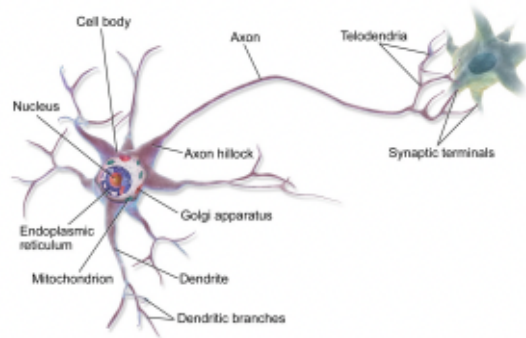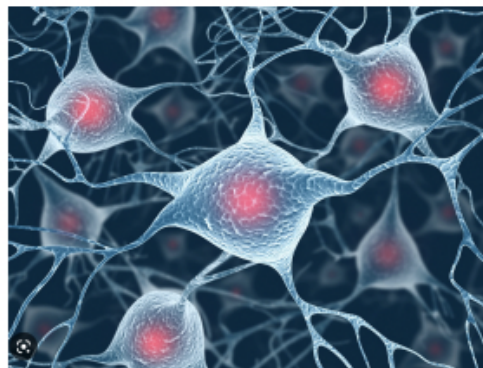CAS CS 505
Lec 9

Deep Learning: Artificial Neural Networks: Design and Implementation; Workflow for
Classification Tasks; Evaluating Classifiers

1. Introduction to Deep Learning
    a. Deep Learning refers to Supervised Learning using an Artificial Neural Network,
       which has the following features:
    b. It is a network/graph of small computation units called artificial neurons, loosely
       modeled on the neurons in our brains, which send signals to each other. The
       signals are floating-point numbers.
    c. The network is typically organized in layers: the first layer is the input layer, the
       last is the output layer, and others are called hidden layers.
    d. The input layer is array/vector of floats, and the output layer produces an array of
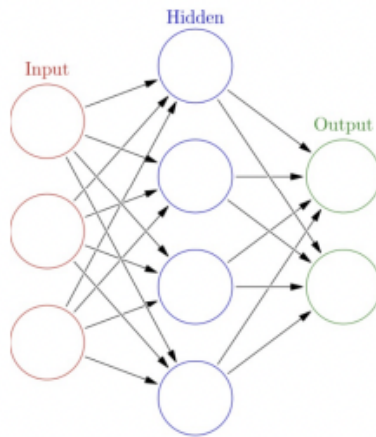       floats. Thus, the network computes a function from vectors to vectors
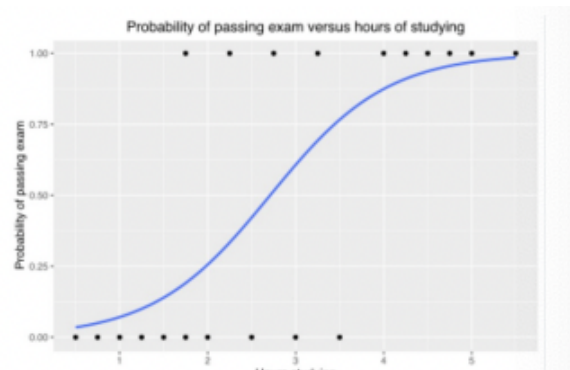
    e. 

    f. 
    g. Each neuron:
        i.   Is connected to each neuron in the previous layer, and each connection has
             a weight or parameter which determines the strength of the signal
             (importance of this input to the neuron);
        ii.  Performs logistic regression, using the sigmoid or other non-linear
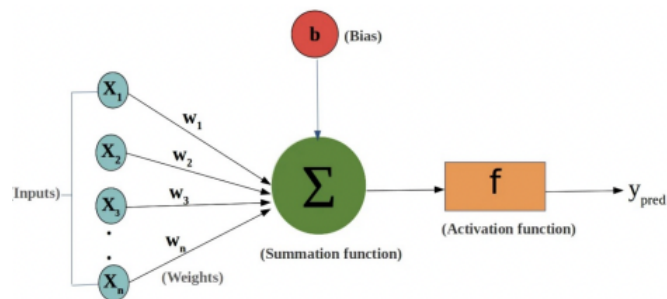             activation function.

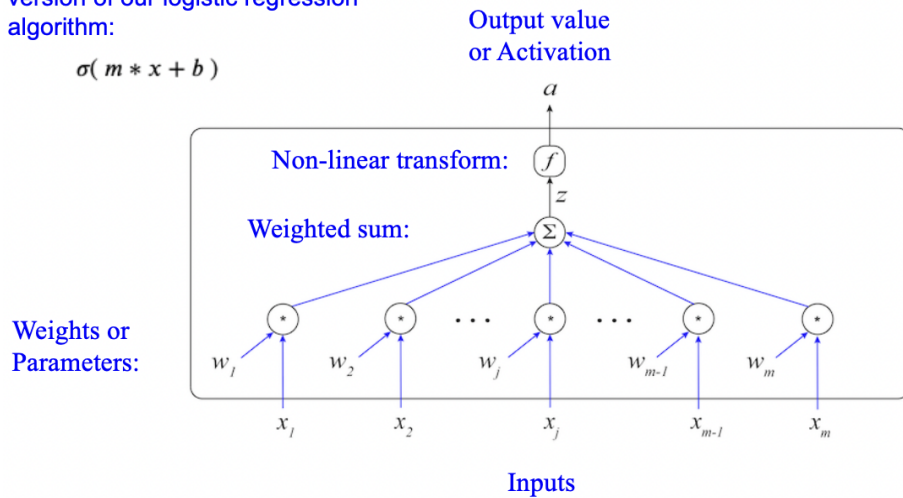h. Gradient descent is used to learn the weights to minimize some cost function on the outputs.

i.



j.

Probability of passing exam versus hours of studying



k.



l. Features of artificial neural networks:
   i. Additional layers may perform data aggregation (e.g., convolution and pooling), dropout (forgets some of the weights), or other kinds of data manipulation (e.g., softmax = transforming the output into a probability distribution).
   ii. In a feedforward network, the network transforms an array of floats through the layers into another array of floats; in a sequence model, the inputs and outputs are sequences of vectors; and recurrent layers have cyclical connections which act as memory.
   iii. BERT, GPT, and other large networks learn to pay attention to complex patterns in the input sequence (e.g., words in a sentence).

Output value
or Activation

$a$

Non-linear transform: $f$

$z$

Weighted sum: $\Sigma$

Weights or
Parameters:

$w_1$   $w_2$   $w_j$   $w_{m-1}$   $w_m$

$x_1$   $x_2$   $x_j$   $x_{m-1}$   $x_m$

Inputs
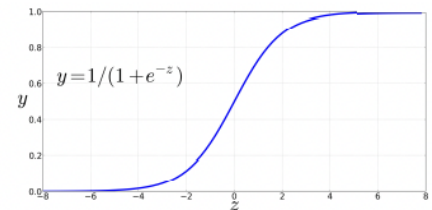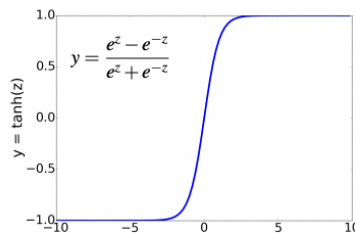
m.

n.  One possible activation function f is the sigmoid from logistic regression:

$$y = \sigma(z) = \frac{1}{1+e^{-z}}$$

$y = 1/(1+e^{-z})$

i.
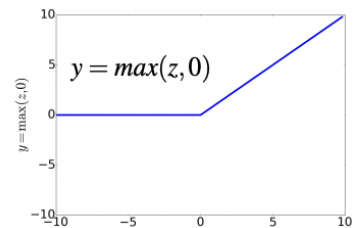
o.  Both non-linear activation functions besides sigmoid are more often used

$y = \dfrac{e^z - e^{-z}}{e^z + e^{-z}}$

$y = max(z,0)$

**Hyperbolic Tangent (Tanh)**            **Rectified Linear Unit (Relu)**

i.

ii.  Hyperbolic tangent works better than sigmoid in recurrent calls

p.  A small amount of Linear Algebra can be used to compactly specify the logistic
regression performed by a neuron.

q.  Inputs and weights/parameters are just vectors:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \qquad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \end{bmatrix}$$

$a$

$f$

$z$

$\Sigma$

$w_1$   $w_2$   $w_j$   $w_{m-1}$   $w_m$

$x_1$   $x_2$   $x_j$   $x_{m-1}$   $x_m$

i.

r.  Often written as tuples

$$\mathbf{x} = (x_1, x_2, \ldots x_m). \qquad \mathbf{w} = (w_1, w_2, \ldots w_m).$$

i.

s.  Then the neuron performs logistic regression by performing a dot product of inputs x and weights w,

$$\mathbf{w} \cdot \mathbf{x} = (w_1 \cdot x_1) + \cdots + (w_k \cdot x_k) = \mathbf{w}^T \mathbf{x}$$

i.

t.  and then applying the activation function f:

$$\sigma(m * x + b) \qquad\qquad a = f(\mathbf{w}^T \mathbf{x}). \qquad\qquad a$$

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_m \end{bmatrix} \qquad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_m \end{bmatrix}$$



i.

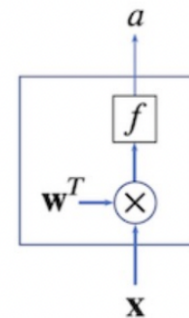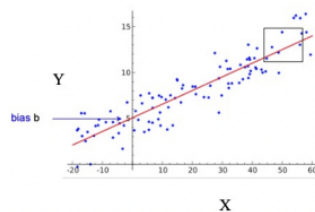u.  A small but important detail: Each neuron has a bias term, because they are simply doing logistic regression! (each neuron will have a bias)

$$\sigma(m * x + b * 1.0)$$



i.                                                 → you always have a '1.0' to match 'b' in matrix multiplication

v.  This connection has a weight (the value b) but is not connected to the inputs; it serves to scale the output, just as b does in linear regression.

w.  A convenient way to encode this is to assume that all input vectors to all neurons contain a constant 1.0 value:

$$\mathbf{x} = \begin{bmatrix} x_1 \\ x_2 \\ \vdots \\ x_{m-1} \\ 1.0 \end{bmatrix} \qquad \mathbf{w} = \begin{bmatrix} w_1 \\ w_2 \\ \vdots \\ w_{m-1} \\ b \end{bmatrix}$$



$w_m = b$

i.

x. The simplest network is simply a row of such neurons, where
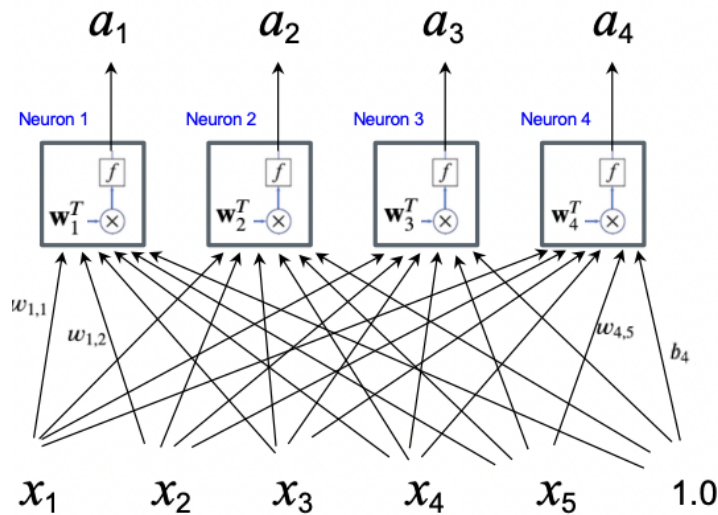    i.    Each has its own weight/parameter vector, but
    ii.   Shares the same input vector

## Activations (outputs)



Inputs

y.
    i.    They can work together to learn different thing about the inputs
    ii.   Each will do logistic regression using its weights on the input and output
          an activation (each separately)
    iii.  Vector-to-vector function
    iv.   It works simultaneously (everything in one layer in principle)
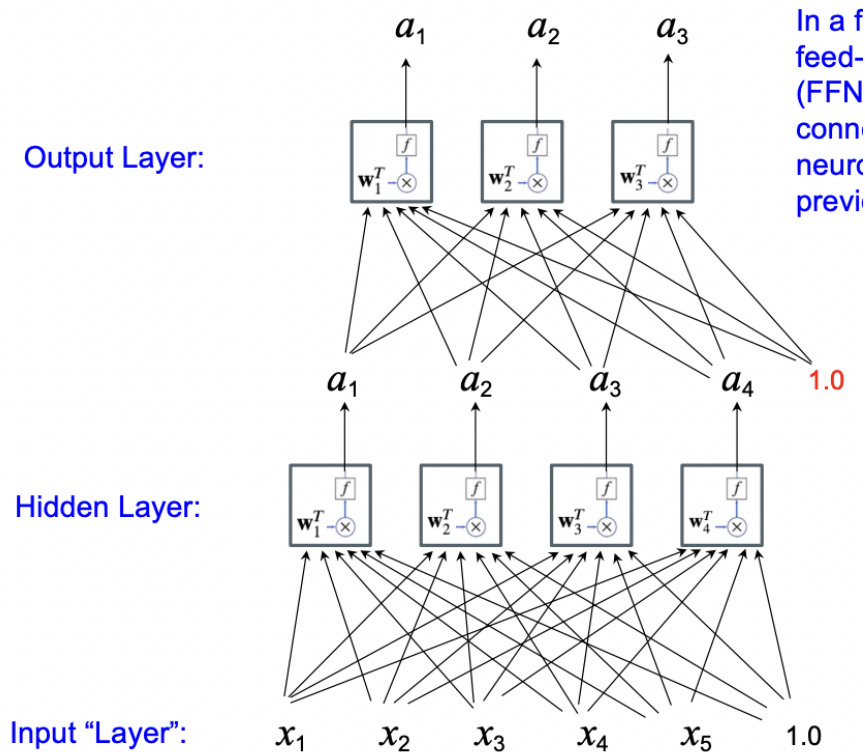z. The inputs and outputs for a layer are vectors

$$\mathbf{a} \; = \; (a_1, a_2, a_3, a_4)$$

    i.

aa. The weights/parameters for a layer form a matrix

$$W\,\mathbf{x} = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} & w_{1,4} & w_{1,5} & b_1 \\ w_{2,1} & w_{2,2} & w_{2,3} & w_{2,4} & w_{2,5} & b_2 \\ w_{3,1} & w_{3,2} & w_{3,3} & w_{3,4} & w_{3,5} & b_3 \\ w_{4,1} & w_{4,2} & w_{4,3} & w_{4,4} & w_{4,5} & b_4 \end{bmatrix} \times \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_2 \\ x_4 \\ x_5 \\ 1.0 \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^T\mathbf{x} \\ \mathbf{w}_2^T\mathbf{x} \\ \vdots \\ \mathbf{w}_n^T\mathbf{x} \end{bmatrix}.$$

$$\mathbf{x} = (x_1, x_2, x_3, x_4, x_5, 1.0)$$

    i.
    ii.   The ws are weights that are learned

In a f
feed-
(FFN
conn
neurc
previ

Output Layer:

Hidden Layer:

Input "Layer":     $x_1$   $x_2$   $x_3$   $x_4$   $x_5$   1.0

bb.

    i.    In a fully connected feed-forward network (FFNN), each layer is connected to each neuron or input in the previous layer

    ii.    Notice that each neuron has a bias weight

    iii.    Each one of them is doing classification but they are working together

    iv.    The bigger the network, there are problems introduced

        1.  Cost of multiplication

        2.  Searching in the highly dimensional space

        3.  Vanishing weights/parameters $\rightarrow$ the probabilities are so small that they disappear/vanish (Same thing can happen here)

        4.  Some numbers can get very small or big (happens on deeper networks but not usually in wider networks)
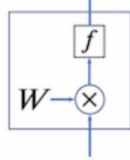
cc. Again, a little bit of linear algebra can make this a whole lot simpler:

$$W\mathbf{x} = \begin{bmatrix} w_{1,1} & w_{1,2} & w_{1,3} & w_{1,4} & w_{1,5} & b_1 \\ w_{2,1} & w_{2,2} & w_{2,3} & w_{2,4} & w_{2,5} & b_2 \\ w_{3,1} & w_{3,2} & w_{3,3} & w_{3,4} & w_{3,5} & b_3 \\ w_{4,1} & w_{4,2} & w_{4,3} & w_{4,4} & w_{4,5} & b_4 \end{bmatrix} \times \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_2 \\ x_4 \\ x_5 \\ 1.0 \end{bmatrix} = \begin{bmatrix} \mathbf{w}_1^T\mathbf{x} \\ \mathbf{w}_2^T\mathbf{x} \\ \vdots \\ \mathbf{w}_n^T\mathbf{x} \end{bmatrix}$$

    i.

$$\begin{bmatrix} w_{3,1} & w_{3,2} \\ w_{4,1} & w_{4,2} \end{bmatrix}$$

Activations (outputs)

$$\mathbf{a} = (a_1, a_2, a_3, a_4)$$



$$\mathbf{a} = f(W\mathbf{x})$$

$$\mathbf{x} = (x_1, x_2, x_3, x_4, x_5, 1.0)$$

Inputs

ii.

dd. Our FFNN:



$a_1$  $a_2$  $a_3$

$a_1$  $a_2$  $a_3$  $a_4$  1.0

$x_1$  $x_2$  $x_3$  $x_4$  $x_5$  1.0
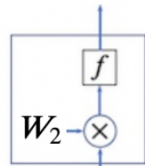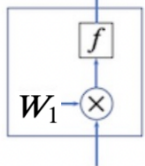
i.

$$\mathbf{a}^2 = (a_1^2, a_2^2, a_3^2)$$



$$\mathbf{a}^1 = (a_1^1, a_2^1, a_3^1, a_4^1) + [1.0]$$



$$\mathbf{x} = (x_1, x_2, x_3, x_4, x_5, 1.0)$$

ii.  $$a^2 = f(W_2(f(W_2\mathbf{x}) + [1.0])$$

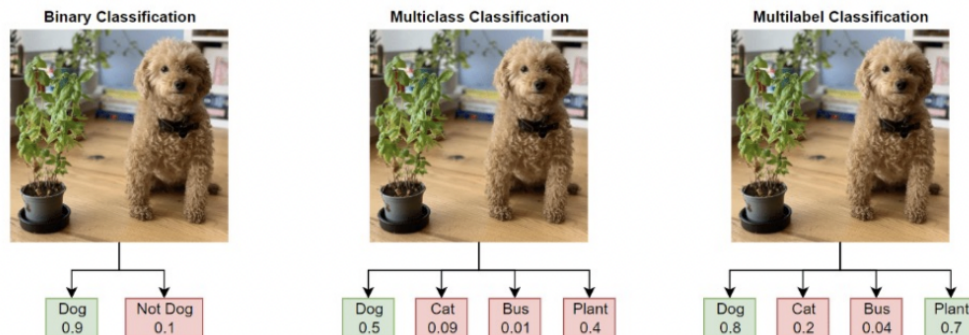1. Do non-linear activation function every layer

2. Classification Methods: Supervised ML
    a. Input:
        i. a fixed set of classes C = {c1, c2,…, cJ}
        ii. a randomly-permuted set of labeled documents (d1,c1),….,(dn,cn) split into
            1. a training set (d1,c1),….,(dm,c)
            2. a testing set dm+1,….,dn (labels withheld)
    b. Output:
        i. A classifier $\gamma : d \rightarrow c$ trained the training set
        ii. The testing set with labels calculated by $\gamma$
        iii. Test results (confusion matrix, metrics, etc.)
3. Classification: Binary, Multiclass and Multilabel
    a. In Binary Classification, we have 2 labels, and we must choose one; often this is phrased as "something" or "not something" (spam or ham, misinformation or not, etc.)
    b. In Multiclass Classification, we have more than 2 labels, and our task is to assign a single label to each sample:
    c. In Multilabel Classification, we have more than 2 labels, and our task is to assign any appropriate labels (not just one):
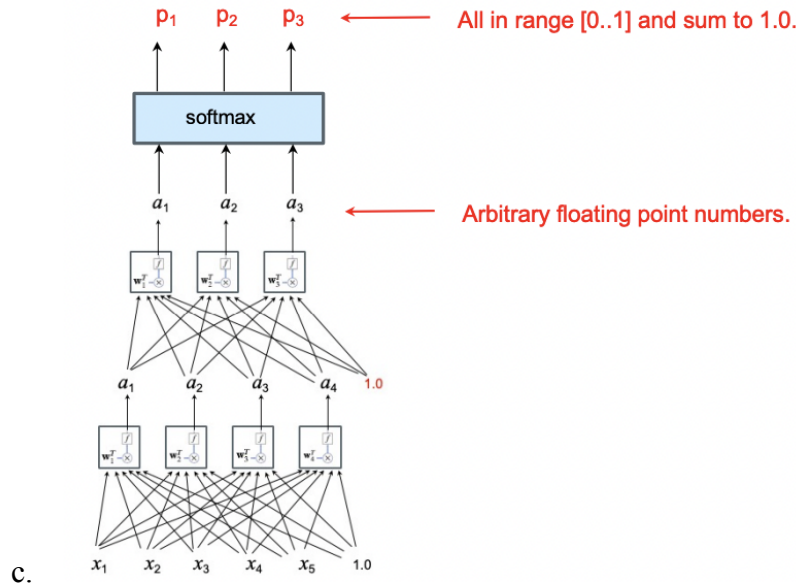


    d.
        i. Multilabel has a threshold $\rightarrow$ anything above the threshold is in the picture
4. Classification Methods: Supervised ML
    a. There are many different kinds of classifiers for labeled data
        i. Naïve Bayes
        ii. Logistic regression
        iii. Neural networks
5. Classification with Deep learning
    a. All of these types of classification are typically implemented by having the network output a probability distribution on all the labels.
    b. To convert the output values into a distribution, we used a generalization of the sigmoid called the softmax

All in range [0..1] and sum to 1.0.

Arbitrary floating point numbers.

c.

    i. No guarantee that the numbers coming out are good probability distributions (might not be between 0 and 1 and add to 1 and etc.)

    ii. Use softmax

    iii. Take any sequence of floating numbers and turn it into a probability distribution (it will maintain the order but scale the value)

    iv. Softmax exaggeratest the differences

6. Introduction to Deep Learning

    a. Softmax = a generalization of sigmoid which scales k numbers into a probability distribution

$$s(z) = \frac{e^z}{e^z + 1}$$

        i.

    b. For a vector z of dimensionality k, the softmax is

$$\text{softmax}(z) = \left[ \frac{\exp(z_1)}{\sum_{i=1}^{k} \exp(z_i)}, \frac{\exp(z_2)}{\sum_{i=1}^{k} \exp(z_i)}, ..., \frac{\exp(z_k)}{\sum_{i=1}^{k} \exp(z_i)} \right]$$

        i.

        ii. Easy to manipulate, eazy to take derivatives, and etc.

    c. Example

```
5  A = [1.2, 0.5, 3.2, 2.2, 1.9]
6  softmax(A)

array([0.07343397, 0.03646623, 0.54260772, 0.19961422, 0.14787785])
```

```
1  A = [1.2, 2.3]
2  softmax(A)

array([0.24973989, 0.75026011])
```

        i.

            1. 3.2 is not that far away from 2.2 but the softmax value is quite far away (emphasizing the larger ones)

7. Data Format for Deep Learning Classifier
   a. The training set consists of a matrix of features X and a vector of labels Y.
   b. For each input observation $x(i)$ , we have a vector of features $[x1, x2, ... , xn]$. Feature j for input $x(i)$ is $x_j$, more precisely $x_j(i)$.
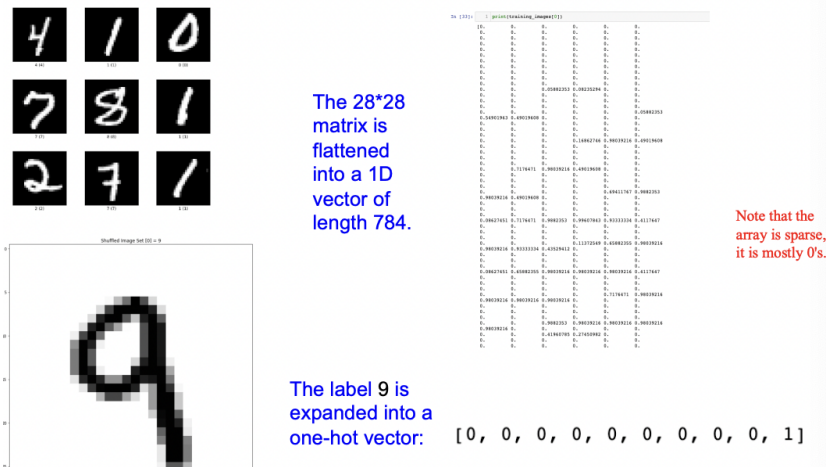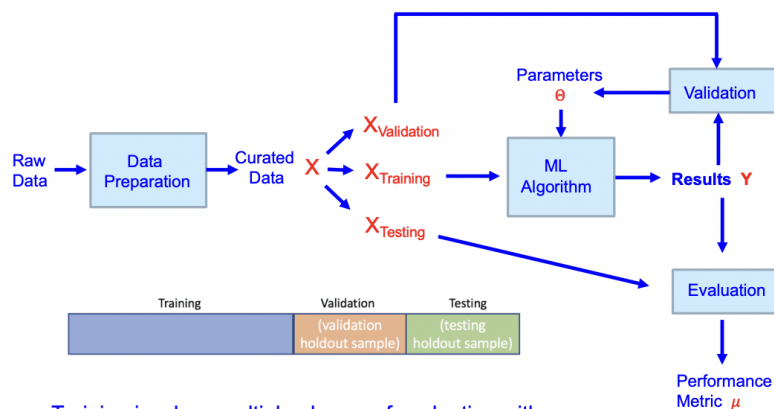


   c.

8. Multiclass Classification
   a. We will consider the MNIST database of handwritten digits: this is the "Hello World" of deep learning.
   b. The MNIST digit database consists of 70,000 28x28 BW pixel images, stored as 28*28=784 floating-point numbers in the range [0..1]; labels are integers 0..9:



The 28*28 matrix is flattened into a 1D vector of length 784.

Note that the array is sparse, it is mostly 0's.

The label 9 is expanded into a one-hot vector: [0, 0, 0, 0, 0, 0, 0, 0, 0, 1]

   c.

9. Recall: Supervised Machine Learning Workflow



Training involves multiple phases of evaluation with a validation set to find optimal values for the hyperparameters.

   a.

10. Cost/Loss Function for Classification
    a. One more detail! What is the cost (loss) function used with the output of a classifier?

Label: [0, 0, 0, 0, 0, 0, 0, 0, 0, 1]

Network output: $[p_0, p_1, p_2, p_3, p_4, p_5, p_6, p_7, p_8, p_9]$

What is the cost?

softmax

Output Layer: 10 neurons

Hidden Layer: 128 neurons

Input "Layer": 784 numbers: $[ x_0 \ x_1 \ x_2 \ x_3 \quad \dots \quad x_{783} ]$

Image Features

    b.
    c. We need to compare two probability distributions

Label: [0, 0, 0, 0, 0, 0, 0, 0, 0, 1] $\mathbf{y}$

Network output: $[p_0, p_1, p_2, p_3, p_4, p_5, p_6, p_7, p_8, p_9]$ $\mathbf{\hat{y}}$

to measure how different they are. The function used to do this is called the Cross-Entropy Loss:

$$\mathbf{\hat{y}} = softmax(\mathbf{a}) \qquad CE(\mathbf{\hat{y}}, \mathbf{y}) = -\frac{1}{n}\sum_{i=1}^{n} y_i \log(\hat{y}_i)$$

    d. One great feature of this cost function is that the derivative of the softmax plus CE Loss function is absurdly simple: it is just the difference of the two distributions

$$\frac{\partial CE(\hat{y}_i, y_i)}{\partial z_i} = \hat{y}_i - y_i$$

        i.