

## MATH 382W Lec 24

Short Survey on other methods we didn't have time for.

### • Deep Learning

- Image Recognition using CNN's (soon)

- Transformers (not covered)

- Machine Teaching: Assume you can generate  $\vec{x}_*$ 's from a div. <sup>its more expensive</sup> cheap  $\vec{y}_*$  is a begin with  $\mathbb{D} \xrightarrow{A} g$ . Then repeat  $\vec{e}$ . Which  $|e_i|$ 's are large?

Generate  $\langle \vec{x}_*, y_* \rangle \approx \vec{x}_i$ 's that produced large  $|e_i|$ 's

Let  $\mathbb{D} = \left\{ \begin{bmatrix} X \\ \vec{x}_* \end{bmatrix}, \begin{bmatrix} Y \\ y_* \end{bmatrix} \right\} \xrightarrow{A} g$ . Repeat until model has low  $|e_i|$ 's

for all  $\vec{x}$  locations or until you run out of resources.

- Reinforcement Learning: same as machine teaching except you generate the  $\vec{x}$  entries by interacting with a dynamic <sup>system which</sup> then gives you  $y$ 's.

For instance: playing a video game.  $y = \#$  of points you get or beat the level?  $x$ 's are everything happening on the screen, and the record of your "moves".

### • Unsupervised Learning

You only see  $X$ ; there is no  $\vec{y}$ . Typical problems,

- Clustering the units into groups  $G_1, G_2, \dots, G_K \subseteq \{\vec{x}_1, \vec{x}_2, \dots, \vec{x}_n\}$   
<sup>hierarchical clustering, K-means clustering</sup>

- Anomaly detection. if the  $\vec{x}_i$ 's are realizations from a data generating process, which  $\vec{x}_i$ 's don't seem to be following the process?

- Dimension reduction: instead of  $p$  features, can you squeeze the same amount of information out of  $p' < p$  dimensions?

Principal Components Analysis (PCA) and Factor Analysis

K-means clustering. Consider data  $X$  which is prenormalized (all col means = 0 and std = 1).

- (1) Pick the # of clusters  $K$ , and distance function  $d$  (default is Euclidean)
- (2) Randomly assign  $K$  of the  $n$  units to be the "centroids"  $\vec{x}_1, \vec{x}_2, \dots, \vec{x}_K$   
 $G_1 = \{\vec{x}_1\}, G_2 = \{\vec{x}_2\}, \dots, G_K = \{\vec{x}_K\}$
- (3) For each remaining  $n-K$  units, compare the distance to each of the  $K$  centroids i.e. for all  $i$

$$d_1 = \|\vec{x}_i - \vec{x}_1\|^2$$

$$d_2 = \|\vec{x}_i - \vec{x}_2\|^2$$

⋮

$$d_K = \|\vec{x}_i - \vec{x}_K\|^2$$

- (4) Now assign the  $i$ th unit to the cluster which minimizes the distance to its centroid i.e.  $\vec{x}_i$  is added to group  $G_{k^*}$  where  $k^* = \argmin \{d_k\}_{k=1}^K$

- (5) Recompute centroids as

$$\vec{x}_{k^*} = \frac{1}{|G_{k^*}|} \sum_{\vec{x} \in G_{k^*}} \vec{x}$$

- (6) For all units, repeat steps 3-6 until "convergence" (where no membership changes for any unit in their group)

- (7) Calculate <sup>within</sup> total dist. for each cluster  $D_k := \sum_{\vec{x} \in G_k} d(\vec{x}, \vec{x}_k)$  <sup>e.g.</sup> then add all  $D_k$  distances for all clusters to get a fit metric  $\sum_{k=1}^K D_k$

- (8) Repeat steps 1-7 with different starting pts and report best starting pt

- (9) Repeat steps 1-7 for  $K$  values  $\in \{1, 2, \dots, K_{max}\}$

- (10) select optimal  $K$  value visually using "the bend" as a guide ("scree plot")

# Artificial Neural Networks / ANN's

3

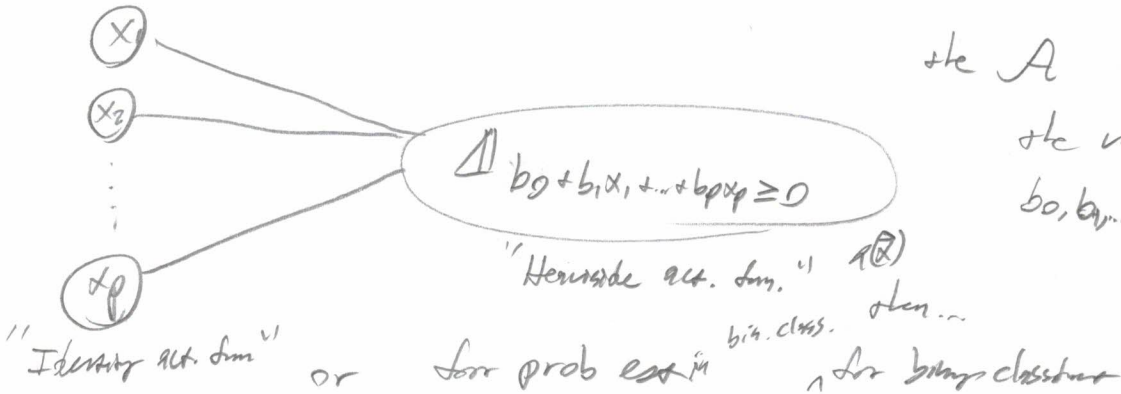
## Neural Networks / Neural "Nets" / NN's

"activation"

Another perspective on our models. Layer:  
For binary classification:

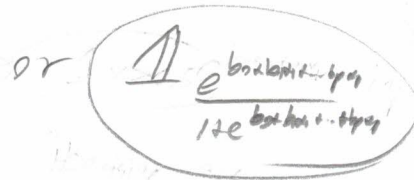
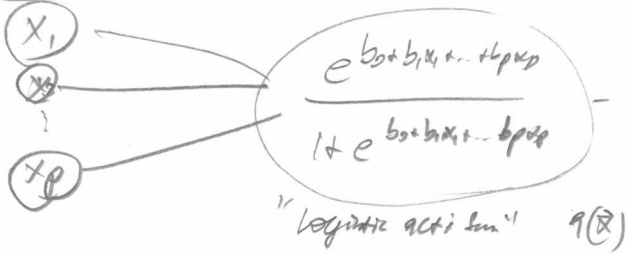
a set of  $L$  "neurons",  
which take input from a  
vector

"Input Layer" "Output Layer" ( $L=1$ )



input layer

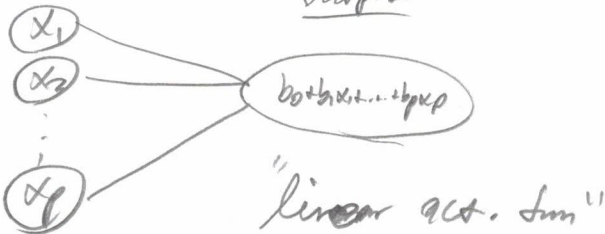
output layer / hidden layer



Input

or for regression

Output



Generally, activation functions are generalized linear functions

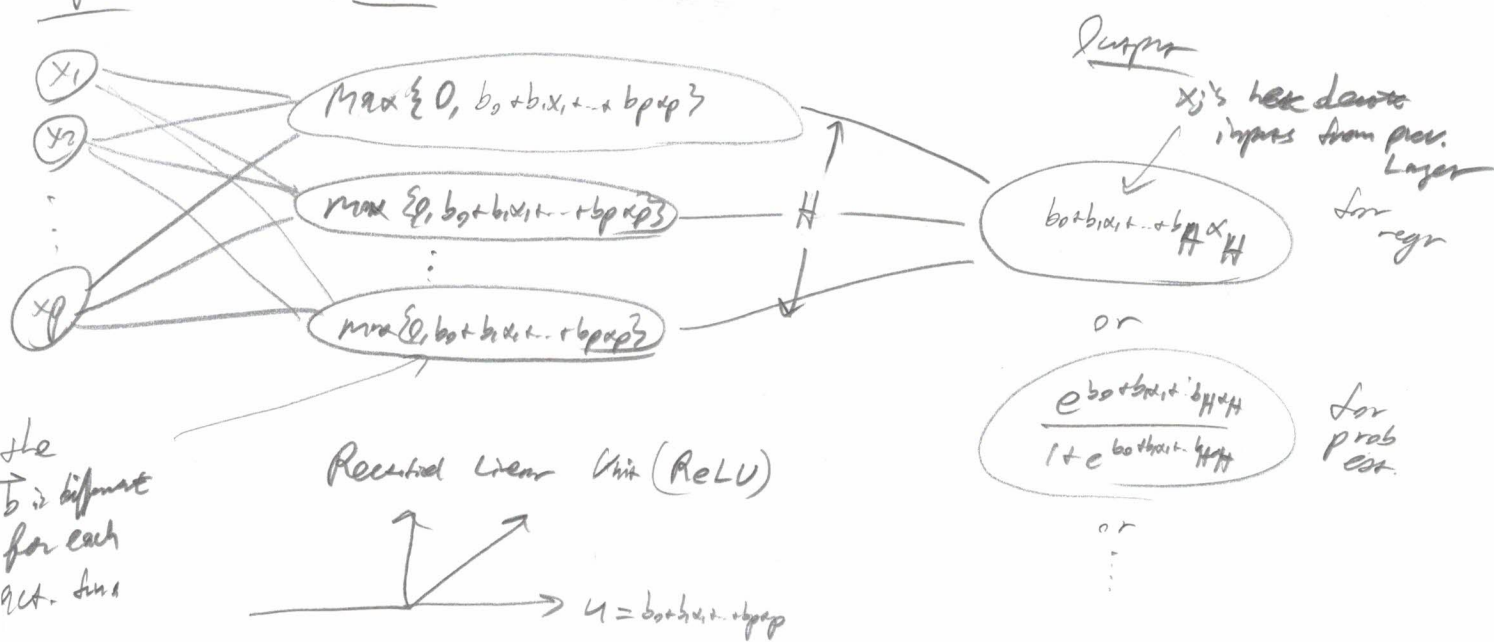
and can be written as  $g(u)$  where  $u = b_0 + b_1x_1 + \dots + b_px_p$  (like 1st function)

If so, how can NN's ever fit interactions? Just wait...

Network Topology: arrangement of Layers in a NN.  
 Multilayer Networks: has at least one hidden layer

Figure

Hidden



This Hidden Layer is an "alternative representation" of the input  $\vec{x}$ .  
 It is also "fully connected" meaning each act. fun.'s use all  $x_1, \dots, x_p$ .

ANN's are called universal approximators i.e. they can fit any function to arbitrary precision. In 2017, they proved fully connected NN's using ReLU's with  $p+1$  fully connected hidden layers and final output.

width: how many hidden layers (0)

height/depth: how many act. fun.'s per layer  $\{H_1, H_2, \dots, H_L\}$

You can create a network of arbitrary width, depth and specification of activation functions. How to fit all the b's?

First specify objective function to minimize  $L(\vec{y}, \vec{\hat{y}})$  or  $L(\vec{y}, \vec{\hat{p}})$

Regression: SSE, Prob Est. log loss:  $-\sum x_i \ln(\hat{p}_i) + (1-y_i) \ln(1-\hat{p}_i)$

generally, NN's are trained to do prob. est. since this obj. function is differentiable.



② Now, pick a random starting point to all  $w$ 's. E.g.  $U(-1,1)$  <sup>iid draws from</sup> [5]

① "Feed the data forward": do all calculations for all  $i=1, \dots, n$  for every activation function in all hidden layers and outer layer. Here, weights are fixed.

② Compute the gradients of all activation functions. How?

The entire network can be expressed as a function  $g(\vec{x})$  <sup>collection of  $H_{D-1}$  functions</sup>  
 can think its composite  

$$= g_D(\vec{w}_D \vec{g}_{D-1}(\vec{w}_{D-1} \vec{g}_{D-2}(\dots$$
  

$$g_1(\vec{w}_1 \vec{x})))$$
 <sup>collection of  $H_0$  function</sup>  
 We seek to minimize 
$$L = \sum_{i=1}^n L(y_i, g(\vec{x}_i))$$
  
 So we compute 
$$\vec{\nabla} L = \begin{bmatrix} \frac{\partial L}{\partial w_{D,1}} & \dots & \frac{\partial L}{\partial w_{D,H_D}} \\ \frac{\partial L}{\partial w_{D-1,1}} & \dots & \frac{\partial L}{\partial w_{D-1,H_{D-1}}} \\ \vdots & & \vdots \\ \frac{\partial L}{\partial w_{1,1}} & \dots & \frac{\partial L}{\partial w_{1,H_1}} \end{bmatrix}$$

And then we set  $\vec{w}_{D,t} = \vec{w}_{D,t-1} + \eta \vec{\nabla} L$  i.e. gradient descent  
 They've iterated the last layer  $D$  first!

Now we go to the second to last layer  $D-1$ .

We can take the gradient here as well using the chain rule, then do gradient descent to update the weights

$\vdots$   
 until the first layer. This is why it's called back-propagation (due to the chain rule)

③ Repeat steps 1, 2 until convergence or max # of iterations

As  $D$  and  $H_0$ 's get larger, the # of parameters increases  
 $\Rightarrow$  risk of overfitting. To mitigate overfitting,

- (1) you can use lasso / ridge penalty in each step of back-prop. to ~~push~~ <sup>push</sup> weights smaller
- (2) dropout: randomly set weights to zero during training to force neurons in other weights.
- (3) randomly add noise to layer inputs during training
- (4) keep a select set of data and monitor its performance over iterations

$\Rightarrow$  There are tons of hyperparameters!  $D$ ,  $H_1$ ,  $H_0$ , a function!, objective function,  $\eta$ , regularization hyperparameters!

So far we've duplicated a flexible basis set (such as trees)

$\Rightarrow$  deep neural nets should perform about equally to RF and xgboost. I've found it's rarely the case I think because the # of hyperparameters is so large

$\Rightarrow$  for "tabular" data where  $D = \langle X, Y \rangle$

$\uparrow$   
 a "table", a data frame, or spreadsheet

$\Rightarrow$  RF or xgboost is the way to go!

So why use NN's? Not tabular data e.g. images, text, sound, video 7

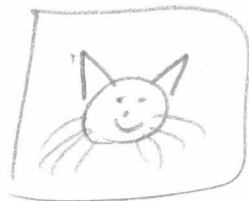
e.g.  
classify  
an  
image

Dog



vs.

Car



$Y=1$

$Y=0$

Why do NN's do well here? Each layer is able to represent the data in a different way. If you can "represent" dog cars vs. cat cars  $\Rightarrow$  single classification!

How do we do this? We first of all need to bundle

image input

RGB  
image

$L \times W \times 3$

usually...

arrays with pixel values  
 $\{0, 1, \dots, 255\}$