

Unconstrained optimization

GECD / MIRI - BarcelonaTech

Lab Assignment Pattern recognition with Single Layer Neural Network (SLNN)

v4.0-02/21

F.-Javier Heredia



UNIVERSITAT POLITÈCNICA DE CATALUNYA
BARCELONATECH

Departament d'Estadística
i Investigació Operativa



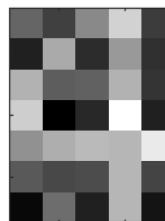
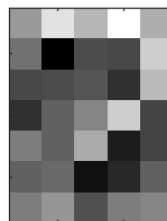
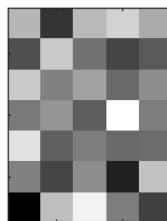
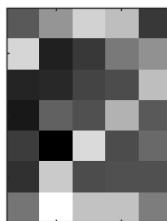
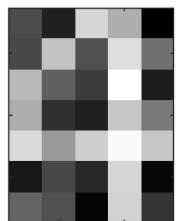
This work is licensed under the Creative Commons Attribution-NonCommercial-NoDerivs 3.0 Unported License. To view a copy of this license, visit <http://creativecommons.org/licenses/by-nc-nd/3.0/>.

Pattern recognition with SLNN

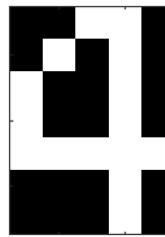
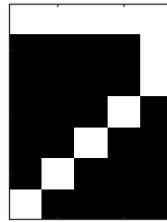
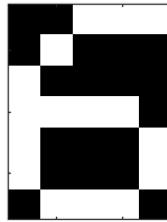
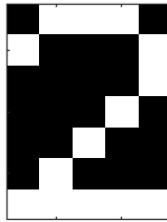
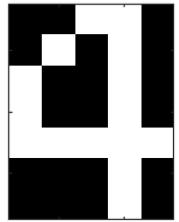
- 1. Presentation.**
- 2. Single Layer Neural Network (SLNN).**
 - Architecture and loss function.
 - Training and testing.
 - Gradient of the loss function.
 - Backtracking line search.
- 3. Pattern recognition with SLNN.**
 - Problem statement and modelling.
 - Generation of training and test data sets.
 - Coding the loss function and gradient.
- 4. Assignment.**
 - Part 1: pattern recognition with GM, CGM and QNM.
 - Part 2: pattern recognition with stochastic gradient.
 - Part 3: computational study of the performance.
 - Report.
- 5. Summary of supporting codes.**

Presentation

- The aim of this project is to develop an application, based on the unconstrained optimization algorithms studied in this course, that allow to recognize the numbers in a sequence of blurred digits:



Original blurred
numbers



Identified by the
neural network

Sequence=42674 ; Identified=42674

- The procedure to achieve that goal will be to formulate a **Single Layer Neural Network** that is going to be **trained to recognize** the different numbers with **First Derivative Optimization methods**.

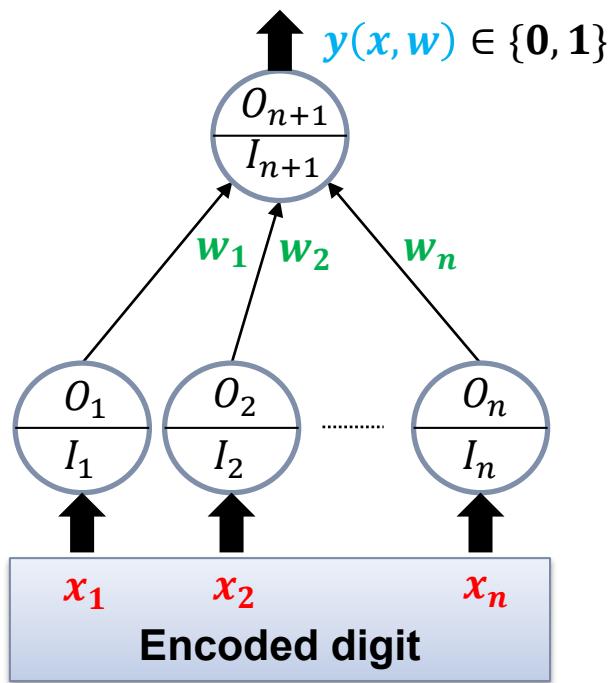
Summary of supporting material

Documents

[uo_nn_v40.pdf](#): this document.

Function/script	Page
<code>uo_nn_main.m</code> : run a single recognition.	link
<code>uo_nn_batch</code> : run a batch of recognitions.	link
<code>function [x,y] = uo_nn_dataset(seed, ncol, target, freq)</code> : generates the dataset <code>x,y</code> .	link
<code>function uo_nn_Xyplot(x,y,w)</code> : plots the dataset <code>x,y</code> . If <code>w</code> is not <code>[]</code> , the plot tells right from wrong predictions of the SLNN.	link
<code>function [alpha,iout] = uo_BLSNW32(f,g,x,d,almax,c1,c2,kBLSmax,epsal)</code> : Algorithm 3.2 of Nocedal & Wright (backtracking line search with SWC and curve fitting).	link

Single layer Neural Network (SLNN): architecture



- Input signal (encoded digit):

$$I_i = \mathbf{x}_i, i = 1, 2, \dots, n ; I_{n+1} = \sum_{i=1}^n \mathbf{w}_i \cdot O_i$$

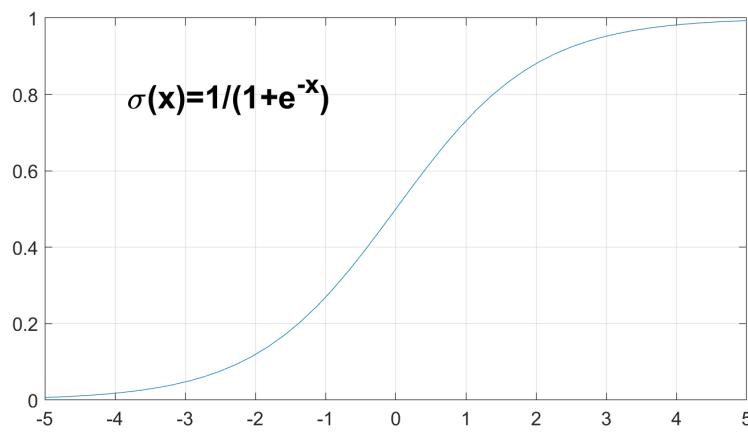
- Activation function (sigmoid function) :

$$O_i = \sigma(I_i) , \quad \sigma(x) = 1/(1 + e^{-x})$$

- Output signal (recognition):

$$\begin{aligned} y(\mathbf{x}, \mathbf{w}) &= \sigma(I_{n+1}) = \sigma\left(\sum_{i=1}^n \mathbf{w}_i O_i\right) = \sigma\left(\sum_{i=1}^n \mathbf{w}_i \cdot \sigma(\mathbf{x}_i)\right) \\ &= \left(1 + e^{-(\sum_{i=1}^n \mathbf{w}_i \cdot \sigma(\mathbf{x}_i))}\right)^{-1} \\ &= \left(1 + e^{-(\sum_{i=1}^n \mathbf{w}_i \cdot (1 + e^{-\mathbf{x}_i})^{-1})}\right)^{-1} \end{aligned}$$

$y(\mathbf{x}, \mathbf{w})$ should be $\begin{cases} = 0 & \text{if } \mathbf{x} \neq \text{target digit} \\ = 1 & \text{if } \mathbf{x} = \text{target digit} \end{cases}$



SLNN: training

- Training data set, size p :

data
model

$$X^{TR} = [x_1^{TR}, x_2^{TR}, \dots, x_p^{TR}] = \begin{bmatrix} x_{1,1}^{TR} & x_{1,2}^{TR} & \cdots & x_{1,p}^{TR} \\ x_{2,1}^{TR} & x_{2,2}^{TR} & \cdots & x_{2,p}^{TR} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1}^{TR} & x_{n,2}^{TR} & \cdots & x_{n,p}^{TR} \end{bmatrix}$$

$$y^{TR} = [y_1^{TR} \quad y_2^{TR} \quad \cdots \quad y_p^{TR}]^T = [1, 0, \dots, 0]^T \text{ (if target=3)}$$

- Mean loss function: for a given (X^{TR}, Y^{TR})

$$L(X^{TR}, y^{TR}) = \min_{w \in \mathbb{R}^n} L(w; X^{TR}, y^{TR}) = \frac{1}{p} \sum_{j=1}^p (y(x_j^{TR}, w) - y_j^{TR})^2$$

- ***L² regularization or weight decay***⁽¹⁾:

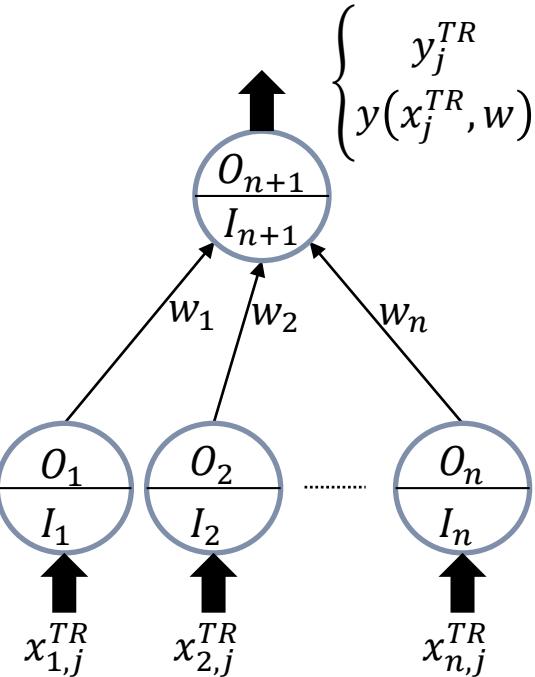
$$\tilde{L}(X^{TR}, y^{TR}, \lambda) = \min_{w \in \mathbb{R}^n} \tilde{L}(w; X^{TR}, y^{TR}, \lambda) = L(w; X^{TR}, y^{TR}) + \lambda \cdot \frac{\|w\|^2}{2}$$

(1): AKA ridge regression or Tikhonov regularization.

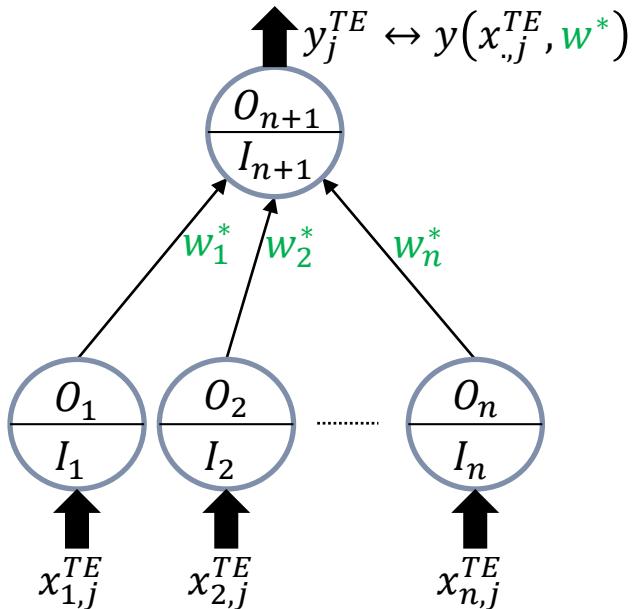
- **Training accuracy (%)**: Accuracy^{TR} = $\frac{100}{p} \cdot \sum_{j=1}^p \delta_{[y(x_j^{TR}, w^*)], \text{round}()}.$

where $\delta_{x,y} = \begin{cases} 1 & \text{if } x = y \\ 0 & \text{if } x \neq y \end{cases}$ (Kronecker delta).

$$w^* = \operatorname{argmin}_{w \in \mathbb{R}^n} \tilde{L}(w; X^{TR}, y^{TR}, \lambda)$$



SLNN : testing



- **Test data set, size q :**

$$X^{TE} = [x_1^{TE}, x_2^{TE}, \dots, x_q^{TE}] = \begin{bmatrix} x_{1,1}^{TE} & x_{1,2}^{TE} & \dots & x_{1,q}^{TE} \\ x_{2,1}^{TE} & x_{2,2}^{TE} & \dots & x_{2,q}^{TE} \\ \vdots & \vdots & \ddots & \vdots \\ x_{n,1}^{TE} & x_{n,2}^{TE} & \dots & x_{n,q}^{TE} \end{bmatrix}$$

$$y^{TE} = [y_1^{TE} \quad y_2^{TE} \quad \dots \quad y_q^{TE}]^T$$

- **Test accuracy (%):**

$$\text{Accuracy}^{TE} = \frac{100}{p} \cdot \sum_{j=1}^p \delta_{[y(x_j^{TE}, \mathbf{w}^*)], y_j^{TE}}$$

- **Overfitting:** if $\text{Accuracy}^{TR} \gg \text{Accuracy}^{TE}$

SLNN : gradient (1/2)

- Mean loss function with regularization (objective function):

$$\tilde{L}(w; X^{TR}, y^{TR}, \lambda) = \frac{1}{p} \sum_{j=1}^p (y(x_j^{TR}, w) - y_j^{TR})^2 + \frac{\lambda}{2} \cdot \sum_{i=1}^n w_i^2$$

- Gradient:

$$\frac{\partial \tilde{L}(w; X^{TR}, y^{TR}, \lambda)}{\partial w_i} = \frac{1}{p} \sum_{j=1}^p 2 \cdot (y(x_j^{TR}, w) - y_j^{TR}) \cdot \frac{\partial y(x_j^{TR}, w)}{\partial w_i} + \lambda \cdot w_i \quad (1)$$

with

$$y(x_j^{TR}, w) = \left(1 + e^{-\left(\sum_{i=1}^n w_i \cdot (1 + e^{-x_{i,j}^{TR}})^{-1} \right)} \right)^{-1} \quad (2)$$

SLNN : gradient (2/2)

- Let us find $\partial y(x_j^{TR}, w) / \partial w_i$:

$$\begin{aligned}
 \frac{\partial y(x_j^{TR}, w)}{\partial w_i} &= \frac{\partial}{\partial w_i} \left(\frac{1}{1 + e^{-\left(\sum_{i=1}^n w_i \cdot (1 + e^{-x_{i,j}^{TR}})^{-1}\right)}} \right)^{-1} = \\
 &= \underbrace{-y(x_j^{TR}, w)^2}_{-\left(1 + e^{-\left(\sum_{i=1}^n w_i \cdot (1 + e^{-x_{i,j}^{TR}})^{-1}\right)}\right)^{-2}} \cdot \underbrace{\left(y(x_j^{TR}, w)^{-1} - 1\right)}_{e^{-\left(\sum_{i=1}^n w_i \cdot (1 + e^{-x_{i,j}^{TR}})^{-1}\right)}} \\
 &\quad \cdot \left(-\left(1 + e^{-x_{i,j}^{TR}}\right)^{-1}\right) = y(x_j^{TR}, w)^2 \cdot \left(y(x_j^{TR}, w)^{-1} - 1\right) \cdot \left(1 + e^{-x_{i,j}^{TR}}\right)^{-1} \\
 &= y(x_j^{TR}, w) \cdot \left(1 - y(x_j^{TR}, w)\right) \cdot \left(1 + e^{-x_{i,j}^{TR}}\right)^{-1}
 \end{aligned}$$

Therefore:

$$\frac{\partial \tilde{L}(w; X^{TR}, y^{TR}, \lambda)}{\partial w_i} = \frac{1}{p} \sum_{j=1}^p 2 \cdot \left(y(x_j^{TR}, w) - y_j^{TR} \right) \cdot y(x_j^{TR}, w) \cdot \left(1 - y(x_j^{TR}, w)\right) \cdot \left(1 + e^{-x_{i,j}^{TR}}\right)^{-1} + \lambda \cdot w_i$$

SLNN: backtracking line search

- The backtracking line search algorithm Alg.BLS cannot handle conveniently the SLNN problem. As a consequence we need to introduce two modifications in the computation of the line search:
 - The maximum step length cannot be a constant for every iteration. Instead, it must be updated dynamically using information of the local behaviour of f near the iterated point at each iteration, using some of the formulas (N&W page 58):

$$\alpha_1^{max} = \alpha^{k-1} \frac{\nabla f^{k-1 T} d^{k-1}}{\nabla f^k T d^k}; \quad \alpha_2^{max} = \frac{2(f^k - f^{k-1})}{\nabla f^k T d^k}.$$

- A BLS based on interpolations must be used (see N&W 3.4), as the one proposed in Alg 3.2 and 3.3 of N&W, implemented in function uo_BLSNW32:

```
function [alpha,iout] =
    uo_BLSNW32(f,g,x,d,almax,c1,c2,kBLSmax,epsal)
```

where **f**, **g**, **d**, **x**, **almax**, **c1**, **c2** are as usual, **iout=0** if the procedure succeeds and:

- ❖ **kBLSmax** is the maximum number of iterations of the BLS algorithm: if exceeded, the algorithm stops with **iout=1**.
- ❖ **epsal** is the minimum variation between two consecutive reductions of α^k , meaning that the algorithm will stop with **iout=2** whenever $|\alpha^{k+1} - \alpha^k| < \text{epsal}$.

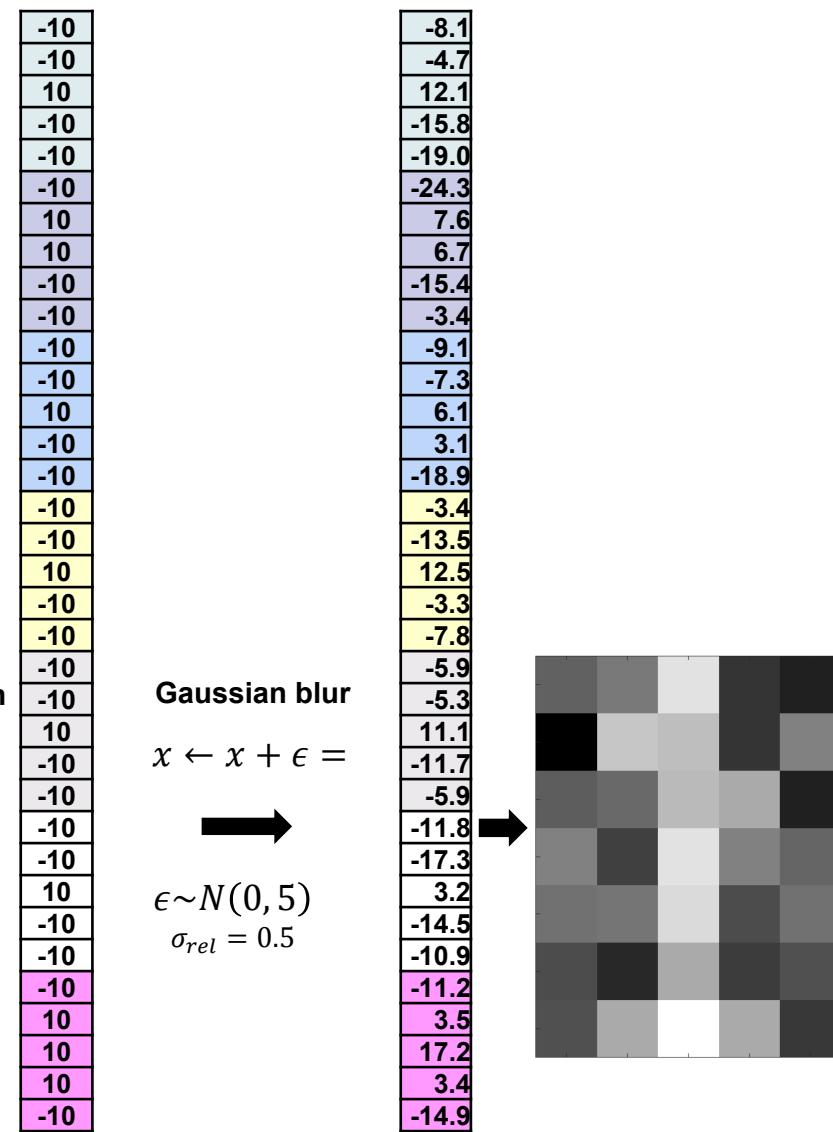
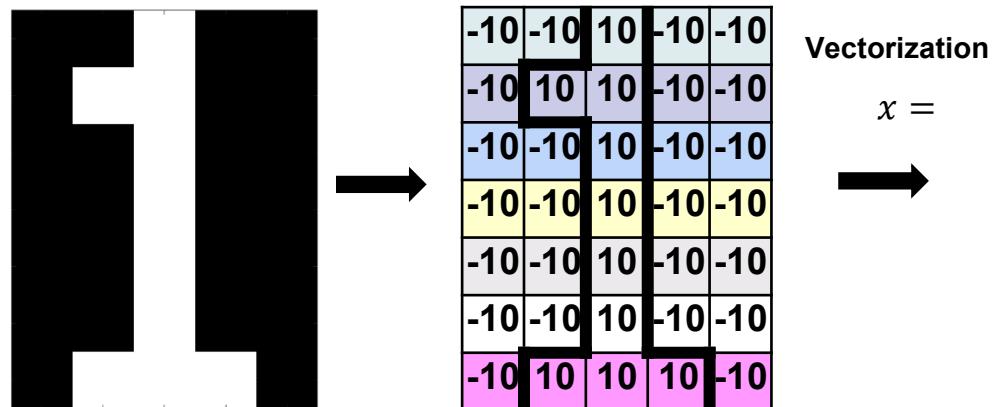


Pattern recognition with SLNN (1/2)

- We are going to use the SLNN to solve a problem of pattern recognition over a small 7x5 pixels matrix picturing the 10 digits:

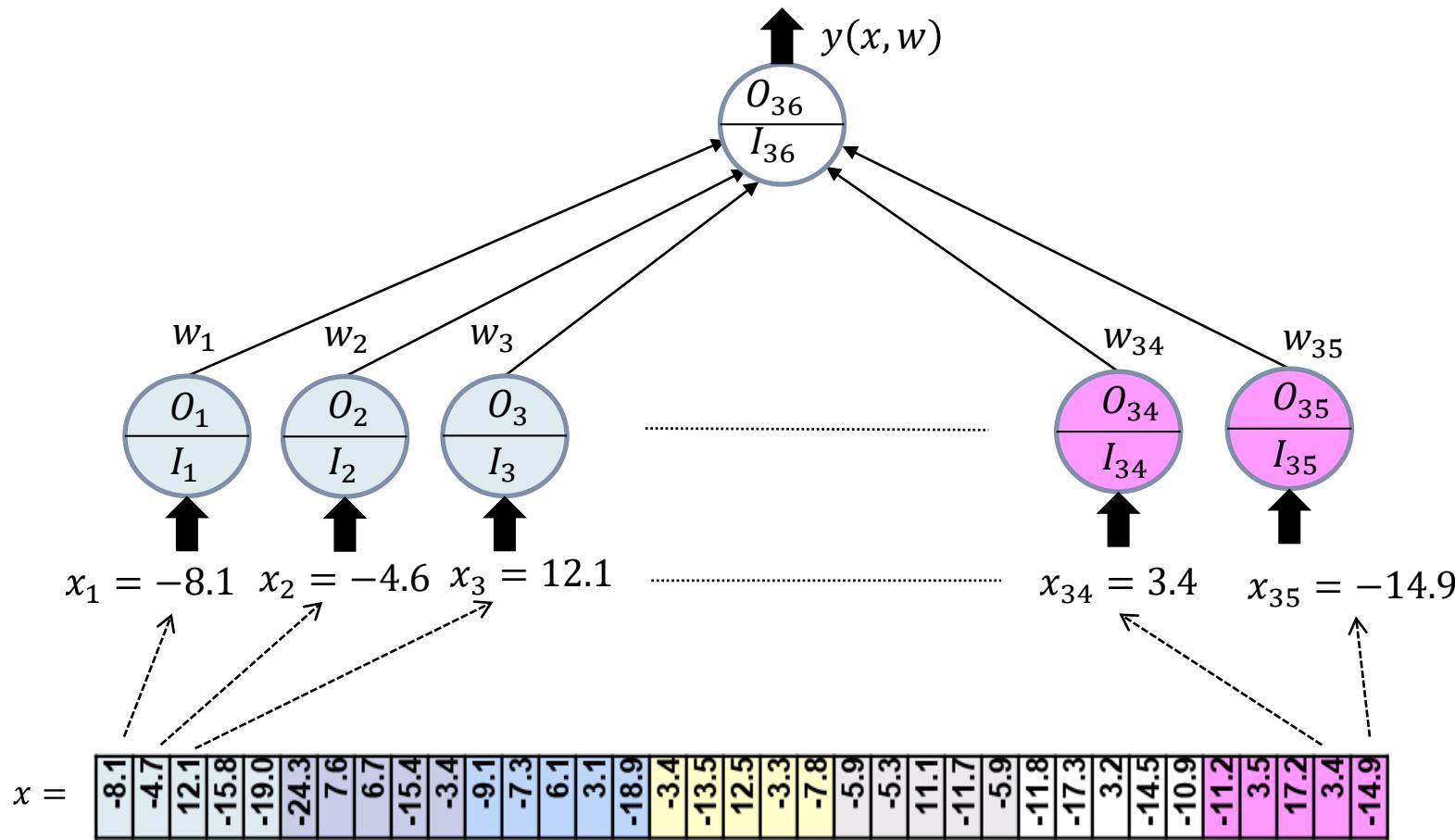


- To obtain the input data of the SLNN x , each white pixel is assigned with a value of 10 and each black pixel with a value of -10 then vectorized and blurred with a Gaussian noise with $\mu = 0$ and $\sigma = \sigma_{rel} \cdot 10$.



Pattern recognition with SLNN (2/2)

- The resulting vectorised and blurred digit x are going to be taken as the inputs of a SLNN:



Training and test data set (1/2)

- The objective of the SLNN is to recognize a set of target numbers, `num_target`, for instance `num_target = [1 3 5 7 9]` will recognize the odd numbers between 0 and 9.
- To this end, the training and test data sets:

$$X^{TR} = [x_1^{TR}, x_2^{TR}, \dots, x_p^{TR}] \equiv \text{xtr}(1:35, 1:\text{tr_p}) \text{ and } y^{TR} \equiv \text{ytr}(1:\text{tr_p})$$

$$X^{TE} = [x_1^{TE}, x_2^{TE}, \dots, x_p^{TE}] \equiv \text{xte}(1:35, 1:\text{te_q}) \text{ and } y^{TE} \equiv \text{yte}(1:\text{te_q})$$

must be generated with the help of function

```
function [x,y] = uo_nn_dataset(seed, number, target, freq)
```

This function will generate a dataset, where:

- x, y** are the generated data sets (`xtr`, `ytr` or `xte`, `yte`).
- seed** is the seed for the Matlab random numbers generator. The numbers in the dataset are randomly choosed, guaranteeing a frequency of the digits in `target` close to `freq`. The value σ_{rel} for each digit is also randomly selected within the range [0.25, 1].
- number** number of columns/elements of array **x/y**.
- target** is the set of digits to be identified (one or several).
- freq** is the frequency of the digits `target` in the data ser. For instance, if `target=[1 2]` and `freq=0.5`, the digits 1 and 2 will be, approximately, half of the total digits in the data set **x**. If `freq≤0.0`, every digit is (approx.) equally represented in the dataset: **this is the value to be used when generating the test dataset X^{TE} , y^{TE}** .



Training and test data set (2/2)

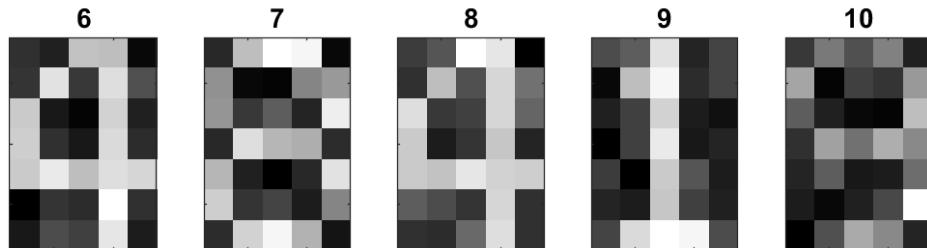
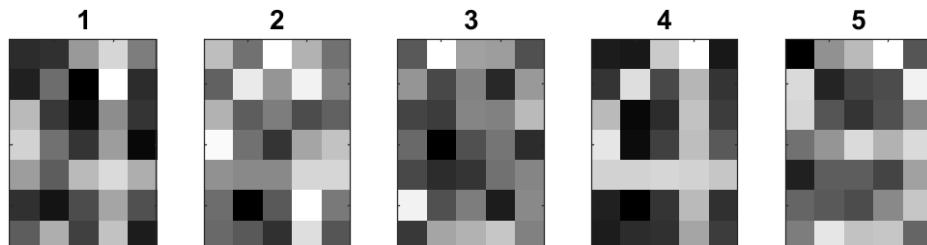
- Exemple: let seed=1234, number=10, target=[4], freq=0.5 then:

and the graphical representation:

```
>> uo nn Xyplot(X,y,[])
```

- **function** `uo_nn_Xyplot(X,y,w)`
plots a set of vectorised digits, and the
recognition brought by a vector w :

- **x** an array of vectorised digits.
 - **y** associated output of the SLNN.
 - **w** vector of weights w (optional).



Loss function and its gradient (1/2)

- Let $\mathbf{x}_{tr}, \mathbf{y}_{tr}$ be the training data set:

```
[Xtr,ytr] = uo_nn_dataset(tr_seed, tr_p, num_target, tr_freq);
```

- Defining the row vector of predictions $y(X^{TR}, w)$ and the sigmoid matrix of inputs $\sigma(X^{TR})$ as

$$y(X^{TR}, w) \stackrel{\text{def}}{=} [y(x_1^{TR}, w), \dots, y(x_p^{TR}, w)]; \quad \sigma(X^{TR}) = \begin{bmatrix} \sigma(x_{11}^{TR}) & \dots & \sigma(x_{1p}^{TR}) \\ \vdots & \ddots & \vdots \\ \sigma(x_{n1}^{TR}) & \dots & \sigma(x_{np}^{TR}) \end{bmatrix}$$

the value of the loss function \tilde{L} and its gradient $\nabla \tilde{L}$ can be expressed as

$$\tilde{L}(w; X^{TR}, y^{TR}, \lambda) = \frac{1}{p} \|y(X^{TR}, w) - y^{TR}\|^2 + \lambda \frac{\|w\|^2}{2}$$

$$\nabla \tilde{L}(w; X^{TR}, y^{TR}, \lambda) = \frac{1}{p} 2\sigma(X^{TR}) \left((y(X^{TR}, w) - y^{TR}) \circ y(X^{TR}, w) \circ (1 - y(X^{TR}, w)) \right)^T + \lambda w$$

where \circ stands for the **element-wise** (or **Hadamard**) product ($\begin{bmatrix} u \\ v \end{bmatrix} \circ \begin{bmatrix} x \\ y \end{bmatrix} = \begin{bmatrix} ux \\ vy \end{bmatrix}$).

The Matlab code using the **element-wise operators** “./” and “.*” is :

σ	<code>sig = @(Xtr) 1./(1+exp(-Xtr));</code>
y	<code>y = @(Xtr,w) sig(w'*sig(Xtr));</code>
\tilde{L}	<code>L = @(w,Xtr,ytr) (norm(y(Xtr,w)-ytr)^2)/size(ytr,2)+ (la*norm(w)^2)/2;</code>
$\nabla \tilde{L}$	<code>gL = @(w,Xtr,ytr) (2*sig(Xtr)*((y(Xtr,w)-ytr).*y(Xtr,w) .* (1-y(Xtr,w)))')/size(ytr,2)+la*w;</code>



Assignment.

There are **three** different parts in this assignment:

- **Part 1:** to develop a code **uo_nn_solve** for the **recognition** of some specific target with GM and QNM).
- **Part 2:** to extend code **uo_nn_solve** with the **stochastic gradient method (SGM)**.
- **Part 3:** to conduct a **comparation of the performance** of the three abovementioned methods, in terms of **global and local convergence** as well as **recognition accuracy**.

Part 1: pattern recognition

- Matlab script to solve the recognition of a given `num_target`

`uo_nn_main.m` : recognition of `num_target` digits.

```

clear;
%
% Parameters for dataset generation
%
num_target =[1];
tr_freq    = .5;
tr_p        = 250;
te_q        = 250;
tr_seed    = 123456;
te_seed    = 789101;
%
% Parameters for optimization
%
la = 0.0;                                % L2 regularization.
epsG = 10^-6; kmax = 10000;                % Stopping criterium.
ils=3; ialmax = 2; kmaxBLS=30; epsal=10^-3;c1=0.01; c2=0.45; % Linesearch (ils=3->uo_BLSNW32)
isd = 1; icg = 2; irc = 2 ; nu = 1.0;       % Search direction.
sg_seed = 565544; sg_a10 = 2; sg_be = 0.3; sg_ga = 0.01;      % SGM iteration.
sg_emax = kmax; sg_ebest = floor(0.01*sg_emax);                 % SGM stopping condition.
%
% Optimization
%
t1=clock;
[Xtr,ytr,wo,fo,tr_acc,Xte,yte,te_acc,niter,tex]=uo_nn_solve(num_target,tr_freq,tr_seed,tr_p,te_seed,te_q,la,epsG,kmax,ils,ialmax,kmaxBLS,epsal,c1,c2,isd,sg_a10,sg_be,sg_ga,sg_emax,sg_ebest,sg_seed,icg,irc,nu);
t2=clock;
fprintf(' wall time = %6.1d s.\n', etime(t2,t1));

```

Part 1: pattern recognition

- Function **uo_nn_solve** called inside **uo_nn_main.m** must solve the instance corresponding to a particular combination of parameters . The actions to be taken inside this function are:
 - i. To generate the training data set (X^{TR}, y^{TR}) and the test dataset (X^{TE}, y^{TE}) with function **uo_nn_dataset.m**.
 - ii. To find the value of w^* minimizing $\tilde{L}(w; X^{TR}, y^{TR}, \lambda)$ **with your own optimization routines developed during the course.**
 - iii. To calculate Accuracy^{TR} and Accuracy^{TE}.

Part 1: example 1, num_target=[1]

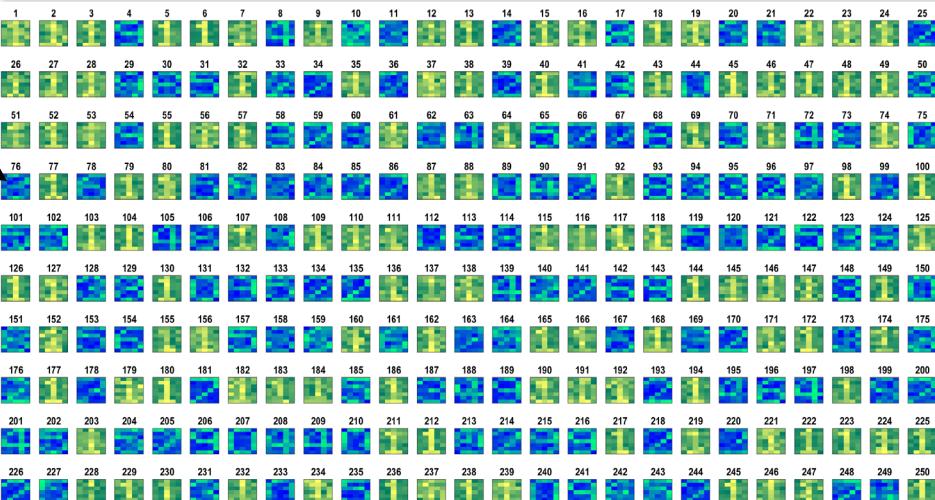
```
[uo_nn_solve] ::::::::::::::::::::
[uo_nn_solve] Pattern recognition with neural networks (OM/GCED).
[uo_nn_solve] 11-Feb-2021 19:34:47
[uo_nn_solve] ::::::::::::::::::::
[uo_nn_solve] Training data set generation.
[uo_nn_solve] num_target = 1
[uo_nn_solve] tr_freq   = 0.50
[uo_nn_solve] tr_p      = 250
[uo_nn_solve] tr_seed   = 123456
[uo_nn_solve] Test data set generation.
[uo_nn_solve] te_freq   = 0.00
[uo_nn_solve] te_q      = 250
[uo_nn_solve] te_seed   = 789101
[uo_nn_solve] Optimization
[uo_nn_solve] L2 reg. lambda = 0.0000
[uo_nn_solve] epsG= 1.0e-06, kmax= 10000
[uo_nn_solve] ils= 3, ialmax= 2, kmaxBLS= 30, epsBLS= 1.0e-03,
[uo_nn_solve] c1= 0.01, c2= 0.45, isd= 1
[uo_nn_solve] w0=[0]
k      al    iW      g'*d      f      ||g|| 
1     1.00e+00  2   -1.18e-01  2.50e-01  3.44e-01
2     3.44e+00  0   -5.81e-02  1.50e-01  2.41e-01
3     5.22e+01  0   -4.37e-03  3.59e-02  6.61e-02
.
.
.
42    3.09e+04  0   -9.68e-12  1.08e-06  3.11e-06
43    2.18e+05  2   -1.47e-12  9.17e-07  1.21e-06
44    6.95e+04  0   -1.75e-12  6.74e-07  1.32e-06
45    6.11e-07
k      al    iW      g'*d      f      ||g|| 
wo=[ -1.1e+00,-1.6e+00,-2.3e-01,-2.4e+00,-1.4e+00
-1.3e+00,+1.5e+00,+2.7e+00,+1.3e-01,-1.4e+00
-1.7e+00,-9.3e-02,+3.0e+00,-7.9e-01,-7.6e-01
-2.1e+00,-1.7e+00,+2.7e+00,-2.4e+00,-1.4e-01
-1.8e+00,-4.8e-01,+2.0e+00,-1.4e+00,-2.0e+00
-9.6e-01,-1.2e+00,+2.2e+00,-1.1e+00,-3.6e-01
-2.5e-01,-2.1e-02,+2.7e-01,+7.4e-01,-7.4e-01
]
[uo_nn_solve] Accuracy.
[uo_nn_solve] tr_accuracy = 100.0
[uo_nn_solve] te_accuracy = 100.0
wall time = 1.6e-01 s.
```

```
>> uo_nn_Xyplot(wo,0,[])
```

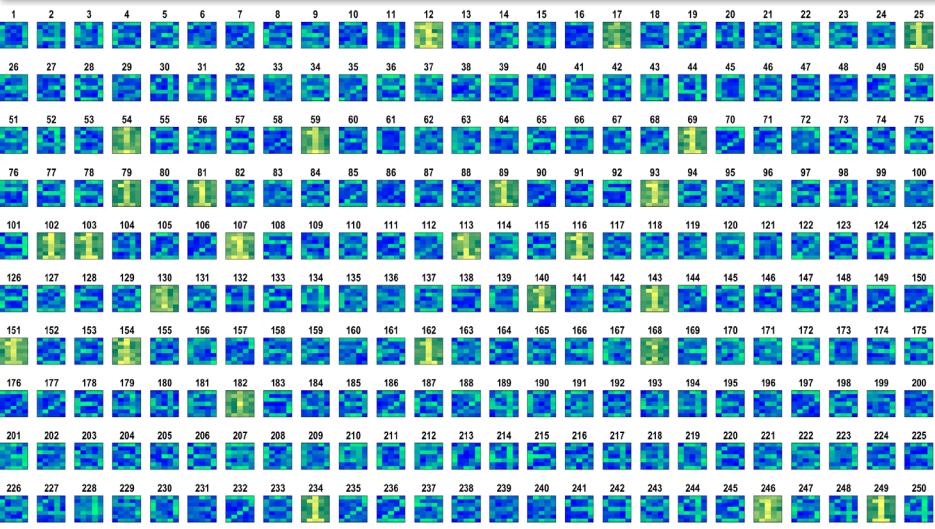
Rigth positive

Rigth negative

```
>> uo_nn_Xyplot(Xtr,ytr,wo)
```



```
>> uo_nn_Xyplot(Xte,yte,wo)
```



Part 1: example 2, num_target=[3]

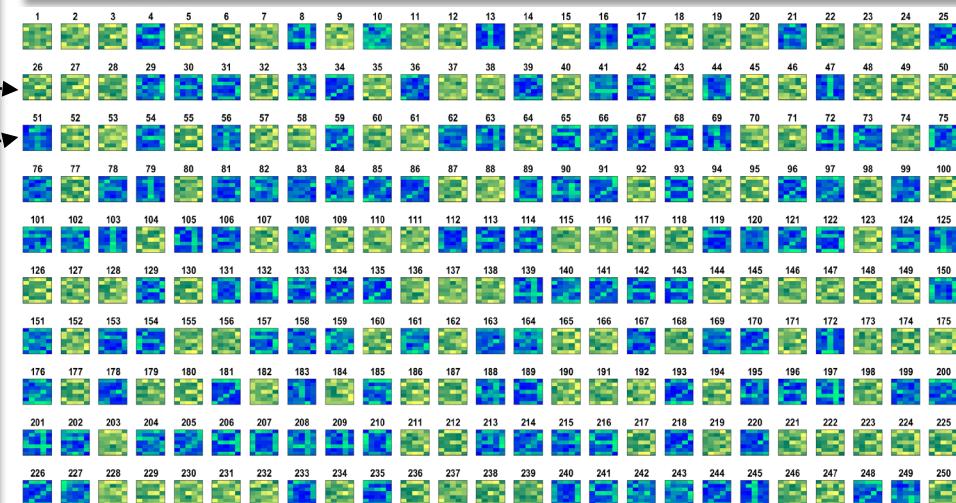
```
[uo_nn_solve] ::::::::::::::::::::
[uo_nn_solve] Pattern recognition with neural networks (OM/GCED).
[uo_nn_solve] 11-Feb-2021 20:03:11
[uo_nn_solve] ::::::::::::::::::::
[uo_nn_solve] Training data set generation.
[uo_nn_solve] num_target = 3
[uo_nn_solve] tr_freq = 0.50
[uo_nn_solve] tr_p = 250
[uo_nn_solve] tr_seed = 123456
[uo_nn_solve] Test data set generation.
[uo_nn_solve] te_freq = 0.00
[uo_nn_solve] te_q = 250
[uo_nn_solve] te_seed = 789101
[uo_nn_solve] Optimization
[uo_nn_solve] L2 reg. lambda = 0.0000
[uo_nn_solve] epsG= 1.0e-06, kmax= 10000
[uo_nn_solve] ils= 3, ialmax= 2, kmaxBLS= 30, epsBLS= 1.0e-03,
[uo_nn_solve] c1= 0.01, c2= 0.45, isd= 1
[uo_nn_solve] w0=[0]
[uo_nn_solve] k al iW g'*d f ||g||
[uo_nn_solve] 1 1.00e+00 2 -3.92e-02 2.50e-01 1.98e-01
[uo_nn_solve] 23 1.41e+01 0 -7.56e-05 1.80e-02 8.70e-03
[uo_nn_solve] 45 1.07e+01 0 -2.36e-05 1.00e-02 4.86e-03
[uo_nn_solve] 2069 9.54e+03 0 -2.50e-12 3.58e-06 1.58e-06
[uo_nn_solve] 2091 1.01e+04 0 -2.22e-12 3.33e-06 1.49e-06
[uo_nn_solve] 2113 1.07e+04 0 -1.97e-12 3.09e-06 1.40e-06
[uo_nn_solve] 2128 2.94e-06 9.93e-07
[uo_nn_solve] k al iW g'*d f ||g||
[uo_nn_solve] wo=[ -8.1e+00,+1.5e+00,-1.1e+01,-3.4e+00,-1.2e+00
[uo_nn_solve] +5.9e+00,-1.2e+00,+5.3e+00,-2.3e+00,-2.0e+00
[uo_nn_solve] -1.0e+01,-1.6e+00,-1.0e+00,-6.1e+00,-1.2e-01
[uo_nn_solve] -2.4e+00,-1.4e+01,+7.9e+00,+1.5e+00,-3.6e+00
[uo_nn_solve] -1.0e+01,+6.0e+00,-1.1e+01,-4.3e-01,-6.9e-01
[uo_nn_solve] +8.7e+00,-9.0e-01,-6.6e+00,+7.8e-01,+5.3e+00
[uo_nn_solve] +1.0e+01,-7.7e-01,+5.6e+00,+4.0e+00,-7.8e+00
[uo_nn_solve] ]
[uo_nn_solve] Accuracy.
[uo_nn_solve] tr_accuracy = 100.0
[uo_nn_solve] te_accuracy = 96.0
wall time = 2.5e+00 s.
```

```
>> uo_nn_Xyplot(wo,0,[])
```

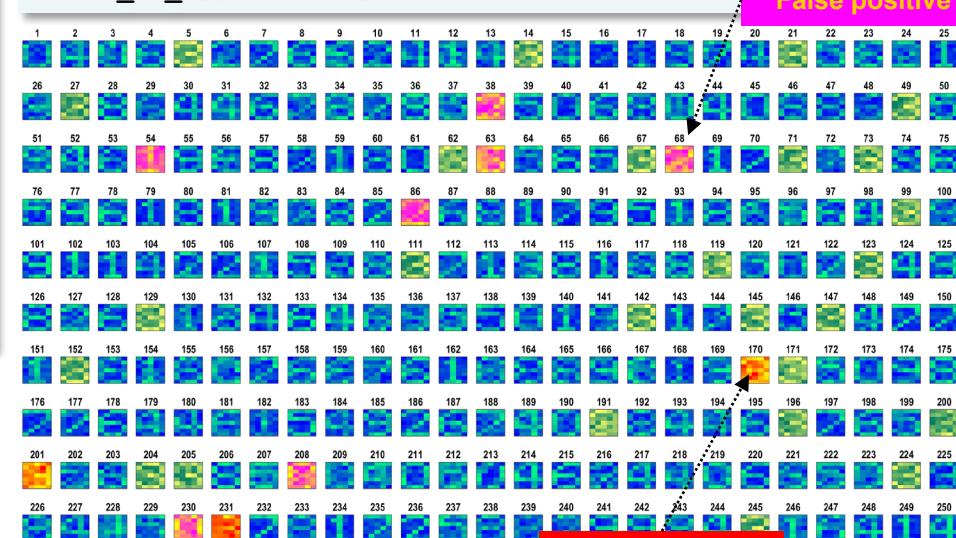
Righ positive

Righ negative

```
>> uo_nn_Xyplot(Xtr,ytr,wo)
```



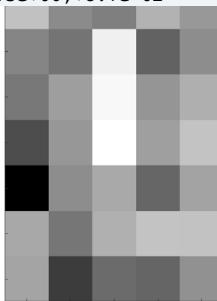
```
>> uo_nn_Xyplot(Xte,yte,wo)
```



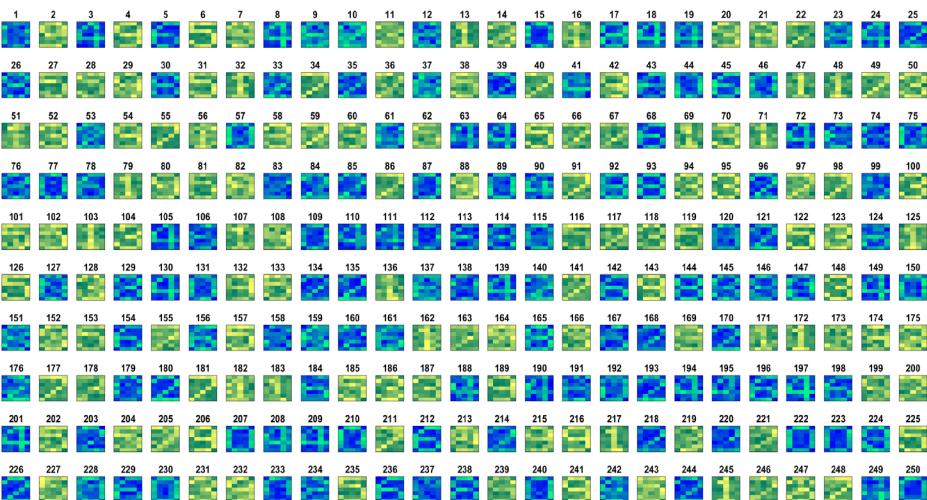
Part 1: example 3, num_target=[1 3 5 7 9]

```
[uo_nn_solve] ::::::::::::::::::::: 
[uo_nn_solve] Pattern recognition with neural networks (OM/GCED). 
[uo_nn_solve] 11-Feb-2021 20:14:05 
[uo_nn_solve] ::::::::::::::::::::: 
[uo_nn_solve] Training data set generation. 
[uo_nn_solve] num_target = 1 3 5 7 
[uo_nn_solve] tr_freq = 0.50 
[uo_nn_solve] tr_p = 250 
[uo_nn_solve] tr_seed = 123456 
[uo_nn_solve] Test data set generation. 
[uo_nn_solve] te_freq = 0.00 
[uo_nn_solve] te_q = 250 
[uo_nn_solve] te_seed = 789101 
[uo_nn_solve] Optimization 
[uo_nn_solve] L2 reg. lambda = 0.0000 
[uo_nn_solve] epsG= 1.0e-06, kmax= 10000 
[uo_nn_solve] ils= 3, ialmax= 2, kmaxBLS= 30, epsBLS= 1.0e-03, 
[uo_nn_solve] c1= 0.01, c2= 0.45, isd= 1 
[uo_nn_solve] wo=[0] 
[uo_nn_solve] k al iW g'*d f ||g|| 
[uo_nn_solve] 1 1.00e+00 2 -3.41e-02 2.50e-01 1.85e-01 
[uo_nn_solve] 11 3.36e+00 0 -1.56e-03 4.89e-02 3.95e-02 
[uo_nn_solve] ..... 
[uo_nn_solve] 1821 1.56e+04 0 -2.31e-12 4.02e-06 1.52e-06 
[uo_nn_solve] 1841 1.49e+04 0 -1.90e-12 3.70e-06 1.38e-06 
[uo_nn_solve] 1861 1.53e+04 0 -1.71e-12 3.43e-06 1.31e-06 
[uo_nn_solve] 1881 1.55e+04 0 -1.57e-12 3.17e-06 1.25e-06 
[uo_nn_solve] 1901 1.55e+04 0 -1.47e-12 2.94e-06 1.21e-06 
[uo_nn_solve] 1914 1.55e+04 0 -1.37e-12 2.79e-06 9.90e-07 
[uo_nn_solve] k al iW g'*d f ||g|| 
[uo_nn_solve] wo=[ 
[uo_nn_solve] +6.1e+00,-7.3e-01,-2.3e+00,+4.1e+00,+9.4e-01 
[uo_nn_solve] -1.5e+00,-3.6e+00,+1.2e+01,-6.0e+00,-4.5e-01 
[uo_nn_solve] -3.1e+00,+1.6e+00,+1.3e+01,+5.3e-01,+3.9e+00 
[uo_nn_solve] -8.5e+00,+8.0e-01,+1.4e+01,+1.8e+00,+6.4e+00 
[uo_nn_solve] -1.8e+01,-5.4e-01,+3.1e+00,-5.4e+00,+2.3e+00 
[uo_nn_solve] +2.9e+00,-3.4e+00,+4.0e+00,+6.4e+00,+6.3e+00 
[uo_nn_solve] +2.4e+00,-1.1e+01,-4.8e+00,-5.5e+00,+8.7e-02 
[uo_nn_solve] 1 
[uo_nn_solve] Accuracy. 
[uo_nn_solve] tr_accuracy = 100.0 
[uo_nn_solve] te_accuracy = 95.6 
wall time = 2.2e+00 s.
```

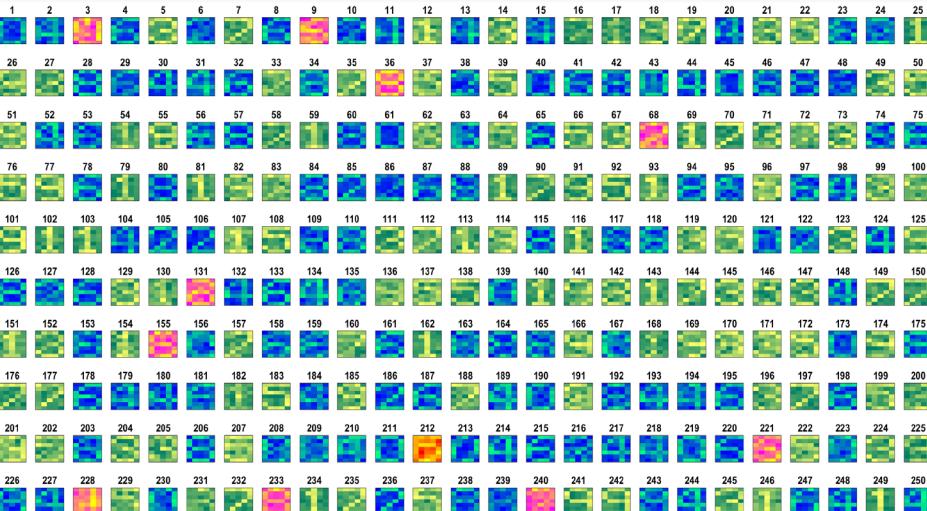
>> uo_nn_Xyplot(wo,0,[])



>> uo_nn_Xyplot(Xtr,ytr,wo)



>> uo_nn_Xyplot(Xte,yte,wo)



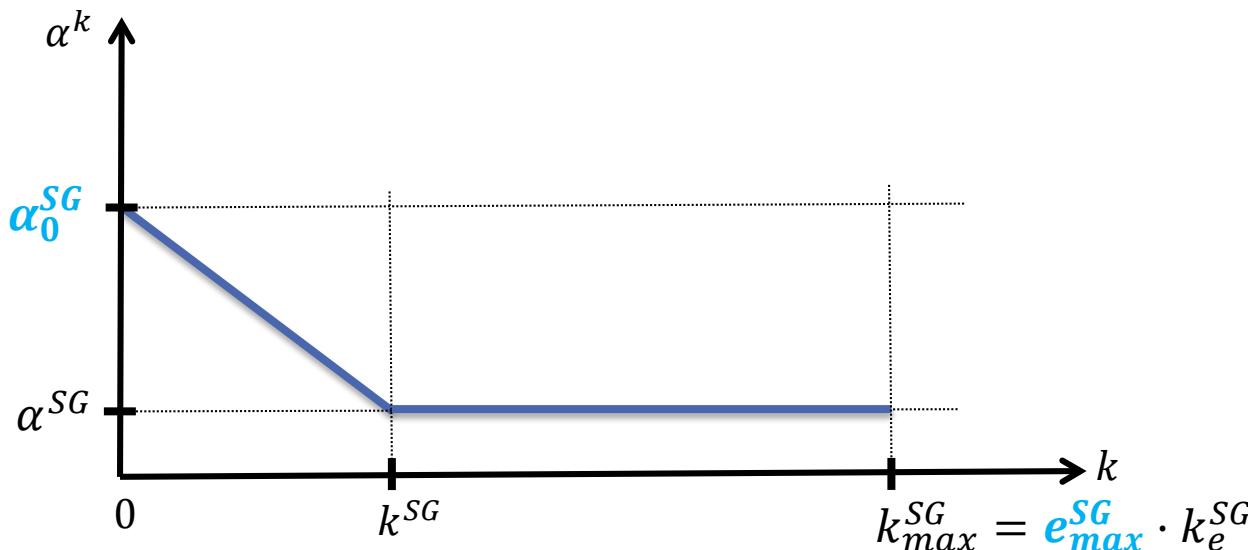
Part 2: Stochastic Gradient (1/6)

- **Stochastic Gradient Method (SGM).** (<https://www.deeplearningbook.org/>, chapter 8)
 - **Principle:** it is possible to obtain an unbiased estimate of the gradient by taking the average gradient on a minibatch of m examples drawn i.i.d. from the training dataset.
 - **Procedure:** given w^0 and $\alpha_0^{SG} > 0$, $\beta^{SG}, \gamma^{SG} \in [0,1]$ and $e_{max}^{SG} > 0$. $k \leftarrow 0$.
 1. Take a random permutation of the training data set: $\mathbf{P} \leftarrow \{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_p\}$
 2. Define the size of the minibatch: $\mathbf{m} := \lfloor \gamma^{SG} \cdot p \rfloor (\ll p)$.
 3. Define a minibatch of observations associated to the first \mathbf{m} indexs of \mathbf{P} :
 - Indexes of the minibatch: $\mathbf{S} \leftarrow \{\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_m\} \equiv \{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_m\}$.
 - Observations of the minibatch:
$$\begin{cases} X_{\mathbf{S}}^{TR} \leftarrow [x_{\mathbf{s}_1}^{TR}, x_{\mathbf{s}_2}^{TR}, \dots, x_{\mathbf{s}_m}^{TR}] \\ y_{\mathbf{S}}^{TR} \leftarrow [y_{\mathbf{s}_1}^{TR} \quad y_{\mathbf{s}_2}^{TR} \quad \dots \quad y_{\mathbf{s}_m}^{TR}]^T \end{cases}$$
 4. Estimate the gradient search direction : $d^k \leftarrow -\nabla \tilde{L}(w^k; X_{\mathbf{S}}^{TR}, y_{\mathbf{S}}^{TR}, \lambda)$
 5. Update the parametres: $w^k \leftarrow w^k + \alpha^k d^k$ with **learning rate** α^k . $k \leftarrow k + 1$
 6. Repeat steps 3 to 5 $k_e^{SG} := [\mathbf{p}/\mathbf{m}]$ iterations, taking at every iteration the next bunch of \mathbf{m} indexes of \mathbf{P} until every observation is used (**epoch**).
 7. Generate a new permutation \mathbf{P} and repeat the steps 2 to 6 for a number of e_{max}^{SG} epochs (the total number of iterations being $k_{max}^{SG} := e_{max}^{SG} \cdot k_e^{SG}$).

Part 2: Stochastic Gradient (2/6)

- **Learning rate α^k :** linear decay until iteration k^{SG} and then a constant value till the last iteration k_{max}^{SG} :

$$\alpha^k \leftarrow \begin{cases} \left(1 - \frac{k}{k^{SG}}\right) \alpha_0^{SG} + \frac{k}{k^{SG}} \alpha^{SG} & \text{if } k \leq k^{SG}, \\ \alpha^{SG} & \text{if } k > k^{SG} \end{cases}, \quad \begin{cases} \alpha^{SG} := 0,01 \cdot \alpha_0^{SG} \\ k^{SG} := \lfloor \beta^{SG} \cdot k_{max}^{SG} \rfloor \end{cases}$$



Part 2: Stochastic Gradient (3/6)

- **Stopping criterion:** the usual stopping condition $\|\nabla f(x^k)\| \approx 0$ is no longer suitable, as the SGM does not compute the true gradient $\nabla f(x)$. One of the most commonly regularization stopping condition is the ***early stopping condition***. The rationale of this stopping criterion is:
 1. To check the value of the loss function for the **test data set** after each epoch e :

$$\tilde{L}^{TE} \leftarrow \tilde{L} \left(w^{e \cdot k_e^{SG}}, \mathbf{X}^{TE}, \mathbf{y}^{TE}, \lambda \right).$$

2. Every time the new value \tilde{L}^{TE} improves the best value so far (\tilde{L}_{best}^{TE}), \tilde{L}^{TE} and the associated solution $w^{e \cdot k_e^{SG}}$ is saved :

$$\tilde{L}^{TE} < \tilde{L}_{best}^{TE} \stackrel{\text{def}}{=} \min_{j=1, \dots, e-1} \left\{ \tilde{L} \left(w^{j \cdot k_e^{SG}}, \mathbf{X}^{TE}, \mathbf{y}^{TE}, \lambda \right) \right\} \Rightarrow \begin{cases} \tilde{L}_{best}^{TE} \leftarrow \tilde{L}^{TE} \\ w^* \leftarrow w^{e \cdot k_e^{SG}} \end{cases}$$

3. If the value of \tilde{L}_{best}^{TE} has not been improved (updated) for the last e_{best}^{SG} epochs, the algorithm is terminated, and the solution of the last update w^* is declared optimal.

Part 2: Stochastic Gradient (4/6)

Algorithm SGM: Stochastic Gradient Method.

```

 $w^* \leftarrow \text{SGM}(w^0, \lambda, \tilde{L}, \nabla \tilde{L}, \mathbf{X}^{TR}, \mathbf{y}^{TR}, \mathbf{X}^{TE}, \mathbf{y}^{TE}, \alpha_0^{SG}, \beta^{SG}, \gamma^{SG}, e_{max}^{SG}, e_{best}^{SG})$ 
 $p := \text{columns}(\mathbf{X}^{TR});$ 
 $\mathbf{m} := \lfloor \gamma^{SG} \cdot p \rfloor; k_e^{SG} := \lceil p/m \rceil; k_{max}^{SG} := e_{max}^{SG} \cdot k_e^{SG}; e := 0; s := 0; \tilde{L}_{best}^{TE} := +\infty; k := 0;$ 
While  $e \leq e_{max}^{SG}$  and  $s < e_{best}^{SG}$  do
     $\mathbf{P} \leftarrow \{\mathbf{p}_1, \mathbf{p}_2, \dots, \mathbf{p}_p\};$ 
    For  $i = 0, 1, \dots, [p/m - 1]$  do
         $\mathbf{S} \leftarrow \{\mathbf{s}_1, \mathbf{s}_2, \dots, \mathbf{s}_m\} \equiv \{\mathbf{p}_{i \cdot m + 1}, \mathbf{p}_{i \cdot m + 2}, \dots, \mathbf{p}_{\min\{(i+1) \cdot m, p\}}\};$ 
         $\mathbf{X}_s^{TR} \leftarrow [x_{\mathbf{s}_1}^{TR}, x_{\mathbf{s}_2}^{TR}, \dots, x_{\mathbf{s}_m}^{TR}]; \mathbf{y}_s^{TR} \leftarrow [y_{\mathbf{s}_1}^{TR} \quad y_{\mathbf{s}_2}^{TR} \quad \dots \quad y_{\mathbf{s}_m}^{TR}]^T;$ 
         $d^k \leftarrow -\nabla \tilde{L}(w^k; \mathbf{X}_s^{TR}, \mathbf{y}_s^{TR}, \lambda);$ 
         $\alpha^k \leftarrow \begin{cases} \left(1 - \frac{k}{k^{SG}}\right) \alpha_0^{SG} + \frac{k}{k^{SG}} \alpha^{SG} & \text{if } k \leq k^{SG} \\ \alpha^{SG} & \text{if } k > k^{SG} \end{cases}, \quad \begin{cases} \alpha^{SG} := 0,01 \cdot \alpha_0^{SG} \\ k^{SG} := \lfloor \beta^{SG} \cdot k_{max}^{SG} \rfloor \end{cases}$ 
         $w^{k+1} \leftarrow w^k + \alpha^k d^k; k \leftarrow k + 1;$ 
    End For
     $e \leftarrow e + 1; \tilde{L}^{TE} \leftarrow \tilde{L}(w^k, \mathbf{X}^{TE}, \mathbf{y}^{TE}, \lambda);$ 
    If  $\tilde{L}^{TE} < \tilde{L}_{best}^{TE}$  then ( $\tilde{L}_{best}^{TE} \leftarrow \tilde{L}^{TE}$ ,  $w^* \leftarrow w^k$ ,  $s \leftarrow 0$ ) else  $s \leftarrow s + 1$  End If
End While
End SGM

```

Part 2: Stochastic Gradient (5/6)

- Script **uo_nn_solve.m** with **isd = 7** now applies the SGM:

uo_nn_main.m : recognition of num_target digits.

```

clear;
%
% Parameters for dataset generation
%
num_target =[1];
tr_freq    = .5;
tr_p       = 250;
te_q       = 250;
tr_seed    = 123456;
te_seed    = 789101;
%
% Parameters for optimization
%
la = 0.0;                                     % L2 regularization.
epsG = 10^-6; kmax = 10000;                   % Stopping criterium.
ils=3; ialmax = 2; kmaxBLS=30; epsal=10^-3;c1=0.01; c2=0.45; % Linesearch.
isd = 7; icg = 2; irc = 2 ; nu = 1.0;        % Search direction.
sg_seed = 565544; sg_a10 = 2; sg_be = 0.3; sg_ga = 0.01; % SGM iteration.
sg_emax = kmax; sg_ebest = floor(0.01*sg_emax); % SGM stopping condition.
%
% Optimization
%
t1=clock;
global iheader; iheader = 1;
[Xtr,ytr,wo,tr_acc,Xte,yte,te_acc,niter,tex]=uo_nn_solve(num_target,tr_freq,tr_seed,tr_p,te_seed,te_q,la,epsG,kmax,ils,ialmax,kmaxBLS,epsal,c1,c2,isd,sg_a10,sg_be,sg_ga,sg_emax,sg_ebest,sg_seed,icg,irc,nu);
t2=clock;
fprintf(' wall time = %6.1d s.\n', etime(t2,t1));

```

Part 2: Stochastic Gradient (6/6)

```
[uo_nn_solve] :::::::::::::::::::::
[uo_nn_solve] Pattern recognition with neural networks (OM/GCED).
[uo_nn_solve] 26-Apr-2021 16:45:24
[uo_nn_solve] :::::::::::::::::::::
[uo_nn_solve] Training data set generation.
[uo_nn_solve] num_target = 3
[uo_nn_solve] tr_freq   = 0.50
[uo_nn_solve] tr_p      = 250
[uo_nn_solve] tr_seed   = 123456
[uo_nn_solve] Test data set generation.
[uo_nn_solve] te_freq   = 0.00
[uo_nn_solve] te_q      = 250
[uo_nn_solve] te_seed   = 789101
[uo_nn_solve] Optimization
[uo_nn_solve] L2 reg. lambda = 0.0000
[uo_nn_solve] epsG= 1.0e-06, kmax= 10000
[uo_nn_solve] ilS= 3, ialmax= 2, kmaxBLS= 30, epsBLS= 1.0e-03,
[uo_nn_solve] c1= 0.01, c2= 0.45, isd= 7
[uo_nn_solve] sg_a10= 2.00, sg_be= 0.3, sg_ge= 0.01
[uo_nn_solve] sg_emax= 10000, sg_ebest= 100
[uo_nn_solve] w0=[0]
[uo_nn_solve] k    al   iW     g'*d      f    ||g|| 
[uo_nn_solve] 1  2.00e+00 0 -8.34e-02 2.50e-01 1.98e-01
[uo_nn_solve] 129 2.00e+00 0 -1.25e-03 3.43e-02 8.60e-02
[uo_nn_solve] 257 2.00e+00 0 +1.07e-05 1.12e-02 1.41e-02
[uo_nn_solve] 385 2.00e+00 0 +1.91e-07 7.41e-03 3.10e-03
[uo_nn_solve] 513 2.00e+00 0 -4.22e-05 7.31e-03 1.08e-02
[uo_nn_solve] 641 2.00e+00 0 -1.34e-04 1.07e-02 1.99e-02
...
[uo_nn_solve] 12673 1.93e+00 0 -2.21e-08 3.56e-04 1.78e-04
[uo_nn_solve] 12751           4.40e-04 1.60e-03
[uo_nn_solve] k    al   iW     g'*d      f    ||g|| 
[uo_nn_solve] wo=[ -5.6e-01,-1.7e-02,-1.3e+00,-2.8e-01,-1.2e+00
[uo_nn_solve] +9.3e-01,-1.3e-01,-2.1e-01,-9.6e-03,+3.3e-01
[uo_nn_solve] -2.8e+00,-6.4e-01,-6.1e-01,-8.9e-01,-1.8e-01
[uo_nn_solve] -6.2e-01,-3.1e+00,+1.6e+00,+5.7e-01,-1.8e+00
[uo_nn_solve] -2.1e+00,+4.4e-01,-2.1e+00,-2.5e-01,+2.6e-01
[uo_nn_solve] +2.4e+00,-9.4e-01,-4.9e-01,+2.4e-01,+7.6e-01
[uo_nn_solve] -5.0e-01,-3.6e-01,+2.1e-01,+2.8e-01,-6.9e-01
[uo_nn_solve] ]
[uo_nn_solve] Accuracy.
[uo_nn_solve] tr_accuracy = 98.8
[uo_nn_solve] te_accuracy = 98.8
wall time = 2.8e-01 s.

>> uo_nn_Xyplot(wo,0,[])

```

Righ positive

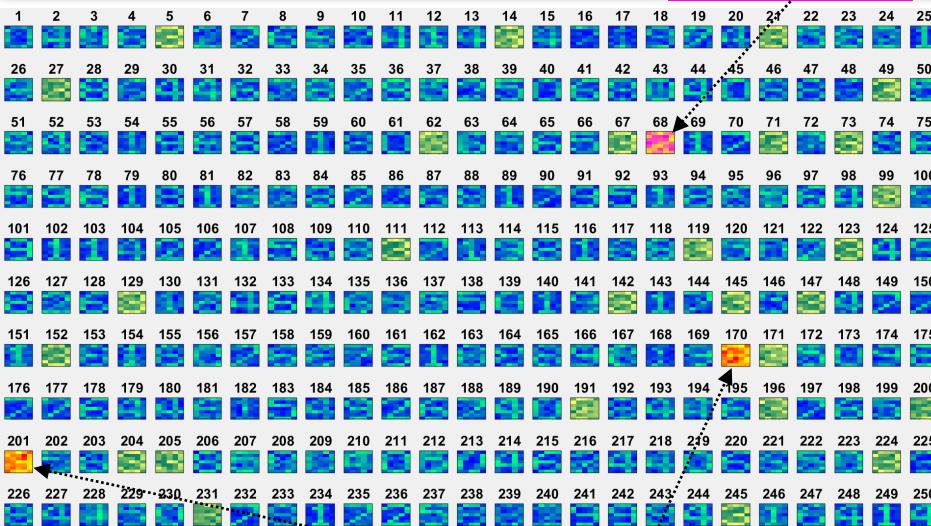
Righ negative



```
>> uo_nn_Xyplot(Xtr,ytr,wo)
```

```
>> uo_nn_Xyplot(Xte,yte,wo)
```

False positive



False negative



Part 3: computational study (1/3)

- In this third part we want to conduct a series of computational experiments to study:
 - i. How the regularization parameter λ affects the results.
 - ii. The relative performance of the different algorithms (GM, QNM, SGM)
- To this end, an instance of the SLNN problem must be solved:
 - For every one of the individual **digits, 0 to 9**.
 - For every value of the regularization parameter $\lambda \in \{0.0, 0.01, 0.1\}$.
 - For every optimization algorithm: **GM, QNM** and **SGM**.That makes a total of **$10 \times 3 \times 3 = 90$ instances to be solved**.
- We will use a small size dataset with

```
tr_p = 250; te_q = 250; tr_freq = 0.5;
```
- Max. number of iterations for GM and QNM will be **kmax = 1000**. To do a fair comparison with SGM, **sg_emax = kmax**, so that every method is allowed to spend a similar computational effort (number of gradient evaluations).

Part 3: computational study (2/3)

- To organize the computational experiments you can use function `uo_nn_batch.m`.

`uo_nn_batch.m` : run a batch of SLNN instances.

```

clear;
% Parameters.
tr_seed = 123456; te_seed = 789101;
tr_p = 250; te_q = 250; tr_freq = 0.5;                                % Datasets generation
epsG = 10^-6; kmax = 1000;                                              % Stopping condition.
ils=3; ialmax = 1; kmaxBLS=10; epsal=10^-3; c1=0.01; c2=0.45;          % Linesearch.
icg = 2; irc = 2 ; nu = 1.0;                                            % Search direction.
sg_seed = 565544; sg_a10 = 2; sg_be = 0.3; sg_ga = 0.01;                % SGM iteration.
sg_emax = kmax; sg_ebest = floor(0.01*sg_emax);                          % SGM stopping condition.
% Optimization
global iheader; iheader = 1;
csvfile = strcat('uo_nn_batch_',num2str(tr_seed),'-',num2str(te_seed),'.csv');
fileID = fopen(csvfile , 'w');
t1=clock;
for num_target = [1:10]
    for la = [0.0, 0.01, 0.1]
        for isd = [1,3,7]
[Xtr,ytr,wo,fo,tr_acc,Xte,yte,te_acc,niter,tex]=uo_nn_solve(num_target,tr_freq,tr_seed,tr_p,te_seed,te_q,la,
,epsG,kmax,ils,ialmax,kmaxBLS,epsal,c1,c2,isd,sg_a10,sg_be,sg_ga,sg_emax,sg_ebest,sg_seed,icg,irc,nu);
        if iheader == 1
            fprintf(fileID,'num_target;      la; isd;  niter;      tex; tr_acc; te_acc;      L*\n');
        end
        fprintf(fileID,'      %i; %7.4f;    %i; %6i; %7.4f;  %5.1f;  %5.1f;  %8.2e;\n',
mod(num_target,10), la, isd, niter, tex, tr_acc, te_acc, fo);
        iheader=0;
    end
end
end
t2=clock; fprintf(' wall time = %6.1d s.\n', etime(t2,t1));
fclose(fileID);

```



Part 3: computational study (3/3)

- The outcome of `uo_nn_batch.m` is the **log file** `uo_nn_batch_tr_seed-te_seed.csv` with the following content:

num_target;	la;	isd;	niter;	tex;	tr_acc;	te_acc;	L*;
1;	0.0000;	1;	75;	0.1055;	100.0;	100.0;	5.72e-07;
1;	0.0000;	3;	6;	0.0099;	100.0;	100.0;	2.91e-50;
1;	0.0000;	7;	125001;	1.5164;	100.0;	100.0;	1.41e-05;
1;	0.0100;	1;	58;	0.0955;	100.0;	100.0;	2.76e-02;
1;	0.0100;	3;	45;	0.1026;	100.0;	100.0;	2.76e-02;
1;	0.0100;	7;	3751;	0.0464;	96.4;	99.6;	6.74e-02;
.....							
0;	0.1000;	1;	36;	0.0645;	99.6;	98.4;	1.45e-01;
0;	0.1000;	3;	20;	0.0520;	99.6;	98.4;	1.45e-01;
0;	0.1000;	7;	2751;	0.0336;	53.6;	92.4;	6.37e-01;

The data included in this file will be the base for the analysis of the performance of the different optimization algorithms that you have to perform in this project.

Report (1/3)

- 1) **Study of the convergence:** first, we are going to study the global and local convergence of the three algorithms only in terms of the objective function \tilde{L} :
- a) **Global convergence:** analyse the global convergence of every algorithm and how this global convergence property depends on the value of the regularization parameter λ . Which combination algorithm- λ gives the best results in terms of global converge? In particular, discuss the application to the SGM of the conditions for global convergence.
 - b) **Local convergence:**
 - i. Compare the speed of convergence of the three algorithms in terms of the execution time and number of iterations.
 - ii. Analyse how the speed of convergence of the three algorithms depend on the value of λ and try to find an explanation for the observed dependence, if any.
 - iii. Analyse the running time per iteration (`tex/niter`) and try to find an explanation for the different values among the three algorithms.
 - c) Finally, according to the previous study, discuss the general performance of the three algorithms in terms of the observed local and global convergence and justify which is the most efficient combination λ -algorithm for the minimization of \tilde{L} .

Report (2/3)

- 2) **Study of the recognition accuracy:** now, we are going to analyse the recognition accuracy of the SLNN, $Accuracy^{TE}$ (`te_acc`) for the different algorithmic options with a more realistic dataset with:

```
tr_p = 20000; te_q = tr_p/10; tr_freq = 0.0;
```

Run the training process for the ten digits with these parameters and every algorithm, GM, QNM, SGM, using, for each one of them, the value of λ showing the best $Accuracy^{TE}$ in the results of the previous section 1- a) with the smaller dataset (`tr_p = 250`).

- a) Analyse the results. Based on this analyse, tell if there is a method that clearly outperforms the others in terms of training speed and recognition accuracy.
- b) Does the best combination λ -algorithm with respect to the maximization of $Accuracy^{TE}$ coincide with best combination with respect to the minimization of \tilde{L} observed in the previous study 1-c)? Discuss the reasons for the discrepancy, in case the combinations do not coincide.

Report (3/3)

- This assignment must be done in **groups of two**. Use a value of **tr_seed** and **te_seed** based on your student's ID number. You must upload to Atenea a file with the name **surname-student-1_surname-student-2.zip** containing:
 - A report (.pdf file) with your answers to the different sections of tasks 1) “Convergence” and 2) “Accuracy”.
 - The report must have a cover with the name of the two students and the values of **tr_seed**, **te_seed** and **sg_seed**
 - The source of all the codes used to do the assignment. The codes must be self-contained and must allow the replication of all the results included in the report.
- The mark will take into account formal quality, understandability and reading easiness, as long as the evidences provided to support the results of the analysis performed in the project (numerical values, tables, plots, ...).

Grading criteria

	Points
1) Study of the convergence:	
a) Global convergence.	1,5
b) Local convergence.	2,5
c) Discussion.	1,0
2) Study of the accuracy for large problems.	
a) Best method for large problems.	2,0
b) Discussion $Accuracy^{TE} - \tilde{L}$	2,0
3) Overall quality of the report.	1,0