

Deep Learning

7 Cross-Entropy Loss

Dr. Konda Reddy Mopuri
Dept. of Artificial Intelligence
IIT Hyderabad
Jan-May 2023

Classification



① Dataset looks like $(x_n, y_n) \in \mathcal{R}^D \times \{1, 2, \dots, C\}, n = 1, 2, \dots, N$

Classification



- ① Dataset looks like $(x_n, y_n) \in \mathcal{R}^D \times \{1, 2, \dots, C\}, n = 1, 2, \dots, N$
- ② We don't generally ($C > 2$) regress the target (Why not?)

Classification



- ① Dataset looks like $(x_n, y_n) \in \mathcal{R}^D \times \{1, 2, \dots, C\}, n = 1, 2, \dots, N$
- ② We don't generally ($C > 2$) regress the target (Why not?)
- ③ In other words, we don't prefer MSE loss for learning

Classification Intuition



- ① Target label y is one-hot encoded

Classification Intuition



- ① Target label y is one-hot encoded
- ② It converts y to a pmf (\mathbf{p}) (e.g., $y_n = 2 \rightarrow \{0, 1, 0, 0\}$ and $y_n = 3 \rightarrow \{0, 0, 1, 0\}$ when $C = 4$)

Classification Intuition



- ① Target label y is one-hot encoded
- ② It converts y to a pmf (\mathbf{p}) (e.g., $y_n = 2 \rightarrow \{0, 1, 0, 0\}$ and $y_n = 3 \rightarrow \{0, 0, 1, 0\}$ when $C = 4$)
- ③ Hence, the DNN's prediction should also be a pmf (\mathbf{q})

Classification Intuition



- ① Target label y is one-hot encoded
- ② It converts y to a pmf (\mathbf{p}) (e.g., $y_n = 2 \rightarrow \{0, 1, 0, 0\}$ and $y_n = 3 \rightarrow \{0, 0, 1, 0\}$ when $C = 4$)
- ③ Hence, the DNN's prediction should also be a pmf (\mathbf{q})
- ④ Loss function should compare \mathbf{p} and \mathbf{q}

Very very brief discussion on related Information Theory



- ① Information contained in an event x can be computed given the probability of that event $P(x)$

Very very brief discussion on related Information Theory



- ① Information contained in an event x can be computed given the probability of that event $P(x)$
- ② Higher the $P(x)$, lesser is the information (less 'surprising')

Very very brief discussion on related Information Theory



- ① Information contained in an event x can be computed given the probability of that event $P(x)$
- ② Higher the $P(x)$, lesser is the information (less 'surprising')
- ③ Hence, the information can be calculated as $I(x) = -\log_2(P(x))$

Very very brief discussion on related Information Theory



- ① Information contained in an event x can be computed given the probability of that event $P(x)$
- ② Higher the $P(x)$, lesser is the information (less 'surprising')
- ③ Hence, the information can be calculated as $I(x) = -\log_2(P(x))$
- ④ This is also the number of bits required to encode x

Very very brief discussion on related Information Theory



- ① Entropy is the number of bits required to encode a randomly chosen message (x) from a probability distribution $p(x)$

Very very brief discussion on related Information Theory



- ① Entropy is the number of bits required to encode a randomly chosen message (x) from a probability distribution $p(x)$
- ② Expected amount of information in an event drawn from that distribution $H(X) = \mathbb{E}_{x \sim p}[I(x)]$

Very very brief discussion on related Information Theory



① One message x needs $-\log(P(x))$ bits

Very very brief discussion on related Information Theory



- ① One message x needs $-\log(P(x))$ bits
- ② There are multiple messages with associated probabilities \rightarrow entropy
$$H(X) = -\sum P(x) \cdot \log_2(P(x))$$

Very very brief discussion on related Information Theory



- ① One message x needs $-\log(P(x))$ bits
- ② There are multiple messages with associated probabilities \rightarrow entropy
$$H(X) = -\sum P(x) \cdot \log_2(P(x))$$
- ③ $H(p) = -\sum_i p_i \cdot \log_2(p_i)$

Very very brief discussion on related Information Theory



- ① One message x needs $-\log(P(x))$ bits
- ② There are multiple messages with associated probabilities \rightarrow entropy
$$H(X) = -\sum P(x) \cdot \log_2(P(x))$$
- ③
$$H(p) = -\sum_i p_i \cdot \log_2(p_i)$$
- ④ Skewed distribution has less entropy, uniform/balanced distribution has more entropy

Very very brief discussion on related Information Theory

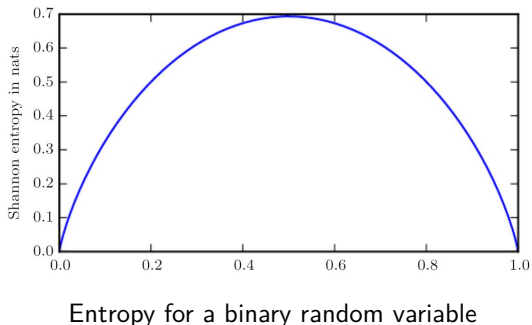


Figure credits Goodfellow et al. 2016

Very very brief discussion on related Information Theory



- ① Cross-entropy $H(p, q)$ is the average number of bits required to encode the messages from a source distribution p when encoded with a different model q

Very very brief discussion on related Information Theory



- ① Cross-entropy $H(p, q)$ is the average number of bits required to encode the messages from a source distribution p when encoded with a different model q
- ② $H(p, q) = - \sum_i p_i \cdot \log_2(q_i)$

Very very brief discussion on related Information Theory



- ① Cross-entropy $H(p, q)$ is the average number of bits required to encode the messages from a source distribution p when encoded with a different model q
- ② $H(p, q) = - \sum_i p_i \cdot \log_2(q_i)$
- ③ Note that cross-entropy is not symmetric metric, i.e., $H(p, q) \neq H(q, p)$

Very very brief discussion on related Information Theory



- ① Cross-entropy $H(p, q)$ is the average number of bits required to encode the messages from a source distribution p when encoded with a different model q
- ② $H(p, q) = - \sum_i p_i \cdot \log_2(q_i)$
- ③ Note that cross-entropy is not symmetric metric, i.e., $H(p, q) \neq H(q, p)$
- ④ Cross-entropy between a distribution and itself ($H(p, p)$) gives the entropy of the distribution $H(p)$

Very very brief discussion on related Information Theory



- ① KL-Divergence : average number of **extra** bits required to represent a message with distribution q instead of p

Very very brief discussion on related Information Theory



- ① KL-Divergence : average number of **extra** bits required to represent a message with distribution q instead of p
- ② $H(p, q) = H(p) + KL(p||q)$ where $KL(p||q) = \sum p_i \cdot \log\left(\frac{p_i}{q_i}\right)$

Cross-entropy as a loss function



- ① Widely used in classification problems (e.g. logistic regression, NNs)

Cross-entropy as a loss function



- ① Widely used in classification problems (e.g. logistic regression, NNs)
- ② Each label is converted into a distribution with 1 and 0s (one-hot encoding)

Cross-entropy as a loss function



- ① Widely used in classification problems (e.g. logistic regression, NNs)
- ② Each label is converted into a distribution with 1 and 0s (one-hot encoding)
- ③ Model predicts the probabilities that sample belongs to different classes

Cross-entropy as a loss function



- ① Random variable is the sample

Cross-entropy as a loss function



- ① Random variable is the sample
- ② Events are the classes

Cross-entropy as a loss function



- ① Random variable is the sample
- ② Events are the classes
- ③ Target distribution (or, groundtruth) is one-hot encoding p , and model predicts a distribution q

Cross-entropy as a loss function



- ① Random variable is the sample
- ② Events are the classes
- ③ Target distribution (or, groundtruth) is one-hot encoding p , and model predicts a distribution q
- ④ `torch.nn.CrossEntropyLoss(q, p)` (PyTorch takes predicted distribution as the first argument)

Softmax



- ① Typically last layer in the DNN classifier is linear (without a nonlinearity)

Softmax



- ① Typically last layer in the DNN classifier is linear (without a nonlinearity)
- ② Predicts the confidences to each class (may not lie in $[0, 1]$)

Softmax



- ① Typically last layer in the DNN classifier is linear (without a nonlinearity)
- ② Predicts the confidences to each class (may not lie in $[0, 1]$)
- ③ But, we need probabilities

Softmax



- ① Typically last layer in the DNN classifier is linear (without a nonlinearity)
- ② Predicts the confidences to each class (may not lie in $[0, 1]$)
- ③ But, we need probabilities
- ④ Softmax operation
 - squashes the predicted confidences to lie in $[0, 1]$

Softmax



- ① Typically last layer in the DNN classifier is linear (without a nonlinearity)
- ② Predicts the confidences to each class (may not lie in $[0, 1]$)
- ③ But, we need probabilities
- ④ Softmax operation
 - squashes the predicted confidences to lie in $[0, 1]$
 - make them probabilities (i.e. sum to 1)

Softmax

$$\textcircled{1} (\alpha_1, \alpha_2, \dots, \alpha_C) \rightarrow \left(\frac{e^{\alpha_1}}{\sum_i e^{\alpha_i}}, \frac{e^{\alpha_2}}{\sum_i e^{\alpha_i}}, \dots, \frac{e^{\alpha_C}}{\sum_i e^{\alpha_i}} \right)$$

Softmax

$$\textcircled{1} (\alpha_1, \alpha_2, \dots, \alpha_C) \rightarrow \left(\frac{e^{\alpha_1}}{\sum_i e^{\alpha_i}}, \frac{e^{\alpha_2}}{\sum_i e^{\alpha_i}}, \dots, \frac{e^{\alpha_C}}{\sum_i e^{\alpha_i}} \right)$$

$$\textcircled{2} (\alpha_1, \alpha_2, \dots, \alpha_C) \rightarrow (q_1, q_2, \dots, q_C)$$

Cross-entropy



- ① Target distribution p has 1 at the position of correct label and 0 at rest of the components

Cross-entropy



- ① Target distribution p has 1 at the position of correct label and 0 at rest of the components
- ② $H(p, q) = -\sum p_i \cdot \log(q_i) = -\log(q_c)$, where c is the groundtruth class of the sample

Cross-entropy



- ① Target distribution p has 1 at the position of correct label and 0 at rest of the components
- ② $H(p, q) = -\sum p_i \cdot \log(q_i) = -\log(q_c)$, where c is the groundtruth class of the sample
- ③ The cross-entropy loss is
 - small when the model predicts high probability to the groundtruth class ($q_c \approx 1$)

Cross-entropy



- ① Target distribution p has 1 at the position of correct label and 0 at rest of the components
- ② $H(p, q) = -\sum p_i \cdot \log(q_i) = -\log(q_c)$, where c is the groundtruth class of the sample
- ③ The cross-entropy loss is
 - small when the model predicts high probability to the groundtruth class ($q_c \approx 1$)
 - large if the model assigns smaller probability for the groundtruth class ($q_c \approx 0$)

Variations/Additions



- ① BCE: Binary Cross Entropy Loss

Variations/Additions



- ① BCE: Binary Cross Entropy Loss
- ② Softmax with temperature