

# Deep Network 2

softmax

$$z \rightarrow \left( \frac{e^{z_1}}{\sum_j e^{z_j}}, \dots, \frac{e^{z_k}}{\sum_j e^{z_j}} \right)$$

converts input into probability error.

$$P_r[Y=y | X=x]$$

If some coordinate  $j$  of  $z$  dominates others, then softmax close to  $e^j$

Cross entropy loss  $p, q \in \Delta_d := \{v \in \mathbb{R}^k : v \geq 0, \sum_i v_i = 1\}$

$$H(q, p) := - \sum_{i=1}^k p_i \ln q_i$$

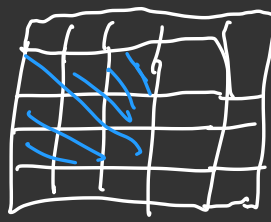
In pytorch:  $\text{Lce}(\hat{y}, y) = H(\text{softmax}(\hat{y}), e_y)$

$$= - \sum_{i=1}^k (e_y)_i \ln(\text{softmax}(\hat{y})_i) = - \ln \left( \frac{\sum_i \exp(\hat{y}_i)}{\exp(\hat{y}_y)} \right) = -\hat{y}_y + \ln(\sum_i \exp(\hat{y}_i))$$

$$\begin{bmatrix} 0 & 0 & \dots & 1 & \dots & 0 & 0 \end{bmatrix}$$

$\uparrow$   
y-th

## Max pooling



Each filter takes the maximum.

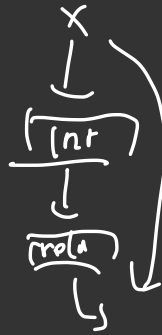
## Batch Normalization

Standardize output

$$x \rightarrow \frac{x - E(x)}{\text{std dev}(x)} \cdot \gamma + \beta$$

## Skip connections

$$z \rightarrow z + f_i(z)$$



# Gradient Descent and backpropagation

$$\nabla (Cy x^T w) = \nabla (C p(n)) = \frac{\partial L(p(n))}{\partial p(n)} \cdot \frac{\partial p(n)}{\partial w}$$

$$p_1 := \sigma_1(w_1 x)$$

$$p_{i+1} := \sigma_{i+1}(w_{i+1} p_i)$$

$$p_L := C y \sigma_L(w_L p_{L-1})$$

$$\nabla_w (Cy f(x; w)) = \frac{\partial p_L}{\partial p_i} \cdot \frac{\partial p_i}{\partial w_i}, \text{ where } \frac{\partial p_L}{\partial p_i} = \frac{\partial p_L}{\partial p_{i+1}} \cdot \frac{\partial p_{i+1}}{\partial p_i}$$

**Forward pass:** compute  $p_L$  by  $p_1$  and more.

**Backward pass:** compute  $w_L$  down to  $w_1$  via inductive chain rule

**vanishing / exploding gradients**

$$\nabla w_i (Cy f(x; w)) = \frac{\partial p_L}{\partial p_i} \cdot \frac{\partial p_i}{\partial w_i} = \prod_{j=i}^{L-1} \frac{\partial p_{j+1}}{\partial p_j}$$

With many-layered networks, this computation easily explodes and vanishes

## Initialization

$W_i \in \mathbb{R}^{d_i \times d_{i-1}} \sim \mathcal{N}(0, 1/d_{i-1})$  per coordinate.

Learning rates  $\eta \in \{0.001, 0.01, 0.1\}$

## Data Augmentation

Besides original data, generate non data,  
with some computation