Stationary Points of Shallow Neural Networks with Quadratic Activation Function

Eren C. Kızıldağ, joint work with David Gamarnik (MIT) and Ilias Zadik (NYU)

arXiv:1912.01599

MIT MLTea

July 28, 2020



Overview

- Intro and Motivation
- 2 Main Results: Optimization Landscape
- Main Results: Initialization
- 4 Main Results: Generalization



(a) Neural networks are extremely powerful in practice: natural language processing, image recognition, image classification, speech recognition, game playing, etc.

- (a) Neural networks are extremely powerful in practice: natural language processing, image recognition, image classification, speech recognition, game playing, etc.
- (b) A mathematical understanding of this success is largely missing.



- (a) Neural networks are extremely powerful in practice: natural language processing, image recognition, image classification, speech recognition, game playing, etc.
- (b) A mathematical understanding of this success is largely missing.
- (c) For instance, training is NP-hard (Blum and Rivest [89]). However, the gradient descent has great empirical success.

- (a) Neural networks are extremely powerful in practice: natural language processing, image recognition, image classification, speech recognition, game playing, etc.
- (b) A mathematical understanding of this success is largely missing.
- (c) For instance, training is NP-hard (Blum and Rivest [89]). However, the gradient descent has great empirical success.
- (d) Our main motivation is to provide further insights for these networks. In particular, we focus on both training (through the lens of optimization landscape), initialization, and generalization aspects.

(a) Shallow architecture, one hidden layer of width $m \in \mathbb{Z}^+$, quadratic activation $(\sigma(x) = x^2)$. Data $X \in \mathbb{R}^d$ has i.i.d. centered sub-Gaussian coordinates. Realizable model: labels are generated by a teacher network with planted weights $W^* \in \mathbb{R}^{m \times d}$.



- (a) Shallow architecture, one hidden layer of width $m \in \mathbb{Z}^+$, quadratic activation $(\sigma(x) = x^2)$. Data $X \in \mathbb{R}^d$ has i.i.d. centered sub-Gaussian coordinates. Realizable model: labels are generated by a teacher network with planted weights $W^* \in \mathbb{R}^{m \times d}$.
- (b) The network computes, for each $X \in \mathbb{R}^d$, the label $f(W^*; X) = \sum_{1 \leq j \leq m} \langle W_j^*, X \rangle^2$ (which is $\|W^*X\|_2^2$). Here, $W_j^* \in \mathbb{R}^d$ is the j^{th} row of W^* . Also, $\operatorname{rank}(W^*) = d$.

- (a) Shallow architecture, one hidden layer of width $m \in \mathbb{Z}^+$, quadratic activation $(\sigma(x) = x^2)$. Data $X \in \mathbb{R}^d$ has i.i.d. centered sub-Gaussian coordinates. Realizable model: labels are generated by a teacher network with planted weights $W^* \in \mathbb{R}^{m \times d}$.
- (b) The network computes, for each $X \in \mathbb{R}^d$, the label $f(W^*; X) = \sum_{1 \leq j \leq m} \langle W_j^*, X \rangle^2$ (which is $\|W^*X\|_2^2$). Here, $W_j^* \in \mathbb{R}^d$ is the j^{th} row of W^* . Also, $\operatorname{rank}(W^*) = d$.
- (c) Generate $X_i \in \mathbb{R}^d$, $1 \le i \le N$ i.i.d. Let $Y_i = f(W^*; X_i)$, $1 \le i \le N$ be the associated label. **The empirical risk** for any "candidate" $W \in \mathbb{R}^{m \times d}$:

$$\widehat{\mathcal{L}}(W) \triangleq \frac{1}{N} \sum_{1 \leq i \leq N} (Y_i - f(W^*; X_i))^2.$$

 $\widehat{\mathcal{L}}(\cdot)$ has a highly non-convex nature.



- (a) Shallow architecture, one hidden layer of width $m \in \mathbb{Z}^+$, quadratic activation $(\sigma(x) = x^2)$. Data $X \in \mathbb{R}^d$ has i.i.d. centered sub-Gaussian coordinates. Realizable model: labels are generated by a teacher network with planted weights $W^* \in \mathbb{R}^{m \times d}$.
- (b) The network computes, for each $X \in \mathbb{R}^d$, the label $f(W^*; X) = \sum_{1 \leq j \leq m} \langle W_j^*, X \rangle^2$ (which is $\|W^*X\|_2^2$). Here, $W_j^* \in \mathbb{R}^d$ is the j^{th} row of W^* . Also, $\operatorname{rank}(W^*) = d$.
- (c) Generate $X_i \in \mathbb{R}^d$, $1 \le i \le N$ i.i.d. Let $Y_i = f(W^*; X_i)$, $1 \le i \le N$ be the associated label. **The empirical risk** for any "candidate" $W \in \mathbb{R}^{m \times d}$:

$$\widehat{\mathcal{L}}(W) \triangleq \frac{1}{N} \sum_{1 \leq i \leq N} (Y_i - f(W^*; X_i))^2.$$

- $\widehat{\mathcal{L}}(\cdot)$ has a highly non-convex nature.
- (d) Goal of the learner. Solve the empirical risk minimization problem $\min_{W \in \mathbb{R}^{m \times d}} \widehat{\mathcal{L}}(W)$ and understand its generalization ability, as measured by the population risk

$$\mathcal{L}(W) \triangleq \mathbb{E}[(f(W;X) - f(W^*;X))^2].$$



(a) Such shallow architectures with planted weights and Gaussian input data have been explored previously in literature: Du et al. [17], Li & Yuan [17], Tian [17], Zhong et al. [17], Soltanolkotabi [17], Brutzkus & Globerson [17], ...

- (a) Such shallow architectures with planted weights and Gaussian input data have been explored previously in literature: Du et al. [17], Li & Yuan [17], Tian [17], Zhong et al. [17], Soltanolkotabi [17], Brutzkus & Globerson [17], ...
- (b) Quadratic networks, in particular, were considered by Du and Lee [18], Soltanolkotabi, Javanmard, and Lee [18].

- (a) Such shallow architectures with planted weights and Gaussian input data have been explored previously in literature: Du et al. [17], Li & Yuan [17], Tian [17], Zhong et al. [17], Soltanolkotabi [17], Brutzkus & Globerson [17], ...
- (b) Quadratic networks, in particular, were considered by Du and Lee [18], Soltanolkotabi, Javanmard, and Lee [18].
- (c) Quadratic activation function is, admittedly, stylized. However,

- (a) Such shallow architectures with planted weights and Gaussian input data have been explored previously in literature: Du et al. [17], Li & Yuan [17], Tian [17], Zhong et al. [17], Soltanolkotabi [17], Brutzkus & Globerson [17], ...
- (b) Quadratic networks, in particular, were considered by Du and Lee [18], Soltanolkotabi, Javanmard, and Lee [18].
- (c) Quadratic activation function is, admittedly, stylized. However,
 - Block of quadratic activations can be stacked together to approximate deeper networks with sigmoid activations (Livni et al. [14]).

- (a) Such shallow architectures with planted weights and Gaussian input data have been explored previously in literature: Du et al. [17], Li & Yuan [17], Tian [17], Zhong et al. [17], Soltanolkotabi [17], Brutzkus & Globerson [17], ...
- (b) Quadratic networks, in particular, were considered by Du and Lee [18], Soltanolkotabi, Javanmard, and Lee [18].
- (c) Quadratic activation function is, admittedly, stylized. However,
 - Block of quadratic activations can be stacked together to approximate deeper networks with sigmoid activations (Livni et al. [14]).
 - They serve as a second order approximation of general nonlinear activations (Venturi et al. [18]).

- (a) Such shallow architectures with planted weights and Gaussian input data have been explored previously in literature: Du et al. [17], Li & Yuan [17], Tian [17], Zhong et al. [17], Soltanolkotabi [17], Brutzkus & Globerson [17], ...
- (b) Quadratic networks, in particular, were considered by Du and Lee [18], Soltanolkotabi, Javanmard, and Lee [18].
- (c) Quadratic activation function is, admittedly, stylized. However,
 - Block of quadratic activations can be stacked together to approximate deeper networks with sigmoid activations (Livni et al. [14]).
 - They serve as a second order approximation of general nonlinear activations (Venturi et al. [18]).
- (d) Study of quadratic networks is an attempt to gain further insights on more complex networks.



Overview

- Intro and Motivation
- Main Results: Optimization Landscape
- Main Results: Initialization
- 4 Main Results: Generalization

Theorem (Gamarnik, K.; and Zadik, 2020)

Let $X_i \in \mathbb{R}^d$, $1 \le i \le N$ be i.i.d. data with centered i.i.d. sub-Gaussian coordinates; and $Y_i = f(W^*; X_i)$ be the associated label. Then, under certain technical assumptions (in particular if $N = d^{O(1)}$), it holds that with high probability.

$$\min_{\substack{W \in \mathbb{R}^{m \times d} \\ \operatorname{rank}(W) \leq d-1}} \widehat{\mathcal{L}}(W) = \min_{\substack{W \in \mathbb{R}^{m \times d} \\ \operatorname{rank}(W) \leq d-1}} \frac{1}{N} \sum_{1 \leq i \leq N} (Y_i - f(W; X_i))^2 \geq \frac{1}{2} C_5 \sigma_{\min}(W^*)^4.$$

Here, $C_5 > 0$ is a constant determined only by the (conditional) moments of data distribution.

Theorem (Gamarnik, K.; and Zadik, 2020)

Let $X_i \in \mathbb{R}^d$, $1 \le i \le N$ be i.i.d. data with centered i.i.d. sub-Gaussian coordinates; and $Y_i = f(W^*; X_i)$ be the associated label. Then, under certain technical assumptions (in particular if $N = d^{O(1)}$), it holds that with high probability,

$$\min_{\substack{W \in \mathbb{R}^{m \times d} \\ \operatorname{rank}(W) \leq d-1}} \widehat{\mathcal{L}}(W) = \min_{\substack{W \in \mathbb{R}^{m \times d} \\ \operatorname{rank}(W) \leq d-1}} \frac{1}{N} \sum_{1 \leq i \leq N} (Y_i - f(W; X_i))^2 \geq \frac{1}{2} C_5 \sigma_{\min}(W^*)^4.$$

Here, $C_5 > 0$ is a constant determined only by the (conditional) moments of data distribution.

(a) An energy barrier in the landscape of $\widehat{\mathcal{L}}(\cdot)$: if $\operatorname{rank}(W) < d$, then $\widehat{\mathcal{L}}(W)$ is bounded away from zero by an explicit constant. Analogue result for the population risk $\mathcal{L}(W)$.

Theorem (Gamarnik, K.; and Zadik, 2020)

Let $X_i \in \mathbb{R}^d$, $1 \le i \le N$ be i.i.d. data with centered i.i.d. sub-Gaussian coordinates; and $Y_i = f(W^*; X_i)$ be the associated label. Then, under certain technical assumptions (in particular if $N = d^{O(1)}$), it holds that with high probability,

$$\min_{\substack{W \in \mathbb{R}^{m \times d} \\ \operatorname{rank}(W) \leq d-1}} \widehat{\mathcal{L}}(W) = \min_{\substack{W \in \mathbb{R}^{m \times d} \\ \operatorname{rank}(W) \leq d-1}} \frac{1}{N} \sum_{1 \leq i \leq N} (Y_i - f(W; X_i))^2 \geq \frac{1}{2} C_5 \sigma_{\min}(W^*)^4.$$

Here, $C_5>0$ is a constant determined only by the (conditional) moments of data distribution.

- (a) An energy barrier in the landscape of $\widehat{\mathcal{L}}(\cdot)$: if $\operatorname{rank}(W) < d$, then $\widehat{\mathcal{L}}(W)$ is bounded away from zero by an explicit constant. Analogue result for the population risk $\mathcal{L}(W)$.
- (b) Tight up to a multiplicative constant. Sub-Gaussianity not essential: $\mathbb{P}(|X_i(j)| > t) \le \exp(-Ct^{\alpha})$ type tail behavior is ok.

Theorem (Gamarnik, K.; and Zadik, 2020)

Theorem (Gamarnik, K.; and Zadik, 2020)

Suppose $\operatorname{rank}(W) = d$, and $\nabla_W \widehat{\mathcal{L}}(W) = 0$. Then, $\widehat{\mathcal{L}}(W) = 0$. If, furthermore, number N of data is $N \geq d(d+1)/2$ then $W = QW^*$ for an orthogonal matrix $Q \in \mathbb{R}^{m \times m}$.

(a) An analogue holds also for the population risk.



Theorem (Gamarnik, K.; and Zadik, 2020)

- (a) An analogue holds also for the population risk.
- (b) Corresponding losses $\widehat{\mathcal{L}}(\cdot)$ and $\mathcal{L}(\cdot)$ admit no full-rank saddle points.

Theorem (Gamarnik, K.; and Zadik, 2020)

- (a) An analogue holds also for the population risk.
- (b) Corresponding losses $\widehat{\mathcal{L}}(\cdot)$ and $\mathcal{L}(\cdot)$ admit no full-rank saddle points.
- (c) Landscape of $(\widehat{\mathcal{L}}(\cdot))$ and $\mathcal{L}(\cdot)$ has fairly benign properties below the aforementioned energy barrier (**recall:** below the barrier no rank-deficient $W \in \mathbb{R}^{m \times d}$, w.h.p.).

Theorem (Gamarnik, K.; and Zadik, 2020)

- (a) An analogue holds also for the population risk.
- (b) Corresponding losses $\widehat{\mathcal{L}}(\cdot)$ and $\mathcal{L}(\cdot)$ admit no full-rank saddle points.
- (c) Landscape of $(\widehat{\mathcal{L}}(\cdot))$ and $\mathcal{L}(\cdot)$ has fairly benign properties below the aforementioned energy barrier (**recall:** below the barrier no rank-deficient $W \in \mathbb{R}^{m \times d}$, w.h.p.).
- (d) **Next.** Benign landscape \implies Convergence of gradient descent.



Theorem (Gamarnik, **K.**; and Zadik, 2020)

- Running gradient descent (with a certain step size) generates a full-rank, ϵ -approximate stationary point $W \in \mathbb{R}^{m \times d}$ (namely $\|\nabla \widehat{\mathcal{L}}(W)\|_F \leq \epsilon$) in time $\operatorname{poly}(\epsilon^{-1}, d)$.
- For this W, $\widehat{\mathcal{L}}(W) \leq C\epsilon\sigma_{\min}(W^*)^{-2}\mathrm{poly}(d)$, $\mathcal{L}(W) \leq C'\epsilon\sigma_{\min}(W^*)^{-1}\mathrm{poly}(d)$; and $\|W^TW (W^*)^TW^*\|_F \leq C''\epsilon^{\frac{1}{2}}\sigma_{\min}(W^*)^{-1}\mathrm{poly}(d)$. C, C', C'' > 0 constants.

Theorem (Gamarnik, K.; and Zadik, 2020)

- Running gradient descent (with a certain step size) generates a full-rank, ϵ -approximate stationary point $W \in \mathbb{R}^{m \times d}$ (namely $\|\nabla \widehat{\mathcal{L}}(W)\|_F \leq \epsilon$) in time $\operatorname{poly}(\epsilon^{-1}, d)$.
- For this W, $\widehat{\mathcal{L}}(W) \leq C\epsilon\sigma_{\min}(W^*)^{-2}\mathrm{poly}(d)$, $\mathcal{L}(W) \leq C'\epsilon\sigma_{\min}(W^*)^{-1}\mathrm{poly}(d)$; and $\|W^TW (W^*)^TW^*\|_F \leq C''\epsilon^{\frac{1}{2}}\sigma_{\min}(W^*)^{-1}\mathrm{poly}(d)$. C, C', C'' > 0 constants.
- (a) If the initialization is "nice", GD finds an approx. stat. point W in polynomial time.

Theorem (Gamarnik, K.; and Zadik, 2020)

- Running gradient descent (with a certain step size) generates a full-rank, ϵ -approximate stationary point $W \in \mathbb{R}^{m \times d}$ (namely $\|\nabla \widehat{\mathcal{L}}(W)\|_F \leq \epsilon$) in time $\operatorname{poly}(\epsilon^{-1}, d)$.
- For this W, $\widehat{\mathcal{L}}(W) \leq C\epsilon\sigma_{\min}(W^*)^{-2}\mathrm{poly}(d)$, $\mathcal{L}(W) \leq C'\epsilon\sigma_{\min}(W^*)^{-1}\mathrm{poly}(d)$; and $\|W^TW (W^*)^TW^*\|_F \leq C''\epsilon^{\frac{1}{2}}\sigma_{\min}(W^*)^{-1}\mathrm{poly}(d)$. C, C', C'' > 0 constants.
- (a) If the initialization is "nice", GD finds an approx. stat. point W in polynomial time.
- (b) For this W, the weights W^TW are (uniformly) close to planted weights $(W^*)^TW^*$, thus has well-controlled generalization error.

Theorem (Gamarnik, K.; and Zadik, 2020)

- Running gradient descent (with a certain step size) generates a full-rank, ϵ -approximate stationary point $W \in \mathbb{R}^{m \times d}$ (namely $\|\nabla \widehat{\mathcal{L}}(W)\|_F \leq \epsilon$) in time $\operatorname{poly}(\epsilon^{-1}, d)$.
- For this W, $\widehat{\mathcal{L}}(W) \leq C\epsilon\sigma_{\min}(W^*)^{-2}\mathrm{poly}(d)$, $\mathcal{L}(W) \leq C'\epsilon\sigma_{\min}(W^*)^{-1}\mathrm{poly}(d)$; and $\|W^TW (W^*)^TW^*\|_F \leq C''\epsilon^{\frac{1}{2}}\sigma_{\min}(W^*)^{-1}\mathrm{poly}(d)$. C, C', C'' > 0 constants.
- (a) If the initialization is "nice", GD finds an approx. stat. point W in polynomial time.
- (b) For this W, the weights W^TW are (uniformly) close to planted weights $(W^*)^TW^*$, thus has well-controlled generalization error.
- (c) From a technical point: control the condition number of a certain matrix with i.i.d. rows consisting of tensorized data $X_i^{\otimes 2}$. Uses results from a very recent work of us analyzing the spectrum of expected covariance matrices of tensorized data.

Remarks.

(a) An energy barrier separating rank-deficient points. Below this barrier, only full-rank points survive.

- (a) An energy barrier separating rank-deficient points. Below this barrier, only full-rank points survive.
- (b) Any full-rank stationary point of the risk is globally optimal. Thus, the landscape below the barrier is actually benign. No spurios stationary points within the set of full-rank matrices.

- (a) An energy barrier separating rank-deficient points. Below this barrier, only full-rank points survive.
- (b) Any full-rank stationary point of the risk is globally optimal. Thus, the landscape below the barrier is actually benign. No spurios stationary points within the set of full-rank matrices.
- (c) Gradient descent, provided **initialized properly**, "approximately" minimizes the empirical risk, and recovers the planted weights in polynomial time.

- (a) An energy barrier separating rank-deficient points. Below this barrier, only full-rank points survive.
- (b) Any full-rank stationary point of the risk is globally optimal. Thus, the landscape below the barrier is actually benign. No spurios stationary points within the set of full-rank matrices.
- (c) Gradient descent, provided **initialized properly**, "approximately" minimizes the empirical risk, and recovers the planted weights in polynomial time.
- (d) On a technical level. Covering number and concentration arguments. The proximity above is established using a novel concentration result of us for matrices with i.i.d. rows obtained from tensorized data $X_i^{\otimes 2}$.

- (a) An energy barrier separating rank-deficient points. Below this barrier, only full-rank points survive.
- (b) Any full-rank stationary point of the risk is globally optimal. Thus, the landscape below the barrier is actually benign. No spurios stationary points within the set of full-rank matrices.
- (c) Gradient descent, provided initialized properly, "approximately" minimizes the empirical risk, and recovers the planted weights in polynomial time.
- (d) **On a technical level.** Covering number and concentration arguments. The proximity above is established using a novel concentration result of us for matrices with i.i.d. rows obtained from tensorized data $X_i^{\otimes 2}$.
- (e) **Next.** "How to initialize properly?"

Overview

- Intro and Motivation
- Main Results: Optimization Landscape
- Main Results: Initialization
- 4 Main Results: Generalization

(a) **Recall.** GD is successful provided initialized properly.



12 / 18

- (a) **Recall.** GD is successful provided initialized properly.
- (b) **Focus.** Initialization in the context of randomly generated planted $W^* \in \mathbb{R}^{m \times d}$:

- (a) **Recall.** GD is successful provided initialized properly.
- (b) **Focus.** Initialization in the context of randomly generated planted $W^* \in \mathbb{R}^{m \times d}$:
 - Networks with random weights is an active research area, define initial loss landscape.

- (a) **Recall.** GD is successful provided initialized properly.
- (b) **Focus.** Initialization in the context of randomly generated planted $W^* \in \mathbb{R}^{m \times d}$:
 - Networks with random weights is an active research area, define initial loss landscape.
 - Closely related to random feature methods (Rahimi & Recht [09]).

- (a) **Recall.** GD is successful provided initialized properly.
- (b) **Focus.** Initialization in the context of randomly generated planted $W^* \in \mathbb{R}^{m \times d}$:
 - Networks with random weights is an active research area, define initial loss landscape.
 - Closely related to random feature methods (Rahimi & Recht [09]).
 - Well approximate dynamical systems (Gonon et al. [20])

- (a) **Recall.** GD is successful provided initialized properly.
- (b) **Focus.** Initialization in the context of randomly generated planted $W^* \in \mathbb{R}^{m \times d}$:
 - Networks with random weights is an active research area, define initial loss landscape.
 - Closely related to random feature methods (Rahimi & Recht [09]).
 - Well approximate dynamical systems (Gonon et al. [20])
 - Also studied in the context of extreme learning machine (Huang et al. [06]), and in random matrix theory (Pennington & Worah [17]).

- (a) **Recall.** GD is successful provided initialized properly.
- (b) **Focus.** Initialization in the context of randomly generated planted $W^* \in \mathbb{R}^{m \times d}$:
 - Networks with random weights is an active research area, define initial loss landscape.
 - Closely related to random feature methods (Rahimi & Recht [09]).
 - Well approximate dynamical systems (Gonon et al. [20])
 - Also studied in the context of extreme learning machine (Huang et al. [06]), and in random matrix theory (Pennington & Worah [17]).
- (c) Intuition.

- (a) **Recall.** GD is successful provided initialized properly.
- (b) **Focus.** Initialization in the context of randomly generated planted $W^* \in \mathbb{R}^{m \times d}$:
 - Networks with random weights is an active research area, define initial loss landscape.
 - Closely related to random feature methods (Rahimi & Recht [09]).
 - Well approximate dynamical systems (Gonon et al. [20])
 - Also studied in the context of extreme learning machine (Huang et al. [06]), and in random matrix theory (Pennington & Worah [17]).
- (c) Intuition.
 - Value of loss function is determined by spectrum of $W^TW (W^*)^TW^*$ and data moments.

- (a) **Recall.** GD is successful provided initialized properly.
- (b) **Focus.** Initialization in the context of randomly generated planted $W^* \in \mathbb{R}^{m \times d}$:
 - Networks with random weights is an active research area, define initial loss landscape.
 - Closely related to random feature methods (Rahimi & Recht [09]).
 - Well approximate dynamical systems (Gonon et al. [20])
 - Also studied in the context of extreme learning machine (Huang et al. [06]), and in random matrix theory (Pennington & Worah [17]).
- (c) Intuition.
 - Value of loss function is determined by spectrum of $W^TW (W^*)^TW^*$ and data moments.
 - Spectrum of Wishart matrices $(W^*)^T W^*$ are tightly concentrated. Semicircle law by Bai & Yin [88,93].

- (a) **Recall.** GD is successful provided initialized properly.
- (b) **Focus.** Initialization in the context of randomly generated planted $W^* \in \mathbb{R}^{m \times d}$:
 - Networks with random weights is an active research area, define initial loss landscape.
 - Closely related to random feature methods (Rahimi & Recht [09]).
 - Well approximate dynamical systems (Gonon et al. [20])
 - Also studied in the context of extreme learning machine (Huang et al. [06]), and in random matrix theory (Pennington & Worah [17]).
- (c) Intuition.
 - Value of loss function is determined by spectrum of $W^TW (W^*)^TW^*$ and data moments.
 - Spectrum of Wishart matrices $(W^*)^T W^*$ are tightly concentrated. Semicircle law by Bai & Yin [88,93].
 - Thus the spectrum of $W^TW (W^*)^TW^*$ can be controlled by choosing W appropriately.





Theorem (Gamarnik, K.; and Zadik, 2020)

Suppose $W^* \in \mathbb{R}^{m \times d}$ has centered i.i.d. entries with unit variance, finite fourth moment. Data $X_i \in \mathbb{R}^d$, $1 \leq i \leq N$ has i.i.d. centered sub-Gaussian coordinates. Initialize W_0 so that $W_0^T W_0 = mI_{d \times d}$. Then, with high probability,

$$\widehat{\mathcal{L}}(W_0) < \frac{1}{2}C_5\sigma_{\mathsf{min}}(W^*)^4,$$

provided $m > Cd^2$, for a large constant C > 0.

Theorem (Gamarnik, K.; and Zadik, 2020)

Suppose $W^* \in \mathbb{R}^{m \times d}$ has centered i.i.d. entries with unit variance, finite fourth moment. Data $X_i \in \mathbb{R}^d$, $1 \leq i \leq N$ has i.i.d. centered sub-Gaussian coordinates. Initialize W_0 so that $W_0^T W_0 = mI_{d \times d}$. Then, with high probability,

$$\widehat{\mathcal{L}}(W_0) < rac{1}{2}C_5\sigma_{\mathsf{min}}(W^*)^4,$$

provided $m > Cd^2$, for a large constant C > 0.

(a) A deterministic initialization beats the energy barrier, provided the network is sufficiently overparametrized $(m = \Omega(d^2))$. Based on the semicircle law.



13 / 18

Theorem (Gamarnik, K.; and Zadik, 2020)

Suppose $W^* \in \mathbb{R}^{m \times d}$ has centered i.i.d. entries with unit variance, finite fourth moment. Data $X_i \in \mathbb{R}^d$, $1 \leq i \leq N$ has i.i.d. centered sub-Gaussian coordinates. Initialize W_0 so that $W_0^T W_0 = m I_{d \times d}$. Then, with high probability,

$$\widehat{\mathcal{L}}(W_0) < \frac{1}{2}C_5\sigma_{\min}(W^*)^4,$$

provided $m > Cd^2$, for a large constant C > 0.

- (a) A deterministic initialization beats the energy barrier, provided the network is sufficiently overparametrized $(m = \Omega(d^2))$. Based on the semicircle law.
- (b) An analogous result for the population risk. In the case of W^* having i.i.d. standard normal entries, non-asymptotic guarantees are available.

Overview

- Intro and Motivation
- Main Results: Optimization Landscape
- Main Results: Initialization
- Main Results: Generalization



Main question. "What is the smallest number of samples required to claim that small empirical risk also controls the generalization error?"



Main question. "What is the smallest number of samples required to claim that small empirical risk also controls the generalization error?"

Theorem (Gamarnik, K.; and Zadik, 2020)

Let $X_i \in \mathbb{R}^d$, $1 \le i \le N$ be the data (not necessarily random). $S \triangleq \{A \in \mathbb{R}^{d \times d} : A^T = A\}$.

Main question. "What is the smallest number of samples required to claim that small empirical risk also controls the generalization error?"

Theorem (Gamarnik, K.; and Zadik, 2020)

Let $X_i \in \mathbb{R}^d$, $1 \le i \le N$ be the data (not necessarily random). $S \triangleq \{A \in \mathbb{R}^{d \times d} : A^T = A\}$.

• Suppose $\operatorname{span}(X_iX_i^T: 1 \leq i \leq N) = \mathcal{S}$, and $\widehat{m} \in \mathbb{N}$ arbitrary. Then, for any $W \in \mathbb{R}^{\widehat{m} \times d}$ "interpolating" the data $(f(W; X_i) = f(W^*; X_i), 1 \leq i \leq N)$, $W^TW = (W^*)^TW^*$. Thus, W generalizes well: $\mathcal{L}(W) = 0$.

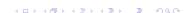
15 / 18

Main question. "What is the smallest number of samples required to claim that small empirical risk also controls the generalization error?"

Theorem (Gamarnik, K.; and Zadik, 2020)

Let $X_i \in \mathbb{R}^d$, $1 \le i \le N$ be the data (not necessarily random). $S \triangleq \{A \in \mathbb{R}^{d \times d} : A^T = A\}$.

- Suppose $\operatorname{span}(X_iX_i^T:1\leq i\leq N)=\mathcal{S}$, and $\widehat{m}\in\mathbb{N}$ arbitrary. Then, for any $W\in\mathbb{R}^{\widehat{m}\times d}$ "interpolating" the data $(f(W;X_i)=f(W^*;X_i),\ 1\leq i\leq N)$, $W^TW=(W^*)^TW^*$. Thus, W generalizes well: $\mathcal{L}(W)=0$.
- Suppose $\operatorname{span}(X_iX_i^T: 1 \leq i \leq N) \subsetneq \mathcal{S}$. Then for any $\widehat{m} \in \mathbb{N}$, there exists a $W \in \mathbb{R}^{\widehat{m} \times d}$ such that while W interpolates the data $(f(W; X_i) = f(W^*; X_i))$ for every i, $W^TW \neq (W^*)^TW^*$. In particular, $\mathcal{L}(W) > 0$ (where \mathcal{L} is defined w.r.t. any jointly continuous distribution on \mathbb{R}^d).







Remarks.

(a) If $\operatorname{span}(X_iX_i^T: 1 \le i \le N) = \mathcal{S}$, then any minimizer W of the empirical risk $\widehat{\mathcal{L}}(\cdot)$ has necessarily zero generalization error.



- (a) If $\operatorname{span}(X_iX_i^T: 1 \le i \le N) = \mathcal{S}$, then any minimizer W of the empirical risk $\widehat{\mathcal{L}}(\cdot)$ has necessarily zero generalization error.
- (b) The condition is not retrospective in manner: the condition $\operatorname{span}(X_iX_i^T:1\leq i\leq N)=\mathcal{S}$ can be checked ahead of optimization task.

- (a) If span(X_iX_i^T: 1 ≤ i ≤ N) = S, then any minimizer W of the empirical risk Â(·) has necessarily zero generalization error.
 (b) The condition is not retrospective in manner: the condition span(X:X^T: 1 ≤ i ≤ N) = S
- (b) The condition is not retrospective in manner: the condition $\operatorname{span}(X_iX_i^T:1\leq i\leq N)=\mathcal{S}$ can be checked ahead of optimization task.
- (c) No randomness. A purely geometric, necessary and sufficient geometric condition.

- (a) If span(X_iX_i^T: 1 ≤ i ≤ N) = S, then any minimizer W of the empirical risk Â(·) has necessarily zero generalization error.
 (b) The condition is not retrospective in manner: the condition span(X:X^T: 1 ≤ i ≤ N) = S
- (b) The condition is not retrospective in manner: the condition $\operatorname{span}(X_iX_i^T:1\leq i\leq N)=\mathcal{S}$ can be checked ahead of optimization task.
- (c) No randomness. A purely geometric, necessary and sufficient geometric condition.
- (d) In case W has non-zero but small training loss $(\widehat{\mathcal{L}}(W))$, earlier results allow bounding $\|W^TW (W^*)^TW^*\|_F$, and $\mathcal{L}(W)$.



- (a) If span(X_iX_i^T: 1 ≤ i ≤ N) = S, then any minimizer W of the empirical risk Â(·) has necessarily zero generalization error.
 (b) The condition is not retrospective in manner; the condition span(X:X^T: 1 ≤ i ≤ N) = S
- (b) The condition is not retrospective in manner: the condition $\operatorname{span}(X_iX_i^T:1\leq i\leq N)=\mathcal{S}$ can be checked ahead of optimization task.
- (c) No randomness. A purely geometric, necessary and sufficient geometric condition.
- (d) In case W has non-zero but small training loss $(\widehat{\mathcal{L}}(W))$, earlier results allow bounding $\|W^TW (W^*)^TW^*\|_F$, and $\mathcal{L}(W)$.
- (e) The parameter $\widehat{m} \in \mathbb{N}$: the interpolating network need not have the same width m. Provided the span condition holds, **any** interpolant, albeit being over-parametrized, generalize well (i.e. in this case, overparametrization does **not** hurt generalization).



- (a) If span(X_iX_i^T: 1 ≤ i ≤ N) = S, then any minimizer W of the empirical risk L̂(·) has necessarily zero generalization error.
 (b) The condition is not retrospective in manner: the condition span(X_iX^T: 1 ≤ i ≤ N) = S
- (b) The condition is not retrospective in manner: the condition $\operatorname{span}(X_iX_i^T:1\leq i\leq N)=\mathcal{S}$ can be checked ahead of optimization task.
- (c) No randomness. A purely geometric, necessary and sufficient geometric condition.
- (d) In case W has non-zero but small training loss $(\widehat{\mathcal{L}}(W))$, earlier results allow bounding $\|W^TW (W^*)^TW^*\|_F$, and $\mathcal{L}(W)$.
- (e) The parameter $\widehat{m} \in \mathbb{N}$: the interpolating network need not have the same width m. Provided the span condition holds, **any** interpolant, albeit being over-parametrized, generalize well (i.e. in this case, overparametrization does **not** hurt generalization).
- (f) **Theorem.** As soon as $N \ge d(d+1)/2$, random data $X_i \in \mathbb{R}^d$, $1 \le i \le N$ enjoys $\operatorname{span}(X_i X_i^T : 1 \le i \le N) = \mathcal{S}$, with probability one.

Sample Complexity Bound for Planted Network.



Sample Complexity Bound for Planted Network.

Theorem (Gamarnik, K.; and Zadik, 2020)

Let $X_i \in \mathbb{R}^d$, $1 \le i \le N$ be i.i.d. with a jointly continuous distribution (on \mathbb{R}^d). Suppose $W^* \in \mathbb{R}^{m \times d}$ with $\operatorname{rank}(W^*) = d$ and the labels $Y_i = f(W^*; X_i) = \sum_{1 \le i \le m} \langle W_i^*, X_i \rangle^2$.

- Suppose $N \ge d(d+1)/2$, and $\widehat{m} \in \mathbb{N}$. Then, with probability one over X_i , $1 \le i \le N$ the following holds: if $f(W; X_i) = f(W^*; X_i)$, $1 \le i \le N$, then $f(W; x) = f(W^*; x)$ for every $x \in \mathbb{R}^d$.
- Suppose X_i has centered i.i.d. coordinates with variance μ_2 and (finite) fourth moment μ_4 , and N < d(d+1)/2. Then, there exists a $W \in \mathbb{R}^{m \times d}$ such that while $\widehat{\mathcal{L}}(W) = 0$ (namely $f(W; X_i) = f(W^*; X_i)$ for $1 \le i \le N$),

$$\mathcal{L}(W) \ge \min\{\mu_4 - \mu_2^2, 2\mu_2^2\} \sigma_{\min}(W^*)^4.$$



Sample Complexity Bound for Planted Network.

Theorem (Gamarnik, K.; and Zadik, 2020)

Let $X_i \in \mathbb{R}^d$, $1 \le i \le N$ be i.i.d. with a jointly continuous distribution (on \mathbb{R}^d). Suppose $W^* \in \mathbb{R}^{m \times d}$ with $\operatorname{rank}(W^*) = d$ and the labels $Y_i = f(W^*; X_i) = \sum_{1 \le i \le m} \langle W_i^*, X_i \rangle^2$.

- Suppose $N \ge d(d+1)/2$, and $\widehat{m} \in \mathbb{N}$. Then, with probability one over X_i , $1 \le i \le N$ the following holds: if $f(W; X_i) = f(W^*; X_i)$, $1 \le i \le N$, then $f(W; x) = f(W^*; x)$ for every $x \in \mathbb{R}^d$.
- Suppose X_i has centered i.i.d. coordinates with variance μ_2 and (finite) fourth moment μ_4 , and N < d(d+1)/2. Then, there exists a $W \in \mathbb{R}^{m \times d}$ such that while $\widehat{\mathcal{L}}(W) = 0$ (namely $f(W; X_i) = f(W^*; X_i)$ for $1 \le i \le N$),

$$\mathcal{L}(W) \ge \min\{\mu_4 - \mu_2^2, 2\mu_2^2\} \sigma_{\min}(W^*)^4.$$

The lower bound in second part coincides with the (earlier) energy barrier.



Thank you!

