

of the input which is in class 6. linear model: it now has Dinput dimensions (has not changed)

Kontput dimensions (has changed)

Rather than have a scalar output for each input,

index of a we have one ucight bector we among K classes

he now have a vector output for each input.

and one bias scalar be for each class of the tasky total # of desses K dimensions

total

of (lodel output is the vector === input vector of dimension D 2 = W = + b vector of dinertion K dimension K of dimensions

of input dimensions (features)

logistic function

T(2)= 1 (+e-t)

softnak output, which

is a vector, is the

Linear model

for multiclass

For RELLE :

setting.

Achiration function Recall in the binary classification setting, we used Here in the nulhillass classification setting, we use the softmax function: (multivariable generalitation of The logistic function)

yh = softmax (2,,..., 2,)k = retri for one dass ke 1.. k index of # 5/ a class clusses $\frac{1}{y} = \text{Softmax}(\frac{1}{z}) =$ Output of the linear model w/ a softmax adjudion Function

K dimensions The components of \overline{z}' , values $\overline{z}_1, \dots, \overline{z}_k$ are called <u>logits</u> (inputs to the) Properties of the softmax: · outputs are all possitive nie because they . outputs all sum to 1 car be interpreted . If one of the bogits the is much larger than the other bogits, as probabilités (each component & of the

becker of

probability of the input Softmax (2)= Softmax (2, ..., 2k) beig in class le) ~ (0, 0, ..., 0, 1, 0, ... 0) In this case the softmax behaves like an argmax. the target (one-hot vector) Cross-entropy for multiclass classification - k" component
of the model $\vec{y} = Softmax(\vec{z})$ $\mathcal{L}_{CE}(\vec{y}, \vec{t})$ = - $\vec{\Sigma}$ t_{R} t_{R} t_{R} t_{R} t_{R} t_{R} t_{R} t_{R} vectors $= -\vec{t} (\log \vec{y})$

we apply the

log elementuise to the (cechon) model output. Here again, we can combine the cross-entropy loss function and the softmax activation function inho the Softmax-cross-entropy function

Gradient descent updale rule: Requires computing now a vector