• To explain the clip feature extraction to loss computation with a toy example, we will cover each key stage.

1) Image & text feature extraction.

2) Similarity computation.

3) loss computation.

○ lets use a sample batch of ③ image-text pair for clarity.

---
① Toy Example set up.
---

• We have a batch of ③ image-text pairs.

- $I_1$ matches $T_1$
- $I_2$ matches $T_2$
- $I_3$ matches $T_3$

Step 1  Extract embeddings for images:

$$I_e = \begin{bmatrix} 0.8 & 0.2 & 0.1 \\ 0.1 & 0.9 & 0.2 \\ 0.3 & 0.1 & 0.9 \end{bmatrix} \begin{matrix} I_1 \\ I_2 \\ I_3 \end{matrix}$$

• each row corresponds to the embedding of an image. (dimensionality = 3)

2. Text Encoding: Extract embedding of

$$T_e = \begin{bmatrix} 0.9 & 0.1 & 0.0 \\ 0.2 & 0.8 & 0.0 \\ 0.1 & 0.2 & 0.9 \end{bmatrix} \begin{matrix} T_1 \\ T_2 \\ T_3 \end{matrix}$$

each row is corresponds to a text embedding (d=3)
description

STEP1

we compute the similarity b/w all image-text pairs (not just aligned pairs). cosine similarity is given by

$$\text{similarity}(I_i, T_j) = \frac{I_i \cdot T_j^T}{\|I_i\| \cdot \|T_j\|}$$

(just to make it a dot product b/w the image & text)

For simplicity, assume embeddings are normalized ($i.e \|I_i\| = \|T_j\| = 1$). Then similarity is just a dot product

$$\text{logits} = I_e \cdot T_i^T = \begin{matrix} & \text{Image} \rightarrow \\ \begin{matrix} I_1 \\ I_2 \\ I_3 \end{matrix} & \begin{bmatrix} 0.8 & 0.2 & 0.1 \\ 0.1 & 0.9 & 0.2 \\ 0.3 & 0.1 & 0.9 \end{bmatrix} \end{matrix} \begin{matrix} \begin{matrix} T_1 & T_2 & T_3 \\ \end{matrix} \\ \begin{bmatrix} 0.9 & 0.2 & 0 \\ 0.1 & 0.8 & 0.2 \\ 0.0 & 0.0 & 0.9 \end{bmatrix} \end{matrix}$$

$3\times3$    $3\times3$

why? $e^{(z)}$   why?

logit mean raw output $z = Wx+b$ before apply softmax/sigmoid

$$= \begin{bmatrix} 0.74 & 0.26 & 0.20 \\ 0.29 & 0.74 & 0.26 \\ 0.30 & 0.20 & 0.85 \end{bmatrix}$$

each row corresponds to an image, and each column corresponds to a text

Step3 apply softmax → Convert logits into Probabilities

To calculate probabilities, apply softmax to each row (image-to-text direction) & each column (text-to-image direction)

Image-to-text direction softmax

For row $I_1$ ( 0.74, 0.26, 0.20)

$$Softmax(\bar{I_1}) = \left( \frac{e^{0.74}}{e^{0.74}+e^{0.26}+e^{0.20}} \ , \ \frac{e^{0.26}}{e^{0.7}+e^{0.26}+e^{0.20}} \ , \ \frac{e^{0.2}}{e^{0.7}+e^{0.6}+e^{0.2}} \right)$$

$$Softmax(I_1) = [\ 0.492, \ 0.260, \ 0.248\ ]$$

$$Softmax(I_1) = [0.258, 0.486, 0.258]$$

$$softmax(I_2) = [0.250, 0.249, 0.501]$$

$$Softmax(I_2) = [0.26, 0.492, 0.248]$$

$$Softmax(I_2) = [0.265, 0.241, 0.494]$$

Text-to-Image softmax

$$softmax(t_1) = [0.46, 0.264, 0.271]$$
$$softmax(t_2) = [0.258, 0.486, 0.258]$$
$$softmax(t_3) = [0.250, 0.249, 0.501]$$

(4)

Step 4    Cross entropy

The ground Truth label [0, 1, 2]

Image — to — text loss (Loss 1):

For $I_1$ (label=0)    [1, 0, 0]

$loss(I_1) = -log(0.492) = 0.709$

for $I_2$ (label = 1) [0, 1, 0]

$loss(I_L) = -log(0.492) = 0.709$

for $I_3$ (label=2) [0, 0, 1]

$loss(I_3) = -log(0.494) = 0.705$

np.op(1)

Average loss1 $= \dfrac{0.709 + 0.709 + 0.705}{3} = 0.708$

cross entropy

Text — to — Image loss (Loss 2):

$$loss = -\sum_{i=1}^{n} y_i \log \hat{y}_i$$

For $t_1$ (label=0): [1 0 0]

$loss(T_1) = -log(0.465) = 0.766$

For $T_2$ (label = 1) [0 1 0]

$loss(T_2) = -log(0.486) = 0.721$

For $T_3$ (label=2) [0 0 1]

$loss(T_3) = -log(0.51) = 0.691$

$loss_2 = \dfrac{0.766 + 0.721 + 0.691}{3}$

Avg. $loss_2 = 0.726$

Step 5

$$loss = \frac{Loss1 + loss2}{2} = \frac{0.708 + 0.726}{2}$$

$$= 0.717$$

Summary:

- Feature extraction: creates embedding/numeric values for image & text

- cosine similarity: b/w embedding from logits matrix.

- Softmax: convert logits into probabilities.

- Cross entropy: penalizes misalignments, encouraging matching imag-text pair

# Temperature Scaling

Image 1: A clear picture of a Cat ( Label: Cat )

Image 2: A clear picture of a Dog ( Label: Dog )

Image 3: A blurry picture of a Cat ( Label : Cat )

Image 4: A blurry Picture of a Dog ( label : Dog )

## RAW Similarity Score

| | | cat $T_1$ | dog $T_2$ |
|---|---|---|---|
| clear cat | $I_1$ | 0.9 | 0.1 |
| clear Dog | $I_2$ | 0.2 | 0.8 |
| blurry cat | $I_3$ | 0.6 | 0.4 |
| blurry Dog | $I_4$ | 0.3 | 0.5 |

Adjusting CONFIDENCE LEVEL

(A) ✗

→ ( as the model is very confident in its prediction & may not learn effectively from other, more ambiguous images.)

### w/o temperature Scaling

we apply softmax act & we get

| | | cat $T$ | dog $T_2$ |
|---|---|---|---|
| | | (very confident) | |
| clear cat | $I_1$ | 0.9 | 0.1 |
| clear Dog | $I_2$ | 0.2 | (0.8) (very confident) |
| blurry cat | $I_3$ | 0.6 | 0.4 |
| blurry Dog | $I_4$ | 0.3 | 0.5 |

## Applying temperature scaling

Now, let's apply temperature scaling $T = 2$

Similarly $(I, t)/2$

| | $T_1$ | $T_2$ |
|---|---|---|
| $I_1$ | 0.45 | 0.05 |
| $I_2$ | 0.1 | 0.4 |
| $I_3$ | 0.3 | 0.2 |
| $I_4$ | 0.15 | 0.25 |

- New probabilities after scaling

original RBN softdet
(logits)

| 0.9 | 0.1 |
|---|---|
| 0.2 | 0.8 |
| 0.6 | 0.4 |
| 0.3 | 0.8 |

| | $T_1$ | $T_2$ |
|---|---|---|
| $I_1$ | 0.6 | 0.4 |
| $I_2$ | 0.425 | 0.575 |
| $I_3$ | 0.525 | 0.425 |
| $I_4$ | 0.475 | 0.525 |

By applying a higher temperature (like $T = 2$)
the model's confidence is tempered, resulting in more balanced
probability. By adjusting a probability of 0.95 for cat + 0.60
indicating that there is still some uncertainty about classification

⇒ uncertain → increase loss → gradient high,
↳ Weights update to learn complex features.

⇒ more **robust** feature learned.

⇒ clear & unclear image. both
cases feature learned.

⇒ learn feature in both clear
& unclear image. → baking
learning.