

Master in Artificial Intelligence

Advanced Human Language Technologies Word Classification

Classification
Task with
Neural
Networks

Classification
Tasks in NLP

Conclusions



Outline

- 1 Classification Task with Neural Networks
 - Classification setup and notation
 - Softmax Classifier
 - Softmax with trainable Word Vectors
 - Neural Networks
- 2 Classification Tasks in NLP
 - Named Entity Recognition (NER)
 - Word-window Classification
 - Stochastic Gradient Descent
 - Other considerations
 - Beyond Word-window Classification
- 3 Conclusions

Classification
Task with
Neural
Networks

Classification
Tasks in NLP

Conclusions

Outline

- 1 Classification Task with Neural Networks
 - Classification setup and notation
 - Softmax Classifier
 - Softmax with trainable Word Vectors
 - Neural Networks
- 2 Classification Tasks in NLP
 - Named Entity Recognition (NER)
 - Word-window Classification
 - Stochastic Gradient Descent
 - Other considerations
 - Beyond Word-window Classification
- 3 Conclusions

Classification
Task with
Neural
Networks

Classification setup
and notation

Classification
Tasks in NLP

Conclusions

Classification setup and notation

- Generally we have a training dataset consisting of samples
 - $\{x_i, y_i\}_{i=1}^N$
- x_i are inputs, e.g. words (indices or vectors), sentences, documents, etc
 - Dimension d .
- y_i are labels (one of C classes) we try to predict, for example:
 - classes: sentiment, named entities, buy/sell decision
 - other words
 - later: multi-word sequences

Classification setup and notation (II)

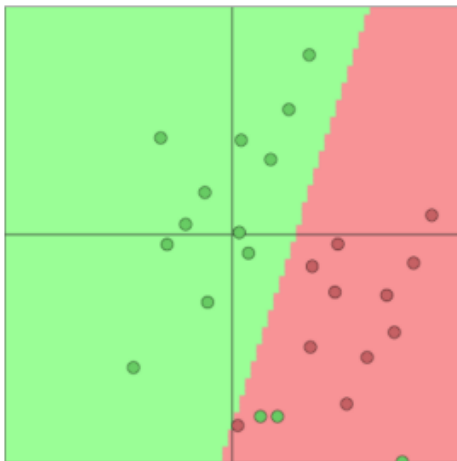


Figure: Simple illustration case: Fixed 2D word vectors to classify. Using softmax/logistic regression. Linear decision boundary.

Classification
Task with
Neural
Networks

Classification setup
and notation

Classification
Tasks in NLP

Conclusions

Outline

1 Classification Task with Neural Networks

- Classification setup and notation
- **Softmax Classifier**
- Softmax with trainable Word Vectors
- Neural Networks

2 Classification Tasks in NLP

- Named Entity Recognition (NER)
- Word-window Classification
- Stochastic Gradient Descent
- Other considerations
- Beyond Word-window Classification

3 Conclusions

Classification
Task with
Neural
Networks

Softmax Classifier

Classification
Tasks in NLP

Conclusions

Softmax Classifier

- Training Data:
- Traditional ML approach:
 - train (i.e. set) softmax/logistic regression weights $W \in \mathbb{R}^{C \times d}$ to determine a decision boundary (hyperplane)
- Method: For each x , predict:

$$p(y|x; \theta) = \frac{e^{(W_y \times x)}}{\sum_{c=1}^C e^{(W_c \times x)}}$$

Softmax Classifier (II)

$$p(y|x; \theta) = \frac{e^{(W_y \times x)}}{\sum_{c=1}^C e^{(W_c \times x)}}$$

We can tease apart the prediction function into two steps:

- 1 Take the y^{th} row of W and multiply that row with x :
 $W_y \times x = \sum W_{y_i} x_{i_{I=1}} = f_y$ Compute all f_c for $c = 1, \dots, c$
- 2 Apply softmax function to get the normalised probability:

$$p(y|x; \theta) = \frac{e^{f_y}}{\sum_{c=1}^C e^{f_c}} = \text{softmax}(f_y)$$

Cross-entropy loss

- For each training example (x, y) , our objective is to maximise the probability of the correct class y
- This is equivalent to minimising the negative log probability of that class:

$$-\log p(y|x; \theta) = -\log\left(\frac{e^{f_y}}{\sum_{c=1}^C e^{f_c}}\right)$$

- Using log probability converts our objective function to sums, which is easier to work with on paper and in implementation.

Cross-entropy loss (II)

- Concept of “cross entropy” is from information theory
- Let the true probability distribution be p
- Let our computed model probability be q
- The cross entropy is:

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

- Assuming a ground truth (or true or gold or target) probability distribution that is 1 at the right class and 0 everywhere else: $p = [0, \dots, 0, 1, 0, \dots, 0]$ then:
- Because of one-hot p , the only term left is the negative log probability of the true class

Cross-entropy loss (III)

- Cross entropy loss function over full dataset $x_i, y_{i=1}^N$

$$J(\theta) = \frac{1}{N} \cdot \sum_{i=1}^N -\log\left(\frac{e^{f_{y_i}}}{\sum_{c=1}^C e^{f_c}}\right)$$

- In general:

$$\theta = \begin{bmatrix} W_1 \\ \dots \\ W_C \end{bmatrix} = W \in \mathbb{R}^{C \cdot d}$$

- So we only update the decision boundary via:

$$\nabla J(\theta) = \begin{bmatrix} \nabla W_1 \\ \dots \\ \nabla W_C \end{bmatrix} \in \mathbb{R}^{C \cdot d}$$

Outline

1 Classification Task with Neural Networks

- Classification setup and notation
- Softmax Classifier
- **Softmax with trainable Word Vectors**
- Neural Networks

2 Classification Tasks in NLP

- Named Entity Recognition (NER)
- Word-window Classification
- Stochastic Gradient Descent
- Other considerations
- Beyond Word-window Classification

3 Conclusions

Softmax with trainable Word Vectors

- Commonly in NLP deep learning:
 - We learn both W and word vectors x
 - We learn both conventional parameters and representations
 - The word vectors re-represent one-hot vectors (move them around in an intermediate layer vector space) for easy classification with a (linear) softmax classifier

$$\nabla_{\theta} J(\theta) = \begin{bmatrix} \nabla W_1 \\ \dots \\ \nabla W_d \\ \nabla x_{word_1} \\ \dots \\ \nabla x_{word_v} \end{bmatrix} \in \mathbb{R}^{C \cdot d + V \cdot d}$$

! But $V \cdot d$ is big!

Outline

1 Classification Task with Neural Networks

- Classification setup and notation
- Softmax Classifier
- Softmax with trainable Word Vectors
- **Neural Networks**

2 Classification Tasks in NLP

- Named Entity Recognition (NER)
- Word-window Classification
- Stochastic Gradient Descent
- Other considerations
- Beyond Word-window Classification

3 Conclusions

Neural Network Classifier

- Softmax (\approx logistic regression) alone not very powerful
- Softmax gives only linear decision boundaries This can be quite limiting: Unhelpful when a problem is complex
- Solution: Neural Networks can learn much more complex functions and nonlinear decision boundaries

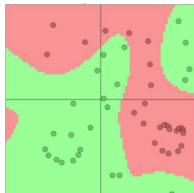
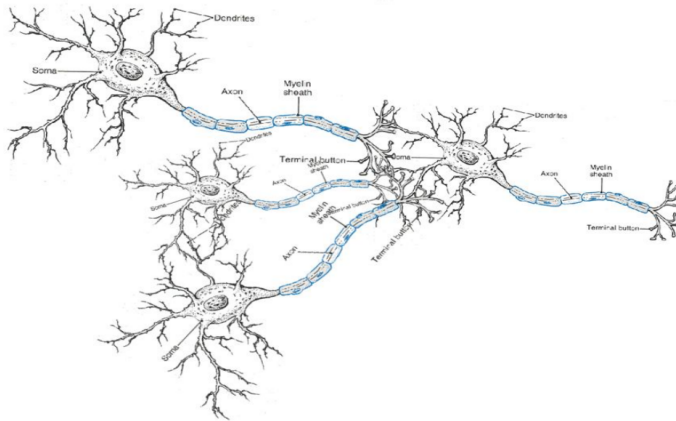


Figure: Non-linear decision boundary

Neural Computation



Classification
Task with
Neural
Networks

Neural Networks

Classification
Tasks in NLP

Conclusions

A Neuron

A neuron can be a binary logistic regression unit

- f = nonlinear activation function (e.g. sigmoid), w = weights, b = bias, h = hidden, x = inputs

$$h_{w,b}(x) = f(w^T \cdot x + b)$$

$$f(z) = \frac{1}{1 + e^{-z}}$$

- b = We can have an “always on” feature, which gives a class prior, or separate it out, as a bias term
- w , b are the parameters of this neuron i.e., this logistic regression model

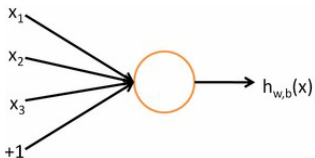
A Neuron

Classification
Task with
Neural
Networks

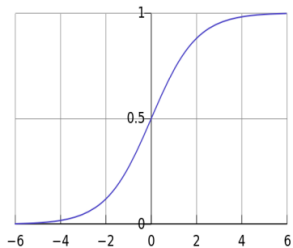
Neural Networks

Classification
Tasks in NLP

Conclusions



(a) Single Neuron



(b) Sigmoid

Neural Network

- A neural network = running several logistic regressions at the same time
- If we feed a vector of inputs through a bunch of logistic regression functions, then we get a vector of outputs ...

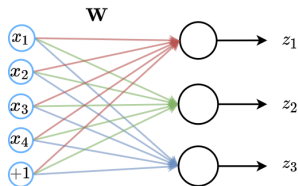


Figure: Neural Network with 3 neurons

But we don't have to decide ahead of time what variables these logistic regressions are trying to predict!

Neural Network (II)

... which we can feed into another logistic regression function

It is the loss function that will direct what the intermediate hidden variables should be, so as to do a good job at predicting the targets for the next layer, etc.

And if we add more layers... Before we know it, we have a multi-layer neural network....

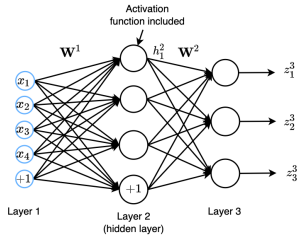


Figure: Multi-layer Neural Network

Neural Network (III)

In a Multi-layer Perceptron (MLP)

$$h(c_1|x; \theta) = \sigma(z) = \frac{1}{1 + e^{-z}}$$

z is no longer linear

$$h(c_k|x; \theta) = \frac{e^{z_k}}{\sum_j e^{z_j}}$$

Then:

$$h_1^2 = f(W_{11}^1 \cdot x_1 + W_{12}^1 \cdot x_2 + W_{13}^1 \cdot x_3 + b_1^1)$$

$$h_2^2 = f(W_{21}^1 \cdot x_1 + W_{22}^1 \cdot x_2 + W_{23}^1 \cdot x_3 + b_2^1)$$

The activation function is applied element-wise

$$f([z_1^2, z_2^2, z_3^2]) = [f(z_1^2), f(z_2^2), f(z_3^2)]$$

The need of Non-linearity

- Without non-linearities, deep neural networks can't do anything more than a linear transform
- Extra layers could just be compiled down into a single linear transform: $W^1 \cdot W^2 \cdot x = W \cdot x$
- More layers approximate more complex functions

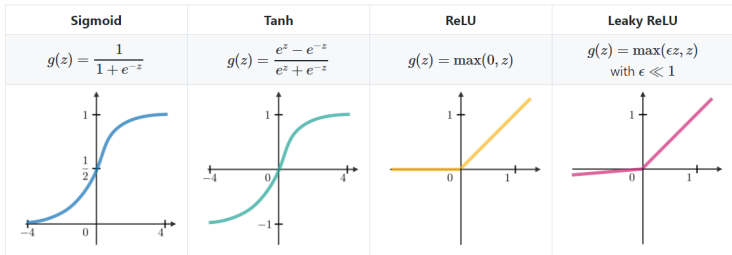


Figure: Common activation functions

You can "play" with them in the TensorFlow Playground

Outline

- 1 Classification Task with Neural Networks
 - Classification setup and notation
 - Softmax Classifier
 - Softmax with trainable Word Vectors
 - Neural Networks
- 2 Classification Tasks in NLP
 - Named Entity Recognition (NER)
 - Word-window Classification
 - Stochastic Gradient Descent
 - Other considerations
 - Beyond Word-window Classification
- 3 Conclusions

Classification
Task with
Neural
Networks

Classification
Tasks in NLP

Conclusions

Outline

- 1 Classification Task with Neural Networks
 - Classification setup and notation
 - Softmax Classifier
 - Softmax with trainable Word Vectors
 - Neural Networks
- 2 Classification Tasks in NLP
 - **Named Entity Recognition (NER)**
 - Word-window Classification
 - Stochastic Gradient Descent
 - Other considerations
 - Beyond Word-window Classification
- 3 Conclusions

Classification
Task with
Neural
Networks

Classification
Tasks in NLP
Named Entity
Recognition (NER)

Conclusions

Named Entity Recognition

Classification
Task with
Neural
Networks

Classification
Tasks in NLP

Named Entity
Recognition (NER)

Conclusions

- A named entity is anything that can be referred to with a proper name: a person, a location, an organization
- The task of named entity recognition (NER) is to find spans of text that constitute proper names and tag the type of the entity
- **PER** (person), **LOC** (location), **ORG** (organization) and others (See CONLL and Ontonotes)

Named Entity Recognition (II)

- The task: Find and classify entities in the text:

The *European Commission* [**ORG**] said on Thursday it disagreed with *German* [**MISC**] advice. Only *France* [**LOC**] and *Britain* [**LOC**] backed *Fischler* [**PER**]’s proposal. “What we have to be extremely careful of is how other countries are going to take Germany’s lead,” *Welsh National Farmers’ Union* [**ORG**] (*NFU* [**ORG**]) chairman *John Lloyd Jones* [**PER**] said on *BBC* [**ORG**] radio.

- Simple approach:
 - We predict entities by classifying words in context and then extracting entities as word sub-sequences (token by token).
 - Classes can be encoded using **IO** (O, ORG, PER, ...); **BIO** (O, B-ORG, I-ORG, B-PER, I-PER, ...); ...

Named Entity Recognition (III)

- Possible purposes:
 - Tracking mentions of particular entities in documents
 - For question answering, answers are usually named entities
 - A lot of wanted information is really associations between named entities
 - The same techniques can be extended to other slot-filling classifications
- Often followed by Named Entity Linking/Canonicalization into Knowledge Base

Challenges in NER

- Hard to work out boundaries of entity
Ex: In "First National Bank Donates 2 Vans to Future School of Fort Smith"
? Is the first entity "First National Bank" or "National Bank"
- Hard to know if something is an entity
? Is there a school called "Future School" or is it a future school?
- Hard to know class of unknown/novel entity:
Ex: ...To find out more about Zig Ziglar and read features by other Creators Syndicate writers and...
? What class is "Zig Ziglar"? (A person.)
- Entity class is ambiguous and depends on context
Ex: ...where Larry Ellison and Charles Schwab can live discreetly amongst wooded estates. And...
? "Charles Schwab" is PER not ORG here!

Outline

- 1 Classification Task with Neural Networks
 - Classification setup and notation
 - Softmax Classifier
 - Softmax with trainable Word Vectors
 - Neural Networks
- 2 Classification Tasks in NLP
 - Named Entity Recognition (NER)
 - **Word-window Classification**
 - Stochastic Gradient Descent
 - Other considerations
 - Beyond Word-window Classification
- 3 Conclusions

Classification
Task with
Neural
Networks

Classification
Tasks in NLP
Word-window
Classification

Conclusions

Word-window Classification

- Idea: classify a word in its context window of neighboring words.
 - Ex: “Museums in Paris are amazing”
to classify whether or not the center word “Paris” is a named-entity
- For example, Named Entity Classification of a word in context:
 - Person, Location, Organization, None
- A simple way to classify a word in context might be to average the word vectors in a window and to classify the average vector
 - Problem: that would **lose position information**

Word-window Classification (II)

- Train softmax classifier to classify a center word by taking the concatenation of words surrounding in a window

Ex: Classify “Paris” in the context of this sentence with window length 2:

... museums in Paris are amazing ...



$$X_{\text{window}} = [x_{\text{museums}} \quad x_{\text{in}} \quad x_{\text{Paris}} \quad x_{\text{are}} \quad x_{\text{amazing}}]^T$$

- Resulting vector $w_{\text{window}} = x \in \mathbb{R}^{5 \cdot d}$, a column vector!

Word-window Classification (III)

- With $x = x_{window}$ we can use the softmax classifier

$$p(y|x; \theta) = \frac{e^{z_y}}{\sum_j e^{z_j}} = \frac{e^{W_y \cdot x}}{\sum_j e^{W_j \cdot x}}$$

- With cross-entropy loss:

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N -\log\left(\frac{e^{z_{y_i}}}{\sum_{j=1}^C e^{z_j}}\right)$$

- How do you update the word vectors?
 - Short answer: Just take derivatives and optimize

Word-window Classification - Binary Logistic Classifier

Classification
Task with
Neural
Networks

Classification
Tasks in NLP

Word-window
Classification

Conclusions

- Train logistic classifier on hand-labeled data to classify center word {yes / no} for each class based on a concatenation of word vectors in a window

Ex: Classify “Paris” as +/– location in context of sentence with window length 2:

... museums in Paris are amazing ...



$$X_{\text{window}} = [x_{\text{museums}} \quad x_{\text{in}} \quad x_{\text{Paris}} \quad x_{\text{are}} \quad x_{\text{amazing}}]^T$$

Word-window Classification - Binary Logistic Classifier (II)

- We do supervised training and want high score if it's a location

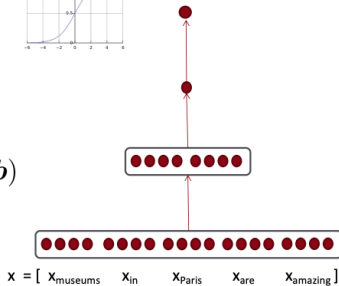
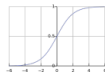
$$J_t(\theta) = \sigma(s) = \frac{1}{1 + e^{-s}}$$

predicted model
probability of class

$$s = \mathbf{u}^T \mathbf{h}$$

$$\mathbf{h} = f(\mathbf{W}\mathbf{x} + \mathbf{b})$$

\mathbf{x} (input)



Outline

- 1 Classification Task with Neural Networks
 - Classification setup and notation
 - Softmax Classifier
 - Softmax with trainable Word Vectors
 - Neural Networks
- 2 Classification Tasks in NLP
 - Named Entity Recognition (NER)
 - Word-window Classification
 - **Stochastic Gradient Descent**
 - Other considerations
 - Beyond Word-window Classification
- 3 Conclusions

Classification
Task with
Neural
Networks

Classification
Tasks in NLP
Stochastic Gradient
Descent

Conclusions

Stochastic Gradient Descent

- Update equation gradient descent:

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

$$J(\theta) = \frac{1}{N} \sum_{i=1}^N -\log\left(\frac{e^{z_{y_i}}}{\sum_{j=1}^C e^{z_j}}\right)$$

- Update equation stochastic gradient descent (SGD):

$$\theta^{new} = \theta^{old} - \alpha \cdot \nabla_{\theta} J_i(\theta; x_i, y_i)$$

- 1 Randomly shuffle dataset
 - 2 For every training sample (i) in the dataset- i apply the update rule
- We can also update the parameter every minibatch, which means a few number of samples.

Gradients

- Given a function with 1 output and n inputs:

$$f(x) = f(x_1, \mathbf{x}_2, \dots, x_n)$$

- Its gradient is a vector of partial derivatives with respect to each input: $\frac{\partial f}{\partial x} = [\frac{\partial f}{\partial x_1}, \frac{\partial f}{\partial x_2}, \dots, \frac{\partial f}{\partial x_n}]$
- Now given a function f with m outputs and n inputs, its Jacobian is:

$$\frac{\partial f}{\partial x} = \begin{bmatrix} \frac{\partial f_1}{\partial x_1}(\mathbf{x}) & \frac{\partial f_1}{\partial x_2}(\mathbf{x}) & \dots & \frac{\partial f_1}{\partial x_n}(\mathbf{x}) \\ \frac{\partial f_2}{\partial x_1}(\mathbf{x}) & \frac{\partial f_2}{\partial x_2}(\mathbf{x}) & \dots & \frac{\partial f_2}{\partial x_n}(\mathbf{x}) \\ \vdots & \vdots & \ddots & \vdots \\ \frac{\partial f_m}{\partial x_1}(\mathbf{x}) & \frac{\partial f_m}{\partial x_2}(\mathbf{x}) & \dots & \frac{\partial f_m}{\partial x_n}(\mathbf{x}) \end{bmatrix}$$

Gradients (II)

- Let's find $\frac{\partial s}{\partial b}$ ¹

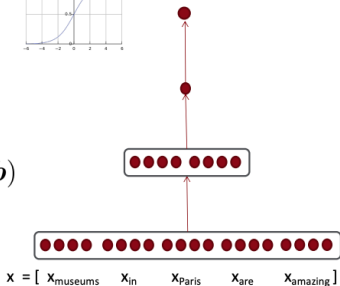
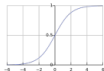
$$J_t(\theta) = \sigma(s) = \frac{1}{1 + e^{-s}}$$

predicted model
probability of class

$$s = \mathbf{u}^T \mathbf{h}$$

$$\mathbf{h} = f(\mathbf{W}\mathbf{x} + \mathbf{b})$$

\mathbf{x} (input)



¹In actuality, we care about the gradient of the loss J_i but we will compute the gradient of the score for simplicity

Gradients (III)

- We apply the chain rule

Ex: Derivative of s respect to b :

$$s = u^T \cdot h \quad h = f(z) \quad z = W \cdot x + b$$

$$\frac{\partial s}{\partial b} = \frac{\partial s}{\partial h} \cdot \frac{\partial h}{\partial z} \cdot \frac{\partial z}{\partial b}$$

Computational Graph

- Software represents our neural net equations as a graph
 - Source nodes: inputs
 - Interior nodes: operations
 - Edges pass along result of the operation

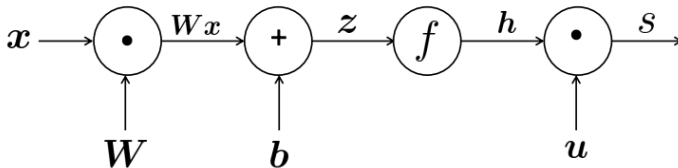


Figure: Forward Pass

Computational Graph (II)

- Then do the backward pass

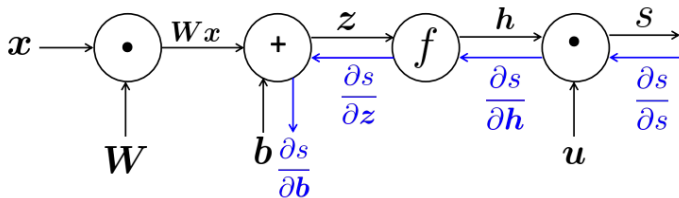
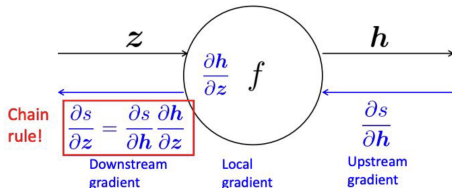


Figure: Backpropagation

Computational Graph (III)

- Backpropagation in a single node:
 - Node receives an “upstream gradient”
 - Goal is to pass on the correct “downstream gradient”
 - Each node has a local gradient
 - The gradient of its output with respect to its input



Outline

- 1 Classification Task with Neural Networks
 - Classification setup and notation
 - Softmax Classifier
 - Softmax with trainable Word Vectors
 - Neural Networks
- 2 Classification Tasks in NLP
 - Named Entity Recognition (NER)
 - Word-window Classification
 - Stochastic Gradient Descent
 - **Other considerations**
 - Beyond Word-window Classification
- 3 Conclusions

Classification
Task with
Neural
Networks

Classification
Tasks in NLP
Other considerations

Conclusions

Faster Activation Functions

- LReLU (Leaky Rectified Linear Unit - Leaky ReLU):
 - Modification ReLU that avoids the "dying ReLU" problem, where neurons stop firing due to a zero output
 - Introduces a small, non-zero slope for negative inputs
($f(x) = \max(\alpha \cdot x, x)$)
 - Can avoid the vanishing gradient problem, which can occur when using sigmoid or other saturating activation functions
 - Allows a small, non-zero gradient when the input is negative, which can prevent the gradient from becoming too small
 - This can lead to faster convergence and better accuracy in some cases.
- ELU (Exponential Linear Unit):
 - Avoids the "dying ReLU" problem and has a smooth output
- SELU (Scaled Exponential Linear Unit):
 - Self-normalizing activation function that can significantly improve the performance

Parameter Initialization

- Proper initialization of model parameters is crucial for effective training and **convergence**. Popular approaches include:
 - Random: e.g., uniform or normal distribution
 - **He**: scaled version of random initialization, designed for ReLU activations
 - **Xavier**: Scaled version of random initialization
 - Designed for sigmoid/tanh activations that have a linear region
 - Sets the variance of the weights to $Var(W_i) = \frac{2}{n_{in} + n_{out}}$, where n_{in} is the number of input neurons and n_{out} is the number of output neurons
 - **Glorot**: Combination of He and Xavier
 - Pre-trained **word-embeddings**: Using pre-trained word-embeddings, such as GloVe or Word2Vec, to initialize the embedding layer of the model
- In general, we initialize the weights to small random values and biases to 0 in the hidden layers.

Optimizers: SGD

- SGD is a commonly used optimizer for neural network training
 - The method iteratively adjusts the model's parameters by computing the **gradient** of the **loss function** with respect to the parameters for a **randomly selected** sample (**stochastic**) of the training data.
- **Simple and efficient.**
- However, getting good results often requires **hand-tuning** the **learning rate**
 - Learning rate determines the **step size** that the optimizer takes to update the weights and biases
 - Inappropriate values can cause the optimizer to converge too slowly or too quickly

Adaptive Optimization Algorithms

- They **scale the learning rate** of each parameter based on the accumulated gradient history
- This provides a **per-parameter learning rate** that can perform well in settings with high curvature, noisy gradients, and sparse data
- Popular adaptive optimizers include:
 - Adagrad: divides the learning rate by the sum of the squares of past gradients
 - RMSprop: exponentially decays the average of past squared gradients to normalize the learning rate
 - Adam: combines the benefits of Adagrad and RMSprop by using both first and second moments of past gradients
 - SparseAdam: similar to Adam, but optimized for sparse gradients
- Each optimizer has its own strengths and weaknesses

Learning Rate

- In NLP models, the learning rate plays a crucial role in training and convergence.
- Learning rate determines the size of the step the optimizer takes in the direction of the negative gradient to update the weights and biases of the model.
 - High LR can cause the model to overshoot the optimal point and **diverge**
 - Low LR can result in the model taking too long to converge or getting stuck in **local minima**
- A while a low learning rate can result in the model taking too long to converge or getting stuck in local minima
- NLP models can benefit from using:
 - Learning rate **schedules**
 - **Adaptive** optimization algorithms (previous slide)
- **Fine-tuning pre-trained models** for downstream NLP tasks may require using a smaller learning rate than for training the original model

Regularization

- Regularization (largely) prevents overfitting when we have a lot of features (or later a very powerful/deep model)
- L1 regularization: adds the sum of absolute values of weights to the loss function

$$L_{reg} = L + \lambda \sum_{i=1}^n |w_i|$$

- L2 regularization: adds the sum of squares of weights to the loss function

$$L_{reg} = L + \lambda \sum_{i=1}^n w_i^2$$

Regularization (II)

Classification
Task with
Neural
Networks

Classification
Tasks in NLP

Other considerations

Conclusions

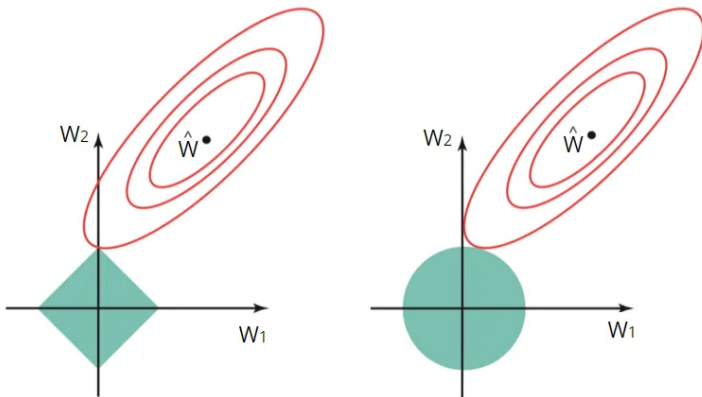
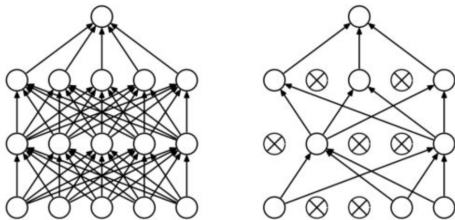


Figure: Representation of the effect of L1 (left) and L2 (right) Regularization. Red lines represent local minima. The red area represents optimal values for the regularization term.

Regularization (III)

- Dropout: randomly sets a fraction of the units to zero during training

$$y = \begin{cases} x & \text{with probability } 1 - p \\ 0 & \text{with probability } p \end{cases}$$



Outline

- 1 Classification Task with Neural Networks
 - Classification setup and notation
 - Softmax Classifier
 - Softmax with trainable Word Vectors
 - Neural Networks
- 2 Classification Tasks in NLP
 - Named Entity Recognition (NER)
 - Word-window Classification
 - Stochastic Gradient Descent
 - Other considerations
 - **Beyond Word-window Classification**
- 3 Conclusions

Classification
Task with
Neural
Networks

Classification
Tasks in NLP

Beyond
Word-window
Classification

Conclusions

Convolution Filters for Word Classification

- Word window classification processes each window independently, not leveraging the fact that windows are sharing words.
- We can extend the word window concept to flexible-sized windows i.e., convolution filters, by applying the filters to regions of input space, represented by:

$$f_k(w_1, w_2, \dots, w_n) \rightarrow l$$

where f_k is the filter and w 's are the words in the window and l is the label.

- Similar to word-window, convolution filters also classify a word based on its local neighbors but with overlappable and variable-sized (1 to n) windows.

Convolution Filters Vs Word Window Classification

- Unlike word-window where each window size requires different parameters, convolution filters share parameters across windows, reducing the number of parameters required and making the model more efficient.
- Convolution filters can handle words near the start and end of a sentence by padding the input.
- While the simple sliding-window approach considers only local context and assumes independence between windows, convolutions consider more global, unbounded context via pooling, represented by:

$$y = g\left(\sum_{i=1}^m f_k(w_{i:i+n})\right)$$

where g is the global pooling function, f_k is the filter, w 's are the words, i is the start of the window

Convolution Operation

- Each filter slides over the input sequence, applying the same linear transformation to each window it covers, represented by:

$$y_i = f_k(w_{i:i+n}; \theta)$$

where θ represents the shared parameters.

- The output (after all filters are applied) is then combined and passed to the next layer of the neural network.
- This operation enables the neural network to maintain a low number of parameters while still learning hierarchical features.

Pooling Mechanism

- After convolutions are applied, it's common to use pooling operations to aggregate information.
- Pooling helps to reduce dimensionality and control overfitting.
- Max-pooling is a commonly used technique, where the maximum value over a certain window (or the whole sequence) is taken.

$$Y_i = g(y_i, y_{i+1}, \dots, y_{i+m})$$

where Y_i represents the pooled output and y 's are the convolution outputs and m is the pooling window.

Convolution Filters As Feature Detectors

- Convolution filters can learn to detect local features, similar to how they detect lines or edges in image processing.
- In the context of NLP, these locally detected features might represent significant n-grams or phrases.
- For instance, a filter might learn to detect negations or other semantic and syntactic features.
- The process can be represented as:

$$f(w_{i:i+n}) = \max(0, w_{i:i+n} * \theta + b)$$

Multi-Layer Convolutions

- Multi-Layer convolutional neural networks (CNNs) stack multiple convolutional layers (sometimes interspersed with pooling layers) to learn an increasingly abstract hierarchy of features.
- In NLP tasks such as Named Entity Recognition (NER), a multi-layer CNN might first learn low-level lexical features (e.g., character or word sequences), then more abstract syntactic/semantic features (e.g., part-of-speech patterns or name-entity type patterns).
- With each layer, the model learns more complex and abstract features, often leading to improved performance.

Multi-Layer Convolutions (II)

- For a multi-layer CNN, the output of the first layer becomes the input to the second layer:

$$Y^1[i] = W^1 * X[i : s^1 : i + k^1] + b^1$$

$$Y^2[j] = W^2 * Y^1[j : s^2 : j + k^2] + b^2$$

- W^1 and W^2 are the weights, b^1 and b^2 are the biases
- X is the input sequence
- i and j represent the index of the input or output signal
- The $*$ operator represents the convolution operation
- The stride s determines the step size for moving the filter window along the input sequence.
- The kernel size k represents the size of the filter window. It can be adjusted to capture different scales of information.

Outline

- 1 Classification Task with Neural Networks
 - Classification setup and notation
 - Softmax Classifier
 - Softmax with trainable Word Vectors
 - Neural Networks
- 2 Classification Tasks in NLP
 - Named Entity Recognition (NER)
 - Word-window Classification
 - Stochastic Gradient Descent
 - Other considerations
 - Beyond Word-window Classification
- 3 Conclusions

Classification
Task with
Neural
Networks

Classification
Tasks in NLP

Conclusions

Conclusions

Classification
Task with
Neural
Networks

Classification
Tasks in NLP

Conclusions

- Classification Tasks can successfully be addressed with Neural Networks because they are able to capture non-linearity
- Named Entity Recognition can be addressed as a Classification Task
- Deep Learning Computation is complex and full of tricks and details