

Olá,

Conforme foi pedido, estudei os principais arquivos dos dois repositórios do projeto e, com base na análise e o meu breve conhecimento, elaborei um esboço de plano de trabalho com foco nas melhorias estruturais e em estratégias iniciais de aprendizado de máquina. Acredito que esses três tópicos já representam um bom ponto de partida para aprimorar o processamento dos dados de solo no Repositório SoilData.

---

## 1. Organização e Funções no Código

Não tenho muita experiência na área, mas acredito que, se pegarmos partes do código e transformarmos em funções — em vez de deixá-lo em blocos sequenciais — podemos facilitar muito a leitura, a manutenção e até mesmo reutilizar essas funções em outros contextos.

**Exemplos de tarefas:**

- `verificar_duplicatas()`
- `corrigir_coordenadas()`
- `validar_limites_geograficos()`
- `padronizar_taxonomia()`

Essa abordagem deixará o código mais limpo, mais organizado e muito mais fácil de revisar.

---

## 2. Comentários e Documentação no Código

Além das melhorias estruturais, acredito que se adicionarmos comentários mais explicativos ao longo do código, principalmente em funções e blocos com múltiplas etapas. Nem sempre conseguimos bater o olho e entender diretamente o que determinada função ou trecho faz, então, um comentário explicando melhor ainda o que estamos fazendo pode facilitar bastante tanto o aprendizado quanto a manutenção do código.

---

## 3. Adaptações com Aprendizado de Máquina – Random Forest (Acredito que seria um pouco difícil)

Será que não poderíamos utilizar o Random Forest para prever valores ausentes com base nas correlações entre as variáveis do solo (pH, carbono, etc.), em vez de métodos mais simples como média ou mediana??

O Random Forest é robusto contra outliers e funciona bem mesmo quando há colinearidade entre variáveis (ou seja, quando algumas variáveis estão fortemente correlacionadas entre si, como pH e teor de alumínio). Isso é importante em dados ambientais, onde muitas variáveis do solo estão relacionadas.

O algoritmo **missForest** é baseado em **Random Forest** e realiza a imputação de forma iterativa: ele começa preenchendo os dados faltantes com estimativas simples (como média), e depois vai refinando as previsões em ciclos sucessivos, treinando modelos melhores a cada rodada, até os resultados estabilizarem.”

---

Acredito que com esses três tópicos iniciais já seria um bom começo para se trabalhar. Como disse ao senhor, não tenho tanto conhecimento, mas pelo que estudei nesses últimos três dias, proponho aplicar esses “ajustes”, acredito que os dois primeiros não serão tão difíceis, já o terceiro terá uma complexidade maior.

**Atenciosamente, Jefferson Korte Junior.**