

# Resumo Detalhado do Plano Estratégico para Análise de Vagas de Emprego no Brasil com Big Data e Python

---

Este resumo tem como foco apresentar o **Plano de Ação** detalhado, mantendo a concisão, mas incorporando os passos práticos e as ferramentas específicas necessárias para a implementação do projeto.

## 1. Plano de Ação: Aquisição de Dados

---

A estratégia de aquisição deve ser implementada em duas frentes, utilizando o **Python** como linguagem principal para automação.

Ação	Ferramentas Chave	Detalhe de Implementação
<b>Coleta de Dados Governamentais</b>	requests , BeautifulSoup , Pandas , Apache Airflow	Desenvolver scripts Python ( requests , BeautifulSoup ) para navegar nos portais do MTE (CAGED, RAIS) e IBGE (PNAD) e automatizar o download dos arquivos (CSV, XLSX). Utilizar o Apache Airflow para agendar e monitorar a execução mensal/anual desses downloads.
<b>Web Scraping de Vagas</b>	requests , BeautifulSoup , Scrapy , Selenium	<b>Para Prova de Conceito:</b> Utilizar requests e BeautifulSoup para sites mais simples. <b>Para Escala:</b> Implementar Scrapy para projetos de scraping mais complexos e escaláveis. Usar Selenium apenas para lidar com conteúdo dinâmico (JavaScript). <b>Atenção:</b> Implementar delays e rotatividade de User-Agents para evitar bloqueios.
<b>Armazenamento Inicial</b>	Amazon S3 ou Google Cloud Storage (GCS)	Mover todos os dados brutos (CSV, XLSX, HTML) para um <b>Data Lake</b> (S3/GCS) em seu formato original, garantindo durabilidade e estrutura de pastas lógica (ex: /raw/<fonte>/<ano>/<mes>/ ).

## 2. Arquitetura de Big Data e Ferramentas

A arquitetura é baseada em camadas, garantindo escalabilidade e separação de responsabilidades.

Camada	Função Prática	Ferramentas Recomendadas
Ingestão	Orquestrar e automatizar o fluxo de dados das fontes para o Data Lake.	<b>Apache Airflow</b> (Orquestração), <b>Python</b> ( <code>requests</code> , <code>Scrapy</code> ).
Armazenamento	Armazenar dados brutos e processados para consumo analítico.	<b>Data Lake</b> ( <code>S3</code> / <code>GCS</code> ) para <i>Raw Data</i> . <b>Data Warehouse</b> ( <code>BigQuery</code> / <code>Redshift</code> ) para <i>Curated Data</i> (dados limpos e estruturados).
Processamento	Executar o pipeline ETL (Transformação) em grandes volumes de dados.	<b>PySpark</b> (Processamento Distribuído), <b>Pandas</b> (Transformação em memória), <b>Python</b> ( <code>NumPy</code> ).
Análise/ML	Construção de modelos preditivos e análises estatísticas.	<b>Python</b> ( <code>Scikit-learn</code> , <code>TensorFlow</code> ), <b>Jupyter Notebooks</b> (Exploração Interativa).
Visualização	Apresentar insights em dashboards e relatórios.	<b>BI Tools</b> ( <code>Tableau</code> , <code>Power BI</code> ) ou <b>Frameworks Web</b> ( <code>Flask</code> / <code>Django</code> com bibliotecas Python).

### 3. Pipeline de Processamento (ETL)

---

O sucesso do projeto depende da qualidade da etapa de Transformação (T), que deve ser implementada com foco em limpeza, padronização e enriquecimento.

Etapa	Ação Detalhada	Ferramentas Chave
Limpeza de Dados	Tratar valores ausentes, remover duplicatas e corrigir erros de formato (ex: datas, números).	Pandas (funções <code>drop_duplicates()</code> , <code>fillna()</code> ), PySpark para grandes volumes.
Padronização	Normalizar termos e categorias (ex: nomes de cargos, setores, localização).	<b>NLP Básico:</b> Regex e dicionários de mapeamento em Python. <b>Exemplo:</b> Padronizar "Desenvolvedor", "Dev" e "Progr." para " <b>Desenvolvedor</b> ".
Enriquecimento (NLP)	Extrair entidades e habilidades das descrições de vagas.	<b>Bibliotecas NLP:</b> spaCy ou NLTK. <b>Ação:</b> Utilizar <i>Named Entity Recognition (NER)</i> para identificar e classificar habilidades técnicas e comportamentais.
Carga (Load)	Carregar os dados limpos e estruturados no Data Warehouse.	Conectores de PySpark ou Pandas para o Data Warehouse escolhido (ex: <code>psycopg2</code> para PostgreSQL, conectores nativos para BigQuery/Redshift).

## 4. Indicadores Chave de Desempenho (KPIs) e Entregáveis

Os resultados devem ser mensuráveis e focados em insights acionáveis para o mercado de trabalho.

KPI	Propósito	Exemplo de Análise
Tendências de Demanda	Identificar onde o mercado de trabalho está crescendo.	Vagas por setor (CNAE) e região geográfica.
Habilidades em Alta	Mapear a lacuna de competências.	Frequência das <i>hard skills</i> (Python, SQL) e <i>soft skills</i> (Comunicação, Liderança) nas descrições.
Análise Salarial	Fornecer referências de remuneração.	Faixas salariais médias e medianas por cargo, senioridade e localização.
Tempo de Preenchimento (TTF)	Medir a liquidez do mercado.	Tempo médio entre a publicação e a remoção de uma vaga.

## 5. Abordagem de Desenvolvimento

O projeto deve ser conduzido sob a metodologia **Ágil (Agile)**, com foco em entregas rápidas e um **Mínimo Produto Viável (MVP)** que cubra a coleta de dados governamentais e a análise exploratória básica.

### Próximos Passos Essenciais (MVP):

- Configuração do Ambiente:** Instalação das bibliotecas Python ( `requests` , `pandas` , `beautifulsoup` , `scikit-learn` ).
- Desenvolvimento do Script de Aquisição:** Criar o script para download e processamento inicial dos dados do CAGED/RAIS.
- Estruturação do Data Lake:** Definir a estrutura de pastas no S3/GCS para o armazenamento dos dados brutos.
- Análise Exploratória Inicial:** Utilizar `Jupyter Notebooks` para as primeiras visualizações (ex: `Matplotlib` , `Seaborn` ) sobre as tendências de emprego formal.

**Autor:** Manus AI **Data:** Outubro de 2025 **Baseado em:** Plano Estratégico para Análise de Vagas de Emprego no Brasil com Big Data e Python (63 páginas)