

DCMMC

A NOTE IN MACHINE LEARNING

TABLE OF CONTENTS

| | | |
|-------|--|----|
| | <i>About this Note</i> | 6 |
| 1 | <i>NLP with DL</i> | 7 |
| 1.1 | <i>Word Vector</i> | 7 |
| 1.1.1 | <i>Word2vec</i> | 7 |
| 1.1.2 | <i>HW1</i> | 10 |
| 1.1.3 | <i>GloVe</i> | 11 |
| 1.1.4 | <i>Word sense ambiguity</i> | 12 |
| 1.2 | <i>Math Backgrounds</i> | 12 |
| 1.2.1 | <i>Data Preprocessing</i> | 14 |
| 1.2.2 | <i>Parameter Initialization</i> | 14 |
| 1.2.3 | <i>Optimizer</i> | 15 |
| 1.2.4 | <i>Regularization</i> | 16 |
| 1.2.5 | <i>Practice: Named Entity Recognition</i> | 17 |
| 1.2.6 | <i>HW2</i> | 17 |
| 1.3 | <i>Dependency Parser</i> | 18 |
| 1.3.1 | <i>Neural Dependency Parsing</i> | 19 |
| 1.3.2 | <i>HW3</i> | 19 |
| 1.4 | <i>Language Modeling and Recurrent Neural Networks</i> | 21 |
| 1.5 | <i>Seq2Seq and Attention</i> | 24 |
| 1.5.1 | <i>HW4</i> | 27 |
| 1.5.2 | <i>Tips for Research</i> | 29 |

| | | |
|-------|--|----|
| 1.6 | <i>Contextual Word Representations and Pretraining</i> | 29 |
| 1.7 | <i>Attention Model and Question Answering</i> | 30 |
| 1.8 | <i>ConvNets for NLP</i> | 31 |
| 1.9 | <i>Subword models</i> | 34 |
| 1.9.1 | <i>Byte Pair Encoding</i> | 34 |
| 1.9.2 | <i>FastText embeddings</i> | 34 |
| 1.9.3 | <i>HW5</i> | 34 |
| 2 | <i>Bibliography</i> | 36 |

LIST OF FIGURES

| | | |
|------|---|----|
| 1.1 | An example of word analogy of man:woman :: king:? | 7 |
| 1.2 | A demo of the window size and $p(w_o w_c)$ | 8 |
| 1.3 | An example of co-occurrence matrix with window size of 1 | 11 |
| 1.4 | An example of the conditional probabilities and their ratio in GloVe paper. | 11 |
| 1.5 | An exmaple of mean subtration. | 14 |
| 1.6 | An example of normalization. | 14 |
| 1.7 | A picture of momentum. | 15 |
| 1.8 | Principle of RNN | 21 |
| 1.9 | The repeating module in an LSTM contains four interacting layers. | 23 |
| 1.10 | Residual connections. | 24 |
| 1.11 | Dense Net. | 24 |
| 1.12 | Alignment from french to english translation. | 24 |
| 1.13 | Training phase for NMT. | 25 |
| 1.14 | Seq2Seq with attention. | 26 |
| 1.15 | Architecture of BiDAF | 31 |
| 1.16 | Simple ConvNet for NLP classification. | 32 |

LIST OF TABLES

ABOUT THIS NOTE

Note that most of symbols in this note are vector, matrix, or tensor. Strictly speaking, we should write them as bold to differ from scalars. But for simplification, most bolds of them are ignored in this note.

Also, \times in superscript will leads into overflow in this latex source code. Therefore, all \times are replaced by $*$ in superscripts.

When there is no possibility for confusion, we write the probability $Pr(W = w)$ where W is the random variable and w is the specific value to the shorthands $P(w)$.

TODO(DCMMC)...

Learning Objectives:

- Word Vector
- Calculus Review
- RNN & Language Model
- Seq2Seq & Attention
- ConvNet for NLP
- Transformer

1.1 Word Vector

Arguably the most simple word vector, i.e., **one-hot vector**: an $\mathbb{R}^{|V| \times 1}$ vector with one 1 and the rest 0s. Note that these one-hot vectors are **orthogonal** (i.e., no similarity/relationship) and V is a very big vocabulary ($\sim 500k$ words for english).

Another idea: **distributional representation** in modern statistical NLP. A word's meaning is given by the words that frequently appear close-by. Using some N -dim ($N \ll |V|$) space is sufficient to encode all semantics of our language into a dense vector. Once we get the word embedding matrix where each column is a word vector, we can query the word vector from one-hot representation by treating it as **lookup table** instead of using matrix product.

To evaluate word vectors, there are two fold: *intrinsic* (directly used, e.g. word analogies/similarity) and *extrinsic* (indirectly used in real task, e.g. Q&A). Word vector analogies for $a : b :: c : d$ is calculated by cosine similarity as example shown in Fig. 1.1:

$$d = \arg \max_i \frac{(x_b - x_a + x_c)^\top x_i}{\|x_b - x_a + x_c\|} \quad (1.1)$$

If we have hundreds of millions of words, it's okay to start the vectors *randomly*. If there is a *small* ($\leq 100,000$) training data set, it's best to just treat the pre-trained word vectors as *fixed*. In the other hand, if there is a large dataset, then we can gain by **fine tuning** of the word vectors.

1.1.1 Word2vec

Two families of models: **Skip-gram** and **Continuous Bag of Words**.

Dependencies: Machine Learning Basic

In traditional NLP (before 2013), words are regarded as discrete symbols (**localist** representation) and cannot capture similarity. One-hot vector is an example.

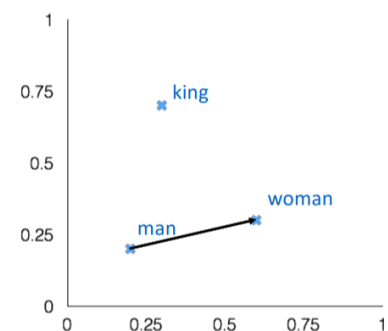


Figure 1.1: An example of word analogy of man:woman :: king:?

Idea of **Skip-gram** (predicting context words by a given center word) in Word2vec¹:

¹ Mikolov et al. 2013

- a large corpus of text T with a vocabulary V
- every word is represented by a vector $w \in \mathbb{R}^d$ and start off as a random vector
- use the (cosine) similarity of the word vectors for c (center word) and o (context/outside word) to calculate the probability of o given c : $p(w_o|w_c)$
- adjusting the word vectors to maximize the probability

The conditional probability is calculated by the **softmax** (normalize to probability distribution) of **cosine** similarity (review dot product: $a \cdot b = |a||b| \cos \langle a, b \rangle$). Note that the visualization of word vectors utilizes 2D projection (e.g. PCA) that will lose huge information.

Why we use two vectors per word? Make it simpler to calculate the gradient of loss function. Because the center word would be one of the choices for the context word and thus squared terms are imported. Average both vectors at the end is the final word vector.

$$p(w_o|w_c) = \frac{\exp(u_o^\top v_c)}{\sum_{w \in V} (u_w^\top v_c)} \quad (1.2)$$

where v_c denotes the center word vector of w when w is used as a center word in the formula, and u_w denotes the context word vector of w as the similar way. A demo of the window size and conditional probability is shown in Fig. 1.2.

The objective function (a.k.a loss or cost function) is given by the (average) negative log likelihood (abbr. **NLL**). The parameters of the model are adjusted by minimizing the loss function $J(\theta)$ or maximizing the likelihood. This is, give a high probability estimate to those words that occur in the context and low probability to those don't typically occur in the context.

$$\begin{aligned} \arg \max_{\theta} L(\theta) &= \prod_{c=1}^T p(w_{c-m}, \dots, w_{c-1}, w_{c+1}, \dots, w_{c+m} | w_c; \theta) \\ &= \prod_{c=1}^T \prod_{\substack{-m \leq j \leq m \\ o=j+c \\ o \neq c}} p(w_o | w_c; \theta) \\ &\Downarrow \\ \arg \min_{\theta} -\frac{1}{T} \log L(\theta) &= -\frac{1}{T} \sum_{c=1}^T \sum_{\substack{-m \leq j \leq m \\ o=j+c \\ o \neq c}} \log p(w_o | w_c; \theta) \end{aligned}$$

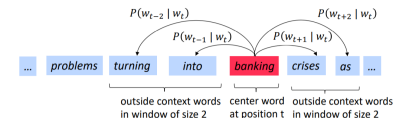


Figure 1.2: A demo of the window size and $p(w_o|w_c)$

$$= -\frac{1}{T} \sum_{c=1}^T \sum_{\substack{-m \leq j \leq m \\ o=j+c \\ o \neq c}} \left(u_o^\top v_c - \log \sum_{w \in V} \exp(u_w^\top v_c) \right) \quad (1.3)$$

where m is the window size, $\theta \in \mathbb{R}^{2d|V|}$ represents all model parameters. And we assume that $p(\cdot|w_c)$ are **i.i.d.**

We use **gradient descent** (i.e. averaged gradient of all samples/windows) to optimize the loss function. Note that stochastic (one sample/window with noisy estimates of the gradients) or mini-batch (a subset of samples/windows with size powered of 2 such as 64) gradient descent methods are useful to prevent over-fitting and train for large dataset. Calculating the gradient of the loss function is trivial:

The properties of log and arg max (arg min) used in Eq. 1.3 are VERY useful. $\exp(\cdot)$ ensures anything positive.

$$\begin{aligned} \frac{\partial J}{\partial v_c} &= -\frac{1}{T} \sum_{c=1}^T \sum_{\substack{-m \leq j \leq m \\ o=j+c \\ o \neq c}} \left(u_o - \sum_{x \in V} \frac{\exp(u_x^\top v_c) u_x}{\sum_w \exp(u_w^\top v_c)} \right) \\ &= -\frac{1}{T} \sum_{c=1}^T \sum_{\substack{-m \leq j \leq m \\ o=j+c \\ o \neq c}} \left(u_o - \sum_{x \in V} p(w_x|w_c) \cdot u_x \right) \end{aligned} \quad (1.4)$$

$$\frac{\partial J}{\partial u_o} = -\frac{1}{T} \sum_{c=1}^T \sum_{\substack{-m \leq j \leq m \\ o=j+c \\ o \neq c}} (v_c - p(w_o|w_c)) \quad (1.5)$$

Iteratively update equation (naïve version) is given by:

$$\theta^{new} = \theta^{old} - \alpha \nabla_{\theta} J(\theta) \quad (1.6)$$

where α is the learning size (step size) such as 10^{-3} .

Note that the summation over $|V|$ ($\sum_{x \in V}$) is very expensive to compute! For every training step, instead of looping over the entire vocabulary, we can just sample several negative examples!

negative sampling: train binary logistic regression instead.

$p(D = 1|w_o, w_c)$ denotes the probability when (w_o, w_c) came from the same window of the corpus data, and $p(D = 0|w_o, \tilde{w}_o)$ is the probability given (w_o, \tilde{w}_o) did not come from the same window (i.e. noisy/invalid pair). Randomly sample a bunch of noise words from the **unigram distribution** raised to the power of 3/4: $p(w) = U(w)^{3/4}/Z$, where $U(w)$ is the counts for every unique words (i.e. unigram) and Z is the normalization term.

To avoid high frequency effect of words such as **of** and **the**, one simple way is just lop off the first biggest component in the word vector. The unigram with power of 3/4 in word2vec is also a trick to handle the effect, where it decrease how often you sample very common words and increase how often you sample rare words.

The objective function is also come from NLL:

$$J(\theta) = -\frac{1}{T} \sum_{c=1}^T \sum_{\substack{-m \leq j \leq m \\ o=j+c \\ o \neq c}} \left(\log \sigma(u_o^\top v_c) + \sum_{j \sim p(w)} \left[\log \sigma(-u_j^\top v_c) \right] \right) \quad (1.7)$$

where **sigmoid** function is $\sigma(x) = \frac{1}{1+e^{-x}}$ which can be seen as the 1D (binary) version of softmax and used to output the probability, and k is the number of negative samples such as 5 and 15. Note that according to the symmetric property of sigmoid function we get: $P(D=0|\tilde{w}_j, w_c) = 1 - P(D=1|\tilde{w}_j, w_c) = \sigma(-u_j^\top v_c)$.

Continuous Bag of Words (CBOW): predict center word from (bag of) context words. Similar to Skip-gram, the objective function is formulated as:

Although word2vec model is fairly simple and clean, there are actually many tricks which aren't particularly theoretical.

$$J = -\frac{1}{T} \sum_{c=1}^T \log P(w_c | w_{c-m}, \dots, w_{c-1}, w_{c+1}, \dots, w_{c+m}) \quad (1.8)$$

$$= -\frac{1}{T} \sum_{c=1}^T \log p(v_c | \hat{u}) \quad (1.9)$$

$$= -\frac{1}{T} \sum_{c=1}^T \log \text{softmax}_c(v_c^\top \hat{u}) \quad (1.10)$$

$$= -\frac{1}{T} \sum_{c=1}^T (v_c^\top \hat{u} - \log \sum_{j=1}^{|V|} \exp(v_j^\top \hat{u})) \quad (1.11)$$

where $\hat{u} = \frac{1}{2m} \sum_{\substack{-m \leq j \leq m \\ o=j+c \\ o \neq c}} u_o$

Although word2vec can capture complex patterns beyond word similarity, it has inefficient usage of statistics (i.e. rely on sampling rather than directly use counts of words).

1.1.2 HW1

A simple intro to co-occurrence matrix, SVD, cosine similarity, and some applications (e.g. word analogy) of word2vec.

1.1.3 GloVe

Co-occurrence matrix $X \in \mathbb{R}^{|V| \times |V|}$ with window size k . Fig. 1.3 shows an example. Note that such matrix is extremely sparse and very high dimensional, and the dimensions of the matrix change very often as new words are added very frequently and corpus changes in size. We can perform SVD on X to reduce the dimensionality to $25 \sim 1000$ -dim. In addition, there are some hacks to X that transform the raw count introduced by ²: (1) set upper bound (e.g. 100) or just ignore them all for the counts of too frequent words, (2) ramped windows that count closer words more. (3) use Pearson correlations instead of counts. Note that they made some interesting observation in their word vector that the verb (e.g. swim) and the corresponding doer (e.g. swimmer) pairs are roughly *linear components* (e.g. $v_{swimmer} - v_{swim} = k(v_{driver} - v_{drive})$).

TODO(DCMC)...SVD

Although the aforementioned conventional method has disproportionate importance given to large counts and mainly only capture word similarity, it enjoys the fast training and efficient usage of statistics. GloVe (**G**lobal **V**ector) ³ combines the advantages from both of this conventional method (global count matrix factorization) and the DL-based methods (local context window methods) such as word2vec. It captures global corpus statistics directly.

Some notations: X_{ij} tabulate the number of times word j occurs in the context of word i , $X_i = \sum_k X_{ik}$ is the number of times any word appears in the context of word i i.e., the normalization denominator. $P_{ij} = P(j|i) = X_{ij}/X_i$ is the probability that word j appear in the context of word i . The crucial insight is that the *ratios* of co-occurrence probabilities as shown in Fig. 1.4 to encode meaning components. We'd like to leverage the word vectors w_i, w_j, \tilde{w}_k to represent such ratio: $F(w_i, w_j, \tilde{w}_k) = P_{ik}/P_{jk}$, where \tilde{w} is a separate *context* word vector for various *probe words* k , instead of the word vector w (similar to center word vector in skip-gram).

We can select a unique choice of F by enforcing a few desiderata (i.e. restrictions). To fit the demand of the *linear components* and the output *scalar* value, in addition to the *homomorphism* between the groups $(\mathbb{R}, -)$ and (\mathbb{R}^+, \div) (i.e., $F(i, j) = P_{ik}/P_{jk} = 1/F(j, i) = P_{jk}/P_{ik}$), we can derivate that $F(w_i, w_j, \tilde{w}_k) = F((w_i - w_j)^\top \tilde{w}_k) = F(w_i^\top \tilde{w}_k)/F(w_j^\top \tilde{w}_k) = P_{ik}/P_{jk}$. Therefore, $F = \exp, w_i^\top \tilde{w}_k = \log(P_{ik}) = \log(X_{ik}) - \log(X_i)$. Note that the symmetry property of co-occurrence: $X_{ik} = X_{ki}$. We add two biases to restore the symmetry: $w_i^\top \tilde{w}_k + b_i + \tilde{b}_k = \log(X_{ik})$, where we can analogy that $b_i + \tilde{b}_j = \log X_i$.

More details, the relationship to the "global skip-gram" and the

1. I enjoy flying.
2. I like NLP.
3. I like deep learning.

The resulting counts matrix will then be:

| | I | like | enjoy | deep | learning | NLP | flying | . |
|----------|---|------|-------|------|----------|-----|--------|---|
| I | 0 | 2 | 1 | 0 | 0 | 0 | 0 | 0 |
| like | 2 | 0 | 0 | 1 | 0 | 1 | 0 | 0 |
| enjoy | 1 | 0 | 0 | 0 | 0 | 0 | 1 | 0 |
| deep | 0 | 1 | 0 | 0 | 1 | 0 | 0 | 0 |
| learning | 0 | 0 | 0 | 1 | 0 | 0 | 0 | 1 |
| NLP | 0 | 1 | 0 | 0 | 0 | 0 | 0 | 1 |
| flying | 0 | 0 | 1 | 0 | 0 | 0 | 0 | 1 |
| . | 0 | 0 | 0 | 0 | 1 | 1 | 1 | 0 |

Figure 1.3: An example of co-occurrence matrix with window size of 1

² Rohde et al. 2005

³ Pennington et al. 2014

| Probability and Ratio | $k = solid$ | $k = gas$ | $k = water$ | $k = fashion$ |
|-----------------------|----------------------|----------------------|----------------------|----------------------|
| $P(k ice)$ | 1.9×10^{-4} | 6.6×10^{-5} | 3.0×10^{-3} | 1.7×10^{-3} |
| $P(k steam)$ | 2.2×10^{-5} | 7.8×10^{-4} | 2.2×10^{-3} | 1.8×10^{-5} |
| $P(k ice)/P(k steam)$ | 8.9 | 8.5×10^{-2} | 1.36 | 0.96 |

Figure 1.4: An example of the conditional probabilities and their ratio in GloVe paper.

complexity refer to the original GloVe paper ⁴.

$$w_i \cdot w_j = \log P(i|j) \quad (1.12)$$

$$w_x \cdot (w_a - w_b) = \log \frac{P(x|a)}{P(x|b)} \quad (1.13)$$

⁴ Pennington et al. 2014

To handle the ill-defined log function when its argument be 0 (its common that $X_{ij} = 0$), the authors use the factorized log: $\log(X_{ik}) \rightarrow \log(1 + X_{ik})$.

Therefore, the ratios of co-occurrence probabilities is the **log-bilinear model with vector differences**. The final objective function is *weighted least squares* (MSE) for this regression problem.

$$J = \sum_{i,j=1}^V f(X_{ij}) \left(w_i^\top \tilde{w}_j + b_i + \tilde{b}_j - \log X_{ij} \right)^2 \quad (1.14)$$

where weighted function (is also a hyperparamter) is:

$$f(x) = \begin{cases} \left(\frac{x}{x_{max}} \right)^\alpha & \text{if } x < x_{max} \\ 1 & \text{otherwise} \end{cases} \quad (1.15)$$

where $x_{max} = 100, \alpha = 3/4$ (*empirical value*).

1.1.4 Word sense ambiguity

Because most words have lots of meanings. One crude way ⁵ is to cluster word windows around words, retrain with each word assigned to multiple different clusters bank₁, bank₂, etc. Another method ⁶ is weighted sum of different senses of a word reside in a linear superposition, e.g.:

⁵ Huang et al. 2012

⁶ Arora et al. 2018

$$v_{\text{pike}} = \alpha_1 v_{\text{pike}_1} + \alpha_2 v_{\text{pike}_2} + \alpha_3 v_{\text{pike}_3} \quad (1.16)$$

where $\alpha_i = \frac{f_i}{\sum_{j=1}^3 f_j}$ for frequency f .

The result is counterintuitive very well, because of the idea from *sparse* coding you can actually separate out the senses.

1.2 Math Backgrounds

For **multi-class classification** problem, **NLL** (negative likelihood loss) is the objective function of **Maximum Likelihood Estimate** (abbr, MLE):

$$J(\theta) = - \sum_i \log p(y = y_i^{true} | x_i; \theta) \quad (1.17)$$

cross entropy (distance measure) between (discrete) distribution p and q is more convenient way:

$$H(p, q) = - \sum_{c=1}^C p(c) \log q(c) \quad (1.18)$$

However, in the multi-class (with single label) setting, the $p(c)$ is the **ground truth distribution** which has the *one-hot* style (**empirical distribution**), i.e. $p = [0, \dots, 0, 1, 0, \dots, 0]$ where 1 at the right class and 0 everywhere else. Therefore, the **cross entropy** in the multi-class classification is *equal* to the NLL.

A simple k -class model example is **dense layer** with *softmax*:

$$p(y|x; \theta) = \text{softmax}(W_2 f(W_1 x + b)) \quad (1.19)$$

where $\theta = [W_1, b, W_2]^\top$ are the parameters, $x \in \mathbb{R}^m$, $W_1 \in \mathbb{R}^{n \times m}$, $b \in \mathbb{R}^n$, $W_2 \in \mathbb{R}^{k \times n}$, $f(\cdot)$ is a kind of simple activate (non-linear) function to provide non-linearity, such as $\text{ReLU}(x) = \max(0, x)$. The visualization of neural network refer to ⁷.

The **Jacobian Matrix** (generalization of the gradient) of function $f(x) : \mathbb{R}^n \rightarrow \mathbb{R}^m$ is a $m \times n$ matrix: $\left(\frac{\partial f}{\partial x}\right)_{ij} = \frac{\partial f_i}{\partial x_j}$.

Supposed that we have a function $g(f(x))$, $f : \mathbb{R} \rightarrow \mathbb{R}^2$, $g : \mathbb{R}^2 \rightarrow \mathbb{R}^2$, we can compute the partial derivative of g w.r.t x by **chain rule**:

$$\frac{\partial g}{\partial x} = \begin{bmatrix} \frac{\partial g_1}{\partial f_1} \frac{\partial f_1}{\partial x} + \frac{\partial g_1}{\partial f_2} \frac{\partial f_2}{\partial x} \\ \frac{\partial g_2}{\partial f_1} \frac{\partial f_1}{\partial x} + \frac{\partial g_2}{\partial f_2} \frac{\partial f_2}{\partial x} \end{bmatrix} \quad (1.20)$$

It is the same as multiplying the two Jacobians:

$$\frac{\partial g}{\partial x} = \frac{\partial g}{\partial f} \frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial g_1}{\partial f_1} & \frac{\partial g_1}{\partial f_2} \\ \frac{\partial g_2}{\partial f_1} & \frac{\partial g_2}{\partial f_2} \end{bmatrix} \begin{bmatrix} \frac{\partial f_1}{\partial x} \\ \frac{\partial f_2}{\partial x} \end{bmatrix} \quad (1.21)$$

There are some useful identities:

- $\frac{\partial x}{\partial x} = I$
- $\frac{\partial Wx}{\partial x} = W$, $\frac{\partial u^\top x}{\partial x} = u^\top$
- $\frac{\partial x^\top W}{\partial x} = W^\top$
- For elemenwise function $f(x)$: $\frac{\partial f}{\partial x} = \text{diag}(f'(x))$
- $\frac{\partial \theta^\top (W \cdot h)}{\partial W} = \theta h^\top$ where $\theta \in \mathbb{R}^{D_\theta \times 1}$, $W \in \mathbb{R}^{D_\theta \times D_h}$, $h \in \mathbb{R}^{D_h \times 1}$

⁷ ConvNetJS: <https://cs.stanford.edu/people/karpathy/convnetjs/demo/classify2d.html>

$\frac{dg_1}{dy} = \frac{\partial g_1}{\partial y_1} + \frac{\partial g_1}{\partial y_2}$ is the relationship of the full differential and the partial differential.

- For cross entropy loss: $J(\mathbf{h}) = -\mathbf{y}^\top \log(\hat{\mathbf{y}}) = -\mathbf{y}^\top \log \text{softmax}(\mathbf{h})$ (\mathbf{y} is one-hot vector) is: $\frac{\partial J}{\partial \mathbf{h}} = (\hat{\mathbf{y}} - \mathbf{y})^\top$

We can use **backward propagation** (reversed of the *topological sort*) and *re-use* intermediate nodes to reduce complexity in the *computation graph*.

Other machine learning basic concepts are: **regularization** (e.g. L2) to prevent **overfitting**, vectorization to parallelization, (non-linear) **activation function** (e.g. sigmoid, tanh, (leaky) ReLU), parameter initialization (e.g. Xavier), **Optimizer** (e.g. RMSprop, Adam), learning rate.

1.2.1 Data Preprocessing

- Mean Subtraction (Shifting)**: Shifting all data so that they have zero mean as shown in Fig. 1.5. Formally, $\mathbf{x}^{(i)} \leftarrow \mathbf{x}^{(i)} - \mathbb{E}[\mathbf{x}^{(i)}]$ for sample i , where \mathbb{E} indicates mean.
- Normalization (Scaling)**: Scale every input feature dimension to have similar ranges of magnitudes, as shown in Fig. 1.6. This is useful since input features are often measured in different "units", but we often want to initially consider all features as equally important. Formally, $x_j^{(i)} \leftarrow \frac{x_j^{(i)}}{\sigma_i(x_j^{(i)})}$ for feature j in sample i where $\sigma(\cdot)$ is the standard variance over $x_j^{(0)}, \dots, x_j^{(N)}$.
- Whitening**: Whitening converts the data to have an identity covariance matrix - that is, features become uncorrelated and have a variance of 1. In the specific, we can divide the principal components achieved from PCA by the square roots of their eigenvalues (singular value).

Compared with activations such as sigmoid and tanh, ReLU does not *saturate* even for larger values. Note that ReLU is not *differentiable* at 0, where we can use *sub-derivatives* in implementation with a certain value such as 0 or 1.

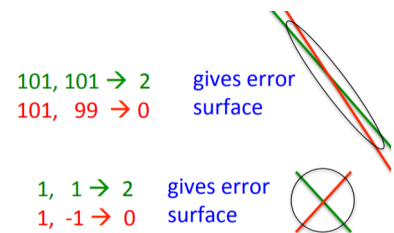


Figure 1.5: An example of mean subtraction.

1.2.2 Parameter Initialization

If two hidden units have exactly the same bias and exactly the same incoming and outgoing weights (e.g. different channels for learning different features in the same convolutional layer), they will always get exactly the same gradient. So we should initialize the weights to small random values.

A good starting strategy is to initialize the weights to small random numbers of **normal distribution** with the mean around 0. Xavier et al. ⁸ suggest that for sigmoid and tanh activation units, it's better for the weight matrix $W \in \mathbb{R}^{n^{(l+1)} \times n^{(l)}}$ to be initialized

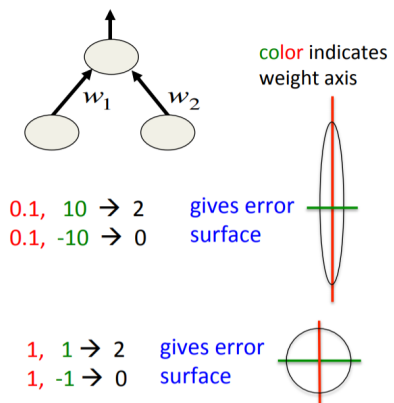


Figure 1.6: An example of normalization.

⁸ Glorot and Bengio 2010

randomly with a **uniform distribution**:

$$W \sim U \left[-\sqrt{\frac{6}{n^{(l)} + n^{(l+1)}}}, \sqrt{\frac{6}{n^{(l)} + n^{(l+1)}}} \right] \quad (1.22)$$

where $n^{(l)}, n^{(l+1)}$ are also called **fan-in** and **fan-out**. Besides, bias units are initialized to 0.

1.2.3 Optimizer

To avoid a diverging loss (too large learning step) or unconverging (too small learning step), there are some learning strategies.

Annealing: start off with a high learning rate to approach a minimum quickly, after several iterations, the learning rate is reduced in some way under a more fine-grained scope.

- Exponential decay: $\alpha(t) = \alpha_0 e^{-kt}$ where α_0 is initial learning rate.
- Decrease over time: $\alpha(t) = (\alpha_0 \tau) / \max(t, \tau)$ where τ denotes the time at which the learning rate should start reducing.

Momentum (a picture of it can be seen in Fig.1.7) based methods:

- AdaGrad: $m \leftarrow m + (\nabla_{\theta} J(\theta))^2, \theta \leftarrow \theta - \alpha \nabla_{\theta} J(\theta) \odot (\sqrt{m} + \epsilon)^{-1}$ where $\odot, (\cdot)^{-1}, (\cdot)^2, \sqrt{\cdot}$ are all element-wise operators, and ϵ is a very small value such as 10^{-8} to prevent **arithmetic underflow**. It leads to that parameters that have not been updated much in the past are likelier to have higher learning rates now.
- RMSprop: $m \leftarrow \beta m + (1 - \beta) (\nabla_{\theta} J(\theta))^2, \theta \leftarrow \theta - \alpha \nabla_{\theta} J(\theta) \odot (\sqrt{m} + \epsilon)^{-1}$ where β is the decay rate with default value 0.9. Unlike AdaGrad, its updates do not become monotonically smaller.
- Adam⁹: $m \leftarrow \beta_1 m + (1 - \beta_1) \nabla_{\theta} J(\theta), v \leftarrow \beta_2 v + (1 - \beta_2) (\nabla_{\theta} J(\theta))^2, \hat{m} = m / (1 - \beta_1^t), \hat{v} = v / (1 - \beta_2^t), \theta \leftarrow \theta - \alpha \hat{m} / (\sqrt{\hat{v}} + \epsilon)$, where $/$ is also a element-wise operator, hyperparameters $\beta_1 = 0.9, \beta_2 = 0.999 \in [0, 1)$. \hat{m}, \hat{v} are the bias-corrected m, v , and they indicate a rolling average of the gradients and a rolling average of the magnitudes of the gradients, respectively. In addition, m, v are all initialized to 0 Adam is like a combination of RMSProp and momentum.

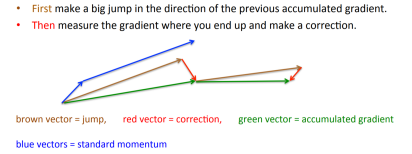


Figure 1.7: A picture of momentum.

For the implementation of momentum such as RMSprop, there is an interesting small trick: use $m = m - (1 - \beta)(m - (\nabla_{\theta} J(\theta))^2)$ instead of $m = \beta m + (1 - \beta) (\nabla_{\theta} J(\theta))^2$. In such way, we need calculate only one multiplication, compared with original two multiplications.

1.2.4 Regularization

1. Dropout

During training, **dropout**¹⁰ randomly disables units in the hidden layer by a mask vector drawn from Bernoulli distribution where each entry is 0 with probability p_{drop} and 1 with probability $(1 - p_{\text{drop}})$:

$$\text{Dropout Layer: } d_i \sim \text{Bernoulli}(1 - p_{\text{drop}}), \hat{h}^{(t)} = \frac{1}{1 - p_{\text{drop}}} d \odot h^{(t)} \quad (1.23)$$

¹⁰ Srivastava et al. 2014

where \odot is element-wise product, $d \in \{0, 1\}^{D_h}$, $h^{(t)} \in \mathbb{R}^{D_h}$.

If the expected output of a neuron during testing is far different as it was during training, the magnitude of the outputs could be radically different, and the behavior of the network is no longer well-defined. Therefore, all the parameters should be divided by retain rate $1 - p$ (blue part in above formula), so that $\mathbb{E}_{p_{\text{drop}}}[\hat{h}_{\text{drop}}]_i = h_i$. If we do not such correction in training, we should multiply $1 - p$ to all related parameters.

If we use dropout in testing, the result is *unstable* (vary from every testing) because of the dropout is drawn from Bernoulli distribution. Therefore, we should apply dropout only during training but not during evaluation or testing.

A motivation of dropout is to reduce complex *co-adaptations* among the hidden units which comes from the superiority of sexual reproduction compared with asexual. The criterion for natural selection may be the mix-ability of genes. A gene which rely only on a *small* number of other genes is able to work well with another random set of genes that could come from another parent's genes in sexual reproduction. Hence, it will more robust against noises and increases the chance of a new gene improving the fitness of an individual.

2. Batch Normalization

Although **batch normalization**¹¹ is like normalization used in data preprocessing with N be the mini-batch size instead of the dataset size, it is inserted between hidden layers to normalize the output of last layer. It leads to achieve the fixed distributions of inputs that would remove the ill effects of the internal covariate shift. *Internal Covariate Shift* is defined as the change in the distribution of network activations due to the change in network parameters during training.

¹¹ Ioffe and Szegedy 2015

The batch normalization in training is defined as follows:

$$\text{BN}(h_i) = \gamma \frac{h_i - \mu_B}{\sqrt{\sigma_B^2 + \epsilon}} + \beta \quad (1.24)$$

where γ, β are **trainable parameters**, μ_B, σ_B^2 are the mean and variance over the mini-batch as the way of normalization for data preprocessing. Since the mean subtraction will *ignore* the learned bias which may be useful to the model. The trainable γ and β are used to correct them and make the BN layer trainable. They ensure that the batch normalization inserted in the network can represent the identity transform.

However, when testing, we cannot use mini-batch in most time. We instead feed one sample into the model. Therefore, we leverage m training mini-batches to perform **unbiased estimates** of them:

$$\begin{aligned}\mathbb{E}[h_i] &\leftarrow \mathbb{E}_B[\mu_B] \\ \text{Var}[h_i] &\leftarrow \frac{m}{m-1} \mathbb{E}_B[\sigma_B^2]\end{aligned}$$

Note that in many implementations, the above estimation is replaced with the way like the momentum used in RMSprop. More details refer to the source code, e.g. Keras. For the motivation, the use of BatchNorm makes models much less sensitive to parameter initialization.

1.2.5 Practice: Named Entity Recognition

To find and classify words as entities (e.g. location, or organization) in text. One simple idea is that train softmax classifier to classify a center word by taking *concatenation* of word vectors surrounding it in a window (*word window*)¹². To perform NER of location, we need (unnormalized) score for each window, and make *true window's* (i.e. location in the center) score larger and other *corrupt window's* score lower. The model is formulated as:

¹² Collobert and Weston 2008

$$s = W_2 f(W_1 x + b) \quad (1.25)$$

The objective function (*max-margin loss*) is:

$$J = \max(0, s_c - (s - 1)) \quad (1.26)$$

where s and s_c is the score of true window and corrupt window. It ensures each window with an NER location at its center should have a score +1 higher than any window without a location at its center.

1.2.6 HW2

Gradient calculation and implementation of word2vec.

1. Written: Understanding word2vec

$$\begin{aligned}
(a) \quad \hat{y}_o &= P(O = o | C = c) \\
(b) \quad \frac{\partial J}{\partial \mathbf{v}_c} &= (\hat{\mathbf{y}} - \mathbf{y})^\top \mathbf{U}^\top \\
(c) \quad \frac{\partial J}{\partial \mathbf{U}} &= \mathbf{v}_c (\hat{\mathbf{y}} - \mathbf{y})^\top \\
(d) \quad \sigma(\mathbf{x}) &= \frac{1}{1 + \exp(-\mathbf{x})}, \quad \frac{d\sigma(\mathbf{x})}{d\mathbf{x}} = \text{diag}(\sigma(x_i)(1 - \sigma(x_i))) \\
(e) \quad \frac{\partial J}{\partial \mathbf{v}_c} &= \sum_k \sigma(\mathbf{u}_k^\top \mathbf{v}_c) \mathbf{u}_k^\top - (1 - \sigma(\mathbf{u}_o^\top \mathbf{v}_c)) \mathbf{u}_o^\top \\
\frac{\partial J}{\partial \mathbf{u}_o} &= (\sigma(\mathbf{u}_o^\top \mathbf{v}_c) - 1) \mathbf{v}_c^\top \\
\frac{\partial J}{\partial \mathbf{u}_k} &= \sigma(\mathbf{u}_k^\top \mathbf{v}_c) \mathbf{v}_c^\top \\
(f) \quad (i) \quad \frac{\partial J}{\partial \mathbf{U}} &= \sum_o \mathbf{v}_c (\hat{\mathbf{y}}_o - \mathbf{y}_o)^\top \\
(ii) \quad \frac{\partial J}{\partial \mathbf{v}_c} &= \sum_o (\hat{\mathbf{y}}_o - \mathbf{y}_o)^\top \mathbf{U}^\top \\
(iii) \quad \frac{\partial J}{\partial \mathbf{v}_w} &= \mathbf{0}
\end{aligned}$$

2 Coding: Implementing word2vec

Note that \mathbf{U}, \mathbf{V} in the handout are the matrices whose i -th column is the n -dimensional embedded vector for word w_i . However, in the codes of HW2, all the centerWordVectors and outsideVectors are as rows.

Use shape convention to check the result.

1.3 Dependency Parser

Two views of linguistic structure: (1) constituency (i.e., phrase structure grammar, or context-free grammar) (2) Dependency structure. Dependence parse trees (single root with optional fake root, acyclic) use binary asymmetric relations which depicted as typed arrows going from *head* to *dependent*. Note that the natural language is ambiguity.

Basic transition-based dependency parser¹³ with stack $\sigma = [\text{ROOT}]$, buffer $\beta = w_1, \dots, w_n$, set of dependency arcs $A = \emptyset$, and a set of actions (*transitions*) based on the above 3-tuple:

1. Shift: $\sigma, w_i | \beta, A \Rightarrow \sigma | w_i, \beta, A$
2. Left-Arc reduction: $\sigma | w_i | w_j, \beta, A \Rightarrow \sigma | w_j, \beta, A \cup \{r(w_j, w_i)\}$
3. Right-Arc reduction: $\sigma | w_i | w_j, \beta, A \Rightarrow \sigma | w_i, \beta, A \cup \{r(w_i, w_j)\}$

where $r(w_j, w_i)$ denotes w_i is the dependency of w_j (e.g. nsubj(ate \rightarrow I)), $|$ and \cup stand for concatenate. The finish state is: $\sigma = [w], \beta =$

¹³ Nivre 2003

∅. How to select (search) the best choice among the exponential size of different possible parse trees is the problem. In 1960s, they use *dynamic programming algorithms* ($\mathcal{O}(n^3)$). In paper ¹⁴, the authors predict each action by a discriminative classifier (e.g. SVM classifier) which is more efficient but the accuracy is fractionally below the state-of-the-art.

¹⁴ Nivre 2003

1.3.1 Neural Dependency Parsing

Compared with traditional sparse feature-based discriminative dependency parsers, the work by ¹⁵ utilizes **feedforward neural network model** with simple **dense layers** and the softmax layer to predict each transition. The input features with embedding dimension d are:

Chen and Manning 2014

1. $x^w \in \mathbb{R}^{d*N_w}$: The top 3 words on the stack and buffer $s_1, s_2, s_3, b_1, b_2, b_3$; the first and second leftmost / rightmost children of the top two words on the stack $lc_1(s_i), rc_1(s_i), lc_2(s_i), rc_2(s_i), i = 1, 2$; the leftmost of leftmost / rightmost of rightmost children of the top two words on the stack $lc_1(lc_1(s_i)), rc_1(rc_1(s_i)), i = 1, 2$; In total, $N_w = 18$.
2. $x^t \in \mathbb{R}^{d*N_t}$: The corresponding POS (Part-of-speech, e.g. noun, verb, adjective) tags for S_{word} , $N_t = 18$.
3. $x^l \in \mathbb{R}^{d*N_l}$: The corresponding arc labels of words, excluding those 6 words on the stack/buffer, $N_l = 12$.

The predicted class is the one of transitions (i.e. shift, left/right arc reduction): $p = \text{softmax}(W_2 f(W_1^w x^w + W_1^t x^t + W_1^l x^l + b_1) + b_2)$, where $f(\cdot)$ is the activation function (e.g. ReLU, or x^3). The number of class is 3 when untyped reductions or $T * 2 + 1$ when typed reductions (e.g. left-arc reduction with type *nsubj*).

Note that we use a special **NULL** token for non-existent elements: when the stack and buffer are empty or dependents have not been assigned yet.

Here are mainly four types of parsing error:

- **Prepositional Phrase Attachment Error**
- **Verb Phrase Attachment Error**
- **Modifier Attachment Error**
- **Coordination Attachment Error**

1.3.2 HW3

1. Machine Learning & Neural Networks

- (a) Adam

- (i) Because $\beta = 0.9$, most of the final gradients (m) come from the past (90%). Even if current gradients are varying much from previous, it only occupy $1 - \beta_1 = 0.1$ of the final gradients.
 - (ii) Parameters that have not been updated much in the past are likelier to have higher learning rates.
- (b) Dropout
- (i) If the expected output of a neuron during testing is far different as it was during training, the magnitude of the outputs could be radically different, and the behavior of the network is no longer well-defined. Thus, all the parameters should divided by retain rate $1 - p$, so that $\mathbb{E}_{p_{\text{drop}}}[\mathbf{h}_{\text{drop}}]_i = h_i$.
 - (ii) If we use dropout in testing, the result is unstable because of the dropout is drawn from Bernoulli distribution.

2. Neural Transition-Based Dependency Parsing

- (a) The remaindered configuration is:

| Stack | Buffer | New dependency | Transition |
|--------------------------------|-----------------------|--------------------------------|------------|
| [ROOT, parsed, this] | [sentence, correctly] | | SHIFT |
| [ROOT, parsed, this, sentence] | [correctly] | | SHIFT |
| [ROOT, parsed, sentence] | [correctly] | sentence \rightarrow this | LEFT-ARC |
| [ROOT, parsed] | [correctly] | parsed \rightarrow sentence | RIGHT-ARC |
| [ROOT, parsed, correctly] | [] | | SHIFT |
| [ROOT, parsed] | [] | parsed \rightarrow correctly | RIGHT-ARC |
| [ROOT] | [] | ROOT \rightarrow parsed | RIGHT-ARC |

(b) $2n$

(e) dev UAS: 88.56, test UAS: 89.07

(f) Refer to [Stanford CoreNLP](#) to visualize the results of dependencies parses (also may contain errors).

- (i) Verb Phrase Attachment Error, wedding \rightarrow fearing should be heading \rightarrow fearing.
- (ii) Coordination Attachment Error, makes \rightarrow rescue should be rush \rightarrow rescue.
- (iii) Prepositional Phrase Attachment Error, named \rightarrow Midland should be guy \rightarrow Midland.
- (iv) Modifier Attachment Error, elements \rightarrow most should be crucial \rightarrow most.

Note that in the source code, the restriction that the version of PyTorch must be 1.0.0 is meaningless and thus I remove it.

1.4 Language Modeling and Recurrent Neural Networks

Language Modeling: given a sequence of words $x^{(1)}, \dots, x^{(t)}$, compute the probability distribution of the next word at $x^{(t+1)}$:

$$P(x^{(t+1)} | x^{(1)}, \dots, x^{(t)}) \quad (1.27)$$

The joint probability of a text is:

$$P(x^{(1)}, \dots, x^{(T)}) = \prod_{t=1}^T P(x^{(t)} | x^{(t-1)}, \dots, x^{(1)}) \quad (1.28)$$

n -gram is a chunk of n consecutive words: unigram, bigram, trigram, 4-gram, ... n -gram language model is based on a simplifying assumption: $x^{(t+1)}$ depends only on the preceding $n - 1$ words with i.i.d.:

$$\begin{aligned} P(x^{(t+1)} | x^{(1)}, \dots, x^{(t)}) &= P(x^{(t+1)} | x^{(t)}, \dots, x^{(t-n+2)}) \\ &= \frac{P(x^{(t+1)}, x^{(t)}, \dots, x^{(t-n+2)})}{P(x^{(t)}, \dots, x^{(t-n+2)})} \end{aligned}$$

where the n -gram and $(n-1)$ -gram probabilities are calculated by *counting*. There are some *sparsity problems* with the above n -gram models such as the numerator or denominator is zero. Some tricks such as *smoothing* (add small δ to the count) and *backoff* (e.g. 4-gram backoff to 3-gram) are proposed to solve them.

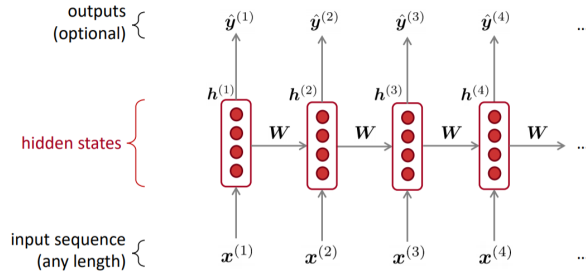


Figure 1.8: Principle of RNN

To process *variable length sequential input* such as text, **Recurrent Neural Network** (RNN) is introduced. As the principle of RNN shown in Fig. 1.8: *repeat* (i.e. **unfold** or unroll) the same RNN cell for each time-step but with different input and previous **hidden state**. A vanilla RNN for language modeling is:

$$\begin{aligned} h^{(t)} &= \sigma(W_h h^{(t-1)} + W_x x^{(t)} + b_1) \\ \hat{y} &= P(x^{(t)} | x^{(t-1)}, \dots, x^{(1)}) \\ &= \text{softmax}(U h^{(t)} + b_2) \end{aligned}$$

Note that for n -gram, increasing n makes sparsity problems worse. Typically $n \leq 5$.

where $\sigma(\cdot)$ is the activation function, and $\mathbf{h}^{(0)}$ is the initial (random or zero) hidden state. The gradient w.r.t. the weight matrix is the *sum* of the gradients w.r.t. each time it appears using **back-propagation through time** (BPTT, just as same as normal back-prop). And the **evaluation metric** for language modeling is *perplexity* which is equal to the exponential of the cross-entropy losses:

$$\begin{aligned} \text{perplexity} &= \prod_{t=1}^T \left(\frac{1}{P_{LM}(\mathbf{x}^{(t+1)} | \mathbf{x}^{(t)}, \dots, \mathbf{x}^{(1)})} \right)^{1/T} \\ &= \exp \left(\frac{1}{T} \sum_{t=1}^T -\log \hat{y}^{(t)} \right) \end{aligned} \quad (1.29)$$

There are some other applications of RNN: part-of-speech tagging, named entity recognition, sentence classification, text generator, encoder module, etc. The final feature can be the final hidden state or element-wise max/mean of all hidden states. Using chain rule, we get $\frac{\partial J^{(n)}}{\partial \mathbf{h}^{(1)}} = \frac{\partial J^{(n)}}{\partial \mathbf{h}^{(n)}} \times \prod_{i=2}^n \frac{\partial \mathbf{h}^{(n)}}{\partial \mathbf{h}^{(i-1)}} = \frac{\partial J^{(n)}}{\partial \mathbf{h}^{(n)}} \times \prod_{i=2}^n \sigma' \circ \mathbf{W}_h$. For a large n and small \mathbf{W}_h , it's easy to encounter the vanishing gradient problem. In overall, the *vanilla* RNN has these disadvantages: (1) recurrent computation is slow (2) hard to access long-term information (**long-term dependencies**) due to *gradient vanish* and *gradient explosion*.

We can formalize the above vanishing intuitions according to ¹⁶. Let \mathbf{W}_h have the **eigenvalues** $\lambda_1, \dots, \lambda_n$ such that $|\lambda_1| > |\lambda_2| > \dots > |\lambda_n|$ and the corresponding (left) eigenvectors $\mathbf{q}_1^\top, \dots, \mathbf{q}_n^\top$ which have unit norms: $\mathbf{q}_i^\top \mathbf{W}_h = \lambda_i \mathbf{q}_i$. We can rewrite the gradients $\frac{\partial J^{(n)}}{\partial \mathbf{h}^{(n)}} = \sum_{i=1}^N c_i \mathbf{q}_i^\top$ where $c_i = 0$ for $i < j$ and $c_j \neq 0$. Thus, the overall gradient is:

¹⁶ Pascanu et al. 2013

$$\begin{aligned} \frac{\partial J^{(n)}}{\partial \mathbf{h}^{(1)}} &= \frac{\partial J^{(n)}}{\partial \mathbf{h}^{(n)}} \times \prod_{i=2}^n \sigma' \circ \mathbf{W}_h \\ &= \sum_{i=1}^N c_i \mathbf{q}_i^\top (\text{diag}(\sigma'))^{n-1} (\mathbf{W}_h)^{n-1} \\ &= \sum_{i=1}^N c_i \mathbf{q}_i^\top (\mathbf{W}_h)^{n-1} (\text{diag}(\sigma'))^{n-1} \\ &= c_j \lambda_j^{n-1} \mathbf{q}_j^\top (\text{diag}(\sigma'))^{n-1} + \lambda_j^{n-1} \sum_{i=j+1}^N c_i \left(\frac{\lambda_i}{\lambda_j} \right)^{n-1} \mathbf{q}_i^\top (\text{diag}(\sigma'))^{n-1} \\ &\approx c_j \lambda_j^{n-1} (\sigma')^{n-1} \mathbf{q}_j^\top \end{aligned}$$

where $\frac{\lambda_i}{\lambda_j} < 1$, and for large n we have $\lim_{n \rightarrow \infty} \left(\frac{\lambda_i}{\lambda_j} \right)^{n-1} = 0$. Therefore, if $\forall j, \sigma' < \frac{1}{\lambda_j}$ then we get vanishing gradient. Note that,

$\sup(\text{sigmoid}') = \frac{1}{4}$, $\sup(\text{ReLU}') = 1$. Thus, the largest eigenvalue $\lambda_1 < \frac{1}{\sigma'}$ will lead to vanishing.

For avoid gradient explosion, one simple method is *gradient clipping*: $\hat{g} \leftarrow \frac{\text{threshold}}{\|g\|} g$ if $\|g\| \geq \text{threshold}$. As for fixing vanishing gradient, many RNN variants are introduced such as **Long Short-Term Memory** (LSTM)¹⁷ and **Gated Recurrent Unit** (GRU)¹⁸. LSTM uses two *separated* memories: *hidden state* $h^{(t)}$ for *short-term* information and *cell state* $c^{(t)}$ for *long-term* information. There are three *gates* performed in each LSTM *cell*:

$$\text{Forget gate: } f^{(t)} = \sigma(W_f h^{(t-1)} + U_f x^{(t)} + b_f)$$

$$\text{Input gate: } i^{(t)} = \sigma(W_i h^{(t-1)} + U_i x^{(t)} + b_i)$$

$$\text{Output gate: } o^{(t)} = \sigma(W_o h^{(t-1)} + U_o x^{(t)} + b_o)$$

$$\text{New cell content: } \tilde{c}^{(t)} = \tanh(W_c h^{(t-1)} + U_c x^{(t)} + b_c)$$

$$\text{Cell state: } c^{(t)} = f^{(t)} \circ c^{(t-1)} + i^{(t)} \circ \tilde{c}^{(t)}$$

$$\text{Hidden state: } h^{(t)} = o^{(t)} \circ \tanh c^{(t)}$$

where σ is sigmoid function, \circ is element-wise product, and \tanh used in hidden state (the last formula) is to provide non-linearity and normalize $c^{(t)}$ to $(0, 1)$. The structure of LSTM is shown in Fig.1.9 made by [colah's blog](#).

Those three gates (forget, input, output) enable the abilities of erase, read and write for LSTM. Each element of the gates are between 1 (open) and 0 (close). While The LSTM architecture makes it *easier* for the RNN to preserve long-distance dependencies, it does not *guarantee* that there is no vanishing/exploding gradient.

In the other hand, GRU combines input and forget gate into *update* gate, and add new *reset* gate to select useful part of previous hidden state to compute new state content. While there is no conclusive evidence that GRU consistently performs better than LSTM or vice versa, GRU is computed more efficient due to fewer parameters.

$$\text{Update gate: } u^{(t)} = \sigma(W_u h^{(t-1)} + U_u x^{(t)} + b_u)$$

$$\text{Reset gate: } r^{(t)} = \sigma(W_r h^{(t-1)} + U_r x^{(t)} + b_r)$$

$$\text{New hidden state content: } \tilde{h}^{(t)} = \tanh(W_h(r^{(t)} \circ h^{(t-1)}) + U_h x^{(t)} + b_h)$$

$$\text{Hidden state: } h^{(t)} = (1 - u^{(t)}) \circ h^{(t-1)} + u^{(t)} \circ \tilde{h}^{(t)}$$

The vanishing gradient problem appears not only in RNNs, but also for most all other neural networks including MLP (dense layers) and CNNs, especially for deep ones. One solution is add more *direct connections* between future apart layers to allow gradients flow more easier. For example, **Residual connections** (aka.

¹⁷ Hochreiter and Schmidhuber 1997

¹⁸ Cho et al. 2014

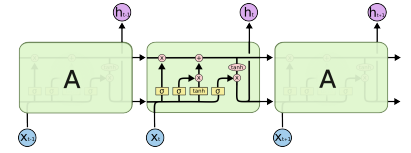


Figure 1.9: The repeating module in an LSTM contains four interacting layers.

ResNet¹⁹ or skip-connections) is shown in Fig.1.10 where an identity skips two layers. Another example is **Dense connections** (aka. DenseNet²⁰) which directly connect every layers to every layers where the output of each layer will **concatenate** the input as presented in Fig.1.11. **Highway connections** is inspired from the gates of LSTM and similar to residual connections, where the identity connect and the transformation layer is controlled by a dynamic gate.

Apart from the above RNNs, there are other important RNN architectures: **Bidirectional RNNs** and **Multi-layer RNNs** (aka. stacked RNNs). The definition of bidirectional RNNs is given by:

$$\text{Forward RNN: } \vec{h}^{(t)} = \overrightarrow{\text{RNN}}_{FW}(\vec{h}^{(t-1)}, x^{(t)})$$

$$\text{Backward RNN: } \overleftarrow{h}^{(t)} = \overleftarrow{\text{RNN}}_{BW}(\overleftarrow{h}^{(t-1)}, x^{(t)})$$

$$h^{(t)} = [\vec{h}^{(t)}; \overleftarrow{h}^{(t)}]$$

While LSTM became the dominant approach between 2013 to 2016, *Transformer* is the state-of-the-art now in 2019. More descriptions refer to Section 1.6. And another notable fact is that RNNs parallelize badly and are slow compared with CNNs.

1.5 Seq2Seq and Attention

Pre-neural machine translation: (1) Rule-based bilingual dictionary in 1950s (2) Statistical machine translation from 1990s to 2010s. More formally, statistical machine translation from *source language* x to *target language* y is given by:

$$\begin{aligned} \arg \max_y P(y|x) &= \arg \max_y \frac{P(x, y)}{P(x)} \\ &= \arg \max_y \frac{P(x|y)P(y)}{P(x)} \\ &= \arg \max_y P(x|y)P(y) \end{aligned}$$

where $P(y|x)$ is broke down to two components according to **Baye's rule**: $P(x|y)$ is the translation model which learn from parallel (bilingual) data to model the fidelity of words and phrases whether x is well- or ill-formed, $P(y)$ is the language model which learn from monolingual data of target language to model the fluency of the whole sentence regardless of their connection to the French.

In practice, we further consider **alignment** (word-level correspondence) because there are one-to-many, many-to-one, many-to-many, and even no counterpart apart from one-to-one

¹⁹ He et al. 2016

²⁰ Huang et al. 2017

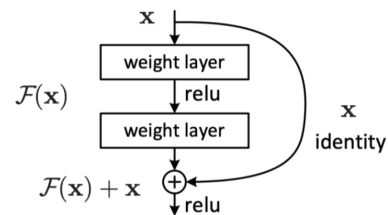


Figure 1.10: Residual connections.

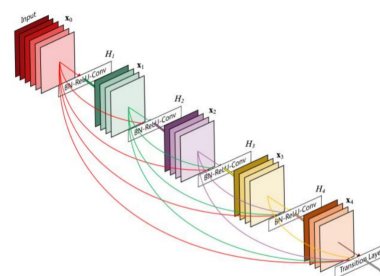


Figure 1.11: Dense Net.

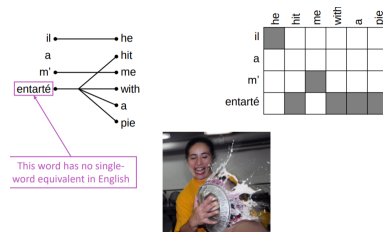


Figure 1.12: Alignment from french to english translation.

reflection relations. One example of one-to-many (entarté) and no counterpart (a) is shown in Fig.1.12. More examples refer to the original paper ²¹. Therefore, we add alignment to the model: $P(x, a|y)$.

²¹ Brown et al. 1993

The core idea of Seq2Seq model is using two RNN to construct an *encode-decode* architecture. At test time, first we feed source sentence (with embedding) into the encoder RNN, then we use the last hidden state of the encode RNN as the initial state of the decoder RNN as a conditional language model. The output word at position t is given by:

$$w_t = \text{softmax}(\text{RNN}(h^{(t-1)}, W_e \arg \max \text{softmax}(h^{(t-1)})))$$

where W_e is the embedding table of the target language and the first input $x_1 = \langle \text{START} \rangle$ is a special token and repeatedly output until output $\langle \text{END} \rangle$. Note that there are two different embedding lookup table for source and target language. When training, we need provide parallel dataset whose samples consist of bilingual sentences. As for training, the diagram is represented in Fig.1.13.

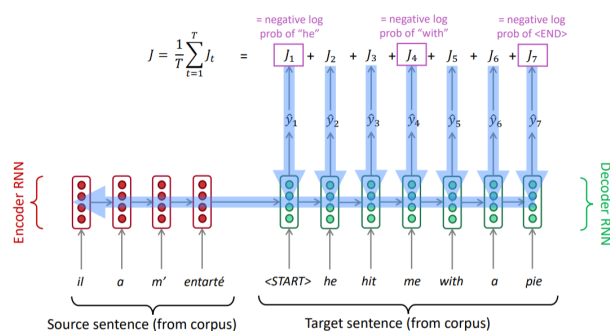


Figure 1.13: Training phase for NMT.

However, the aforementioned (greedy) decoding has no way to undo decisions. This is, if one of the output words are wrong, all the follow-up outputs are also wrong. **Beam search** decoding is utilized to fix this problem: on each step of decoder, keep track of the k (in practice around 5 to 10) most probable partial translation (*hypotheses*, path, or branch). The target is the path with the largest cumulative log probabilities with shortest one (i.e. average). Apart from all paths reach $\langle \text{END} \rangle$ as the stop sign, we can set some pre-defined cutoff for maximum number of timesteps or finished paths.

Although NMT is much simple and less human engineering effort while achieves better performance, NMT is difficult to control (i.e. specify rules or guidelines) and less interpretable which leads to hard to debug.

The popular evaluation metric for MT is **BLEU** (Bilingual Evaluation Understudy). Its calculation is based on n-gram precision

plus a penalty for too-short translations. Note that BLEU is useful but imperfect because there are many valid translations which has low n -gram overlap with the ground truth translation. NMT outperforms SMT quickly, but there are still many difficulties remain:

- Out-of-vocabulary words.
- Domain mismatch between train and test data.
- Maintaining context over longer text (long-term dependencies).
- Low-resource language pairs (few-shot learning).
- Using common sense is still hard.

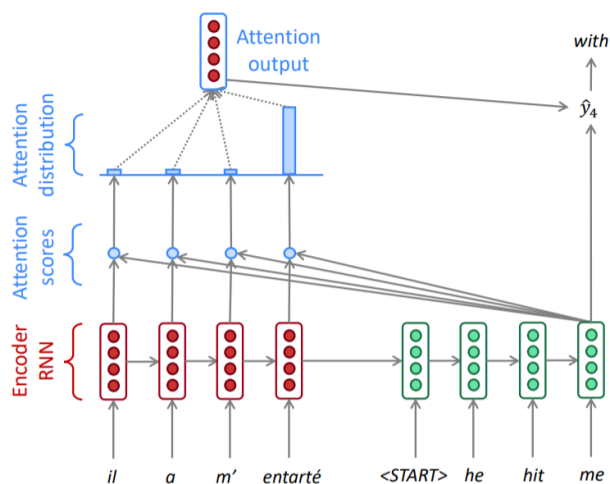


Figure 1.14: Seq2Seq with attention.

We notice that only the last hidden state from encoder RNN to represent the whole source sentence may be the information bottleneck. Thus we introduce Seq2Seq with **attention**: on each step of the decoder, we use the attention distribution (like soft version of alignment) for each hidden states of encoder RNN to take a weighted sum of the encoder hidden states as the current decoder hidden state. Note that sometimes the input of decoder RNN will concatenate the previous translated word vector and the previous attention output. The diagram can be seen in Fig.1.14. In addition, attention in here helps with vanishing gradient problem by providing shortcut to faraway states, and also provides

some interpretability. More formally, the attention for Seq2Seq is:

$$\mathbf{e}^{(t)} = [\mathbf{s}_t^\top \mathbf{h}_1, \dots, \mathbf{s}_t^\top \mathbf{h}_n] \in \mathbb{R}^n$$

$$\text{Attention distribution: } \boldsymbol{\alpha}^{(t)} = \text{softmax}(\mathbf{e}^{(t)})$$

$$\text{Attention output: } \mathbf{a}^{(t)} = \sum_{i=1}^n \alpha_i^{(t)} \mathbf{h}_i \in \mathbb{R}^h$$

$$\text{Attention decoder hidden state: } [\mathbf{a}^{(t)}; \mathbf{s}_t] \in \mathbb{R}^{2h} \quad (1.30)$$

where encoder hidden states $\mathbf{h}_1, \dots, \mathbf{h}_n \in \mathbb{R}^h$, and decoder hidden state at timestep t is $\mathbf{s}_t \in \mathbb{R}^h$.

More general definition of attention: given a set of vector *values*, and a vector *query*, attention is to compute a weighted sum of the values, dependent on the query. We can find that attention is a way to obtain a *fixed-size* representation of an arbitrary set of representations (e.g. sequential features). There are many ways to obtain query vector: dot product $e_i = \mathbf{s}^\top \mathbf{h}_i$, multiplicative $e_i = \mathbf{s}^\top \mathbf{W} \mathbf{h}_i$, additive $e_i = \mathbf{v}^\top \tanh(\mathbf{W}_1 \mathbf{h}_i + \mathbf{W}_2 \mathbf{s})$, where $\mathbf{v}, \mathbf{W}, \mathbf{W}_1, \mathbf{W}_2$ are trainable parameters.

As for *out-of-vocabulary* (OOV) problem in NMT, we can increase the size of vocabulary. However, large vocabulary leads to huge computations in the softmax layer. To solve it, we can use *hierarchical softmax* (tree-structured vocabulary). We can also use character-based ²² (or sub-word) model instead of word-based to handle OOV problem. But pure character-level seq2seq is very slow (as longer sequence got through LSTM, see HW5).

²² Kim et al. 2016

1.5.1 HW4

1. Neural Machine Translation with RNNs

Refer to the source codes for detail. Note that `pack_padded_sequence` is often used to reduce unnecessary computations of padded elements for variable length sequence. In implementation, I found that for reshaping the `last_hidden` with shape (2, b, h) to shape (b, 2h), I had to use `torch.cat((last_hidden[0, :], last_hidden[1, :]), axis=1)` instead of `.permute(1, 0, 2).view(b, -1)`. Common used torch function include: `torch.squeeze`, `torch.unsqueeze`, `torch.split`, `torch.cat`, `torch.stack`, `torch.Tensor.view`, `torch.Tensor.permute`, `torch.Tensor.size`.

(i) Corpus BLEU score when testing: 22.68.

(j) Dot product attention is simple but requires $\mathbf{s}_t, \mathbf{h}_i$ have the same shape. Additive attention could learn more complicated features but need more computations.

2. Analyzing NMT Systems

(a)

- (i) Error: duplication of translations ("favorite"), Reason: model limitations in preventing duplications, Solution: add penalty on duplicated translations.
- (ii) Error: the superlative comparison is mistranslated (most instead of more), Reason: could be because specific linguistic construct of comparison in Spanish, Solution: train with more comparison examples or add penalty for comparison examples.
- (iii) Reason: out-of-vocabulary words of names, Reason: the model do not consider out-of-vocabulary words, Solution: keep the out-of-vocabulary the same according to named entity recognition or attention output, or just increase the size of vocabulary.
- (iv) Error: "go around" is mistranslated to "go back," "block" is mistranslated to "apple", Reason: "apple" and "block" are same in Spanish "manzana", "go around" and "go back" are also similar in Spanish, Solution: train with more examples about such polysemic words.
- (v) Error: "teacher" is mistranslated to "women" for Spanish "profesores", Reason: training dataset may contain the above mistakes, Solution: correct the dataset and add rules to correct this error.
- (vi) Error: "hectareas" (1/4 acres) is mistranslated to "acres", Reason: there is no such measure unit in English, Solution: identify measure unit in the source sentence and perform some rules.

(b)

(c) BLEU

- (i) $p_1 = 0.5477, p_2 = 0.6325, c_2$, agree.
- (ii) $p_1 = 0.4484, p_2 = 0.2589, c_1$, disagree.
- (iii) There are often many valid translations for one sentence.
- (iv) Pros and cons of BLEU:
 - Pros: 1. BLEU has frequently been reported as correlating well with human judgement. 2. few human effort with inexpensive computations.
 - Cons: 1. A good translation can get a poor BLEU score because it has low n-gram overlap with the human (reference) translation. 2. It does not consider meaning and sentence structure, etc.

Note that despite the HW taking sentence-level BLEU, BLEU is used as a **corpus-level measure**.

1.5.2 Tips for Research

Two basic start points:

- Start with a (domain) problem of interest and try to find good/better ways to address it than currently known/used.
- Start with a technical approach of interest, and work out good ways to extend or improve it or new ways to apply it.
- Experimental strategy: start from simple model and work incrementally on models, initially run on a tiny dataset (e.g. 8 examples), make sure get 100% accuracy on testing, then run your model on a large dataset that still get score to 100% on training (it's okay to overfitting), and try to regularize the model.

Recommendation sources: ACL, NeurIPS, ICML, ICLR, arXiv ([arXiv-sanity](#)), [PapersWithCode](#) (leaderboards for various tasks).

1.6 Contextual Word Representations and Pretraining

The context-independent word representations (such as Word2Vec and GloVe) have limitations: the same word will have a ton of meanings in different contexts, but has a fixed word representation in these models. Instead of a fixed embedding table, we can (unsupervisedly) pre-train a neural network model and freeze it as a word representation model (e.g., $f(w, c(w))$).

TagLM ²³ ("Pre-ELMo"): pre-train a Neural Language Modeling on large unlabeled corpus and get the output as the contextual word representation. Use both (context-independent) word embeddings (e.g., Word2Vec) and (frozen) LM embeddings in the downstream sequence tagging model (NER task) w/ small task-labeled data.

²³ [Peters et al. 2017](#)

CoVe ²⁴: Use NMT as the pre-training task. But, the results aren't as strong as the simpler NLM.

²⁴ [McCann et al. 2017](#)

ELMo ²⁵: the word representation is the (learnable) weighted sum of hidden features of each layer in stacked BiLSTM.

²⁵ [Peters et al. 2018](#)

Transformer ²⁶: scaled dot-product self-attention (QKV formula, the softmax scores are derived from the same source as the target sequence), multihead (like channels in CNN), position embedding (because self-attention has not position-aware), good parallelization than RNN (sequential because of the unroll operation), feed-forward layer (like 1×1 convolution), residual connection (propagate the position information to higher layers), layer normalization, masked attention. Key resource:

²⁶ [Vaswani et al. 2017](#)

<http://nlp.seas.harvard.edu/2018/04/03/attention.html>. Other applications²⁷ of the transformer are shown in the sidenote.

BERT²⁸ (Bidirectional Encoder Representations from Transformers): two tasks: masked (bidirectional) LM and next sentence prediction (NSP, but is useless according to follow up papers because this task is misconstrued and too easy).

1.7 Attention Model and Question Answering

Compared with previous versions, Stanford Question Answering Dataset (SQuAD) 2.0 has many unanswerable questions. Each question has 3 gold answers, and there are two evaluation metrics:

- Exact match (EM): 1/o accuracy on whether match one of 3 answers.
- F1: take system (predict) and gold answers as bag of words to calculate F1 scores.

where both metrics ignore punctuation and articles. Disadvantages: the answer is only sub-span of the paragraph (the answers are directly appeared in the question), without indirect answers, e.g., yes/no, counting, and implicit why.

A simple and successful architecture is the Stanford Attentive Reader²⁹ which likes the attentive Seq2Seq model. For paragraph $[p_1, \dots, p_m]$ and question $[q_1, \dots, q_n]$ with GloVe embedding, the index of start token and end token for the sub-span answer is:

$$q = [\overrightarrow{\text{LSTM}}(\overrightarrow{h}^{(m-1)}, q_{m-1}); \overleftarrow{\text{LSTM}}(\overleftarrow{h}^{(2)}, q_2)]$$

$$\tilde{p}_i = \text{Bi-LSTM}(p_i)$$

$$\text{Start token: } \text{softmax}_i(q^\top W_s \tilde{p}_i)$$

$$\text{End token: } \text{softmax}_i(q^\top W_e \tilde{p}_i)$$

Compared with SAR, Stanford Attentive Reader++³⁰ mainly leverages these updates: (1) use weighted sum (attention) of all hidden states in the encoder of question, (2) multi-layers Bi-LSTM (3) add more features in the encoder of the paragraph, including linguistic features such as POS&NER tags&term frequency, exact match (whether the word appears in the question), aligned question embedding ($\sum_j \text{softmax}_{j'}(\alpha(p_i) \cdot \alpha(q_{j'})) \times q_j$ where α is dense layer with ReLU activation and p_i, q_j are embedding of the paragraph and the question word).

Bi-Directional Attention Flow (DiDAF)³¹ is a state-of-the-art (SoTA) method in 2018 and 2017. It comes up with word and char embedding, and the main innovation is attention flow layer

²⁷ Parmar et al. 2018; Huang et al. 2019

²⁸ Devlin et al. 2019

Transformer for Image: due to the self-similarity (or translational equivariance) in images, we can limit the receptive field of the self-attention by masking with diagonal matrix (e.g., tridiagonal matrix). Transformer for Music Generation: decrease memory usage for storing the relative position embeddings of $L \times L$ to L by padding, reshape, and slice.

²⁹ Chen et al. 2016

³⁰ Chen et al. 2017

³¹ Seo et al. 2016

which consists of Query2Context (i.e. Question2Paragraph) and Context2Query. C2Q attention is defined as:

$$\text{Similarity matrix: } s_{ij} = w_s^\top [p_i; q_j; p_i \circ q_j] \in \mathbb{R} \quad (1.31)$$

$$\alpha_i = \text{softmax}(s_{i,:}) \quad (1.32)$$

$$a_i = \sum_j \alpha_{ij} q_j \in \mathbb{R}^{2h} \quad (1.33)$$

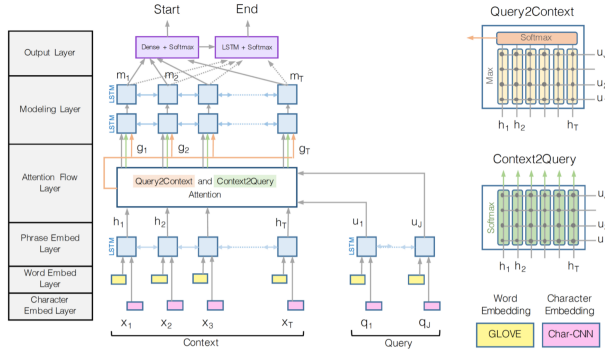


Figure 1.15: Architecture of BiDAF

Q2C attention (the weighted sum of the most important words in the context w.r.t the query) is given by:

$$m_i = \max_j s_{ij} \quad (1.34)$$

$$\beta = \text{softmax}(m) \quad (1.35)$$

$$p' = \sum_i \beta_i p_i \quad (1.36)$$

The BiDAF layer output of each paragraph position is: $b_i = [p_i; a_i; p_i \circ a_i; p_i \circ p'] \in \mathbb{R}^{8h}$. The overall architecture is presented in Fig.1.15. The start token comes from the output of dense layer with softmax built upon the concatenation of the BiDAF output and the output M_1 of modelling layer (2-layer Bi-LSTM) whose input is BiDAF output. And the end token is the dense layer with softmax upon the concatenation of the BiDAF output and the output M_2 of another modelling layer whose input is M_1 .

1.8 ConvNets for NLP

Convolutions are classically used to extract (position-invariant) features from images by leveraging fixed-size *filter (kernel)* weights which are used to perform element-wise product with certain window (patch) of the whole image. *Zero padding* ("same" padding) is often used to keep the length. The kernel window slides with a stride (e.g. 1) over the whole sentence. One more variant about

convolutions is *dilation convolution* (skip conv) which can access more wide *receptive field* with smaller filter. Same as the *channels* (aka. feature maps) in CNNs, we can also define the dimensions of word embedding as channels like 3 channels (RGB) of images. Each channel may indicate one feature (e.g. words about food or sport). *Pooling* (e.g. max, average) over time (aka. global pooling) is utilized to aggregate the information of the whole sentence. *k*-max (global) pooling is also a common technique used in ConvNet for NLP. Besides, local pooling is more common used in CNN circumstances. 1×1 convolution (aka. Network-in-Network connection) is a famous technique that can be seen as a fully connection linear layer across channels. It can be also (often) used to map from many channels to fewer channels to reduce computation while keeping performance.

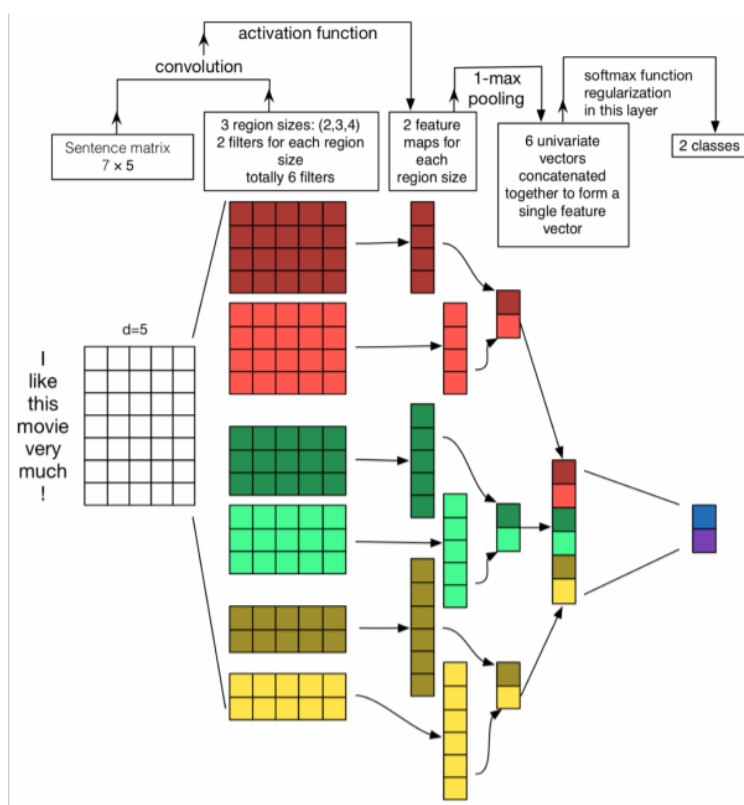


Figure 1.16: Simple ConvNet for NLP classification.

Above is an example ³² of single layer CNN for sentence classification. Fig.1.16 illustrates the architecture. Multiple single convolution layer with different kernel size and global max pooling are used. Specifically, there are 100 feature maps for each of kernel size 3, 4, and 5. The input is two copies of pre-trained word vector (word2vec-300). Keep one as 'static' (no gradient), and the

³² Kim 2014

other 'dynamic' will be updated during training. The 'static' word vector gives unseen words appeared in the test set a better chance of being interpreted correctly. There are several ways of handling these two channels, most common is to simply average them before using in a CNN. The other method is to double the length of the CNN filters. The output of max poolings are concatenated into a vector, which is then fed into a dense layer with softmax to perform classification. Dropout and L2 regularization are used.

VD-CNN³³ for text classification is to build a ResNet like very deep CNN which has 29 depth and many shortcut (residual) connections. Note that 29 depth is not as deep as ResNet whose depths are typically 34, 50, and 101. The researchers found that, for more deep one like 47 layers, the results are fraction worse than the one with 29 layers. The input text sequence will pad into fixed length (1024).

³³ [Conneau et al. 2017](#)

Quasi-recurrent neural network takes the best and parallelizable parts of RNNs and CNNs. It achieves 3 ~ 16 times faster than LSTM. One Quasi-RNN layer consists of a normal convolution layer with left same padding to learn new hidden content z , forget gate f , optional output gate o , and optional input gate i , and three dynamic average pooling (non-parametric) variants to perform RNN-like unfolding.

New hidden contents: $Z = \tanh(W_z * X)$

Forget gates: $F = \sigma(W_f * X)$

Optional output gates: $O = \sigma(W_o * X)$

Optional input gates: $I = \sigma(W_i * X)$

f -pooling: $h_t = f_t \odot h_{t-1} + (1 - f_t) \odot z_t$

fo -pooling:

$c_t = f_t \odot c_{t-1} + (1 - f_t) \odot z_t$

$h_t = o_t \odot c_t$

ifo -pooling:

$c_t = f_t \odot c_{t-1} + i_t \odot z_t$

$h_t = o_t \odot c_t$

where the hidden state h and the cell state c are initialized to zeros, $*$ indicates convolution operator, and \odot denotes element-wise product. A variant dropout called zoneout is adopt in the forget gate channels of stacked QRNN. Dense connection layers from DenseNet are also extended in stacked QRNN. Note that QRNN often need deeper network to get as good performance as LSTM.

1.9 Subword models

Although what our mouths produce is continuous space, phonology posits phonetics the sound stream that consists of smaller distinctive categorical units: phonemes. In linguistics, morphemes (subwords) is the smallest semantic unit that has meanings, e.g. unfortunately can be divided into tree structure of morephemes using *recursive neural networks* ³⁴:

³⁴ Luong et al. 2013

[[un [[fortun(e)]_{ROOT} ate]_{STEM}] _{STEM} ly]_{WORD}

There are many languages that do not segment words or has many separated and joined words, e.g. Chinese, French. Moreover, subword models solve the OOV problem.

1.9.1 Byte Pair Encoding

Like gzip compression. Replace bytes with character ngrams. Merge together several bytes that are common sequences into a new byte. Rather than splitting the word into pieces, the sentence piece ³⁵ tokenizes from the raw text. For example, [Hello] [-Wor] [ld] [.]

³⁵ <https://github.com/google/sentencepiece>

BERT uses a variant of the wordpiece model:

1. (Relatively) common words directly in the vocabulary.
2. Other words are built from wordpieces. e.g., hypatia = h, ##yp, ##ati, ##a (non-initial word pieces are start with ##)

1.9.2 FastText embeddings

Proposed by the authors of Word2Vec. Represent word as char n-grams augmented with boundary symbols and as whole word. For example, where = <wh, whe, her, ere, re>, where.

1.9.3 HW5

Enhance the NMT model with character-based CNN encoder and LSTM decoder.

1. Character-based convolutional encoder for NMT ³⁶

³⁶ Kim et al. 2016

(a) Characters are significantly less diverse than words. So, characters have less semantic information. Besides, we will use CharCNN to map $e_{\text{char}} \times m_{\text{word}}$ to e_{word} .

(b)

$$\begin{aligned}\text{\#parameters(char-based)} &= V_{\text{char}} \times e_{\text{char}} + k \times e_{\text{char}} \times e_{\text{word}} + e_{\text{word}} + (2e_{\text{word}} + 2) \times e_{\text{word}} \\ &= 200,640\end{aligned}$$

$$\begin{aligned}\text{\#parameters(word-based)} &= V_{\text{word}} \times e_{\text{word}} \\ &\approx 12,800,000\end{aligned}$$

The parameters of char-based embedding are much less than word-based. However, the char-based embedding need more computation resources.

(c) The computational efficiency of CNN is better than RNN. Although CNN works worse than RNN on sequential information due to the sliding window, CNN has good ability to capture sequential information inside the kernel window (e.g. $k = 5$). Also, the parameters of CNN is less than LSTM in this scenario.

(d) If the position of objects is not important, Max Pooling seems to be the better choice. If it is, it seems that better results can be achieved with Average Pooling. For me, I don't think one of the poolings has any significant advantage over the other.

(l) The test q1 result is 99.29792465574434 BLEU.

2. Character-based LSTM decoder for NMT

The key idea of this ³⁷ is that when word-level decoder produces an <UNK> token, we can run character-level decoder to generate the rare and out-of-vocabulary target words. CharD-ecoderLSTM is trained like LSTM-based LM but on character-level of every word. If the target word is *music*, the input seq is [<START>, m, u, s, i, c], and the target seq is [m, u, s, i, c, <END>].

³⁷ Luong and Manning 2016

The test q2 result is 99.29792465574434 BLEU.

The final score of the full test is 24.660168424600386 (which is a lot more than the score required for full points).

- Arora, S., Li, Y., Liang, Y., Ma, T., and Risteski, A. (2018). Linear algebraic structure of word senses, with applications to polysemy. *Transactions of the Association for Computational Linguistics*, 6(o):483–495.
- Brown, P. F., Della Pietra, S. A., Della Pietra, V. J., and Mercer, R. L. (1993). The mathematics of statistical machine translation: Parameter estimation. *Computational Linguistics*, 19(2):263–311.
- Chen, D., Bolton, J., and Manning, C. D. (2016). A thorough examination of the CNN/daily mail reading comprehension task. In *Proceedings of the 54th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 2358–2367, Berlin, Germany. Association for Computational Linguistics.
- Chen, D., Fisch, A., Weston, J., and Bordes, A. (2017). Reading Wikipedia to answer open-domain questions. In *Proceedings of the 55th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 1870–1879, Vancouver, Canada. Association for Computational Linguistics.
- Chen, D. and Manning, C. (2014). A fast and accurate dependency parser using neural networks. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 740–750, Doha, Qatar. Association for Computational Linguistics.
- Cho, K., van Merriënboer, B., Gulcehre, C., Bahdanau, D., Bougares, F., Schwenk, H., and Bengio, Y. (2014). Learning phrase representations using RNN encoder–decoder for statistical machine translation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1724–1734, Doha, Qatar. Association for Computational Linguistics.

- Collobert, R. and Weston, J. (2008). A unified architecture for natural language processing: Deep neural networks with multitask learning. In *Proceedings of the 25th International Conference on Machine Learning, ICML '08*, page 160–167, New York, NY, USA. Association for Computing Machinery.
- Conneau, A., Schwenk, H., Barrault, L., and Lecun, Y. (2017). Very deep convolutional networks for text classification. In *Proceedings of the 15th Conference of the European Chapter of the Association for Computational Linguistics: Volume 1, Long Papers*, pages 1107–1116, Valencia, Spain. Association for Computational Linguistics.
- Devlin, J., Chang, M.-W., Lee, K., and Toutanova, K. (2019). Bert: Pre-training of deep bidirectional transformers for language understanding.
- Glorot, X. and Bengio, Y. (2010). Understanding the difficulty of training deep feedforward neural networks. In Teh, Y. W. and Titterton, M., editors, *Proceedings of the Thirteenth International Conference on Artificial Intelligence and Statistics*, volume 9 of *Proceedings of Machine Learning Research*, pages 249–256, Chia Laguna Resort, Sardinia, Italy. PMLR.
- He, K., Zhang, X., Ren, S., and Sun, J. (2016). Deep residual learning for image recognition. In *2016 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 770–778.
- Hochreiter, S. and Schmidhuber, J. (1997). Long short-term memory. *Neural Computation*, 9(8):1735–1780.
- Huang, C.-Z. A., Vaswani, A., Uszkoreit, J., Simon, I., Hawthorne, C., Shazeer, N., Dai, A. M., Hoffman, M. D., Dinculescu, M., and Eck, D. (2019). Music transformer. In *International Conference on Learning Representations*.
- Huang, E., Socher, R., Manning, C., and Ng, A. (2012). Improving word representations via global context and multiple word prototypes. In *Proceedings of the 50th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers)*, pages 873–882, Jeju Island, Korea. Association for Computational Linguistics.
- Huang, G., Liu, Z., v. d. Maaten, L., and Weinberger, K. Q. (2017). Densely connected convolutional networks. In *2017 IEEE Conference on Computer Vision and Pattern Recognition (CVPR)*, pages 2261–2269.

- Ioffe, S. and Szegedy, C. (2015). Batch normalization: Accelerating deep network training by reducing internal covariate shift. In *Proceedings of the 32nd International Conference on International Conference on Machine Learning - Volume 37, ICML'15*, page 448–456. JMLR.org.
- Kim, Y. (2014). Convolutional neural networks for sentence classification. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1746–1751, Doha, Qatar. Association for Computational Linguistics.
- Kim, Y., Jernite, Y., Sontag, D., and Rush, A. M. (2016). Character-aware neural language models. In *Proceedings of the Thirtieth AAAI Conference on Artificial Intelligence, AAAI'16*, page 2741–2749. AAAI Press.
- Kingma, D. P. and Ba, J. (2014). Adam: A method for stochastic optimization. *arXiv preprint arXiv:1412.6980*.
- Luong, M.-T. and Manning, C. D. (2016). Achieving open vocabulary neural machine translation with hybrid word-character models.
- Luong, M.-T., Socher, R., and Manning, C. D. (2013). Better word representations with recursive neural networks for morphology. In *Proceedings of the Seventeenth Conference on Computational Natural Language Learning*, pages 104–113.
- McCann, B., Bradbury, J., Xiong, C., and Socher, R. (2017). Learned in translation: Contextualized word vectors. In Guyon, I., Luxburg, U. V., Bengio, S., Wallach, H., Fergus, R., Vishwanathan, S., and Garnett, R., editors, *NeurIPS*, volume 30.
- Mikolov, T., Sutskever, I., Chen, K., Corrado, G. S., and Dean, J. (2013). Distributed representations of words and phrases and their compositionality. In Burges, C. J. C., Bottou, L., Welling, M., Ghahramani, Z., and Weinberger, K. Q., editors, *NeurIPS*, pages 3111–3119.
- Nivre, J. (2003). An efficient algorithm for projective dependency parsing. In *Proceedings of the Eighth International Conference on Parsing Technologies*, pages 149–160, Nancy, France.
- Parmar, N., Vaswani, A., Uszkoreit, J., Kaiser, L., Shazeer, N., Ku, A., and Tran, D. (2018). Image transformer. In Dy, J. and Krause, A., editors, *Proceedings of the 35th International Conference on Machine Learning*, volume 80 of *Proceedings of Machine Learning Research*, pages 4055–4064. PMLR.

- Pascanu, R., Mikolov, T., and Bengio, Y. (2013). On the difficulty of training recurrent neural networks. In *Proceedings of the 30th International Conference on International Conference on Machine Learning - Volume 28, ICML'13*, page III-1310-III-1318. JMLR.org.
- Pennington, J., Socher, R., and Manning, C. (2014). Glove: Global vectors for word representation. In *Proceedings of the 2014 Conference on Empirical Methods in Natural Language Processing (EMNLP)*, pages 1532-1543, Doha, Qatar. Association for Computational Linguistics.
- Peters, M. E., Ammar, W., Bhagavatula, C., and Power, R. (2017). Semi-supervised sequence tagging with bidirectional language models. In *ACL (1)*, pages 1756-1765.
- Peters, M. E., Neumann, M., Iyyer, M., Gardner, M., Clark, C., Lee, K., and Zettlemoyer, L. (2018). Deep contextualized word representations. In *NAACL*, pages 2227-2237.
- Rohde, D. L., Gonnerman, L. M., and Plaut, D. C. (2005). An improved model of semantic similarity based on lexical co-occurrence.
- Seo, M., Kembhavi, A., Farhadi, A., and Hajishirzi, H. (2016). Bidirectional attention flow for machine comprehension.
- Srivastava, N., Hinton, G., Krizhevsky, A., Sutskever, I., and Salakhutdinov, R. (2014). Dropout: A simple way to prevent neural networks from overfitting. *Journal of Machine Learning Research*, 15:1929-1958.
- Vaswani, A., Shazeer, N., Parmar, N., Uszkoreit, J., Jones, L., Gomez, A. N., Kaiser, L. u., and Polosukhin, I. (2017). Attention is all you need. In *NeurIPS*, volume 30.