

Assignment 2, Problem 1 (Tensorflow Softmax)

In this question, we will implement a linear classifier with loss function

$$J_{\text{softmax-CE}}(\mathbf{W}) = CE(\mathbf{y}, \text{softmax}(\mathbf{x}\mathbf{W})) \quad (1)$$

Here the rows of \mathbf{x} are feature vectors. We will use `tensorflow`'s automatic differentiation capability to fit this model to provided data.

Solution to Assignment 2, Problem 1 (a)

James Le

Implement the softmax function using `tensorflow` in `q1_softmax.py`. Remember that,

$$\text{softmax}(\mathbf{x})_i = \frac{e^{x_i}}{\sum_j e^{x_j}} \quad (2)$$

NOTE: that you may **NOT** use `tf.nn.softmax` or related build-in functions. You can run basic (non-exhaustive tests) by running: `python q1_softmax.py`.

Implementation in `numpy`:

```
import numpy as np

def softmax(x):
    c = np.max(x, axis=x.ndim - 1, keepdims=True)
    #for numerical stability
    y = np.sum(np.exp(x - c), axis=x.ndim - 1, keepdims=True)
    x = np.exp(x - c)/y
    return x
```

The implementation in `tensorflow`:

```
import numpy as np
import tensorflow as tf

def softmax(x):
    """
```

Compute the softmax function in tensorflow.

You might find the tensorflow functions `tf.exp`, `tf.reduce_max`, `tf.reduce_sum`, `tf.expand_dims` useful. (Many solutions are possible, so you may not need to use all of these functions). Recall also that many common tensorflow operations are sugared (e.g. `x * y` does a tensor multiplication if `x` and `y` are both tensors). Make sure to implement the numerical stability fixes as in the previous homework!

Args:

`x`: `tf.Tensor` with shape `(n_samples, n_features)`. Note feature vectors are represented by row-vectors. (For simplicity, no need to handle 1-d input as in the previous homework)

Returns:

`out`: `tf.Tensor` with shape `(n_sample, n_features)`. You need to construct this tensor in this problem.

"""

YOUR CODE HERE

```
log_c = tf.reduce_max(x,
                      reduction_indices=[len(x.get_shape()) - 1],
                      keep_dims=True)
y      = tf.reduce_sum(tf.exp(x - log_c),
                      reduction_indices=[len(x.get_shape()) - 1],
                      keep_dims=True)
out     = tf.exp(x - log_c) / y
### END YOUR CODE
```

```
return out
```

Solution to Assignment 2, Problem 1 (b)*James Le*

Implement the cross-entropy loss using `tensorflow` in `q1_softmax.py`. Remember that:

$$CE(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^{N_c} y_i \log(\hat{y}_i) \quad (3)$$

where $\mathbf{y} \in \mathbb{R}^5$ is a one-hot label vector and N_c is the number of classes. **Note:** that you may **NOT** use `tensorflow`'s built-in cross-entropy functions for this question. You can run basic (non-exhaustive tests) by running `python q1_softmax.py`.

```
import numpy as np
```

Solution to Assignment 2, Problem 1 (c)*James Le*

Carefully study the `Model` class in `model.py`. Briefly explain the purpose of placeholder variables and feed dictionaries in `tensorflow` computations. Fill in the implementation for the `add_placeholders`, `create_feed_dict` in `q1_classifier.py`.

HINT: that configuration variables are stored in `Config` class. You will need to use these configuration variables in the code.

```
import numpy as np
```

Solution to Assignment 2, Problem 1 (d)

James Le

Implement the transformoin for a softmax classifier in function `add_model` in `q1_classifier.py`. Add cross-entropy loss in function `add_loss_op` in the same file. Use the implementation form earlier pars of the problem, **NOT** `tensorflow` built-ins.

```
import numpy as np
```

Solution to Assignment 2, Problem 1 (e)*James Le*

Fill in the implementation for `add_training_op` in `q1_classifier.py`. Explain how tensorflow's automatic differentiation removes the need for us to define gradients explicitly. Verify that your model is able to fit to synthetic data by running `python q1_classifier.py` and make sure that the tests pass.

HINT: Make sure to use the learning rate specified in `Config`.

```
import numpy as np
```

Assignment 2, Problem 2 (Deep Networks for Named Entity Recognition (NER))

In this section, we'll get to practice backpropagation and training deep networks to attack the task of **Named Entity Recognition (NER)**: predicting whether a given word, in context, represents one of four categories:

- Person (PER)
- Organization (ORG)
- Location (LOC)
- Miscellaneous (MISC)

We formulate this as a 5-class classification problem, using the four above classes and a null-class (0) for words that do not represent a named entity (**NOTE**: most words fall into this category).

The model is a 1-hidden-layer neural network, with an additional representation layer similar to what you saw with `word2vec`. Rather than averaging or sampling, here we explicitly represent context as a 'window' consisting of a word concatenated with its immediate neighbours:

$$\mathbf{x}^{(t)} = [\mathbf{x}_{t-1}\mathbf{L}, \mathbf{x}_t\mathbf{L}, \mathbf{x}_{t+1}\mathbf{L}] \in \mathbb{R}^{3d} \quad (4)$$

where the input \mathbf{x}_{t-1} , \mathbf{x}_t , \mathbf{x}_{t+1} are one-hot row vectors into an embedding matrix $\mathbf{L} \in \mathbb{R}^{|V| \times d}$, which each row \mathbf{L}_i is the vector for a particular word $i = \mathbf{x}_t$. We then compute our prediction as:

$$1 = \quad (5)$$

$$2 = \quad (6)$$

And evaluate by cross-entropy loss,

$$J(\theta) = CE(\mathbf{y}, \hat{\mathbf{y}}) = - \sum_{i=1}^d y_i \log(\hat{y}_i) \quad (7)$$

To compute the loss for the training set, we sum (or average) this $J(\theta)$ as computed with respect to each training example.

For this problem, we let $d = 50$ be the length of our word vectors, which are concatenated into a window of width $3 \times 50 = 150$. The hidden layer has dimensions of 100, and the output layer \hat{y} has dimension of 5.

Solution to Assignment 2, Problem 2 (a)*James Le*

Compute the gradients of $J(\theta)$ with respect to all the model parameters:

$$\frac{\partial J}{\partial \mathbf{U}} \quad \frac{\partial J}{\partial \mathbf{b}_2} \quad \frac{\partial J}{\partial \mathbf{W}} \quad \frac{\partial J}{\partial \mathbf{b}_1} \quad \frac{\partial J}{\partial \mathbf{L}_i} \quad (8)$$

where,

$$\mathbf{U} \in \mathbb{R}^{100 \times 5} \quad \mathbf{b}_2 \in \mathbb{R}^5 \quad \mathbf{W} \in \mathbb{R}^{150 \times 100} \quad \mathbf{b}_1 \in \mathbb{R}^{100} \quad \mathbf{L}_i \in \mathbb{R}^{50} \quad (9)$$

In the spirit of backpropagation, you should express the derivative of activation functions (\tanh , **softmax**) in terms of their function values (as with sigmoid in Assignment 1). This identity may be helpful:

$$\tanh(z) = 2\sigma(2z) - 1 \quad (10)$$

Furthermore, you should express the gradients by using an ‘error vector’ propagated back to each layer; this just amounts to putting parentheses around factors in the chain rule, and will greatly simplify your analysis. All resulting gradients should have simple, closed-form expressions in terms of matrix operations. (**Hint:** you’ve already done most of the work here as part of Assignment 1.)

Solution to Assignment 2, Problem 2 (b)

James Le

To avoid parameters from exploding or becoming highly correlated, it is helpful to augment our cost function with a Gaussian prior: this tends to push parameter weights closer to zero, without constraining their direction, and often leads to classifiers with better generalization ability.

If we maximize **log-likelihood** (as with the cross-entropy loss above), then the Gaussian prior becomes a quadratic term¹ (L2 regularization):

$$J_{\text{reg}}(\theta) = \frac{\lambda}{2} \left[\sum_{i,j} \mathbf{W}_{ij}^2 + \sum_{i',j'} \mathbf{U}_{i'j'}^2 \right] \quad (11)$$

Update your gradients from part (a) to include the additional term in this loss function (i.e. compute $\frac{\partial J_{\text{full}}}{\partial \mathbf{W}}$, etc.).

¹Optional (**not graded**): The interested reader should prove that this is indeed the maximum-likelihood objective when we let $\mathbf{W}_{ij} \equiv \mathcal{N}(0, 1/\lambda)$ for all i, j .

Solution to Assignment 2, Problem 2 (c)

James Le

In order to avoid neurons becoming too correlated and ending up in poor local minima, it is often helpful to randomly initialize parameters. One of the most frequent initializations used is called **Xavier** initialization².

Given a matrix \mathbf{A} of dimensions $m \times n$, select values \mathbf{A}_{ij} uniformly from $[-\epsilon, \epsilon]$, where

$$\epsilon = \frac{\sqrt{6}}{\sqrt{m+n}} \quad (12)$$

Implement the initialization for use in `xavier_weight_init` in `q2_initialization.py` and use it for the weights \mathbf{W} and \mathbf{U} .

²This is also referred to as **Glorot** initialization and was initially described in <http://jmlr.org/proceedings/papers/v9/glorot10a/glorot10a.pdf>.

Solution to Assignment 2, Problem 2 (d)

James Le

In `q2_NER.py` implement the **NER** window model by filling in the appropriate sections. The gradients you derived in (a) and (b) will be computed for you automatically, showing the benefits that automatic differentiation can provide for rapid prototyping.

Run `python q2_NER.py` to evaluate your model's performance on the **dev** set, and compute predictions on the **test** data (make sure turn off debug settings when doing final evaluation). **Note:** the test set has only dummy labels; we'll compare your predictions against ground truth after you submit.

Deliverables:

- Working implementation of the **NER** window model in `q2_NER.py`.
 - In your write-up, *briefly* state the optimal hyperparameters you found for your model: **regularization**, **dimensions**, **learning rate** (including time-varying, such as **annealing**), **SGC batch size**, etc. Report the performance of your model on the **validation** set. You should be able to get a validation loss below **0.2**.
 - List of predicted labels for the **test** set, one per line, in the file `q2 test.predicted`.
 - **HINT:** When debugging, set `max_epochs = 1`. Pass the keyword argument `debug = True` to the call to `load_data` in the `__init__` method.
 - **HINT:** The code should run within 15 minutes on a GPU and 1 hour on a CPU.
-

Assignment 2, Problem 3 (Recurrent Neural Networks: Language Modeling)

In this section, you'll implement your first recurrent neural network (**RNN**) for building a language model.

Language modeling is a central task in NLP, and language models can be found at the heart of speech recognition, machine translation, and many other systems. Given $\mathbf{x}_1, \dots, \mathbf{x}_t$, a language model predicts the following word \mathbf{x}_{t+1} by modeling:

$$P(\mathbf{x}_{t+1} = v_j | \mathbf{x}_1, \dots, \mathbf{x}_t) \quad (13)$$

where v_j is a word in the vocabulary.

Your job is to implement a recurrent neural network language model, which uses feedback information in the hidden layer to model the 'history' x_t, x_{t-1}, \dots, x_1 . Formally, the model³ is, for $t \in \{1, \dots, n-1\}$:

$$\mathbf{e}^{(t)} = \mathbf{x}^{(t)} \mathbf{L} \quad (14)$$

$$\mathbf{h}^{(t)} = \sigma(\mathbf{h}^{(t-1)} \mathbf{H} + \mathbf{e}^{(t)} \mathbf{I} + \mathbf{b}_1) \quad (15)$$

$$\hat{\mathbf{y}}^{(t)} = \text{softmax}(\mathbf{h}^{(t)} \mathbf{U} + \mathbf{b}_2) \quad (16)$$

$$\bar{P}(\mathbf{x}^{(t)} = \mathbf{v}_j | \mathbf{x}_{(t)}, \dots, \mathbf{x}_{(1)}) = \hat{\mathbf{y}}_j^{(t)} \quad (17)$$

$$(18)$$

where $\mathbf{h}^0 = \mathbf{h}_0 \in \mathbb{R}^{D_h}$ is some initialization vector for the hidden layer and $\mathbf{x}^{(t)} \mathbf{L}$ is the product of \mathbf{L} with the one-hot row-vector $\mathbf{x}^{(t)}$ representing index of the current word. The parameters are:

$$\mathbf{L} \in \mathbb{R}^{|V| \times d} \quad (19)$$

where \mathbf{L} is the embedding matrix, \mathbf{I} input word representation matrix, \mathbf{H} the hidden transformation matrix, and \mathbf{U} is the output word representation matrix, \mathbf{b}_1 and \mathbf{b}_2 are biases (d is the embedding dimension, $|V|$ is the vocabulary size, and D_h is the hidden layer dimension).

The output vector $\hat{\mathbf{y}}^{(t)} \in \mathbb{R}^{|V|}$ is a probability distribution over the vocabulary, and we optimize the (unregularized) cross-entropy loss:

$$J(\theta) = CE(\mathbf{y}^{(t)}, \hat{\mathbf{y}}^{(t)}) = - \sum_{i=1}^{|V|} \mathbf{y}_i^{(t)} \log(\hat{\mathbf{y}}_i^{(t)}) \quad (20)$$

³This model is adapted from a paper by Toma Mikolov, et al. from 2010: http://www.fit.vutbr.cz/research/groups/speech/publi/2010/mikolov_interspeech2010_IS100722.pdf

where $\mathbf{y}^{(t)}$ is the one-hot vector corresponding to the target word (which here is equal to x_{t+1}). As in question 2, this is a point-wise loss, and we sum (or average) the cross-entropy loss across all examples in a sequence, across all sequences⁴ in the dataset in order to evaluate model performance.

⁴We use the tensor function `sequence_loss` to do this.

Solution to Assignment 2, Problem 3 (a)

James Le

Conventionally, when reporting performance of a language model, we evaluate on **perplexity**, which is defined as:

$$PP \tag{21}$$

i.e. the inverse probability of the correct word, according to the model distribution \bar{P} . Show how you can derive perplexity from the cross-entropy loss (**HINT**: remember that $y^{(t)?}$ is one-hot!), and thus argue that minimizing the (arithmetic) mean cross-entropy loss will also minimize the (geometric) mean perplexity across the training set.⁵

For a vocabulary of $|V|$ words, what would you expect perplexity to be if your model predictions were completely random? Compute the corresponding cross-entropy loss $|V| = 2000$ and $|V| = 10,000$ and keep this in mind as a baseline.

⁵**NOTE**: this should be a very short problem - not too perplexing!

Solution to Assignment 2, Problem 3 (b)

James Le

As you did in question 2, compute the gradients with respect to all the model parameters at a single point in time t :

$$\frac{\partial J(t)}{\partial \mathbf{U}} \quad (22)$$

where $\mathbf{L}_{\mathbf{x}^{(t)}}$ is the column of \mathbf{L} corresponding to the current word $\mathbf{x}^{(t)}$, and $\big|_{(t)}$ denotes the gradient for the appearance of that parameter at time t . Equivalently, $\mathbf{h}^{(t-1)}$ is taken to be fixed, and you need not backpropagate to earlier timesteps just yet - you'll do this in 3 (c).

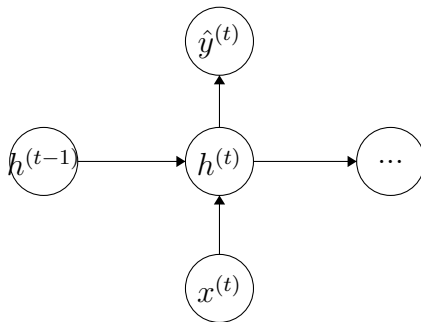
Additionally compute the derivative with respect to the **previous** hidden layer value:

$$\frac{\partial J(t)}{\partial \mathbf{h}^{(t-1)}} \quad (23)$$

Solution to Assignment 2, Problem 3 (c)

James Le

Below is a sketch of the network at a single timestep:



Draw the ‘unrolled’ network for 3 timesteps, and compute the backpropagation through time gradients:

$$\frac{\partial J(t)}{\partial \mathbf{L}_{\mathbf{x}(t-1)}} \quad (24)$$

where $|_{(t-1)}$ denotes the gradient for the appearance of that parameter at time $(t-1)$. Because parameters are used multiple times in feed-forward computation, we need to compute the gradient for each time they appear.

You should use the backpropagation rules from Lecture 5⁶ to express these derivatives in terms of error term

$$\delta^{(t-1)} = \frac{\partial}{\partial} \quad (25)$$

computed in the previous part. Doing so will allow for re-use of the expressions for $(t-2)$, $(t-3)$, and so on.

NOTE: that the true gradient with respect to a training example requires us to run backpropagation all the way back to $t = 0$. In practice, however, we generally truncate this and only backpropagate for a fixed number of $\tau \approx 3 - 5$ timesteps.

⁶<http://cs224d.stanford.edu/lectures/CS224d-Lecture5.pdf>

Solution to Assignment 2, Problem 3 (d)

James Le

Given $\mathbf{h}^{(t-1)}$, how many operations are required to perform one step of forward propagation to compute $J^{(t)}(\theta)$? How about backpropagation for a single step in time? For τ steps in time? Express your answer in big-O notation in terms of the dimensions d , D_h and $|V|$ (Equation 19). What is the slow step?

Solution to Assignment 2, Problem 3 (e)

James Le

Implement the above model in `q3.RNNLM.py`. Data loaders and other starter code are provided. Follow the directions in the code to understand which parts need to be filled in. Running `python q3.RNNLM.py` will run the model. Note that you may **NOT** use built-in tensorflow functions such as those in the `rnn.cell` module.

Train a model on the `ptb-train` data, consisting of the first 20 sections of the WSJ corpus of the Penn Treebank. As in question 2, you should tune your model to maximize generalization performance (minimize cross-entropy loss) on the `dev` set. We'll evaluate your model on an unseen but similar set of sentences.

Deliverables:

- In your writeup, include the best hyperparameters you used (train schedule, number of iterations, learning rate, backprop timesteps), and your perplexity score on the `ptb-dev` set. You should be able to get validation perplexity below 175.
 - Include your saved model parameters for your best model; we'll use these to test your model.
 - **HINT:**When debugging, set `max_epochs = 1`. Pass the keyword argument `debug = True` to the call to `load_data` in the `__init__` method.
 - **HINT:**On a GPU this code should run quickly (below 30 minutes). On a CPU the code may take up to 4 hours to run.
-

Solution to Assignment 2, Problem 3 (f)

James Le

The networks that you've seen in Assignment 1 and in question 2 of this assignment are discriminative models: they take data, and make a prediction. The **RNNLM** model you've just implemented is a **generative** model, in that it actually models the distribution of the **data** sequence x_1, \dots, x_n . This means that not only can we use it to evaluate the likelihood of a sentence, but we can actually use it to generate one!

After training, in `q3_RNNLM.py`, implement the `generate_text()` function. This should run the RNN forward in time, beginning with the index for the start token `<eos>`, and sampling a new word x_{t+1} from the distribution \hat{y}_{t+1} at each timestep. Then feed this word in as input at the next step, and repeat until the model emits an end token (index of `</eos>`).

Deliverables:

- Include 2-3 generated sentences in your write up. See if you can generate something humorous!

COMPLETELY OPTIONAL, (not graded): If you want to experiment further with language models you're welcome to load up your own texts and train on them - sometimes the results can be entertaining. See <http://kingjamesprogramming.tumblr.com/> for a great one⁷ trained on a mix of the King James Bible and the Structure and Interpretation of Computer Programs.

⁷This one just uses a simple **n-gram Markov model**, but there's no reason an RNNLM can't compete!