

Spin models and ML

1 Intro

Bla

2 Nishimori identity also known as Bayes rule

The main contribution of this section is the introduction of a very convenient notation that will be used throughout these notes.

Definition 2.1

Let (\mathbf{X}, \mathbf{Y}) be a couple of random variables on a polish space. We denote $\langle \cdot \rangle$ the expectation with respect to $\mathbb{P}(\mathbf{X} = \cdot | \mathbf{Y})$ and \mathbb{E} the expectation with respect to (\mathbf{X}, \mathbf{Y}) .

A replica denoted by \mathbf{x} is a random variable sampled (given \mathbf{Y}) from the distribution $\mathbb{P}(\mathbf{X} = \cdot | \mathbf{Y})$, independently of every other random variables. For $k \geq 1$, we denote k replicas as $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}$ being k i.i.d. samples (given \mathbf{Y}) from the distribution $\mathbb{P}(\mathbf{X} = \cdot | \mathbf{Y})$, independently of every other random variables.

Formally, we defined random variables $(\mathbf{X}, \mathbf{Y}, \mathbf{x}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)})$ with distribution $P_{(\mathbf{X}, \mathbf{Y})} \otimes P_{(\mathbf{X} | \mathbf{Y})} \otimes \dots \otimes P_{(\mathbf{X} | \mathbf{Y})}$ and a notation $\langle \cdot \rangle$ such that for all continuous bounded function f :

$$\langle f(\mathbf{X}, \mathbf{Y}, \mathbf{x}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}) \rangle \stackrel{\text{def}}{=} \mathbb{E} \left[f(\mathbf{X}, \mathbf{Y}, \mathbf{x}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}) | \mathbf{Y} \right].$$

The basic properties of conditional expectation gives:

$$\begin{aligned} \langle \mathbf{X} \rangle &= \mathbb{E}[\mathbf{X} | \mathbf{Y}] \\ \langle f(\mathbf{X})g(\mathbf{Y}) \rangle &= \langle f(\mathbf{X}) \rangle g(\mathbf{Y}) \\ \mathbb{E} \langle \mathbf{X} \rangle &= \mathbb{E}[\mathbf{X}]. \end{aligned}$$

Moreover, the independence assumption ensures that

$$\langle f(\mathbf{X})g(\mathbf{x}) \rangle = \langle f(\mathbf{X}) \rangle \langle g(\mathbf{x}) \rangle.$$

Thanks to Bayes rule, we have $P_{(\mathbf{X}, \mathbf{Y})} = P_{\mathbf{Y}} \otimes P_{(\mathbf{X} | \mathbf{Y})}$ so that (\mathbf{X}, \mathbf{Y}) has the same law as (\mathbf{x}, \mathbf{Y}) and

$$\langle \mathbf{x} \rangle = \mathbb{E}[\mathbf{x} | \mathbf{Y}] = \langle \mathbf{X} \rangle$$

This argument extends easily and $(\mathbf{Y}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)})$ has the same law as $(\mathbf{Y}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k-1)}, \mathbf{X})$, so that we get the following convenient reformulation:

Proposition 2.1 (Nishimori identity)

With the notations introduced in Definition 2.1, we have for all continuous bounded function f

$$\langle f(\mathbf{Y}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}) \rangle = \langle f(\mathbf{Y}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k-1)}, \mathbf{X}) \rangle.$$

We will now use the Nishimori identity in the particular case of two random vectors \mathbf{X} and \mathbf{Y} that live respectively in \mathbb{R}^n and \mathbb{R}^m . In this setting, we have useful relations:

$$\langle \mathbf{X}^T \mathbf{x} \rangle = \langle \mathbf{X} \rangle^T \langle \mathbf{x} \rangle = \|\langle \mathbf{x} \rangle\|^2,$$

where the first equality follows from the independence assumption and the second one from Nishimori identity. Moreover, we have

$$\mathbb{E} [\mathbf{X}^T \langle \mathbf{x} \rangle] = \mathbb{E} [\langle \mathbf{X}^T \mathbf{x} \rangle] = \mathbb{E} [\|\langle \mathbf{x} \rangle\|^2],$$

so that we obtain the classical result:

Proposition 2.2

$\langle \mathbf{x} \rangle$ is the orthogonal projection of \mathbf{X} on $L^2(\mathbf{Y})$:

$$\mathbb{E} [(\mathbf{X} - \langle \mathbf{x} \rangle)^T \langle \mathbf{x} \rangle] = 0.$$

As a direct consequence, we get another classical result about the problem of estimating \mathbf{X} when observing \mathbf{Y} where the performance of an estimator $\hat{\theta}$ (i.e. a measurable function of the observations \mathbf{Y}) is given by its Mean Squared Error $\text{MSE}(\hat{\theta}) = \mathbb{E} \|\mathbf{X} - \hat{\theta}(\mathbf{Y})\|^2$.

Definition 2.2

For a signal \mathbf{X} and observations \mathbf{Y} , the Minimum Mean Squared Error is defined by

$$\text{MMSE} \stackrel{\text{def}}{=} \min_{\hat{\theta}} \text{MSE}(\hat{\theta}) = \mathbb{E} [\|\mathbf{X} - \langle \mathbf{x} \rangle\|^2] = \mathbb{E} [\|\mathbf{X}\|^2] - \mathbb{E} [\langle \mathbf{x}^T \mathbf{X} \rangle],$$

where the minimum is taken over all measurable function $\hat{\theta}$ of the observations \mathbf{Y} . The optimal estimator (in term of Mean Squared Error) is the posterior mean of the signal \mathbf{X} given \mathbf{Y} .

Proof. Using $\langle \mathbf{X}^T \hat{\theta}(\mathbf{Y}) \rangle = \langle \mathbf{X} \rangle^T \hat{\theta}(\mathbf{Y})$, we have

$$\mathbb{E} [\|\mathbf{X} - \hat{\theta}(\mathbf{Y})\|^2] = \mathbb{E} [\|\mathbf{X}\|^2] - 2\mathbb{E} [\langle \mathbf{X} \rangle^T \hat{\theta}(\mathbf{Y})] + \mathbb{E} [\|\hat{\theta}(\mathbf{Y})\|^2]$$

and

$$\mathbb{E} [\|\langle \mathbf{X} \rangle - \hat{\theta}(\mathbf{Y})\|^2] = \mathbb{E} [\|\langle \mathbf{X} \rangle\|^2] - 2\mathbb{E} [\langle \mathbf{X} \rangle^T \hat{\theta}(\mathbf{Y})] + \mathbb{E} [\|\hat{\theta}(\mathbf{Y})\|^2] \geq 0,$$

we deduce

$$\begin{aligned} \mathbb{E} [\|\mathbf{X} - \hat{\theta}(\mathbf{Y})\|^2] &\geq \mathbb{E} [\|\mathbf{X}\|^2] - \mathbb{E} [\|\langle \mathbf{X} \rangle\|^2] \\ &= \mathbb{E} [\|\mathbf{X}\|^2] - 2\mathbb{E} [\langle \mathbf{X} \rangle^T \langle \mathbf{X} \rangle] + \mathbb{E} [\|\langle \mathbf{X} \rangle\|^2] \\ &= \mathbb{E} [\|\mathbf{X} - \langle \mathbf{x} \rangle\|^2] = \mathbb{E} [\|\mathbf{X} - \langle \mathbf{x} \rangle\|^2], \end{aligned}$$

and the last equality follows from Proposition 2.2. \square

Example 2.1. *In some particular cases, explicit computations can be done. The additive Gaussian scalar channel with Gaussian prior is a very simple case. We define*

$$Y = \sqrt{\lambda}X + Z, \quad (2.1)$$

where $Z \sim \mathcal{N}(0, 1)$ and X is sampled from a distribution $P_0 = \mathcal{N}(0, 1)$ over \mathbb{R} , independently of Z . We can compute the densities as $p(x) = \mathcal{N}(x|0, 1)$, $p(y|x) = \mathcal{N}(y|\sqrt{\lambda}x, 1)$ so that $p(y) = \mathcal{N}(y|0, 1 + \lambda)$ and the law of the replica is also Gaussian with $p(x|y) = \mathcal{N}\left(x|\frac{\sqrt{\lambda}}{1+\lambda}y, \frac{1}{1+\lambda}\right)$. In particular, we have

$$\langle x \rangle = \frac{\sqrt{\lambda}}{1 + \lambda} Y = \frac{\lambda}{1 + \lambda} X + \frac{\sqrt{\lambda}}{1 + \lambda} Z, \text{ and, } \mathbb{E}[\langle Xx \rangle] = \frac{\lambda}{1 + \lambda}.$$

In particular, we have $\text{MMSE}(\lambda) = \frac{1}{1+\lambda}$.

3 Bayesian inference with Gaussian additive noise

We start with the following model:

$$\mathbf{Y} = \sqrt{\lambda} \mathbf{X} + \mathbf{Z}, \quad (3.1)$$

where the signal \mathbf{X} is sampled according to some probability distribution P_X over \mathbb{R}^n , and where the noise $\mathbf{Z} = (Z_1, \dots, Z_n) \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$ is independent from \mathbf{X} . The parameter $\lambda \geq 0$ plays the role of a signal-to-noise ratio. We assume that P_X admits a finite second moment: $\mathbb{E}\|\mathbf{X}\|^2 < \infty$. No other assumption is made on P_X . In particular, the fact that the noise is i.i.d. in the dimension n will be crucial in the computation below but no such assumption is made on the signal which is allowed to have any complicated structure component-wise.

Given the observation channel (3.1), the goal of the statistician is to estimate \mathbf{X} given the observations \mathbf{Y} . We assume to be in the “Bayes-optimal” setting, where the statistician knows all the parameters of the inference model, that is the prior distribution P_X and the signal-to-noise ratio λ . As will be clear below, the Mean Squared Error (MSE) is the natural measure of performance when dealing with Gaussian noise.

Following Definition 2.2, we see that a natural object to study is the posterior distribution of \mathbf{X} . We have for the joint distribution:

$$dP_{(X,Y)}(\mathbf{x}, \mathbf{y}) \propto dP_X(\mathbf{x}) e^{\frac{\|\mathbf{y} - \sqrt{\lambda}\mathbf{x}\|^2}{2}} d\mathbf{y}$$

By Bayes rule, the posterior distribution of \mathbf{X} given \mathbf{Y} is

$$dP(\mathbf{x} | \mathbf{Y}) = \frac{1}{\mathcal{Z}(\lambda, \mathbf{Y})} e^{H_{\lambda, \mathbf{Y}}(\mathbf{x})} dP_X(\mathbf{x}), \quad (3.2)$$

where

$$H_{\lambda, \mathbf{Y}}(\mathbf{x}) = \sqrt{\lambda} \mathbf{x}^\top \mathbf{Y} - \frac{\lambda}{2} \|\mathbf{x}\|^2 = \sqrt{\lambda} \mathbf{x}^\top \mathbf{Z} + \lambda \mathbf{x}^\top \mathbf{X} - \frac{\lambda}{2} \|\mathbf{x}\|^2.$$

Definition 3.1

$H_{\lambda, \mathbf{Y}}$ is called the Hamiltonian¹ and the normalizing constant

$$\mathcal{Z}(\lambda, \mathbf{Y}) = \int dP_X(\mathbf{x}) e^{H_{\lambda, \mathbf{Y}}(\mathbf{x})}$$

is called the partition function.

Expectations with respect the posterior distribution (3.2) will be denoted by the Gibbs brackets $\langle \cdot \rangle_\lambda$:

$$\langle f(\mathbf{x}) \rangle_\lambda = \mathbb{E}[f(\mathbf{X})|\mathbf{Y}] = \frac{1}{\mathcal{Z}(\lambda, \mathbf{Y})} \int dP_X(\mathbf{x}) f(\mathbf{x}) e^{H_{\lambda, \mathbf{Y}}(\mathbf{x})},$$

for any measurable function f such that $f(\mathbf{X})$ is integrable.

Definition 3.2

$F(\lambda) = \mathbb{E} \log \mathcal{Z}(\lambda, \mathbf{Y})$ is called the free energy². It is related to the mutual information between \mathbf{X} and \mathbf{Y} by

$$F(\lambda) = \frac{\lambda}{2} \mathbb{E} \|\mathbf{X}\|^2 - I(\mathbf{X}; \mathbf{Y}). \quad (3.3)$$

Proof. The mutual information $I(\mathbf{X}; \mathbf{Y})$ is defined as the Kullback-Leibler divergence between $P_{(\mathbf{X}, \mathbf{Y})}$, the joint distribution of (\mathbf{X}, \mathbf{Y}) and $P_X \otimes P_Y$ the product of the marginal distributions of \mathbf{X} and \mathbf{Y} . $P_{(\mathbf{X}, \mathbf{Y})}$ is absolutely continuous with respect to $P_X \otimes P_Y$ with Radon-Nikodym derivative:

$$\frac{dP_{(\mathbf{X}, \mathbf{Y})}}{dP_X \otimes P_Y}(\mathbf{X}, \mathbf{Y}) = \frac{\exp\left(-\frac{1}{2}\|\mathbf{Y} - \sqrt{\lambda}\mathbf{X}\|^2\right)}{\int \exp\left(-\frac{1}{2}\|\mathbf{Y} - \sqrt{\lambda}\mathbf{x}\|^2\right) dP_X(\mathbf{x})}.$$

Therefore

$$\begin{aligned} I(\mathbf{X}; \mathbf{Y}) &= \mathbb{E} \log \left(\frac{dP_{(\mathbf{X}, \mathbf{Y})}}{dP_X \otimes P_Y}(\mathbf{X}, \mathbf{Y}) \right) \\ &= -\mathbb{E} \log \int dP_X(\mathbf{x}) \exp \left(\sqrt{\lambda} \mathbf{x}^\top \mathbf{Y} - \sqrt{\lambda} \mathbf{X}^\top \mathbf{Y} - \frac{\lambda}{2} \|\mathbf{x}\|^2 + \frac{\lambda}{2} \|\mathbf{X}\|^2 \right) \\ &= -F(\lambda) + \sqrt{\lambda} \mathbb{E} \mathbf{X}^\top \mathbf{Y} - \frac{\lambda}{2} \mathbb{E} \|\mathbf{X}\|^2 \\ &= -F(\lambda) + \frac{\lambda}{2} \mathbb{E} \|\mathbf{X}\|^2, \end{aligned}$$

since $\sqrt{\lambda} \mathbf{X}^\top \mathbf{Y} = \lambda \|\mathbf{X}\|^2 + \sqrt{\lambda} \mathbf{X}^\top \mathbf{Z}$. □

We state now two basic properties of the MMSE. A more detailed analysis can be found in [?, ?].

Proposition 3.1

$\lambda \mapsto \text{MMSE}(\lambda)$ is non-increasing over $\mathbb{R}_{\geq 0}$. Moreover

- $\text{MMSE}(0) = \mathbb{E} \|\mathbf{X} - \mathbb{E}[\mathbf{X}]\|^2$,
- $\text{MMSE}(\lambda) \xrightarrow{\lambda \rightarrow +\infty} 0$.

¹According to the physics convention, this should be minus the Hamiltonian, since a physical system tries to minimize its energy. However, we chose here to remove it for simplicity.

²This is in fact minus the free energy, but we chose to remove the minus sign for simplicity.

Proof. Let $0 < \lambda_2 < \lambda_1$. Define $\Delta_1 = \lambda_1^{-1}$, $\Delta_2 = \lambda_2^{-1}$ and

$$\begin{cases} \mathbf{Y}_1 = \mathbf{X} + \sqrt{\Delta_1} \mathbf{Z}_1 \\ \mathbf{Y}_2 = \mathbf{X} + \sqrt{\Delta_1} \mathbf{Z}_1 + \sqrt{\Delta_2 - \Delta_1} \mathbf{Z}_2, \end{cases}$$

where $\mathbf{X} \sim P_X$ is independent from $\mathbf{Z}_1, \mathbf{Z}_2 \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, \text{Id}_n)$. Now, by independence between $(\mathbf{X}, \mathbf{Y}_1)$ and \mathbf{Z}_2 we have

$$\begin{aligned} \text{MMSE}(\lambda_1) &= \mathbb{E} \|\mathbf{X} - \mathbb{E}[\mathbf{X}|\mathbf{Y}_1]\|^2 = \mathbb{E} \|\mathbf{X} - \mathbb{E}[\mathbf{X}|\mathbf{Y}_1, \mathbf{Z}_2]\|^2 = \mathbb{E} \|\mathbf{X} - \mathbb{E}[\mathbf{X}|\mathbf{Y}_1, \mathbf{Y}_2]\|^2 \\ &\leq \mathbb{E} \|\mathbf{X} - \mathbb{E}[\mathbf{X}|\mathbf{Y}_2]\|^2 = \text{MMSE}(\lambda_2). \end{aligned}$$

Next, notice that

$$\text{MMSE}(\lambda_1) = \mathbb{E} \|\mathbf{X} - \mathbb{E}[\mathbf{X}|\mathbf{Y}_1]\|^2 \leq \mathbb{E} \|\mathbf{X} - \mathbb{E}[\mathbf{X}]\|^2 = \text{MMSE}(0). \quad (3.4)$$

This shows that the MMSE is non-increasing on $\mathbb{R}_{\geq 0}$. The first point is obvious while the second follows from:

$$0 \leq \text{MMSE}(\lambda) = \mathbb{E} \|\mathbf{X} - \mathbb{E}[\mathbf{X}|\mathbf{Y}]\|^2 \leq \mathbb{E} \|\mathbf{X} - \frac{1}{\sqrt{\lambda}} \mathbf{Y}\|^2 = \frac{n}{\lambda} \xrightarrow{\lambda \rightarrow +\infty} 0.$$

□

Proposition 3.2

$\lambda \mapsto \text{MMSE}(\lambda)$ is continuous over $\mathbb{R}_{\geq 0}$.

The proof of Proposition 3.2 can be found in Appendix 8.1.

We present now the very useful “I-MMSE” relation from [?]. This relation was previously known (under a different formulation) as “de Bruijn identity” see [?, Equation 2.12].

Proposition 3.3

For all $\lambda \geq 0$,

$$\frac{\partial}{\partial \lambda} I(\mathbf{X}; \mathbf{Y}) = \frac{1}{2} \text{MMSE}(\lambda) \quad \text{and} \quad F'(\lambda) = \frac{1}{2} \mathbb{E} \langle \mathbf{x}^\top \mathbf{X} \rangle_\lambda = \frac{1}{2} (\mathbb{E} \|\mathbf{X}\|^2 - \text{MMSE}(\lambda)). \quad (3.5)$$

F thus is a convex, differentiable, non-decreasing, and $\frac{1}{2} \mathbb{E} \|\mathbf{X}\|^2$ -Lipschitz function over $\mathbb{R}_{\geq 0}$. If P_X is not a Dirac mass, then F is strictly convex.

TBC...

Proposition 3.3 is proved in Appendix ???. Proposition 3.3 reduces the computation of the MMSE to the computation of the free energy. This will be particularly useful because the free energy F is much easier to handle than the MMSE.

We end this section with the simplest model of the form (3.1), namely the additive Gaussian scalar channel:

$$Y = \sqrt{\lambda}X + Z, \quad (3.6)$$

where $Z \sim \mathcal{N}(0, 1)$ and X is sampled from a distribution P_0 over \mathbb{R} , independently of Z . The corresponding free energy and the MMSE are respectively

$$\psi_{P_0}(\lambda) = \mathbb{E} \log \int dP_0(x) e^{\sqrt{\lambda}Yx - \lambda x^2/2} \quad \text{and} \quad \text{MMSE}_{P_0}(\lambda) = \mathbb{E} \left[(X - \mathbb{E}[X|Y])^2 \right]. \quad (3.7)$$

The study of this simple inference channel will be very useful in the following, because we will see that the inference problems that we are going to study enjoy asymptotically a “decoupling principle” that reduces them to scalar channels like (3.6).

Let us compute the mutual information and the MMSE for particular choices of prior distributions:

Example 3.1 (Gaussian prior: $P_0 = \mathcal{N}(0, 1)$). *In that case $\mathbb{E}[X|Y]$ is simply the orthogonal projection of X on Y :*

$$\mathbb{E}[X|Y] = \frac{\mathbb{E}[XY]}{\mathbb{E}[Y^2]}Y = \frac{\sqrt{\lambda}}{1+\lambda}Y.$$

One deduces $\text{MMSE}_{P_0}(\lambda) = \frac{1}{1+\lambda}$. Using (3.5), we get $I(X; Y) = \frac{1}{2} \log(1+\lambda)$ and $\psi_{P_0}(\lambda) = \frac{1}{2}(\lambda - \log(1+\lambda))$.

Remark 3.1 (Worst-case prior). *Let P_0 be a probability distribution on \mathbb{R} with unit second moment $\mathbb{E}_{P_0}[X^2] = 1$. By considering the estimator $\hat{x} = \frac{\sqrt{\lambda}}{1+\lambda}Y$, one obtain $\text{MMSE}_{P_0}(\lambda) \leq \frac{1}{1+\lambda}$. We conclude:*

$$\sup_{P_0} \text{MMSE}_{P_0}(\lambda) = \frac{1}{1+\lambda} \quad \text{and} \quad \inf_{P_0} \psi_{P_0}(\lambda) = \frac{1}{2}(\lambda - \log(1+\lambda)),$$

where the supremum and infimum are both over the probability distributions that have unit second moment. The standard normal distribution $P_0 = \mathcal{N}(0, 1)$ achieves both extrema.

Example 3.2 (Rademacher prior: $P_0 = \frac{1}{2}\delta_{+1} + \frac{1}{2}\delta_{-1}$). *We compute $\psi_{P_0}(\lambda) = \mathbb{E} \log \cosh(\sqrt{\lambda}Z + \lambda) - \frac{\lambda}{2}$ and $I(X; Y) = \lambda - \mathbb{E} \log \cosh(\sqrt{\lambda}Z + \lambda)$. The I-MMSE relation gives*

$$\begin{aligned} \frac{1}{2} \text{MMSE}(\lambda) &= \frac{\partial}{\partial \lambda} I(X; Y) = 1 - \mathbb{E} \left[\left(\frac{1}{2\sqrt{\lambda}}Z + 1 \right) \tanh(\sqrt{\lambda}Z + \lambda) \right] \\ &= 1 - \mathbb{E} \tanh(\sqrt{\lambda}Z + \lambda) - \frac{1}{2} \mathbb{E} \tanh'(\sqrt{\lambda}Z + \lambda) \\ &= \frac{1}{2} - \mathbb{E} \tanh(\sqrt{\lambda}Z + \lambda) + \frac{1}{2} \mathbb{E} \tanh^2(\sqrt{\lambda}Z + \lambda) \end{aligned}$$

where we used Gaussian integration by parts. Since by the Nishimori property $\mathbb{E}\langle xX \rangle_\lambda = \mathbb{E}\langle x \rangle_\lambda^2$, one has $\mathbb{E} \tanh(\sqrt{\lambda}Z + \lambda) = \mathbb{E} \tanh^2(\sqrt{\lambda}Z + \lambda)$ and therefore $\text{MMSE}(\lambda) = 1 - \mathbb{E} \tanh(\sqrt{\lambda}Z + \lambda)$.

4 Bayes-optimal inference

We introduce in this chapter some general properties of Bayes-optimal inference, that will be used repeatedly in the sequel. Let us first define what we mean by *Bayes-optimal inference*.

We consider a statistical problem where we would like to recover a signal vector $\mathbf{X} \in \mathbb{R}^n$ from some observations $\mathbf{Y} \in \mathbb{R}^m$. We assume that (\mathbf{X}, \mathbf{Y}) is drawn from some probability distribution μ over $\mathbb{R}^n \times \mathbb{R}^m$. Given a performance criterion, a Bayes-optimal estimator (or simply Bayes estimator) is an estimator of \mathbf{X} given \mathbf{Y} that achieves the best performance for this criterion. For instance if we measure the performance of an estimator $\hat{\mathbf{x}}$ by its mean square error $\text{MSE}(\hat{\mathbf{x}}) = \mathbb{E}\|\mathbf{X} - \hat{\mathbf{x}}(\mathbf{Y})\|^2$, then the Bayes-optimal estimator is simply the posterior mean $\hat{\mathbf{x}}^{\text{Bayes}}(\mathbf{Y}) = \mathbb{E}[\mathbf{X}|\mathbf{Y}]$.

The goal of this chapter is to present some general properties of Bayes-optimal estimators. In Section 5 we introduce what we will call (according to the statistical physics terminology) the “Nishimori identity” which is nothing more than a rewriting of Bayes rule. In Section 6 we will study the links between various natural performance metrics for estimators and show that they are in some sense equivalent. In Sections 3 we analyse the special case where $\mathbf{Y} = \sqrt{\lambda}\mathbf{X} + \mathbf{Z}$, where $\lambda \geq 0$ and \mathbf{Z} is some Gaussian noise. This is the starting point of the study of the “spiked” matrix and tensor models. Finally we consider in Section 7 a simple example to illustrate the tools of this chapter.

5 The Nishimori identity

In order to analyze Bayes-optimal estimators, we will need to examine the posterior distribution of \mathbf{X} given \mathbf{Y} . To do so we will often consider i.i.d. samples $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}$ from the posterior distribution $P(\cdot|\mathbf{Y})$, independently of everything else. Such samples are called replicas. The (obvious) identity below (which is simply Bayes rule) is named after the works of Nishimori [?, ?] on “gauge-symmetric” spin glasses. It states that the planted solution \mathbf{X} behaves like a replica.

Proposition 5.1 (*Nishimori identity*)

Let (\mathbf{X}, \mathbf{Y}) be a couple of random variables on a polish space. Let $k \geq 1$ and let $\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}$ be k i.i.d. samples (given \mathbf{Y}) from the distribution $\mathbb{P}(\mathbf{X} = \cdot | \mathbf{Y})$, independently of every other random variables. Let us denote $\langle \cdot \rangle$ the expectation with respect to $\mathbb{P}(\mathbf{X} = \cdot | \mathbf{Y})$ and \mathbb{E} the expectation with respect to (\mathbf{X}, \mathbf{Y}) . Then, for all continuous bounded function f

$$\mathbb{E}\langle f(\mathbf{Y}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)}) \rangle = \mathbb{E}\langle f(\mathbf{Y}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k-1)}, \mathbf{X}) \rangle.$$

Proof. It is equivalent to sample the couple (\mathbf{X}, \mathbf{Y}) according to its joint distribution or to sample first \mathbf{Y} according to its marginal distribution and then to sample \mathbf{X} conditionally to \mathbf{Y} from its conditional distribution $\mathbb{P}(\mathbf{X} = \cdot | \mathbf{Y})$. Thus the $(k+1)$ -tuple $(\mathbf{Y}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k)})$ is equal in law to $(\mathbf{Y}, \mathbf{x}^{(1)}, \dots, \mathbf{x}^{(k-1)}, \mathbf{X})$. \square

6 Performance measure and optimal estimators

We consider two random vectors \mathbf{X} and \mathbf{Y} that live respectively in \mathbb{R}^n and \mathbb{R}^m . We assume (for simplicity) that $\|\mathbf{X}\| = 1$ almost surely. As explained above, given the observations \mathbf{Y} , our goal is to estimate \mathbf{X} with an estimator $\hat{\mathbf{x}}(\mathbf{Y})$. In order to evaluate the performance of such an estimator, what criterion should we take?

The probably most natural way to characterize the performance of $\hat{\mathbf{x}}$ is by its mean-squared error:

$$\text{MSE}(\hat{\mathbf{x}}) = \mathbb{E}\|\mathbf{X} - \hat{\mathbf{x}}(\mathbf{Y})\|^2.$$

By Pythagorean theorem, we know that the optimal estimator which respect to this metric is the posterior mean $\hat{\mathbf{x}}(\mathbf{Y}) = \mathbb{E}[\mathbf{X}|\mathbf{Y}]$ which achieves the minimal mean square error:

$$\text{MMSE}(\mathbf{X}|\mathbf{Y}) \stackrel{\text{def}}{=} \mathbb{E}\|\mathbf{X} - \mathbb{E}[\mathbf{X}|\mathbf{Y}]\|^2. \quad (6.1)$$

However, the MSE is not always an appropriate criterion. Indeed in many cases it is only possible to recover \mathbf{X} up to its sign: think for instance of the Spiked Wigner Model $\mathbf{Y} = \mathbf{X}\mathbf{X}^\top + \text{Noise}$ with $\mathbf{X} \sim \text{Unif}(\mathbb{S}^{n-1})$. In such case, $\mathbb{E}[\mathbf{X}|\mathbf{Y}] = 0$: the best estimator in term of MSE does not even depend on the observations \mathbf{Y} !

For this kind of problems one should rather consider the correlation (also known as cosine similarity) *in absolute value* between $\hat{\mathbf{x}}$ and the signal \mathbf{X} :

$$\sup_{\|\hat{\mathbf{x}}\|=1} \mathbb{E}\left[|(\hat{\mathbf{x}}(\mathbf{Y}); \mathbf{X})|\right] \quad \text{or} \quad \sup_{\|\hat{\mathbf{x}}\|=1} \mathbb{E}\left[(\hat{\mathbf{x}}(\mathbf{Y}); \mathbf{X})^2\right], \quad (6.2)$$

where $(\cdot; \cdot)$ denotes the Euclidean inner product and where the suprema are taken over all estimators $\hat{\mathbf{x}} : \mathbb{R}^m \rightarrow \mathbb{S}^{n-1}$.

Let us introduce some notations. We will use the Gibbs Bracket $\langle \cdot \rangle$ to write expectations with respect to the posterior distribution of \mathbf{X} given \mathbf{Y} :

$$\langle f(\mathbf{x}) \rangle = \mathbb{E}[f(\mathbf{X})|\mathbf{Y}],$$

for all measurable function f such that $f(\mathbf{X})$ is integrable. In particular, we will be interested by the $n \times n$ positive semi-definite (random) matrix:

$$\mathbf{M} \stackrel{\text{def}}{=} \langle \mathbf{x}\mathbf{x}^\top \rangle = \mathbb{E}[\mathbf{X}\mathbf{X}^\top|\mathbf{Y}]. \quad (6.3)$$

\mathbf{M} is the Bayes-optimal estimator (in terms of mean square error) for estimating the matrix $\mathbf{X}\mathbf{X}^\top$:

$$\text{MMSE}(\mathbf{X}\mathbf{X}^\top|\mathbf{Y}) = \mathbb{E}\|\mathbf{X}\mathbf{X}^\top - \mathbb{E}[\mathbf{X}\mathbf{X}^\top|\mathbf{Y}]\|^2. \quad (6.4)$$

An easy computation gives $\text{MMSE}(\mathbf{X}\mathbf{X}^\top|\mathbf{Y}) = 1 - \mathbb{E}[\text{Tr}(\mathbf{M}^2)]$. The matrix \mathbf{M} is also related to the second quantity of (6.2) through its largest eigenvalue $\lambda_{\max}(\mathbf{M})$:

Lemma 6.1

$$\sup_{\|\hat{\mathbf{x}}\|=1} \mathbb{E}\left[(\hat{\mathbf{x}}(\mathbf{Y}); \mathbf{X})^2\right] = \mathbb{E}\left[\lambda_{\max}(\mathbf{M})\right]$$

and the optimal estimator for this metric is a unit eigenvector of \mathbf{M} associated to its largest eigenvalue $\lambda_{\max}(\mathbf{M})$.

Proof. Let $\hat{\mathbf{x}}$ be an estimator of \mathbf{X} . By the Nishimori identity (Proposition 5.1), we have

$$\mathbb{E}[(\hat{\mathbf{x}}(\mathbf{Y}); \mathbf{X})^2] = \mathbb{E}[\hat{\mathbf{x}}(\mathbf{Y})^\top \mathbf{X} \mathbf{X}^\top \hat{\mathbf{x}}(\mathbf{Y})] = \mathbb{E}[\langle \hat{\mathbf{x}}(\mathbf{Y})^\top \mathbf{x} \mathbf{x}^\top \hat{\mathbf{x}}(\mathbf{Y}) \rangle] = \mathbb{E}[\hat{\mathbf{x}}(\mathbf{Y})^\top \mathbf{M} \hat{\mathbf{x}}(\mathbf{Y})],$$

the lemma follows. \square

Lemma 6.1 tells us that the top unit eigenvector $\hat{\mathbf{v}}$ of \mathbf{M} maximizes $\mathbb{E}[(\hat{\mathbf{x}}(\mathbf{Y}); \mathbf{X})^2]$. In the following, we will show that under a simple condition (that will hold for the models we consider in this manuscript), the estimator $\hat{\mathbf{v}}$ is “asymptotically optimal” (in the limit of large dimension) for the two metrics (6.2) and $\lambda_{\max} \hat{\mathbf{v}} \hat{\mathbf{v}}^\top$ is optimal for the estimation of $\mathbf{X} \mathbf{X}^\top$ in terms of mean square error.

To introduce this condition and the asymptotic limit, we need to consider to a sequence of inference problems. We assume that for all $n \geq 1$ we have two random vectors $\mathbf{X}_{[n]}$ and $\mathbf{Y}_{[n]}$ respectively in \mathbb{S}^{n-1} and \mathbb{R}^{m_n} , for some sequence $(m_n)_{n \geq 1}$. Our goal is again to estimate $\mathbf{X}_{[n]}$ from the observation of $\mathbf{Y}_{[n]}$ when n is very large: we would like for instance to compute the limits of (6.2) and (6.4) as $n \rightarrow \infty$. Moreover, we would like to know which estimators are “asymptotically optimal”, i.e. whose performance reach in the $n \rightarrow \infty$ limit the optimal one. In the following, in order to simplify the notations, we will write \mathbf{X} and \mathbf{Y} instead of $\mathbf{X}_{[n]}$ and $\mathbf{Y}_{[n]}$.

Proposition 6.1

Let us denote by G_n the posterior distribution of \mathbf{X} given \mathbf{Y} . Notice that G_n is a random probability distributions on \mathbb{S}^{n-1} . Assume that there exists $q \in [0, 1]$ such that for $\mathbf{x}^{(1)}, \mathbf{x}^{(2)} \stackrel{i.i.d.}{\sim} G_n$ we have

$$\left| (\mathbf{x}^{(1)}; \mathbf{x}^{(2)}) \right| \xrightarrow[n \rightarrow \infty]{(d)} q. \quad (6.5)$$

Then $\text{Tr}(\mathbf{M}^2) \xrightarrow[n \rightarrow \infty]{(d)} q^2$ and

$$\lambda_{\max}(\mathbf{M}) \xrightarrow[n \rightarrow \infty]{(d)} q.$$

Proof. Let us compute $\text{Tr}(\mathbf{M}^2) = \text{Tr}(\langle \mathbf{x} \mathbf{x}^\top \rangle \langle \mathbf{x} \mathbf{x}^\top \rangle) = \langle (\mathbf{x}^{(1)}; \mathbf{x}^{(2)})^2 \rangle \xrightarrow[n \rightarrow \infty]{} q^2$, by assumption.

If $q = 0$, then the result is obvious since $\lambda_{\max}(\mathbf{M})^2 \leq \text{Tr}(\mathbf{M}^2)$.

Notice that $\text{Tr}(\mathbf{M}^3) \leq \lambda_{\max}(\mathbf{M}) \text{Tr}(\mathbf{M}^2)$, so it suffices to show that

$$\text{Tr}(\mathbf{M}^3) = \langle (\mathbf{x}^{(1)}; \mathbf{x}^{(2)})(\mathbf{x}^{(2)}; \mathbf{x}^{(3)})(\mathbf{x}^{(3)}; \mathbf{x}^{(1)}) \rangle \xrightarrow[n \rightarrow \infty]{} q^3,$$

for $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)} \stackrel{i.i.d.}{\sim} G_n$. This follows from Lemma 6.2 below. \square

Lemma 6.2

Under the assumptions of Proposition 6.1, we have for $\mathbf{x}^{(1)}, \mathbf{x}^{(2)}, \mathbf{x}^{(3)} \stackrel{i.i.d.}{\sim} G_n$,

$$(\mathbf{x}^{(1)}; \mathbf{x}^{(2)})(\mathbf{x}^{(2)}; \mathbf{x}^{(3)})(\mathbf{x}^{(3)}; \mathbf{x}^{(1)}) \xrightarrow[n \rightarrow \infty]{(d)} q^3.$$

Lemma 6.2 will be proved in Appendix ???. From Proposition 6.1 we deduce the main result of this section:

Proposition 6.2

Let $\hat{\mathbf{v}}$ be a leading unit eigenvector of \mathbf{M} (which is defined by (6.3)). Under the assumptions of Proposition 6.1, we have

$$|(\hat{\mathbf{v}}; \mathbf{X})| \xrightarrow[n \rightarrow \infty]{(d)} \sqrt{q}. \quad (6.6)$$

Further $\lim_{n \rightarrow \infty} \text{MMSE}(\mathbf{X}\mathbf{X}^\top | \mathbf{Y}) = 1 - q^2$,

$$\lim_{n \rightarrow \infty} \sup_{\|\hat{\mathbf{x}}\|=1} \mathbb{E}[|(\hat{\mathbf{x}}(\mathbf{Y}); \mathbf{X})|] = \sqrt{q}, \quad \text{and} \quad \lim_{n \rightarrow \infty} \sup_{\|\hat{\mathbf{x}}\|=1} \mathbb{E}[(\hat{\mathbf{x}}(\mathbf{Y}); \mathbf{X})^2] = q. \quad (6.7)$$

Proof. Let us abbreviate $\lambda_{\max} \stackrel{\text{def}}{=} \lambda_{\max}(\mathbf{M})$. By Lemma 6.1 and Proposition 6.1 we have

$$\mathbb{E}[(\hat{\mathbf{v}}; \mathbf{X})^2] = \mathbb{E}[\lambda_{\max}] \xrightarrow[n \rightarrow \infty]{} q. \quad (6.8)$$

Hence if $q = 0$, the Proposition follows easily. Assume now that $q > 0$. Using Pythagorean Theorem and the Nishimori identity (Proposition 5.1) we get

$$\begin{aligned} \mathbb{E}\|\mathbf{X}^{\otimes 4} - \lambda_{\max}^2 \hat{\mathbf{v}}^{\otimes 4}\|^2 &\geq \mathbb{E}\|\mathbf{X}^{\otimes 4} - \langle \mathbf{x}^{\otimes 4} \rangle\|^2 \\ &= 1 + \mathbb{E}[\langle \langle \mathbf{x}^{\otimes 4} \rangle; \langle \mathbf{x}^{\otimes 4} \rangle \rangle] - 2\mathbb{E}[(\mathbf{X}^{\otimes 4}; \langle \mathbf{x}^{\otimes 4} \rangle)] \\ &= 1 - \mathbb{E}\left[\left\langle (\mathbf{x}^{(1)}; \mathbf{x}^{(2)})^4 \right\rangle\right] \xrightarrow[n \rightarrow \infty]{} 1 - q^4, \end{aligned}$$

where the last limit follows from the assumption (6.5). Since

$$\mathbb{E}\|\mathbf{X}^{\otimes 4} - \lambda_{\max}^2 \hat{\mathbf{v}}^{\otimes 4}\|^2 = 1 + \mathbb{E}[\lambda_{\max}^4] - 2\mathbb{E}[\lambda_{\max}^2(\mathbf{X}; \hat{\mathbf{v}})^4],$$

using Proposition 6.1, we deduce (recall that we assumed $q > 0$) that $\limsup_{n \rightarrow \infty} \mathbb{E}[(\hat{\mathbf{v}}; \mathbf{X})^4] \leq q^2$. Together with (6.8) this gives that $|(\hat{\mathbf{v}}; \mathbf{X})| \xrightarrow[n \rightarrow \infty]{} \sqrt{q}$.

The next point is a consequence of Proposition 6.1 because $\text{MMSE}(\mathbf{X}\mathbf{X}^\top | \mathbf{Y}) = 1 - \mathbb{E}[\text{Tr}(\mathbf{M}^2)]$. To prove (6.7) simply notice that

$$\mathbb{E}[|(\hat{\mathbf{v}}; \mathbf{X})|]^2 \leq \sup_{\|\hat{\mathbf{x}}\|=1} \mathbb{E}[|(\hat{\mathbf{x}}(\mathbf{Y}); \mathbf{X})|]^2 \leq \sup_{\|\hat{\mathbf{x}}\|=1} \mathbb{E}[(\hat{\mathbf{x}}(\mathbf{Y}); \mathbf{X})^2] = \mathbb{E}[\lambda_{\max}],$$

which proves (6.7) using (6.6) and Proposition 6.1. \square

From Proposition 6.2, we deduce that the estimator $\widehat{\mathbf{A}} \stackrel{\text{def}}{=} \lambda_{\max}(\mathbf{M}) \hat{\mathbf{v}} \hat{\mathbf{v}}^\top$ achieves asymptotically the minimal mean square error for the estimation of $\mathbf{X}\mathbf{X}^\top$:

$$\lim_{n \rightarrow \infty} \mathbb{E}\|\mathbf{X}\mathbf{X}^\top - \widehat{\mathbf{A}}\|^2 = 1 - q^2.$$

Remark 6.1. For simplicity we assumed in this section that $\|\mathbf{X}\|^2 = 1$ almost surely. However we will need to work in the next chapters under a slightly weaker condition, namely $\|\mathbf{X}\|^2 \xrightarrow[n \rightarrow \infty]{} 1$. It is not difficult to modify the proofs of this section to see that Lemma 6.1, Proposition 6.1, Lemma 6.2 and Proposition 6.2 still hold, provided that $\|\mathbf{X}\|^2 \xrightarrow[n \rightarrow \infty]{} 1$ for the Wasserstein distance of order 4 (i.e. $\|\mathbf{X}\|^2 \xrightarrow[n \rightarrow \infty]{} 1$ in distribution and $\mathbb{E}\|\mathbf{X}\|^8 \xrightarrow[n \rightarrow \infty]{} 1$).

7 A warm-up: the “needle in a haystack” problem

In order to illustrate the results seen in the previous sections, we study now a very simple inference model. Let (e_1, \dots, e_{2^n}) be the canonical basis of \mathbb{R}^{2^n} . Let $\sigma_0 \sim \text{Unif}(\{1, \dots, 2^n\})$ and define $\mathbf{X} = e_{\sigma_0}$ (i.e. \mathbf{X} is chosen uniformly over the canonical basis of \mathbb{R}^{2^n}). Suppose here that we observe:

$$\mathbf{Y} = \sqrt{\lambda n} \mathbf{X} + \mathbf{Z},$$

where $\mathbf{Z} = (Z_1, \dots, Z_{2^n}) \stackrel{\text{i.i.d.}}{\sim} \mathcal{N}(0, 1)$, independently from σ_0 . The goal here is to estimate \mathbf{X} or equivalently to find σ_0 . The posterior distribution reads:

$$\begin{aligned} \mathbb{P}(\sigma_0 = \sigma | \mathbf{Y}) &= \mathbb{P}(\mathbf{X} = e_\sigma | \mathbf{Y}) = \frac{1}{\mathcal{Z}_n(\lambda)} 2^{-n} \exp \left(\sqrt{\lambda n} e_\sigma^\top \mathbf{Y} - \frac{\lambda n}{2} \|e_\sigma\|^2 \right) \\ &= \frac{1}{\mathcal{Z}_n(\lambda)} 2^{-n} \exp \left(\sqrt{\lambda n} Z_\sigma + \lambda n \mathbb{1}(\sigma = \sigma_0) - \frac{\lambda n}{2} \right), \end{aligned}$$

where $\mathcal{Z}_n(\lambda)$ is the partition function

$$\mathcal{Z}_n(\lambda) = \frac{1}{2^n} \sum_{\sigma=1}^{2^n} \exp \left(\sqrt{\lambda n} Z_\sigma + \lambda n \mathbb{1}(\sigma = \sigma_0) - \frac{\lambda n}{2} \right).$$

We will be interested in computing the free energy $F_n(\lambda) = \frac{1}{n} \mathbb{E} \log \mathcal{Z}_n(\lambda)$ in order to deduce then the minimal mean squared error using the I-MMSE relation (3.5) presented in the previous section.

Although its simplicity, this model is interesting for many reasons. First, it is one of the simplest statistical model for which one observes a phase transition. Second it is the “planted” analog of the random energy model (REM) introduced in statistical physics by Derrida [?, ?], for which the free energy reads $\frac{1}{n} \mathbb{E} \log \sum_{\sigma} \frac{1}{2^n} \exp(\sqrt{\lambda n} Z_\sigma)$. Third, as we will see in Section ??, this model correspond to the “large order limit” of a rank-one tensor estimation model.

We start by computing the limiting free energy:

Theorem 7.1

$$\lim_{n \rightarrow \infty} F_n(\lambda) = \begin{cases} 0 & \text{if } \lambda \leq 2 \log 2, \\ \frac{\lambda}{2} - \log(2) & \text{if } \lambda \geq 2 \log 2. \end{cases}$$

Proof. Using Jensen’s inequality

$$\begin{aligned} F_n(\lambda) &\leq \frac{1}{n} \mathbb{E} \log \mathbb{E} [\mathcal{Z}_n(\lambda) | \sigma_0, Z_{\sigma_0}] = \frac{1}{n} \mathbb{E} \log \left(1 - \frac{1}{2^n} + e^{\sqrt{\lambda n} Z_{\sigma_0} + \frac{\lambda n}{2} - \log(2)n} \right) \\ &\leq \frac{1}{n} \mathbb{E} \log \left(1 + e^{\frac{\lambda n}{2} - \log(2)n} \right) + \sqrt{\frac{\lambda}{n}} \xrightarrow{n \rightarrow \infty} \begin{cases} 0 & \text{if } \lambda \leq 2 \log(2), \\ \frac{\lambda}{2} - \log(2) & \text{if } \lambda \geq 2 \log(2). \end{cases} \end{aligned}$$

F_n is non-negative since $F_n(0) = 0$ and F_n is non-decreasing. We have therefore $F_n(\lambda) \xrightarrow{n \rightarrow \infty} 0$ for all $\lambda \in [0, 2 \log(2)]$. We have also, by only considering the term $\sigma = \sigma_0$:

$$F_n(\lambda) \geq \frac{1}{n} \mathbb{E} \log \left(\frac{e^{\sqrt{\lambda n} Z_{\sigma_0} + \frac{\lambda n}{2}}}{2^n} \right) = \frac{\lambda}{2} - \log(2).$$

We obtain therefore that $F_n(\lambda) \xrightarrow{n \rightarrow \infty} \frac{\lambda}{2} - \log(2)$ for $\lambda \geq 2 \log(2)$. \square

Using the I-MMSE relation (3.5), we deduce the limit of the minimum mean squared error $\text{MMSE}_n(\lambda) = \min_{\hat{\theta}} \mathbb{E} \|\mathbf{X} - \hat{\theta}(\mathbf{Y})\|^2$:

$$\text{MMSE}_n(\lambda) = \mathbb{E} \|\mathbf{X}\|^2 - 2F'_n(\lambda) = 1 - 2F'_n(\lambda).$$

F_n is a convex function of λ , thus (see Proposition ??) its derivative converges to the derivative of its limit at each λ at which the limit is differentiable, i.e. for all $\lambda \in (0, +\infty) \setminus \{2 \log(2)\}$. We obtain therefore that for all $\lambda > 0$,

- if $\lambda < 2 \log(2)$, then $\text{MMSE}_n(\lambda) \xrightarrow{n \rightarrow \infty} 1$: one can not recover \mathbf{X} better than a random guess.
- if $\lambda > 2 \log(2)$, then $\text{MMSE}_n(\lambda) \xrightarrow{n \rightarrow \infty} 0$: one can recover \mathbf{X} perfectly.

Of course, the result we obtain here is (almost) trivial since the maximum likelihood estimator

$$\hat{\sigma}(\mathbf{Y}) = \arg \max_{1 \leq \sigma \leq 2^n} Y_\sigma$$

of σ_0 is easy to analyze. Indeed, $\max_{\sigma} Z_\sigma \simeq \sqrt{2 \log(2)n}$ with high probability so that the maximum likelihood estimator recovers perfectly the signal for $\lambda > 2 \log(2)$ with high probability.

8 Appendix

8.1 Proof of Proposition 3.2

We start by proving that MMSE is continuous at $\lambda = 0$. Let $\lambda \geq 0$ and consider $\mathbf{Y}, \mathbf{X}, \mathbf{Z}$ as given by (3.1). By dominated convergence one has almost surely that

$$\mathbb{E}[\mathbf{X}|\mathbf{Y}] = \frac{\int dP_X(\mathbf{x}) \mathbf{x} e^{-\frac{1}{2} \|\sqrt{\lambda} \mathbf{x} - \mathbf{Y}\|^2}}{\int dP_X(\mathbf{x}) e^{-\frac{1}{2} \|\sqrt{\lambda} \mathbf{x} - \mathbf{Y}\|^2}} \xrightarrow{\lambda \rightarrow 0} \mathbb{E}[\mathbf{X}].$$

Then by Fatou's Lemma we get

$$\liminf_{\lambda \rightarrow 0} \text{MMSE}(\lambda) \geq \mathbb{E} \left[\liminf_{\lambda \rightarrow 0} \|\mathbf{X} - \mathbb{E}[\mathbf{X}|\mathbf{Y}]\|^2 \right] = \mathbb{E} \|\mathbf{X} - \mathbb{E}[\mathbf{X}]\|^2.$$

Combining this with the bound $\text{MMSE}(\lambda) \leq \mathbb{E} \|\mathbf{X} - \mathbb{E}[\mathbf{X}]\|^2$ gives $\text{MMSE}(\lambda) \xrightarrow{\lambda \rightarrow 0} \mathbb{E} \|\mathbf{X} - \mathbb{E}[\mathbf{X}]\|^2$. This proves that the MMSE is continuous at $\lambda = 0$.

Let us now prove that the MMSE is continuous on $\mathbb{R}_{\geq 0}^*$. We need here a technical lemma:

Lemma 8.1

For all $\lambda > 0$, $p \geq 1$

$$\mathbb{E} \|\mathbf{X} - \langle \mathbf{x} \rangle_\lambda\|^{2p} \leq \frac{2^p (2p!)}{\lambda^p p!} n^{p+1}.$$

Proof. We reproduce here the proof from [?], Proposition 5. We start with the equality

$$\sqrt{\lambda}(\mathbf{X} - \langle \mathbf{x} \rangle_\lambda) = \sqrt{\lambda}\mathbf{X} - \mathbb{E}[\sqrt{\lambda}\mathbf{X}|\mathbf{Y}] = \mathbf{Y} - \mathbf{Z} - \mathbb{E}[\mathbf{Y} - \mathbf{Z}|\mathbf{Y}] = \mathbb{E}[\mathbf{Z}|\mathbf{Y}] - \mathbf{Z}.$$

We have therefore

$$\mathbb{E}\|\mathbf{X} - \langle \mathbf{x} \rangle_\lambda\|^{2p} = \frac{1}{\lambda^p} \mathbb{E}\|\mathbb{E}[\mathbf{Z}|\mathbf{Y}] - \mathbf{Z}\|^{2p} \leq \frac{2^{2p-1}}{\lambda^p} \mathbb{E}[\|\mathbb{E}[\mathbf{Z}|\mathbf{Y}]\|^{2p} + \|\mathbf{Z}\|^{2p}] \leq \frac{2^{2p}}{\lambda^p} \mathbb{E}\|\mathbf{Z}\|^{2p}.$$

It remains to bound

$$\mathbb{E}\|\mathbf{Z}\|^{2p} \leq n^p \mathbb{E} \left[\sum_{i=1}^n Z_i^{2p} \right] = n^{p+1} \frac{(2p)!}{2^p p!}.$$

□

Let $\lambda_0 > 0$. The family of random variables $(\|\mathbf{X} - \langle \mathbf{x} \rangle_\lambda\|^2)_{\lambda \geq \lambda_0}$ is bounded in L^2 by Lemma 8.1 and is therefore uniformly integrable. The function $\lambda \mapsto \|\mathbf{X} - \langle \mathbf{x} \rangle_\lambda\|^2$ is continuous on $[\lambda_0, +\infty)$, the uniform integrability ensures then that $\text{MMSE} : \lambda \mapsto \mathbb{E}\|\mathbf{X} - \langle \mathbf{x} \rangle_\lambda\|^2$ is continuous over $[\lambda_0, +\infty)$. This is valid for all $\lambda_0 > 0$: we conclude that MMSE is continuous over $(0, +\infty)$.