of neurones every day. Not so in computing machines.
4. Neural processing is highly distributed, seeming to show equipotentiality and recruitment of neural machinery as needed. Less true in computers, which have rigid hardware functional specification (e.g. memory vs. ALU).
5. The elementary ‘cycle time’ of neurones (i.e. their time constant) is of the order of one millisecond, whereas silicon gate times are on the order of one nanosecond, i.e., a million times faster.
6. The number of elementary functional units in brains ( $10^{11}$  neurones,  $10^{15}$  synapses) exceeds by a factor of more than 1,000 those in computers.
7. Brains are able to tolerate ambiguity, to learn on the basis of very impoverished and disorganized data (e.g. to learn a grammar from random samples of poorly structured, unsystematic, natural language); whereas much more precisely formatted and structured data and rules are required by computers.
8. Communications media in neural tissue involve stochastic codes; those in computing machines are deterministic and formally structured.
9. Brains do not seem to have formal symbolic or numerical representations, unlike computers (although humans certainly can perform symbolic manipulation, if only very slowly by comparison). In general, it appears that those tasks for which we humans have *cognitive penetrance* (i.e. an understanding of how we do them, like mental arithmetic), we are not very efficient at, in comparison to machines; but those tasks at which we excel, and machines can hardly perform at all (e.g. adaptive behaviour in unpredictable or novel environments, or face recognition, or language acquisition), are tasks for which we have virtually no cognitive penetrance (ability to explain how we do them).

## Exercise 2

1. Illustrate how stochasticity can be used in artificial neural networks to solve, at least in an asymptotic sense, problems that would otherwise be intractable. Name at least two such stochastic engines; describe the role of stochasticity in each; and identify the kinds of problems that such artificial neural devices seem able to solve.
2. Illustrate the evidence for stochasticity in natural nervous systems, and comment on the role that it might play in neurobiological function. What is the case supporting John von Neumann's deathbed prediction that it might be a computational engine for the nervous system, rather than just random noise? Describe at least one experiment involving neural tissue in support of this theory.

## Model Answer – Exercise 2

1. Stochasticity offers an opportunity to explore very large problem spaces in search of a solution or a globally optimal match by *blind variation and selective retention*. Random variation is analogous to “temperature,” or perturbations in state whose average variance specifies a relationship between entropy and energy that is analogous to temperature. The use of random variation ensures that (in a statistical sense) all corners of the state space can be represented and/or explored. This is an approach to solving NP-complete problems that relies upon asymptotic convergence of expected values, rather than deterministic convergence of an algorithm upon the solution. Its prime disadvantages are (i) very slow operation; and (ii) no guarantee of finding the optimal solution.

Two examples of stochastic engines: Simulated Annealing; and Genetic Algorithms. Role of stochasticity in SA: temperature that declines according to a specific annealing schedule, representing random jumps through state space but with declining average amplitude, so that improvements are more likely, but traps are avoided in the long-term. Role of stochasticity in GA’s: mutations of the genotype, with those that increase fitness being retained. Type of problem approached with SA: the Travelling Salesman Problem. Type of problem approached with GA: Monte Carlo combinatorial optimization.

2. Sequences of nerve action potentials are stochastic time-series whose random structure resembles, to first order, a variable-rate Poisson process. The inter-arrival time distributions tend to be exponentials, and the counting distributions tend to be gamma distributions.

von Neumann’s prediction that stochasticity may play an important role in neurobiological function is supported by the fact that seemingly *identical* visual stimuli can generate very different spike sequences from the same neurone in successive presentations, as though possibly different hypotheses were being “entertained” about the pattern and compared with the responses from other neurones. A second argument is that if noise were disadvantageous, then it should quickly have been eliminated in evolution (i.e. Nature could easily have evolved less “noisy” membranes than those with the electrophysiological properties of nerve cells). One set of experiments supporting the hypothesis are those of Optican and Richmond, in which a Principle Components Analysis (PCA) of the response sequences of neurones in the infero-temporal (IT) lobe of macaque monkeys were recorded while they looked at various orthogonal visual patterns (2D Walsh functions). It was found that there are systematic eigenfunctions (shared among large populations of IT neurones) of spike train variation that seem to form temporal-modulation codes for spatial patterns. The conclusion was that much more than just the “mean firing rate” of neurones matters, and that these higher moments of (otherwise seemingly random variation in firing) were in fact responsible for about two-thirds of all the information being encoded.

### Exercise 3

Discuss how neural operators that encode, analyze, and represent image structure in natural visual systems can be implemented in artificial neural networks. Include the following issues:

- receptive field structure
- adaptiveness
- perceptual learning
- hierarchies of tuning variables in successive layers
- the introduction of new signal processing dimensions and of non-linearities in successive layers
- wavelet codes for extracting pattern information in highly compressed form
- self-similarity of weighting functions
- associative memory or content-addressable memory for recognizing patterns such as faces and eliciting appropriate response sequences

### Model Answer – Exercise 3

- Receptive Field Concept: a linear combination of image pixels is taken by some neurone, with weights (either positive or negative) to produce a sum which determines the output response of the neurone. The Receptive Field constitutes that region of visual space in which information can directly influence the neurone. Its distribution of weights determines primarily the functionality of the neurone. These elements of natural nervous systems are the standard elements of Artificial Neural Networks (ANN's).
- Adaptiveness: the summation weights over the receptive field can be adaptive, controlled by higher-order neural processes, which may be hormonal or neuro-peptides in the case of natural nervous systems. In ANN's, the standard model for adaptiveness is the ADELIN (Adaptive Linear Combiner), and involves global feedback control over all gain parameters in the network.
- Perceptual learning involves the modification of synaptic strengths or other gain factors in response to visual experience, such as the learning of a particular face. In natural visual systems, the almost real-time modification of receptive field properties of neurones has been observed, depending upon other stimulation occurring in (possibly remote) parts of visual space.
- Hierarchies of tuning variables: In the retina and the lateral geniculate nucleus, the spatial tuning variables for visual neurones are primarily size and center-surround structure. But in the visual cortex, the new tuning variable of orientation selectivity is introduced. Another one is stereoscopic selectivity (disparity tuning). At still higher levels in the infero-temporal cortex, still more abstract selectivities emerge, such as neurones tuned to detect faces and to be responsive even to particular aspects of facial expression, such as the gaze.
- Beyond the primary visual cortex, all neurones have primarily non-linear response selectivities. But up to the level of “simple cells” in V1, many response properties can be described as linear, or quasi-linear.
- Cortical receptive field structure of simple cells can be described by a family of 2D wavelets, which have five primary degrees-of-freedom: (1) and (2) the X,Y coordinates of the neurone's receptive field in visual space; (3) the size of its receptive field; (4) its orientation of modulation between excitatory and inhibitory regions; and (5) its phase, or symmetry. These wavelet properties generate complete representations for image structure, and moreover do so in highly compressed form because the wavelets serve as decorrelators.
- To a good approximation, the receptive field profiles of different visual neurones in this family are *self-similar* (related to each other by dilation, rotation, and translation).
- Many neurones show the property of associative recall (or content addressability), in the sense that even just very partial information such as a small portion of an occluded face seems able to suffice to generate the full response of that neurone when presented with the entire face. This idea has been exploited in “Hopfield networks” of artificial neurones for fault-tolerant and content-addressable visual processing and recognition.

#### Exercise 4

Explain the concepts of the ‘curse of dimensionality’ and ‘intrinsic dimensionality’ in the context of pattern recognition. Discuss why models based on linear combinations of fixed basis functions of the form

$$y(\mathbf{x}) = \sum_j w_j \phi_j(\mathbf{x})$$

suffer from the curse of dimensionality, and explain how neural networks, which use adaptive basis functions, overcome this problem.

#### Model Answer – Exercise 4

The term *curse of dimensionality* refers to a range of phenomena whereby certain pattern recognition techniques require quantities of training data which increase exponentially with the number of input variables. Regression models consisting of linear combinations of fixed basis functions  $\phi_i(\mathbf{x})$  are prone to this problem. As a specific example, suppose that each of the input variables is divided into a number of intervals, so that the value of a variable can be specified approximately by saying in which interval it lies. By increasing the number of divisions along each axis we could increase the precision with which the input variables can be described. This leads to a division of the whole input space into a large number of cells. Let us now choose the basis function  $\phi_j(\mathbf{x})$  to be zero everywhere except within the  $j$ th cell (over which it is assumed to be constant). Suppose we are given data points in each cell, together with corresponding values for the output variable. If we are given a new point in the input space, we can determine a corresponding value for  $y$  by finding which cell the point falls in, and then returning the average value of  $y$  for all of the training points which lie in that cell. We see that, if each input variable is divided into  $M$  divisions, the total number of cells is  $M^d$  and this grows *exponentially* with the dimensionality of the input space. Since each cell must contain at least one data point, this implies that the quantity of training data needed to specify the mapping also grows exponentially. Although the situation can be improved somewhat by better choices for the basis functions, the underlying difficulties remain as long as the basis functions are chosen independently of the problem being solved. For most real problems, however, the input variables will have significant (often non-linear) correlations, so that the data does not fill the input space uniformly, but rather is confined (approximately) to a lower-dimensional manifold whose dimensionality is called the *intrinsic dimensionality* of the data set. Furthermore, the output value may have significant dependence only on particular directions within this manifold. If the basis functions depend on adjustable parameters, they can adapt to the position and shape of the manifold and to the dependence of the output variable(s) on the inputs. The number of such basis functions required to learn a suitable input-output function will depend primarily on the complexity of the non-linear mapping, and not simply on the dimensionality of the input space.

#### Exercise 5

By using the example of polynomial curve fitting through noisy data, explain the concept of generalization. You should include a discussion of the role of model complexity, an explanation of why there is an optimum level of complexity for a given data set, and a discussion of how you would expect this optimum complexity to depend on the size of the data set.



## Model Answer – Exercise 5

The goal in solving a pattern recognition problem is to achieve accurate predictions for new data. This is known as *generalization*, and it is important to understand that good performance on the training does not necessarily imply good generalization (although poor performance on the training data will almost certainly result in equally poor generalization). A practical method of assessing generalization is to partition the available data into a training set and a separate validation set. The training set is used to optimize the parameters of the model, while the validation set is used to assess generalization performance. An important factor governing generalization is the complexity (or flexibility) of the model. In the case of a polynomial, the complexity is governed by the order of the polynomial as this controls the number of adaptive parameters (corresponding to the coefficients in the polynomial). Polynomials of order  $M$  include polynomials of all orders  $M' < M$  as special cases (obtained by setting the corresponding coefficients to zero), so an increase in the order of the polynomial will never result in an increase in training set error, and will often result in a decrease. Naively we might expect that the same thing would hold true also for the error measured on the validation set. However, in practice this is not the case since we must deal with a data set of finite size in which the data values are noisy. Suppose we use a polynomial to fit a data set in which the output variable has a roughly quadratic dependence on the input variable, and where the data values are corrupted by noise. A linear (first order) polynomial will give a poor fit to the training data, and will also have poor generalization, since it is unable to capture the non-linearity in the input-output mapping. A quadratic (second-order) polynomial will give a smaller error on both training and validation sets. Polynomials of higher order will give still smaller training set error, and will even give zero error on the training set if the number of coefficients equals the number of training data points. However, they do so at the expense of fitting the noise on the data as well as its underlying trend. Since the noise component of the validation data is independent of that on the training data, the result is an increase in the validation set error. Figure 1 shows a schematic illustration of the error of a trained model measured with respect to the training set and also with respect to an independent validation set.

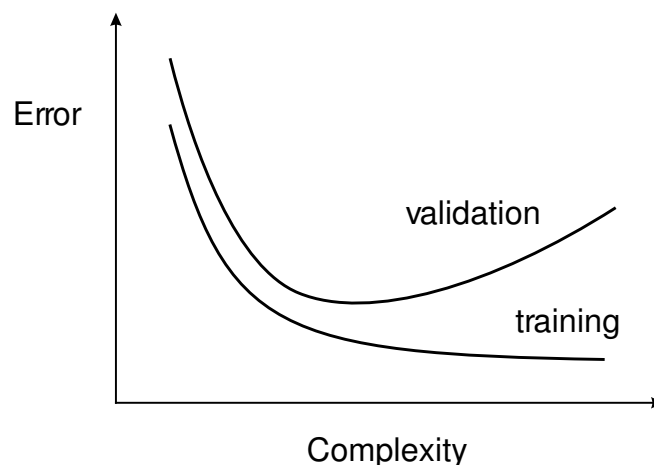


Figure 1: Schematic illustration of the behaviour of training set and validation set error versus model complexity.

## Exercise 6

1.

Many classes of artificial neural networks learn from data by forming a lower dimensional parametric representation, or mapping, that resembles the process of curve-fitting. Explain this idea in reference to least-squares error minimisation or statistical regression.

Explain why increasing the complexity of a model may cause a neural network's error in the training phase to become smaller and smaller, but its generalisation in the validation phase to become worse and worse. How would you expect the optimal complexity of a neural network model to depend on the amount of data?

2.

Answer each of the following short questions:

1. What is the approximate capacitance of nerve cell membrane, in microFarads per  $\text{cm}^2$ , and what functional parameters of neural activity are determined by this?
2. Approximately what range of voltages does a nerve cell membrane move through during the course of generating a neural impulse, and what determines this range?
3. What is the role of positive feedback in nerve impulse generation?
4. From which organ does the retina develop embryologically, and to what cells elsewhere in the body are the retinal photoreceptors most closely related?
5. What causes the refractory deadtime of about 1 msec after each nerve impulse, and what is its consequence?

## Model Answer – Exercise 6

1.

A neural network that learns a model for a process or a data set is estimating a small(ish) set of parameters that describe its main trends and dependencies. This is similar to the process of fitting a curve, such as a polynomial function, to data by the methods of statistical regression analysis or of least-squares (finding that set of parameters that minimizes the sum of the squared deviations between the model and the available data). Although an exact solution can be found for the coefficients of models such as polynomials which produce a system of linear equations in the same number of unknowns, more generally we wish to fit non-linear parameters that do not lead to systems of linear equations. Training such models requires slow iterative adjustment of all parameters while seeking to minimize an error function by gradient descent. This is basically how neural network training procedures work.

The parameters of a neural network model can be set by the minimisation of an appropriate error function. However, the goal of training is not to give good performance on the training data, but instead to give the best performance (in terms of smallest error) on independent, unseen data drawn from the same distribution as the training data. The ability of a neural network to give good results on unseen data is called generalisation. If all of the available data were used in the training set and the model were made too large in order to fit all of it well, then it may become “over-trained” on this one set and will not generalize well to other sample sets. A simpler model trained just on a subset of the available data may actually show better generalisation. More complex models can fit the training data better, but do so at the expense of fitting the noise on the data as well as its underlying trend. Since the noise component of the validation data is independent of that on the training data, the result is an increase in the validation set error. Therefore the optimal model complexity would be somewhere between these two cases, and it should grow only slowly with the amount of data from which the neural net is trying to learn.

2.

1. The capacitance of nerve cell membrane is approximately 10,000 microFarads per  $\text{cm}^2$ . This largely determines the nerve cell membrane time-constant and hence the neural response speed, as well as the velocity of nerve impulse propagation.
2. When generating a neural impulse, the voltage across a nerve cell membrane moves from about -40mV to +70mV (the inside relative to the outside of the neuron). The Nernst equilibrium potential determines the resting potential of -40mV based on the relative concentrations of  $\text{K}^+$ ,  $\text{Na}^+$ , and  $\text{Cl}^-$  ions. The net positive polarisation of the membrane when  $\text{Na}^+$  ions have flowed into the cell down their osmotic gradient, due to the voltage-dependent opening of  $\text{Na}^+$  conductance channels in the membrane, balanced by the outward flow of  $\text{K}^+$  ions, determines the peak spike potential of about +70mV.
3. Nerve impulse generation is based on the existence of voltage-dependent conductances for  $\text{Na}^+$  and  $\text{K}^+$  ions, with differing time constants, and osmotic gradients of different sign. Positive feedback exists because the more the trans-membrane voltage rises, the lower the

resistance becomes (the greater the conductance becomes) for  $\text{Na}^+$  ions. This causes still more  $\text{Na}^+$  ions to flow in, and the voltage to rise still more. The positive feedback process halts only when  $\text{Na}^+$  equilibrium has been reached. (Then the opposite flow of  $\text{K}^+$ , with a slower time constant, brings the voltage down again until voltage-dependent conductance levels are restored, and ion pumps can restore the original ion separations and concentration gradients.)

4. The retina develops embryologically from brain tissue, of which the eye itself is a collapsed ventricle. Retinal photoreceptors are derived from hair cells, as indeed are the transducers along the basilar membrane which are basis for hearing, and similarly the somatosensory receptors.
5. The need to restore ionic equilibrium by the action of ion pumps is the reason for the 1 msec refractory deadtime after each nerve impulse. Collapse of ion-specific resistances during the nerve impulse allowed  $\text{Na}^+$  and  $\text{K}^+$  ions to flow down their osmotic gradients until reaching equilibrium with charge gradients. Restoring these to resting levels limits the highest frequency of nerve impulse generation to a few 100 Hz.

## Exercise 7

**1.**

In a Hopfield neural network configured as an associative memory, with all of its weights trained and fixed, what three possible behaviours may occur over time in configuration space as the net continues to iterate in response to a given input?

How many stable content-addressable memories would you expect a fully connected Hopfield network consisting of 100 neurons to be capable of storing?

What property of those memory patterns would make it most probable that you could successfully train the network to store the maximum number, and why?

**2.**

Explain how 5 independent dimensions of visual processing are multiplexed together into the 3 available spatial dimensions of neural tissue, by the structure of the cubic millimeter hypercolumns in the brain's visual cortex.

**3.** The retina is often regarded as an image capture device; but it has about 100 million input sensors (photoreceptors) yet only 1 million output fibres (optic nerve axons). What are some implications of this 100-to-1 ratio of input channels to output channels?

**4.** Provide some statistics and arguments supporting the proposition that: "Connectivity is the basic computational principle in the brain."

## Model Answer – Exercise 7

1.

The Hopfield network (i) may have reached a stable state or attractor, from which it will make no further changes; (ii) it may be caught in a limit cycle in which it will continue to oscillate indefinitely between states as the iterations continue; or (iii) it may wander around chaotically in state space without reaching either any periodic or stable states.

A fully connected Hopfield network consisting of 100 neurons should be capable of storing about 15 stable content-addressable memories.

If the stored patterns are orthogonal to each other, or as nearly so as possible (i.e. if the inner product projections of the memories onto each other were minimal, or zero), this would maximize the number that could be stored.

2.

Neurons' orientation selectivity, size or spatial frequency selectivity, and alternation between inputs coming from the left and right eyes for stereo disparity, are embedded locally within each cubic millimeter of brain tissue in the visual cortex. (Receptive field size varies from the superficial to the deep layers; orientation selectivity varies tangentially to the cortical surface, spanning  $2\pi$  radians over 1 mm; and ocular dominance slabs of about 0.5mm width alternate between inputs from the left and right eyes in the orthogonal tangential direction.) These "hypercolumn modules" of  $1\text{mm}^3$  neural machinery themselves migrate systematically across the global 2D visual field, in their own receptive field map positions. Hence 3 dimensions of neural selectivity mappings are embedded locally within the global 2D retinotopic mapping for position, spanning 5 dimensions altogether.

3.

The fact that the retina has about 100 times more input channels (photoreceptors) than output channels (fibres in the optic nerve) indicates that far from being just an image transfer device, the retina actively processes, encodes, and summarizes the spatio-temporal and chromatic structure of the dynamical scene. What reaches the brain is already an abstracted description in terms of several dimensions of image processing and analysis. Indeed, it is worth remembering that the eye itself develops embryologically from a collapsed ventricle of the brain; the retina should be regarded as a piece of the brain.

4.

Synapses in the brain far outnumber neurons: on average there are about 10,000 synapses per neuron. In certain sagittal brain slices, one sees white matter almost completely dominating grey matter. (White matter consists of the myelinated axons that connect neurons; grey matter consists of neuron cell bodies themselves.) Learning and memory formation are believed to be based on the formation of new synapses amongst groups of neurons; drugs that interfere with the synthesis of new connections produce amnesia from the time of the drug's application. There are some 30 different known neurotransmitters, and more than a dozen distinct types of synapses. With some  $10^{15}$  synaptic connections in a human brain, many modifiable based on experiences, it seems reasonable to say that connectivity is the basic computational principle in the brain.

## Exercise 8

Consider a network having a vector  $\mathbf{x}$  of inputs and a single output  $y(\mathbf{x})$  in which the output value represents the posterior probability of class membership for a binary classification problem. The error function for such a network can be written

$$E = - \sum_{n=1}^N \{t_n \ln y(\mathbf{x}_n) + (1 - t_n) \ln(1 - y(\mathbf{x}_n))\}$$

where  $t_n \in \{0, 1\}$  is the target value corresponding to input pattern  $\mathbf{x}_n$ , and  $N$  is the total number of patterns. In the limit  $N \rightarrow \infty$ , the average error per data point takes the form

$$E = - \iint \{t \ln y(\mathbf{x}) + (1 - t) \ln(1 - y(\mathbf{x}))\} p(t|\mathbf{x})p(\mathbf{x}) dt d\mathbf{x}. \quad (1)$$

By functional differentiation of (1), show that, for a network model with unlimited flexibility, the minimum of  $E$  occurs when

$$y(\mathbf{x}) = \int tp(t|\mathbf{x}) dt$$

so that the network output represents the conditional average of the target data. Next consider a network with multiple outputs representing posterior probabilities for several mutually exclusive classes, for which the error function is given by

$$E = - \sum_{n=1}^N \sum_{k=1}^K t_{nk} \ln y_k(\mathbf{x}_n)$$

where the target values have a 1-of- $K$  coding scheme so that  $t_{nk} = \delta_{kl}$  for a pattern  $n$  from class  $l$ . (Here  $\delta_{kl} = 1$  if  $k = l$  and  $\delta_{kl} = 0$  otherwise.) Write down the average error per data point in the limit  $N \rightarrow \infty$ , and by functional differentiation show that the minimum occurs when the network outputs are again given by the conditional averages of the corresponding target variables. Hint: remember that the network outputs are constrained to sum to unity, so express the outputs in terms of the softmax activation function

$$y_k = \frac{\exp(a_k)}{\sum_l \exp(a_l)}$$

and perform the differentiation with respect to the  $\{a_k(\mathbf{x})\}$ .

## Model Answer – Exercise 8

If we start from the error function in the form

$$E = - \iint \{t \ln y(\mathbf{x}) + (1 - t) \ln(1 - y(\mathbf{x}))\} p(t|\mathbf{x})p(\mathbf{x}) dt d\mathbf{x}$$

and set the functional derivative with respect to  $y(\mathbf{x})$  equal to zero we obtain

$$\frac{\delta E}{\delta y(\mathbf{x})} = - \int \left\{ \frac{t - y(\mathbf{x})}{y(\mathbf{x})(1 - y(\mathbf{x}))} \right\} p(t|\mathbf{x})p(\mathbf{x}) dt.$$

Provided  $p(\mathbf{x}) \neq 0$  we can solve for  $y(\mathbf{x})$  to give

$$y(\mathbf{x}) = \int t p(t|\mathbf{x}) dt$$

where we have used the normalization property  $\int p(t|\mathbf{x}) dt = 1$  for the conditional distribution. For the case of multiple classes, we can again take the limit  $N \rightarrow \infty$ , so that the error function per data point becomes

$$E = - \iint t \ln\{y(\mathbf{x})\} p(t|\mathbf{x})p(\mathbf{x}) dt d\mathbf{x}.$$

To evaluate the functional derivatives with respect to  $a_k(\mathbf{x})$  we first note that, for the softmax activation function

$$\frac{\partial y_k}{\partial a_m} = \frac{\partial}{\partial a_m} \left\{ \frac{\exp(a_k)}{\sum_l \exp(a_l)} \right\} = y_k \delta_{km} - y_k y_m$$

where  $\delta_{km} = 1$  if  $k = m$  and  $\delta_{km} = 0$  otherwise. Hence we obtain

$$\begin{aligned} \frac{\delta E}{\delta a_m(\mathbf{x})} &= \sum_{k=1}^K \int \frac{\delta E}{\delta y_k(\mathbf{x}')} \frac{\delta y_k(\mathbf{x}')}{\delta a_m(\mathbf{x})} d\mathbf{x}' \\ &= - \int \{t_k - y_k(\mathbf{x})\} p(t_k|\mathbf{x})p(\mathbf{x}) dt_k. \end{aligned}$$

Assuming  $p(\mathbf{x}) \neq 0$  we can solve for  $y_k(\mathbf{x})$  to obtain

$$y_k(\mathbf{x}) = \int t_k p(t_k|\mathbf{x}) dt_k$$

which again is the conditional average of the target data, conditioned on the input vector.



## Exercise 9

Consider a cancer screening application, based on medical images, in which only 1 person in 1000 in the population to be screened has cancer. Suppose that a neural network has been trained on a data set consisting of equal numbers of ‘cancer’ and ‘normal’ images, and that the outputs of the network represent the corresponding posterior probabilities  $P(C|\mathbf{x})$  and  $P(N|\mathbf{x})$  where  $C \equiv$  ‘cancer’ and  $N \equiv$  ‘normal’. Assume that the loss in classifying ‘cancer’ as ‘normal’ is 500 times larger than the loss in classifying ‘normal’ as ‘cancer’ (with no loss for correct decisions). Explain clearly how you would use the outputs of the network to assign a new image to one of the classes so as to minimize the expected (i.e. average) loss. If you were also permitted to reject some fraction of the images, explain what the reject criterion would be in order again to minimize the expected loss.

### Model Answer – Exercise 9

The prior probabilities of cancer and normal are  $P(C) = 10^{-3}$  and  $P(N) = 1 - 10^{-3}$  respectively. From Bayes’ theorem we know that the posterior probabilities  $P(C|\mathbf{x})$  and  $P(N|\mathbf{x})$  are proportional to the artificial prior probabilities of  $\hat{P}(C) = \hat{P}(N) = 0.5$  used to train the network. Hence we can find the posterior probabilities  $\tilde{P}(C|\mathbf{x})$  and  $\tilde{P}(N|\mathbf{x})$  corresponding to the true priors by dividing by the old priors and multiplying by the new ones, so that

$$\tilde{P}(C|\mathbf{x}) \propto P(C|\mathbf{x}) \frac{10^{-3}}{0.5}$$

and similarly for  $\tilde{P}(N|\mathbf{x})$ . Normalizing and cancelling the factors of 0.5 we then obtain

$$\tilde{P}(C|\mathbf{x}) = \frac{P(C|\mathbf{x})}{P(C|\mathbf{x}) + P(N|\mathbf{x})(10^3 - 1)}$$

with an analogous expression for  $\tilde{P}(N|\mathbf{x})$ .

We can now use these corrected posterior probabilities to find the decision rule for minimum expected loss. If an input pattern  $\mathbf{x}$  is assigned to class  $C$  then the expected loss will be

$$\tilde{P}(N|\mathbf{x})$$

while if the pattern is assigned to class  $N$  the expected loss will be

$$500\tilde{P}(C|\mathbf{x}).$$

Thus to make a minimum expected loss decision we simply evaluate both of these expressions for the new value of  $\mathbf{x}$  and assign the corresponding image to the class for which the expected loss is smaller.

Finally, to reject some fraction of the images we choose a threshold value  $\theta$  and reject an image if the value of the expected loss, when the image is assigned to the class having the smaller loss, is greater than  $\theta$ . By changing the value of  $\theta$  we can control the overall fraction of images which will be rejected, so that by increasing  $\theta$  we reject fewer images, while if  $\theta = 0$  all of the images will be rejected.

## Exercise 10

In Computer Science, a fundamental distinction has classically been erected between computing and communications. The former creates, requires, or manipulates data, and the latter moves it around. But in living neural systems, this distinction is less easy to establish; a given neurone performs both functions by generating nerve impulses, and it is not clear where to draw the distinction between processing and communication. Still more so with artificial neural networks, where the entire essence of computing is modeled as just changes in connectivity. Flesh out and discuss this issue. Would you argue that some of the limitations of efforts in artificial intelligence have been the result of such a spurious dichotomy?

## Model Answer – Exercise 10

Computing by logical and arithmetic operations is based upon deterministic rules which are guaranteed (in a proper algorithm) to lead to a solution by a sequence of state transitions. Apart from moving bits to and from registers or storage locations, the pathways of communications and their properties are classically not part of the analysis. In wet neural systems, there are no (or few) known rules which could be described as “formal,” and little or nothing appears to be deterministic. Rather, stochasticity appears to be the best description of membrane properties, signalling events, and overall neural activity. Similarly, the connectivity between and among neurones is not based upon precise blueprints, but rather upon connectivity matrices which are probabilistic both in their wiring and in their connection strengths. It may be that signalling, or communications, among neurones are the essence of wet neural computing, rather than any distinct processing rules which in any way resemble a sequence of instructions. In artificial neural nets, this general view is the basic strategy for learning and problem-solving. Connectivity is everything. Learning occurs by adaptive modification of connection strengths, often following rules which are primarily probabilistic and rarely even remotely formal. The events which underlie neural network computing are analog, or graded, rather than discrete states and transitions among states. Finally, an influential view today is that “physics is computation,” meaning that the laws of nature underlying dynamical systems, energy minimization, and entropy flows over time, may be the only way to implement the sorts of computations that are required which cannot readily be reduced to mere symbol manipulation. If the tasks which require solution in artificial intelligence (e.g. vision, or learning) are formally “intractable,” as is generally accepted, then this observation may well account for the failure of AI largely to deliver on its promises. Implementing the ill-posed problems of AI instead as optimization problems or as stochastic explorations of huge-dimensional solution spaces may be the key strategy behind wet neural computing, and may indeed be the only way forward for AI.

## Exercise 11

Explain the mechanisms and computational significance of nerve impulse generation and transmission. Include the following aspects:

1. Equivalent electrical circuit for nerve cell membrane.
2. How different ion species flow across the membrane, in terms of currents, capacitance, conductances, and voltage-dependence. (Your answer can be qualitative.)
3. Role of positive feedback and voltage-dependent conductances.
4. The respect in which a nerve impulse is a mathematical catastrophe.
5. Approximate time-scale of events, and the speed of nerve impulse propagation.
6. What happens when a propagating nerve impulse reaches an axonal branch.
7. What would happen if two impulses approached each other from opposite directions along a single nerve fibre and collided.
8. How linear operations like integration in space and time can be combined in dendritic trees with logical or Boolean operations such as AND, OR, NOT, and veto.
9. Whether “processing” can be distinguished from “communications,” as it is for artificial computing devices.
10. Respects in which stochasticity in nerve impulse time-series may offer computational opportunities that are absent in synchronous deterministic logic.

## Model Answer – Exercise 11

1. Equivalent Circuit diagram:

(see lecture notes)

2. A nerve cell membrane can be modelled in terms of electrical capacitance  $C$  and several conductances (or resistances  $R$ ) specific to particular ions which carry charge across the membrane. These currents  $I$  into and out of the nerve cell affect its voltage  $V$  in accordance with Ohm's Law for resistance ( $I = V/R$ ) and the law for current flow across a capacitor ( $I = C \frac{dV}{dt}$ ). The charge-carrying ionic species are sodium ( $Na^+$ ), potassium ( $K^+$ ), and chloride ( $Cl^-$ ) ions.
3. The crucial element underlying nerve impulse generation is the fact that the conductances (resistances) for  $Na^+$  and  $K^+$  are not constant, but voltage-dependent. Moreover, these two voltage-dependent conductances have different time courses (time constants). The more the voltage across a nerve cell rises due to  $Na^+$  ions flowing into it, the lower the resistance for  $Na^+$  becomes. ( $Na^+$  current continues to flow until its osmotic concentration gradient is in equilibrium with the voltage gradient.) This positive feedback process causes a voltage spike, which is a nerve impulse.
4. Since the positive feedback process is unstable, causing the voltage to climb higher and higher, faster and faster, it can be described as a catastrophe, rather like an explosion. Once combustion starts in a small corner of a keg of dynamite, matters just get worse and worse. The positive climb of voltage only stops when, on a slower time-scale,  $K^+$  begins to flow in the opposite direction. Once the trans-membrane voltage falls below its threshold, the resting state of ionic concentrations can be restored by ion pumps and the catastrophe (the nerve impulse) is over.
5. The process described above is complete within about 2 milliseconds. There is a refractory period for restoration of ionic equilibrium that lasts for about 1 millisecond, so the fastest frequency of nerve impulses is about 300 Hz. The speed of nerve impulse propagation down an excitable myelinated axon, by saltatory spike propagation, can reach 100 meters per second in warm-blooded vertebrates.
6. A nerve impulse reaching an axonal branch would normally go down both paths, unless vetoed at either one by other shunting synapses. Some axonal branches are "steerable" by remote signalling.
7. The two approaching nerve impulses would annihilate each other when they collided. The minimum refractory period of excitable nerve membrane prevents the impulses from being able to pass through each other, as they would if they were pulses propagating in a linear medium such as air or water or the aether.
8. The linear components of nerve cells (i.e. their capacitance and any non-voltage-dependent resistances) behave as linear integrators, providing linear (but leaky) summation of currents over space and time. However, the fundamentally non-linear interactions at synapses can implement logical operations such as AND, OR, NOT, and veto. The basic factor which

underlies this “logico-linear” combination of signal processing is the mixture of excitable (“logical”) and non-excitable (“linear”) nerve cell membranes.

9. It is very difficult to distinguish between processing and communications in living nervous tissue. The generation and propagation of nerve impulses is the basis for both. A steerable axonal branch particularly illustrates the impossibility of making such a distinction.
10. Stochasticity in nerve impulse time-series may provide a means to search very large spaces for solutions (e.g. to pattern recognition problems) in a way resembling “simulated annealing.” Evolutionary computing (blind variation and selective retention of states) can be the basis of learning in neural networks, and stochasticity may provide the blind variation.

## Exercise 12

Explain why probability theory plays a central role in neural computation. Discuss how the problem of classification can be expressed in terms of the estimation of a probability distribution.

Explain what is meant by a *likelihood function* and by the concept of *maximum likelihood*.

Consider a neural network regression model which takes a vector  $\mathbf{x}$  of input values and produces a single output  $y = f(\mathbf{x}, \mathbf{w})$  where  $\mathbf{w}$  denotes the vector of all adjustable parameters ('weights') in the network. Suppose that the conditional probability distribution of the target variable  $t$ , given an input vector  $\mathbf{x}$ , is a Gaussian distribution of the form

$$p(t|\mathbf{x}, \mathbf{w}) = \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{\{t - f(\mathbf{x}, \mathbf{w})\}^2}{2\sigma^2}\right)$$

where  $\sigma^2$  is the variance parameter.

Given a data set of input vectors  $\{\mathbf{x}_n\}$ , and corresponding target values  $\{t_n\}$ , where  $n = 1, \dots, N$ , write down an expression for the likelihood function, assuming the data points are independent. Hence show that maximization of the likelihood (with respect to  $\mathbf{w}$ ) is equivalent to minimization of a sum-of-squares error function.

## Model Answer – Exercise 12

Neural computation deals with problems involving real-world data and must therefore address the issue of *uncertainty*. The uncertainty arises from a variety of sources including noise on the data, mislabelled data, the natural variability of data sources, and overlapping class distributions. Probability theory provides a consistent framework for the quantification of uncertainty, and is unique under a rather general set of axioms.

The goal in classification is to predict the class  $\mathcal{C}_k$  of an input pattern, having observed a vector  $\mathbf{x}$  of features extracted from that pattern. This can be achieved by estimating the conditional probabilities of each class given the input vector, i.e.  $P(\mathcal{C}_k|\mathbf{x})$ . The optimal decision rule, in the sense of minimising the average number of mis-classifications, is obtained by assigning each new  $\mathbf{x}$  to the class having the largest posterior probability.

The likelihood function, for a particular probabilistic model and a particular observed data set, is defined as the probability of the data set given the model, viewed as a function of the adjustable parameters of the model. Maximum likelihood estimates the parameters to be those values for which the likelihood function is maximized. It therefore gives the parameter values for which the observed data set is the most probable.

Since the data points are assumed to be independent, the likelihood function is given by the product of the densities evaluated for each data point

$$\begin{aligned}\mathcal{L}(\mathbf{w}) &= \prod_{n=1}^N p(t_n|\mathbf{x}_n, \mathbf{w}) \\ &= \prod_{n=1}^N \frac{1}{(2\pi\sigma^2)^{1/2}} \exp\left(-\frac{\{t_n - f(\mathbf{x}_n, \mathbf{w})\}^2}{2\sigma^2}\right).\end{aligned}$$

Following the standard convention, we can define an error function by the negative logarithm of the likelihood

$$E(\mathbf{w}) = -\ln \mathcal{L}(\mathbf{w}).$$

Since the negative logarithm is a monotonically decreasing function, maximization of  $\mathcal{L}(\mathbf{w})$  is equivalent to minimization of  $E(\mathbf{w})$ . Hence we obtain

$$E(\mathbf{w}) = \frac{1}{2\sigma^2} \sum_{n=1}^N \{t_n - f(\mathbf{x}_n, \mathbf{w})\}^2 + \frac{N}{2} \ln(2\pi\sigma^2)$$

which, up to an additive constant independent of  $\mathbf{w}$  and a multiplicative constant also independent of  $\mathbf{w}$ , is the sum-of-squares error function.

### **Exercise 13**

Give brief explanations of the following terms:

- (a) the curse of dimensionality,
- (b) the Perceptron,
- (c) error back-propagation,
- (d) generalisation,
- (e) loss matrix.



### Model Answer – Exercise 13

(a) *The curse of dimensionality.*

Many simple models used for pattern recognition have the unfortunate property that the number of adaptive parameters in the model grows rapidly, sometimes exponentially, with the number of input variables (i.e. with the dimensionality of the input space). Since the size of the data set must grow with the number of parameters, this leads to the requirement for excessively large data sets, as well as increasing the demands on computational resources. An important class of such models is based on linear combinations of fixed, non-linear basis functions. The worst aspects of the curse of dimensionality in such models can be alleviated, at the expense of greater computational complexity, by allowing the basis functions themselves to be adaptive.

(b) *The Perceptron.*

The Perceptron is a simple neural network model developed in the 1960s by Rosenblatt. He built hardware implementations of the Perceptron, and also proved that the learning algorithm is guaranteed to find an exact solution in a finite number of steps, provided that such a solution exists. The limitations of the Perceptron were studied mathematically by Minsky and Papert.

(c) *Error back-propagation.*

Neural networks consisting of more than one layer of adaptive connections can be trained by error function minimisation using gradient-based optimisation techniques. In order to apply such techniques, it is necessary to evaluate the gradient of the error function with respect to the adaptive parameters in the network. This can be achieved using the chain rule of calculus which leads to the error back-propagation algorithm. The name arises from the graphical interpretation of the algorithm in terms of a backwards propagation of error signals through the network.

(d) *Generalisation.*

The parameters of a neural network model can be determined through minimisation of an appropriate error function. However, the goal of training is not to give good performance on the training data, but instead to give the best performance (in terms of smallest error) on independent, unseen data drawn from the same distribution as the training data. The capability of the model to give good results on unseen data is termed generalisation.

(e) *Loss matrix.*

In many classification problems, different misclassification errors can have different consequences and hence should be assigned different penalties. For example, in a medical screening application the cost of predicting that a patient is normal when in fact they have cancer is much more serious than predicting they have cancer when in fact they are healthy. This effect can be quantified

using a loss matrix consisting of penalty values for each possible combination of true class versus predicted class. The elements on the leading diagonal correspond to correct decisions and are usually chosen to be zero. A neural network model can be used to estimate the posterior probabilities of class membership for a given input vector. Simple decision theory then shows that if these posterior probabilities are weighted by the appropriate elements of the loss matrix, then selection of the largest weighted probability represents the optimal classification strategy in the sense of minimising the average loss.

## Exercise 14

**1.**

A competitive Kohonen neural network forms feature maps which can be regarded as performing dimensionality reduction. Explain this.

Is training time normally faster, or slower, in a supervised neural network compared with an unsupervised one? What is the major disadvantage inherent in the use of supervised neural networks?

What class of neural network can be used to overcome the mathematical difficulties caused by the use of non-orthogonal sensory and motor representations?

**2.** Give three examples of biological sensory or motor control systems that seem to rely on the use of non-orthogonal coordinates.

Explain why this creates a problem in the computational evaluation and simulation of such systems, and discuss whether or not you think this issue matters in the function of the actual neurobiological systems.

**3.**

Give four examples of neural activity having a fundamentally quantal structure, in the sense that signals or events are quantised into discrete packages rather than being continuous.

For purposes of understanding neurobiological computation, what can be learned from studying the brain's failures, either as the consequences of specific forms of trauma or in normal function as revealed in the systematic visual illusions?

## Model Answer – Exercise 14.

1.

A Kohonen topological feature map finds the organisation of relationships among patterns. Incoming patterns are classified by the units that they activate in the competitive layer, and similarities among patterns are mapped into closeness relationships in this competitive layer. In a Kohonen feature map, a chain of competitive units can span a pattern space of two or more dimensions by competing with each other to represent specific local neighbourhoods in the ( $n$ -dimensional) input space (i.e. particular combinations of input appearing at the  $n$  input entry points). Because of the neighbourhood that each competitive unit in the chain has won the right to represent, the chain becomes a space-filling curve that can cover, e.g., 2 or 3 or more dimensions. In this sense the feature map achieves dimensionality reduction.

The training (or learning) time is significantly less in supervised neural nets, because with the desired answers known and with the relationship between parameters and the error function known, direct feedback can be provided to the parameters or connection strengths within the network so as to optimise its ability to represent the data or the desired pattern classes. Learning in unsupervised networks occurs much more slowly because their variation is blind and the origin of their errors may not be clear. However, the major disadvantage of supervised neural networks is that (besides requiring a supervisor) they are less free in exploring solution spaces and are less able to discover and extract hidden structure in the training data.

A relaxation network can overcome the mathematical problems associated with non-orthogonal sensory and motor representations.

2.

Non-orthogonal representations underlie the following biological sensory or motor control systems (any three from this list):

- The semi-circular canals of the vestibular system, which have accelerometers measuring angular acceleration to provide a sense of balance, are inclined in planes at about  $105^\circ$
- The absorption spectra of visual pigments are non-orthogonal to each other
- Spatial visual receptive field profiles are non-orthogonal (their inner product projections onto each other are non-zero)
- Attachments of skeletal musculature involve non-orthogonal control coordinates
- The muscles that control eye movements are attached on non-orthogonal axes

Non-orthogonal coordinate systems are difficult to work with mathematically and computationally, because effects are not independent of each other. Necessary control signals become much more complicated, and the meaning of sensory coding variables becomes more difficult to

interpret than if the coding primitives were orthogonal to each other. However, none of this seems to matter in neurobiological systems, perhaps because they use maps instead of numerical representations.

### 3.

Examples of quantal neural structure (any four from this list):

- Nerve impulse generation (all-or-none spikes)
- Synaptic vesicle release (sacculs containing 500,000 molecules of neurotransmitter, all released as one packet)
- Post-synaptic “bumps” that can be observed in the trans-membrane voltage due to the arrival of individual molecules of neurotransmitter
- Trans-membrane voltage quantal fluctuations from gating of individual ion channels
- The reliable sensitivity of dark-adapted retinal photoreceptors to individual photons of light: a person can actually “see” single photons, or dim flashes containing fewer than about a dozen absorbed photons across the retina
- “*Eigengrau*” percepts due to the thermal isomerisation of individual visual pigment molecules, nonetheless perceivable as distinct flashes
- The quantised degrees-of-freedom in the receptive field profiles of neurons and therefore in the dimensions of visual codes

The study of neurological trauma to the brain gives clues to its modular organisation, and specialisation of function. In particular, its fault tolerance when the damage occurs gradually (e.g. a slowly growing tumour as opposed to a sudden injury) suggests that other parts of the brain can be “recruited” to perform the tasks of the damaged parts. It also reveals the degree to which specific functions (like linguistic abilities) are associated with specific brain regions. Finally, systematic visual illusions (such as the universal geometrical distortions) may reveal how our visual algorithms actually work, e.g. by mechanisms of short-range competition, long-range cooperation, and adaptive resonance.

## Exercise 15

(A)

Explain the key ideas of a Hopfield artificial neural network for content-addressable, associative memory. In explaining how memories are stored and retrieved, be sure to define the notions of:

- configuration space
- connectivity matrix
- stable attractor
- basin of attraction
- network capacity, and its dependence on the number of “neurones”

(B)

1. Marshall as many lines of evidence as you can to support the view that in human vision “what you see is your own ‘graphics,’ rather than the retinal image as faithfully recorded by photoreceptors in the eye.” Explain the significance of this observation for vision theory and for machine vision.
  
2. Suppose you were trying to design a machine vision system based as closely as possible upon human vision. Would you aim to design in the visual illusions that nearly all people “see” as well? (These include the distortions of geometrical form, angle and relative length illusions, etc.) If such properties emerged as unintended consequences of your vision design, would you consider them to be features, or bugs?

## Model Answer – Exercise 15

(A)

A Hopfield content-addressable, associative memory allows patterns to be stored in such a way that each one can be retrieved by using just a portion, or a corrupted version, of it as its “address.” The memory functions dynamically by detecting an association between an input pattern and a stored memory, which it then retrieves in its entirety.

The network consists of  $N$  “neurons” that are fully connected to each other. Their connections have weights of adjustable strengths, defining an  $N \times N$  connectivity matrix. This matrix of weights is adjusted during a learning phase that stores the patterns in memory. Each of the  $N$  neurons has a binary state, and the  $N$ -dimensional space of their states is called the configuration space. A stored “memory” is just a particular point in this  $N$ -dimensional configuration space. The training phase consists in setting all the  $N \times N$  weights such that a particular memory is a stable state consistent with all the other stable states, so that once the network reaches such a state, no neurons will change their state further. Then, whenever a new input pattern first appears (implemented by re-setting the states of the  $N$  neurons), the collective state of the network will evolve by each neurone following a threshold rule: The  $N - 1$  inputs that it receives from all the other neurons (each of whose state is either a  $+1$  or a  $-1$ ), each multiplied by the corresponding weight for that connection, are all added together and compared against a threshold. If this inner product (sum) exceeds the threshold, then the neurone’s state is set to  $+1$ ; else it is set to  $-1$ . This process reiterates for all the  $N$  neurons until a stable state is reached, when no neurone changes its state anymore. This collective state is called an stable attractor. Each stable attractor represents one memory. The range of different input states that will all converge eventually to a particular such stable attractor is called the basin of attraction for that memory. Because the states within this basin of attraction represents similar, but corrupted (or partial) versions of the stable state, a Hopfield memory is capable of “completing” a memory from just parts of it serving as a trigger, rather like happens in human memory. Hence it is deemed content-addressable and associative. It is also capable of overcoming noise and corruption of an input pattern, provided that the input remains within the basin of attraction of the attractor. The network capacity of a Hopfield network is about  $(0.15)N$ , where  $N$  is the number of neurons. This capacity would be greater if only orthogonal patterns were stored.

....Continued....

(B)

1. The purpose of vision is not to reproduce faithfully the 2D retinal image, but rather to construct an internal 3D model of the surrounding world and of the 3D objects that populate it. In this regard, vision is a kind of “inverse graphics.” (In graphics, a 3D world model is projected into a 2D screen image; in vision, just the reverse must be accomplished.) Evidence that what we see is our own “graphics:”

- Perceptual size invariance: all the faces in a lecture theatre appear to be roughly the same size. Yet in terms of the 2D retinal image, some faces may actually be 20 or 30 times larger than others, depending on their distance.
- Motion compensation: every eye movement causes the retinal image to shift ballistically, in tremendous sweeps and jerks. Yet the world appears to be stable.
- Homogeneity of spatial resolution: high visual resolution exists only within the fovea (the central 1 or 2 degrees). Yet the world appears to have uniform resolution everywhere.
- Homogeneity of colour signals: the colour-sensitive photoreceptors (“cones”) exist only near the fovea; outside the fovea, only black and white information is transduced. Yet the world appears uniformly full of colour.
- Invisibility of the silhouette of the retinal blood vessels: the retinal image is corrupted by a dense tree of opaque blood vessels, lying in front of the actual photoreceptors. Yet we discount their shadows, and “fill-in” the gaps with graphics. Similarly, the blind spot is filled in; we are not aware of this large hole.
- Colour constancy: a banana continues to look yellow, whether reddish light or bluish light is shown on it. We see the spectral reflectance properties of the pigmented surface, rather than the actual wavelength mix received on the retina.
- Interpreted scene properties: we perceptually interpret events like occlusion, or motion and rotation in depth, in terms of solid 3D objects and configurations.

One implication of these observations for machine vision is that we should cast the problems of vision as “inverse problems,” for which the goal is to invert the processes of graphics (projections, ray-tracing, radiosity, etc) to try to arrive at the unique 3D world situation that could produce such a 2D signal as the one received on the retina.

2. The geometrical illusions (distortions of form, angle, length, shape) may be inevitable consequences of trying to achieve such goals. They may reveal, for example, the existence of near-range competition and far-range cooperation in the processing of orientation information. Mere fidelity to the retinal image should *not* be a primary goal of vision, as it might be in a sound reproduction system. The goal is to *understand* the world; who cares if this *changes* its appearance (inverting the famous maxim of Karl Marx about changing the world versus understanding it...)



## Exercise 16

(A)

1. Define “generalisation” in neural networks that learn from training data, and then are tested on new data. Why should not *all* the available data be used in the training set?
2. Draw a simple connectivity diagram that illustrates the idea of lateral inhibition in a competitive neural network.
3. With another diagram showing plots for input and output, illustrate how lateral inhibition in such a competitive network sharpens any input signal by effectively amplifying its first derivative.
4. What class of multi-layer neural network can be used to overcome the mathematical difficulties caused by intrinsic non-orthogonality of representation in many sensory and control systems?

(B)

The study of neurological trauma to the brain gives clues about its modular organisation and specialisations of function, which may reveal some computational principles.

1. What two fundamental principles of brain function did Karl Lashley’s neurological research seem to reveal?
2. What are generally the differences between recovery prospects after a sudden brain trauma (in which all the damage is done at once), versus the same damage done more gradually (e.g. by a growing tumour)?
3. What mechanism might explain this difference?
4. Comment on its possible computational significance in terms of fault-tolerance, circuit adaptability and flexibility.
5. Describe two different types of language-related disorders that can result from trauma to Broca’s area or Wernicke’s area, and comment on the computational inferences we might draw concerning language processing and linguistic representation.

## Model Answer – Exercise 16

(A)

1. The parameters of a neural network model can be set via the minimisation of an appropriate error function. However, the goal of training is not to give good performance on the training data, but instead to give the best performance (in terms of smallest error) on independent, unseen data drawn from the same distribution as the training data. The ability of a neural network to give good results on unseen data is termed generalisation. If all of the available data were used in the training set and the model were made too large in order to fit all of it well, then it may become “over-trained” on this one set and will not generalize well to other sample sets. A simpler model trained just on a subset of the available data may actually show better generalisation.
2. (Diagram page 124)
3. (Diagram page 125)
4. A relaxation network.

(B)

1. The two principles of brain function that emerged from Karl Lashley’s neurological research were:
  - The Law of Mass Action – the idea that much of the brain seems to participate in the performance of many tasks.
  - The Law of Equipotentiality – the idea that all brain tissue is basically the same, and has the *potential* to subserve any type of brain task if it is recruited to do so.
2. The brain is far more able to overcome, and to compensate for, an insult that develops or spreads slowly, than one caused by sudden trauma.
3. It is thought that the reason is recruitment: other non-injured tissue that may currently not be utilised is recruited to take over the functions served by the damaged tissue, provided there is time for this process to occur.
4. The implication is that brain circuits can train each other. Brain tissue that is involved in the performance of a specific task can transfer its acquired architecture or configuration to other tissue, during a gradual insult (like a spreading neoplasm). This is a remarkable form of adaptability and circuit flexibility, which constitutes a type of fault-tolerance that is unknown in computer hardware.

5. Any two examples from the following list:

**Broca's aphasia:** disorders of speech production. Patients may be unable to read aloud, but can read to themselves and comprehend written material perfectly.

**Wernicke's aphasia:** disorders of speech understanding. Such aphasics may be able to read sentences aloud but unable to comprehend the material.

**echolalia:** a strong tendency to repeat exactly what is heard.

**conduction phasia:** Fluent spontaneous speech, normal comprehension, but inability to repeat what is heard.

**alexia:** Lost ability to comprehend written language.

**paralexia:** Words are read as other words (substitution).

**dyslexia:** Confounded letter combinations when reading.

Computational significance: such disorders imply a remarkable specialisation of function for the different aspects of language. There is evidence that even the different grammatical parts of speech (such as verbs, nouns, prepositions) are stored in different specific places, because highly localised brain traumas can cause the loss of a particular such grammatical category.

## Exercise 17

(A)

Give evidence supporting the view that the main computational load that has shaped the evolution of the human brain is “social computation,” with sexual success being the ultimate measure of the value of an algorithm or neural design feature. What implications does this have for:

- The cognitive skills and perceptual faculties that have been selected for in brain evolution, as contrasted with the goals which are the traditional focus of AI.
- The design of face recognition algorithms, which aim to interpret facial expression, gesture, and intent, as well as gender and identity.
- The construction of the theory that other persons have minds, too.
- Models of action, planning, and interaction between self and others.

Comment on whether this “social computation” view of human brain evolution implies that brain science is less relevant to the goals of computer science than is usually thought.

(B)

Answer any 5 of the following 7 short questions:

1. Roughly what is the total number of neurones in the human brain?
2. Roughly what is the total number of synapses in the human brain? How does this compare with the total number of stars in our galaxy, and with the total number of galaxies in the known universe?
3. Why is nerve impulse propagation described as “saltatory,” and what purposes are achieved by this method of signalling?
4. What is the approximate speed of nerve impulse propagation in warm-blooded animals, in meters/sec?
5. Why is “white matter” white, what cells are responsible for this, and what purpose do they serve?
6. Name the three principal ions involved in nerve membrane current flows, and identify which two of them transit through voltage-dependent conductances.
7. What causes the refractory deadtime of about 1 msec after each nerve impulse?

## Model Answer – Exercise 17

(A)

The neural faculties that have been selected for in human brain evolution are those which optimised sexual success and thereby gene transmission; that is how evolution works. These skills focus on: being charming and seductive; outwitting one's rivals; deceiving others; being able to figure out what others are thinking; out-guessing their intentions and contradicting their plans; manipulating others' desires; forming and maintaining alliances; in short, the whole set of social skills cultivated and achieved by ambitious and upwardly mobile humans. (The time scale of evolution is such that no significant brain evolution can have occurred in the short recent period since the dawn of civilisation.)

Superior execution of the above social tasks leads to higher status in a power hierarchy in social primates, and therefore to better mating privileges. When the ratio of brain size to body size is compared among various primate species, it is most highly correlated with the characteristic size of social groups for that species. (For purposes of this issue, it does not matter which factor is regarded as 'cause' and which as 'effect.')

From smallest to largest, some typical examples of primate social group sizes are: orangutans (3); baboons (20); and chimpanzees (100). It appears that the need to compute and monitor social factors, keeping track of one's con-specifics and of one's own place among them, monitoring others' gaze and actions and modelling their intentions, were important factors determining relative brain size and brain evolution.

Other specific forms of evidence include: the existence of brain areas specialized for processing facial information, including expression, gesture, and intent, as well as facial identity; neurons discovered that are sensitive to the gaze direction of other primates, and especially to eye contact; and evoked brain potentials (electroencephalograms) showing signals sensitive to social power hierarchy (relative position of self and other) in both human and non-human primates.

If such forms of "social computation" are indeed the major computational load of the human brain,

- The traditional goals of AI have not been formulated with a corresponding focus.
- Face processing algorithms, if aimed at simulating human performance with faces, should focus much more on reading social intent and expression than mere identity.
- A crucial step was the evolution of the perspective that other people have minds, too. This (and indeed consciousness itself) would not have been necessary in evolution were it not for the demands of social computation.
- The notion of agency (that others are agents, with their own plans and intentions), and that two agents might need to interact, either in conflict or in cooperation, is perhaps the computational origin of human social psychology.

These perspectives highlight a divergence between the computational goals that have shaped the brain, and those of machine computation. Thus they may limit the current mutual relevance between these two subjects of study. If we wish to "make machines more like humans, rather than humans more like machines," then some consideration and reformulation of goals along these lines may be necessary.

(B)

1. There are roughly  $10^{11}$  neurones in the human brain.
2. There are roughly  $10^{15}$  synapses in the human brain, i.e. about 10,000 times more than the total number of stars in our galaxy, and about 10,000 times more than the total number of galaxies in the known universe.
3. Nerve impulse propagation is described as “saltatory” because when propagating down a single axon, it “jumps” between Nodes of Ranvier which are the isolated patches of excitable membrane separated by about 0.1 mm. The axonal membrane between the Nodes of Ranvier is insulated, and contains no excitable patches. The purposes served by this arrangement are (i) greater speed of impulse propagation, and (ii) less energy consumption by the process, since the dissipative trans-membrane current flows are limited to discrete patches.
4. The approximate speed of nerve impulse propagation is 100 meters/sec in warm-blooded animals.
5. “White matter” consists of the myelinated axons which are the communications channels between neurones. It is white because of myelin insulation, which prevents current leakage and short-circuits. Each axonal fibre is wrapped around many times by insulating Schwann cells.
6. The three principal ions involved in nerve membrane current flows are sodium ( $Na^+$ ), potassium ( $K^+$ ), and chloride ( $Cl^-$ ) ions. Of these three,  $Na^+$  and  $K^+$  transit through voltage-dependent conductances.
7. The time required to restore ionic equilibrium to the resting state concentrations of these three ion species is about 1 msec. Therefore most neurones cannot generate nerve impulses at frequencies exceeding 1 KHz.