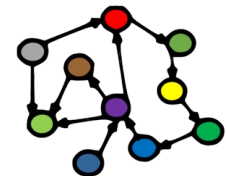


**Welcome to INFO216:
Knowledge Graphs
Spring 2022**

**Andreas L Opdahl
<Andreas.Opdahl@uib.no>**

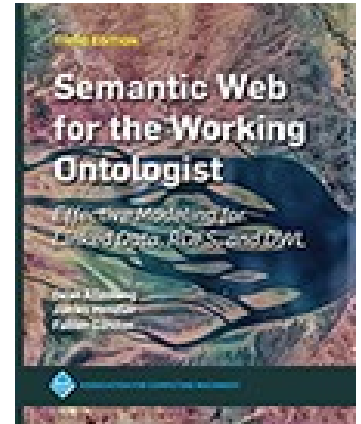
Session 12: Graph embeddings

- Themes:
 - *KGs and machine learning (ML)*
 - *what are embeddings?*
 - word embeddings
 - how to find and use them
 - other types of embeddings
 - *what are graph embeddings?*
 - how to find them...
 - ...and what to use them for



Readings

- Material at <http://wiki.uib.no/info216>:
 - Introduction to Machine Learning
 - Introduction to Word Embeddings
 - Introduction to Knowledge Graph Embeddings
- Supplementary (links in the wiki):
 - Mikolov et al's original word2vec paper
 - Bordes et al's original TransE paper
 - TorchKGE documentation (for the labs):
 - <https://torchkge.readthedocs.io/en/latest/index.html>

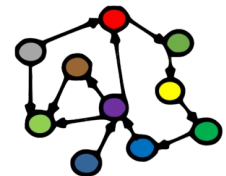


THE KNOWLEDGE GRAPH
COOKBOOK
RECIPES THAT WORK



ANDREAS BLUMAUER
AND HELMUT NAGY

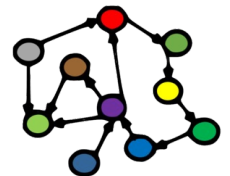
1st edition, 2020



KGs and Machine Learning (ML)

What are the connections?

- *Knowledge graphs are well matched with machine learning!*
- Preparing inputs to ML (varying origins, formats, modalities...)
 - also managing outputs from ML
- Infusing world knowledge into ML
 - common sense knowledge, world knowledge (domain and general), ...
- As a native ML technique



A micro-introduction to machine learning (ML)

- Sole purpose: to be able to understand and use KG embeddings
- *How to make computers do useful things based on examples (training data)*

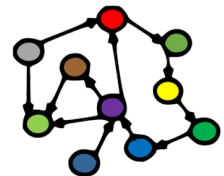
Supervised learning:

- training materials comprise input-output value pairs as examples

- *Unsupervised learning:*

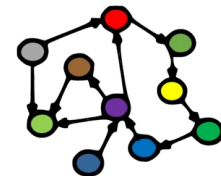
- training materials comprise only input examples

- Several other variants: *semi-supervised, reinforcement learning, ...*
- Learning KG (and other) embeddings is *unsupervised*



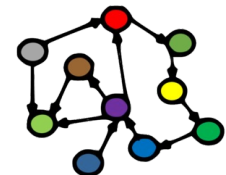
Train, evaluate, and test

- Training examples can be split in three:
 - *training data* are used to train the model
 - *validation data* are used to optimise hyper-parameters and monitor progress
 - *test data* are used only for final evaluation
 - 60%-20%-20% or 80%-10%-10% split is common
 - also minimum requirements for test examples
- *k-fold cross-validation*:
 - training and validation data are split in k folds
 - $k-1$ folds are used for training, 1 for validation
 - repeated k times for each validation fold
 - finally, the measures are averaged



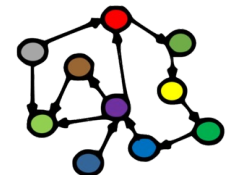
Epochs and batches

- We can go through the training data many times
 - each time is an *epoch*
- We can go through the training examples in groups
 - each group is called a *batch*
- Each example creates a *loss*
- So:
 - training consists of many epochs
 - each epoch consists of many batches
 - each batch consists of many training examples
 - each training example creates a loss
 - after each batch, steps are taken to minimise future loss



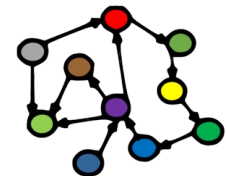
Evaluation measures

- Results without ranking:
 - *accuracy (A)*: ratio of correct results
 - there are lots of others:
 - precision (P), recall (R), $F1 = 2PR/(P+R)$, ...
- Ranked results:
 - *Hit@n*: number of correct results in the “top n”, e.g., Hit@10
 - *Mean Rank*: average rank of the correct results
 - *Mean Reciprocal Rank (MRR)*: average inverse rank of the correct results, example:
 - the correct results have rank 1, 3, 28
 - $MRR = (1/1 + 1/3 + 1/28) / 3$
- Other measures for other data types, e.g., time series data



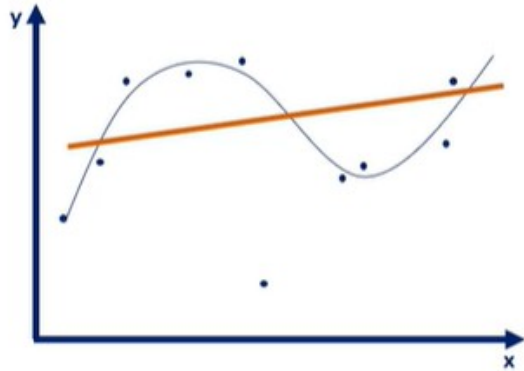
Under- and overfitting

- Underfitting:
 - we have not trained for long enough, too few epochs
 - there is more to learn from the training data
 - high loss, weak validation measures
- Overfitting:
 - we have trained for too long, too many epochs
 - the model has specialised on the training data
 - low loss, weak validation measures



Underfitting and overfitting

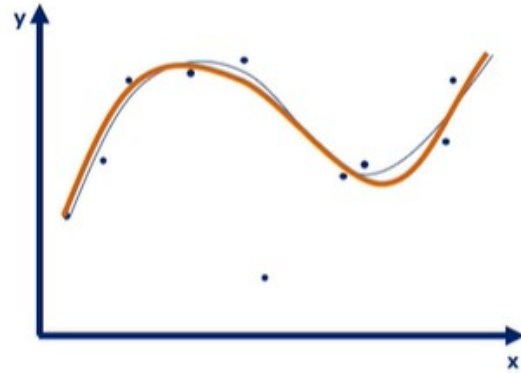
An **underfitted** model



Doesn't capture any logic

- High loss
- Low accuracy

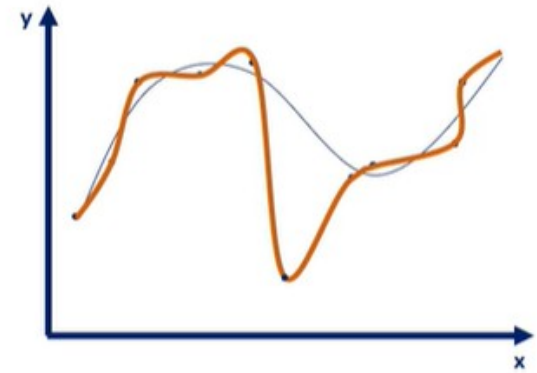
A **good** model



Captures the underlying logic of the dataset

- Low loss
- High accuracy

An **overfitted** model



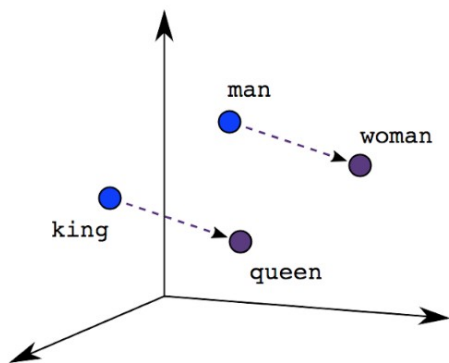
Captures all the noise, thus "missed the point"

- Low loss
- Low accuracy

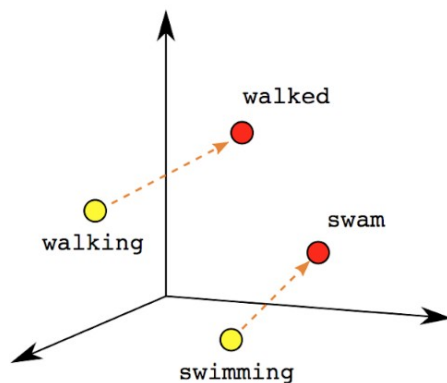
What are
embeddings?

How can we represent the meaning of words?

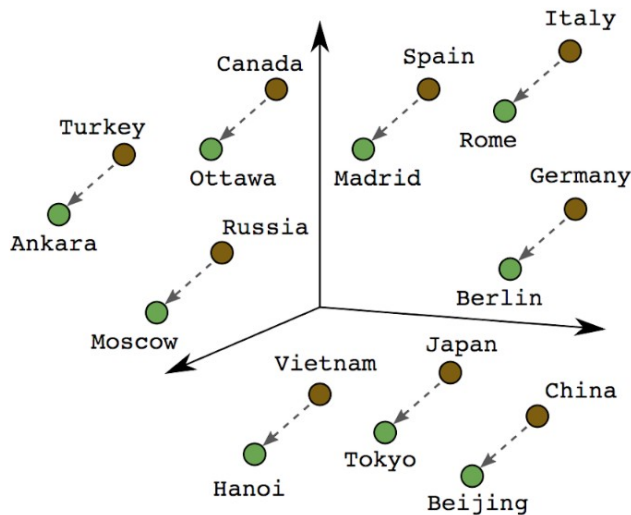
- By designation (e.g., textual descriptions in a dictionary)
- As nodes in a network (e.g., in WordNet or a knowledge graph)
- Formally (e.g., adding axioms to)
- *As vectors in a latent semantic space!*



Male-Female



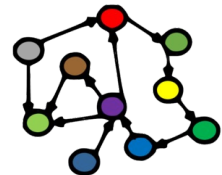
Verb Tense



Country-Capital

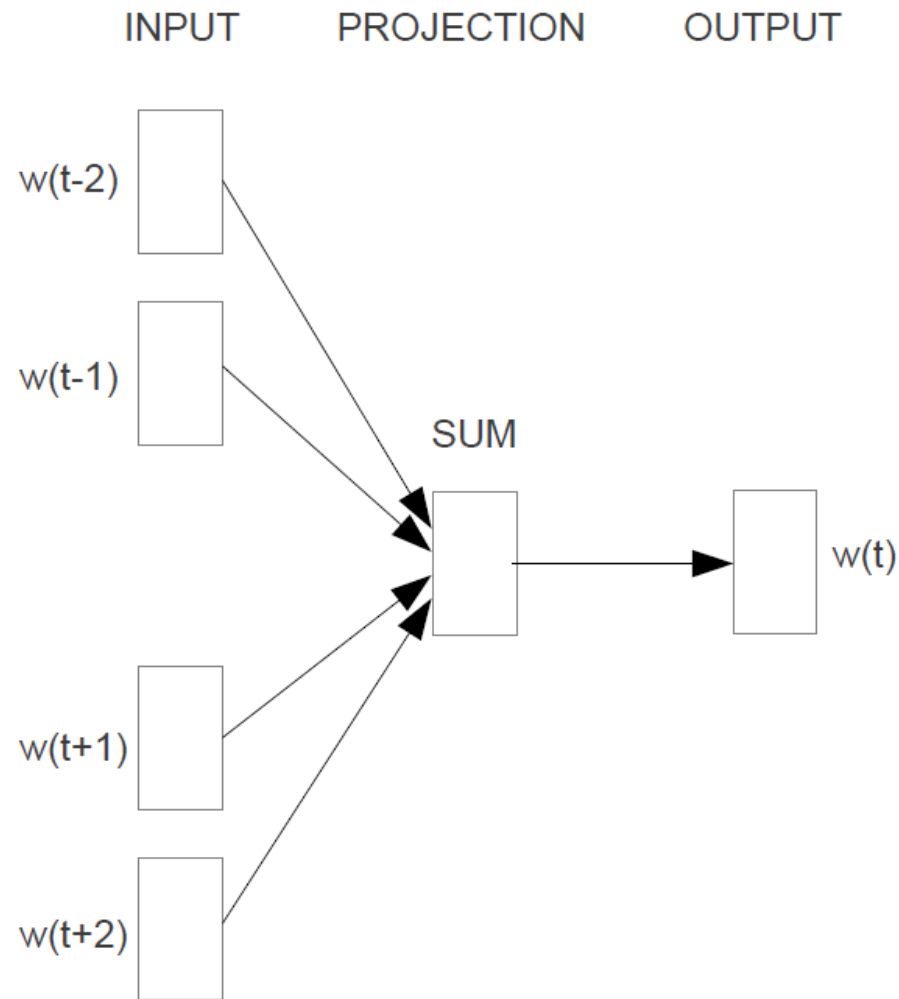
How can we represent the meaning of words?

- By designation (e.g., textual descriptions in a dictionary)
- As nodes in a network (e.g., in WordNet or a knowledge graph)
- Formally (e.g., adding axioms to)
- *As vectors in a latent semantic space!*
 - [0.01 0.62 0.03 ... 0.41]
 - similar words are close to one another
 - relative positions between words can be systematic
 - [Paris] – [France] + [Italy] \approx [Rome]
 - distances between words can represent relations
 - [J. K. Rowling] + [influenced by] \approx [J. R. R. Tolkien]
- Important use: as inputs to deep neural networks that process NL text



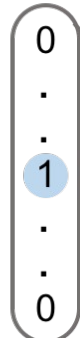
How to learn the vectors?

- CBOW (Continuous Bag of Words):
 - part of *word2vec*
 - neural network with one hidden layer
 - trained on large corpus of NL text (1.6 billion words)
 - input examples: sentences with one word missing
 - expected output: the missing word
 - the weights in the neural network are used as word vectors
- Also: Skip-gram, GloVe, FastText, ...
- Ubiquitous as *inputs to deep neural networks that process NL text*

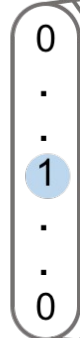


CBOW

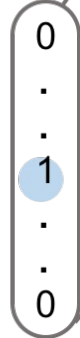
Input



X



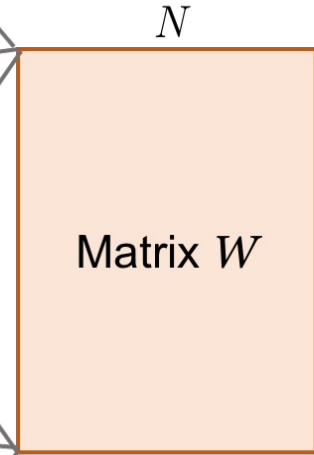
X



X

Example dimensions:

- $V = 10000$
- $N = 300$



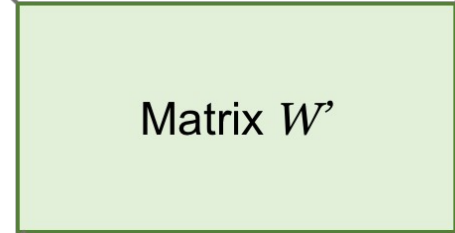
Matrix W

$V =$
avg

Hidden

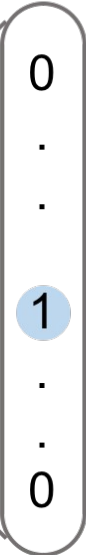


X



Matrix W'

Output softmax

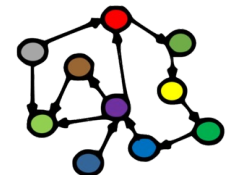


N-dimension vector
(**Average** of vectors of
all input words)

After training, W'
consists of V word
vectors with dimension N

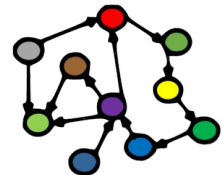
Word similarity

- Extremely powerful and much used, *but be careful*
- The distributional hypothesis:
 - “words that occur in the same contexts tend to have similar meanings” (Harris 1954)
 - hence, word similarity can be measured in terms of vector similarity
 - *this is not true*
 - synonyms will often appear close to the same words
 - but so will many antonyms (“love”, “hate”)
 - syntagmatic similarity:
the words are able to combine in sentences with the same other words
 - paradigmatic similarity:
the words can be substituted with one another



Other types of embeddings

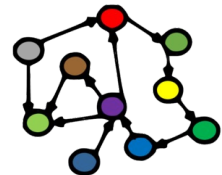
- The idea has caught on:
 - phrase embeddings (“baseball bat”, “linear algebra”, ...)
 - word piece embeddings ([lin-] + [-ear], [al-] + [-ge-]+ [-bra])
- Contextual embeddings (ELMo):
 - how to deal with words that are
 - homonymous (different words that look/sound the same)
 - polysemous (same word form has several meanings)
 - words have different embeddings in different neighbourhoods
- Sentence and paragraph embeddings:
 - transformer models with attention
 - BERT and descendants, e.g., S-BERT
- *Graph embeddings!*



What are
graph embeddings?

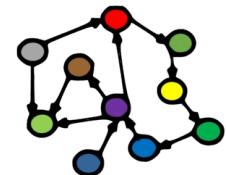
How can we represent the meaning of graphs?

- By designation (e.g., textual descriptions) of nodes and edges
- By URIs defined in open KGs and standard vocabularies
- Formally (e.g., using description logic)
- *As vectors in a latent semantic space!*
 - node vectors
 - edge vectors
 - graph vectors



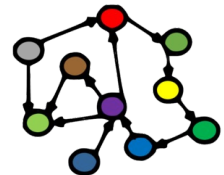
What can we do with graph embeddings?

- Graph completion and validation:
 - node classification: given a node which type should it have?
 - link prediction: between two given nodes, should there be an edge?
 - relation prediction: given two nodes, which edge type should link them?
 - triple classification: given two nodes and an edge, is the triple correct?
- Graph (or sub-graph) classification:
 - what type of entity/situation/event does the graph represent?
 - which class does the graph represent?
- Input to deep networks:
 - perhaps in combination with text, images, ...
 - “early or late fusion”
 - deep “multi-path” networks ((example))



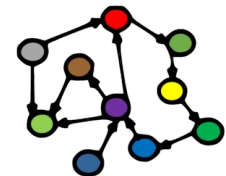
How to learn the vectors?

- Early and simple example:
 - *Deepwalk* (2014)
- Algorithm:
 - 1) drop a marker randomly onto a graph node
 - 2) let the marker traverse the graph randomly along edges for n steps
 - additional parameters can guide traversal
 - 3) treat each resulting walk of n nodes as a sentence of n words
 - 4) feed a corpus of n -node walks into CBOW or similar
- *Instead of a vector for each word, this produces a vector for each node*
- Limitations:
 - all relations are equal
 - sampling may not fully exploit graph structure



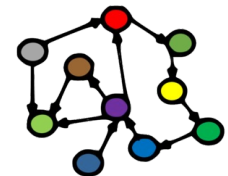
Translational embeddings (TransE)

- The *translational property*:
 - if $(h, r, t) \in KG$, then $[h] + [r] \approx [t]$
- Approach:
 - start out with random vectors for nodes and edges
 - repeat:
 - for each $(h, r, t) \in KG$, generate corrupted (h', r, t') *not* in KG (either h' or t' is changed)
 - adjust vectors to
 - minimise $dist([h] + [r], [t])$
 - maximise $dist([h'] + [r], [t'])$
 - loss is $L = \gamma + dist([h] + [r], [t]) - dist([h'] + [r], [t'])$

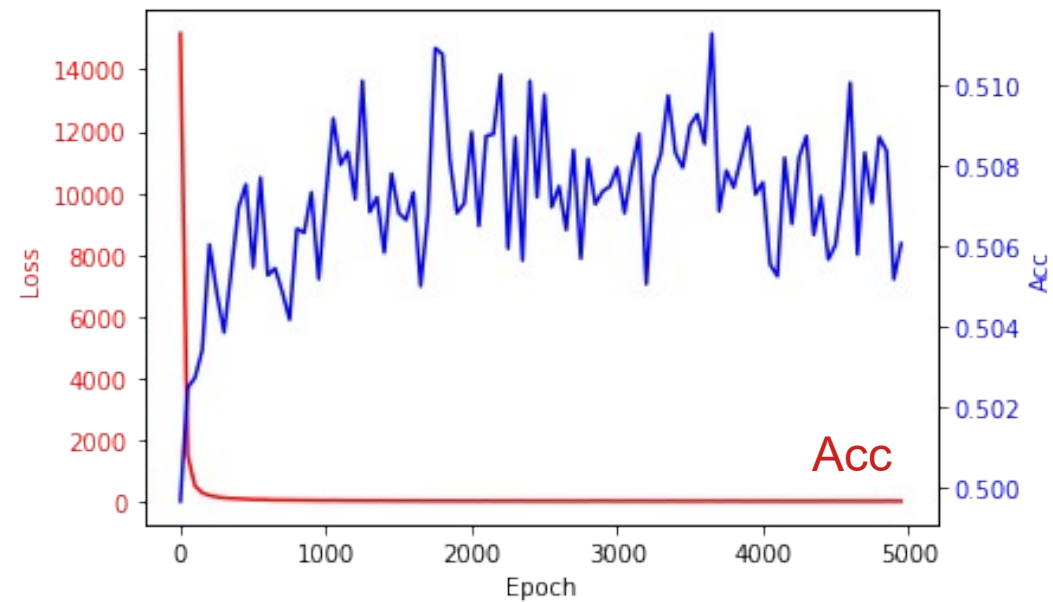
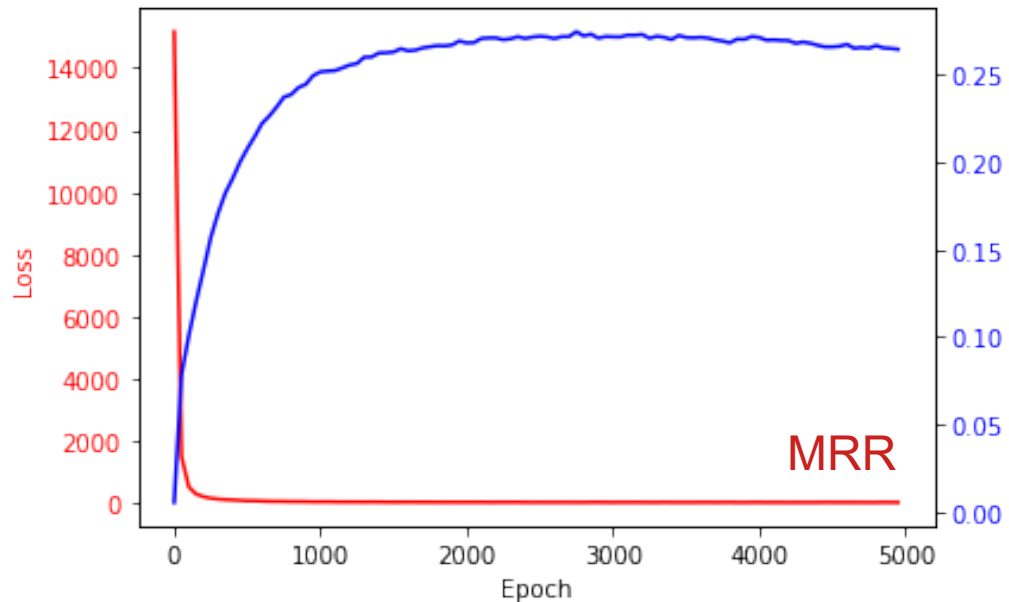


Evaluation

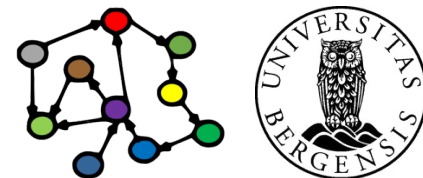
- *Link prediction:*
 - $h + r \approx$ which t ?
 - MRR (not reciprocal), Mean Rank, Hit@n (@10).
 - filtered and raw variants
- *Relation prediction:*
 - $h - t \approx$ which r ?
 - MRR (not reciprocal), Mean Rank, Hit@n (@10).
 - filtered and raw variants
- *Relation classification:*
 - are (h, t, r) and (h', t, r') in KG?
 - accuracy (A)



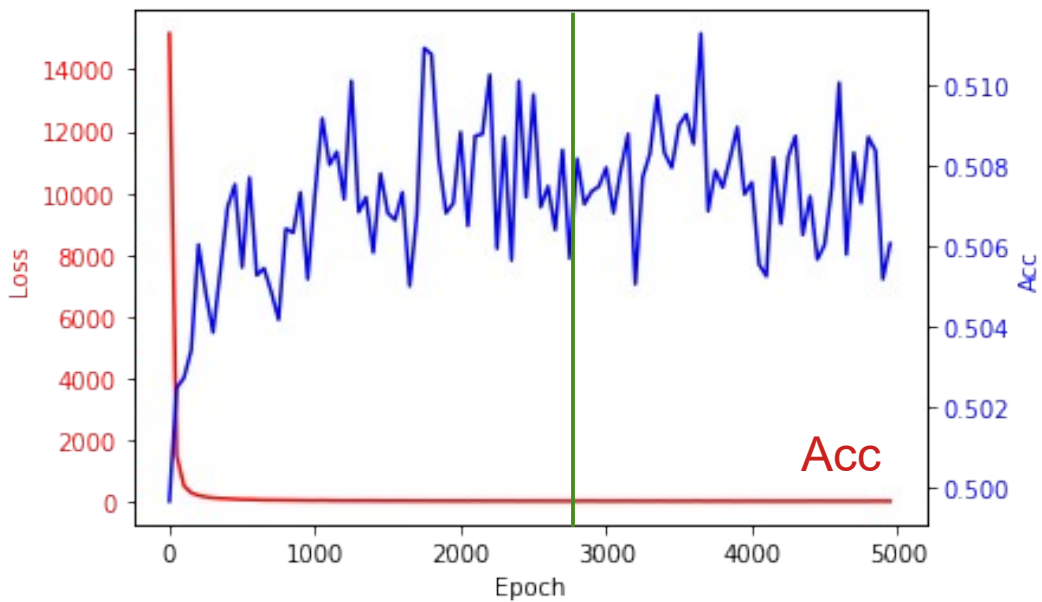
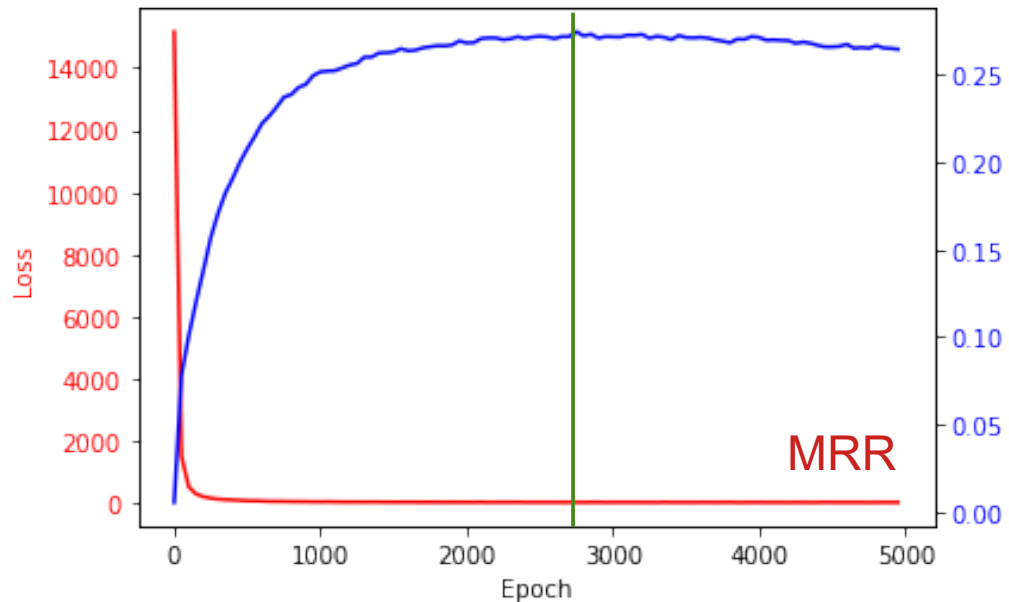
Learning curves



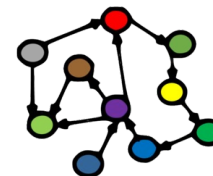
TransE on FB15k237 with 5000 epochs



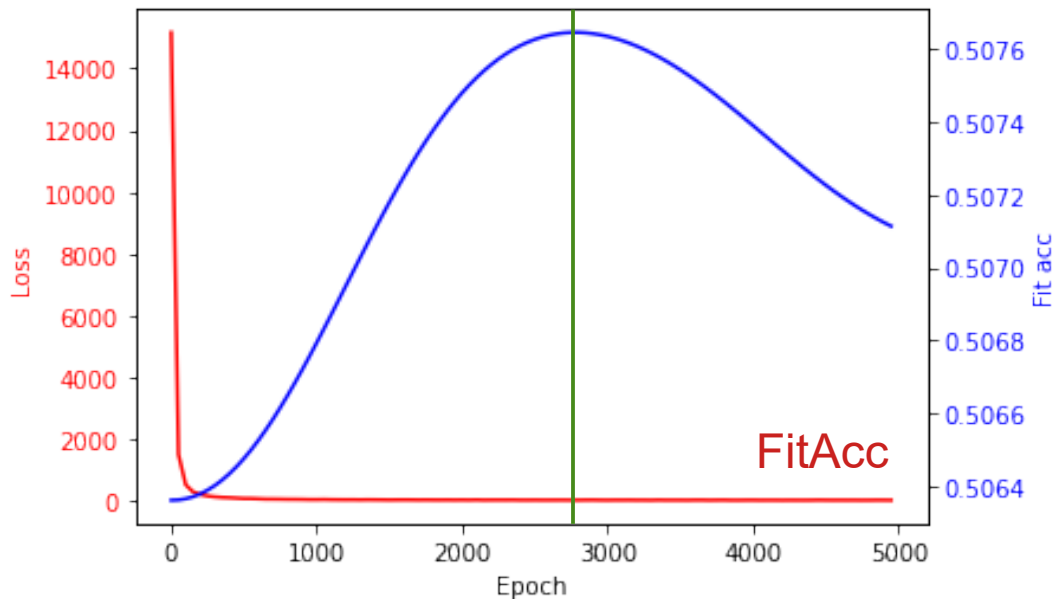
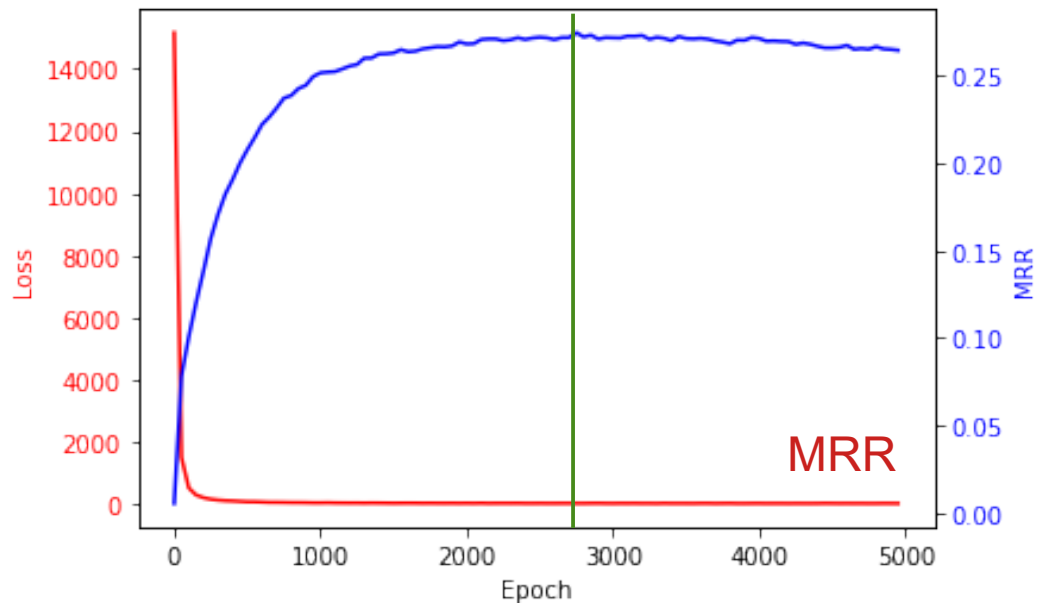
Learning curves



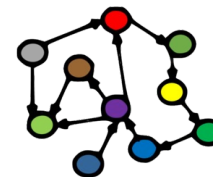
TransE on FB15k237 with 5000 epochs



Learning curves



TransE on FB15k237 with 5000 epochs

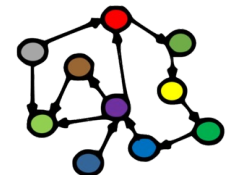


Datasets and pre-trained models

- Datasets:
 - Freebase extract (FB15k)
 - WordNet synsets (WN)
 - both have problems with training/validation/test overlap:

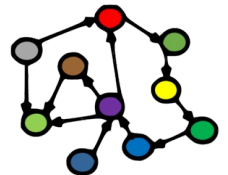
DATA SET	WN	FB15K	FB1M
ENTITIES	40,943	14,951	1×10^6
RELATIONSHIPS	18	1,345	23,382
TRAIN. EX.	141,442	483,142	17.5×10^6
VALID EX.	5,000	50,000	50,000
TEST EX.	5,000	59,071	177,404

- use FB15k237 and WN18RR instead
- Pre-trained models:
 - for example TransE already trained on FB15k237



Limitations

- *TransE* is powerful and simple, but has limitations:
 - works best for 1-1 relations
 - trained on corrupted (h', r, t) and (h, r, t') variants, but never (h, r', t)
 - therefore bad on relation prediction
 - several derivations:
 - *TransH, TransR, TransD, TorusE, ...*
 - more recent developments:
 - *Graph Neural Networks (GNNs)*
 - e.g., *Graph Convolutional Networks (CGNs)*
 - combines ideas from:
 - Convolutional Neural Networks (CNNs)
 - big graph databases



Next week:
Knowledge Engineering /
Wrapping Up