These Lectures are based on
two excellent, free on-line books
(see links on Deep materials page)

- Neural Networks and Deep Learning
  by Michael Nielsen

- Deep Learning by Goodfellow, Bengio
  and Courville

# MACHINE LEARNING BIG PICTURE

## Lecture ML9

training data $\longrightarrow$ in batches

Machine
$M_Z$
depends on
Parameters $Z$

$\longrightarrow$ Compare output
to correct answer

adjust.
$M$ to improve

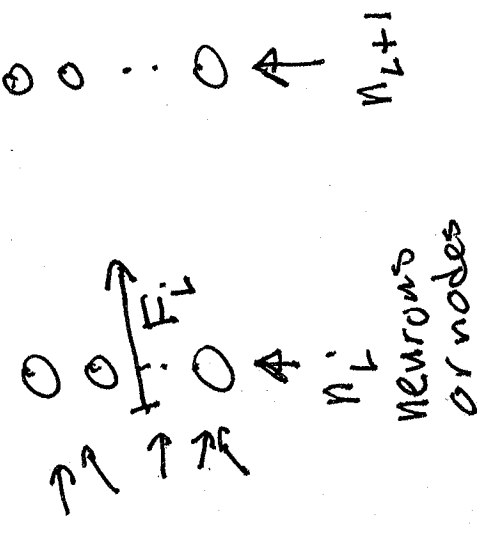DATA is sent in as mini-batches, after each
mini batch $M_Z$ is adjusted. After all the mini batches
have been "learned" the final machine $M_{Z_{final}}$
is run on test data to see how often it
yields correct answers.

— Many possible structures for MM, etc.

— We first cover the big math picture for a push forward, fully connected net (deep)

— Each layer is represented by a function

$$F_i: \text{1+ inputs } \bar{x} \in \mathbb{R}^{n_i} \text{ and outputs } \bar{@2} \in \mathbb{R}^{n_{L+1}}$$



— $F_i(x, A_i, \bar{b}_L)$ depends on the parameters

- $A_i$ an $(n_{L+1} \times n_i)$ matrix of weights so
  $$A: \mathbb{R}^{n_i} \to \mathbb{R}^{n_{L+1}}$$

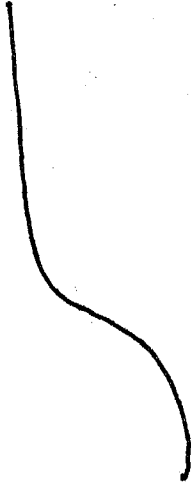- $b_L$ an $(n_{L+1} \times 1)$ vector $=$ the bias

- an activation function $\nabla$

$$F(x, A_i, b_i) = \nabla(A_L x + \vec{b_L})$$

The activation function has variou forms

• step function



1 ———

• Sigmoid

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



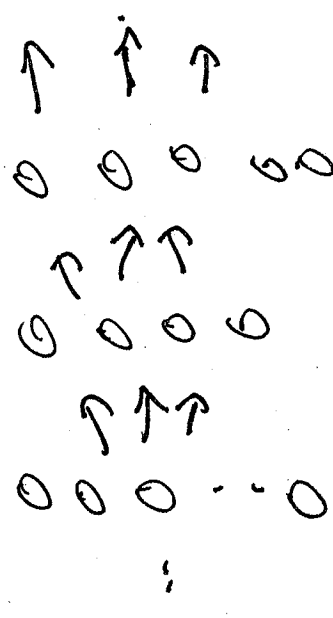• ramp or ReLu

$$\nabla(z) = max(z, 0)$$

- The activation function is vectorized, ie, it acts on each component of a vector

- So $\sigma(z_1, \ldots, z_n) = (\sigma(z_1), \ldots, \sigma(z_n))$

And $F_L : \mathbb{R}^{n_L} \to \mathbb{R}^{n_{L+1}}$

- Putting the pieces together from many layers

$$\begin{array}{cccc} 0 & 0 & 0 & \nearrow \\ 0 & 0 & \nearrow & \nearrow \\ 0 & \nearrow & \nearrow & \rightarrow \\ 0 & \nearrow & \nearrow & \rightarrow \\ & \cdots & 0 & 0 \\ 0 & & 0 & 0 \end{array}$$

$$G_\eta = F_K \circ F_{K-1} \circ \cdots \circ F_1, \quad \eta = (A_1, \vec{b}_1, A_2, \vec{b}_2, \ldots, A_K, \vec{b}_K)$$

all the weights and bias together ( lots of parameters!).

- Now we train the machine with training data $\vec{x}_1, ..., \vec{x}_n$ which are correctly characterized as $\vec{y}_1, ..., \vec{y}_n$.

- Now throw in all the training data and construct the cost or objective or error function

$$\Phi_m (\vec{x}_1, ..., \vec{x}_n) = \frac{1}{N} \sum_{i=1}^{N} | G_m (\vec{x}_i) - y_i |^2$$

(This is simplest, least squares version. More sophisticated versions later)

- We treat this as a function of $m$ and use an optimization $\Theta$ routine to diminish $\Phi_m$ to $\Phi_{m'}$

- Repeat with $m \to m'$.

- In practise, a radomn susset of training data is thrown in, $M_\theta$ is adjusted, then another mini batch, etc.

○ The goal is to get a $M_\theta$ given by $G_{M_\theta}$ to get a $M_\theta$ that generalizes ie. works well on @test that generalizes ie. works well on @test data that is not in the trainingset

. A big issue is how much to optimize $M$ for just the trainging set. Don't want the machine to memorize the training data and not generalize to other test data

- This is called over-fitting.

● Now the question is why $F_{M_\theta}$ takes this form?

○ This is connected to why it is called a <u>neural net</u>

- It will be easier to understand
the form of a neural net after
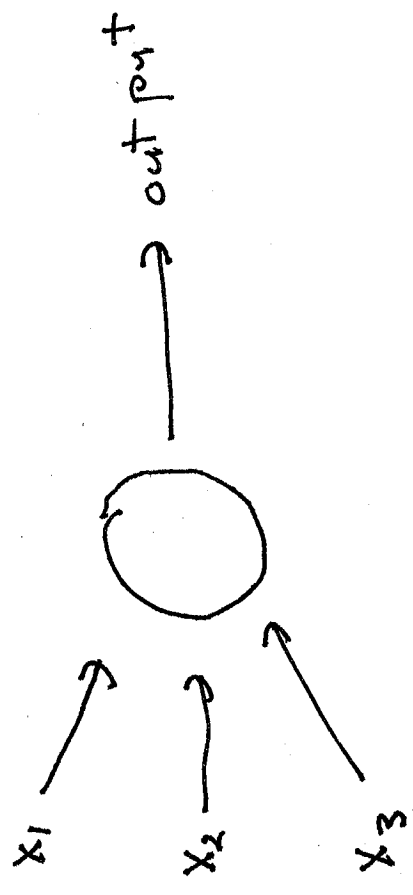we know a little about actual neurons

- Video on youtube by Marc Dingman

Lecture MLb →

Artificial neurons and
a single layer net

• We first describe a simple artificial
neuron called "the Perceptron".



$x_1$
$x_2$
$x_3$
→ output

• The output is zero or one (fire or don't fire)
weights

• The neuron uses the input using weights
$w_1, w_2$ and $w_3$

• A threshold $-b$ is set (minus sign explained later)
also called the bias

- Rule: Output is Zero if

$$w_1x_1 + w_2x_2 + w_3x_3 \leq -b$$

output is one if

$$w_1x_1 + w_2x_2 + w_3x_3 > -b$$

- Example: you are trying to decide whether to do your math #W tonight
  - $x_1 =$ how close is the due date
  - $x_2 =$ how long is the HW
  - $x_3 =$ what your friends are doing tonight

- you weigh up these various factors and make a decision $0 = no, 1 = yes$.
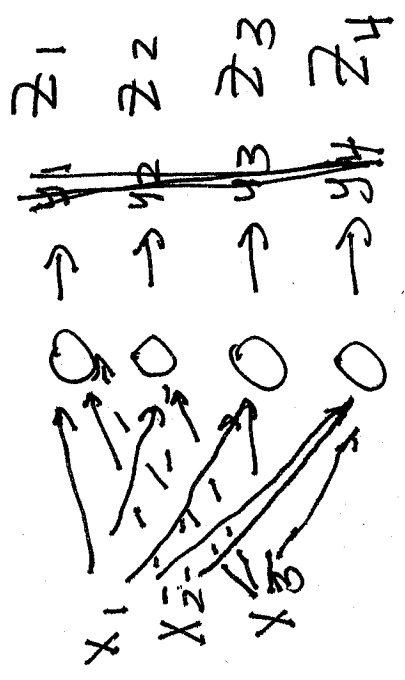
- We want to express the decision process more succinctly.

○ Let $\nabla(z) = 0$    $z \leq 0$    activation function
            $= 1$    $z > 0$

○  Let $F(x) = \nabla\left(\vec{w}^T \vec{x} + b\right)$    with $\vec{w}^T = [w_1, w_2, w_3]$

$$\vec{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$$

Then $F(x) = 0 \Rightarrow$ no
            $= 1 \Rightarrow$ yes.

○ Now we want to model a more complicated
   decision or classification problem using
   multiple perceptron

$x_1 \to$ ○ $\to y_1$  $z_1$

$x_2 \to$ ○ $\to y_2$  $z_2$

$x_3 \to$ ○ $\to y_3$  $z_3$

$\to$ ○ $\to y_4$  $z_4$

- each to the $k$'th perceptron has weights $\vec{w}_k$ and threshold or bias $b_k$

  we get for each $k$    $k = 1, \ldots, m$   $m = \#$ of neurons

  $$z_{,k} = \nabla\left(\vec{w}_k^T \vec{x} + b_k\right)$$

- We want to combine all these into a matrix form

  Let $W = \begin{bmatrix} \vec{w}_1 & \cdots & \vec{w}_m \end{bmatrix}$,   $\vec{b} = \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix}$

- then $W^T \vec{x} + \vec{b} = \begin{bmatrix} \vec{w}_1^T \\ \vdots \\ \vec{w}_m^T \end{bmatrix} \vec{x} + \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix}$

  $$= \begin{bmatrix} \vec{w}_1^T \vec{x} + b_1 \\ \vdots \\ \vec{w}_m^T \vec{x} + b_m \end{bmatrix} = \begin{bmatrix} z_1 \\ z_2 \\ \vdots \\ z_m \end{bmatrix}$$

- The last step is to vectorize the activation $\nabla$

$$\nabla \begin{bmatrix} z_1 \\ \vdots \\ z_m \end{bmatrix} = \begin{bmatrix} \nabla(z_1) \\ \vdots \\ \nabla(z_m) \end{bmatrix}$$

- So letting $A = W^T$, the weight matrix, our one layer machine is described by

$$F(\vec{x}) = \nabla\left(A\vec{x} + \vec{b}\right)$$

- We now change our point of view on the one layer of perceptions and treat it as a learning machine.

Lecture MLC

Learning and Multiple Layers

- So we return to the characterization problem

   So $\vec{x}_1, \vec{x}_p$ are inputs with correct
   outputs $\vec{y}_1, \vec{y}_2, \ldots, \vec{y}_p$
   
   $A \uparrow$  $\vec{b} \uparrow$ (the weights and bias)

- We fix a value of the parameters (the weights and bias)
   and for each $\vec{x}_i$ input the machine outputs

$$F(\vec{x}_i) = \nabla(A\vec{x}_i + \vec{b})$$

○ For each $i$, this has least squares error

$$E_i = \|(A\vec{x}_i + \vec{b}) - \vec{y}_i\|^2$$

● For the total error over all inputs we average these

$$\Phi(A, \vec{b}) = \frac{1}{p} \sum_{i=1}^{p} E_i = \frac{1}{p} \sum_{i=1}^{p} \|(A\vec{x}_i + \vec{b}) - \vec{y}_i\|^2$$

- Now we optimize, i.e. find $A_f$ and $\vec{b}_f$ which

MINIMIZE $\boxed{\Phi(A,B)} \; \Phi(A,\vec{b})$

- We then declare our final machine to

be $\otimes \nabla(A_f \vec{x} + \vec{b}_f)$

- Note the similarity to least squares

and polynomial fitting

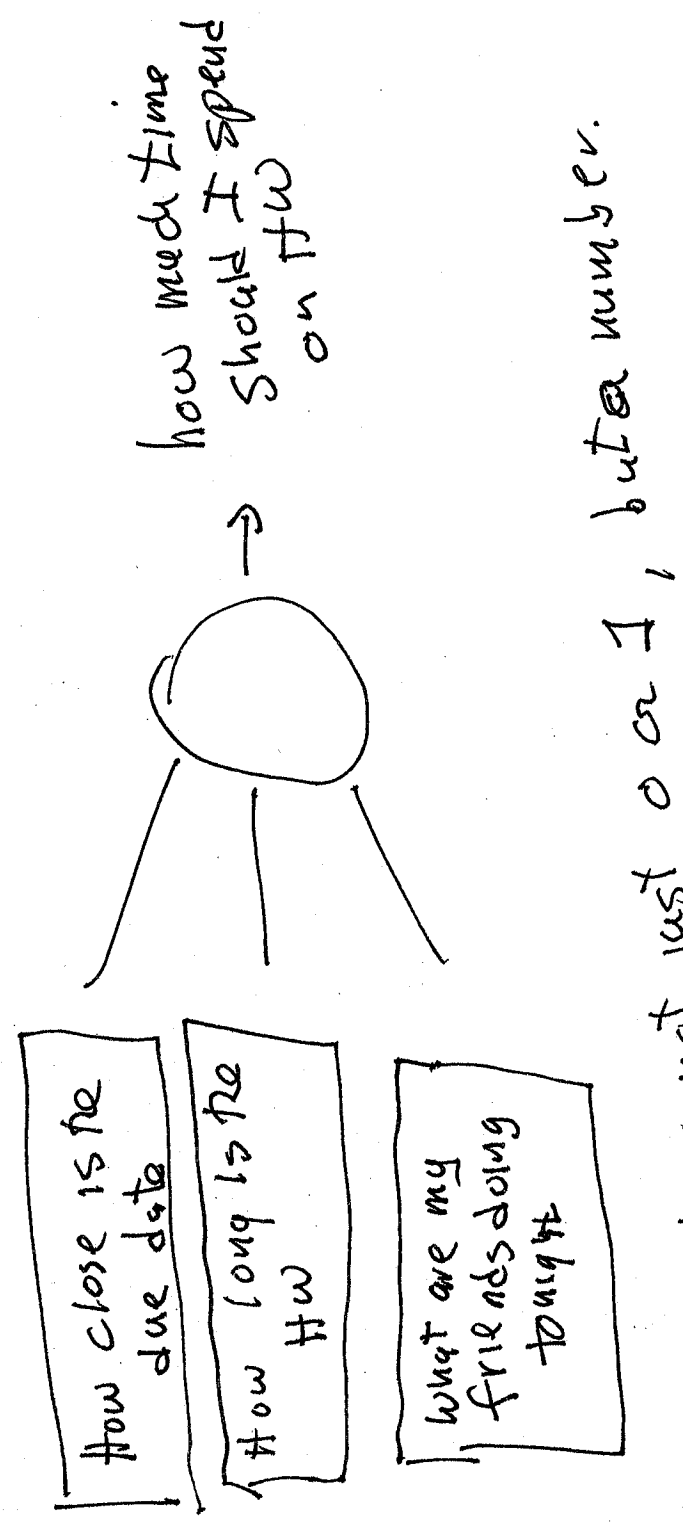- Looking more closely, how do we

optimize?

- usual thing is to differentiate $\Phi$

with respect to $A$ and $b$, etc.

- But $\nabla$ is not differentiable.

o Another issue with T is that it is restrictive, binary output

Back to the #w example



how much time should I spend on #w →

How close is the due date

How long is the #w
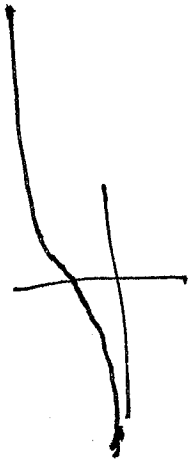
what are my friends doing tonight

The output is not just 0 or 1, but a number.

So we probably want $\sigma$ at least continuous
or maybe differentiable. We want small changes
in the parameter to yield small changes in the output

- The sigmoid $\sigma(z) = \dfrac{1}{1+e^{-z}}$    note: $\sigma(z) \to 1$ as $z \to \infty$

  $\sigma(z) \to 0$ as $z \to -\infty$

  is nice and differentiable
  but computationally expensive

- The ramp $\sigma(z) = MAX(z, 0)$

  is continuous, it's "derivative"
  is  ____    which is not so bad

  is  ____  0

  and is computationally tame.

We don't need to specify which $\sigma$ for now,
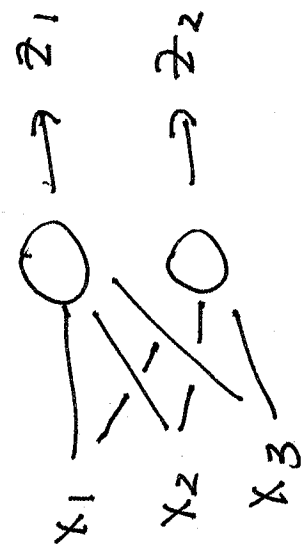but eventually /we will/ $f_X/f_{he}$ the activation function $\sigma$

For Now, let $\sigma$ be the sigmoid for theoretical ease

Let's study optimization or learning for one level

$$F(x, A, b) = \nabla(A_{\#} x + b_{\#})$$

Let's have three inputs and 2 neurons



$$x_1 \quad \longrightarrow z_1$$
$$x_2 \quad \longrightarrow z_2$$
$$x_3$$

So
$$z_1 = \nabla\left( (w_{11} x_1 + w_{12} x_2 + w_{13} x_3 + b_1) \right)$$

$$z_2 = \nabla\left( (w_{21} x_1 + w_{22} x_2 + w_{23} x_3 + b_2) \right)$$

and if $z_0 +$ Rue values are $y_1$ and $y_2$ we have

$$\Phi = \left( \nabla(w_{11} x_1 + w_{12} x_2 + w_{13} x_3 + b_1) - y_1 \right)^2$$
$$+ \left( \nabla(w_{21} x_1 + w_{22} x_2 + w_{23} x_3 + b_2) - y_2 \right)^2$$

/2

- we want to minimize the error $\Phi$ as a function of the parameters i.e. the w's and b's

- So we treat $\Phi$ as a function of these and compute $\nabla \Phi$ and find critical points and see if they are loc max/min or saddles

- For example, $\nabla \Phi = \begin{bmatrix} \dfrac{\partial \Phi}{\partial \omega_{11}}, \ldots, \dfrac{\partial \Phi}{\partial \omega_{28}}, \ldots, \dfrac{\partial \Phi}{\partial b_1}, \dfrac{\partial \Phi}{\partial b_2} \end{bmatrix}$

with

$$\frac{\partial \Phi}{\partial \omega_{11}} = \sigma'( \text{same argument}) \cdot x_1$$

$$\frac{\partial \Phi}{\partial b_1} = \sigma'(\text{same argument}) \cdot 1$$

by the chain rule

- This is complicated for just this simple one layer but we need many layers with many neurons and maybe thousands of spasronal of parameters.

- So we need new idea

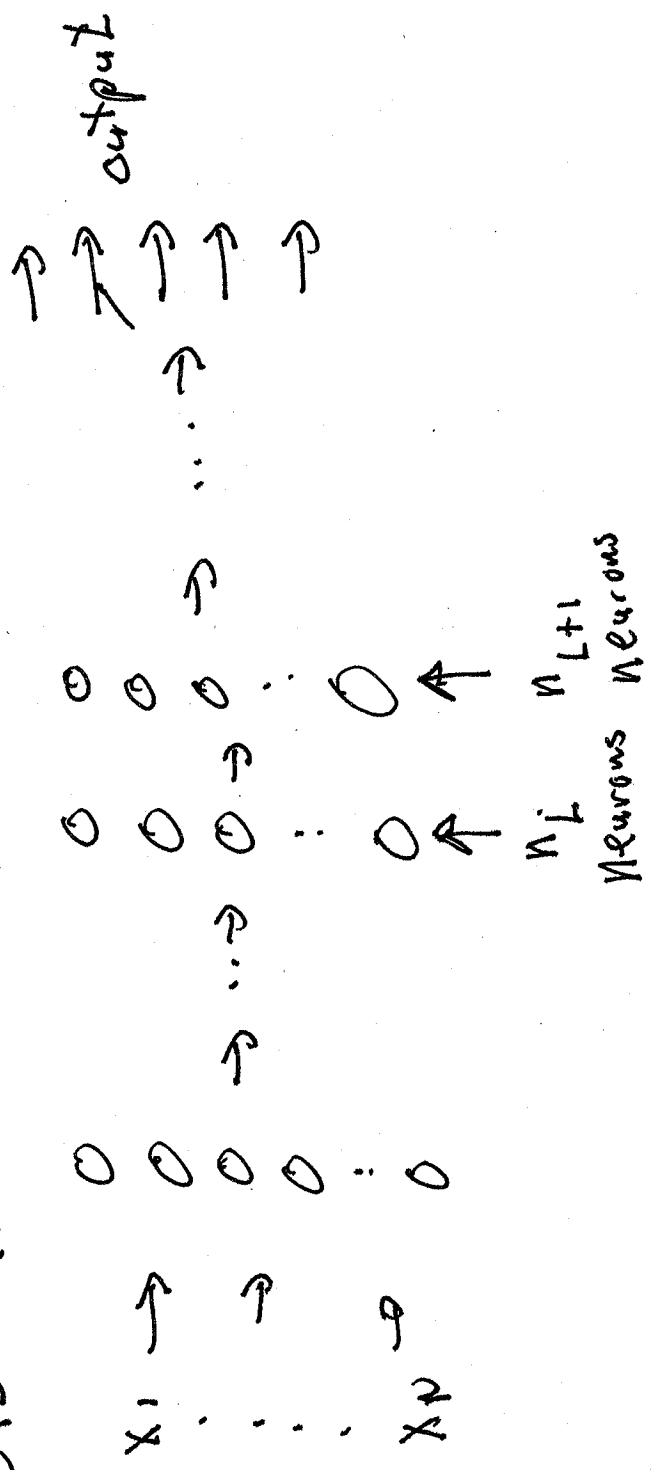(1) A better optimization scheme namely
**Gradient Descent**

(2) A clever way of computing $\nabla \Phi$ when
there are many layers

We will cover each of these in more detail
in later lectures

- Now to finish the Introduction we describe
multiple layers — This is the "Deep" in deep
learning

- one way to Think of This is decision
making in stages

• For example, first you decide how much time to allot to your math HW tonight, then you decide how what order to fit it in with your other HW.



$$x_1 \quad \cdots \quad x_n$$

$n_L$
Neurons

$n_{L+1}$
Neurons

output

• Each layer is given by a function

$$\vec{F}_i \left( \vec{x}, A_i, \vec{b}_i \right) = \nabla \left( A_i \vec{x} + \vec{b}_i \right)$$

PIX

- The layers act sequentially

$$F_1 \text{ Then } F_2 \text{ Then } F_3 \dots F_L$$

- Mathematically, This is a composition (recall it is written in the reverse order)

$$F = F_L \circ F_{L-1} \circ \dots \circ F_2 \circ F_1$$

- Then the least squares error is

$$\Phi = \frac{1}{2} \sum_{i=1}^{W} \| F(x_i) - y_i \|^2 \text{ and}$$

A_i and b_i. So $\nabla \Phi$

it depends on all the A_i and b_i.

to compute $\nabla \Phi$ is a Chore

- This net is called "Feed forward" Since Information just flows in one direction

$$\text{Input} \longrightarrow \text{output}$$