



# AI Pitch Deck Analysis Model: Methodology and Insights Report

**Version: 1.0**

**Date: 2025-03-30**

**Author: Gitanshu Vaghasiya**

## 1. Executive Summary

This report details the methodology behind the development of an AI-powered Pitch Deck Analysis Model. The primary objective was to create a tool capable of automatically evaluating startup pitch decks (provided as PDFs) to generate an objective score (0-100) and qualitative feedback (strengths, weaknesses, suggestions). The model focuses on six key sections: Problem, Solution, Market Size, Business Model, Financials (Projections), and Team. The core technology relies on Google's Gemini Large Language Models (LLMs) – specifically a vision model ( `gemini-1.5-pro-latest` ) for robust Optical Character Recognition (OCR) from potentially image-based PDFs and a text model ( `gemini-1.5-flash-latest` ) for analysis, scoring, and feedback generation. The system includes strategies to handle API rate limits and is presented through a basic Flask web interface for user interaction. The model demonstrated plausible results when tested on a known early-stage pitch deck (Uber, 2008), correctly identifying key strengths and weaknesses typical for that stage.

## 2. Introduction

Reviewing startup pitch decks is a critical but often time-consuming task for investors, accelerators, and internal reviewers. Manual review can be subjective and prone to inconsistency. The goal of this project was to leverage the power of modern LLMs to automate the preliminary analysis of pitch decks, providing a standardized evaluation framework.

The objective was to develop an application that:

1. Accepts pitch decks in PDF format.
2. Reliably extracts text content, even from image-based slides using OCR.
3. Identifies and isolates key sections relevant to startup evaluation.
4. Scores each key section based on predefined, objective criteria using an LLM.
5. Calculates an overall weighted score reflecting the pitch's strength across core areas.
6. Generates actionable, qualitative feedback summarizing strengths and suggesting improvements.
7. Presents the analysis through a simple user interface.

This tool aims to augment, not replace, human judgment by providing a quick, structured first pass of a pitch deck's content.

## 3. Methodology

The development followed a structured pipeline, broken down into the following key stages:

### 3.1. Data Input & Handling

- **Input:** The model accepts standard PDF files as input via a web upload form.
- **Challenge:** PDFs vary widely; some contain selectable text, while others (especially exports from presentation software or scans) are essentially collections of images.
- **Solution:** The system was designed to handle both types, prioritizing OCR for robustness, assuming many real-world decks might be image-based.

Uploaded files are temporarily stored for processing and deleted afterward.

## 3.2. Text Extraction (OCR Focus)

- **Necessity:** Standard PDF text extraction libraries (like `PyPDF2`) fail on image-based PDFs. Therefore, Optical Character Recognition (OCR) is essential.
- **Approach:** A multi-step OCR process using Gemini Vision was implemented:
  1. **PDF to Images:** The `pdf2image` library (leveraging the external `Poppler` utility) converts each page of the input PDF into a high-resolution image (JPEG format).
  2. **Image-to-Text (OCR):** Each page image is sent individually to the Google Gemini Vision API (`gemini-1.5-pro-latest`). A specific prompt instructs the model to extract all visible text.
  3. **Rate Limit Handling:** Vision API calls can be rate-limited. To mitigate this, two strategies were employed:
    - **Inter-Page Delay:** A configurable delay (`INTER_PAGE_DELAY`) is introduced between processing each page to avoid hitting requests-per-minute limits.
    - **Exponential Backoff:** The API call function includes retry logic that waits progressively longer (`RATE_LIMIT_BACKOFF_MULTIPLIER`) if a rate limit error (HTTP 429) is encountered.
- **Output:** A list of strings, where each string contains the extracted text for the corresponding page. Empty strings are returned for pages where OCR fails or is blocked.

## 3.3. Text Preprocessing

- Extracted text (from OCR or potentially direct extraction) undergoes basic cleaning:
  - Conversion to lowercase.
  - Removal of excessive whitespace (multiple spaces, tabs, newlines replaced with single spaces).
  - Stripping leading/trailing whitespace.

- This standardization helps improve the consistency of input for subsequent LLM analysis steps.

### 3.4. Section Identification

- **Objective:** Identify which page(s) correspond to the six predefined `TARGET_SECTIONS` : Problem, Solution, Market Size, Business Model, Financial Projections, Team.
- **Approach:**
  1. The extracted text from all pages (potentially truncated for brevity) is formatted with page number indicators.
  2. This formatted text is sent to the Gemini Text API ( `gemini-1.5-flash-latest` ).
  3. The prompt specifically instructs the LLM to identify pages corresponding *only* to the `TARGET_SECTIONS` and return the mapping as a JSON object.
  4. The returned JSON is parsed, and the results are validated to ensure section names are correct and page numbers are within range. Pages are stored as 0-based indices.

### 3.5. Content Aggregation

- Using the validated section-to-page mapping, the preprocessed text content for all pages belonging to a specific section is concatenated. This creates a single block of text representing the content of each identified target section.

### 3.6. Scoring Model

- **Objective:** Evaluate the quality of each identified target section based on predefined metrics.
- **Approach:**
  1. For each identified section with aggregated content:
    - The section's text is sent to the Gemini Text API.
    - The prompt includes detailed, specific `SCORING_CRITERIA` for *that particular section*.

- The LLM is instructed to provide a score (integer 0-100) and a brief justification (2-3 sentences) based *only* on the provided criteria and text.
  - The output is requested strictly in JSON format ( `{"score": ..., "justification": ...}` ).
2. Robust JSON parsing is used to extract the score and justification from the LLM's response. Default values (score 0, error message) are assigned if parsing fails or the LLM call is unsuccessful.
- **Overall Score Calculation:**
    - A weighted average is calculated using the scores obtained for the `TARGET_SECTIONS` .
    - The predefined `SECTION_WEIGHTS` (summing to 100) determine the contribution of each section.
    - The score is normalized based on the total weight of the target sections that were successfully identified and scored (even if the score was 0), ensuring a consistent 0-100 scale reflecting the quality of the evaluated components.

### 3.7. Strength & Weakness Analysis (Feedback Generation)

- **Objective:** Provide high-level qualitative feedback.
- **Approach:**
  1. A summary containing the scores and justifications for the identified `TARGET_SECTIONS` is compiled.
  2. This summary is sent to the Gemini Text API.
  3. The prompt asks the LLM to act as an analyst, synthesize the provided information, and identify:
    - 2-3 key strengths (based on high scores/positive justifications).
    - 2-3 critical weaknesses (based on low scores/negative justifications).
    - Actionable suggestions for each weakness.

4. The output is requested strictly in JSON format ( `{"strengths": [...], "weaknesses": [...]}` ).
5. The parsed lists of strengths and weaknesses are included in the final results.

### 3.8. Technology Stack

- **Programming Language:** Python 3
- **Core AI:** Google Gemini API ( `google-generativeai` library)
  - `gemini-1.5-pro-latest` (or similar vision model) for OCR.
  - `gemini-1.5-flash-latest` (or similar text model) for Section ID, Scoring, Feedback.
- **PDF Processing:** `pdf2image` (requires `Poppler` system dependency)
- **Image Handling:** `Pillow`
- **Web Framework:** Flask
- **Configuration:** `python-dotenv`
- **Utilities:** `os` , `re` , `json` , `time` , `logging`

### 3.9. Deployment & UI

- A basic web interface was developed using Flask.
- It provides a simple form to upload a PDF file.
- Upon submission, it triggers the analysis pipeline defined in the `PitchAnalyzer` class.
- Results (overall score, section scores/justifications, feedback) are displayed on a separate results page using HTML templates.
- Basic error handling and user feedback (flash messages, loading indicator) are included.
- The code is structured using OOP ( `PitchAnalyzer` class) and modular Flask blueprints ( `config.py` , `analyzer/core.py` , `routes/main.py` , `run.py` ) for better organization.

## 4. Implementation Details & Challenges

- **LLM Selection:** Gemini was chosen for its strong multimodal capabilities (Vision API for OCR) and capable text models. The use of separate models for vision and text optimizes for cost and speed where possible ( `gemini-1.5-flash` being faster/cheaper for text tasks). A key trade-off is the potential cost/speed compared to local OCR (Tesseract), but Gemini Vision may offer higher accuracy on complex layouts.
- **Prompt Engineering:** Crafting effective prompts was crucial for each LLM task (OCR, Section ID, Scoring, Feedback). Prompts needed to be clear, provide sufficient context and examples, and strictly enforce the desired output format (especially JSON) to ensure reliable parsing. This was an iterative process.
- **OCR Robustness & Rate Limits:** Handling image-based PDFs reliably required implementing the `pdf2image` → Gemini Vision pipeline. The primary challenge here was managing API rate limits, necessitating the `INTER_PAGE_DELAY` and exponential backoff strategies, which significantly impacts processing time for longer decks.
- **JSON Output Enforcement:** LLMs do not always perfectly adhere to output format instructions. Robust parsing logic (using regex to find JSON blocks within potential surrounding text or markdown) was implemented to handle variations in LLM responses.
- **Code Structure:** Refactoring the initial script into a class ( `PitchAnalyzer` ) and a modular Flask application improved code organization, reusability, and maintainability.

## 5. Results & Evaluation (Illustrative)

The model was tested using the well-known early Uber pitch deck (circa 2008).

- **OCR:** Successfully extracted text from the image-based slides, enabling further analysis.
- **Section Identification:** Correctly identified the pages corresponding to the target sections (Problem, Solution, Market Size, etc.).
- **Scoring:** Generated an overall score of **51/100**. Individual section scores highlighted strengths in Problem (70) and Solution (65) clarity, but weaknesses in Market Size (45), Business Model (40), Go-to-Market (N/A - not a target

section, but would be low), Financials (30), and Traction (45 - due to lack of quantifiable metrics).

- **Feedback:** The generated feedback accurately reflected these findings, praising the clear problem definition while pointing out the lack of specific data, detailed financial projections, and concrete go-to-market plans. The suggestions provided were relevant for improving such an early-stage deck.

This result suggests the model provides a plausible and critical assessment aligned with the known characteristics of that specific deck at that point in time.

## 6. Limitations

- **Subjectivity:** While criteria are defined, LLM interpretation still involves a degree of subjectivity inherent in language understanding. Scores might vary slightly between runs or models.
- **Financial Audit:** The model does *not* perform a deep financial audit. It checks for the presence, clarity, and apparent realism of projections based on stated criteria, but cannot verify calculations or complex financial soundness.
- **Visual Nuance:** The OCR extracts text but does not deeply analyze the *design*, layout quality, or intricate details within complex charts/diagrams beyond basic text recognition.
- **Contextual Understanding:** While powerful, LLMs may struggle with highly niche market contexts or extremely novel technological concepts not well-represented in their training data.
- **Bias:** LLMs can reflect biases present in their training data. The analysis might inadvertently favor certain presentation styles or business models.
- **OCR Errors:** OCR is not perfect and errors can propagate into the analysis.
- **Scope:** The analysis is limited to the 6 predefined sections; other important aspects (e.g., detailed competitive analysis visuals, exit strategy) might be present in the deck but are not explicitly scored in this configuration.

## 7. Future Work



- **Refined Prompts & Criteria:** Continuously iterate on scoring criteria and LLM prompts based on feedback and testing with more diverse decks.
- **Configurability:** Allow users to customize section weights or even scoring criteria via the UI.
- **Advanced Visual Analysis:** Explore integrating models capable of interpreting charts and graphs more deeply (beyond simple OCR).
- **Comparative Analysis:** Store results and allow comparison across different pitch decks.
- **Alternative Extraction:** Add optional fallback to faster text extraction (e.g., `PyPDF2`) if a PDF is text-based, using OCR only when necessary.
- **UI Enhancements:** Improve the user interface with progress indicators, result visualization, and user accounts.
- **Error Handling:** Implement more granular error handling and reporting to the user.
- **Benchmarking:** Evaluate performance against human expert scoring on a larger set of diverse decks.

## 8. Conclusion

The AI Pitch Deck Analysis Model successfully demonstrates the feasibility of using Large Language Models (specifically Google Gemini) for automated, structured evaluation of startup pitch decks. By combining multimodal OCR with text analysis capabilities, the tool can process diverse PDF formats and provide a valuable preliminary assessment, including a quantitative score and qualitative feedback focused on core business areas. While limitations exist, particularly concerning deep financial validation and inherent LLM subjectivity, the model serves as a powerful tool to augment human review, increase efficiency, and provide consistent, data-driven insights for initial pitch screening. Future iterations can further enhance its accuracy, scope, and user experience.