



Paris, le 18/01/2022

## Offre de stage

*Sujet : Apprentissage dynamique de caractéristiques par champs aléatoires conditionnels pour la détection des sentiments dans les conversations téléphoniques*

*Possibilité de poursuivre sur une thèse*

### La Chaire Data Science and Artificial Intelligence for Digitized Industry and Services

Portée par Florence d'Alché-Buc, enseignante-chercheur dans le département Image, Données, Signal de Télécom ParisTech, la chaire DSAI réunit cinq partenaires industriels : Airbus Defence & Space, Engie, Idemia, Safran et Valeo Finance. Son objectif général est de développer, en liaison étroite entre les Parties, une formation et une recherche de niveau international.

Ses quatre principaux axes de recherche sont :

1. Analyse et prévision de séries temporelles (Predictive Analytics on Time Series) ;
2. Exploitation de données hétérogènes, massives et partiellement étiquetées (Exploiting Large Scale and Heterogeneous, Partially Labelled Data) ;
3. Apprentissage pour une prise de décision robuste et fiable (Learning for Trusted and Robust Decision) ;
4. Apprentissage dans un environnement dynamique (Learning through Interactions with a Changing Environment).

## Description du stage

### Encadrement

Ekhine Irurozki, Ons Jelassi, Pascal Bianchi

### Lieu et dates du stage

Telecom Paris, 19 Place Marguerite Perey, 91120 Palaiseau

Date de début du stage : 2022

### Équipe(s) d'accueil de la thèse

Equipe Signal, Statistique et Apprentissage (S<sup>2</sup>A)

### Mots clés

Machine learning; combinatorial optimization;

### Sujet détaillé

The motivation of this proposal is to bridge the gap between two major areas in applied mathematics: Deep Learning (DL) and Combinatorial Optimization (CO). In Combinatorial Optimization (CO) the goal is to find an optimal object (i.e., graph) among a finite set of such objects (i.e., all the possible graphs of n nodes). One of the best-known problems is the Traveling salesman problem, TSP. In the TSP, the input to the problem is a set of cities and the distances among them (represented by a graph) and the goal is to find the tour (a permutation of cities) that travels to all the cities in which the total distance is minimized. In a regular computer, a problem with 50k cities takes more than 20 years! Since huge problems of this kind arise in many areas (for example network planning or logistics) the most common approach when solving CO problems is to obtain a 'good' solution in a 'reasonable' time along with an estimation/bound on how far our solution is from the optimal solution. Deep learning is the hottest topic in research and industry in applied mathematics. There are amazing results for real-valued data, image, and text data and are state-of-the-art in all these domains. However, there are no successful, ad-hoc results for combinatorial data except for graphs with being Graph Neural Networks [8].

**Literature gap and goal** Currently, there exist several approaches for modeling combinatorial data. For example, permutations are modeled using complex network structures such as attention mechanisms [1, 9] or by real-valued vectors via the argsort function, thus using an infinite cardinality space to model a finite cardinality space. The goals of this internship are (1) to propose transformations of combinatorial data into different spaces that are more adequate for NN, and (2) to implement efficient NN algorithms for combinatorial data with theoretical guarantees. The tools used for the transformations will be bijections among combinatorial structures, embeddings, and probabilistic modeling.

The proposal is divided into two tasks: In the first stage, the intern will review the different transformations, enumerations, and properties of them, including but not limited to Birkhoff theorem for doubly-stochastic matrices, sorting algorithms, different injections among planar graphs, lattice paths, permutations, sequences such as Narayana, Catalan or Stirling numbers,...[2, 3, 6, 7, 8]. A theoretical analysis of the convergence properties for existing learning algorithms is in order.

As a second step, the intern will code an end-to-end NN using the embedding studied in the first task. The intern will use benchmarking problems and compare the results to baseline methods and state-of-the-art methods [1, 4, 5, 9], elaborating the results in terms of quality of the solution, transferability and time performance. For the simulations, he/she will have access to computational resources

This proposal describes some first steps in a hot area of research in which both methodological and practical lines follow naturally.

## Profil du candidat

Student having master 2 research

Statistical learning, bases of probability and interest in combinatorics and discrete algebra

Good level of programming

Good command of English

## Candidatures

A envoyer à irurozki@telecom-paris.fr

- Curriculum Vitae

- Lettre de motivation personnalisée expliquant l'intérêt du candidat sur le sujet (directement dans le corps du mail)

- Relevés de notes des années précédentes

- Contact d'une personne de référence

Les candidatures incomplètes ne seront pas examinées.

## Références

- [1] I. Bello, H. Pham, Q. V. Le, M. Norouzi, and S. Bengio. Neural combinatorial optimization with reinforcement learning. In 5th International Conference on Learning Representations, ICLR 2017 - Workshop Track Proceedings, 2019.
- [2] J. Ceberio, E. Irurozki, A. Mendiburu, and J. Lozano. A review of distances for the Mallows and Generalized Mallows estimation of distribution algorithms. Computational Optimization and Applications, 62(2), 2015.
- [3] E. Irurozki, B. Calvo, and J. Lozano. Learning probability distributions over permutations by means of Fourier coefficients, volume 6657 LNAI. 2011.
- [4] E. Irurozki and M. López-Ibáñez. Unbalanced Mallows Models for Optimizing Expensive Black-Box Permutation Problems. In The Genetic and Evolutionary Computation Conference (GECCO21), Lille, France, jul 2021.
- [5] M. Malagon, E. Irurozki, and J. Ceberio. Alternative Representations for Codifying Solutions in Permutation-Based Problems. In 2020 IEEE Congress on Evolutionary Computation, CEC 2020 - Conference Proceedings, 2020.
- [6] A. Seshadri, S. Ragain, and J. Ugander. Learning Rich Rankings. In Neural Information Processing System (NeurIPS), volume 5, pages 1–12, 2020.
- [7] N. J. A. Sloane. On-Line Encyclopedia of Integer Sequences, <http://oeis.org/>, 2009.
- [8] P. Veličković, G. Cucurull, A. Casanova, A. Romero, P. Lio, and Y. Bengio. Graph attention networks. arXiv preprint arXiv:1710.10903, 2017.
- [9] O. Vinyals, M. Fortunato, and N. Jaitly. Pointer networks. In Advances in Neural Information Processing Systems, 2015.