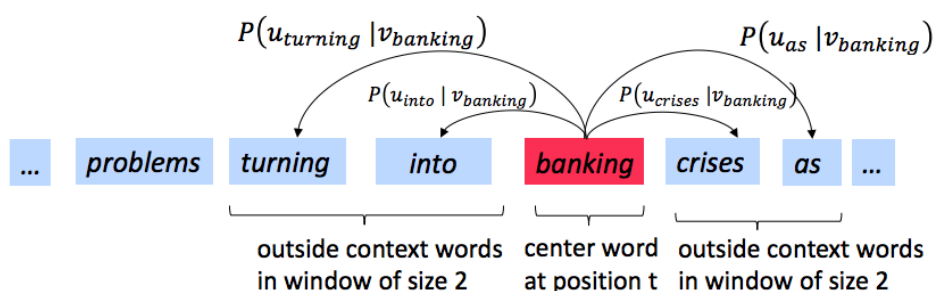


1 手写作业：理解 word2vec

还记得吗？word2vec 背后的核心想法：“一个单词的意义取决于其搭配”，具体来说，对于一个中心词 c ，它被前后特定长度的上下文围绕。我们将这个上下文窗口中的单词定义为外部词(‘outside words’)。举个例子，在图 1 中，上下文窗口长度为 2，中心词 c 是 ‘banking’，外部词就是 ‘turning’, ‘into’, ‘crises’ 和 ‘as’。



Skip-gram word2vec 的目标是学习一个概率分布 $P(O|C)$ 。特殊一点的话，给定一个特定的单词 o 和一个特定的单词 c ，我们想要预测 $P(O=o|C=c)$ ，即单词 o 是单词 c 的外部词的概率(单词 o 落在 c 的上下文窗口里)。我们将通过在一组向量点积上求 softmax 函数对这个概率建模。

$$P(O=o|C=c) = \frac{\exp(\mathbf{u}_o^\top \mathbf{v}_c)}{\sum_{w \in \text{Vocab}} \exp(\mathbf{u}_w^\top \mathbf{v}_c)} \quad (1)$$

对每一个单词，我们学习向量 \mathbf{u} 和 \mathbf{v} ，其中 \mathbf{u}_o 是代表外部词 o 的外部向量， \mathbf{v}_c 是代表中心词的中心向量。我们将这些参数存储在两个矩阵 U 和 V 当中。 U 的各列是所有的外部向量 \mathbf{u}_w ， V 的各列是所有的中心向量 \mathbf{v}_w 。 U 和 V 都包含词汇表中每一个单词 w 对应的向量。

还记得吗？对于一对单词 c 和 o ，损失值由以下表达式给出：

$$J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U}) = -\log P(O = o | C = c) \quad (2)$$

我们可以将这个损失视为真实分布 \mathbf{y} 和预测分布 $\hat{\mathbf{y}}$ 的交叉熵损失，对于一个特定的中心词 c 和一个特定的外部词 o 。在此处， \mathbf{y} 和 $\hat{\mathbf{y}}$ 都是长度等于词汇表中单词数的向量。而且，这些向量的第 k 个分量表示第 k 个单词成为给定中心词 c 的外部词的条件概率。真正的经验分布 \mathbf{y} 是一个独热向量，在真正的外部词对应的分量上的值为 1，其他值为 0。预测分布 $\hat{\mathbf{y}}$ 是方程(1)中的模型给出的概率分布 $P(O|C=c)$ 。

注意：在整个作业中，当计算导数时，请使用课程中复习到的方法(不要使用泰勒级数近似)

(a) 请证明方程 2 中的朴素 softmax 损失和 \mathbf{y} - $\hat{\mathbf{y}}$ 交叉熵损失等价(注意 \mathbf{y} 和 $\hat{\mathbf{y}}$ 都是向量，而 $\hat{\mathbf{y}}_o$ 是标量)

$$-\sum_{w \in \text{Vocab}} \mathbf{y}_w \log(\hat{\mathbf{y}}_w) = -\log(\hat{\mathbf{y}}_o).$$

My answer: $\mathbf{y}_w = 1$ if $w = o$ else $\mathbf{y}_w = 0$

(b)

① 请计算 $J_{\text{naive-softmax}}(\mathbf{v}_c, o, \mathbf{U})$ 关于 \mathbf{v}_c 的偏导数，结果以向量形式表示

$$\frac{\partial J_{naive-softmax}(\mathbf{v}_c, \mathbf{o}, U)}{\partial \mathbf{v}_c} = - \left[\mathbf{u}_o - \sum_{x \in Vocab} \frac{\exp(\mathbf{u}_x^T \mathbf{v}_c)}{\sum_{w \in Vocab} \exp(\mathbf{u}_w^T \mathbf{v}_c)} \mathbf{u}_x \right]$$

$$= -[\mathbf{u}_o - E[\mathbf{u}_w]]$$

② 什么情况下计算出的梯度等于 0？

My answer: $\mathbf{u}_o = E[\mathbf{u}_w]$ ，即模型对外部词向量的估计值等于外部词向量的实际值时计算出的梯度等于 0。

③ 你找到的梯度是两项的差，当词向量 \mathbf{v}_c 减去梯度时请对梯度中的每一项如何改进词向量进行解释说明

My answer: 负梯度的值为 $\mathbf{u}_o - E[\mathbf{u}_w]$ ，表示参数的更新将使得模型对外部词向量的估计更接近真实观测到的外部词向量

④ 在许多使用词向量的下游应用中，使用的是 L2 正则化向量而非它们的原始形式。现在，假定你想要将词组分类成正例或者反例。

什么时候 L2 正则化向量会从下游任务提取有效信息？什么时候不会？提示：考虑当单词 $x \neq y$ ，存在标量 α 使得 $\mathbf{u}_x = \alpha \mathbf{u}_y$ 的情况

My answer: 当不同单词对应的词向量共线时，L2 正则化后不同的单词将对应相同的词向量表示，此时无法从下游任务中提取有效信息；而当这种情况不存在或在小部分单词上存在时可以从下游任务中提取有效信息。

(c) 请计算 $J_{naive-softmax}(v_c, o, U)$ 关于每一个外部词向量 u_w 的偏导数。

这将会会有两种情况:当 $w=o$ 时,对真正的外部词向量,当 $w \neq o$ 时,对所有其他单词。请写出关于 y , \hat{y} 和 v_c 的答案。在这个子部分,你可能会用到这些项中特定的元素,比如 y_1, y_2, \dots , 请注意 u_w 是一个向量而 y_1, y_2, \dots 是标量。

My answer:

$$\frac{\partial J_{naive-softmax}(v_c, o, U)}{\partial u_w} = \begin{cases} -[u_w - y_w v_c] & (w = o) \\ -[-y_w v_c] & (w \neq o) \end{cases}$$

$$y_w = \frac{\exp(u_w^T v_c)}{\sum_{w \in Vocab} \exp(u_w^T v_c)}$$

(d) 请写出 $J_{naive-softmax}(v_c, o, U)$ 关于 U 的偏导数。请将答案分解成各个列向量的组合。这里不需要用到具体的导数,只要矩阵形式的答案即可。

My answer:

$$\frac{\partial J_{naive-softmax}(v_c, o, U)}{\partial U} = v_c \hat{y}^T - u_o y$$

(e) Leaky ReLU 激活函数通过方程 4 和图 2 给出,

$$f(x) = \max(\alpha x, x) \quad (4)$$

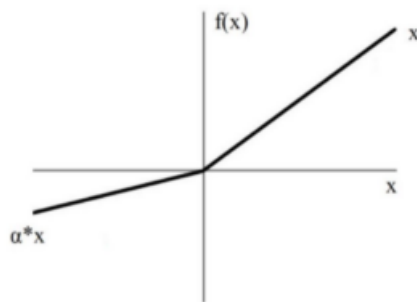


Figure 2: Leaky ReLU

x 是一个标量, $0 < \alpha < 1$, 请计算 $f(x)$ 关于 x 的导数, 你可以忽略在 0 处导数未定义的情况。

My answer:

$$\frac{df}{dx} = \begin{cases} 1 & (x > 0) \\ \alpha & (x < 0) \end{cases}$$

(f) sigmoid 函数由方程 5 给出:

$$\sigma(x) = \frac{1}{1 + e^{-x}} = \frac{e^x}{e^x + 1} \quad (5)$$

请计算 $\sigma(x)$ 关于 x 的导数, x 是一个标量。请写出关于 $\sigma(x)$ 的答案。

My answer:

$$\sigma'(x) = \frac{e^{-x}}{(1 + e^{-x})^2} = \sigma(x) (1 - \sigma(x))$$

(g) 现在我们考虑负采样损失, 这是朴素 softmax 损失的一个选项。

假定从词汇表中抽取 K 个负样本。为了方便表示我们记为

w_1, w_2, \dots, w_K , 并将外部词向量记为 $\mathbf{u}_{w_1}, \mathbf{u}_{w_2}, \dots, \mathbf{u}_{w_K}$ 。对这个问

题, 假定 K 个负样本相异, 即 $i \neq j$, 有 $w_i \neq w_j$, 对 $i, j \in \{1, 2, \dots, K\}$ 。

请注意 $o \notin \{w_1, w_2, \dots, w_K\}$ 。对于一个中心词 c 和一个外部词 o ，负采样损失函数通过下式给出：

$$J_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \quad (6)$$

- ① 请重复(b), (c)，计算 $J_{\text{neg-sample}}$ 分别关于 \mathbf{v}_c ， \mathbf{u}_o 和第 s 个负样本 \mathbf{u}_{w_s} 的偏导数。答案关于 \mathbf{v}_c ， \mathbf{u}_o 和 \mathbf{u}_{w_s} ， $s \in \{1, 2, \dots, K\}$ 。注意：你应该使用(f)中的答案帮助计算这里必需的梯度。

My answer:

$$J_{\text{neg-sample}} = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c))$$

$$\frac{\partial J_{\text{neg-sample}}}{\partial \mathbf{v}_c} = -[1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)]\mathbf{u}_o + \sum_{s=1}^K [1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)]\mathbf{u}_{w_s}$$

$$\frac{\partial J_{\text{neg-sample}}}{\partial \mathbf{u}_o} = -[1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)]\mathbf{v}_c$$

$$\frac{\partial J_{\text{neg-sample}}}{\partial \mathbf{u}_{w_s}} = [1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)]\mathbf{v}_c$$

- ② 在课程中，我们学到了反向传播的一种高效实现，这种实现复用了之前计算过的偏导数。你在多大程度上能复用上面计算出的偏导数以最小化计算的冗余程度？请写出关于 $U_{o, \{w_1, \dots, w_K\}} = [\mathbf{u}_o, -\mathbf{u}_{w_1}, \dots, -\mathbf{u}_{w_K}]$ (一个由外部词向量按列放置形成的矩阵) 和 $\mathbf{1}$ (一个形状为 $(K+1) \times 1$ 的全 1 向量) 的答案，除 $U_{o, \{w_1, \dots, w_K\}}$ 和 $\mathbf{1}$ 之外的其他项也能在答案中使用。

My answer:

$$U_{o, \{w_1, \dots, w_K\}} = [u_o, -u_{w_1}, \dots, -u_{w_K}]$$

$$\frac{\partial J_{neg-sample}}{\partial U_{o, \{w_1, \dots, w_K\}}} = -\mathbf{v}_c \cdot [1^T - \sigma(\mathbf{v}_c^T U_{o, \{w_1, \dots, w_K\}})]$$

$$\frac{\partial J_{neg-sample}}{\partial \mathbf{v}_c} = -U_{o, \{w_1, \dots, w_K\}} [1 - \sigma(U_{o, \{w_1, \dots, w_K\}}^T \mathbf{v}_c)]$$

③ 用一句话描述为什么这个损失函数比朴素 softmax 损失在计算上更高效

My answer:

负采样损失函数只在 K 个样本上计算点积，而朴素 softmax 损失函数需要在全部词向量上计算点积。

到目前为止，我们已经研究了通过变量复用和采样近似 softmax 以实现更快的梯度下降。但请注意这些优化的一部分在现代 GPU 上是不必要的，并且在某种程度上，是当时计算资源有限这一背景的产物。

(h) 现在我们将重复之前的练习，但是不使用 K 个采样词相异这一假设。假定从词汇表中抽取 K 个负采样词，为了方便表示我们记为 w_1, w_2, \dots, w_K ，并将外部词向量记为 $\mathbf{u}_{w_1}, \mathbf{u}_{w_2}, \dots, \mathbf{u}_{w_K}$ 。在这一问题中，你将不再假设这些词相异。换言之，当 $i \neq j$ 时，仍然可能有 $w_i = w_j$ ，请注意 $o \notin \{w_1, \dots, w_K\}$ 。对一个中心词 c 和一个外部词 o，负采样损失函数由下式给出：

$$J_{\text{neg-sample}}(\mathbf{v}_c, o, \mathbf{U}) = -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \quad (7)$$

请计算 $J_{\text{neg-sample}}$ 关于 \mathbf{u}_{w_s} 的偏导数，请写出关于 \mathbf{v}_c 和 \mathbf{u}_{w_s} 的答案， $s \in \{1, 2, \dots, K\}$ 。提示：将损失函数的和式分解为两个部分的和：一部分和基于所有等于 w_s 的采样词，另一部分和基于所有不等于 w_s 的采样词。为方便表示，你可以在求和符号下写“等于”和“不等于”的条件。

My answer:

$$\begin{aligned} J_{\text{neg-sample}} &= -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{s=1}^K \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \\ &= -\log(\sigma(\mathbf{u}_o^\top \mathbf{v}_c)) - \sum_{w_s} \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) - \sum_{\neq w_s} \log(\sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)) \end{aligned}$$

$$\frac{\partial J_{\text{neg-sample}}}{\partial \mathbf{u}_{w_s}} = \sum_{w_s} [1 - \sigma(-\mathbf{u}_{w_s}^\top \mathbf{v}_c)] \mathbf{v}_c$$

(i) 假定中心词 $c = w_t$ ，上下文窗口是

$[w_{t-m}, \dots, w_{t-1}, w_t, w_{t+1}, \dots, w_{t+m}]$ ，其中 m 是上下文窗口的大小。还记得吗？在 skip-gram 版本的 word2vec 中整个上下文窗口的总损失为

$$J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, \mathbf{U}) = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} J(\mathbf{v}_c, w_{t+j}, \mathbf{U})$$

在这里， $J(\mathbf{v}_c, w_{t+j}, \mathbf{U})$ 代表关于中心词 $c = w_t$ 和外部词 w_{t+j} 的任意一个损失项， $J(\mathbf{v}_c, w_{t+j}, \mathbf{U})$ 可以是 $J_{\text{naive-softmax}}(\mathbf{v}_c, w_{t+j}, \mathbf{U})$ 或者 $J_{\text{neg-sample}}(\mathbf{v}_c, w_{t+j}, \mathbf{U})$ ，取决于你的实现。

请写出以下三个偏导数：

- (i) $\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial U}$
- (ii) $\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial \mathbf{v}_c}$
- (iii) $\frac{\partial J_{\text{skip-gram}}(\mathbf{v}_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial \mathbf{v}_w}$ when $w \neq c$

请写出关于 $\frac{\partial J(\mathbf{v}_c, w_{t+j}, U)}{\partial U}$ 和 $\frac{\partial J(\mathbf{v}_c, w_{t+j}, U)}{\partial \mathbf{v}_c}$ 的答案。这非常简单，每个答案应该只有一行。

My answer:

$$\frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial U} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(v_c, w_{t+j}, U)}{\partial U}$$

$$\frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial \mathbf{v}_c} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(v_c, w_{t+j}, U)}{\partial \mathbf{v}_c}$$

$$\frac{\partial J_{\text{skip-gram}}(v_c, w_{t-m}, \dots, w_{t+m}, U)}{\partial \mathbf{v}_w} = \sum_{\substack{-m \leq j \leq m \\ j \neq 0}} \frac{\partial J(v_c, w_{t+j}, U)}{\partial U} \mathbf{y}_j \text{ (a little puzzle)}$$

一旦你完成了上述所有练习：假定你通过(a)-(c)计算出了 $J(v_c, w_{t+j}, U)$ 关于 U 和 V 所有模型参数的导数，你就已经计算出所有损失函数关于所有参数的导数，你就做好了实现 word2vec 的准备。