

Continuous Optimization

3667 Recall from Section 1.1.4 that training a machine learning model often
 3668 boils down to finding a good set of parameters. The notion of “good” is
 3669 determined by the objective function or the probabilistic model, which we
 3670 will see examples of in the second part of this book. Given an objective
 3671 function finding the best value is done using optimization algorithms. We
 3672 will assume in this chapter that our objective function is differentiable
 3673 (see Chapter 5), hence we have access to a gradient at each location in
 3674 the space to help us find the optimum value. By convention most objective
 3675 functions in machine learning are intended to be minimized, that is the
 3676 best value is the minimum value. Intuitively finding the best value is like
 3677 finding the valleys of the objective function, and the gradients point us
 3678 uphill. The idea is to move downhill (opposite to the gradient) and hope
 3679 to find the deepest point.

Consider the function in Figure 7.1. The function has a *global minimum*
 around the value $x = -4.5$ which has the objective function value of

Since we consider data and models in \mathbb{R}^D the optimization problems we face are *continuous* optimization problems, as opposed to *combinatorial* optimization problems for discrete variables. global minimum

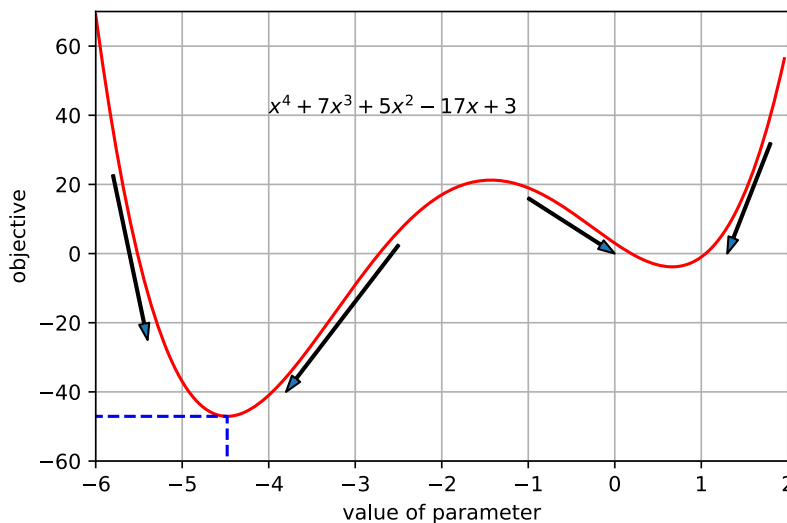


Figure 7.1 Example objective function. Gradients are indicated by arrows, and the global minimum is indicated by the dashed blue line.

local minimum
Stationary points
are points that have
zero gradient.

around -47 . Since the function is “smooth” the gradients can be used to help find the minimum by indicating whether we should take a step to the right or left. This assumes that we are in the correct bowl, as there exists another *local minimum* around the value $x = 0.7$. Recall that we can solve for all the stationary points of a function by calculating its derivative and setting it to zero. Let

$$\ell(x) = x^4 + 7x^3 + 5x^2 - 17x + 3. \quad (7.1)$$

Its gradient is given by

$$\frac{d\ell(x)}{dx} = 4x^3 + 21x^2 + 10x - 17. \quad (7.2)$$

Since this is a cubic equation, it has three solutions when set to zero. Two of them are minima and one is a maximum (around $x = -1.4$). Recall that to check whether a stationary point is a minimum or maximum we need to take the derivative a second time and check whether the second derivative is positive or negative at the stationary point.

$$\frac{d^2\ell(x)}{dx^2} = 12x^2 + 42x + 10 \quad (7.3)$$

3680 By substituting our visually estimated values of $x = -4.5, -1.4, 0.7$ we
3681 will observe that as expected the middle point is a maximum ($\frac{d^2\ell(x)}{dx^2} < 0$)
3682 and the other two stationary points are minimums.

In fact according to
the Abel-Ruffini
theorem, also
known as Abel's
impossibility
theorem, there is
general no algebraic
solution for
polynomials of
degree 5 or more.

3683 Note that we have avoided analytically solving for values of x in the pre-
3684 vious discussion, although for low order polynomials such as the above we
3685 could. In general, we are unable to find analytic solutions, and hence we
3686 need to start at some value, say $x_0 = -10$ and follow the gradient. The
3687 gradient indicates that we should go right, but not how far (this is called
3688 the step size). Furthermore, if we had started at the right side (e.g. $x_0 = 0$)
3689 the gradient would have led us to the wrong minimum. Figure 7.1 illus-
3690 trates the fact that for $x > -1$, the gradient points towards the minimum
3691 on the right of the figure, which has a larger objective value.

convex functions

3692 We will see in Section 7.3 a class of functions called convex functions
3693 that do not exhibit this tricky dependency on the starting point of the
3694 optimization algorithm. For *convex functions* all local minima are global
3695 minimum. It turns out that many machine learning objective functions
3696 are designed such that they are convex, and we will see an example in
3697 Chapter 10.

3698 The discussion in this Chapter so far was about a one dimensional func-
3699 tion, where we are able to visualize the ideas of gradients, descent direc-
3700 tions and optimal values. In the rest of this chapter we develop the same
3701 ideas in high dimensions. Unfortunately we can only visualize the con-
3702 cepts in one dimension, but some concepts do not generalize directly to
3703 higher dimensions, therefore some care needs to be taken when reading.

7.1 Optimization using Gradient Descent

We now consider the problem of solving for the minimum of a real-valued function

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad (7.4)$$

where $f : \mathbb{R}^d \rightarrow \mathbb{R}$ is an objective function that captures the machine learning problem at hand. We assume that our function f is differentiable, and we are unable to analytically find a solution in closed form.

Gradient descent is a first-order optimization algorithm. To find a local minimum of a function using gradient descent, one takes steps proportional to the negative of the gradient of the function at the current point. Recall from Chapter 5 that the gradient points in the direction of the steepest ascent and it is orthogonal to the contour lines of the function we wish to optimize.

Consider the quadratic function $f(x) = (x - 1)^2$ in Figure 7.2. We begin our search for the minimum value at $x_0 = -0.7$, where the subscript 0 indicates the initial estimate. Observe that the gradient $\left. \frac{d}{dx} f(x) \right|_{x=-0.7}$ points toward the right hand side. This means that we should take a step to the right (increasing x), which makes intuitive sense. However it is unclear how big a step to take, and the example algorithm overshoots the (unknown) minimum to $x_1 = 2.5$. Now the gradient points to the left, and we take a step to the left, and so on.

Let us consider multivariate functions. Imagine a surface (described by the function $f(\mathbf{x})$) with a ball starting at a particular location \mathbf{x}_0 . When the ball is released, it will move downhill in the direction of steepest descent. Gradient descent exploits the fact that $f(\mathbf{x}_0)$ decreases fastest if one moves from \mathbf{x}_0 in the direction of the negative gradient $-(\nabla f)(\mathbf{x}_0)^\top$ of f at \mathbf{x}_0 . We assume in this book that the functions are differentiable, and refer the reader to more general settings in Section 7.4. Then, if

$$\mathbf{x}_1 = \mathbf{x}_0 - \gamma((\nabla f)(\mathbf{x}_0))^\top \quad (7.5)$$

for a small *step size* $\gamma \geq 0$ then $f(\mathbf{x}_1) \leq f(\mathbf{x}_0)$. Note that we use the transpose for the gradient since otherwise the dimensions will not work out.

This observation allows us to define a simple gradient-descent algorithm: If we want to find a local optimum $f(\mathbf{x}_*)$ of a function $f : \mathbb{R}^n \rightarrow \mathbb{R}$, $\mathbf{x} \mapsto f(\mathbf{x})$, we start with an initial guess \mathbf{x}_0 of the parameters we wish to optimize and then iterate according to

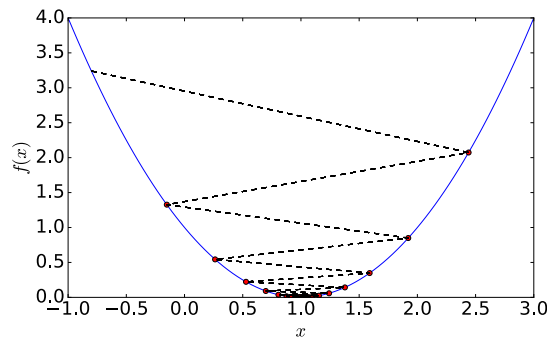
$$\mathbf{x}_{i+1} = \mathbf{x}_i - \gamma_i((\nabla f)(\mathbf{x}_i))^\top. \quad (7.6)$$

For suitable step size γ_i , the sequence $f(\mathbf{x}_0) \geq f(\mathbf{x}_1) \geq \dots$ converges to a local minimum.

Remark. Gradient descent can be relatively slow close to the minimum:

We use the convention of row vectors for gradients.

Figure 7.2 Gradient descent can lead to zigzagging and slow convergence.



3728 Its asymptotic rate of convergence is inferior to many other methods. Us-
 3729 ing the ball rolling down the hill analogy, when the surface is a long thin
 3730 valley the problem is poorly conditioned. For poorly conditioned convex
 3731 problems, gradient descent increasingly ‘zigzags’ as the gradients point
 3732 nearly orthogonally to the shortest direction to a minimum point, see
 3733 Fig. 7.2. \diamond

7.1.1 Step size

3734
 3735 The step size is also
 3736 called the learning
 3737 rate

3735 As mentioned earlier, choosing a good step size is important in gradient
 3736 descent. If the step size is too small, gradient descent can be slow. If the
 3737 step size is chosen too large, gradient descent can overshoot, fail to con-
 3738 verge, or even diverge. We will discuss the use of momentum in the next
 3739 section. It is a method that smoothes out erratic behavior of gradient up-
 3740 dates and dampens oscillations.

3741 Adaptive gradient methods rescale the step size at each iteration, de-
 3742 pending on local properties of the function. There are two simple heuris-
 3743 tics (Toussaint, 2012):

- 3744 • When the function value increases after a gradient step, the step size
 3745 was too large. Undo the step and decrease the step size.
- 3746 • When the function value decreases the step could have been larger. Try
 3747 to increase the step size.

3748 Although the “undo” step seems to be a waste of resources, using this
 3749 heuristic guarantees monotonic convergence.

Example (Solving a Linear Equation System)

When we solve linear equations of the form $\mathbf{Ax} = \mathbf{b}$, in practice we solve $\mathbf{Ax} - \mathbf{b} = \mathbf{0}$ approximately by finding \mathbf{x}_* that minimizes the squared error

$$\|\mathbf{Ax} - \mathbf{b}\|^2 = (\mathbf{Ax} - \mathbf{b})^\top (\mathbf{Ax} - \mathbf{b}) \quad (7.7)$$

if we use the Euclidean norm. The gradient of (7.7) with respect to \mathbf{x} is

$$\nabla_{\mathbf{x}} = 2(\mathbf{A}\mathbf{x} - \mathbf{b})^\top \mathbf{A}. \quad (7.8)$$

3750 We can use this gradient directly in a gradient descent algorithm. How-
 3751 ever for this particular special case, it turns out that there is an analytic
 3752 solution, which can be found by setting the gradient to zero. We can see
 3753 that this analytic solution is given by $\mathbf{A}\mathbf{x} = \mathbf{b}$. We will see more on solving
 3754 squared error problems in Chapter 9.

3755

3756

3757 *Remark.* When applied to the solution of linear systems of equations $\mathbf{A}\mathbf{x} =$
 3758 \mathbf{b} gradient descent may converge slowly. The speed of convergence of gra-
 3759 dient descent is dependent on the condition number $\kappa = \frac{\sigma(\mathbf{A})_{\max}}{\sigma(\mathbf{A})_{\min}}$, which is
 3760 the ratio of the maximum to the minimum singular value of \mathbf{A} . The con-
 3761 dition number essentially measures the ratio of the most curved direction
 3762 versus the least curved direction, which corresponds to our imagery that
 3763 poorly conditioned problems are long thin valleys: they are very curved in
 3764 one direction, but very flat in the other. Instead of directly solving $\mathbf{A}\mathbf{x} = \mathbf{b}$,
 3765 one could instead solve $\mathbf{P}^{-1}(\mathbf{A}\mathbf{x} - \mathbf{b}) = 0$, where \mathbf{P} is called the pre-
 3766 conditioner. The goal is to design \mathbf{P}^{-1} such that $\mathbf{P}^{-1}\mathbf{A}$ has a better condition
 3767 number, but at the same time \mathbf{P}^{-1} is easy to compute. For further infor-
 3768 mation on gradient descent, pre-conditioning and convergence we refer
 3769 to (Boyd and Vandenberghe, 2004, Chapter 9). \diamond

3770

7.1.2 Gradient Descent with Momentum

Gradient descent with momentum (Rumelhart et al., 1986) is a method that introduces an additional term to remember what happened in the previous iteration. This memory dampens oscillations and smoothes out the gradient updates. Continuing the ball analogy, the momentum term emulates the phenomenon of a heavy ball which is reluctant to change directions. The idea is to have a gradient update with memory to implement a moving average. The momentum-based method remembers the update $\Delta\mathbf{x}_i$ at each iteration i and determines the next update as a linear combination of the current and previous gradients

$$\mathbf{x}_{i+1} = \mathbf{x}_i - \gamma_i((\nabla f)(\mathbf{x}_i))^\top + \alpha\Delta\mathbf{x}_i \quad (7.9)$$

$$\Delta\mathbf{x}_i = \mathbf{x}_i - \mathbf{x}_{i-1} = -\gamma_{i-1}((\nabla f)(\mathbf{x}_{i-1}))^\top, \quad (7.10)$$

3771 where $\alpha \in [0, 1]$. Sometimes we will only know the gradient approxi-
 3772 mately. In such cases the momentum term is useful since it averages out
 3773 different noisy estimates of the gradient. One particularly useful way to
 3774 obtain an approximate gradient is using a stochastic approximation, which
 3775 we discuss next.

Goh (2017) wrote an intuitive blog post on gradient descent with momentum.

7.1.3 Stochastic Gradient Descent

Computing the gradient can be very time consuming. However, often it is possible to find a “cheap” approximation of the gradient. Approximating the gradient is still useful as long as it points in roughly the same direction as the true gradient.

Stochastic gradient
descent

Stochastic gradient descent (often shortened to SGD) is a stochastic approximation of the gradient descent method for minimizing an objective function that is written as a sum of differentiable functions. The word stochastic here refers to the fact that we acknowledge that we do not know the gradient precisely, but instead only know a noisy approximation to it. By constraining the probability distribution of the approximate gradients, we can still theoretically guarantee that SGD will converge.

In machine learning given $n = 1, \dots, N$ data points, we often consider objective functions which are the sum of the losses L_n incurred by each example n . In mathematical notation we have the form

$$L(\boldsymbol{\theta}) = \sum_{n=1}^N L_n(\boldsymbol{\theta}) \quad (7.11)$$

where $\boldsymbol{\theta}$ is the vector of parameters of interest, i.e., we want to find $\boldsymbol{\theta}$ that minimizes L . An example from regression (Chapter 9), is the negative log-likelihood, which is expressed as a sum over log-likelihoods of individual examples,

$$L(\boldsymbol{\theta}) = - \sum_{n=1}^N \log p(y_n | \mathbf{x}_n, \boldsymbol{\theta}) \quad (7.12)$$

where $\mathbf{x}_n \in \mathbb{R}^D$ are the training inputs, y_n are the training targets and $\boldsymbol{\theta}$ are the parameters of the regression model.

Standard gradient descent, as introduced previously, is a “batch” optimization method, i.e., optimization is performed using the full training set by updating the vector of parameters according to

$$\boldsymbol{\theta}_{i+1} = \boldsymbol{\theta}_i - \gamma_i (\nabla L(\boldsymbol{\theta}_i))^\top = \boldsymbol{\theta}_i - \gamma_i \sum_{n=1}^N (\nabla L_n(\boldsymbol{\theta}_i))^\top \quad (7.13)$$

for a suitable stepsize parameter γ_i . Evaluating the sum-gradient may require expensive evaluations of the gradients from all individual functions L_n . When the training set is enormous and/or no simple formulas exist, evaluating the sums of gradients becomes very expensive.

Consider the term $\sum_{n=1}^N (\nabla L_n(\boldsymbol{\theta}_i))$ in (7.13) above: we can reduce the amount of computation by taking a sum over a smaller set of L_n . In contrast to batch gradient descent, which uses all L_n for $n = 1, \dots, N$, we randomly choose a subset of L_n for mini-batch gradient descent. In the extreme case, we randomly select only a single L_n to estimate the gradient.

Why should one consider using an approximate gradient? A major reason is practical implementation constraints, such as the size of CPU/GPU memory or limits on computational time. We can think of the size of the subset used to estimate the gradient in the same way that we thought of the size of a sample when estimating empirical means 6.4.1. In practice, it is good to keep the size of the mini-batch as large as possible. Large mini-batches reduce the variance in the parameter update. Furthermore large mini-batches take advantage of highly optimized matrix operations in vectorized implementations of the cost and gradient. However when we choose the mini-batch size, we need to make sure it fits into CPU/GPU memory. Typical mini-batch sizes are 64, 128, 256, 512, 1024, which depends on the way computer memory is laid out and accessed.

This often leads to more stable convergence since the gradient estimator is less noisy.

Remark. When the learning rate decreases at an appropriate rate, and subject to relatively mild assumptions, stochastic gradient descent converges almost surely to local minimum (Bottou, 1998). \diamond

If we keep the mini-batch size small, the noise in our gradient estimate will allow us to get out of some bad local optima, which we may otherwise get stuck in.

Stochastic gradient descent is very effective in large-scale machine learning problems (Bottou et al., 2018), such as training deep neural networks on millions of images (Dean et al., 2012), topic models (Hoffman et al., 2013), reinforcement learning (Mnih et al., 2015) or training large-scale Gaussian process models (Hensman et al., 2013; Gal et al., 2014).

7.2 Constrained Optimization and Lagrange Multipliers

In the previous section, we considered the problem of solving for the minimum of a function

$$\min_{\mathbf{x}} f(\mathbf{x}) \quad (7.14)$$

where $f : \mathbb{R}^D \rightarrow \mathbb{R}$.

In this section we have additional constraints. That is for real valued functions $g_i : \mathbb{R}^D \rightarrow \mathbb{R}$ for $i = 1, \dots, m$ we consider the constrained optimization problem

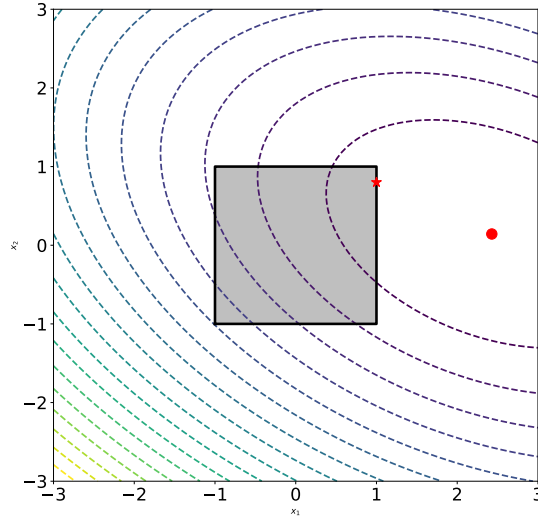
$$\begin{aligned} \min_{\mathbf{x}} \quad & f(\mathbf{x}) \\ \text{subject to} \quad & g_i(\mathbf{x}) \leq 0 \quad \text{for all } i = 1, \dots, m \end{aligned} \quad (7.15)$$

It is worth pointing out that the functions f and g_i could be non-convex in general, and we will consider the convex case in the next section.

One obvious, but not very practical, way of converting the constrained problem (7.15) into an unconstrained one is to use an indicator function

$$J(\mathbf{x}) = f(\mathbf{x}) + \sum_{i=1}^m \mathbf{1}(g_i(\mathbf{x})) \quad (7.16)$$

Figure 7.3
Illustration of constrained optimization. The unconstrained problem (indicated by the contour lines) has a minimum on the right side (indicated by the circle). The box constraints ($-1 \leq x \leq 1$ and $-1 \leq y \leq 1$) require that the optimal solution are within the box, resulting in an optimal value indicated by the star.



where $\mathbf{1}(z)$ is an infinite step function

$$\mathbf{1}(z) = \begin{cases} 0 & \text{if } z \leq 0 \\ \infty & \text{otherwise} \end{cases}. \quad (7.17)$$

3827 This gives infinite penalty if the constraint is not satisfied, and hence
 3828 would provide the same solution. However, this infinite step function is
 3829 equally difficult to optimize. We can overcome this difficulty by introduc-
 Lagrange multipliers 3830 ing *Lagrange multipliers*. The idea of Lagrange multipliers is to replace the
 3831 step function with a linear function.

Lagrangian We associate to problem (7.15) the *Lagrangian* by introducing the La-
 grange multipliers λ_i corresponding to each inequality constraint respec-
 tively (Boyd and Vandenberghe, 2004, Chapter 4).

$$\begin{aligned} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) &= f(\mathbf{x}) + \sum_{i=1}^m \lambda_i g_i(\mathbf{x}) \\ &= f(\mathbf{x}) + \boldsymbol{\lambda}^\top \mathbf{g}(\mathbf{x}) \end{aligned} \quad (7.18)$$

3832 where in the last line we have a concatenated all constraints $g_i(\mathbf{x})$ into a
 3833 vector $\mathbf{g}(\mathbf{x})$, and all the Lagrange multipliers into a vector $\boldsymbol{\lambda} \in \mathbb{R}^m$.

3834 We now introduce the idea of Lagrangian duality. In general, duality
 3835 in optimization is the idea of converting an optimization problem in one
 3836 set of variables \mathbf{x} (called the primal variables), into another optimization
 3837 problem in a different set of variables $\boldsymbol{\lambda}$ (called the dual variables). We
 3838 introduce two different approaches to duality: in this section we discuss

3839 Lagrangian duality, and in Section 7.3.3 we discuss Legendre-Fenchel du-
 3840 ality.

Theorem 7.1. *The problem in (7.15)*

$$\begin{aligned} & \min_{\mathbf{x}} f(\mathbf{x}) \\ & \text{subject to } g_i(\mathbf{x}) \leq 0 \quad \text{for all } i = 1, \dots, m \end{aligned}$$

is known as the primal problem, corresponding to the primal variables \mathbf{x} .
 The associated Lagrangian dual problem is given by

primal problem
 Lagrangian dual
 problem

$$\max_{\boldsymbol{\lambda} \in \mathbb{R}^m} \mathcal{D}(\boldsymbol{\lambda}) \tag{7.19}$$

$$\text{subject to } \boldsymbol{\lambda} \geq 0, \tag{7.20}$$

3841 where $\boldsymbol{\lambda}$ are the dual variables and $\mathcal{D}(\boldsymbol{\lambda}) = \min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$.

Proof Recall that the difference between $J(\mathbf{x})$ in (7.16) and the La-
 grangian in (7.18) is that we have relaxed the indicator function to a
 linear function. Therefore when $\lambda \geq 0$, the Lagrangian $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ is a lower
 bound of $J(\mathbf{x})$. Hence the maximum of $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ with respect to $\boldsymbol{\lambda}$ is $J(\mathbf{x})$

$$J(\mathbf{x}) = \max_{\boldsymbol{\lambda} \geq 0} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}). \tag{7.21}$$

Recall that the original problem was minimising $J(\mathbf{x})$,

$$\min_{\mathbf{x} \in \mathbb{R}^d} \max_{\boldsymbol{\lambda} \geq 0} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}). \tag{7.22}$$

By the minimax inequality (Boyd and Vandenberghe, 2004) it turns out
 that, for any function swapping the order of the minimum and maximum
 above results in a smaller value.

$$\min_{\mathbf{x} \in \mathbb{R}^d} \max_{\boldsymbol{\lambda} \geq 0} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) \geq \max_{\boldsymbol{\lambda} \geq 0} \min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) \tag{7.23}$$

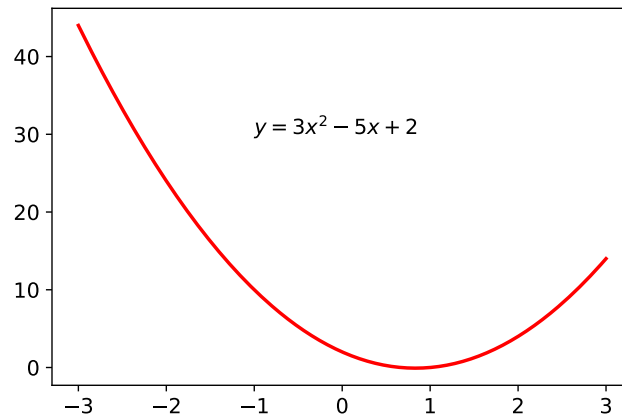
3842 This is also known as *weak duality*. Note that the inner part of the right
 3843 hand side is the dual objective function $\mathcal{D}(\boldsymbol{\lambda})$ and the theorem follows. \square

weak duality

3844 In contrast to the original optimization problem which has constraints,
 3845 $\min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ is an unconstrained optimization problem for a given
 3846 value of $\boldsymbol{\lambda}$. If solving $\min_{\mathbf{x} \in \mathbb{R}^d} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ is easy, then the overall problem
 3847 is easy to solve. The reason is that the outer problem (maximization over
 3848 $\boldsymbol{\lambda}$) is a maximum over a set of affine functions, and hence is a concave
 3849 function, even though $f(\cdot)$ and $g_i(\cdot)$ may be non-convex. The maximum
 3850 of a concave function can be efficiently computed.

3851 Assuming $f(\cdot)$ and $g_i(\cdot)$ are differentiable, we find the Lagrange dual
 3852 problem by differentiating the Lagrangian with respect to \mathbf{x} and setting
 3853 the differential to zero and solving for the optimal value. We will discuss
 3854 two concrete examples in Section 7.3.1 and 7.3.2, where $f(\cdot)$ and $g_i(\cdot)$
 3855 are convex.

Figure 7.4 Example of a convex function.



Remark (Equality constraints). Consider (7.15) with additional equality constraints

$$\begin{aligned} \min_{\mathbf{x}} f(\mathbf{x}) & \qquad (7.24) \\ \text{subject to } g_i(\mathbf{x}) & \leq 0 \quad \text{for all } i = 1, \dots, m \\ h_j(\mathbf{x}) & = 0 \quad \text{for all } j = 1, \dots, n \end{aligned}$$

3856 We can model equality constraints by replacing them with two inequality
3857 constraints. That is for each equality constraint $h_j(\mathbf{x}) = 0$ we equivalently
3858 replace it by two constraints $h_j(\mathbf{x}) \leq 0$ and $h_j(\mathbf{x}) \geq 0$. It turns out that
3859 the resulting Lagrange multipliers are then unconstrained.

3860 Therefore we constrain the Lagrange multipliers corresponding to the
3861 inequality constraints in (7.24) to be non-negative, and leave the La-
3862 grange multipliers corresponding to the equality constraints unconstrained.

3863

◇

7.3 Convex Optimization

3864

3865 We focus our attention of a particularly useful class of optimization prob-
3866 lems, where we can guarantee global optimality. When $f(\cdot)$ is a convex
3867 function, and when the constraints involving $g(\cdot)$ and $h(\cdot)$ are convex sets,
3868 this is called a *convex optimization problem*. In this setting, we have *strong*
3869 *duality*: The optimal solution of the dual problem is the same as the opti-
3870 mal solution of the primal problem. The distinction between *convex func-*
3871 *tions* and *convex sets* are often not strictly presented in machine learning
3872 literature, but one can often infer the implied meaning from context.

3873 Convex functions are functions such that a straight line between any
3874 two points of the function lie above the function. Figure 7.1 shows a non-
3875 convex function and Figure 7.2 shows a convex function. Another convex
3876 function is shown in Figure 7.4.

convex optimization
problem
strong duality
convex functions
convex sets

Definition 7.2. A function $f : \mathbb{R}^D \rightarrow \mathbb{R}$ is *convex* if for all \mathbf{x}, \mathbf{y} in the domain of f , and for any scalar θ with $0 \leq \theta \leq 1$, we have

$$f(\theta\mathbf{x} + (1 - \theta)\mathbf{y}) \leq \theta f(\mathbf{x}) + (1 - \theta)f(\mathbf{y}) \tag{7.25}$$

Technically, the domain of the function f must also be a convex set.

The constraints involving $g(\cdot)$ and $h(\cdot)$ above truncate functions at a scalar value, resulting in sets. Another relation between convex functions and convex sets is to consider the set obtained by “filling in” a convex function. A convex function is a bowl like object, and we imagine pouring water into it to fill it up. This resulting filled in set, called the epigraph of the convex function, is a convex set. Convex sets are sets such that a straight line connecting any two elements of the set lie inside the set. Figure 7.5 and Figure 7.6 illustrates convex and nonconvex sets respectively.

Definition 7.3. A set C is *convex* if for any $\mathbf{x}, \mathbf{y} \in C$ and for any scalar θ with $0 \leq \theta \leq 1$, we have

$$\theta\mathbf{x} + (1 - \theta)\mathbf{y} \in C \tag{7.26}$$

Figure 7.5 Example of a convex set



Figure 7.6 Example of a nonconvex set



Remark. We can check that a function or set is convex from first principles by recalling the definitions. In practice we often rely on operations that preserve convexity to check that a particular function or set is convex. Although the details are vastly different, this is again the idea of closure that we introduced in Chapter 2 for vector spaces. \diamond

Example

A nonnegative weighted sum of convex functions is convex. Observe that if f is a convex function, and $\alpha \geq 0$ is a nonnegative scalar, then the function αf is convex. We can see this by multiplying α to both sides of equation in Definition 7.2, and recalling that multiplying a nonnegative number does not change the inequality.

If f_1 and f_2 are convex functions, then we have by the definition

$$f_1(\theta\mathbf{x} + (1 - \theta)\mathbf{y}) \leq \theta f_1(\mathbf{x}) + (1 - \theta)f_1(\mathbf{y}) \tag{7.27}$$

$$f_2(\theta\mathbf{x} + (1 - \theta)\mathbf{y}) \leq \theta f_2(\mathbf{x}) + (1 - \theta)f_2(\mathbf{y}). \tag{7.28}$$

Summing up both sides gives us

$$\begin{aligned} & f_1(\theta\mathbf{x} + (1 - \theta)\mathbf{y}) + f_2(\theta\mathbf{x} + (1 - \theta)\mathbf{y}) \\ & \leq \theta f_1(\mathbf{x}) + (1 - \theta)f_1(\mathbf{y}) + \theta f_2(\mathbf{x}) + (1 - \theta)f_2(\mathbf{y}) \end{aligned} \tag{7.29}$$

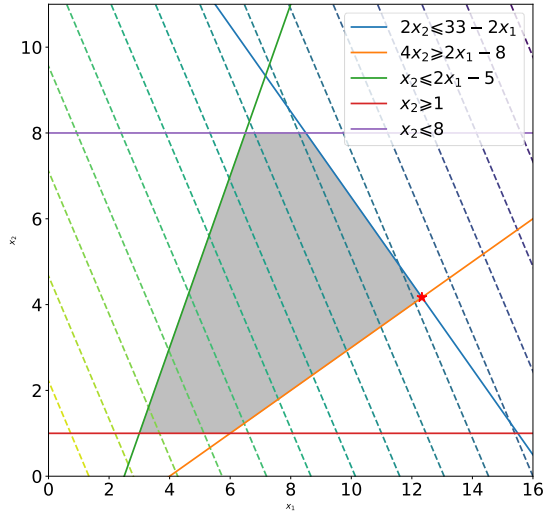
where the right hand side can be rearranged to

$$\theta(f_1(\mathbf{x}) + f_2(\mathbf{x})) + (1 - \theta)(f_1(\mathbf{y}) + f_2(\mathbf{y})) \tag{7.30}$$

completing the proof that the sum of convex functions is convex.

Combining the two facts above, we see that $\alpha f_1(\mathbf{x}) + \beta f_2(\mathbf{x})$ is convex for $\alpha, \beta \geq 0$. This closure property can be extended using a similar argument for nonnegative weighted sums of more than two convex functions.

Figure 7.7
 Illustration of a linear program. The unconstrained problem (indicated by the contour lines) has a minimum on the right side. The optimal value given the constraints are shown by the star.



3901

3902

Jensen's inequality

3903 *Remark.* The inequality defining convex functions, see 7.25, is sometimes
 3904 called *Jensen's inequality*. In fact a whole class of inequalities for taking
 3905 nonnegative weighted sums of convex functions are all called Jensen's
 3906 inequality. \diamond

3907

7.3.1 Linear Programming

Consider the special case when all the functions above are linear, that is

$$\min_{\mathbf{x} \in \mathbb{R}^d} \mathbf{c}^\top \mathbf{x} \tag{7.31}$$

subject to $\mathbf{Ax} \leq \mathbf{b}$

Linear programs are
 one of the most
 widely used
 approaches in
 industry.

3908 where $\mathbf{A} \in \mathbb{R}^{m \times d}$ and $\mathbf{b} \in \mathbb{R}^m$. This is known as a *linear program*. It has d
 3909 variables and m linear constraints.

Example

3910
 3911 An example of a linear program is illustrated in Figure 7.7, which has
 3912 two variables. The objective function is linear, resulting in linear contour
 3913 lines. The constraint set in standard form is translated into the legend. The
 3914 optimal value must lie in the shaded (feasible) region, and is indicated by
 3915 the star.
 3916

$$\min_{\mathbf{x} \in \mathbb{R}^2} - \begin{bmatrix} 5 \\ 3 \end{bmatrix}^\top \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (7.32)$$

$$\text{subject to } \begin{bmatrix} 2 & 2 \\ 2 & -4 \\ -2 & 1 \\ 0 & -1 \\ 0 & 1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \leq \begin{bmatrix} 33 \\ 8 \\ 5 \\ -1 \\ 8 \end{bmatrix} \quad (7.33)$$

3917

3918

3919

The Lagrangian is given by

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = \mathbf{c}^\top \mathbf{x} + \boldsymbol{\lambda}^\top (\mathbf{A}\mathbf{x} - \mathbf{b})$$

where $\boldsymbol{\lambda} \in \mathbb{R}^m$ is the vector of non-negative Lagrange multipliers. It is easier to see what is going on by rearranging the terms corresponding to \mathbf{x} .

$$\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) = (\mathbf{c} + \mathbf{A}^\top \boldsymbol{\lambda})^\top \mathbf{x} - \boldsymbol{\lambda}^\top \mathbf{b}$$

Taking the derivative of $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ with respect to \mathbf{x} and setting it to zero gives us

$$\mathbf{c} + \mathbf{A}^\top \boldsymbol{\lambda} = \mathbf{0}.$$

Therefore the dual Lagrangian is $\mathcal{D}(\boldsymbol{\lambda}) = -\boldsymbol{\lambda}^\top \mathbf{b}$. Recall we would like to maximize $\mathcal{D}(\boldsymbol{\lambda})$. This is traditionally presented as minimising the negative of it. In addition to the constraint due to the derivative of $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ being zero, we also have the fact that $\boldsymbol{\lambda} \geq \mathbf{0}$, resulting in the following dual optimization problem

$$\min_{\boldsymbol{\lambda} \in \mathbb{R}^m} \mathbf{b}^\top \boldsymbol{\lambda} \quad (7.34)$$

$$\text{subject to } \mathbf{c} + \mathbf{A}^\top \boldsymbol{\lambda} = \mathbf{0}$$

$$\boldsymbol{\lambda} \geq \mathbf{0}.$$

3920

3921

3922

3923

This is also a linear program, but with m variables. We have the choice of solving the primal (7.31) or the dual (7.34) program depending on whether m or d is larger. Recall that d is the number of variables and m is the number of constraints in the primal linear program.

3924

7.3.2 Quadratic Programming

Consider when the objective function is a convex quadratic function, and the constraints are affine,

$$\min_{\mathbf{x} \in \mathbb{R}^d} \frac{1}{2} \mathbf{x}^\top \mathbf{Q}\mathbf{x} + \mathbf{c}^\top \mathbf{x} \quad (7.35)$$

subject to $\mathbf{Ax} \leq \mathbf{b}$

3925 where $\mathbf{A} \in \mathbb{R}^{m \times d}$, $\mathbf{b} \in \mathbb{R}^m$ and $\mathbf{c} \in \mathbb{R}^d$. The square symmetric matrix $\mathbf{Q} \in$
 3926 $\mathbb{R}^{d \times d}$ is positive definite, and therefore the objective function is convex.
 3927 This is known as a *quadratic program*. Observe that it has d variables and
 3928 m linear constraints.

3929 Example

3930 An example of a quadratic program is illustrated in Figure 7.3, which has
 3931 two variables. The objective function is quadratic with a positive semidefinite
 3932 matrix \mathbf{Q} , resulting in elliptical contour lines. The optimal value must
 3933 lie in the shaded (feasible) region, and is indicated by the star.
 3934

$$\min_{\mathbf{x} \in \mathbb{R}^2} \frac{1}{2} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix}^\top \begin{bmatrix} 2 & 1 \\ 1 & 4 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} + \begin{bmatrix} 5 \\ 3 \end{bmatrix}^\top \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \quad (7.36)$$

$$\text{subject to } \begin{bmatrix} 1 & 0 \\ -1 & 0 \\ 0 & 1 \\ 0 & -1 \end{bmatrix} \begin{bmatrix} x_1 \\ x_2 \end{bmatrix} \leq \begin{bmatrix} 1 \\ 1 \\ 1 \\ 1 \end{bmatrix} \quad (7.37)$$

3935

3936

3937

The Lagrangian is given by

$$\begin{aligned} \mathcal{L}(\mathbf{x}, \boldsymbol{\lambda}) &= \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + \mathbf{c}^\top \mathbf{x} + \boldsymbol{\lambda}^\top (\mathbf{Ax} - \mathbf{b}) \\ &= \frac{1}{2} \mathbf{x}^\top \mathbf{Q} \mathbf{x} + (\mathbf{c} + \mathbf{A}^\top \boldsymbol{\lambda})^\top \mathbf{x} - \boldsymbol{\lambda}^\top \mathbf{b} \end{aligned}$$

where again we have rearranged the terms. Taking the derivative of $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ with respect to \mathbf{x} and setting it to zero gives

$$\mathbf{Q} \mathbf{x} + (\mathbf{c} + \mathbf{A}^\top \boldsymbol{\lambda}) = 0$$

Assuming that \mathbf{Q} is invertible, we get

$$\mathbf{x} = -\mathbf{Q}^{-1}(\mathbf{c} + \mathbf{A}^\top \boldsymbol{\lambda}) \quad (7.38)$$

Substituting (7.38) into the primal Lagrangian $\mathcal{L}(\mathbf{x}, \boldsymbol{\lambda})$ we get the dual Lagrangian

$$\mathcal{D}(\boldsymbol{\lambda}) = -\frac{1}{2}(\mathbf{c} + \mathbf{A}^\top \boldsymbol{\lambda})^\top \mathbf{Q}^{-1}(\mathbf{c} + \mathbf{A}^\top \boldsymbol{\lambda}) - \boldsymbol{\lambda}^\top \mathbf{b}$$

Therefore the dual optimization problem is given by

$$\min_{\boldsymbol{\lambda} \in \mathbb{R}^m} \frac{1}{2}(\mathbf{c} + \mathbf{A}^\top \boldsymbol{\lambda})^\top \mathbf{Q}^{-1}(\mathbf{c} + \mathbf{A}^\top \boldsymbol{\lambda}) + \boldsymbol{\lambda}^\top \mathbf{b} \quad (7.39)$$

$$\text{subject to } \boldsymbol{\lambda} \geq 0. \quad (7.40)$$

We will see an application of Quadratic Programming in machine learning in Chapter 10.

7.3.3 Legendre-Fenchel Transform and Convex Conjugate

Let us revisit the idea of duality, which we saw in Section 7.2, without considering constraints. One useful fact about convex sets is that a convex set can be equivalently described by its supporting hyperplanes. A hyperplane is called a *supporting hyperplane* of a convex set if it intersects the convex set and the convex set is contained on just one side of it. Recall that for a convex function, we can fill it up to obtain the epigraph which is a convex set. Therefore we can also describe convex functions in terms of their supporting hyperplanes. Furthermore observe that the supporting hyperplane just touches the convex function, and is in fact the tangent to the function at that point. And recall that the tangent of a function $f(\mathbf{x})$ at a given point \mathbf{x}_0 is the evaluation of the gradient of that function at that point $\left. \frac{d^2 f(\mathbf{x})}{d\mathbf{x}^2} \right|_{\mathbf{x}=\mathbf{x}_0}$. In summary, because convex sets can be equivalently described by its supporting hyperplanes, convex functions can be equivalently described by a function of their gradient. The *Legendre transform* formalizes this concept.

supporting
hyperplane

We begin with the most general definition which unfortunately has a counterintuitive form, and look at special cases to try to relate the definition to the intuition above. The *Legendre-Fenchel transform* is a transformation (in the sense of a Fourier transform) from a convex differentiable function $f(\mathbf{x})$ to a function that depends on the tangents $s(\mathbf{x}) = \nabla_{\mathbf{x}} f(\mathbf{x})$. It is worth stressing that this is a transformation of the function $f(\cdot)$ and not the variable \mathbf{x} or the function evaluated at a value. The Legendre-Fenchel transform is also known as the convex conjugate (for reasons we will see soon) and is closely related to duality.

Legendre transform
Physics students are often introduced to the Legendre transform as relating the Lagrangian and the Hamiltonian in classical mechanics. Legendre-Fenchel transform

Definition 7.4. The *convex conjugate* of a function $f : \mathbb{R}^D \rightarrow \mathbb{R}$ is a function f^* defined by

convex conjugate

$$f^*(\mathbf{s}) = \sup_{\mathbf{x} \in \mathbb{R}^D} \langle \mathbf{s}, \mathbf{x} \rangle - f(\mathbf{x}) \quad (7.41)$$

Note that the convex conjugate definition above does not need the function f to be convex nor differentiable. In the definition above, we have used a general inner product (Section 3.2) but in the rest of this section we will consider the standard dot product between finite dimensional vectors ($\langle \mathbf{s}, \mathbf{x} \rangle = \mathbf{s}^\top \mathbf{x}$) to avoid too many technical details.

To understand the above definition in a geometric fashion, consider an nice simple one dimensional convex and differentiable function, for example $f(x) = x^2$. Note that since we are looking at a one dimensional problem, hyperplanes reduce to a line. Consider a line $y = sx + c$. Recall that we are able to describe convex functions by their supporting hyper-

This derivation is easiest to understand by drawing the reasoning as it progresses.

planes, so let us try to describe this function $f(x)$ by its supporting lines. Fix the gradient of the line $s \in \mathbb{R}$ and for each point $(x_0, f(x_0))$ on the graph of f , find the minimum value of c such that the line still intersects $(x_0, f(x_0))$. Note that the minimum value of c is the place where a line with slope s “just touches” the function $f(x) = x^2$. The line passing through $(x_0, f(x_0))$ with gradient s is given by

$$y - f(x_0) = s(x - x_0) \quad (7.42)$$

The y -intercept of this line is $-sx_0 + f(x_0)$. The minimum of c for which $y = sx + c$ intersects with the graph of f is therefore

$$\inf_{x_0} -sx_0 + f(x_0). \quad (7.43)$$

3970 The convex conjugate above is by convention defined to be the negative
3971 of this. The reasoning in this paragraph did not rely on the fact that we
3972 chose a one dimensional convex and differentiable function, and holds for
3973 $f : \mathbb{R}^D \rightarrow \mathbb{R}$ which is nonconvex and non differentiable.

The classical Legendre transform is defined on convex differentiable functions in \mathbb{R}^D

However convex differentiable functions such as the example $f(x) = x^2$ is a nice special case, where there is no need for the supremum, and there is a one to one correspondence between a function and its Legendre transform. For a convex differentiable function, we know that at x_0 the tangent touches $f(x_0)$, therefore

$$f(x_0) = sx_0 + c \quad (7.44)$$

Rearranging to get an expression for $-c$

$$-c = sx_0 - f(x_0). \quad (7.45)$$

Note that $-c$ changes with x_0 and therefore with s , which is why we can think of it as a function of s , which we call $f^*(s)$.

$$f^*(s) = sx - f(x). \quad (7.46)$$

3974 The conjugate function has nice properties, for example for convex
3975 functions, applying the Legendre transform again gets us back to the origi-
3976 nal function. In the same way that the slope of $f(x)$ is s , the slope of $f^*(s)$
3977 is x . The following two examples show common uses of convex conjugates
3978 in machine learning.

Example

To illustrate the application of convex conjugates, consider the quadratic function based on a positive definite matrix $\mathbf{K} \in \mathbb{R}^{n \times n}$. We denote the primal variable to be $\mathbf{y} \in \mathbb{R}^n$ and the dual variable to be $\boldsymbol{\alpha} \in \mathbb{R}^n$.

$$f(\mathbf{y}) = \frac{\lambda}{2} \mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y} \quad (7.47)$$

Applying Definition 7.4, we obtain the function

$$f^*(\alpha) = \sup_{\mathbf{y} \in \mathbb{R}^n} \langle \mathbf{y}, \alpha \rangle - \frac{\lambda}{2} \mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}. \tag{7.48}$$

Observe that the function is differentiable, and hence we can find the maximum by taking the derivative and with respect to \mathbf{y} setting it to zero.

$$\frac{\partial [\langle \mathbf{y}, \alpha \rangle - \frac{\lambda}{2} \mathbf{y}^\top \mathbf{K}^{-1} \mathbf{y}]}{\partial \mathbf{y}} = (\alpha - \lambda \mathbf{K}^{-1} \mathbf{y})^\top \tag{7.49}$$

and hence when the gradient is zero we have $\mathbf{y} = \frac{1}{\lambda} \mathbf{K} \alpha$. Substituting into (7.48) yields

$$f^*(\alpha) = \frac{1}{\lambda} \alpha^\top \mathbf{K} \alpha - \frac{\lambda}{2} \left(\frac{1}{\lambda} \mathbf{K} \alpha \right)^\top \mathbf{K}^{-1} \left(\frac{1}{\lambda} \mathbf{K} \alpha \right) \tag{7.50}$$

$$= \frac{1}{2\lambda} \alpha^\top \mathbf{K} \alpha. \tag{7.51}$$

3979
3980
3981

Example

In machine learning we often use sums of functions, for example the objective function of the training set includes a sum of the losses for each example in the training set. In the following, we derive the convex conjugate of a sum of losses $\ell(t)$, where $\ell : \mathbb{R} \rightarrow \mathbb{R}$. This also illustrates the application of the convex conjugate to the vector case. Let $\mathcal{L}(\mathbf{t}) = \sum_{i=1}^n \ell_i(t_i)$,

$$\mathcal{L}^*(\mathbf{z}) = \sup_{\mathbf{t} \in \mathbb{R}^n} \langle \mathbf{z}, \mathbf{t} \rangle - \sum_{i=1}^n \ell_i(t_i) \tag{7.52}$$

$$= \sup_{\mathbf{t} \in \mathbb{R}^n} \sum_{i=1}^n z_i t_i - \ell_i(t_i) \quad \text{definition of dot product} \tag{7.53}$$

$$= \sum_{i=1}^n \sup_{t \in \mathbb{R}} z_i t_i - \ell_i(t_i) \tag{7.54}$$

$$= \sum_{i=1}^n \ell_i^*(z_i) \quad \text{definition of conjugate} \tag{7.55}$$

3982
3983
3984
3985
3986
3987
3988
3989

Recall that in Section 7.2 we derived a dual optimization problem using Lagrange multipliers. Furthermore for convex optimization problems we have strong duality, that is the solutions of the primal and dual problem match. The Fenchel-Legendre transform described here also can be used to derive a dual optimization problem. Furthermore then the function is

3990 convex and differentiable, the supremum is unique. To further investigate
 3991 the relation between these two approaches, let us consider a linear equal-
 3992 ity constrained convex optimization problem.

Example

Let $f(\mathbf{y})$ and $g(\mathbf{x})$ be convex functions, and \mathbf{A} a real matrix of appropriate dimensions such that $\mathbf{Ax} = \mathbf{y}$. Then

$$\min_{\mathbf{x}} f(\mathbf{Ax}) + g(\mathbf{x}) = \min_{\mathbf{Ax}=\mathbf{y}} f(\mathbf{y}) + g(\mathbf{x}). \quad (7.56)$$

By introducing the Lagrange multiplier \mathbf{u} for the constraints $\mathbf{Ax} = \mathbf{y}$,

$$\min_{\mathbf{Ax}=\mathbf{y}} f(\mathbf{y}) + g(\mathbf{x}) = \min_{\mathbf{x}, \mathbf{y}} \max_{\mathbf{u}} f(\mathbf{y}) + g(\mathbf{x}) + (\mathbf{Ax} - \mathbf{y})^\top \mathbf{u} \quad (7.57)$$

$$= \max_{\mathbf{u}} \min_{\mathbf{x}, \mathbf{y}} f(\mathbf{y}) + g(\mathbf{x}) + (\mathbf{Ax} - \mathbf{y})^\top \mathbf{u} \quad (7.58)$$

where the last step of swapping max and min is due to the fact that $f(\mathbf{y})$ and $g(\mathbf{x})$ are convex functions. By splitting up the dot product term and collecting \mathbf{x} and \mathbf{y} ,

$$\max_{\mathbf{u}} \min_{\mathbf{x}, \mathbf{y}} f(\mathbf{y}) + g(\mathbf{x}) + (\mathbf{Ax} - \mathbf{y})^\top \mathbf{u} \quad (7.59)$$

$$= \max_{\mathbf{u}} \left[\min_{\mathbf{y}} -\mathbf{y}^\top \mathbf{u} + f(\mathbf{y}) \right] + \left[\min_{\mathbf{x}} (\mathbf{Ax})^\top \mathbf{u} + g(\mathbf{x}) \right] \quad (7.60)$$

$$= \max_{\mathbf{u}} \left[\min_{\mathbf{y}} -\mathbf{y}^\top \mathbf{u} + f(\mathbf{y}) \right] + \left[\min_{\mathbf{x}} \mathbf{x}^\top \mathbf{A}^\top \mathbf{u} + g(\mathbf{x}) \right] \quad (7.61)$$

For general inner products, \mathbf{A}^\top is replaced by the adjoint \mathbf{A}^* .

Recall the convex conjugate (Definition 7.4) and the fact that dot products are symmetric,

$$\max_{\mathbf{u}} \left[\min_{\mathbf{y}} -\mathbf{y}^\top \mathbf{u} + f(\mathbf{y}) \right] + \left[\min_{\mathbf{x}} \mathbf{x}^\top \mathbf{A}^\top \mathbf{u} + g(\mathbf{x}) \right] \quad (7.62)$$

$$= \max_{\mathbf{u}} -f^*(\mathbf{u}) - g^*(-\mathbf{A}^\top \mathbf{u}). \quad (7.63)$$

Therefore we have shown that

$$\min_{\mathbf{x}} f(\mathbf{Ax}) + g(\mathbf{x}) = \max_{\mathbf{u}} -f^*(\mathbf{u}) - g^*(-\mathbf{A}^\top \mathbf{u}). \quad (7.64)$$

3993

3994

3995

3996

3997

3998

3999

4000

The Legendre-Fenchel conjugate turns out to be quite useful for machine learning problems that can be expressed as convex optimization problems. In particular for convex loss functions that apply independently to each example, the conjugate loss is a convenient way to derive a dual problem. We will see such an example in Chapter 10.

7.4 Further Reading

Continuous optimization is an active area of research, and we do not try to provide a comprehensive account of recent advances.

From a gradient descent perspective, there are two major weaknesses which each have their own set of literature. The first challenge is the fact that gradient descent is a first order algorithm, and does not use information about the curvature of the surface. When there are long valleys, the gradient points perpendicularly to the direction of interest. Conjugate gradient methods avoid the issues faced by gradient descent by taking previous directions into account (Shewchuk, 1994). Second order methods such as Newton methods use the Hessian to provide information about the curvature. Quasi-Newton methods such as L-BFGS try to use cheaper computational methods to approximate the Hessian (Nocedal and Wright, 2006).

The second challenge are non-differentiable functions. Gradient methods are not well defined when there are kinks in the function. In these cases, *subgradient methods* can be used (Shor, 1985). For further information and algorithms for optimizing non-differentiable functions, we refer to the book by Bertsekas (1999).

Modern applications of machine learning often mean that the size of datasets prohibit the use of batch gradient descent, and hence stochastic gradient descent is the current workhorse of large scale machine learning methods. Recent surveys of the literature include (Hazan, 2015; Bottou et al., 2018).

For duality and convex optimization, the book by Boyd and Vandenberghe (Boyd and Vandenberghe, 2004) includes lectures and slides online. A more mathematical treatment is provided by Bertsekas (2009). Convex optimization is based upon convex analysis, and the reader interested in more foundational results about convex functions is referred to Hiriart-Urruty and Lemaréchal (2001); Rockafellar (1970); Borwein and Lewis (2006). Legendre-Fenchel transforms are also covered in the above books on convex analysis, but more beginner friendly presentations are available at Zia et al. (2009); Gonçalves (2014).

Exercises

7.1 Consider the univariate function

$$f(x) = x^3 + 2x^2 + 5x - 3.$$

Find its stationary points and indicate whether they are maximum, minimum or saddle points.

7.2 Consider the update equation for stochastic gradient descent (Equation (7.13)).

Write down the update when we use a mini-batch size of one.

7.3 Express the following optimization problem as a standard linear program in

matrix notation

$$\max_{\mathbf{x} \in \mathbb{R}^2, \xi \in \mathbb{R}} \mathbf{p}^\top \mathbf{x} + \xi$$

4039 subject to the constraints that $\xi \geq 0$, $x_0 \leq 0$ and $x_1 \leq 3$.

7.4 The hinge loss (which is the loss used by the Support Vector Machine) is given by

$$L(\alpha) = \max\{0, 1 - \alpha\}$$

If we are interested in applying gradient methods such as L-BFGS, and do not want to resort to subgradient methods, we need to smooth the kink in the hinge loss. Compute the convex conjugate of the hinge loss $L^*(\beta)$ where β is the dual variable. Add a ℓ_2 proximal term, and compute the conjugate of the resulting function

$$L^*(\beta) + \frac{\gamma}{2}\beta^2$$

4040 where γ is a given hyperparameter.