

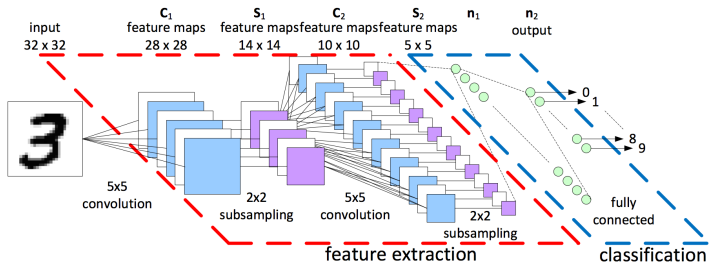
A Mean Field View of the Landscape of Two-Layers Neural Networks

Song Mei

Stanford University

November 14, 2018

Joint work with Andrea Montanari and Phan-Minh Nguyen



Empirical surprise [Zhang, Bengio, Hardt, Recht, Vinyals, 2016]

- ▶ Overparameterized regime.
- ▶ Efficiently fit all the data.
- ▶ Generalize well.

Empirical surprise

- ▶ Overparameterized regime.
- ▶ Efficiently fit all the data.
- ▶ Generalize well.

Questions

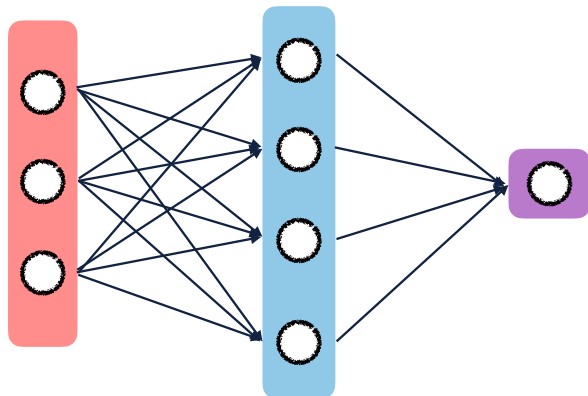
- ▶ Why can complex neural network be **optimized** efficiently?
- ▶ Why does overparameterization not harm **generalization**?

Two-layers neural networks

Input layer

Hidden layer

Output layer



Two-layers neural networks

▶ Parameter: $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N) \in \mathbb{R}^{N \times D}$.

▶ Prediction:

$$\hat{y}(\boldsymbol{x}; \boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_i).$$

▶ An example: $\boldsymbol{\theta}_i = (\boldsymbol{a}_i, \boldsymbol{w}_i)$, $\sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_i) = \boldsymbol{a}_i \sigma(\langle \boldsymbol{x}, \boldsymbol{w}_i \rangle)$.

▶ Data distribution: $(\boldsymbol{x}, y) \sim \mathbb{P}_{\boldsymbol{x}, y}$.

▶ Risk function:

$$R_N(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{x}, y} \left[\left(y - \frac{1}{N} \sum_{i=1}^N \sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_i) \right)^2 \right].$$

Two-layers neural networks

- ▶ Parameter: $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N) \in \mathbb{R}^{N \times D}$.
- ▶ Prediction:

$$\hat{y}(\boldsymbol{x}; \boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_i).$$

- ▶ An example: $\boldsymbol{\theta}_i = (\boldsymbol{a}_i, \boldsymbol{w}_i)$, $\sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_i) = \boldsymbol{a}_i \sigma(\langle \boldsymbol{x}, \boldsymbol{w}_i \rangle)$.
- ▶ Data distribution: $(\boldsymbol{x}, y) \sim \mathbb{P}_{\boldsymbol{x}, y}$.
- ▶ Risk function:

$$R_N(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{x}, y} \left[\left(y - \frac{1}{N} \sum_{i=1}^N \sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_i) \right)^2 \right].$$

Two-layers neural networks

▶ Parameter: $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N) \in \mathbb{R}^{N \times D}$.

▶ Prediction:

$$\hat{y}(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \sigma_*(\mathbf{x}; \boldsymbol{\theta}_i).$$

▶ An example: $\boldsymbol{\theta}_i = (\mathbf{a}_i, \mathbf{w}_i)$, $\sigma_*(\mathbf{x}; \boldsymbol{\theta}_i) = \mathbf{a}_i \sigma(\langle \mathbf{x}, \mathbf{w}_i \rangle)$.

▶ Data distribution: $(\mathbf{x}, y) \sim \mathbb{P}_{\mathbf{x}, y}$.

▶ Risk function:

$$R_N(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, y} \left[\left(y - \frac{1}{N} \sum_{i=1}^N \sigma_*(\mathbf{x}; \boldsymbol{\theta}_i) \right)^2 \right].$$

Two-layers neural networks

▶ Parameter: $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N) \in \mathbb{R}^{N \times D}$.

▶ Prediction:

$$\hat{y}(\boldsymbol{x}; \boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_i).$$

▶ An example: $\boldsymbol{\theta}_i = (\boldsymbol{a}_i, \boldsymbol{w}_i)$, $\sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_i) = \boldsymbol{a}_i \sigma(\langle \boldsymbol{x}, \boldsymbol{w}_i \rangle)$.

▶ Data distribution: $(\boldsymbol{x}, y) \sim \mathbb{P}_{\boldsymbol{x}, y}$.

▶ Risk function:

$$R_N(\boldsymbol{\theta}) = \mathbb{E}_{\boldsymbol{x}, y} \left[\left(y - \frac{1}{N} \sum_{i=1}^N \sigma_*(\boldsymbol{x}; \boldsymbol{\theta}_i) \right)^2 \right].$$

Two-layers neural networks

▶ Parameter: $\boldsymbol{\theta} = (\boldsymbol{\theta}_1, \dots, \boldsymbol{\theta}_N) \in \mathbb{R}^{N \times D}$.

▶ Prediction:

$$\hat{y}(\mathbf{x}; \boldsymbol{\theta}) = \frac{1}{N} \sum_{i=1}^N \sigma_*(\mathbf{x}; \boldsymbol{\theta}_i).$$

▶ An example: $\boldsymbol{\theta}_i = (\mathbf{a}_i, \mathbf{w}_i)$, $\sigma_*(\mathbf{x}; \boldsymbol{\theta}_i) = \mathbf{a}_i \sigma(\langle \mathbf{x}, \mathbf{w}_i \rangle)$.

▶ Data distribution: $(\mathbf{x}, y) \sim \mathbb{P}_{\mathbf{x}, y}$.

▶ Risk function:

$$R_N(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, y} \left[\left(y - \frac{1}{N} \sum_{i=1}^N \sigma_*(\mathbf{x}; \boldsymbol{\theta}_i) \right)^2 \right].$$

Two-layers neural networks

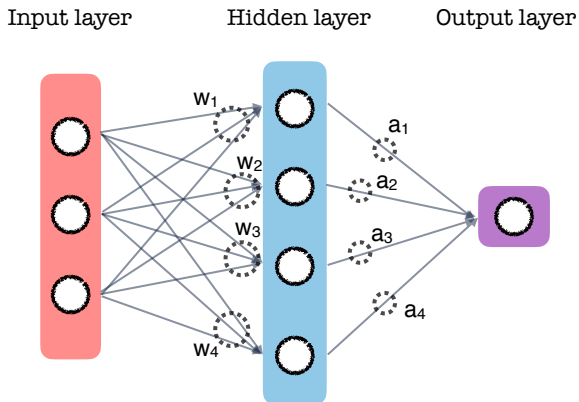


Figure: $\theta_i = (a_i, w_i)$.

Related literatures (before 2018)

$$R_N(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, y} \left[\left(y - \frac{1}{N} \sum_{j=1}^N \sigma_*(\mathbf{x}; \boldsymbol{\theta}_j) \right)^2 \right].$$

- ▶ Landscape analysis: [Soudry, Carmon, 2016], [Freeman, Bruna, 2016], [Ge, Lee, Ma, 2017], [Soltanolkotabi, Javanmard, Lee, 2017], [Zhong, Song, Jain, Bartlett, Dhillon, 2017]...
- ▶ Optimization dynamics: [Tian, 2017], [Soltanolkotabi, 2017], [Li, Yuan, 2017]...
- ▶ Generalization: [Bartlett, Foster, Telgarsky, 2017], [Neyshabur, Bhojanapalli, McAllester, Srebro, 2017]...

Related literatures (before 2018)

$$R_N(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, y} \left[\left(y - \frac{1}{N} \sum_{j=1}^N \sigma_*(\mathbf{x}; \boldsymbol{\theta}_j) \right)^2 \right].$$

- ▶ Landscape analysis: [Soudry, Carmon, 2016], [Freeman, Bruna, 2016], [Ge, Lee, Ma, 2017], [Soltanolkotabi, Javanmard, Lee, 2017], [Zhong, Song, Jain, Bartlett, Dhillon, 2017]...
- ▶ Optimization dynamics: [Tian, 2017], [Soltanolkotabi, 2017], [Li, Yuan, 2017]...
- ▶ Generalization: [Bartlett, Foster, Telgarsky, 2017], [Neyshabur, Bhojanapalli, McAllester, Srebro, 2017]...

Related literatures (before 2018)

$$R_N(\boldsymbol{\theta}) = \mathbb{E}_{\mathbf{x}, y} \left[\left(y - \frac{1}{N} \sum_{j=1}^N \sigma_*(\mathbf{x}; \boldsymbol{\theta}_j) \right)^2 \right].$$

- ▶ Landscape analysis: [Soudry, Carmon, 2016], [Freeman, Bruna, 2016], [Ge, Lee, Ma, 2017], [Soltanolkotabi, Javanmard, Lee, 2017], [Zhong, Song, Jain, Bartlett, Dhillon, 2017]...
- ▶ Optimization dynamics: [Tian, 2017], [Soltanolkotabi, 2017], [Li, Yuan, 2017]...
- ▶ Generalization: [Bartlett, Foster, Telgarsky, 2017], [Neyshabur, Bhojanapalli, McAllester, Srebro, 2017]...

Overparameterization: what happens for large N ?

- ▶ [Bengio, et. al, 2006]. Expand the square

$$R_N(\boldsymbol{\theta}) = \mathbb{E}[y^2] + \frac{2}{N} \sum_{i=1}^N V(\boldsymbol{\theta}_i) + \frac{1}{N^2} \sum_{i,j=1}^N U(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j),$$

where

$$\begin{aligned} V(\boldsymbol{\theta}_i) &= -\mathbb{E}[y\sigma_*(\mathbf{x}; \boldsymbol{\theta}_i)], \\ U(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) &= \mathbb{E}[\sigma_*(\mathbf{x}; \boldsymbol{\theta}_i)\sigma_*(\mathbf{x}; \boldsymbol{\theta}_j)]. \end{aligned}$$

- ▶ R_N depends on $(\boldsymbol{\theta}_i)_{i \leq N}$ through $\rho_N = (1/N) \sum_{i=1}^N \delta_{\boldsymbol{\theta}_i}$.
- ▶ Motivate us to define $R(\rho)$, $\rho \in \mathcal{P}(\mathbb{R}^D)$,

$$R(\rho) = \mathbb{E}[y^2] + 2 \int V(\boldsymbol{\theta})\rho(d\boldsymbol{\theta}) + \int U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)\rho(d\boldsymbol{\theta}_1)\rho(d\boldsymbol{\theta}_2).$$

Overparameterization: what happens for large N ?

- ▶ [Bengio, et. al, 2006]. Expand the square

$$R_N(\boldsymbol{\theta}) = \mathbb{E}[y^2] + \frac{2}{N} \sum_{i=1}^N V(\boldsymbol{\theta}_i) + \frac{1}{N^2} \sum_{i,j=1}^N U(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j),$$

where

$$\begin{aligned} V(\boldsymbol{\theta}_i) &= -\mathbb{E}[y\sigma_*(\mathbf{x}; \boldsymbol{\theta}_i)], \\ U(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) &= \mathbb{E}[\sigma_*(\mathbf{x}; \boldsymbol{\theta}_i)\sigma_*(\mathbf{x}; \boldsymbol{\theta}_j)]. \end{aligned}$$

- ▶ R_N depends on $(\boldsymbol{\theta}_i)_{i \leq N}$ through $\rho_N = (1/N) \sum_{i=1}^N \delta_{\boldsymbol{\theta}_i}$.
- ▶ Motivate us to define $R(\rho)$, $\rho \in \mathcal{P}(\mathbb{R}^D)$,

$$R(\rho) = \mathbb{E}[y^2] + 2 \int V(\boldsymbol{\theta})\rho(d\boldsymbol{\theta}) + \int U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)\rho(d\boldsymbol{\theta}_1)\rho(d\boldsymbol{\theta}_2).$$

Overparameterization: what happens for large N ?

- ▶ [Bengio, et. al, 2006]. Expand the square

$$R_N(\boldsymbol{\theta}) = \mathbb{E}[y^2] + \frac{2}{N} \sum_{i=1}^N V(\boldsymbol{\theta}_i) + \frac{1}{N^2} \sum_{i,j=1}^N U(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j),$$

where

$$\begin{aligned} V(\boldsymbol{\theta}_i) &= -\mathbb{E}[y\sigma_*(\mathbf{x}; \boldsymbol{\theta}_i)], \\ U(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j) &= \mathbb{E}[\sigma_*(\mathbf{x}; \boldsymbol{\theta}_i)\sigma_*(\mathbf{x}; \boldsymbol{\theta}_j)]. \end{aligned}$$

- ▶ R_N depends on $(\boldsymbol{\theta}_i)_{i \leq N}$ through $\rho_N = (1/N) \sum_{i=1}^N \delta_{\boldsymbol{\theta}_i}$.
- ▶ Motivate us to define $R(\rho)$, $\rho \in \mathcal{P}(\mathbb{R}^D)$,

$$R(\rho) = \mathbb{E}[y^2] + 2 \int V(\boldsymbol{\theta})\rho(d\boldsymbol{\theta}) + \int U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2)\rho(d\boldsymbol{\theta}_1)\rho(d\boldsymbol{\theta}_2).$$

Overparameterization: what happens for large N ?

Correspondence $R_N(\boldsymbol{\theta}) = R((1/N) \sum_{i=1}^N \delta_{\boldsymbol{\theta}_i})$, where

$$R_N(\boldsymbol{\theta}) = \mathbb{E}[\mathbf{y}^2] + \frac{2}{N} \sum_{i=1}^N V(\boldsymbol{\theta}_i) + \frac{1}{N^2} \sum_{i,j=1}^N U(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j),$$

$$R(\rho) = \mathbb{E}[\mathbf{y}^2] + 2 \int V(\boldsymbol{\theta}) \rho(d\boldsymbol{\theta}) + \int U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \rho(d\boldsymbol{\theta}_1) \rho(d\boldsymbol{\theta}_2).$$

Overparameterization: what happens for large N ?

Correspondence $R_N(\boldsymbol{\theta}) = R((1/N) \sum_{i=1}^N \delta_{\boldsymbol{\theta}_i})$, where

$$R_N(\boldsymbol{\theta}) = \mathbb{E}[\mathbf{y}^2] + \frac{2}{N} \sum_{i=1}^N V(\boldsymbol{\theta}_i) + \frac{1}{N^2} \sum_{i,j=1}^N U(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j),$$

$$R(\rho) = \mathbb{E}[\mathbf{y}^2] + 2 \int V(\boldsymbol{\theta}) \rho(d\boldsymbol{\theta}) + \int U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \rho(d\boldsymbol{\theta}_1) \rho(d\boldsymbol{\theta}_2).$$

What is the relationship of minimum value of R_N and R ?

Overparameterization: what happens for large N ?

Correspondence $R_N(\boldsymbol{\theta}) = R((1/N) \sum_{i=1}^N \delta_{\boldsymbol{\theta}_i})$, where

$$R_N(\boldsymbol{\theta}) = \mathbb{E}[y^2] + \frac{2}{N} \sum_{i=1}^N V(\boldsymbol{\theta}_i) + \frac{1}{N^2} \sum_{i,j=1}^N U(\boldsymbol{\theta}_i, \boldsymbol{\theta}_j),$$

$$R(\rho) = \mathbb{E}[y^2] + 2 \int V(\boldsymbol{\theta}) \rho(d\boldsymbol{\theta}) + \int U(\boldsymbol{\theta}_1, \boldsymbol{\theta}_2) \rho(d\boldsymbol{\theta}_1) \rho(d\boldsymbol{\theta}_2).$$

What is the relationship of minimum value of R_N and R ?

Lemma

If $\int U(\boldsymbol{\theta}, \boldsymbol{\theta}) \rho_{\text{opt}}(d\boldsymbol{\theta}) < K$, then

$$\inf_{\rho} R(\rho) \leq \inf_{\boldsymbol{\theta}} R_N(\boldsymbol{\theta}) \leq \inf_{\rho} R(\rho) + \frac{K}{N}.$$

How to optimize $R(\rho)$?

[Bengio, et. al, 2006] proposed to optimize over ρ

$$R(\rho) = \mathbb{E}[y^2] + 2 \int V(\theta)\rho(d\theta) + \int U(\theta_1, \theta_2)\rho(d\theta_1)\rho(d\theta_2).$$

Exponential bases functions to discretize ρ !

How to optimize $R(\rho)$?

[Bengio, et. al, 2006] proposed to optimize over ρ

$$R(\rho) = \mathbb{E}[y^2] + 2 \int V(\theta) \rho(d\theta) + \int U(\theta_1, \theta_2) \rho(d\theta_1) \rho(d\theta_2).$$

Exponential bases functions to discretize ρ !

[This work]: run SGD on θ , and give a scaling limit dynamics for ρ .

SGD and distributional dynamics (DD)

- ▶ SGD for θ^k , with $(\mathbf{x}_k, \mathbf{y}_k) \sim \mathbb{P}_{\mathbf{x}, \mathbf{y}}, i \in [N]$,

$$\theta_i^{k+1} = \theta_i^k - 2s_k \nabla_{\theta_i} \ell(\mathbf{x}_k, \mathbf{y}_k; \theta^k). \quad (\text{SGD})$$

SGD and distributional dynamics (DD)

- ▶ SGD for θ^k , with $(\mathbf{x}_k, \mathbf{y}_k) \sim \mathbb{P}_{\mathbf{x}, \mathbf{y}}$, $i \in [N]$,

$$\theta_i^{k+1} = \theta_i^k - 2s_k \nabla_{\theta_i} \ell(\mathbf{x}_k, \mathbf{y}_k; \theta^k). \quad (\text{SGD})$$

- ▶ Claim: $s_k = \varepsilon \xi(k\varepsilon)$, $k = t/\varepsilon$, $N \rightarrow \infty$, $\varepsilon \rightarrow 0$:

$$\hat{\rho}_k^{(N)} \equiv \frac{1}{N} \sum_{i=1}^N \delta_{\theta_i^k} \Rightarrow \rho_t.$$

SGD and distributional dynamics (DD)

- ▶ SGD for θ^k , with $(\mathbf{x}_k, \mathbf{y}_k) \sim \mathbb{P}_{\mathbf{x}, \mathbf{y}}$, $i \in [N]$,

$$\theta_i^{k+1} = \theta_i^k - 2s_k \nabla_{\theta_i} \ell(\mathbf{x}_k, \mathbf{y}_k; \theta^k). \quad (\text{SGD})$$

- ▶ Claim: $s_k = \varepsilon \xi(k\varepsilon)$, $k = t/\varepsilon$, $N \rightarrow \infty$, $\varepsilon \rightarrow 0$:

$$\hat{\rho}_k^{(N)} \equiv \frac{1}{N} \sum_{i=1}^N \delta_{\theta_i^k} \Rightarrow \rho_t.$$

- ▶ Distributional dynamics (DD) for ρ_t ,

$$\partial_t \rho_t(\theta) = 2\xi(t) \nabla_{\theta} \cdot (\rho_t(\theta) \nabla_{\theta} \Psi(\theta; \rho_t)), \quad (\text{DD})$$

where

$$\Psi(\theta; \rho) = \frac{\delta R(\rho)}{\delta \rho(\theta)} = V(\theta) + \int U(\theta, \theta') \rho(d\theta').$$

More precisely

Assumption

(i) σ_* bounded; (ii) $\nabla_{\theta}\sigma_*(\mathbf{x}; \boldsymbol{\theta})$ sub-Gaussian; (iii) $\nabla V, \nabla U$ bdd. Lipschitz.

Theorem (M., Montanari, Nguyen, 2018)

Let $(\boldsymbol{\theta}_i^0)_{i \leq N} \sim_{iid} \rho_0$. Then, $\forall f$ bounded Lipschitz:

$$\sup_{t \leq T} \left| \frac{1}{N} \sum_{i=1}^N f(\boldsymbol{\theta}_i^{[t/\varepsilon]}) - \int f(\boldsymbol{\theta}) \rho_t(\boldsymbol{\theta}) \right| \leq K e^{KT} \text{err}_{N,D}(z),$$

where

$$\text{err}_{N,D}(z) \equiv \sqrt{\frac{1}{N} \vee \varepsilon} \cdot \left[\sqrt{D \vee \log \frac{N}{\varepsilon}} + z \right],$$

with probability at least $1 - 4e^{-z^2/2}$.

More precisely

Assumption

(i) σ_* bounded; (ii) $\nabla_{\theta} \sigma_*(\mathbf{x}; \boldsymbol{\theta})$ sub-Gaussian; (iii) $\nabla V, \nabla U$ bdd. Lipschitz.

Theorem (M., Montanari, Nguyen, 2018)

Let $(\boldsymbol{\theta}_i^0)_{i \leq N} \sim_{iid} \rho_0$. Then, $\forall f$ bounded Lipschitz:

$$\sup_{t \leq T} \left| \frac{1}{N} \sum_{i=1}^N f(\boldsymbol{\theta}_i^{\lfloor t/\varepsilon \rfloor}) - \int f(\boldsymbol{\theta}) \rho_t(\boldsymbol{\theta}) \right| \leq K e^{KT} \text{err}_{N,D}(z),$$

where

$$\text{err}_{N,D}(z) \equiv \sqrt{\frac{1}{N} \vee \varepsilon} \cdot \left[\sqrt{D \vee \log \frac{N}{\varepsilon}} + z \right],$$

with probability at least $1 - 4e^{-z^2/2}$.

Number of neurons, sample size, and dimensions

$$\text{err}_{N,D} \asymp \sqrt{\frac{D}{N}} \vee (D\varepsilon).$$

N : number of neurons; D : feature dimension; ε : stepsize.

► Small if $N \gg D$, $\varepsilon \ll 1/D$. This is very practical!

Number of neurons, sample size, and dimensions

$$\text{err}_{N,D} \asymp \sqrt{\frac{D}{N}} \vee (D\varepsilon).$$

N : number of neurons; D : feature dimension; ε : stepsize.

- ▶ Small if $N \gg D$, $\varepsilon \ll 1/D$. This is very practical!

An animation

- ▶ A specific model (classifying two Gaussians), fix dimension $d = 80$, $N = 200$.
- ▶ Run SGD and solve PDE.
- ▶ Track the norm statistics of the weights,
 $\hat{\mu}_k(s) = (1/N) \sum_{i=1}^N \delta_{\|\theta_i^k\|_2}$ versus $\mu_t(s) = \rho_t(\{\|\theta\|_2 = s\})$.

An animation

- ▶ A specific model (classifying two Gaussians), fix dimension $d = 80$, $N = 200$.
- ▶ Run SGD and solve PDE.
- ▶ Track the norm statistics of the weights,
 $\hat{\mu}_k(s) = (1/N) \sum_{i=1}^N \delta_{\|\theta_i^k\|_2}$ versus $\mu_t(s) = \rho_t(\{\|\theta\|_2 = s\})$.

An animation

- ▶ A specific model (classifying two Gaussians), fix dimension $d = 80$, $N = 200$.
- ▶ Run SGD and solve PDE.
- ▶ Track the norm statistics of the weights,
 $\hat{\mu}_k(s) = (1/N) \sum_{i=1}^N \delta_{\|\theta_i^k\|_2}$ versus $\mu_t(s) = \rho_t(\{\|\theta\|_2 = s\})$.

Message

Approximately $(1/N) \sum_{i=1}^N \delta_{\theta_i^k} \approx \rho_t$, where

$$\theta_i^{k+1} = \theta_i^k - 2s_k \nabla_{\theta_i} \ell(\mathbf{x}_k, \mathbf{y}_k; \theta^k), \quad i \in [N], \quad (\text{SGD})$$

$$\partial_t \rho_t(\theta) = 2\xi(t) \nabla_{\theta} \cdot (\rho_t(\theta) \nabla_{\theta} \Psi(\theta; \rho_t)). \quad (\text{DD})$$

Overparameterization $N \rightarrow \infty$ does not affect the limiting dynamics, and therefore

- ▶ Overparameterization does not slow down convergence!
- ▶ Overparameterization does not affect generalization!

Message

Approximately $(1/N) \sum_{i=1}^N \delta_{\theta_i^k} \approx \rho_t$, where

$$\theta_i^{k+1} = \theta_i^k - 2s_k \nabla_{\theta_i} \ell(\mathbf{x}_k, \mathbf{y}_k; \theta^k), \quad i \in [N], \quad (\text{SGD})$$

$$\partial_t \rho_t(\theta) = 2\xi(t) \nabla_{\theta} \cdot (\rho_t(\theta) \nabla_{\theta} \Psi(\theta; \rho_t)). \quad (\text{DD})$$

Overparameterization $N \rightarrow \infty$ does not affect the limiting dynamics, and therefore

- ▶ Overparameterization does not slow down convergence!
- ▶ Overparameterization does not affect generalization!

Message

Approximately $(1/N) \sum_{i=1}^N \delta_{\theta_i^k} \approx \rho_t$, where

$$\theta_i^{k+1} = \theta_i^k - 2s_k \nabla_{\theta_i} \ell(\mathbf{x}_k, \mathbf{y}_k; \theta^k), \quad i \in [N], \quad (\text{SGD})$$

$$\partial_t \rho_t(\theta) = 2\xi(t) \nabla_{\theta} \cdot (\rho_t(\theta) \nabla_{\theta} \Psi(\theta; \rho_t)). \quad (\text{DD})$$

Overparameterization $N \rightarrow \infty$ does not affect the limiting dynamics, and therefore

- ▶ Overparameterization does not slow down convergence!
- ▶ Overparameterization does not affect generalization!

What is this?

$$\partial_t \rho_t = \nabla_{\theta} \cdot \left(\rho_t \nabla_{\theta} \Psi(\theta; \rho_t) \right).$$

Existence and uniqueness: [Sznitman, 1991].

- ▶ Physics: nonlinear transport equation describing motions of particles with pairwise interaction (mean field approach).
- ▶ Math: Gradient flow of $R(\rho)$...
- ▶ ... in the metric space $(\mathcal{P}(\mathbb{R}^D), W_2)$.
- ▶ [Jordan, Kinderlehrer, Otto, 1998; Ambrosio, Gigli, Savaré, 2006; Carrillo, McCann, Villani, 2013; ...]

What is this?

$$\partial_t \rho_t = \nabla_{\theta} \cdot \left(\rho_t \nabla_{\theta} \Psi(\theta; \rho_t) \right).$$

Existence and uniqueness: [Sznitman, 1991].

- ▶ Physics: nonlinear transport equation describing motions of particles with pairwise interaction (mean field approach).
- ▶ Math: Gradient flow of $R(\rho)$...
- ▶ ... in the metric space $(\mathcal{P}(\mathbb{R}^D), W_2)$.
- ▶ [Jordan, Kinderlehrer, Otto, 1998; Ambrosio, Gigli, Savaré, 2006; Carrillo, McCann, Villani, 2013; ...]

What is this?

$$\partial_t \rho_t = \nabla_{\theta} \cdot \left(\rho_t \nabla_{\theta} \Psi(\theta; \rho_t) \right).$$

Existence and uniqueness: [Sznitman, 1991].

- ▶ Physics: nonlinear transport equation describing motions of particles with pairwise interaction (mean field approach).
- ▶ Math: Gradient flow of $R(\rho)$...
- ▶ ... in the metric space $(\mathcal{P}(\mathbb{R}^D), W_2)$.
- ▶ [Jordan, Kinderlehrer, Otto, 1998; Ambrosio, Gigli, Savaré, 2006; Carrillo, McCann, Villani, 2013; ...]

What is this?

$$\partial_t \rho_t = \nabla_{\theta} \cdot \left(\rho_t \nabla_{\theta} \Psi(\theta; \rho_t) \right).$$

Existence and uniqueness: [Sznitman, 1991].

- ▶ Physics: nonlinear transport equation describing motions of particles with pairwise interaction (mean field approach).
- ▶ Math: Gradient flow of $R(\rho)$...
- ▶ ... in the metric space $(\mathcal{P}(\mathbb{R}^D), W_2)$.
- ▶ [Jordan, Kinderlehrer, Otto, 1998; Ambrosio, Gigli, Savaré, 2006; Carrillo, McCann, Villani, 2013; ...]

What is this?

$$\partial_t \rho_t = \nabla_{\theta} \cdot \left(\rho_t \nabla_{\theta} \Psi(\theta; \rho_t) \right).$$

Existence and uniqueness: [Sznitman, 1991].

- ▶ Physics: nonlinear transport equation describing motions of particles with pairwise interaction (mean field approach).
- ▶ Math: Gradient flow of $R(\rho)$...
- ▶ ... in the metric space $(\mathcal{P}(\mathbb{R}^D), W_2)$.
- ▶ [Jordan, Kinderlehrer, Otto, 1998; Ambrosio, Gigli, Savaré, 2006; Carrillo, McCann, Villani, 2013; ...]

What is gradient flow in a metric space?

Euclidean space gradient descent

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \varepsilon \nabla F(\mathbf{x}^k).$$

More insightfully

$$\begin{aligned}\mathbf{x}^{k+1} &= \arg \min_{\mathbf{x}} \left\{ F(\mathbf{x}^k) + \langle \nabla F(\mathbf{x}^k), \mathbf{x} - \mathbf{x}^k \rangle + \frac{1}{2\varepsilon} \|\mathbf{x} - \mathbf{x}^k\|_2^2 \right\} \\ &\approx \arg \min_{\mathbf{x}} \left\{ F(\mathbf{x}) + \frac{1}{2\varepsilon} \|\mathbf{x} - \mathbf{x}^k\|_2^2 \right\}\end{aligned}$$

What is gradient flow in a metric space?

Euclidean space gradient descent

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \varepsilon \nabla F(\mathbf{x}^k).$$

More insightfully

$$\begin{aligned}\mathbf{x}^{k+1} &= \arg \min_{\mathbf{x}} \left\{ F(\mathbf{x}^k) + \langle \nabla F(\mathbf{x}^k), \mathbf{x} - \mathbf{x}^k \rangle + \frac{1}{2\varepsilon} \|\mathbf{x} - \mathbf{x}^k\|_2^2 \right\} \\ &\approx \arg \min_{\mathbf{x}} \left\{ F(\mathbf{x}) + \frac{1}{2\varepsilon} \|\mathbf{x} - \mathbf{x}^k\|_2^2 \right\}\end{aligned}$$

What is gradient flow in a metric space?

Euclidean space gradient descent

$$\mathbf{x}^{k+1} = \mathbf{x}^k - \varepsilon \nabla F(\mathbf{x}^k).$$

More insightfully

$$\begin{aligned} \mathbf{x}^{k+1} &= \arg \min_{\mathbf{x}} \left\{ F(\mathbf{x}^k) + \langle \nabla F(\mathbf{x}^k), \mathbf{x} - \mathbf{x}^k \rangle + \frac{1}{2\varepsilon} \|\mathbf{x} - \mathbf{x}^k\|_2^2 \right\} \\ &\approx \arg \min_{\mathbf{x}} \left\{ F(\mathbf{x}) + \frac{1}{2\varepsilon} \|\mathbf{x} - \mathbf{x}^k\|_2^2 \right\} \end{aligned}$$

Use this as definition!

Metric space (M, d)

$$\mathbf{x}^{k+1} = \arg \min_{\mathbf{x}} \left\{ F(\mathbf{x}) + \frac{1}{2\varepsilon} d(\mathbf{x}, \mathbf{x}^k)^2 \right\}.$$

Wasserstein space $(\mathcal{P}(\mathbb{R}^D), W_2)$

$$\rho_{t+\varepsilon} = \arg \min_{\mathcal{P}(\mathbb{R}^D)} \left\{ R(\rho) + \frac{1}{2\varepsilon} W_2(\rho, \rho_t)^2 \right\}$$

$$\begin{aligned} \partial_t \rho_t &= \nabla_{\theta} \cdot \left(\rho_t \nabla_{\theta} \Psi(\theta; \rho_t) \right), \\ \Psi(\theta; \rho) &\equiv \frac{\delta R}{\delta \rho}(\theta). \end{aligned}$$

Related work

Last few months

- ▶ Rotskoff, Vanden-Eijnden [arXiv:1805.00915](https://arxiv.org/abs/1805.00915)
- ▶ Sirignano, Spiliopoulos [arXiv:1805.01053](https://arxiv.org/abs/1805.01053)
- ▶ Chizac, Bach [arXiv:1805.09545](https://arxiv.org/abs/1805.09545)

Does distributional dynamics converge?

Gradient flow minimizing $R(\rho)$,

$$\partial_t \rho_t(\theta) = 2\xi(t) \nabla_{\theta} \cdot (\rho_t(\theta) \nabla_{\theta} \Psi(\theta; \rho_t)). \quad (\text{DD})$$

- ▶ Does distributional dynamics converge to minimizers?
- ▶ In general, no; sometimes, yes.

In the following

- ▶ Concrete examples with convergence.
- ▶ A general convergence result for noisy SGD.

Does distributional dynamics converge?

Gradient flow minimizing $R(\rho)$,

$$\partial_t \rho_t(\boldsymbol{\theta}) = 2\xi(t) \nabla_{\boldsymbol{\theta}} \cdot (\rho_t(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}; \rho_t)). \quad (\text{DD})$$

- ▶ Does distributional dynamics converge to minimizers?
- ▶ In general, no; sometimes, yes.

In the following

- ▶ Concrete examples with convergence.
- ▶ A general convergence result for noisy SGD.

Does distributional dynamics converge?

Gradient flow minimizing $R(\rho)$,

$$\partial_t \rho_t(\boldsymbol{\theta}) = 2\xi(t) \nabla_{\boldsymbol{\theta}} \cdot (\rho_t(\boldsymbol{\theta}) \nabla_{\boldsymbol{\theta}} \Psi(\boldsymbol{\theta}; \rho_t)). \quad (\text{DD})$$

- ▶ Does distributional dynamics converge to minimizers?
- ▶ In general, no; sometimes, yes.

In the following

- ▶ Concrete examples with convergence.
- ▶ A general convergence result for noisy SGD.

Concrete examples

Simplest example requiring more than one neuron

With probability 1/2: $y = +1$, $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma_+)$,

With probability 1/2: $y = -1$, $\mathbf{x} \sim \mathcal{N}(\mathbf{0}, \Sigma_-)$.

$$\Sigma_{\pm} = \begin{bmatrix} \tau_{\pm}^2 \mathbf{I}_{s_0} & \mathbf{0} \\ \mathbf{0} & \mathbf{I}_{d-s_0} \end{bmatrix}.$$

Invariant under $\mathcal{O}(s_0) \times \mathcal{O}(d - s_0) \Rightarrow$ Reduced PDE.

Activation functions

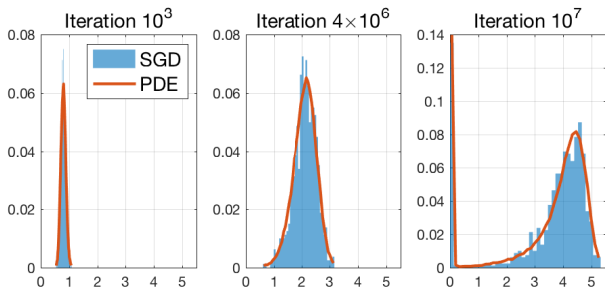
Simple $\theta_i = \mathbf{w}_i \in \mathbb{R}^d$ (no offset, no scaling weights)

$$\sigma_{\star}(\mathbf{x}; \theta_i) = \sigma(\langle \mathbf{x}, \mathbf{w}_i \rangle).$$

ReLU $\theta_i = (a_i, b_i, \mathbf{w}_i) \in \mathbb{R}^{d+2}$

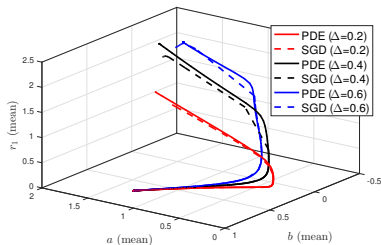
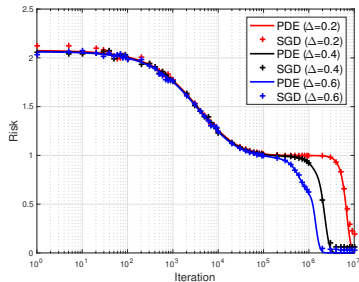
$$\sigma_{\star}(\mathbf{x}; \theta_i) = a_i \max(\langle \mathbf{x}, \mathbf{w}_i \rangle + b_i, 0).$$

Distributional dynamics



- ▶ $s_0 = d = 40$, $N = 800$, $\tau_+^2 = 1.8$, $\tau_-^2 = 0.2$.
- ▶ Simple activation.
- ▶ Histogram: empirical results. Cont. lines: PDE solutions.

Evolution of the risk



- ▶ $d = 320$, $s_0 = 60$, $N = 800$, $\tau_{\pm}^2 = 1 \pm \Delta$.
- ▶ ReLU activation.

Classifying anisotropic Gaussians: analysis

Assumption

(i) $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ truncated ReLU; (ii) $s_0 = \gamma d$, $\gamma \in (0, 1)$ fixed; (iii) $\bar{\rho}_0 \in \mathcal{P}(\mathbb{R}_+)$ has bounded density and $R(\rho_0) < 1$.

Theorem (M., Montanari, Nguyen, 2018)

For $T \geq T_0$, $d \geq d_0$, $N \geq C_0 d \log d$ (T_0, d_0, C_0 depend on $(\eta, \bar{\rho}_0, \Delta)$), consider SGD initialized with $(\theta_i^0)_{i \leq N} \sim_{i.i.d} \bar{\rho}_0 \times \text{Unif}(\mathbb{S}^{d-1})$ and step size $\varepsilon \leq 1/(C_0 d)$. Then, for any $k \in [T/\varepsilon, 10T/\varepsilon]$, whp

$$R_N(\theta^k) \leq \inf_{\theta \in \mathbb{R}^{d \times N}} R_N(\theta) + \eta.$$

- ▶ Learning from $k = O(1/\varepsilon) = O(d)$ samples.
- ▶ Independent of number of neurons $N \geq O(d \log d)$.

Classifying anisotropic Gaussians: analysis

Assumption

(i) $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ truncated ReLU; (ii) $s_0 = \gamma d$, $\gamma \in (0, 1)$ fixed; (iii) $\bar{\rho}_0 \in \mathcal{P}(\mathbb{R}_+)$ has bounded density and $R(\rho_0) < 1$.

Theorem (M., Montanari, Nguyen, 2018)

For $T \geq T_0$, $d \geq d_0$, $N \geq C_0 d \log d$ (T_0, d_0, C_0 depend on $(\eta, \bar{\rho}_0, \Delta)$), consider SGD initialized with $(\theta_i^0)_{i \leq N} \sim_{iid} \bar{\rho}_0 \times \text{Unif}(\mathbb{S}^{d-1})$ and step size $\varepsilon \leq 1/(C_0 d)$. Then, for any $k \in [T/\varepsilon, 10T/\varepsilon]$, whp

$$R_N(\theta^k) \leq \inf_{\theta \in \mathbb{R}^{d \times N}} R_N(\theta) + \eta.$$

- ▶ Learning from $k = O(1/\varepsilon) = O(d)$ samples.
- ▶ Independent of number of neurons $N \geq O(d \log d)$.

Classifying anisotropic Gaussians: analysis

Assumption

(i) $\sigma : \mathbb{R} \rightarrow \mathbb{R}$ truncated ReLU; (ii) $s_0 = \gamma d$, $\gamma \in (0, 1)$ fixed; (iii) $\bar{\rho}_0 \in \mathcal{P}(\mathbb{R}_+)$ has bounded density and $R(\rho_0) < 1$.

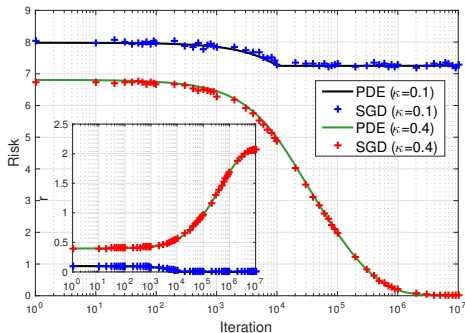
Theorem (M., Montanari, Nguyen, 2018)

For $T \geq T_0$, $d \geq d_0$, $N \geq C_0 d \log d$ (T_0, d_0, C_0 depend on $(\eta, \bar{\rho}_0, \Delta)$), consider SGD initialized with $(\theta_i^0)_{i \leq N} \sim_{iid} \bar{\rho}_0 \times \text{Unif}(\mathbb{S}^{d-1})$ and step size $\varepsilon \leq 1/(C_0 d)$. Then, for any $k \in [T/\varepsilon, 10T/\varepsilon]$, whp

$$R_N(\theta^k) \leq \inf_{\theta \in \mathbb{R}^{d \times N}} R_N(\theta) + \eta.$$

- ▶ Learning from $k = O(1/\varepsilon) = O(d)$ samples.
- ▶ Independent of number of neurons $N \geq O(d \log d)$.

Predicting failure



- ▶ $s_0 = d = 320$, $N = 800$, $\tau_+^2 = 1.5$, $\tau_-^2 = 0.5$.
- ▶ Non-monotone activation.
- ▶ Two different initialization ($\kappa =$ initialization variance).

Predicting failure

- ▶ SGD does not necessarily converge to global min.
- ▶ Can we fix it?

Noisy stochastic gradient descent

Regularized noisy SGD

SGD

$$\theta_i^{k+1} = \theta_i^k - 2s_k \nabla_{\theta_i} \ell(x_k, y_k; \theta^k)$$

Distributional dynamics

$$\partial_t \rho_t(\theta) = 2\xi(t) \nabla_{\theta} \cdot (\rho_t(\theta) \nabla_{\theta} \Psi(\theta; \rho_t))$$

Regularized noisy SGD

SGD with $(g_i^k)_{i \leq N, k \geq 0} \sim_{iid} \mathcal{N}(\mathbf{0}, \mathbf{I})$,

$$\theta_i^{k+1} = (1 - 2\lambda s_k) \theta_i^k - 2s_k \nabla_{\theta_i} \ell(x_k, y_k; \theta^k) + \sqrt{s_k / \beta} g_i^k.$$

Distributional dynamics **with diffusion term**

$$\partial_t \rho_t(\theta) = 2\xi(t) \nabla_{\theta} \cdot (\rho_t(\theta) \nabla_{\theta} \Psi_{\lambda}(\theta; \rho_t)) + \beta^{-1} \Delta_{\theta} \rho_t(\theta).$$

Regularized noisy SGD

SGD with $(g_i^k)_{i \leq N, k \geq 0} \sim_{iid} \mathcal{N}(\mathbf{0}, \mathbf{I})$,

$$\theta_i^{k+1} = (1 - 2\lambda s_k) \theta_i^k - 2s_k \nabla_{\theta_i} \ell(x_k, y_k; \theta^k) + \sqrt{s_k / \beta} g_i^k.$$

Distributional dynamics **with diffusion term**

$$\partial_t \rho_t(\theta) = 2\xi(t) \nabla_{\theta} \cdot (\rho_t(\theta) \nabla_{\theta} \Psi_{\lambda}(\theta; \rho_t)) + \beta^{-1} \Delta_{\theta} \rho_t(\theta).$$

Theorem

Same approximation theorem: noisy SGD \leftrightarrow PDE.

Gradient flow interpretation

$$F_{\beta,\lambda}(\rho) = \frac{1}{2}R(\rho) + \frac{\lambda}{2} \int \|\boldsymbol{\theta}\|_2^2 \rho(d\boldsymbol{\theta}) - \beta^{-1} \text{Ent}(\rho),$$
$$\text{Ent}(\rho) = - \int \rho(\boldsymbol{\theta}) \log \rho(\boldsymbol{\theta}) d\boldsymbol{\theta}.$$

- ▶ Distributional dynamics is the gradient flow of $F_{\beta,\lambda}(\rho)$...
- ▶ ... in Wasserstein metric space.

[Jordan, Kinderlehrer, Otto, 1998]

Convergence of DD

Theorem (M., Montanari, Nguyen, 2018)

Assume V, U, ρ_0 “sufficiently” regular. If ρ_t is a solution of DD, then $F_{\beta, \lambda}(\rho_t)$ is non-increasing:

$$\partial_t F_{\beta, \lambda}(\rho_t) = - \int \left\| \nabla \left(\Psi_{\lambda}(\boldsymbol{\theta}; \rho_t) - \frac{1}{\beta} \log \rho_t(\boldsymbol{\theta}) \right) \right\|_2^2 \rho_t(d\boldsymbol{\theta}) \leq 0.$$

In particular, there exists a unique fixed point ρ_{\star} of $F_{\beta, \lambda}$ satisfies

$$\rho_{\star}(\boldsymbol{\theta}) = \frac{1}{Z_{\star}(\beta, \lambda)} \exp\{-\beta \Psi_{\lambda}(\boldsymbol{\theta}; \rho_{\star})\}.$$

Moreover, as $t \rightarrow \infty$, $\rho_t \rightarrow \rho_{\star}$.

Generalized the analysis of [Carrillo, McCann, Villani, 2013].

Key remark

▶

$$\rho_*(\theta) = \frac{1}{Z_*(\beta, \lambda)} \exp\{-\beta\Psi(\theta; \rho_*)\}.$$

is the stationery equation for

$$F_{\beta, \lambda}(\rho) = \frac{1}{2}R(\rho) + \frac{\lambda}{2} \int \|\theta\|_2^2 \rho(d\theta) - \beta^{-1}\text{Ent}(\rho).$$

- ▶ $F_{\beta, \lambda}(\cdot)$ is strongly convex.
- ▶ The fixed point is unique!

General convergence for noisy SGD

Theorem (M., Montanari, Nguyen, 2018)

Assumptions of previous theorem. Initialization $(\theta_i^0)_{i \leq N} \sim_{iid} \rho_0$. Then there exists $\beta_0 = \beta_0(D, U, V, \eta)$, such that, for $\beta \geq \beta_0$, there exists $T = T(D, U, V, \beta, \eta)$ such that for any $k \in [T/\varepsilon, 10T/\varepsilon]$, $N \geq C_0 D \log D$, $\varepsilon \leq 1/(C_0 D)$, we have, whp

$$R_{\lambda, N}(\theta^k) \leq \inf_{\theta \in \mathbb{R}^{D \times N}} R_{\lambda, N}(\theta) + \eta.$$

General convergence for noisy SGD

Theorem (M., Montanari, Nguyen, 2018)

Assumptions of previous theorem. Initialization $(\theta_i^0)_{i \leq N} \sim_{iid} \rho_0$. Then there exists $\beta_0 = \beta_0(D, U, V, \eta)$, such that, for $\beta \geq \beta_0$, there exists $T = T(D, U, V, \beta, \eta)$ such that for any $k \in [T/\varepsilon, 10T/\varepsilon]$, $N \geq C_0 D \log D$, $\varepsilon \leq 1/(C_0 D)$, we have, whp

$$R_{\lambda, N}(\theta^k) \leq \inf_{\theta \in \mathbb{R}^{D \times N}} R_{\lambda, N}(\theta) + \eta.$$

- ▶ For general distribution $(x, y) \sim \mathbb{P}_{x, y}$!
- ▶ Convergence time depends on D , but not on N !

Conclusion

Conclusion

Correspondence

- ▶ Two layer neural networks.
- ▶ Dynamics of particles with pairwise interactions.
- ▶ Gradient flow in measure spaces.

Conclusion

Correspondence

- ▶ Two layer neural networks.
- ▶ Dynamics of particles with pairwise interactions.
- ▶ Gradient flow in measure spaces.

Partially explained the optimization/generalization surprise.