

Reinforcement Learning and Optimal Control

by

Dimitri P. Bertsekas

Massachusetts Institute of Technology

Chapter 3 *Parametric Approximation*

DRAFT

This is Chapter 3 of the draft textbook “Reinforcement Learning and Optimal Control.” The chapter represents “work in progress,” and it will be periodically updated. It more than likely contains errors (hopefully not serious ones). Furthermore, its references to the literature are incomplete. Your comments and suggestions to the author at dimitrib@mit.edu are welcome. The date of last revision is given below.

December 14, 2018

Parametric Approximation

Contents

3.1. Approximation Architectures	p. 2
3.1.1. Linear and Nonlinear Feature-Based Architectures . .	p. 2
3.1.2. Training of Linear and Nonlinear Architectures . .	p. 7
3.1.3. Incremental Gradient and Newton Methods	p. 9
3.2. Neural Networks	p. 21
3.2.1. Training of Neural Networks	p. 24
3.2.2. Multilayer and Deep Neural Networks	p. 26
3.3. Sequential Dynamic Programming Approximation	p. 29
3.4. Q-factor Parametric Approximation	p. 31
3.5. Notes and Sources	p. 33

Clearly, for the success of approximation in value space, it is important to select a class of lookahead functions \tilde{J}_k that is suitable for the problem at hand. In the preceding chapter we discussed several methods for choosing \tilde{J}_k based mostly on problem approximation, including rollout. In this chapter we discuss an alternative approach, whereby \tilde{J}_k is chosen to be a member of a parametric class of functions, including neural networks, with the parameters “optimized” or “trained” by using some algorithm.

3.1 APPROXIMATION ARCHITECTURES

The starting point for the schemes of this chapter is a class of functions $\tilde{J}_k(x_k, r_k)$ that for each k , depend on the current state x_k and a vector $r_k = (r_{1,k}, \dots, r_{m,k})$ of m “tunable” scalar parameters, also called *weights*. By adjusting the weights, one can change the “shape” of \tilde{J}_k so that it is a reasonably good approximation to the true optimal cost-to-go function J_k . The class of functions $\tilde{J}_k(x_k, r_k)$ is called an *approximation architecture*, and the process of choosing the parameter vectors r_k is commonly called *training* or *tuning* the architecture.

The simplest training approach is to do some form of semi-exhaustive or semi-random search in the space of parameter vectors and adopt the parameters that result in best performance of the associated one-step lookahead controller (according to some criterion). More systematic approaches are based on numerical optimization, such as for example a least squares fit that aims to match the cost approximation produced by the architecture to a “training set,” i.e., a large number of pairs of state and cost values that are obtained through some form of sampling process. Throughout this chapter we will focus on this latter approach.

3.1.1 Linear and Nonlinear Feature-Based Architectures

There is a large variety of approximation architectures, based for example on polynomials, wavelets, discretization/interpolation schemes, neural networks, and others. A particularly interesting type of cost approximation involves *feature extraction*, a process that maps the state x_k into some vector $\phi_k(x_k)$, called the *feature vector* associated with x_k at time k . The vector $\phi_k(x_k)$ consists of scalar components

$$\phi_{1,k}(x_k), \dots, \phi_{m,k}(x_k),$$

called *features*. A feature-based cost approximation has the form

$$\tilde{J}_k(x_k, r_k) = \hat{J}_k(\phi_k(x_k), r_k),$$

where r_k is a parameter vector and \hat{J}_k is some function. Thus, the cost approximation depends on the state x_k through its feature vector $\phi_k(x_k)$.

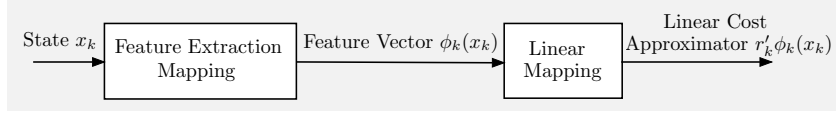


Figure 3.1.1 The structure of a linear feature-based architecture. We use a feature extraction mapping to generate an input $\phi_k(x_k)$ to a linear mapping defined by a weight vector r_k .

Note that we are allowing for different features $\phi_k(x_k)$ and different parameter vectors r_k for each stage k . This is necessary if the state space changes over time. On the other hand, for stationary problems with a long or infinite horizon, where the state space does not change with k , it is common to use the same features and parameters for all stages. The subsequent discussion can easily be adapted to infinite horizon methods, as we will discuss in Chapter 4.

Features are often handcrafted, based on whatever human intelligence, insight, or experience is available, and are meant to capture the most important characteristics of the current state. There are also systematic ways to construct features, including the use of neural networks, which we will discuss shortly. In this section, we provide a brief and selective discussion of architectures, and we refer to the specialized literature (e.g., Bertsekas and Tsitsiklis [BeT96], Bishop [Bis95], Haykin [Hay09], Sutton and Barto [SuB18]), and the author's [Ber12], Section 6.1.1, for more detailed presentations.

One idea behind using features is that the optimal cost-to-go functions J_k may be complicated nonlinear mappings, so it is sensible to try to break their complexity into smaller, less complex pieces. In particular, if the features encode much of the nonlinearity of J_k , we may be able to use a relatively simple architecture \hat{J}_k to approximate J_k . For example, with a well-chosen feature vector $\phi_k(x_k)$, a good approximation to the cost-to-go is often provided by *linearly* weighting the features, i.e.,

$$\tilde{J}_k(x_k, r_k) = \hat{J}_k(\phi_k(x_k), r_k) = \sum_{\ell=1}^m r_{\ell,k} \phi_{\ell,k}(x_k) = r'_k \phi_k(x_k), \quad (3.1)$$

where $r_{\ell,k}$ and $\phi_{\ell,k}(x_k)$ are the ℓ th components of r_k and $\phi_k(x_k)$, respectively, and $r'_k \phi_k(x_k)$ denotes the inner product of r_k and $\phi_k(x_k)$, viewed as column vectors of \mathbb{R}^m (see Fig. 3.1.1).

This is called a *linear feature-based architecture*, and the scalar parameters $r_{\ell,k}$ are also called *weights*. The approximating function $\tilde{J}_k(x_k, r_k)$ can be viewed as a member of the subspace spanned by the features $\phi_{\ell,k}(x_k)$, $\ell = 1, \dots, m$, which for this reason are also referred to *basis functions*. In addition to their simplicity, linear architectures have the advantage that they admit simpler training algorithms than their nonlinear counterparts. We provide a few examples, where for simplicity we drop the index k .

Example 3.1.1 (Piecewise Constant Approximation)

Suppose that the state space is partitioned in subsets S_1, \dots, S_m , so that every state belongs to one and only one subset. Let the ℓ th feature be defined by membership to the set S_ℓ , i.e.,

$$\phi_\ell(x) = \begin{cases} 1 & \text{if } x \in S_\ell, \\ 0 & \text{if } x \notin S_\ell. \end{cases}$$

Consider the architecture

$$\tilde{J}(x, r) = \sum_{\ell=1}^m r_\ell \phi_\ell(x),$$

where r is the vector consists of the m scalar parameters r_1, \dots, r_m . It can be seen that $\tilde{J}(x, r)$ is the piecewise constant function that has value r_ℓ for all states within the set S_ℓ .

The piecewise constant approximation is an example of a linear feature-based architecture that involves exclusively *local features*. These are features that take a nonzero value only for a relatively small subset of states. Thus a change of a single weight causes a change of the value of $\tilde{J}(x, r)$ for relatively few states x . At the opposite end we have linear feature-based architectures that involve *global features*. These are features that take nonzero values for a large number of states. The following is an example.

Example 3.1.2 (Polynomial Approximation)

An important case of linear architecture is one that uses polynomial basis functions. Suppose that the state consists of n components x^1, \dots, x^n , each taking values within some range of integers. For example, in a queueing system, x^i may represent the number of customers in the i th queue, where $i = 1, \dots, n$. Suppose that we want to use an approximating function that is quadratic in the components x^i . Then we can define a total of $1 + n + n^2$ basis functions that depend on the state $x = (x^1, \dots, x^n)$ via

$$\phi_0(x) = 1, \quad \phi_i(x) = x^i, \quad \phi_{ij}(x) = x^i x^j, \quad i, j = 1, \dots, n.$$

A linear approximation architecture that uses these functions is given by

$$\tilde{J}(x, r) = r_0 + \sum_{i=1}^n r_i x^i + \sum_{i=1}^n \sum_{j=i}^n r_{ij} x^i x^j,$$

where the parameter vector r has components r_0 , r_i , and r_{ij} , with $i = 1, \dots, n$, $j = k, \dots, n$. Indeed, any kind of approximating function that is polynomial in the components x^1, \dots, x^n can be constructed similarly.

A more general polynomial approximation may be based on some other known features of the state. For example, we may start with a feature vector

$$\phi(x) = (\phi_1(x), \dots, \phi_m(x))',$$

and transform it with a quadratic polynomial mapping. In this way we obtain approximating functions of the form

$$\tilde{J}(x, r) = r_0 + \sum_{i=1}^m r_i \phi_i(x) + \sum_{i=1}^m \sum_{j=k}^m r_{ij} \phi_i(x) \phi_j(x),$$

where the parameter r has components r_0 , r_i , and r_{ij} , with $i, j = 1, \dots, m$. This can also be viewed as a linear architecture that uses the basis functions

$$w_0(x) = 1, \quad w_i(x) = \phi_i(x), \quad w_{ij}(x) = \phi_i(x) \phi_j(x), \quad i, j = 1, \dots, m.$$

Here are two different types of examples of feature-based architectures involving games.

Example 3.1.3 (Tetris)

Let us revisit the game of tetris, which we discussed in Example 1.3.4. We can model the problem of finding an optimal playing strategy as a finite horizon problem with a very large horizon.

In Example 1.3.4 we viewed as state the pair of the board position x and the shape of the current falling block y . We viewed as control, the horizontal positioning and rotation applied to the falling block. However, the DP algorithm can be executed over the space of x only, since y is an uncontrollable state component. The optimal cost-to-go function is a vector of huge dimension (there are 2^{200} board positions in a “standard” tetris board of width 10 and height 20). However, it has been successfully approximated in practice by low-dimensional linear architectures.

In particular, the following features have been proposed in [BeI96]: the heights of the columns, the height differentials of adjacent columns, the wall height (the maximum column height), the number of holes of the board, and the constant 1 (the unit is often included as a feature in cost approximation architectures, as it allows for a constant shift in the approximating function). These features are readily recognized by tetris players as capturing important aspects of the board position.† There are a total of 22 features for a “standard” board with 10 columns. Of course the $2^{200} \times 22$ matrix of feature values cannot be stored in a computer, but for any board position, the corresponding row of features can be easily generated, and this is sufficient for implementation of the associated approximate DP algorithms. For recent works involving approximate DP methods and the preceding 22 features, see [Sch13] and [GGS13], which reference several other related papers.

† The use of feature-based approximate DP methods for the game of tetris was first suggested in the paper [TsV96], which introduced just two features (in addition to the constant 1): the wall height and the number of holes of the board. Most studies have used the set of features described here, but other sets of features have also been used; see [ThS09] and the discussion in [GGS13].

Example 3.1.4 (Computer Chess)

Computer chess programs that are based on feature-based architectures have been available for many years, and are still used widely (they have been upstaged in the mid-2010s by alternative types of chess programs that are based on neural network-based techniques that will be discussed later). These programs are based on approximate DP for minimax problems,[†] a feature-based parametric architecture, and multistep lookahead. The computer chess training methodology, however, is qualitatively different from the parametric approximation methods that we consider in this book.

In particular, with few exceptions, the training of chess architectures has been done with ad hoc hand-tuning techniques (as opposed to some form of optimization). Moreover, the features have traditionally been hand-crafted based on chess-specific knowledge (as opposed to automatically generated through a neural network or some other method). Indeed, it has been argued that the success of chess programs in outperforming the best humans, can be more properly attributed to the brute-force calculating power of modern computers, which provides long and accurate lookahead, than to the ability of simulation-based algorithmic approaches, which can learn powerful playing strategies that humans have difficulty conceiving or executing. For this reason, computer chess provides a distinct approximate DP paradigm, which may be better suited for some practical contexts where the limiting resource is raw computational power rather than innovative and sophisticated algorithms.

The fundamental paper on which all computer chess programs are based was written by one of the most illustrious modern-day applied mathematicians, C. Shannon [Sha50]. Shannon proposed a limited lookahead of a few moves and evaluating the end positions by means of a “scoring function” (in our terminology this plays the role of a cost function approximation). The scoring function may involve, for example, the calculation of a numerical value for each of a set of major features of a position that chess players easily recognize (such as material balance, mobility, pawn structure, and other positional factors), together with a method to combine these numerical values into a single score.

We may view the scoring function as a hand-crafted feature-based architecture for evaluating a chess position/state (cf. Fig. 3.1.2). In our earlier notation, it is a function of the form

$$\tilde{J}(x, r) = \hat{J}(\phi(x), r),$$

which maps a position x into a cost-to-go approximation $\hat{J}(\phi(x), r)$, a score whose value depends on the feature vector $\phi(x)$ of the position and a parameter vector r .[‡]

[†] We have not discussed DP for minimax problems and two-player games, but the ideas of approximation in value space apply to these contexts as well; see [Ber17] for a discussion that is focused on computer chess.

[‡] Because of the nature of the chess game, $\tilde{J}(x, r)$ does not depend critically on time. The duration of the game is unknown and so is the horizon of the

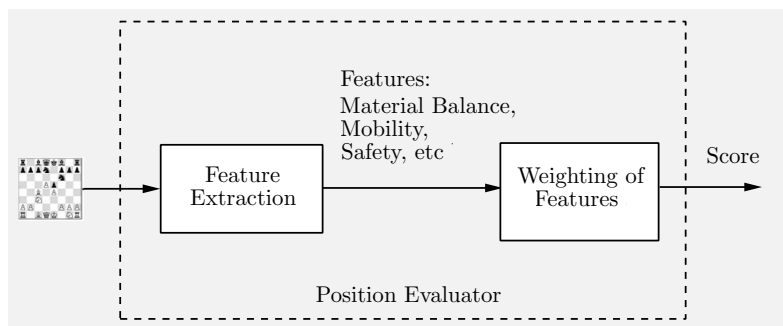


Figure 3.1.2 A feature-based architecture for computer chess.

In most computer chess programs, the features are weighted linearly, i.e., the architecture $\tilde{J}(x, r)$ that is used for limited lookahead is linear [cf. Eq. (3.1)]. In many cases, the weights are determined manually, by trial and error based on experience. However, in some programs, the weights are determined with pattern recognition techniques by training using examples of grandmaster play, i.e., by adjustment to bring the play of the program as close as possible to the play of chess grandmasters. This is a technique that applies more broadly in artificial intelligence; see Tesauro [Tes89], [Tes01].

Unfortunately, in many situations an adequate set of features is not available, so it is important to have methods that construct features automatically, to supplement whatever features may already exist. Indeed, there are architectures that do not rely on the knowledge of good features. One of the most popular is neural networks, which we will describe in the Section 3.2. Some of these architectures involve training that constructs simultaneously both the feature vectors $\phi_k(x_k)$ and the parameter vectors r_k that weigh them.

3.1.2 Training of Linear and Nonlinear Architectures

The process of choosing the parameter vector r of a parametric architecture $\tilde{J}(x, r)$ is generally referred to as training. The most common type of training is based on a *least squares* optimization, also known as *least squares regression*. Here a set of state-cost training pairs (x^s, β^s) , $s = 1, \dots, q$, called the *training set*, is collected and r is determined by solving the

problem. We are dealing essentially with an infinite horizon minimax problem, whose termination time is finite but unknown, similar to the stochastic shortest path problems to be discussed in Chapter 4. Still, however, chess programs often use features and weights that depend on the phase of the game (opening, middlegame, or endgame). Moreover the programs include specialized knowledge, such as opening and endgame databases. In our subsequent discussion we will ignore such possibilities.

problem

$$\min_r \sum_{s=1}^q (\tilde{J}(x^s, r) - \beta^s)^2. \quad (3.2)$$

Thus r is chosen to minimize the error between sample costs β^s and the architecture-predicted costs $\tilde{J}(x^s, r)$ in a least squares sense.

The cost function of the training problem (3.2) is generally nonconvex, which may pose challenges, since there may exist multiple local minima. However, for a linear architecture the cost function is convex quadratic, and the training problem admits a closed-form solution. In particular, for the linear architecture

$$\tilde{J}(x, r) = r' \phi(x_k),$$

the problem becomes

$$\min_r \sum_{s=1}^q (r' \phi(x^s) - \beta^s)^2.$$

By setting the gradient of the quadratic objective to 0 yields the linear equation

$$\sum_{s=1}^q \phi(x^s) (r' \phi(x^s) - \beta^s) = 0,$$

or

$$\sum_{s=1}^q \phi(x^s) \phi(x^s)' r = \sum_{s=1}^q \phi(x^s) \beta^s.$$

Thus by matrix inversion we obtain the minimizing parameter vector

$$\hat{r} = \left(\sum_{s=1}^q \phi(x^s) \phi(x^s)' \right)^{-1} \sum_{s=1}^q \phi(x^s) \beta^s. \quad (3.3)$$

(Since sometimes the inverse above may not exist, an additional quadratic in r , called a *regularization* function, is added to the least squares objective to deal with this, and also to help with other issues to be discussed later.)

Thus a linear architecture has the important advantage that the training problem can be solved exactly and conveniently with the formula (3.3) (of course it may be solved by any other algorithm that is suitable for linear least squares problems). By contrast, if we use a nonlinear architecture, such as a neural network, the associated least squares problem is nonquadratic and also nonconvex. Despite this fact, through a combination of sophisticated implementation of special gradient algorithms, called *incremental*, and powerful computational resources, neural network methods have been successful in practice as we will discuss in Section 3.2.

3.1.3 Incremental Gradient and Newton Methods

We will now digress to discuss special methods for solution of the least squares training problem (3.2), assuming a parametric architecture that is differentiable in the parameter vector. This methodology is properly viewed as a subject in nonlinear programming and iterative algorithms, and as such it can be studied independently of the approximate DP methods of this book. Thus the reader who has already some exposure to the subject may skip to the next section, and return later as needed.

Incremental methods have a rich theory, and our presentation in this section is brief, focusing primarily on implementation and intuition. Some surveys and book references are provided at the end of the chapter, which include a more detailed treatment, including convergence analysis. Since, we want to cover problems that are more general than the specific least squares training problem (3.2), we will adopt a broader formulation and notation that are standard in nonlinear programming.

Incremental Gradient Methods

We view the training problem (3.2) as a special case of the minimization of a sum of component functions

$$f(x) = \sum_{i=1}^m f_i(x), \quad (3.4)$$

where each f_i is a differentiable scalar function of the n -dimensional vector x (this is the parameter vector). Thus we use the more common symbols x and m in place of r and q , respectively, and we replace the least squares terms $\tilde{J}(x^s, r)$ in the training problem (3.2) with the generic terms $f_i(x)$.

The (ordinary) gradient method for problem (3.4) has the form[†]

$$x^{k+1} = x^k - \gamma^k \nabla f(x^k) = x^k - \gamma^k \sum_{i=1}^m \nabla f_i(x^k), \quad (3.5)$$

where γ^k is a positive stepsize parameter. Incremental gradient methods have the general form

$$x^{k+1} = x^k - \gamma^k \nabla f_{i_k}(x^k), \quad (3.6)$$

where i_k is some index from the set $\{1, \dots, m\}$, chosen by some deterministic or randomized rule. Thus a single component function f_{i_k} is used at

[†] We use standard calculus notation for gradients; see, e.g., [Ber16], Appendix A. In particular, $\nabla f(x)$ denotes the n dimensional vector whose components are the first partial derivatives $\partial f(x)/\partial x_i$ of f with respect to the components x_1, \dots, x_n of the vector x .

each iteration of the incremental method (3.6), with great economies in gradient calculation cost over the ordinary gradient method (3.5), particularly when m is large.

The method for selecting the index i_k of component to be iterated on at iteration k is important for the performance of the method. Three common rules are:

- (1) A *cyclic order*, whereby the indexes are taken up in the fixed deterministic order $1, \dots, m$, so that i_k is equal to $(k \text{ modulo } m) \text{ plus } 1$. A contiguous block of iterations involving f_1, \dots, f_m in this order and exactly once is called a *cycle*.
- (2) A *uniform random order*, whereby the index i_k chosen randomly by sampling over all indexes with a uniform distribution, independently of the past history of the algorithm. This rule may perform better than the cyclic rule in some circumstances.
- (3) A *cyclic order with random reshuffling*, whereby the indexes are taken up one by one within each cycle, but their order after each cycle is reshuffled randomly (and independently of the past). This rule is used widely in practice, particularly when the number of components m is modest, for reasons to be discussed later.

Note that in the cyclic case, it is essential to include all components in a cycle, for otherwise some components will be sampled more often than others, leading to a bias in the convergence process. Similarly, it is necessary to sample according to the uniform distribution in the random order case.

Focusing for the moment on the cyclic rule, we note that the motivation for the incremental gradient method is faster convergence: we hope that far from the solution, a single cycle of the method will be as effective as several (as many as m) iterations of the ordinary gradient method (think of the case where the components f_i are similar in structure). Near a solution, however, the incremental method may not be as effective.

To be more specific, we note that there are two complementary performance issues to consider in comparing incremental and nonincremental methods:

- (a) *Progress when far from convergence.* Here the incremental method can be much faster. For an extreme case take m large and all components f_i identical to each other. Then an incremental iteration requires m times less computation than a classical gradient iteration, but gives exactly the same result, when the stepsize is scaled to be m times larger. While this is an extreme example, it reflects the essential mechanism by which incremental methods can be much superior: far from the minimum a single component gradient will point to “more or less” the right direction, at least most of the time; see the following example.

- (b) *Progress when close to convergence.* Here the incremental method can be inferior. In particular, the ordinary gradient method (3.5) can be shown to converge with a constant stepsize under reasonable assumptions, see e.g., [Ber16], Chapter 1. However, the incremental method requires a diminishing stepsize, and its ultimate rate of convergence can be much slower.

This type of behavior is illustrated in the following example.

Example 3.1.5

Assume that x is a scalar, and that the problem is

$$\begin{aligned} &\text{minimize} && f(x) = \frac{1}{2} \sum_{i=1}^m (c_i x - b_i)^2 \\ &\text{subject to} && x \in \mathbb{R}, \end{aligned}$$

where c_i and b_i are given scalars with $c_i \neq 0$ for all i . The minimum of each of the components $f_i(x) = \frac{1}{2}(c_i x - b_i)^2$ is

$$x_i^* = \frac{b_i}{c_i},$$

while the minimum of the least squares cost function f is

$$x^* = \frac{\sum_{i=1}^m c_i b_i}{\sum_{i=1}^m c_i^2}.$$

It can be seen that x^* lies within the range of the component minima

$$R = \left[\min_i x_i^*, \max_i x_i^* \right],$$

and that for all x outside the range R , the gradient

$$\nabla f_i(x) = c_i(c_i x - b_i)$$

has the same sign as $\nabla f(x)$ (see Fig. 3.1.3). As a result, when outside the region R , the incremental gradient method

$$x^{k+1} = x^k - \gamma^k c_{i_k} (c_{i_k} x^k - b_{i_k})$$

approaches x^* at each step, provided the stepsize γ^k is small enough. In fact it can be verified that it is sufficient that

$$\gamma^k \leq \min_i \frac{1}{c_i^2}.$$

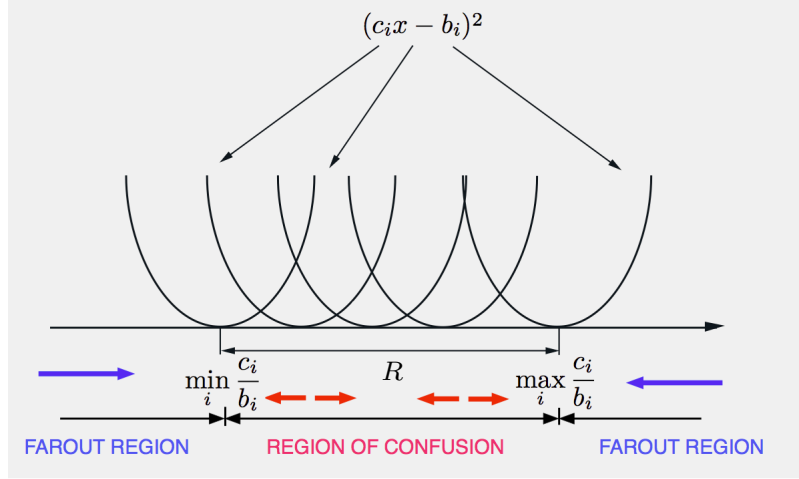


Figure 3.1.3. Illustrating the advantage of incrementalism when far from the optimal solution. The region of component minima

$$R = \left[\min_i x_i^*, \max_i x_i^* \right],$$

is labeled as the “region of confusion.” It is the region where the method does not have a clear direction towards the optimum. The i th step in an incremental gradient cycle is a gradient step for minimizing $(c_i x - b_i)^2$, so if x lies outside the region of component minima $R = [\min_i x_i^*, \max_i x_i^*]$, (labeled as the “farout region”) and the stepsize is small enough, progress towards the solution x^* is made.

However, for x inside the region R , the i th step of a cycle of the incremental gradient method need not make progress. It will approach x^* (for small enough stepsize γ^k) only if the current point x^k does not lie in the interval connecting x_i^* and x^* . This induces an oscillatory behavior within the region R , and as a result, the incremental gradient method will typically not converge to x^* unless $\gamma^k \rightarrow 0$. By contrast, the ordinary gradient method, which takes the form

$$x^{k+1} = x^k - \gamma \sum_{i=1}^m c_i (c_i x^k - b_i),$$

can be verified to converge to x^* for any constant stepsize γ with

$$0 < \gamma \leq \frac{1}{\sum_{i=1}^m c_i^2}.$$

However, for x outside the region R , a full iteration of the ordinary gradient method need not make more progress towards the solution than a single step of

the incremental gradient method. In other words, with comparably intelligent stepsize choices, *far from the solution (outside R), a single pass through the entire set of cost components by incremental gradient is roughly as effective as m passes by ordinary gradient.*

The discussion of the preceding example relies on x being one-dimensional, but in many multidimensional problems the same qualitative behavior can be observed. In particular, a pass through the i th component f_i by the incremental gradient method can make progress towards the solution in the region where the component gradient $\nabla f_{i_k}(x^k)$ makes an angle less than 90 degrees with the cost function gradient $\nabla f(x^k)$. If the components f_i are not “too dissimilar”, this is likely to happen in a region of points that are not too close to the optimal solution set. This behavior has been verified in many practical contexts, including the training of neural networks (cf. the next section), where incremental gradient methods have been used extensively, frequently under the name *backpropagation methods*.

Stepsize Choice and Diagonal Scaling

The choice of the stepsize γ^k plays an important role in the performance of incremental gradient methods. On close examination, it turns out that the direction used by the method differs from the gradient direction by an error that is proportional to the stepsize, and for this reason a diminishing stepsize is essential for convergence to a local minimum of f (convergence to a local minimum is the best we hope for since the cost function may not be convex).

However, it turns out that a peculiar form of convergence also typically occurs for a constant but sufficiently small stepsize. In this case, the iterates converge to a “limit cycle”, whereby the i th iterates ψ_i within the cycles converge to a different limit than the j th iterates ψ_j for $i \neq j$. The sequence $\{x^k\}$ of the iterates obtained at the end of cycles converges, except that the limit obtained *need not* be a minimum of f , even when f is convex. The limit tends to be close to the minimum point when the constant stepsize is small (see Section 2.4 of [Ber16] for analysis and examples). In practice, it is common to use a constant stepsize for a (possibly prespecified) number of iterations, then decrease the stepsize by a certain factor, and repeat, up to the point where the stepsize reaches a prespecified floor value.

An alternative possibility is to use a diminishing stepsize rule of the form

$$\gamma^k = \min \left\{ \gamma, \frac{\beta_1}{k + \beta_2} \right\},$$

where γ , β_1 , and β_2 are some positive scalars. There are also variants of the incremental gradient method that use a constant stepsize throughout, and generically converge to a stationary point of f at a linear rate. In one

type of such method the degree of incrementalism gradually diminishes as the method progresses (see [Ber97a]). Another incremental approach with similar aims, is the aggregated incremental gradient method, which will be discussed later in this section.

Regardless of whether a constant or a diminishing stepsize is ultimately used, to maintain the advantage of faster convergence when far from the solution, the incremental method must use a much larger stepsize than the corresponding nonincremental gradient method (as much as m times larger so that the size of the incremental gradient step is comparable to the size of the nonincremental gradient step).

One possibility is to use an adaptive stepsize rule, whereby the stepsize is reduced (or increased) when the progress of the method indicates that the algorithm is oscillating because it operates within (or outside, respectively) the region of confusion. There are formal ways to implement such stepsize rules with sound convergence properties (see [Tse98], [MYF03], [GOP15a]).

The difficulty with stepsize selection may also be addressed with *diagonal scaling*, i.e., using a stepsize γ_i^k that is different for each of the components x_i of x . Second derivatives, or approximations thereof, can be very effective for this purpose; see also the discussion on incremental Newton methods later in this section.

Stochastic Gradient Descent

Incremental gradient methods are related to methods that aim to minimize an expected value

$$f(x) = E\{F(x, w)\},$$

where w is a random variable, and $F(\cdot, w) : \Re^n \mapsto \Re$ is a differentiable function for each possible value of w . The *stochastic gradient method* for minimizing f is given by

$$x^{k+1} = x^k - \gamma^k \nabla_x F(x^k, w^k), \quad (3.7)$$

where w^k is a sample of w and $\nabla_x F$ denotes gradient of F with respect to x . This method has a rich theory and a long history, and it is strongly related to the classical algorithmic field of *stochastic approximation*; see the books [BeT96], [KuY03], [Spa03], [Mey07], [Bor08], [BPP13]. The method is also often referred to as *stochastic gradient descent*, particularly in the context of machine learning applications.

If we view the expected value cost $E\{F(x, w)\}$ as a weighted sum of cost function components, we see that the stochastic gradient method (3.7) is related to the incremental gradient method

$$x^{k+1} = x^k - \gamma^k \nabla f_{i_k}(x^k) \quad (3.8)$$

for minimizing a finite sum $\sum_{i=1}^m f_i$, when randomization is used for component selection. An important difference is that the former method involves stochastic sampling of cost components $F(x, w)$ from a possibly infinite population, under some probabilistic assumptions, while in the latter the set of cost components f_i is predetermined and finite. However, it is possible to view the incremental gradient method (3.8), with uniform randomized selection of the component function f_i (i.e., with i_k chosen to be any one of the indexes $1, \dots, m$, with equal probability $1/m$, and independently of preceding choices), as a stochastic gradient method.

Despite the apparent similarity of the incremental and the stochastic gradient methods, the view that the problem

$$\begin{aligned} \text{minimize} \quad & f(x) = \sum_{i=1}^m f_i(x) \\ \text{subject to} \quad & x \in \mathbb{R}^n, \end{aligned} \tag{3.9}$$

can simply be treated as a special case of the problem

$$\begin{aligned} \text{minimize} \quad & f(x) = E\{F(x, w)\} \\ \text{subject to} \quad & x \in \mathbb{R}^n, \end{aligned}$$

is questionable.

One reason is that once we convert the finite sum problem to a stochastic problem, we preclude the use of methods that exploit the finite sum structure, such as the incremental aggregated gradient methods to be discussed in the next subsection. Another reason is that the finite-component problem (3.9) is often genuinely deterministic, and to view it as a stochastic problem at the outset may mask some of its important characteristics, such as the number m of cost components, or the sequence in which the components are ordered and processed. These characteristics may potentially be algorithmically exploited. For example, with insight into the problem's structure, one may be able to discover a special deterministic or partially randomized order of processing the component functions that is superior to a uniform randomized order. On the other hand analysis indicates that in the absence of problem-specific knowledge that can be exploited to select a favorable deterministic order, a uniform randomized order (each component f_i chosen with equal probability $1/m$ at each iteration, independently of preceding choices) sometimes has superior worst-case complexity; see [NeB00], [NeB01], [BNO03], [Ber15a], [WaB16].

Finally, let us note the popular hybrid technique, which reshuffles randomly the order of the cost components after each cycle. Practical experience and recent analysis [GOP15c] indicates that it has somewhat better performance to the uniform randomized order when m is large. One possible reason is that random reshuffling allocates exactly one computation slot

to each component in an m -slot cycle, while uniform sampling allocates one computation slot to each component *on the average*. A nonzero variance in the number of slots that any fixed component gets within a cycle, may be detrimental to performance. While it seems difficult to establish this fact analytically, a justification is suggested by the view of the incremental gradient method as a gradient method with error in the computation of the gradient. The error has apparently greater variance in the uniform sampling method than in the random reshuffling method, and heuristically, if the variance of the error is larger, the direction of descent deteriorates, suggesting slower convergence.

Incremental Aggregated Gradient Methods

Another incremental algorithm is the *incremental aggregated gradient method*, which has the form

$$x^{k+1} = x^k - \gamma^k \sum_{\ell=0}^{m-1} \nabla f_{i_{k-\ell}}(x^{k-\ell}), \quad (3.10)$$

where f_{i_k} is the new component function selected for iteration k .[†] In the most common version of the method the component indexes i_k are selected in a cyclic order [$i_k = (k \text{ modulo } m) + 1$]. Random selection of the index i_k has also been suggested.

From Eq. (3.10) it can be seen that the method computes the gradient incrementally, one component per iteration. However, in place of the single component gradient $\nabla f_{i_k}(x^k)$, used in the incremental gradient method (3.6), it uses the sum of the component gradients computed in the past m iterations, which is an approximation to the total cost gradient $\nabla f(x^k)$.

The idea of the method is that by aggregating the component gradients one may be able to reduce the error between the true gradient $\nabla f(x^k)$ and the incrementally computed approximation used in Eq. (3.10), and thus attain a faster asymptotic convergence rate. Indeed, it turns out that under certain conditions the method exhibits a linear convergence rate, just like in the nonincremental gradient method, without incurring the cost of a full gradient evaluation at each iteration (a strongly convex cost function and with a sufficiently small constant stepsize are required for this; see [Ber16], Section 2.4.2). This is in contrast with the incremental gradient method (3.6), for which a linear convergence rate can be achieved only at the expense of asymptotic error, as discussed earlier.

A disadvantage of the aggregated gradient method (3.10) is that it requires that the most recent component gradients be kept in memory, so that when a component gradient is reevaluated at a new point, the

[†] In the case where $k < m$, the summation in Eq. (3.10) should go up to $\ell = k$, and the stepsize should be replaced by a corresponding larger value.

preceding gradient of the same component is discarded from the sum of gradients of Eq. (3.10). There have been alternative implementations that ameliorate this memory issue, by recalculating the full gradient periodically and replacing an old component gradient by a new one. More specifically, instead of the gradient sum

$$s^k = \sum_{\ell=0}^{m-1} \nabla f_{i_{k-\ell}}(x^{k-\ell}),$$

in Eq. (3.10), these methods use

$$\tilde{s}^k = \nabla f_{i_k}(x^k) - \nabla f_{i_k}(\tilde{x}^k) + \sum_{\ell=0}^{m-1} \nabla f_{i_{k-\ell}}(\tilde{x}^k),$$

where \tilde{x}^k is the most recent point where the full gradient has been calculated. To calculate \tilde{s}_k one only needs to compute the difference of the two gradients

$$\nabla f_{i_k}(x^k) - \nabla f_{i_k}(\tilde{x}^k)$$

and add it to the full gradient $\sum_{\ell=0}^{m-1} \nabla f_{i_{k-\ell}}(\tilde{x}^k)$. This bypasses the need for extensive memory storage, and with proper implementation, typically leads to small degradation in performance. However, periodically calculating the full gradient when m is very large can be a drawback. Another potential drawback of the aggregated gradient method is that for a large number of terms m , one hopes to converge within the first cycle through the components f_i , thereby reducing the effect of aggregating the gradients of the components.

Incremental Newton Methods

We will now consider an incremental version of Newton's method for unconstrained minimization of an additive cost function of the form

$$f(x) = \sum_{i=1}^m f_i(x),$$

where the functions f_i are convex and twice continuously differentiable.[†]

[†] We will denote by $\nabla^2 f(x)$ the $n \times n$ Hessian matrix of f at x , i.e., the matrix whose (i, j) th component is the second partial derivative $\partial^2 f(x) / \partial x^i \partial x^j$. A beneficial consequence of assuming convexity of f_i is that the Hessian matrices $\nabla^2 f_i(x)$ are positive semidefinite, which facilitates the implementation of the algorithms to be described. On the other hand, the algorithmic ideas of this section may also be adapted for the case where f_i are nonconvex.

The ordinary Newton's method, at the current iterate x^k , obtains the next iterate x^{k+1} by minimizing over x the quadratic approximation/second order expansion

$$\tilde{f}(x; x^k) = \nabla f(x^k)'(x - x^k) + \frac{1}{2}(x - x^k)'\nabla^2 f(x^k)(x - x^k),$$

of f at x^k . Similarly, the incremental form of Newton's method minimizes a sum of quadratic approximations of components of the general form

$$\tilde{f}_i(x; \psi) = \nabla f_i(\psi)'(x - \psi) + \frac{1}{2}(x - \psi)'\nabla^2 f_i(\psi)(x - \psi), \quad (3.11)$$

as we will now explain.

As in the case of the incremental gradient method, we view an iteration as a cycle of m subiterations, each involving a single additional component f_i , and its gradient and Hessian at the current point within the cycle. In particular, if x^k is the vector obtained after k cycles, the vector x^{k+1} obtained after one more cycle is

$$x^{k+1} = \psi_{m,k},$$

where starting with $\psi_{0,k} = x^k$, we obtain $\psi_{m,k}$ after the m steps

$$\psi_{i,k} \in \arg \min_{x \in \mathbb{R}^n} \sum_{\ell=1}^i \tilde{f}_\ell(x; \psi_{\ell-1,k}), \quad i = 1, \dots, m, \quad (3.12)$$

where \tilde{f}_ℓ is defined as the quadratic approximation (3.11). If all the functions f_i are quadratic, it can be seen that the method finds the solution in a single cycle.[†] The reason is that when f_i is quadratic, each $f_i(x)$ differs from $\tilde{f}_i(x; \psi)$ by a constant, which does not depend on x . Thus the difference

$$\sum_{i=1}^m f_i(x) - \sum_{i=1}^m \tilde{f}_i(x; \psi_{i-1,k})$$

is a constant that is independent of x , and minimization of either sum in the above expression gives the same result.

[†] Here we assume that the m quadratic minimizations (3.12) to generate $\psi_{m,k}$ have a solution. For this it is sufficient that the first Hessian matrix $\nabla^2 f_1(x^0)$ be positive definite, in which case there is a unique solution at every iteration. A simple possibility to deal with this requirement is to add to f_1 a small positive regularization term, such as $\frac{\epsilon}{2}\|x - x^0\|^2$. A more sound possibility is to lump together several of the component functions (enough to ensure that the sum of their quadratic approximations at x^0 is positive definite), and to use them in place of f_1 . This is generally a good idea and leads to smoother initialization, as it ensures a relatively stable behavior of the algorithm for the initial iterations.

It is important to note that the computations of Eq. (3.12) can be carried out efficiently. For simplicity, let us assume that $\tilde{f}_1(x; \psi)$ is a positive definite quadratic, so that for all i , $\psi_{i,k}$ is well defined as the unique solution of the minimization problem in Eq. (3.12). We will show that the incremental Newton method (3.12) can be implemented in terms of the incremental update formula

$$\psi_{i,k} = \psi_{i-1,k} - D_{i,k} \nabla f_i(\psi_{i-1,k}), \quad (3.13)$$

where $D_{i,k}$ is given by

$$D_{i,k} = \left(\sum_{\ell=1}^i \nabla^2 f_\ell(\psi_{\ell-1,k}) \right)^{-1}, \quad (3.14)$$

and is generated iteratively as

$$D_{i,k} = \left(D_{i-1,k}^{-1} + \nabla^2 f_i(\psi_{i-1,k}) \right)^{-1}. \quad (3.15)$$

Indeed, from the definition of the method (3.12), the quadratic function $\sum_{\ell=1}^{i-1} \tilde{f}_\ell(x; \psi_{\ell-1,k})$ is minimized by $\psi_{i-1,k}$ and its Hessian matrix is $D_{i-1,k}^{-1}$, so we have

$$\sum_{\ell=1}^{i-1} \tilde{f}_\ell(x; \psi_{\ell-1,k}) = \frac{1}{2}(x - \psi_{\ell-1,k})' D_{i-1,k}^{-1} (x - \psi_{\ell-1,k}) + \text{constant}.$$

Thus, by adding $\tilde{f}_i(x; \psi_{i-1,k})$ to both sides of this expression, we obtain

$$\begin{aligned} \sum_{\ell=1}^i \tilde{f}_\ell(x; \psi_{\ell-1,k}) &= \frac{1}{2}(x - \psi_{\ell-1,k})' D_{i-1,k}^{-1} (x - \psi_{\ell-1,k}) + \text{constant} \\ &\quad + \frac{1}{2}(x - \psi_{i-1,k})' \nabla^2 f_i(\psi_{i-1,k}) (x - \psi_{i-1,k}) + \nabla f_i(\psi_{i-1,k})' (x - \psi_{i-1,k}). \end{aligned}$$

Since by definition $\psi_{i,k}$ minimizes this function, we obtain Eqs. (3.13)-(3.15).

The recursion (3.15) for the matrix $D_{i,k}$ can often be efficiently implemented by using convenient formulas for the inverse of the sum of two matrices. In particular, if f_i is given by

$$f_i(x) = h_i(a_i'x - b_i),$$

for some twice differentiable convex function $h_i : \mathbb{R} \mapsto \mathbb{R}$, vector a_i , and scalar b_i , we have

$$\nabla^2 f_i(\psi_{i-1,k}) = \nabla^2 h_i(\psi_{i-1,k}) a_i a_i',$$

and the recursion (3.15) can be written as

$$D_{i,k} = D_{i-1,k} - \frac{D_{i-1,k} a_i' D_{i-1,k}}{\nabla^2 h_i(\psi_{i-1,k})^{-1} + a_i' D_{i-1,k} a_i};$$

this is the well-known Sherman-Morrison formula for the inverse of the sum of an invertible matrix and a rank-one matrix.

We have considered so far a single cycle of the incremental Newton method. Similar to the case of the incremental gradient method, we may cycle through the component functions multiple times. In particular, we may apply the incremental Newton method to the extended set of components

$$f_1, f_2, \dots, f_m, f_1, f_2, \dots, f_m, f_1, f_2, \dots$$

The resulting method asymptotically resembles an incremental gradient method with diminishing stepsize of the type described earlier. Indeed, from Eq. (3.14)], the matrix $D_{i,k}$ diminishes roughly in proportion to $1/k$. From this it follows that the asymptotic convergence properties of the incremental Newton method are similar to those of an incremental gradient method with diminishing stepsize of order $O(1/k)$. Thus its convergence rate is slower than linear.

To accelerate the convergence of the method one may employ a form of restart, so that $D_{i,k}$ does not converge to 0. For example $D_{i,k}$ may be reinitialized and increased in size at the beginning of each cycle. For problems where f has a unique nonsingular minimum x^* [one for which $\nabla^2 f(x^*)$ is nonsingular], one may design incremental Newton schemes with restart that converge linearly to within a neighborhood of x^* (and even superlinearly if x^* is also a minimum of all the functions f_i , so there is no region of confusion). Alternatively, the update formula (3.15) may be modified to

$$D_{i,k} = \left(\beta_k D_{i-1,k}^{-1} + \nabla^2 f_\ell(\psi_{i,k}) \right)^{-1}, \quad (3.16)$$

by introducing a parameter $\beta_k \in (0, 1)$, which can be used to accelerate the practical convergence rate of the method; cf. the discussion of the incremental Gauss-Newton methods.

Incremental Newton Method with Diagonal Approximation

Generally, with proper implementation, the incremental Newton method is often substantially faster than the incremental gradient method, in terms of numbers of iterations (there are theoretical results suggesting this property for stochastic versions of the two methods; see the end-of-chapter references). However, in addition to computation of second derivatives, the incremental Newton method involves greater overhead per iteration due to matrix-vector calculations in Eqs. (3.13), (3.15), and (3.16), so it is suitable only for problems where n , the dimension of x , is relatively small.

A way to remedy in part this difficulty is to approximate $\nabla^2 f_i(\psi_{i,k})$ by a diagonal matrix, and recursively update a diagonal approximation of $D_{i,k}$ using Eqs. (3.15) or (3.16). In particular, one may set to 0 the off-diagonal components of $\nabla^2 f_i(\psi_{i,k})$. In this case, the iteration (3.13) becomes a diagonally scaled version of the incremental gradient method, and involves comparable overhead per iteration (assuming the required diagonal second derivatives are easily computed or approximated). As an additional scaling option, one may multiply the diagonal components with a stepsize parameter that is close to 1 and add a small positive constant (to bound them away from 0). Ordinarily this method is easily implementable, and requires little experimentation with stepsize selection.

3.2 NEURAL NETWORKS

There are several different types of neural networks that can be used for a variety of tasks, such as pattern recognition, classification, image and speech recognition, and others. We focus here on our finite horizon DP context, and the role that neural networks can play in approximating the optimal cost-to-go functions J_k . As an example within this context, we may first use a neural network to construct an approximation to J_{N-1} . Then we may use this approximation to approximate J_{N-2} , and continue this process backwards in time, to obtain approximations to all the optimal cost-to-go functions J_k , $k = 1, \dots, N-1$, as we will discuss in more detail in Section 3.3.

To describe the use of neural networks in finite horizon DP, let us consider the typical stage k , and for convenience drop the index k ; the subsequent discussion applies to each value of k separately. We consider parametric architectures $\tilde{J}(x, v, r)$ of the form

$$\tilde{J}(x, v, r) = r' \phi(x, v) \quad (3.17)$$

that depend on two parameter vectors v and r . Our objective is to select v and r so that $\tilde{J}(x, v, r)$ approximates some cost function that can be sampled (possibly with some error). The process is to collect a training set that consists of a large number of state-cost pairs (x^s, β^s) , $s = 1, \dots, q$, and to find a function $\tilde{J}(x, v, r)$ of the form (3.17) that matches the training set in a least squares sense, i.e., (v, r) minimizes

$$\sum_{s=1}^q (\tilde{J}(x^s, v, r) - \beta^s)^2.$$

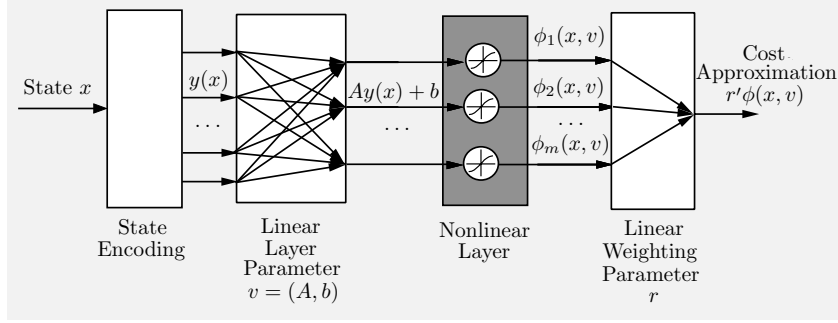


Figure 3.2.1 A perceptron consisting of a linear layer and a nonlinear layer. It provides a way to compute features of the state, which can be used for cost function approximation. The state x is encoded as a vector of numerical values $y(x)$, which is then transformed linearly as $Ay(x) + b$ in the linear layer. The m scalar output components of the linear layer, become the inputs to nonlinear functions that produce the m scalars $\phi_\ell(x, v) = \sigma((Ay(x) + b)_\ell)$, which can be viewed as features that are in turn linearly weighted with parameters r_ℓ .

We postpone for later the question of how the training pairs (x^s, β^s) are generated, and how the training problem is solved.[†] Notice the different roles of the two parameter vectors here: v parametrizes $\phi(x, v)$, which in some interpretation may be viewed as a feature vector, and r is a vector of linear weighting parameters for the components of $\phi(x, v)$.

A neural network architecture provides a parametric class of functions $\tilde{J}(x, v, r)$ of the form (3.17) that can be used in the optimization framework just described. The simplest type of neural network is the *single layer perceptron*; see Fig. 3.2.1. Here the state x is encoded as a vector of numerical values $y(x)$ with components $y_1(x), \dots, y_n(x)$, which is then transformed linearly as

$$Ay(x) + b,$$

where A is an $m \times n$ matrix and b is a vector in \mathbb{R}^m .[‡] This transformation is called the *linear layer* of the neural network. We view the components of A and b as parameters to be determined, and we group them together into the parameter vector $v = (A, b)$.

[†] The least squares training problem used here is based on *nonlinear regression*. This is a classical method for approximating the expected value of a function with a parametric architecture, and involves a least squares fit of the architecture to simulation-generated samples of the expected value. We refer to machine learning and statistics textbooks for more discussion.

[‡] The method of encoding x into the numerical vector $y(x)$ is problem-dependent, but it is important to note that some of the components of $y(x)$ could be some known interesting features of x that can be designed based on problem-specific knowledge.

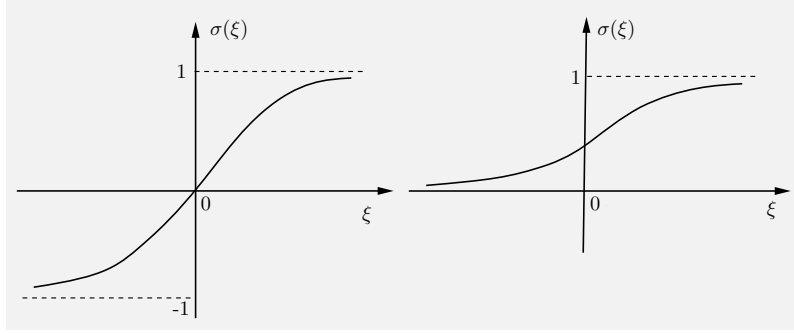


Figure 3.2.2 Some examples of sigmoid functions. The hyperbolic tangent function is on the left, while the logistic function is on the right.

Each of the m scalar output components of the linear layer,

$$(Ay(x) + b)_\ell, \quad \ell = 1, \dots, m,$$

becomes the input to a nonlinear differentiable function σ that maps scalars to scalars. Typically σ is differentiable and monotonically increasing. A simple and popular possibility is the *rectified linear unit*, which is simply the function $\max\{0, \xi\}$, “rectified” to a differentiable function by some form of smoothing operation; for example $\sigma(\xi) = \ln(1 + e^\xi)$. Other functions, used since the early days of neural networks, have the property

$$-\infty < \lim_{\xi \rightarrow -\infty} \sigma(\xi) < \lim_{\xi \rightarrow \infty} \sigma(\xi) < \infty;$$

see Fig. 3.2.2. Such functions are called *sigmoids*, and some common choices are the *hyperbolic tangent* function

$$\sigma(\xi) = \tanh(\xi) = \frac{e^\xi - e^{-\xi}}{e^\xi + e^{-\xi}},$$

and the *logistic* function

$$\sigma(\xi) = \frac{1}{1 + e^{-\xi}}.$$

In what follows, we will ignore the character of the function σ (except for differentiability), and simply refer to it as a “nonlinear unit” and to the corresponding layer as a “nonlinear layer.”

At the outputs of the nonlinear units, we obtain the scalars

$$\phi_\ell(x, v) = \sigma((Ay(x) + b)_\ell), \quad \ell = 1, \dots, m.$$

One possible interpretation is to view $\phi_\ell(x, v)$ as features of x , which are linearly combined using weights r_ℓ , $\ell = 1, \dots, m$, to produce the final output

$$\tilde{J}(x, v, r) = \sum_{\ell=1}^m r_\ell \phi_\ell(x, v) = \sum_{\ell=1}^m r_\ell \sigma((Ay(x) + b)_\ell).$$

Note that each value $\phi_\ell(x, v)$ depends on just the ℓ th row of A and the ℓ th component of b , not on the entire vector v . In some cases this motivates placing some constraints on individual components of A and b to achieve special problem-dependent “handcrafted” effects.

3.2.1 Training of Neural Networks

Given a set of state-cost training pairs (x^s, β^s) , $s = 1, \dots, q$, the parameters of the neural network A , b , and r are obtained by solving the training problem

$$\min_{A, b, r} \sum_{s=1}^q \left(\sum_{\ell=1}^m r_\ell \sigma((Ay(x^s) + b)_\ell) - \beta^s \right)^2. \quad (3.18)$$

Note that the cost function of this problem is generally nonconvex, so there may exist multiple local minima.

In practice it is common to augment the cost function of this problem with a *regularization* function, such as a quadratic in the parameters A , b , and r . This is customary in least squares problems in order to make the problem easier to solve algorithmically. However, in the context of neural network training, regularization is primarily important for a different reason: it helps to avoid *overfitting*, which refers to a situation where a neural network model matches the training data very well but does not do as well on new data. This is a well known difficulty in machine learning, which may occur when the number of parameters of the neural network is relatively large (roughly comparable to the size of the training set). We refer to machine learning and neural network textbooks for a discussion of algorithmic questions regarding regularization and other issues that relate to the practical implementation of the training process. In any case, the training problem (3.18) is an unconstrained nonconvex differentiable optimization problem that can in principle be addressed with any of the standard gradient-type methods. Significantly, it is well-suited for the incremental methods discussed in Section 3.1.3.

Let us now consider a few issues regarding the neural network formulation and training process just described:

- (a) The first issue is to select a method to solve the training problem (3.18). While we can use any unconstrained optimization method that is based on gradients, in practice it is important to take into account the cost function structure of problem (3.18). The salient characteristic of this cost function is that it is the sum of a potentially very large number q of component functions. This makes the computation of the cost function value of the training problem and/or its gradient very costly. For this reason the incremental methods of Section 3.1.3

are universally used for training.[†] Experience has shown that these methods can be vastly superior to their nonincremental counterparts in the context of neural network training.

The implementation of the training process has benefited from experience that has been accumulated over time, and has provided guidelines for scaling, regularization, initial parameter selection, and other practical issues; we refer to books on neural networks such as Bishop [Bis95], Goodfellow, Bengio, and Courville [GBC16], and Haykin [Hay09] for related accounts. Still, incremental methods can be quite slow, and training may be a time-consuming process. Fortunately, the training is ordinarily done off-line, in which case computation time may not be a serious issue. Moreover, in practice the neural network training problem typically need not be solved with great accuracy.

- (b) Another important question is how well we can approximate the optimal cost-to-go functions J_k with a neural network architecture, assuming we can choose the number of the nonlinear units m to be as large as we want. The answer to this question is quite favorable and is provided by the so-called *universal approximation theorem*. Roughly, the theorem says that assuming that x is an element of a Euclidean space X and $y(x) \equiv x$, a neural network of the form described can approximate arbitrarily closely (in an appropriate mathematical sense), over a closed and bounded subset $S \subset X$, any piecewise continuous function $J : S \mapsto \mathbb{R}$, provided the number m of nonlinear units is sufficiently large. For proofs of the theorem at different levels of generality, we refer to Cybenko [Cyb89], Funahashi [Fun89], Hornik, Stinchcombe, and White [HSW89], and Leshno et al. [LLP93]. For intuitive explanations we refer to Bishop ([Bis95], pp. 129-130) and Jones [Jon90].
- (c) While the universal approximation theorem provides some assurance about the adequacy of the neural network structure, it does not predict how many nonlinear units we may need for “good” performance in a given problem. Unfortunately, this is a difficult question to even pose precisely, let alone to answer adequately. In practice, one is reduced to trying increasingly larger values of m until one is convinced that satisfactory performance has been obtained for the task at hand. Experience has shown that in many cases the number of re-

[†] The incremental methods are valid for an arbitrary order of component selection within the cycle, but it is common to randomize the order at the beginning of each cycle. Also, in a variation of the basic method, we may operate on a batch of several components at each iteration, called a *minibatch*, rather than a single component. This has an averaging effect, which reduces the tendency of the method to oscillate and allows for the use of a larger stepsize; see the end-of-chapter references.

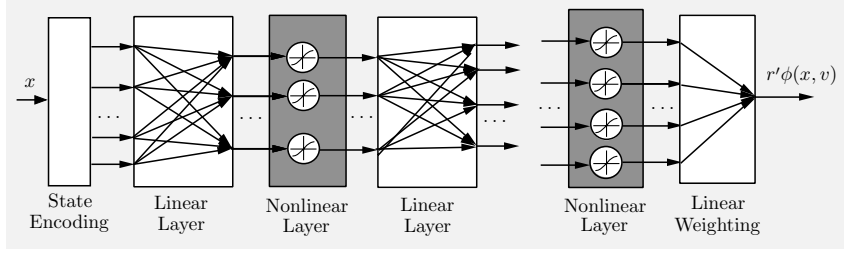


Figure 3.2.3 A neural network with multiple layers. Each nonlinear layer constructs the inputs of the next linear layer.

quired nonlinear units and corresponding dimension of A can be very large, adding significantly to the difficulty of solving the training problem. This has given rise to many suggestions for modifications of the neural network structure. One possibility is to concatenate multiple single layer perceptrons so that the output of the nonlinear layer of one perceptron becomes the input to the linear layer of the next, as we will now discuss.

3.2.2 Multilayer and Deep Neural Networks

An important generalization of the single layer perceptron architecture involves a concatenation of multiple layers of linear and nonlinear functions; see Fig. 3.2.3. In particular the outputs of each nonlinear layer become the inputs of the next linear layer. In some cases it may make sense to add as additional inputs some of the components of the state x or the state encoding $y(x)$.

There are a few questions to consider here. The first has to do with the reason for having multiple nonlinear layers, when a single one is sufficient to guarantee the universal approximation property. Here are some qualitative explanations:

- (a) If we view the outputs of each nonlinear layer as features, we see that the multilayer network provides a sequence of features, where each set of features in the sequence is a function of the preceding set of features in the sequence [except for the first set of features, which is a function of the encoding $y(x)$ of the state x]. Thus the network produces a hierarchy of features. In the context of specific applications, this hierarchical structure can be exploited in order to specialize the role of some of the layers and to enhance particular characteristics of the state.
- (b) Given the presence of multiple linear layers, one may consider the possibility of using matrices A with a particular sparsity pattern, or other structure that embodies special linear operations such as

convolution. When such structures are used, the training problem often becomes easier, because the number of parameters in the linear layers is drastically decreased.

Note that while in the early days of neural networks practitioners tended to use few nonlinear layers (say one to three), more recently a lot of success in certain problem domains (including image and speech processing, as well as approximate DP) has been achieved with so called *deep neural networks*, which involve a considerably larger number of layers. In particular, the use of deep neural networks, in conjunction with Monte Carlo tree search, has been an important factor for the success of the computer programs AlphaGo and AlphaZero, which perform better than the best humans in the games of Go and chess; see the papers by Silver et al. [SHM16], [SHS17]. By contrast, Tesauro's backgammon program and its descendants (cf. Section 2.4.2) do not require multiple nonlinear layers for good performance at present.

Training and Backpropagation

Let us now consider the training problem for multilayer networks. It has the form

$$\min_{v,r} \sum_{s=1}^q \left(\sum_{\ell=1}^m r_{\ell} \phi_{\ell}(x^s, v) - \beta^s \right)^2,$$

where v represents the collection of all the parameters of the linear layers, and $\phi_{\ell}(x, v)$ is the ℓ th feature produced at the output of the final nonlinear layer. Various types of incremental gradient methods, which modify the weight vector in the direction opposite to the gradient of a single sample term

$$\left(\sum_{\ell=1}^m r_{\ell} \phi_{\ell}(x^s, v) - \beta^s \right)^2,$$

can also be applied here. They are the methods most commonly used in practice. An important fact is that the gradient with respect to v of each feature $\phi_{\ell}(x, v)$ can be efficiently calculated using a special procedure known as *backpropagation*. This is just an intelligent way of applying the chain rule of differentiation, as we will explain now.

Multilayer perceptrons can be represented compactly by introducing certain mappings to describe the linear and the sigmoidal layers. In particular, let L_1, \dots, L_{m+1} denote the matrices representing the linear layers; that is, at the output of the 1st linear layer we obtain the vector $L_1 x$ and at the output of the k th linear layer ($k > 1$) we obtain $L_k \xi$, where ξ is the output of the preceding sigmoidal layer. Similarly, let $\Sigma_1, \dots, \Sigma_m$ denote the mappings representing the sigmoidal layers; that is, when the input of the k th sigmoidal layer ($k > 1$) is the vector y with components $y(j)$, we

obtain at the output the vector $\Sigma_k y$ with components $\sigma(y(j))$. The output of the multilayer perceptron is

$$F(L_1, \dots, L_{m+1}, x) = L_{m+1} \Sigma_m L_m \cdots \Sigma_1 L_1 x.$$

The special nature of this formula has an important computational consequence: the gradient (with respect to the weights) of the squared error between the output and a desired output y ,

$$E(L_1, \dots, L_{m+1}) = \frac{1}{2} (y - F(L_1, \dots, L_{m+1}, x))^2,$$

can be efficiently calculated using a special procedure known as *backpropagation*, which is just an intelligent way of using the chain rule.[†] In particular, the partial derivative of the cost function $E(L_1, \dots, L_{m+1})$ with respect to $L_k(i, j)$, the ij th component of the matrix L_k , is given by

$$\frac{\partial E(L_1, \dots, L_{m+1})}{\partial L_k(i, j)} = -e' L_{m+1} \bar{\Sigma}_m L_m \cdots L_{k+1} \bar{\Sigma}_k I_{ij} \Sigma_{k-1} L_{k-1} \cdots \Sigma_1 L_1 x, \quad (3.19)$$

where e is the error vector

$$e = y - F(L_1, \dots, L_{m+1}, x),$$

$\bar{\Sigma}_n$, $n = 1, \dots, m$, is the diagonal matrix with diagonal terms equal to the derivatives of the sigmoidal functions σ of the n th hidden layer evaluated at the appropriate points, and I_{ij} is the matrix obtained from L_k by setting all of its components to 0 except for the ij th component which is set to 1. This formula can be used to obtain efficiently all of the terms needed in the partial derivatives (3.19) of E using a two-step calculation:

- (a) Use a forward pass through the network to calculate sequentially the outputs of the linear layers

$$L_1 x, L_2 \Sigma_1 L_1 x, \dots, L_{m+1} \Sigma_m L_m \cdots \Sigma_1 L_1 x.$$

This is needed in order to obtain the points at which the derivatives in the matrices $\bar{\Sigma}_n$ are evaluated, and also in order to obtain the error vector $e = y - F(L_1, \dots, L_{m+1}, x)$.

- (b) Use a backward pass through the network to calculate sequentially the terms

$$e' L_{m+1} \bar{\Sigma}_m L_m \cdots L_{k+1} \bar{\Sigma}_k$$

[†] The name *backpropagation* is used in several different ways in the neural networks literature. For example feedforward neural networks of the type shown in Fig. 3.2.3 are sometimes referred to as *backpropagation networks*. The somewhat abstract derivation of the backpropagation formulas given here comes from Section 3.1.1 of the book [BeT96].

in the derivative formulas (3.19), starting with $e' L_{m+1} \bar{\Sigma}_m$, proceeding to $e' L_{m+1} \bar{\Sigma}_m L_m \bar{\Sigma}_{m-1}$, and continuing to $e' L_{m+1} \bar{\Sigma}_m \cdots L_2 \bar{\Sigma}_1$.

As a final remark, we mention that the ability to simultaneously extract features and optimize their linear combination is not unique to neural networks. Other approaches that use a multilayer architecture have been proposed (see the survey by Schmidhuber [Sch15]), and they admit similar training procedures based on appropriately modified forms of backpropagation. An example of an alternative multiplayer architecture approach is the Group Method for Data Handling (GMDH), which is principally based on the use of polynomial (rather than sigmoidal) nonlinearities. The GMDH approach was investigated extensively in the Soviet Union starting with the work of Ivakhnenko in the late 60s; see e.g., [Iva68]. It has been used in a large variety of applications, and its similarities with the neural network methodology have been noted (see the survey by Ivakhnenko [Iva71], and the large literature summary at the web site <http://www.gmdh.net>). Most of the GMDH research relates to inference-type problems, and we are unaware of any application to approximate DP. However, this may be a fruitful area of investigation, since in some applications it may turn out that polynomial nonlinearities are more suitable than sigmoidal or rectified linear unit nonlinearities.

3.3 SEQUENTIAL DYNAMIC PROGRAMMING APPROXIMATION

Let us describe a popular approach for training an approximation architecture $\tilde{J}_k(x_k, r_k)$ for a finite horizon DP problem. The parameter vectors r_k are determined sequentially, with an algorithm known as *fitted value iteration*, starting from the end of the horizon, and proceeding backwards as in the DP algorithm: first r_{N-1} then r_{N-2} , and so on. The algorithm samples the state space for each stage k , and generates a large number of states x_k^s , $s = 1, \dots, q$. It then determines sequentially the parameter vectors r_k to obtain a good “least squares fit” to the DP algorithm.

In particular, each r_k is determined by generating a large number of sample states and solving a least squares problem that aims to minimize the error in satisfying the DP equation for these states at time k . At the typical stage k , having obtained r_{k+1} , we determine r_k from the least squares problem

$$r_k \in \arg \min_r \sum_{s=1}^q \left(\tilde{J}_k(x_k^s, r) - \min_{u \in U_k(x_k^s)} E \left\{ g(x_k^s, u, w_k) + \tilde{J}_{k+1}(f_k(x_k^s, u, w_k), r_{k+1}) \right\} \right)^2$$

where x_k^s , $i = 1, \dots, q$, are the sample states that have been generated for the k th stage. Since r_{k+1} is assumed to be already known, the complicated minimization term in the right side of this equation is the known scalar

$$\beta_k^s = \min_{u \in U_k(x_k^s)} E \left\{ g(x_k^s, u, w_k) + \tilde{J}_{k+1}(f_k(x_k^s, u, w_k), r_{k+1}) \right\},$$

so that r_k is obtained as

$$r_k \in \arg \min_r \sum_{s=1}^q (\tilde{J}_k(x_k^s, r) - \beta_k^s)^2. \quad (3.20)$$

The algorithm starts at stage $N - 1$ with the minimization

$$r_{N-1} \in \arg \min_r \sum_{s=1}^q \left(\tilde{J}_{N-1}(x_{N-1}^s, r) - \min_{u \in U_{N-1}(x_{N-1}^s)} E \left\{ g_{N-1}(x_{N-1}^s, u, w_{N-1}) + g_N(f_{N-1}(x_{N-1}^s, u, w_{N-1})) \right\} \right)^2$$

and ends with the calculation of r_0 at time $k = 0$.

In the case of a linear architecture, where the approximate cost-to-go functions are

$$\tilde{J}_k(x_k, r_k) = r_k' \phi_k(x_k), \quad k = 0, \dots, N - 1,$$

the least squares problem (3.20) greatly simplifies, and admits the closed form solution

$$r_k = \left(\sum_{s=1}^q \phi_k(x_k^s) \phi_k(x_k^s)' \right)^{-1} \sum_{s=1}^q \beta_k^s \phi_k(x_k^s);$$

cf. Eq. (3.3). For a nonlinear architecture such as a neural network, incremental gradient algorithms may be used.

An important implementation issue is how to select the sample states x_k^s , $s = 1, \dots, q$, $k = 0, \dots, N - 1$. In practice, they are typically obtained by some form of Monte Carlo simulation, but the distribution by which they are generated is important for the success of the method. In particular, it is important that the sample states are “representative” in the sense that they are visited often under a nearly optimal policy. More precisely, the frequencies with which various states appear in the sample should be roughly proportional to the probabilities of their occurrence under an optimal policy. This point will be discussed later in Chapter 4, in the context of infinite horizon problems, and the notion of “representative” state will be better quantified in probabilistic terms.

Aside from the issue of selection of the sampling distribution that we have just described, a difficulty with fitted value iteration arises when the horizon is very long. In this case, however, the problem is often stationary, in the sense that the system and cost per stage do not change as time progresses. Then it may be possible to treat the problem as one with an infinite horizon and bring to bear additional methods for training approximation architectures; see the discussion in Chapter 4.

3.4 Q-FACTOR PARAMETRIC APPROXIMATION

We will now consider an alternative form of approximation in value space. It involves approximation of the optimal Q-factors of state-control pairs (x_k, u_k) at time k , without an intermediate approximation of cost-to-go functions. The optimal Q-factors are defined by

$$Q_k(x_k, u_k) = E \left\{ g_k(x_k, u_k, w_k) + J_{k+1}(f_k(x_k, u_k, w_k)) \right\}, \quad k = 0, \dots, N-1, \quad (3.21)$$

where J_{k+1} is the optimal cost-to-go function for stage $k+1$. Thus $Q_k(x_k, u_k)$ is the cost attained by using u_k at state x_k , and subsequently using an optimal policy.

As noted in Section 1.2, the DP algorithm can be written as

$$J_k(x_k) = \min_{u \in U_k(x_k)} Q_k(x_k, u), \quad (3.22)$$

and by using this equation, we can write Eq. (3.21) in the following equivalent form that relates Q_k with Q_{k+1} :

$$Q_k(x_k, u_k) = E \left\{ g_k(x_k, u_k, w_k) + \min_{u \in U_{k+1}(f_k(x_k, u_k, w_k))} Q_{k+1}(f_k(x_k, u_k, w_k), u) \right\}. \quad (3.23)$$

This suggests that in place of the Q-factors $Q_k(x_k, u_k)$, we may use Q-factor approximations and Eq. (3.23) as the basis for suboptimal control.

We can obtain such approximations by using methods that are similar to the ones we have considered so far (parametric approximation, enforced decomposition, certainty equivalent control, etc). Parametric Q-factor approximations $\tilde{Q}_k(x_k, u_k, r_k)$ often involve a feature vector that depends on the state, or depends on both the state and the control. In particular, a linear feature-based architecture may have the form

$$\tilde{Q}_k(x_k, u_k, r_k) = r_k(u_k)' \phi_k(x_k), \quad (3.24)$$

where $r_k(u_k)$ is a separate weight vector for each control u_k . Another type of linear feature-based architecture is

$$\tilde{Q}_k(x_k, u_k, r_k) = r_k' \phi_k(x_k, u_k), \quad (3.25)$$

where r_k is a weight vector that is independent of u_k . The architecture (3.24) is suitable for problems with a relatively small number of control options at each stage. In what follows, we will focus on the architecture (3.25), but the discussion with few modifications, also applies to the architecture (3.24).

We may adapt the fitted value iteration approach of the preceding section to compute sequentially the parameter vectors r_k in Q-factor parametric approximations, starting from $k = N - 1$. This algorithm is based on Eq. (3.23), with r_k obtained by solving least squares problems similar to the ones of Eq. (3.20). As an example, the parameters r_k of the architecture (3.25) are computed sequentially by collecting sample state-control pairs (x_k^s, u_k^s) , $s = 1, \dots, q$, and solving the linear least squares problems

$$r_k \in \arg \min_r \sum_{s=1}^q (r' \phi_k(x_k^s, u_k^s) - \beta_k^s)^2, \quad (3.26)$$

where

$$\beta_k^s = E \left\{ g_k(x_k^s, u_k^s, w_k) + \min_{u \in U_{k+1}(f_k(x_k^s, u_k^s, w_k))} r'_{k+1} \phi_k(f_k(x_k^s, u_k^s, w_k), u) \right\}. \quad (3.27)$$

Thus r_k is obtained through a least squares fit that aims to minimize the squared errors in satisfying Eq. (3.23). Note that the solution of the least squares problem (3.26) can be obtained in closed form as

$$r_k = \left(\sum_{s=1}^q \phi_k(x_k^s, u_k^s) \phi_k(x_k^s, u_k^s)' \right)^{-1} \sum_{s=1}^q \beta_k^s \phi_k(x_k^s, u_k^s);$$

[cf. Eq. (3.3)]. Once r_k has been computed, the one-step lookahead control $\tilde{\mu}_k(x_k)$ is obtained on-line as

$$\tilde{\mu}_k(x_k) \in \arg \min_{u \in U_k(x_k)} \tilde{Q}_k(x_k, u, r_k), \quad (3.28)$$

without the need to calculate any expected value. This latter property is a primary incentive for using Q-factors in approximate DP, particularly when there are tight constraints on the amount of on-line computation that is possible in the given practical setting.

Having obtained the one-step lookahead policy $\tilde{\pi} = \{\tilde{\mu}_0, \dots, \tilde{\mu}_{N-1}\}$, a further possibility is to approximate it with a parametric architecture. This is *approximation in policy space built on top of approximation in value space*; see the discussion of Section 2.1.3. The idea here is to simplify even further the on-line computation of the suboptimal controls by avoiding the minimization of Eq. (3.28).

Finally, let us note an alternative to computing Q-factor approximations, which is motivated by the potential benefit of using Q-factor

differences in contexts involving approximation. In this method, called *advantage updating*, instead of computing and comparing $Q_k(x_k, u_k)$ for all $u_k \in U_k(x_k)$, we compute

$$A_k(x_k, u_k) = Q_k(x_k, u_k) - \min_{u_k \in U_k(x_k)} Q_k(x_k, u_k).$$

The function $A_k(x_k, u_k)$ can serve just as well as $Q_k(x_k, u_k)$ for the purpose of comparing controls, but may have a much smaller range of values than $Q_k(x_k, u_k)$. In the absence of approximations, advantage updating is clearly equivalent to selecting controls by comparing their Q-factors. However, when approximation is involved, using Q-factor differences can be important, because $Q_k(x_k, u_k)$ may embody sizable quantities that are independent of u , and may interfere with algorithms such as the fitted value iteration (3.26)-(3.27). This question is discussed further and is illustrated with an example in the neuro-dynamic programming book [BeT96], Section 6.6.2.

3.5 NOTES AND SOURCES

Our discussion of approximation architectures, neural networks, and training has been limited, and aimed just to provide the connection with approximate DP and a starting point for further exploration. The literature on the subject is vast, and the textbooks mentioned in the references to Chapter 1 provide detailed accounts and many references in addition to the ones given in Sections 3.1.3 and 3.2.1.

There are two broad directions of inquiry in parametric architectures:

- (1) The design of architectures, either in a general or a problem-specific context. Research in this area has intensified following the increase in popularity of deep neural networks.
- (2) The algorithms for training of neural networks as well as linear architectures.

Both of these issues have been extensively investigated and research relating to these and other related issues is continuing.

Incremental algorithms are supported by substantial theoretical analysis, which addresses issues of convergence, rate of convergence, stepsize selection, and component order selection. It is beyond our scope to cover this analysis, and we refer to the book [BeT96] and paper [BeT00] by Bertsekas and Tsitsiklis, and the recent survey by Bottou, Curtis, and Nocedal [BCN18] for theoretically oriented accounts. The author's surveys [Ber10] and [Ber15b], and convex optimization and nonlinear programming textbooks [Ber15a], [Ber16a], collectively contain an extensive account of theoretical and practical issues regarding incremental methods, including the Kaczmarz, incremental gradient, incremental subgradient, incremental

aggregated gradient, incremental Newton, and incremental Gauss-Newton methods.

Fitted value iteration has interesting properties, and at times exhibits interesting behavior, such as instability over a long horizon. We will discuss this behavior in Section 4.4 in the context of infinite horizon problems. Advantage updating was proposed by Baird [Bai93], [Bai94].