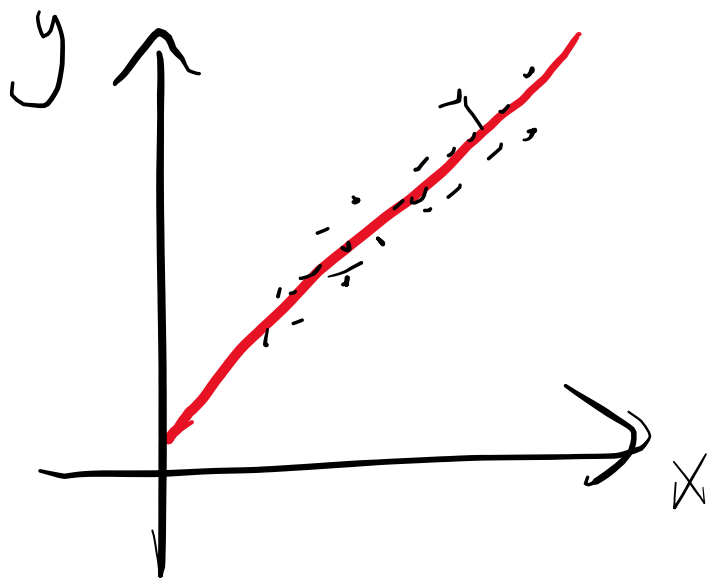


- Please join Slack workspace and channel
- Please leave your video on if you can.

A closer look at regression



in particular linear least squares regression.

$$L = \sum_{i=1}^N (y_i - f(x_i; w))^2$$

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial b} = 0.$$

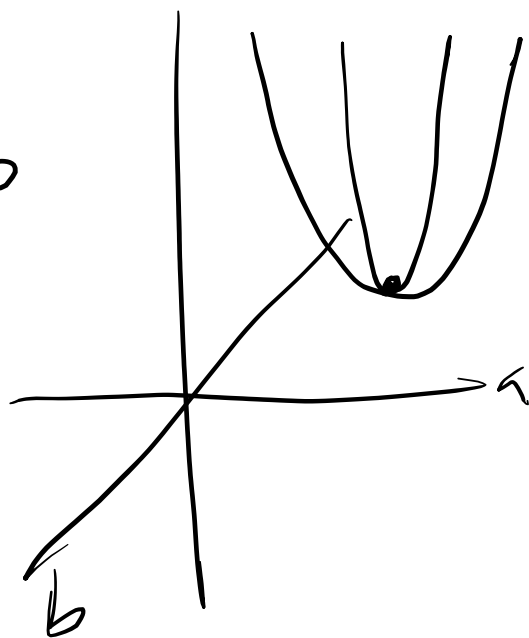
$$f(x) = ax + b$$

Mean squared error  
(MSE)

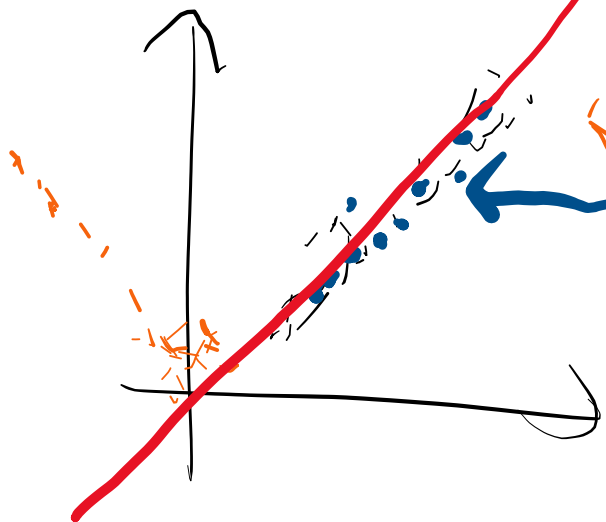
$L$  is convex wrt  $a, b$

$\exists$  a global min wrt  $a, b$

$$\frac{\partial L}{\partial a} = \frac{\partial L}{\partial b} = 0 \text{ trivial to solve}$$



can use to predict on new data



$f(x_j)$   $x_j \in$  new data

if  $x_j$  is previous domain  $\rightarrow$  interpolation ✓

if new data is outside range  $\longrightarrow$  extrapolation ~~X~~  
DANGEROUS

---

In ML: loss or objective function  
usually trying to minimize wrt parameters  
in general,  $L$  is not convex, minimum cannot  
be found analytically,  $f$  is not linear,  $x$  or  $y$   
are multivariate

In regression, what is special about MSE loss?

General, formal reason: Maximum Likelihood Estimation

maximize:  $L = P(\text{data} | \text{model})$

$-\log p$

"log likelihood"

usually

model:  $P(y | x_i, a, b) = N e^{-\frac{(y - (ax+b))^2}{2\sigma^2}}$

each data pt is iid

$$P(\text{data} | \text{model}) = \prod_{i=1}^N P(y_i | x_i, a, b)$$

$$= \exp(-MSE)$$

MLE:

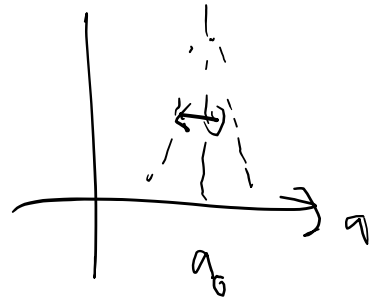
1. "consistency"

as data  $\rightarrow \infty$

estimated parameters  $\rightarrow$  true parameters  
( $a, b$ )

2. "efficiency"

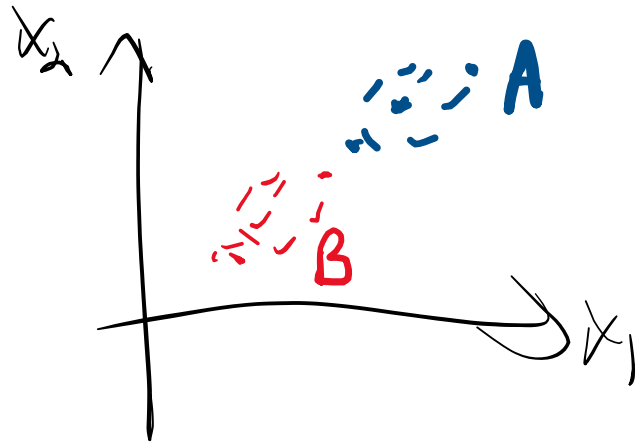
as data  $\rightarrow \infty$  min variance estimator



# • (Binary) Classification

iid = independent & identically dist'd

- What should loss be under MLE?



$f(x_j; w) = \text{prob. that } x \text{ is type A.}$

$$= P(A|x) = 1 - P(B|x)$$

$$P(\text{data} | \text{model}) = \prod_{i \in A} P(A|x_i) \prod_{j \in B} P(B|x_j)$$

$$= \prod_{i \in A} f(x_i; w) \prod_{j \in B} (1 - f(x_j; w))$$

$$L = -\log p = -\sum_{i \in A} \log f(x_i; w) - \sum_{i \in B} \log (1 - f(x_i; w))$$

"log loss"  
"binary cross entropy"

$$= -\sum_{i \in \text{data}} (y_i \log f_i + (1 - y_i) \log (1 - f_i))$$

$y_i = \begin{cases} 1 & (A) \\ 0 & (B) \end{cases}$

---

What family of fns to use?  $f(x; w)$

popular options  
in HGP

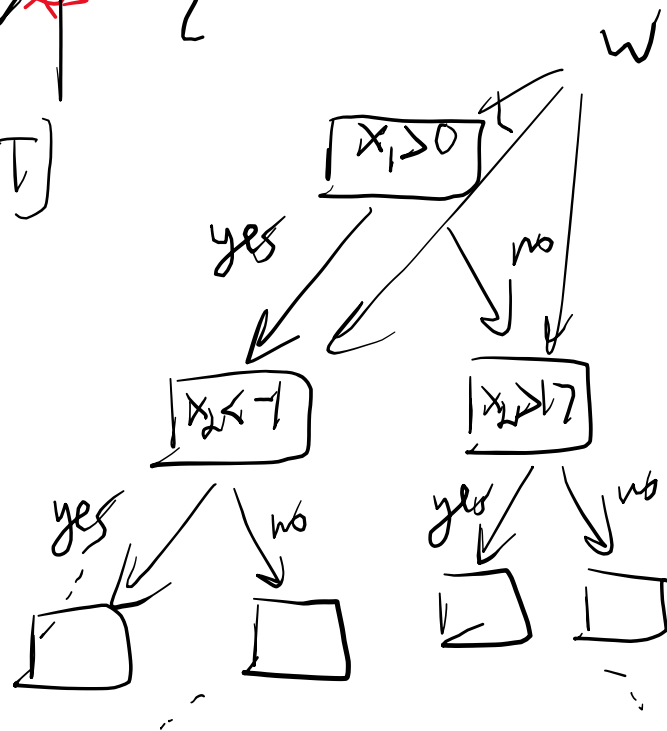
1.  $f(x; w) = f_0(x)$

"high level feature"



$$f_0(x) = r$$

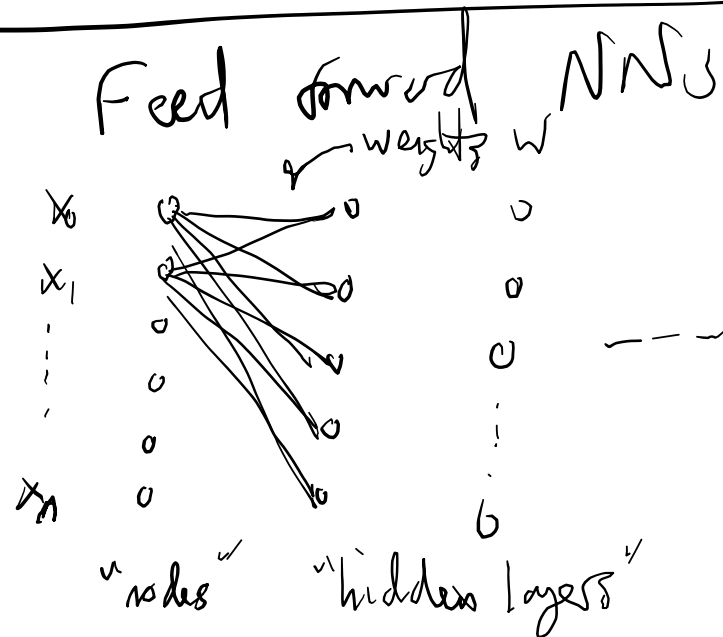
2. (Boosted) Decision Tree (BDT)





### 3. Neural Networks

- Biologically inspired
- Scale well to higher dimensions
- Expressive ("Universal Approx Thm")
- Generalizes unreasonably well.



"fully connected"  
"dense" NN

(DNN)

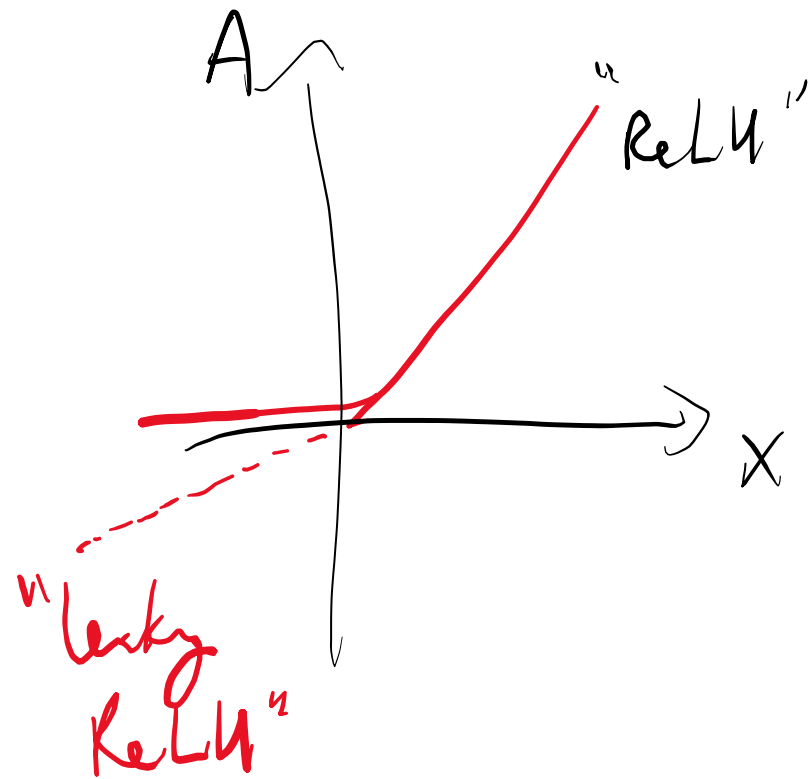
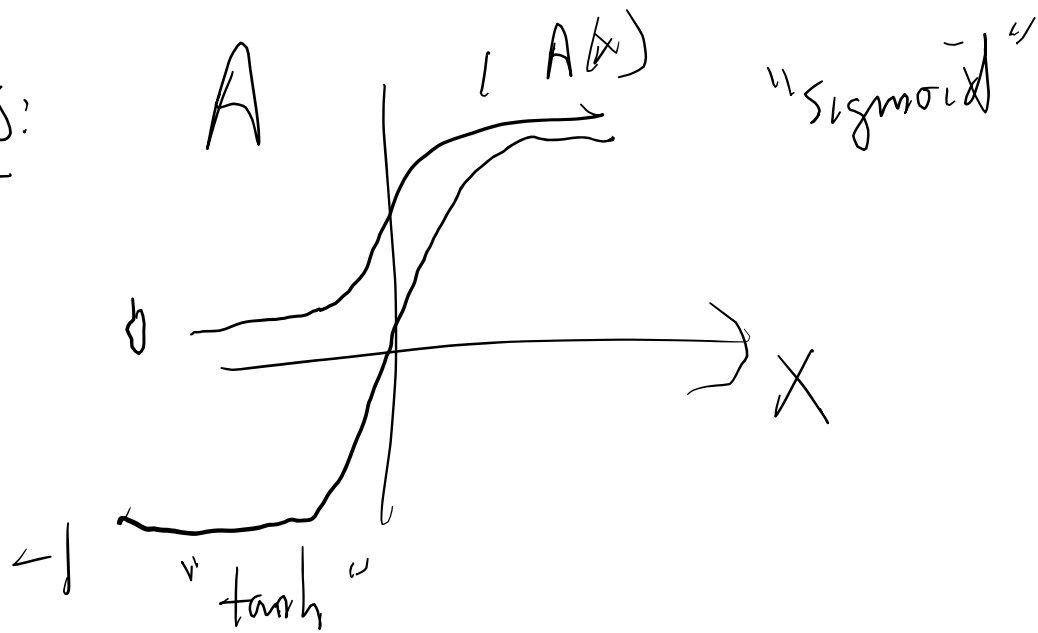
weights biases

output layer

$$\vec{x}^{(i)} = A([w^{(i)}], \vec{x}^{(i-1)} + b^{(i)})$$

elementwise activation fn.

Examples:



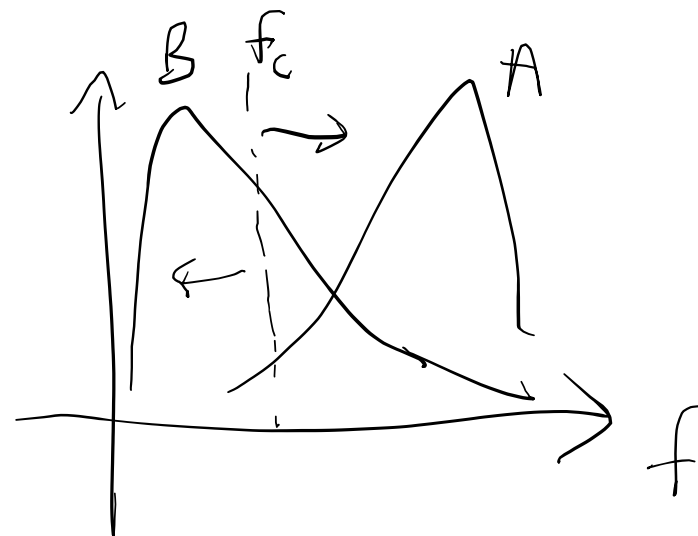
$$A(x) = \begin{pmatrix} A(x_0) \\ A(x_1) \\ \vdots \\ A(x_n) \end{pmatrix}$$

• Motives for binary classification w/

• Suppose we've learned  $f(x; w)$ . How do we measure / characterize its performance?

1.  $L[f(x; w)]$  loss fn is ultimate metric  
however it's not practically useful / interpretable

2.  $f(x; w) > f_c$  A  
 $< f_c$  B  
"working point"



Fraction of class A passing cut : "true positive rate" tpr  
 "signal efficiency"

" " B " " : "false positive rate" fpr  
 "background efficiency"

