

IN5550: Neural Methods in
Natural Language Processing
Sub-lecture 3.4
Computation graphs

Andrey Kutuzov

University of Oslo

7 February 2023





- 1 Computation graphs
- 2 Neural network toolkits
- 3 Next week trailer



Neural networks are just differentiable parameterized functions:



Neural networks are just differentiable parameterized functions:

- ▶ so we can use the same gradient-based optimization that we already saw (SGD, etc)



Neural networks are just differentiable parameterized functions:

- ▶ so we can use the same gradient-based optimization that we already saw (SGD, etc)
- ▶ But we have to **calculate the gradient**
 - ▶ potentially error prone and needs to be repeated for every change in model configuration/number of hidden layers/activation function/loss function, etc



Neural networks are just differentiable parameterized functions:

- ▶ so we can use the same gradient-based optimization that we already saw (SGD, etc)
- ▶ But we have to **calculate the gradient**
 - ▶ potentially error prone and needs to be repeated for every change in model configuration/number of hidden layers/activation function/loss function, etc
- ▶ The good news is that we can use **reverse-mode automatic differentiation (i.e. back-propagation)** to automatically calculate the gradients



Neural networks are just differentiable parameterized functions:

- ▶ so we can use the same gradient-based optimization that we already saw (SGD, etc)
- ▶ But we have to **calculate the gradient**
 - ▶ potentially error prone and needs to be repeated for every change in model configuration/number of hidden layers/activation function/loss function, etc
- ▶ The good news is that we can use **reverse-mode automatic differentiation (i.e. back-propagation)** to automatically calculate the gradients
- ▶ This means that instead of spending lots of your time deriving gradients by hand, you can use all that time to explore model options.



- ▶ These days, deep neural networks are usually trained using the **computation graph abstraction**:



- ▶ These days, deep neural networks are usually trained using the **computation graph abstraction**:
- ▶ representation of the process of computing a mathematical expression;



- ▶ These days, deep neural networks are usually trained using the **computation graph abstraction**:
- ▶ representation of the process of computing a mathematical expression;
- ▶ allows to **construct arbitrary deep and complex network architectures**;



- ▶ These days, deep neural networks are usually trained using the **computation graph abstraction**:
- ▶ representation of the process of computing a mathematical expression;
- ▶ allows to **construct arbitrary deep and complex network architectures**;
- ▶ computation is broken down into separate operations;

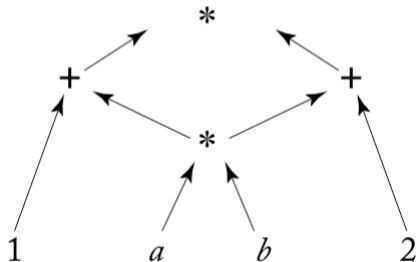


- ▶ These days, deep neural networks are usually trained using the **computation graph abstraction**:
- ▶ representation of the process of computing a mathematical expression;
- ▶ allows to **construct arbitrary deep and complex network architectures**;
- ▶ computation is broken down into separate operations;
- ▶ each operation and variable is a node in a graph:

Computation graphs



- ▶ These days, deep neural networks are usually trained using the **computation graph abstraction**:
- ▶ representation of the process of computing a mathematical expression;
- ▶ allows to **construct arbitrary deep and complex network architectures**;
- ▶ computation is broken down into separate operations;
- ▶ each operation and variable is a node in a graph:

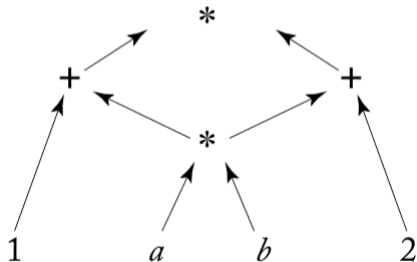


directed acyclic graph (DAG)

Computation graphs



- ▶ These days, deep neural networks are usually trained using the **computation graph abstraction**:
- ▶ representation of the process of computing a mathematical expression;
- ▶ allows to **construct arbitrary deep and complex network architectures**;
- ▶ computation is broken down into separate operations;
- ▶ each operation and variable is a node in a graph:



directed acyclic graph (DAG)

Question: **what does this graph produce when $a = 2, b = 1$?**



Training basics

- ▶ Similar to linear classifiers, deep NNs are **trained by gradient calculation** (SGD, etc)



Training basics

- ▶ Similar to linear classifiers, deep NNs are **trained by gradient calculation** (SGD, etc)
- ▶ Two main stages:



Training basics

- ▶ Similar to linear classifiers, deep NNs are **trained by gradient calculation** (SGD, etc)
- ▶ Two main stages:
 1. **Forward pass** (compute predictions for given inputs);
 - ▶ move forward along the graph from input to output and compute the loss.



Training basics

- ▶ Similar to linear classifiers, deep NNs are **trained by gradient calculation** (SGD, etc)
- ▶ Two main stages:
 1. **Forward pass** (compute predictions for given inputs);
 - ▶ move forward along the graph from input to output and compute the loss.
 2. **Backward pass** (compute gradients with respect to the loss)
 - ▶ move backwards along the graph from the resulting loss through all the layers to the input.



Training basics

- ▶ Similar to linear classifiers, deep NNs are **trained by gradient calculation** (SGD, etc)
- ▶ Two main stages:
 1. **Forward pass** (compute predictions for given inputs);
 - ▶ move forward along the graph from input to output and compute the loss.
 2. **Backward pass** (compute gradients with respect to the loss)
 - ▶ move backwards along the graph from the resulting loss through all the layers to the input.
- ▶ How exactly they are performed is dependent on which nodes are there in the graph.



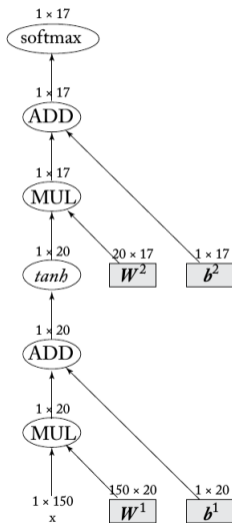
Training basics

- ▶ Similar to linear classifiers, deep NNs are **trained by gradient calculation** (SGD, etc)
- ▶ Two main stages:
 1. **Forward pass** (compute predictions for given inputs);
 - ▶ move forward along the graph from input to output and compute the loss.
 2. **Backward pass** (compute gradients with respect to the loss)
 - ▶ move backwards along the graph from the resulting loss through all the layers to the input.
- ▶ How exactly they are performed is dependent on which nodes are there in the graph.
- ▶ Operations, parameters and variables are nodes in a DAG.



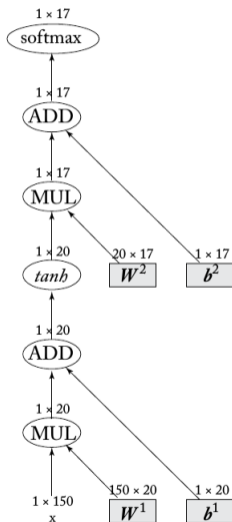
Training basics

- ▶ Similar to linear classifiers, deep NNs are **trained by gradient calculation** (SGD, etc)
- ▶ Two main stages:
 1. **Forward pass** (compute predictions for given inputs);
 - ▶ move forward along the graph from input to output and compute the loss.
 2. **Backward pass** (compute gradients with respect to the loss)
 - ▶ move backwards along the graph from the resulting loss through all the layers to the input.
- ▶ How exactly they are performed is dependent on which nodes are there in the graph.
- ▶ Operations, parameters and variables are nodes in a DAG.
- ▶ The whole graph should be **differentiable**.



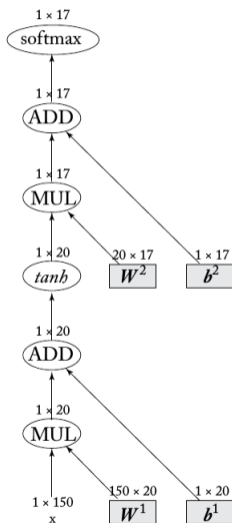
Think over before the Q&A session / lab session:

- How many total layers are there in this computation graph?



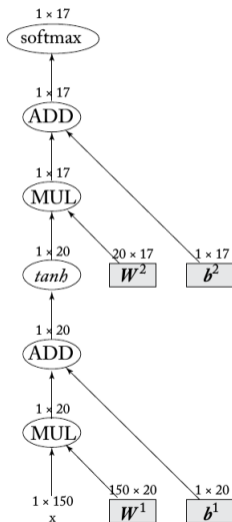
Think over before the Q&A session / lab session:

- ▶ How many total layers are there in this computation graph?
- ▶ How many hidden layers?



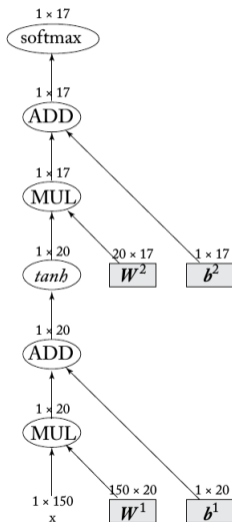
Think over before the Q&A session / lab session:

- ▶ How many total layers are there in this computation graph?
- ▶ How many hidden layers?
- ▶ What is the shape of its output?



Think over before the Q&A session / lab session:

- ▶ How many total layers are there in this computation graph?
- ▶ How many hidden layers?
- ▶ What is the shape of its output?
- ▶ **How can you describe what this architecture does?**



Think over before the Q&A session / lab session:

- ▶ How many total layers are there in this computation graph?
- ▶ How many hidden layers?
- ▶ What is the shape of its output?
- ▶ **How can you describe what this architecture does?**
- ▶ Any NLP task for which it can be used?



Forward pass

- ▶ Go over all nodes starting from the input;



Forward pass

- ▶ Go over all nodes starting from the input;
- ▶ apply the operations sequentially;



Forward pass

- ▶ Go over all nodes starting from the input;
- ▶ apply the operations sequentially;
- ▶ produce the output up to the final scalar **loss-node**.



Forward pass

- ▶ Go over all nodes starting from the input;
- ▶ apply the operations sequentially;
- ▶ produce the output up to the final scalar **loss-node**.

Loss calculation

- ▶ Loss functions are as a rule the same as with linear models;
- ▶ **cross-entropy** is arguably the dominant one: $L(\hat{y}, y) = -y \log(\hat{y}) - (1 - y)\log(1 - \hat{y})$
- ▶ NB: the output should be squashed by sigmoid or like.



Forward pass

- ▶ Go over all nodes starting from the input;
- ▶ apply the operations sequentially;
- ▶ produce the output up to the final scalar **loss-node**.

Loss calculation

- ▶ Loss functions are as a rule the same as with linear models;
- ▶ **cross-entropy** is arguably the dominant one: $L(\hat{y}, y) = -y \log(\hat{y}) - (1 - y) \log(1 - \hat{y})$
- ▶ NB: the output should be squashed by sigmoid or like.
- ▶ ...however, new loss functions are also being introduced:
 - ▶ e.g., Von Mises-Fisher loss [Kumar and Tsvetkov, 2018], etc.



Backward pass

- ▶ In simple linear classifiers, we compute the **partial derivatives of model parameters with respect to the loss...**



Backward pass

- ▶ In simple linear classifiers, we compute the **partial derivatives of model parameters with respect to the loss...**
- ▶ ...in multi-layer networks it is the same...



Backward pass

- ▶ In simple linear classifiers, we compute the **partial derivatives of model parameters with respect to the loss...**
- ▶ ...in multi-layer networks it is the same...
- ▶ ...but the number of parameters can be crazy big.



Backward pass

- ▶ In simple linear classifiers, we compute the **partial derivatives of model parameters with respect to the loss...**
- ▶ ...in multi-layer networks it is the same...
- ▶ ...but the number of parameters can be crazy big.
- ▶ Gradient is computed using a special algorithm:
- ▶ **back-propagation a.k.a reverse differentiation** [Rumelhart et al., 1986];



Backward pass

- ▶ In simple linear classifiers, we compute the **partial derivatives of model parameters with respect to the loss...**
- ▶ ...in multi-layer networks it is the same...
- ▶ ...but the number of parameters can be crazy big.
- ▶ Gradient is computed using a special algorithm:
- ▶ **back-propagation a.k.a reverse differentiation** [Rumelhart et al., 1986];
- ▶ propagating loss values backwards through all layers...
- ▶ ...religiously computing derivatives at each.



Backward pass

- ▶ *Back-propagation* is very efficient in terms of computation speed, but **very error-prone** if done manually...



Backward pass

- ▶ *Back-propagation* is very efficient in terms of computation speed, but **very error-prone** if done manually...

Neural network gradients

Andrew Ng

Backpropagation intuition (optional)

- ▶ ... fortunately, **automatic tools for gradient computation with *backprop*** do exist and are well-developed!



1 Computation graphs

2 Neural network toolkits

3 Next week trailer



Software libraries

No lack of open-source deep learning toolkits:



Software libraries

No lack of open-source deep learning toolkits:

- ▶ **Static graph construction** (compile then run):
 - ▶ **TensorFlow** by Google [Abadi et al., 2015] (<https://tensorflow.org>)



Software libraries

No lack of open-source deep learning toolkits:

- ▶ **Static graph construction** (compile then run):
 - ▶ **TensorFlow** by Google [Abadi et al., 2015] (<https://tensorflow.org>)
- ▶ **Dynamic graph construction** (run on the fly, a new graph for each sample):
 - ▶ **PyTorch** by Facebook [Paszke et al., 2017] (<https://pytorch.org/>)



Software libraries

No lack of open-source deep learning toolkits:

- ▶ **Static graph construction** (compile then run):
 - ▶ **TensorFlow** by Google [Abadi et al., 2015] (<https://tensorflow.org>)
- ▶ **Dynamic graph construction** (run on the fly, a new graph for each sample):
 - ▶ **PyTorch** by Facebook [Paszke et al., 2017] (<https://pytorch.org/>)
 - ▶ (**TensorFlow 2** added support for dynamic graphs as well)

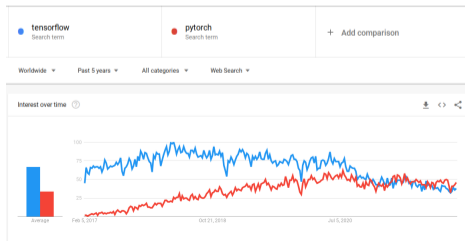


Software libraries

No lack of open-source deep learning toolkits:

- ▶ **Static graph construction** (compile then run):
 - ▶ **TensorFlow** by Google [Abadi et al., 2015] (<https://tensorflow.org>)
- ▶ **Dynamic graph construction** (run on the fly, a new graph for each sample):
 - ▶ **PyTorch** by Facebook [Paszke et al., 2017] (<https://pytorch.org/>)
 - ▶ (**TensorFlow 2** added support for dynamic graphs as well)

Every tech giant wants to introduce its own 'definitive deep learning library'. The field is very competitive.





This course

- ▶ IN5550 is currently based on **PyTorch** library...



This course

- ▶ IN5550 is currently based on **PyTorch** library...
- ▶ ...since it is more flexible and arguably easier to learn.



This course

- ▶ IN5550 is currently based on **PyTorch** library...
- ▶ ...since it is more flexible and arguably easier to learn.



- ▶ We are going to use well-established architectures...



This course

- ▶ IN5550 is currently based on **PyTorch** library...
- ▶ ...since it is more flexible and arguably easier to learn.



- ▶ We are going to use well-established architectures...
- ▶ ...so everything can be done via **PyTorch**.



This course

- ▶ IN5550 is currently based on **PyTorch** library...
- ▶ ...since it is more flexible and arguably easier to learn.



- ▶ We are going to use well-established architectures...
- ▶ ...so everything can be done via **PyTorch**.
- ▶ Makes it easy to use deep learning elements in your code.
- ▶ <https://pytorch.org/>



- 1 Computation graphs
- 2 Neural network toolkits
- 3 Next week trailer**



Dense representations of textual features

- ▶ One-hot VS dense encodings of linguistic entities.



Dense representations of textual features

- ▶ One-hot VS dense encodings of linguistic entities.
- ▶ Making the discrete world continuous.



Dense representations of textual features




- ▶ One-hot VS dense encodings of linguistic entities.
- ▶ Making the discrete world continuous.
- ▶ Preliminaries for **word embeddings**.
 - ▶ Yay! *Word2vec*!



Dense representations of textual features

- ▶ One-hot VS dense encodings of linguistic entities.
- ▶ Making the discrete world continuous.
- ▶ Preliminaries for **word embeddings**.
 - ▶ Yay! *Word2vec*!
- ▶ Stay tuned!

References I

-  Abadi, M., Agarwal, A., Barham, P., Brevdo, E., Chen, Z., Citro, C., Corrado, G. S., Davis, A., Dean, J., Devin, M., Ghemawat, S., Goodfellow, I., Harp, A., Irving, G., Isard, M., Jia, Y., Jozefowicz, R., Kaiser, L., Kudlur, M., Levenberg, J., Mané, D., Monga, R., Moore, S., Murray, D., Olah, C., Schuster, M., Shlens, J., Steiner, B., Sutskever, I., Talwar, K., Tucker, P., Vanhoucke, V., Vasudevan, V., Viégas, F., Vinyals, O., Warden, P., Wattenberg, M., Wicke, M., Yu, Y., and Zheng, X. (2015).
TensorFlow: Large-scale machine learning on heterogeneous systems.
Software available from [tensorflow.org](https://www.tensorflow.org).
-  Kumar, S. and Tsvetkov, Y. (2018).
Von mises-fisher loss for training sequence to sequence models with continuous outputs.
In International Conference on Learning Representations.
-  Paszke, A., Gross, S., Chintala, S., and Chanan, G. (2017).
Pytorch: Tensors and dynamic neural networks in Python with strong GPU acceleration.



Rumelhart, D. E., Hinton, G. E., and Williams, R. J. (1986).
Learning representations by back-propagating errors.
Nature, 323(6088):533.