# Less supervision?

L101: Machine Learning for Language Processing Andreas Vlachos



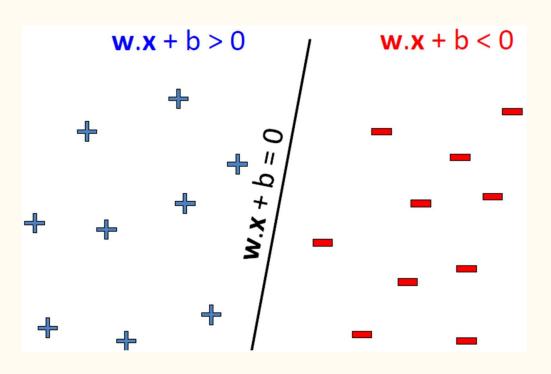
## So far: supervised learning

#### Why?

- Training models using real outputs for the task you want to solve has better chances of success
- Even if we didn't need labeled data for training, we still want to know how good our model is before deploying it (and not tune hyperparams on the test data!)

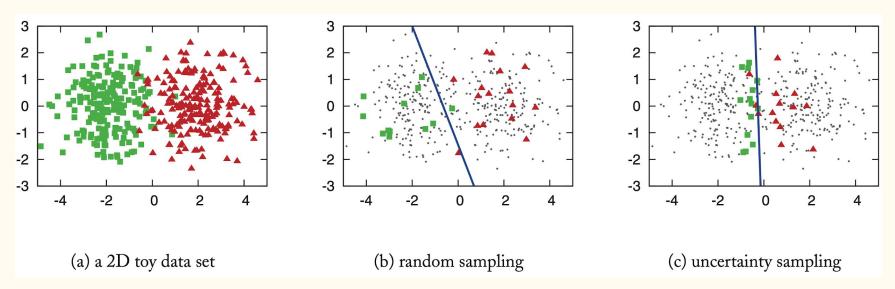
But it needs labeled data for training, how can we help ourselves given that we can only find unlabeled data for free (usually!)

#### Do we need all that labeled data for training?



- Are all instances equally useful in learning a model?
- Let the algorithm decide!
  - Typically select a few instances to label, update model, repeat
- Active learning in education refers to the student asking questions

#### Active learning works



- Savings against randomly selected training data can be impressive
  - Even if you are <u>fine-tuning BERT</u>
- In fact <u>less can be better</u> (sometimes)
- But this is not always the case (more later)

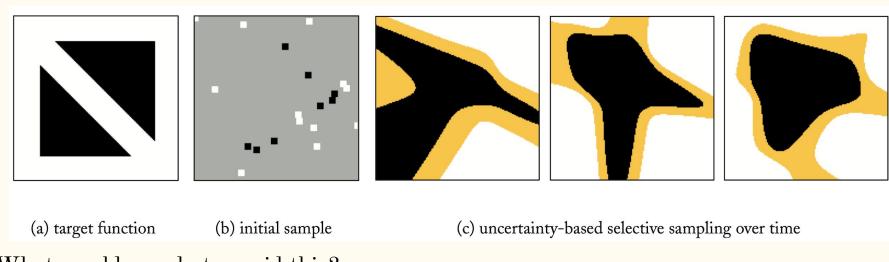
#### Active learning setups

- Pool-based
  - we have a fixed pool of unlabeled data and iterate through it
  - o most common: we get some unlabelled instances and need to build a model
- Streaming/online:
  - decide whether to ask for a label now or never
  - o ask a customer if the transaction is fraudulent, if the e-mail is spam, etc.
- Constructing instances:
  - generate instances for labeling
  - very rare in NLP; need very good generative models
  - o now seems within reach
- Feature-based active learning
  - o <u>label features instead of instances</u>

#### How to choose informative instances?

- Uncertainty based sampling
  - Least confident: pick the instances with the lowest score by the model for any label
  - Entropy in the label distribution (for probabilistic models)
- Query by committee: train a few models and select the instances where they disagree the most
- Meta-learning: learn a model to select the most useful instances (<u>pool</u>, <u>stream</u>)
- <u>Discriminate</u> in favour of instances that don't look like the labeled data
- <u>Bayesian Active Learning by Disagreement</u>: combine uncertainty with informativity on model parameters
  - Selecting instances in batches is a common practical consideration

#### Things can go wrong



What would you do to avoid this?

- Make sure your initial sample is representative?
- Select instances at random too?

#### When not to active learn

- If we need data for evaluation
  - o it should be obvious that AL biases the data, <u>not guaranteed to be</u> <u>representative of the task</u>
  - need to approach the data selection differently (<u>active testing</u>)
- If we don't think the model(s) we have can give good estimates of uncertainty
- If we are planning to change models later
  - Data selected by one model <u>can be worse than random</u> for another

### Bandit learning

Some times obtaining a complete label for an instance is impossible: e.g. for a search query we need to check all webpages

#### Repeat:

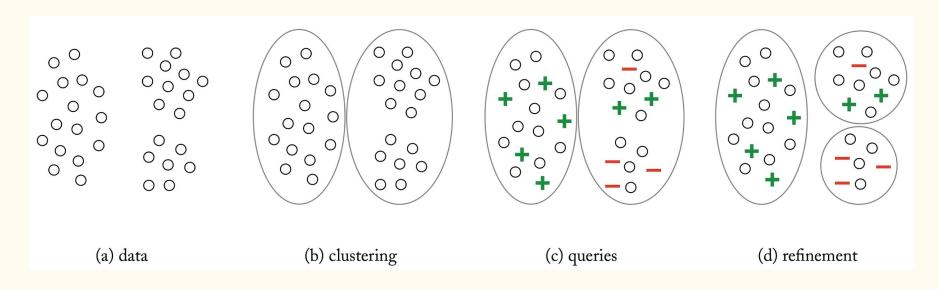
- pick the most promising handle
  (= label)
- get a reward for that handle
- Update using the reward



The key here is how to define promising: a balance of exploration/exploitation Can be adaptive: first explore, then exploit

See <u>Kreutzer et al. (2017)</u> for an application to NMT

### Can we use the unlabeled data to help?



Yes, with clustering!

Hierarchical sampling: prefer instances from clusters that are impure

#### Clustering

Many good algorithms:

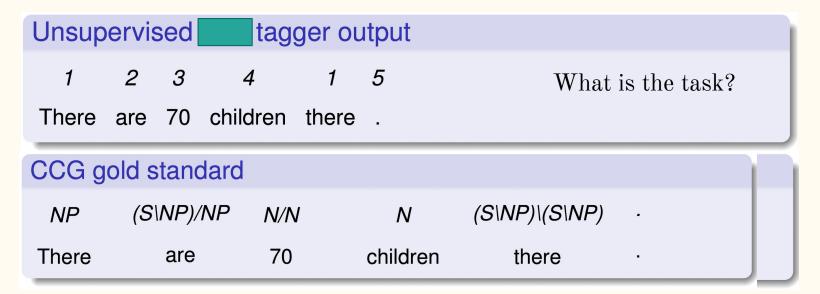
- K-means (++)
- Gaussian mixture models
- Spectral clustering
- Topic models can be thought of as soft clustering

They are great to explore unlabeled data and learn about theirs properties, but

Hard to evaluate...

<u>Clustering: Science or Art</u>: "the major obstacle is the difficulty in evaluating a clustering algorithm without taking into account the context"

### Unsupervised NLP evaluation

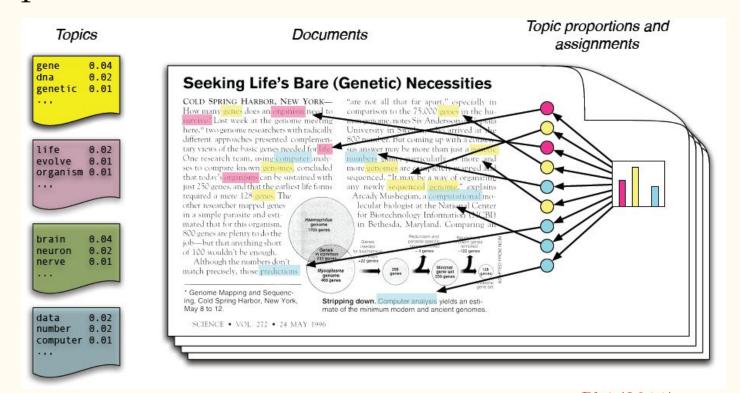


Why PoS tagging and not CCG super-tagging?

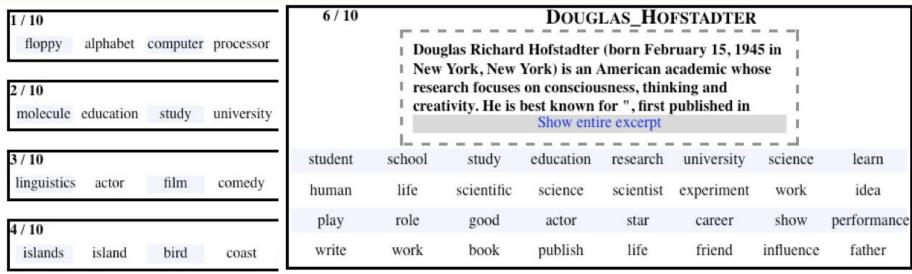
So no point in doing unsupervised learning?

No, but you need to put it in context to evaluate it

#### Topic models



### Topic modeling evaluation (<u>reading tea leaves</u>)



- Word intrusion: Pick a topic, take top words, throw in a low-prob word for the topic, ask humans to spot the intruder
- Topic intrusion: Pick a document, take its top topics, throw in a low-prob topic for the document, ask humans to spot the intruder
- Avoid automatic evaluation if you can

#### But isn't unsupervised learning common?

Yes! But often this is not what unsupervised in the ML sense (no labeled data, e.g. in clustering no cluster info, in topic models no topic info):

- Some supervision gets in through dictionaries, mapping to labels, etc.
- Or (distant) supervision was readily available on the web
- Some dev set was most likely used (hopefully not the test set!)
- Nothing wrong with this; in fact it is a great way to solve tasks
- But if we want labels as output we need to provide them to the model

### What about language modeling?

- Supervised or unsupervised learning?
- For me it is supervised, but we can harvest data for it at will
- More data beats better model

- The main application of LMs for a while was to score outputs from MT, ASR, etc.
- Nowadays?

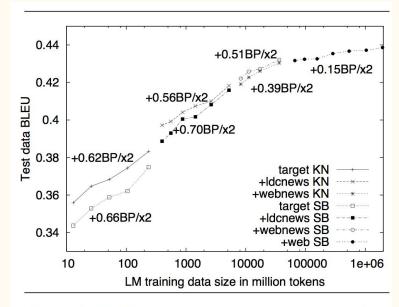


Figure 5: BLEU scores for varying amounts of data using Kneser-Ney (KN) and Stupid Backoff (SB).

Brants et al. (2007)

#### Word Embeddings

- Once a by-product of learning RNN-based language models, now a goal in itself
- That's the key insight of the <u>word2vec</u> paper: stop worrying about trying to build a language model, focus on the embeddings
- Supervised or unsupervised?
- Unsupervised: LM has a supervised training objective, but we don't have gold standard embeddings
  - This is also why their evaluation is difficult
  - Often done in context: input for supervised models (e.g. <u>BERT</u>)

#### Language models as models for task X?

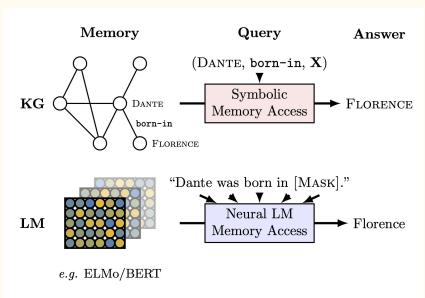


Figure 1: Querying knowledge bases (KB) and language models (LM) for factual knowledge.

- Yes! But usually we need to add a bit task-specific supervision (fine-tuning)
- If our task can be modelled as an LM, we can take advantage of a lot of data and pre-trained models
- If we are building task-specific models, we'd better improve on it!
- If we developing a task/dataset, make sure a LM can't solve it (easily), e.g. don't use an LM (alone) to construct it

#### Prompting large language models

#### Few-shot

In addition to the task description, the model sees a few examples of the task. No gradient updates are performed.

```
Translate English to French: 

sea otter => loutre de mer 

peppermint => menthe poivrée

plush girafe => girafe peluche

cheese => 

prompt
```

Brown et al. (2020)

- Look ma, no training, just a few examples!
- From feature engineering, to architecture, and now to prompt engineering?
- Not so fast:
  - Needs a lot of dev data for model/ prompt selection (<u>Perez et al., 2021</u>)
  - Fine-tuning is more reliable and not necessarily more expensive (<u>Logan IV</u> et al., 2021)
- Descriptions/instructions in zero-shot approaches can also be seen as prompts (Aly et al. 2021, Wei et al. 2022)

### Bibliography

Active learning book (Burr Settles, where a lot of images where taken from)

Bandit learning book (Slivkins)

Foundations of LLMs (Xiao and Zhu, 2025)