# Stock Embeddings

## Learning Distributed Representations for Financial Assets
## ( Dolphin et al. (2022))

Paper presentation for
Statistical Physics Reading Group

2023-05-03

The University of Edinburgh

By:
Martin Lellep
(http://lellep.xyz)

# Overview

- ML primer
- Introduction
- Maths
- Implementation
- Use cases

# ML primer  (ML := machine learning)

This is a ML project!

What you need for ML project:

- Problem to solve

- Data

- Algorithm:
  - Prediction algorithm
  - Loss fct
  - Optimiser

# Introduction

## Problem to solve:

Turn an asset $a_i$ like "Apple Inc" into a vector $\underline{e}_{Apple} \in \mathbb{R}^N$. Generally, $\underline{e}_{a_i} \in \mathbb{R}^N$ should be dense, i.e. all its entries should be non-zero & thereby be used.

E.g. $\underline{e}_{a_i}$ should <span style="color:red">not</span> be a unit vector (which are also called one-hot encodings).

## Data:

Prices of 500 largest US companies against time, i.e. $P_t^{a_i}$.

Specifically, this set of data is called S&P 500.

# Inspiration

Comes from "Natural language processing" (NLP), e.g. "continuous bag of words" approach.

There, words are turned into vectors ("word2vec" by [Mikolov 2013], [Pennington 2014]),

e.g. $\underline{e}_{hi} = (0.7, 0.3)^T$ &

$$\underline{e}_{hello} = (0.8, 0.35)^T \; \&$$

$$\underline{e}_{bye} = (-0.1, 0.5)^T.$$

Useful for similarity,

$$sim(\underline{e}_i, \underline{e}_j) = \frac{\underline{e}_i \cdot \underline{e}_j}{\|\underline{e}_i\| \cdot \|\underline{e}_j\|} \sim \text{cosine of angle between } \underline{e}_i \, \& \, \underline{e}_j$$
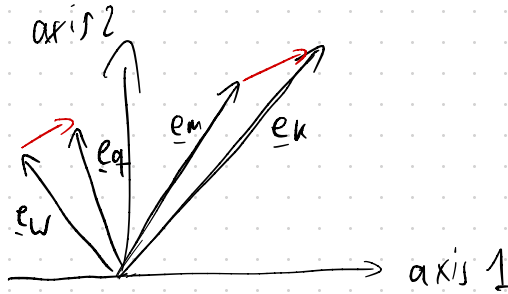
So that
$$sim(\underline{e}_{hi}, \underline{e}_{hello}) = 0.992 \quad \Leftarrow \text{similar}$$
$$sim(\underline{e}_{hi}, \underline{e}_{bye}) = 0.206 \quad \Leftarrow \text{not so similar}$$

& arithmetics [Mikolov 2013]:

$$\underline{e}_{king} - \underline{e}_{man} + \underline{e}_{woman} \approx \underline{e}_{queen}$$



Also useful in many downstream ML tasks that operate on language.

E.g. Transformer models that power large language models like GPT models, & hence also ChatGPT, use word (or to be precise: token) embeddings [Vaswani (2017) NIPS].

## Framework & prediction algo

o $\mathcal{U} = \{ a_1, ..., a_{|\mathcal{U}|} \}$   ~ asset universe

o $P_{a_i} = \{ p_0^{a_i}, ..., p_T^{a_i} \}$   ~ prices

$\quad\rightarrow\quad \underline{r}_{a_i} = \{ r_1^{a_i}, ..., r_T^{a_i} \}$   w/

$$r_t^{a_i} = \frac{p_t^{a_i} - p_{t-1}^{a_i}}{p_{t-1}^{a_i}} \quad \sim \text{returns}$$

o Construct context dataset by selecting $\forall a_i$ & $t$   $S(a_i, t)$ the $C$ assets $a_j$ that minimise $| r_t^{a_i} - r_t^{a_j} |$.

$\rightarrow$ one obtains dataset of size $|\mathcal{U}| \times T$

o $\underline{x}_{a_i}$ ~ one-hot encoded vector, i.e. is components $(\underline{x}_{a_i})_j = \delta_{ij}$

- Define embedding matrix

$$\underline{\underline{W}} = \begin{pmatrix} - & e_1 & - \\ & \vdots & \\ - & e_{|u|} & - \end{pmatrix} \in \mathbb{R}^{|u| \times N}$$

- Given $S(a_i, t)$, compute hidden state as mean embedding

$$\underline{h} = \underline{\underline{W}}^T \left( \frac{1}{c} \sum_{i=1}^{c} \underline{x}_{a_{j_i}} \right) \quad \text{with} \quad a_{j_i} \in S(a_i, t)$$

- Lastly compute neural network (NN) prediction as

$$p(\text{target} \mid \text{context}) = \text{softmax} \left( \underline{\underline{W}} \, \underline{h} \right)$$

with $\left( \text{softmax} \left( \underline{z} \right) \right)_i = \dfrac{e^{z_i}}{\sum_{j=1}^{k} e^{z_j}}$ for $i = 1, \dots, K$

& $\underline{z} = (z_1, \dots, z_k) \in \mathbb{R}^K$, i.e. softmax produces a discrete probability distribution.

Two possible noise reduction strategies:

o Noise red 1: different weighting

o Noise red. 2: Exclude $S(a_i, t)$ if $r_t^{a_i}$ statistically common

# Loss fct

Loss function defines optimisation objective. Here we want to predict target asset given context to then obtain $\underline{\underline{W}}$ which stores the embeddings.

Consider a training sample $S(a_i, t)$ of target asset $a_i$ at time $t$.

We now compare the predicted probability over assets from the NN,

$$P_p\left(a_j \mid S(a_i, t), \theta\right) \in \mathbb{R}^{|u|}$$

$\uparrow$ predicted

$\uparrow$ NN parameters $\underline{\underline{W}}$ to optimise

to the ground truth,

$$P_t\left(a_j \mid S(a_i, t)\right) = \left(\delta_{ik}\right)_{k=1,\dots,|u|} \in \mathbb{R}^{|u|}.$$

$\uparrow$ true

The loss function is chosen to be the <span style="color:red">categorical crossentropy</span> so that the loss of a single sample is computed as:

$$\ell(p_t, p_p) = -\sum_a p_t(a) \cdot \log p_p(a)$$

Motivation behind categorical crossentropy loss: Use Kullback-Leibler (KL) divergence

$$KL(p_t, p_p) = \sum_a p_t(a) \log\left(\frac{p_t}{p_p}\right)$$

$$\underset{\substack{\log\left(\frac{M}{N}\right) \\ = \log M - \log N}}{\curvearrowright} = \sum_a p_t(a) \cdot \left[\log p_t - \log p_p\right]$$

$$= \underbrace{\sum_a p_t(a) \log p_t(a)}_{= \text{const} = 0} - \sum_a p_t(a) \log p_p(a)$$

$= \text{const} = 0$
$\Rightarrow$ omitted for optimisation 😊.

$=: H(p_t, p_p)$
$- H(p_t)$

# Optimisation

Optimise $\underline{\underline{W}}$ to minimise loss to approach $P_t$ with $P_p$, which is actually only a proxy to obtain $\underline{e}_{a_i} \; \forall \, a_i$.

Basic idea is called <span style="color:red">stochastic gradient descent (SGD)</span>. Consider loss over batch:

$$L(\theta) = \frac{1}{b} \sum_{i=1}^{b} \ell \left( p_t(a_j | D_i)), P_p(a_j | D_i, \theta) \right)$$

w/ $b \hat{=}$ batch size, which is a hyper-parameter that is to be tuned manually, and $D_i$ the i-th sample from data batch $D$. Side note: Joining all batches, one obtains the full training dataset $\mathcal{D}$,

$$\mathcal{D} = [ \underbrace{S(a_1, t_1), S(a_1, t_2), \ldots}_{D} \;_{\cdots}\; \underbrace{S(a_{|n|}, T)}_{D} ].$$

Lastly, the optimisation step is performed as

$$\theta_{new} = \theta_{old} - \eta \nabla_\theta L(\theta) ,$$

w/ $\eta$ as learning rate that is yet another hyperparameter that is to be tuned manually.

Side note 1: There are much more sophisticated optimisers than SGD. Most people use those more sophisticated ones.

Side note 2: An algorithm called "back propagation" is used to compute $\nabla_\theta L(\theta)$ for NNs of almost arbitrary topology.

# Implementation

Implemented to compute one batch,
that is then used for optimisation
step, simultaneously.

# Use cases & benchmarks

Use embeddings as measure for similarity
between assets; traditionally, corre-
lations are used for that.

Quality of embeddings:

- Neighbours; Table I
- Arithmetics; Table II
- Visualisation; Fig. 3

Potential use:

- Construct hedged portfolio: Fig. 4

$$\left( \begin{array}{l} \text{Tables \& Figs refer to} \\ \text{Dolphin et al. (2022).} \end{array} \right)$$