

SENSORIMOTOR LEARNING

**CoSMo 2013
Kingston, Ontario**

DOUGLAS TWEED

1 CONTROL

Here we will study how neural networks in the brain learn to control the body. Our focus is on the brain, but we will discuss the challenges facing *any* system that pursues its aims by interacting with the world. Therefore many of the concepts will apply also to other goal-directed systems, such as robots, immune responses, and gene networks.

State Dynamics

Any control problem starts with an agent of some sort embedded in a world which obeys rules of its own that the agent has to work with and can't change. For instance, the neural networks that steer your arm can't change the geometry of the arm or the laws of physics.

These unchangeable rules are summed up in a *state equation*, which describes how a vector of interest called the *state*, \mathbf{x} , evolves under the influence of a *command* vector \mathbf{u} issued by the agent:

$$\dot{\mathbf{x}}(t) = \dot{\mathbf{x}}(\mathbf{x}(t), \mathbf{u}(t)) \quad (1)$$

That is, $\dot{\mathbf{x}}(t)$, the rate of change of \mathbf{x} at time t , depends on the current values of \mathbf{x} and \mathbf{u} . In arm control, for example, \mathbf{x} might be the positions and velocities of all the arm joints, and \mathbf{u} might be the firing rates of the motor-neurons to the arm muscles. So the agent's job is to issue commands, \mathbf{u} , that achieve its aims, given the rules laid out by the state equation.

More generally, $\dot{\mathbf{x}}$ may also depend on other things besides \mathbf{x} and \mathbf{u} . For instance your arm dynamics change if you pick up a heavy suitcase. We will ignore that complication, but it is easy to incorporate using the principles we will learn.

We will use the word *controller*, rather than *agent*, for the thing that issues the commands \mathbf{u} , to emphasize that for us it is usually a set of neurons rather than a whole organism. The controller receives information about \mathbf{x} , and about some target state \mathbf{x}^* , and uses it to compute \mathbf{u} . For instance \mathbf{x}^* might be a location you want to reach to. The function the controller uses to compute \mathbf{u} 's from \mathbf{x} 's and \mathbf{x}^* 's is called its *policy*. We can expand formula (1) to spell out that \mathbf{u} depends on \mathbf{x} and \mathbf{x}^* :

$$\dot{\mathbf{x}}(t) = \dot{\mathbf{x}}(\mathbf{x}(t), \mathbf{u}(\mathbf{x}(t), \mathbf{x}^*(t))) \quad (2)$$

In these notes I will use bold lower-case letters to represent vectors, and I will add parentheses to refer to the functions that give rise to those vectors, e.g. $\mathbf{u}()$ is a policy, and \mathbf{u} and $\mathbf{u}(\mathbf{x}, \mathbf{x}^*)$ are commands. For brevity I will usually write \mathbf{u} rather than $\mathbf{u}(\mathbf{x}(t), \mathbf{x}^*(t))$, but it is understood that \mathbf{u} depends on \mathbf{x} and \mathbf{x}^* , which vary with time.

Biological Dynamics

Engineers study special cases of state dynamics that are easy to control. One is the *linear* form, $\dot{\mathbf{x}} = \mathbf{A}\mathbf{x} + \mathbf{B}\mathbf{u}$, but linear systems are rare in biology. A special case that does apply to biology is the *u-affine* form,

$$\dot{\mathbf{x}} = \mathbf{f}(\mathbf{x}) + \mathbf{G}(\mathbf{x})\mathbf{u} \quad (3)$$

This equation is widely applicable because $f()$ and $G()$ can be any smooth functions.

Biological dynamics often involve a *configuration*, meaning a vector q with fewer elements than the state x and with the property that q and its time-derivatives determine x . For example in reaching, q is all the angles of the arm joints, and x is these angles and their derivatives, $x = (q^T, \dot{q}^T)^T$. If we can rewrite the state dynamics this way

$$\ddot{q} = f(x) + G(x)u \quad (4)$$

then we say they are in *control-canonical form*.

Simulating Dynamics

Most differential equations, including most state equations, can be solved only by numerical integration. The simplest method is Euler's, where we solve $\dot{x} = f(x)$ by breaking time into discrete instants Δt apart, and setting

$$x(t + \Delta t) = x(t) + f(x(t))\Delta t \quad (5)$$

We will use this method in all our control simulations, with $\Delta t = 0.01$ s.

```
>> run EULER_INTEGRATION.m
```

Desired Dynamics

A controller needs a policy that will achieve its aims, but what are those aims? Often they are to make the state x obey some desired dynamics, e.g. to snap to a target, or track a moving object.

Desired dynamics should be feasible. They shouldn't ask the controller to whip the hand to a target supersonically, but should specify an achievable motion. Computationally, this kind of path planning isn't trivial. In the brain, suitable paths are likely wired in by natural selection, and then improved by learning. We will consider that kind of learning in Section 5.

For now we will consider an engineering approach to path planning, which is to choose *linear* desired dynamics, e.g. in the arm example we might choose

$$\ddot{\mathbf{q}}^* = c_1 \dot{\mathbf{q}} + c_0 (\mathbf{q} - \mathbf{q}^*) \quad (6)$$

where the c 's are constants. How do we choose c 's that give the behavior we want, e.g. driving \mathbf{q} to \mathbf{q}^* ? We define the *Hurwitz polynomial*

$$s^2 - c_1 s^1 - c_0 s^0, \text{ i.e. } s^2 - c_1 s - c_0 \quad (7)$$

and choose our c 's so the real parts of all the roots of this polynomial are less than 0. Then the dynamics in (6) are said to be *Hurwitz*, and can be proven to drive \mathbf{q} to \mathbf{q}^* .

Having chosen our desired dynamics, how do we find a policy that will make them come true, i.e. make $\ddot{\mathbf{q}} = \ddot{\mathbf{q}}^*$? If the state dynamics have control-canonical form (4), and the matrix $\mathbf{G}(\mathbf{x})$ has a right-inverse, $\mathbf{G}(\mathbf{x})^+$, then a suitable policy is

$$\mathbf{u}(\mathbf{x}) = \mathbf{G}(\mathbf{x})^+ [\ddot{\mathbf{q}}^*(\mathbf{x}) - \mathbf{f}(\mathbf{x})] \quad (8)$$

To verify that this policy works, substitute the right-hand side for \mathbf{u} in (4) and simplify to get $\ddot{\mathbf{q}} = \ddot{\mathbf{q}}^*$. This

policy, with a Hurwitz $\ddot{\mathbf{q}}^*$, is called *feedback linearization*.

```
>> run FEEDBACK_LINEARIZATION.m
```

Policy (8) assumes the controller knows the state dynamics, as defined by $f()$ and $\mathbf{G}()$. Control networks in the brain probably do have approximations to their $f()$ and $\mathbf{G}()$ wired in at birth, thanks to natural selection. But as the body grows, its dynamics change. So neural controllers must *learn* $f()$ and $\mathbf{G}()$.

2 ELEMENTS OF LEARNING

Neural controllers must learn their state dynamics and policies. To explore the principles in a simple setting, we will consider a task where a learner deduces a *target function* $y()$ based on examples of x 's and $y(x)$'s. This task is called *supervised learning*; the x 's are *inputs*; the y 's are *desired outputs*, and each ordered pair (x, y) is an *example*. Usually the learner can't deduce $y()$ precisely but seeks an estimate $\hat{y}()$ that *approximates* $y()$.

Models

Any learner has a *model*, which is the set of all the functions $\hat{y}()$ it is willing to consider as approximations to $y()$. For instance it might have a *linear* model, which permits only linear functions, of the form

$$\hat{y} = \mathbf{W}\mathbf{x} \tag{9}$$

Here \mathbf{W} is a matrix of numbers called *parameters* or *weights*, which the controller adjusts to improve its estimates. So learning comes down to finding the \mathbf{W} that makes $\hat{y}()$ as close as possible to $y()$.

More generally the model may have the form

$$\hat{y} = \hat{y}(\mathbf{x}, \mathbf{W}) \tag{10}$$

where $\hat{y}()$ may be nonlinear in \mathbf{x} or \mathbf{W} or both. In the brain, the parameters are (probably mainly) the weights (i.e. the strengths) of the synapses.

Risk

The aim of learning is to compute the \mathbf{W} that gives us the best approximation. But what do we mean by *best*? One approach that seems well suited to the brain is to say the best approximation is the one that minimizes the *least-squares risk*, which is the expected value of the squared error in our approximation,

$$R(\mathbf{W}) = \mathbb{E}(\|\hat{\mathbf{y}}(\mathbf{x}, \mathbf{W}) - \mathbf{y}\|^2) \quad (11)$$

i.e. the average squared error over *all possible* examples (\mathbf{x}, \mathbf{y}) .

Optimization

Learning is a matter of finding the weights \mathbf{W} that optimize (i.e. minimize) the risk, R . How might the brain do this? Here we will look at the main methods of optimization from math and computer science. Then in Sections 3 and 4 we will apply them to neural networks. For simplicity we will focus for now on the case where \mathbf{W} is a row vector, \mathbf{w} .

Optimization methods are classed according to the derivatives of R they use: none at all, first derivatives only, or first and second derivatives.

Derivative-free Optimization

If we can compute $R()$ but not its derivatives then the best we can do is try different \mathbf{w} 's and remember which worked best. We make an initial guess \mathbf{w} , compute $R(\mathbf{w})$, then try another weight \mathbf{w}' close to \mathbf{w} . If $R(\mathbf{w}') < R(\mathbf{w})$ then \mathbf{w}' becomes our new \mathbf{w} ; otherwise we reject \mathbf{w}' . Then we repeat, choosing a new \mathbf{w}' near our cur-

rent \mathbf{w} and so on. As we keep doing this, $R(\mathbf{w})$ gradually decreases.

Gradient Descent

If we know the first derivative of R with respect to \mathbf{w} — written $\nabla_{\mathbf{w}}R$ or simply $R_{\mathbf{w}}$ or ∇R and called the *gradient* — then we can move toward a minimum by *gradient descent*. That is, we choose some \mathbf{w} as a starting point, then step to

$$\mathbf{w}_{\text{new}} = \mathbf{w} - \eta \nabla R(\mathbf{w})^T \quad (12)$$

(where η is some positive number), and repeat. $\nabla R(\mathbf{w})^T$ points in the direction of steepest increase of R at the point \mathbf{w} , and so by stepping in the opposite direction, we *decrease* our risk.

Unfortunately this method doesn't tell us how *far* to step, i.e. how big to make η . If we choose η too small then we may approach the minimum extremely slowly. If we make η too large then we may overstep the minimum and increase our risk. Sometimes we can choose good step sizes using the second derivative of R , written $\nabla^2 R$ and called the *Hessian*.

Newton's Method

If we know the Hessian then we can often find minima quickly by *Newton's method*. In this approach we exploit the fact that any smooth function, such as R , in any small region near a minimum, is shaped roughly like the bottom of a paraboloid, i.e. like a quadratic surface $Q(\mathbf{w}) = a + \mathbf{w}\mathbf{b} + \mathbf{w}\mathbf{C}\mathbf{w}^T$, where a is a scalar, \mathbf{b} a column vector, and \mathbf{C} a symmetric, positive-semidefinite matrix.

Clearly $\nabla Q(\mathbf{w}) = \mathbf{b} + 2\mathbf{C}\mathbf{w}^\top$, $\nabla^2 Q = 2\mathbf{C}$, and therefore

$$\nabla Q(\mathbf{w}) = \mathbf{b} + \nabla^2 Q \mathbf{w}^\top \quad (13)$$

We are at the minimum if ∇Q is zero. It follows that, given any \mathbf{w} , the minimum is at $\mathbf{w} + \Delta\mathbf{w}$ if

$$[\nabla^2 Q] \Delta\mathbf{w}^\top = -\nabla Q(\mathbf{w}) \quad (14)$$

because then $\nabla Q(\mathbf{w} + \Delta\mathbf{w}) = \mathbf{b} + \nabla^2 Q(\mathbf{w} + \Delta\mathbf{w})^\top$ [by (13)] = $\mathbf{b} + \nabla^2 Q \mathbf{w}^\top - \nabla Q(\mathbf{w})$ [by (14)] = $\mathbf{b} + \nabla^2 Q \mathbf{w}^\top - \mathbf{b} - \nabla^2 Q \mathbf{w}^\top$ [by (13)] = $\mathbf{0}$.

This formula motivates Newton's method of minimizing R , which is to choose a \mathbf{w} , step to $\mathbf{w}_{\text{new}} = \mathbf{w} + \Delta\mathbf{w}$, where $\Delta\mathbf{w}$ is a solution to

$$[\nabla^2 R] \Delta\mathbf{w}^\top = -\nabla R(\mathbf{w}) \quad (15)$$

and repeat.

If R is a perfect, concave-up paraboloid then a $\Delta\mathbf{w}$ fitting (15) will take us to the minimum in a single step. If R is *roughly* parabolic near \mathbf{w} then we will need a series of steps, though fewer than with gradient descent. But if R isn't even roughly parabolic near \mathbf{w} — if \mathbf{w} is far from a minimum — then Newton's method may fling us into the wild blue yonder.

Another issue is that we may need a lot of effort and memory to compute the Hessian, e.g. if \mathbf{w} has 1000 elements then $\nabla^2 R$ has a million.

Next we will see how these 3 classes of optimization methods may be implemented in the brain.

3 ONLINE LEARNING

The brain learns *online*, meaning it receives a stream of examples (\mathbf{x}_s , y_s), one after another, and adjusts its weights after each one. We will assume the examples are drawn from a probability distribution (unknown to the learning algorithm). So \mathbf{x} and y are random variables. For convenience we will assume they have means of zero.

We will start with methods that approximate the relation between \mathbf{x} and y with linear functions, $\hat{y} = \mathbf{W}\mathbf{x}$. These methods are worth studying because, as we will see, they generalize to learn not just linear but any smooth function.

Again we will focus on the case where \mathbf{W} is a row vector, \mathbf{w} . Our goal is to find the \mathbf{w} that minimizes the risk,

$$R = E(\|\mathbf{w}\mathbf{x} - y\|^2) \quad (16)$$

For now we will assume our neurons are linear: each receives a column vector of inputs \mathbf{x} and multiplies it by a row vector of weights, \mathbf{w} . For instance if the neuron has 1000 incoming synapses then \mathbf{x} has 1000 elements. Each of the 1000 incoming signals x_i gets multiplied by the weight of its synapse, w_i , and the sum of the 1000 products is the cell's firing rate, or signal, $\hat{y} = \mathbf{w}\mathbf{x}$. We call it \hat{y} because the cell is trying to estimate the desired output y .

Weight Perturbation

This is a derivative-free method. A neuron starts with a random weight vector \mathbf{w}_0 , and perturbs it, trying out

nearby \mathbf{w} 's to see if they work better. As each example input \mathbf{x}_s arrives, our neuron computes its estimate $\hat{y}_s = \mathbf{w}_{s-1} \mathbf{x}_s$ and its error $\tilde{y}_s = \hat{y}_s - y_s$. Then it chooses a small random perturbation \mathbf{w}_{pert} , computes a new estimate $\hat{y}'_s = (\mathbf{w}_{s-1} + \mathbf{w}_{pert}) \mathbf{x}_s$ and error \tilde{y}'_s . If the new error is smaller than the first then the neuron accepts the perturbed weights, i.e. $\mathbf{w}_s = \mathbf{w}_{s-1} + \mathbf{w}_{pert}$. Otherwise it adjusts in the opposite direction, $\mathbf{w}_s = \mathbf{w}_{s-1} - \mathbf{w}_{pert}$.

`>> run WEIGHT_PERTURBATION.m`

Node Perturbation

This is a faster derivative-free method, where the neuron perturbs its signal (its firing rate) rather than its synapses. When input \mathbf{x}_s arrives, the neuron fires at rate $\hat{y}_s = \mathbf{w}_{s-1} \mathbf{x}_s$, and computes its error. Then it alters its firing rate slightly, to $\hat{y}'_s = \hat{y}_s + \hat{y}_{pert}$, and computes its new error. Based on those 2 errors, it computes an appropriate adjustment to \mathbf{w} . That computation makes node perturbation more complicated than weight perturbation, but faster.

`>> run NODE_PERTURBATION.m`

Risk Derivatives

To apply gradient descent or Newton's method we need the derivatives of the risk R . Clearly $\nabla R = \nabla E((\mathbf{w}\mathbf{x} - y)^T(\mathbf{w}\mathbf{x} - y)) = E(2\mathbf{x}(\mathbf{w}\mathbf{x} - y))$ [because E is linear and so we can switch E and ∇ and apply the chain rule] = $E(2\mathbf{x}(\mathbf{x}^T\mathbf{w}^T - y)) = 2E(\mathbf{x}\mathbf{x}^T\mathbf{w}^T - \mathbf{x}\mathbf{y})$, i.e.

$$\nabla R = 2[E(\mathbf{x}\mathbf{x}^T)\mathbf{w}^T - E(\mathbf{x}\mathbf{y})] \quad (17)$$

And by similar reasoning,

$$\nabla^2 R = 2E(\mathbf{xx}^\top) \quad (18)$$

This matrix $E(\mathbf{xx}^\top)$ is the *covariance* of \mathbf{x} .

The expected values in (17) and (18) are usually unknown, but we can estimate them based on the current example, \mathbf{x}_s and y_s . To begin with, we can estimate the covariance and cross-covariance this way:

$$\hat{E}(\mathbf{xx}^\top) = \mathbf{x}_s \mathbf{x}_s^\top, \quad \hat{E}(\mathbf{xy}) = \mathbf{x}_s y_s \quad (19)$$

And therefore

$$\hat{\nabla}R_s = 2\left(\mathbf{x}_s \mathbf{x}_s^\top \mathbf{w}^\top - \mathbf{x}_s y_s\right) \quad (20)$$

$$\hat{\nabla}^2 R_s = 2\mathbf{x}_s \mathbf{x}_s^\top \quad (21)$$

Online Gradient Descent

To find the formula for gradient descent in this online linear setting, we plug (20) into (12) to get

$$\mathbf{w}_s = \mathbf{w}_{s-1} - \eta [\mathbf{w}_{s-1} \mathbf{x}_s - y_s] \mathbf{x}_s^\top \quad (22)$$

In the case where the output is a vector, \mathbf{y} , rather than a scalar, the formula is the same except that \mathbf{w} becomes an n_y -by- n_x matrix \mathbf{W} . This formula is also called *least-mean square* learning, or *LMS*.

$$\mathbf{W}_s = \mathbf{W}_{s-1} - \eta [\mathbf{W}_{s-1} \mathbf{x}_s - \mathbf{y}_s] \mathbf{x}_s^\top \quad (23)$$

For a more-compact expression, we define an error vector

$$\tilde{\mathbf{y}}_s = \hat{\mathbf{y}}_s - \mathbf{y}_s = \mathbf{W}_{s-1} \mathbf{x}_s - \mathbf{y}_s \quad (24)$$

and write LMS as

$$\Delta \mathbf{W} = -\eta \tilde{\mathbf{y}} \mathbf{x}^T \quad (25)$$

LMS will drive $\tilde{\mathbf{y}}$ to $\mathbf{0}$, so long as η is not too large. But if η is too small then the convergence is unbearably slow. So η must be chosen carefully.

`>> run LMS.m`

Online Newton

To apply Newton's method, from (15), we have to find a $\Delta \mathbf{w}$ such that $[\hat{\nabla}^2 R_s] \Delta \mathbf{w}^T = -\hat{\nabla} R_s(\mathbf{w})$, i.e. $[\mathbf{x}_s \mathbf{x}_s^T] \Delta \mathbf{w}^T = \mathbf{x}_s \mathbf{x}_s^T \mathbf{w}^T - \mathbf{x}_s \mathbf{y}_s$. A little algebra will confirm that $\Delta \mathbf{w} = -[\mathbf{w} \mathbf{x}_s - \mathbf{y}_s] [\mathbf{x}_s^T \mathbf{x}_s]^{-1} \mathbf{x}_s^T$ does the job. Hence we have this formula for online Newton, also known as *normalized least-mean-square* learning, *NLMS*:

$$\Delta \mathbf{W} = -\tilde{\mathbf{y}} \mathbf{x}^T / (\mathbf{x}^T \mathbf{x}) \quad (26)$$

(when $\mathbf{x} = \mathbf{0}$ we set $\Delta \mathbf{W} = \mathbf{0}$). NLMS learns faster than LMS, and doesn't require us to choose an η .

`>> run NLMS.m`

A drawback of NLMS is that it doesn't converge as well as LMS to a steady \mathbf{W} in cases where only rough fits are possible. If there is a single \mathbf{W}^* that yields a close fit to all the data then NLMS will find it, but if there is no such all-purpose \mathbf{W}^* then NLMS will keep adjusting \mathbf{W} forever, always minimizing the error for just the current input, never finding a \mathbf{W} that is a good global compro-

mise. It helps to introduce a damping factor η between 0 and 1 into the learning rule, like this:

$$\Delta \mathbf{W} = -\eta \tilde{\mathbf{y}} \mathbf{x}^+ \quad (27)$$

to make it less aggressive in adjusting \mathbf{W} to suit the current input.

Recursive Least Squares

NLMS learns faster than LMS but not optimally fast. For optimal speed, we need a more complex algorithm called *recursive least squares*, RLS, also known as an online *Gaussian process*. We won't derive it here, but you can see its equations in the code file RLS.m.

`>> run RLS.m`

4 LEARNING NONLINEARITIES

So far, all our methods have learned *linear* approximations, $\hat{y} = \mathbf{W}\mathbf{x}$, but the brain needs algorithms that work with nonlinear models

$$\hat{y} = \hat{y}(\mathbf{x}, \mathbf{W}) \quad (28)$$

Backpropagation

The backprop algorithm trains layered networks of nonlinear neurons. Networks of this form are capable of *universal approximation*, i.e. you can approximate any smooth function as closely as you like if you have a big enough network. Backprop adjusts the network weights by gradient descent.

>> run BACKPROPAGATION.m

Linear-in-the-parameters Learning

To learn nonlinear functions $y(\mathbf{x})$, we needn't make \hat{y} nonlinear in both \mathbf{x} and \mathbf{W} , as in (28). We can make \hat{y} nonlinear in \mathbf{x} but *linear* in \mathbf{W} ,

$$\hat{y} = \mathbf{W}\phi(\mathbf{x}) \quad (29)$$

This way, we can learn nonlinear functions of \mathbf{x} but use linear learning algorithms. This approach is called *linear-in-the-parameters* learning.

The vector $\phi(\mathbf{x})$ in (29) is called the *feature vector*. Its elements $\phi_i(\mathbf{x})$, called *features*, are functions of \mathbf{x} , usually nonlinear. Common choices are Gaussians, polyno-

mials, and hyperbolic tangents. Like layered networks, models of the form (29) are universal approximators.

By analogy with (24), our error is

$$\tilde{\mathbf{y}} = \mathbf{W}\varphi(\mathbf{x}) - \mathbf{y} \quad (30)$$

For linear-in-the-parameters learning we simply run any of the algorithms from Section 3 using $\varphi(\mathbf{x})$ in place of \mathbf{x} .

For example, suppose we want to learn the function $\mathbf{x}^T \mathbf{x}$ over the domain $x_1, x_2 \in \{-1, 1\}$, as in Backpropagation.m, but now using RLS with Gaussian features, i.e. features of the form $\varphi_i(\mathbf{x}) = \exp(-\gamma(\mathbf{x} - \mathbf{x}_i)^2)$. The vector \mathbf{x}_i is called the *center* of the Gaussian. We have to guess where to put these centers, and what value of γ will work best, and how many features we need. The resulting algorithm is much faster than backprop.

```
>> run RLS_GAUSSIAN.m
```

How might linear-in-the-parameters learning work in the brain? It could run on a layered network with linear neurons in the output layer. The upstream neurons are nonlinear, and the signals coming from the second-last layer of cells serve as features $\varphi_i(\mathbf{x})$. Only the synapses from this second-last layer to the output layer are adjusted by learning, because only they are linearly related to the network's output. All upstream synapses are frozen. They may be set by natural selection to yield good features $\varphi_i(\mathbf{x})$, but they don't change with learning, just as the widths and centers of the Gaussians don't change in RLS_Gaussian.m. So linear-in-the-parameters networks are full of synapses that may be far from optimal but can't be improved.

The Curse

With linear-in-the-parameters methods we can't improve our features during learning, so we have to choose good ones at the outset. Can we just choose a huge number of them, and hope there will be some good ones in the mix? In fact that *is* part of the method. But a brute-force approach — just choosing a lot of features randomly — has its limits, thanks to a theorem called the *curse of dimensionality*, which says the number of random features we need grows exponentially with the dimensionality of the input \mathbf{x} .

Linear versus Nonlinear

Linear-in-the-parameters methods learn faster than nonlinear ones and they may be more robust because they hinge on fewer choices. With RLS, for instance, we choose only the number and type of features, e.g. Gaussians of a certain width. But with backprop we have to choose the type of neuron, the number of layers, the number of neurons in each layer, the η 's for all the neurons, and the initial values of the weights. If we choose wrong then learning fails.

On the other hand, linear-in-the-parameters networks have a lot of non-optimal synapses they can't improve whereas nonlinear methods adjust all the synapses. Nonlinear methods still suffer the curse of dimensionality — the number of adjustable parameters they need grows exponentially with the dimensions of \mathbf{x} — but they make better use of the parameters they have.

In light of all this, the brain may mix linear and nonlinear methods — fast linear learning for output-layer synapses, and slower backprop or perturbation upstream.

5 OPTIMAL CONTROL

In optimal control we try to find desired dynamics and policies that minimize certain costs. Typically we have some non-negative scalar quantity we want to keep small, such as our rate of fuel consumption or the distance to a target, and we call that quantity the *cost rate*, r . This r may depend on \mathbf{x} , \mathbf{u} , and other things such as the location of a target, but for simplicity we will assume the target is always **0** and r is determined by \mathbf{x} or \mathbf{u} or both — $r(\mathbf{x}, \mathbf{u})$.

The thing we want to minimize is the *cost*, J , which is the time-integral of r ,

$$J = \int_t^\infty r(\mathbf{x}(\tau), \mathbf{u}(\mathbf{x}(\tau))) d\tau \quad (31)$$

Cost-to-go

The cost (31) varies as a function of the initial state $\mathbf{x}(t)$. When we want to emphasize its dependence on \mathbf{x} , we call it the *cost-to-go*, and write

$$J^n(\mathbf{x}(t)) = \int_t^\infty r(\mathbf{x}(\tau), \mathbf{u}^n(\mathbf{x}(\tau))) d\tau \quad (32)$$

In other words $J^n(\mathbf{x}(t))$ is the total cost we will accumulate from now on if we are currently (i.e. at time t) in state \mathbf{x} and we choose our commands using the policy $\mathbf{u}^n()$. The superscript n 's emphasize that the cost-to-go depends on the policy. Any one system has many pos-

sible policies, and (32) shows the cost-to-go for one of them, $\mathbf{u}^n()$.

Our aim is to find the *optimal policy* \mathbf{u}^* (\cdot), which is the one with the smallest cost-to-go $J(\mathbf{x})$ for all \mathbf{x} in state space. This smallest J is the *optimal cost-to-go*

$$J^*(\mathbf{x}(t)) = \int_t^\infty r(\mathbf{x}(\tau), \mathbf{u}^*(\mathbf{x}(\tau))) d\tau \quad (33)$$

Finding \mathbf{u}^* (\cdot) is challenging because the cost is an integral through time, and it is hard to work out how a command at any one moment will affect the cost later on. But we can simplify things with 2 concepts: the principle of optimality and Hamiltonians.

Principle of Optimality

This says that every piece of an optimal trajectory is itself optimal. For any stretch of time, t to $t + \Delta t$, and any state $\mathbf{x}(t)$, the optimal policy \mathbf{u}^* (\cdot) always yields a smaller value than does any other policy \mathbf{u} (\cdot) for this integral:

$$\begin{aligned} & \int_t^{t+\Delta t} \nabla J^*(\mathbf{x}(\tau)) \dot{\mathbf{x}}(\mathbf{x}(\tau), \mathbf{u}(\mathbf{x}(\tau))) + r(\mathbf{x}(\tau), \mathbf{u}(\mathbf{x}(\tau))) d\tau \\ &= \int_t^{t+\Delta t} \nabla J^* \dot{\mathbf{x}} + r d\tau \\ &= \int_t^{t+\Delta t} J^* d\tau + \int_t^{t+\Delta t} r d\tau \end{aligned} \quad (34)$$

To see why, note that in the last line of (34), the first integral is the decrease in J^* achieved by the policy \mathbf{u} (\cdot) between times t and $t + \Delta t$, and the second integral is the price paid — the cost that \mathbf{u} (\cdot) accumulated to get this reduction in J^* . Clearly, if \mathbf{u} (\cdot) achieved a smaller sum than \mathbf{u}^* (\cdot), then we could use \mathbf{u} (\cdot) from time t to Δt , and \mathbf{u}^* (\cdot) thereafter, and so accumulate less cost from

time t to ∞ than we would with $\mathbf{u}^*(t)$ alone. But this is a contradiction because $\mathbf{u}^*(t)$ is optimal. Hence $\mathbf{u}^*(t)$ minimizes (34).

Hamiltonians

As we have seen, $\mathbf{u}^*(t)$ minimizes the integral (34) over any time interval Δt , no matter how brief. That means it must minimize the integrand at each moment, i.e. the optimal command \mathbf{u}^* is the one that minimizes

$$\nabla J^*(\mathbf{x}) \dot{\mathbf{x}}(\mathbf{x}, \mathbf{u}) + r(\mathbf{x}, \mathbf{u}) \quad (35)$$

The quantity (35) is called the *Hamiltonian* of the system, $H^*(\mathbf{x}, \mathbf{u})$.

Moreover, differentiating (33) yields $\dot{J}^* = -r$ (where the minus sign appears because we differentiate with respect to the left rather than the right side of the integration domain $[t, \infty)$). Adding r to both sides and expanding terms, we get

$$\nabla J^*(\mathbf{x}) \dot{\mathbf{x}}(\mathbf{x}, \mathbf{u}^*(\mathbf{x})) + r(\mathbf{x}, \mathbf{u}^*(\mathbf{x})) = 0 \quad (36)$$

i.e. $H^*(\mathbf{x}, \mathbf{u}^*) = 0$. From (35) and (36) we see that \mathbf{u}^* both minimizes and zeroes the Hamiltonian. So for all \mathbf{x} ,

$$\min_{\mathbf{u} \in U} \{ \nabla J^*(\mathbf{x}) \dot{\mathbf{x}}(\mathbf{x}, \mathbf{u}) + r(\mathbf{x}, \mathbf{u}) \} = 0 \quad (37)$$

where U is the set of admissible commands.

Finally, it can be shown that no other function can substitute for $J^*(t)$ in (37) and still make it true; i.e. $J^*(t)$ is the *unique* solution to

$$\min_{\mathbf{u} \in U} \{ \nabla J(\mathbf{x}) \dot{\mathbf{x}}(\mathbf{x}, \mathbf{u}) + r(\mathbf{x}, \mathbf{u}) \} = 0 \quad (38)$$

which is called the *Hamilton-Jacobi-Bellman equation*.

So if we can find the $J()$ that solves the HJB equation then we have the optimal cost-to-go J^* (\cdot), and we get $\mathbf{u}^*(\cdot)$ as well because it is defined by (37). Unfortunately, in most problems the HJB equation can't be solved analytically. We have to look for *approximate* solutions.

Solving HJB Equations by Learning

There are many approaches, but most are versions of the following scheme. First we choose a policy $\mathbf{u}(\cdot)$ and a scalar function $\hat{J}(\cdot)$ to be our initial estimates of $\mathbf{u}^*(\cdot)$ and $J^*(\cdot)$. This initial $\mathbf{u}(\cdot)$ may be far from $\mathbf{u}^*(\cdot)$, but must at least have a finite cost-to-go. Often a suitable $\mathbf{u}(\cdot)$ can be found by feedback linearization.

This $\mathbf{u}(\cdot)$ and $\hat{J}(\cdot)$ imply an estimate of the Hamiltonian,

$$\hat{H}(\mathbf{x}, \mathbf{u}) = \nabla \hat{J}(\mathbf{x}) \dot{\mathbf{x}}(\mathbf{x}, \mathbf{u}(\mathbf{x})) + r(\mathbf{x}, \mathbf{u}(\mathbf{x})) \quad (39)$$

Our $\mathbf{u}(\cdot)$, $\hat{J}(\cdot)$ and $\hat{H}(\cdot)$ won't have the properties of the real $\mathbf{u}^*(\cdot)$, $J^*(\cdot)$ and $H^*(\cdot)$ as laid out in (37), i.e. our $\mathbf{u}(\cdot)$ won't yield commands \mathbf{u} that minimize $\hat{H}(\mathbf{x}, \mathbf{u})$, and the minimum of $\hat{H}(\mathbf{x}, \mathbf{u})$ won't be 0 for all \mathbf{x} . But we can gradually adjust $\mathbf{u}(\cdot)$ and $\hat{J}(\cdot)$ to get closer and closer to those properties. That is, we alternate 2 operations: Step 1, adjust $\hat{J}(\cdot)$ to bring \hat{H} closer to 0, i.e. to shrink \hat{H}^2 ; Step 2, adjust $\mathbf{u}(\cdot)$ to lower \hat{H} .

For example, we may choose an estimate $\hat{J}(\mathbf{x})$ of the form $\mathbf{w}\phi(\mathbf{x})$ where \mathbf{w} is a row-vector of weights and $\phi(\mathbf{x})$ is a feature vector. Then

$$\hat{H} = \mathbf{w} \nabla \phi \dot{\mathbf{x}} + r = \mathbf{w} \dot{\phi} + r \quad (40)$$

and

$$\partial \hat{H}^2 / \partial \mathbf{w} = 2\hat{H}\dot{\phi} \quad (41)$$

So in Step 1 we can adjust \mathbf{w} to shrink \hat{H}^2 by the LMS formula, $\Delta \mathbf{w} = -\eta \hat{H}\dot{\phi}^\top$. Similarly, we choose a parametrized $\mathbf{u}() — \mathbf{u}(\mathbf{x}, \mathbf{W})$ — so in Step 2 we can adjust \mathbf{W} to lower \hat{H} by gradient descent, $\Delta \mathbf{W} = -\eta_u [\partial \hat{H} / \partial \mathbf{W}]^\top$. There is an example in the code file HG_TO_1DOF.m.

```
>> run HG_TO_1DOF.m
```