

Otimizando Projetos de Dados: Controle de Versionamento para Grandes Datasets e Código

Olá! Bem-vindo(a) à nossa newsletter semanal sobre Ciência de Dados! Nesta edição, vamos explorar um tema crucial para qualquer analista ou cientista de dados: o controle de versionamento de dados e código. Sabemos que tanto a manipulação quanto a análise de dados exigem ferramentas que nos ajudem a manter tudo organizado e versionado adequadamente. Hoje, vamos entender as limitações do Git em projetos de grande escala e apresentar soluções que podem ser a chave para um fluxo de trabalho mais eficiente.

Embora não haja um limite de tamanho de arquivo formalmente imposto pelo Git, o desempenho pode degradar significativamente com arquivos individuais que excedam alguns gigabytes. No que diz respeito a repositórios no GitHub, existe um limite de tamanho individual de arquivos de 100 MB. Portanto, projetos de dados que exigem a manipulação e armazenamento de grandes volumes de dados frequentemente enfrentam dificuldades, pois, como mencionado, o Git não lida eficientemente com esses arquivos volumosos.

Para resolver esse problema, existem ferramentas como o DVC (Data Version Control) e o Git LFS (Large File Storage).

O DVC é uma ferramenta open-source que estende o Git para gerenciar versões de dados de maneira eficiente. Ele permite que você armazene grandes arquivos de dados de forma externa, mantendo o controle de versão dentro do repositório Git. Dessa forma, você pode manter um histórico completo das mudanças nos datasets sem sobrecarregar o seu repositório Git. O DVC integra-se perfeitamente com outros sistemas de armazenamento, como Amazon S3, Google Drive e Azure Blob Storage, oferecendo flexibilidade e escalabilidade para seus projetos.

Já o Git LFS, que também ajuda a gerenciar arquivos grandes dentro de um repositório Git, substitui os arquivos grandes no seu repositório com ponteiros textuais, mantendo os arquivos grandes armazenados separadamente. O DVC oferece funcionalidades adicionais específicas para projetos de dados, como pipelines de processamento de dados e métricas de performance.

Para mais informações sobre DVC, acesse: <https://dvc.org/> e para mais informações sobre Git LFS, acesse: <https://git-lfs.com/>.

A implementação dessas ferramentas pode transformar a maneira como você gerencia e versiona seus dados e código, permitindo uma colaboração mais eficiente e um controle mais rigoroso sobre as mudanças realizadas ao longo do tempo.

Esperamos que você tenha achado esta abordagem útil. Seja você um iniciante no mundo da Ciência de Dados ou um profissional experiente, entender as nuances entre diferentes ferramentas de versionamento de dados pode abrir novas possibilidades em seus projetos. Se você tiver alguma dúvida ou sugestão, não hesite em entrar em contato. Queremos garantir que nossos conteúdos sejam sempre relevantes e úteis para você.

Fique atento à nossa próxima edição, onde continuaremos a explorar ferramentas, técnicas e dicas valiosas para aprimorar suas habilidades em Ciência de Dados. Até lá, continue explorando e aprendendo! Um abraço,

Prof. Dr. Dilermando Piva Jr
Coordenador de Ciência de Dados para Negócios / Fatec Votorantim
E-mail: f301.cdn@fatec.sp.gov.br