
The Optimal Noise in Noise-Contrastive Learning Is Not What You Think (Supplementary Material)

Omar Chehab¹

Alexandre Gramfort¹

Aapo Hyvärinen²

¹Université Paris-Saclay, Inria, CEA, Palaiseau, France

²Department of Computer Science, University of Helsinki, Finland

1 VISUALIZATIONS OF THE MSE LANDSCAPE

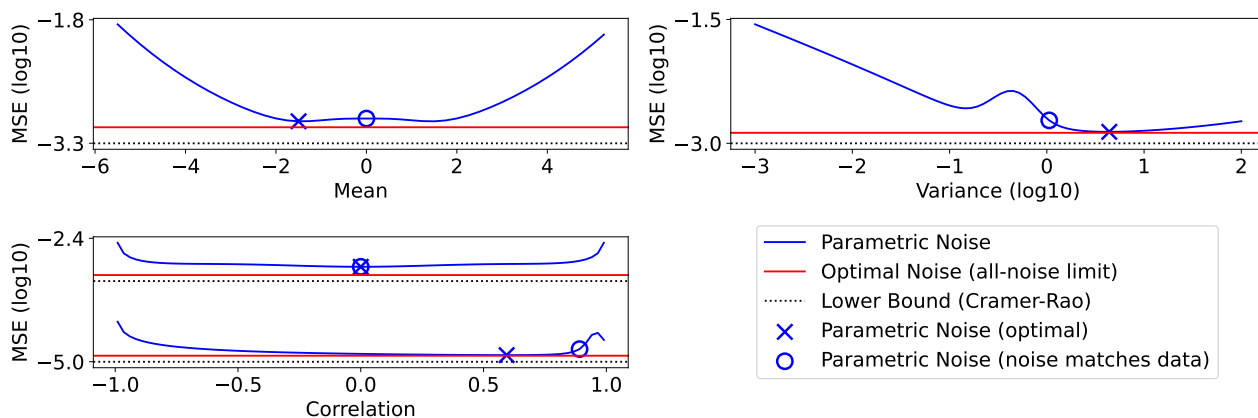


Figure 1: MSE vs. the noise parameter. Top left panel for model (i), Gaussian mean; Top right panel for model (ii), Gaussian variance; Bottom left for model (iii), Gaussian correlation.

We provide visualizations of the MSE landscape of the NCE estimator, when the noise is constrained within a parametric family containing the data.

We draw attention to the two local minima symmetrically placed to the left and to the right of the Gaussian mean. This corroborates the indeterminacies observed in this paper (Conjecture on limit of zero noise), as to where the optimal noise should place its mass for this estimation problem.

2 INTRACTABILITY OF THE 1D GAUSSIAN CASE

Suppose the data distribution p_d is a one-dimensional standardized zero-mean Gaussian. The model and noise distributions are of the same family, parameterized by mean and/or variance (we write these together in one model):

$$p_\theta(x) = \frac{1}{\sqrt{2\pi\alpha}} e^{-\frac{1}{2} \frac{(x-\mu)^2}{\alpha}}, \quad p_n(x) = \frac{1}{\sqrt{2\pi\beta}} e^{-\frac{1}{2} \left(\frac{x-\pi}{\beta}\right)^2} \quad x \in \mathbb{R}$$

We can write out the relevant functions, evaluated at $\alpha = 1, \mu = 0$ as the 2D score:

$$\mathbf{g}(x) = \left(\begin{array}{c} \partial_\mu \log p_\theta \\ \partial_\alpha \log p_\theta \end{array} \right) \Big|_{\mu=0, \alpha=1} = \left(\begin{array}{c} x \\ -1 + x^2 \end{array} \right)$$

and its ‘‘pointwise covariance’’: $\mathbf{g}(x)\mathbf{g}(x)^\top = \left(\begin{array}{cc} x^2 & -x + x^3 \\ -x + x^3 & x^4 - x^2 + 1 \end{array} \right)$

In the following, we consider estimation of variance only. i.e. only the second term in \mathbf{m} and the second diagonal term in the Fisher information matrix I . Now we can compute the generalized score mean m and mean of square I as they intervene in the MSE formula for Noise-Contrastive Estimation:

$$m = \int g(x)(1 - D(x))p(x)dx$$

which gives

$$m = -\frac{1}{2\sqrt{2\pi}} \int \left(e^{-\frac{x^2}{2}} \frac{1}{1 + \frac{1}{\nu}\sqrt{\beta}e^{-\frac{x^2}{2}(1-\frac{1}{\beta})}} \right) dx + \frac{1}{2\sqrt{2\pi}} \int x^2 \left(e^{-\frac{x^2}{2}} \frac{1}{1 + \frac{1}{\nu}\sqrt{\beta}e^{-\frac{x^2}{2}(1-\frac{1}{\beta})}} \right) dx$$

and

$$I = \int g(x)^2(1 - D(x))p(x)dx$$

which gives

$$I = \frac{1}{4\sqrt{2\pi}} \int x^4 \left(e^{-\frac{x^4}{2}} \frac{1}{1 + \frac{1}{\nu}\sqrt{\beta}e^{-\frac{x^2}{2}(1-\frac{1}{\beta})}} \right) dx - \frac{1}{2\alpha^3\sqrt{2\pi}} \int x^2 \left(e^{-\frac{x^2}{2}} \frac{1}{1 + \frac{1}{\nu}\sqrt{\beta}e^{-\frac{x^2}{2}(1-\frac{1}{\beta})}} \right) dx + \frac{1}{4\sqrt{2\pi}} \int \left(e^{-\frac{x^2}{2}} \frac{1}{1 + \frac{1}{\nu}\sqrt{\beta}e^{-\frac{x^2}{2}(1-\frac{1}{\beta})}} \right) dx$$

We see that even in a simple 1D Gaussian setting, evaluating the asymptotic MSE of the Noise-Contrastive Estimator is untractable in closed-form, given the integrals in I , where the integrand includes the product of a Gaussian density with the logistic function compounded by the Gaussian density, further multiplied by monomials. While here we considered the case of variance, the intractability is seen even in the case of the mean. Optimizing the asymptotic MSE with respect to β and π (noise distribution) or ν (identifiable to the noise proportion) yields similarly intractable integrals.

3 OPTIMAL NOISE PROPORTION WHEN THE NOISE DISTRIBUTION MATCHES THE DATA DISTRIBUTION: PROOF

We wish to minimize the MSE given by

$$\text{MSE}_{\text{NCE}}(T, \nu, p_n) = \frac{\nu + 1}{T} \text{tr}(\mathbf{I}^{-1} - \frac{\nu + 1}{\nu} (\mathbf{I}^{-1} \mathbf{m} \mathbf{m}^\top \mathbf{I}^{-1}))$$

when $p_n = p_d$. In that case,

$$D(\mathbf{x}) = \frac{p_d}{p_d + \nu p_n}(\mathbf{x}) = \frac{p_d}{p_d + \nu p_d}(\mathbf{x}) = \frac{1}{1 + \nu}$$

and the integrals involved become

$$\begin{aligned} \mathbf{m} &= \int \mathbf{g}(\mathbf{x})(1 - D(\mathbf{x}))p(\mathbf{x})d\mathbf{x} \\ &= \frac{\nu}{1 + \nu} \int \mathbf{g}(\mathbf{x})p(\mathbf{x})d\mathbf{x} \\ &= 0 \end{aligned}$$

given the score has zero mean, and

$$\begin{aligned} \mathbf{I} &= \int \mathbf{g}(\mathbf{x})\mathbf{g}(\mathbf{x})^\top (1 - D(\mathbf{x}))p(\mathbf{x})d\mathbf{x} \\ &= \frac{\nu}{1 + \nu} \int \mathbf{g}(\mathbf{x})\mathbf{g}(\mathbf{x})^\top p(\mathbf{x})d\mathbf{x} \\ &= \frac{\nu}{1 + \nu} \mathbf{I}_F . \end{aligned}$$

The objective function thus reduces to

$$\text{MSE}_{\text{NCE}}(T, \nu, p_n) = \frac{\nu + 1}{T} \text{tr}(\mathbf{I}^{-1}) = \frac{(\nu + 1)^2}{\nu T} \text{tr}(\mathbf{I}_F^{-1}) \propto \frac{(\nu + 1)^2}{\nu} .$$

The derivative with respect to ν is proportional to $\frac{1}{\nu^2} - 1$ and is null when $\nu = 1$ so when the noise proportion is 50%.

Note that in that case where $p_n = p_d$, we can compare the MSE achieved by NCE (using T_d data samples and T_n noise samples) with the MSE achieved by MLE (using T_d data samples):

$$\frac{\text{MSE}_{\text{NCE}}(T, \nu, p_n)}{\text{MSE}_{\text{MLE}}(T_d)} = \frac{\frac{(\nu+1)^2}{\nu T} \text{tr}(\mathbf{I}_F^{-1})}{\frac{1}{T_d} \text{tr}(\mathbf{I}_F^{-1})} = \frac{\frac{(\nu+1)^2}{\nu T} \text{tr}(\mathbf{I}_F^{-1})}{\frac{\nu+1}{T} \text{tr}(\mathbf{I}_F^{-1})} = 1 + \frac{1}{\nu}$$

which is known from [Gutmann and Hyvärinen, 2012, Pihlaja et al., 2010].

4 OPTIMAL NOISE FOR ESTIMATING A PARAMETER: PROOFS

We here prove the theorem and conjecture for the optimal noise distribution in three limit cases $\nu \rightarrow 0$ (all data samples), $\nu \rightarrow \infty$ (all noise samples), and $\frac{p_d}{p_n}(\cdot) = 1 + \epsilon(\cdot)$ as $\epsilon(\cdot) \rightarrow 0$ (noise distribution is an infinitesimal perturbation of the data distribution).

The goal is to optimize the $\text{MSE}_{\text{NCE}}(T, \nu, p_n)$ with respect to the noise distribution p_n , where

$$\text{MSE}_{\text{NCE}}(T, \nu, p_n) = \frac{\nu + 1}{T} \text{tr}(\mathbf{I}^{-1} - \frac{\nu + 1}{\nu} (\mathbf{I}^{-1} \mathbf{m} \mathbf{m}^\top \mathbf{I}^{-1})) \quad (1)$$

where the integrals

$$\begin{aligned} \mathbf{m} &= \int \mathbf{g}(\mathbf{x})(1 - D(\mathbf{x}))p(\mathbf{x})d\mathbf{x} \\ \mathbf{I} &= \int \mathbf{g}(\mathbf{x})\mathbf{g}(\mathbf{x})^\top (1 - D(\mathbf{x}))p(\mathbf{x})d\mathbf{x} \end{aligned}$$

depend non-linearly on p_n via the optimal discriminator:

$$1 - D(\mathbf{x}) = \frac{\nu p_n(\mathbf{x})}{p_d(\mathbf{x}) + \nu p_n(\mathbf{x})}$$

The general proof structure is:

- Perform a Taylor expansion of $1 - D(\mathbf{x})$ in the $\nu \rightarrow 0$ or $\nu \rightarrow \infty$ limit
- Plug into the integrals \mathbf{m} , \mathbf{I} and evaluate them (up to a certain order)
- Perform a Taylor expansion of \mathbf{I}^{-1} (up to a certain order)
- Evaluate the MSE_{NCE} (up to a certain order)

- Optimize the MSE_{NCE} w.r.t. p_n
- Compute the MSE gaps at optimality

Theorem 1 In either of the following two limits:

- (i) the noise distribution is a (infinitesimal) perturbation of the data distribution $\frac{p_d}{p_n} = 1 + \epsilon(x)$;
- (ii) in the limit of all noise samples $\nu \rightarrow \infty$;

the noise distribution minimizing asymptotic MSE is

$$p_n^{\text{opt}}(\mathbf{x}) \propto p_d(\mathbf{x}) \|\mathbf{I}_F^{-1} \mathbf{g}(\mathbf{x})\| . \quad (2)$$

Proof: case where $\nu \rightarrow \infty$.

We start with a change of variables $\gamma = \frac{1}{\nu} \rightarrow 0$ to bring us to a zero-limit.

The MSE in terms of our new variable $\gamma = \frac{1}{\nu}$ can be written as:

$$\text{MSE}_{\text{NCE}}(T, \gamma, p_n) = \frac{\gamma + 1}{\gamma T} \text{tr}(\mathbf{I}^{-1}) - \frac{(\gamma + 1)^2}{T\gamma} \text{tr}(\mathbf{I}^{-1} \mathbf{m} \mathbf{m}^\top \mathbf{I}^{-1}) \quad (3)$$

$$= \left(\gamma^{-1} T^{-1} + \gamma^0 T^{-1} \right) \text{tr}(\mathbf{I}^{-1}) - \left(\gamma^{-1} T^{-1} + \gamma^0 2T^{-1} + \gamma^1 T^{-1} \right) \text{tr}(\mathbf{I}^{-1} \mathbf{m} \mathbf{m}^\top \mathbf{I}^{-1}) \quad (4)$$

Given the term up until γ^{-1} in the MSE, we will use Taylor expansions up to order 2 throughout the proof, in anticipation that the MSE will be expanded until order 1.

- Taylor expansion of the discriminator

$$1 - D(\mathbf{x}) = \frac{\nu p_n(\mathbf{x})}{p_d(\mathbf{x}) + \nu p_n(\mathbf{x})} = \frac{1}{1 + \gamma \frac{p_d}{p_n}(\mathbf{x})} = 1 - \gamma \frac{p_d}{p_n}(\mathbf{x}) + \gamma^2 \frac{p_d^2}{p_n^2}(\mathbf{x}) + o(\gamma^2)$$

- Evaluating the integrals \mathbf{m} , \mathbf{I}

$$\begin{aligned} \mathbf{m} &= \int \mathbf{g}(\mathbf{x}) p_d(\mathbf{x}) \left(1 - D(\mathbf{x}) \right) d\mathbf{x} = \int \mathbf{g}(\mathbf{x}) p_d(\mathbf{x}) \left(1 - \gamma \frac{p_d}{p_n}(\mathbf{x}) + \gamma^2 \frac{p_d^2}{p_n^2}(\mathbf{x}) + o(\gamma^2) \right) d\mathbf{x} \\ &= \mathbf{m}_F - \gamma \mathbf{a} + \gamma^2 \mathbf{b} + o(\gamma^2) \end{aligned} \quad (5)$$

where \mathbf{m}_F is the Fisher-score mean of the (possibly unnormalized) model and we use shorthand notations \mathbf{a} and \mathbf{b} for the remaining integrals:

$$\begin{aligned} \mathbf{m}_F &= \int \mathbf{g}(\mathbf{x}) p_d(\mathbf{x}) d\mathbf{x} = 0 \\ \mathbf{a} &= \int \mathbf{g}(\mathbf{x}) \frac{p_d^2}{p_n}(\mathbf{x}) d\mathbf{x} \\ \mathbf{b} &= \int \mathbf{g}(\mathbf{x}) \frac{p_d^3}{p_n^2}(\mathbf{x}) d\mathbf{x} . \end{aligned}$$

Similarly,

$$\begin{aligned} \mathbf{I} &= \int \mathbf{g}(\mathbf{x}) \mathbf{g}(\mathbf{x})^\top p_d(\mathbf{x}) \left(1 - D(\mathbf{x}) \right) d\mathbf{x} = \int \mathbf{g}(\mathbf{x}) \mathbf{g}(\mathbf{x})^\top p_d(\mathbf{x}) \left(1 - \gamma \frac{p_d}{p_n}(\mathbf{x}) + \gamma^2 \frac{p_d^2}{p_n^2}(\mathbf{x}) + o(\gamma^2) \right) d\mathbf{x} \\ &= \mathbf{I}_F - \gamma \mathbf{A} + \gamma^2 \mathbf{B} + o(\gamma^2) \end{aligned}$$

where the Fisher-score covariance (Fisher information) is \mathbf{I}_F and we use shorthand notations \mathbf{A} and \mathbf{B} for the remaining integrals:

$$\begin{aligned}\mathbf{I}_F &= \int \mathbf{g}(\mathbf{x})\mathbf{g}(\mathbf{x})^\top p_d(\mathbf{x})d\mathbf{x} \\ \mathbf{A} &= \int \mathbf{g}(\mathbf{x})\mathbf{g}(\mathbf{x})^\top \frac{p_d^2}{p_n}(\mathbf{x})d\mathbf{x} \\ \mathbf{B} &= \int \mathbf{g}(\mathbf{x})\mathbf{g}(\mathbf{x})^\top \frac{p_d^3}{p_n^2}(\mathbf{x})d\mathbf{x} .\end{aligned}$$

- Taylor expansion of \mathbf{I}^{-1}

$$\begin{aligned}\mathbf{I}^{-1} &= \left(\mathbf{I}_F - \gamma\mathbf{A} + \gamma^2\mathbf{B} + \circ(\gamma^2) \right)^{-1} \\ &= \left(\mathbf{I}_F(\mathbf{Id} - \gamma\mathbf{I}_F^{-1}\mathbf{A} + \gamma^2\mathbf{I}_F^{-1}\mathbf{B}) + \circ(\gamma^2) \right)^{-1} \\ &= \mathbf{I}_F^{-1} \left(\mathbf{Id} - \gamma\mathbf{I}_F^{-1}\mathbf{A} + \gamma^2\mathbf{I}_F^{-1}\mathbf{B} \right)^{-1} + \circ(\gamma^2) \\ &= \mathbf{I}_F^{-1} \left(\mathbf{Id} + \gamma\mathbf{I}_F^{-1}\mathbf{A} + \gamma^2((\mathbf{I}_F^{-1}\mathbf{A})^2 - \mathbf{I}_F^{-1}\mathbf{B}) + \circ(\gamma^2) \right) + \circ(\gamma^2) \\ &= \mathbf{I}_F^{-1} + \gamma\mathbf{I}_F^{-2}\mathbf{A} + \gamma^2(\mathbf{I}_F^{-1}(\mathbf{I}_F^{-1}\mathbf{A})^2 - \mathbf{I}_F^{-2}\mathbf{B}) + \circ(\gamma^2)\end{aligned}\tag{6}$$

- Evaluating the MSE_{NCE}

$$\mathbf{I}^{-1}\mathbf{m}\mathbf{m}^\top\mathbf{I}^{-1} = \mathbf{I}_F^{-1}\mathbf{m}_F\mathbf{m}_F^\top\mathbf{I}_F^{-1}\gamma^2(\mathbf{I}_F^{-1}\mathbf{a}\mathbf{a}^\top\mathbf{I}_F^{-1} + \mathbf{I}_F^{-2}\mathbf{A}\mathbf{m}_F\mathbf{m}_F^\top\mathbf{I}_F^{-2}\mathbf{A}) + \circ(\gamma^2)$$

by plugging in the Taylor expansions of \mathbf{I}^{-1} and \mathbf{m} and retaining only terms up to the second order. Hence, the second term of the MSE without the trace is

$$\begin{aligned}& \left(\gamma^{-1}T^{-1} + \gamma^0 2T^{-1} + \gamma^1 T^{-1} \right) \mathbf{I}^{-1}\mathbf{m}\mathbf{m}^\top\mathbf{I}^{-1} \\ &= \gamma^{-1}\frac{1}{T}(\mathbf{I}_F^{-1}\mathbf{m}_F\mathbf{m}_F^\top\mathbf{I}_F^{-1}) + \gamma^0\frac{2}{T}(\mathbf{I}_F^{-1}\mathbf{m}_F\mathbf{m}_F^\top\mathbf{I}_F^{-1}) + \\ & \quad \gamma^1\frac{1}{T}(\mathbf{I}_F^{-1}\mathbf{m}_F\mathbf{m}_F^\top\mathbf{I}_F^{-1} + \mathbf{I}_F^{-1}\mathbf{a}\mathbf{a}^\top\mathbf{I}_F^{-1} + \mathbf{I}_F^{-2}\mathbf{A}\mathbf{m}_F\mathbf{m}_F^\top\mathbf{I}_F^{-2}\mathbf{A}) + \circ(\gamma)\end{aligned}$$

and the first term of the MSE without the trace is

$$\begin{aligned}& \left(\gamma^{-1}T^{-1} + \gamma^0 T^{-1} \right) (\mathbf{I}^{-1}) \\ &= \left(\gamma^{-1}T^{-1} + \gamma^0 T^{-1} \right) \left(\mathbf{I}_F^{-1} + \gamma\mathbf{I}_F^{-2}\mathbf{A} + \gamma^2(\mathbf{I}_F^{-1}(\mathbf{I}_F^{-1}\mathbf{A})^2 - \mathbf{I}_F^{-2}\mathbf{B}) + \circ(\gamma^2) \right) \\ &= \gamma^{-1}\frac{1}{T}\mathbf{I}_F^{-1} + \gamma^0\frac{1}{T}(\mathbf{I}_F^{-2}\mathbf{A} + \mathbf{I}_F^{-1}) + \gamma^1\frac{1}{T}[\mathbf{I}_F^{-1}(\mathbf{I}_F^{-1}\mathbf{A})^2 - \mathbf{I}_F^{-2}\mathbf{B} + \mathbf{I}_F^{-2}\mathbf{A}] + \circ(\gamma) .\end{aligned}$$

Subtracting the second term from the first term and applying the trace, we finally write the MSE:

$$\text{MSE}_{\text{NCE}} = \text{tr} \left(\gamma^{-1}\frac{1}{T}(\mathbf{I}_F^{-1} - \mathbf{I}_F^{-1}\mathbf{m}_F\mathbf{m}_F^\top\mathbf{I}_F^{-1}) + \gamma^0\frac{1}{T}(\mathbf{I}_F^{-2}\mathbf{A} + \mathbf{I}_F^{-1} - 2\mathbf{I}_F^{-1}\mathbf{m}_F\mathbf{m}_F^\top\mathbf{I}_F^{-1}) \right) + \circ(\gamma)\tag{7}$$

- Optimize the MSE_{NCE} w.r.t. p_n

To optimize w.r.t. p_n , we need only keep the two first orders of the MSE_{NCE} , which depends on p_n only via the term $\text{tr}(\mathbf{I}_F^{-2}\mathbf{A}) = \int \|\mathbf{I}_F^{-1}\mathbf{g}(\mathbf{x})\|_{\frac{p_d^2}{p_n}}^2 d\mathbf{x}$. Hence, we need to optimize

$$J(p_n) = \frac{1}{T} \int \|\mathbf{I}_F^{-1}\mathbf{g}(\mathbf{x})\|_{\frac{p_d^2}{p_n}}^2 d\mathbf{x}\tag{8}$$

with respect to p_n . We compute the variational (Fréchet) derivative together with the Lagrangian of the constraint $\int p_n(\mathbf{x}) = 1$ (with λ denoting the Lagrangian multiplier) to obtain

$$\delta_{p_n} J = -\|I_F^{-1} \mathbf{g}\|^2 \frac{p_d^2}{p_n^2} + \lambda . \quad (9)$$

Setting this to zero and taking into account the non-negativity of p_n gives

$$p_n(\mathbf{x}) = \|I_F^{-1} \mathbf{g}(x)\| p_d(\mathbf{x}) / Z \quad (10)$$

where $Z = \int \|I_F^{-1} \mathbf{g}(x)\| p_d(\mathbf{x}) d\mathbf{x}$ is the normalization constant. This is thus the optimal noise distribution, as a first-order approximation.

- Compute the MSE gaps at optimality

Plugging this optimal p_n into the formula of MSE_{NCE} and subtracting the Cramer-Rao MSE (which is a lower bound for a normalized model), we get:

$$\begin{aligned} \Delta_{\text{opt}} \text{MSE}_{\text{NCE}} &= \text{MSE}_{\text{NCE}}(p_n = p_n^{\text{opt}}) - \text{MSE}_{\text{Cramer-Rao}} \\ &= \frac{1}{T} \left(\int \|I_F^{-1} \boldsymbol{\psi}\| p_d \right)^2 . \end{aligned}$$

This is interesting to compare with the case where the noise distribution is the data distribution, which gives

$$\begin{aligned} \Delta_{\text{data}} \text{MSE}_{\text{NCE}} &= \text{MSE}_{\text{NCE}}(p_n = p_d) - \text{MSE}_{\text{Cramer-Rao}} \\ &= \frac{1}{T} \int \|I_F^{-1} \boldsymbol{\psi}\|^2 p_d \end{aligned}$$

where the squaring is in a different place. In fact, we can compare these two quantities by the Cauchy-Schwartz inequality, or simply the fact that

$$\begin{aligned} \Delta \text{MSE}_{\text{NCE}} &= \Delta_{\text{data}} \text{MSE}_{\text{NCE}} - \Delta_{\text{opt}} \text{MSE}_{\text{NCE}} \\ &= \text{MSE}_{\text{NCE}}(p_n = p_d) - \text{MSE}_{\text{NCE}}(p_n = p_n^{\text{opt}}) \\ &= \frac{1}{T} \text{Var}_{X \sim p_d} \{ \|I_F^{-1} \mathbf{g}(\mathbf{X})\| \} \end{aligned}$$

This implies that the two MSEs, when the noise distribution is either p_n^{opt} or p_d , can be equal only if $\|I_F^{-1} \mathbf{g}(\cdot)\|$ is constant in the support of p_d . This does not seem to be possible for any reasonable distribution.

Proof: case where $p_n \approx p_d$

We consider the limit case where $\frac{p_d}{p_n}(\mathbf{x}) = 1 + \epsilon(\mathbf{x})$ with $|\epsilon(\mathbf{x}) - 0| < \epsilon_{\max} \quad \forall \mathbf{x}$.

Note that in order to use Taylor expansions for terms containing $\epsilon(\mathbf{x})$ in an integral, we assume for any integrand $h(\mathbf{x})$ that $\int h(\mathbf{x})\epsilon(\mathbf{x})d\mathbf{x} \approx \epsilon \int h(\mathbf{x})d\mathbf{x}$, where ϵ would be a constant.

- Taylor expansion of the discriminator

$$\begin{aligned} 1 - D(\mathbf{x}) &= \frac{\nu p_n(\mathbf{x})}{p_d(\mathbf{x}) + \nu p_n(\mathbf{x})} = \frac{1}{1 + \frac{1}{\nu} + \frac{p_d}{p_n}(\mathbf{x})} = \frac{1}{1 + \frac{1}{\nu} + \frac{1}{\nu}\epsilon(\mathbf{x})} \\ &= \frac{\nu}{1 + \nu}\epsilon^0(\mathbf{x}) - \frac{\nu}{(1 + \nu)^2}\epsilon^1(\mathbf{x}) + \frac{\nu}{(1 + \nu)^3}\epsilon^2(\mathbf{x}) + o(\epsilon^2) \end{aligned}$$

- Evaluating the integrals \mathbf{m} , \mathbf{I}

$$\begin{aligned} \mathbf{m} &= \int \mathbf{g}(\mathbf{x})p_d(\mathbf{x})\left(1 - D(\mathbf{x})\right)d\mathbf{x} \\ &= \int \mathbf{g}(\mathbf{x})p_d(\mathbf{x})\left(\frac{\nu}{1 + \nu}\epsilon^0(\mathbf{x}) - \frac{\nu}{(1 + \nu)^2}\epsilon^1(\mathbf{x}) + \frac{\nu}{(1 + \nu)^3}\epsilon^2(\mathbf{x}) + o(\epsilon^2)\right)d\mathbf{x} \\ &= \frac{\nu}{1 + \nu}\mathbf{m}_F - \frac{\nu}{(1 + \nu)^2}\mathbf{a}(\epsilon) + \frac{\nu}{(1 + \nu)^3}\mathbf{b}(\epsilon^2) + o(\epsilon^3) \end{aligned}$$

where the Fisher-score mean \mathbf{m}_F is null and we use shorthand notations \mathbf{a} and \mathbf{b} for the remaining integrals:

$$\begin{aligned} \mathbf{m}_F &= \int \mathbf{g}(\mathbf{x})p_d(\mathbf{x})d\mathbf{x} \\ \mathbf{a}(\epsilon) &= \int \mathbf{g}(\mathbf{x})p_d\epsilon(\mathbf{x})d\mathbf{x} \\ \mathbf{b}(\epsilon^2) &= \int \mathbf{g}(\mathbf{x})p_d\epsilon^2(\mathbf{x})d\mathbf{x} . \end{aligned}$$

Similarly,

$$\begin{aligned} \mathbf{I} &= \int \mathbf{g}(\mathbf{x})\mathbf{g}(\mathbf{x})^\top p_d(\mathbf{x})\left(1 - D(\mathbf{x})\right)d\mathbf{x} \\ &= \int \mathbf{g}(\mathbf{x})\mathbf{g}(\mathbf{x})^\top p_d(\mathbf{x})\left(\frac{\nu}{1 + \nu}\epsilon^0(\mathbf{x}) - \frac{\nu}{(1 + \nu)^2}\epsilon^1(\mathbf{x}) + \frac{\nu}{(1 + \nu)^3}\epsilon^2(\mathbf{x}) + o(\epsilon^2)\right)d\mathbf{x} \\ &= \frac{\nu}{1 + \nu}\mathbf{I}_F - \frac{\nu}{(1 + \nu)^2}\mathbf{A}(\epsilon) + \frac{\nu}{(1 + \nu)^3}\mathbf{B}(\epsilon^2) + o(\epsilon^3) \end{aligned}$$

where the Fisher-score covariance (Fisher information) is \mathbf{I}_F and we use shorthand notations \mathbf{A} and \mathbf{B} for the remaining integrals:

$$\begin{aligned} \mathbf{I}_F &= \int \mathbf{g}(\mathbf{x})\mathbf{g}(\mathbf{x})^\top p_d(\mathbf{x})d\mathbf{x} \\ \mathbf{A}(\epsilon) &= \int \mathbf{g}(\mathbf{x})\mathbf{g}(\mathbf{x})^\top p_d\epsilon(\mathbf{x})d\mathbf{x} \\ \mathbf{B}(\epsilon^2) &= \int \mathbf{g}(\mathbf{x})\mathbf{g}(\mathbf{x})^\top p_d\epsilon^2(\mathbf{x})d\mathbf{x} . \end{aligned}$$

- Taylor expansion of \mathbf{I}^{-1}

$$\begin{aligned} \mathbf{I}^{-1} &= \left(\frac{\nu}{1 + \nu}\mathbf{I}_F - \frac{\nu}{(1 + \nu)^2}\mathbf{A}(\epsilon) + \frac{\nu}{(1 + \nu)^3}\mathbf{B}(\epsilon^2) + o(\epsilon^3)\right)^{-1} \\ &= \frac{1 + \nu}{\nu}\mathbf{I}_F^{-1} + \frac{1}{\nu}\mathbf{I}_F^{-2}\mathbf{A}(\epsilon) + \frac{\nu}{1 + \nu}\mathbf{I}_F^{-2}(\mathbf{I}_F^{-1}\mathbf{A}^2(\epsilon) - \mathbf{B}(\epsilon^2)) + o(\epsilon^3) \end{aligned}$$

- Evaluating the MSE_{NCE}

$$\mathbf{I}^{-1} \mathbf{m} \mathbf{m}^\top \mathbf{I}^{-1} = \mathbf{I}_F^{-1} \mathbf{m}_F \mathbf{m}_F^\top \mathbf{I}_F^{-1} + \frac{1}{(1+\nu)^2} \left(\mathbf{I}_F^{-2} \mathbf{A}(\epsilon) \mathbf{m}_F \mathbf{m}_F^\top \mathbf{I}_F^{-2} \mathbf{A}(\epsilon) + \mathbf{I}_F^{-1} \mathbf{a}(\epsilon) \mathbf{a}(\epsilon)^\top \mathbf{I}_F^{-1} \right) + \circ(\epsilon^3)$$

by plugging in the Taylor expansions of \mathbf{I}^{-1} and \mathbf{m} and retaining only terms up to the second order. Finally, the MSE becomes:

$$\begin{aligned} \text{MSE}_{\text{NCE}}(T, \nu, p_n) &= \frac{\nu+1}{T} \text{tr} \left(\mathbf{I}^{-1} - \frac{\nu+1}{\nu} (\mathbf{I}^{-1} \mathbf{m} \mathbf{m}^\top \mathbf{I}^{-1}) \right) \\ &= \text{tr} \left(\frac{(1+\nu)^2}{T\nu} (\mathbf{I}_F^{-1} - \mathbf{I}_F^{-1} \mathbf{m}_F \mathbf{m}_F^\top \mathbf{I}_F^{-1}) + \frac{1+\nu}{T\nu} \mathbf{I}_F^{-2} \mathbf{A}(\epsilon) + \right. \\ &\quad \left. \frac{1}{T\nu} (\mathbf{I}_F^{-3} \mathbf{A}^2(\epsilon) - \mathbf{I}_F^{-2} \mathbf{B}(\epsilon^2) - \mathbf{I}_F^{-1} \mathbf{a}(\epsilon) \mathbf{a}(\epsilon)^\top \mathbf{I}_F^{-1} - \mathbf{I}_F^{-2} \mathbf{A}(\epsilon) \mathbf{m}_F \mathbf{m}_F^\top \mathbf{I}_F^{-2} \mathbf{A}(\epsilon)) \right) + \circ(\epsilon^3) \end{aligned}$$

- Optimize the MSE_{NCE} w.r.t. p_n

To optimize w.r.t. p_n , we need only keep the MSE_{NCE} up to order 1, which depends on p_n only via the term

$$\text{tr}(\mathbf{I}_F^{-2} \mathbf{A}(\epsilon)) = \text{tr} \left(\mathbf{I}_F^{-2} \left(\int \mathbf{g}(\mathbf{x}) \mathbf{g}(\mathbf{x})^\top \frac{p_d^2}{p_n}(\mathbf{x}) d\mathbf{x} - \mathbf{I}_F \right) \right)$$

. where we unpacked p_n from $\epsilon = \frac{p_d}{p_n} - 1$. Hence, we need to optimize

$$J(p_n) = \frac{1}{T} \int \|\mathbf{I}_F^{-1} \mathbf{g}(\mathbf{x})\|^2 \frac{p_d^2}{p_n}(\mathbf{x}) d\mathbf{x} \quad (11)$$

with respect to p_n . This was already done in the all-noise limit $\nu \rightarrow \infty$ and yielded

$$p_n(\mathbf{x}) = \|\mathbf{I}_F^{-1} \mathbf{g}(\mathbf{x})\| p_d(\mathbf{x}) / Z \quad (12)$$

where $Z = \int \|\mathbf{I}_F^{-1} \mathbf{g}(\mathbf{x})\| p_d(\mathbf{x}) d\mathbf{x}$ is the normalization constant. This is thus the optimal noise distribution, as a first-order approximation.

In the third case, the limit of all data, we have the following conjecture:

Conjecture 1 *In case (iii), the limit of all data samples $\nu \rightarrow 0$, the optimal noise distribution is such that it is all concentrated at the set of those ξ which are given by*

$$\begin{aligned} \arg \max_{\xi} p_d(\xi) \text{tr} \left((\mathbf{g}(\xi) \mathbf{g}(\xi)^\top)^{-1} \right)^{-1} \\ \text{s.t. } \mathbf{g}(\xi) = \text{constant} \end{aligned} \quad (13)$$

Informal and heuristic “proof”:

We have the $\text{MSE}_{\text{NCE}}(T, \nu, p_n) = \frac{\nu+1}{T} \text{tr}(\mathbf{I}^{-1} - \frac{\nu+1}{\nu} (\mathbf{I}^{-1} \mathbf{m} \mathbf{m}^\top \mathbf{I}^{-1}))$.

Given the term up until ν^{-1} in the MSE, we will use Taylor expansions up to order 2 throughout the proof, in anticipation that the MSE will be expanded until order 1.

Note that in this no noise limit, the assumption made by Gutmann and Hyvärinen (2012) that p_n is non-zero whenever p_d is nonzero is not true for this optimal p_n , which reduces the rigour of this analysis. (This we denote by heuristic approximation 1.)

- Taylor expansion of the discriminator

$$1 - D(\mathbf{x}) = \frac{\nu p_n(\mathbf{x})}{p_d(\mathbf{x}) + \nu p_n(\mathbf{x})} = \frac{1}{1 + \frac{1}{\nu} \frac{p_d}{p_n}(\mathbf{x})} = \nu \frac{p_n}{p_d}(\mathbf{x}) - \nu^2 \frac{p_n^2}{p_d^2}(\mathbf{x}) + o(\nu^2)$$

- Evaluating the integrals \mathbf{m} , \mathbf{I}

$$\begin{aligned} \mathbf{m} &= \int \mathbf{g}(\mathbf{x}) p_d(\mathbf{x}) (1 - D(\mathbf{x})) d\mathbf{x} = \int \mathbf{g}(\mathbf{x}) p_d(\mathbf{x}) \left(\nu \frac{p_n}{p_d}(\mathbf{x}) - \nu^2 \frac{p_n^2}{p_d^2}(\mathbf{x}) + o(\nu^2) \right) d\mathbf{x} \\ &= \nu \mathbf{m}_n - \nu^2 \mathbf{b} + o(\nu^2) \end{aligned}$$

where

$$\begin{aligned} \mathbf{m}_n &= \int \mathbf{g}(\mathbf{x}) p_n(\mathbf{x}) d\mathbf{x} \\ \mathbf{b} &= \int \mathbf{g}(\mathbf{x}) \frac{p_n^2}{p_d}(\mathbf{x}) d\mathbf{x} . \end{aligned}$$

Similarly,

$$\begin{aligned} \mathbf{I} &= \int \mathbf{g}(\mathbf{x}) \mathbf{g}(\mathbf{x})^\top p_d(\mathbf{x}) (1 - D(\mathbf{x})) d\mathbf{x} = \int \mathbf{g}(\mathbf{x}) \mathbf{g}(\mathbf{x})^\top p_d(\mathbf{x}) \left(\nu \frac{p_n}{p_d}(\mathbf{x}) - \nu^2 \frac{p_n^2}{p_d^2}(\mathbf{x}) + o(\nu^2) \right) d\mathbf{x} \\ &= \nu \mathbf{I}_n - \nu^2 \mathbf{B} + o(\nu^2) \end{aligned}$$

where the Fisher-score covariance (Fisher information) is \mathbf{I}_F and we use shorthand notations A and B for the remaining integrals:

$$\begin{aligned} \mathbf{I}_n &= \int \mathbf{g}(\mathbf{x}) \mathbf{g}(\mathbf{x})^\top p_n(\mathbf{x}) d\mathbf{x} \\ \mathbf{B} &= \int \mathbf{g}(\mathbf{x}) \mathbf{g}(\mathbf{x})^\top \frac{p_n^2}{p_d}(\mathbf{x}) d\mathbf{x} . \end{aligned}$$

- Taylor expansion of \mathbf{I}^{-1}

$$\begin{aligned}
\mathbf{I}^{-1} &= \left(\nu \mathbf{I}_n - \nu^2 \mathbf{B} + \circ(\nu^2) \right)^{-1} \\
&= \left(\nu \mathbf{I}_n (\mathbf{Id} - \nu \mathbf{I}_n^{-1} \mathbf{B}) + \circ(\nu^2) \right)^{-1} \\
&= \nu^{-1} \mathbf{I}_n^{-1} \left(\mathbf{Id} + \nu \mathbf{I}_n^{-1} \mathbf{B} + \nu^2 (\mathbf{I}_n^{-1} \mathbf{B})^2 + \nu^3 (\mathbf{I}_n^{-1} \mathbf{B})^3 + \circ(\nu^3) \right) + \circ(\nu^2) \\
&= \nu^{-1} \mathbf{I}_n^{-1} + \nu^0 \mathbf{I}_n^{-2} \mathbf{B} + \nu^1 \mathbf{I}_n^{-1} (\mathbf{I}_n^{-2} \mathbf{B})^2 + \nu^2 \mathbf{I}_n^{-1} (\mathbf{I}_n^{-2} \mathbf{B})^3 + \circ(\nu^2)
\end{aligned}$$

- Evaluating the MSE_{NCE}

$$\begin{aligned}
\mathbf{I}^{-1} \mathbf{m} \mathbf{m}^\top \mathbf{I}^{-1} &= \\
&\nu^0 (\mathbf{I}_n^{-1} \mathbf{m}_n \mathbf{m}_n^\top \mathbf{I}_n^{-1}) + \nu^2 (\mathbf{I}_n^{-1} \mathbf{b} \mathbf{b}^\top \mathbf{I}_n^{-1} + \mathbf{I}_n^{-2} \mathbf{B} \mathbf{m}_n \mathbf{m}_n^\top \mathbf{I}_n^{-2} \mathbf{B}) + \circ(\nu^2)
\end{aligned}$$

by plugging in the Taylor expansions of \mathbf{I}^{-1} and \mathbf{m} and retaining only terms up to the second order. Hence, the second term of the MSE without the trace is

$$\begin{aligned}
&\left(\nu^1 T^{-1} + \nu^0 2T^{-1} + \nu^{-1} T^{-1} \right) \mathbf{I}^{-1} \mathbf{m} \mathbf{m}^\top \mathbf{I}^{-1} \\
&= \left(\nu^1 T^{-1} + \nu^0 2T^{-1} + \nu^{-1} T^{-1} \right) \left(\nu^0 (\mathbf{I}_n^{-1} \mathbf{m}_n \mathbf{m}_n^\top \mathbf{I}_n^{-1}) + \nu^2 (\mathbf{I}_n^{-1} \mathbf{b} \mathbf{b}^\top \mathbf{I}_n^{-1} + \mathbf{I}_n^{-2} \mathbf{B} \mathbf{m}_n \mathbf{m}_n^\top \mathbf{I}_n^{-2} \mathbf{B}) + \circ(\nu^2) \right) \\
&= \nu^{-1} \frac{1}{T} (\mathbf{I}_n^{-1} \mathbf{m}_n \mathbf{m}_n^\top \mathbf{I}_n^{-1}) + \nu^0 \frac{1}{T} (2 \mathbf{I}_n^{-1} \mathbf{m}_n \mathbf{m}_n^\top \mathbf{I}_n^{-1}) + \\
&\nu^1 \frac{1}{T} (\mathbf{I}_n^{-1} \mathbf{b}_n \mathbf{b}_n^\top \mathbf{I}_n^{-1} + \mathbf{I}_n^{-2} \mathbf{B} \mathbf{m}_n \mathbf{m}_n^\top \mathbf{I}_n^{-2} \mathbf{B} + \mathbf{I}_n^{-1} \mathbf{m}_n \mathbf{m}_n^\top \mathbf{I}_n^{-1}) + \circ(\nu)
\end{aligned}$$

and the first term of the MSE without the trace is

$$\begin{aligned}
&\left(\nu^0 T^{-1} + \nu^1 T^{-1} \right) \text{tr}(\mathbf{I}^{-1}) \\
&= \left(\nu^0 T^{-1} + \nu^1 T^{-1} \right) \left(\nu^{-1} \mathbf{I}_n^{-1} + \nu^0 \mathbf{I}_n^{-2} \mathbf{B} + \nu^1 \mathbf{I}_n^{-1} (\mathbf{I}_n^{-2} \mathbf{B})^2 + \nu^2 \mathbf{I}_n^{-1} (\mathbf{I}_n^{-2} \mathbf{B})^3 + \circ(\nu^2) \right) \\
&= \nu^{-1} \frac{1}{T} \mathbf{I}_n^{-1} + \nu^0 \frac{1}{T} (\mathbf{I}_n^{-2} \mathbf{B} + \mathbf{I}_n^{-1}) + \nu^1 \frac{1}{T} [\mathbf{I}_n^{-1} (\mathbf{I}_n^{-1} \mathbf{B})^2 + \mathbf{I}_n^{-2} \mathbf{B}] + \circ(\nu) .
\end{aligned}$$

Subtracting the second term from the first term and applying the trace, we finally write the MSE:

$$\begin{aligned}
\text{MSE}_{\text{NCE}} &= \text{tr} \left(\nu^{-1} \frac{1}{T} (\mathbf{I}_n^{-1} - \mathbf{I}_n^{-1} \mathbf{m}_n \mathbf{m}_n^\top \mathbf{I}_n^{-1}) + \nu^0 \frac{1}{T} (\mathbf{I}_n^{-2} \mathbf{B} + \mathbf{I}_n^{-1} - 2 \mathbf{I}_n^{-1} \mathbf{m}_n \mathbf{m}_n^\top \mathbf{I}_n^{-1}) + \right. \\
&\left. \nu^1 \frac{1}{T} [\mathbf{I}_n^{-1} (\mathbf{I}_n^{-1} \mathbf{B})^2 + \mathbf{I}_n^{-2} \mathbf{B} - \mathbf{I}_n^{-1} \mathbf{b}_n \mathbf{b}_n^\top \mathbf{I}_n^{-1} - \mathbf{I}_n^{-2} \mathbf{B} \mathbf{m}_n \mathbf{m}_n^\top \mathbf{I}_n^{-2} \mathbf{B} - \mathbf{I}_n^{-1} \mathbf{m}_n \mathbf{m}_n^\top \mathbf{I}_n^{-1}] + \circ(\nu) \right) .
\end{aligned}$$

Rewriting $\mathbf{I}_n^{-1} = \mathbf{I}_n^{-1} \mathbf{I}_n \mathbf{I}_n^{-1}$, using the circular invariance of the trace operator and stopping at order ν^0 , we get:

$$\begin{aligned}
\text{MSE}_{\text{NCE}} &= \nu^{-1} \frac{1}{T} \langle \mathbf{I}_n^{-2}, \mathbf{I}_n - \mathbf{m}_n \mathbf{m}_n^\top \rangle + \nu^0 \frac{1}{T} \langle \mathbf{I}_n^{-2}, \mathbf{B} + \mathbf{I}_n - 2 \mathbf{m}_n \mathbf{m}_n^\top \rangle + \circ(1) \\
&= \nu^{-1} \frac{1}{T} \langle \mathbf{I}_n^{-2}, \text{Var}_{N \sim p_n} \mathbf{g}(\mathbf{N}) \rangle + \nu^0 \frac{1}{T} \langle \mathbf{I}_n^{-2}, \mathbf{B} + \mathbf{I}_n - 2 \mathbf{m}_n \mathbf{m}_n^\top \rangle + \circ(1) . \tag{14}
\end{aligned}$$

- Optimize the MSE_{NCE} w.r.t. p_n

Looking at the above MSE, the dominant term of order ν^{-1} is $\langle \mathbf{I}_n^{-2}, \text{Var}_{N \sim p_n} \mathbf{g}(\mathbf{N}) \rangle \geq 0$ is minimized when it is 0, that is, when \mathbf{g} is constant in the support of p_n . Typically this means that p_n is concentrated on a set of zero measure. In the 1D case, such case is typically the Dirac delta $p_n = \delta_z$, or a distribution with two deltas in case of symmetrical \mathbf{g} .

We can plug this in the terms of the next order ν^0 , which remain to be minimized:

$$\begin{aligned}\langle \mathbf{I}_n^{-2}, \mathbf{B} + \mathbf{I}_n - 2\mathbf{m}_n\mathbf{m}_n^T \rangle &= \langle \mathbf{I}_n^{-2}, \mathbf{B} - \mathbf{I}_n + 2\mathbf{I}_n - 2\mathbf{m}_n\mathbf{m}_n^T \rangle \\ &= \langle \mathbf{I}_n^{-2}, \mathbf{B} - \mathbf{I}_n + 2\text{Var}_{N \sim p_n} \mathbf{g}(\mathbf{N}) \rangle \\ &= \langle \mathbf{I}_n^{-2}, \mathbf{B} - \mathbf{I}_n \rangle\end{aligned}$$

given we chose p_n so that the variance is 0.

The integrands of \mathbf{B} and \mathbf{I} respectively involve p_n^2 and p_n . Because p_n is concentrated on a set of zero measure (Dirac-like), the term in \mathbf{B} dominates the term in \mathbf{I} . This is because if we consider the p_n as the limit of a sequence of some proper pdf's, the value of the pdf gets infinite in the support of that pdf in the limit, and thus p_n^2 is infinitely larger than p_n . Hence we are left with $\langle \mathbf{I}_n^{-2}, \mathbf{B} \rangle$.

The integral with respect to p_n simplifies to simply evaluating the $\mathbf{g}(\mathbf{x})\mathbf{g}(\mathbf{x})^\top / p_d(\mathbf{x})$ the support of p_n . Since we know that $\mathbf{g}(\mathbf{x})$ is constant in that set, the main question is whether p_d is constant in that set as well. Here, we heuristically assume that it is; this is intuitively appealing in many cases, if not necessarily true. (This we denote by heuristic approximation 2.)

Thus, we have

$$\int \mathbf{g}(\mathbf{x})\mathbf{g}(\mathbf{x})^\top \frac{\delta_z^2(\mathbf{x})}{p_d} d\mathbf{x} \approx c \mathbf{g}(\mathbf{z})\mathbf{g}(\mathbf{z})^\top \frac{1}{p_d(\mathbf{z})}$$

for some constant c taking into account the effect of squaring of p_n (it is ultimately infinite, but the reasoning is still valid in any sequence going to the limit.)

Next we make heuristic approximation 3: we neglect any problems of inversion of singular, rank 1 matrices (note this is not a problem in the 1D case), and further obtain

$$\langle \mathbf{I}_n^{-2}, \mathbf{B} \rangle \approx \text{tr} \left((\mathbf{g}(\mathbf{z})\mathbf{g}(\mathbf{z})^\top)^{-1} \mathbf{g}(\mathbf{z})\mathbf{g}(\mathbf{z})^\top \frac{1}{p_d(\mathbf{z})} (\mathbf{g}(\mathbf{z})\mathbf{g}(\mathbf{z})^\top)^{-1} \right) \approx \frac{1}{p_d(\mathbf{z})} \text{tr} \left((\mathbf{g}(\mathbf{z})\mathbf{g}(\mathbf{z})^\top)^{-1} \right). \quad (15)$$

Minimizing this term is equivalent to the following maximization setup (still applying heuristic approximation 3):

$$\arg \max_{\xi} p_d(\xi) \text{tr} \left((\mathbf{g}(\xi)\mathbf{g}(\xi)^\top)^{-1} \right)^{-1}.$$

Those points z obtained by the above condition are the best candidates for p_n to concentrate its mass on.

We arrived this result by making three heuristic approximations as explained above; we hope to be able to remove some of them in future work.

Numerically, evaluating the optimal noise in the all-data limit requires computing a weight $w(\mathbf{x}) = \text{tr} \left((\mathbf{g}(\xi)\mathbf{g}(\xi)^\top)^{-1} \right)^{-1}$ that is intractable in dimensions bigger than 1, due to the singularity of the rank 1 matrix. We can avoid this numerically by introducing an (infinitesimal) perturbation $\epsilon > 0$ which removes the singularity problem:

$$\begin{aligned}w_\epsilon(\xi) &= \text{tr} \left((\mathbf{g}(\xi)\mathbf{g}(\xi)^\top + \epsilon \text{Id})^{-1} \right)^{-1} \\ &= \text{tr} \left(\epsilon^{-1} \text{Id} - \frac{1}{\epsilon^2 + \epsilon \mathbf{g}(\xi)^\top \text{Id} \mathbf{g}(\xi)} \mathbf{g}(\xi)\mathbf{g}(\xi)^\top \right)^{-1} \quad \text{by the Sherman-Morrison formula} \\ &= \left(\epsilon^{-1} d - \frac{1}{\epsilon^2 + \epsilon \|\mathbf{g}(\xi)\|^2} \|\mathbf{g}(\xi)\|^2 \right)^{-1} \\ &= \left(\epsilon^{-1} (d-1) + \epsilon^0 \frac{1}{\|\mathbf{g}(\xi)\|^2} + \epsilon^1 \frac{-1}{\|\mathbf{g}(\xi)\|^4} + O(\epsilon^2) \right)^{-1} \quad \text{by Taylor expansion} \\ &= \epsilon \frac{1}{d-1} + \epsilon^2 \frac{-1}{\|\mathbf{g}(\xi)\|^2 (d-1)^2} + \epsilon^3 \frac{(2-d)}{\|\mathbf{g}(\xi)\|^4 (d-1)^3} + O(\epsilon^4) \quad \text{by further Taylor expansion}\end{aligned}$$

where we go up to order 3 to ensure the weight $w_\epsilon(\xi)$ is positive. Finally, we can approximate the arg max operator with its relaxation $\text{soft arg max}^\epsilon(x) = \frac{e^{\frac{x}{\epsilon}}}{\int e^{\frac{x}{\epsilon}} dx}$, so that

$$p_n(\mathbf{x}) \approx \text{soft arg max}^{\epsilon_1} (p_d(\mathbf{x})w_{\epsilon_2}(\mathbf{x}))$$

where $(\epsilon_1, \epsilon_2) \in (\mathbb{R}_+^*)^2$ are two hyperparameters taken close to zero.

5 OPTIMAL NOISE FOR ESTIMATING A DISTRIBUTION: PROOFS

So far, we have optimized hyperparameters (such as the noise distribution) so that the reduce the uncertainty of the *parameter* estimation, measured by the Mean Squared Error $\mathbb{E}[\|\hat{\theta}_T - \theta^*\|^2] = \frac{1}{T_d} \text{tr}(\Sigma)$.

Sometimes, we might wish to reduce the uncertainty of the *distribution* estimation, which we can measure using the Kullback-Leibler (KL) divergence $\mathbb{E}[\mathcal{D}_{\text{KL}}(p_d, p_{\hat{\theta}_T})]$.

We can specify this error, by using the Taylor expansion of the estimated $\hat{\theta}_T$ near optimality, given in Gutmann and Hyvärinen [2012]:

$$\hat{\theta}_T - \theta^* = z + O(\|\hat{\theta}_T - \theta^*\|^2) \quad (16)$$

where $z \sim \mathcal{N}(0, \frac{1}{T_d} \Sigma)$ and Σ is the asymptotic variance matrix.

We can similarly take the Taylor expansion of the KL divergence with respect to its second argument, near optimality:

$$\begin{aligned} J(\hat{\theta}_T) &:= \mathcal{D}_{\text{KL}}(p_d, p_{\hat{\theta}_T}) \\ &= J(\theta^*) + \langle \nabla_{\theta} J(\theta^*), \hat{\theta}_T - \theta^* \rangle + \frac{1}{2} \langle (\hat{\theta}_T - \theta^*), \nabla_{\theta}^2 J(\theta^*) (\hat{\theta}_T - \theta^*) \rangle + O(\|\hat{\theta}_T - \theta^*\|^3) \\ &= J(\theta^*) + \langle \nabla_{\theta} J(\theta^*), \hat{\theta}_T - \theta^* \rangle + \frac{1}{2} \|\hat{\theta}_T - \theta^*\|_{\nabla_{\theta}^2 J(\theta^*)}^2 + O(\|\hat{\theta}_T - \theta^*\|^3) \end{aligned}$$

Note that some simplifications occur:

- $J(\theta^*) = \mathcal{D}_{\text{KL}}(p_{\theta^*}, p_{\theta^*}) = 0$
- $\nabla_{\theta} J(\theta^*) = 0$ as the gradient the KL divergence at θ^* is the mean of the (negative) Fisher score, which is null.
- $\nabla_{\theta}^2 J(\theta^*) = \mathbf{I}_F$

Plugging in the estimation error 16 into the distribution error yields:

$$\begin{aligned} J(\hat{\theta}_T) &= \frac{1}{2} \left\| z + O(\|\hat{\theta}_T - \theta^*\|^2) \right\|_{\mathbf{I}_F}^2 + O(\|\hat{\theta}_T - \theta^*\|^3) \\ &= \frac{1}{2} \left(\|z\|_{\mathbf{I}_F}^2 + 2 \langle z, O(\|\hat{\theta}_T - \theta^*\|^2) \rangle_{\mathbf{I}_F} + \|O(\|\hat{\theta}_T - \theta^*\|^2)\|_{\mathbf{I}_F}^2 \right) + O(\|\hat{\theta}_T - \theta^*\|^3) \\ &= \frac{1}{2} \|z\|_{\mathbf{I}_F}^2 + O(\|\hat{\theta}_T - \theta^*\|^2) \end{aligned}$$

by truncating the Taylor expansion to the first order. Hence up to the first order, the expectation yields:

$$\begin{aligned} \mathbb{E}[\mathcal{D}_{\text{KL}}(p_d, p_{\hat{\theta}_T})] &= \frac{1}{2} \mathbb{E}[\|z\|_{\mathbf{I}_F}^2] = \frac{1}{2} \mathbb{E}[z^T \mathbf{I}_F z] = \frac{1}{2} \mathbb{E}[\text{tr}(z^T \mathbf{I}_F z)] = \frac{1}{2} \mathbb{E}[\text{tr}(\mathbf{I}_F z z^T)] \\ &= \frac{1}{2} \text{tr}(\mathbf{I}_F \mathbb{E}[z z^T]) = \frac{1}{2} \text{tr}(\mathbf{I}_F \text{Var}[z]) = \frac{1}{2T_d} \text{tr}(\mathbf{I}_F \Sigma) \end{aligned}$$

Note that this is a general and known result which is applicable beyond the KL divergence: for any divergence, the 0th order term is null as it measures the divergence between the data distribution and itself, the 1st order term is null in expectation if the estimator $\hat{\theta}_T$ is asymptotically unbiased, which leaves an expected error given by the 2nd-order term $\frac{1}{2T_d} \text{tr}(\nabla^2 J \Sigma)$ where J is the chosen divergence. Essentially, one would replace the Fisher Information above, which is the Hessian for a forward-KL divergence, by the Hessian for a given divergence.

Finding the optimal noise that minimizes the distribution error means minimizing $\frac{1}{T_d} \text{tr}(\Sigma \mathbf{I}_F)$. Contrast that with the optimal noise that minimizes the parameter estimation error (asymptotic variance) $\frac{1}{T_d} \text{tr}(\Sigma)$. We can reprise each of the three limit cases from the previous proofs, and derive novel optimal noise distributions:

Theorem 2 *In the two limit cases of Theorem 1, the noise distribution minimizing the expected Kullback-Leibler divergence is given by*

$$p_n^{\text{opt}}(\mathbf{x}) \propto p_d(\mathbf{x}) \|\mathbf{I}_F^{-\frac{1}{2}} \mathbf{g}(\mathbf{x})\| . \quad (17)$$

Proof: case of $\nu \rightarrow \infty$

We recall the asymptotic variance $\frac{1}{T_d} \Sigma$ in the all-noise limit is given by equation 7 at the first order and without the trace. Multiplying by I_F introduces no additional dependency in p_n , hence we retain the only term dependent that was dependent on p_n , $\mathbf{I}_F^{-2} \int \mathbf{g}(\mathbf{x}) \mathbf{g}(\mathbf{x})^\top \frac{p_d^2(\mathbf{x})}{p_n} d\mathbf{x}$, multiply it with I_F and take the trace. This yields the following cost to minimize:

$$J(p_n) = \frac{1}{T} \int \|\mathbf{I}_F^{-\frac{1}{2}} \mathbf{g}(\mathbf{x})\|^2 \frac{p_d^2(\mathbf{x})}{p_n} d\mathbf{x} \quad (18)$$

with respect to p_n . As in previous proofs, we compute the variational (Fréchet) derivative together with the Lagrangian of the constraint $\int p_n(\mathbf{x}) = 1$ (with λ denoting the Lagrangian multiplier) to obtain

$$\delta_{p_n} J = -\|\mathbf{I}_F^{-\frac{1}{2}} \mathbf{g}\|^2 \frac{p_d^2}{p_n^2} + \lambda . \quad (19)$$

Setting this to zero and taking into account the non-negativity of p_n gives

$$p_n(\mathbf{x}) = \|\mathbf{I}_F^{-\frac{1}{2}} \mathbf{g}(\mathbf{x})\| p_d(\mathbf{x}) / Z \quad (20)$$

where $Z = \int \|\mathbf{I}_F^{-\frac{1}{2}} \mathbf{g}(\mathbf{x})\| p_d(\mathbf{x}) d\mathbf{x}$ is the normalization constant. This is thus the optimal noise distribution, as a first-order approximation.

In the third case, the limit of all data, we have the following conjecture:

Conjecture 2 *In the limit of Conjecture 1 the noise distribution minimizing the expected Kullback-Leibler divergence is such that it is all concentrated at the set of those ξ which are given by*

$$\begin{aligned} \arg \max_{\xi} p_d(\xi) \text{tr} \left((\mathbf{g}(\xi) \mathbf{g}(\xi)^\top)^{-\frac{1}{2}} \right)^{-1} \\ \text{s.t. } \mathbf{g}(\xi) = \text{constant} \end{aligned} \quad (21)$$

Proof: case of $\nu \rightarrow 0$

By the same considerations, we can obtain the optimal noise that minimizes the asymptotic error in distribution space in the all-data limit, using equation 15 with a multiplication by I_F inside the the trace. This leads to the result.

6 NUMERICAL VALIDATION OF THE PREDICTED DISTRIBUTION ERROR

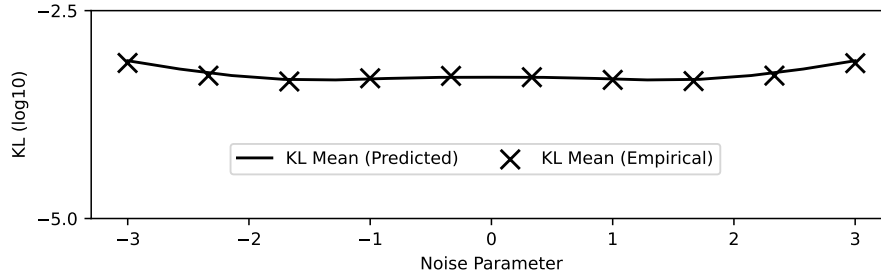


Figure 2: KL vs. the noise parameter (Gaussian Mean). The noise proportion is fixed at 50%.

We here numerically validate our formulae predicting the asymptotic estimation error in distribution space $\mathcal{D}_{\text{KL}}(p_d, p_{\hat{\theta}_{\text{NCE}}})$, when the noise is constrained within a parametric family containing the data; here, the model is a one-dimensional centered Gaussian with unit variance, parameterized by its mean.

References

- Michael Gutmann and Aapo Hyvärinen. Noise-contrastive estimation of unnormalized statistical models, with applications to natural image statistics. *Journal of Machine Learning Research*, 13:307–361, 2012.
- Miika Pihlaja, Michael Gutmann, and Aapo Hyvärinen. A family of computationally efficient and simple estimators for unnormalized statistical models. In *UAI*, 2010.