

Homework 1

Problem 1

Given a D-dimensional vector \mathbf{x} (x_k is the k-th component of \mathbf{x} , and $1 \leq k \leq D$):

$$\text{softmax}(\mathbf{x})_j = \frac{e^{x_j}}{\sum_{k=1}^D e^{x_k}}.$$

Then,

$$\text{softmax}(\mathbf{x} + c)_j = \frac{e^{x_j + c}}{\sum_{k=1}^D e^{x_k + c}} = \frac{e^c e^{x_j}}{\sum_{k=1}^D e^c e^{x_k}} = \frac{e^c e^{x_j}}{e^c \sum_{k=1}^D e^{x_k}} = \frac{e^{x_j}}{\sum_{k=1}^D e^{x_k}} = \text{softmax}(\mathbf{x})_j$$

Since we have shown that for all components j $\text{softmax}(\mathbf{x} + c)_j = \text{softmax}(\mathbf{x})_j$, then it follows that $\text{softmax}(\mathbf{x} + c) = \text{softmax}(\mathbf{x})$.

Problem 2

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

$$\begin{aligned} \frac{\partial}{\partial x} \sigma(x) &= -(1 + e^{-x})^{-2} * (e^{-x} * -1) = \frac{e^{-x}}{(1 + e^{-x})^2} = \frac{1}{1 + e^{-x}} * \frac{e^{-x}}{1 + e^{-x}} = \frac{1}{1 + e^{-x}} * \left(\frac{1 + e^{-x}}{1 + e^{-x}} - \frac{1}{1 + e^{-x}} \right) \\ &= \sigma(x) * (1 - \sigma(x)) \end{aligned}$$

$\sigma'(x) = \sigma(x) * (1 - \sigma(x))$
--

b) $\frac{\partial}{\partial \theta} CE = \frac{\partial}{\partial \theta} \left[-\sum_i y_i * \log \frac{e^{\theta_i}}{\sum_{k=1}^D e^{\theta_k}} \right] = \frac{\partial}{\partial \theta} \left[-\sum_i y_i * (\theta_i - \log \sum_{k=1}^D e^{\theta_k}) \right]$

When $j = i$:

$$\frac{\partial}{\partial \theta} CE_j = - \left(y_i * 1 - \frac{1}{\sum_{k=1}^D e^{\theta_k}} * e^{\theta_i} \right) = y_i * \left(\frac{e^{\theta_i}}{\sum_{k=1}^D e^{\theta_k}} - 1 \right) = y_i * (\hat{y}_i - 1)$$

$y_i = 1$ since y is the 1-hot label, so:

$$\frac{\partial}{\partial \theta} CE_j = \hat{y}_j - 1 \text{ when } j = i.$$

When $j \neq i$:

$$\frac{\partial}{\partial \theta} CE_j = - \left(y_i * 0 - \frac{1}{\sum_{k=1}^D e^{\theta_k}} * e^{\theta_j} \right) = \frac{e^{\theta_j}}{\sum_{k=1}^D e^{\theta_k}} = \hat{y}_j$$

Thus,

$\frac{\partial}{\partial \theta} CE = \hat{\mathbf{y}} - \mathbf{y}.$
--

- c) Let's define $\theta = hW_2 + b_2$, $\hat{y} = \text{softmax}(\theta)$, $h = \text{sigmoid}(a)$, and $a = xW_1 + b_1$.
Using the chain rule:

$$\frac{\partial J}{\partial x} = \left(\left(\frac{\partial J}{\partial \theta} * \frac{\partial \theta}{\partial h} \right) \circ \frac{\partial h}{\partial a} \right) * \frac{\partial a}{\partial x}$$

Here \circ represents element-wise product and $*$ represents dot product. We take the element-wise product with $\frac{\partial h}{\partial a}$ because h is an element-wise function.

From (b) we know that for cross-entropy cost J , $\frac{\partial J}{\partial \theta} = \hat{y} - y$.

$$\frac{\partial \theta}{\partial h} = W_2^T$$

$$\frac{\partial h}{\partial a} = h \circ (1 - h)$$

$$\frac{\partial a}{\partial x} = W_1^T$$

Thus,

$$\frac{\partial J}{\partial x} = \left((\hat{y} - y) * W_2^T \right) \circ h \circ (1 - h) * W_1^T.$$

We know that, since J is a scalar, $\frac{\partial J}{\partial x}$ must have the same dimensions as x . Let's do some dimensional analysis on our answer above to make sure this is the case.

$\hat{y} - y$ has dimensions $1 \times D_y$

W_2 has dimensions $H \times D_y$

h and $1 - h$ both have dimensions $1 \times H$

W_1 has dimensions $D_x \times H$

Thus $\frac{\partial J}{\partial x}$ has dimensions $(1 \times D_y) * (D_y \times H) \circ (1 \times H) * (H \times D_x) = (1 \times H) * (H \times D_x) = 1 \times D_x$, which is the same dimensions as x .

- d) The parameters are W_1 (D_x by H), W_2 (H by D_y), b_1 (1 by H), b_2 (1 by D_y).

$$\text{Total \# of params} = HD_x + HD_y + H + D_y$$

Problem 3

a) We are assuming cross-entropy cost, so for a given predicted word at index i:

$$J = -\log P(\text{word}_i | \hat{\mathbf{r}}, \mathbf{w}) = -\log \frac{\exp(\mathbf{w}_i^T \hat{\mathbf{r}})}{\sum_{j=1}^{|V|} \exp(\mathbf{w}_j^T \hat{\mathbf{r}})} = -\mathbf{w}_i^T \hat{\mathbf{r}} + \log \sum_{j=1}^{|V|} \exp(\mathbf{w}_j^T \hat{\mathbf{r}})$$

$$\begin{aligned} \frac{\partial J}{\partial \hat{\mathbf{r}}} &= -\hat{\mathbf{r}} + \frac{1}{\sum_{j=1}^{|V|} \exp(\mathbf{w}_j^T \hat{\mathbf{r}})} * \sum_{k=1}^{|V|} \exp(\mathbf{w}_k^T \hat{\mathbf{r}}) * \mathbf{w}_k = -\hat{\mathbf{r}} + \sum_{k=1}^{|V|} \frac{\exp(\mathbf{w}_k^T \hat{\mathbf{r}}) * \mathbf{w}_k}{\sum_{j=1}^{|V|} \exp(\mathbf{w}_j^T \hat{\mathbf{r}})} = \\ &= -\hat{\mathbf{r}} + \sum_{k=1}^{|V|} P(\text{word}_k | \hat{\mathbf{r}}, \mathbf{w}) * \mathbf{w}_k \end{aligned}$$

$$\boxed{\frac{\partial J}{\partial \hat{\mathbf{r}}} = -\mathbf{w}_i + \sum_{k=1}^{|V|} P(\text{word}_k | \hat{\mathbf{r}}, \mathbf{w}) * \mathbf{w}_k}$$

b) $J = -\log P(\text{word}_i | \hat{\mathbf{r}}, \mathbf{w}) = -\log \frac{\exp(\mathbf{w}_i^T \hat{\mathbf{r}})}{\sum_{j=1}^{|V|} \exp(\mathbf{w}_j^T \hat{\mathbf{r}})} = -\mathbf{w}_i^T \hat{\mathbf{r}} + \log \sum_{j=1}^{|V|} \exp(\mathbf{w}_j^T \hat{\mathbf{r}})$

When $j = i$:

$$\frac{\partial J}{\partial \mathbf{w}_j} = -\hat{\mathbf{r}} + \frac{1}{\sum_{k=1}^{|V|} \exp(\mathbf{w}_k^T \hat{\mathbf{r}})} * \exp(\mathbf{w}_j^T \hat{\mathbf{r}}) * \hat{\mathbf{r}}$$

When $j \neq i$:

$$\frac{\partial J}{\partial \mathbf{w}_j} = \frac{1}{\sum_{k=1}^{|V|} \exp(\mathbf{w}_k^T \hat{\mathbf{r}})} * \exp(\mathbf{w}_j^T \hat{\mathbf{r}}) * \hat{\mathbf{r}}$$

We have:

$$\frac{\partial J}{\partial \mathbf{w}_j} = -\hat{\mathbf{r}} + \frac{\exp(\mathbf{w}_j^T \hat{\mathbf{r}}) * \hat{\mathbf{r}}}{\sum_{k=1}^{|V|} \exp(\mathbf{w}_k^T \hat{\mathbf{r}})} \text{ when } j = i$$

$$\frac{\partial J}{\partial \mathbf{w}_j} = \frac{\exp(\mathbf{w}_j^T \hat{\mathbf{r}}) * \hat{\mathbf{r}}}{\sum_{k=1}^{|V|} \exp(\mathbf{w}_k^T \hat{\mathbf{r}})} \text{ when } j \neq i$$

c)

$$J = -\log \sigma(\mathbf{w}_i^T \hat{\mathbf{r}}) - \sum_{k=1}^K \log \sigma(-\mathbf{w}_k^T \hat{\mathbf{r}})$$

$$\frac{\partial J}{\partial \hat{\mathbf{r}}} = -\frac{1}{\sigma(\mathbf{w}_i^T \hat{\mathbf{r}})} * \sigma(\mathbf{w}_i^T \hat{\mathbf{r}}) * (1 - \sigma(\mathbf{w}_i^T \hat{\mathbf{r}})) * \mathbf{w}_i - \sum_{k=1}^K \frac{1}{\sigma(-\mathbf{w}_k^T \hat{\mathbf{r}})} * \sigma(-\mathbf{w}_k^T \hat{\mathbf{r}}) * (1 - \sigma(-\mathbf{w}_k^T \hat{\mathbf{r}})) (-\mathbf{w}_k)$$

$$\frac{\partial J}{\partial \hat{\mathbf{r}}} = \mathbf{w}_i * (\sigma(\mathbf{w}_i^T \hat{\mathbf{r}}) - 1) + \sum_{k=1}^K \mathbf{w}_k (1 - \sigma(-\mathbf{w}_k^T \hat{\mathbf{r}}))$$

When $j = i$:

$$\frac{\partial J}{\partial \mathbf{w}_j} = -\frac{1}{\sigma(\mathbf{w}_j^T \hat{\mathbf{r}})} * \sigma(\mathbf{w}_j^T \hat{\mathbf{r}}) (1 - \sigma(\mathbf{w}_j^T \hat{\mathbf{r}})) \hat{\mathbf{r}}$$

$$\frac{\partial J}{\partial \mathbf{w}_j} = \hat{\mathbf{r}} (\sigma(\mathbf{w}_j^T \hat{\mathbf{r}}) - 1)$$

When $j \neq i$:

$$\frac{\partial J}{\partial \mathbf{w}_j} = -\frac{1}{\sigma(-\mathbf{w}_j^T \hat{\mathbf{r}})} * \sigma(-\mathbf{w}_j^T \hat{\mathbf{r}}) (1 - \sigma(-\mathbf{w}_j^T \hat{\mathbf{r}})) (-\hat{\mathbf{r}})$$

$$\frac{\partial J}{\partial \mathbf{w}_j} = -\hat{\mathbf{r}} (\sigma(-\mathbf{w}_j^T \hat{\mathbf{r}}) - 1)$$

We have:

$$\frac{\partial J}{\partial \hat{\mathbf{r}}} = \mathbf{w}_i * (\sigma(\mathbf{w}_i^T \hat{\mathbf{r}}) - 1) + \sum_{k=1}^K \mathbf{w}_k (1 - \sigma(-\mathbf{w}_k^T \hat{\mathbf{r}}))$$

$$\frac{\partial J}{\partial \mathbf{w}_j} = \hat{\mathbf{r}} (\sigma(\mathbf{w}_j^T \hat{\mathbf{r}}) - 1) \text{ when } j = i$$

$$\frac{\partial J}{\partial \mathbf{w}_j} = -\hat{\mathbf{r}} (\sigma(-\mathbf{w}_j^T \hat{\mathbf{r}}) - 1) \text{ when } j \neq i, j \text{ is in } \{1, \dots, K\}$$

$$\frac{\partial J}{\partial \mathbf{w}_j} = 0 \text{ when } j \neq i \text{ and } j \text{ is not in } \{1, \dots, K\}$$

d)

$$J_{\text{skip-gram}} = \sum_{-c \leq j \leq c, j \neq 0} F(v'_{w_{i+j}}, v_{w_i})$$

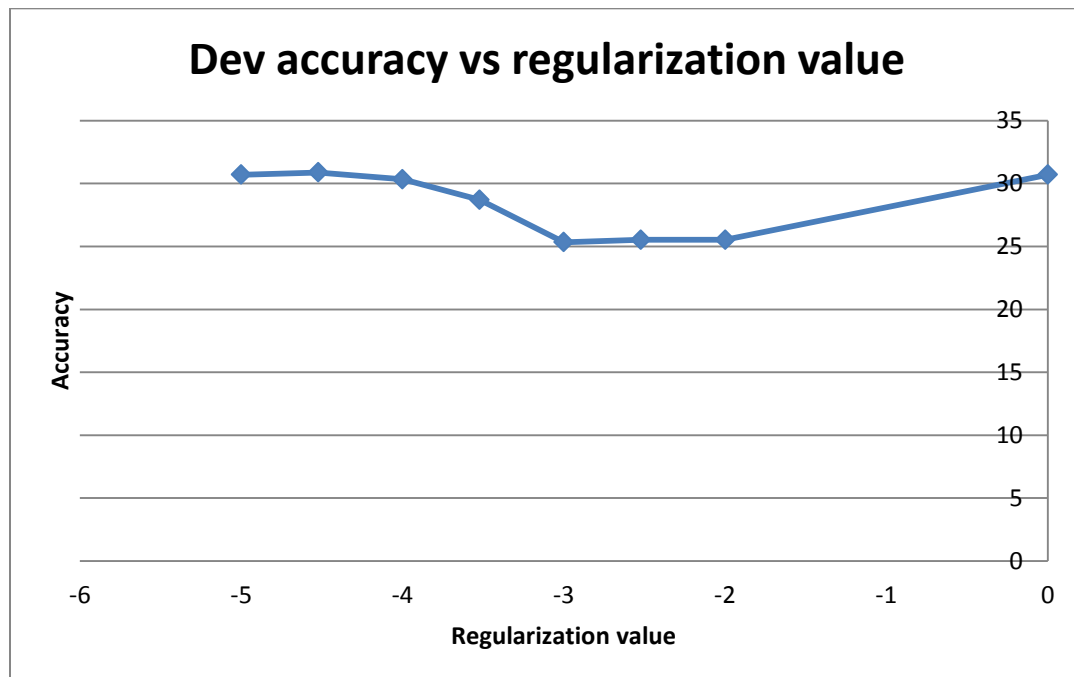
Given the previous parts we know how to compute the gradient of F with respect to any word vector.

$$\frac{\partial J}{\partial v_{w_i}} = \sum_{-c \leq j \leq c, j \neq 0} \frac{\partial F(v'_{w_{i+j}}, v_{w_i})}{\partial v_{w_i}}$$

$$\frac{\partial J}{\partial v'_{w_{i+k}}} = \frac{\partial F(v'_{w_{i+k}}, v_{w_i})}{\partial v'_{w_{i+k}}} \text{ for } -c \leq k \leq c, k \neq 0$$

Problem 4

- a) We want to introduce regularization because it helps to prevent overfitting. Regularization penalizes more highly parameters with larger values, so it forces the parameters to be closer to 0. This reduces model complexity and helps to reduce overfitting.
- b)



We see that as the regularization value increases, the dev accuracy first slightly increases, then decreases significantly, then increases again. (we note that the rightmost point of the curve corresponds to log 0 which we set to be 0 here). We achieve the best dev accuracy with a regularization value of 0.0003, after which the accuracy drops significantly but then again starts to increase, though to values smaller than the best accuracy.