# State-of-the-Art GenAI and LLMs Creative Projects, with Solutions

GAN started with seed = 103 (X axis = epoch, Y axis = loss function)

GAN started with seed = 102 (X axis = epoch, Y axis = loss function)

d_history          g_history

# Preface

With 23 top projects, 96 subprojects, and over 6000 lines of Python code, this vendor-neutral coursebook is a goldmine for any analytic professional or AI/ML engineer interested in developing superior GenAI or LLM enterprise apps using ground-breaking technology.

## Key Features

- 23 enterprise-grade projects, 96 subprojects all with solutions, 6000 lines of Python code, with real-world datasets and case studies.
- Ground-breaking technology covering RAG, multi-LLMs, knowledge graphs, synthetic data generation (tabular, geospatial, time series), generative adversarial network (GANs), model evaluation, adaptive loss functions, and more.
- Learn how to create fast, memory-frugal apps from scratch with explainable AI, to deliver better results.

## Target Audience

This vendor-neutral coursebook is a goldmine for any analytic professional or AI/ML engineer interested in developing superior GenAI or LLM enterprise apps using ground-breaking technology. This is not another book discussing the same topics that you learn in bootcamps, college classes, Coursera, or at work. Instead, the focus is on implementing solutions that address and fix the main problems encountered in current applications. Using foundational redesign rather than patches such as prompt engineering to fix backend design flaws.

You will learn how to quickly implement from scratch applications actually used by Fortune 100 companies, outperforming OpenAI and the likes by several order of magnitudes, in terms of quality, speed, memory requirements, costs, interpretability (explainable AI), security, latency, and training complexity.

This book is also an invaluable resource to instructors and professors teaching AI, or for corporate training. Also, it is useful to prepare for job interviews or to build a robust portfolio. And for hiring managers, there are plenty of original interview questions. The amount of Python code accompanying the solutions is considerable, using a vast array of libraries as well as home-made implementations showing the inner workings and improving existing black-box algorithms. By itself, this book constitutes a solid introduction to Python and scientific programming. The code is also on GitHub.

Knowledge of Python and machine learning concepts is required to fully understand the use cases and code examples. However, the first chapter starts from scratch: getting a Python environment. Mathematics are used only when needed, rarely going beyond a first course in calculus, linear algebra and statistics. Advanced concepts, when useful, are introduced and explained in simple English.

## What you will learn

- How to create fast, memory-frugal apps from scratch with explainable AI, to deliver better results.
- Solving a problem from start to finish, with Python: Learn how to choose the best algorithms, fine-tune, evaluate, and how to produce great visualizations including data animations (videos) to efficiently convey insights to stakeholders.
- How to use dozens of Python libraries in a business-like setting, including your home-made, and from the author. Covering anything from smart crawling, deep neural networks, scientific computing, model evaluation, NLP, and machine learning.

Figure 1: The first three steps in the main xLLM project

## Approach

Each project starts with a summary of the problem and an introduction to the concepts needed for completion. In several cases, projects do not consist of writing code from scratch, but instead, to understand the code offered by the author and build on top of it, and/or improve the speed. Each project is broken down into multiple steps of increasing scope and complexity. Two glossaries, a large index, and modern references are included. Considerable efforts have been mode to avoid jargon and advanced mathematics, using simple English as much as possible. Yet, all the terminology you need to know is introduced when needed.

## Description

With tutorials, enterprise-grade projects with solutions, and real-world case studies, this coursebook covers state-of-the-art material on GenAI, generative adversarial networks (GAN), data synthetization, and much more, in a compact format with many recent references. It includes deep dives into innovative, ground-breaking AI technologies such a xLLMs (extreme LLMs), invented by the author.

The focus is on explainable AI with faster, simpler, high-performance, automated algorithms. For instance: NoGAN, new evaluation metrics, xLLM (self-tuned multi-LLM based on taxonomies with application to clustering and predictive analytics), variable-length embeddings, generating observations outside the training set range, fast probabilistic vector search, or Python-generated SQL queries. I also discuss alternatives to traditional methods, for instance to synthesize geo-spatial data or music.

## Author Bio

Vincent Granville is a pioneering GenAI scientist and machine learning expert, co-founder of Data Science Central (acquired by a publicly traded company), Chief AI Scientist at MLTechniques.com and GenAItechLab.com, former VC-funded executive, author, and patent owner. Vincent's past corporate experience includes Microsoft, Visa, Wells Fargo, eBay, NBC, and CNET.

Vincent is also a former post-doc at Cambridge University. He published in Journal of Number Theory, Journal of the Royal Statistical Society (Series B), and IEEE Transactions on Pattern Analysis and Machine Intelligence. He is the author of multiple books, including "Synthetic Data and Generative AI" (Elsevier). Vincent lives in Washington state, and enjoys doing research on stochastic processes, dynamical systems, experimental math and probabilistic number theory.

# Contents