

MAP:

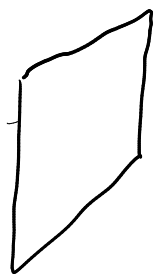
maximize  
w/ respect  
to  
 $\vec{w}$

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

↑  
model parameters

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

example: machine learning classification

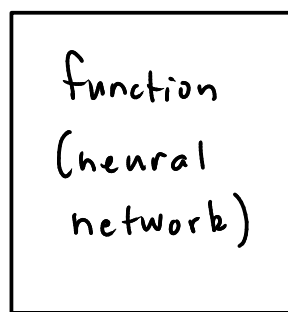


input:  
image (100x100 px  
x 3 colors)

2 categories  
(i.e. cat or dog)



input  
vector  
 $\vec{x}$   
(length  
 $3 \times 10^4$ )



prob.  
vector

$$\vec{p} = \begin{pmatrix} p_1 \\ p_2 \end{pmatrix}$$

$p_i$  = prob. that  
input is  
category  $i$

Neural network: dense neural net (DNN)  
or multilayer perceptron

Series of operations:

"one layer"  
of net

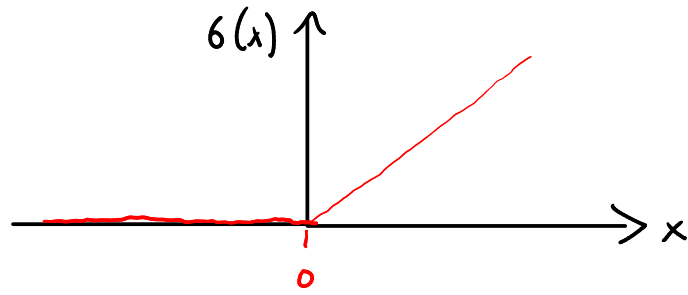
$$G \left( \begin{matrix} A^{(1)} & \vec{x} \\ \downarrow & \downarrow \\ \text{matrix} & \text{vector} \end{matrix} + \begin{matrix} \vec{b}^{(1)} \\ \downarrow \\ \text{vector} \end{matrix} \right) = \vec{z}^{(1)}$$

(could be rectangular)

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

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

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



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

o o o

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



same dim  
as # categories  
(here: 2)

Overall process  
depends on  
components:

$A^{(i)}, \vec{b}^{(i)} \quad i=1, \dots, n$

final step: "softmax"

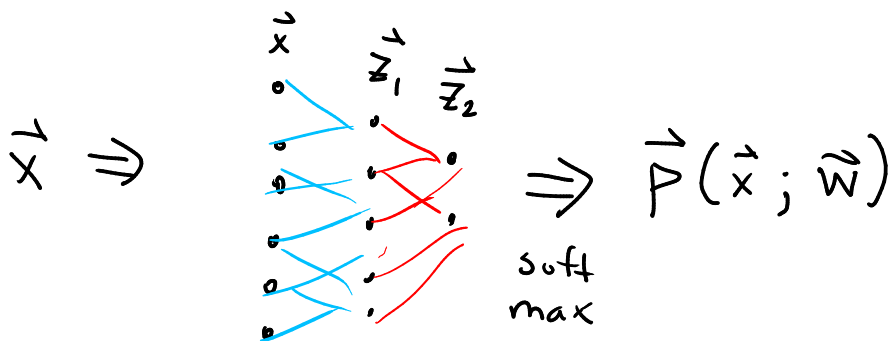
convert  $\vec{z}^{(n)}$  into prob.  $\vec{p}$

⇒ vector of model  
params  $\vec{w}$

$$p_i = \frac{e^{-z_i^{(n)}}}{e^{-z_1^{(n)}} + e^{-z_2^{(n)}} + \dots}$$

$0 \leq p_i \leq 1$

$\sum_i p_i = 1$



universal approx. theorems:

for DNNs if # layers + size of layers is large enough

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

---

Bayesian interpretation:

|                       |  |
|-----------------------|--|
| dataset $\mathcal{D}$ | "training data" for net  |
| "image"               | label:   |
| $\vec{x}^{(1)}$       | $l_1 = \begin{cases} 1 & \text{cat} \\ 2 & \text{dog} \end{cases}$ |
| $\vdots$              |  |
| $\vec{x}^{(N)}$       | $l_N$  |

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

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

first step: calculate  $\mathcal{P}(\mathcal{D} \mid \vec{w})$

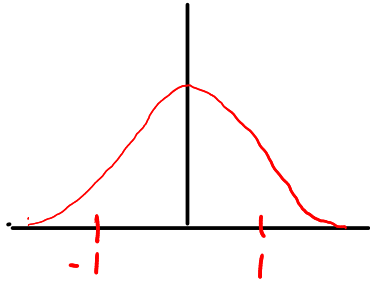
$$\mathcal{P}(\mathcal{D} \mid \vec{w}) = \prod_{i=1}^N P_{l_i}(\vec{x}^{(i)}; \vec{w}) \leq 1$$

= 1 when you have a perfect classifier

want to maximize

$$P(\vec{w} | D) = \frac{P(D | \vec{w}) P(\vec{w})}{P(D)} \quad \text{prior}$$

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



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

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

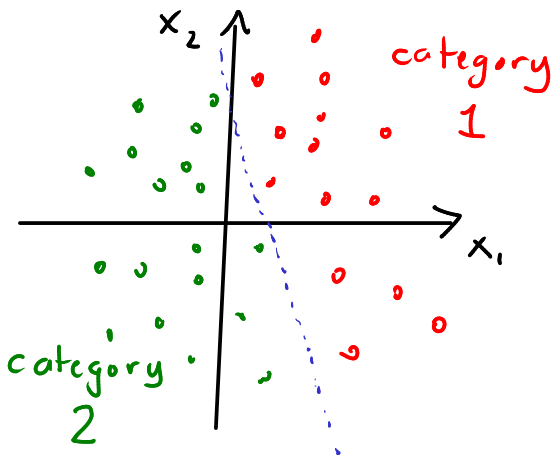
trick: to maximize  $P(\vec{w} | D)$  we take  $-\log$  of it & minimize:

loss func.

$$\begin{aligned} \mathcal{L}(\vec{w}) &= -\log P(\vec{w} | D) \\ &\downarrow \\ &= -\log P(D | \vec{w}) - \log P(\vec{w}) + \underbrace{\log P(D)}_{\text{constant (ignore)}} \end{aligned}$$

we can minimize w/ gradient descent

prob. 2 of PS #1



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

$$P_1(\vec{x}; \vec{w}) = \text{prob. that } \vec{x} \text{ is category 1}$$

$$p_2(\vec{x}; \vec{w}) = 1 - p_1(\vec{x}; \vec{w})$$