

Maximum Mean Discrepancy Gradient Flow

Michael Arbel ¹ Anna Korba ¹ Adil Salim ² Arthur Gretton ¹

¹Gatsby Computational Neuroscience Unit, UCL, London

²Visual Computing Center, KAUST, Saudi Arabia

October 15, 2019

Outline

- ▶ General problem
- ▶ Wasserstein gradient of the MMD
- ▶ A Criterion for global convergence
- ▶ A noise-injection algorithm for better convergence

General problem

Finite dimensional **non-convex** optimization:

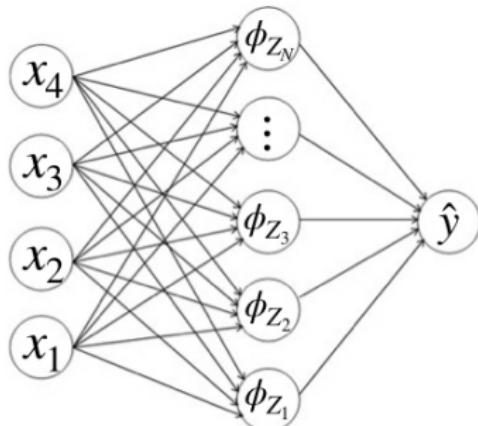
$$\min_{Z^1, \dots, Z^N \in \mathcal{Z}} \mathcal{L} \left(\frac{1}{N} \sum_{i=1}^N \delta_{Z^n} \right)$$

General problem

Finite dimensional **non-convex** optimization:

$$\min_{Z^1, \dots, Z^N \in \mathcal{Z}} \mathcal{L} \left(\frac{1}{N} \sum_{i=1}^N \delta_{Z^i} \right)$$

$(x, y) \sim data$



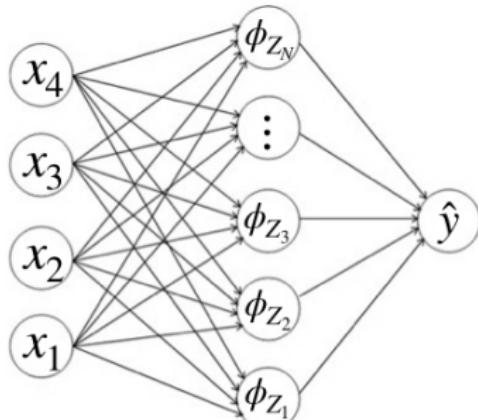
$$\min_{Z_1, \dots, Z_N} \mathbb{E}_{data} \left[\|y - \frac{1}{N} \sum_{i=1}^N \phi_{Z_i}(x)\|^2 \right]$$

General problem

Finite dimensional **non-convex** optimization:

$$\min_{Z^1, \dots, Z^N \in \mathcal{Z}} \mathcal{L} \left(\frac{1}{N} \sum_{i=1}^N \delta_{Z_i^n} \right)$$

$(x, y) \sim data$



► Optimization using gradient descent GD:

$$Z_{t+1}^n = Z_t^n - \gamma \nabla_{Z^n} \mathcal{L} \left(\frac{1}{N} \sum_{i=1}^N \delta_{Z_t^n} \right)$$

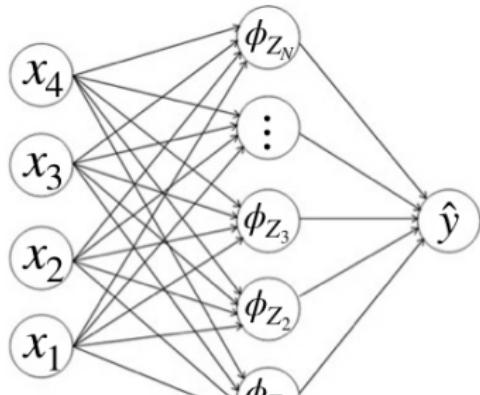
$$\min_{Z_1, \dots, Z_N} \mathbb{E}_{data} \left[\|y - \frac{1}{N} \sum_{i=1}^N \phi_{Z_i}(x)\|^2 \right]$$

General problem

Finite dimensional **non-convex** optimization:

$$\min_{Z^1, \dots, Z^N \in \mathcal{Z}} \mathcal{L} \left(\frac{1}{N} \sum_{i=1}^N \delta_{Z^i} \right)$$

$(x, y) \sim data$



$$\min_{Z_1, \dots, Z_N} \mathbb{E}_{data} \left[\|y - \frac{1}{N} \sum_{i=1}^N \phi_{Z_i}(x)\|^2 \right]$$

- ▶ Optimization using gradient descent GD:

$$Z_{t+1}^n = Z_t^n - \gamma \nabla_{Z^n} \mathcal{L} \left(\frac{1}{N} \sum_{i=1}^N \delta_{Z_t^i} \right)$$

- ▶ Hard to describe the dynamics of GD!

General problem

Infinite dimensional **non-convex** optimization:

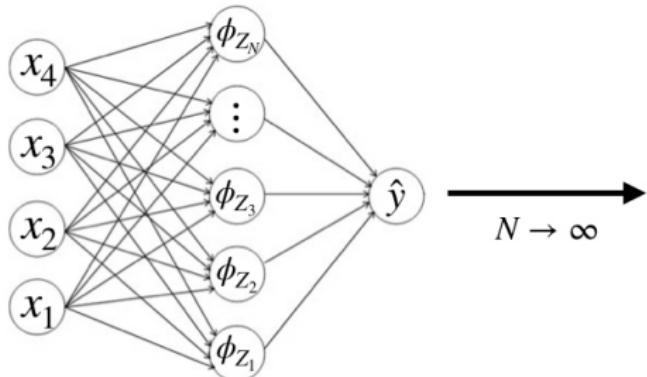
$$\min_{Z^1, \dots, Z^N \in \mathcal{Z}} \mathcal{L} \left(\frac{1}{N} \sum_{i=1}^N \delta_{Z^n} \right) \quad \xrightarrow{N \rightarrow \infty} \quad \min_{\nu \in \mathcal{P}} \mathcal{L}(\nu)$$

General problem

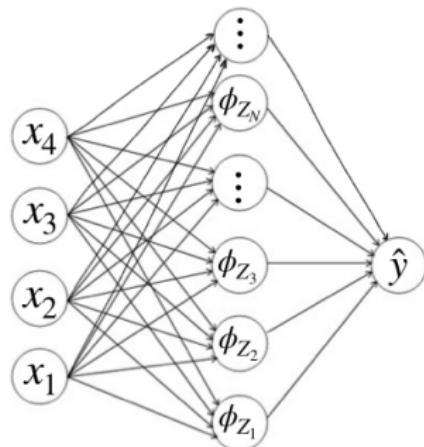
Infinite dimensional **non-convex** optimization:

$$\min_{Z^1, \dots, Z^N \in \mathcal{Z}} \mathcal{L} \left(\frac{1}{N} \sum_{i=1}^N \delta_{Z^i} \right) \xrightarrow{N \rightarrow \infty} \min_{\nu \in \mathcal{P}} \mathcal{L}(\nu)$$

$(x, y) \sim data$



$$\min_{Z_1, \dots, Z_N} \mathbb{E}_{data} [\|y - \frac{1}{N} \sum_{i=1}^N \phi_{Z_i}(x)\|^2] \xrightarrow{N \rightarrow \infty} \min_{\nu \in \mathcal{P}} \mathbb{E}_{data} [\|y - \mathbb{E}_{Z \sim \nu} [\phi_Z(x)]\|^2]$$

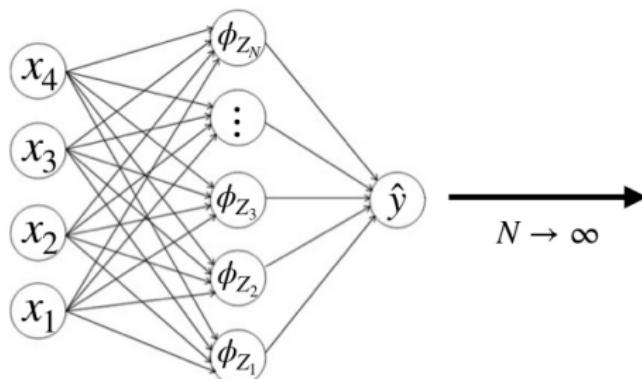


General problem

- Global convergence of Gradient descent¹ when $N \rightarrow \infty$ and $\phi_Z(x)$ of the form:

$$\phi_Z(x) = w g_\theta(x), \quad Z = (w, \theta)$$

$(x, y) \sim data$



$N \rightarrow \infty$

$$\min_{Z_1, \dots, Z_N} \mathbb{E}_{data} \left[\|y - \frac{1}{N} \sum_{i=1}^N \phi_{Z_i}(x)\|^2 \right] \xrightarrow{N \rightarrow \infty} \min_{\nu \in \mathcal{P}} \mathbb{E}_{data} [\|y - \mathbb{E}_{Z \sim \nu} [\phi_Z(x)]\|^2]$$

¹[Rotskoff and Vanden-Eijnden, 2018, Chizat and Bach, 2018]

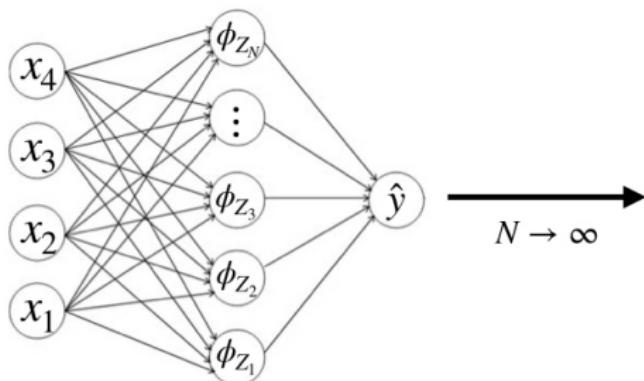
General problem

- ▶ Global convergence of Gradient descent¹ when $N \rightarrow \infty$ and $\phi_Z(x)$ of the form:

$$\phi_Z(x) = w g_\theta(x), \quad Z = (w, \theta)$$

- ▶ Interested in more general form for $\phi_Z(x)$.

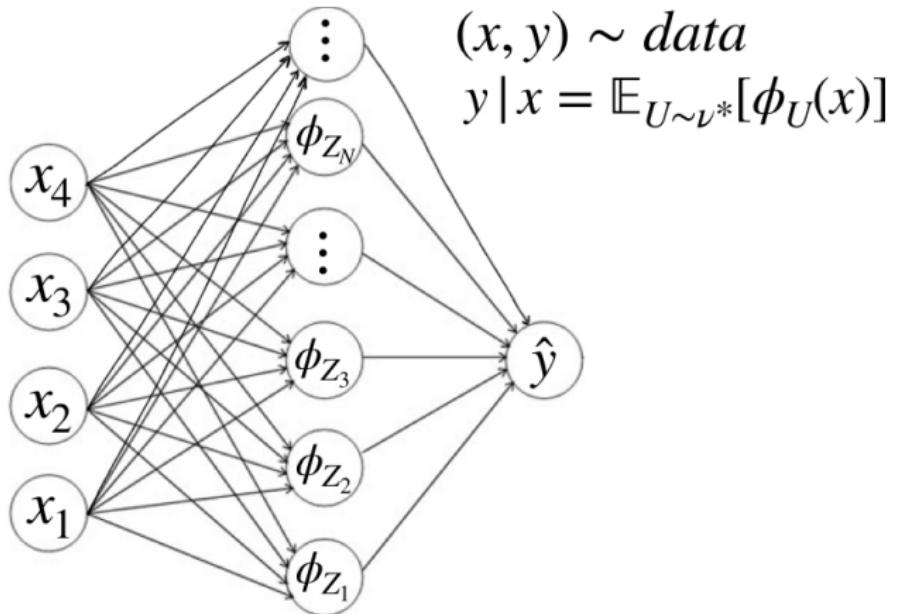
$(x, y) \sim data$



$$\min_{Z_1, \dots, Z_N} \mathbb{E}_{data} \left[\|y - \frac{1}{N} \sum_{i=1}^N \phi_{Z_i}(x)\|^2 \right] \xrightarrow{N \rightarrow \infty} \min_{\nu \in \mathcal{P}} \mathbb{E}_{data} [\|y - \mathbb{E}_{Z \sim \nu} [\phi_Z(x)]\|^2]$$

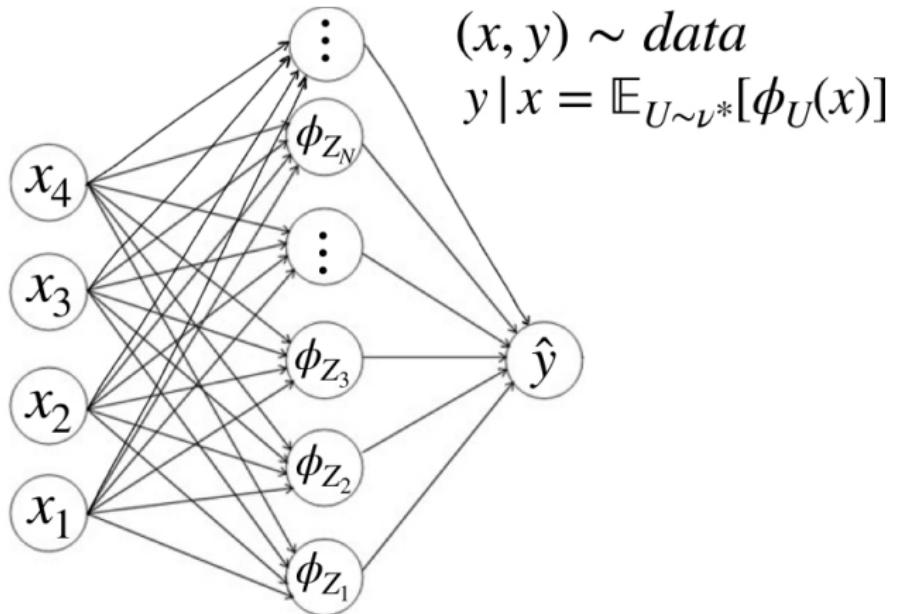
¹[Rotskoff and Vanden-Eijnden, 2018, Chizat and Bach, 2018]

General problem



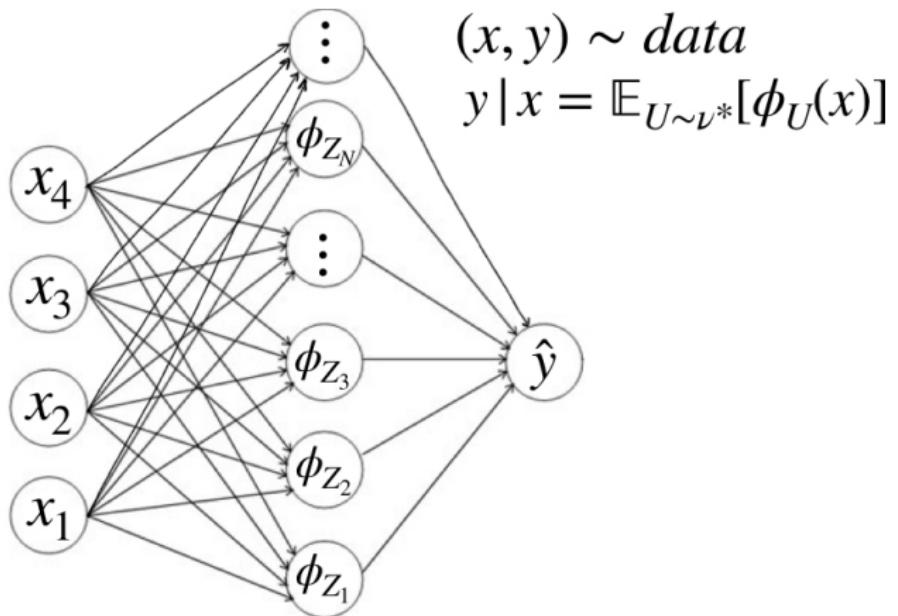
$$\min_{\nu \in \mathcal{P}} \mathbb{E}_{data} [\|y - \mathbb{E}_{Z \sim \nu} [\phi_Z(x)]\|^2]$$

General problem



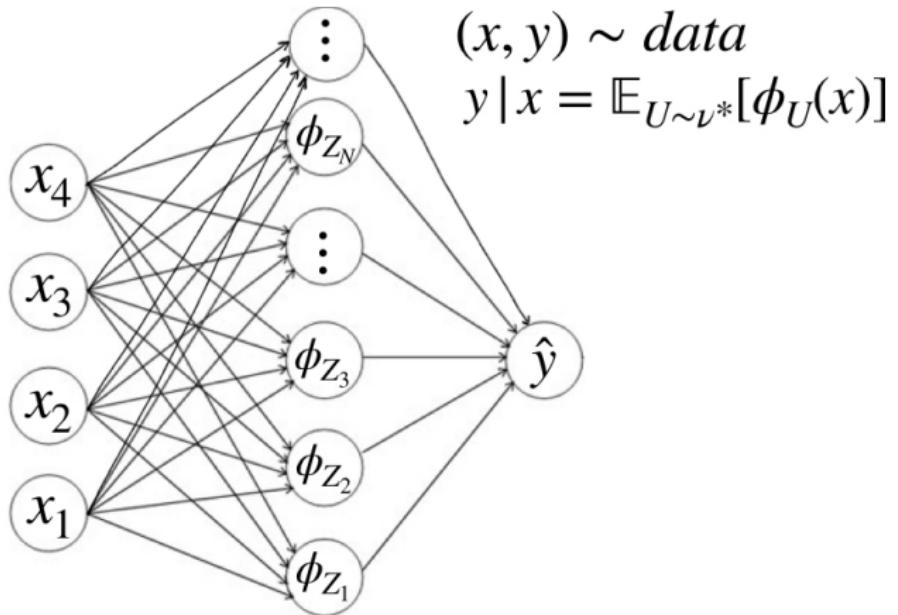
$$\min_{\nu \in \mathcal{P}} \mathbb{E}_{data} [\|\mathbb{E}_{U \sim \nu^*}[\phi_U(x)] - \mathbb{E}_{Z \sim \nu}[\phi_Z(x)]\|^2]$$

General problem



$$\min_{\nu \in \mathcal{P}} MMD^2(\nu^*, \nu)$$

General problem



$$\min_{\nu \in \mathcal{P}} MMD^2(\nu^*, \nu)$$

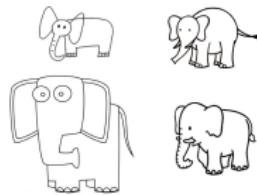
$$k(Z, Z') = \mathbb{E}_{data}[\phi_Z(x)\phi_{Z'}(x)]$$

The Maximum Mean Discrepancy²

Consider samples from two distributions ν^* and ν_0 .



$$U^m \sim \nu^*$$

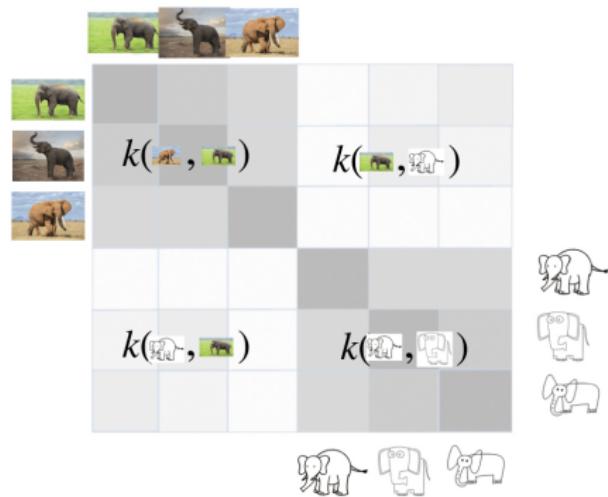


$$Z^n \sim \nu_0$$

²[Gretton et al., 2012]

The Maximum Mean Discrepancy²

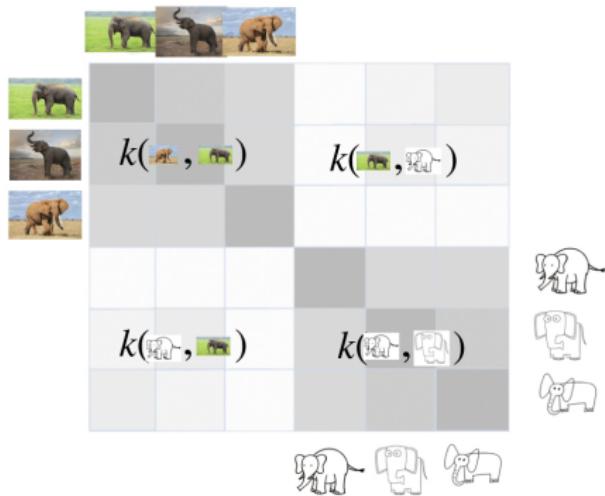
Compute a similarity matrix using a kernel k



²[Gretton et al., 2012]

The Maximum Mean Discrepancy²

Compute a similarity matrix using a kernel k

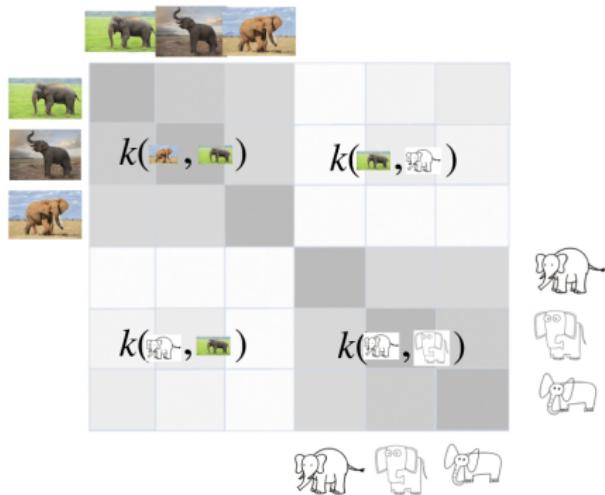


$$\widehat{MMD}^2 = \frac{1}{n(n-1)} \sum_{n,n'} k(\text{brown}, \text{green}) + \frac{1}{n(n-1)} \sum_{n,n'} k(\text{brown}, \text{tan}) - \frac{2}{n^2} \sum_{n,n'} k(\text{brown}; \text{green})$$

²[Gretton et al., 2012]

The Maximum Mean Discrepancy²

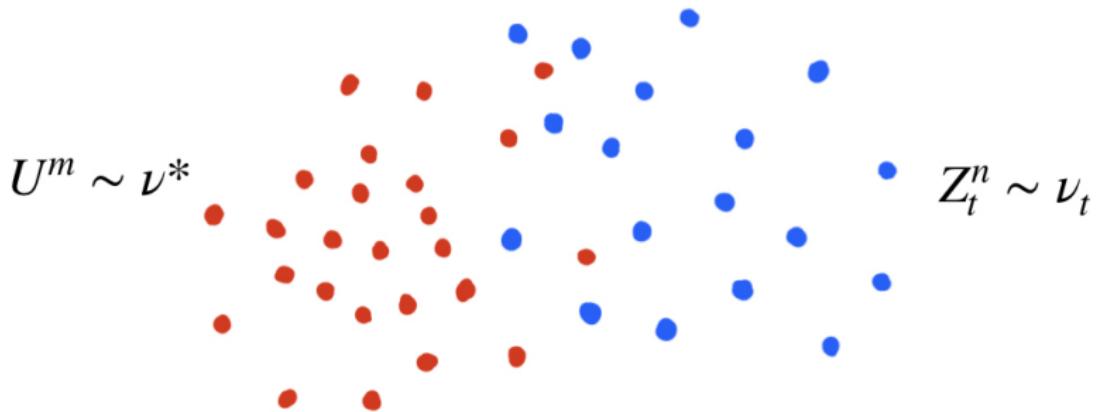
Compute a similarity matrix using a kernel k



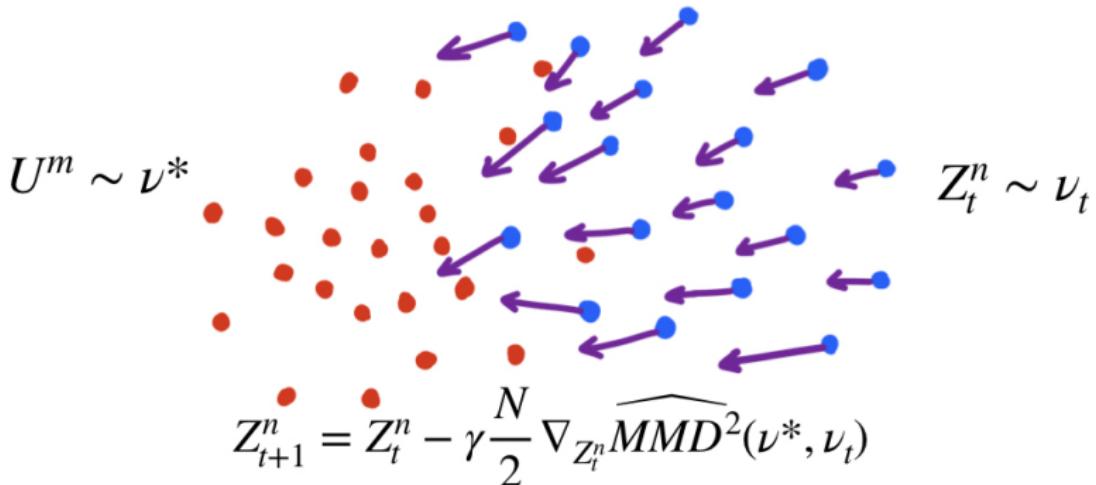
$$MMD^2(\nu^*, \nu_0) = \mathbb{E}_{\substack{U \sim \nu^* \\ U' \sim \nu^*}} [k(U, U')] + \mathbb{E}_{\substack{Z \sim \nu_0 \\ Z' \sim \nu_0}} [k(Z, Z')] - 2 \mathbb{E}_{\substack{U \sim \nu^* \\ Z' \sim \nu_0}} [k(U, Z)]$$

²[Gretton et al., 2012]

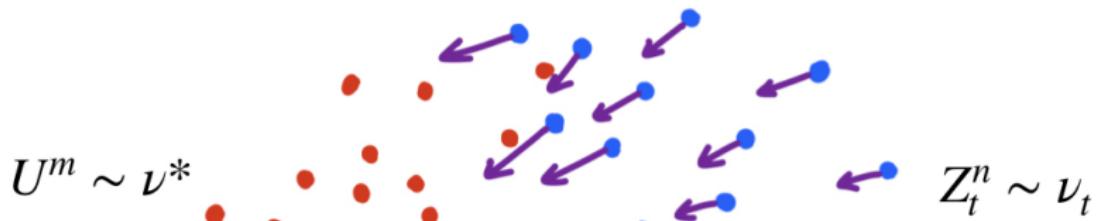
Gradient descent



Gradient descent



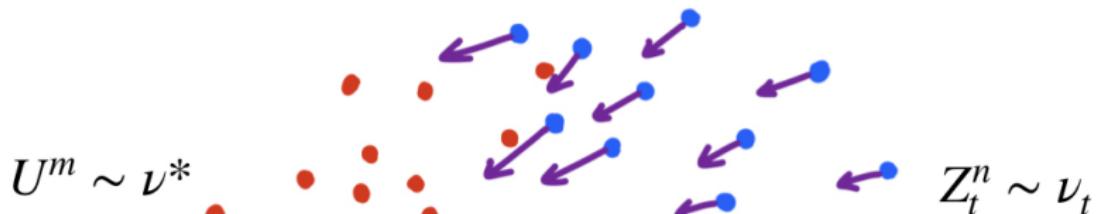
Gradient descent



$$Z_{t+1}^n = Z_t^n - \gamma \frac{N}{2} \nabla_{Z_t^n} \widehat{MMD}^2(\nu^*, \nu_t)$$

$$\frac{N}{2} \nabla_{Z_t^n} \widehat{MMD}^2(\nu^*, \nu_t) = \frac{1}{N-1} \sum_{n' \neq n} \partial_1 k(Z_t^n, Z_t^{n'}) - \frac{1}{M} \sum_{m=1}^M \partial_1 k(Z_t^n, U^m)$$

Gradient descent



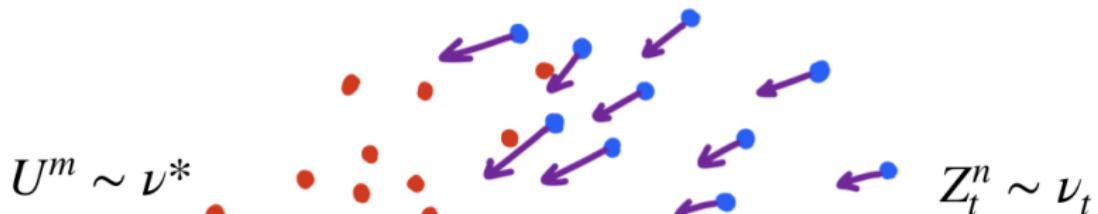
$$Z_{t+1}^n = Z_t^n - \gamma \frac{N}{2} \nabla_{Z_t^n} \widehat{\text{MMD}}^2(\nu^*, \nu_t)$$

$$\frac{N}{2} \nabla_{Z_t^n} \widehat{\text{MMD}}^2(\nu^*, \nu_t) = \frac{1}{N-1} \sum_{n' \neq n} \partial_1 k(Z_t^n, Z_t^{n'}) - \frac{1}{M} \sum_{m=1}^M \partial_1 k(Z_t^n, U^m)$$

$$\nabla_{Z_t} \left(\mathbb{E}_{Z' \sim \nu_t} [k(Z_t, Z')] - \mathbb{E}_{U \sim \nu^*} [k(Z_t, U)] \right)$$

$\downarrow N, M \rightarrow \infty$

Gradient descent



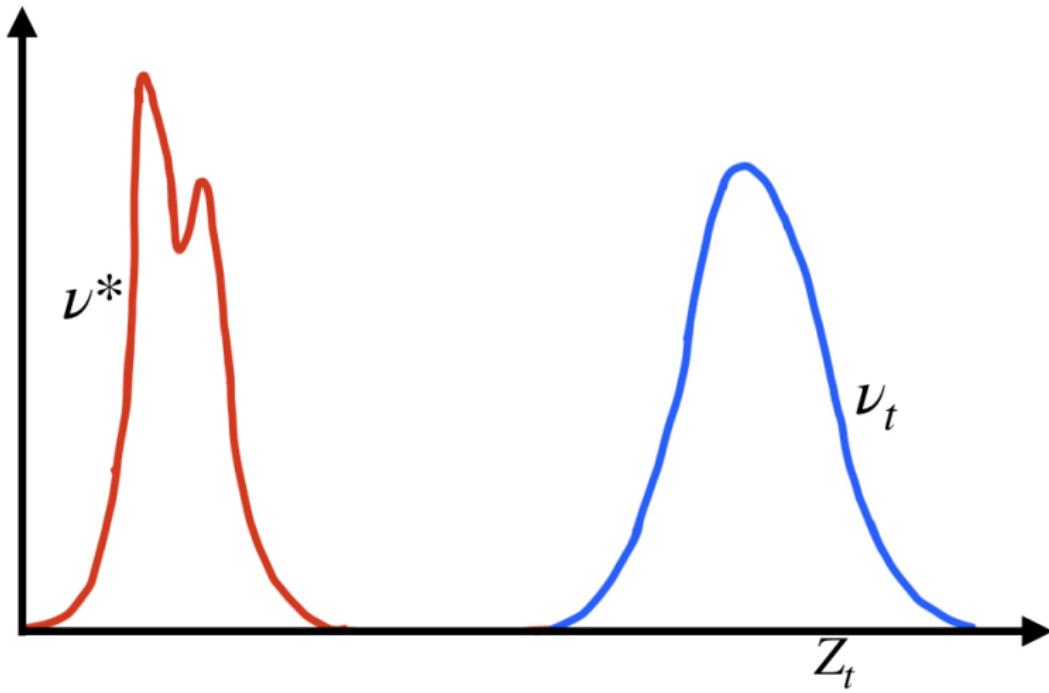
$$Z_{t+1}^n = Z_t^n - \gamma \frac{N}{2} \nabla_{Z_t^n} \widehat{\text{MMD}}^2(\nu^*, \nu_t)$$

$$\frac{N}{2} \nabla_{Z_t^n} \widehat{\text{MMD}}^2(\nu^*, \nu_t) = \frac{1}{N-1} \sum_{n' \neq n} \partial_1 k(Z_t^n, Z_t^{n'}) - \frac{1}{M} \sum_{m=1}^M \partial_1 k(Z_t^n, U^m)$$

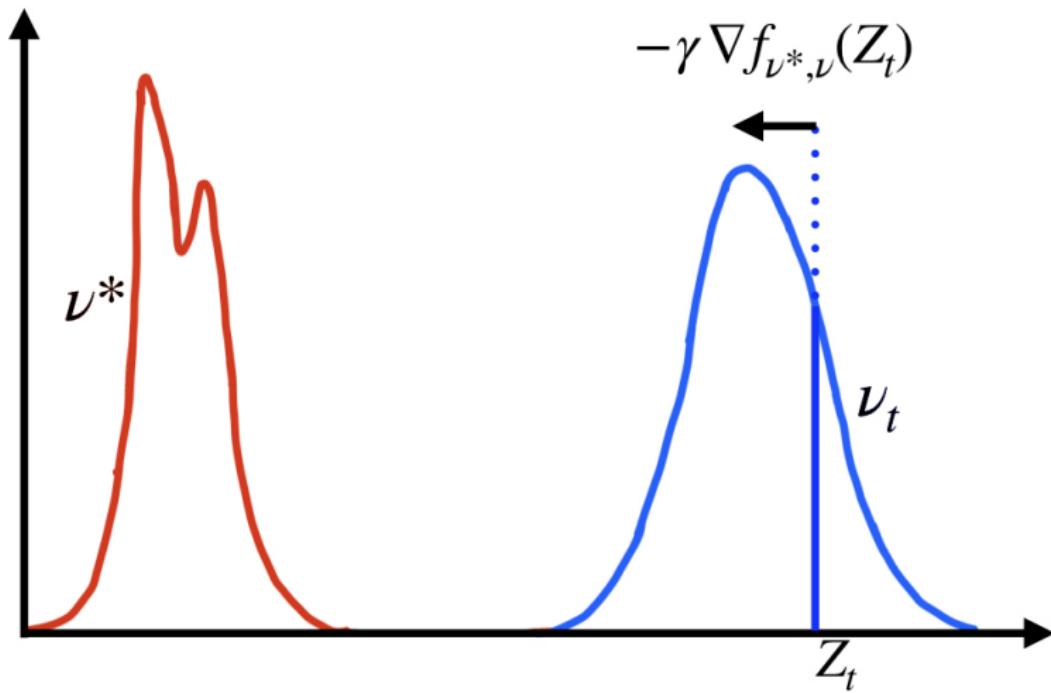
$\downarrow N, M \rightarrow \infty$

$$\underbrace{\nabla_{Z_t} \left(\mathbb{E}_{Z' \sim \nu_t} [k(Z_t, Z')] - \mathbb{E}_{U \sim \nu^*} [k(Z_t, U)] \right)}_{f_{\nu^*, \nu_t}(Z_t)}$$

Wasserstein gradient descent

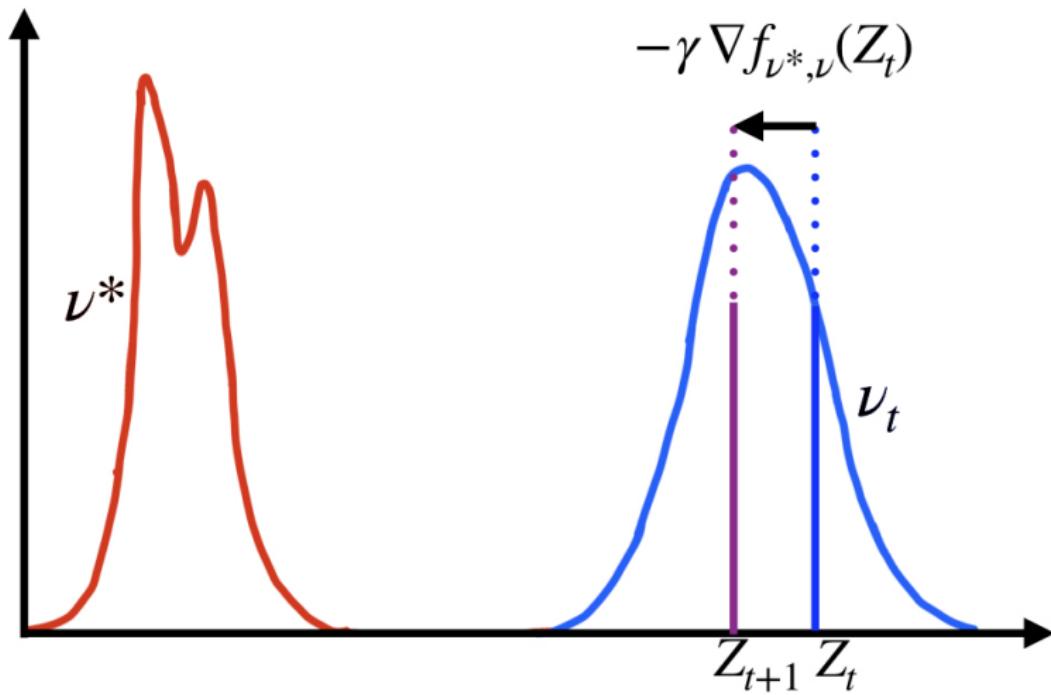


Wasserstein gradient descent



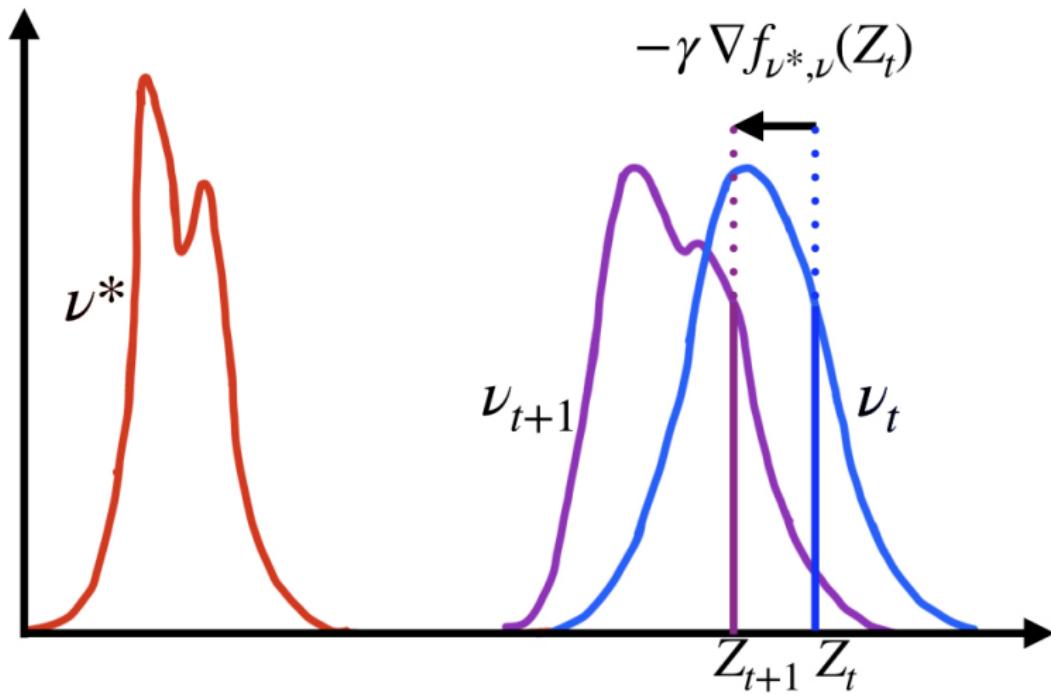
$$f_{\nu^*, \nu_t}(Z_t) = \mathbb{E}_{Z' \sim \nu_t}[k(Z_t, Z')] - \mathbb{E}_{U \sim \nu^*}[k(Z_t, U)]$$

Wasserstein gradient descent



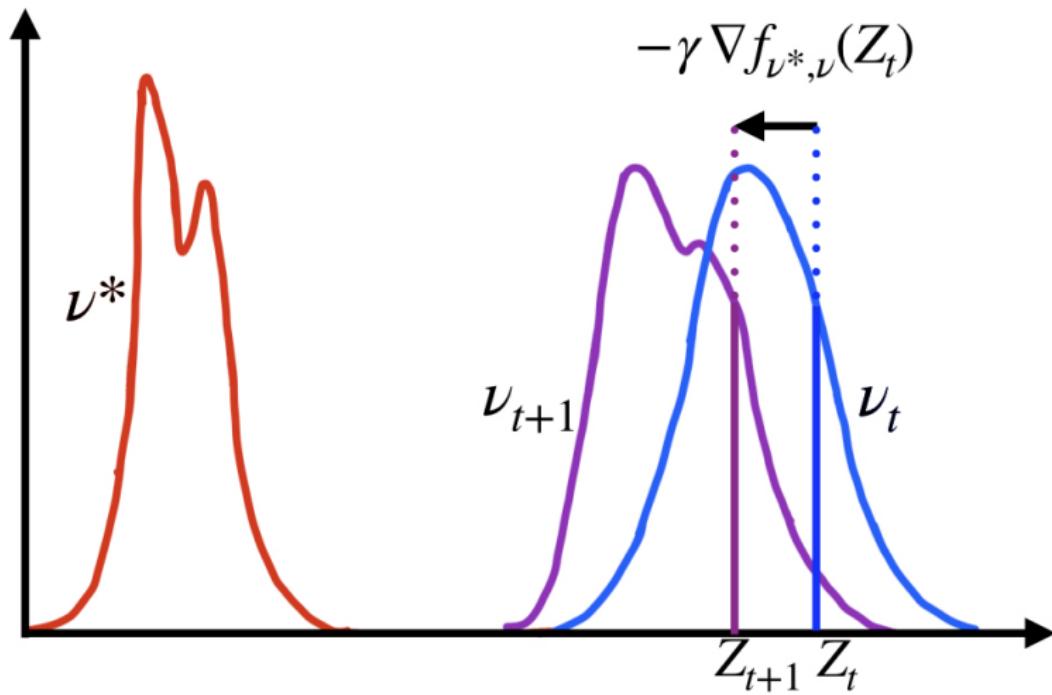
$$f_{\nu^*, \nu_t}(Z_t) = \mathbb{E}_{Z' \sim \nu_t}[k(Z_t, Z')] - \mathbb{E}_{U \sim \nu^*}[k(Z_t, U)]$$

Wasserstein gradient descent



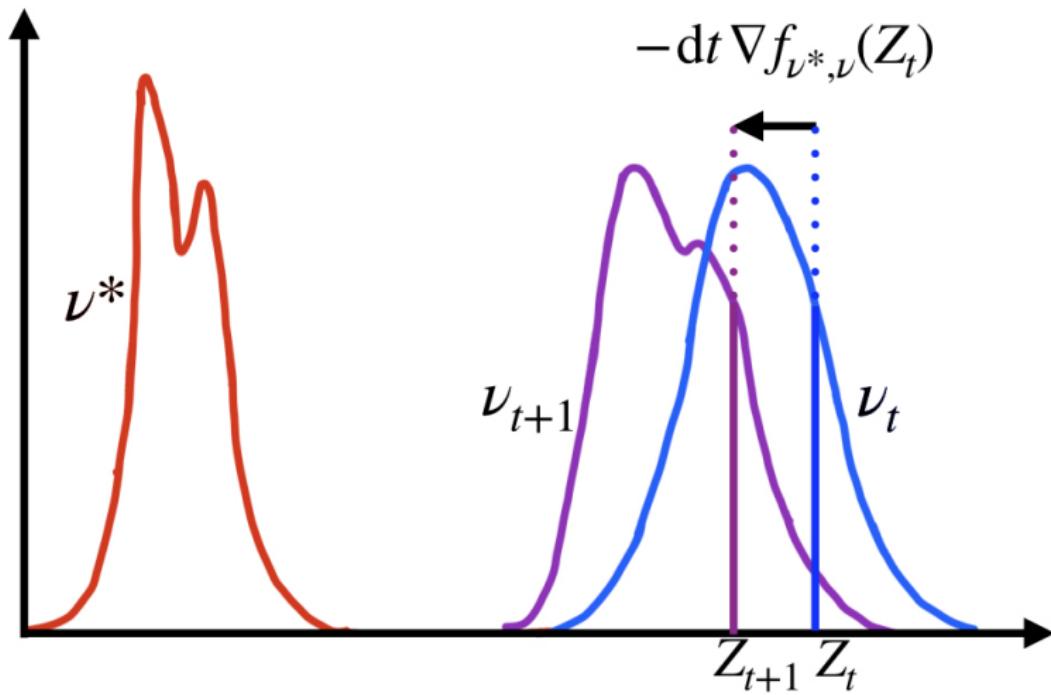
$$f_{\nu^*, \nu_t}(Z_t) = \mathbb{E}_{Z' \sim \nu_t}[k(Z_t, Z')] - \mathbb{E}_{U \sim \nu^*}[k(Z_t, U)]$$

Wasserstein gradient descent



$$Z_{t+1} = Z_t - \gamma \nabla_{Z_t} f_{\nu^*, \nu_t}(Z_t), \quad Z_t \sim \nu_t$$

Wasserstein gradient descent



$$dZ_t = - dt \nabla_{Z_t} f_{\nu^*, \nu_t}(Z_t), \quad Z_t \sim \nu_t$$

Wasserstein gradient flow

Continuity Equation

- Continuous time equation: Mc-Kean Vlasov dynamics

$$\frac{dZ_t}{dt} = -\nabla_{Z_t} f_{\nu^*, \nu_t}(Z_t), \quad Z_t \sim \nu_t$$

³[Otto, 2001, Villani, 2004, Ambrosio et al., 2004]

Continuity Equation

- ▶ Continuous time equation: Mc-Kean Vlasov dynamics

$$\frac{dZ_t}{dt} = -\nabla_{Z_t} f_{\nu^*, \nu_t}(Z_t), \quad Z_t \sim \nu_t$$

- ▶ Equivalent to a PDE in ν_t :

$$\partial_t \nu_t = \operatorname{div}(\nu_t \nabla f_{\nu^*, \nu_t})$$

³[Otto, 2001, Villani, 2004, Ambrosio et al., 2004]

Continuity Equation

- ▶ Continuous time equation: Mc-Kean Vlasov dynamics

$$\frac{dZ_t}{dt} = -\nabla_{Z_t} f_{\nu^*, \nu_t}(Z_t), \quad Z_t \sim \nu_t$$

- ▶ Equivalent to a PDE in ν_t :

$$\partial_t \nu_t = \text{div}(\nu_t \nabla f_{\nu^*, \nu_t})$$

- ▶ Interpretation as a gradient flow in probability space³:

$$\partial_t \nu_t = -\nabla_{\nu_t} \mathcal{L}(\nu_t) \quad \mathcal{L}(\nu) := \frac{1}{2} MMD^2(\nu^*, \nu)$$

³[Otto, 2001, Villani, 2004, Ambrosio et al., 2004]

A criterion for convergence

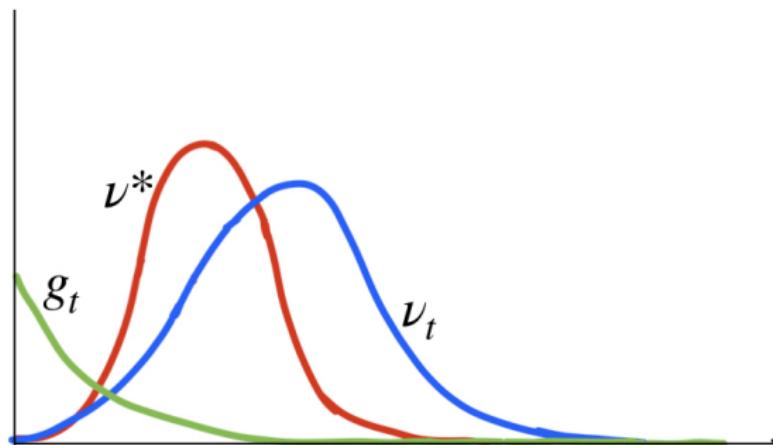
- ▶ Define the Negative Sobolev distance:

$$S(\nu^* | \nu_t) = \sup_{g, \mathbb{E}_{Z \sim \nu_t} [\|\nabla g(Z)\|^2] \leq 1} |\mathbb{E}_{Z \sim \nu_t} [g(Z)] - \mathbb{E}_{U \sim \nu^*} [g(U)]|$$

A criterion for convergence

- ▶ Define the Negative Sobolev distance:

$$S(\nu^* | \nu_t) = \sup_{g, \mathbb{E}_{Z \sim \nu_t} [\|\nabla g(Z)\|^2] \leq 1} |\mathbb{E}_{Z \sim \nu_t} [g(Z)] - \mathbb{E}_{U \sim \nu^*} [g(U)]|$$



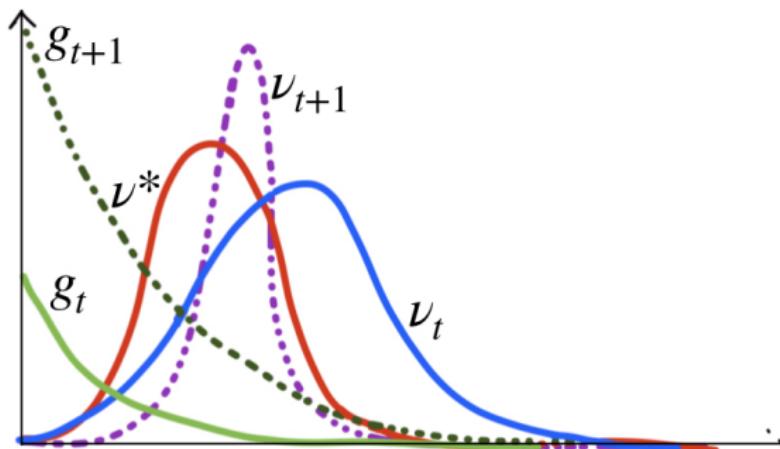
$$\int \|\nabla g_t(x)\|^2 d\nu_t(x) \leq 1$$

$$S(\nu^* | \nu_t) = |\mathbb{E}_{Z \sim \nu_t} [g_t(Z)] - \mathbb{E}_{U \sim \nu^*} [g_t(U)]|$$

A criterion for convergence

- ▶ Define the Negative Sobolev distance:

$$S(\nu^* | \nu_t) = \sup_{g, \mathbb{E}_{Z \sim \nu_t} [\|\nabla g(Z)\|^2] \leq 1} |\mathbb{E}_{Z \sim \nu_t} [g(Z)] - \mathbb{E}_{U \sim \nu^*} [g(U)]|$$

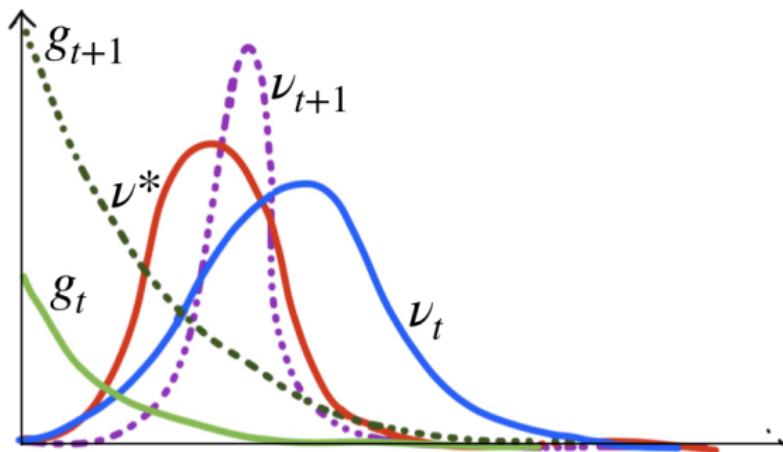


$$S(\nu^* | \nu_{t+1}) - S(\nu^* | \nu_t) > 0$$

A criterion for convergence

- ▶ Define the Negative Sobolev distance:

$$S(\nu^* | \nu_t) = \sup_{g, \mathbb{E}_{Z \sim \nu_t} [\|\nabla g(Z)\|^2] \leq 1} |\mathbb{E}_{Z \sim \nu_t} [g(Z)] - \mathbb{E}_{U \sim \nu^*} [g(U)]|$$



$$S(\nu^* | \nu_{t+1}) - S(\nu^* | \nu_t) > C$$

A criterion for convergence

- ▶ Assume that $S(\nu^* | \nu_t) \leq C$ for all t , then for γ small enough

$$MMD^2(\nu^*, \nu_t) \leq \frac{1}{MMD^2(\nu^*, \nu_0) + 8\gamma C^{-1}t}$$

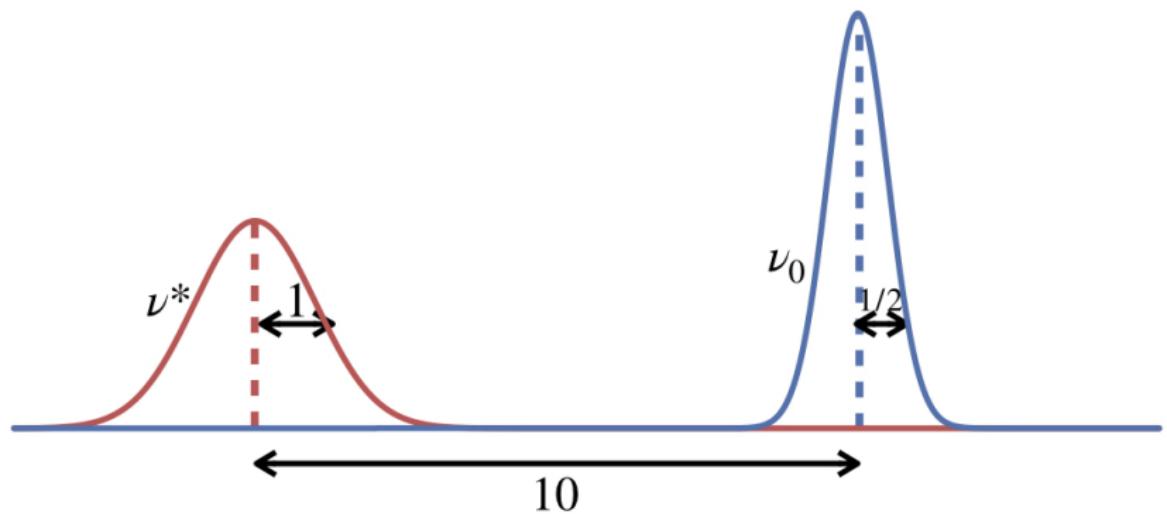
A criterion for convergence

- ▶ Assume that $S(\nu^* | \nu_t) \leq C$ for all t , then for γ small enough

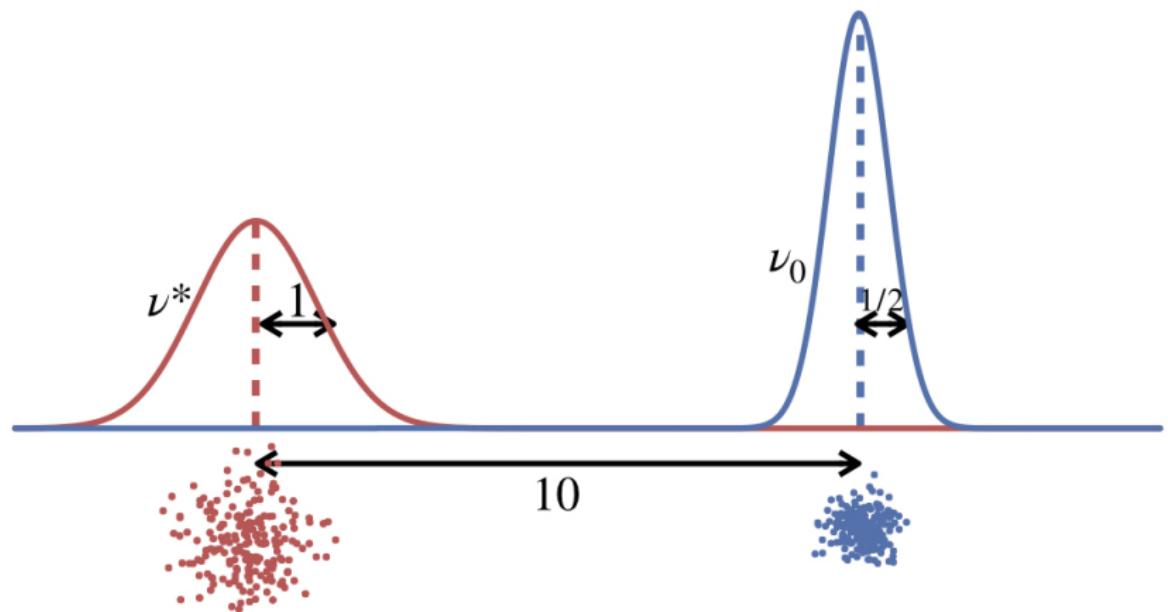
$$MMD^2(\nu^*, \nu_t) \leq \frac{1}{MMD^2(\nu^*, \nu_0) + 8\gamma C^{-1}t}$$

- ▶ Depends on the whole sequence ν_t : Hard to verify in general
- ▶ Can be checked for simple examples.

Convergence: Failure case

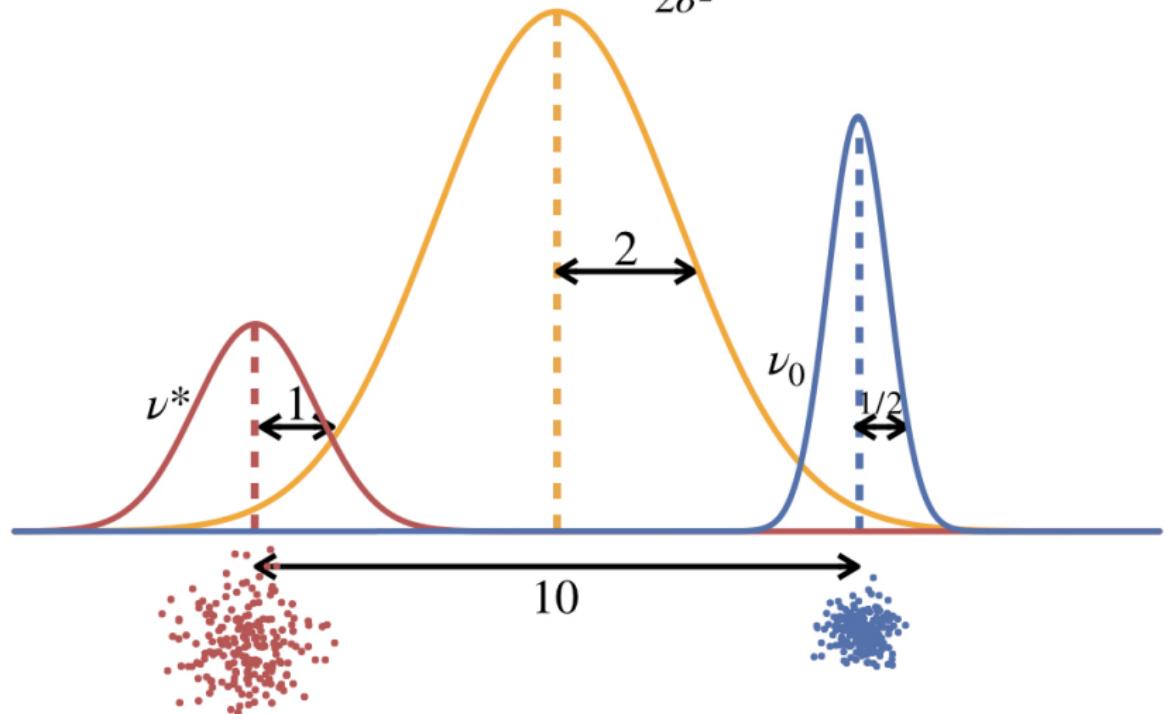


Convergence: Failure case



Convergence: Failure case

$$k(z, z') = \exp\left(-\frac{\|z - z'\|^2}{2\sigma^2}\right)$$



Convergence: Failure case

Noise Injection

- ▶ Idea: Evaluate $\nabla f_{\nu^*, \nu_t}$ outside of the support of ν_t to get a better signal!

⁴[Chaudhari et al., 2017, Hazan et al., 2015]

⁵[Mei et al., 2018]

Noise Injection

- ▶ Idea: Evaluate $\nabla f_{\nu^*, \nu_t}$ outside of the support of ν_t to get a better signal!
- ▶ Sample $u_t \sim \mathcal{N}(0, 1)$ and β_t is the noise level:

$$Z_{t+1} = Z_t - \gamma \nabla f_t(Z_t + \beta_t u_t); \quad Z_t \sim \nu_t$$

⁴[Chaudhari et al., 2017, Hazan et al., 2015]

⁵[Mei et al., 2018]

Noise Injection

- ▶ Idea: Evaluate $\nabla f_{\nu^*, \nu_t}$ outside of the support of ν_t to get a better signal!
- ▶ Sample $u_t \sim \mathcal{N}(0, 1)$ and β_t is the noise level:

$$Z_{t+1} = Z_t - \gamma \nabla f_t(Z_t + \beta_t u_t); \quad Z_t \sim \nu_t$$

- ▶ Similar to *continuation methods*⁴, but extended to interacting particles.

⁴[Chaudhari et al., 2017, Hazan et al., 2015]

⁵[Mei et al., 2018]

Noise Injection

- ▶ Idea: Evaluate $\nabla f_{\nu^*, \nu_t}$ outside of the support of ν_t to get a better signal!
- ▶ Sample $u_t \sim \mathcal{N}(0, 1)$ and β_t is the noise level:

$$Z_{t+1} = Z_t - \gamma \nabla f_t(Z_t + \beta_t u_t); \quad Z_t \sim \nu_t$$

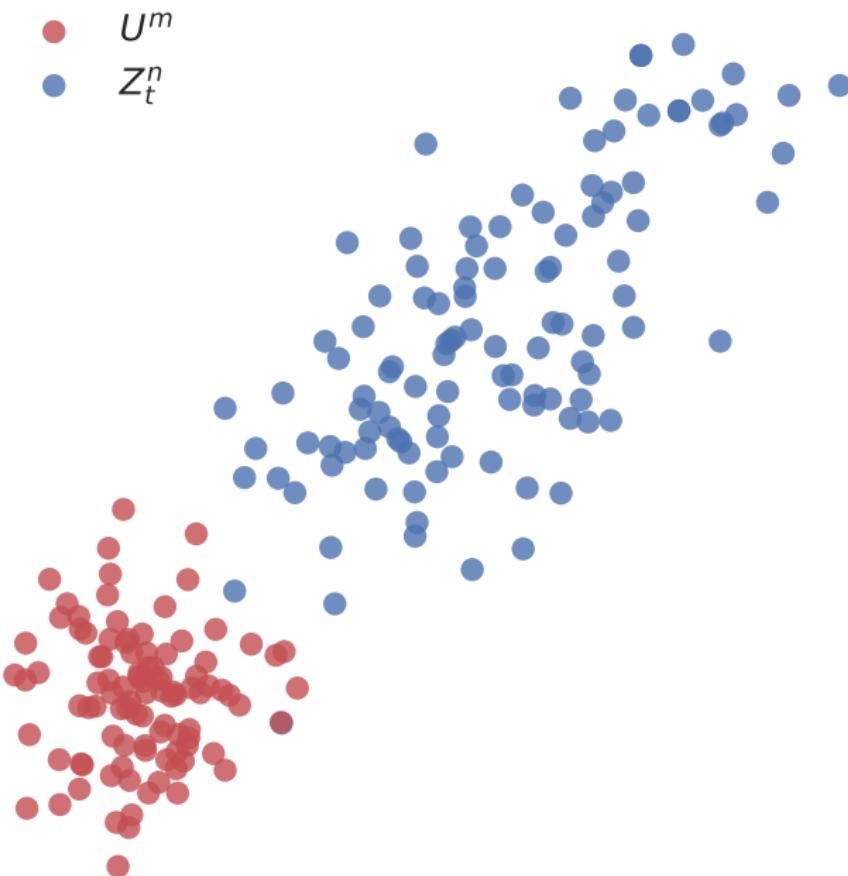
- ▶ Similar to *continuation methods*⁴, but extended to interacting particles.
- ▶ Different from entropic regularization⁵

$$Z_{t+1} = Z_t - \gamma \nabla f_{\nu^*, \nu_t}(Z_t) + \beta_t u_t$$

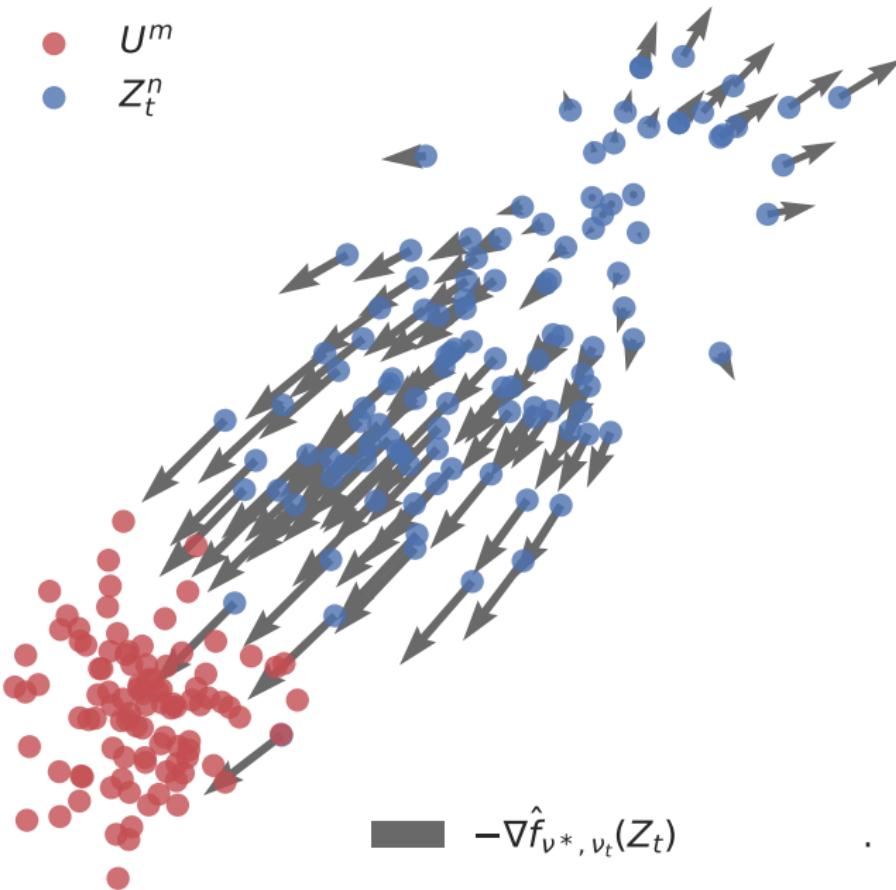
⁴[Chaudhari et al., 2017, Hazan et al., 2015]

⁵[Mei et al., 2018]

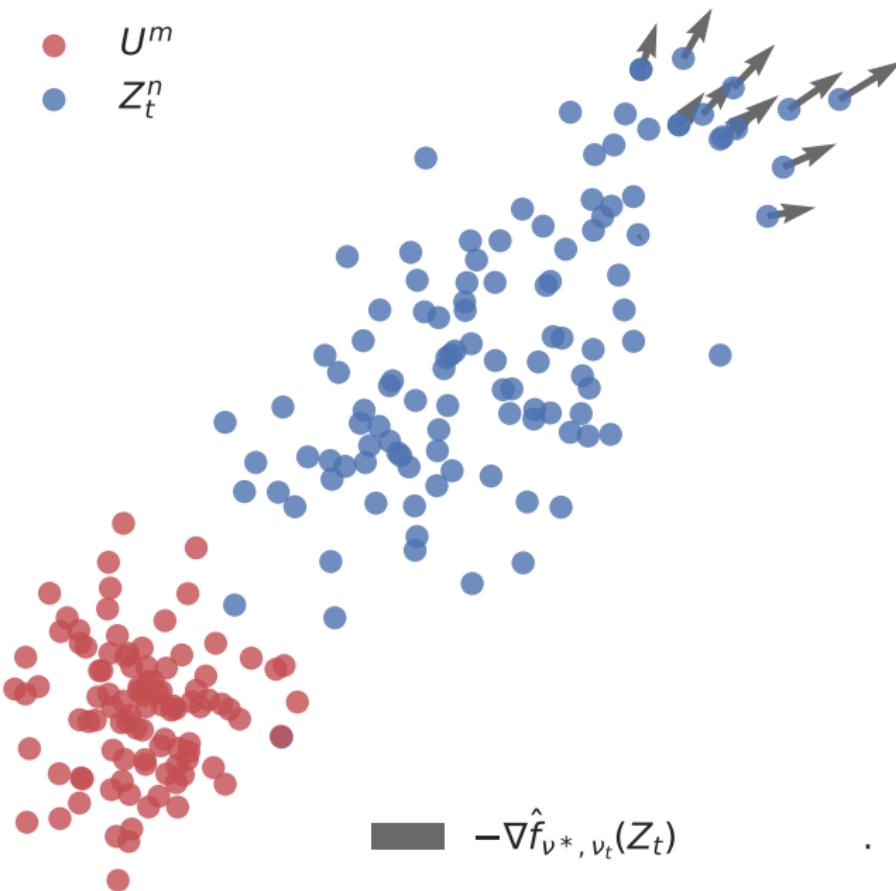
Noise Injection



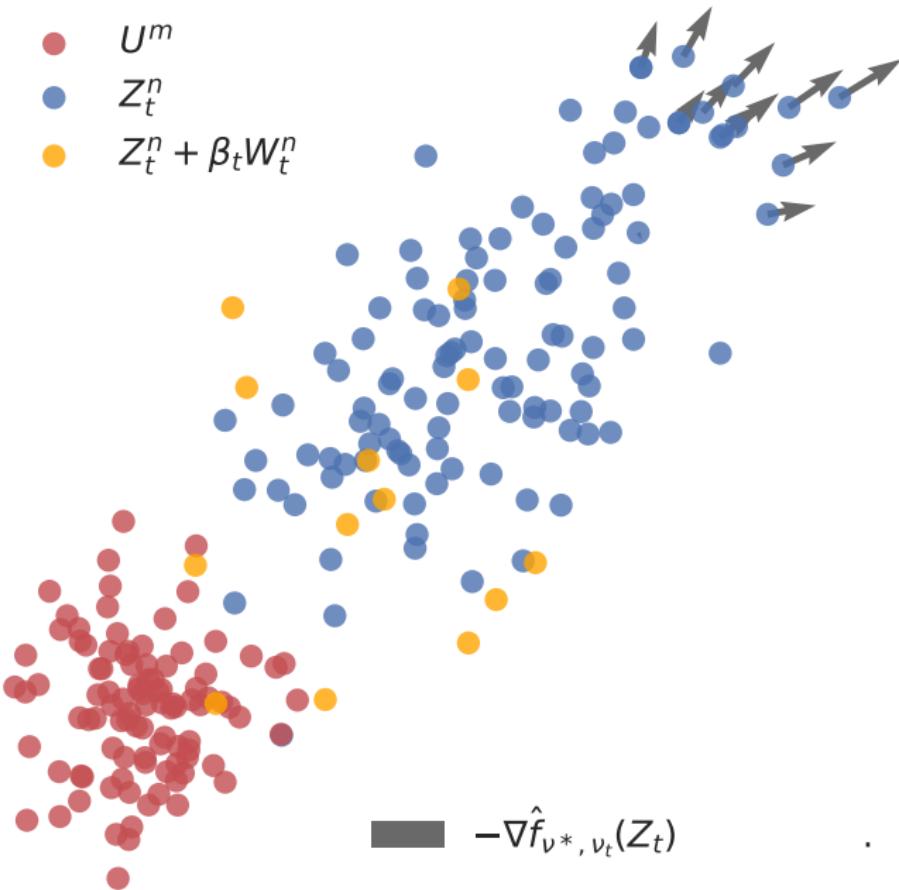
Noise Injection



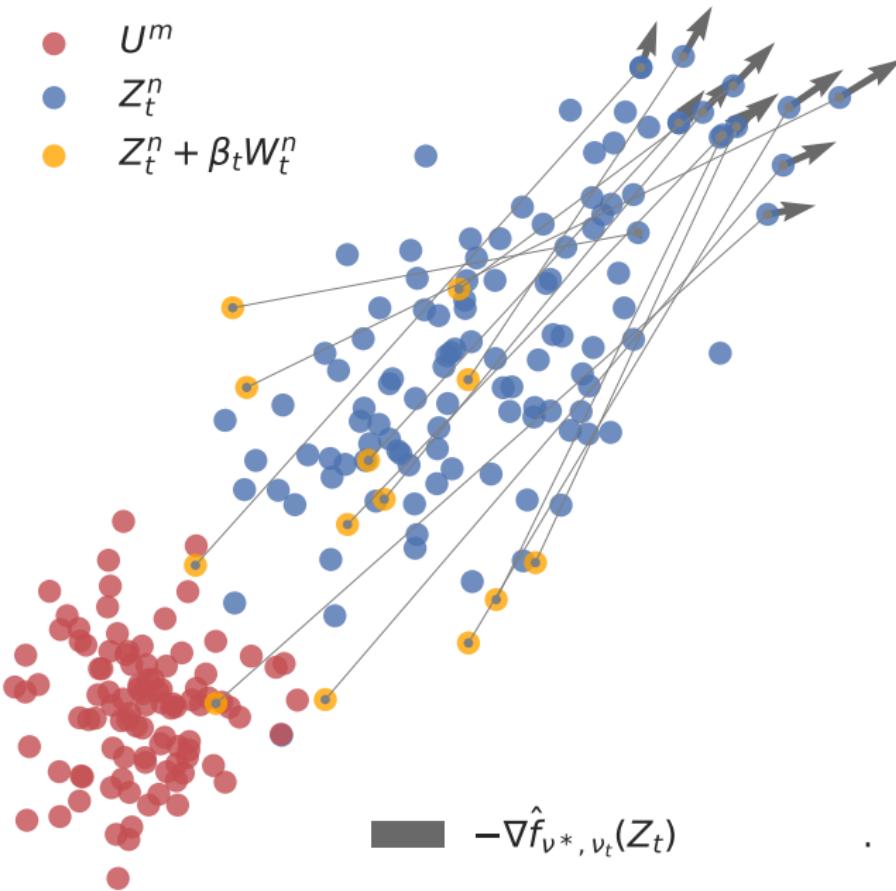
Noise Injection



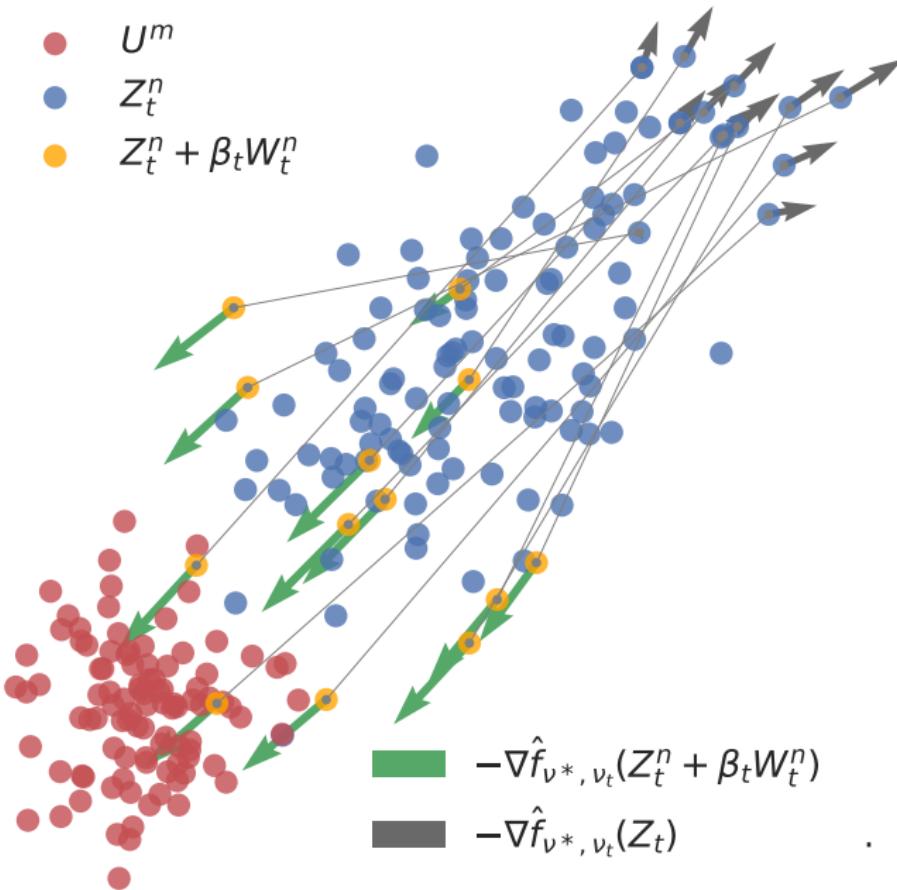
Noise Injection



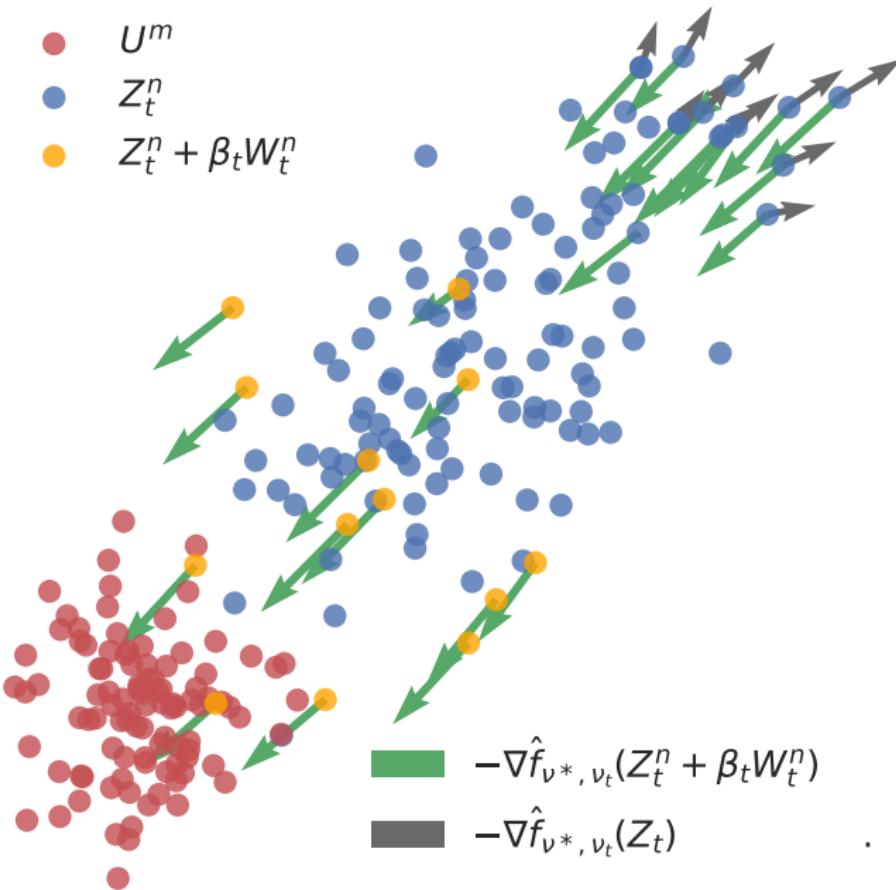
Noise Injection



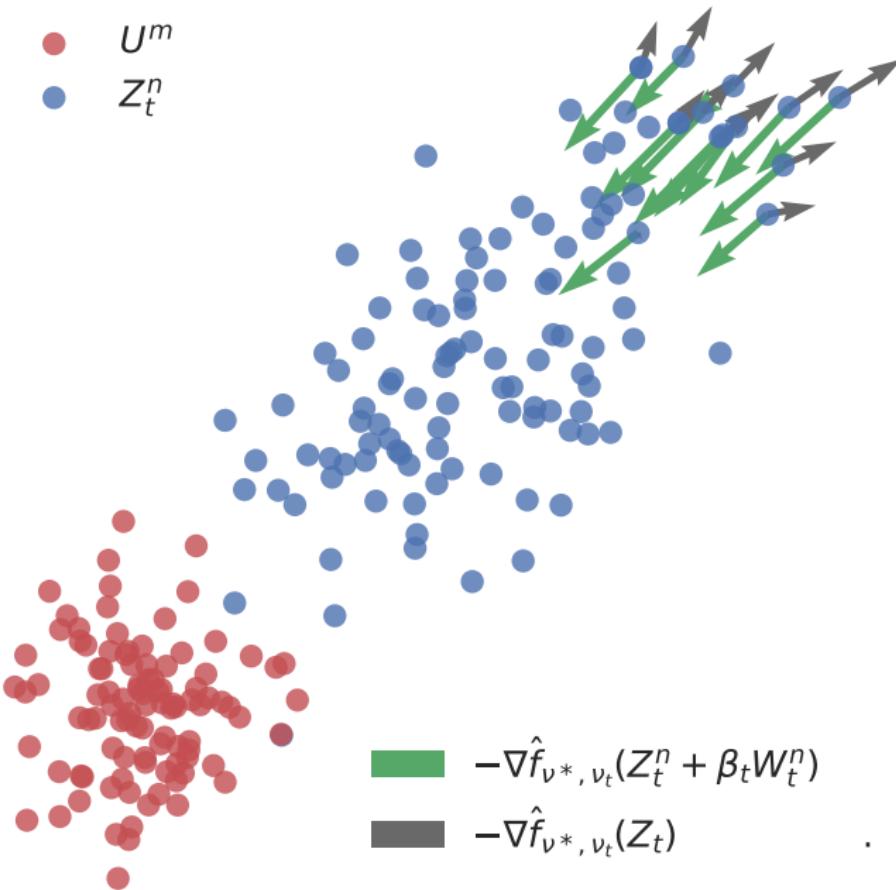
Noise Injection



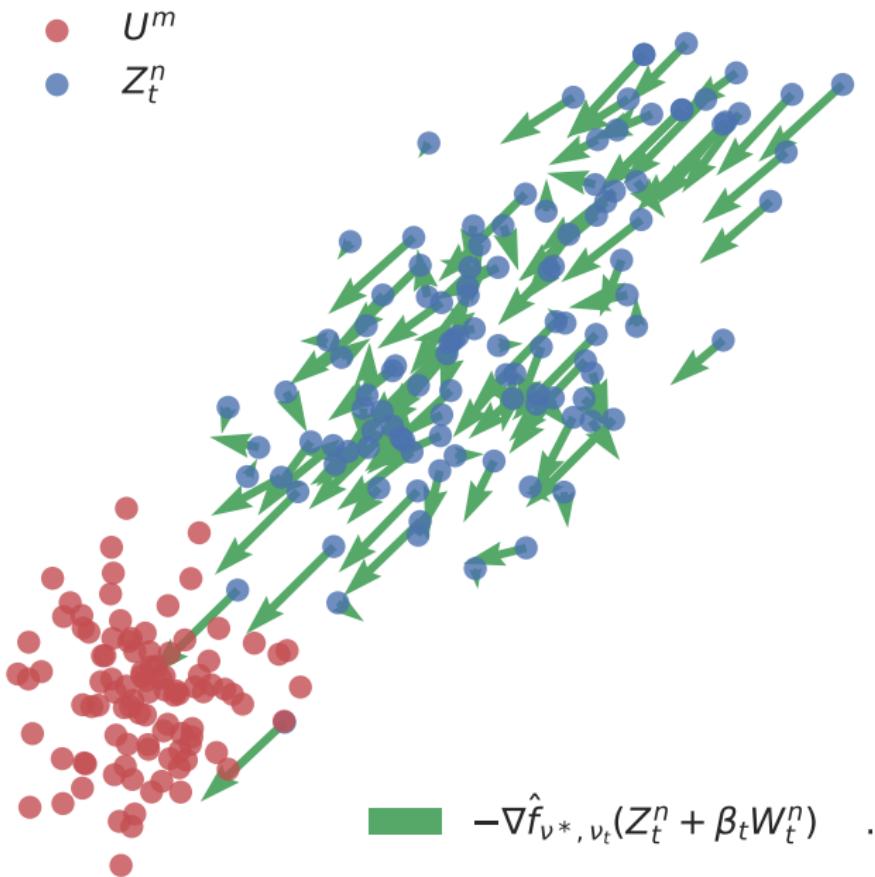
Noise Injection



Noise Injection



Noise Injection

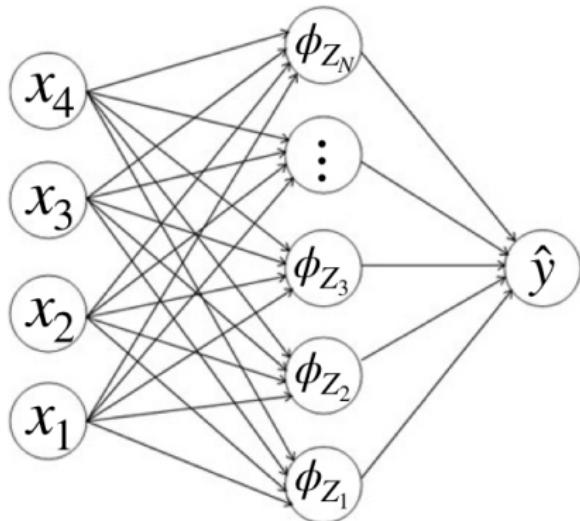


Noise Injection: Experiments

Noise Injection: Experiments

Noise Injection: Experiments

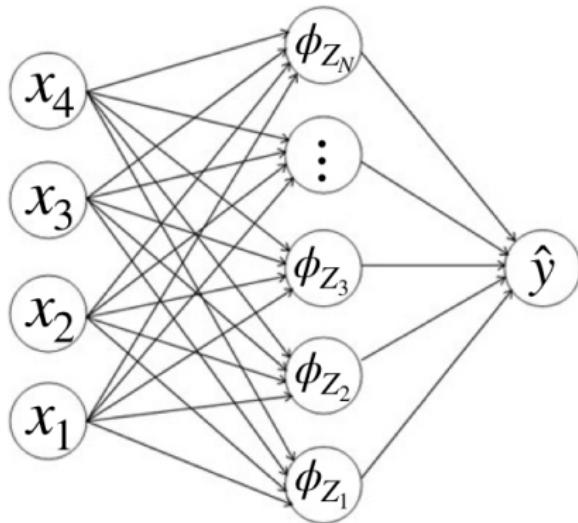
$(x, y) \sim data$



$$\min_{Z_1, \dots, Z_N} \mathbb{E}_{data} \left[\left\| \frac{1}{M} \sum_m \phi_{U^m}(x) - \frac{1}{N} \sum_{n=1}^N \phi_{Z^n}(x) \right\|^2 \right]$$

Noise Injection: Experiments

$$(x, y) \sim data$$

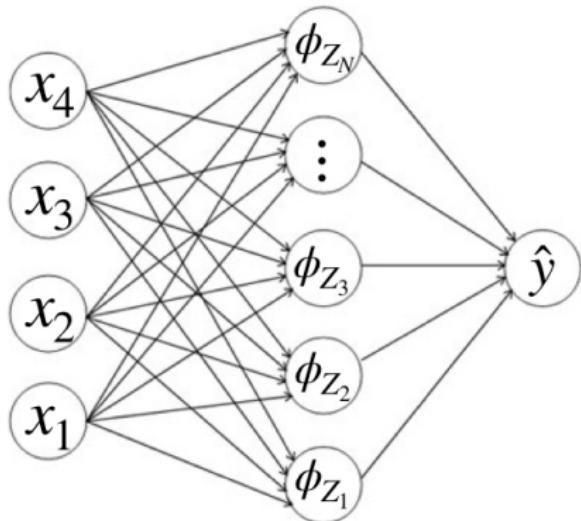


$$\min_{Z_1, \dots, Z_N} MMD^2(\nu^*, \frac{1}{N} \sum_{n=1}^N \delta_{Z^n})$$

$$k(Z, Z') = \mathbb{E}_{data}[\phi_Z(x)\phi_{Z'}(x)]$$

Noise Injection: Experiments

$(x, y) \sim data$



$$\min_{Z_1, \dots, Z_N} \mathbb{E}_{data} \left[\left\| \frac{1}{M} \sum_m^M \phi_{U^m}(x) - \frac{1}{N} \sum_{n=1}^N \phi_{Z^n}(x) \right\|^2 \right]$$

$$\hat{k}(Z, Z') = \frac{1}{B} \sum_{b=1}^B \phi_Z(x_b) \phi_{Z'}(x_b)$$

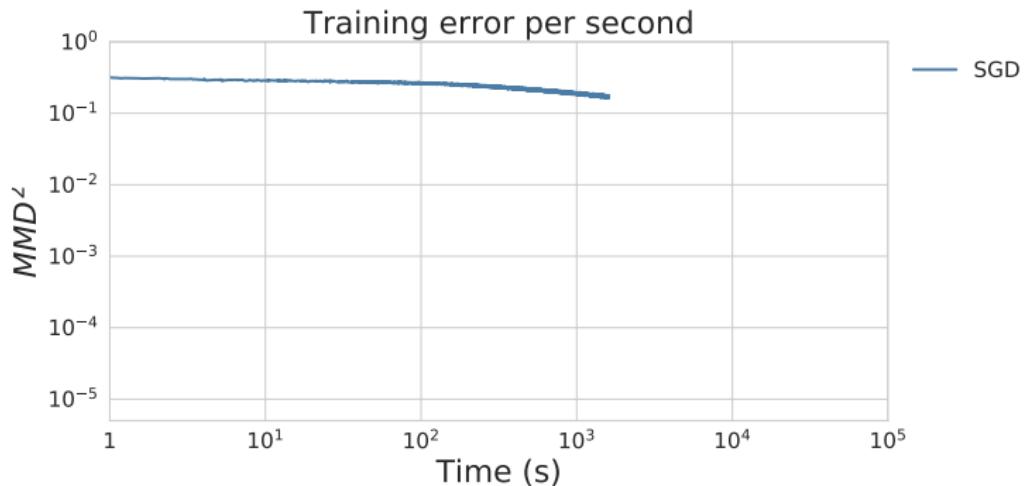
Noise Injection: Experiments

Methods:

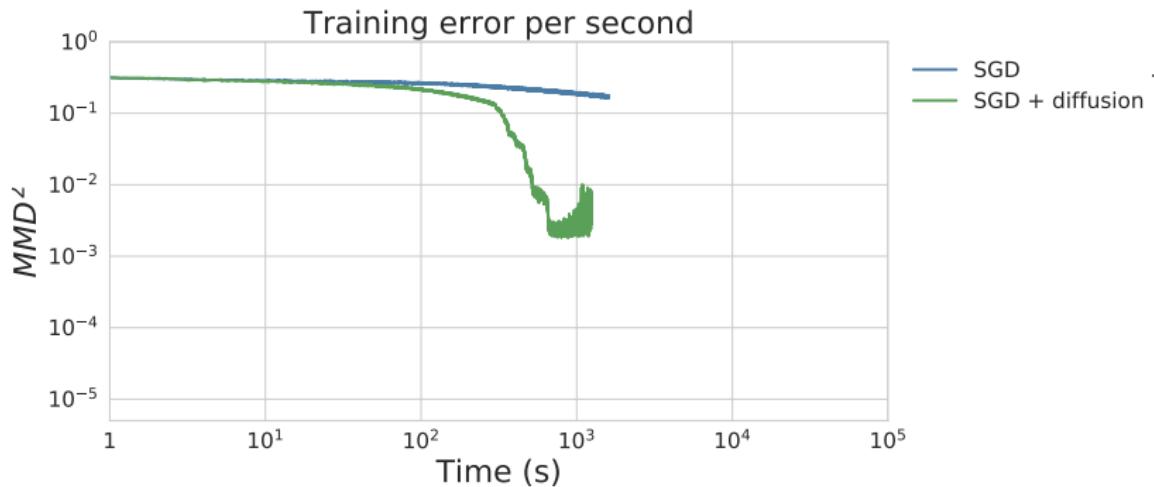
- ▶ SGD
- ▶ SGD + Noise injection
- ▶ SGD + diffusion
- ▶ KSD⁶: SGD using the Negative Sobolev distance
 $\nu \mapsto S(\nu^*|\nu)$ as a loss function: also minimizes the MMD.

⁶[Mroueh et al., 2019]

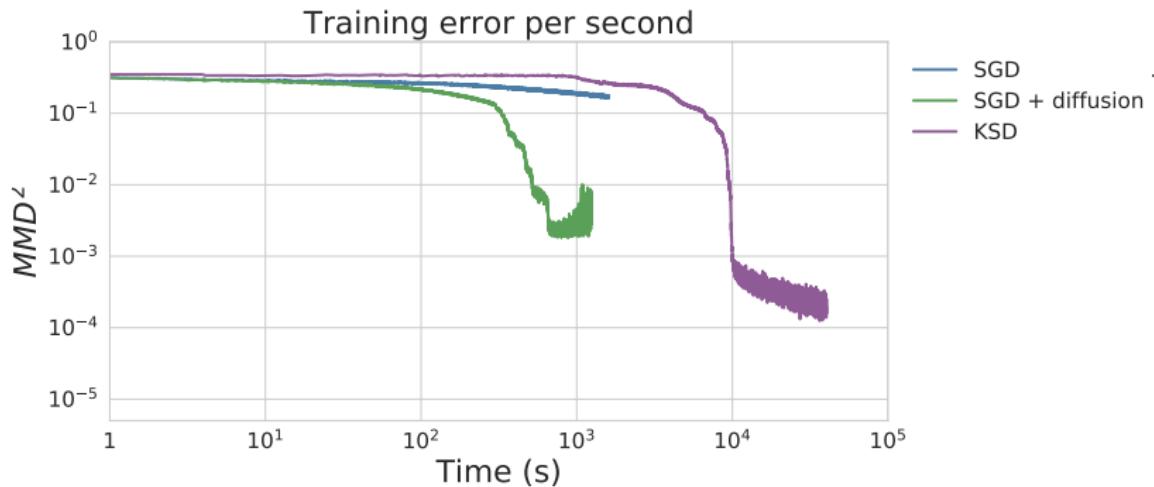
Noise Injection: Experiments



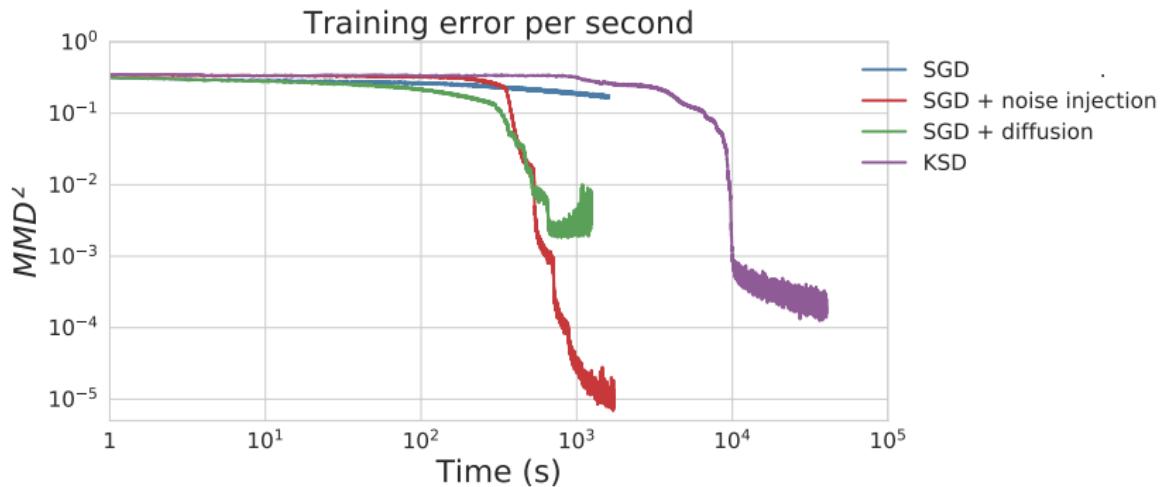
Noise Injection: Experiments



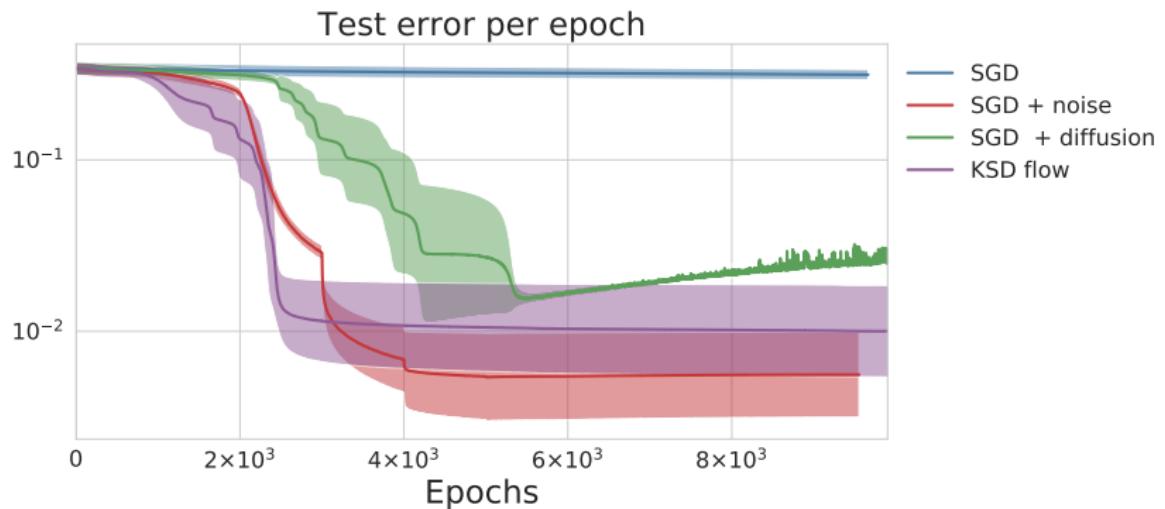
Noise Injection: Experiments



Noise Injection: Experiments

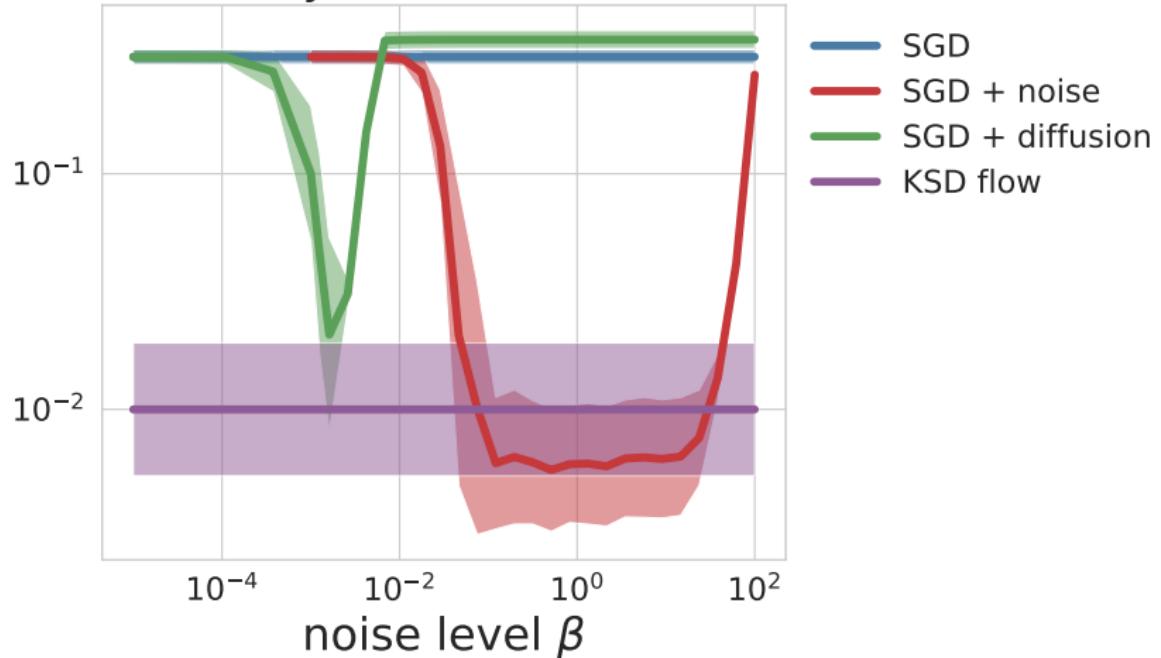


Noise Injection: Experiments



Noise Injection: Experiments

Sensitivity to noise (Test error)



Conclusion

Contributions:

- ▶ Provided a convergence criterion for the Wasserstein gradient descent.
- ▶ Proposed an extension to the noise injection algorithm for interacting particles and showed its effectiveness on simple examples.

Future work:

- ▶ A criterion for convergence that is independent from the whole optimization trajectory.
- ▶ Stronger guarantees for the convergence of the noise injection algorithm.

Thank you!

 Ambrosio, L., Gigli, N., and Savaré, G. (2004). Gradient flows with metric and differentiable structures, and applications to the Wasserstein space.

Atti della Accademia Nazionale dei Lincei. Classe di Scienze Fisiche, Matematiche e Naturali. Rendiconti Lincei. Matematica e Applicazioni, 15(3-4):327–343.

 Chaudhari, P., Oberman, A., Osher, S., Soatto, S., and Carlier, G. (2017).

Deep Relaxation: partial differential equations for optimizing deep neural networks.

arXiv:1704.04932 [cs, math].

arXiv: 1704.04932.

 Chizat, L. and Bach, F. (2018).

On the Global Convergence of Gradient Descent for Over-parameterized Models using Optimal Transport.

 Gretton, A., Sejdinovic, D., Strathmann, H., Balakrishnan, S., Pontil, M., Fukumizu, K., and Sriperumbudur, B. K. (2012).

Optimal kernel choice for large-scale two-sample tests

Noise Injection: Theory

Tradeoff for β_t

- ▶ Large β_t : μ_{t+1} not a descent direction anymore:
 $MMD^2(\nu^*, \mu_{t+1}) > MMD^2(\nu^*, \mu_t)$

Noise Injection: Theory

Tradeoff for β_t

- ▶ Large β_t : μ_{t+1} not a descent direction anymore:
 $MMD^2(\nu^*, \mu_{t+1}) > MMD^2(\nu^*, \mu_t)$
- ▶ Small β_t : Back to the failure mode: $\nabla f_t(X_t + \beta_t u_t) \simeq 0$.

Noise Injection: Theory

Tradeoff for β_t

- ▶ Large β_t : μ_{t+1} not a descent direction anymore:
 $MMD^2(\nu^*, \mu_{t+1}) > MMD^2(\nu^*, \mu_t)$
- ▶ Small β_t : Back to the failure mode: $\nabla f_t(X_t + \beta_t u_t) \simeq 0$.

Need β_t such that:

$$MMD^2(\nu^*, \mu_{t+1}) - MMD^2(\nu^*, \mu_t) \leq C\gamma \mathbb{E}_{\substack{X_t \sim \mu_t \\ U_t \sim \mathcal{N}(0,1)}} [\|\nabla f_t(X_t + \beta_t U_t)\|^2]$$

and:

$$\sum_i^t \beta_i^2 \rightarrow \infty$$

Noise Injection: Theory

Tradeoff for β_t

- ▶ Large β_t : μ_{t+1} not a descent direction anymore:
 $MMD^2(\nu^*, \mu_{t+1}) > MMD^2(\nu^*, \mu_t)$
- ▶ Small β_t : Back to the failure mode: $\nabla f_t(X_t + \beta_t u_t) \simeq 0$.

Need β_t such that:

$$MMD^2(\nu^*, \mu_{t+1}) - MMD^2(\nu^*, \mu_t) \leq C\gamma \mathbb{E}_{\substack{X_t \sim \mu_t \\ U_t \sim \mathcal{N}(0,1)}} [\|\nabla f_t(X_t + \beta_t U_t)\|^2]$$

and:

$$\sum_i^t \beta_i^2 \rightarrow \infty$$

Then

$$MMD^2(\nu^*, \nu_t) \leq MMD^2(\nu^*, \nu_0) e^{-C\gamma \sum_i^t \beta_i^2}$$