One model to rule them all

Achieve state-of-the-art performance with less data than usually on multiple tasks by **transfer learning**





Agenda





01 Machine Learning

De-mystifying something a lot are afraid of



Illustration by Freepik Stories







How is this possible?

It's all statistics

Complex words, simple idea



It's all statistics

Complex words, simple idea



Example: Rectified Linear Unit (ReLU)

The most popular activation function of Neural Networks



The Relationship

Artificial Intelligence

The Relationship



The Relationship



So, what?





More about Neural Networks

Standard



Train from scratch **Random** initiation

Standard

Train from scratch **Random** initiation



Actually remembering things

1

Illustration by Freepik Stories

Reasons to use Transfer Learning

- Low-resource tasks
- Improve performance
- Reduce training-time













Take a pre-trained network **Fine-tune**



WC

Ω

維

BBC NEWS



Practical end-result



In practice, transfer learning has often been shown to **achieve similar performance** compared to a non-pretrained model **with 10x fewer examples** (Howard and Ruder, 2018).







Major Themes (NLP)
Language Models Words to **Pre-training** Word-In-Context Major Themes **Pretraining vs** Shallow to Deep (NLP) Target Task

Words to words-in-context

Context matters!

Words to words-in-context

Context matters!

I think this is a very good movie

Words to words-in-context

Context matters!

I don't think this is a very good movie

Understand a language Contextual Representations Versatile

Understand a language Contextual Representations Versatile

I think this is a very good movie \rightarrow **Positive**

Understand a language **Contextual Representations** Versatile I think this is a very novie \rightarrow **Positive**

Understand a language Contextual Representations Versatile

I think this is a very [MASK] movie

Shallow to Deep

Quickly models have grown deeper.

2016: 2-3 layers



Shallow to Deep

Quickly models have grown deeper.

2016: 2-3 layers 2019: ~60 layers (24 Transformers Blocks, BERT/GPT2)



Pretraining vs Target Task

Pre-training similar to target task

Sentence Representation is not useful for word-level predictions, and v.v.

How do we use "Transfer Learning"?

Feature Extraction

Fine-Tuning

Feature Extraction



¥

Ystads kommun

Feature Extraction



How to fine-tune?





¥

Input 1 Hidden 4 Hidden 1 Output 1 Input 2 Ystads kommun Input 3 Output 2 Hidden 6 Hidden 3 Input 4



¥

Input 1 Hidden 1 Hidden 4 Output 1 Input 2 Hidden 5 Ystads kommun Input 3 Output 2 Hidden 6 Hidden 3 Input 4



Fine-tune progressively in time



Lower Learning Rate

Fine-tune: progressively in intensity



Lower Learning Rate

Fine-tune: progressively in intensity



Lower Learning Rate

Fine-tune: progressively in intensity





Fine-tune: progressively vs. a pretrained model



Trade-offs: Working with pre-trained models

Feature Extraction

Fine-tuning

- Slower
- Space-efficient

- Better for dissimilar tasks

Getting more signal

Sequential adaptation

A related task with more data?

- 1. Fine-tune on related task
- 2. Fine-tune on target task

Improves if limited data (<u>Phang et al.</u>, <u>2018</u>) Improves sample efficiency on target task (<u>Yogatama et al.</u>, 2019)

Multi-task fine-tuning

Fine-tune jointly on related task Language Model is a good choice (even w/o pre-train)

Led to improvement in multiple target tasks (Liu et al., 2019, Wang et al., 2019)

Semi-supervised learning

- More consistent predictions by perturbing unlabelled examples
 - E.g. noise, masking or data-augmentation

Ensembling

An ensemble of models

- Different hyper-params
- Different pre-trained models
- Different target-tasks

Distilling

- Large model distilled into a small model
- Distilled model is a lot simpler

(Tang et al., 2019)



Extremes

	Date of original paper	Energy consumption (kWh)	Carbon footprint (lbs of CO2e)	Cloud compute cost (USD)
Transformer (65M parameters)	Jun, 2017	27	26	\$41-\$140
Transformer (213M parameters)	Jun, 2017	201	192	\$289-\$981
ELMo	Feb, 2018	275	262	\$433-\$1,472
BERT (110M parameters)	Oct, 2018	1,507	1,438	\$3,751-\$12,571
Transformer (213M parameters) w/ neural architecture search	Jan, 2019	656,347	626,155	\$942,973-\$3,201,722
GPT-2	Feb, 2019	-	-	\$12,902-\$43,008

GPT3 - Estimated 4,600,000 USD

Downstream Applications

- Hubs ("blackbox")
- Checkpoints
- Third party library

Applications for us

- Routing issues for municipalities or companies
- Data Understanding
 - Sentiment
 - Keywords
 - Tagging
 - Severity Ranking
- & much more

Final Words

There exists no "magic wand" Be knowledgeable Never forget baselines

Thanks!

Do you have any questions?

hampus.londogard@afry.com +46 733 673 179






Live Example

Transformer on Swedish: Colaboratory (google.com)