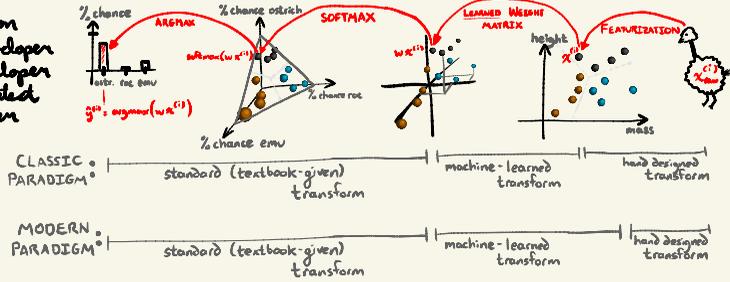


SOFTMAX @ 5 LEVELS OF SOPHISTICATION

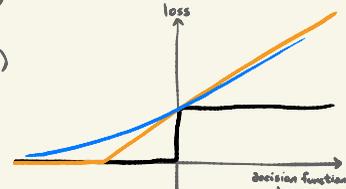
SAM'S EXPLAINERS
2022-06-07

- A. lay person
- B. MFBT developer
- C. FBR developer
- D. deep architect
- E. researcher



C. FBR developer SOFTMAX is... a computationally convenient interface w/ likelihood loss

convex, so easy to optimize
bounded gradients (despite miracle sensitivity)
tight surrogate for 0-1 loss



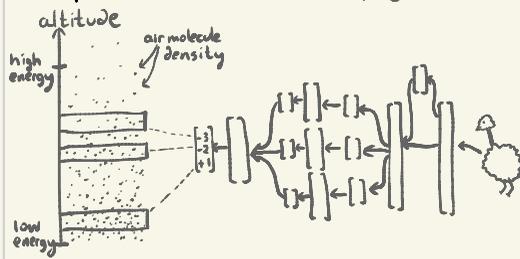
$$\text{decision func} = (-y^{(i)} w \cdot x^{(i)})$$

trespass

$$\text{hinge loss} = \max(0, 1 + \dots)$$

$$\text{softplus loss} = \log(1 + \exp(\dots)) = \log\left(\frac{\exp(\dots)}{1 + \exp(\dots)}\right) = \log \text{Softmax}\left[\begin{matrix} \dots \\ 0 \end{matrix}\right]$$

D. deep architect SOFTMAX is... a physics-inspired final layer for deep nets



(START) FOUR SCORE AND ? ?

factory

A. lay person SOFTMAX is... how we pick 'the' best answer while being fair about ties

models uncertainty: key to learning to calibrated predictions

DISCRETE \leftrightarrow CONTINUOUS

allows continuous methods to solve discrete tasks

SOFTMAX

$$\begin{bmatrix} \text{die?} \\ \text{not die?} \end{bmatrix} \leftarrow \text{ANN} \begin{bmatrix} 0.99 \\ 0.02 \end{bmatrix} \leftarrow \text{SM} \begin{bmatrix} 3.5 \\ -1.0 \end{bmatrix}$$

B. MFBT developer SOFTMAX is... how we normalize real-valued scores to probabilities

$$\text{softmax} : \mathbb{R}^k \rightarrow \{(p_i : 0 \leq p_i \leq 1) \in \mathbb{R}^k : (\forall 0 \leq k \leq \text{osp}_i) \& (\sum_i p_i = 1)\} \in \mathbb{R}^k$$

$$(z_i : 0 \leq i \leq k) \mapsto \left(\frac{\exp(z_i)}{\sum_j \exp(z_j)} : 0 \leq i \leq k \right)$$

$$z = \text{np.matmul}(w, x[i])$$

$$e = \text{np.exp}(z)$$

$$n = \text{np.sum}(e)$$

$$p = e/n$$

$$dl_dp = (-1./p) \cdot y[i]$$

$$dl_dn = -np.dot(dl_dp, e/n^2)$$

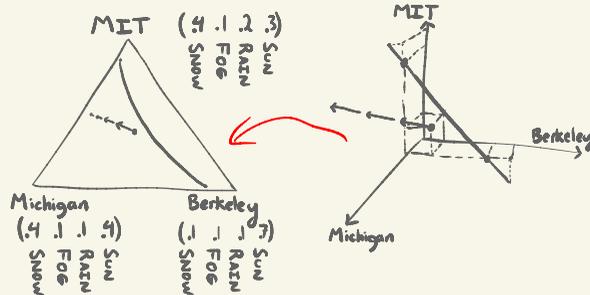
$$dl_de = dl_dp/n + dl_dn$$

$$dl_dz = dl_de \cdot e$$

$$dl_dw = \text{np.outer}(dl_dz, x[i])$$

one-hot

E. researcher SOFTMAX is... the canonical e-flat parameterization of the simplex



MIT Mich Berk
(.33 .33 .33)
observe SNOW
= (.4 .4 .1)
(.13 .13 .03)
(.44 .44 .11)
observe SNOW
= (.4 .4 .1)
(.17 .17 .01)
(.48 .48 .03)
observe SNOW
= (.4 .4 .1)
(.19 .19 .03)
(.46 .46 .08)