

14.1 Review and overview

In the last lecture, we introduced the concept of *algorithmic regularization* (or *implicit regularization*) and demonstrated it for the problem of overparametrized linear regression. We also began applying it to a non-convex case where we assume that the ground-truth model is sparse. Consider the following quadratic parametrization:

$$y = \langle \beta^* \odot \beta^*, x \rangle, \quad (14.1)$$

where β^* is an r -sparse vector. We write $\beta^* = \mathbf{1}_S$ where S is the set of indices on which β^* takes nonzero value. Since β^* is r -sparse, we have that $|S| \leq r$. Our goal is to show that with $\Omega(r^2)$ samples, gradient descent on $\hat{L}(\beta)$ starting from small initialization can converge to β^* . This is expressed more formally by the following theorem.

Theorem 14.1. *Let c be a sufficiently large universal constant. Suppose $n \geq cr^2 \log^2(d)$ and $\alpha \leq 1/d^c$, then when $\frac{\log(d/\alpha)}{\eta} \lesssim T \lesssim \frac{1}{\eta\sqrt{d\alpha}}$, we have*

$$\|\beta^T \odot \beta^T - \beta^* \odot \beta^*\|_2^2 \leq O(\alpha\sqrt{d}). \quad (14.2)$$

The main intuition for proving Theorem 14.1 is to establish strong uniform convergence on $S_r = \{\beta : \|\beta\|_0 \leq O(r)\}$, i.e. gradient descent on the population loss \hat{L} is sufficiently “similar” to gradient descent on the empirical loss $\hat{L}(\beta)$.

As a warm-up, we began to prove the following result on convergence for gradient descent on $L(\beta)$:

Theorem 14.2. *For sufficiently small η , gradient descent on $L(\beta)$ converges to β^* in $\Theta\left(\frac{\log(1/(\epsilon\alpha))}{\eta}\right)$ iterations with ϵ -error in ℓ_2 -distance.*

In this lecture, we will wrap up the proof of Theorem 14.2 and then prove Theorem 14.1 in the case of $r = 1$.

14.2 Warm-up: Gradient descent on population loss

The analysis of Theorem 14.2 can be broken into two cases depending on where coordinate i of β^* is nonzero or not, i.e. $i \in S$ or $i \notin S$. Refer to the previous lecture for the proof in the case of $i \in S$. Here, we will prove the following lemma for $i \notin S$, then explain informally how the lemma implies Theorem 14.2.

Lemma 14.3. *For all $i \notin S$, when $t \leq 1/(10\eta\alpha^2)$ we have that $\beta_i^t \leq 2\alpha$.*

Proof. For a coordinate $i \notin S$, the gradient descent update for this problem becomes

$$\beta_i^{t+1} = [\beta^t - \eta(\beta^t \odot \beta^t - \beta^* \odot \beta^*) \odot \beta^t]_i \quad (14.3)$$

$$= \beta_i^t - \eta(\beta_i^t \cdot \beta_i^t) \cdot \beta_i^t \quad (\text{since } \beta_i^* = 0 \ \forall i \notin S) \quad (14.4)$$

$$= \beta_i^t - \eta(\beta_i^t)^3. \quad (14.5)$$

Since our initialization β^0 was small, the update to these coordinates will be even smaller because $(\beta_i^t)^3$ is small. We can prove the desired claim using strong induction. Suppose $\beta_i^s \leq 2\alpha$ for all $s \leq t$ and $i \notin S$, and that $t+1 \leq 1/(10\eta\alpha^2)$. Then, for all $s \leq t$,

$$\beta_i^{s+1} = (1 - \eta(\beta_i^s)^2)\beta_i^s \quad (14.6)$$

$$\leq (1 + \eta(\beta_i^s)^2)\beta_i^s \quad (14.7)$$

$$\leq (1 + 4\eta\alpha^2)\beta_i^s. \quad (\text{since } \beta_i^s \leq 2\alpha) \quad (14.8)$$

With strong induction, we can repeatedly apply this gradient update starting from $t = 0$ to obtain

$$\beta_i^{t+1} \leq \beta_0 \cdot (1 + 4\eta\alpha^2)^t \quad (14.9)$$

$$\leq \beta_0(1 + 4\eta\alpha^2)^{\frac{1}{10\eta\alpha^2}} \quad (14.10)$$

$$\leq \beta_0 \exp\left(\frac{4\eta\alpha^2}{10\eta\alpha^2}\right) \quad (14.11)$$

$$= \beta_0 \cdot e^{2/5} \quad (14.12)$$

$$\leq 2\alpha, \quad (14.13)$$

which completes the inductive proof of the claim. □

Informally, what we observe here is that when $i \notin S$, β_i stays small and will take many iterations before it even gets to 2α , which is close to 0 since α is chosen to be small. Last lecture, we saw that when $i \in S$, β_i converges to 1. Putting these two together, we can be sure that gradient descent on $L(\beta)$ converges to β^* .

14.3 Main result: Gradient descent on empirical loss

The main part of this lecture focuses on analyzing gradient descent on the empirical risk $\hat{L}(\beta)$. Analyzing the empirical risk is more complicated than analyzing the population risk, so we focus on the case when β^* is 1-sparse, i.e. $r = 1$. (When $r > 1$, the main idea is the same but requires some more advanced analysis techniques.) We restate the case of $r = 1$ in the following theorem.

Theorem 14.4. *Suppose $\eta \geq \tilde{\Omega}(1)$. Then, gradient descent on $\hat{L}(\beta)$ with $t = \Theta\left(\frac{\alpha \log(1/\delta)}{\eta}\right)$ steps satisfies*

$$\|\beta^t \odot \beta^t - \beta^* \odot \beta^*\|_2^2 \leq \tilde{O}\left(\frac{1}{\sqrt{n}}\right). \quad (14.14)$$

Remark 14.5. Note that Theorem 14.4 is a slightly weaker version of Theorem 14.1 for $r = 1$, since the bound on the RHS depends on the number of examples and not the initialization α . In Theorem 14.1, we could take α as small as we like to drive the bound to zero; we cannot do this for Theorem 14.4.

In the proof of this theorem, we will make use of the *restricted isometry property* (r, δ) -RIP (refer to last lecture) as well as the following lemma.

Lemma 14.6. *Suppose $x^{(1)}, x^{(2)}, \dots, x^{(n)}$ satisfy the (r, δ) -RIP condition. Then, $\forall v, w$ such that $\|v\|_0 \leq r$ and $\|w\|_0 \leq r$, we have that*

$$\left| \frac{1}{n} \sum_{i=1}^n \langle x^{(i)}, v \rangle \langle x^{(i)}, w \rangle - \langle v, w \rangle \right| = \left| v^T \left(\frac{1}{n} \sum_{i=1}^n x^{(i)} (x^{(i)})^\top \right) w - \langle v, w \rangle \right| \quad (14.15)$$

$$\leq 4\delta \|v\|_2 \cdot \|w\|_2. \quad (14.16)$$

Corollary 14.7. Taking $w = e_1, \dots, e_d$ in Lemma 14.6, we can conclude that

$$\left\| \frac{1}{n} \sum_{i=1}^n \langle x^{(i)}, v \rangle x^{(i)} - v \right\|_\infty = \left\| \left(\frac{1}{n} \sum_{i=1}^n x^{(i)} (x^{(i)})^\top \right) v - v \right\|_\infty \quad (14.17)$$

$$\leq 4\delta \|v\|_2. \quad (14.18)$$

Specifically, in the Gaussian case when $x^{(1)} \dots x^{(n)} \stackrel{\text{iid}}{\sim} \mathcal{N}(0, I_{d \times d})$ and when $n \geq \tilde{\Omega}(r/\delta^2)$, with high probability the data satisfy the (r, δ) -RIP condition. It follows that when $r = 1$ and $\delta = \tilde{O}(1/\sqrt{n})$, the data are $(1, \delta)$ -RIP. With this, we proceed to prove Theorem 14.4 over the next few subsections.

14.3.1 Computing the gradient update $\nabla \hat{L}(\beta)$

WLOG, assume that $\beta^\star = e_1$. We can decompose the gradient descent iterate β^t as

$$\beta^t = r_t \cdot e_1 + \zeta_t, \quad (14.19)$$

where $\zeta_t \perp e_1$. The idea is to prove convergence to β^\star by showing that (i) $r_t \rightarrow 1$ as $t \rightarrow \infty$, and (ii) $\|\zeta_t\|_\infty \leq O(\alpha)$ for $t \leq \tilde{O}(1/\eta)$. In other words, the *signal* r_t converges quickly to 1 while the *noise* ζ_t remains small for some number of initial iterations. One may be concerned that it is possible for the noise to amplify after many iterations, but we will not have to worry about this scenario if we can guarantee that β^t converges to β^\star first.

We can compute the gradient of $\hat{L}(\beta^t)$ as follows. Since $y^{(i)} = \langle \beta^\star \odot \beta^\star, x^{(i)} \rangle$ and $\beta^t = r_t e_1 + \zeta_t =$

$$r_t \beta^* + \zeta_t,$$

$$\nabla \widehat{L}(\beta^t) = \frac{1}{n} \sum_{i=1}^n (\langle \beta^t \odot \beta^t, x^{(i)} \rangle - y^{(i)}) x^{(i)} \odot \beta^t \quad (14.20)$$

$$= \frac{1}{n} \sum_{i=1}^n (\langle \beta^t \odot \beta^t - \beta^* \odot \beta^*, x^{(i)} \rangle) x^{(i)} \odot \beta^t \quad (14.21)$$

$$= \frac{1}{n} \sum_{i=1}^n \langle r_t^2 \beta^* \odot \beta^* + \zeta_t \odot \zeta_t - \beta^* \odot \beta^*, x^{(i)} \rangle x^{(i)} \odot \beta^t \quad (14.22)$$

$$= \frac{1}{n} \sum_{i=1}^n \underbrace{\langle (r_t^2 - 1) \beta^* \odot \beta^* + \zeta_t \odot \zeta_t, x^{(i)} \rangle}_{m_t} x^{(i)} \odot \beta^t. \quad (14.23)$$

To simplify the analysis, we can rearrange some of the terms that are part of the gradient. Define m_t such that $\nabla \widehat{L}(\beta^t) = m_t \odot \beta^t$. Also, let $X = \frac{1}{n} \sum_{i=1}^n x^{(i)} (x^{(i)})^\top$. Then,

$$m_t = \frac{1}{n} \sum_{i=1}^n \left\langle (r_t^2 - 1) \beta^* \odot \beta^* + \zeta_t \odot \zeta_t, x^{(i)} \right\rangle x^{(i)} \quad (14.24)$$

$$= \left(\frac{1}{n} \sum_{i=1}^n x^{(i)} (x^{(i)})^\top \right) (r_t^2 - 1) \cdot (\beta^* \odot \beta^*) + \left(\frac{1}{n} \sum_{i=1}^n x^{(i)} (x^{(i)})^\top \right) (\zeta_t \odot \zeta_t) \quad (14.25)$$

$$= \underbrace{X(r_t^2 - 1) \cdot (\beta^* \odot \beta^*)}_{\text{part of } u_t} + \underbrace{X(\zeta_t \odot \zeta_t)}_{v_t}. \quad (14.26)$$

Now, define $u_t := (r_t^2 - 1)(\beta^* \odot \beta^*) - X(r_t^2 - 1)(\beta^* \odot \beta^*)$ and $v_t := X(\zeta_t \odot \zeta_t)$. Then we can rewrite the gradient as

$$\nabla \widehat{L}(\beta^t) = m_t \odot \beta^t = [(r_t^2 - 1) \beta^* \odot \beta^* - u_t + v_t] \odot \beta^t. \quad (14.27)$$

Our goal is to show that both u_t and v_t are small, so that $\nabla \widehat{L}(\beta^t)$ is close to its population version $\nabla L(\beta^t)$. Observe that X appears in both u_t and v_t . This matrix is challenging to deal with mathematically because it does not have full rank (because $n < d$). Instead, we rely on the RIP condition to reason about the behavior of X : the idea is that X behaves like the identity for sparse vector multiplication. Applying Corollary 14.7, we can bound $\|u_t\|_\infty$ as

$$\|u_t\|_\infty \leq 4\delta \|(r_t^2 - 1) \beta^* \odot \beta^*\|_2 \leq 4\delta \|\beta^* \odot \beta^*\|_2 \leq 4\delta. \quad (14.28)$$

(In the second inequality, we assume that $|r_t| < 1$. We can do this because r_t starts out at α which is small; if $r_t \geq 1$, then we are already in the regime where gradient descent has converged.) We can bound $\|v_t\|_\infty$ in a similar manner: since Corollary 14.7 implies $\|v_t - \zeta_t \odot \zeta_t\|_\infty \leq 4\delta \|\zeta_t \odot \zeta_t\|_2$,

$$\|v_t\|_\infty \leq \|\zeta_t \odot \zeta_t\|_\infty + 4\delta \|\zeta_t \odot \zeta_t\|_2 \quad (\text{by the triangle inequality}) \quad (14.29)$$

$$\leq \|\zeta_t\|_\infty^2 + 4\delta \|\zeta_t \odot \zeta_t\|_1 \quad (\text{since } \zeta_t \text{ very small}) \quad (14.30)$$

$$= \|\zeta_t\|_\infty^2 + 4\delta \|\zeta_t\|_2^2. \quad (14.31)$$

Note that the size of v_t depends on the size of the noise ζ_t . Thus, by bounding ζ_t (e.g. with a small initialization), we can ensure that v_t is also small. (Ensuring bounds on u_t is more difficult because it depends only on δ .) In the next subsection, we analyze the growth of ζ_t and r_t .

14.3.2 Dynamics analysis of ζ_t

First, we analyze the dynamics of the noise ζ_t , which we want to ensure does not grow too fast.

Lemma 14.8. *For all $t \leq 1/(c\eta\delta)$ with sufficiently large constant c , we have*

$$\|\zeta_t\|_\infty \leq 2\alpha, \quad \|\zeta_t\|_2^2 \leq \frac{1}{2}, \quad \text{and} \quad \|\zeta_{t+1}\|_\infty \leq (1 + O(\eta\delta)) \|\zeta_t\|_\infty. \quad (14.32)$$

Note that this result is weaker than what we were able to show for the population gradient (exponential growth with a small fixed rate), but we will ultimately show that the growth of the signal will be even faster.

Proof. Recall that the empirical gradient (14.27) is $\nabla \hat{L}(\beta) = [(r_t^2 - 1)\beta^\star \odot \beta^\star - u_t + v_t] \odot \beta^t$. Hence, the gradient update to β^t is

$$\beta^{t+1} = \beta^t - \eta [(r_t^2 - 1)\beta^\star \odot \beta^\star - u_t + v_t] \odot \beta^t \quad (14.33)$$

$$= \underbrace{\beta^t - \eta (r_t^2 - 1)\beta^\star \odot \beta^\star \odot \beta^t}_{\text{GD update for population loss}} - \eta (-u_t + v_t) \odot \beta^t. \quad (14.34)$$

Recall that ζ_{t+1} is simply β^{t+1} except for the first coordinate (where it has a zero instead of r_{t+1}), i.e. ζ_{t+1} is the projection of β^{t+1} onto the subspace orthogonal to e_1 . Hence,

$$\zeta_{t+1} = (I - e_1 e_1^\top) \beta^{t+1} \quad (14.35)$$

$$= (I - e_1 e_1^\top) \cdot \beta^t - \eta (I - e_1 e_1^\top) (v_t - u_t) \odot \beta^t \quad (\text{by (14.34), second term} = 0) \quad (14.36)$$

$$= \zeta_t - \eta [(I - e_1 e_1^\top) (v_t - u_t) \odot (I - e_1 e_1^\top) \beta^t] \quad (\text{by distribution law for } \odot) \quad (14.37)$$

$$= \zeta_t - \underbrace{\eta [(I - e_1 e_1^\top) (v_t - u_t)]}_{\rho_t} \odot \zeta_t. \quad (14.38)$$

If we define ρ_t such that $\zeta_{t+1} = \zeta_t - \eta \rho_t \odot \zeta_t$, then the growth of ζ_t is dictated by the size of ρ_t . We can bound this as

$$\|\zeta_{t+1}\|_\infty \leq (1 + \eta \|\rho_t\|_\infty) \|\zeta_t\|_\infty. \quad (14.39)$$

Now, we will prove the lemma by using strong induction on t . Suppose that the first two pieces of (14.32) hold for all iterations up to t . We can show that

$$\|\rho_t\|_\infty \leq \|u_t\|_\infty + \|v_t\|_\infty \quad (14.40)$$

$$\leq 4\delta + \|\zeta_t\|_\infty^2 + 4\delta \|\zeta_t\|_2^2 \quad (\text{by (14.28) and (14.31)}) \quad (14.41)$$

$$\leq 4\delta + (2\alpha)^2 + 4\delta \cdot \frac{1}{2} \quad (\text{by the inductive hypothesis}) \quad (14.42)$$

$$\leq 8\delta, \quad (14.43)$$

where the last step holds because we can take α to be arbitrarily small (e.g. $\alpha \leq \text{poly}(1/n) \leq O(\delta)$). Plugging this into (14.39), we have

$$\|\zeta_{t+1}\|_\infty \leq (1 + 8\eta\delta) \|\zeta_t\|_\infty = (1 + O(\eta\delta)) \|\zeta_t\|_\infty, \quad (14.44)$$

which proves the third piece of the lemma. Using this piece, we can show that

$$\|\zeta_{t+1}\|_\infty \leq (1 + 8\eta\delta)^{t+1} \|\zeta_0\|_\infty \leq (1 + 8\eta\delta)^{1/(c\eta\delta)} \cdot \alpha \leq 2\alpha \quad (14.45)$$

for a sufficiently large constant c , which proves the second piece. Finally, we show that

$$\|\zeta_{t+1}\|_2^2 \leq (1 + 8\eta\delta)^{t+1} \|\zeta_0\|_2^2 \leq (1 + 8\eta\delta)^{1/(c\eta\delta)} \cdot \alpha \sqrt{d} \leq \frac{1}{2}, \quad (14.46)$$

if $\alpha \leq \frac{1}{n^{O(1)}}$, which proves the first piece. □

14.3.3 Dynamics analysis of r_t

Next, we analyze the dynamics of the signal r_t , which we want to show converges to 1.

Lemma 14.9. *For all $t \leq 1/(c\eta\delta)$ with sufficiently large constant c , we have that*

$$r_{t+1} = (1 + \eta(1 - r_t^2))r_t + O(\eta\delta)r_t.$$

Note that the first term on the RHS is r_{t+1} during gradient descent on the population loss, and the second term captures the error.

Proof. Recall that the gradient descent update from the empirical gradient (14.27) is

$$\beta^{t+1} = \beta^t - \eta[(r_t^2 - 1)\beta^\star \odot \beta^\star - u_t + v_t] \odot \beta_t. \quad (14.47)$$

We have that

$$r_{t+1} = \langle \beta^{t+1}, e_1 \rangle \quad (14.48)$$

$$= \langle \beta^t, e_1 \rangle - \eta(r_t^2 - 1)\langle \beta^t, e_1 \rangle - \eta\langle v_t - u_t, e_1 \rangle \langle \beta^t, e_1 \rangle \quad (14.49)$$

$$= r_t - \eta(r_t^2 - 1)r_t - \eta\langle v_t - u_t, e_1 \rangle r_t \quad (14.50)$$

$$= (1 + \eta(1 - r_t^2))r_t + \eta\langle u_t - v_t, e_1 \rangle r_t \quad (14.51)$$

so all we need to do is bound the second term as follows:

$$|\eta\langle v_t - u_t, e_1 \rangle r_t| \leq \eta \cdot r_t \|v_t - u_t\|_\infty \quad (14.52)$$

$$\leq \eta \cdot r_t \cdot 8\delta \quad (\text{by (14.43)}) \quad (14.53)$$

$$= O(\eta\delta) \cdot r_t. \quad (14.54)$$

□

14.3.4 Putting it all together

Finally, we return to the proof of Theorem 14.4. By Lemma 14.9, we know that as long as $r_t \leq 1/2$ it will grow exponentially fast, since

$$r_{t+1} \geq \left(1 + \eta(1 - r_t^2) - O(\eta\delta)\right) \cdot r_t \geq \left(1 + \frac{\eta}{2}\right) \cdot r_t. \quad (14.55)$$

This implies that at some $t_0 = O\left(\frac{\log(1/\alpha)}{\eta}\right)$, we'll observe $r_{t_0} > 1/2$ for the first time. Consider what happens after this point.

- When $1/2 < r_t \leq 1$, we have that

$$1 - r_{t+1} \leq 1 - r_t - \eta(1 - r_t^2)r_t + O(\eta\delta) \cdot r_t \quad (14.56)$$

$$\leq 1 - r_t - \frac{\eta(1 - r_t^2)}{2} + O(\eta\delta) \quad (14.57)$$

$$\leq 1 - r_t - \frac{\eta(1 - r_t)}{2} + O(\eta\delta) \quad (14.58)$$

$$= \left(1 - \frac{\eta}{2}\right)(1 - r_t) + O(\eta\delta). \quad (14.59)$$

Thus, we can achieve $1 - r_{t+1} \leq 2 \cdot \frac{O(n\delta)}{\eta/2} = O(\delta)$ in $\Theta\left(\frac{\log(1/\delta)}{\eta}\right)$ iterations.

- When $r_t > 1$, we can show in a similar manner that

$$r_{t+1} - 1 \leq (1 - \eta)(r_t - 1) + O(\eta\delta) \leq O(\delta), \quad (14.60)$$

implying that r_t remains very close to 1 after the same order of iterations.

This completes the proof of Theorem 14.4, bounding the number of iterations needed for gradient descent on the empirical loss to converge to β^* . \square