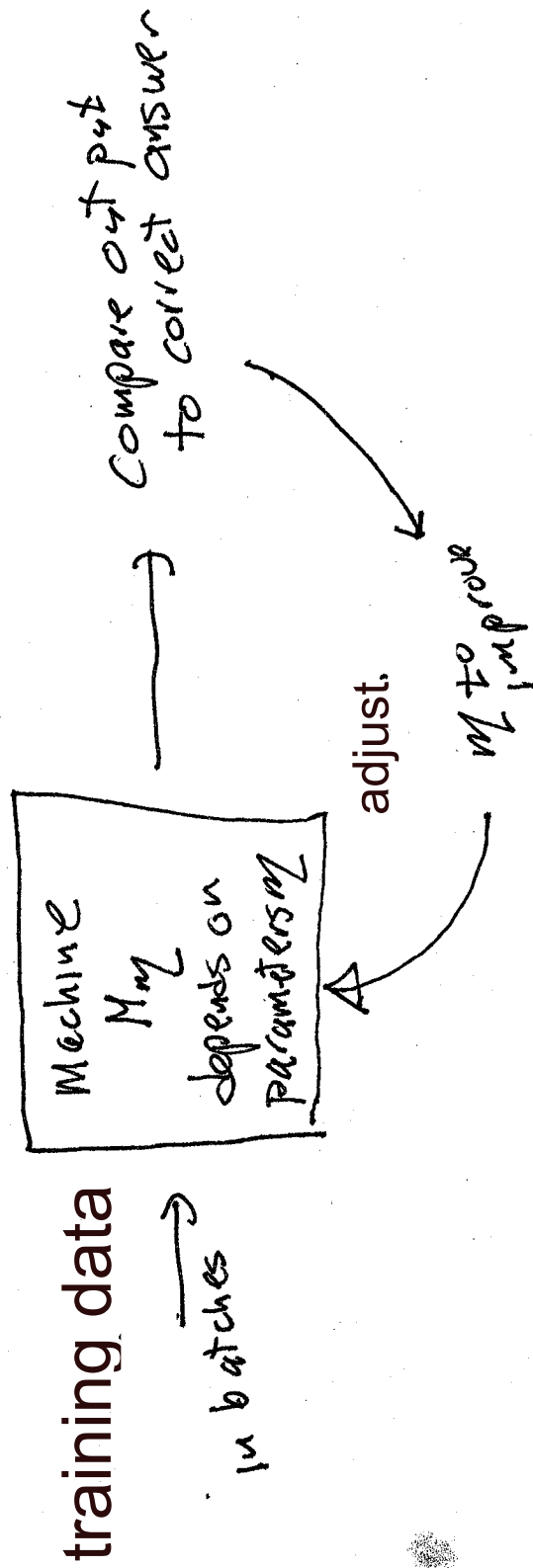


These Lectures are based on
two excellent, free on-line books
(links on syllabus page)

- Neural Networks and Deep Learning
by Michael Nielsen
- Deep Learning by Goodfellow, Bengio
and Courville.

MACHINE LEARNING BIG PICTURE

Lecture ML9



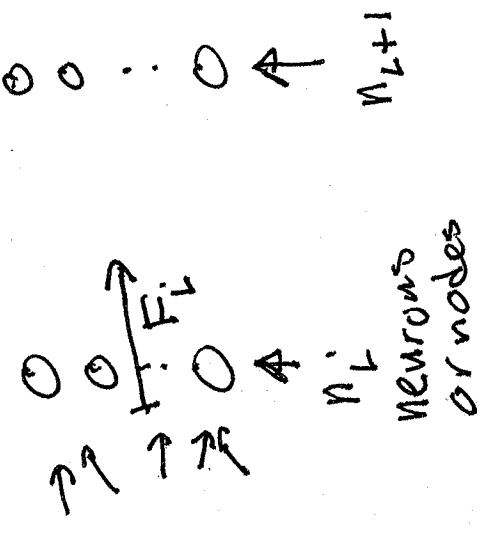
DATA is sent in as MINI-batches, after each mini batch θ is adjusted. After all the mini batches have been "learned" the final machine θ_{final} is run on test data to see how often it yields correct answers.

- Many possible structures for MLP, etc.
- We first cover the big math picture for a push forward, fully connected net

deep,

- Each layer is represented by a function $\vec{z} \in \mathbb{R}^{n_{l+1}}$

F_l : it inputs $\vec{x} \in \mathbb{R}^{n_l}$ and outputs $\vec{z} \in \mathbb{R}^{n_{l+1}}$



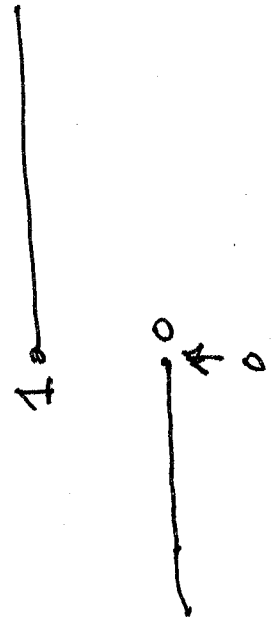
- $F_l(x, A_l, \vec{b}_l)$ depends on the parameters

- A_l an $(n_{l+1} \times n_l)$ matrix of weights so
- $A_l: \mathbb{R}^{n_l} \rightarrow \mathbb{R}^{n_{l+1}}$
- \vec{b}_l an (n_{l+1}) vector = the bias
- an activation function σ

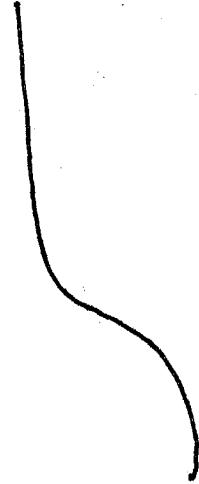
$$F(x, A_i, b_i) = \sigma(A_i x + b_i)$$

The activation function has various forms

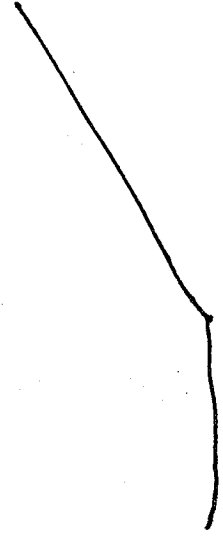
- Step function



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



- Ramp or ReLU $\sigma(z) = \max(z, 0)$

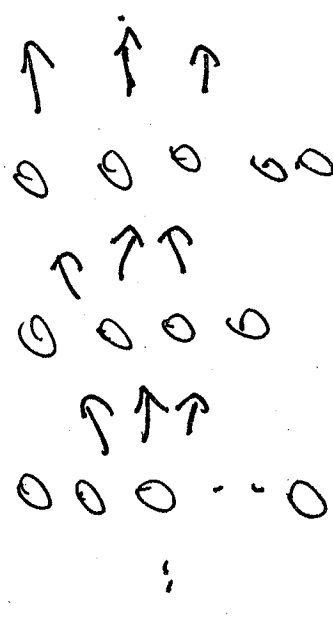


- The activation function is vectorized, i.e. it acts on each component of a vector

- so $\sigma(z_1, \dots, z_n) = (\sigma(z_1), \dots, \sigma(z_n))$

and $F_L: \mathbb{R}^{n_L} \rightarrow \mathbb{R}^{n_{L+1}}$

- Putting the pieces together from many layers



$$G_M = F_k \circ F_{k-1} \circ \dots \circ F_1, \quad \eta = (A_1, \vec{b}_1, A_2, \vec{b}_2, \dots, A_k, \vec{b}_k)$$

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

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

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

$$\Phi_m(\vec{x}_1, \dots, \vec{x}_n) = \frac{1}{N} \sum_{l=1}^N |G_m(\vec{x}_l) - g_l|^2$$

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

• We treat this as a function of m and use an optimization routine to diminish Φ_m to Φ_m^*

• Repeat with $m \rightarrow m'$

- In practice, a random subset of training data is thrown in, m is adjusted, then another minibatch, etc.

- The goal is to get a M_m given by G_m that generalizes i.e. works well on test data that is not in the training set

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

- This is called over-fitting. F_m takes

- Now the question is why F_m takes this form?

- This is connected to why it is called a neural net

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

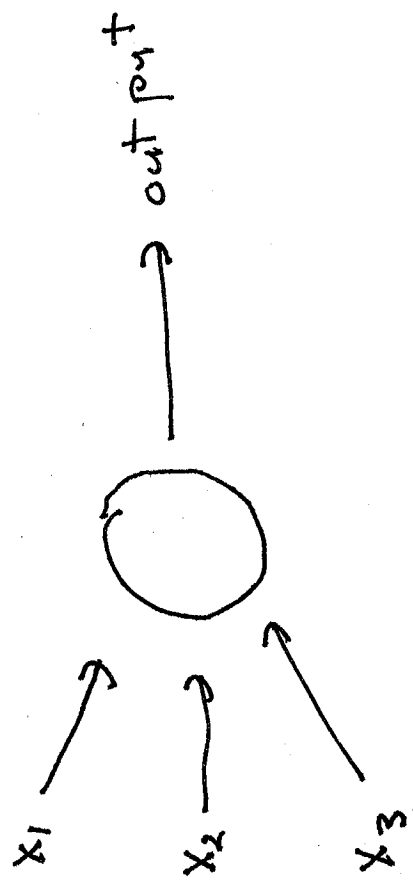
- Video on Youtube by Marc Dingman

<https://www.youtube.com/watch?v=6qS83wD29PY>

Lecture MLb →

Artificial neurons and a single layer net

~~17~~
We first describe a simple artificial
neuron called "ne perceptron"



- The out put is zero or one (fire or don't fire)
weights the input using weights
- The neuron weights w_1 , w_2 and w_3
 w_3 (minus sign explained later)
- A threshold $-b$ is set (minus sign explained later)
also called the bias

- Rule! output is zero if

$$w_1 x_1 + w_2 x_2 + w_3 x_3 \leq -b$$

output is one if

$$w_1 x_1 + w_2 x_2 + w_3 x_3 > -b$$

whether

Example: You are trying to decide whether to do your math HW tonight

- $x_1 =$ how close is the due date
- $x_2 =$ how long is the HW
- $x_3 =$ what your friends are doing tonight

various factors and

• You weigh up these various factors and make a decision 0=no, 1=yes.

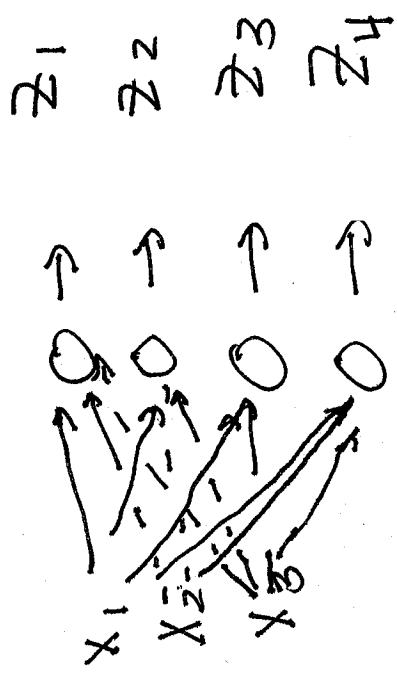
• We want to express the decision process

more succinctly.

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

Let $F(x) = \sigma(\tilde{w}^T \tilde{x} + b)$ with $\tilde{w}^T = [w_1, w_2, w_3]$
 $\tilde{x} = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \end{bmatrix}$
 then $F(x) = 0 \Rightarrow$ NO
 $\quad \quad \quad = 1 \Rightarrow$ YES.

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



• ~~each~~ k 'th perceptron has weights \vec{w}_k and threshold or bias b_k

we get for each k $K = 1, \dots, M = \# \text{ neurons}$

$$z_k = \sigma(\vec{w}_k^T \vec{x} + b_k)$$

• We want to combine all these into a matrix form

$$\vec{z} = \begin{bmatrix} z_1 \\ \vdots \\ z_m \end{bmatrix} = \vec{w} + \begin{bmatrix} b_1 \\ \vdots \\ b_m \end{bmatrix}$$

$$\text{Let } \vec{w} = \begin{bmatrix} \vec{w}_1 \\ \vdots \\ \vec{w}_m \end{bmatrix}$$

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

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

• The last step is to vectorize the activation τ III

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

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

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

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

12*

Lecture MLC

Learning and Multiple Layers

• So we return to the characterization problem

with correct

so $\vec{x}_1, \dots, \vec{x}_p$ are inputs with correct outputs $\vec{y}_1, \vec{y}_2, \dots, \vec{y}_p$

A

b

↑ weights and bias

• We fix a value of the parameters (the machine outputs

and for each \vec{x}_L input the

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

least squares error

• For each i , this has

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

• For the total error over all inputs we average these

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

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

$$\Phi(A, \vec{b})$$

minimize

• We then declare our final machine to be $\nabla(A_f \vec{x} + \vec{b}_f)$

• Note the similarity to least squares and polynomial fitting

and poly more closely, how do we

looking more closely, how do we optimize?

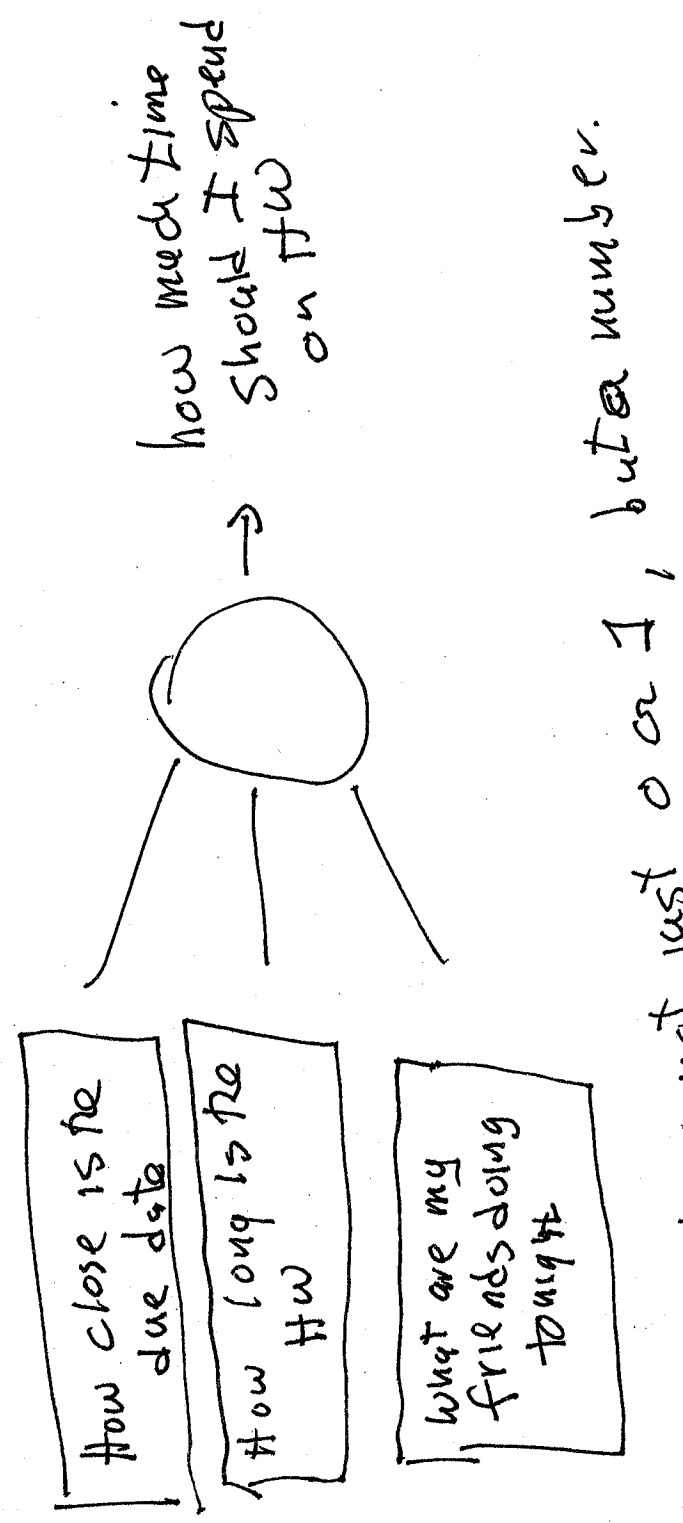
$$\Phi$$

• Usual thing is to differentiate with respect to A and \vec{b} , etc.

• But ∇ is not differentiable.

Another issue with Γ is that it is restrictive, binary output

BACK TO THE HW example

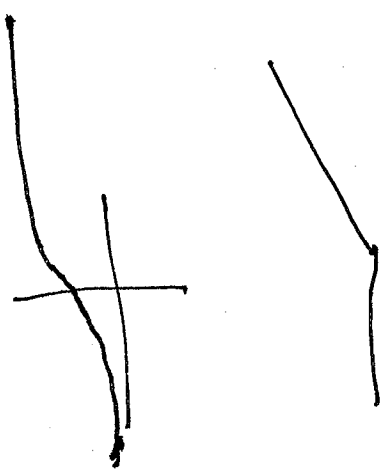


The output is not just 0 or 1, but a number.

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

• The sigmoid $\sigma(z) = \frac{1}{1 + e^{-z}}$

note! $\sigma(z) \rightarrow 1$ as $z \rightarrow \infty$
 $\sigma(z) \rightarrow 0$ as $z \rightarrow -\infty$



is nice and differentiable but computationally expensive

• The ramp $\sigma(z) = \max(z, 0)$

is continuous, it's "derivative"

is 1 which is not so bad

and is computationally tame.

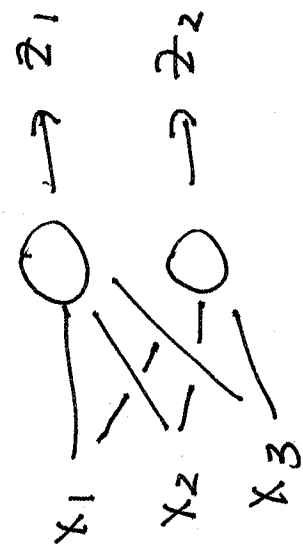
For now, let σ be the sigmoid for theoretical ease

Let's study optimization or learning for one level

$$F(x, A, b) = \sigma(Ax + b)$$

Three inputs and 2 neurons

Let's have



$$z_1 = \sigma(w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + b_1)$$

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

we have

and if $R_e + R_{ne}$ values are y_1 and y_2

$$\Phi = \left(\sigma(w_{11}x_1 + w_{12}x_2 + w_{13}x_3 + b_1) - y_1 \right)^2$$

$$+ \left(\sigma(w_{21}x_1 + w_{22}x_2 + w_{23}x_3 + b_2) - y_2 \right)^2$$

$$/ 2$$

- We want to minimize the error Φ as a function of the parameters w 's and b 's
- So we treat Φ as a function of these and compute $\nabla\Phi$ and find critical points
- and see if they are local max, min or saddle
- For example, $\nabla\Phi = \left[\frac{\partial\Phi}{\partial w_{11}}, \dots, \frac{\partial\Phi}{\partial w_{23}}, \frac{\partial\Phi}{\partial b_1}, \dots, \frac{\partial\Phi}{\partial b_2} \right]$

with

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

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

(σ (same argument) - y_{-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 parameters.

18

• So we need new ideas

(1) A better optimization scheme namely
Gradient Descent $\nabla \Phi$ when

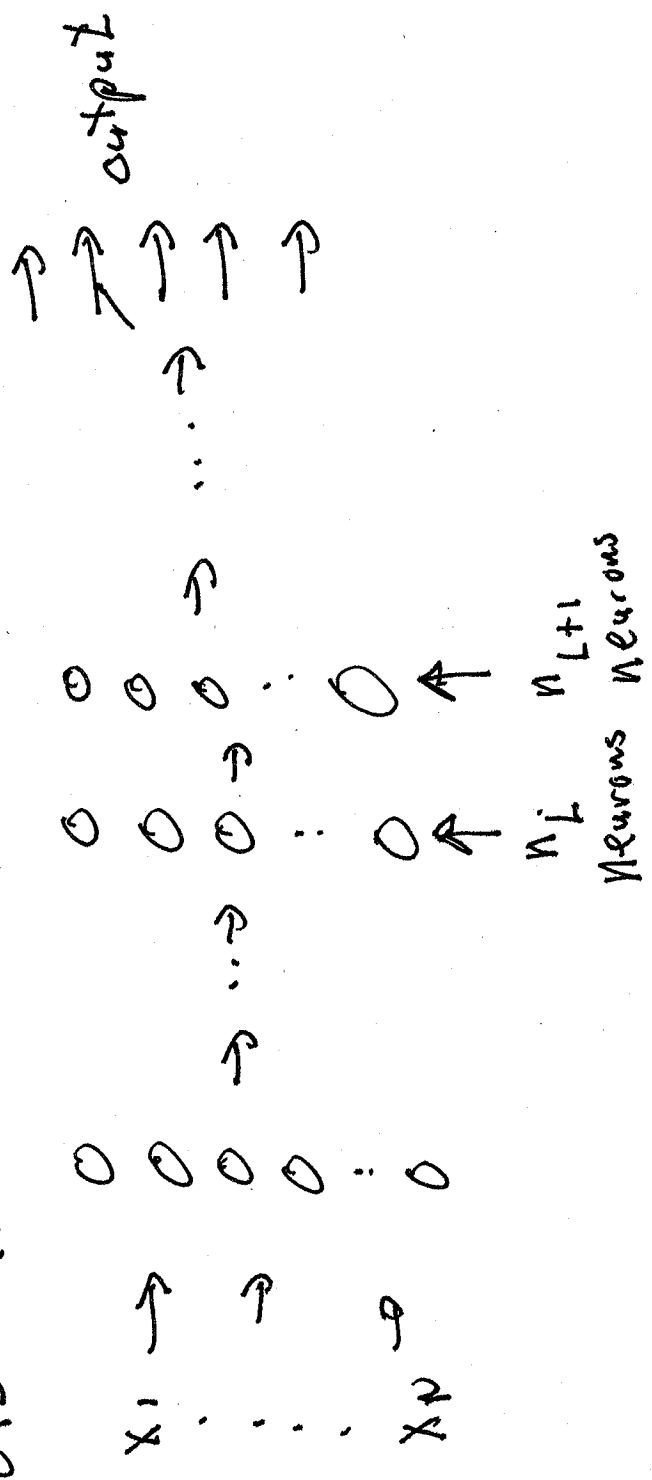
(2) A clever way of computing $\nabla \Phi$ when
there are many layers

• 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 what order to fit it in with your other HW.



Each layer is given by a function

$$F_i(\vec{x}, A_i, \vec{b}_i) = \sigma(A_i x + \vec{b}_i)$$

PIX

- The layers act sequentially

F_1 then F_2 then $F_3 \dots F_L$

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

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

- The least squares error is $\sum_{i=1}^N \|F(x_i) - y_i\|^2$ and

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

it depends on all the A_i and so $\nabla \Phi$ is a chore to compute since

- This net is called "Feed forward"

Information just flows in one direction
input \rightarrow output