

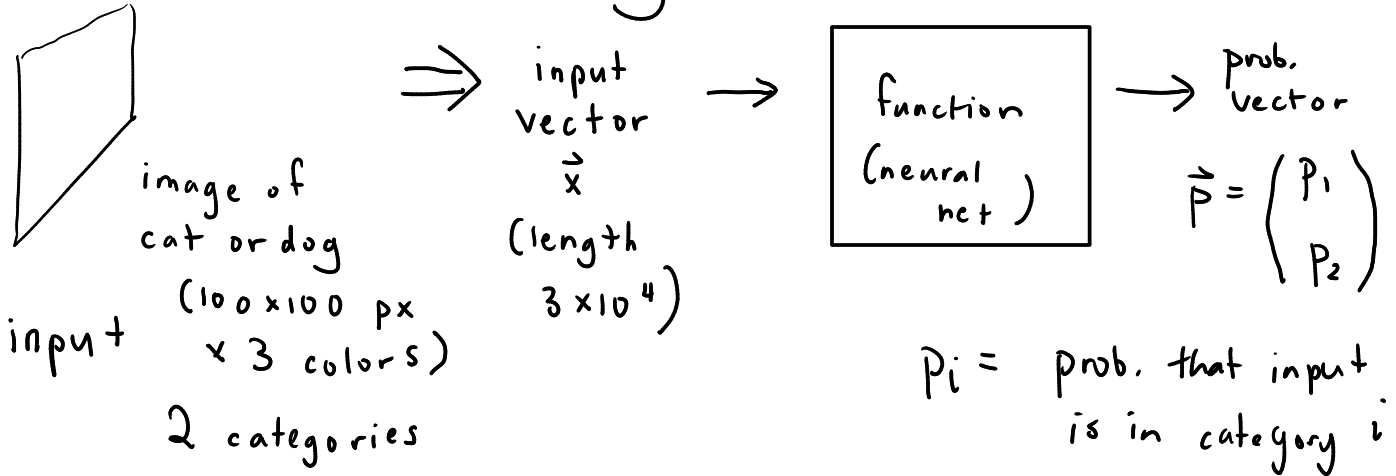
MAP : maximize  
w/ respect  
to  $\vec{w}$

$$P(\vec{w} | \mathcal{D}) = \frac{P(\mathcal{D} | \vec{w}) P(\vec{w})}{P(\mathcal{D})}$$

↑  
model  
params

⇒ find "best" set of params  $\vec{w}^*$

example: machine learning classification



What does the neural net func. look like?

one example: dense neural net (DNN)  
or multilayer perceptron

to find output  $\vec{P}$  we do a series of operations:

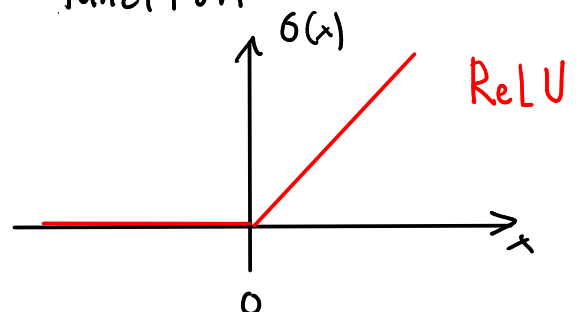
$$\vec{z}^{(1)} = \sigma \left( \underset{\substack{\downarrow \\ \text{matrix} \\ \text{(could be rectangular)}}}{A^{(1)}} \vec{x} + \underset{\substack{\downarrow \\ \text{vector}}}{b^{(1)}} \right)$$

one "layer" of net

$\sigma(x)$  = nonlinear "threshold" function

example:  $\sigma(x) = \max(0, x)$

$$\sigma(\vec{x}) = (\sigma(x_1), \sigma(x_2), \dots)$$



next layer:  $\vec{z}^{(2)} = \sigma \left( A^{(2)} \vec{z}^{(1)} + \vec{b}^{(2)} \right)$

...

$\vec{z}^{(n)} = \sigma \left( A^{(n)} \vec{z}^{(n-1)} + \vec{b}^{(n)} \right)$       n-th layer

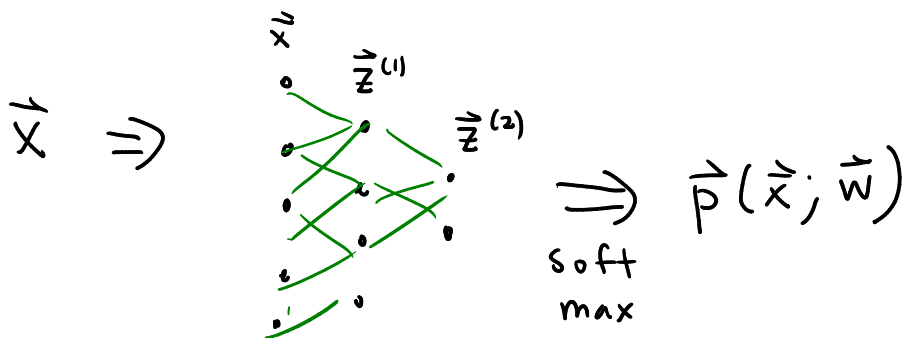
↓  
same dim. as  
# categories

overall func. depends  
all the components  
of  $A^{(i)}, \vec{b}^{(i)}$   $i=1, \dots, n$

⏟  
⇒ vector of params  $\vec{w}$

final step: convert  $\vec{z}^{(n)}$  into prob.  $\vec{p}$

"softmax": 
$$p_i = \frac{e^{-z_i^{(n)}}}{e^{-z_1^{(n)}} + e^{-z_2^{(n)}} + \dots}$$
       $\sum_i p_i = 1$   
 $p_i \geq 0$



universal approx. theorems:

in the case of DNNs if # layers

+ size of layers is large enough

⇒ you can approx. any func. that  
takes input from finite domain

Bayesian interpretation:

dataset  $\mathcal{D}$  : "training" data for net

"image"

$$\vec{x}^{(1)}$$

$$\vec{x}^{(2)}$$

$\vdots$

$$\vec{x}^{(N)}$$

label:

$$l_1 = \begin{cases} 1 \\ 2 \end{cases} \text{ category}$$

$$l_2$$

$\vdots$

$$l_N$$

$$\mathcal{D} = \{ (\vec{x}^{(i)}, l_i) \mid i=1, \dots, N \}$$

interested in  $\mathcal{P}(\vec{w} \mid \mathcal{D}) \Rightarrow$  first step is calculate likelihood

$$\mathcal{P}(\mathcal{D} \mid \vec{w})$$

$$\mathcal{P}(\mathcal{D} \mid \vec{w}) = \prod_{i=1}^N P_{l_i}(\vec{x}^{(i)}; \vec{w}) \leq 1$$

$= 1$  when perfect classifier

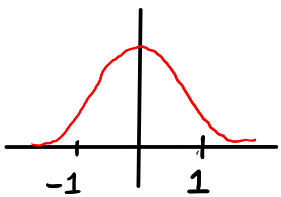
want to maximize  $\mathcal{P}(\vec{w} \mid \mathcal{D}) = \frac{\mathcal{P}(\mathcal{D} \mid \vec{w}) \mathcal{P}(\vec{w})}{\mathcal{P}(\mathcal{D})} \rightarrow$  "prior"?

"regularization": assume net params are not too big (to prevent blow-ups)

$$\mathcal{P}(w_i) = \frac{1}{\sqrt{2\pi}} e^{-w_i^2/2} \text{ Gaussian of width 1}$$

$$\mathcal{P}(\vec{w}) = \prod_j \mathcal{P}(w_j)$$

trick: to maximize  $\mathcal{P}(\vec{w} \mid \mathcal{D})$  we take  $-\log$  of it & minimize



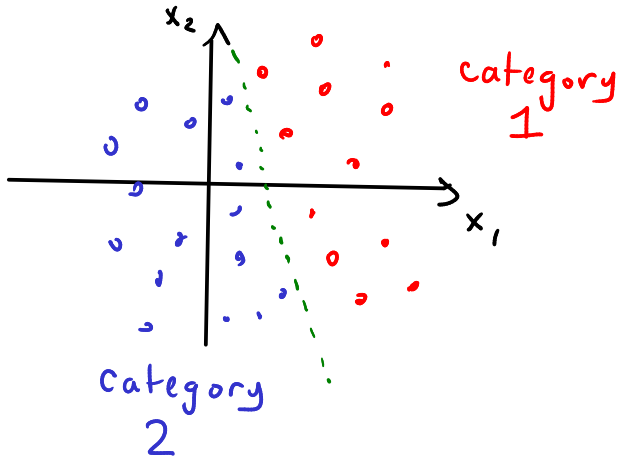
$$\mathcal{L}(\vec{w}) = -\log \mathcal{P}(\vec{w} | \mathcal{D})$$

$$= -\log \mathcal{P}(\mathcal{D} | \vec{w}) - \log \mathcal{P}(\vec{w}) + \log \mathcal{P}(\mathcal{D})$$

const.  
↓  
ignore

⇒ minimize  $\mathcal{L}(\vec{w})$  w/ respect to  $\vec{w}$   
(i.e. gradient descent)

prob.  
2  
of  
PS # 1



$$\vec{x} = (x_1, x_2)$$

$\mathcal{P}_1(\vec{x}; \vec{w})$   
= prob. that  $\vec{x}$  is  
in category 1

$$\mathcal{P}_2(\vec{x}; \vec{w}) = 1 - \mathcal{P}_1(\vec{x}; \vec{w})$$