

# Math 124 - Programming for Mathematical Applications

UC Berkeley, Spring 2023

## Project 5 - Neural Networks for Character Recognition

Due Friday, April 28

In [1]: `1 using PyPlot, Random, LinearAlgebra, Optim # Packages needed`

### Description

In this project, you will implement an artificial neural network for *Optical Character Recognition* (OCR). We will use a so-called *Multilayer Perceptron* (MLP) with a single hidden layer.

### Preliminaries

First we define the characters that we will use to train the network. Each character will be an image of size 7-by-6, represented as a vector of 42 values, and we will use the 4 characters "MATH". In the network, these characters will be encoded using two output variables  $y = (y_1, y_2)$ , where  $y = (0, 0)$  represents "M",  $y = (0, 1)$  represents "A", and so on. The code below defines these images and returns the following variables:

- `training` , size 42-by-4 array containing the 42 pixel-values for each of the 4 characters
- `target` , size 2-by-4 array containing the desired output for each of the 4 characters
- `mapstr` , the string "MATH" which is the characters that each `target` output corresponds to

```

In [2]: 1 charstr = ""
2         000000 000000 000000 00..00
3         000000 000000 000000 00..00
4         0.00.0 00..00 ..00.. 00..00
5         0.00.0 00..00 ..00.. 000000
6         0....0 000000 ..00.. 00..00
7         0....0 00..00 ..00.. 00..00
8         0... 0 00..00 ..00.. 00..00
9         ""
10
11 training = reshape(collect(charstr), :, 7)
12 training = Int.(training[[1:6;9:14;17:22;25:30],:] .== '0')
13 training = reshape(training', 7*6, 4)
14 target = [0 0; 0 1; 1 0; 1 1]'
15 mapstr = "MATH";

```

We also define the plotting function below, which takes an array `images` with 42 rows and one column for each image, and shows the images in a grid:

```

In [3]: 1 function plot_chars(images)
2         gray()
3         n_images = size(images,2)
4         for j = 1:n_images
5             subplot(ceil(Int, n_images/4), 4, j)
6             im = 1 .- reshape(images[:,j], 7, 6)
7             imshow(im); axis("off");
8         end
9     end
10 plot_chars(training)

```



The image displays the word "MATH" in a large, bold, black, pixelated font. Each letter is composed of a 7x6 grid of pixels, where black pixels represent the letter and white pixels represent the background. The letters are arranged horizontally with consistent spacing between them.

## Problem 1 - Generating noisy test characters

To test our trained OCR code, we need noisy or perturbed characters. Here we will artificially produce these by modifying the true character images in `training`.

Write a function `make_testdata(training)` which returns an array `testdata` of size 42-by-20. The first 4 images (columns) of `testdata` are identical copies of the `training` array. Generate the next 4 images by randomly choosing 2 pixels in each training image, and flip their values (that is, 0 becomes 1, 1 becomes 0). For the next 4 images you choose  $2 \times 2 = 4$  random pixels to flip, then  $2 \times 3 = 6$  pixels, and finally  $2 \times 4 = 8$  pixels. This gives a total of 16 new perturbed images, and 20 columns total including the original images.

Plot your test data using the commands below. This should show a 5-by-4 array of successively worse letters MATH in each row.

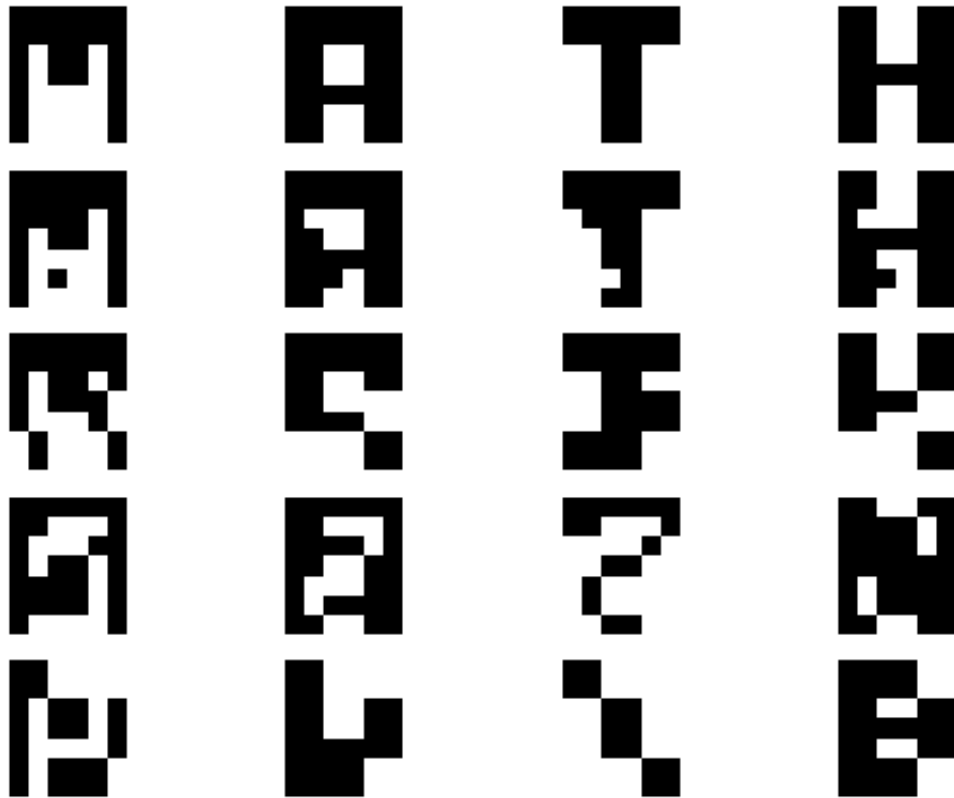
```

In [4]: 1 function make_testdata(training)
2
3     function pixel_mod(dim1, dim2, images)
4         images2 = zeros{Int64, size(images, 1), size(images, 2)}
5         range1 = 7 - dim1 + 1
6         range2 = 6 - dim2 + 1
7         row1 = rand(collect(1:range1))
8         col1 = rand(collect(1:range2))
9         row2, col2 = 0, 0
10
11         while true
12             row2 = rand(collect(1:range1))
13             col2 = rand(collect(1:range2))
14             if row2 != row1 || col2 != col1
15                 break
16             end
17         end
18
19         for k = 1:size(images, 2)
20             A = reshape(images[:, k], 7, 6)
21             for i = 1:dim1
22                 for j = 1:dim2
23
24                     A[row1 + i - 1, col1 + j - 1] = 1 - A[row1 + i
25                     A[row2 + i - 1, col2 + j - 1] = 1 - A[row2 + i
26                 end
27             end
28             images2[:, k] = reshape(A, 1, 42)
29         end
30
31         return images2
32
33     end
34
35     return training2 = [training pixel_mod(1, 1, training) pixel_m
36
37 end

```

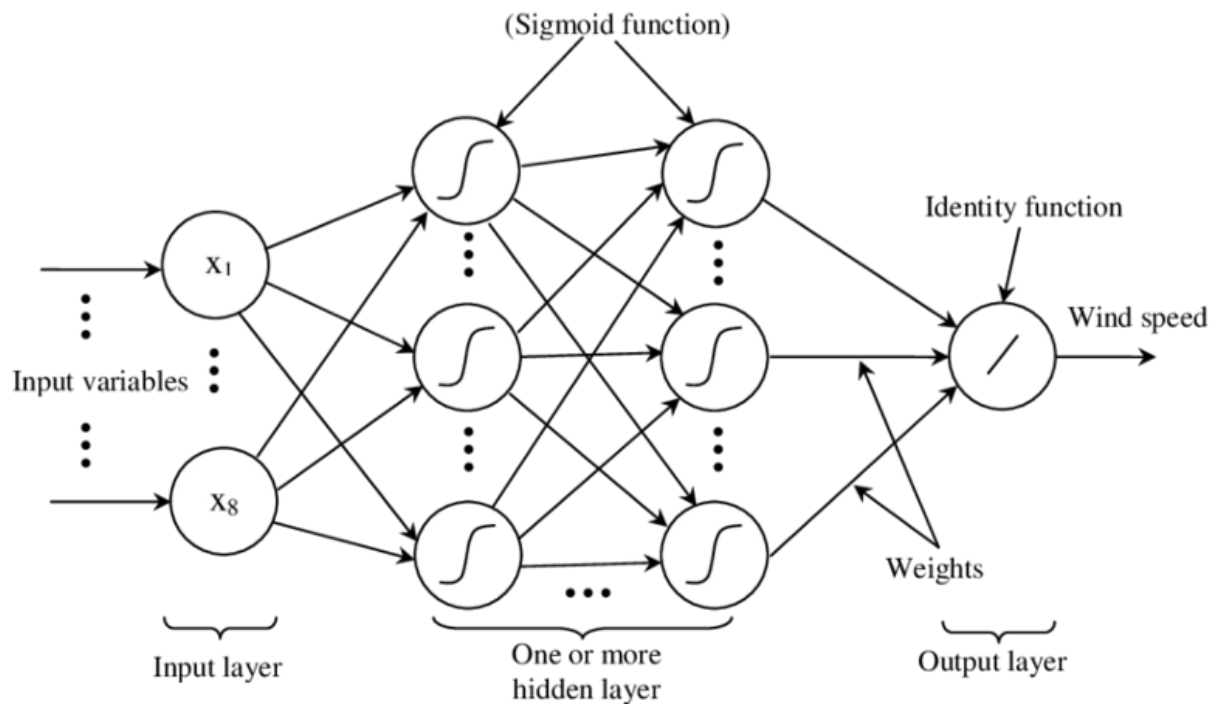
Out[4]: make\_testdata (generic function with 1 method)

```
In [5]: 1 testdata = make_testdata(training);  
        2 plot_chars(testdata)
```



## Machine learning using a multilayer perceptron

Our network is illustrated below. The input  $x = (x_1, \dots, x_{42})$  is in our case an image represented as a vector of length 42, and the output  $y = (y_1, y_2)$  is a vector of length 2, with the predicted character encoded as described above.



(by V. Palaniappan, <https://medium.com/engineer-quant/multilayer-perceptron-4453615c4337> (<https://medium.com/engineer-quant/multilayer-perceptron-4453615c4337>))

Here, the *sigmoid function* is defined by

$$\sigma(x) = \frac{1}{1 + e^{0.5-x}}$$

The *weights* (that is, the arrows between the layers) can be represented in a linear form as a matrix multiplication. The goal of the training is to determine these weights, that is, the entries of the matrices.

We will use a hidden layer with 10 neurons. This means that the first layer can be written as a matrix-vector product  $Vx$  for some weight matrix  $V$  of size 10-by-42. Next we apply the sigmoid function (element-wise) to get  $r = \sigma(Vx)$ , apply another weight matrix  $W$  of size 2-by-10, and finally apply the sigmoid function again to predict the output  $y$ . In total this gives  $10 \cdot 42 + 2 \cdot 10 = 440$  weights that need to be determined during the training. The entire network can be written in a compact form as

$$y(x) = \sigma(W\sigma(Vx))$$

We will train our network by solving the following optimization problem: Find  $V \in \mathbb{R}^{10 \times 42}$ ,  $W \in \mathbb{R}^{2 \times 10}$  that minimize the following so-called *misfit function*

$$E(V, W) = \frac{1}{2} \sum_{j=1}^4 \|y(x_j) - t_j\|_2^2$$

Here,  $x_j$  is the  $j$ th training vector (from the `training` array),  $t_j$  is the desired target output for the input  $x_j$  (from the `target` array), and  $y(x)$  is the network described above.

## Stochastic Gradient Descent training

Finding the true gradient of  $E$  is somewhat involved, and instead we will use the so-called Stochastic Gradient Descent method. At each iteration, we pick:

- A random number  $j \in \{1, 2, 3, 4\}$  to decide which training vector  $x_j$  and target output  $t_j$  to use
- A random number  $k \in \{1, 2\}$  to decide which component of the vector  $y(x_j) - t_j$  to consider

We then only compute the gradient direction for the part of  $E(V, W)$  that contribute to output  $k$  for training vector  $j$ . The gradient calculations get simplified with these assumptions, and they can be computed with the following sequence of computations:

$r = \sigma(Vx_j)$	column-vector (matrix-vector product)
$y = \sigma(W_k \cdot r)$	scalar (dot-product, $W_k$ is $k$ th row of $W$ )
$q = (y - t_{kj})y(1 - y)$	scalar ( $t_{kj}$ is component $k$ of $t_j$ )
$u = W_k^T r(1 - r)$	column-vector (all products are element-wise)
$\nabla_{W_k} E(V, W) = qr^T$	row-vector (scalar times row-vector)
$\nabla_V E(V, W) = qux_j^T$	matrix (scalar times outer product)

Using these gradients, we update  $V$  and the  $k$ th row of  $W$  using standard gradient descent iterations for a given stepsize  $\alpha$ .

## Problem 2 - Training using SGD

Implement a function

```
function train_sgd(; maxiter=10000, rate=1)
```

which initializes  $V, W$  to normal-distributed random numbers and performs `maxiter` stochastic gradient descent iterations as described above. Use the step-size  $\alpha$  set to `rate` (the so-called *learning rate*). The function returns  $V, W$ , and can be tested using the code below.

```

In [6]: 1 maxiter=10000
        2 rate=1
        3 function train_sgd(; maxiter=10000, rate=1)
        4
        5     α = rate
        6     V = reshape(randn(420), 10, 42)
        7     W = reshape(randn(20), 2, 10)
        8
        9     for i = 1:α:maxiter
       10
       11         j = rand(collect(1:4))
       12         k = rand([1, 2])
       13
       14         σ(x) = 1 / (1 + exp(0.5 - x))
       15         r = σ.(V * training[:,j])
       16         y = σ.(W[k,:] * r)
       17         q = (y - target[k,j])*y*(1 - y)
       18         u = W[k,:] .* r .* (1 .- r)
       19
       20         ΔW = q .* transpose(r)
       21         ΔV = q .* u[:, :] * transpose(training[:,j][:,:])
       22         W[k,:] -= α * transpose(ΔW)
       23         V -= α * ΔV
       24
       25     end
       26
       27     return V, W
       28
       29 end
       30

```

Out[6]: train\_sgd (generic function with 1 method)

```

In [7]: 1 V,W = train_sgd()

```

Out[7]: ([-0.10825595586026325 0.0331458569540207 ... -2.048336177478499 -1.097  
966823695898; 0.05614314715294488 -0.10207318114433704 ... -0.669606574  
1167333 -1.0773834045219517; ... ; 0.6626871883600354 0.009559912284699  
055 ... 0.18650223741770477 1.0283653526404901; 0.7374185610854761 -1.1  
658514061187704 ... 0.6262096445010441 1.2309072812497177], [4.71618517  
5211911 -1.627875174615892 ... -1.144231978254841 -2.0399458752900785;  
-0.859484338806269 1.6736657349852249 ... -0.7219365928937113 0.9720185  
138070203])



### Problem 3 - Predict output for noisy characters

We can now apply the trained network  $V, W$  to the noisy characters in `testdata`, and see how well they match. We do this by simply computing  $y(x_j)$  for each image  $x_j$  in `testdata`, rounding the outputs  $y_1, y_2$  to integers, and mapping the four possibilities to the letters in `mapstr`.

Write a function `predict(testdata, V, W)` which performs these operations, and returns a 5-by-4 character array with the predicted characters. The ideal output would be 5 rows of 'M', 'A', 'T', 'H', but we do not of course expect to be perfect for the highly noisy cases.

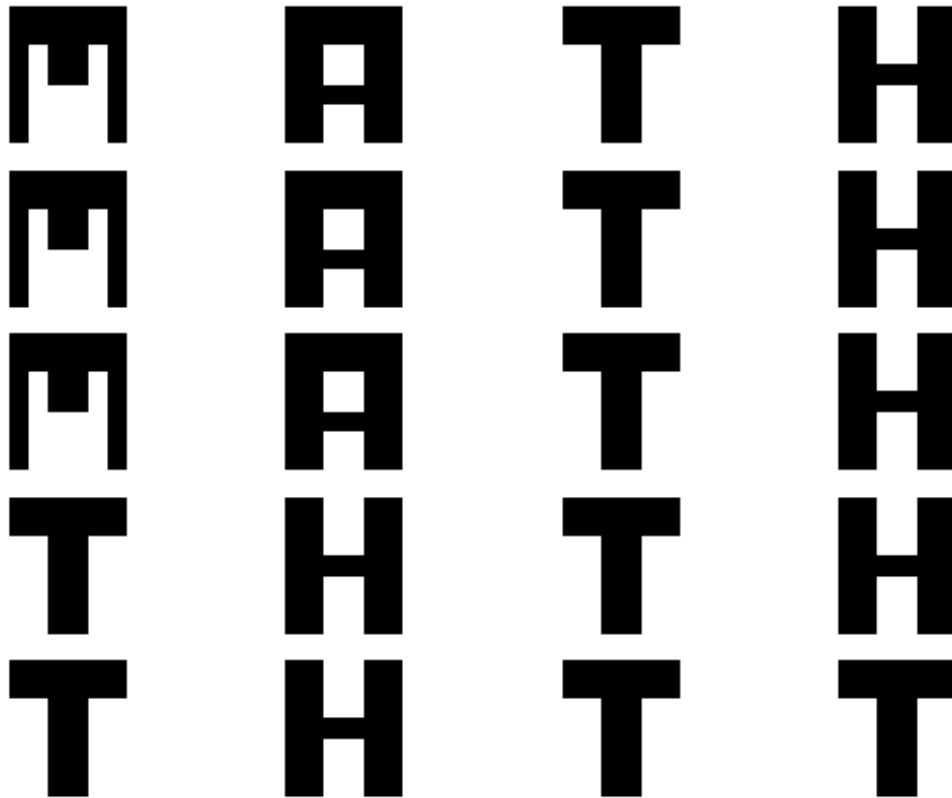
```

In [8]: 1 function predict(testdata, V, W)
        2
        3  $\sigma(x) = 1 / (1 + \exp(0.5 - x))$ 
        4 min = -Inf
        5 newdata = zeros(Int64, size(testdata, 1), size(testdata,2))
        6
        7 for i = 1:size(testdata,2)
        8
        9     y =  $\sigma.(W*\sigma.(V*testdata[:,i]))$ 
       10
       11     for j = 1:2
       12         if abs(y[j] - floor(y[j])) ≤ abs(y[j] - ceil(y[j]))
       13             y[j] = floor(y[j])
       14         else
       15             y[j] = ceil(y[j])
       16         end
       17     end
       18
       19     if y[1] == 0
       20         if y[2] == 0
       21             newdata[:,i] = training[:,1]
       22         else
       23             newdata[:,i] = training[:,2]
       24         end
       25     elseif y[2] == 0
       26         newdata[:,i] = training[:,3]
       27     else
       28         newdata[:,i] = training[:,4]
       29     end
       30 end
       31
       32 return plot_chars(newdata)
       33
       34 end

```

Out[8]: predict (generic function with 1 method)

```
In [9]: 1 plot_chars(testdata)
        2 predict(testdata, V, W)
```



#### Problem 4 - True gradient descent using the `Optim` package

Since the `Optim` package can perform automatic differentiation, we can try to solve the optimization problem using true gradient descent. Write a function `train_optim()` which does this and returns  $V$ ,  $W$  like before. Use the solver `GradientDescent()`, and remember to set `autodiff=:forward`.

One tricky thing with this is that `optimize` expects a single input vector  $x$ , but we have two matrices  $V$ ,  $W$ . This is a quite common problem, and you can get around it by converting the  $420+20=440$  numbers in  $V$ ,  $W$  to a vector. In the objective function, you also have to convert back to matrices in order to evaluate the misfit function  $E(V, W)$ .

Run your function, then predict the output using the resulting network  $V$ ,  $W$  with the code below.

```

In [10]: 1 using Optim
          2 function train_optim()
          3     X = randn(440)
          4      $\sigma(x) = 1 / (1 + \exp(0.5 - x))$ 
          5     E(x) = 0.5 * sum([norm( $\sigma$ .(reshape(x[421:end], 2, 10) *  $\sigma$ .(resh
          6     res = optimize(E, X, GradientDescent(); autodiff=:forward)
          7     return reshape(res.minimizer[1:420], 10, 42), reshape(res.mini
          8 end

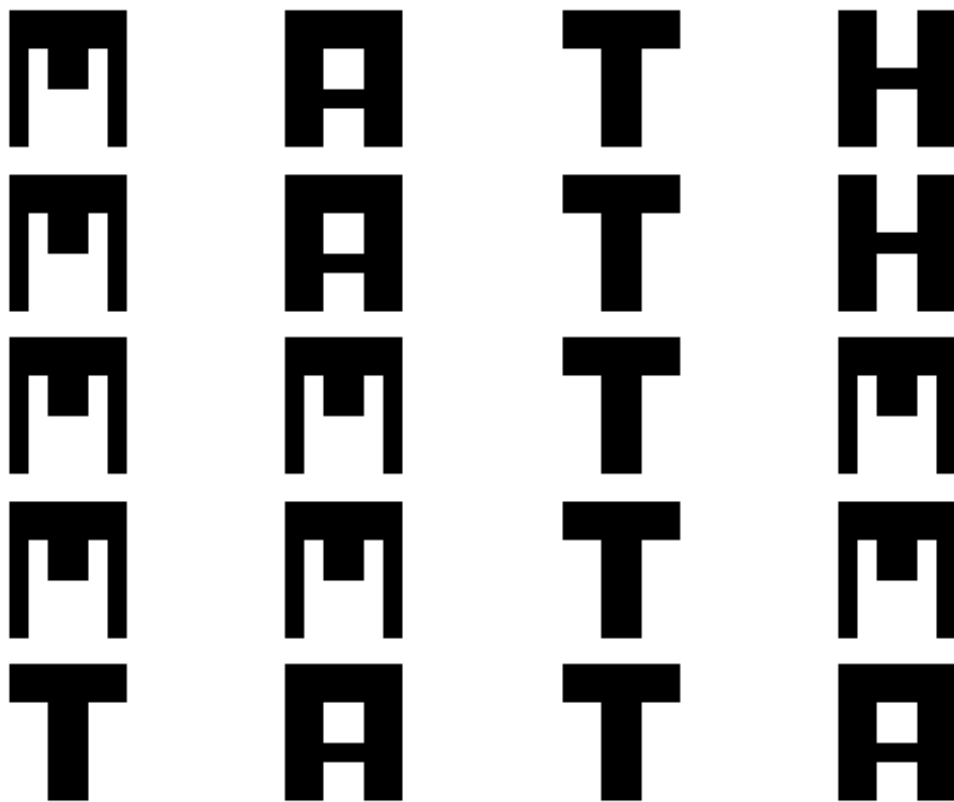
```

Out[10]: train\_optim (generic function with 1 method)

```

In [11]: 1 plot_chars(testdata)
          2 V,W = train_optim()
          3 predict(testdata, V, W)

```



In [ ]: 1