

Lecture 21: Derivation of AMP II

Lecturer: Song Mei

Scriber: Alexander Tsigler

Proof reader: Alexander Tsigler

1 MP for LASSO

Previously we have defined the following Message Passing (MP) algorithm for a general Markov random field:

Definition 1 (Message passing algorithm). For each k , $\{m_{i \rightarrow a}^k, v_{i \rightarrow a}^k, \hat{m}_{a \rightarrow i}^k, \hat{v}_{a \rightarrow i}^k\}$ are called “beliefs”, which are real values. Define $\rho_{i \rightarrow a}^k(x_i)$ and $\hat{\rho}_{a \rightarrow i}^k(x_i)$ as the densities of $N(m_{i \rightarrow a}^k, v_{i \rightarrow a}^k)$ and $N(\hat{m}_{a \rightarrow i}^k, \hat{v}_{a \rightarrow i}^k)$ respectively, i.e.,

$$\rho_{i \rightarrow a}^k(x_i) = \frac{1}{\sqrt{2\pi v_{i \rightarrow a}^k}} \exp\left\{-\frac{(x_i - m_{i \rightarrow a}^k)^2}{2v_{i \rightarrow a}^k}\right\},$$

$$\hat{\rho}_{a \rightarrow i}^k(x_i) = \frac{1}{\sqrt{2\pi \hat{v}_{a \rightarrow i}^k}} \exp\left\{-\frac{(x_i - \hat{m}_{a \rightarrow i}^k)^2}{2\hat{v}_{a \rightarrow i}^k}\right\}.$$

Given initialization $\{m_{i \rightarrow a}^0, v_{i \rightarrow a}^0, \hat{m}_{a \rightarrow i}^0, \hat{v}_{a \rightarrow i}^0\}$, compute

$$\hat{\gamma}_{a \rightarrow i}^k(x_i) \propto \int \psi_a(x_{\partial a}) \prod_{j \in \partial a \setminus i} \rho_{j \rightarrow a}^k(x_j) dx_{\partial a \setminus i}, \quad (1)$$

$$\gamma_{i \rightarrow a}^{k+1}(x_i) \propto \psi_i(x_i) \prod_{b \in \partial i \setminus a} \hat{\rho}_{b \rightarrow i}^k(x_i), \quad (2)$$

and update the messages as

$$(\hat{m}_{a \rightarrow i}^k, \hat{v}_{a \rightarrow i}^k) = \text{mean and variance of } \hat{\gamma}_{a \rightarrow i}^k(x_i), \quad (3)$$

$$(m_{i \rightarrow a}^k, v_{i \rightarrow a}^k) = \text{mean and variance of } \gamma_{i \rightarrow a}^k(x_i). \quad (4)$$

Finally, we extract the marginal as

$$\gamma_i^{k+1}(x_i) \propto \psi_i(x_i) \prod_{b \in \partial i} \hat{\rho}_{b \rightarrow i}^k(x_i). \quad (5)$$

A very important property of this algorithm is that it is “non-backtracking”, meaning that the message that a sends to i does not depend on the message that i sent to a in the previous iteration and vice versa. Note, however, that the expression for the final output includes the product over the whole ∂i .

1.1 LASSO with finite temperature

In this lecture we apply that procedure to the LASSO problem:

Example 1 (LASSO with temperature β). Consider

$$\mu_\beta(x) = \prod_{a=1}^n \underbrace{\exp\left\{-\frac{(y_a - \langle A_a, x \rangle_2)^2}{2}\right\}}_{\psi_a(x_{\partial a})} \prod_{i=1}^d \underbrace{\exp(-\beta \lambda |x_i|)}_{\psi_i(x_i)}.$$

Here $V = [d]$ and $F = [n]$. The BP update rule is given by

$$\begin{aligned}\hat{\mu}_{a \rightarrow i}^k(x_i) &\propto \int_{\mathbb{R}^{d-1}} \prod_{j \neq i} dx_j \exp\left\{-\frac{(y_a - \langle A_a, x \rangle)^2}{2}\right\} \prod_{j \neq i} \mu_{j \rightarrow a}^k(x_j), \\ \mu_{i \rightarrow a}^{k+1}(x_i) &\propto \exp(-\beta\lambda|x_i|) \prod_{b \neq a} \hat{\mu}_{b \rightarrow i}^k(x_i),\end{aligned}$$

and the extracted marginal is

$$\mu_i^{k+1}(x_i) \propto \exp(-\beta\lambda|x_i|) \prod_{b \in [n]} \hat{\mu}_{b \rightarrow i}^k(x_i).$$

We need to understand what are the means and variances of distributions $\gamma_{a \rightarrow i}^k$ and $\hat{\gamma}_{i \rightarrow a}^k$. Let's start with $\gamma_{a \rightarrow i}^k$. Plugging in the expressions for LASSO in the MP we get

$$\begin{aligned}\gamma_{i \rightarrow a}^{k+1}(x_i) &\propto \psi_i(x_i) \prod_{b \in \partial i \setminus a} \hat{\rho}_{b \rightarrow i}^k(x_i) \\ &\propto \exp(-\beta\lambda|x_i|) \times \exp\left\{-\sum_{b \neq a} \frac{(x_i - \hat{m}_{b \rightarrow i}^k)^2}{2\hat{v}_{b \rightarrow i}^k}\right\} \\ &\propto \exp\left\{-\beta\left(\frac{(x_i - \theta_{i \rightarrow a}^k)^2}{2\zeta_{i \rightarrow i}^k} + \lambda|x_i|\right)\right\},\end{aligned}$$

where $\theta_{i \rightarrow a}^k$ and $\zeta_{i \rightarrow i}^k$ are mean and variance of the Gaussian distribution with density, proportional to

$$\exp\left\{-\sum_{b \neq a} \frac{(x_i - \hat{m}_{b \rightarrow i}^k)^2}{2\beta\hat{v}_{b \rightarrow i}^k}\right\}.$$

The exact expressions are the following

$$\beta(\zeta_{i \rightarrow a}^k)^{-1} = \sum_{b \neq a} (\hat{v}_{b \rightarrow i}^k)^{-1}, \quad (6)$$

$$\beta(\zeta_{i \rightarrow a}^k)^{-1} \theta_{i \rightarrow a}^k = \sum_{b \neq a} (\hat{v}_{b \rightarrow i}^k)^{-1} \hat{m}_{b \rightarrow i}^k, \quad (7)$$

$$m_{i \rightarrow a}^{k+1} = \mathbb{E}_{x_i \sim \pi(\beta, \lambda, \theta_{i \rightarrow a}^k, \zeta_{i \rightarrow a}^k)}[x_i], \quad (8)$$

$$v_{i \rightarrow a}^{k+1} = \text{Var}_{x_i \sim \pi(\beta, \lambda, \theta_{i \rightarrow a}^k, \zeta_{i \rightarrow a}^k)}[x_i], \quad (9)$$

$$(10)$$

where

$$\pi(\beta, \lambda, \theta, \zeta) \propto \exp\left\{-\beta\left(\frac{(x - \theta)^2}{2\zeta} + \lambda|x|\right)\right\}. \quad (11)$$

The expression for $\hat{\gamma}_{a \rightarrow i}^k$ turns out to be a bit simpler:

$$\hat{\gamma}_{a \rightarrow i}^k(x) \propto \int_{\mathbb{R}^{d-1}} \underbrace{\exp\left\{\frac{\beta}{2}(y_a - \langle A_a, x \rangle)^2\right\}}_{\psi_a(x_{\partial a})} \times \underbrace{\exp\left\{-\sum_{j \neq i} \frac{(x_j - m_{j \rightarrow a}^k)^2}{2v_{j \rightarrow a}^k}\right\}}_{\prod_{j \in \partial a \setminus i} \rho_{j \rightarrow a}^k(x_j)} dx_{\setminus i}.$$

This is a Gaussian integration, which gives a Gaussian density.

$$\hat{\gamma}_{a \rightarrow i}^k(x) \propto \exp\left\{-\sum_{j \neq i} \frac{(x_j - \hat{m}_{a \rightarrow i}^{k+1})^2}{2\hat{v}_{a \rightarrow i}^{k+1}}\right\}$$

Because of that the result is explicit: mean and variance are affine functions of means and variances of incoming messages correspondingly.

$$A_{ai} \cdot \hat{m}_{a \rightarrow i}^k = y_a - \sum_{j \neq i} A_{aj} m_{j \rightarrow a}^k, \quad (12)$$

$$A_{ai}^2 \cdot \hat{v}_{a \rightarrow i}^k = \sum_{j \neq i} A_{aj}^2 v_{j \rightarrow a}^k + \frac{1}{\beta}. \quad (13)$$

1.2 LASSO with infinite temperature

The formulas above give the AMP algorithm with finite β . This is sufficient in Bayesian setting with the corresponding prior. However, for the vanilla LASSO we need to send β to infinity. Recall the definition of π in Equation 11. Let's write out the following for convenience:

$$\frac{d}{dx} \left(\frac{(x - \theta)^2}{2\zeta} + \lambda|x| \right) = \begin{cases} \zeta^{-1}(x - \theta - \lambda\zeta), & x < 0, \\ \zeta^{-1}(x - \theta + \lambda\zeta), & x > 0. \end{cases}$$

By Laplace approximation we have for expectation

$$\begin{aligned} \lim_{\beta \rightarrow \infty} \mathbb{E}_{X \sim \pi(\beta, \lambda, \theta, \zeta)}[X] &= \arg \min_x \left\{ \frac{(x - \theta)^2}{2\zeta} + \lambda|x| \right\} \\ &= \begin{cases} \theta + \lambda\zeta, & \text{if } \theta < -\lambda\zeta, \\ \theta - \lambda\zeta, & \text{if } \theta > \lambda\zeta, \\ 0, & \text{else;} \end{cases} \\ &= \eta(\theta; \lambda\zeta) \text{ --- soft-threshold.} \end{aligned}$$

When it comes to variance, we look at the second derivative at the minimum:

$$\frac{d^2}{dx^2} \left(\frac{(x - \theta)^2}{2\zeta} + \lambda|x| \right)_{x=\eta(\theta; \lambda\zeta)}.$$

We see that that second derivative is ζ^{-1} if $|\theta| > \lambda\zeta$, but the function is not twice differentiable at the minimum for $|\theta| \leq \lambda\zeta$, i.e. the second derivative is infinite. Thus we obtain

$$\begin{aligned} \lim_{\beta \rightarrow \infty} \text{Var}_{X \sim \pi(\beta, \lambda, \theta, \zeta)}[X] \times \beta &= \begin{cases} \zeta, & \text{if } |\theta| > \lambda\zeta, \\ 0, & \text{if } |\theta| < \lambda\zeta; \end{cases} \\ &= \zeta \cdot \eta'(\theta; \lambda\zeta). \end{aligned}$$

In view of these asymptotic relations, we **change notation**:

$$\begin{aligned} m_{i \rightarrow a}^k &\leftarrow \lim_{\beta \rightarrow \infty} m_{i \rightarrow a}^k, & \hat{m}_{i \rightarrow a}^k &\leftarrow \lim_{\beta \rightarrow \infty} \hat{m}_{a \rightarrow i}^k \\ v_{i \rightarrow a}^k &\leftarrow \lim_{\beta \rightarrow \infty} \beta v_{i \rightarrow a}^k, & \hat{v}_{i \rightarrow a}^k &\leftarrow \lim_{\beta \rightarrow \infty} \beta \hat{v}_{a \rightarrow i}^k \end{aligned}$$

Plugging this into equations 6–9, 12, 13 gives

$$\begin{aligned}
(\zeta_{i \rightarrow a}^k)^{-1} &= \sum_{b \neq a} (\hat{v}_{b \rightarrow i}^k)^{-1}, \\
(\zeta_{i \rightarrow a}^k)^{-1} \theta_{i \rightarrow a}^k &= \sum_{b \neq a} (\hat{v}_{b \rightarrow i}^k)^{-1} \hat{m}_{b \rightarrow i}^k, \\
m_{i \rightarrow a}^{k+1} &= \eta(\theta_{i \rightarrow a}^k; \lambda \zeta_{i \rightarrow a}^k), \\
v_{i \rightarrow a}^{k+1} &= \zeta_{i \rightarrow a}^k \cdot \eta'(\theta_{i \rightarrow a}^k; \lambda \zeta_{i \rightarrow a}^k), \\
A_{ai} \cdot \hat{m}_{a \rightarrow i}^k &= y_a - \sum_{j \neq i} A_{aj} m_{j \rightarrow a}^k, \\
A_{ai}^2 \cdot \hat{v}_{a \rightarrow i}^k &= \sum_{j \neq i} A_{aj}^2 v_{j \rightarrow a}^k + 1.
\end{aligned}$$

So far we treated matrix \mathbf{A} as a given design matrix. To further analyze the AMP for LASSO we will need to impose some assumptions on it. Eventually we would like to have $A_{ai} \sim_{i.i.d.} \mathcal{N}(0, 1/n)$. However, for now we assume $A_{ai} \sim_{i.i.d.} \text{Unif}(\{\pm 1/\sqrt{n}\})$, so we can replace A_{ai}^2 by $1/n$. The intuition for this is that by universality the distribution of A_{ai} doesn't matter as long we match the first and the second moments. An alternative explanation is that the objects of interest (such as $\sum_{j \neq i} A_{aj}^2 v_{j \rightarrow a}^k$) should concentrate around their expectations, and therefore only the moments of A_{ai} matter.

Let's make one more slight change of variables:

$$\begin{aligned}
z_{a \rightarrow i}^k &\equiv A_{ai} \cdot \hat{m}_{a \rightarrow i}^k = y_a - \sum_{j \neq i} A_{aj} m_{j \rightarrow a}^k, \\
\tau_{a \rightarrow i}^k &\equiv A_{ai}^2 \hat{v}_{a \rightarrow i}^k = \sum_{j \neq i} A_{aj}^2 v_{j \rightarrow a}^k + 1.
\end{aligned}$$

Our MP equations become:

$$\theta_{i \rightarrow a}^k = \frac{\sum_{b \neq a} A_{bi} z_{b \rightarrow i}^k / \tau_{b \rightarrow i}^k}{\frac{1}{n} \sum_{b \neq a} 1 / \tau_{b \rightarrow i}^k}, \quad (14)$$

$$m_{i \rightarrow a}^{k+1} = \eta(\theta_{i \rightarrow a}^k; \lambda \zeta_{i \rightarrow a}^k), \quad (15)$$

$$z_{a \rightarrow i}^k = y_a - \sum_{j \neq i} A_{aj} m_{j \rightarrow a}^k, \quad (16)$$

$$\zeta_{i \rightarrow a}^k = \left(\frac{1}{n} \sum_{b \neq a} 1 / \tau_{b \rightarrow i}^k \right)^{-1}, \quad (17)$$

$$v_{i \rightarrow a}^{k+1} = \zeta_{i \rightarrow a}^k \cdot \eta'(\theta_{i \rightarrow a}^k; \lambda \zeta_{i \rightarrow a}^k), \quad (18)$$

$$\tau_{a \rightarrow i}^k = \frac{1}{n} \sum_{j \neq i} v_{j \rightarrow a}^k + 1. \quad (19)$$

These give the Message Passing algorithm for LASSO. Note that since we have $i \in V$, $a \in F$, $|V| = n$ and $|F| = d$, the number of messages is of order dn .

2 From message passing to approximate message passing.

Our goal in this section is to simplify MP ($2 \times n \times d$ # of messages). After simplification there are going to be only d messages. First we will introduce the main idea that leads to such simplification in Section 2.1. However, the straightforward application of this idea does not lead to a good approximation. We will then show how to relax it and make it work in Section 2.2

2.1 Crude approximation of MP

Observe that the source of such a large number of messages is a non-backtracking property. For example, we only need to track the dependence of $\theta_{i \rightarrow a}^k$ on a because we don't want to account for the message from a in its update rule. A simple idea that comes to mind is to give up on this property and to hope that adding one extra term to each summation will not influence much. This leads to the following crude approximation:

$$\begin{aligned}\theta_i^k &= \frac{\sum_b A_{bi} z_b^k / \tau_b^k}{\frac{1}{n} \sum_b 1 / \tau_b^k}, \\ m_i^{k+1} &= \eta(\theta_i^k; \lambda \zeta_i^k), \\ z_a^k &= y_a - \sum_j A_{aj} m_j^k, \\ \zeta_i^k &= \left(\frac{1}{n} \sum_b 1 / \tau_b^k \right)^{-1}, \\ v_i^{k+1} &= \zeta_i^k \cdot \eta'(\theta_i^k, \lambda \zeta_i^k), \\ \tau_a^k &= \frac{1}{n} \sum_j v_j^k + 1,\end{aligned}$$

From the last equation we see that $\tau_a^k =: \tau^k$ — does not depend on a . This propagates to the equation for ζ_i^k and gives $\zeta_i^k = \tau^k =: \zeta^k$ — does not depend on i . τ^k also cancels in the equation for θ_i^k . After these simplifications we only get 3 equations:

$$\mathbf{m}^{k+1} = \eta(\mathbf{m}^k + \mathbf{A}^\top \mathbf{z}^k; \lambda \zeta^k) \in \mathbb{R}^d, \quad (20)$$

$$\mathbf{z}^k = \mathbf{y} - \mathbf{A} \mathbf{m}^k \text{ (without Onsager term)} \in \mathbb{R}^n, \quad (21)$$

$$\zeta^{k+1} = \zeta^k \times \frac{1}{n} \sum_{i \in [n]} \eta'(\mathbf{m}^k + \mathbf{A} \mathbf{z}^k; \lambda \zeta^k)_i + 1 \in \mathbb{R}. \quad (22)$$

It turns out that the approximation above is too crude, so it needs a correction. The equations for ζ , v and τ (the variables that correspond to variances of the messages) are fine. The equations for θ , m , z (variables that correspond to means) are off by $O(1)$. We cannot drop non-backtracking property in equations for those variables.

2.2 Derivation of AMP from MP

As we mentioned before, we only want to drop the non-backtracking property in equations for ζ , v and τ . When we replace $\tau_{i \rightarrow a}^k$ by τ^k , $\zeta_{i \rightarrow a}^k$ by ζ^k and v_i^k by v_k , the equations 17, 18 and 19 can be combined into a single update rule

$$\zeta^{k+1} = \zeta^k \times \frac{1}{n} \sum_{i \in [n]} \eta'(\mathbf{m}^k + \mathbf{A} \mathbf{z}^k; \lambda \zeta^k)_i + 1.$$

So far, the equations 14, 15 and 16 become

$$\begin{aligned}\theta_{i \rightarrow a}^k &= \sum_{b \neq a} A_{bi} z_{b \rightarrow i}^k, \\ m_{i \rightarrow a}^k &= \eta(\theta_{i \rightarrow a}^k; \lambda \zeta^k), \\ z_{a \rightarrow i}^k &= y_a - \sum_{j \neq i} A_{aj} m_{j \rightarrow a}^k,\end{aligned}$$

(equations for m and z remain unchanged yet, in the equation for θ we cancelled τ^k and plugged in $(n-1)/n \approx 1$ in the denominator).

Now we would like to deal with the variables θ , m , z . The idea to completely eliminate non-backtracking property does not work for these variables, so let's try to relax it a little bit. For example, instead of fully dropping the dependence of $\theta_{i \rightarrow a}$ on a , let's say that this dependence results in a small correction to the update rule. More precisely, we would like to introduce the following decomposition:

$$\theta_{i \rightarrow a}^k = \theta_i^k + \delta\theta_{i \rightarrow a}^k, \quad (23)$$

$$m_{i \rightarrow a}^k = m_i^k + \delta m_{i \rightarrow a}^k, \quad (24)$$

$$z_{a \rightarrow i}^k = z_a^k + \delta z_{a \rightarrow i}^k, \quad (25)$$

$$(26)$$

where $\delta\theta$, δm and δz are terms of smaller order of magnitude than θ , m and z correspondingly.

Defining $\delta\theta_{i \rightarrow a}^k$ is rather straightforward:

$$\theta_{i \rightarrow a}^k = \underbrace{\sum_b A_{bi} z_{b \rightarrow i}^k}_{\theta_i^k} - \underbrace{A_{ai} z_{a \rightarrow i}^k}_{\delta\theta_{i \rightarrow a}^k}.$$

Note that because of CLT we can think of θ_i^k as of something of order $O(1)$, while $\delta\theta_{i \rightarrow a}^k$ is $O(1/\sqrt{n})$. So it is indeed a small correction. The update rules are as follows:

$$\begin{cases} \theta_i^k = \sum_b A_{bi} z_{b \rightarrow i}^k = \sum_b A_{bi} z_b^k + \sum_b A_{bi} \delta z_{b \rightarrow i}^k, \\ \delta\theta_{i \rightarrow a}^k = -A_{ai} z_{a \rightarrow i}^k = -A_{ai} (z_i^k + \delta z_{a \rightarrow i}^k). \end{cases}$$

The derivation for m is a bit more involved because we need to push small correction $\delta\theta_{i \rightarrow a}^k$ through non-linear function η :

$$\begin{aligned} m_{i \rightarrow a}^{k+1} &= \underbrace{\eta(\theta_i^k; \lambda \zeta^k)}_{m_i^{k+1}} + \underbrace{\eta'(\theta_i^k; \lambda \zeta^k) \cdot \delta\theta_{i \rightarrow a}^k}_{\delta m_{i \rightarrow a}^{k+1}}, \\ \begin{cases} m_i^{k+1} &= \eta(\theta_i^k; \lambda \zeta^k), \\ \delta m_{i \rightarrow a}^{k+1} &= \eta'(\theta_i^k; \lambda \zeta^k) \cdot \delta\theta_{i \rightarrow a}^k. \end{cases} \end{aligned}$$

Finally, we derive the decomposition for z . It is quite similar to that for θ :

$$\begin{aligned} z_{a \rightarrow i}^k &= y_a - \underbrace{\sum_b A_{bj} m_{j \rightarrow a}^k}_{z_a^k} + \underbrace{A_{ai} m_{i \rightarrow a}^k}_{\delta z_{a \rightarrow i}^k} \\ \begin{cases} z_a^k &= y_a - \sum_j A_{aj} m_{j \rightarrow a}^k = y_a - \sum_j A_{aj} m_j^k - \sum_j A_{aj} \delta m_{j \rightarrow a}^k, \\ \delta z_{a \rightarrow i}^k &= A_{ai} m_{i \rightarrow a}^k = A_{ai} (m_i^k + \delta m_{i \rightarrow a}^k). \end{cases} \end{aligned}$$

So far we didn't actually change much: by introducing the small correction terms we didn't decrease the number of updates, but only emphasized which quantities in those updates may be negligible. Now we are in a good shape to make some simplifying approximations. We have already stated that simply neglecting $\delta\theta_{i \rightarrow a}^k$, $\delta m_{i \rightarrow a}^{k+1}$ and $\delta z_{a \rightarrow i}^k$ in the updates of θ_i^k , m_i^{k+1} and z_a^k is too crude. So let's try to neglect those quantities in their own update rules, i.e. let's take

$$\begin{aligned} \delta\theta_{i \rightarrow a}^k &= -A_{ai} z_i^k, \\ \delta m_{i \rightarrow a}^{k+1} &= \eta'(\theta_i^k; \lambda \zeta^k) \cdot \delta\theta_{i \rightarrow a}^k, \\ \delta z_{a \rightarrow i}^k &= A_{ai} m_i^k. \end{aligned}$$

We get

$$\begin{aligned}
m_i^{k+1} &= \eta(\theta_i^k; \lambda \zeta^k), \\
\theta_i^k &= \sum_b A_{bi} z_b^k + \left(\sum_b A_{bi}^2 \right) m_i^k \\
&\approx m_i^k + \sum_b A_{bi} z_b^k, \\
z_a^k &= y_a - \sum_j A_{aj} m_j^k + \sum_j A_{aj}^2 \partial \eta(\theta_j^k; \lambda \zeta^k) z_a^{k-1} \\
&\approx y_a - \sum_j A_{aj} m_j^k + \frac{1}{n} \sum_j \partial \eta(\theta_j^k; \lambda \zeta^k) z_a^{k-1},
\end{aligned}$$

where we also substituted $A_{ai}^2 \rightarrow 1/n$ for every a, i .

Our final system of update rules is

$$\begin{aligned}
\mathbf{m}^{k+1} &= \eta(\mathbf{m}^k + \mathbf{A}^\top \mathbf{z}^k; \lambda \zeta^k) && \in \mathbb{R}^d, \\
\mathbf{z}^k &= \mathbf{y} - \mathbf{A} \mathbf{m}^k + \underbrace{\left[\frac{1}{n} \sum_j \partial \eta(\theta_j^k; \lambda \zeta^k) \right] \mathbf{z}^{k-1}}_{\text{Onsager term}} && \in \mathbb{R}^n, \\
\zeta^{k+1} &= \zeta^k \times \frac{1}{n} \sum_{i \in [n]} \eta'(\mathbf{m}^k + \mathbf{A} \mathbf{z}^k; \lambda \zeta^k)_i + 1 && \in \mathbb{R}.
\end{aligned}$$

It is interesting to compare it to our crude system (Equations 20, 21, 22). The only correction appears in the update rule for \mathbf{z}^k . This correction term is called "Onsager term".