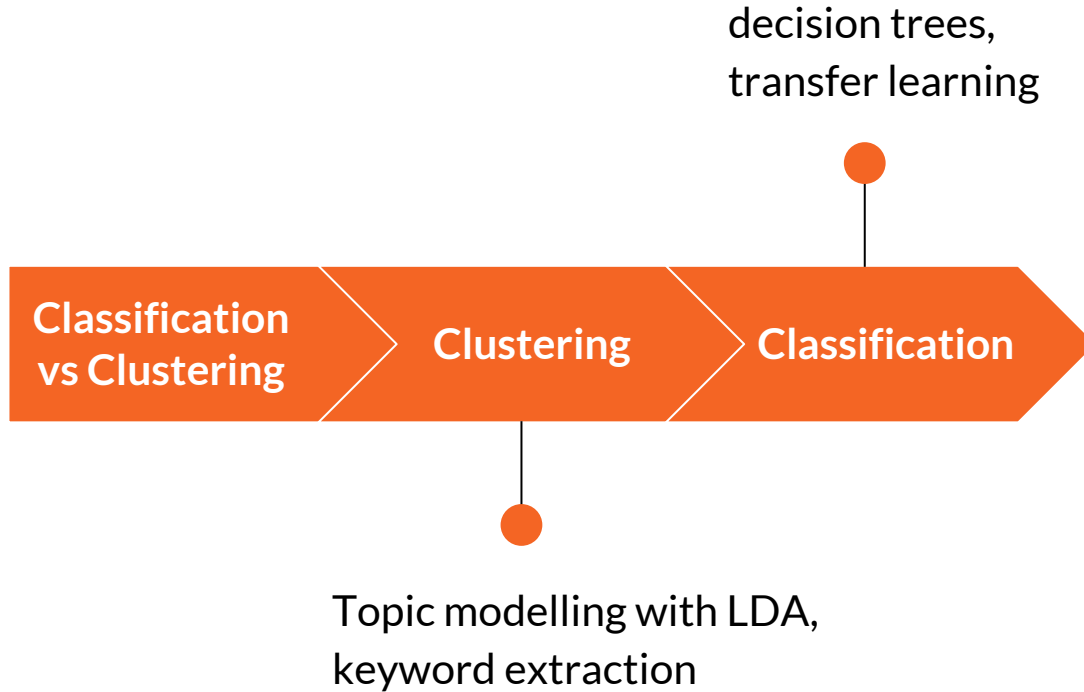

Clustering & Classification

Chris Swart Data Scientist @ Resolver

15/10/2018 - mcubed.london

Overview



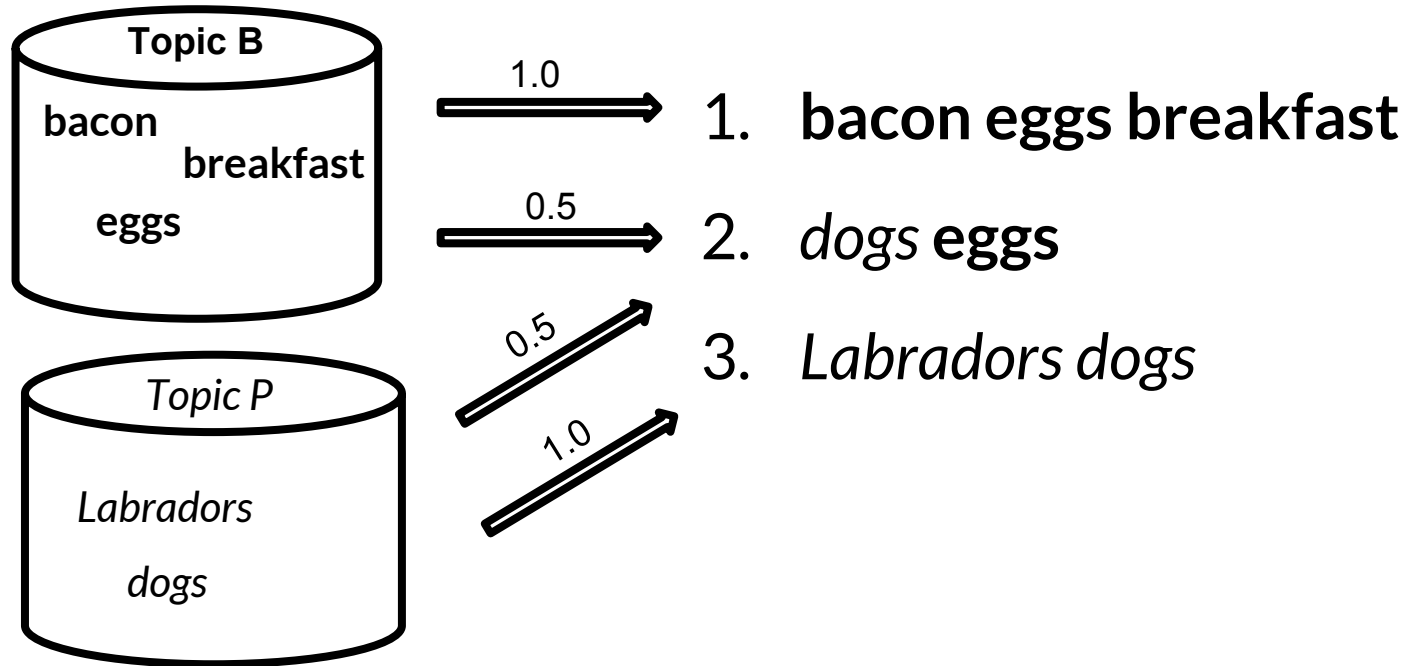
Clustering vs Classification

- Discover patterns based on **unlabelled** data
 - **Easier to gather data** as we don't need labels
 - **Less clear how to interpret**
 - More flexible structure
 - Utilise patterns based on a **labels** you have
 - **More expensive to gather data**
 - **Interpretable**
-

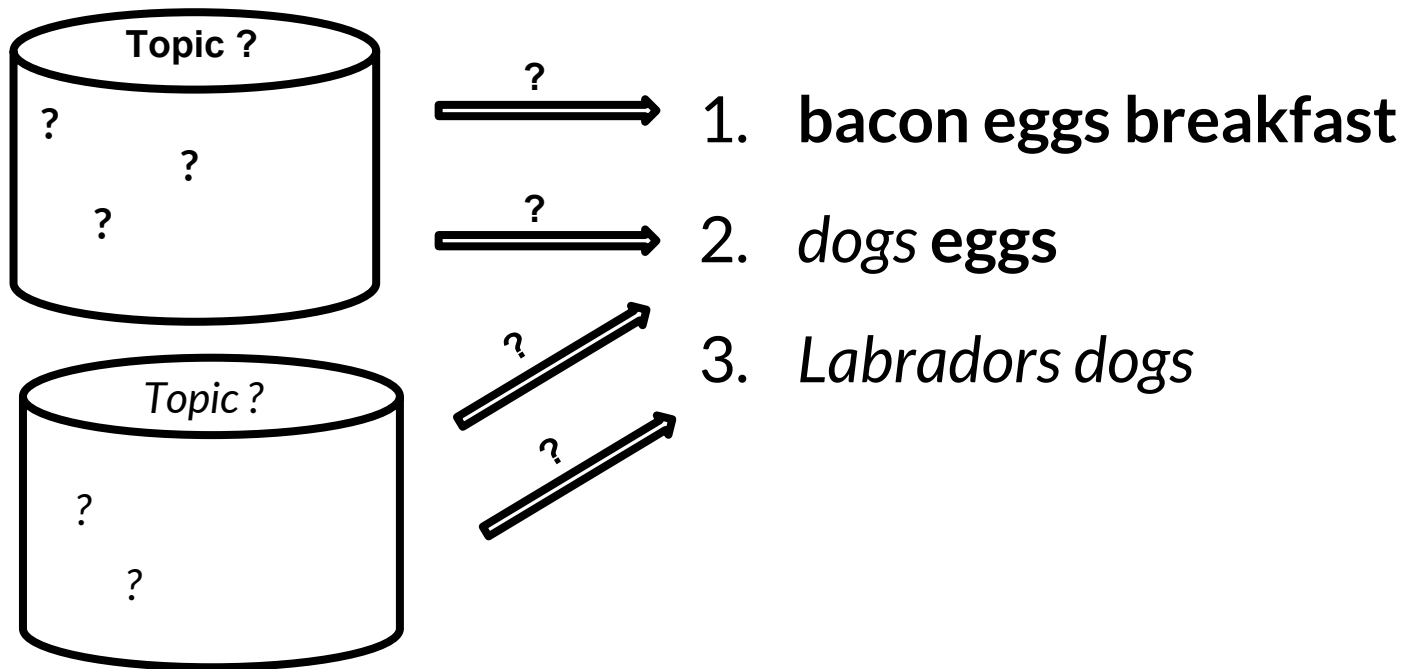
Topic modelling with LDA

1. I had **bacon** and **eggs** for **breakfast**.
 2. *Labradors* are cute *dogs*.
 3. My *dogs* do not like **eggs**.
-

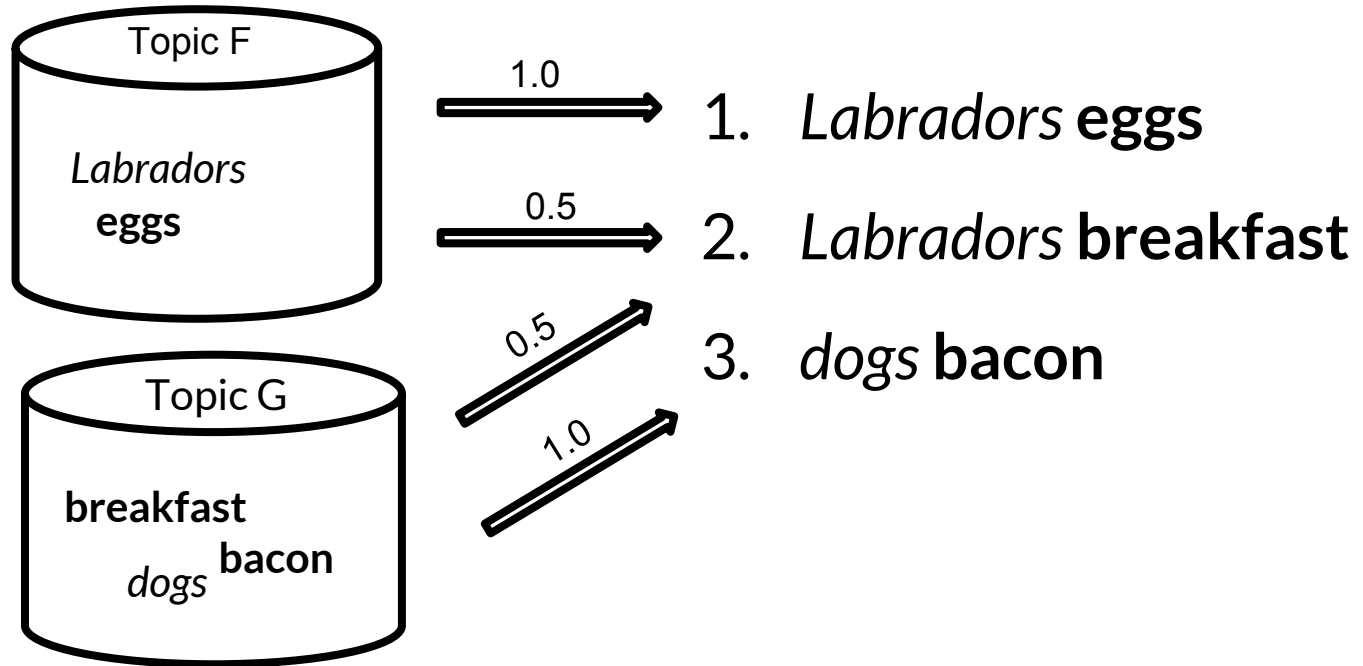
Generative process with LDA



Inference with LDA



Generative process with LDA



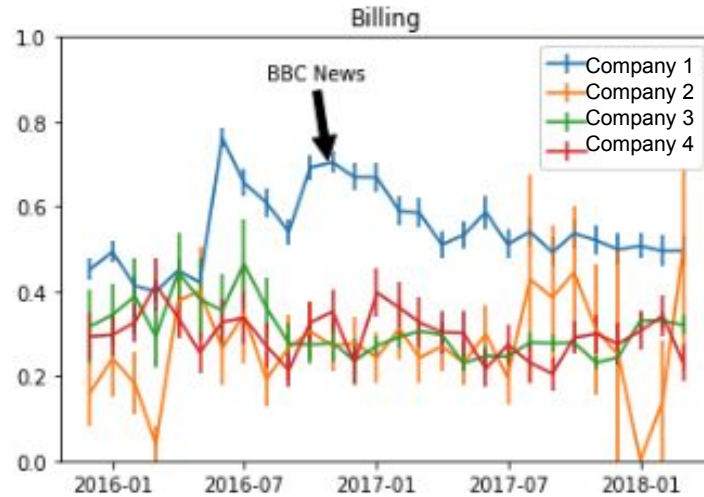
Evaluating LDA

- Word intrusion [eggs, breakfast, Labradors, bacon]
 - Consolidate topics and annotate documents
-

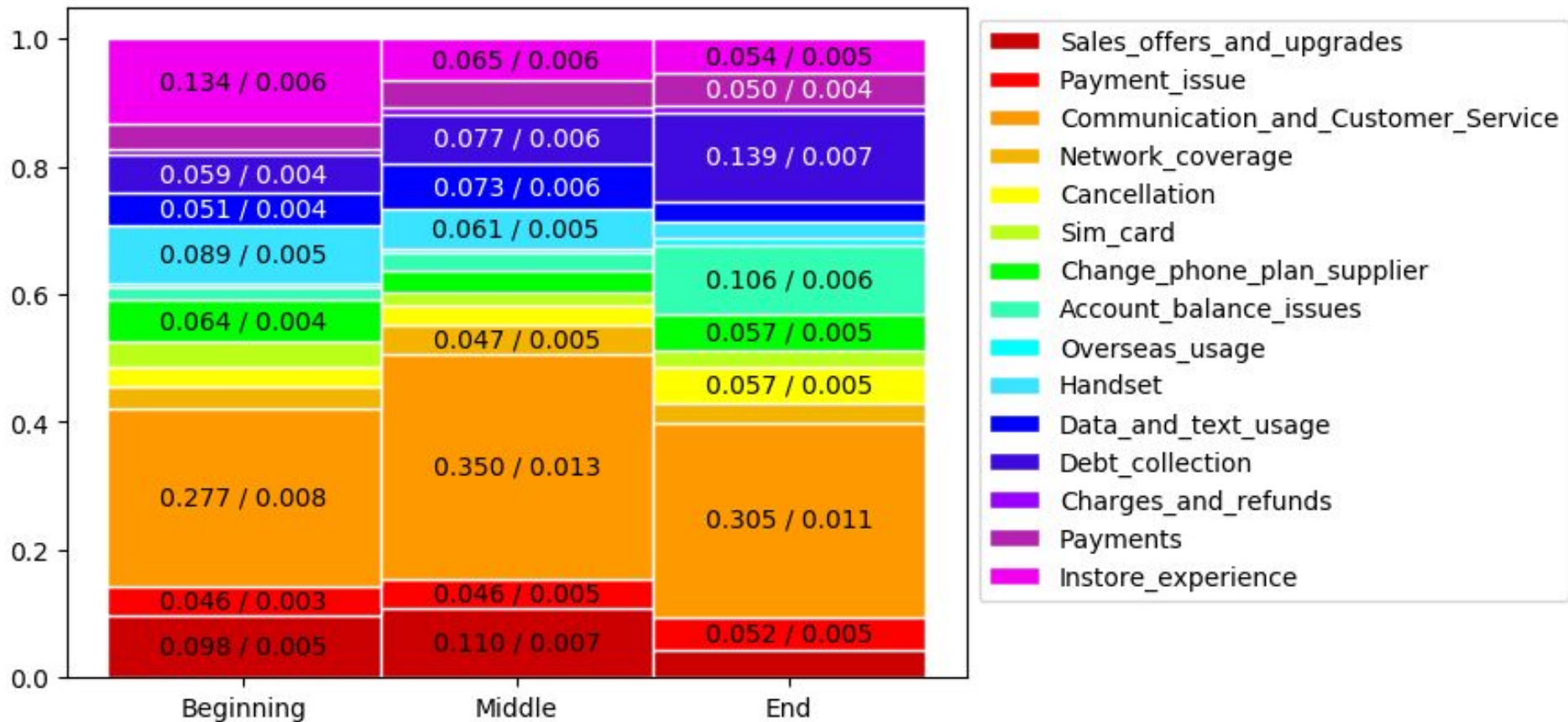
Bank1 performance compared to industry



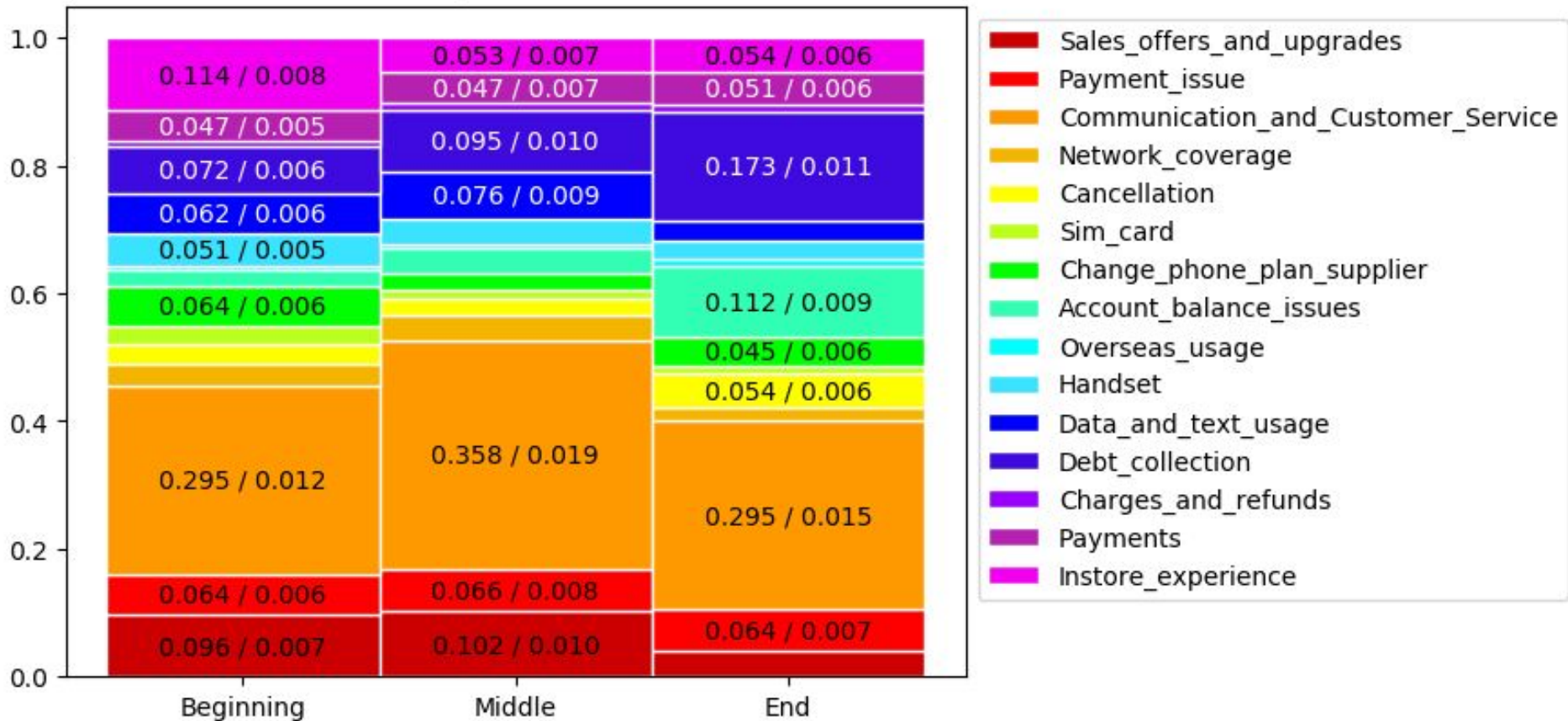
Case study - Telco1



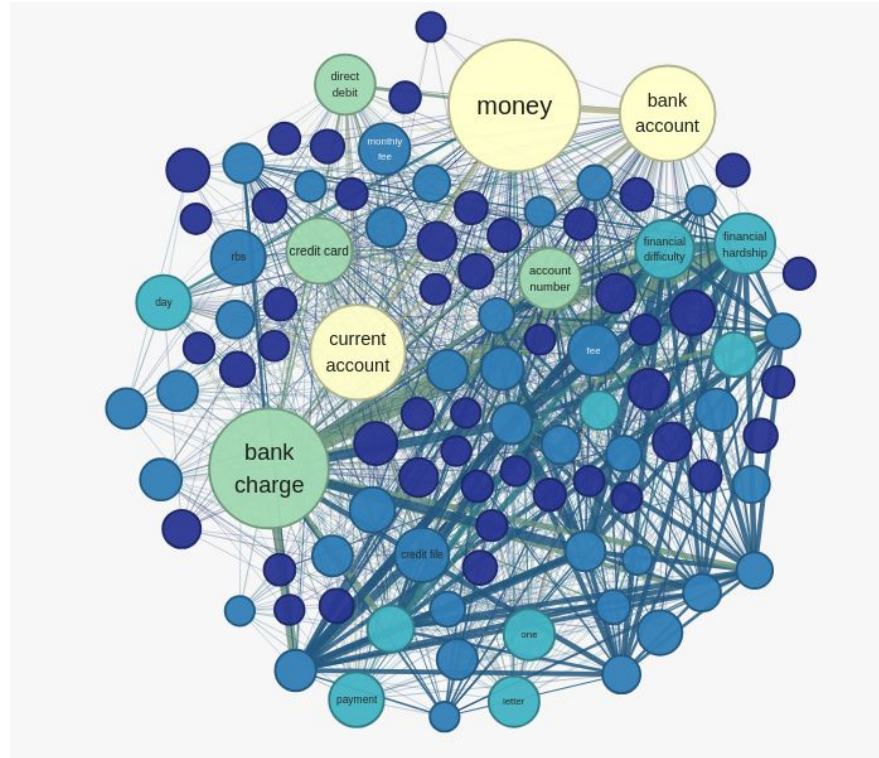
Industry user journey



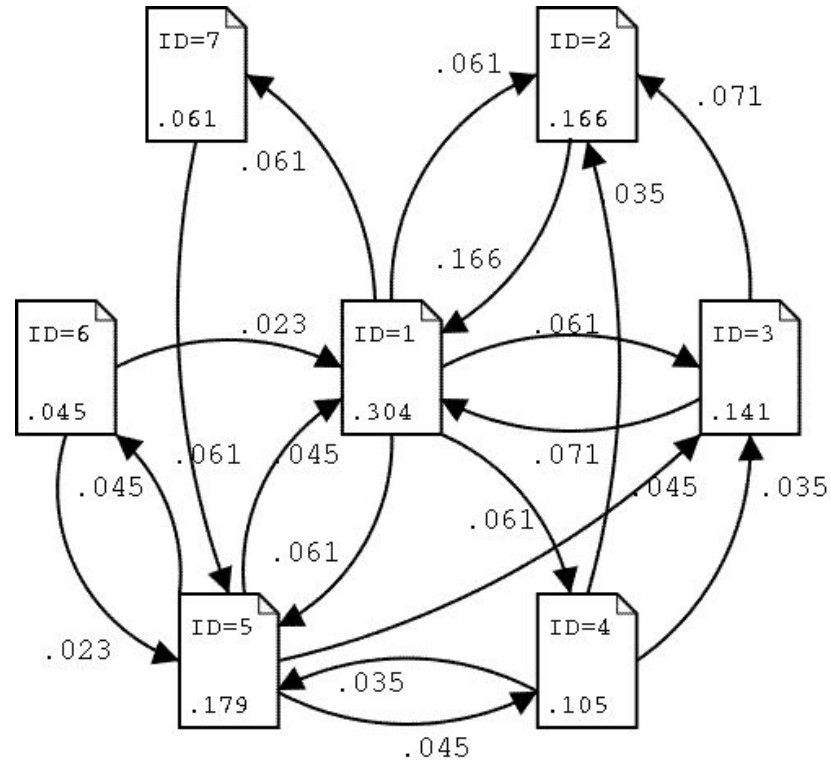
Telco1 user journey



Keyword extraction with TextRank



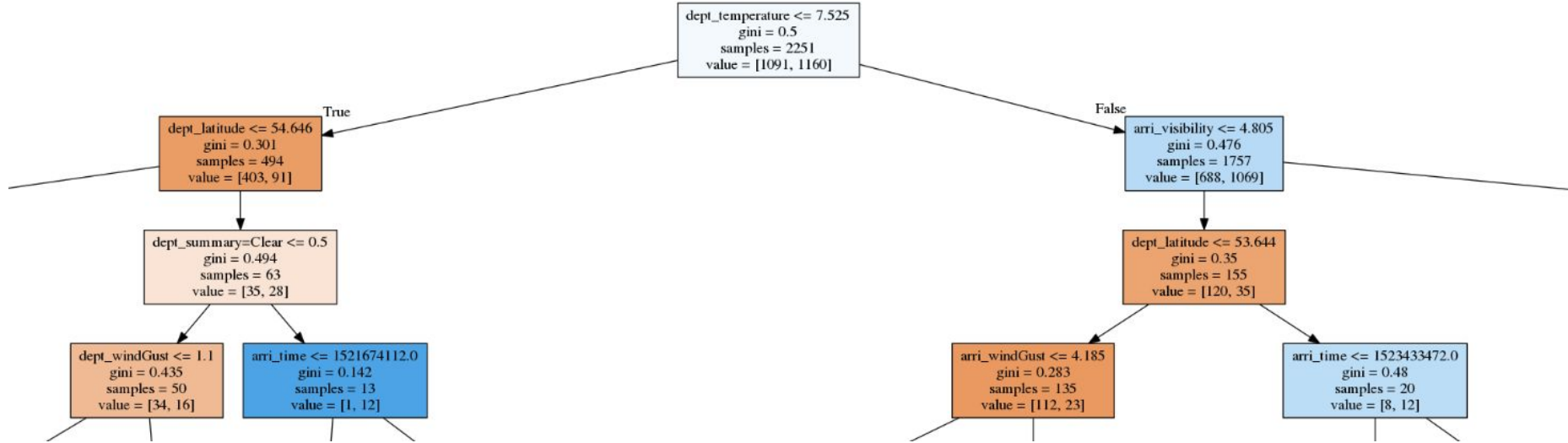
Pagerank + Text data = TextRank

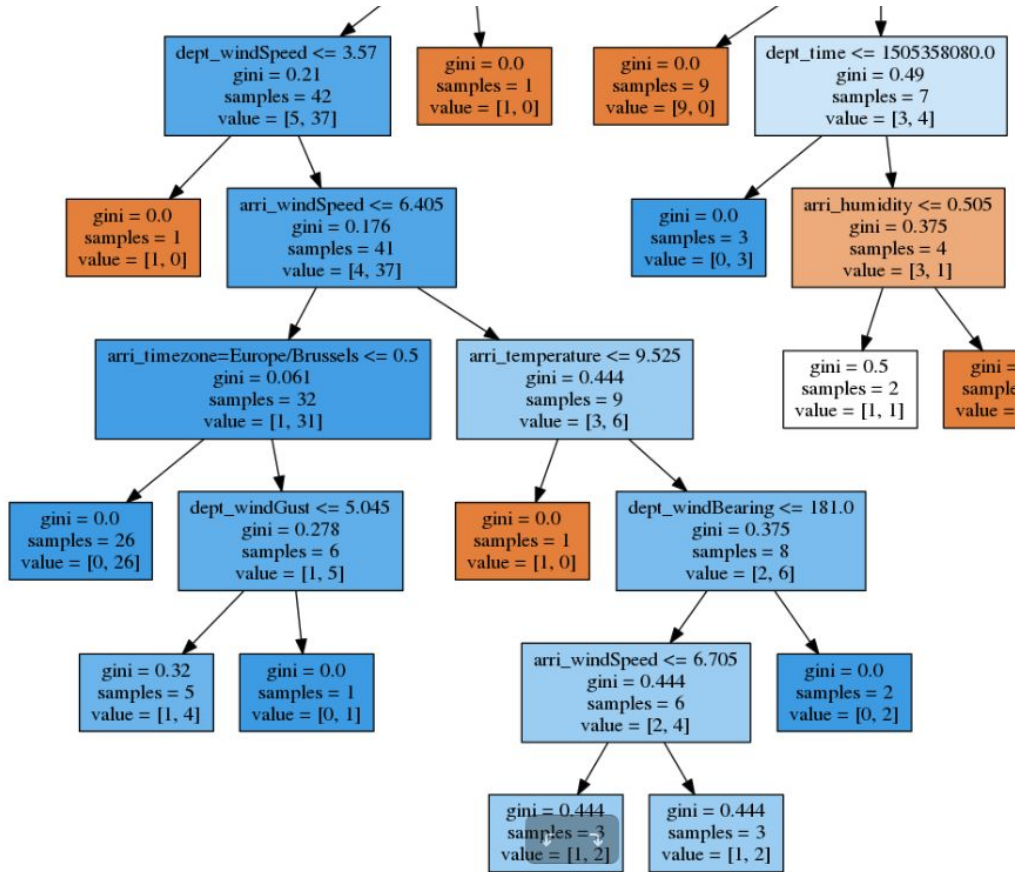


Evaluating keyword extraction

- Supervised ranking problem
 - User feedback
-

Decision trees





Decision forests

- Split the data
- Build a decision tree on each split
- Average their outputs



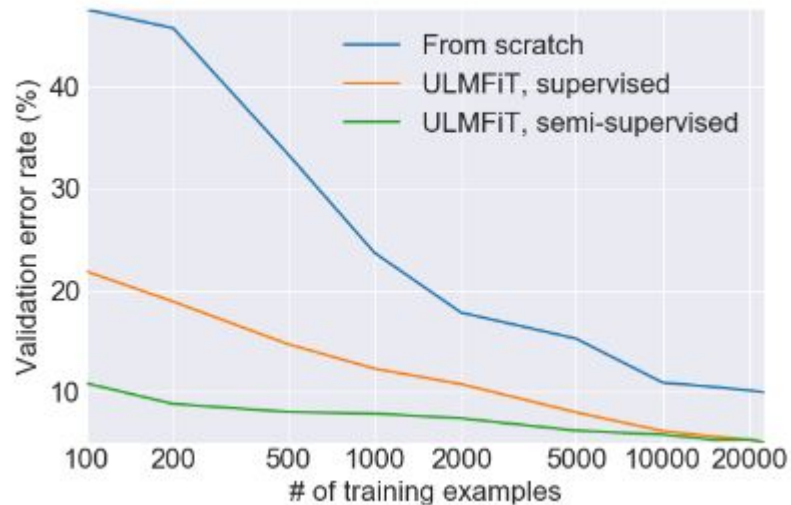
Evaluating decision forests

- Out of Bag scores
- Train / test / validation split

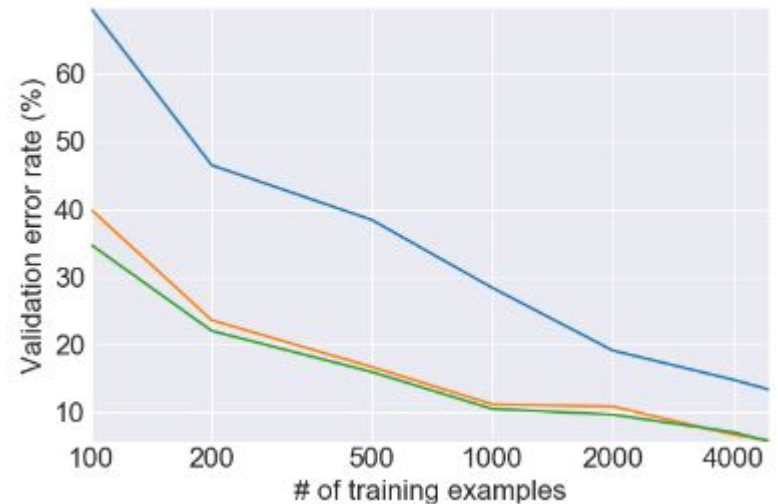


Transfer learning

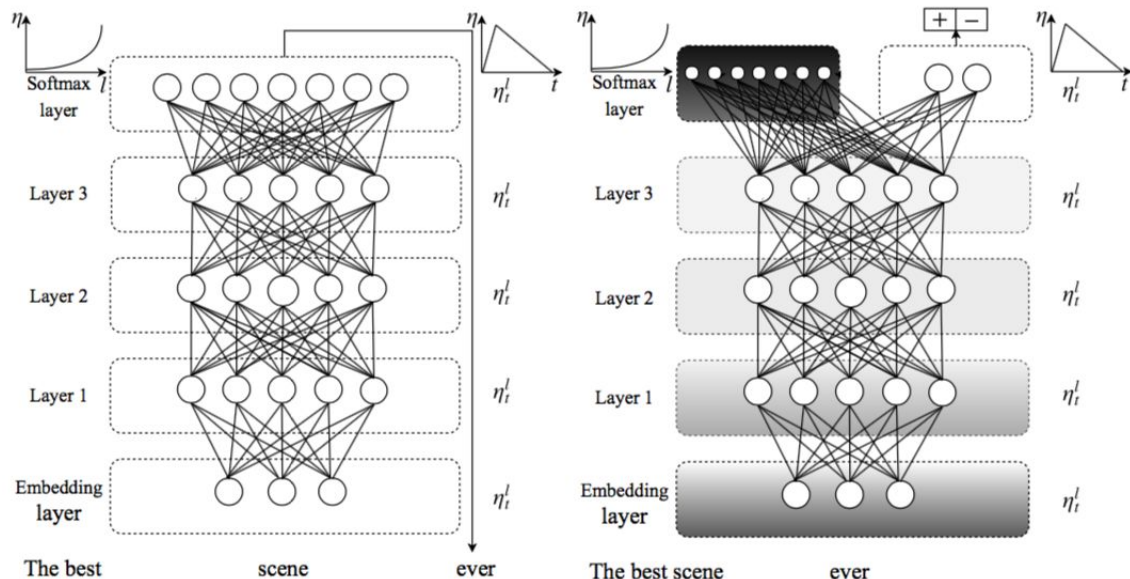
Movie reviews



Topic extraction of non-complaints

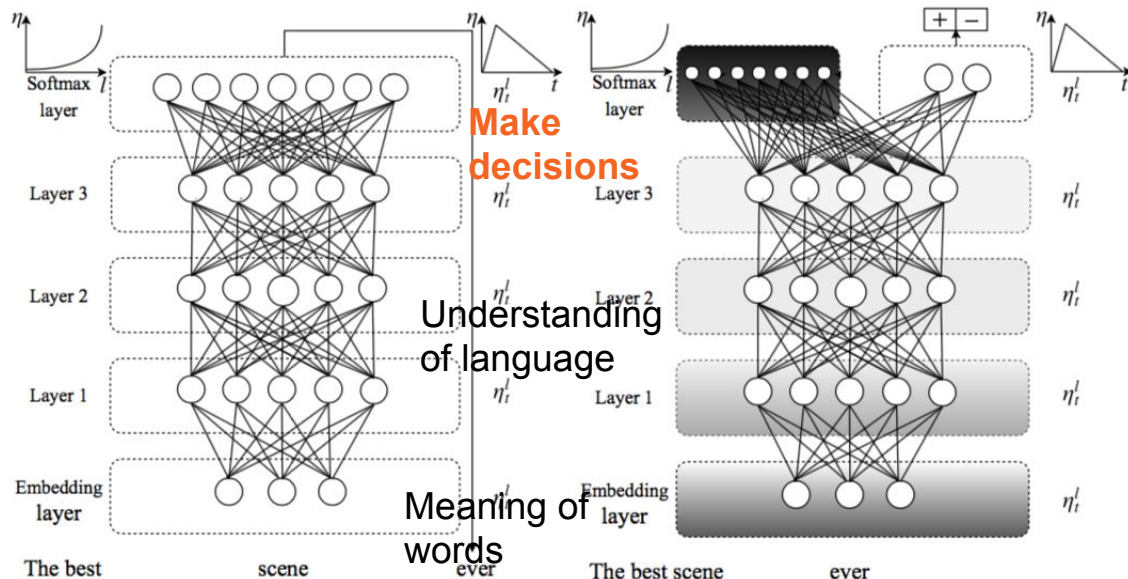


How does it work?



- The more knowledge the Neural Network has before it starts the less it has to work out.
- Here we train a “language model” first.

How does it work?



- The more knowledge the Neural Network has before it starts the less it has to work out.
- Here we train a “language model” first.

References

- Chang, Jonathan, et al. "[Reading tea leaves: How humans interpret topic models.](#)" Advances in neural information processing systems. 2009.
 - Howard, Jeremy, and Sebastian Ruder. "[Universal language model fine-tuning for text classification.](#)" Proceedings of the 56th Annual Meeting of the Association for Computational Linguistics (Volume 1: Long Papers). Vol. 1. 2018.
-

Questions

If you have any questions about the talk you can email me at chris@resolver.co.uk or you can message me on twitter [swartchris8](https://twitter.com/swartchris8)

Thanks
