**PUCP**

**SEMANA INTERNACIONAL**

**CURSOS INTERNACIONALES 2022-2**

**FORMULARIO DE POSTULACIÓN**

| | |
|---|---|
| Nombre del curso | Seminario Internacional de Ingenieria -Introduction to Large Language Models and Agents- |
| Código del curso | |
| Facultad | Ciencias e Ingenieria |
| Nombre del/de la profesor/a | Ronald Cardenas Acosta |
| Mini bio del/de la profesora *(1 párrafo que señale credenciales académicas, institución extranjera de filiación, campo de especialización, entre otros aspectos relevantes)* | Ronald received his B.S. degree in Mechatronics from the Department of Mechanical Engineering, Universidad Nacional de Ingenieria, Peru, in 2015. Later he received his M.Sc. in Language and Communication Technologies from Charles University in Prague in 2019. He obtained his Ph.D. degree from the Centre for Doctorate Training in Natural Language Processing programme at the University of Edinburgh, in 2024. Since then, he joined Huawei's London Research Centre as a Research Scientist. His reseach includes agentic generative AI with a focus on human preference alignment and reasoning. |
| Foto del/de la profesor/a |  |
| Sumilla del curso | During the course development, students will learn the basic inner workings of the transformer, the technology behind Large Language Models (LLMs). The course will cover a comprehensive view of language modeling and other NLP applications, as well as the latest techniques to implement LLM-based |

| | |
|---|---|
| | solutions. |
| Descripcion del curso | Foundation models, and language models in particular, are at the centre of the current AI revolution.<br>Hence, the understanding of the architecture and optimization techniques of these neural networks has become paramount not only for researchers but AI practicioners in general.<br>This course aims to introduce key concepts crucial for the effective use of these technologies. |
| Metodología<br><br>*(señale el número total de horas sincrónicas y, si estas fueran menos de 16, explique cómo se llevarán a cabo las horas asincrónicas)* | The course will consist of<br><br>- Four synchronous classes, 3 hours/each, 12 hours total;<br><br>- Office hour session, 4 hours total, during which students can ask questions about the assignments. |
| Evaluación<br><br>*(señale la fórmula y tipo de evaluación/es)* | Evaluation consists of two coding assignments.<br>Assignments<br>A1. Implementing Multi-Head Attention, Perplexity<br>A2. Evaluating Zero-shot, ICL, CoT<br>Final grade:<br>Grade = (A1 + A2) / 2 |
| Horario<br><br>*(señale días y horas de clase en la semana previa al inicio oficial del ciclo 2022-2; eventualmente, puede haber una hora más en la semana 1)* | <u>Classes (4), 3h/each</u><br>1. March, 17 (Monday)    - 15:00-18:00<br>2. March, 18 (Tuesday)    - 15:00-18:00<br>3. March, 19 (Wednesday) - 15:00-18:00<br>4. March, 20 (Thursday)  - 15:00-18:00<br><u>Office hours,  4 hours</u><br>5. March, 21 (Friday)  - 14:00-18:00 |
| Público objetivo<br><br>*(señale qué especialidades o facultades tienen alumnos/as potencialmente interesados en el curso)* | El curso está dirigido a estudiantes de ingeniería Informatica o carreras afines (Telecomunicaciones, Electrónica, Informática, Mecánica, etc). Se recomienda haber llevado cursos en programacion, estadistica y probabilidad. |
| Número de vacantes<br><br>*(señale el número de vacantes* | 25 |

**DIRECCIÓN**
**ACADÉMICA DE**
**RELACIONES**
**INSTITUCIONALES**
SECCIÓN DE INTERNACIONALIZACIÓN ACADÉMICA

**PUCP**

| | |
|---|---|
| *que se ofrecerá para el curso)* | |
| Contenido | # PROGRAMA ANALÍTICO<br><br>## UNIT 1: Basic Concepts<br>- Statistics and Probabilities<br>- Machine Learning for Language<br>- Text clasification, clustering<br>- Neural Networks and Backpropagation<br><br>## UNIT 2: Language Modeling<br>- N-gram models<br>- Word Vectors: skip-gram, glove<br>- Convolutional neural networks<br>- Recurrent neural networks<br>- Transformers<br><br>## UNIT 3: Training and Fine-tuning<br>- Supervised Finetuning (SFT)<br>- Reinforcement Learning from Human Feedback (RLHF)<br>- Direct Preference Optimization (DPO)<br>- Parameter-efficient Finetuning (PEFT)<br>- Prompt Engineering<br>   * In-context learning (ICL)<br>   * Chain of Thought (CoT)<br><br>## UNIT 4: LLM applications<br>- LLM interfaces and resources<br>- NLP tasks<br>  * Sentiment Analysis<br>  * Question Answering<br>  * Summarization<br>  * Code generation<br>- Tool calling<br>- Agents<br>- Semantic Search<br>- Retrieval Augmented Generation |
| Nombre y correo electrónico de la persona de contacto en la facultad | evillota@pucp.edu.pe |