# Training day 10 report

## Supervised Machine Learning and OpenAI: End-to-End Workflow – Detailed Summary

This notebook outlines a complete process for applying supervised learning principles to a code generation task using OpenAI's language models. It emphasizes how prompt-completion pairs can be used to train or prompt an AI system to generate code, optimize it, and even debug it. The project demonstrates an end-to-end pipeline, starting from data preparation to evaluation, with each step focused on practical implementation in a machine learning workflow. The following is a comprehensive theoretical summary of each section.

### Dataset Preparation for Fine-Tuning

The first step in any supervised learning pipeline is data preparation. The task here is code generation, so the dataset consists of pairs of natural language prompts and the corresponding Python code completions. For example, a prompt might ask to write a function to add two numbers, and the completion would be the correct implementation of that function. This structure mirrors a typical supervised learning dataset where the model is expected to learn the mapping from inputs to outputs.

The dataset is constructed manually using a small number of examples. These examples are then converted into a structured format using a data manipulation library, where each row represents a prompt and its corresponding completion. This format is essential for training because it clearly defines what the input and expected output are for the model. The dataset is saved in a standard format that can be used later for both training and testing.

### Train-Test Split and JSONL Conversion

Before training or fine-tuning a model, it is necessary to divide the data into training and testing subsets. A typical split involves using eighty percent of the data for training and the remaining twenty percent for testing. This separation allows for the evaluation of the model's performance on unseen data after training.

Since OpenAI models expect data in a specific format for fine-tuning, the training and testing datasets are converted to JSON Lines (JSONL) format. This format stores each example as a separate JSON object on a new line. For fine-

tuning with OpenAI, each object must include a sequence of messages, where one message is from the user and one is from the assistant. This conversational structure supports the chat-based format used by models like GPT.

Each prompt becomes a user message, and each completion becomes the assistant's reply. This formatting step is critical because improperly formatted data can result in errors during the fine-tuning process or suboptimal performance during inference.

## Environment Setup and API Configuration

To use the OpenAI API, necessary libraries such as OpenAI's Python package and dotenv are installed. The dotenv package is used to securely load environment variables, such as the API key, from a local file. This practice ensures that sensitive information is not hardcoded into scripts.

The OpenAI client is then initialized using the retrieved API key. This setup is a required step for any project that interacts with OpenAI's models through their API. Once the client is configured, it can be used to send prompts to the model and receive responses.

## Prompt Engineering for Code Generation

Prompt engineering involves crafting effective inputs that elicit the desired response from a language model. In this context, the goal is to generate Python code. A prompt is designed to ask the model to perform a specific task, such as writing a function to reverse a string.

A function is defined to send such prompts to the OpenAI model and retrieve the completion. This demonstrates how to use the OpenAI chat completion endpoint to perform code generation. The model is provided with a system message that defines its role as a helpful assistant and a user message that contains the actual prompt.

The response from the model typically includes the requested code along with a brief explanation. This part of the workflow illustrates how pretrained language models can be used directly for supervised tasks without additional fine-tuning, simply by using well-constructed prompts.

## Code Optimization with Language Models

Beyond basic generation, language models can also be used to optimize code. In this case, an example function is deliberately written in a redundant or

inefficient way. A prompt is then constructed to ask the model to simplify and optimize the code.

The response includes a revised version of the function that achieves the same functionality in a more efficient or readable manner. This application shows how language models can assist in software development tasks such as refactoring, which traditionally require manual effort and domain knowledge.

Prompting the model with optimization tasks introduces a level of automation in improving code quality. This can be particularly useful for software engineers looking to save time during development or for educators who want to teach coding best practices.

## Debugging with Language Models

Another valuable use of language models in code-related tasks is debugging. A snippet of buggy code is provided as input, along with a prompt asking the model to find and fix the error. In this example, the code attempts to divide by zero, which causes a runtime error.

The model responds with a corrected version of the function that includes a check for division by zero. This step demonstrates the model's ability to reason about runtime behavior and suggest modifications that make the code more robust.

Debugging assistance is an important feature of intelligent coding tools. It can help beginners understand common errors and improve the overall reliability of programs by suggesting safe coding practices.

## Evaluation of Generated Code

To validate that the generated or optimized code functions as expected, the notebook includes a mechanism to evaluate the code. This is done by extracting code blocks from the model's response and running them within a controlled environment.

Specific test cases are defined to check the output of the generated functions. For instance, a function that reverses a string is tested with several inputs to ensure that it behaves correctly. If the function passes all test cases, it confirms that the model's output is both syntactically and logically correct.

Automated evaluation is crucial when working with code generation tasks, as it provides an objective measure of performance. Without testing, it is difficult to verify whether the model is producing correct and useful code.

**Summary of the Workflow**

The notebook provides a structured approach to using large language models for supervised code generation tasks. It begins with the construction of a dataset, formatted for fine-tuning, and progresses through environment setup, prompt engineering, code generation, optimization, debugging, and automated evaluation.

This approach demonstrates that large language models can be integrated into traditional supervised learning workflows either by fine-tuning on custom data or by leveraging prompt engineering techniques. It also shows how different types of model interactions—generation, refinement, error correction—can be accomplished using simple input-output patterns.

The techniques in this notebook can be extended to larger datasets, more complex tasks, and different domains. Although the examples used are simple, they provide a foundation for building powerful applications in intelligent code assistance, automated documentation, and software development education.