

# The Large and Small Noise Regimes in Diffusion Models

Yi Zhang

2024/10/10

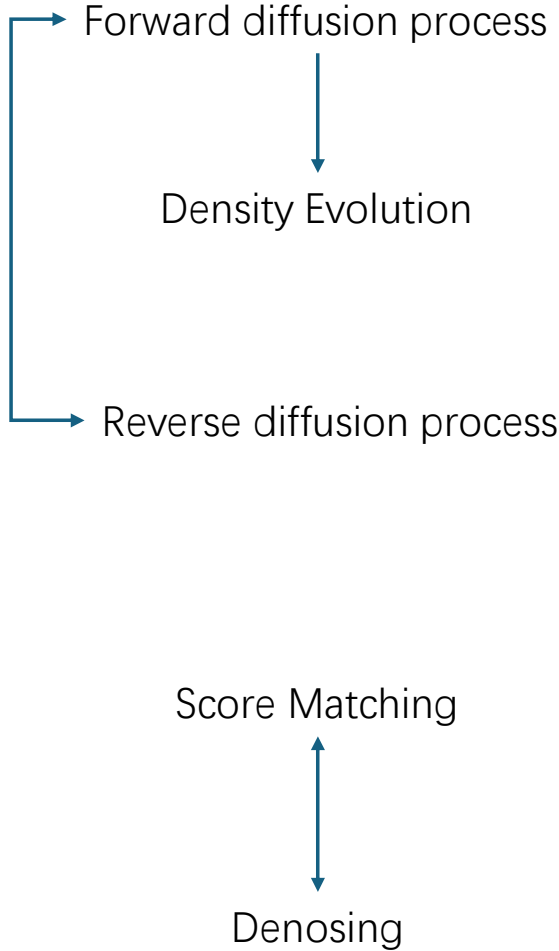
HKU IDS

Shah, K., Chen, S., & Klivans, A. (2023). Learning mixtures of gaussians using the DDPM objective. *Advances in Neural Information Processing Systems*, 36, 19636-19649.

# Outline

- Denoising for Score Matching
  - Preliminary
  - Data settings MoG and network structure
- Large Noise Regime: Connection to power iteration
  - Approximate DDPM loss gradient using Taylor's Expansion
  - One-step DDPM loss (population) gradient decent can be approximated by a matrix by Taylor's Expansion.
- Small Noise Regime: Connection to EM algorithm
  - One-step DDPM loss (population) gradient decent can be (partially) seen as a gradient EM step.
  - The remaining terms are very small.
- A trade-off between learning in large and small noise regime
  - Collapse error induced by determinant sampler
  - Analysis

- Denoising for score matching



### 1.3 Preliminaries

**Diffusion models.** Throughout the paper, we use either  $q$  or  $q_0$  to denote the data distribution and  $X$  or  $X_0$  to denote the corresponding random variable on  $\mathbb{R}^d$ . The two main components in diffusion models are the *forward process* and the *reverse process*. The forward process transforms samples from the data distribution into noise, for instance via the *Ornstein-Uhlenbeck (OU) process*:

$$dX_t = -X_t dt + \sqrt{2} dW_t \quad \text{with} \quad X_0 \sim q_0,$$

where  $(W_t)_{t \geq 0}$  is a standard Brownian motion in  $\mathbb{R}^d$ . We use  $q_t$  to denote the law of the OU process at time  $t$ . Note that for  $X_t \sim q_t$ ,

$$X_t = \exp(-t)X_0 + \sqrt{1 - \exp(-2t)}Z_t \quad \text{with} \quad X_0 \sim q_0, \quad Z_t \sim \mathcal{N}(0, \text{Id}).$$

The reverse process then transforms noise into samples, thus performing generative modeling. Ideally, this could be achieved by running the following stochastic differential equation for some choice of terminal time  $T$ :

$$dX_t^{\leftarrow} = \{X_t^{\leftarrow} + 2\nabla_x \ln q_{T-t}(X_t^{\leftarrow})\} dt + \sqrt{2} dW_t \quad \text{with} \quad X_0^{\leftarrow} \sim q_T,$$

where now  $W_t$  is the reversed Brownian motion. In this reverse process, the iterate  $X_t^{\leftarrow}$  is distributed according to  $q_{T-t}$  for every  $t \in [0, T]$ , so that the final iterate  $X_T^{\leftarrow}$  is distributed according to the data distribution  $q_0$ . The function  $\nabla_x \ln q_t$  is called the *score function*, and because it depends on  $q$  which is unknown, in practice one estimates it by minimizing the *score matching loss*

$$\min_{s_t} \mathbb{E}_{X_t \sim q_t} [\|s_t(X_t) - \nabla_x \ln q_t(X_t)\|^2]. \quad (4)$$

A standard calculation (see e.g. Appendix A of [CCL<sup>+</sup>23b]) shows that this is equivalent to minimizing the *DDPM objective* in which one wants to predict the noise  $Z_t$  from the noisy observation  $X_t$ , i.e.

$$\min_{s_t} L_t(s_t) = \mathbb{E}_{X_0, Z_t} \left[ \left\| s_t(X_t) + \frac{Z_t}{\sqrt{1 - \exp(-2t)}} \right\|^2 \right]. \quad (5)$$

# MoG Settings

We begin by describing in greater detail the algorithm we analyze in this work. For the sake of intuition, in this overview we will focus on the case of mixtures of two Gaussians ( $K = 2$ ) where the centers are well-separated and symmetric about the origin, that is, the data distribution is given by

$$q = \frac{1}{2}\mathcal{N}(\mu^*, \text{Id}) + \frac{1}{2}\mathcal{N}(-\mu^*, \text{Id}). \quad (1)$$

At the end of the overview, we briefly discuss the key challenges for handling smaller separation and general  $K$ .

Two gaussians with symmetric means  
and identity covariance matrix

# Construct a single-layer “MLP” with Tanh Activation

**Lemma 4.** *The score function for distribution  $q_t$ , for any  $t > 0$ , is given by*

$$\nabla_x \ln q_t(x) = \sum_{i=1}^K w_{i,t}^*(x) \mu_{i,t}^* - x, \quad \text{where} \quad w_{i,t}^*(x) = \frac{\exp(-\|x - \mu_{i,t}^*\|^2/2)}{\sum_{j=1}^K \exp(-\|x - \mu_{j,t}^*\|^2/2)}.$$

Calculate Score Function

*For a mixture of two Gaussians, the score function simplifies to*

$$\nabla_x \log q_t(x) = \tanh(\mu_t^{*\top} x) \mu_t^* - x, \quad \text{where} \quad \mu_t^* \triangleq \mu^* \exp(-t)$$

See Appendix A for the calculation.

Recall that  $\nabla_x \log q_t(x)$  is the minimizer for the score-matching objective given in Eq. (4). Therefore, we parametrize our student network architecture similarly to the optimal score function. Our student architecture for mixtures of  $K$  Gaussians is

Construct network structure

$$s_{\theta_t}(x) = \sum_{i=1}^K w_{i,t}(x) \mu_{i,t} - x, \quad \text{where} \quad w_{i,t}(x) \triangleq \frac{\exp(-\|x - \mu_{i,t}\|^2/2)}{\sum_{j=1}^K \exp(-\|x - \mu_{j,t}\|^2/2)} \quad (9)$$

$$\mu_{i,t} \triangleq \mu_i \exp(-t).$$

where  $\theta_t = \{\mu_{1,t}, \mu_{2,t}, \dots, \mu_{K,t}\}$  denotes the set of parameters at the noise scale  $t$ . For mixtures of two Gaussians, we simplify the student architecture as follows:

$$s_{\theta_t}(x) = \tanh(\mu_t^\top x) \mu_t - x, \quad \text{where} \quad \mu_t \triangleq \mu \exp(-t).$$

$\mu$  is the learnable parameters

# Derivation

$$\tanh(\mu^{*T} x) = w_{1,t}^*(x) + w_{2,t}^*(x)$$

Tanh comes from the summation of the belongingness of  $x$  to the 2 Gaussians.

## A.2 Derivation of score function

*Proof of Lemma 4.* For mixtures of  $K$  Gaussians in the form of Eq. (6), the score function at time  $t$  is given by

$$\begin{aligned} \nabla \log q_t(x) &= - \frac{\sum_{i=1}^K e^{-\frac{\|x - \mu_{i,t}^*\|^2}{2}} (x - \mu_{i,t}^*)}{\sum_{j=1}^K e^{-\frac{\|x - \mu_{j,t}^*\|^2}{2}}} \\ &= \sum_{i=1}^K w_{i,t}^*(x) \mu_{i,t}^* - x \quad \text{where} \quad w_{i,t}^*(x) = \frac{e^{-\frac{\|x - \mu_{i,t}^*\|^2}{2}}}{\sum_{j=1}^K e^{-\frac{\|x - \mu_{j,t}^*\|^2}{2}}}. \end{aligned}$$

For mixtures of two Gaussians in the form of Eq. (7), the score function is given by

$$\begin{aligned} \nabla \log q_t(x) &= w_{1,t}^*(x) \mu_{1,t}^* + w_{2,t}^*(x) \mu_{2,t}^* - x \\ &= \boxed{w_{1,t}^*(x) \mu^* - (1 - w_{1,t}^*(x)) \mu^*} - x \\ &= (2w_{1,t}^*(x) - 1) \mu^* - x \end{aligned} \tag{A.1}$$

By simplifying  $w_{1,t}^*(x)$ , we obtain

$$\begin{aligned} w_{1,t}^*(x) &= \frac{1}{1 + \exp(\frac{\|x - \mu^*\|^2}{2} - \frac{\|x + \mu^*\|^2}{2})} \\ &= \frac{1}{1 + \exp(-2\mu^{*\top} x)} \\ &= \sigma(2\mu^{*\top} x) \end{aligned} \tag{A.2}$$

where  $\sigma(\cdot)$  denotes the sigmoid function. Using Eq. (A.2) in Eq. (A.1), we obtain

$$\nabla \log q_t(x) = \boxed{\tanh(\mu^{*\top} x)} \mu^* - x.$$

□



# Recall: EM step for MoG Modeling

**Expectation-Maximization (EM) algorithm.** The EM algorithm is composed of two steps: the E-step and the M-step. For mixtures of Gaussians, the E-step computes the expected log-likelihood based on the current mean parameters and the M-step maximizes this expectation to find a new estimate of the parameters.

**Fact 5** (See e.g., [DTZ17, YYS17, KC20] for more details). *When  $X$  is the mixture of  $K$  Gaussian and  $\{\mu_1, \mu_2, \dots, \mu_K\}$  are current estimates of the means, the population EM update for all  $i \in \{1, 2, \dots, K\}$  is given by*

$$\text{M step} \rightarrow \mu'_i = \frac{\mathbb{E}_X[w_i(X)X]}{\mathbb{E}_X[w_i(X)]}, \quad \text{where } w_i(X) = \frac{\exp(-\|X - \mu_i\|^2/2)}{\sum_{j=1}^K \exp(-\|X - \mu_j\|^2/2)}.$$

Latent Variable

E step

$q = q_0 = \frac{1}{2}\mathcal{N}(\mu^*, \text{Id}) + \frac{1}{2}\mathcal{N}(-\mu^*, \text{Id}).$  (7)

*The EM update for mixtures of two Gaussians given in Eq. (7) simplifies to*

$$\mu' = \mathbb{E}_{X \sim \mathcal{N}(\mu^*, \text{Id})}[\tanh(\mu^\top X)X].$$

An analogous version of the EM algorithm, called the gradient EM algorithm, takes a gradient step in the direction of the M-step instead of optimizing the objective in the M-step fully.

# Derivation from the cited paper DTZ17

**Analysis of Population EM for Mixtures of Two Gaussians.** To elucidate the optimization features of the algorithm and avoid analytical distractions arising due to sampling error, it has been standard practice in the literature of theoretical analyses of EM to consider the “population version” of the algorithm, where the EM iterations are performed assuming access to infinitely many samples from a distribution  $p_{\mu_1, \mu_2}$  as above. With infinitely many samples, we can identify the mean,  $\frac{\mu_1 + \mu_2}{2}$ , of  $p_{\mu_1, \mu_2}$ , and re-parametrize the density around the mean as follows:

$$p_{\mu}(\mathbf{x}) = 0.5 \cdot \mathcal{N}(\mathbf{x}; \mu, \Sigma) + 0.5 \cdot \mathcal{N}(\mathbf{x}; -\mu, \Sigma). \quad (1.1)$$

We first study the convergence of EM when we perform iterations with respect to the parameter  $\mu$  of  $p_{\mu}(\mathbf{x})$  in (1.1). Starting with an initial guess  $\lambda^{(0)}$  for the unknown mean vector  $\mu$ , the  $t$ -th iteration of EM amounts to the following update:

$$\lambda^{(t+1)} = M(\lambda^{(t)}, \mu) \triangleq \frac{\mathbb{E}_{\mathbf{x} \sim p_{\mu}} \left[ \frac{0.5 \mathcal{N}(\mathbf{x}; \lambda^{(t)}, \Sigma)}{p_{\lambda^{(t)}}(\mathbf{x})} \mathbf{x} \right]}{\mathbb{E}_{\mathbf{x} \sim p_{\mu}} \left[ \frac{0.5 \mathcal{N}(\mathbf{x}; \lambda^{(t)}, \Sigma)}{p_{\lambda^{(t)}}(\mathbf{x})} \right]}, \quad (1.2)$$

where we have compacted both the E- and M-step of EM into one update.

The intuition behind the EM update formula is as follows. First, we take expectations with respect to  $\mathbf{x} \sim p_{\mu}$  because we are studying the population version of EM, hence we assume access to infinitely many samples from  $p_{\mu}$ . For each sample  $\mathbf{x}$ , the ratio  $\frac{0.5 \mathcal{N}(\mathbf{x}; \lambda^{(t)}, \Sigma)}{p_{\lambda^{(t)}}(\mathbf{x})}$  is our belief, at step  $t$ , that  $\mathbf{x}$  was sampled from the first Gaussian component of  $p_{\mu}$ , namely the one for which our current estimate of its mean vector is  $\lambda^{(t)}$ . (The complementary probability is our present belief that  $\mathbf{x}$  was sampled from the other Gaussian component.) Given these beliefs for all vectors  $\mathbf{x}$ , the update (1.2) is the result of the M-step of EM. Intuitively, our next guess  $\lambda^{(t+1)}$  for the mean vector of the first Gaussian component is a weighted combination over all samples  $\mathbf{x} \sim p_{\mu}$  where the weight of every  $\mathbf{x}$  is our belief that it came from the first Gaussian component.

## 2 Preliminary Observations

In this section we illustrate some simple properties of the EM update (1.2) and simplify the formula. First, it is easy to see that plugging in the values  $\lambda \in \{-\mu, \mathbf{0}, \mu\}$  into  $M(\lambda, \mu)$  results into

$$M(-\mu, \mu) = -\mu \quad ; \quad M(\mathbf{0}, \mu) = \mathbf{0} \quad ; \quad M(\mu, \mu) = \mu. \quad (2)$$

In particular, for all  $\mu$ , these values are certainly fixed points of the EM iteration. Next, we rewrite  $M(\lambda, \mu)$  as follows:

$$M(\lambda, \mu) = \frac{\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)} \left[ \frac{0.5 \mathcal{N}(\mathbf{x}; \lambda, \Sigma)}{p_{\lambda}(\mathbf{x})} \mathbf{x} \right] + \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(-\mu, \Sigma)} \left[ \frac{0.5 \mathcal{N}(\mathbf{x}; \lambda, \Sigma)}{p_{\lambda}(\mathbf{x})} \mathbf{x} \right]}{\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)} \left[ \frac{0.5 \mathcal{N}(\mathbf{x}; \lambda, \Sigma)}{p_{\lambda}(\mathbf{x})} \right] + \frac{1}{2} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(-\mu, \Sigma)} \left[ \frac{0.5 \mathcal{N}(\mathbf{x}; \lambda, \Sigma)}{p_{\lambda}(\mathbf{x})} \right]}.$$

It is easy to observe that by symmetry this simplifies to

$$M(\lambda, \mu) = \frac{\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)} \left[ \frac{\frac{1}{2} \mathcal{N}(\mathbf{x}; \lambda, \Sigma) - \frac{1}{2} \mathcal{N}(\mathbf{x}; -\lambda, \Sigma)}{\frac{1}{2} \mathcal{N}(\mathbf{x}; \lambda, \Sigma) + \frac{1}{2} \mathcal{N}(\mathbf{x}; -\lambda, \Sigma)} \mathbf{x} \right]}{\frac{1}{2} \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)} \left[ \frac{\frac{1}{2} \mathcal{N}(\mathbf{x}; \lambda, \Sigma) + \frac{1}{2} \mathcal{N}(\mathbf{x}; -\lambda, \Sigma)}{\frac{1}{2} \mathcal{N}(\mathbf{x}; \lambda, \Sigma) + \frac{1}{2} \mathcal{N}(\mathbf{x}; -\lambda, \Sigma)} \right]} = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)} \left[ \frac{\mathcal{N}(\mathbf{x}; \lambda, \Sigma) - \mathcal{N}(\mathbf{x}; -\lambda, \Sigma)}{\mathcal{N}(\mathbf{x}; \lambda, \Sigma) + \mathcal{N}(\mathbf{x}; -\lambda, \Sigma)} \mathbf{x} \right].$$

Simplifying common terms in the density functions  $\mathcal{N}(\mathbf{x}; \lambda, \Sigma)$ , we get that

$$M(\lambda, \mu) = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)} \left[ \frac{\exp(\lambda^T \Sigma^{-1} \mathbf{x}) - \exp(-\lambda^T \Sigma^{-1} \mathbf{x})}{\exp(\lambda^T \Sigma^{-1} \mathbf{x}) + \exp(-\lambda^T \Sigma^{-1} \mathbf{x})} \mathbf{x} \right].$$

We thus get the following expression for the EM iteration

$$M(\lambda, \mu) = \mathbb{E}_{\mathbf{x} \sim \mathcal{N}(\mu, \Sigma)} [\tanh(\lambda^T \Sigma^{-1} \mathbf{x}) \mathbf{x}]. \quad (2)$$

$\tanh(\mu^{*T} \mathbf{x}) = w_{1,t}^*(\mathbf{x}) + w_{2,t}^*(\mathbf{x})$  Tanh also comes from the summation of the belief (belongingness of  $\mathbf{x}$  to the 2 Gaussians)



# Large Noise Regime: Power Iteration

Gradient Approximation



Power Iteration



Angle Contraction

**Part I: Analysis of high noise regime and connection to power iteration.** We show that in the large noise regime, the negative gradient  $-\nabla L_t(s_t)$  is well-approximated by  $2\mu_t^* \mu_t^{*\top} \mu_t - 3\|\mu_t\|^2 \mu_t$ . Recall that this result is the key to showing the resemblance between gradient descent and power iteration. Concretely, we show the following lemma:

**Lemma 8** (See Lemma C.3 for more details). *For  $t = O(\log d)$ , the gradient descent update on the DDPM objective  $L_t(s_t)$  can be approximated with  $2\mu_t^* \mu_t^{*\top} \mu_t - 3\|\mu_t\|^2 \mu_t$ :*

$$\left\| (-\nabla L_t(s_t)) - \left( 2\mu_t^* \mu_t^{*\top} \mu_t - 3\|\mu_t\|^2 \mu_t \right) \right\| \leq \text{poly}(1/d).$$

From Lemma 8, it immediately follows that  $\mu'_t$ , the result of taking a single gradient step starting from  $\mu_t$ , is well-approximated by the result of taking a single step of power iteration for a matrix whose leading eigenvector is  $\mu_t^*$ :

$$\mu'_t = \mu_t - \eta \nabla L_t(s_t) \approx (\text{Id}(1 - 3\eta\|\mu_t\|^2) + 2\mu_t^* \mu_t^{*\top}) \mu_t.$$

The second key element is to show that as a consequence of the above power iteration update, the gradient descent converges in *angular distance* to the leading eigenvector. Concretely, we show the following lemma:

**Lemma 9** (Informal, see Lemma C.5 for more details). *Suppose  $\mu'_t$  is the iterate after one step of gradient descent on the DDPM objective from  $\mu_t$ . Denote the angle between  $\mu_t$  and  $\mu_t^*$  to be  $\theta$  and between  $\mu'_t$  and  $\mu_t^*$  to be  $\theta'$ . In this case, we show that*

$$\tan \theta' = \max(\kappa_1 \tan \theta, \kappa_2),$$

where  $\kappa_1 < 1$  and  $\kappa_2 \leq 1/\text{poly}(d)$ .

Note  $\tan \theta' < \tan \theta$  implies that  $\theta' < \theta$  or equivalently  $\langle \hat{\mu}'_t, \hat{\mu}_t^* \rangle > \langle \hat{\mu}_t, \hat{\mu}_t^* \rangle$ . Thus, the above lemma shows that by taking a gradient step in the DDPM objective, the angle between  $\mu_t$  and  $\mu_t^*$  decreases. By iterating this, we obtain the following lemma:

$\mu_t^*$  is the first eigenvector, because it is the first eigenvector of  $\text{Id}$  and  $\mu_t^* \mu_t^{*T}$

# Calculate Loss Gradient by standard Calculus

A standard calculation (see e.g. Appendix A of [CCL<sup>+</sup>23b]) shows that this is equivalent to minimizing the *DDPM objective* in which one wants to predict the noise  $Z_t$  from the noisy observation  $X_t$ , i.e.

$$\min_{s_t} L_t(s_t) = \mathbb{E}_{X_0, Z_t} \left[ \left\| s_t(X_t) + \frac{Z_t}{\sqrt{1 - \exp(-2t)}} \right\|^2 \right]. \quad (5)$$

$$s_{\theta_t}(x) = \tanh(\mu_t^\top x) \mu_t - x, \quad \text{where } \mu_t \triangleq \mu \exp(-t).$$

*Proof of Lemma C.2.* By calculating the negative gradient of the DDPM objective in Eq. (5), we obtain

$$\begin{aligned} -\nabla_{\mu_t} L_t(s_{\mu_t}) &= -\mathbb{E}_{X_0, Z_t} \left[ (\tanh(\mu_t^\top X_t) I + \tanh'(\mu_t^\top X_t) X_t \mu_t^\top) (s_{\mu_t}(X_t) + \frac{Z_t}{\beta_t}) \right] \\ &= -\mathbb{E} \left[ (\tanh(\mu_t^\top X_t) I + \tanh'(\mu_t^\top X_t) X_t \mu_t^\top) (\tanh(\mu_t^\top X_t) \mu_t - X_t + \frac{Z_t}{\beta_t}) \right] \\ &= \mathbb{E} \left[ -\tanh^2(\mu_t^\top X_t) \mu_t - \tanh(\mu_t^\top X_t) \tanh'(\mu_t^\top X_t) X_t \|\mu_t\|^2 + \tanh(\mu_t^\top X_t) X_t \right. \\ &\quad \left. + \tanh'(\mu_t^\top X_t) \mu_t^\top X_t X_t - \tanh(\mu_t^\top X_t) \frac{Z_t}{\beta_t} - \tanh'(\mu_t^\top X_t) X_t \mu_t^\top \frac{Z_t}{\beta_t} \right] \end{aligned}$$

Take derivative on  $\mu_t$   
and simply expand

# Simplify Loss Gradient by Stein's Lemma

Because  $Z_t$  and  $X_t$  are dependent random variable, analyzing them at the same time is difficult.

## F Additional proofs

### F.1 Proof of Lemma C.2

*Proof of Lemma C.2.* By calculating the negative gradient of the DDPM objective in Eq. (5), we obtain

$$\begin{aligned} -\nabla_{\mu_t} L_t(s_{\mu_t}) &= -\mathbb{E}_{X_0, Z_t}[(\tanh(\mu_t^\top X_t)I + \tanh'(\mu_t^\top X_t)X_t\mu_t^\top)(s_{\mu_t}(X_t) + \frac{Z_t}{\beta_t})] \\ &= -\mathbb{E}[(\tanh(\mu_t^\top X_t)I + \tanh'(\mu_t^\top X_t)X_t\mu_t^\top)(\tanh(\mu_t^\top X_t)\mu_t - X_t + \frac{Z_t}{\beta_t})] \\ &= \mathbb{E}[-\tanh^2(\mu_t^\top X_t)\mu_t - \tanh(\mu_t^\top X_t)\tanh'(\mu_t^\top X_t)X_t\|\mu_t\|^2 + \tanh(\mu_t^\top X_t)X_t \\ &\quad + \tanh'(\mu_t^\top X_t)\mu_t^\top X_t X_t - \tanh(\mu_t^\top X_t)\frac{Z_t}{\beta_t} - \tanh'(\mu_t^\top X_t)X_t\mu_t^\top \frac{Z_t}{\beta_t}] \end{aligned}$$

By simplifying the gradient terms involving  $Z_t$  by the Stein's identity as in Lemma F.1 and plugging it back in the gradient, we obtain

$$\begin{aligned} -\nabla_{\mu_t} L_t(s_{\mu_t}) &= \mathbb{E}\left[\left(\tanh(\mu_t^\top X_t) - \tanh(\mu_t^\top X_t)\tanh'(\mu_t^\top X_t)\|\mu_t\|^2 + \tanh'(\mu_t^\top X_t)\mu_t^\top X_t\right)X_t\right] \\ &\quad - \mu_t - \mathbb{E}\left[\tanh''(\mu_t^\top X_t)\|\mu_t\|^2 X_t\right] - \mathbb{E}\left[\tanh'(\mu_t^\top X_t)\mu_t\right] \\ &= \mathbb{E}\left[\left(\tanh(\mu_t^\top X_t) - 0.5\tanh''(\mu_t^\top X_t)\|\mu_t\|^2 + \tanh'(\mu_t^\top X_t)\mu_t^\top X_t\right)X_t\right] \\ &\quad - \mu_t - \mathbb{E}\left[\tanh'(\mu_t^\top X_t)\mu_t\right] \end{aligned}$$

Observe that  $(\tanh(\mu^\top x) - \frac{1}{2}\tanh''(\mu^\top x)\|\mu\|^2 + \tanh'(\mu^\top x)\mu^\top x)x$  and  $\tanh'(\mu^\top x)$  are even functions and  $X_t$  is a symmetric distribution, therefore, for any even function  $f$ , we can write  $\mathbb{E}_{X_t}[f(X_t)] = \frac{1}{2}\mathbb{E}_{X_t \sim \mathcal{N}(\mu_t^\top, I_d)}[f(X_t)] + \frac{1}{2}\mathbb{E}_{X_t \sim \mathcal{N}(-\mu_t^\top, I_d)}[f(X_t)] = \mathbb{E}_{X_t \sim \mathcal{N}(\mu_t^\top, I_d)}[f(X_t)]$ . Applying this property of the even function on the gradient update, we obtain the result.  $\square$

**Lemma F.1.** When random variable  $X_t = \alpha_t X_0 + \beta_t Z_t$  where  $Z_t \sim \mathcal{N}(0, I)$ ,  $\alpha_t = \exp(-t)$  and  $\beta_t = \sqrt{1 - \exp(-2t)}$ , then for any  $t > 0$ , the following two equations hold.

$$\mathbb{E}_{X_0, Z_t}\left[\tanh(\mu_t^\top X_t)\frac{Z_t}{\beta_t} + \tanh^2(\mu_t^\top X_t)\mu_t\right] = \mu_t$$

$$\mathbb{E}_{X_0, Z_t}\left[\tanh'(\mu_t^\top X_t)\frac{\mu_t^\top Z_t}{\beta_t} X_t\right] = \mathbb{E}_{X_0, Z_t}\left[\tanh''(\mu_t^\top X_t)\|\mu_t\|^2 X_t + \tanh'(\mu_t^\top X_t)\mu_t\right]$$

*Proof.* Applying Stein's lemma on the first term, we get the first equation of the statement in the Lemma.

$$\begin{aligned} \mathbb{E}_{X_0, Z_t}\left[\tanh(\mu_t^\top X_t)\frac{Z_t}{\beta_t}\right] &= \mathbb{E}_{X_0, Z_t}\left[\tanh(\mu_t^\top(\alpha_t X_0 + \beta_t Z_t))\frac{Z_t}{\beta_t}\right] \xrightarrow{\text{Stein's lemma}} \mathbb{E}_{X_0, Z_t}\left[\tanh'(\mu_t^\top X_t)\mu_t\right] \\ &= \mathbb{E}_{X_0, Z_t}\left[\left(1 - \tanh^2(\mu_t^\top X_t)\right)\mu_t\right] \end{aligned}$$

For the second term, we have

$$\begin{aligned} \mathbb{E}\left[\tanh'(\mu_t^\top X_t)\frac{\mu_t^\top Z_t}{\beta_t} X_t\right] &= \mathbb{E}\left[\tanh'(\mu_t^\top X_t)\frac{\mu_t^\top Z_t}{\beta_t}\alpha_t X_0\right] + \mathbb{E}\left[\tanh'(\mu_t^\top X_t)\mu_t^\top Z_t Z_t\right] \\ &= \sum_{i=1}^d \mathbb{E}\left[\alpha_t X_0 \tanh'(\mu_t^\top X_t)\frac{\mu_t(i)Z_t(i)}{\beta_t}\right] + \mathbb{E}\left[\tanh'(\mu_t^\top X_t)\mu_t\right] + \mathbb{E}\left[\tanh''(\mu_t^\top X_t)\mu_t^\top Z_t \beta_t \mu_t\right] \\ &= \sum_{i=1}^d \mathbb{E}\left[\alpha_t X_0 \tanh''(\mu_t^\top X_t)\mu_t(i)\mu_t(i)\right] + \mathbb{E}\left[\tanh'(\mu_t^\top X_t)\mu_t\right] + \mathbb{E}\left[\tanh''(\mu_t^\top X_t)\mu_t^\top Z_t \beta_t \mu_t\right] \end{aligned}$$

where the second equality follows from the Stein's lemma on the  $\mathbb{E}[\tanh'(\mu_t^\top X_t)\mu_t^\top Z_t Z_t]$  and the last equality follows from the Stein's lemma on  $\mathbb{E}[\alpha_t X_0 \tanh''(\mu_t^\top X_t)\mu_t(i)Z_t(i)]$ . Applying Stein's inequality on the  $\mathbb{E}\left[\tanh''(\mu_t^\top X_t)\mu_t^\top Z_t \beta_t \mu_t\right]$ , we obtain

$$\begin{aligned} &= \mathbb{E}\left[\alpha_t X_0 \tanh''(\mu_t^\top X_t)\|\mu_t\|^2\right] + \mathbb{E}\left[\tanh'(\mu_t^\top X_t)\mu_t\right] + \sum_{i=1}^d \beta_t \mu_t \mathbb{E}\left[\tanh'''(\mu_t^\top X_t)\mu_t(i)\beta_t \mu_t(i)\right] \\ &= \mathbb{E}\left[X_t \tanh''(\mu_t^\top X_t)\|\mu_t\|^2\right] - \mathbb{E}\left[\beta_t Z_t \tanh''(\mu_t^\top X_t)\|\mu_t\|^2\right] + \mathbb{E}\left[\tanh'(\mu_t^\top X_t)\mu_t\right] \\ &\quad + \beta_t^2 \|\mu_t\|^2 \mu_t \mathbb{E}\left[\tanh'''(\mu_t^\top X_t)\right] \\ &= \mathbb{E}\left[X_t \tanh''(\mu_t^\top X_t)\|\mu_t\|^2\right] + \mathbb{E}\left[\tanh'(\mu_t^\top X_t)\mu_t\right]. \end{aligned}$$

I don't know why summation here

$\square$

After Stein's Lemma, the 2<sup>th</sup> order and  $Z_t$  are canceled



# Approximate loss gradient by Taylor Expansion

**Lemma C.3.** For any noise scale  $t > t'$  and number of samples  $n > n'$  where  $t' \lesssim \log d$  and  $n' = \Theta(\frac{d^4 B^3}{\varepsilon^2})$ , with high probability, the negative gradient of the diffusion model objective  $L_t(s_t)$  can be approximated by  $2\mu_t^* \mu_t^{*\top} \mu_t - 3\|\mu_t\|^2 \mu_t$ . More precisely, given independent samples  $\{x_{i,t}\}_{i=1,\dots,n}$  from  $q_t$  generated using noise vectors  $\{z_{i,t}\}_{i=1,\dots,n}$  sampled from  $\mathcal{N}(0, \text{Id})$ , we have

$$\left\| -\nabla \left( \frac{1}{n} \sum_{i=1}^n L_t(s_{\mu_t}(x_{i,t}, z_{i,t})) \right) - (2\mu_t^* \mu_t^{*\top} \mu_t - 3\|\mu_t\|^2 \mu_t) \right\| \leq 250\sqrt{d}\|\mu_t\|^5 + 10\|\mu_t\|^3 \|\mu_t^*\|^2 + \varepsilon.$$

*Proof.* Recall that the population gradient update on the DDPM objective is given by

$$\begin{aligned} -\nabla L_t(s_{\mu_t}) &= \mathbb{E}_{x \sim \mathcal{N}(\mu_t^*, \text{Id})} [\tanh(\mu_t^\top x) x - \frac{1}{2} \tanh''(\mu_t^\top x) \|\mu_t\|^2 x + \tanh'(\mu_t^\top x) \mu_t^\top x x] \\ &\quad - \mu_t - \mathbb{E}_{x \sim \mathcal{N}(\mu_t^*, \text{Id})} [\tanh'(\mu_t^\top x) \mu_t] \quad \xrightarrow{\text{Stein's Lemma}} \\ &= \mathbb{E}_{x \sim \mathcal{N}(\mu_t^*, \text{Id})} [\tanh(\mu_t^\top x) x - \frac{1}{2} \tanh''(\mu_t^\top x) \|\mu_t\|^2 x + \tanh'(\mu_t^\top x) \mu_t^\top x \mu_t^* \\ &\quad + \tanh''(\mu_t^\top x) \mu_t^\top x \mu_t] - \mu_t, \end{aligned}$$

where the last equality follows from the Stein's lemma on  $\mathbb{E}_{x \sim \mathcal{N}(\mu_t^*, \text{Id})} [\tanh'(\mu_t^\top x) \mu_t^\top x x]$ , as  $\mathbb{E}_{x \sim \mathcal{N}(\mu_t^*, \text{Id})} [\tanh'(\mu_t^\top x) \mu_t^\top x x] = \mathbb{E}_{x \sim \mathcal{N}(\mu_t^*, \text{Id})} [\tanh'(\mu_t^\top x) \mu_t^\top x \mu_t^* + \tanh'(\mu_t^\top x) \mu_t + \tanh''(\mu_t^\top x) \mu_t^\top x \mu_t]$ .

Using Taylor's theorem, we know that

$$\tanh(\mu_t^\top x) = \mu_t^\top x - \frac{2}{3}(\mu_t^\top x)^3 + O(\xi(x)^5) \quad \text{where } \xi(x) \in [0, \mu_t^\top x]$$

$$\Rightarrow \tanh(\mu_t^\top x) x = \mu_t^\top x x - \frac{2}{3}(\mu_t^\top x)^3 x + O(\xi(x)^5 x) \quad \text{It does not follow big O definition?}$$

$$\Rightarrow \left\| \mathbb{E}_{x \sim \mathcal{N}(\mu_t^*, \text{Id})} [\tanh(\mu_t^\top x) x] - \mathbb{E}_{x \sim \mathcal{N}(\mu_t^*, \text{Id})} [\mu_t^\top x x - \frac{2}{3}(\mu_t^\top x)^3 x] \right\| \leq \|\mathbb{E}[\xi(x)^5 x]\| \lesssim \sqrt{d}\|\mu_t\|^5$$

where the last inequality follows from  $\|\mathbb{E}[\xi(x)^5 x]\| \leq \mathbb{E}[\|\mu_t^\top x\|^5 \|x\|] \leq (\mathbb{E}[\|\mu_t^\top x\|^{10}])^{1/2} (\mathbb{E}[\|x\|^2])^{1/2} \lesssim$

$\|\mu_t\|^5 \sqrt{d + \|\mu_t^*\|^2} \lesssim \sqrt{d}\|\mu_t\|^5$ . Similarly, using Taylor's theorem, we get

Holds for large noise regime ( $\mu_t^*$  is small)

$$\begin{aligned} \tanh''(\mu_t^\top x) &= -2\mu_t^\top x + O(\xi(x)^3) \quad \text{where } \xi(x) \in [0, \mu_t^\top x] \\ \Rightarrow \tanh''(\mu_t^\top x) \left( -\frac{1}{2}\|\mu_t\|^2 x + \mu_t^\top x \mu_t \right) &= \left( -2\mu_t^\top x + O(\xi(x)^3) \right) \left( -\frac{1}{2}\|\mu_t\|^2 x + \mu_t^\top x \mu_t \right) \\ \Rightarrow \left\| \mathbb{E}[\tanh''(\mu_t^\top x) \left( -\frac{1}{2}\|\mu_t\|^2 x + \mu_t^\top x \mu_t \right)] - \mathbb{E} \left[ -2\mu_t^\top x \left( -\frac{1}{2}\|\mu_t\|^2 x + \mu_t^\top x \mu_t \right) \right] \right\| \\ &\leq \left\| -\frac{1}{2}\|\mu_t\|^2 \mathbb{E}_{x \sim \mathcal{N}(\mu_t^*, \text{Id})} [O(\xi(x)^3) x] + \mathbb{E}_{x \sim \mathcal{N}(\mu_t^*, \text{Id})} [O(\xi(x)^3) \mu_t^\top x \mu_t] \right\| \\ &\leq \frac{1}{2}\|\mu_t\|^2 \mathbb{E}[\|\mu_t^\top x\|^3 \|x\|] + \|\mu_t\| \mathbb{E}[\|\mu_t^\top x\|^4] \\ &\leq \frac{1}{2}\|\mu_t\|^2 \sqrt{\mathbb{E}[\|\mu_t^\top x\|^6] \mathbb{E}[\|x\|^2]} + \|\mu_t\| \mathbb{E}[\|\mu_t^\top x\|^4] \\ &\leq 10\|\mu_t\|^5 \sqrt{d} + 6\|\mu_t\|^5 \end{aligned}$$

Holds for large noise regime (x almost follows a standard gaussian)

$$\begin{aligned} \text{Why it is 10 and 6.} \\ \tanh'(\mu_t^\top x) &= 1 - (\mu_t^\top x)^2 + O(\xi(x)^4) \quad \text{where } \xi(x) \in [0, \mu_t^\top x] \\ \Rightarrow \tanh'(\mu_t^\top x) \mu_t^\top x \mu_t^* &= \mu_t^\top x \mu_t^* - (\mu_t^\top x)^3 \mu_t^* + O(\xi(x)^4 \mu_t^\top x \mu_t^*) \quad \text{where } \xi(x) \in [0, \mu_t^\top x] \\ \Rightarrow \left\| \mathbb{E}[\tanh'(\mu_t^\top x) \mu_t^\top x \mu_t^*] - \mathbb{E}[\mu_t^\top x \mu_t^* - (\mu_t^\top x)^3 \mu_t^*] \right\| &\leq \left\| \mathbb{E}[\xi(x)^4 (\mu_t^\top x) \mu_t^*] \right\| \\ &\leq \mathbb{E}[\|\mu_t^\top x\|^5] \|\mu_t^*\| \lesssim \|\mu_t^*\| \|\mu_t\|^5 \end{aligned}$$

Holds for large noise regime (x almost follows a standard gaussian)

# Approximate loss gradient by Taylor Expansion

Sum up all Taylor approximations

$$\begin{aligned}
 & \mathbb{E}_{x \sim \mathcal{N}(\mu_t^*, \text{Id})} [xx^\top \mu_t (1 + \|\mu_t\|^2) - \frac{2}{3} (\mu_t^\top x)^3 x - 2\mu_t (\mu_t^\top x)^2 + \mu_t^\top x \mu_t^* - (\mu_t^\top x)^3 \mu_t^*] \\
 &= (I + \mu_t^* \mu_t^{*\top}) \mu_t (1 + \|\mu_t\|^2) - \frac{5}{3} \mathbb{E}[(\mu_t^\top x)^3 \mu_t^*] + \mu_t^* \mu_t^{*\top} \mu_t - 4 \mathbb{E}[\mu_t (\mu_t^\top x)^2] \quad \text{Stein's Lemma} \\
 &= (I + \mu_t^* \mu_t^{*\top}) \mu_t (1 + \|\mu_t\|^2) - \frac{5\mu_t^*}{3} ((\mu_t^\top \mu_t^*)^3 + 3(\mu_t^\top \mu_t^*) \|\mu_t\|^2) \quad \text{Gaussian moments} \\
 &\quad + \mu_t^* \mu_t^{*\top} \mu_t - 4\mu_t (\|\mu_t\|^2 + (\mu_t^\top \mu_t^*)^2) \\
 &= \mu_t^* \mu_t^{*\top} \mu_t (2 - 4\|\mu_t\|^2) + \mu_t (1 - 3\|\mu_t\|^2) - \frac{5\mu_t^* (\mu_t^\top \mu_t^*)^3}{3} - 4\mu_t (\mu_t^\top \mu_t^*)^2
 \end{aligned}$$

Approximate Taylor approximation by another polynomials

$$\begin{aligned}
 & \| -\nabla L_t(s_{\mu_t}) - (2\mu_t^* \mu_t^{*\top} \mu_t - 3\|\mu_t\|^2 \mu_t) \| \\
 & \leq \left\| -\nabla L_t(s_{\mu_t}) - \mathbb{E}[xx^\top \mu_t (1 + \|\mu_t\|^2) - \frac{2}{3} (\mu_t^\top x)^3 x - 2\mu_t (\mu_t^\top x)^2 + \mu_t^\top x \mu_t^* - (\mu_t^\top x)^3 \mu_t^* - \mu_t] \right\| \\
 & \quad + \left\| \mathbb{E}[xx^\top \mu_t (1 + \|\mu_t\|^2) - \frac{2}{3} (\mu_t^\top x)^3 x - 2\mu_t (\mu_t^\top x)^2 + \mu_t^\top x \mu_t^* - (\mu_t^\top x)^3 \mu_t^* - \mu_t] \right. \\
 & \quad \left. - (2\mu_t^* \mu_t^{*\top} \mu_t - 3\|\mu_t\|^2 \mu_t) \right\| \\
 & \leq \left( 200\sqrt{d}\|\mu_t\|^5 + 10\|\mu_t\|^5 \sqrt{d} + 6\|\mu_t\|^5 + 20\|\mu_t^*\| \|\mu_t\|^5 \right) + 10\|\mu_t\|^3 \|\mu_t^*\|^2 \\
 & \leq 250\sqrt{d}\|\mu_t\|^5 + 10\|\mu_t\|^3 \|\mu_t^*\|^2
 \end{aligned}$$



# Angle Contraction

**Lemma C.5.** *Suppose that the vector  $\mu_t$  satisfies  $|\langle \hat{\mu}_t, \hat{\mu}_t^* \rangle| \geq \frac{1}{2d}$ , and let  $\mu'_t$  denote the iterate resulting from a single empirical gradient step with learning rate  $\eta$  starting from  $\mu_t$ . Suppose that the empirical gradient and the population gradient differ by at most  $\varepsilon$ . Denote the angle between  $\mu_t$  (resp.  $\mu'_t$ ) and  $\mu_t^*$  by  $\theta$  (resp.  $\theta'$ ). Then*

$$\tan \theta' = \max(\kappa_1 \tan \theta, \kappa_2)$$

for

$$\begin{aligned} \kappa_1 &= \frac{1 - 3\eta\|\mu_t\|^2}{1 - 3\eta\|\mu_t\|^2 + \eta(\|\mu_t^*\|^2 - 500\sqrt{d^3}\|\mu_t\|^4 - 20d\|\mu_t\|^2\|\mu_t^*\|^2 - \eta\tilde{\varepsilon})} \ , \\ \kappa_2 &= \frac{500\eta\sqrt{d^3}\|\mu_t\|^4 + 20\eta d\|\mu_t\|^2\|\mu_t^*\|^2 + \eta\tilde{\varepsilon}}{\|\mu_t^*\|^2} \quad \text{and} \quad \tilde{\varepsilon} \lesssim \frac{d\varepsilon}{\|\mu_t\|} \ . \end{aligned}$$

# Angle Contraction

This because in power iteration,  $\mu', \mu, \mu^*$  lie in the same plane.

*Proof.* Define  $\hat{\mu}_t^{*\perp}$  as the orthogonal vector to  $\mu_t^*$  in the plane of  $\mu_t$  and  $\mu_t^*$ . Note that  $\mu'_t$  still lies in this plane, so the orthogonal vector to  $\mu_t^*$  in the plane of  $\mu'_t$  and  $\mu_t^*$  is also given by  $\hat{\mu}_t^{*\perp}$ .

We have

$$\begin{aligned} \tan \theta' &= \frac{\langle \hat{\mu}_t^{*\perp}, \hat{\mu}_t' \rangle}{\langle \hat{\mu}_t^*, \hat{\mu}_t' \rangle} = \frac{\langle \hat{\mu}_t^{*\perp}, \mu_t' \rangle}{\langle \hat{\mu}_t^*, \mu_t' \rangle} \\ &= \frac{\langle \hat{\mu}_t^{*\perp}, \mu_t + \eta F(\mu_t, \mu_t^*) \rangle + \langle \hat{\mu}_t^{*\perp}, -\eta \nabla L_t(s_t) - \eta F(\mu_t, \mu_t^*) \rangle + \eta \varepsilon}{\langle \hat{\mu}_t^*, \mu_t + \eta F(\mu_t, \mu_t^*) \rangle + \langle \hat{\mu}_t^*, -\eta \nabla L_t(s_t) - \eta F(\mu_t, \mu_t^*) \rangle - \eta \varepsilon} \\ &\quad \text{where } F(\mu, \mu^*) = (2\mu_t^* \mu_t^{*\top} \mu_t - 3\|\mu_t\|^2 \mu_t) \\ &\leq \frac{\sigma_2 \langle \hat{\mu}_t^{*\perp}, \mu_t \rangle + \eta \|\nabla L_t(s_t) + F(\mu_t, \mu_t^*)\| + \eta \varepsilon}{\sigma_1 \langle \hat{\mu}_t^*, \mu_t \rangle - \eta \|\nabla L_t(s_t) + F(\mu_t, \mu_t^*)\| - \eta \varepsilon} \end{aligned} \quad (\text{C.1})$$

Because  $\text{Id} + F$  is symmetric, the eigenvectors are orthogonal to each other. Therefore, if  $\mu^*$  is the first eigenvector,  $(\text{Id} + F)\mu^{*T} \leq \sigma_2 \mu^{*T}$

where  $\sigma_1$  and  $\sigma_2$  are the first and second eigenvalues of  $\text{Id} + F(\mu_t, \mu_t^*) = (1 - 3\eta\|\mu_t\|^2)\text{Id} + 2\eta\mu_t^* \mu_t^{*\top}$ , given by

$$\begin{aligned} \sigma_1 &= 1 + \eta(2\|\mu_t^*\|^2 - 3\|\mu_t\|^2) \\ \sigma_2 &= 1 - 3\eta\|\mu_t\|^2. \end{aligned}$$

The last inequality (C.1) follows from the fact that

$$\begin{aligned} \langle \hat{\mu}_t^*, \mu_t + \eta F(\mu_t, \mu_t^*) \rangle &= \hat{\mu}_t^{*\top} ((1 - 3\eta\|\mu_t\|^2)\text{Id} + 2\eta\mu_t^* \mu_t^{*\top}) \mu_t \\ &= \mu_t^\top ((1 - 3\eta\|\mu_t\|^2)\text{Id} + 2\eta\mu_t^* \mu_t^{*\top}) \hat{\mu}_t^* = \sigma_1 \mu_t^\top \hat{\mu}_t^* \end{aligned}$$

because  $\hat{\mu}_t^*$  is the first eigenvector of  $(1 - 3\eta\|\mu_t\|^2)\text{Id} + 2\eta\mu_t^* \mu_t^{*\top}$ . Recall from Lemma C.3 that the

# Angle Contraction

This is because of assumption that

$$\mu_t \text{ satisfies } |\langle \hat{\mu}_t, \hat{\mu}_t^* \rangle| \geq \frac{1}{2d},$$

because  $\hat{\mu}^*$  is the first eigenvector of  $(1 - 3\eta\|\mu_t\|^2)\text{Id} + 2\eta\mu_t^*\mu_t^{*\top}$ . Recall from Lemma C.3 that the deviation between the negative population gradient and the power iteration update  $F(\mu_t, \mu_t^*)$  is bounded by

$$\frac{\|\nabla L_t(s_t) + F(\mu_t, \mu_t^*)\|}{\langle \mu_t, \hat{\mu}_t^* \rangle} \leq \frac{250\eta\sqrt{d}\|\mu_t\|^4 + 10\eta\|\mu_t\|^2\|\mu_t^*\|^2}{\langle \hat{\mu}_t, \hat{\mu}_t^* \rangle} \leq 500\eta\sqrt{d^3}\|\mu_t\|^4 + 20d\eta\|\mu_t\|^2\|\mu_t^*\|^2.$$

Substituting this into Eq. (C.1), we get

$$\begin{aligned} \tan \theta' &\leq \frac{\sigma_2 \langle \hat{\mu}_t^{\perp}, \mu_t \rangle + \eta \|\nabla L_t(s_t) + F(\mu_t, \mu_t^*)\| + \eta \varepsilon}{\langle \hat{\mu}_t^*, \mu_t \rangle (\sigma_1 - 500\eta\sqrt{d^3}\|\mu_t\|^4 - 20d\eta\|\mu_t\|^2\|\mu_t^*\|^2 - \eta\tilde{\varepsilon})} \quad \text{where } \tilde{\varepsilon} \lesssim \frac{d\varepsilon}{\|\mu\|} \\ &\leq \frac{\sigma_2}{\tilde{\sigma}_1} \tan \theta + \frac{1}{\tilde{\sigma}_1} \left( 500\eta\sqrt{d^3}\|\mu\|^4 + 20d\eta\|\mu\|^2\|\mu_t^*\|^2 + \eta\tilde{\varepsilon} \right) \\ &\quad \text{where } \tilde{\sigma}_1 \triangleq \sigma_1 - 500\eta\sqrt{d^3}\|\mu\|^4 - 20d\eta\|\mu\|^2\|\mu_t^*\|^2 - \eta\tilde{\varepsilon} \\ &\leq \left( 1 - \frac{\eta\|\mu_t^*\|^2}{\tilde{\sigma}_1} \right) \frac{\sigma_2}{\tilde{\sigma}_1 - \eta\|\mu_t^*\|^2} \tan \theta + \left( \frac{\eta\|\mu_t^*\|^2}{\tilde{\sigma}_1} \right) \frac{500\eta\sqrt{d^3}\|\mu_t\|^4 + 20d\eta\|\mu_t\|^2\|\mu_t^*\|^2 + \eta\tilde{\varepsilon}}{\eta\|\mu_t^*\|^2} \\ &\leq \max \left( \frac{\sigma_2}{\tilde{\sigma}_1 - \eta\|\mu_t^*\|^2} \tan \theta, \frac{500\eta\sqrt{d^3}\|\mu_t\|^4 + 20\eta d \|\mu_t\|^2 \|\mu_t^*\|^2 + \eta\tilde{\varepsilon}}{\|\mu_t^*\|^2} \right) \end{aligned}$$

where the last inequality uses the fact that convex combinations of two values is less than the maximum of two values.  $\square$

# Low Noise Regime: EM algorithm

As before, we denote  $\mu_t$  as the current iterate and  $\mu'_t$  as the next iterate obtained by performing (population) gradient descent on the DDPM objective with step size  $\eta$ . We upper bound  $\|\mu'_t - \mu_t^*\|$  as follows:

$$\begin{aligned}\|\mu'_t - \mu_t^*\| &= \|\mu_t - \eta \nabla_{\mu_t} L_t(s_{\mu_t}) - \mu_t^*\| \\ &= \left\| (1 - \eta)(\mu_t - \mu_t^*) + \eta \mathbb{E}_{x \sim \mathcal{N}(\mu_t^*, 1)} \left[ \left( \tanh(\mu_t^\top x) - \frac{1}{2} \tanh''(\mu_t^\top x) \|\mu_t\|^2 \right. \right. \right. \\ &\quad \left. \left. \left. + \tanh'(\mu_t^\top x) \mu_t^\top x \right) x \right] - \eta \mathbb{E}_{x \sim \mathcal{N}(\mu_t^*, 1)} [\tanh'(\mu_t^\top x) \mu_t] - \eta \mu_t^* \right\| \\ &\leq (1 - \eta) \|\mu_t - \mu_t^*\| + \eta \left\| \mathbb{E}_{x \sim \mathcal{N}(\mu_t^*, 1)} [\tanh(\mu_t^\top x) x] - \mu_t^* \right\| + \eta \|G(\mu_t, \mu_t^*)\|,\end{aligned}$$

where

One-step of EM

$$G(\mu_t, \mu_t^*) \triangleq \mathbb{E}_{x \sim \mathcal{N}(\mu_t^*, \text{Id})} \left[ -\frac{1}{2} \tanh''(\mu_t^\top x) \|\mu_t\|^2 x + (\tanh'(\mu_t^\top x) \mu_t^\top x) x - \tanh'(\mu_t^\top x) \mu_t \right].$$

It remains to show  $G$  contracts.

Intuition: The first and high order derivative of Tanh decays exponentially with  $\mu$ .

And different from the large noise regime,  $\mu$  is large in small noise regime.

# Contraction of 1-D $G(\mu, \mu^*)$

**Lemma C.8** (One-dimensional version). Let  $\mu, \mu^* > 0$ , and consider  $\mu \in [c, \frac{4\mu^*}{3}]$  for some constant  $c$ . In this one-dimensional case, the function  $G$  specializes to

$$G(\mu, \mu^*) = \mathbb{E}_{x \sim \mathcal{N}(\mu^*, 1)} \left[ -\frac{1}{2} \tanh''(\mu x) \mu^2 x + \tanh'(\mu x) \mu x^2 - \tanh'(\mu x) \mu \right], \quad (\text{C.2})$$

In one 1-D case,  $G$  contracts.

and we have

$$G(\mu, \mu^*) \leq 0.01 |\mu - \mu^*|$$

The proof uses the fact that the function  $G$  only contains first or higher-order derivatives of the tanh function and all the derivatives of tanh decay exponential quickly as  $\mu$  increases. Therefore, when  $\mu$  is at least a constant, we obtain the result. The complete proof of lemma C.8 is given in Appendix F.2.

*Proof of Lemma C.8.* Recall that the gradient update for any  $\mu_t^*$  is given by

$$-\nabla_{\mu_t^*} L_t(s_{\mu_t^*}) = G(\mu_t^*, \mu_t^*) + \eta \mathbb{E}_{x \sim \mathcal{N}(\mu_t^*, 1)} [\tanh(\mu_t^{*\top} x) x] - \eta \mu_t^* \quad (\text{F.1})$$

We know that  $\mathbb{E}_{x \sim \mathcal{N}(\mu_t^*, 1)} [\tanh(\mu_t^{*\top} x) x] = \mu_t^*$  (Eq.(2.1) of [DTZ17]) and  $\nabla_{\mu_t^*} L_t(s_{\mu_t^*}) = 0$  because  $\mu_t^*$  is a stationary point of the regression objective of diffusion model. This implies that  $G(\mu_t^*, \mu_t^*) = 0$  for any  $\mu_t^*$ .

Note that this proof only talks about 1D case therefore, for the purpose of this proof, we use  $a$  to denote  $\mu$  and  $b$  to denote  $\mu^*$ . In 1D, using Mean value theorem, we have

$$\frac{G(a, b) - G(a, a)}{b - a} = \frac{dG(a, \xi)}{d\xi} \text{ for some } \xi \in [a, b] \text{ (if } a < b) \quad (\text{F.2})$$

Using the fact that  $G(a, a) = 0$  in Eq. (F.2), we have

$$|G(a, b)| = \left| \frac{dG(a, \xi)}{d\xi} \right| |b - a| \quad \text{Mean Value Thm}$$

Exponentially decay with  $a$  ( $\mu$ )

Observe that it suffices to prove  $\left| \frac{dG(a, \xi)}{d\xi} \right| \leq 0.01$  to obtain the lemma. By computing the gradient of  $G$ , we obtain

$$\frac{dG(a, \xi)}{d\xi} = \eta \mathbb{E}_{x \sim \mathcal{N}(\xi, 1)} \left[ 2 \tanh'(ax) ax + \tanh''(ax) \left( \frac{-3a^2}{2} + a^2 x^2 \right) - \frac{1}{2} a^3 x \tanh'''(ax) \right]$$



# Conclusion

- If we perfectly build a network to fit MoG score,
  - In large noise regime (small  $\mu_t$ ), the one-step DDPM loss gradient decent can be regarded as power iteration. (angle contraction)
    - Proved by Taylor's expansions and remaining terms are high orders of  $\mu_t$ .
  - In small noise regime (large  $\mu_t$ ), the gradient decent can be regarded as gradient EM algorithm. (amplitude contraction)
    - Proved by derivatives of tanh decays exponentially with  $\mu_t$