

Chapter 9

Implicit/Algorithmic Regularization Effect

One of the miracles of modern deep learning is the phenomenon of “*algorithmic regularization*” (also known as *implicit regularization* or *implicit bias*): although the loss landscape may contain infinitely many global minimizers, many of which do not generalize well, in practice our optimizer (e.g. SGD) tends to recover solutions with good generalization properties.

The focus of this chapter will be to illustrate algorithmic regularization in simple settings. In particular, we first show that gradient descent (with the right initialization) identifies the minimum norm interpolating solution in overparametrized linear regression. Next, we show that for a certain non-convex reparametrization of the linear regression task (where the data is generated from a sparse ground-truth model), gradient descent (again, suitably initialized) approximately recovers a sparse solution (with good generalization). Finally, we discuss algorithmic regularization in the classification setting, and how stochasticity can contribute to algorithmic regularization.

9.1 Implicit regularization effect of zero initialization in overparametrized linear regression ...GD finds the “minimum norm interpolating soln!”

We prove that gradient descent (initialized at the origin) converges to the minimum norm (interpolating solution) (assuming such a solution exists).

Let $X \triangleq [x^{(1)}, \dots, x^{(n)}]^\top \in \mathbb{R}^{n \times d}$ denote our data matrix and $\vec{y} \triangleq [y^{(1)}, \dots, y^{(n)}]^\top \in \mathbb{R}^n$ denote our label vector, where $n < d$. (Assume X is full rank.) Our goal is to find a weight vector β that minimizes our empirical loss function $\hat{L}(\beta) \triangleq \frac{1}{2} \|\vec{y} - X\beta\|_2^2$.

As we are in the overparametrized setting with $n < d$ and X full rank, there exist infinitely many global minimizers that interpolate the data and hence achieve zero loss. In fact, the following lemma shows that the set of global minimizers forms a subspace.

Lemma 9.1. Let X^+ denote the pseudoinverse¹ of X . Then β is a global minimizer if and only if $\beta = X^+ \vec{y} + \zeta$ for some ζ such that $\zeta \perp x_1, \dots, x_n$.

Proof. For any $\beta \in \mathbb{R}^d$, we can decompose it as $\beta = X^+ \vec{y} + \zeta$ for some $\zeta \in \mathbb{R}^d$. Since

$$X\beta = X(X^+ \vec{y} + \zeta) = \vec{y} + X\zeta, \quad \text{If } X\zeta = 0 \text{ do that } \beta \text{ is a global minimizer.} \quad (9.1)$$

β is a global minimizer if and only if $X\zeta = 0$, which happens if and only if $\zeta \perp x_1, \dots, x_n$.

□

¹Since X is full rank, XX^\top is invertible and so we have $X^+ = X^\top (XX^\top)^{-1}$. Note that $XX^+X = X$.

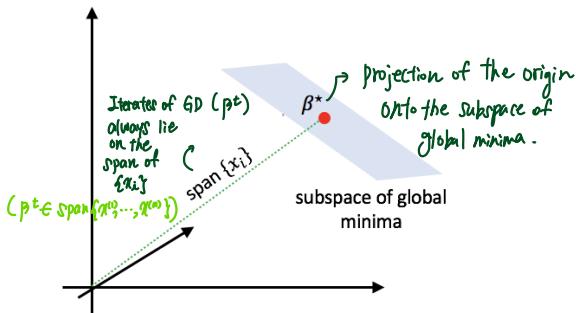


Figure 9.1: Visualization of proof intuition for Theorem 9.3. The solution β^* is the projection of the origin onto the subspace of global minima.

From Lemma 9.1, we can derive an explicit formula for the minimum norm interpolant $\beta^* \triangleq \arg\min_{\beta: \hat{L}(\beta)=0} \|\beta\|_2$.

$$\beta^* \triangleq \arg\min_{\substack{\beta: \hat{L}(\beta)=0 \\ \beta \perp \frac{1}{2}\|\vec{y} - X\beta\|_2^2=0}} \|\beta\|_2$$

minimum norm
interpolating solution.

Corollary 9.2. $\beta^* = X^+ \vec{y}$.

$$X^+ = X^T (X X^T)^{-1}$$

Proof. Take any β such that $\hat{L}(\beta) = 0$, and write $\beta = X^+ \vec{y} + \zeta$. Then from the definition of X^+ and the fact that $X\zeta = 0$ (see the proof of Lemma 9.1), we have

$$\|\beta\|_2^2 = \|X^+ \vec{y}\|_2^2 + \|\zeta\|_2^2 + 2\langle X^+ \vec{y}, \zeta \rangle \quad (9.2)$$

$$\stackrel{X^+ \vec{y} + \zeta}{=} \|X^+ \vec{y}\|_2^2 + \|\zeta\|_2^2 + 2\langle X^T (X X^T)^{-1} \vec{y}, \zeta \rangle \quad (9.3)$$

$$= \|X^+ \vec{y}\|_2^2 + \|\zeta\|_2^2 + 2\langle (X X^T)^{-1} \vec{y}, X \zeta \rangle \quad (9.4)$$

$$= \|X^+ \vec{y}\|_2^2 + \|\zeta\|_2^2 \quad (\text{because } X\zeta = 0) \quad (9.5)$$

$$\geq \|X^+ \vec{y}\|_2^2, \quad \stackrel{\beta = \text{minimum norm soln.}}{\beta(X^+ \vec{y}) = X^+ \vec{y}} \quad (9.6)$$

(with equality if and only if $\zeta = 0$.)

learn β (solution) using GD
with init. β^0
with $\beta^t = \beta^{t-1} - \eta \nabla \hat{L}(\beta^{t-1})$ at iteration t.

□

Now, suppose we learn β using "gradient descent" with initialization β^0 , where at iteration t we set $\beta^t = \beta^{t-1} - \eta \nabla \hat{L}(\beta^{t-1})$ (for some learning rate η). Since $\hat{L}(\beta)$ is convex, we know from standard results in convex optimization that gradient descent will converge to a global minimizer (for a suitably chosen learning rate η (in particular, taking η to be sufficiently small)). Assuming $\beta^0 = 0$, we will in fact recover the minimum norm interpolating solution. β^0 는 수렴하는지 알지 못해. $\beta^0 = 0$ 일 때!

Theorem 9.3. Suppose gradient descent on $\hat{L}(\beta)$ with initialization $\beta^0 = 0$ converges to a solution $\hat{\beta}$ such that $\hat{L}(\hat{\beta}) = 0$. Then $\hat{\beta} = \beta^*$. ("then $\hat{\beta}$ is actually β^* " $\hat{\beta}$ 는 실제로 β^* 이다!)

The main idea of the proof is that the iterates of gradient descent always lie in the span of the $x^{(i)}$'s (see Figure 9.1 for an illustration).

Proof. We first show via induction that $\beta^t \in \text{span}\{x^{(1)}, \dots, x^{(n)}\}$ for all t . For the induction base case, note that $\beta^0 = 0 \in \text{span}\{x^{(1)}, \dots, x^{(n)}\}$. Now suppose $\beta^{t-1} \in \text{span}\{x^{(1)}, \dots, x^{(n)}\}$. Recall that $\beta^t = \beta^{t-1} - \eta \nabla \hat{L}(\beta^{t-1})$. Since left-multiplying any vector by X^T amounts to taking a linear combination of the rows of X , it follows that $\eta \nabla \hat{L}(\beta^{t-1})$ (which is $\eta X^T (X \beta^{t-1} - \vec{y})$) $\in \text{span}\{x^{(1)}, \dots, x^{(n)}\}$ and so $\beta^t = \beta^{t-1} - \eta \nabla \hat{L}(\beta^{t-1}) \in \text{span}\{x^{(1)}, \dots, x^{(n)}\}$. This proves the induction step. $\beta^t \in \text{span}\{x^{(1)}, \dots, x^{(n)}\}$

Next, we show that $\hat{\beta} \in \text{span}\{x^{(1)}, \dots, x^{(n)}\}$ and $\hat{L}(\hat{\beta}) = 0$ implies $\hat{\beta} = \beta^*$. By definition, $\hat{\beta} \in \text{span}\{x^{(1)}, \dots, x^{(n)}\}$ implies $\hat{\beta} = X^T v$ for some $v \in \mathbb{R}^n$. Since $\hat{L}(\hat{\beta}) = 0$, we have $0 = X \hat{\beta} - \vec{y} = X X^T v - \vec{y}$. This implies $v = (X X^T)^{-1} \vec{y}$, and so $\hat{\beta} = X^T v = X^T (X X^T)^{-1} \vec{y} = X^+ \vec{y} = \beta^*$. $\hat{\beta} = X^T v$ $\hat{\beta} = X^+ \vec{y}$ $\hat{\beta} = \beta^*$ \square

$$\begin{aligned} \hat{L}(\beta) &= \frac{1}{2} \|\vec{y} - X\beta\|_2^2 = \frac{1}{2} \|\vec{y} - X\beta^*\|_2^2 && \text{By defn. } X^+ \\ \nabla \hat{L}(\beta^t) &= X^T (X \beta^t - \vec{y}) && 102 \\ \nabla \hat{L}(\beta^{t-1}) &= X^T (X \beta^{t-1} - \vec{y}) \\ \eta \nabla \hat{L}(\beta^{t-1}) &= \eta X^T (X \beta^{t-1} - \vec{y}) \in \text{span}\{x^{(1)}, \dots, x^{(n)}\} \\ \text{The vector } X \beta^{t-1} - \vec{y} &\text{ is left-multiplied by } X^T \\ &\Rightarrow \text{it actually is just a lin. comb of the rows of } X. \end{aligned}$$

9.2 Implicit regularization of small initialization in nonlinear models⁶⁶

We give another example of implicit regularization effect of small initialization in a non-convex version (of the overparametrized linear regression task considered in the previous section). The results in this subsection are largely simplifications of the paper Li et al. [2017] which studies over-parameterized compressed sensing and two-layer neural nets with quadratic activation.)

We assume $x^{(1)}, \dots, x^{(n)}$ iid $\mathcal{N}(0, I_{d \times d})$ and $y^{(i)} = f_{\beta^*}(x^{(i)})$, where the ground truth vector β^* is r -sparse (i.e. $\|\beta^*\|_0 = r$). For simplicity, we assume $\beta_i^* = \mathbf{1}\{i \in S\}$ for some $S \subseteq [d]$ such that $|S| = r$. We again analyze the overparametrized setting, where this time $n \ll d$ but also $n \geq \Omega(r^2)$.

Our goal is to find a weight vector that minimizes our empirical loss function

$$\text{Empirical loss function} \quad \widehat{L}(\beta) \triangleq \frac{1}{4n} \sum_{i=1}^n \underbrace{\left(\underbrace{y^{(i)}}_{\text{obs}} - \underbrace{f_\beta(x^{(i)})}_{\langle \beta \odot \beta, x^{(i)} \rangle} \right)^2}_{\text{Hadamard product: } (\beta \odot \beta)_j \triangleq \beta_j \beta_j \text{ for } j=1, \dots, d}, \quad (9.7)$$

where $f_\beta(x) \triangleq \langle \beta \odot \beta, x \rangle$. (The operation \odot denotes the Hadamard product: for $u, v \in \mathbb{R}^d$, $u \odot v \in \mathbb{R}^d$ is defined by $(u \odot v)_i \triangleq u_i v_i$ for $i = 1, \dots, d$.)

9.2.1 Main results of algorithmic regularization

Note that while f_β is still linear over x , our loss is no longer convex over β . (To see this, suppose $\beta \neq 0$ is a global minimizer. Then we have $\widehat{L}(0) > \widehat{L}(\beta) = \widehat{L}(-\beta)$.) Thus, the effect of algorithmic regularization induced by gradient descent will be much different from the overparametrized linear regression setting.

In the previous setting of linear regression, solutions with low ℓ_2 norm are desirable as they tend to generalize well. In the present setting, we know our ground-truth parameter β^* is sparse. Thus, we want to learn a sparse solution β , avoiding non-sparse solutions that may not generalize well. One approach to finding sparse solutions, called *lasso regression*, is to minimize the ℓ_1 -regularized proxy loss

$$\text{Lasso regression: minimize}_{\beta} \quad \left(\sum_{i=1}^n \left(\langle \theta, x^{(i)} \rangle - y^{(i)} \right)^2 + \lambda \|\theta\|_1 \right) \quad (9.8)$$

with respect to θ where $\theta = \beta \odot \beta$. However, it turns out that we can equivalently learn a sparse solution by running gradient descent from a suitable initialization on the original unregularized loss.)

To be specific, let $\beta^0 = \alpha \mathbf{1} \in \mathbb{R}^d$ be the initialization (where α is a small positive number.) The update rule of gradient descent algorithm is given by $\beta^{t+1} = \beta^t - \eta \nabla \widehat{L}(\beta^t)$. The next theorem shows that when $n = \tilde{\Omega}(r^2)$, gradient descent on $\widehat{L}(\beta)$ converges to β^* .

Theorem 9.4. 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

$$\left\| \beta^\top \odot \beta^\top - \beta^* \odot \beta^* \right\|_2^2 \leq O\left(\alpha \sqrt{d}\right).$$

expectation of
recovery error (test error)

$\beta^* \text{ is } r\text{-sparse}$
 $\|\beta^*\|_0 = r, |S|=r$
small positive number,
initialization $\beta^0 = \alpha \mathbf{1} \in \mathbb{R}^d$
 $\beta^{t+1} = \beta^t - \eta \nabla \widehat{L}(\beta^t)$
universal constant.
 $\beta_i^* = \mathbf{1}\{i \in S\}$
for $i \in [d]$
 $\text{s.t. } |S|=r$
take it as
small inverse polynomial of d
so that the lower bound of T
doesn't change much.
upper bound:
 $\text{depends on } \frac{1}{\alpha}$
which is really big when
 c is really big enough

(Here, T indexes the gradient descent steps.)

We make several remarks about Theorem 9.4 before presenting the proof.

(But we use $\beta^0 = \alpha \mathbf{1} \in \mathbb{R}^d$)

Remark 9.5. In this problem we do not use $\beta^0 = 0$ as the initialization point because $\beta = 0$ is a critical point, that is, $\nabla \widehat{L}(0) = 0$. Note that the lower bound on T depends logarithmically on $1/\alpha$, so we can take α to be a small inverse polynomial on d and the lower bound won't change much. Also, the upper bound depends polynomially on $1/\alpha$ (which is considered very big when c is sufficiently large), so we do not need to use early stopping in a serious way.

Remark 9.6. Theorem 9.4 is a simplified version of Theorem 1.1 in [Li et al., 2018].

Remark 9.7. $\widehat{L}(\beta)$ has many global minima. To see this, observe that the number of parameters is d and the number of constraints to fit all the examples is $O(n)$ because there are only n examples. Recall that for overparameterized model we have $d \gg n$; consequently, there exists many global minima of $\widehat{L}(\beta)$.

Remark 9.8. β^* is the min-norm solution in this case. That is,

$$\text{Ground truth param} \quad \beta^* = \operatorname{argmin}_{\beta} \|\beta\|_2^2 \quad \text{s.t. } \widehat{L}(\beta) = 0. \quad (9.10)$$

Informally, this is because we can view $\beta \odot \beta$ as a vector $\theta \in \mathbb{R}^d$, which leads to $\|\beta\|_2^2 = \|\theta\|_1$. Then in the θ space (and with a little abuse of notation), the optimization problem (9.10) becomes

$$\text{View } \beta \odot \beta \text{ as a vector } \theta \in \mathbb{R}^d, \|\beta\|_2^2 = \|\theta\|_1, \quad \theta^* = \operatorname{argmin}_{\theta} \|\theta\|_1 \quad \text{s.t. } \widehat{L}(\theta) = 0, \quad (9.11)$$

which is a lasso regression, whose solution is sparse.

Remark 9.9. In this non-linear case and the linear case before, gradient descent with small initialization converges to minimum ℓ_2 -norm solution. (Similarly, in the NTK regime, gradient descent converges to a solution that is very close to the initialization.) Therefore, it seems conceivable that GD generally prefers global minima nearest to the initialization. However, we do not have a general theorem for this phenomenon (and the instructor also believes that this is not universally true without other conditions). \downarrow?

9.2.2 Ground work for proof and the restricted isometry property

In this section we prepare the ground work for the proof of Theorem 9.4.

We start by showing several basic properties about $\widehat{L}(\beta)$. Note that for any fixed vector $v \in \mathbb{R}^d$ and $x \in \mathbb{R}^d$ (when x is drawn from $\mathcal{N}(0, I)$), we have $\text{Our empirical loss function}$

$$\mathbb{E} [\langle x, v \rangle^2] = \mathbb{E} [v^\top x x^\top v] = v^\top \mathbb{E} [x x^\top] v = \|v\|_2^2. \quad (9.12)$$

$x^{(1)}, \dots, x^{(n)} \text{ iid } \mathcal{N}(0, \text{Id}_d)$

It follows that

$$L(\beta) = \frac{1}{4} \mathbb{E}_{x \sim \mathcal{N}(0, I)} [(y - \langle \beta \odot \beta, x \rangle)^2] \quad (9.9 \text{ 2nd}) \quad (9.13)$$

$$= \frac{1}{4} \mathbb{E}_{x \sim \mathcal{N}(0, I)} [\langle \beta^* \odot \beta^* - \beta \odot \beta, x \rangle^2] \quad \begin{array}{l} y^{(i)} = f_{\beta^*}(x^{(i)}) \\ y = f_{\beta^*}(x), f_p(x) = \langle p \odot p, x \rangle \end{array} \quad (\text{by definition of } y) \quad (9.14)$$

$$= \frac{1}{4} \|\beta^* \odot \beta^* - \beta \odot \beta\|_2^2 \cdot \mathbb{E} [\langle z, v \rangle^2] = \|v\|_2^2 \quad \begin{array}{l} f_{\beta^*}(x) = \langle \beta^* \odot \beta^*, x \rangle \\ z \text{ is drawn from } \mathcal{N}(0, I) \end{array} \quad (\text{by (9.12)}) \quad (9.15)$$

How close β is to the ground-truth parameter β^* ? (see (9.9)).

Note that (9.15) is the metric that we use to characterize how close β is to the ground-truth parameter β^* (see (9.9)).

In the following lemma we show that $\widehat{L}(\beta) \approx L(\beta)$ by uniform convergence. (Generally speaking, uniform convergence of the loss function for all β requires $n \geq \Omega(d)$ samples, so in our setting (where $n \ll d$) $\widehat{L}(\beta) \approx L(\beta)$ does not always hold.) However, since we assume β^* is sparse, the analysis only requires uniform convergence for sparse vectors. \downarrow .

Lemma 9.10. Assume $n \geq \tilde{\Omega}(r^2)$. (With high probability over the randomness in $x^{(1)}, \dots, x^{(n)}$,) $\forall v$ such that $\|v\|_0 \leq r$ we have $\exists v$ s.t. $\|v\|_1 \leq r$ (where $n \geq \Omega(r^2)$), $\left(\begin{array}{l} \text{(r, \delta)-RIP condition} \\ \text{with } \{x^{(1)}, \dots, x^{(n)}\} \end{array} \right)$

$$(1 - \delta) \|v\|_2^2 \leq \frac{1}{n} \sum_{i=1}^n \langle v, x^{(i)} \rangle^2 \leq (1 + \delta) \|v\|_2^2. \quad (9.16)$$

Lemma 9.10 is a special case of Lemma 2.2 in [Li et al., 2018] so the proof is omitted here. We say the set $\{x^{(1)}, \dots, x^{(n)}\}$ (or $X = [x^{(1)}, \dots, x^{(n)}]$) satisfies (r, δ) -RIP condition (restricted isometric property) if (9.16) holds.

9.10을 만족하는 X : RIP를 만족한다~고 한다.

RIP condition의 consequence:

$\frac{1}{n} \sum_{i=1}^n \langle x^{(i)}, v \rangle^2$ behaves like $\|v\|_2^2$,

104 for approximately low rank N .

$$\text{(9.16)
if } v \text{ s.t. } \|v\|_0 \leq r \text{ (where } n \geq \delta(r^2) \text{) , } \\ (1 - \delta)\|v\|_2^2 \leq \frac{1}{n} \sum_{i=1}^n \langle v, x^{(i)} \rangle^2 \leq (1 + \delta)\|v\|_2^2. \quad \left. \begin{array}{l} \text{(r, } \delta\text{-RIP condition)} \\ \text{with } \{x^{(1)}, \dots, x^{(n)}\} \end{array} \right.$$

By algebraic manipulation, (9.16) is equivalent to

$$\text{(9.17)
} (1 - \delta)\|v\|_2^2 \leq \underbrace{\langle v^\top \left(\frac{1}{n} \sum_{i=1}^n x^{(i)}(x^{(i)})^\top \right) v \rangle}_{\approx 1} \leq (1 + \delta)\|v\|_2^2. \quad \left. \begin{array}{l} \text{IMSE (2)} \\ (1 - \delta)\|v\|_0\|w\|_0 \leq v^\top \left(\frac{1}{n} \sum_{i=1}^n x^{(i)}(x^{(i)})^\top \right) w \leq (1 + \delta)\|v\|_0\|w\|_0 \end{array} \right.$$

In other words, from the point of view of a sparse vector v we have $\sum_{i=1}^n x^{(i)}(x^{(i)})^\top \approx I$. (Note however that $\sum_{i=1}^n x^{(i)}(x^{(i)})^\top$ is not close to $I_{d \times d}$ in other notions of closeness.) For example, $\sum_{i=1}^n x^{(i)}(x^{(i)})^\top$ is not close to $I_{d \times d}$ in spectral norm. Another way to see this is that $\sum_{i=1}^n x^{(i)}(x^{(i)})^\top$ is a $d \times d$ matrix but only has rank $n \ll d$.)

As a result, with the RIP condition we have $\hat{L}(\beta) \approx L(\beta)$ if β is sparse. With more tools we can also get $\nabla \hat{L}(\beta) \approx \nabla L(\beta)$. Let us define the set $S_r = \{\beta : \|\beta\|_0 \leq O(r)\}$, the set where we have uniform convergence of \hat{L} and $\nabla \hat{L}$. Informally, as long as we are in the set S_r , \hat{L} and $\nabla \hat{L}$ have similar behavior to their population counterparts. (Note, on the other hand, that there exists a dense $\beta \notin S_r$ such that $\hat{L}(\beta) = 0$ but $L(\beta) \gg 0$.)

The RIP condition also gives us the following lemma which will be needed for the proof of Theorem 9.4.

Lemma 9.11. 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

RIP condition's consequence:

$\frac{1}{n} \sum_{i=1}^n \langle x^{(i)}, v \rangle \langle x^{(i)}, w \rangle - \langle v, w \rangle$ behave like V ,
for approximately low rank V .
(in 행렬 풍선 (v vector 끝))

$$\left| \frac{1}{n} \sum_{i=1}^n \langle x^{(i)}, v \rangle \langle x^{(i)}, w \rangle - \langle v, w \rangle \right| = \left| v^\top \left(\frac{1}{n} \sum_{i=1}^n x^{(i)}(x^{(i)})^\top \right) w - \langle v, w \rangle \right| \quad \text{(9.18)}$$

$$\leq 4\delta \|v\|_2 \cdot \|w\|_2 \cdot \dots \quad \text{(9.19)}$$

$$\text{range: } 2\delta \|v\|_2 \|w\|_2 \quad \text{(9.19)}$$

Corollary 9.12. Taking $w = e_1, \dots, e_d$ in Lemma 9.11, we can conclude that

$$\left\| \frac{1}{n} \sum_{i=1}^n \langle x^{(i)}, v \rangle x^{(i)} - v \right\|_\infty = \left\| \left(\frac{1}{n} \sum_{i=1}^n x^{(i)}(x^{(i)})^\top \right) v - v \right\|_\infty \quad \text{(9.20)}$$

$$\left\| \frac{1}{n} \sum_{i=1}^n \langle x^{(i)}, v \rangle x^{(i)} - v \right\|_\infty \leq 4\delta \|v\|_2. \quad \left\| \left(\frac{1}{n} \sum_{i=1}^n x^{(i)}(x^{(i)})^\top \right) v - v \right\|_\infty \quad \text{(9.21)}$$

9.2.3 Warm-up for analysis: Gradient descent on population loss

The main intuition for proving Theorem 9.4 is to leverage the uniform convergence when β belongs to the set S_r (see Figure 9.2). Note that the initialization β^0 is not exactly r -sparse, but taking α to be sufficiently small, β^0 is approximately 0-sparse. The proof is decomposed into the following steps:

- Step 1: thm 9.13 (In part 9.23)
 - Gradient descent on $L(\beta)$ converges to β^* without leaving S_r , and
- Step 2: thm 9.14 (In part 9.24)
 - Gradient descent on $\hat{L}(\beta)$ is similar to gradient descent on $L(\beta)$ inside S_r .

Combining the two steps we can show that gradient descent on $\hat{L}(\beta)$ does not leave S_r and converges to β^* .

As a warm up, we prove the following theorem for gradient descent on $L(\beta)$.

Theorem 9.13. For sufficiently small η , gradient descent on $L(\beta)$ converges to β^* (in $\Theta\left(\frac{\log(1/(\epsilon\alpha))}{\eta}\right)$ iteration) (with ϵ -error in ℓ_2 -distance.)

Proof. Since

$$\begin{aligned} & \text{(9.15) } L(\beta) \text{ is } \\ & \text{(9.16) } \nabla L(\beta) = (\beta \odot \beta - \beta^* \odot \beta^*) \odot \beta, \end{aligned} \quad \left. \begin{array}{l} \text{population version } \nabla L(\beta) = (\beta \odot \beta - \beta^* \odot \beta^*) \odot \beta \\ \nabla L(\beta) = (\beta \odot \beta - \beta^* \odot \beta^*) \odot \beta, \end{array} \right. \quad \text{(9.22)}$$

the gradient descent step is

$$\begin{aligned} \beta^{t+1} &= \beta^t - \eta \nabla L(\beta) \\ &= \beta^t - \eta(\beta^t \odot \beta^t - \beta^* \odot \beta^*) \odot \beta^t. \end{aligned}$$

For η not large, GD on $L(\beta)$ converges to β^* in $\Theta\left(\frac{\log(1/\epsilon\alpha)}{\eta}\right)$ iteration with ϵ -error in ℓ_2 -distance.

GD update for population loss:

$$\beta^{t+1} = \beta^t - \eta(\beta^t \odot \beta^t - \beta^* \odot \beta^*) \odot \beta^t$$

$$\begin{aligned} \nabla L(\beta) &= (\beta \odot \beta - \beta^* \odot \beta^*) \odot \beta \\ \beta^{t+1} &= \beta^t - \eta L(\beta) \\ \beta^{t+1} &= \beta^t - \eta(\beta^t \odot \beta^t - \beta^* \odot \beta^*) \odot \beta^t \end{aligned}$$

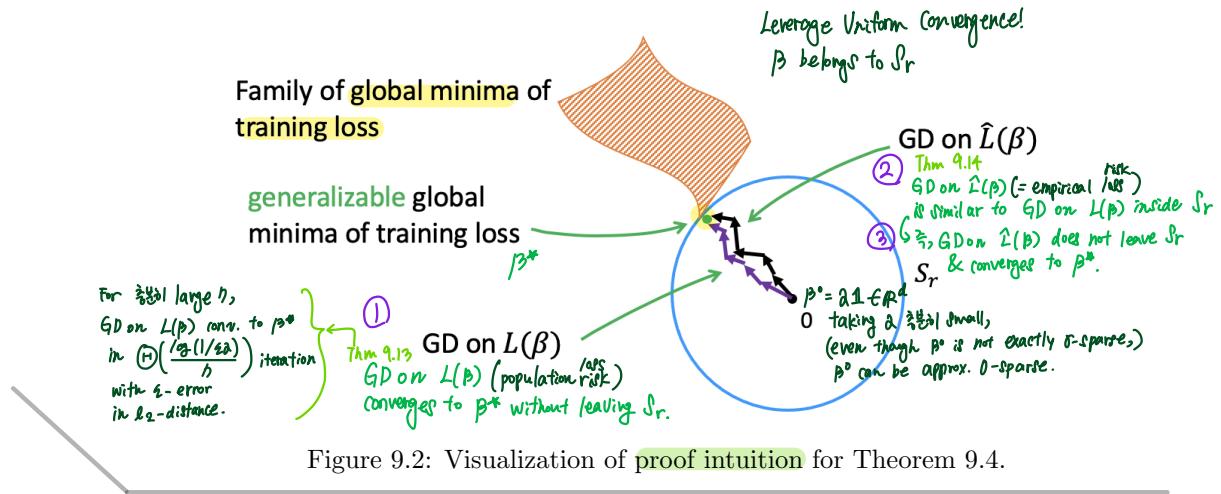


Figure 9.2: Visualization of proof intuition for Theorem 9.4.

Recall that $\beta^* = \mathbf{1}\{i \in S\}$ and $\beta^0 = \alpha \mathbf{1}$, and the update rule above decouples across the coordinates of β^t . Thus, we only need to show that $(\beta_i^* - \beta_i^t) \leq \epsilon$ for the number of iterations stated in the Theorem.

Case 1: $i \in S$. For $i \in S$, the update rule for coordinate i is

$$\nabla L(\beta) = (\beta^0 \otimes \beta - \beta^* \otimes \beta^*) \otimes \beta \quad \beta_i^{t+1} = \beta_i^t - \eta(\beta_i^t \cdot \beta_i^t - 1 \cdot 1) \cdot \beta_i^t \quad (9.24)$$

$$\begin{aligned} \beta^{t+1} &= \beta^t - \eta(\beta^0 \otimes \beta^t - \beta^* \cdot \beta^*) \otimes \beta^t \\ &= \beta^t - \eta[(\beta^t)^2 - 1] \beta^t. \end{aligned} \quad (9.25)$$

Consider the following two cases:

- If $\beta_i^t \leq 1/2$, we have

$$-(\beta_i^t)^2 \geq 1/4 \quad \beta_i^{t+1} = \beta_i^t \left[1 + \eta \left(1 - (\beta_i^t)^2 \right) \right] \quad (9.26)$$

$$\geq \beta_i^t \left(1 + \frac{3}{4}\eta \right). \quad \text{If } \beta_i^{t+1} \text{ grows exponentially.} \quad (9.27)$$

$$(1+\eta)^t \approx e$$

$$(1+\eta)^{t+1} \approx e^c$$

Consequently, β_i^{t+1} grows exponentially, and it takes $\Theta\left(\frac{\log(1/\alpha)}{\eta}\right)$ iterations for β_i^t to grow from α to at least $1/2$.² This will bring us into the second case.

- if $\beta_i^t \geq 1/2$, we have

$$1 - \beta_i^{t+1} = 1 - \beta_i^t + \eta \left[(\beta_i^t)^2 - 1 \right] \beta_i^t \quad (9.28)$$

$$= 1 - \beta_i^t - \eta(1 - \beta_i^t)(1 + \beta_i^t) \beta_i^t \quad (9.29)$$

$$\leq (1 - \beta_i^t) - \eta(1 - \beta_i^t) \beta_i^t \quad (\text{because } 1 + \beta_i^t \geq 1) \quad (9.30)$$

$$= (1 - \beta_i^t)(1 - \eta \beta_i^t) \quad (9.31)$$

$$\leq (1 - \beta_i^t)(1 - \eta/2). \quad (\text{because } \beta_i^t \geq 1/2) \quad (9.32)$$

Therefore it takes $\Theta\left(\frac{\log(1/\epsilon)}{\eta}\right)$ iterations to achieve $1 - \beta_i^t \leq \epsilon$.

Case 2: $i \notin S$. For all $i \notin S$, we claim (informally) that it is sufficient to show that when $t \leq 1/(10\eta\alpha^2)$, $\beta_i^t \leq 2\alpha$. This is because 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.

²This is because $(1 + \eta)^{1/\eta} \approx e$, so $(1 + \eta)^{c/\eta} \approx e^c$.

β_i^t is small when
(less iterations)
(β_i^t doesn't get to 2α)
↓
 β_i^t stays small,
 $\beta_i^t \leq 2\alpha$

it is sufficient to show that
when $t \leq 1/(10\eta\alpha^2)$, $\beta_i^t \leq 2\alpha$

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

$$\nabla L(\beta) = (\beta \odot \beta^t - \beta^* \odot \beta^*) \odot \beta \quad \beta_i^{t+1} = [\beta^t - \eta(\beta^t \odot \beta^t - \beta^* \odot \beta^*) \odot \beta^t]_i = 0 \quad (9.33)$$

$$\beta^{t+1} = \beta^t - \eta L(\beta) \quad \beta_i^{t+1} = \beta_i^t - \eta(\beta_i^t \cdot \beta_i^t) \cdot \beta_i^t \quad \text{(since } \beta_i^* = 0 \forall i \notin S\text{)} \quad (9.34)$$

$$\beta^{t+1} = [\beta^t - \eta(\beta^t \odot \beta^t - \beta^* \odot \beta^*) \odot \beta^t]_i = \beta_i^t - \eta(\beta_i^t)^3 \cdot 0 \quad (9.35)$$

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 $s \leq t$,

$$\text{"It is sufficient to show that when } t \leq 1/(10\eta\alpha^2), \beta_i^t \leq 2\alpha \text{"} \quad \beta_i^{s+1} = (1 - \eta(\beta_i^s)^2)\beta_i^s \quad (9.36)$$

$$\leq (1 + \eta(\beta_i^s)^2)\beta_i^s \quad \text{By induction hypothesis.} \quad (9.37)$$

$$\leq (1 + 4\eta\alpha^2)\beta_i^s. \quad \text{(since } \beta_i^s \leq 2\alpha\text{)} \quad (9.38)$$

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

$$\beta_i^{t+1} \leq \beta_0 \cdot (1 + 4\eta\alpha^2)^t \quad (9.39)$$

$$\leq \beta_0(1 + 4\eta\alpha^2)^{\frac{1}{10\eta\alpha^2}} \quad (9.40)$$

$$\leq \beta_0 \exp\left(\frac{4\eta\alpha^2}{10\eta\alpha^2}\right) \quad (9.41)$$

$$= \beta_0 \cdot e^{2/5} \quad (9.42)$$

$$\leq 2\alpha, \quad \beta_0 = 2\alpha \quad (9.43)$$

which completes the inductive proof of the claim.

Theorem 9.13 \square

9.2.4 Proof of main result: gradient descent on empirical loss

Analyzing gradient descent on the empirical risk $\hat{L}(\beta)$ is more complicated than analyzing gradient descent on 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.)

Note that in our setup, i.e. when $x^{(1)}, \dots, x^{(n)}$ iid $\mathcal{N}(0, I_{d \times d})$ 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. This will allow us to use the lemmas involving the RIP condition for the proof.

We restate the case of $r = 1$ in the following theorem.

Theorem 9.14. 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

$$\|\beta^t \odot \beta^t - \beta^* \odot \beta^*\|_2^2 \leq \tilde{O}\left(\frac{1}{\sqrt{n}}\right). \quad \begin{matrix} \text{Thm 9.14.} \\ \text{GD on } \hat{L}(\beta) \text{ with } t = \Theta\left(\frac{\alpha \log(1/\delta)}{\eta}\right) \text{ satisfies} \\ \|\beta^t \odot \beta^t - \beta^* \odot \beta^*\|_2^2 \leq \tilde{O}\left(\frac{1}{\sqrt{n}}\right). \end{matrix} \quad (9.44)$$

Remark 9.15. Note that Theorem 9.14 is a slightly "weaker" version of Theorem 9.4 for $r = 1$, since the bound on the RHS depends on the number of examples and not the initialization α . (In Theorem 9.4, we could take α as small as we like to drive the bound to zero; we cannot do this for Theorem 9.14.)

We proceed to prove Theorem 9.14 with the follow steps:

1. Computing the gradient update $\nabla \hat{L}(\beta)$,
2. Dynamics analysis of noise ζ_t , \Leftarrow we want it "Does not grow too fast"
3. Dynamics analysis of signal r_t , and \Leftarrow we want.. "Converges to 1"

4. Putting it all together.

(Step 1.)

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

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

$$\beta^* = \mathbf{1} \{ \text{first} \} \quad \beta^t = r_t \cdot e_1 + \zeta_t, \quad \text{signal} \quad \text{noise} \quad (9.45)$$

where $\zeta_t \perp e_1$. The idea is to prove convergence to β^* 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^* first.

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

$$\nabla \hat{L}(\beta^t) = \frac{1}{n} \sum_{i=1}^n (\langle \beta^t \odot \beta^t, x^{(i)} \rangle - y^{(i)}) x^{(i)} \odot \beta^t \quad (9.46)$$

$$= \frac{1}{n} \sum_{i=1}^n (\langle \beta^t \odot \beta^t - \beta^* \odot \beta^*, x^{(i)} \rangle) x^{(i)} \odot \beta^t \quad (9.47)$$

$$= \frac{1}{n} \sum_{i=1}^n \langle r_t^2 \beta^* \odot \beta^* + \zeta_t \odot \zeta_t - \beta^* \odot \beta^*, x^{(i)} \rangle x^{(i)} \odot \beta^t \quad (9.48)$$

$$= \frac{1}{n} \sum_{i=1}^n \underbrace{\left\langle (r_t^2 - 1) \beta^* \odot \beta^* + \zeta_t \odot \zeta_t, x^{(i)} \right\rangle}_{m_t} x^{(i)} \odot \beta^t. \quad (9.49)$$

$$\nabla \hat{L}(\beta^t) = \frac{1}{n} \sum_{i=1}^n \left\langle (r_t^2 - 1) \beta^* \odot \beta^* + \zeta_t \odot \zeta_t, x^{(i)} \right\rangle q^{(i)} \odot p^t$$

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

$$\text{gradient} = m_t \odot \beta^t \quad m_t = \frac{1}{n} \sum_{i=1}^n \left\langle (r_t^2 - 1) \beta^* \odot \beta^* + \zeta_t \odot \zeta_t, x^{(i)} \right\rangle x^{(i)} \quad (9.50)$$

$$= \left(\frac{1}{n} \sum_{i=1}^n x^{(i)} (x^{(i)})^\top \right) (r_t^2 - 1) \cdot (\beta^* \odot \beta^*) + \left(\frac{1}{n} \sum_{i=1}^n x^{(i)} (x^{(i)})^\top \right) (\zeta_t \odot \zeta_t) \quad (9.51)$$

$$\text{gradient } \nabla \hat{L}(\beta^t) = m_t \odot \beta^t \quad = X(r_t^2 - 1) \cdot (\beta^* \odot \beta^*) + X(\zeta_t \odot \zeta_t). \quad (9.52)$$

$$\text{Define } u_t \triangleq (r_t^2 - 1)(\beta^* \odot \beta^*) \quad v_t \triangleq X(\zeta_t \odot \zeta_t) \quad \text{part of } u_t \quad \text{part of } v_t \quad \text{except for...}$$

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

$$\nabla \hat{L}(\beta^t) = m_t \odot \beta^t = [(r_t^2 - 1)(\beta^* \odot \beta^*) - u_t + v_t] \odot \beta^t. \quad (9.53)$$

Our goal is to show that both u_t and v_t are small, so that $\nabla \hat{L}(\beta^t)$ 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 9.12, we can bound $\|u_t\|_\infty$ as

$$\begin{aligned} \|u_t\|_\infty &\leq 4\delta \| (r_t^2 - 1)(\beta^* \odot \beta^*) \|_2 \\ &\leq 4\delta \| (r_t^2 - 1)(\beta^* \odot \beta^*) \|_2 \quad * \text{Goal: } u_t \text{ is small} \\ &\leq 4\delta \| (r_t^2 - 1)(\beta^* \odot \beta^*) \|_2 \quad \|r_t\| < 1 \\ &\leq 4\delta \| (r_t^2 - 1)(\beta^* \odot \beta^*) \|_2 \quad \text{In the second inequality, we assume that } |r_t| < 1. \end{aligned} \quad (9.54)$$

(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

$\|v_t\|_\infty$ bounded, remains very close to 1.

$\|v_t\|_\infty$ in a similar manner: since Corollary 9.12 implies $\|v_t - \zeta_t \odot \zeta_t\|_\infty \leq 4\delta \|\zeta_t \odot \zeta_t\|_2$,

$$\|v_t\|_\infty \leq \|\zeta_t \odot \zeta_t\|_\infty + 4\delta \|\zeta_t \odot \zeta_t\|_2 \quad (\text{by the triangle inequality}) \quad (9.55)$$

$$\begin{aligned} &\leq \|\zeta_t\|_\infty^2 + 4\delta \|\zeta_t \odot \zeta_t\|_2 \\ &= \|\zeta_t\|_\infty^2 + 4\delta \|\zeta_t\|_2^2. \end{aligned} \quad (\text{since } \zeta_t \text{ very small}) \quad (9.56)$$

$$= \|\zeta_t\|_\infty^2 + 4\delta \|\zeta_t\|_2^2. \quad (9.57)$$

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 two subsections, we analyze the growth of ζ_t and r_t .

(Step 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 9.16. For all $t \leq 1/(c\eta\delta)$ with sufficiently large constant c , we have

$$\begin{aligned} \textcircled{1} \quad \|\zeta_t\|_\infty &\leq 2\alpha, \\ \textcircled{2} \quad \|\zeta_t\|_2^2 &\leq \frac{1}{2}, \quad \text{and} \quad \textcircled{3} \quad \|\zeta_{t+1}\|_\infty \leq (1 + O(\eta\delta)) \|\zeta_t\|_\infty. \end{aligned} \quad (9.58)$$

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 r_t will be even faster.

Proof. Recall that the empirical gradient (9.53) is $\nabla \hat{L}(\beta) = [(r_t^2 - 1)\beta^* \odot \beta^* - u_t + v_t] \odot \beta^t$. Hence, the gradient update to β^t is

$$\beta^{t+1} = \beta^t - \eta \nabla \hat{L}(\beta) \quad (9.59)$$

$$\beta^{t+1} = \beta^t - \eta [(r_t^2 - 1)\beta^* \odot \beta^* - u_t + v_t] \odot \beta^t \quad (9.60)$$

GD update for population loss
 $\underbrace{(\beta^{t+1} \text{ for population loss } L(\beta))}_{(9.23)} = 0$

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.

$$\zeta_{t+1} = (I - e_1 e_1^\top) \beta^{t+1} \quad (\text{projection of } \beta^{t+1} \text{ onto the subspace orthogonal to } e_1) \quad (9.61)$$

$$\zeta_{t+1} = (I - e_1 e_1^\top) \cdot \beta^t - \eta (I - e_1 e_1^\top) (v_t - u_t) \odot \beta^t \quad (\text{by (9.60), second term } = 0) \quad (9.62)$$

$$\zeta_{t+1} = \zeta_t - \eta [(I - e_1 e_1^\top) (v_t - u_t) \odot (I - e_1 e_1^\top) \beta^t] \quad (\text{by distribution law for } \odot) \quad (9.63)$$

$$\zeta_{t+1} = \zeta_t - \eta [(I - e_1 e_1^\top) (v_t - u_t)] \odot \zeta_t. \quad (9.64)$$

dense matrices and the cartesian basis vectors $\{e_k\}$

$$(A \odot B) e_k = (A e_k) \odot (B e_k)$$

$$e_k^\top (A \odot B) = (e_k^\top A) \odot (e_k^\top B)$$

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 $\rho_t = (I - e_1 e_1^\top)(v_t - u_t)$

$$\|\zeta_{t+1}\|_\infty \leq (1 + \eta \|\rho_t\|_\infty) \|\zeta_t\|_\infty. \quad (9.65)$$

Now, we will prove the lemma by using strong induction on t . Suppose that the first two pieces of (9.58) hold for all iterations up to t . We can show that

$$\|\rho_t\|_\infty \leq \|u_t\|_\infty + \|v_t\|_\infty \quad (9.66)$$

$$\leq 4\delta + \|\zeta_t\|_\infty^2 + 4\delta \|\zeta_t\|_2^2 \quad (\text{by (9.54) and (9.57)}) \quad (9.67)$$

$$\leq 4\delta + (2\alpha)^2 + 4\delta \cdot \frac{1}{2} \quad (\text{by the inductive hypothesis}) \quad (9.68)$$

$$\leq 8\delta, \quad (\text{arbitrarily small}) \quad (9.69)$$

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 (9.65), we have $\|\zeta_{t+1}\|_\infty \leq (1 + \eta \|\rho_t\|_\infty) \|\zeta_t\|_\infty \leq 8\delta$

$$\|\zeta_{t+1}\|_\infty \leq (1 + 8\eta\delta) \|\zeta_t\|_\infty = (1 + O(\eta\delta)) \|\zeta_t\|_\infty. \quad (9.70)$$

$$\text{③} \quad \text{which proves the third piece of the lemma. Using this piece, we can show that}$$

$$\text{LEMMA 22: } t \leq \frac{1}{C\eta\delta} \text{ & sufficiently large } C$$

$$\|z_{t+1}\|_\infty \leq (1 + 8\eta\delta)^{t+1} \|z_0\|_\infty \leq (1 + 8\eta\delta)^{1/(c\eta\delta)} \cdot \alpha \leq 2\alpha \quad \left. \begin{array}{l} \text{① first} \\ \text{LEMMA 22: } t \leq \frac{1}{C\eta\delta} \text{ & sufficiently large } C \end{array} \right\} \quad (9.71)$$

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

$$\|z_{t+1}\|_2^2 \leq (1 + 8\eta\delta)^{t+1} \|z_0\|_2^2 \leq (1 + 8\eta\delta)^{1/(c\eta\delta)} \cdot \alpha \sqrt{d} \leq \frac{1}{2}, \quad \left. \begin{array}{l} \text{② second} \\ \text{if } \alpha \leq \frac{1}{n^{O(1)}} \text{, which proves the first piece.} \end{array} \right\} \quad (9.72)$$

(Step 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 9.17. For all $t \leq 1/(c\eta\delta)$ with sufficiently large constant c , we have that

$$r_{t+1} = \underbrace{(1 + \eta(1 - r_t^2))r_t}_{\text{RHS in population version}} + \underbrace{O(\eta\delta)r_t}_{\text{error}}. \quad \begin{array}{l} \text{GD of population loss} \\ p^{t+1} = p^t - \eta \nabla L(p) \\ = p^t - \eta(p^t \odot p^t - p^* \odot p^*) \odot p^* \\ p^t = p_{\text{ref}} + p_t \end{array}$$

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 (9.53) is

$$\beta^{t+1} = \beta^t - \eta[(r_t^2 - 1)\beta^* \odot \beta^* - u_t + v_t] \odot \beta_t. \quad (9.73)$$

$\beta^* = e_1$ assumed beforehand

We have that

$$r_{t+1} = \langle \beta^{t+1}, e_1 \rangle \quad (9.74)$$

$$= \langle \beta^t, e_1 \rangle - \eta(r_t^2 - 1) \langle \beta^t, e_1 \rangle - \eta \langle v_t - u_t, e_1 \rangle \langle \beta^t, e_1 \rangle \quad (9.75)$$

$$= r_t - \eta(r_t^2 - 1)r_t - \eta \langle v_t - u_t, e_1 \rangle r_t \quad (9.76)$$

$$= \underbrace{(1 + \eta(1 - r_t^2))r_t}_{\text{RHS in population version}} + \underbrace{\eta \langle u_t - v_t, e_1 \rangle r_t}_{\text{error}} \quad (9.77)$$

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

$$|\eta \langle u_t - v_t, e_1 \rangle r_t| \leq \eta \cdot r_t \|v_t - u_t\|_\infty \max \quad (9.78)$$

$$\leq \eta \cdot r_t \cdot 8\delta \quad (\text{by (9.69)}) \quad (9.79)$$

$$= O(\eta\delta) \cdot r_t. \quad (9.80)$$

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

Putting it all together [Finally], we return to the proof of Theorem 9.14. By Lemma 9.17, we know that as long as $r_t \leq 1/2$ it will grow exponentially fast, since $\eta \geq 0$, GD on $L(p)$ with $t = \Theta(\frac{\log(1/\delta)}{\eta})$ steps satisfies... $\|p^t \odot p^t - p^* \odot p^*\|_2^2 \leq \delta(\frac{1}{\sqrt{n}})$

$$\Theta_{\mathcal{G}} r_{t+1} \geq (1 + \eta(1 - r_t^2) - O(\eta\delta)) \cdot r_t \geq \left(1 + \frac{\eta}{2}\right) \cdot r_t. \quad (9.81)$$

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.

As long as $r_t \leq 1/2$, it $(1 + \eta(1 - r_t^2))r_t$ will grow exponentially fast

Since $r_{t+1} \geq (1 + \eta(1 - r_t^2))r_t$

\Rightarrow At some $t_0 = O\left(\frac{\log(1/\alpha)}{\eta}\right)$,

$r_{t_0} > 1/2$ for the first time.

As long as $r_t \leq 1/2$,
it (more) will grow exponentially fast
since $r_{t+1} \geq (1 + \frac{1}{2})r_t$.
 \Rightarrow At some $t_0 = o(\frac{\log(1/\delta)}{\eta})$,
 $r_{t_0} > 1/2$ for the first time.

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

$$1 - r_{t+1} \leq 1 - r_t - \eta(1 - r_t^2)r_t + O(\eta\delta) \cdot r_t \quad (9.82)$$

$$\leq 1 - r_t - \frac{\eta(1 - r_t^2)r_t}{2} + O(\eta\delta) \quad (9.83)$$

$$\leq 1 - r_t - \frac{\eta(1 - r_t)}{2} + O(\eta\delta) \quad (9.84)$$

$$= \left(1 - \frac{\eta}{2}\right)(1 - r_t) + O(\eta\delta). \quad (9.85)$$

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

$$r_{t+1} - 1 \leq (1 - \eta)(r_t - 1) + O(\eta\delta) \leq O(\delta), \quad (9.86)$$

$r_{t+1} - 1 \leq 0(f) \text{ in } \Theta\left(\frac{\log(1/\delta)}{\eta}\right) \text{ iterations. } \Rightarrow r_t \text{ remains very close to 1.}$

This completes the proof of Theorem 9.14, bounding the number of iterations needed for gradient descent on the empirical loss to converge to β^* . \square

9.3 From small to large initialization: a precise characterization

We have previously discussed how certain initializations of gradient descent converge to minimum-norm solutions. In the sequel, we characterize the effect of initialization more precisely—we will show that in a variant of the settings in Section 9.2, we can precisely compute the corresponding regularizer induced by any initialization.) We will see that when the initialization is small, we obtain the bias towards minimum norm solution (in the parameter space used in optimization), whereas when the initialization is large, we are in the NTK regime (Section 8.4) where the implicit bias is towards the min norm solution under the NTK kernel. The materials in this subsection are simplifications of results in the recent paper Woodworth et al. [2020].

effect of initialization \rightarrow bias in regularizer | ① init small \rightarrow bias towards min norm
② " larger " " " " " under the NTK kernel

9.3.1 Preparation: gradient flow : GD with tiny step learning rate.

To simplify the analysis, we will consider the concept of gradient flow, i.e. gradient descent with an infinitesimal learning rate. This allows us to omit the "second order effect" from the learning rate and simplify the analysis.

We begin by recalling the gradient descent update formula. In our previous description of gradient descent, we indexed the updated parameters by $t = 1, 2, \dots$. Anticipating our generalization to infinitesimal steps, we will index the updated parameters using parentheses (instead of subscripts.) In particular, the standard gradient descent update given a loss function $L(w)$ is

from GD update:
 $w(t+1) = w(t) - \eta \nabla L(w(t)). \quad (9.87)$

? Update: after time step w. size η
If we scale the time by η so that each update (by gradient descent) corresponds to a time step of size η (rather than size 1), the update becomes

each update: after time step $t\eta$.
 $w(t+\eta) = w(t) - \eta \nabla L(w(t)). \quad (9.88)$

time step size $\eta \rightarrow 0$ (infinitesimally small step size)
Taking $\eta \rightarrow 0$ yields a differential equation, which can be thought of as a continuous process rather than discrete updates:
 $\eta \rightarrow dt$

$$w(t + dt) = w(t) - dt \cdot \nabla L(w(t)). \quad (9.89)$$

GD : $w(t+dt) = w(t) - dt \cdot \nabla L(w(t))$

This can also be written as:

$$\text{GD: } w(t+dt) = w(t) - dt \cdot \nabla L(w(t))$$

gradient flow dynamics: $\dot{w}(t) = -\nabla L(w(t))$ with $\dot{w}(t) = \frac{\partial w(t)}{\partial t}$

(9.90)

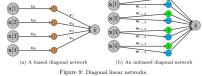
This allows us to ignore the η^2 term (alternatively the (dt^2) term), which will simplify some of the technical details that follow.

\ominus algorithmic reg.

9.3.2 Characterizing the implicit bias of initialization $w = [w_+ \ w_-] \in \mathbb{R}^{2d}$

Our squared parameterization model: $* f(w, x) = \frac{1}{2} (w_{+,-}^T - w_{-,-}^T) x_i = \langle (w_+ - w_-), x \rangle = \langle \theta_w, x \rangle$

The results in this section are slight simplification of the recent paper by Woodworth et al. [2020]. The model is a variant of the one we introduced in (9.7). Recalling that $x^{\odot 2} = x \odot x$, let



$$\begin{aligned} f(w, x) &= \langle w, x \rangle \\ f_w(x) &= (w_+^{\odot 2} - w_-^{\odot 2})^T x. \end{aligned} \quad \begin{aligned} * f_\beta(x) &= (\beta \odot \beta)^T x \\ &\text{Can only represent positive linear combinations of } x. \end{aligned} \quad (9.91)$$

for unbiased diagonal network... use two weights w_+ and w_- . where $w_+, w_- \in \mathbb{R}^d$. Let w denote the concatenation of the two parameter vectors, i.e. $w = (w_+, w_-)$. In (9.7), we defined $f_\beta(x) = (\beta \odot \beta)^T x$; this model can only represent positive linear combinations of x . By contrast, $f_w(x)$ can represent any linear model. Moreover, if we choose our initialization (for w) such that $w_+(0) = w_-(0)$, we obtain $f_{w(0)}(x) \equiv 0$ for all x . Similar to our analysis of the NTK, this initialization will simplify the subsequent derivations.

Remember that $\hat{L}(p) = \frac{1}{4n} \sum_{i=1}^n (y^{(i)} - f_p(x^{(i)}))^2$ where $f_p(x) = \langle p \odot p, x \rangle$

Next, we define the following loss function,

squared loss of the model over a training set $(x_1, y_1), \dots, (x_n, y_n)$

$$\hat{L}(w) = \frac{1}{2} \sum_{i=1}^n (y^{(i)} - f_w(x^{(i)}))^2, \quad (9.92)$$

and consider the initialization

$$w_+(0) = w_-(0) = \alpha \cdot \vec{1} \quad (9.93)$$

where $\vec{1}$ denotes the all-ones vector. The analysis technique still applies to any general initializations as long as all the dimension are initialized to be non-zero, but the the initialization scale is the most important factor, and therefore we chose this simplification for the ease of exposition.

Note that every $w = (w_+, w_-)$ corresponds to a de facto linear function of x . We denote the resulting linear model as θ_w :

$$\begin{aligned} w &= (w_+, w_-) \leftarrow \text{linear function of } x. \\ \theta_w &= w_+^{\odot 2} - w_-^{\odot 2} = (w_+ \odot w_+) - (w_- \odot w_-) \\ f_w(x) &= (w_+^{\odot 2} - w_-^{\odot 2})^T x \\ &= \theta_w^T x \end{aligned} \quad (9.94)$$

Note that $\theta_w^T x = f_w(x)$. (many possible solutions $x \theta = y$)

Let $w(\infty)$ denote the limit of the gradient flow, i.e.

$$w(\infty) = \lim_{t \rightarrow \infty} w(t). \quad (9.95)$$

Then, the converged linear model in the θ space is defined by $\theta_{w(\infty)} = \theta_{w(\infty)}$ —we are interested in understanding its properties. For simplicity, we will omit the ∞ index and refer to this quantity as θ_α . We assume throughout that the limit exists and all corresponding regularity conditions are met.

Let

$$X = \begin{bmatrix} x^{(1)\top} \\ \vdots \\ x^{(n)\top} \end{bmatrix} \in \mathbb{R}^{n \times d} \quad \text{and} \quad \vec{y} = \begin{bmatrix} y^{(1)} \\ \vdots \\ y^{(n)} \end{bmatrix}. \quad (9.96)$$

\exists , solution reached by GD when init. is $w_+(0) = w_-(0) = \lambda w_0$ ($\& w_0 = 1$)

In the sequel, we formally state our result relating the complexity of the solution discovered by gradient flow to the size of the initialization.

Size of initialization $\left(\begin{array}{l} w_+(0) = w_-(0) \\ = \alpha \cdot \vec{1} \end{array} \right)$ \curvearrowleft complexity of the solution discovered by gradient flow

$\hat{\theta}_2$, the GD Solution $\hat{\theta}_2$
(for the squared parametrization model $f(w, x) = \langle \theta w, x \rangle$)
satisfies $X\hat{\theta}_2 = \vec{y}$

Theorem 9.18 (Theorem 1 in Woodworth et al. [2020]). For any $0 < \alpha < \infty$, assume that gradient flow with initialization $w_+(0) = w_-(0) = \alpha \cdot \vec{1}$ converges to a solution that fits the data exactly: $X\hat{\theta}_\alpha = \vec{y}$.³ Then, the solution satisfies the following notion of minimum complexity:

$$\begin{aligned} \text{Our solution, } \hat{\theta}_\alpha &= \underset{\theta}{\operatorname{argmin}} Q_\alpha(\theta) \\ \text{Our linear predictor} &\quad \text{as } \hat{\theta}_\alpha \text{ is "between the two regimes" as } \alpha \rightarrow 0 \text{ or } \alpha \rightarrow \infty \\ \text{s.t. } X\theta &= \vec{y} \end{aligned}$$

$$\hat{\theta}_\alpha = \theta_{w(\alpha)} = (\vec{w}_+(\alpha) \odot \vec{w}_-(\alpha)) - (\vec{w}_-(\alpha) \odot \vec{w}_+(\alpha)) \quad (9.97)$$

$$(9.98)$$

where

and

Q_α :
The implicit regularizer
that biases the
gradient flow solution
towards
one \approx zero-error solution
(out of many possible solutions.)

$$Q_\alpha(\theta) = \alpha^2 \cdot \sum_{i=1}^n q\left(\frac{\theta_i}{\alpha^2}\right)$$

$$q(z) = 2 - \sqrt{4 + z^2} + z \cdot \operatorname{arcsinh}\left(\frac{z}{2}\right) \quad (9.100)$$

$$q(z) = \int_0^z \operatorname{arcsinh}\left(\frac{u}{2}\right) du \uparrow \Theta$$

$$(9.99)$$

In words, Theorem 9.18 claims that $\hat{\theta}_\alpha$ is the minimum complexity solution for the complexity measure Q_α .

Remark 9.19. In particular, when $\alpha \rightarrow \infty$ we have that

$$q(\theta_i/\alpha^2) \asymp \theta_i^2/\alpha^4 \quad \text{when } \alpha \rightarrow \infty, \quad (9.101)$$

and so

$$Q_\alpha(\theta) \asymp \frac{1}{\alpha^2} \|\theta\|_2^2 = \frac{\alpha^2}{\alpha^2} \sum_{i=1}^n \frac{\theta_i^2}{\alpha^4} = \frac{1}{\alpha^2} \sum_{i=1}^n \theta_i^2 = \frac{1}{\alpha^2} \|\theta\|_2^2. \quad (9.102)$$

This means that if $\alpha \rightarrow \infty$ then the complexity measure Q_α is the ℓ_2 -norm, $\|\theta\|_2$. If $\alpha \rightarrow 0$, then the complexity measure becomes

$$q\left(\frac{\theta_i}{\alpha^2}\right) \asymp \frac{|\theta_i|}{\alpha^2} \log\left(\frac{1}{\alpha^2}\right) \quad \text{(by Taylor expansion)} \quad (9.103)$$

and so,

$$Q_\alpha(\theta) \asymp \frac{\|\theta\|_1}{\alpha^2} \log\left(\frac{1}{\alpha^2}\right) = \frac{\alpha^2}{\alpha^2} \sum_{i=1}^n \frac{|\theta_i|}{\alpha^2} \log\left(\frac{1}{\alpha^2}\right) = \sum_{i=1}^n |\theta_i| \log\left(\frac{1}{\alpha^2}\right) = \frac{\|\theta\|_1}{\alpha^2} \log\left(\frac{1}{\alpha^2}\right) \quad (9.104)$$

To summarize, for $\alpha \rightarrow \infty$, the constrained minimization problem we solve in (9.98) yields the minimum ℓ_2 -norm solution of θ (i.e. the ℓ_4 -norm for w). When $\alpha \rightarrow 0$, solving (9.98) yields the minimum ℓ_1 -norm θ (which is the ℓ_2 -norm for w). For $0 < \alpha < \infty$, we obtain some interpolation (of ℓ_1 and ℓ_2 regularization) of the optimum.

Remark 9.20. Note that when $\alpha \rightarrow 0$, the intuition is similar to what we had observed in previous analyses; in particular, the solution is the global minimum closest to the initialization. Note however, that when $\alpha \neq 0$, the solution discovered by gradient descent will not exactly correspond to the solution closest to the initialization.

Chizat et al. **Remark 9.21.** When $\alpha \rightarrow \infty$, we claim that the model optimization is in the neural tangent kernel (NTK) regime. Recall that we had two parameters, (σ, β) , that determined if we could treat the optimization problem as a kernel regression. Further recall that σ denotes the minimum singular value of Φ and β is the Lipschitzness of the gradient. Let us now compute σ and β for large α initializations of our model.

as we increase the scale of init., the dynamics converge to the kernel gradient flow dynamics

³This assumption can likely be proved to be true and thus not required. Here we still include the condition because the original paper Woodworth et al. [2020] assumed it.

④ Kernel regime:

- where it does not change significantly.
- $K(w)$ does not change over the optimization.
- $K_t, K(w_t) \approx K(w)$

⑤ Tangent kernel

$$K_w(x, x') = \langle \nabla_w f(w(x), x), \nabla_w f(w(x'), x') \rangle = \langle \phi_{w(x)}(x), \phi_{w(x')}(x') \rangle$$

where we have our affine model $f(w, x) \approx f_0(x) + \langle w, \phi_{w(x)}(x) \rangle$

For $w_-(0) = w_+(0) = \alpha \vec{1}$, $\nabla f_w(\alpha) = (w_+^{(0)} - w_-^{(0)})^\top \alpha$

$$\nabla f_{w(0)}(x) = 2 \begin{bmatrix} w_+(0) \cdot x \\ -w_-(0) \odot x \end{bmatrix} = 2\alpha \begin{bmatrix} x \\ -x \end{bmatrix} \quad (9.105)$$

by the chain rule. It is clear then that both σ and β linearly depend on α . This implies that

$$\frac{\beta}{\sigma^2} \xrightarrow{\alpha \rightarrow \infty} 0 \quad \text{as } \alpha \rightarrow \infty \quad (9.106)$$

since the denominator is $O(\alpha^2)$, while the numerator is $O(\alpha)$. In particular, the features used in this kernel method are:

$$\phi(x) = \nabla f_{w(0)}(x) = 2\alpha \begin{bmatrix} x \\ -x \end{bmatrix} \quad (9.107)$$

The neural tangent kernel perspective then gives an alternative proof of this complexity minimization result for $\alpha \rightarrow \infty$. In the NTK regime, the solution (to our convex problem) is always the **minimum ℓ_2 -norm** solution for the feature matrix, which in this case equals $\begin{bmatrix} X \\ -X \end{bmatrix}$.

Note that practice tends not to follow the assumptions made here. Often, people either do not use large initializations or do not use infinitesimally small step sizes. But this is a good thing because we do not want to be in the NTK regime; being in the NTK regime implies that we are doing no different or better than just using a kernel method.

We can now prove Theorem 9.18, which is similar to the overparametrized linear regression proof of Theorem 9.3.

This proof follows in two steps:

1. We find an **invariance** maintained by the optimizer. In the overparametrized linear regression proof of Theorem 9.3, we required $\theta \in \text{span}\{x^{(i)}\}$. For this proof, we will use a slightly more complicated invariance.
2. We **characterize the solution** using this invariance. The **invariance**, which depends on α , will tell us which **zero error solution** the optimization converges to.

Note also that all of these conditions only depend upon the **empirically observed samples**. The invariance and minimum is not defined with respect to any population quantities.

Proof. Let

Our solution Θ_α given by f applied to the limit of the ⁽ⁱⁱⁱ⁾ gradient flow dynamics on w .

$$\tilde{X} = [X \quad -X] \in \mathbb{R}^{n \times 2d} \quad \text{and} \quad w(t) = \begin{bmatrix} w_+(t) \\ w_-(t) \end{bmatrix} \in \mathbb{R}^{2d}. \quad (9.108)$$

Then, the **model output on n data points** can be described in matrix notation as follows:

$$\left(\tilde{X} w(t)^{\odot 2} \right) = \left([X \quad -X] \begin{bmatrix} w_+(t)^{\odot 2} \\ w_-(t)^{\odot 2} \end{bmatrix} \right) = \begin{bmatrix} f_{w(t)}(x^{(1)}) \\ \vdots \\ f_{w(t)}(x^{(n)}) \end{bmatrix} \in \mathbb{R}^n. \quad (9.109)$$

Given the loss function,

$$L(w(t)) = \frac{1}{2} \left\| \tilde{X} w(t)^{\odot 2} - \vec{y} \right\|_2^2, \quad (9.110)$$

$$\begin{aligned} \tilde{X} &= [X \quad -X] \in \mathbb{R}^{n \times 2d} && \text{model output on } n \text{ datapoints:} \\ w(t) &= \begin{bmatrix} w_+(t) \\ w_-(t) \end{bmatrix} \in \mathbb{R}^{2d} && \tilde{X} w(t)^{\odot 2} = [X \quad -X] \begin{bmatrix} w_+(t)^{\odot 2} \\ w_-(t)^{\odot 2} \end{bmatrix} = \begin{bmatrix} f_{w(t)}(x^{(1)}) \\ \vdots \\ f_{w(t)}(x^{(n)}) \end{bmatrix} \in \mathbb{R}^n \\ L(w(t)) &= \frac{1}{2} \| \tilde{X} w(t)^{\odot 2} - \vec{y} \|_2^2 && \end{aligned}$$

$\tilde{X} = [X \ -X] \in \mathbb{R}^{n \times 2d}$ model output on n datapoints:
 $w(t) = \begin{bmatrix} w_+(t) \\ w_-(t) \end{bmatrix} \in \mathbb{R}^{2d}$
 $\tilde{X}w(t)^{\odot 2} = \begin{bmatrix} X \ -X \end{bmatrix} \begin{bmatrix} w_+(t)^{\odot 2} \\ w_-(t)^{\odot 2} \end{bmatrix} = \begin{bmatrix} f_{w_+}(x^{(1)}) \\ f_{w_+}(x^{(2)}) \\ \vdots \\ f_{w_+}(x^{(n)}) \end{bmatrix} \in \mathbb{R}^n$

the gradient of $w(t)$ can be computed as

Our gradient flow dynamics:
 $\dot{w}(t) = -\nabla L(w(t))$ (9.111)

$$= -\nabla \left(\left\| \tilde{X}w(t)^{\odot 2} - \vec{y} \right\|_2^2 \right) = -2 \tilde{X}^\top (\tilde{X}w(t)^{\odot 2} - \vec{y}) \odot w(t) \quad (9.112)$$

$$= (2\tilde{X}^\top r(t)) \odot w(t) \quad (\text{chain rule}) \quad (9.113)$$

where $r(t) = \tilde{X}w(t)^{\odot 2} - \vec{y}$ denotes the residual vector. We see that the $\tilde{X}^\top r(t)$ term in (9.113) is reminiscent of linear regression for which it would correspond to the gradient, although the $\odot w(t)$ reminds us that this problem is indeed quadratic.

We cannot directly solve this differential equation, but we claim that

$\downarrow \text{why!}$

$$w(t) = w(0) \odot \exp \left(-2\tilde{X}^\top \int_0^t r(s) ds \right) \quad (\text{exp is applied entry-wise}) \quad (9.114)$$

$w(t) = w(0) \odot \exp(-2\tilde{X}^\top \int_0^t r(s) ds)$

which is not quite a closed form solution of equation 9.113 since $r(s)$ is still a function of $w(t)$. To understand how we obtained this “solution,” we consider a more abstract setting. Suppose that

$\dot{u}(t) = v(t)\dot{u}(t)$
 $u(t) = u(0) \odot \exp \left(\int_0^t v(s) ds \right)$

$$\dot{u}(t) = v(t) \dot{u}(t) \quad (9.115)$$

We can then “solve” this differential equation as follows. Rearranging, we observe that

$$\frac{\dot{u}(t)}{u(t)} = v(t) \quad (9.116)$$

$$\frac{d \log u(t)}{dt} = v(t) \quad (\text{chain rule}) \quad (9.117)$$

$$\log u(t) - \log u(0) = \int_0^t v(s) ds \quad (\text{integration}) \quad (9.118)$$

$\downarrow \log \frac{u(t)}{u(0)}$

$$\frac{u(t)}{u(0)} = \exp \left(\int_0^t v(s) ds \right) \Rightarrow u(t) = u(0) \exp \int_0^t v(s) ds \quad (9.119)$$

In our problem, $u \leftrightarrow w_i$ and $v \leftrightarrow (\tilde{X}^\top r(t))_i$.

We have characterized w , but we want to transform this to a characterization that involves θ . Recall that $w_+(0) = \alpha \vec{1}$ and $w_-(0) = \alpha \vec{1}$ (so that $w(0) = \alpha \vec{1} \in \mathbb{R}^{2d}$). Additionally, we have that $\theta(t) = w_+(t)^{\odot 2} - w_-(t)^{\odot 2}$. We can now apply (9.114) to expand $w(t)$ and simplify. (def. in 9.14)

Note that if we have $\tilde{X}^\top = \begin{bmatrix} X^\top \\ -X^\top \end{bmatrix} \in \mathbb{R}^{2n \times d}$, then for some vector v ,

$$w(t)^{\odot 2} \rightarrow (\exp(-2\tilde{X}^\top v))^{\odot 2} = \begin{bmatrix} \exp(-2X^\top v) \\ \exp(2X^\top v) \end{bmatrix}^{\odot 2} \quad (9.120)$$

$$(w(t) = w(0) \odot \exp(-2\tilde{X}^\top \int_0^t r(s) ds)) \quad = \begin{bmatrix} \exp(-4X^\top v) \\ \exp(4X^\top v) \end{bmatrix}. \quad (9.121)$$

Applying this result for $v = \int_0^t r(s) ds$, we obtain that:

$$\theta(t) = w_+(t)^{\odot 2} - w_-(t)^{\odot 2} \quad (9.122)$$

$\downarrow : w(0)$

$$= \alpha^2 \left[\exp \left(-4X^\top \int_0^t r(s) ds \right) - \exp \left(4X^\top \int_0^t r(s) ds \right) \right] \quad (9.123)$$

$$= 2\alpha^2 \sinh \left(-4X^\top \int_0^t r(s) ds \right). \quad (9.124)$$

$$X\theta_a = \vec{y}$$

$$\theta_a = b_a(X^\top v)$$

Letting $t \rightarrow \infty$, we have that

$$\begin{aligned} \theta(t) &= 2\alpha^2 \sinh \left(-4X^\top \int_0^t r(s) ds \right) \\ \theta_\alpha &= 2\alpha^2 \sinh \left(-4X^\top \int_0^\infty r(s) ds \right). \end{aligned} \quad (9.125)$$

Lastly, we also know

$$X\theta_\alpha = \vec{y} \quad (\because \theta_\alpha \text{ is our global minimum with zero error}) \quad (9.126)$$

since this is the assumption by the theorem (which should can be proven because the optimization should converge to a zero-error solution). We next show that (9.125) and (9.126) are also sufficient conditions for a solution to the constrained optimization problem given by (9.98). In particular, (9.125) and (9.126) correspond to the Karush-Kuhn-Tucker (or KKT) conditions of (9.98).

A KKT condition is an optimality condition for constrained optimization problems. While these conditions can have a variety of formulations and typically one can invoke some off-the-shelf theorems to use them, we can motivate the conditions we encountered by considering the following general optimization program:

$$\begin{array}{ll} \operatorname{argmin} & Q(\theta) \\ \text{s.t.} & X\theta = \vec{y} \end{array} \quad (9.127)$$

$$\theta_a = \operatorname{argmin}_\theta Q_a(\theta) \text{ s.t. } X\theta = \vec{y} \quad (9.128)$$

We say that θ satisfies the (first order) KKT conditions if

$$\begin{cases} \nabla Q(\theta) = X^\top \nu \text{ for some } \nu \in \mathbb{R}^n \\ X\theta = \vec{y} \end{cases} \quad (9.129)$$

$$\text{are } X\theta^* = \vec{y} \text{ and } \nabla Q_a(\theta^*) = X^\top \nu \quad (9.130)$$

$$\text{then, } \nabla Q_a(\theta_a) = \nabla Q_a(b_a(X^\top \nu)) \quad (9.130)$$

$$\text{(a certain function) (a vector) such that } \theta_a = b_a(X^\top \nu) \quad (9.130)$$

More intuitively, we know that optimality implies that there are no first order local improvements that satisfy the constraint (up to first order). Then, consider a perturbation $\Delta\theta$. In order to satisfy the constraint, we must enforce the following:

$$\Delta\theta \perp \text{row-span}\{X\} \quad \text{so} \quad X\Delta\theta = 0 \quad (9.131)$$

So, if we look at $\theta + \Delta\theta$ satisfying the constraint, we can use a Taylor expansion to show that

$$Q(\theta + \Delta\theta) = Q(\theta) + \langle \Delta\theta, \nabla Q(\theta) \rangle \leq Q(\theta) \quad (9.132)$$

because if $\langle \Delta\theta, \nabla Q(\theta) \rangle$ is positive it violates the optimality assumption. In fact, it is very easy to make the sign flip for $\langle \Delta\theta, \nabla Q(\theta) \rangle$ because you can flip $\Delta\theta$ to be the opposite direction. This means that

$$\forall \Delta\theta \perp \text{row-span}\{X\}, \quad \langle \Delta\theta, \nabla Q(\theta) \rangle = 0 \quad (9.133)$$

because if it is negative, you can equivalently flip it to be positive which violates optimality. This means that $Q(\theta) \subseteq \text{row-span}\{X\}$, or $Q(\theta) = X^\top \nu$ for some ν .

Returning to our problem, the KKT condition gives

$$\nabla Q(\theta) = X^\top \nu \quad (9.134)$$

and the invariance gives us

$$\begin{aligned} \theta_\alpha &= 2\alpha^2 \sinh \left(-4X^\top \int_0^\infty r(s) ds \right) \quad (9.125) \\ &= 2\alpha^2 \sinh(-4X^\top v') \end{aligned} \quad (9.135)$$

$$X\theta_a = \vec{y}$$

$$\theta_a = b_a(X^\top v)$$

$$\text{with } b_a(z) = 2\alpha^2 \sinh(z)$$

$$\text{and } \nu = -4 \int_0^\infty r(s) ds$$

$$\theta_a = 2\alpha^2 \sinh \left(X^\top \left(-4 \int_0^\infty r(s) ds \right) \right)$$

$$= \operatorname{arcsinh} \left(\frac{1}{2\alpha^2} \theta \right) \quad (9.135)$$

where we let $v' = \int_0^\infty r(s) ds$ for simplicity. Taking the gradient of Q gives

$$\nabla Q_\alpha(\theta) = \operatorname{arcsinh} \left(\frac{1}{2\alpha^2} \theta \right)$$

$$\theta_a(\theta) = \theta^2 \sum_{i=1}^n q_i \left(\frac{\theta_i}{\alpha^2} \right)$$

$$\text{where } q_i(z) = \frac{2}{z^2+1+z^2} + z \cdot \operatorname{arcsinh}(z/2) \quad (9.137)$$

Plugging in θ_α , we get

$$\begin{aligned} \nabla Q(\theta_\alpha) &= \operatorname{arcsinh} \left(\frac{1}{2\alpha^2} 2\alpha^2 \sinh(-4X^\top v') \right) = -4X^\top v' \\ \nabla Q(\theta_\alpha) &= \operatorname{arcsinh} \left(\frac{1}{2\alpha^2} \theta_\alpha \right) = -4X^\top v' \end{aligned} \quad (9.138)$$

Thus, θ_α satisfies both KKT conditions. Even further, since our optimization problem (9.98) is convex (we do not formally argue this), we conclude that θ_α is a global minimum. \square

9.4 Implicit regularization towards max-margin solutions in classification⁶

We now switch our focus to classification problems. We consider linear models (though these results also apply to nonlinear models with a weaker version of the conclusion). We assume that our data is separable and will prove that gradient descent converges to the max-margin solution. This result holds for any initialization and does not require any additional regularization; we only require the use of gradient descent and the standard logistic loss function. The results in this subsection are originally given by Soudry et al. [2018], and our exposition heavily depends on those in [Ji and Telgarsky, 2018, Telgarsky, 2021].

Assume we have data $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$, where $x^{(i)} \in \mathbb{R}^d$ and $y^{(i)} \in \{\pm 1\}$. We consider the linear model $h_w(x) = w^\top x$ and the cross entropy loss function $\hat{L}(w) = \sum_{i=1}^n \ell(y^{(i)}, h_w(x^{(i)}))$, where $\ell(t) = \log(1 + \exp(-t))$ is the logistic loss.

As we have separable data, there can be multiple global minima, as you can trivially take an infinite number of separators. More formally, there are an infinite number of unit vectors \bar{w} such that $\bar{w}^\top x^{(i)} y^{(i)} > 0$ for all i as one can perturb any strict separator while still maintaining a separation of classes. Then, we can scale the separator to make the loss arbitrarily small—we have that $\hat{L}(\alpha \bar{w}) \rightarrow 0$ as $\alpha \rightarrow \infty$. Thus, informally, for any unit vector \bar{w} that separate the data, $\infty \cdot \bar{w}$ is a global minimum.

We would like to understand which global minimum gradient descent converges to. We will now show that it finds the max-margin solution. Before we can do so, we recall/introduce the following definitions.

$\gamma(w)$ to which global minima does the GD converge? \Rightarrow unit max-margin solution $\gamma(w)$ converge stably ($\gamma(w)$)

Definition 9.22 (Margin). Let $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$ be given data. Assuming $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$ is linearly separable, the margin is defined as

$$\min_{i \in [n]} y^{(i)} w^\top x^{(i)} \quad (\text{with } \bar{w} = \frac{w}{\|w\|_2} \text{ and } \bar{w}^\top x^{(i)} y^{(i)} > 0) \quad (9.139)$$

Definition 9.23 (Normalized Margin). Let $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$ be given data. Assuming $\{(x^{(i)}, y^{(i)})\}_{i=1}^n$ is linearly separable, the normalized margin is defined as

$$\gamma(w) = \frac{\min_{i \in [n]} y^{(i)} w^\top x^{(i)}}{\|w\|_2} \quad * \text{Scale-invariant} \quad \gamma(w) = \frac{\min_{i \in [n]} y^{(i)} w^\top x^{(i)}}{\|w\|_2} \quad \gamma(w) = \frac{\min_{i \in [n]} y^{(i)} w^\top x^{(i)}}{\|w\|_2} \quad (9.140)$$

Definition 9.24 (Max-Margin Solution). Using the normalized margin γ defined in Definition 9.23, we define a max-margin solution as

$$\bar{\gamma} = \max_w \gamma(w) \quad (9.141)$$

and let w^* be the unit-norm maximizer.⁴

\hookrightarrow "direction" of the max-margin solution!

Using these definitions, we claim the following result.

Theorem 9.25. Gradient flow converges to the direction of max-margin solution in the sense that

$$\gamma(w(t)) \rightarrow \bar{\gamma} \text{ as } t \rightarrow \infty \quad \sigma(w(t)) \rightarrow \bar{\sigma} \text{ as } t \rightarrow \infty \quad (9.142)$$

where $w(t)$ is the iterate at time t .

The following observations provide some intuition for Theorem 9.25.

$\hat{L}(w(t)) \rightarrow 0$ for large t , $\hat{L}'(w(t)) \approx 0$

1. $\hat{L}(w(t)) \rightarrow 0$ by a standard optimization argument. Namely, if the objective is monotone decreasing at each iteration, $\hat{L}(w(t)) \approx 0$ for large enough t .
2. Using a Taylor expansion, we can show that $\ell(z) = \log(1 + \exp(-z)) \approx \exp(-z)$ for large z . Thus, logistic loss is close to exponential loss when z is very large.

⁴The normalized margin $\bar{\gamma}$ is scale-invariant. For $c \neq 0$, $\gamma(cw) = \min_{i \in [n]} \frac{y^{(i)} cw^\top x^{(i)}}{\|cw\|_2} = \min_{i \in [n]} \frac{y^{(i)} w^\top x^{(i)}}{\|w\|_2} = \gamma(w)$.

$$\begin{aligned} \ell_n(x) &= \sum_{k=1}^n \left[\frac{(-1)^{k+1}}{k} (x^{-k}) \right] \\ \ell_n(\ln(e^{-z})) &= \sum_{k=1}^n \left[\frac{(-1)^{k+1}}{k} (e^{-z})^{-k} \right] \\ &= e^{-z} - \frac{1}{1!} e^{-2z} + \frac{1}{2!} e^{-3z} - \dots \\ \text{At large } z, \ell_n(\ln(e^{-z})) &\approx e^{-z} \end{aligned}$$

$$\widehat{L}(w) = \sum_{i=1}^n \ell(y^{(i)} w^\top x^{(i)}) \quad \ell(z) = \log(1 + e^{-z})$$

$\widehat{L}(w(t)) \rightarrow 0$ for large t , $\widehat{L}(w(t)) \approx 0$ 根据经验.

3. Using observation 1, we see that $\|w(t)\|_2 \rightarrow \infty$ because if $\|w(t)\|_2$ were instead bounded, then the loss $\widehat{L}(w(t))$ will be bounded below by a constant that is strictly greater than zero, contradicting observation 1. Formally, if $\|w(t)\|_2 \leq B$, then

$$(\because \|w(t)\|_2 \leq B) \text{ finite} \quad |y^{(i)} w^\top x^{(i)}| \leq B \|x^{(i)}\|_2 \quad (9.143)$$

and therefore we get

$$\widehat{L}(w(t)) \geq \sum_{i=1}^n \exp(-B \|x^{(i)}\|_2) > 0. \quad \downarrow \quad \text{因为 } \widehat{L}(w(t)) \xrightarrow[t \rightarrow \infty]{} 0 \text{ 由经验.}$$

$$\ell(z) = \log(1 + e^{-z}) \approx \exp(-z)$$

4. Suppose we have w such that $\|w\|_2 = q$ is very big. Then, using observation 2, we see that

$$\ell(z) = \log(1 + e^{-z}) \approx \exp(-z)$$

$$\widehat{L}(w) = \sum_{i=1}^n \ell(y^{(i)} w^\top x^{(i)}) \quad (9.145)$$

$$\ell(z) \approx e^{-z} \quad \approx \sum_{i=1}^n \exp(-y^{(i)} w^\top x^{(i)}) \quad (9.146)$$

$$\log \widehat{L}(w) \approx \log \sum_{i=1}^n \exp(-y^{(i)} w^\top x^{(i)}) \quad \|\|w\|_2 = q\|.$$

$$(9.147)$$

$$= \log \sum_{i=1}^n \exp(-q y^{(i)} \bar{w}^\top x^{(i)}) \quad \text{unit } (\bar{w} = \frac{w}{\|w\|_2} = \frac{w}{q})$$

$$(9.148)$$

$$\approx \max_{i \in [n]} (-q y^{(i)} \bar{w}^\top x^{(i)}) \quad \text{LSE}$$

$$(9.149)$$

where $\bar{w} = \frac{w}{\|w\|_2}$ and the last step holds because the log of a sum of exponentials (log-sum-exp) is a smooth approximation to the maximum function. To motivate this claim, observe that:

Log-Sum-exp

$$\text{LSE}(x_1, \dots, x_n) = \log(\exp(x_1) + \dots + \exp(x_n))$$

$$\text{take } \exp(\max_i x_i) \leq \sum_{i=1}^n \exp(x_i) \leq n \exp(\max_i x_i)$$

$$\text{by } \max_i \{x_1, \dots, x_n\} \leq \text{LSE}(x_1, \dots, x_n) \leq \max_i \{x_1, \dots, x_n\} + \log n$$

$$\log \sum_{i=1}^n \exp(q u_i) \geq q \max_i u_i$$

$$(9.150)$$

$$\log \sum_{i=1}^n \exp(q u_i) \leq \log(n \exp(q \max_i u_i))$$

$$(9.151)$$

$$= \log n + q \max_i u_i$$

$$(9.152)$$

$$\approx q \max_{i \in [n]} u_i + o(q) \text{ as } q \rightarrow \infty$$

$$(9.153)$$

Thus, minimizing the loss is the same as

minimizing the logistic loss is ...

$$\min_w \log \widehat{L}(w) \quad (9.154)$$

$$\min_w \left(\max_{i \in [n]} -q y^{(i)} \bar{w}^\top x^{(i)} \right)$$

$$\text{with } \bar{w} = \frac{w}{\|w\|_2} = \frac{w}{q}$$

which can be reformulated as

$$\left(\max_w \min_{i \in [n]} q y^{(i)} \bar{w}^\top x^{(i)} \right) \quad \dots \text{the same as maximizing the margin!} \quad (9.155)$$

$$\text{margin: } \min_{i \in [n]} q y^{(i)} \bar{w}^\top x^{(i)}$$

$$\ell(t) = \log(1 + e^{-t})$$

The above observations heuristically demonstrate that minimizing the logistic loss (with gradient descent) is equivalent (in the limit) to maximizing the margin. Below, we prove Theorem 9.25 rigorously for the exponential loss function $\ell(t) = \exp(-t)$, which is nearly the same as the logistic loss.

$$\sigma(w(t)) \rightarrow \bar{\sigma} \text{ as } t \rightarrow \infty$$

Proof of Theorem 9.25. We begin by defining the smooth margin as

$$\text{smooth margin} \triangleq \frac{-\log \hat{L}(w)}{\|w\|_2}$$

$$\begin{aligned} \tilde{\gamma}(w) &\triangleq \frac{-\log \hat{L}(w)}{\|w\|_2} \\ &= \frac{-\log (\sum_{i=1}^n \exp(-y^{(i)} w^\top x^{(i)}))}{\|w\|_2} \quad \text{or.} \end{aligned} \quad (9.157)$$

Note that $\tilde{\gamma}(w)$ approximates $\gamma(w)$ by the log-sum-exp approximation. (To make this precise, recall that $\gamma(w) \geq \tilde{\gamma}(w)$ because $y^{(i)} w^\top x^{(i)} \geq \gamma(w) \|w\|_2$ for all i . (by definition, $\frac{\gamma(w) w^\top x^{(i)}}{\|w\|_2} \geq \gamma(w)$ or.))

Then, since $\gamma(w) \leq \bar{\gamma}$ by definition, it suffices to show that

$$\bar{\sigma} = \max_w \sigma(w) \quad \lim_{t \rightarrow \infty} \tilde{\gamma}(w(t)) = \bar{\gamma}. \quad (9.158)$$

Let $\dot{w}(t) = -\nabla \hat{L}(w(t))$. Then,

$$\left(\dot{w}(t) - \frac{\partial w(t)}{\partial t} \right) \triangleq \frac{\partial}{\partial t} \left(-\log \hat{L}(w(t)) \right) = \langle \nabla (-\log \hat{L}(w(t))), \dot{w}(t) \rangle \quad (9.159)$$

$$= \langle \nabla (-\log \hat{L}(w(t))), -\nabla \hat{L}(w(t)) \rangle \quad \cancel{X} \quad = \left\langle -\frac{\nabla \hat{L}(w(t))}{\hat{L}(w(t))}, \dot{w}(t) \right\rangle \quad \text{or.} \quad (9.160)$$

$$= \frac{\|\nabla \hat{L}(w(t))\|_2^2}{\hat{L}(w(t))} \quad \text{or.} \quad (9.161)$$

$$= \frac{\|\dot{w}(t)\|_2^2}{\hat{L}(w(t))} \geq 0 \quad (9.162)$$

This result tells us that the log loss is decreasing at each infinitesimal step of the gradient flow. By integrating (9.162), we can also evaluate the log loss at time T :

$$-\log \hat{L}(w(T)) = -\log \hat{L}(w(0)) + \int_0^T \frac{\partial}{\partial t} \left(-\log \hat{L}(w(t)) \right) dt \quad (9.163)$$

$$\log \text{loss at time } T = -\log \hat{L}(w(0)) + \int_0^T \frac{\|\dot{w}(t)\|_2^2}{\hat{L}(w(t))} dt. \quad (9.164)$$

While the derivation above tells us how the numerator of (9.156) is changing, we have yet to relate this to the denominator, i.e. the norm of w . Recall that w^* is the direction of the max-margin solution. Then, we have

$$\|\dot{w}(t)\|_2 \geq \langle \dot{w}(t), w^* \rangle \quad \|\dot{w}(t)\|_2 \geq \langle \dot{w}(t), w^* \rangle \geq \bar{\sigma} \cdot \hat{L}(w(t)) \quad (\text{Cauchy-Schwarz}) \quad (9.165)$$

$$= \langle -\nabla \hat{L}(w(t)), w^* \rangle \quad (9.166)$$

$$= \left\langle \sum_{i=1}^n y^{(i)} \exp(-y^{(i)} w^\top x^{(i)}) \cdot x^{(i)}, w^* \right\rangle \quad (9.167)$$

$$= \sum_{i=1}^n y^{(i)} \exp(-y^{(i)} w^\top x^{(i)}) \cdot \langle w^*, x^{(i)} \rangle \quad (9.168)$$

$$\stackrel{\text{By def.}}{\geq} \bar{\gamma} \sum_{i=1}^n \exp(-y^{(i)} w^\top x^{(i)}) \quad (9.169)$$

$$= \bar{\gamma} \cdot \hat{L}(w(t)). \quad (9.170)$$

This shows that $\dot{w}(t)$ is correlated to w^* , and that this correlation depends on $\bar{\gamma}$ and the loss. In addition, $\dot{w}(t)$ is not too small compared to the loss.

$$(9.163) \quad \log_{\text{loss at time } T} -\hat{L}(w(T)) = -\hat{L}(w(0)) + \int_0^T \frac{\|\dot{w}(t)\|_2^2}{\hat{L}(w(t))} dt \geq -\hat{L}(w(0)) + \bar{\gamma} \|w(T)\|_2. \quad (9.164)$$

Next, we substitute (9.165) into the second term of the right-hand-side of (9.163):

$$\int_0^T \frac{\|\dot{w}(t)\|_2^2}{\hat{L}(w(t))} dt \geq \bar{\gamma} \cdot \int_0^T \|\dot{w}(t)\|_2 dt. \quad (9.171)$$

$$\frac{\|\dot{w}(t)\|_2^2}{\hat{L}(w(t))} = \frac{\|\dot{w}(t)\|_2 \|\dot{w}(t)\|_2}{\hat{L}(w(t))} \geq \frac{\|\dot{w}(t)\|_2 \bar{\gamma} \|\dot{w}(t)\|_2}{\hat{L}(w(t))} \geq \bar{\gamma} \cdot \left\| \int_0^T \frac{\dot{w}(t)}{\|\dot{w}(t)\|_2} dt \right\|_2. \quad (9.172)$$

$$= \bar{\gamma} \|w(T)\|_2. \quad (9.173)$$

Applying this bound to the RHS of (9.163), we obtain

$$-\log \hat{L}(w(T)) \geq -\log \hat{L}(w(0)) + \bar{\gamma} \|w(T)\|_2. \quad (9.174)$$

Dividing both sides by $\|w(T)\|_2$,

$$-\frac{\log \hat{L}(w(T))}{\|w(T)\|_2} \geq -\frac{\log \hat{L}(w(0))}{\|w(T)\|_2} + \bar{\gamma}. \quad (9.175)$$

Since $\lim_{T \rightarrow \infty} \|w(T)\|_2 = \infty$, we know that the first term on the RHS of (9.175) goes to 0 in the limit. Thus,

$$\lim_{T \rightarrow \infty} \left(-\frac{\log \hat{L}(w(T))}{\|w(T)\|_2} \right) \geq \bar{\gamma}. \quad (9.176)$$

Recognizing the LHS as the definition of the smooth margin, i.e. (9.156), we conclude that

$$\lim_{T \rightarrow \infty} \tilde{\gamma}(w(T)) \geq \bar{\gamma}. \quad \begin{matrix} \xrightarrow{\text{smooth margin}} \\ \tilde{\gamma}(w) \triangleq \frac{-\log \hat{L}(w)}{\|w\|_2} \end{matrix} \quad (9.177)$$

Meanwhile, since we know that

$$\tilde{\gamma} \geq \gamma(w(T)) \geq \tilde{\gamma}(w(T)), \quad (\text{for sandwich}) \quad (9.178)$$

we conclude by the squeeze theorem that

$$\lim_{T \rightarrow \infty} \gamma(w(T)) = \lim_{T \rightarrow \infty} \tilde{\gamma}(w(T)) = \bar{\gamma}. \quad \text{oo.} \quad (9.179)$$

Thm 9.23 " $\sigma(w(t)) \rightarrow \bar{\gamma}$ as $t \rightarrow \infty$ "

□

9.5 Implicit regularization effect of noise in SGD

In the previous section, we discussed implicit regularization via initialization and the implicit regularization of gradient descent for logistic loss-minimizing classifiers. In the sequel, we will move forward to the implicit regularization effect of SGD noise. Starting from the quadratic case, we analyze how the SGD noise will affect the optimization solution, and present (heuristically) a result for non-quadratic loss functions. In particular, the main (heuristic) results are:

1. On the one dimensional quadratic function, the iterate can be disentangled into a contraction part and a stochastic part, the latter of which is characterized by the Ornstein–Uhlenbeck (OU) process. The noise makes the iterate bounce around the global minimum.
2. On the multi-dimensional quadratic function, the iterate can be disentangled into multiple separate 1-D OU processes. The noise makes the iterate bounce around the global minimum, while the fluctuation is closely related to the shape of the noise.
3. On non-quadratic functions, SGD with *label noise* on empirical loss $\hat{L}(\theta)$ converges to a stationary point of the regularized loss $\hat{L}(\theta) + \lambda \text{tr}(\nabla^2 \hat{L}(\theta))$, which is mainly due to the accumulation of a third order effect.