

Yewno

Data Science Assignment

Thanks for your interest in Yewno. At Yewno, we don't believe in arbitrary, onerous "what is the difference between supervised and unsupervised learning?" type interviews. At work every day, you'll be dealing with a range of challenges, some modeling, some developing, some testing, hopefully all fun. The objective of this assignment is to see how you deal with challenges in a realistic setting, rather than in an artificial one hour interview.

The process goes like this:

1. You: Thoroughly read the exercise below, if you have any questions, email haris@yewno.com.
2. You: Complete the exercise within 3 days of receiving this document.
3. We: Contact you to setup a time to chat about your submission.

We do not expect you to develop the ultimate solution to the problem below. We are however interested in seeing a number of things:

1. How do you approach a problem?
2. How do you manage your work?
3. Are you aware of potential pitfalls your solution might have and could you propose alternative paths?

Introduction

Data and data processing is the foundation of Yewno. With our goal to ingest the world's knowledge, we are working to consume both public and private data sources in both batch and streaming methods. Both data pipelines are built around sets of algorithms that are ran against the datasets to build the Yewno inference engine.

One of the key roles within Yewno will be developing cutting edge machine learning algorithms to extract useful, decision-making, knowledge from raw data. Scalability of the algorithms is a must, and this role will work closely with the data engineering team to tune the algorithms into performant, production-ready systems.

Task

You are given a real-world undirected semantic graph represented via an adjacency matrix wherein links encode similarity between concepts. You are asked to explore the (*emerging*) properties of such a graph (aka *network*) in terms of its embedding on a surface with a given *genus* by using the techniques provided in the references attached [1,2]. You are also asked to infer potential connections between properties of hyperbolic embeddings and semantic information hidden in the network.

You are free to choose the programming language/library you are most comfortable with. Please don't forget to provide us with a brief description of your approach.

Additional question:

1. Is your system scalable w.r.t. network dimension / genus? If not, how would address the scalability (in terms of algorithms, infrastructure, or both)?

When you are finished, send us a link to the code repository - [Github](#) or [BitBucket](#) are great.

Please be sure to **save the outputs of your test** run so we can take a look. Remember, we care for as much about **how you think** about the problem as the code itself! Document the code as needed and be ready to discuss your project. Impress us!

Above all, have fun and reach out if you have any questions. The task is designed to take approximately 1 day to complete.

References

[1] Tomaso Aste, Ruggero Gramatica, and T. Di Matteo. Exploring complex networks via topological embedding on surfaces. Phys. Rev. E 86. 2012

[2] Tomaso Aste, Ruggero Gramatica, and T. Di Matteo. Random and frozen states in complex triangulations. Phil. Mag. 2012