

1

Introduction

This book is about the fundamentals of algorithms for solving *continuous optimization* problems, which involve minimizing functions of multiple real-valued variables, possibly subject to some restrictions or constraints on the values that those variables may take. We focus particularly (though not exclusively) on *convex* problems, and our choice of topics is motivated by relevance to data science. That is, the formulations and algorithms that we discuss are useful in solving problems from machine learning, statistics, and data analysis.

To set the stage for subsequent chapters, the rest of this chapter outlines several paradigms from data science and shows how they can be formulated as continuous optimization problems. We must pay attention to particular properties of these formulations – their smoothness properties and structure – when we choose algorithms to solve them.

1.1 Data Analysis and Optimization

The typical optimization problem in data analysis is to find a model that agrees with some collected data set but also adheres to some structural constraints that reflect our beliefs about what a good model should be. The data set in a typical analysis problem consists of m objects:

$$\mathcal{D} := \{(a_j, y_j), j = 1, 2, \dots, m\}, \quad (1.1)$$

where a_j is a vector (or matrix) of *features* and y_j is an *observation* or *label*. (We can assume that the data has been cleaned so that all pairs (a_j, y_j) , $j = 1, 2, \dots, m$ have the same size and shape.) The data analysis task then consists of discovering a function ϕ such that $\phi(a_j) \approx y_j$ for most $j = 1, 2, \dots, m$. The process of discovering the mapping ϕ is often called “learning” or “training.”

The function ϕ is often defined in terms of a vector or matrix of parameters, which we denote in what follows by x or X (and occasionally by other notation). With these parametrizations, the problem of identifying ϕ becomes a traditional data-fitting problem: *Find the parameters x defining ϕ such that $\phi(a_j) \approx y_j$, $j = 1, 2, \dots, m$ in some optimal sense.* Once we come up with a definition of the term “optimal” (and possibly also with restrictions on the values that we allow to parameters to take), we have an optimization problem. Frequently, these optimization formulations have objective functions of the finite-sum type

$$\mathcal{L}_{\mathcal{D}}(x) := \frac{1}{m} \sum_{j=1}^m \ell(a_j, y_j; x). \quad (1.2)$$

The function $\ell(a, y; x)$ here represents a “loss” incurred for not properly aligning our prediction $\phi(a)$ with y . Thus, the objective $\mathcal{L}_{\mathcal{D}}(x)$ measures the average loss accrued over the entire data set when the parameter vector is equal to x .

Once an appropriate value of x (and thus ϕ) has been learned from the data, we can use it to make predictions about other items of data not in the set \mathcal{D} (1.1). Given an unseen item of data \hat{a} of the same type as a_j , $j = 1, 2, \dots, m$, we predict the label \hat{y} associated with \hat{a} to be $\phi(\hat{a})$. The mapping ϕ may also expose other structures and properties in the data set. For example, it may reveal that only a small fraction of the features in a_j are needed to reliably predict the label y_j . (This is known as *feature selection*.) When the parameter x is a matrix, it could reveal a low-dimensional subspace that contains most of the vectors a_j , or it could reveal a matrix with particular structure (low-rank, sparse) such that observations of X prompted by the feature vectors a_j yield results close to y_j .

The form of the labels y_j differs according to the nature of the data analysis problem.

- If each y_j is a *real number*, we typically have a *regression* problem.
- When each y_j is a *label*, that is, an integer drawn from the set $\{1, 2, \dots, M\}$ indicating that a_j belongs to one of M classes, this is a *classification* problem. When $M = 2$, we have a binary classification problem, whereas $M > 2$ is multiclass classification. (In data analysis problems arising in speech and image recognition, M can be very large, of the order of thousands or more.)
- The labels y_j may not even exist; the data set may contain only the feature vectors a_j , $j = 1, 2, \dots, m$. There are still interesting data analysis problems associated with these cases. For example, we may wish to group

the a_j into clusters (where the vectors within each cluster are deemed to be functionally similar) or identify a low-dimensional subspace (or a collection of low-dimensional subspaces) that approximately contains the a_j . In such problems, we are essentially learning the labels y_j alongside the function ϕ . For example, in a clustering problem, y_j could represent the cluster to which a_j is assigned.

Even after cleaning and preparation, the preceding setup may contain many complications that need to be dealt with in formulating the problem in rigorous mathematical terms. The quantities (a_j, y_j) may contain noise or may be otherwise corrupted, and we would like the mapping ϕ to be robust to such errors. There may be *missing data*: Parts of the vectors a_j may be missing, or we may not know all the labels y_j . The data may be arriving in *streaming* fashion rather than being available all at once. In this case, we would learn ϕ in an *online* fashion.

One consideration that arises frequently is that we wish to avoid *overfitting* the model to the data set \mathcal{D} in (1.1). The particular data set \mathcal{D} available to us can often be thought of as a finite sample drawn from some underlying larger (perhaps infinite) collection of possible data points, and we wish the function ϕ to perform well on the unobserved data points as well as the observed subset \mathcal{D} . In other words, we want ϕ to be not too sensitive to the particular sample \mathcal{D} that is used to define empirical objective functions such as (1.2). One way to avoid this issue is to modify the objective function by adding constraints or penalty terms, in a way that limits the “complexity” of the function ϕ . This process is typically called *regularization*. An optimization formulation that balances fit to the training data \mathcal{D} , model complexity, and model structure is

$$\min_{x \in \Omega} \mathcal{L}_{\mathcal{D}}(x) + \lambda \text{pen}(x), \quad (1.3)$$

where Ω is a set of allowable values for x , $\text{pen}(\cdot)$ is a *regularization function* or *regularizer*, and $\lambda \geq 0$ is a *regularization parameter*. The regularizer usually takes lower values for parameters x that yield functions ϕ with lower complexity. (For example, ϕ may depend on fewer of the features in the data vectors a_j or may be less oscillatory.) The parameter λ can be “tuned” to provide an appropriate balance between fitting the data and lowering the complexity of ϕ : Smaller values of λ tend to produce solutions that fit the training data \mathcal{D} more accurately, while large values of λ lead to less complex models.¹

¹ Interestingly, the concept of overfitting has been reexamined in recent years, particularly in the context of deep learning, where models that perfectly fit the training data are sometimes observed to also do a good job of classifying previously unseen data. This phenomenon is a topic of intense current research in the machine learning community.

The constraint set Ω in (1.3) may be chosen to exclude values of x that are not relevant or useful in the context of the data analysis problem. For example, in some applications, we may not wish to consider values of x in which one or more components are negative, so we could set Ω to be the set of vectors whose components are all greater than or equal to zero.

We now examine some particular problems in data science that give rise to formulations that are special cases of our master problem (1.3). We will see that a large variety of problems can be formulated using this general framework, but we will also see that within this framework, there is a wide range of structures that must be taken into account in choosing algorithms to solve these problems efficiently.

1.2 Least Squares

Probably the oldest and best-known data analysis problem is linear least squares. Here, the data points (a_j, y_j) lie in $\mathbb{R}^n \times \mathbb{R}$, and we solve

$$\min_x \frac{1}{2m} \sum_{j=1}^m (a_j^T x - y_j)^2 = \frac{1}{2m} \|Ax - y\|_2^2, \quad (1.4)$$

where A the matrix whose rows are a_j^T , $j = 1, 2, \dots, m$ and $y = (y_1, y_2, \dots, y_m)^T$. In the preceding terminology, the function ϕ is defined by $\phi(a) := a^T x$. (We can introduce a nonzero intercept by adding an extra parameter $\beta \in \mathbb{R}$ and defining $\phi(a) := a^T x + \beta$.) This formulation can be motivated statistically, as a maximum-likelihood estimate of x when the observations y_j are exact but for independent identically distributed (i.i.d.) Gaussian noise. We can add a variety of penalty functions to this basic least squares problem to impose desirable structure on x and, hence, on ϕ . For example, *ridge regression* adds a squared ℓ_2 -norm penalty, resulting in

$$\min_x \frac{1}{2m} \|Ax - y\|_2^2 + \lambda \|x\|_2^2, \quad \text{for some parameter } \lambda > 0.$$

The solution x of this regularized formulation has less sensitivity to perturbations in the data (a_j, y_j) . The LASSO formulation

$$\min_x \frac{1}{2m} \|Ax - y\|_2^2 + \lambda \|x\|_1 \quad (1.5)$$

tends to yield solutions x that are sparse – that is, containing relatively few nonzero components (Tibshirani, 1996). This formulation performs feature selection: The locations of the nonzero components in x reveal those

components of a_j that are instrumental in determining the observation y_j . Besides its statistical appeal – predictors that depend on few features are potentially simpler and more comprehensible than those depending on many features – feature selection has practical appeal in making predictions about future data. Rather than gathering all components of a new data vector \hat{a} , we need to find only the “selected” features because only these are needed to make a prediction.

The LASSO formulation (1.5) is an important prototype for many problems in data analysis in that it involves a regularization term $\lambda\|x\|_1$ that is non-smooth and convex but has relatively simple structure that can potentially be exploited by algorithms.

1.3 Matrix Factorization Problems

There are a variety of data analysis problems that require estimating a low-rank matrix from some sparse collection of data. Such problems can be formulated as natural extension of least squares to problems in which the data a_j are naturally represented as matrices rather than vectors.

Changing notation slightly, we suppose that each A_j is an $n \times p$ matrix, and we seek another $n \times p$ matrix X that solves

$$\min_X \frac{1}{2m} \sum_{j=1}^m (\langle A_j, X \rangle - y_j)^2, \quad (1.6)$$

where $\langle A, B \rangle := \text{trace}(A^T B)$. Here we can think of the A_j as “probing” the unknown matrix X . Commonly considered types of observations are random linear combinations (where the elements of A_j are selected i.i.d. from some distribution) or single-element observations (in which each A_j has 1 in a single location and zeros elsewhere). A regularized version of (1.6), leading to solutions X that are low rank, is

$$\min_X \frac{1}{2m} \sum_{j=1}^m (\langle A_j, X \rangle - y_j)^2 + \lambda \|X\|_*, \quad (1.7)$$

where $\|X\|_*$ is the nuclear norm, which is the sum of singular values of X (Recht et al., 2010). The nuclear norm plays a role analogous to the ℓ_1 norm in (1.5), where as the ℓ_1 norm favors sparse vectors, the nuclear norm favors low-rank matrices. Although the nuclear norm is a somewhat complex nonsmooth function, it is at least convex so that the formulation (1.7) is also convex. This formulation can be shown to yield a statistically valid solution when the true

X is low rank and the observation matrices A_j satisfy a “restricted isometry property,” commonly satisfied by random matrices but not by matrices with just one nonzero element. The formulation is also valid in a different context, in which the true X is incoherent (roughly speaking, it does not have a few elements that are much larger than the others), and the observations A_j are of single elements (Candès and Recht, 2009).

In another form of regularization, the matrix X is represented explicitly as a product of two “thin” matrices L and R , where $L \in \mathbb{R}^{n \times r}$ and $R \in \mathbb{R}^{p \times r}$, with $r \ll \min(n, p)$. We set $X = LR^T$ in (1.6) and solve

$$\min_{L, R} \frac{1}{2m} \sum_{j=1}^m (\langle A_j, LR^T \rangle - y_j)^2. \quad (1.8)$$

In this formulation, the rank r is “hard-wired” into the definition of X , so there is no need to include a regularizing term. This formulation is also typically much more compact than (1.7); the total number of elements in (L, R) is $(n + p)r$, which is much less than np . However, this function is nonconvex when considered as a function of (L, R) jointly. An active line of current research, pioneered by Burer and Monteiro (2003) and also drawing on statistical sources, shows that the nonconvexity is benign in many situations and that, under certain assumptions on the data (A_j, y_j) , $j = 1, 2, \dots, m$ and careful choice of algorithmic strategy, good solutions can be obtained from the formulation (1.8). A clue to this good behavior is that although this formulation is nonconvex, it is in some sense an approximation to a tractable problem: If we have a complete observation of X , then a rank- r approximation can be found by performing a singular value decomposition of X and defining L and R in terms of the r leading left and right singular vectors.

Some applications in computer vision, chemometrics, and document clustering require us to find factors L and R like those in (1.8) in which all elements are nonnegative. If the full matrix $Y \in \mathbb{R}^{n \times p}$ is observed, this problem has the form

$$\min_{L, R} \|LR^T - Y\|_F^2, \quad \text{subject to } L \geq 0, R \geq 0$$

and is called *nonnegative matrix factorization*.

1.4 Support Vector Machines

Classification via support vector machines (SVM) is a classical optimization problem in machine learning, tracing its origins to the 1960s. Given the input

data (a_j, y_j) with $a_j \in \mathbb{R}^n$ and $y_j \in \{-1, 1\}$, SVM seeks a vector $x \in \mathbb{R}^n$ and a scalar $\beta \in \mathbb{R}$ such that

$$a_j^T x - \beta \geq 1 \quad \text{when } y_j = +1, \quad (1.9a)$$

$$a_j^T x - \beta \leq -1 \quad \text{when } y_j = -1. \quad (1.9b)$$

Any pair (x, β) that satisfies these conditions defines a *separating hyperplane* in \mathbb{R}^n , that separates the “positive” cases $\{a_j \mid y_j = +1\}$ from the “negative” cases $\{a_j \mid y_j = -1\}$. Among all separating hyperplanes, the one that minimizes $\|x\|^2$ is the one that maximizes the *margin* between the two classes – that is, the hyperplane whose distance to the nearest point a_j of either class is greatest.

We can formulate the problem of finding a separating hyperplane as an optimization problem by defining an objective with the summation form (1.2):

$$H(x, \beta) = \frac{1}{m} \sum_{j=1}^m \max(1 - y_j(a_j^T x - \beta), 0). \quad (1.10)$$

Note that the j th term in this summation is zero if the conditions (1.9) are satisfied, and it is positive otherwise. Even if no pair (x, β) exists for which $H(x, \beta) = 0$, a value (x, β) that minimizes (1.2) will be the one that comes as close as possible to satisfying (1.9) in some sense. A term $\lambda \|x\|_2^2$ (for some parameter $\lambda > 0$) is often added to (1.10), yielding the following regularized version:

$$H(x, \beta) = \frac{1}{m} \sum_{j=1}^m \max(1 - y_j(a_j^T x - \beta), 0) + \frac{1}{2} \lambda \|x\|_2^2. \quad (1.11)$$

Note that, in contrast to the examples presented so far, the SVM problem has a nonsmooth loss function and a smooth regularizer.

If λ is sufficiently small, and if separating hyperplanes exist, the pair (x, β) that minimizes (1.11) is the maximum-margin separating hyperplane. The maximum-margin property is consistent with the goals of generalizability and robustness. For example, if the observed data (a_j, y_j) is drawn from an underlying “cloud” of positive and negative cases, the maximum-margin solution usually does a reasonable job of separating other empirical data samples drawn from the same clouds, whereas a hyperplane that passes close to several of the observed data points may not do as well (see Figure 1.1).

Often, it is not possible to find a hyperplane that separates the positive and negative cases well enough to be useful as a classifier. One solution is to transform all of the raw data vectors a_j by some nonlinear mapping ψ and

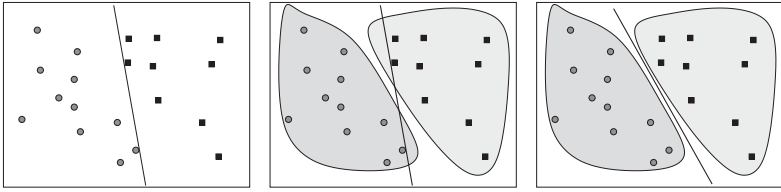


Figure 1.1 Linear support vector machine classification, with the one class represented by circles and the other by squares. One possible choice of separating hyperplane is shown at left. If the training data is an empirical sample drawn from a cloud of underlying data points, this plane does not do well in separating the two clouds (middle). The maximum-margin separating hyperplane does better (right).

then perform the support vector machine classification on the vectors $\psi(a_j)$, $j = 1, 2, \dots, m$. The conditions (1.9) would thus be replaced by

$$\psi(a_j)^T x - \beta \geq 1 \quad \text{when } y_j = +1; \quad (1.12a)$$

$$\psi(a_j)^T x - \beta \leq -1 \quad \text{when } y_j = -1, \quad (1.12b)$$

leading to the following analog of (1.11):

$$H(x, \beta) = \frac{1}{m} \sum_{j=1}^m \max(1 - y_j(\psi(a_j)^T x - \beta), 0) + \frac{1}{2} \lambda \|x\|_2^2. \quad (1.13)$$

When transformed back to \mathbb{R}^m , the surface $\{a \mid \psi(a)^T x - \beta = 0\}$ is nonlinear and possibly disconnected, and is often a much more powerful classifier than the hyperplanes resulting from (1.11).

We note that SVM can also be expressed naturally as a minimization problem over a convex set. By introducing artificial variables, the problem (1.13) (and (1.11)) can be formulated as a convex quadratic program – that is, a problem with a convex quadratic objective and linear constraints. By taking the dual of this problem, we obtain another convex quadratic program, in m variables:

$$\min_{\alpha \in \mathbb{R}^m} \frac{1}{2} \alpha^T Q \alpha - \mathbf{1}^T \alpha \quad \text{subject to } 0 \leq \alpha \leq \frac{1}{\lambda} \mathbf{1}, \quad y^T \alpha = 0, \quad (1.14)$$

where

$$Q_{kl} = y_k y_l \psi(a_k)^T \psi(a_l), \quad y = (y_1, y_2, \dots, y_m)^T, \quad \mathbf{1} = (1, 1, \dots, 1)^T.$$

Interestingly, problem (1.14) can be formulated and solved without explicit knowledge or definition of the mapping ψ . We need only a technique to define the elements of Q . This can be done with the use of a *kernel function* $K: \mathbb{R}^n \times \mathbb{R}^n \rightarrow \mathbb{R}$, where $K(a_k, a_l)$ replaces $\psi(a_k)^T \psi(a_l)$ (Boser et al., 1992; Cortes

and Vapnik, 1995). This is the so-called kernel trick. (The kernel function K can also be used to construct a classification function ϕ from the solution of (1.14).) A particularly popular choice of kernel is the Gaussian kernel:

$$K(a_k, a_l) := \exp\left(-\frac{1}{2\sigma} \|a_k - a_l\|^2\right),$$

where σ is a positive parameter.

1.5 Logistic Regression

Logistic regression can be viewed as a softened form of binary support vector machine classification in which, rather than the classification function ϕ giving a unqualified prediction of the class in which a new data vector a lies, it returns an estimate of the *odds* of a belonging to one class or the other. We seek an “odds function” p parametrized by a vector $x \in \mathbb{R}^n$,

$$p(a; x) := (1 + \exp(a^T x))^{-1}, \quad (1.15)$$

and aim to choose the parameter x in so that

$$p(a_j; x) \approx 1 \quad \text{when } y_j = +1; \quad (1.16a)$$

$$p(a_j; x) \approx 0 \quad \text{when } y_j = -1. \quad (1.16b)$$

(Note the similarity to (1.9).) The optimal value of x can be found by minimizing a negative-log-likelihood function:

$$L(x) := -\frac{1}{m} \left[\sum_{j:y_j=-1} \log(1 - p(a_j; x)) + \sum_{j:y_j=1} \log p(a_j; x) \right]. \quad (1.17)$$

Note that the definition (1.15) ensures that $p(a; x) \in (0, 1)$ for all a and x ; thus, $\log(1 - p(a_j; x)) < 0$ and $\log p(a_j; x) < 0$ for all j and all x . When the conditions (1.16) are satisfied, these log terms will be only *slightly* negative, so values of x that satisfy (1.17) will be near optimal.

We can perform feature selection using the model (1.17) by introducing a regularizer $\lambda \|x\|_1$ (as in the LASSO technique for least squares (1.5)),

$$\min_x -\frac{1}{m} \left[\sum_{j:y_j=-1} \log(1 - p(a_j; x)) + \sum_{j:y_j=1} \log p(a_j; x) \right] + \lambda \|x\|_1, \quad (1.18)$$

where $\lambda > 0$ is a regularization parameter. As we see later, this term has the effect of producing a solution in which few components of x are nonzero,

making it possible to evaluate $p(a; x)$ by knowing only those components of a that correspond to the nonzeros in x .

An important extension of this technique is to *multiclass* (or *multinomial*) logistic regression, in which the data vectors a_j belong to more than two classes. Such applications are common in modern data analysis. For example, in a speech recognition system, the M classes could each represent a *phoneme* of speech, one of the potentially thousands of distinct elementary sounds that can be uttered by humans in a few tens of milliseconds. A multinomial logistic regression problem requires a distinct odds function p_k for each class $k \in \{1, 2, \dots, M\}$. These functions are parametrized by vectors $x_{[k]} \in \mathbb{R}^n$, $k = 1, 2, \dots, M$, defined as follows:

$$p_k(a; X) := \frac{\exp(a^T x_{[k]})}{\sum_{l=1}^M \exp(a^T x_{[l]})}, \quad k = 1, 2, \dots, M, \quad (1.19)$$

where we define $X := \{x_{[k]} \mid k = 1, 2, \dots, M\}$. As in the binary case, we have $p_k(a) \in (0, 1)$ for all a and all $k = 1, 2, \dots, M$ and, in addition, that $\sum_{k=1}^M p_k(a) = 1$. The functions (1.19) perform a “softmax” on the quantities $\{a^T x_{[l]} \mid l = 1, 2, \dots, M\}$.

In the setting of multiclass logistic regression, the labels y_j are vectors in \mathbb{R}^M whose elements are defined as follows:

$$y_{jk} = \begin{cases} 1 & \text{when } a_j \text{ belongs to class } k, \\ 0 & \text{otherwise.} \end{cases} \quad (1.20)$$

Similarly to (1.16), we seek to define the vectors $x_{[k]}$ so that

$$p_k(a_j; X) \approx 1 \quad \text{when } y_{jk} = 1 \quad (1.21a)$$

$$p_k(a_j; X) \approx 0 \quad \text{when } y_{jk} = 0. \quad (1.21b)$$

The problem of finding values of $x_{[k]}$ that satisfy these conditions can again be formulated as one of minimizing a negative-log-likelihood:

$$L(X) := -\frac{1}{m} \sum_{j=1}^m \left[\sum_{\ell=1}^M y_{j\ell} (x_{[\ell]}^T a_j) - \log \left(\sum_{\ell=1}^M \exp(x_{[\ell]}^T a_j) \right) \right]. \quad (1.22)$$

“Group-sparse” regularization terms can be included in this formulation to select a set of features in the vectors a_j , common to each class, that distinguish effectively between the classes.

1.6 Deep Learning

Deep neural networks are often designed to perform the same function as multiclass logistic regression – that is, to classify a data vector a into one of M possible classes, often for large M . The major innovation is that the mapping ϕ from data vector to prediction is now a nonlinear function, explicitly parametrized by a set of structured transformations.

The neural network shown in Figure 1.2 illustrates the structure of a particular neural net. In this figure, the data vector a_j enters at the left of the network, and each box (more often referred to as a “layer”) represents a transformation that takes an input vector and applies a nonlinear transformation of the data to produce an output vector. The output of each operator becomes the input for one or more subsequent layers. Each layer has a set of its own parameters, and the collection of all of the parameters over all the layers comprises our optimization variable. The different shades of boxes here denote the fact that the types of transformations might differ between layers, but we can compose them in whatever fashion suits our application.

A typical transformation, which converts the vector a_j^{l-1} representing output from layer $l - 1$ to the vector a_j^l representing output from layer l , is

$$a_j^l = \sigma(W^l a_j^{l-1} + g^l), \quad (1.23)$$

where W^l is a matrix of dimension $|a_j^l| \times |a_j^{l-1}|$ and g^l is a vector of length $|a_j^l|$. The function σ is a *componentwise* nonlinear transformation, usually called an *activation function*. The most common forms of the activation function σ act independently on each component of their argument vector as follows:

- Sigmoid: $t \rightarrow 1/(1 + e^{-t})$;
- Rectified Linear Unit (ReLU): $t \rightarrow \max(t, 0)$.

Alternative transformations are needed when the input to box l comes from two or more preceding boxes (as in the case for some boxes in Figure 1.2).

The rightmost layer of the neural network (the output layer) typically has M outputs, one for each of the possible classes to which the input (a_j , say) could belong. These are compared to the labels y_{jk} , defined as in (1.20) to indicate which of the M classes that a_j belongs to. Often, a softmax is applied to the



Figure 1.2 A deep neural network, showing connections between adjacent layers, where each layer is represented by a shaded rectangle.

outputs in the rightmost layer, and a loss function similar to (1.22) is obtained, as we describe now.

Consider the special (but not uncommon) case in which the neural net structure is a linear graph of D levels, in which the output for layer $l - 1$ becomes the input for layer l (for $l = 1, 2, \dots, D$) with $a_j = a_j^0$, $j = 1, 2, \dots, m$, and the transformation within each box has the form (1.23). A softmax is applied to the output of the rightmost layer to obtain a set of odds. The parameters in this neural network are the matrix-vector pairs (W^l, g^l) , $l = 1, 2, \dots, D$ that transform the input vector $a_j = a_j^0$ into the output a_j^D of the final layer. We aim to choose all these parameters so that the network does a good job of classifying the training data correctly. Using the notation w for the layer-to-layer transformations, that is,

$$w := (W^1, g^1, W^2, g^2, \dots, W^D, g^D),$$

we can write the loss function for deep learning as

$$L(w) = -\frac{1}{m} \sum_{j=1}^m \left[\sum_{\ell=1}^M y_{j\ell} a_{j,\ell}^D(w) - \log \left(\sum_{\ell=1}^M \exp a_{j,\ell}^D(w) \right) \right], \quad (1.24)$$

where $a_{j,\ell}^D(w) \in \mathbb{R}$ is the output of the ℓ th element in layer D corresponding to input vector a_j^0 . (Here we write $a_{j,\ell}^D(w)$ to make explicit the dependence on the transformations w as well as on the input vector a_j .) We can view multiclass logistic regression as a special case of deep learning with $D = 1$, so that $a_{j,\ell}^1 = W_{\ell,\cdot}^1 a_j^0$, where $W_{\ell,\cdot}^1$ denotes row ℓ of the matrix W^1 .

Neural networks in use for particular applications (for example, in image recognition and speech recognition, where they have been quite successful) include many variants on the basic design. These include restricted connectivity between the boxes (which corresponds to enforcing sparsity structure on the matrices W^l , $l = 1, 2, \dots, D$) and sharing parameters, which corresponds to forcing subsets of the elements of W^l to take the same value. Arrangements of the boxes may be quite complex, with outputs coming from several layers, connections across nonadjacent layers, different componentwise transformations σ at different layers, and so on. Deep neural networks for practical applications are highly engineered objects.

The loss function (1.24) shares with many other applications the finite-sum form (1.2), but it has several features that set it apart from the other applications discussed before. First, and possibly most important, it is *nonconvex* in the parameters w . Second, the total number of parameters in w is usually very large. Effective training of deep learning classifiers typically requires a great deal of data and computation power. Huge clusters of powerful computers –

often using multicore processors, GPUs, and even specially architected processing units – are devoted to this task.

1.7 Emphasis

Many problems can be formulated as in the framework (1.3), and their properties may differ significantly. They might be convex or nonconvex, and smooth or nonsmooth. But there are important features that they all share.

- They can be formulated as functions of *real variables*, which we typically arrange in a vector of length n .
- The functions are continuous. When nonsmoothness appears in the formulation, it does so in a structured way that can be exploited by the algorithm. Smoothness properties allow an algorithm to make good inferences about the behavior of the function on the basis of knowledge gained at nearby points that have been visited previously.
- The objective is often made up in part of a summation of many terms, where each term depends on a single item of data.
- The objective is often a sum of two terms: a “loss term” (sometimes arising from a maximum likelihood expression for some statistical model) and a “regularization term” whose purpose is to impose structure and “generalizability” on the recovered model.

Our treatment emphasizes algorithms for solving these various kinds of problems, with analysis of the convergence properties of these algorithms. We pay attention to complexity guarantees, which are bounds on the amount of computational effort required to obtain solutions of a given accuracy. These bounds usually depend on fundamental properties of the objective function and the data that defines it, including the dimensions of the data set and the number of variables in the problem. This emphasis contrasts with much of the optimization literature, in which global convergence results do not usually involve complexity bounds. (A notable exception is the analysis of interior-point methods (see Nesterov and Nemirovskii, 1994; Wright, 1997)).

At the same time, we try as much as possible to emphasize the practical concerns associated with solving these problems. There are a variety of trade-offs presented by any problem, and the optimizer has to evaluate which tools are most appropriate to use. On top of the problem formulation, it is imperative to account for the time budget for the task at hand, the type of computer on which the problem will be solved, and the guarantees needed for the

solution to be useful in the application that gave rise to the problem. Worst-case complexity guarantees are only a piece of the story here, and understanding the various parameters and heuristics that form part of any practical algorithmic strategy are critical for building reliable solvers.

Notes and References

The softmax operator is ubiquitous in problems involving multiple classes. Given real numbers z_1, z_2, \dots, z_M , we define $p_j = e^{z_j} / \sum_{i=1}^M e^{z_i}$ and note that $p_j \in (0, 1)$ for all j , and $\sum_{j=1}^M p_j = 1$. Moreover, if for some j we have $z_j \gg \max_{i \neq j} z_i$, then $p_j \approx 1$ while $p_i \approx 0$ for all $i \neq j$.

The examples in this chapter are adapted from an article by one of the authors (Wright, 2018).