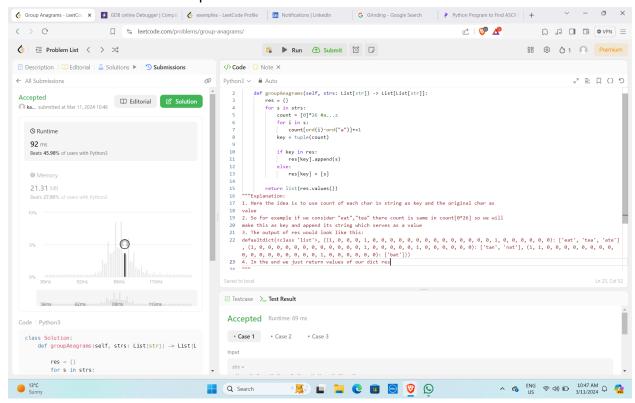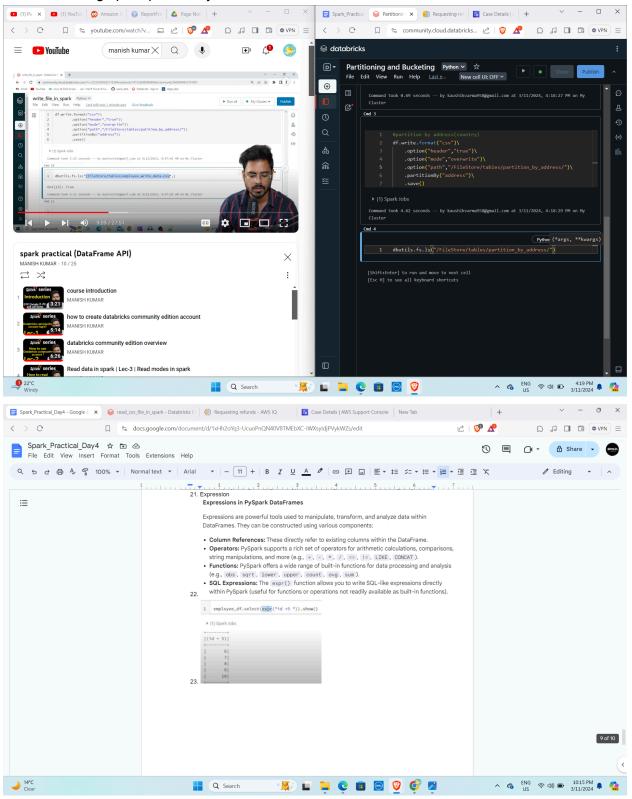Day35 - March 11th 2024

1. Started my day as usual
2. Packed food and headed to library
3. Started marketing my profile for data engineer
4. Solved one medium leetcode problem

# 5. Started learning spark practically

6. Ended my day by solving complex SQL questions from Ankit'YT