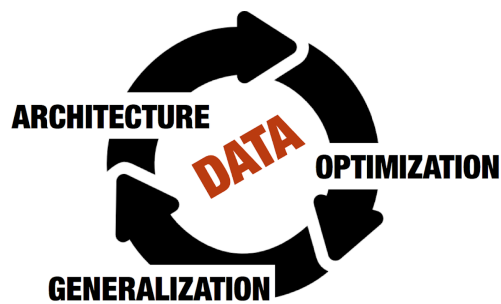


Chapter 9

Background on Optimization, Gradient Descent

We have covered the cliff-notes of the practice of deep learning in the previous eight chapters. It is by no means a complete overview. The practice of deep learning is an enticing, mysterious, and sometimes frustrating. The more time you spend playing with code the more you will learn about deep learning. New ideas are routinely discovered using very simple experiments that each of you is capable of running in your Colab now.

As we discussed, there three main concepts in machine learning. First, the class of functions $f(x; w)$ that you use to make predictions, this is called the hypothesis class or the architecture. Second, the algorithm you use to find the best model in this class of functions that fits your data; this uses tools from optimization theory. Third is the generalization performance of your classifier. Machine Learning is about picking a good hypothesis class, finding the best model within this class and making sure that the model generalizes.



The above process is relatively well-understood for simpler models such as SVMs but the story is quite murky for deep networks. Often in practice, it is never clear which architecture you should pick for your problem (many of you have asked this question in the office hours for instance), training a deep network involves a number of bells and whistles (some of which like Batch-Normalization and Dropout that we have seen) and if at the end of this exercise we get a high validation error, it is unclear how one should change

the parts of the process to improve performance. Disentangling this vicious cycle is what “understanding deep learning” is all about.

Goal Module 2 will develop an understanding of optimization and generalization for more generic machine learning models first. It will end with an insight into understanding their interplay for deep networks. Module 2 has a different flavor, it is more theoretical. Our goal is to grasp the general concepts behind these theoretical results and understand the training process of deep networks better. This will also help us train deep networks much better in practice.

9.1 Convexity

Consider a function $\ell : \mathbb{R}^p \rightarrow \mathbb{R}$ that is convex, i.e., for any w, w' that lie in the domain (which is assumed to be a convex set) of f and any $\lambda \in [0, 1]$ we have

$$\ell(\lambda w + (1 - \lambda)w') \leq \lambda \ell(w) + (1 - \lambda)\ell(w'). \quad (9.1)$$

A function $\ell(w)$ is concave if $-\ell(w)$ is convex. Some examples of convex functions are

- affine functions $Aw + b$, norms $\|w\|_p = (\sum_{i=1}^p |w_i|^p)^{1/p}$, or $\|w\|_\infty = \max_k |w_k|$.
- exponential e^w for $w \in \mathbb{R}$
- powers w^α for $w > 0$ and $\alpha \geq 1$ or $\alpha \leq 0$
- powers of absolute values $|w|^p$ for $w \in \mathbb{R}$ and $p \geq 1$

Strictly convex functions Strictly convex functions have the property that for all $w \neq w'$ in the domain (which is assumed to be convex) and $\lambda \in (0, 1)$

$$\ell(\lambda w + (1 - \lambda)w') < \lambda \ell(w) + (1 - \lambda)\ell(w').$$

First-order condition for convexity If ℓ is differentiable, the definition of convexity in (9.1) is equivalent to the following first-order condition. A differentiable function ℓ with convex domain is convex iff

$$\ell(w') \geq \ell(w) + \langle \nabla \ell(w), w' - w \rangle. \quad (9.2)$$

for all w, w' in the domain. Note that the first-order condition is equivalent to the definition of convexity in (9.1) for differentiable functions. The proof is long but easy; you can see https://www.princeton.edu/aaa/Public/Teaching/ORF523/S16/ORF523_S16_Lec7_gh.pdf for the proof. For strictly convex functions the inequality is strict

$$\ell(w') > \ell(w) + \langle \nabla \ell(w), w' - w \rangle.$$

56 **Monotonicity of the gradient for convex functions** The first-order con-
 57 dition for convexity gives a useful, and equivalent, characterization of the
 58 gradient. Write (9.2) for w, w' in two opposite directions

$$\begin{aligned}\ell(w) &\geq \ell(w') + \langle \nabla \ell(w'), w' - w \rangle \\ \ell(w') &\geq \ell(w) + \langle \nabla \ell(w), w - w' \rangle\end{aligned}$$

59 and add them to get

$$\langle \nabla \ell(w) - \nabla \ell(w'), w - w' \rangle \geq 0. \quad (9.3)$$

60 It is also true that monotonicity of the gradient implies convexity (try to prove
 61 it).

62 **Second-order condition for convexity** If ℓ is twice-differentiable with a
 63 convex domain, it is convex iff

$$\nabla^2 \ell(w) \succeq 0 \quad (9.4)$$

64 for all w in the domain. The symbol \succeq denotes positive semi-definiteness of
 65 the Hessian matrix $\nabla^2 \ell(w)$

$$(\nabla^2 \ell(w))_{ij} = \frac{\partial^2 \ell(w)}{\partial w_i \partial w_j}.$$

66 For strictly convex functions, the inequality in (9.4) is strict, i.e., the Hessian
 67 is positive definite. As an example, the least squares objective $\ell(w) = \|y -$
 68 $Xw\|_2^2$ is convex because

$$\nabla^2 \ell(w) = 2X^\top X$$

69 which is positive definite for any X (why is it positive definite and not just
 70 positive semi-definite?).

71 **Strongly convex functions** A function is strongly convex if there exists an
 72 $m > 0$ such that

$$\ell(w) - \frac{m}{2} \|w\|_2^2 \text{ is convex.} \quad (9.5)$$

73 It is easy to see that strict convexity implies convexity. Since the function
 74 $\ell(w) - m/2 \|w\|^2$ is convex, it satisfies the definition of convexity:

$$\begin{aligned}\ell(\lambda w + (1 - \lambda)w') - \frac{m}{2} \|\lambda w + (1 - \lambda)w'\|^2 \\ \leq \lambda \left(\ell(w) - \frac{m}{2} \|w\|^2 \right) + (1 - \lambda) \left(\ell(w') - \frac{m}{2} \|w'\|^2 \right).\end{aligned} \quad (9.6)$$

75 But

$$\frac{\lambda m}{2} \|w\|^2 + \frac{(1 - \lambda)m}{2} \|w'\|^2 - \frac{m}{2} \|\lambda w + (1 - \lambda)w'\|^2 > 0$$

76 for $\lambda \in (0, 1)$ for all $w \neq w'$ because $\|w\|^2$ is strictly convex which shows
 77 that if we have a strongly convex function ℓ it also satisfies

$$\ell(\lambda w + (1 - \lambda)w') \leq \lambda \ell(w) + (1 - \lambda) \ell(w').$$

In other words, we have

$$\text{strong convexity} \Rightarrow \text{strict convexity} \Rightarrow \text{convexity}.$$

Observe that strongly convexity in (9.6) is a stronger version of Jensen's inequality. Strongly convex functions are easier to optimize for algorithms. It will also always be much easier to prove a result in optimization on strongly convex functions. It is easy to see using the second-order condition for convexity that an m -strongly convex function has

$$\nabla^2 \ell(w) \succeq mI_{p \times p}.$$

We will use the following first-order condition for strongly convex functions often. A function is m -strongly convex if and only if

$$\ell(w') \geq \ell(w) + \langle \nabla \ell(w), w' - w \rangle + \frac{m}{2} \|w' - w\|^2 \quad (9.7)$$

for any w, w' in the domain.

9.2 Introduction to Gradient Descent

In this chapter, we will write $\ell(w)$ to denote the training objective, i.e., if we have a classifier $f(x; w)$ and a dataset $D = \{(x^i, y^i)\}_{i=1, \dots, n}$ of n samples we will denote

$$\ell(w) := \frac{1}{n} \sum_{i=1}^n \ell(w; x^i, y^i).$$

The objective ℓ will always be a function of the entire dataset but we will keep the dependence implicit. Note that the number of samples n is usually quite large in deep learning, so the summation above has a large number of terms on the right-hand side.

Gradient descent is a simple algorithm to minimize $\ell(w)$. Before we study its properties, it will help to refresh the following few facts.

9.2.1 Conditions for optimality

Local and global minima A point w is a local minimum of the function $\ell(w)$ for all w' in a neighborhood of w we have $\ell(w) \leq \ell(w')$. The point is a global minimum of the function ℓ if this condition is true for all w' in the domain, not just the ones in the neighborhood.

Local minima are global minima for convex functions This is easy to see using an argument by contradiction. If w is a local minimum that is not the global minimum, there exists a point w' in the domain such that $\ell(w') < \ell(w)$. The domain of the function is convex, so pick a point $v = \lambda w' + (1 - \lambda)w$ and see that

$$\ell(v) - \ell(w) \leq \lambda(\ell(w') - \ell(w))$$

using the definition of convexity. Since w is only a local minimum, we can pick λ to be small enough that the left hand side is non-negative. This shows that $\ell(w') \geq \ell(w)$ but this means that w is a global minimum and we have a contradiction.

111 **Global minimum is unique for strictly convex functions** If a function is
 112 strictly convex on a convex domain the optimal solution (if it exists) must be
 113 unique. Indeed, if there were two solutions w, w' that were both minimizers
 114 we would have

$$\ell(w) = \ell(w') \leq \ell(w'') \quad \forall w'' \quad (9.8)$$

115 We can now apply the definition of convexity to the point $v = (w + w')/2$ to
 116 get

$$\ell(v) < \frac{1}{2}\ell(w) + \frac{1}{2}\ell(w') = \ell(w).$$

117 which contradicts (9.8). The least-squares objective is strictly convex, so the
 118 solution is unique global minimizer of the objective.

119 **First-order optimality condition** If w is a local minimum of a continuously
 120 differentiable function ℓ , then it satisfies

$$\nabla \ell(w) = 0. \quad (9.9)$$

121 If further ℓ is convex, then $\nabla \ell(w) = 0$ is a sufficient condition for global
 122 optimality from the above discussion.

123 9.2.2 Different types of convergence

124 Let us assume that we have a continuously differentiable convex function ℓ
 125 and let

$$w^* = \underset{w}{\operatorname{argmin}} \ell(w)$$

126 be the global minimizer of this function.

127 We would like to develop an iterative scheme that takes in the initialization
 128 of the weights w^0 and updates them to obtain a sequence

$$w^0, w^1, \dots, w^t, \dots$$

129 Along this sequence we are interested in understanding the

- 130 1. convergence of the function value $\ell(w^t)$ to the minimal value $\ell(w^*)$,
 131 and
- 132 2. convergence of the iterates $\|w^t - w^*\|$.

133 **Descent direction** We are going to perform a sequence of updates given by

$$w^{t+1} = w^t + \eta d^t \quad (9.10)$$

134 where d^t is called the descent direction and the scalar parameter $\eta > 0$ is called
 135 the step-size and determines how far we travel using this descent direction.
 136 Any direction such that

$$\langle \nabla \ell(w^t), d^t \rangle < 0$$

137 is a good descent direction because this leads to a reduction in the value of the
 138 function $\ell(w^{t+1})$ after the weight update. There are numerous ways to pick a
 139 good descent direction. Among the simplest ones is gradient descent which

140 descends along the direction of the negative gradient and thereby performs the
141 following set of updates

$$w^{t+1} = w^t - \eta \nabla \ell(w^t) \quad (9.11)$$

142 given an initial value w^0 . The step-size (also called the learning rate) is chosen
143 by the user. The step-size need not always be fixed, for instance you chose
144 it to be a function of the number of weight updates t in the homework. A
145 good step-size is one that does not overshoot the minimum w^* . For instance,
146 after having chosen a particular descent direction d^t we can compute the best
147 step-size to use at time t by solving

$$\eta^t = \operatorname{argmin}_{\eta \geq 0} \ell(w^t + \eta d^t).$$

148 This is known as line-search in the optimization literature. You may have seen
149 Newton's method

$$w^{t+1} = w^t - (\nabla^2 \ell(w^t))^{-1} \nabla \ell(w^t). \quad (9.12)$$

150 which does not have a user-tuned step-size and further modifies the descent
151 direction to be the product of the inverse Hessian with the gradient.

▲ Draw a picture of overshooting using a large step-size.

❓ Can you think of an algorithm for minimizing a function that does not use the gradient of the function to compute the descent direction?

152 9.3 Convergence rate for gradient descent

153 We will next understand how quickly gradient descent converges to the global
154 minimum. There are two concrete goals of this analysis

- 155 1. to be able to pick the step-size to avoid overshooting without doing
156 line-search, and
- 157 2. characterize how many iterations of gradient descent to run until we are
158 guaranteed to be within some distance of the global minimum.

159 9.3.1 Some assumptions

160 Before we begin, we will make a few simplifying assumptions on the function
161 $\ell(w)$. These are quite typical in optimization and ensure that we are not dealing
162 with pathological functions that make minimizing them arbitrarily hard.

- 163 1. **Lipschitz continuity/bounded gradients** We will assume that ℓ is Lip-
164 schitz continuous

$$|\ell(w) - \ell(w')| \leq B \|w - w'\|_2. \quad (9.13)$$

165 for some $B > 0$. You might also see this condition written as

$$\|\nabla \ell(w)\| \leq B$$

166 for differentiable functions.

167 **2. Smoothness** We will always consider functions such that their gradients
 168 are L -Lipschitz, i.e.,

$$\|\nabla \ell(w) - \nabla \ell(w')\|_2 \leq L \|w - w'\|_2. \quad (9.14)$$

169 If ℓ is twice-differentiable, this is equivalent to assuming

$$\nabla^2 \ell(w) \preceq L I_{p \times p}. \quad (9.15)$$

170 From the Cauchy-Schwarz inequality which states that

$$\langle u, v \rangle \leq \|u\| \|v\|$$

171 for two vectors u, v , we have the following implication of smoothness:

$$\langle \nabla \ell(w) - \nabla \ell(w^*), w - w^* \rangle \leq L \|w - w^*\|^2. \quad (9.16)$$

172 A related concept is called **co-coercivity** of the gradient. The gradient
 173 being L -Lipschitz is equivalent to co-coercivity of the gradient with
 174 parameter $1/L$

$$\frac{1}{L} \|\nabla \ell(w) - \nabla \ell(w')\|^2 \leq \langle \nabla \ell(w) - \nabla \ell(w'), w - w' \rangle. \quad (9.17)$$

175 We can see that co-coercivity implies Lipschitz continuity of the gradi-
 176 ents $\nabla \ell(w)$ using (9.16) and (9.17). The reverse is also true, Lipschitz-
 177 continuity of the gradient implies the Descent Lemma Lemma 1 which
 178 is seen by applying the Descent Lemma twice for the two functions
 179 $g(u) = \ell(u) - \langle \nabla \ell(w'), u \rangle$ and $h(u) = \ell(u) - \langle \nabla \ell(w), u \rangle$.

180 9.3.2 GD for convex functions

181 We begin with the so-called Descent Lemma.

182 **Lemma 1 (Descent Lemma).** For an L -smooth function, we have

$$\ell(w') \leq \ell(w) + \langle \nabla \ell(w), w' - w \rangle + \frac{L}{2} \|w' - w\|^2. \quad (9.18)$$

183 for any two w, w' in the domain.

184 **Proof.** First, you should compare this with the first-order characterization of
 185 convexity

$$\ell(w') \geq \ell(w) + \langle \nabla \ell(w), w' - w \rangle.$$

186 The two conditions can be used to sandwich the value of $\ell(w^{t+1})$ given the
 187 value of $\ell(w^t)$ in gradient descent with room for a quadratic term $\frac{L}{2} \|w' - w\|^2$.
 188 This also gives some intuition as to what L -smooth really means; a large value
 189 of L means that the function ℓ decreases quickly. Let $v = w + \lambda(w' - w)$ and
 190 use Taylor's theorem to see that

$$\ell(w') = \ell(w) + \int_0^1 \langle \nabla \ell(v), w' - w \rangle \, d\lambda \quad (9.19)$$

191 Subtract $\langle \nabla \ell(w), w' - w \rangle$ from both sides to get

$$\ell(w') - \ell(w) - \langle \nabla \ell(w), w' - w \rangle = \int_0^1 \langle \nabla \ell(v) - \nabla \ell(w), w' - w \rangle \, d\lambda.$$

192 Observe that

$$\begin{aligned}
 |\ell(w') - \ell(w) - \langle \nabla \ell(w), w' - w \rangle| &= \left| \int_0^1 \langle \nabla \ell(v) - \nabla \ell(w), w' - w \rangle \, d\lambda \right| \\
 &\leq \int_0^1 |\langle \nabla \ell(v) - \nabla \ell(w), w' - w \rangle| \, d\lambda \\
 &\leq \int_0^1 \|\nabla \ell(v) - \nabla \ell(w)\| \|w' - w\| \, d\lambda \\
 &\leq L \int_0^1 \lambda \|w' - w\|^2 \, d\lambda \\
 &= \frac{L}{2} \|w' - w\|^2.
 \end{aligned}$$

193 This completes the proof after removing the absolute value on the left-hand
 194 side. \square

195 We can use the Descent Lemma twice on two points to w, w' to get (9.16).
 196 Another direct consequence of the Descent Lemma is the following corollary
 197 that relates the value $\ell(w)$ at any point w in the domain to that of the global
 198 minimum.

199 **Corollary 2.** For L -smooth convex function ℓ , if w^* is the global minimizer,
 200 then

$$\frac{1}{2L} \|\nabla \ell(w)\|^2 \leq \ell(w) - \ell(w^*) \leq \frac{L}{2} \|w - w^*\|^2. \quad (9.20)$$

201 **Proof.** Since $\nabla \ell(w^*) = 0$, the right-hand side follows directly from the
 202 Descent Lemma. To get the left-hand side, let us optimize the upper bound in
 203 the Descent Lemma using $w' = w + \lambda v$ with $\|v\| = 1$ as follows

$$\begin{aligned}
 \ell(w^*) &= \inf_{w'} \ell(w') \leq \inf_{w'} \left\{ \ell(w) + \langle \nabla \ell(w), w' - w \rangle + \frac{L}{2} \|w' - w\|^2 \right\} \\
 &= \inf_{\|v\|=1} \inf_{\lambda} \left\{ \ell(w) + \lambda \langle \nabla \ell(w), v \rangle + \frac{L}{2} \lambda^2 \right\} \\
 &= \inf_{\|v\|=1} \left\{ \ell(w) - \frac{1}{2L} (\langle \nabla \ell(w), v \rangle)^2 \right\} \\
 &= \ell(w) - \frac{1}{2L} \|\nabla \ell(w)\|^2.
 \end{aligned}$$

204 \square

205 In other words, the gap between the function values $\ell(w) - \ell(w^*)$ is upper-
 206 bounded by the gap to the minimizer $\frac{L}{2} \|w - w^*\|^2$ and lower-bounded by the
 207 norm of the gradient $\frac{1}{2L} \|\nabla \ell(w)\|^2$.

208 **Lemma 3 (Monotonic progress for gradient descent).** For gradient descent
 209 $w^{t+1} = w^t - \eta \nabla \ell(w^t)$, if we pick the step-size

$$\eta \leq \frac{1}{L} \quad (9.21)$$

210 we have

$$\ell(w^{t+1}) \leq \ell(w^t) - \frac{\eta}{2} \|\nabla \ell(w^t)\|^2. \quad (9.22)$$

211 Further,

$$\ell(w^{t+1}) - \ell(w^*) \leq \frac{1}{2\eta} (\|w^t - w^*\|^2 - \|w^{t+1} - w^*\|^2) \quad (9.23)$$

212 which implies

$$\|w^{t+1} - w^*\|^2 \leq \|w^t - w^*\|^2. \quad (9.24)$$

213 **Proof.** Substitute $\eta \leq 1/L$ in the Descent Lemma and simplify to get (9.22).
 214 The second result is obtained by

$$\begin{aligned} 0 \leq \ell(w^{t+1}) - \ell(w^*) &\leq \ell(w^t) - \ell(w^*) - \frac{\eta}{2} \|\nabla \ell(w^t)\|^2 \\ &\leq \langle \nabla \ell(w^t), w^t - w^* \rangle - \frac{\eta}{2} \|\nabla \ell(w^t)\|^2 \\ &= \frac{1}{2\eta} (\|w^t - w^*\|^2 - \|w^t - w^* - \eta \nabla \ell(w^t)\|^2) \\ &= \frac{1}{2\eta} (\|w^t - w^*\|^2 - \|w^{t+1} - w^*\|^2). \end{aligned}$$

215 Observe that since the left-hand side is positive, the claim in (9.24) is true. \square

216 We have therefore shown that if the step-size is not too large (the smooth-
 217 ness parameter of the function determines how large the step-size can be)
 218 then gradient descent always improves the value of the function with each
 219 iteration (9.22). It also improves the distance of the weights to the global
 220 minimum at each iteration (9.24).

221 **Lemma 4 (Convergence rate for gradient descent, convex function).** For
 222 gradient descent $w^{t+1} = w^t - \eta \nabla \ell(w^t)$ with step-size $\eta < 1/L$, we have

$$\ell(w^{t+1}) - \ell(w^*) \leq \frac{1}{2t\eta} \|w^0 - w^*\|^2. \quad (9.25)$$

223 **Proof.** We sum up the expression in (9.23) for all times t to get

$$\begin{aligned} \sum_{s=1}^t \ell(w^s) - \ell(w^*) &\leq \frac{1}{2\eta} \sum_{s=1}^t (\|w^{s-1} - w^*\|^2 - \|w^s - w^*\|^2) \\ &= \frac{1}{2\eta} (\|w^0 - w^*\|^2 - \|w^t - w^*\|^2) \\ &\leq \frac{1}{2\eta} \|w^0 - w^*\|^2. \end{aligned}$$

224 We know from (9.22) that $\ell(w^t)$ is non-increasing, so we can write

$$\ell(w^t) - \ell(w^*) \leq \frac{1}{t} \sum_{s=1}^t (\ell(w^s) - \ell(w^*)) \leq \frac{1}{2t\eta} \|w^0 - w^*\|^2.$$

225 \square

If we want to find a weights with

$$\ell(w^t) - \ell(w^*) \leq \epsilon$$

for a convex function, we need to run gradient descent for at least

$$\mathcal{O}(1/\epsilon)$$

iterations. This is an important result to remember.

9.3.3 Gradient descent for strongly convex functions

Things are much better if the function we are minimizing is strongly convex. First we have the following lemma for strongly-convex functions which involves a rewriting co-coercivity condition for strongly convex functions.

Lemma 5 (Co-coercivity for strongly convex function). If $\ell(w)$ is m -strongly convex, and L -smooth, then the function $g(w) = \ell(w) - \frac{m}{2}\|w\|^2$ is convex and $L - m$ -smooth. The co-coercivity condition for $\nabla g(w)$ can therefore be re-written as

$$\langle \nabla \ell(w) - \nabla \ell(w'), w - w' \rangle \geq \frac{mL}{m+L} \|w - w'\|^2 + \frac{1}{m+L} \|\nabla \ell(w) - \nabla \ell(w')\|^2. \quad (9.26)$$

Proof. The convexity of $g(w)$ is immediate to see from the definition of strong convexity of $\ell(w)$. Use the monotonicity of the gradient of $g(w)$ to get

$$\begin{aligned} 0 &\leq \langle \nabla g(w) - \nabla g(w'), w - w' \rangle \\ &= \langle \nabla \ell(w) - \nabla \ell(w'), w - w' \rangle - m\|w - w'\|^2 \\ &\leq (L - m)\|w - w'\|^2. \end{aligned}$$

We can now rewrite the co-coercivity condition for $\nabla g(w)$ with the smoothness parameter $L - m$ and simplify to get (9.26). \square

Lemma 6 (Convergence rate of gradient descent for strongly convex functions). For strongly convex functions we have pick a step-size

$$0 < \eta < \frac{2}{m+L}$$

to get

$$\|w^{t+1} - w^*\|^2 \leq \left(1 - \eta \frac{2mL}{m+L}\right) \|w^t - w^*\|^2. \quad (9.27)$$

which gives

$$\|w^t - w^*\|^2 \leq c^t \|w^0 - w^*\|^2 \quad (9.28)$$

where $c = \left(1 - \eta \frac{2mL}{m+L}\right)$.

Proof. We expand the left hand-side in (9.27) to get

$$\begin{aligned} \|w^{t+1} - w^*\|^2 &= \|w^t - \eta \nabla \ell(w^t) - w^*\|^2 \\ &= \|w^t - w^*\|^2 - 2\eta \langle \nabla \ell(w^t), w^t - w^* \rangle + \eta^2 \|\nabla \ell(w^t)\|^2 \\ &\leq \left(1 - \eta \frac{2mL}{m+L}\right) \|w^t - w^*\|^2 + \eta \left(\eta - \frac{2}{m+L}\right) \|\nabla \ell(w^t)\|^2 \\ &\leq \left(1 - \eta \frac{2mL}{m+L}\right) \|w^t - w^*\|^2. \end{aligned}$$

We have substituted the co-coercivity condition from (9.26) for the inner-product with $w' := w^t$ and $w := w^*$ to get the first inequality. This implies that the distance to the global minimum $\|w^t - w^*\|$ decreases multiplicatively; compare this with (9.24) where the progress is additive. The additional assumption of strong convexity therefore means that we are making very quick progress towards the global minimum. We can use this inequality repeatedly for all iterations t to get

$$\|w^t - w^*\|^2 \leq c^t \|w^0 - w^*\|^2$$

where $c = \left(1 - \eta \frac{2mL}{m+L}\right)$. □

Strong convexity enables much faster progress towards the global minimum. If we want $\|w^t - w^*\| \leq \epsilon$ we need

$$\mathcal{O}(\log(1/\epsilon))$$

iterations of gradient descent. This is *much* less than that for a convex function. This is called *linear* convergence because we need a constant number of iterations to reduce the gap to the optimal in half. The naming convention is a bit unusual here but you will see that if we plot $\log\|w^t - w^*\|$ (or $\log(\ell(w^t) - \ell(w^*))$) on the Y-axis and number of iterations t on the X-axis, we get a straight line for gradient descent on strongly-convex functions; you can see this from (9.28).

We say that the convergence rate of gradient descent for non-strongly convex functions is *sub-linear*. The longer we run GD for convex functions, the slower its progress.

Further, if we pick the largest step-size $\eta = 2/(m + L)$ we get

$$c = \left(\frac{\kappa - 1}{\kappa + 1}\right)^2 < 1. \quad (9.29)$$

where $\kappa = L/m$ is the condition number of the Hessian (it is the ratio of the largest eigenvalue and the smallest eigenvalue). Larger the condition number κ , closer to 1 the multiplicative constant c and *slower* the convergence rate of gradient descent.

▲ Plot the convergence rate of gradient descent for convex and strongly-convex functions.

▲ In the optimization literature, an algorithm with

$$\lim_{t \rightarrow \infty} \frac{\ell(w^{t+1}) - \ell(w^*)}{\ell(w^t) - \ell(w^*)} = \rho$$

is said to be sub-linear if $\rho \in (0, 1)$, linear if $\rho = 1$ and super-linear if $\rho = 0$.

A few more points to note

1. The step-size is limited by $m + L$ but the convergence rate depends on $\kappa = L/m$. Smaller the value of c , faster the convergence.
2. Larger the L , smaller the ideal step-size η
3. You can get the upper bound

$$\ell(w^t) - \ell(w^*) \leq \frac{L}{2} \|w^t - w^*\|^2 \leq \frac{c^t L}{2} \|w^0 - w^*\|^2 \quad (9.30)$$

using (9.20).

You will also see the convergence rate written in many papers as

$$\|w^t - w^*\| \leq e^{-4t/\kappa} \|w^0 - w^*\|. \quad (9.31)$$

259 You can get this inequality by using the fact that $1 + x \leq e^x$ in (9.29). We can
 260 use this to pull out the dependence on κ in the convergence rate; for strongly
 261 convex functions, gradient descent requires

$$\mathcal{O}(\kappa \log(1/\epsilon))$$

262 iterations to reach within an ϵ -neighborhood of the global minimum $\ell(w^*)$.
 263 This suggests that smaller the condition number κ fewer the iterations required.

264 We can intuitively understand why convergence of gradient descent is
 265 slower for a large condition number. A large condition number means that
 266 some directions of the objective ℓ are highly curved while some others are
 267 very flat. It is difficult to pick one single scalar step-size in such situations that
 268 makes quick progress along the flat directions but also does not overshoot the
 269 highly curved directions. You might imagine that clever schemes to change the
 270 step-size depending upon the local geometry of the function $\ell(w^t)$ could help
 271 solve this issue and indeed it does, but such adaptive schemes are expensive to
 272 implement computationally. We will see some algorithms that pick different
 273 step-sizes for different weights in Chapter 11.

▲ Draw a picture of this phenomenon for a quadratic objective $\ell(w) = \langle w, Aw \rangle$ for matrices $A \succ 0$ with different condition numbers κ .

274 9.4 Limits on convergence rate of first-order meth- 275 ods

276 It is a powerful and deep result that we cannot do better than a linear conver-
 277 gence rate for optimization methods that only use the gradient of the function
 278 $\ell(w)$. More precisely, for any first-order method, i.e., any method where the
 279 iterate at step t given by w^t is chosen to be

$$w^t \in w^0 + \text{span} \{ \nabla \ell(w^0), \dots, \nabla \ell(w^t) \},$$

280 we have the following theorem by Yurii Nesterov.

281 **Theorem 7 (Nesterov's lower bound).** If $w \in \mathbb{R}^p$, for any $t \leq (p-1)/2$
 282 and every initialization of weights w^0 there exist functions $\ell(w)$ that are
 283 convex, differentiable, L -smooth with finite optimal value $\ell(w^*)$ such that any
 284 first-order method has

$$\ell(w^t) - \ell(w^*) \geq \frac{3}{32} \frac{L \|w^0 - w^*\|^2}{(t+1)^2}.$$

p is the dimension of the parameter space

285 Let us read the statement of the theorem carefully. It states that *fixed* a time
 286 t and initial condition w^0 , we can *find* a convex function $\ell(w)$ such that it takes
 287 gradient descent at least $\mathcal{O}(1/\epsilon^2)$ to reach an ϵ -neighborhood of the optimal
 288 value $\ell(w^*)$. The implication of this theorem is as follows. The convergence
 289 rate $\mathcal{O}(1/\epsilon)$ we obtained for convex functions is not the best rate we can
 290 get. Nesterov's lower bound suggests that there should be gradient-based
 291 algorithms that only require $\mathcal{O}(1/\sqrt{\epsilon})$ iterations. Such methods will be the
 292 topic of the next Chapter.