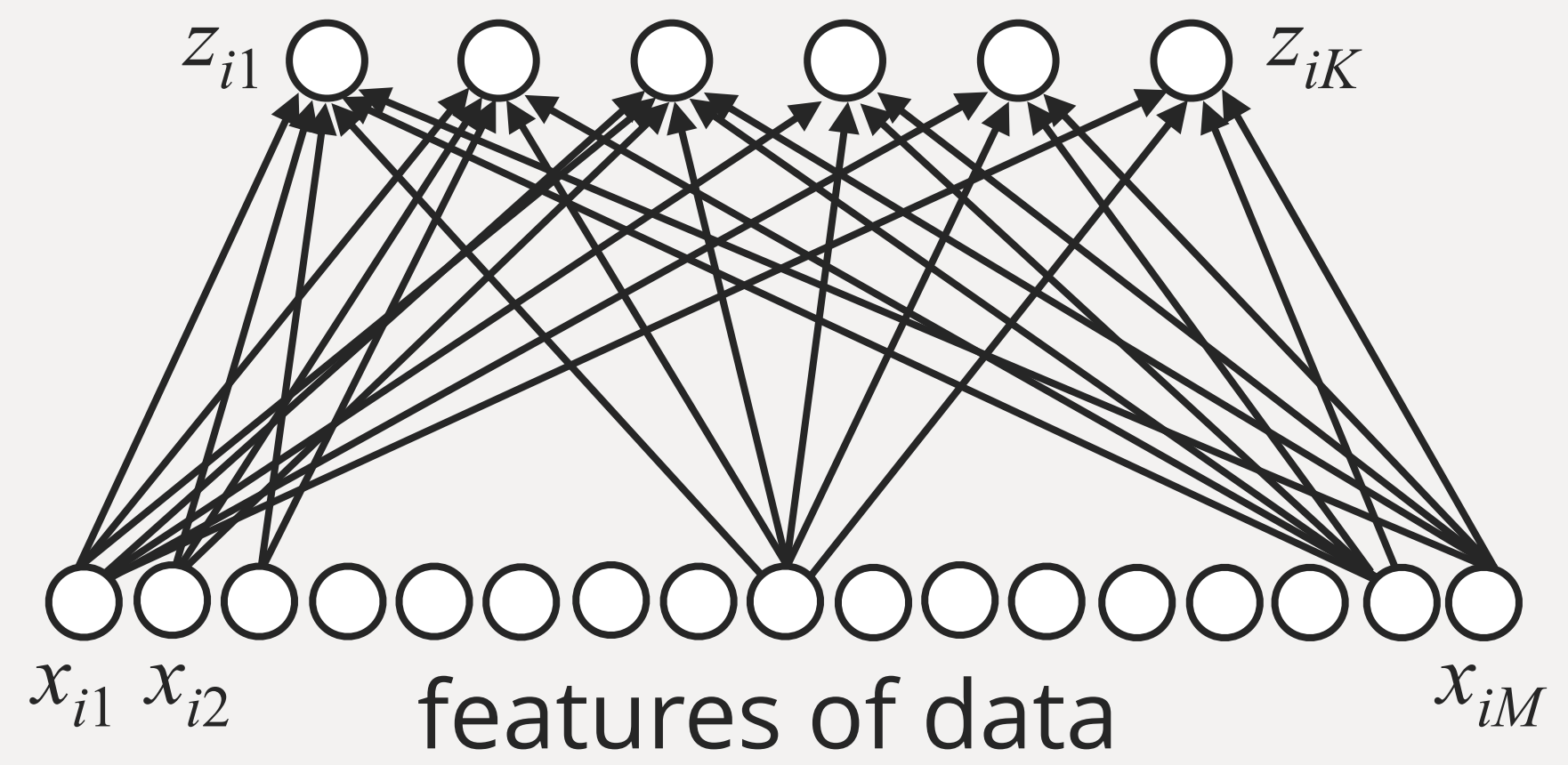




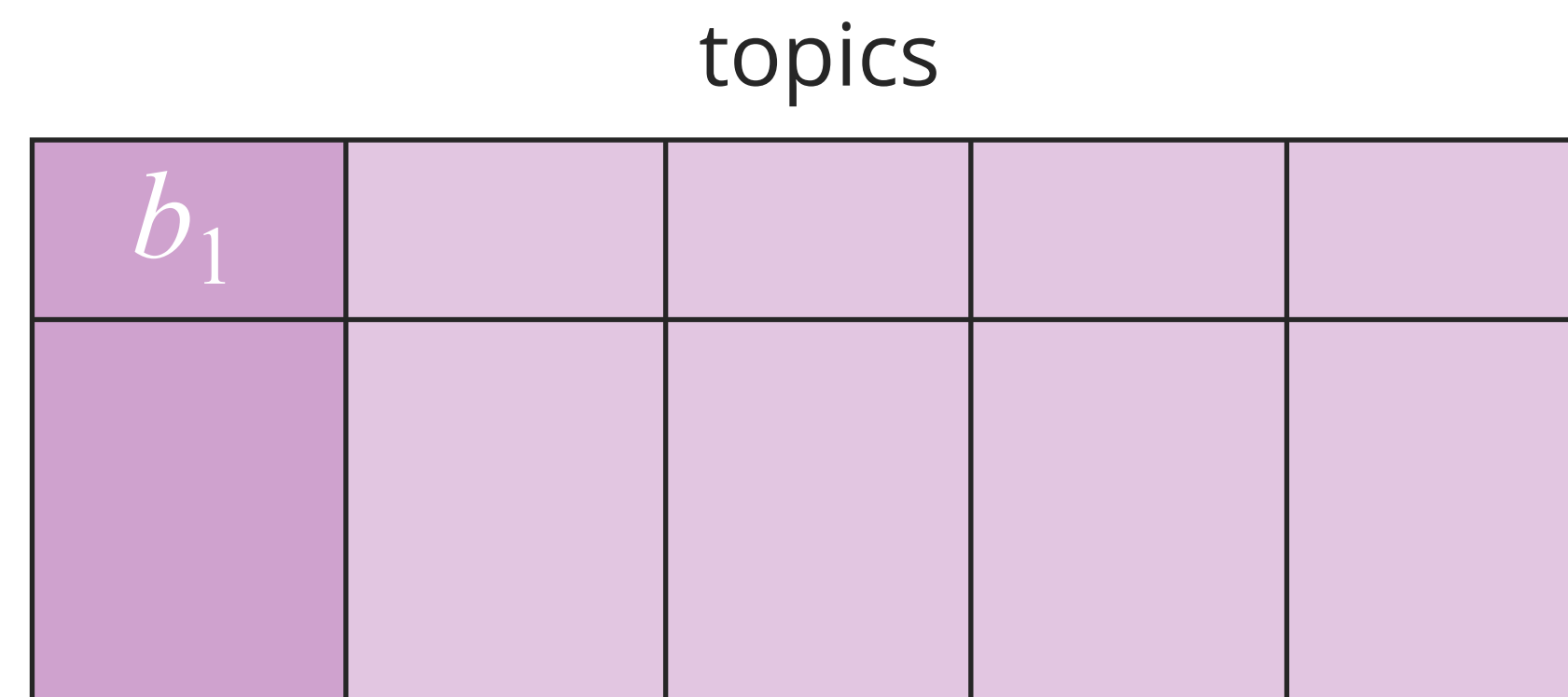
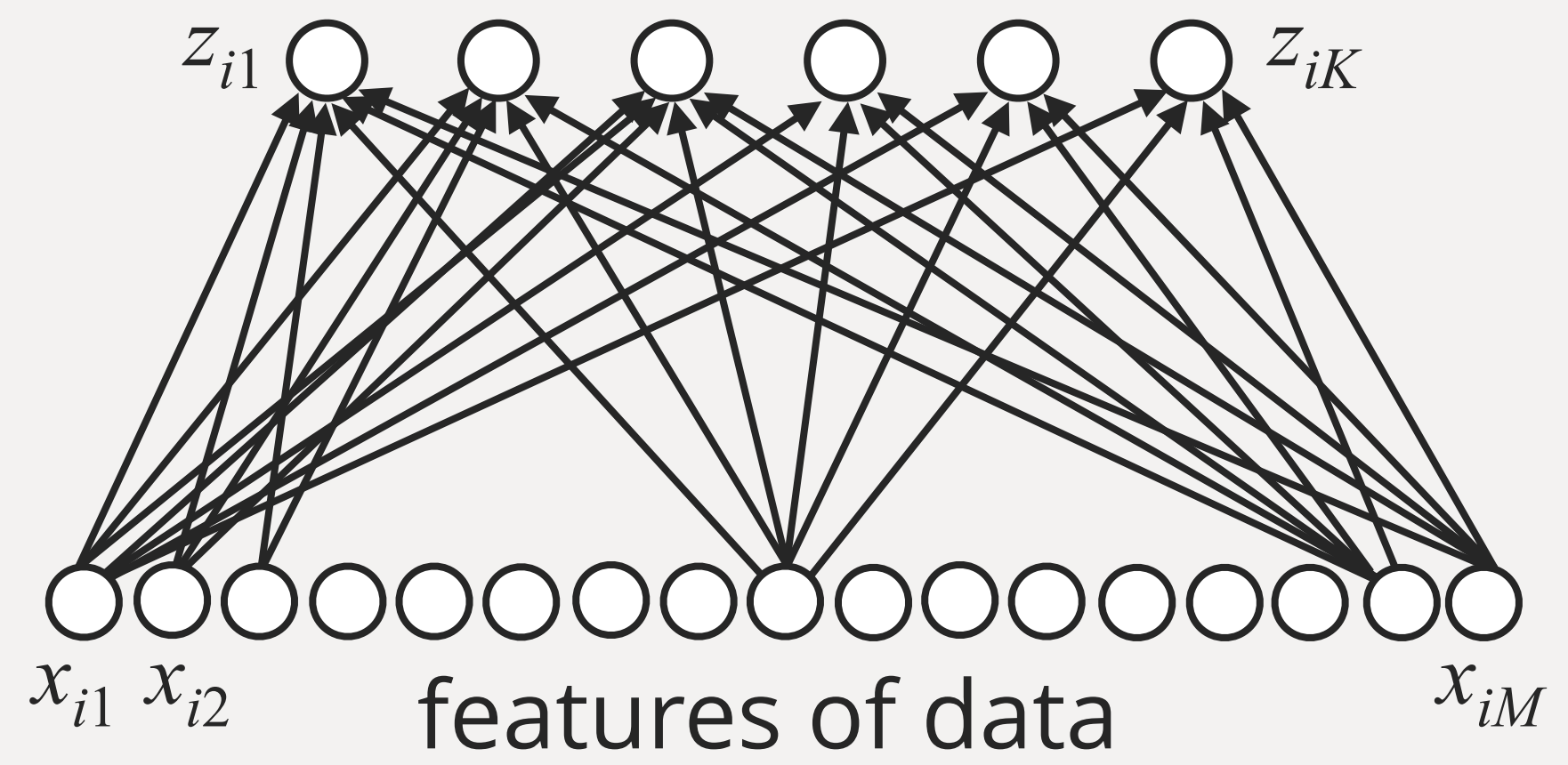
# Interpretation of Multilayer Perceptron



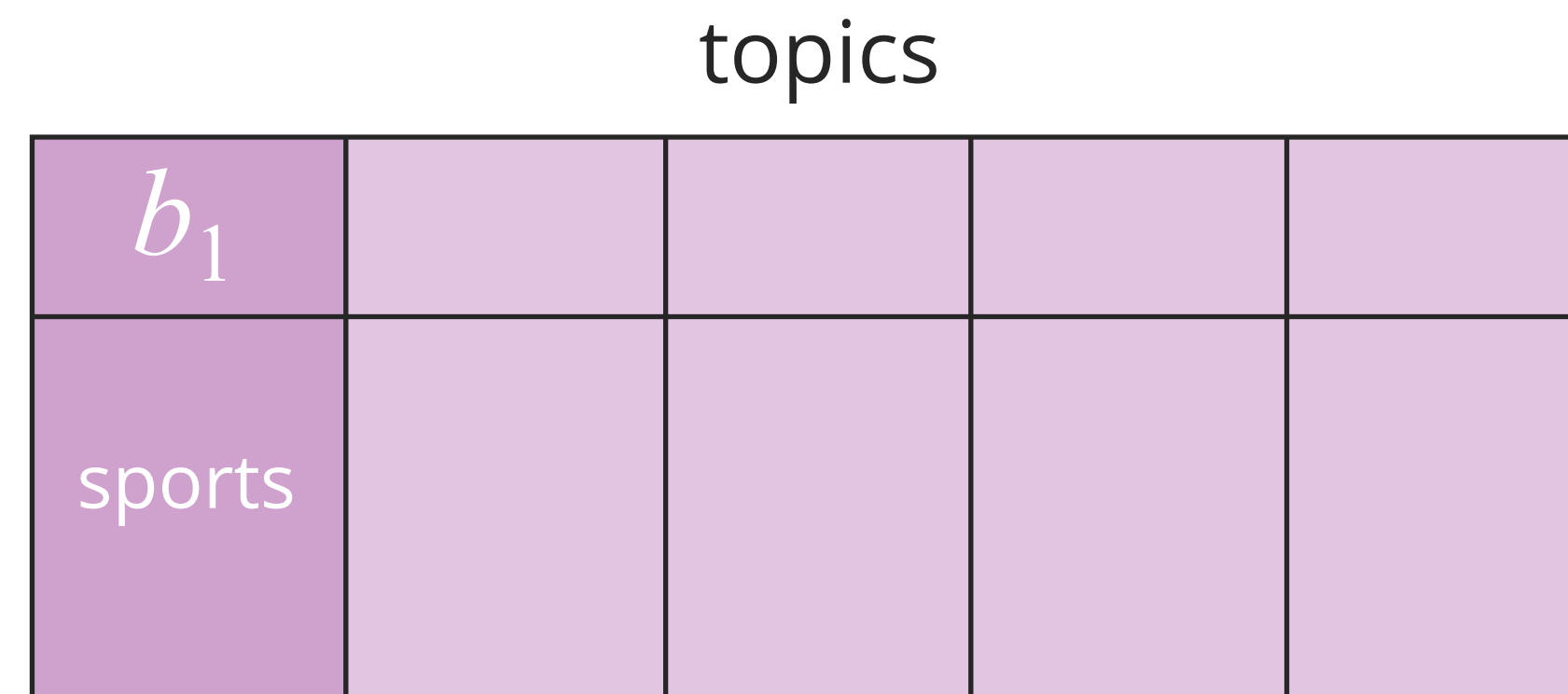
topics

$b_1$	$b_2$	$\cdot$	$\cdot$	$b_K$

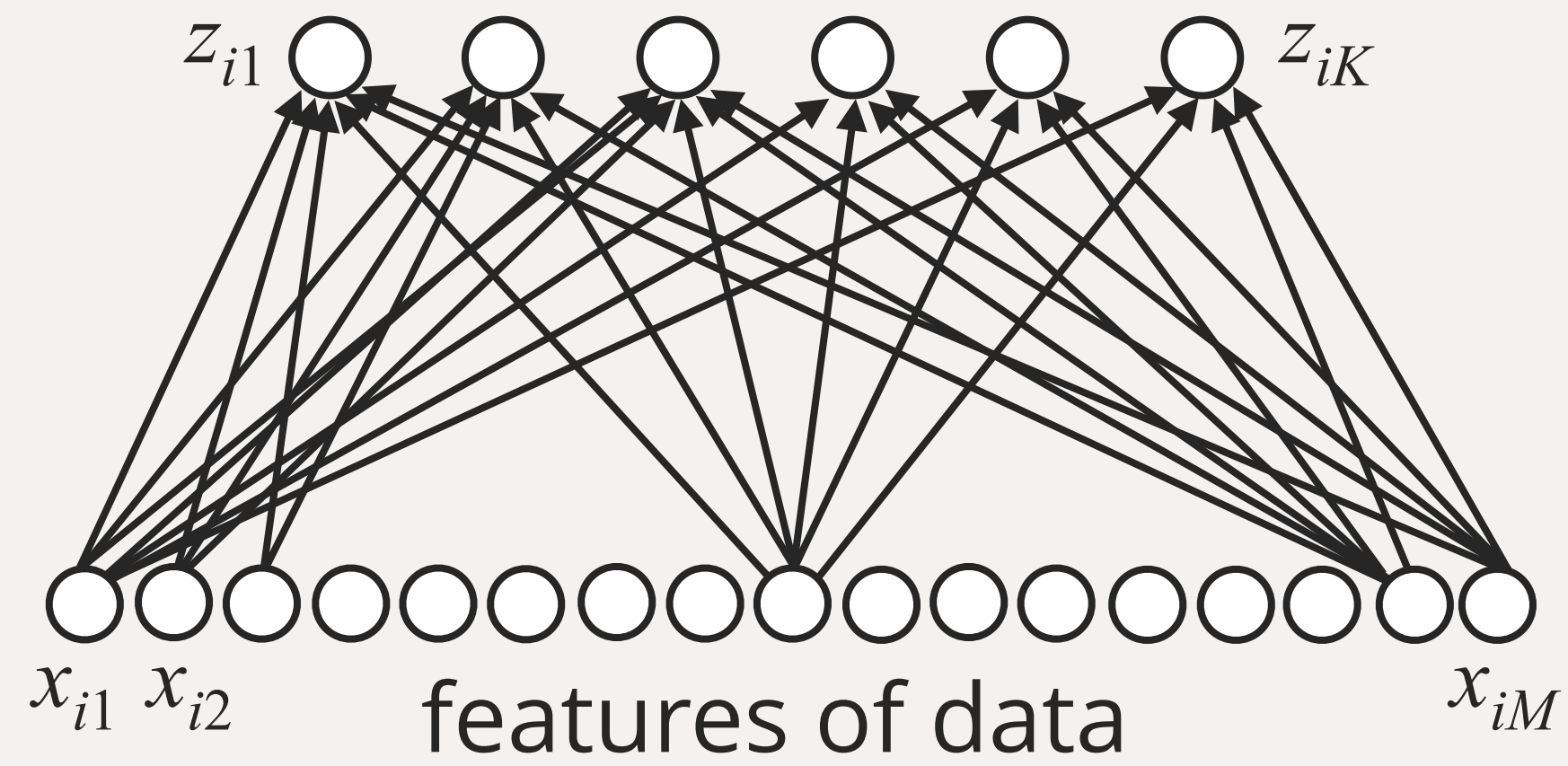
$$\begin{aligned}
 z_{i1} &= b_{01} + x_i \odot b_1 \\
 z_{i2} &= b_{02} + x_i \odot b_2 \\
 &\vdots \\
 z_{iK} &= b_{0K} + x_i \odot b_K
 \end{aligned}$$



$$\begin{aligned}
 z_{i1} &= b_{01} + x_i \odot b_1 \\
 z_{i2} &= b_{02} + x_i \odot b_2 \\
 &\vdots \\
 z_{iK} &= b_{0K} + x_i \odot b_K
 \end{aligned}$$



$$\begin{aligned}
 z_{i1} &= b_{01} + x_i \odot b_1 \\
 z_{i2} &= b_{02} + x_i \odot b_2 \\
 &\vdots \\
 z_{iK} &= b_{0K} + x_i \odot b_K
 \end{aligned}$$

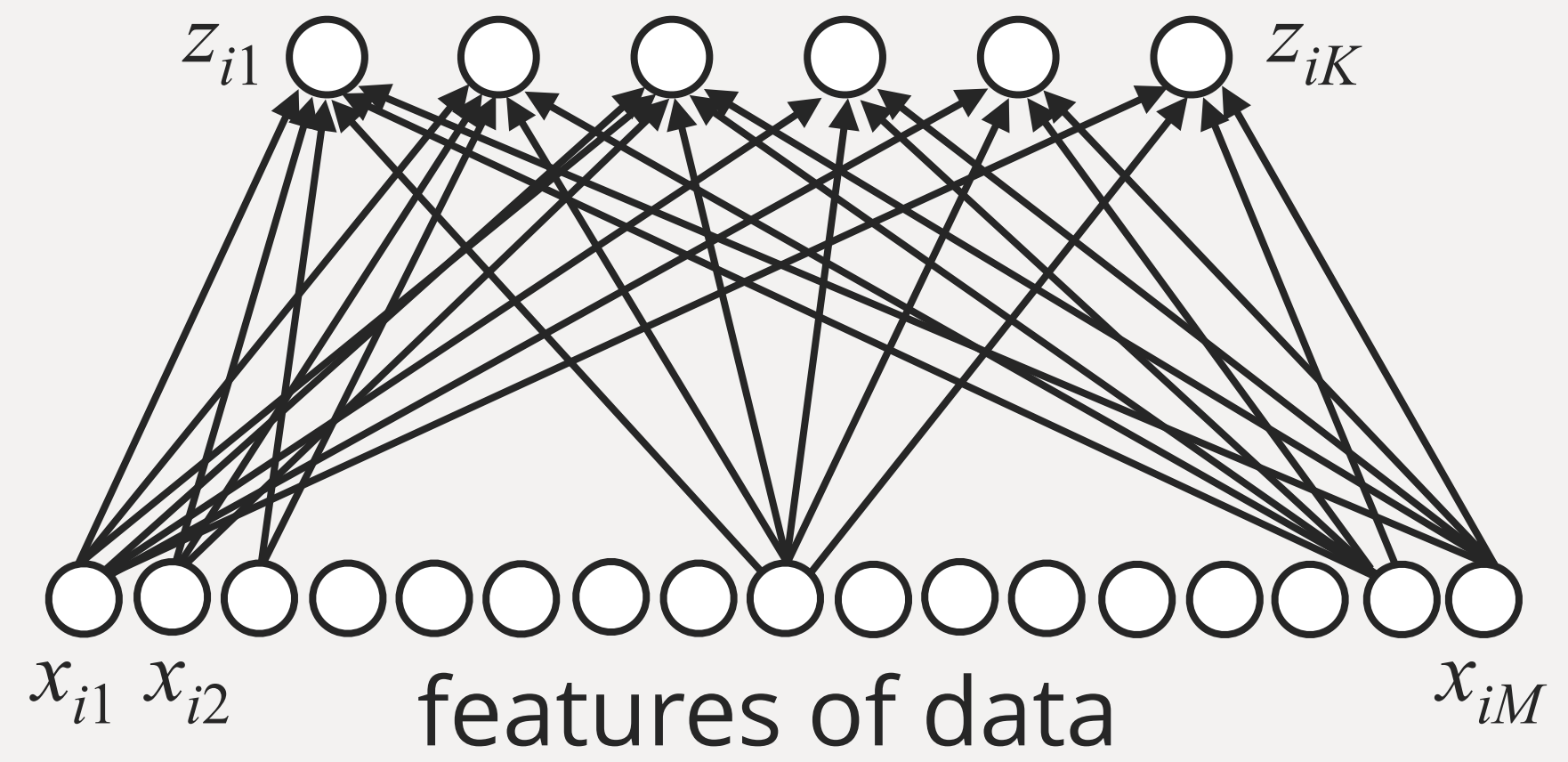


topics

$b_1$	$b_2$			
sports	history			

$$\begin{aligned}
 z_{i1} &= b_{01} + x_i \odot b_1 \\
 z_{i2} &= b_{02} + x_i \odot b_2 \\
 &\quad \vdots \\
 z_{iK} &= b_{0K} + x_i \odot b_K
 \end{aligned}$$

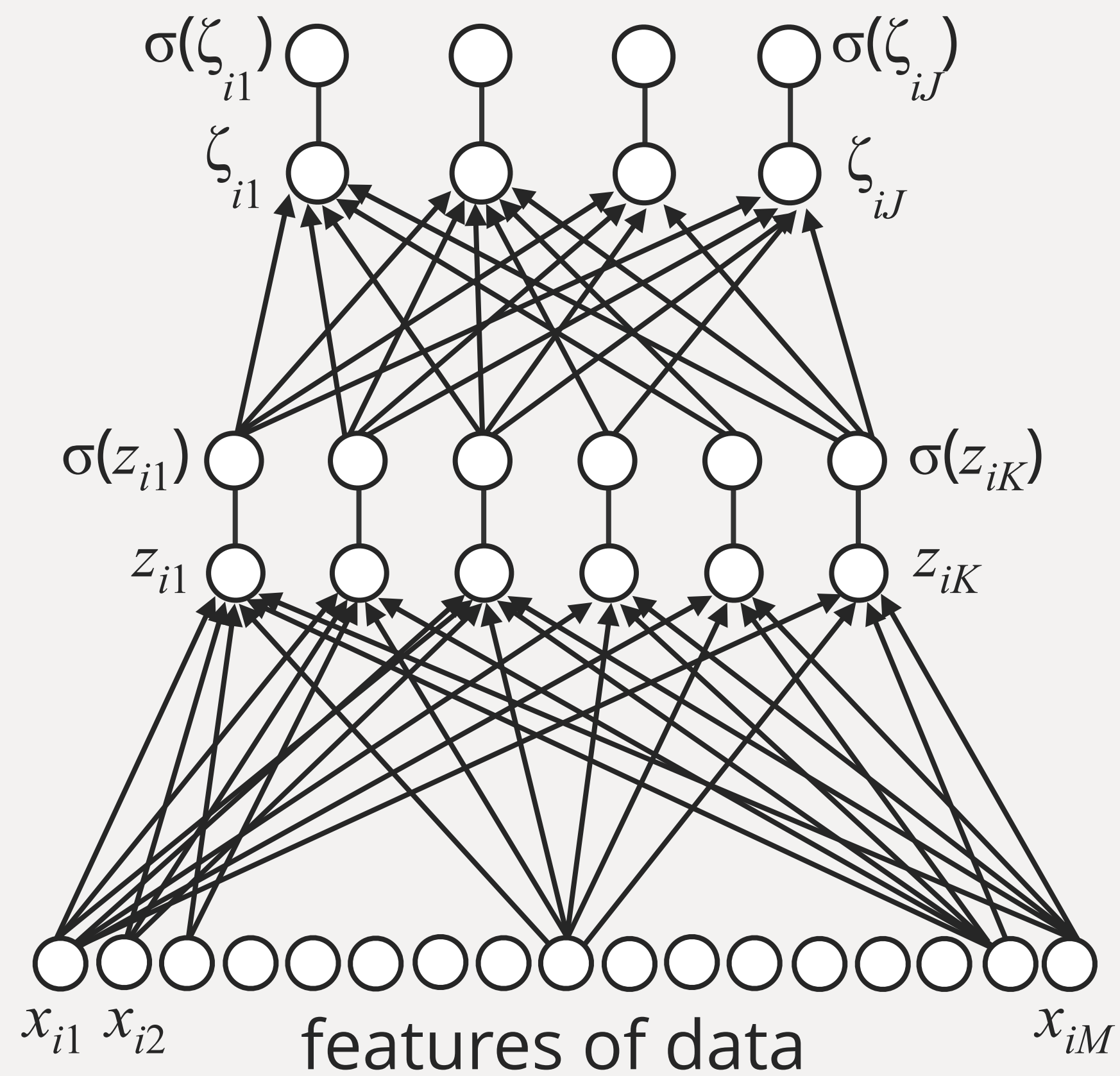




topics

$b_1$	$b_2$	$\cdot$	$\cdot$	$b_K$
sports	history	$\cdot$	$\cdot$	politics

$$\begin{aligned}
 z_{i1} &= b_{01} + x_i \odot b_1 \\
 z_{i2} &= b_{02} + x_i \odot b_2 \\
 &\vdots \\
 z_{iK} &= b_{0K} + x_i \odot b_K
 \end{aligned}$$



meta-topics

$c_1$	$c_2$	•	•	$c_J$

topics

$b_1$	$b_2$	•	•	$b_K$
sports	history	•	•	politics

$$\zeta_{i1} = c_{01} + \sigma(z_i) \odot c_1$$

$$\zeta_{i2} = c_{02} + \sigma(z_i) \odot c_2$$

•

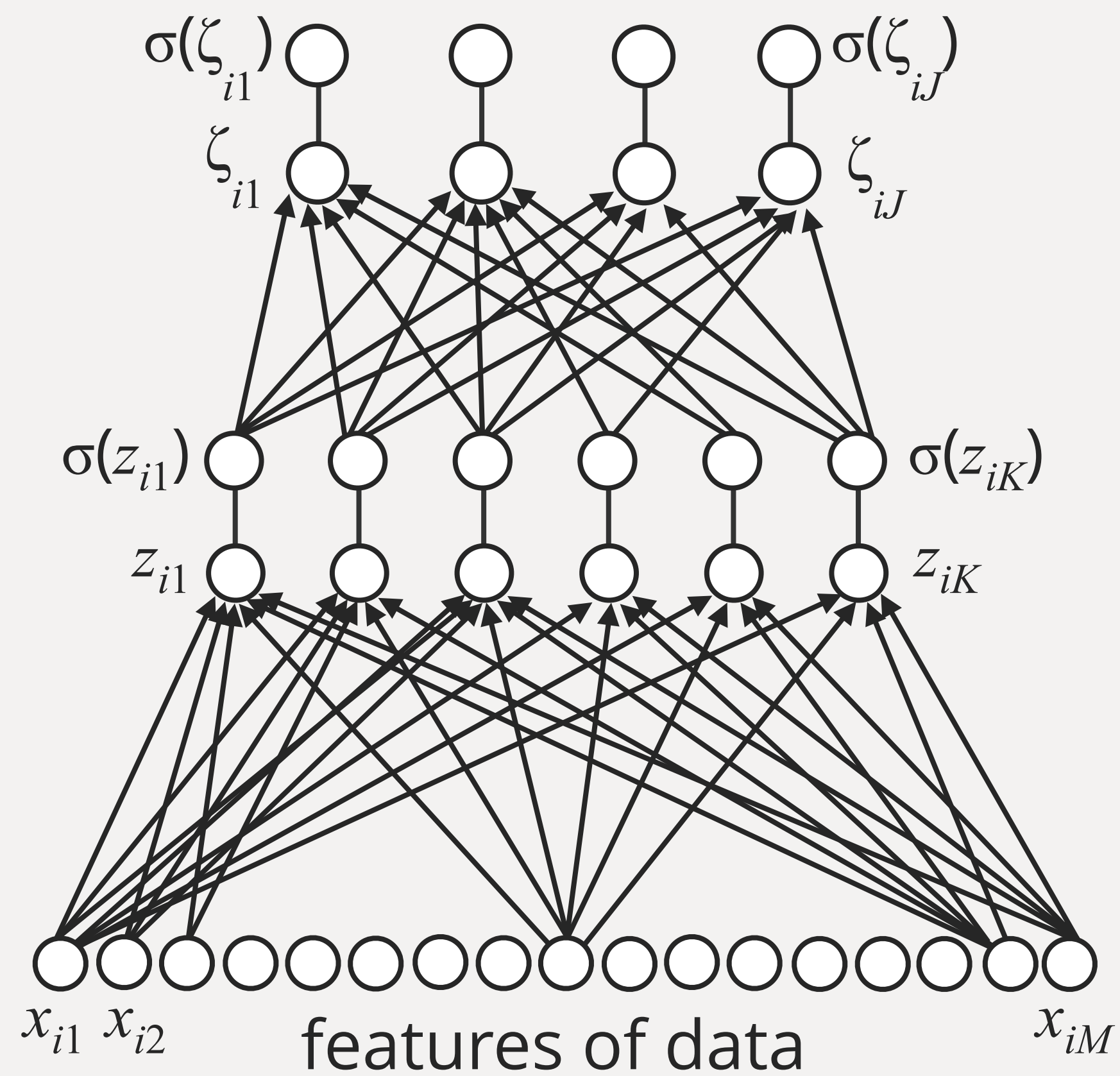
$$\zeta_{iJ} = c_{0J} + \sigma(z_i) \odot c_J$$

$$z_{i1} = b_{01} + x_i \odot b_1$$

$$z_{i2} = b_{02} + x_i \odot b_2$$

•

$$z_{iK} = b_{0K} + x_i \odot b_K$$



meta-topics

$c_1$				
sports + history				

topics

$b_1$	$b_2$	•	•	$b_K$
sports	history	•	•	politics

$$\zeta_{i1} = c_{01} + \sigma(z_i) \odot c_1$$

$$\zeta_{i2} = c_{02} + \sigma(z_i) \odot c_2$$

•

$$\zeta_{iJ} = c_{0J} + \sigma(z_i) \odot c_J$$

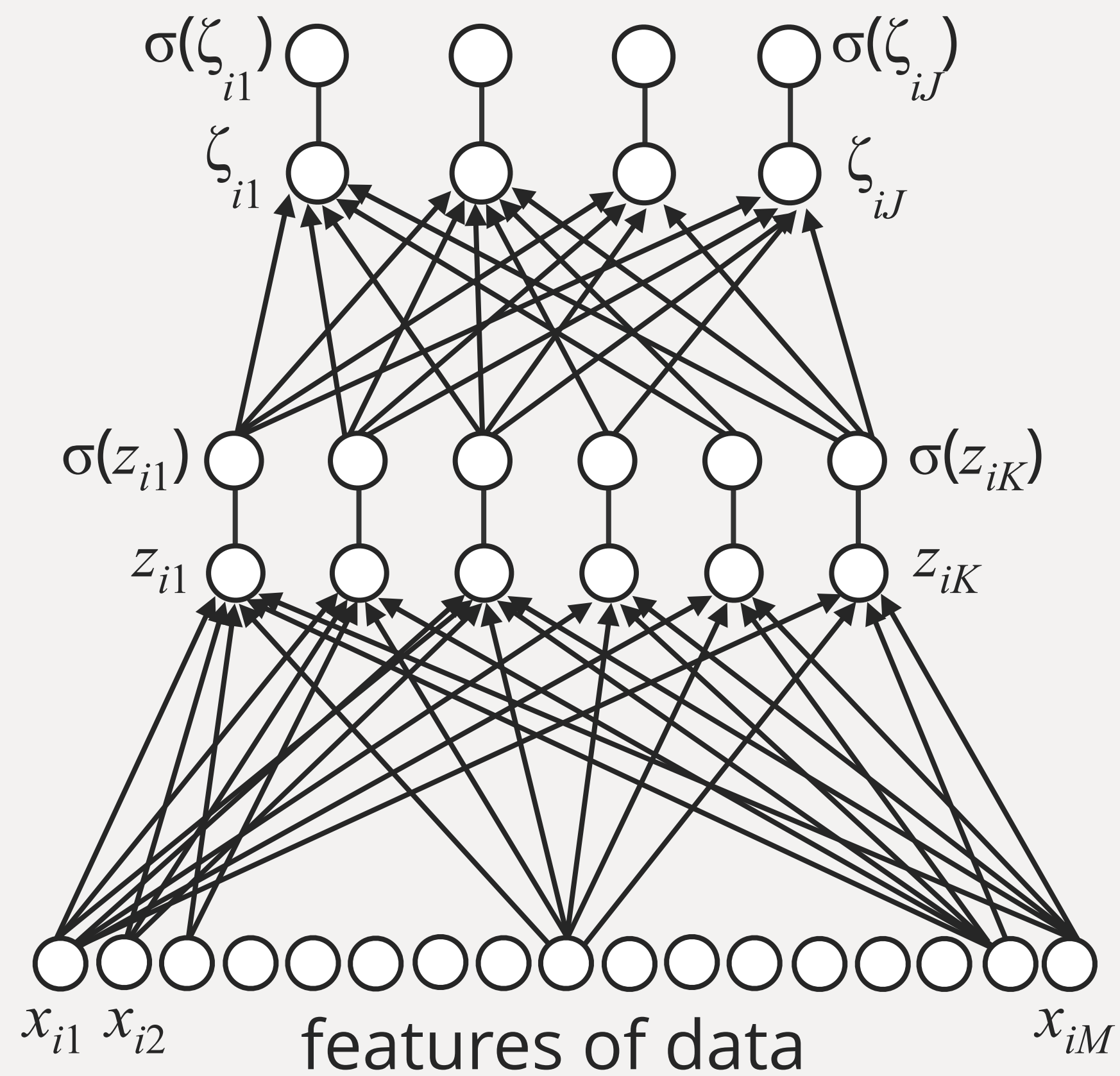
$$z_{i1} = b_{01} + x_i \odot b_1$$

$$z_{i2} = b_{02} + x_i \odot b_2$$

•

$$z_{iK} = b_{0K} + x_i \odot b_K$$





meta-topics

$c_1$	$c_2$			
sports + history	politics + sports			

topics

$b_1$	$b_2$	•	•	$b_K$
sports	history	•	•	politics

$$\zeta_{i1} = c_{01} + \sigma(z_i) \odot c_1$$

$$\zeta_{i2} = c_{02} + \sigma(z_i) \odot c_2$$

⋮

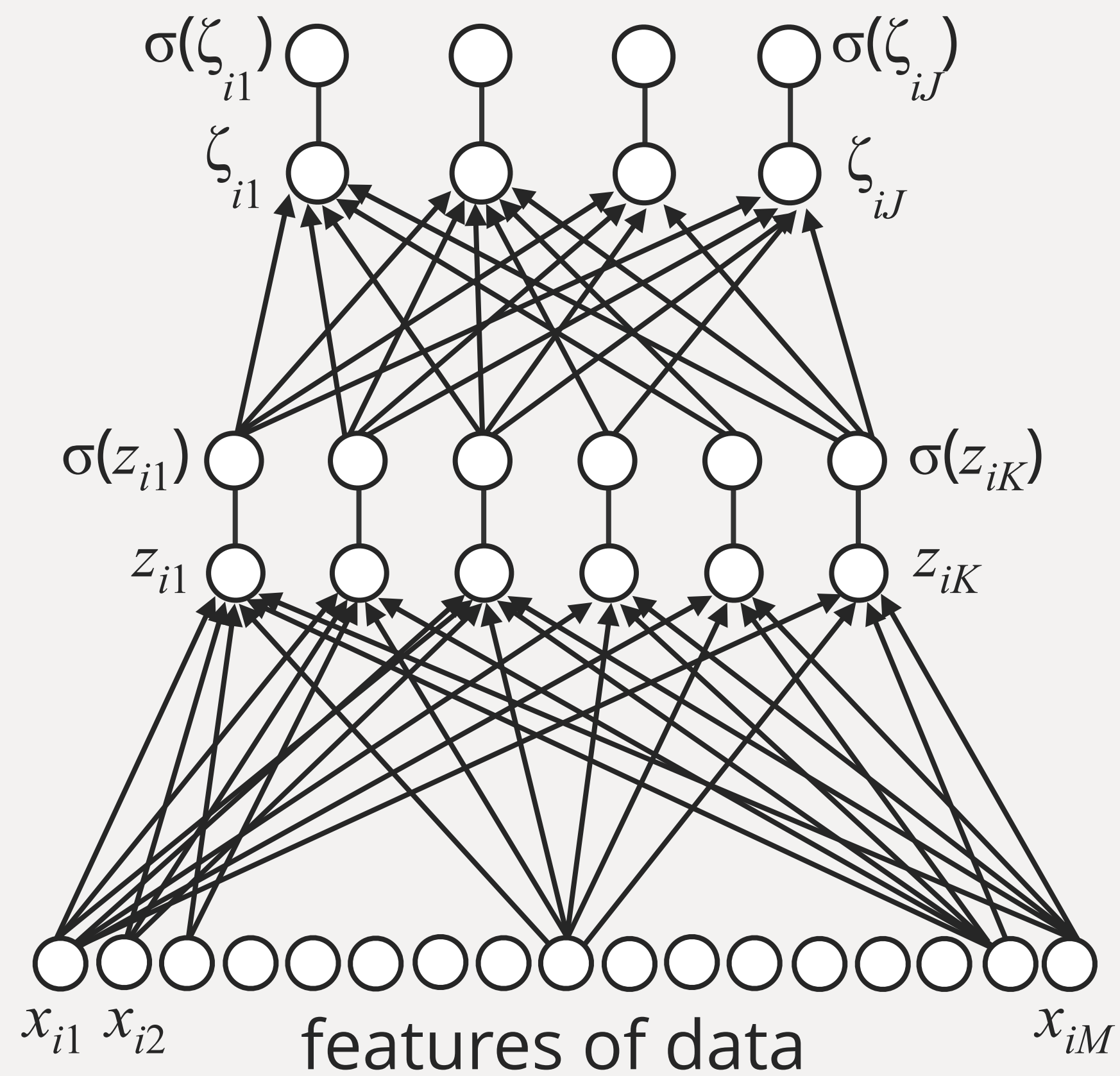
$$\zeta_{iJ} = c_{0J} + \sigma(z_i) \odot c_J$$

$$z_{i1} = b_{01} + x_i \odot b_1$$

$$z_{i2} = b_{02} + x_i \odot b_2$$

⋮

$$z_{iK} = b_{0K} + x_i \odot b_K$$



meta-topics

$c_1$	$c_2$	$\cdot$	$\cdot$	$c_J$
sports + history	politics + sports	$\cdot$	$\cdot$	politics + sports + history

topics

$b_1$	$b_2$	$\cdot$	$\cdot$	$b_K$
sports	history	$\cdot$	$\cdot$	politics

$$\zeta_{i1} = c_{01} + \sigma(z_i) \odot c_1$$

$$\zeta_{i2} = c_{02} + \sigma(z_i) \odot c_2$$

$\vdots$

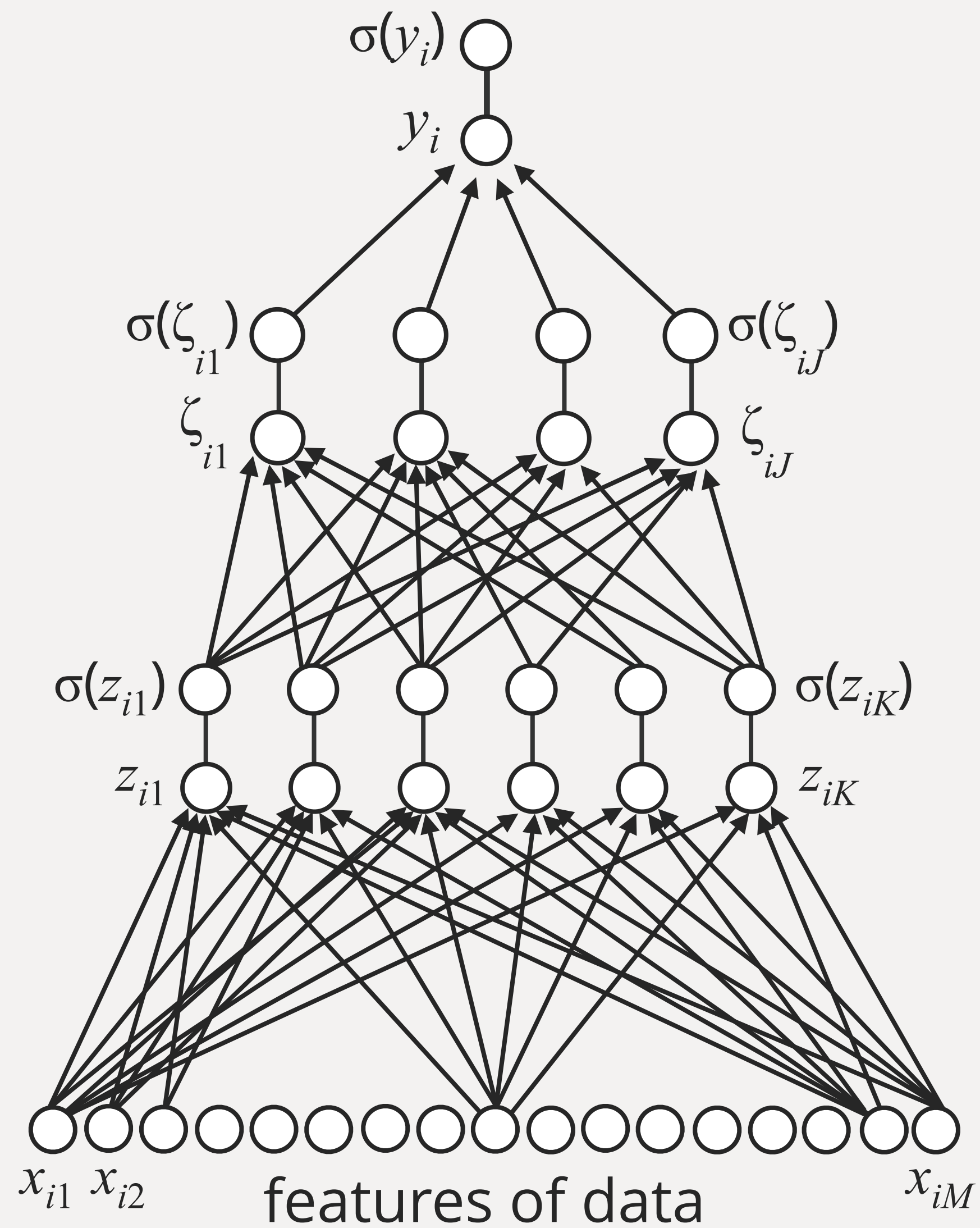
$$\zeta_{iJ} = c_{0J} + \sigma(z_i) \odot c_J$$

$$z_{i1} = b_{01} + x_i \odot b_1$$

$$z_{i2} = b_{02} + x_i \odot b_2$$

$\vdots$

$$z_{iK} = b_{0K} + x_i \odot b_K$$



prediction



$$y_i = d_0 + \sigma(\zeta_i) \odot d$$

meta-topics

$c_1$	$c_2$			$c_J$
sports + history	politics + sports	•	•	politics + sports + history

$$\begin{aligned}\zeta_{i1} &= c_{01} + \sigma(z_i) \odot c_1 \\ \zeta_{i2} &= c_{02} + \sigma(z_i) \odot c_2 \\ &\vdots \\ \zeta_{iJ} &= c_{0J} + \sigma(z_i) \odot c_J\end{aligned}$$

topics

$b_1$	$b_2$			$b_K$
sports	history	•	•	politics

$$\begin{aligned}z_{i1} &= b_{01} + x_i \odot b_1 \\ z_{i2} &= b_{02} + x_i \odot b_2 \\ &\vdots \\ z_{iK} &= b_{0K} + x_i \odot b_K\end{aligned}$$