

## Lecture 2

Minimize error rate

$$\frac{1}{l} \sum_{i=1}^l [a(x_i) \neq y_i] \rightarrow \min_a$$

$y_i \in \{1, \dots, c\}$  - classes

$$p = \begin{pmatrix} p(\hat{y}_i = 1) \\ p(\hat{y}_i = 2) \\ \vdots \\ p(\hat{y}_i = c) \end{pmatrix}$$

$$z_i = \sigma(W_i z_{i-1} + b_i)$$

$$z_n = W_n z_{n-1} + b_n$$

$$z_n \in \underline{\mathbb{R}^c}$$

$$p = \{p_i\}_{i=1}^c$$

$$p_i \geq 0,$$

$$\sum_{i=1}^c p_i = 1$$

Softmax

$$z = (z_1, \dots, z_c)$$

$$p_i = \frac{\exp(z_i)}{\sum_{j=1}^c \exp(z_j)} \geq 0$$

$p$  - probabilities  $\forall i, j \quad z_i > z_j \Rightarrow p_i > p_j$

$z$  - logits

$$\hat{y} \rightarrow 1, \dots, n$$

$$y \rightarrow 1, \dots, n$$

$$\hat{y} \rightarrow \begin{pmatrix} p_1 \\ \vdots \\ p_y \\ \vdots \\ p_c \end{pmatrix} \rightarrow 1$$

$$y \rightarrow \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}_y$$

$$y\text{-target} \in \{1, \dots, c\}$$

$$-\sum_{k=1}^c [k=y] \log p_k \rightarrow \min_{\{W, b\}}$$

||

$$-\log p_y \rightarrow \min \sim p_y \rightarrow \max$$

NLL Negative log-likelihood

$$\{(x_i, y_i)\}_{i=1}^l$$

$$-\sum_{i=1}^l \sum_{k=1}^c [k=y_i] \log p_k^{(i)} \rightarrow \min_{\{W, b\}}$$

$$-\sum_{i=1}^l \log p_{y_i}^{(i)} \rightarrow \min \sim \prod_{i=1}^l p_{y_i}^{(i)} \rightarrow \max$$

$$NLL \text{ loss}(\text{softmax}(z), y) \rightarrow \min$$

$$\log p_i = \log \text{softmax}(z_i) = \log \frac{\exp(z_i)}{\sum_{k=1}^c \exp z_k} =$$

$$= \log \exp(z_i) - \underbrace{\log \left( \sum_{k=1}^c \exp z_k \right)}_{\log \text{sumexp}(z)}$$

$\parallel$   
 $z_i$

$$\log \frac{\exp(z_i - \max_j z_j)}{\sum_{k=1}^c \exp(z_k - \max_j z_j)}$$

Cross entropy loss  $(z, y)$

$\uparrow$   
logits  $\in \mathbb{R}^c$

$$CE(\{p_i\}_{i=1}^c, \{q_i\}_{i=1}^c) = - \sum_{i=1}^c q_i \log p_i$$

$$\hat{y} = \begin{pmatrix} p_1 \\ \vdots \\ p_c \end{pmatrix} = p$$

$$y = \begin{pmatrix} 0 \\ \vdots \\ 1 \\ \vdots \\ 0 \end{pmatrix}_y = q$$

$$NLL = CE = - \sum_{i=1}^c q_i \log p_i \rightarrow \min$$

Dropout

$$z_k \in \mathbb{R}^{d_k} \quad d_k - \text{hidden size}$$

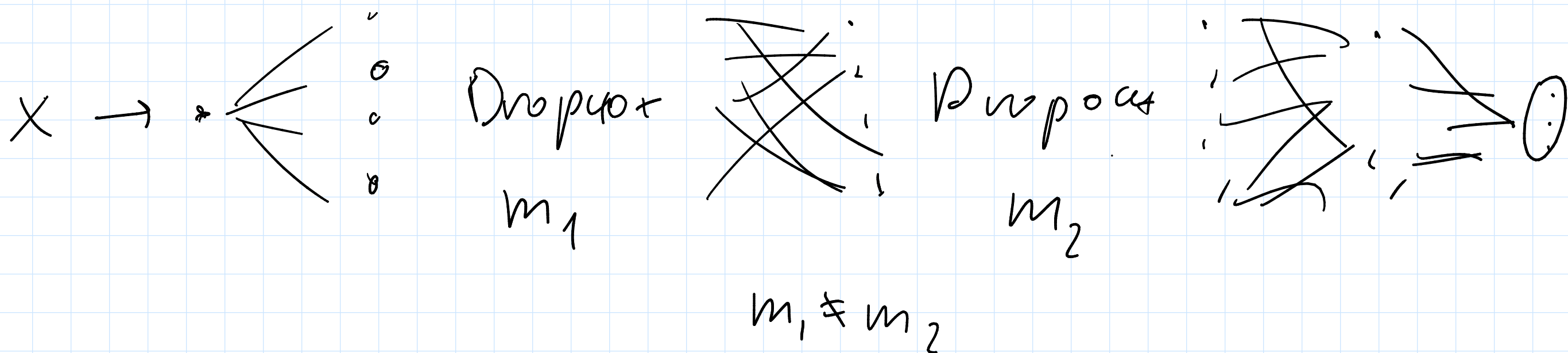
$$z_k = \sigma(W_k z_{k-1} + b_k)$$

$$m \in \{0, 1\}^{d_k} \quad m_i \sim \text{Bernoulli}(1-p)$$

$$p \in [0, 1; 0, 25]$$

$$y = m \odot z_k$$

↑  
element wise product



$x_{\text{next}}$

$m_{1, \text{next}}$

$m_{2, \text{next}}$

$$\mathbb{E}[y_k] = \mathbb{E}[m \odot z_k] = \mathbb{E}[m] \odot z_k =$$

$$y_k \in \mathbb{R}^{d_k}$$

$$= (1-p) \cdot z_k$$

$$\langle w_{k+1}, z_k \rangle \leftarrow 1-p$$

Train mode:

$$m_i \sim \text{Bernoulli}(1-p)$$

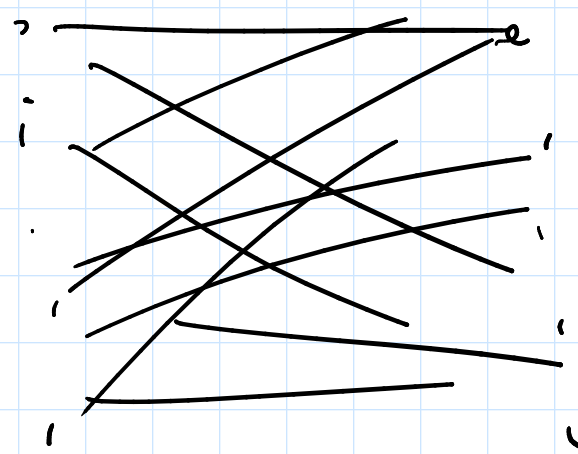
$$y_k = m \odot z_k \frac{1}{1-p}$$

$$\mathbb{E}[y_k] = z_k$$

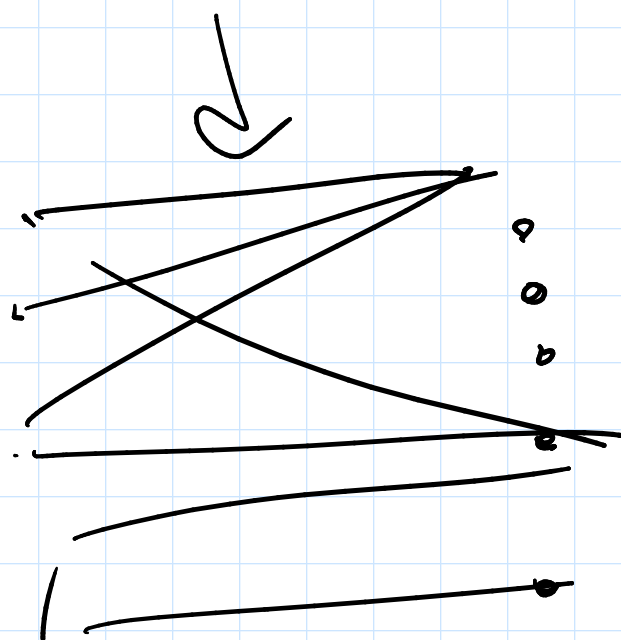
Eval mode:

$$y = z$$

$$m = \begin{pmatrix} 1 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{pmatrix} \begin{matrix} \leftarrow m_1 \\ \\ \\ \leftarrow m_i \\ \\ \\ \leftarrow m_{d_k} \end{matrix}$$



$$m = \begin{pmatrix} 1 \\ \vdots \\ 1 \\ \vdots \end{pmatrix}$$



Batch Normalization

$$(-\infty, +\infty)$$

Train mode

$$Z \in \mathbb{R}^{B \times d_k}$$

B - batch size

$$z_i \in \mathbb{R}^{d_k}$$

$$\{(x_i, y_i)\}$$

$$W \underset{\mathbb{R}^d}{z} + b$$

$$x_1 \ y_1 \rightarrow$$

$$x_2 \ y_2 \rightarrow$$

$$(x_1, x_2, \dots, x_B)$$

$$W \overset{\parallel}{\underset{\leftarrow}{\overset{\rightarrow}{z}}} + (b, \dots, b)$$

$$\mu = \frac{1}{B} \sum_{i=1}^B z_i$$

$$\sigma^2 = \frac{1}{B} \sum_{i=1}^B (z_i - \mu)^2$$

$$\mu, \sigma \in \mathbb{R}^{d_k}$$

$$\hat{z}_i = \frac{z_i - \mu}{\sqrt{\sigma^2 + \epsilon}}$$

element-wise

$\leftarrow$  small

$$\frac{1}{B} \sum \hat{z}_i = 0$$

$$\frac{1}{B} \sum \hat{z}_i^2 \approx 1$$

$$y_i = \underset{\alpha}{\hat{z}_i} \odot \underset{\beta}{w} + b$$

$$w, b \in \mathbb{R}^{d_k}$$

learnable

$$\text{running\_mean} := \text{running\_mean}_{\text{prev}} (1-m) + \mu \cdot m$$

$$\text{running\_var} := \text{running\_var}_{\text{prev}} (1-m) + \sigma^2 m \frac{B}{B-1} \quad m \approx 0.1$$

Eval mode:

$$\hat{z}_i = \frac{z_i - \text{running\_mean}}{\sqrt{\text{running\_var} + \epsilon}}$$

$$\epsilon = \text{const}$$

$$y_i = \hat{z}_i \odot W + b$$

---