

# Deep Message Passing

Around the middle of the XVII century, a young man called Baruch Spinoza was dealing with a problem tougher than yours: finding the essence of God Himself.

## CONTENTS

I. Introduction	2
II. Learning through message passing	2
III. Belief Propagation	2
A. BP updates	2
B. Approximate Message Passing	3
IV. ArgMax layer	4
A. Approach 1: Jensen Inequality	4
B. Approach 2: Jensen again	4
V. Numerical results	5
A. Experiments on Fashion-MNIST	5
1. Good parameters: batch-size 1 and batch-size 128	5
2. Varying batch-size: computational performance	6
B. Experiments on Fashion-MNIST: parameters	8
1. Varying $\rho$	8
2. Varying damping $\psi$	10
3. Varying initial weights	12
4. Varying dataset size	13
5. Varying $maxiters$	14
6. Varying $r$	15
7. Varying the number of layers and the size of the hidden layer	16
C. Experiments on RFM	19

## I. INTRODUCTION

### II. LEARNING THROUGH MESSAGE PASSING

TODO ENTIRE SECTION.

In this deep inference problem, we assume that a signal with prior  $P^{\text{in}}$  is fed to a deep feedforward networks with  $L + 1$  layers of weights  $\mathbf{W}^\ell \in \mathbb{R}^{N_{\ell+1} \times N_\ell}$ ,  $\ell = 0, \dots, L$  and biases  $\mathbf{b}^\ell \in \mathbb{R}^{N_{\ell+1}}$ . The signal is propagated through stochastic neuron layers described by probability distributions  $P^\ell$  conditioned on the preactivations, therefore we have the following Markov chain:

$$\mathbf{x}^0 \sim P^{\text{in}} \quad (\text{II.1})$$

$$\mathbf{x}^{\ell+1} \sim P^{\ell+1} \left( \bullet \mid \mathbf{W}^\ell \mathbf{x}^\ell + \mathbf{b}^\ell \right) \quad \ell = 0, \dots, L \quad (\text{II.2})$$

Only  $\mathbf{y} = \mathbf{x}^{L+1}$  is observed, and the task is to reconstruct the original signal  $\mathbf{x}^0$ . The posterior distribution  $p(\mathbf{x}^{0:L}) = P(\mathbf{x}^{0:L} \mid \mathbf{x}^{L+1} = \mathbf{y})$  reads

$$p(\mathbf{x}^{0:L}) \propto \prod_{\ell=0}^L \prod_{k=1}^{N_{\ell+1}} P_k^{\ell+1} \left( x_k^{\ell+1} \mid \sum_{i=1}^{N_\ell} W_{ki}^\ell x_i^\ell \right) \prod_{k=1}^{N_0} P_k^{\text{in}}(x_k^0), \quad (\text{II.3})$$

Typical channels are given by deterministic elementwise activation function  $f_\ell(z)$  (e.g.  $f_\ell(z) = \text{sign}(z)$  or  $f_\ell(z) = \text{relu}(z) = \max(0, z)$ ), combined with Gaussian additive pre-activation noise with variance  $\sigma^2$ . In such cases we have

$$P_k^\ell(x \mid z) = \int D\xi \delta(x - f_\ell(z + \sigma_k^\ell \xi))$$

We also call  $\alpha_\ell = N_{\ell+1}/N_\ell$  the layer *expansion ratio*.

### III. BELIEF PROPAGATION

#### A. BP updates

The AMP equations have been derived from the BP equation in the Appendix. First we introduce the neuron scalar entropy functions:

$$\varphi_k^0(B, A) = \log \int dx e^{-\frac{1}{2}A^2x^2+Bx} P_k^{\text{in}}(x) \quad (\text{III.1})$$

$$\varphi_k^\ell(B, A, \omega, V) = \log \int dx dz e^{-\frac{1}{2}A^2x^2+Bx} P_k^\ell(x|z) e^{-\frac{(\omega-z)^2}{2V}} \quad \ell = 1, \dots, L \quad (\text{III.2})$$

$$\varphi_k^{L+1}(\omega, V, y) = \log \int dz P_k^{L+1}(y|z) e^{-\frac{(\omega-z)^2}{2V}} \quad (\text{III.3})$$

$$\psi_{ki}^\ell(H, G) = \log \int dw e^{-\frac{1}{2}G^2w^2+Hw} P_{ki}^\ell(w) \quad (\text{III.4})$$

For convenience define  $\varphi_i^{0,t} = \varphi_i^0(B_i^{0,t}, A_i^{0,t})$  and  $\varphi_i^{\ell,t} = \varphi_i^\ell(B_i^{\ell,t}, A_i^{\ell,t}, \omega_i^{\ell-1,t}, V_i^{\ell-1,t})$  and  $\varphi_i^{L+1,t} = \varphi_i^{L+1}(\omega_i^{L,t}, V_i^{L,t}, y_i)$ .

Then, we can decompose the BP update rules in a forward and a backward step.

*Forward pass.* As the initial condition for the iterations, we set to zero the following quantities:  $B_i^{\ell,t=0} = 0$ ,  $A_i^{\ell,t=0} = 0$  and  $g_k^{\ell,t=0} = 0$ . The following iterations hold at time  $t \geq 1$ . In the FORWARD pass, starting from  $\ell = 0$  and up to  $\ell = L$ , we have

$$\hat{x}_{ia \rightarrow k}^{\ell,t} = \partial_B \varphi_{ia \rightarrow k}^\ell \left( B_{ia \rightarrow k}^{\ell,t-1}, A_{ia}^{\ell,t-1}, \omega_{ia}^{\ell-1,t}, V_{ia}^{\ell-1,t} \right) \quad (\text{III.5})$$

$$\Delta_{ia \rightarrow k}^{\ell+1,t} = \partial_B^2 \varphi_{ia \rightarrow k}^{\ell+1,t} \quad (\text{III.6})$$

$$m_{ki \rightarrow a}^{\ell,t} = \partial_H \psi_{ki}^\ell (H_{ki \rightarrow a}^{t-1}, G_{ki}^{t-1}) \quad (\text{III.7})$$

$$\sigma_{ki \rightarrow a}^{\ell,t} = \partial_H^2 \psi_{ki}^\ell (H_{ki \rightarrow a}^{t-1}, G_{ki}^{t-1}) \quad (\text{III.8})$$

$$V_{ka}^{\ell,t} = \sum_i \left( \left( m_{ki \rightarrow a}^{\ell,t} \right)^2 \Delta_{ia \rightarrow k}^{\ell,t} + \Sigma_{ki \rightarrow a}^{\ell,t} (\hat{x}_{ia \rightarrow k}^{\ell,t})^2 + \sigma_{ki \rightarrow a}^{\ell,t} \Delta_{ia \rightarrow k}^{\ell,t} \right) \quad (\text{III.9})$$

$$\omega_{ka \rightarrow i}^{\ell,t} = \sum_{i' \neq i} m_{ki \rightarrow a}^{\ell,t} \hat{x}_{ia \rightarrow k}^{\ell,t} \quad (\text{III.10})$$

Here  $V^\ell$  and  $\omega^\ell$  are computed as a function of the previous layer values  $V^{\ell-1}$  and  $\omega^{\ell-1}$ .

*Backward pass.* In the BACKWARD sweep, starting from  $\ell = L$  and down to  $\ell = 0$ , we have

$$g_{ka \rightarrow i}^{\ell,t} = \partial_\omega \varphi_{ka \rightarrow i}^{\ell+1} \left( B_{ka}^{\ell+1,t}, A_{ka}^{\ell+1,t}, \omega_{ka \rightarrow i}^{\ell,t}, V_{ka}^{\ell,t} \right) \quad (\text{III.11})$$

$$\Gamma_{ka \rightarrow i}^{\ell,t} = -\partial_\omega^2 \varphi_{ka \rightarrow i}^{\ell+1,t} \quad (\text{III.12})$$

$$A_{ia}^{\ell,t} = \sum_k \left( (m_{ki \rightarrow a}^{\ell,t})^2 + \sigma_{ki \rightarrow a}^{\ell,t} \right) \Gamma_{ka \rightarrow i}^{\ell,t} - \sigma_{ki \rightarrow a}^{\ell,t} \left( g_{ka \rightarrow i}^{\ell,t} \right)^2 \quad (\text{III.13})$$

$$B_{ia \rightarrow k}^{\ell,t} = \sum_{k' \neq k} m_{k'i \rightarrow a}^{\ell,t} g_{k'a \rightarrow i}^{\ell,t} \quad (\text{III.14})$$

$$G_{ki}^{\ell,t} = \sum_a \left( (\hat{x}_{ia \rightarrow k}^{\ell,t})^2 + \Delta_{ia \rightarrow k}^{\ell,t} \right) \Gamma_{ka \rightarrow i}^{\ell,t} - \Delta_{ia \rightarrow k}^{\ell,t} \left( g_{ka \rightarrow i}^{\ell,t} \right)^2 \quad (\text{III.15})$$

$$H_{ki \rightarrow a} = \sum_{a' \neq a} \hat{x}_{ia' \rightarrow k}^{\ell,t} g_{ka' \rightarrow i}^{\ell,t} \quad (\text{III.16})$$

Notice that  $A^\ell$  and  $B^\ell$  are computed as a function of the  $A^{\ell+1}, B^{\ell+1}$  of the layer above, with the initial condition given by the output  $\mathbf{x}^{L+1} = \mathbf{y}$  on the top layer.

## B. Approximate Message Passing

*Forward pass.* As the initial condition for the iterations, we set to zero the following quantities:  $B_i^{\ell,t=0} = 0$ ,  $A_i^{\ell,t=0} = 0$  and  $g_k^{\ell,t=0} = 0$ . The following iterations hold at time  $t \geq 1$ . In the FORWARD pass, starting from  $\ell = 0$  and up to  $\ell = L$ , we have

$$\hat{x}_{ia}^{\ell,t} = \partial_B \varphi_{ia}^{\ell,t-} \quad (\text{III.17})$$

$$\Delta_{ia}^{\ell,t} = \partial_B^2 \varphi_{ia}^{\ell,t-} \quad (\text{III.18})$$

$$m_{ki}^{\ell,t} = \partial_H \psi_{ki}^{\ell,t-} \quad (\text{III.19})$$

$$\sigma_{ki}^{\ell,t} = \partial_H^2 \psi_{ki}^{\ell,t-} \quad (\text{III.20})$$

$$V_{ka}^{\ell,t} = \sum_i \left( \left( m_{ki}^{\ell,t} \right)^2 \Delta_{ia}^{\ell,t} + \sigma_{ki}^{\ell,t} (\hat{x}_{ia}^{\ell,t})^2 + \sigma_{ki}^{\ell,t} \Delta_{ia}^{\ell,t} \right) \quad (\text{III.21})$$

$$\omega_{ka}^{\ell,t} = \sum_i m_{ki}^{\ell,t} \hat{x}_{ia}^{\ell,t} + \text{TODO : onsagsize}(x)er \quad (\text{III.22})$$

Here  $V^\ell$  and  $\omega^\ell$  are computed as a function of the previous layer values  $V^{\ell-1}$  and  $\omega^{\ell-1}$ .

*Backward pass.* In the BACKWARD sweep, starting from  $\ell = L$  and up to  $\ell = 0$ , we have

$$g_{ka}^{\ell,t} = \partial_\omega \varphi_{ka}^{\ell+1,t} \quad (\text{III.23})$$

$$\Gamma_{ka}^{\ell,t} = -\partial_\omega^2 \varphi_{ka}^{\ell+1,t} \quad (\text{III.24})$$

$$A_{ia}^{\ell,t} = \sum_k \left( (m_{ki}^{\ell,t})^2 + \sigma_{ki}^{\ell,t} \right) \Gamma_{ka}^{\ell,t} - \sigma_{ki}^{\ell,t} \left( g_{ka}^{\ell,t} \right)^2 \quad (\text{III.25})$$

$$B_{ia}^{\ell,t} = \sum_k m_{ki}^{\ell,t} g_{ka}^{\ell,t} + \text{TODO : onsager} \quad (\text{III.26})$$

$$G_{ki}^{\ell,t} = \sum_a \left( (\hat{x}_{ia}^{\ell,t})^2 + \Delta_{ia} \right) \Gamma_{ka}^{\ell,t} - \Delta_{ia} \left( g_{ka}^{\ell,t} \right)^2 \quad (\text{III.27})$$

$$H_{ki} = \sum_a \hat{x}_{ia}^{\ell,t} g_{ka}^{\ell,t} + \text{TODO : onsager} \quad (\text{III.28})$$

Notice that  $A^\ell$  and  $B^\ell$  are computed as a function of the  $A^{\ell+1}, B^{\ell+1}$  of the layer above, with the initial condition given by the output  $\mathbf{x}^{L+1} = \mathbf{y}$  on the top layer.

#### IV. ARGMAX LAYER

In order to perform multi-class classification, we have perform an argmax operation. Call  $z_k$ , for  $k = 1, \dots, K$ , the Gaussian random variables output of the last layer of the network in correspondence of some input  $\mathbf{x}$ . Assuming the correct label is class  $k$ , the effective partition function corresponding to the output constraint reads

$$Z_{k^*} = \int \prod_k dz_k \mathcal{N}(z_k; \omega_k, V_k) \prod_{k \neq k^*} \theta(z_{k^*} - z_k) \quad (\text{IV.1})$$

$$= \int dz_{k^*} \mathcal{N}(z_{k^*}; \omega_{k^*}, V_{k^*}) \prod_{k \neq k^*} H \left( -\frac{z_{k^*} - \omega_k}{\sqrt{V_k}} \right) \quad (\text{IV.2})$$

This last integral is intractable, therefore we have to resort to approximations.

##### A. Approach 1: Jensen Inequality

Using the Jensen inequality we obtain

$$\phi_{k^*} = \log Z_{k^*} = \log \mathbb{E}_{z \sim \mathcal{N}(\omega_{k^*}, V_{k^*})} \prod_{k \neq k^*} H \left( -\frac{z - \omega_k}{\sqrt{V_k}} \right) \quad (\text{IV.3})$$

$$\geq \sum_{k \neq k^*} \mathbb{E}_{z \sim \mathcal{N}(\omega_{k^*}, V_{k^*})} \log H \left( -\frac{z - \omega_k}{\sqrt{V_k}} \right) \quad (\text{IV.4})$$

Reparameterizing the expectations we have

$$\tilde{\phi}_{k^*} = \sum_{k \neq k^*} \mathbb{E}_{\epsilon \sim \mathcal{N}(0,1)} \log H \left( -\frac{\omega_{k^*} + \epsilon \sqrt{V_{k^*}} - \omega_k}{\sqrt{V_k}} \right) \quad (\text{IV.5})$$

The derivative  $\partial_{\omega_k} \tilde{\phi}_{k^*}$  and  $\partial_{\omega_k}^2 \tilde{\phi}_{k^*}$  that we need can then be estimated by sampling (once?)  $\epsilon$ .

##### B. Approach 2: Jensen again

A further simplification is obtained by applying Jensen inequality again to (IV.5) but in the opposite direction. We have the new effective free energy

$$\tilde{\phi}_{k^*} = \sum_{k \neq k^*} \log \mathbb{E}_{\epsilon \sim \mathcal{N}(0,1)} H \left( -\frac{\omega_{k^*} + \epsilon \sqrt{V_{k^*}} - \omega_k}{\sqrt{V_k}} \right) \quad (\text{IV.6})$$

$$= \sum_{k \neq k^*} \log H \left( -\frac{\omega_{k^*} - \omega_k}{\sqrt{V_k + V_{k^*}}} \right) \quad (\text{IV.7})$$

This gives, for  $k \neq k^*$

$$\partial_{\omega_k} \tilde{\phi}_{k^*} = \begin{cases} -\frac{1}{\sqrt{V_k + V_{k^*}}} GH \left( -\frac{\omega_{k^*} - \omega_k}{\sqrt{V_k + V_{k^*}}} \right) & k \neq k^* \\ \sum_{k' \neq k^*} \frac{1}{\sqrt{V_{k'} + V_{k^*}}} GH \left( -\frac{\omega_{k^*} - \omega_{k'}}{\sqrt{V_{k'} + V_{k^*}}} \right) & k = k^* \end{cases} \quad (\text{IV.8})$$

Notice that  $\partial_{\omega_{k^*}} \tilde{\phi}_{k^*} = -\sum_{k \neq k^*} \partial_{\omega_k} \tilde{\phi}_{k^*}$ . In last formulas we used the definition

$$GH(x) = \frac{G(x)}{H(x)} = \frac{\sqrt{2/\pi}}{\text{erfcx}(x/2)} \quad (\text{IV.9})$$

## V. NUMERICAL RESULTS

In this section we study multi-layer perceptrons (MLPs) on a binary classification task on the Fashion-MNIST dataset (divided arbitrarily in two classes: the first/last 5 classes in the original dataset respectively represent the first/second class in the binary classification task). We perform experiments on the whole dataset and with a MLP with two hidden layers of size 101 (apart from the sections in which we modify the architecture or the dataset size, or otherwise stated).

NB: in all figures the title reports “MNIST” while it should be “Fashion-MNIST”.

$M$  and  $P$  are the same parameter (the dataset size).

### A. Experiments on Fashion-MNIST

We compare the BP family with binary-SGD without biases and without additional parameters for batch normalization. In order to keep the pre-activations of each hidden layer normalized we rescale them by  $\frac{1}{\sqrt{N_{\text{in}}}}$  where  $N_{\text{in}}$  is the size of the previous hidden layer (or the input size in the case of the preactivations afferent to the first hidden layer).

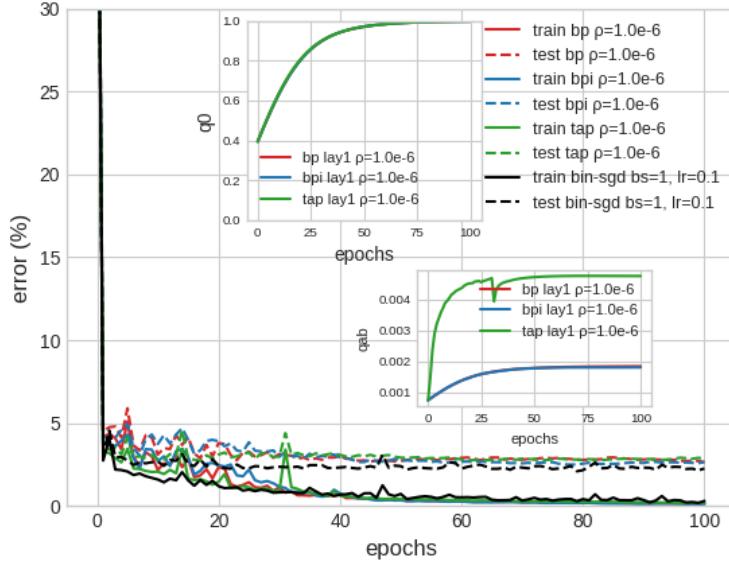
#### 1. Good parameters: batch-size 1 and batch-size 128

We report that it is possible to find values of the hyper-parameters such that the BP family has final train/test error comparable with SGD. In our experiments this holds for generic dataset sizes and batch-sizes, and for generic MLPs (generic depth / hidden layer size).

In summary, up to now the best hyper-parameters (apart from  $P$ , batch-size and  $\rho$  that has to be tuned similarly to the learning rate for SGD) are:

- $\psi = 0.8$
- $\epsilon_{\text{init}} = 1$
- $\text{maxiters} = 1$  ( $r = 0$ )

MNIST 2class, P=60000, K=[784, 501, 501, 501, 1], bs=1,  $\psi=0.8$ , maxit=1,  $\Delta_{init}=1.0$



MNIST 2class, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.8$ , maxit=1,  $\Delta_{init}=1.0$

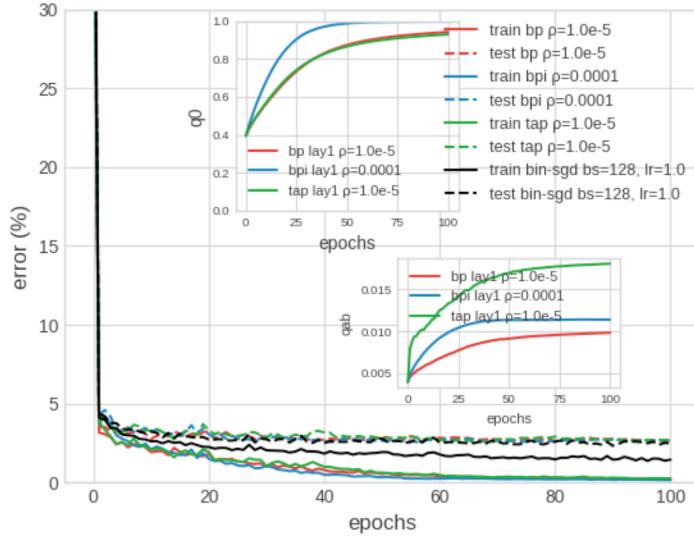


Figure V.1. Left panel: MLP with 3 hidden layers with 501 hidden units each, batch-size=1 on the Fashion-MNIST dataset. Right panel: MLP with 2 hidden layers with 101 hidden units each, batch-size=128 on the Fashion-MNIST dataset. Here we have selected some “good” values of the parameters. **NB: in all figures the upper inset is the self-overlap of the first layer, while the lower inset is the mean overlap of the couples of sub-perceptrons in the first layer.**

## 2. Varying batch-size: computational performance

Here we vary only the batchsize (per completezza farlo per un buon set degli altri iperparametri), in order to compare the performance (time) of BP with binary-SGD (both on GPUs).

The command to reproduce the experiments in this section is:

```
run_experiment(9; M=Int(6e4), batchsize=batchsize, usecuda=true, gpu_id=0, ρ=1+1e-5,
ψ=0.5, lay=lay, epochs=100)
```

with batchsize={1,16,128,1024} and lay={:bp, :bpi, :tap}.

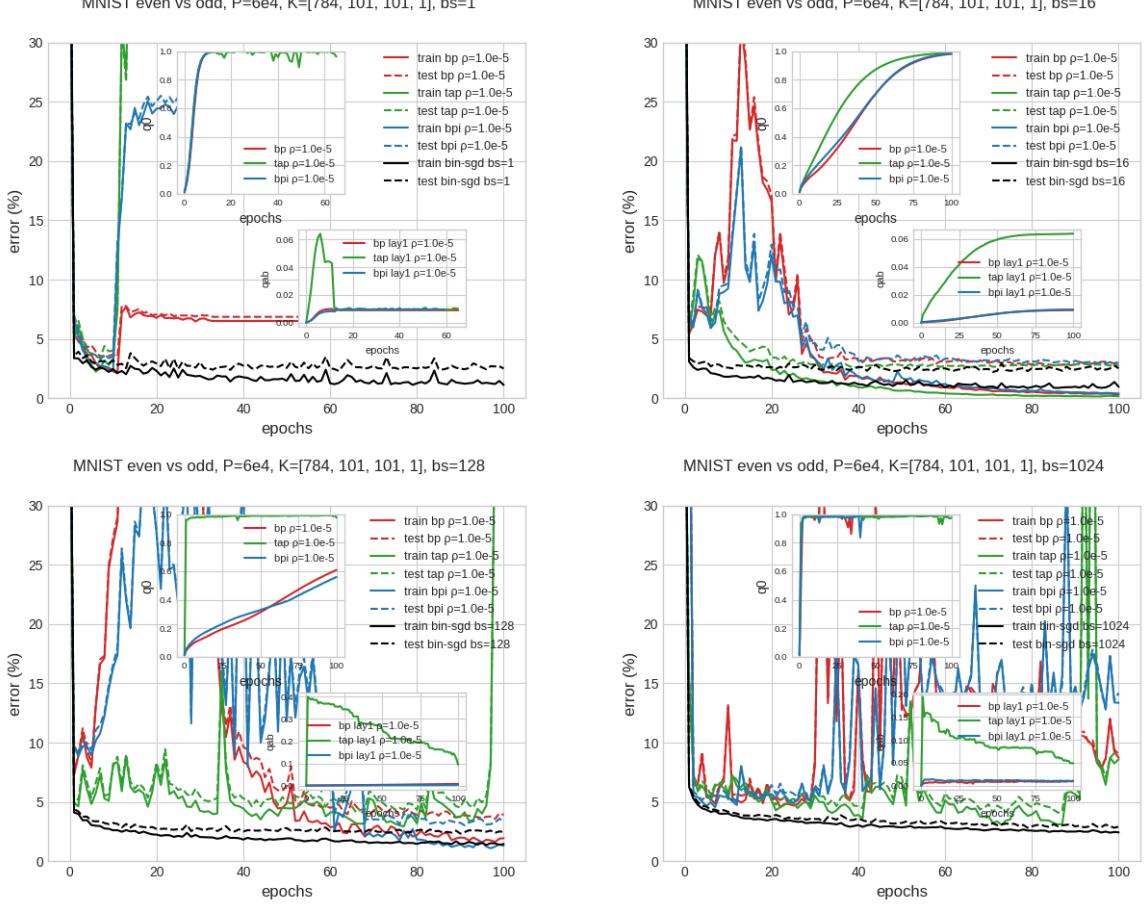


Figure V.2. Comparison of BP, TAP, BPI, SGD varying the batchsize (upper left: bs=1; upper right: bs=16, lower left: bs=128; lower right=1024). The parameter  $\rho - 1$  is fixed in all experiments to  $10^{-5}$ .

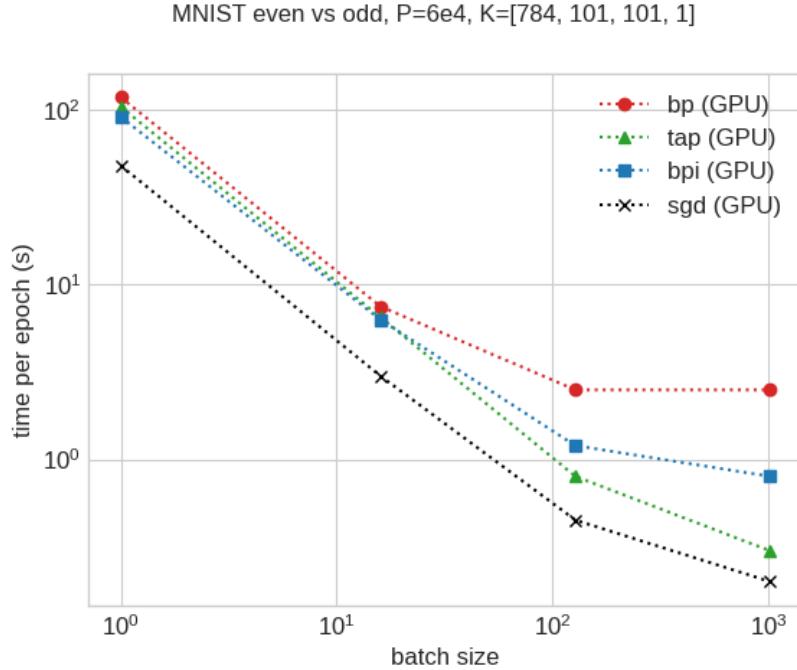


Figure V.3. Algorithms time scaling with the batchsize. The reported time refers to one epoch for each algorithm.

## B. Experiments on Fashion-MNIST: parameters

Here we stick to batchsize 128 and we vary the other hyperparameters. However, we expect that some of the parameters scale with the batchsize (in particular we expect  $\rho - 1 \propto \frac{\text{batchsize}}{\text{data set size}}$ ).

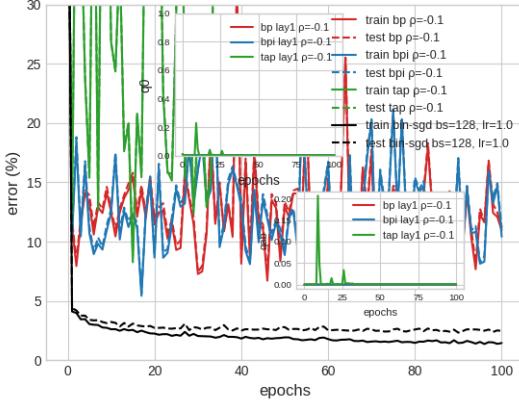
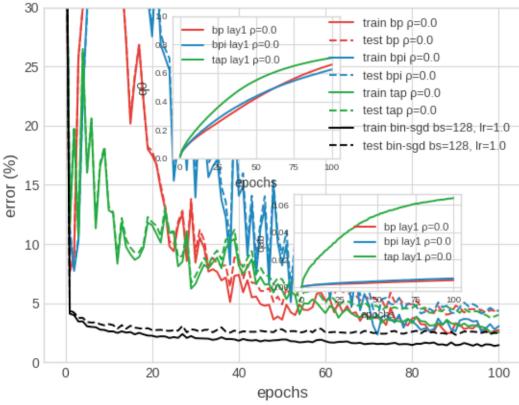
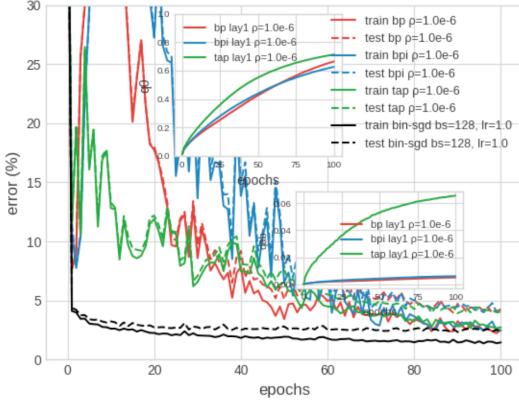
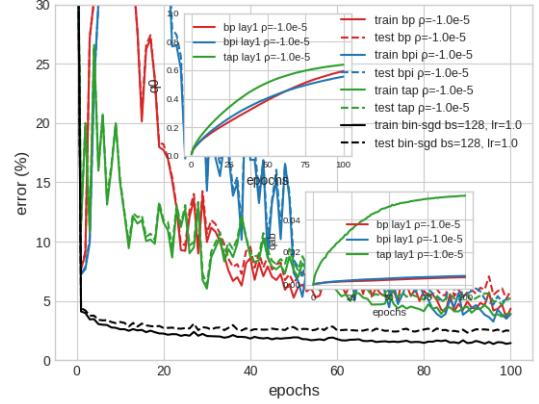
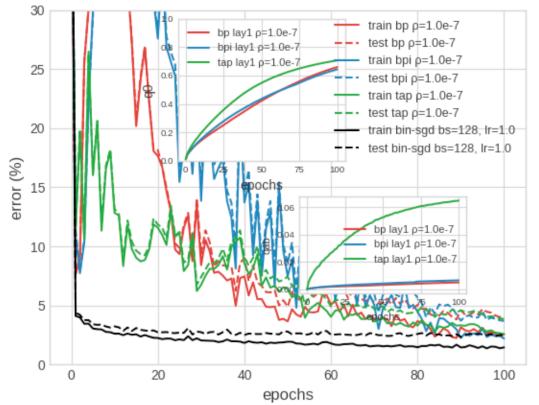
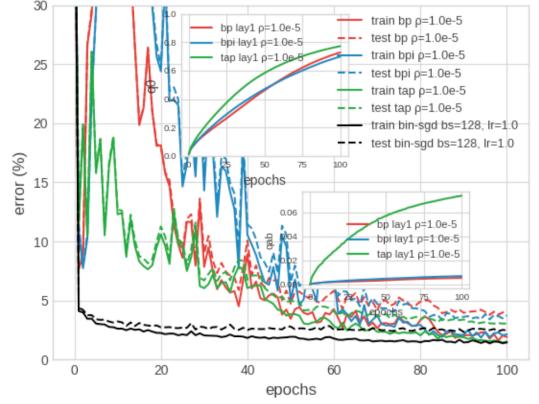
The command to reproduce the experiments is along the lines of:

```
run_experiment(9; usecuda=true, gpu_id=0, epochs=100, lay={:bp, :bpi, :tap}, batchsize=128,
ρ=p[2], ψ=0.5, M=Int(6e4), maxiters=1, r=0., K=[28*28,101,101,1])
```

### 1. Varying $\rho$

$$\rho_s = [-1e-1, -1e-5, 0., 1e-7, 1e-6, 1e-5, 1e-4, 1e-3, 1e-2, 1e-1] .+ 1.$$

For  $bs = 128$  we choose  $\rho - 1 = 10^{-4}$ . The command is the same as the one in the first paragraph of section B except for  $\rho$ .

MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.5$ MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.5$ MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.5$ MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.5$ MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.5$ MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.5$ Figure V.4. MLP with 2 hidden layers with 101 hidden units each, batch-size=128 on the Fashion-MNIST dataset. We vary the parameter  $\rho$ .

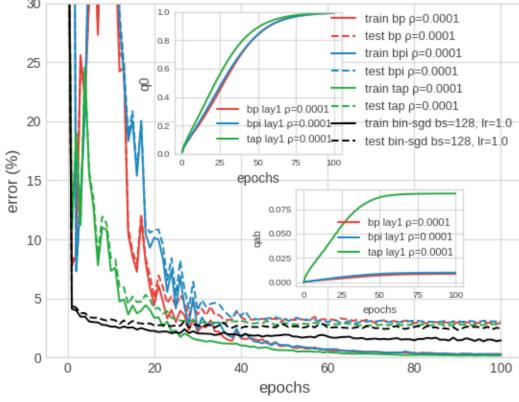
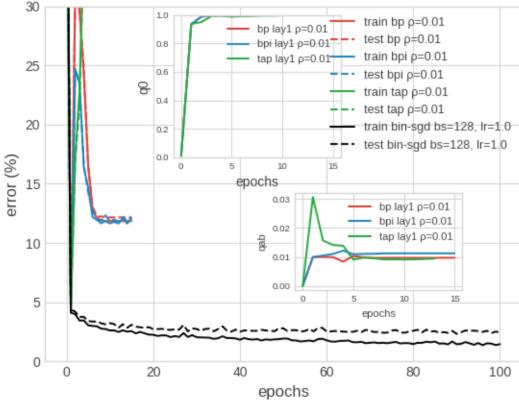
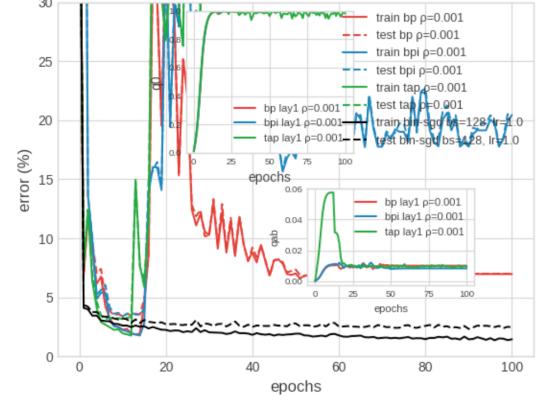
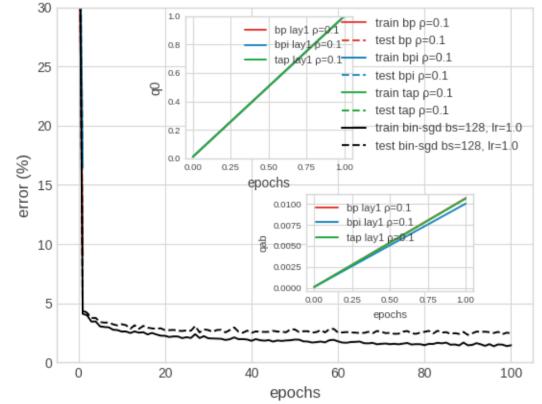
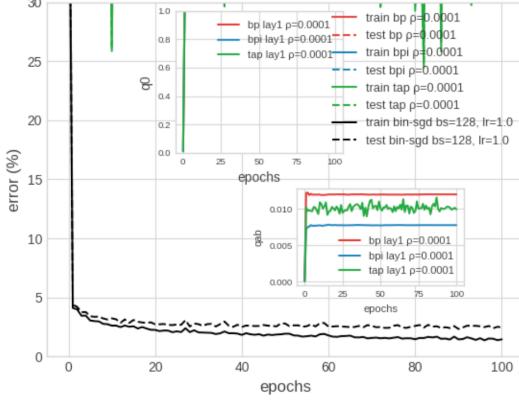
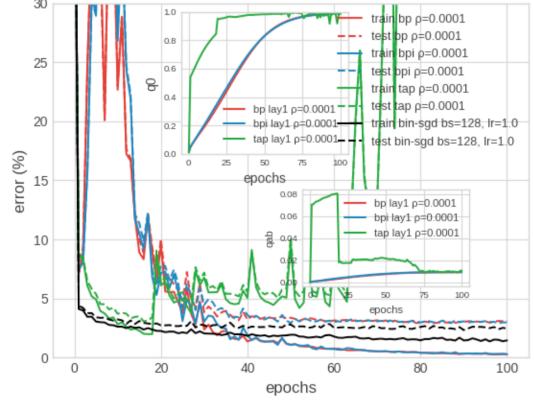
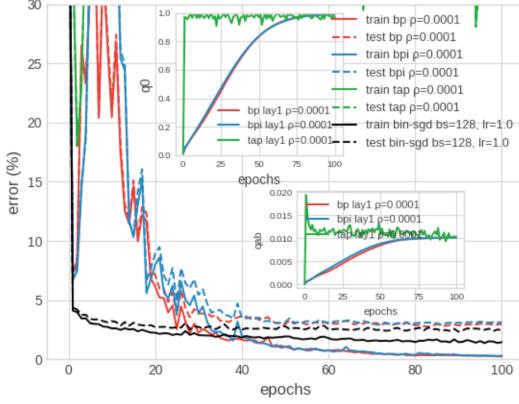
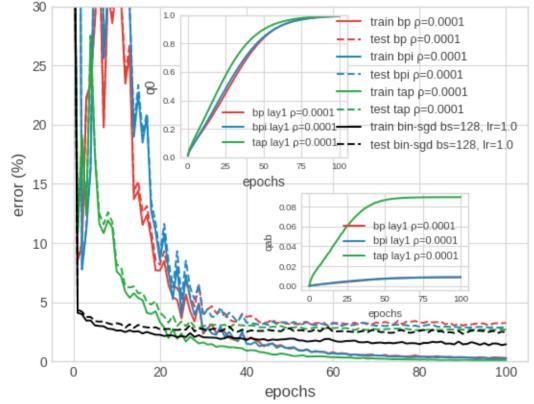
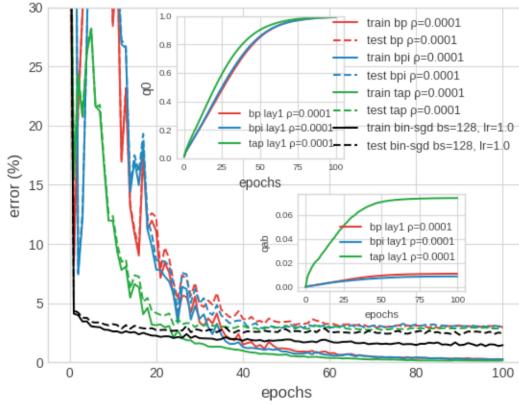
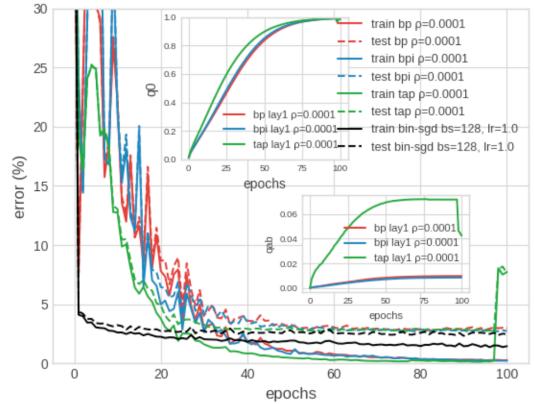
MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.5$ MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.5$ MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.5$ MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.5$ 

Figure V.5. MLP with 2 hidden layers with 101 hidden units each, batch-size=128 on the Fashion-MNIST dataset. We vary the parameter  $\rho$ .

## 2. Varying damping $\psi$

$$\psi_s = [0:0.2:1] \cup [0.9, 0.99, 0.999, 0.9999]$$

For  $bs = 128$  and  $\rho - 1 = 10^{-4}$  we choose  $\psi = 0.8$ . The command is the same as the one in the first paragraph of section B except for  $\rho$  and  $\psi$ .

MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.0$ MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.2$ MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.4$ MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.6$ MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.8$ MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.9$ Figure V.6. MLP with 2 hidden layers with 101 hidden units each, batch-size=128 on the Fashion-MNIST dataset. We vary the parameter  $\psi$ .

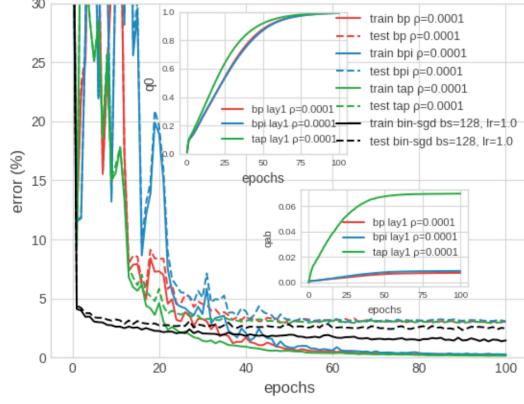
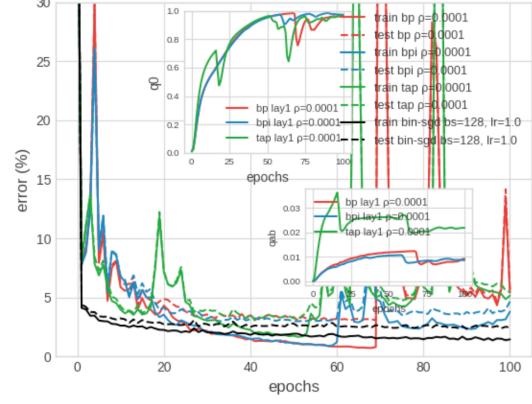
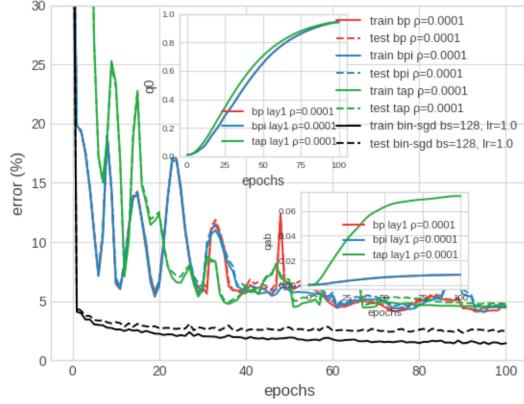
MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.99$ MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.999$ MNIST even vs odd, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.9999$ 

Figure V.7. MLP with 2 hidden layers with 101 hidden units each, batch-size=128 on the Fashion-MNIST dataset. We vary the parameter  $\psi$ .

### 3. Varying initial weights

$\epsilon_{\text{init}} = [0., 1e-3, 1e-2, 1e-1, 5e-1, 1e0]$

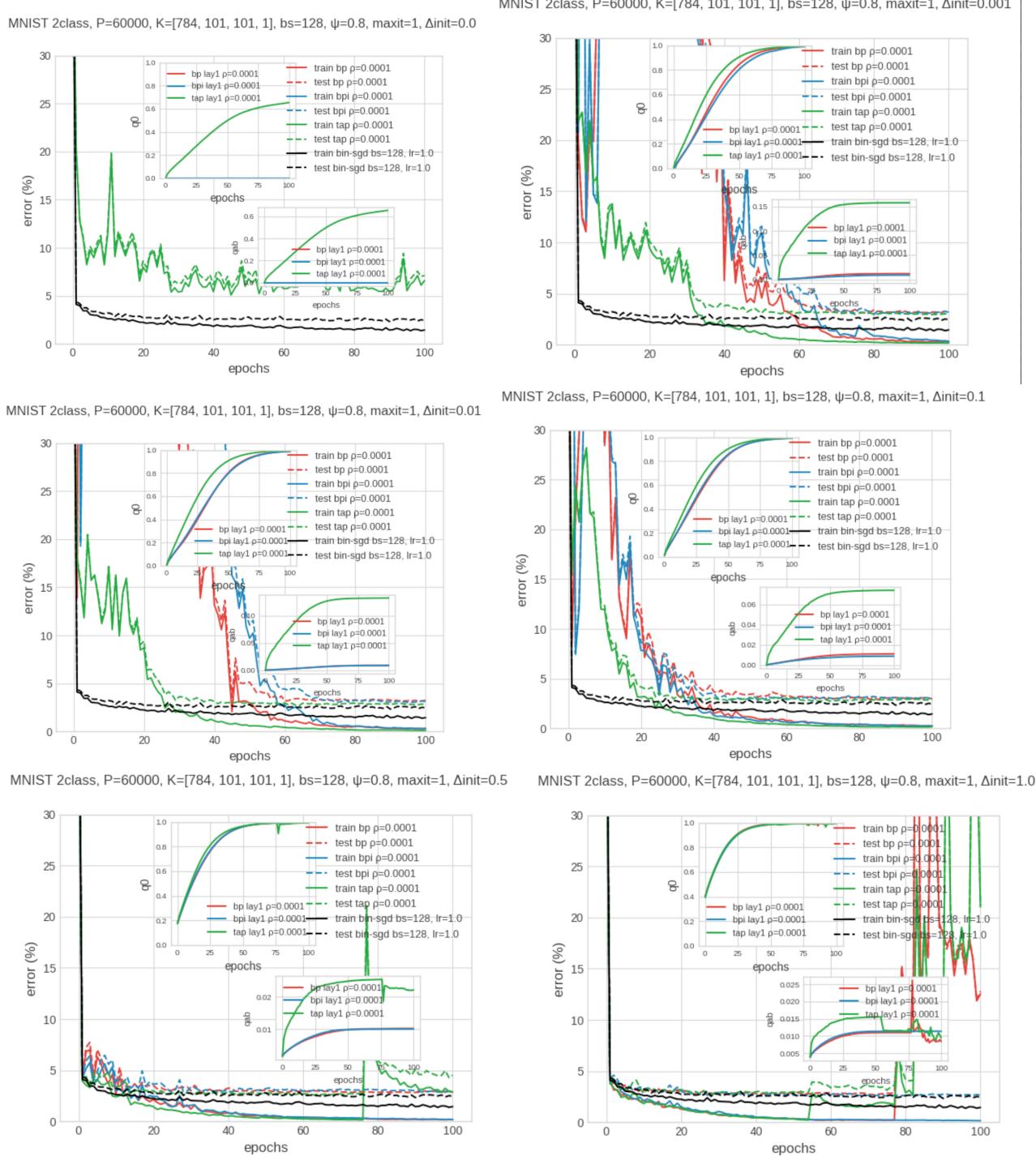


Figure V.8. MLP with 2 hidden layers with 101 hidden units each, batch-size=128 on the Fashion-MNIST dataset. We vary the parameter  $\epsilon\text{init}$ .

#### 4. Varying dataset size

$\text{Ms} = [\text{Int}(1e2), \text{Int}(1e3), \text{Int}(1e4), \text{Int}(6e4)]$   
 $\text{bs} = [\text{Int}(1e0), \text{Int}(1e1), \text{Int}(1e2), \text{Int}(6e2)]$   
 We fix the ratio  $\frac{\text{dataset size}}{\text{batch size}} = \frac{M}{b} = 10^2$ .

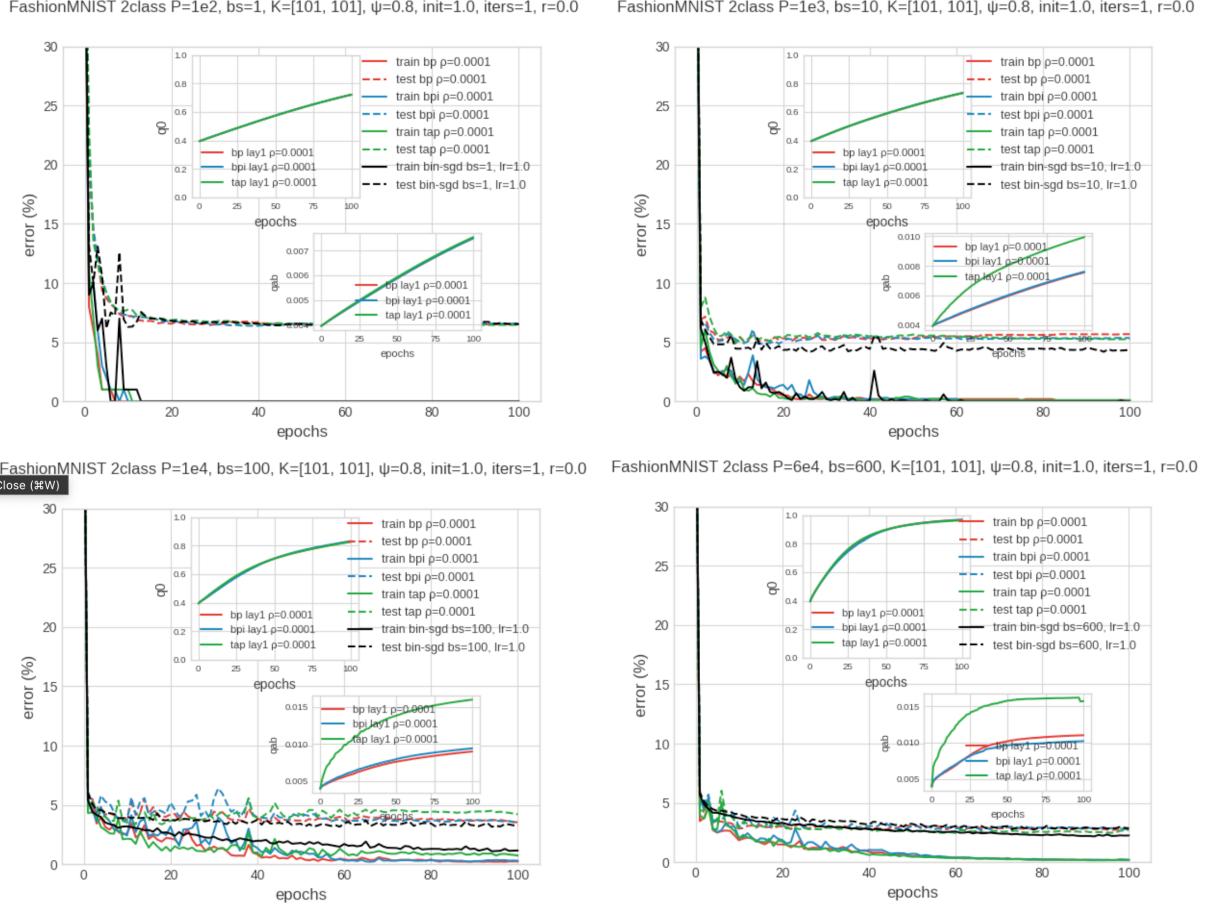
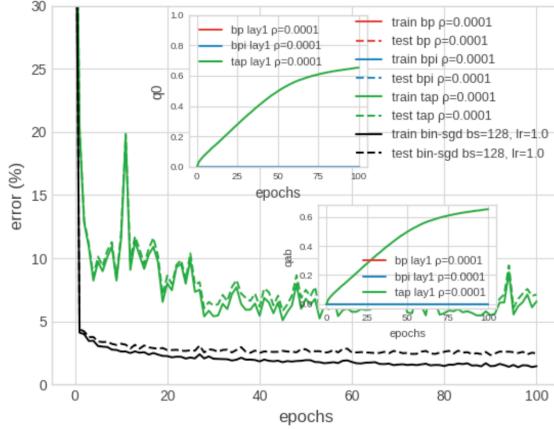


Figure V.9. MLP with 2 hidden layers with 101 hidden units each, batch-size=128 on the Fashion-MNIST dataset. We vary the parameter  $M$  (with  $\frac{M}{b} = 10^2$ ).

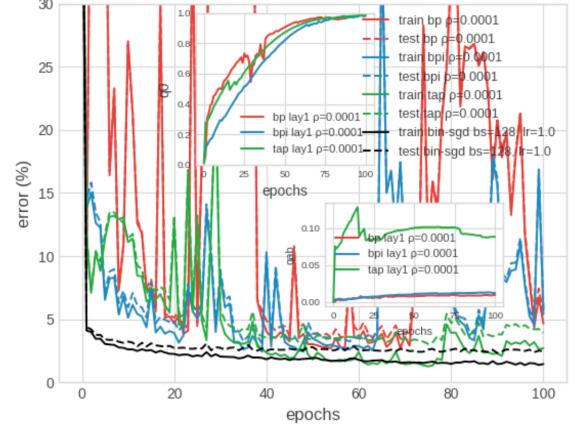
##### 5. Varying maxiters

maxiters = [1, 10, 50, 100]. Here also time is interesting. It is not very favorable in terms of time compared to SGD to choose  $maxiters > 1$ , however it is interesting how many iterations are necessary for convergence.

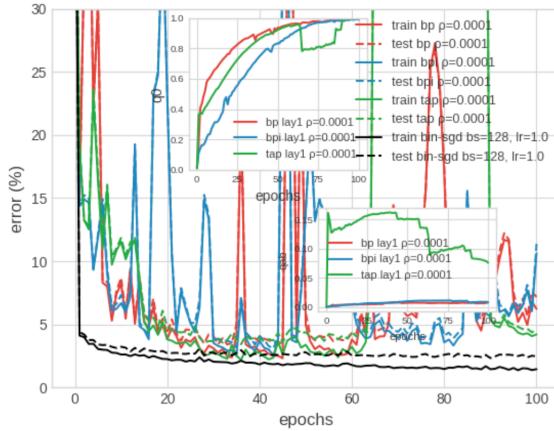
MNIST 2class, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.8$ , maxit=1,  $\Delta\text{init}=0.0$



MNIST 2class, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.8$ , maxit=10,  $\Delta\text{init}=0.0$



MNIST 2class, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.8$ , maxit=50,  $\Delta\text{init}=0.0$



MNIST 2class, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.8$ , maxit=100,  $\Delta\text{init}=0.0$

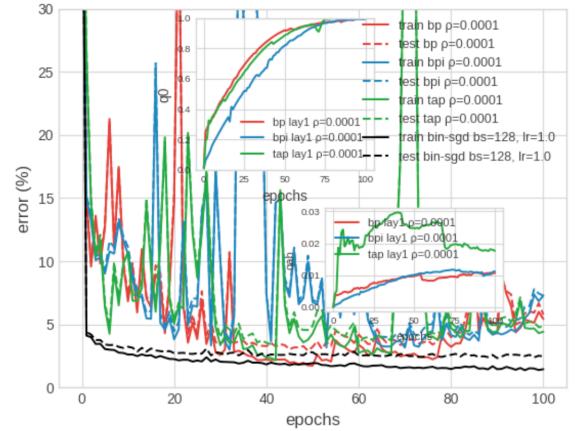


Figure V.10. MLP with 2 hidden layers with 101 hidden units each, batch-size=128 on the Fashion-MNIST dataset. We vary the parameter *maxiters*.

## 6. Varying $r$

$rs = [0:0.2:1.2;]$  (for maxiters=10). Da rifare.

MNIST 2class, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.8$ , maxit=10,  $\Delta\text{init}=0.0$

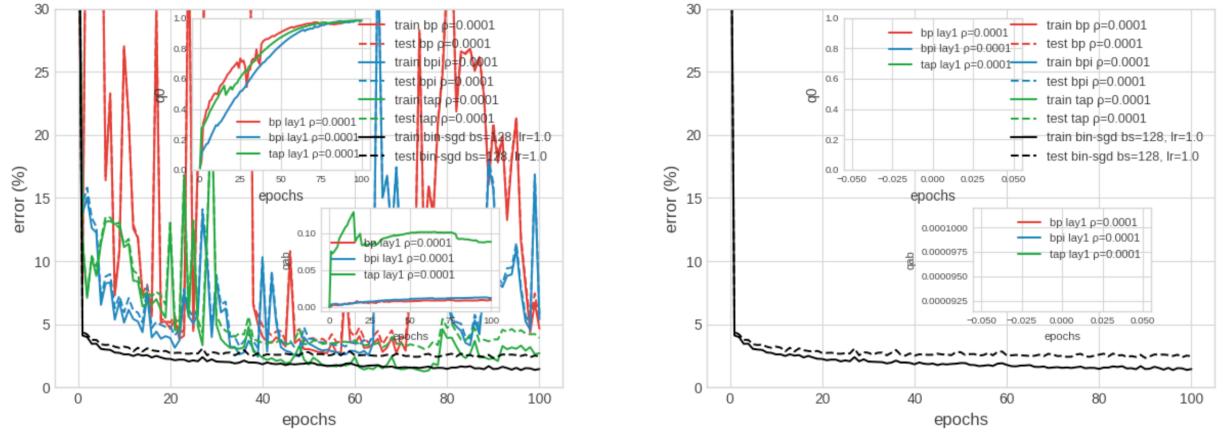
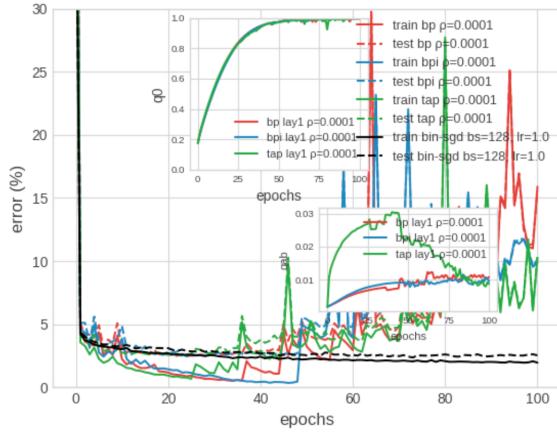


Figure V.11. MLP with 2 hidden layers with 101 hidden units each, batch-size=128 on the Fashion-MNIST dataset. We fix  $\text{maxiters} = 10$  (just to be faster) and we vary the parameter  $r$ . I vari valori di  $r \neq 0$  non hanno proprio funzionato (sembra siano tutti NaN) - gli altri parametri vanno cambiati.

## 7. Varying the number of layers and the size of the hidden layer

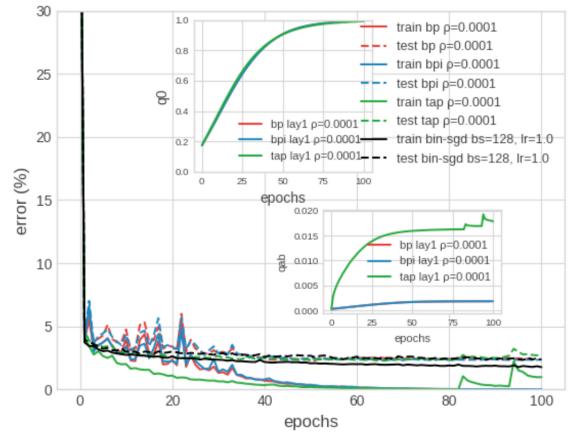
We fix some reasonable values of the hyperparameters (see previous experiments, in particular:  $P = 6e4$ ,  $bs = 128$ ,  $\psi = 0.8$ ,  $\epsilon\text{init} = 1$ ,  $\text{maxiters} = 1$ ,  $r = 0$ ) and want to check if the algorithm converges (and todo the time scaling).

MNIST 2class, P=60000, K=[784, 101, 1], bs=128,  $\psi=0.8$ , maxit=1,  $\Delta_{init}=0.5$



MNIST 2class, P=60000, K=[784, 501, 1], bs=128,  $\psi=0.8$ , maxit=1,  $\Delta_{init}=0.5$

MNIST 2class, P=60000, K=[784, 501, 1], bs=128,  $\psi=0.8$ , maxit=1,  $\Delta_{init}=0.5$



MNIST 2class, P=60000, K=[784, 1001, 1], bs=128,  $\psi=0.8$ , maxit=1,  $\Delta_{init}=0.5$

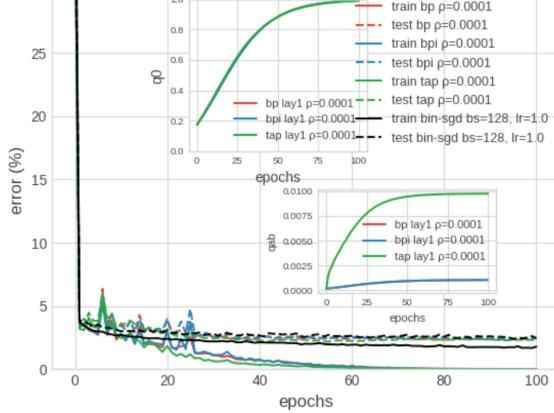
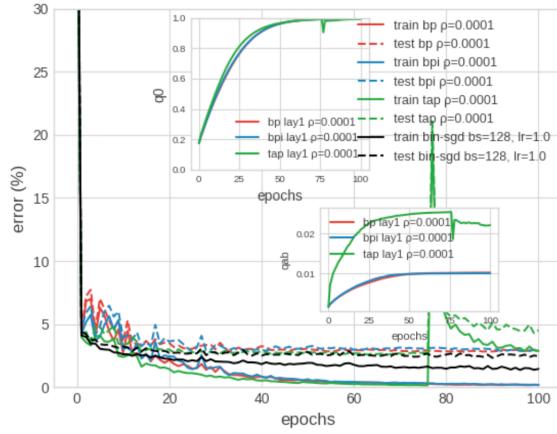
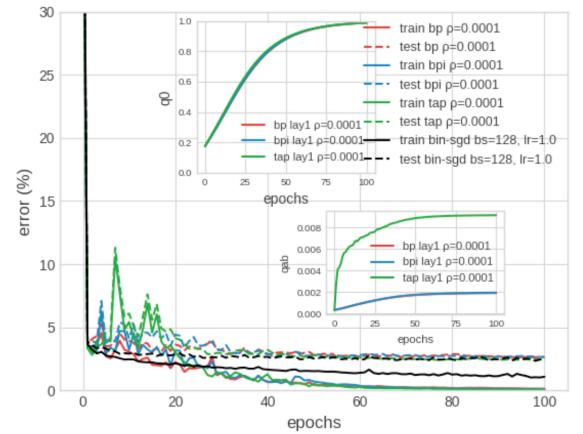


Figure V.12. MLP with 1 hidden layer with 101/501/1001 hidden units each, batch-size=128 on the Fashion-MNIST dataset.

MNIST 2class, P=60000, K=[784, 101, 101, 1], bs=128,  $\psi=0.8$ , maxit=1,  $\Delta\text{init}=0.5$



MNIST 2class, P=60000, K=[784, 501, 501, 1], bs=128,  $\psi=0.8$ , maxit=1,  $\Delta\text{init}=0.5$



MNIST 2class, P=60000, K=[784, 1001, 1001, 1], bs=128,  $\psi=0.8$ , maxit=1,  $\Delta\text{init}=0.5$

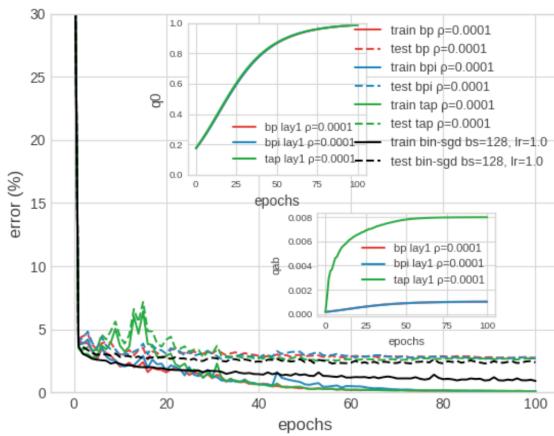
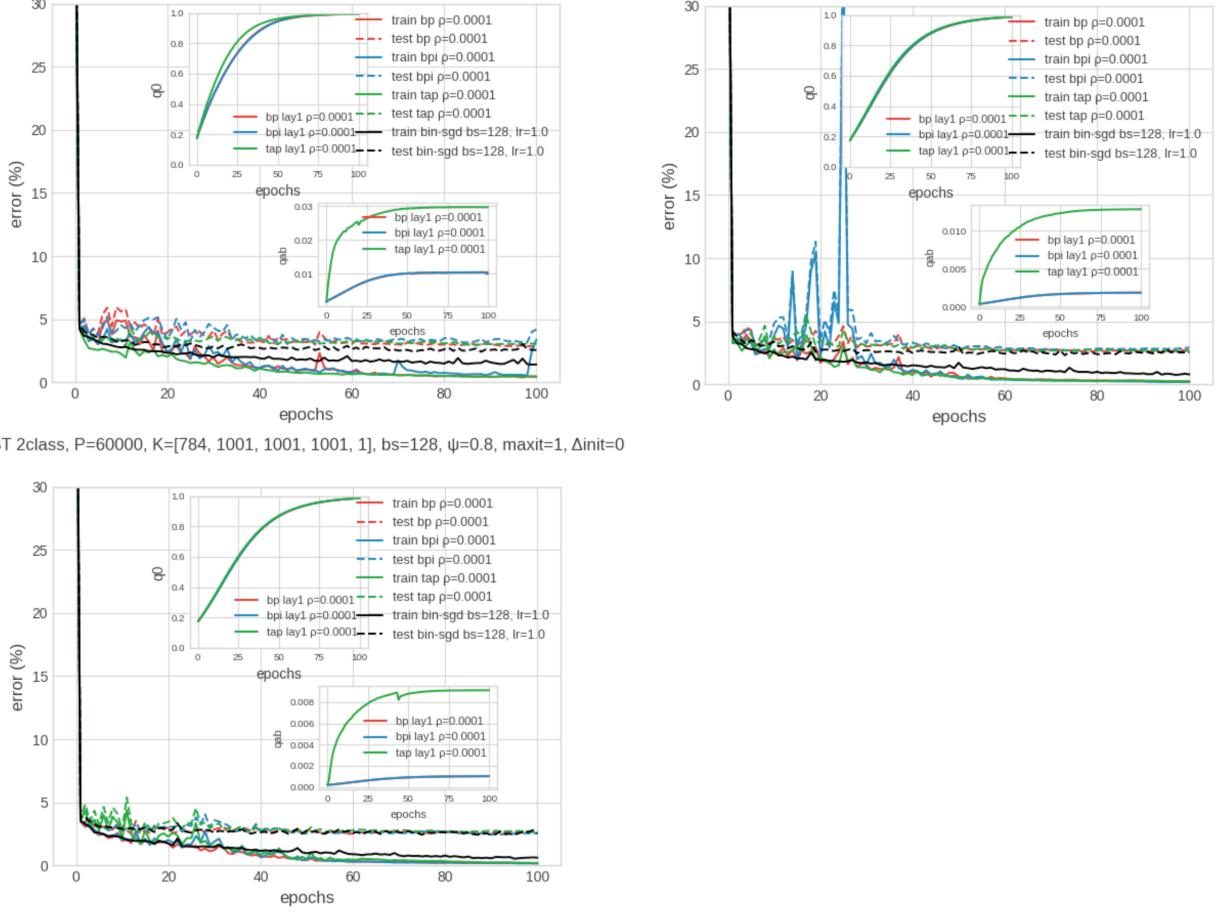


Figure V.13. MLP with 2 hidden layers with 101/501/1001 hidden units each, batch-size=128 on the Fashion-MNIST dataset.

MNIST 2class, P=60000, K=[784, 101, 101, 101, 1], bs=128,  $\psi=0.8$ , maxit=1,  $\Delta\text{init}=0.5$  MNIST 2class, P=60000, K=[784, 501, 501, 501, 1], bs=128,  $\psi=0.8$ , maxit=1,  $\Delta\text{init}=0.5$



Fashion-MNIST, P=60000, K=[784, 1001, 1001, 1001, 1], bs=128,  $\psi=0.8$ , maxit=1,  $\Delta\text{init}=0$

Figure V.14. MLP with 3 hidden layers with 101/501/1001 hidden units each, batch-size=128 on the Fashion-MNIST dataset.

### C. Experiments on RFM

Here we present results with the Random Features Model, concerning in particular the permutation symmetry. In order to investigate the role of the permutation symmetry we present results on the fully connected committee machine (1 hidden layer network learning only the first weight layer, with the weights of the second layer fixed to all ones).

The teacher in the latent space is a perceptron (check).

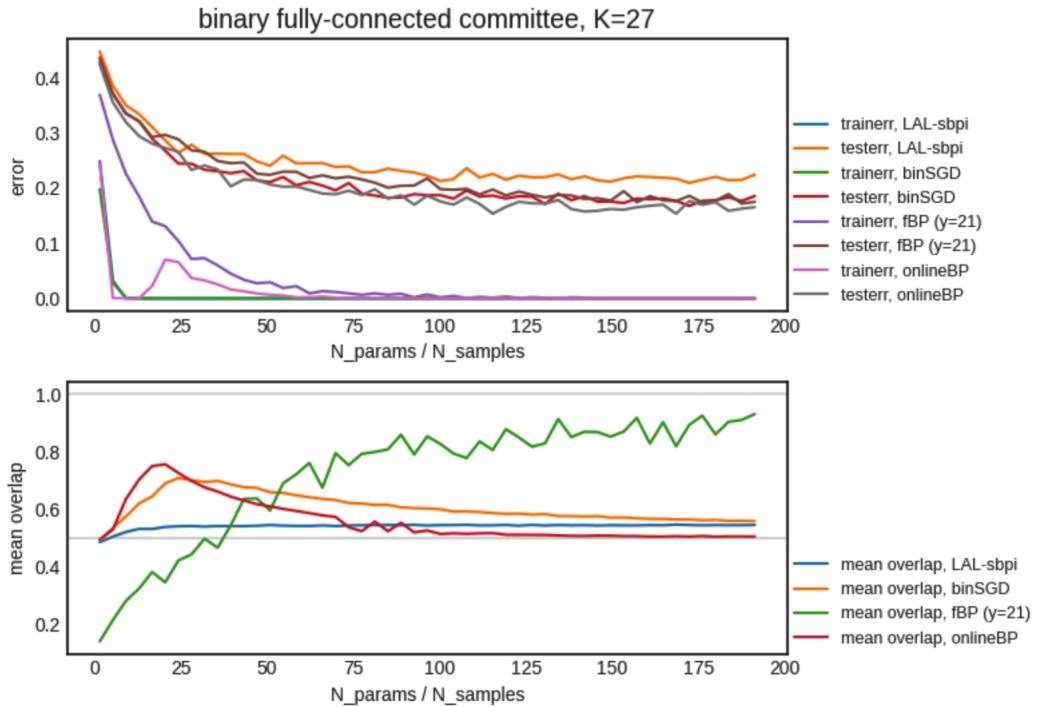


Figure V.15. Overlap varying  $N$  in the RFM, fully connected committee machine with various algorithms.