

# Hindi में परियोजना दस्तावेज़

## परियोजना अवलोकन: ट्विटर से प्रोफाइल निकालना (Author ID-1 के लिए)

यह परियोजना ट्विटर डेटा से प्रोफाइल जानकारी निकालने के लिए डिज़ाइन की गई है। विशेष रूप से, यह Author ID-1 के लिए डेटा निकालने पर केंद्रित है। मुख्य उद्देश्य ट्विटर उपयोगकर्ताओं के डेटा को एक संरचित प्रारूप में संसाधित करना है, जिसमें उनके देश, उनके संबद्धताओं और उनके सहयोगी (co-authors) के बारे में जानकारी शामिल है। यह परियोजना ट्विटर API का उपयोग करके डेटा प्राप्त करती है और फिर इसे प्रोसेसिंग और विश्लेषण के लिए तैयार करने के लिए pandas लाइब्रेरी का उपयोग करती है।

### मुख्य विशेषताएं:

- \* **डेटा स्रोत:** ट्विटर API
- \* **डेटा प्रारूप:** CSV फ़ाइल से इनपुट
- \* **प्रोसेसिंग चरण:**
  - \* ट्विटर API से डेटा प्राप्त करना
  - \* देश कोड का पता लगाना और भरना
  - \* सहयोगी (co-authors) की पहचान करना और उनके देश कोड प्राप्त करना
  - \* कोऑथर्स के जोड़े बनाना और डेटा को विस्तृत करना
  - \* अनावश्यक कॉलम हटाना
- \* **आउटपुट:** एक्सेल फ़ाइल में संसाधित डेटा

## वास्तुकला और डिज़ाइन

यह परियोजना एक बहु-चरणीय प्रक्रिया का पालन करती है, जिसमें निम्नलिखित प्रमुख घटक शामिल हैं:

1. **डेटा अधिग्रहण:** ट्विटर API का उपयोग करके ट्विटर उपयोगकर्ताओं के डेटा को प्राप्त किया जाता है।
2. **डेटा सफाई और मानकीकरण:** प्राप्त डेटा को साफ़ किया जाता है और एक सुसंगत प्रारूप में परिवर्तित किया जाता है। इसमें अनावश्यक वर्णों को हटाना और डेटा प्रकारों को बदलना शामिल है।
3. **देश कोड पहचान:** ओपनएलेक्स API का उपयोग करके प्रत्येक लेखक के देश कोड को निर्धारित किया जाता है।
4. **सहयोगी विश्लेषण:** यह कोड विभिन्न लेखकों के सहयोग को उजागर करने पर ध्यान केंद्रित करता है, उनके द्वारा लिखे गए शोध पत्रों के आधार पर। यह लेखकों के जोड़े बनाने और उनके देशों को निर्दिष्ट करने के लिए itertools लाइब्रेरी का उपयोग करता है।
5. **डेटा भंडारण:** संसाधित डेटा को एक एक्सेल फ़ाइल में संग्रहीत किया जाता है।

यह परियोजना एक रैखिक वर्कफ़्लो का उपयोग करती है, जिसमें प्रत्येक चरण पिछले चरण के आउटपुट पर निर्भर करता है। मुख्य

लॉजिक को pandas लाइब्रेरी और द्विटर API के साथ इंटरैक्ट करने के लिए कोड में लागू किया गया है।

## मुख्य कार्यक्षमताएं

- \* **द्विटर API से डेटा प्राप्त करना:** कोड द्विटर API का उपयोग करके निर्दिष्ट उपयोगकर्ताओं के बारे में जानकारी प्राप्त करने के लिए किया जाता है।
- \* **देश कोड पहचान और भरना:** कोड ओपनएलेक्स API का उपयोग करके प्रत्येक लेखक के देश कोड को निर्धारित करता है और डेटाफ्रेम में गुम मानों को भरता है।
- \* **सहयोगी (co-authors) की पहचान और विश्लेषण:** कोड लेखकों के जोड़े (pairs) की पहचान करता है और उनके देश कोड को एकत्रित करता है। यह सहयोगी विश्लेषण के लिए एक महत्वपूर्ण कदम है।
- \* **डेटा प्रसंस्करण और सफाई:** कोड डेटा को साफ़ करता है, अनावश्यक कॉलम को हटाता है, और डेटा को अंतिम आउटपुट के लिए तैयार करता है।
- \* **डेटा निर्यात:** समीक्षित डेटा को एक एक्सेल फ़ाइल में निर्यात किया जाता है जिसे आगे विश्लेषण या उपयोग के लिए आसानी से एक्सेस किया जा सकता है।

## वर्कफ़्लो और तर्क

1. **डेटा इनपुट:** द्विटर से डेटा CSV फ़ाइल में लोड किया जाता है।
2. **डेटा फ़िल्टरिंग:** आवश्यक कॉलम का चयन किया जाता है।
3. **देश कोड पहचान:** प्रत्येक लेखक के देश कोड को ओपनएलेक्स API का उपयोग करके प्राप्त किया जाता है।
4. **गुम मानों को भरना:** यदि किसी लेखक के लिए देश कोड उपलब्ध नहीं है, तो "अज्ञात" के रूप में चिह्नित किया जाता है।
5. **सहयोगी (co-authors) की पहचान:** लेखकों के जोड़े (pairs) की पहचान की जाती है।
6. **डेटा विस्तार:** सहयोगी (co-authors) के जोड़े को डेटाफ्रेम में विस्तारित किया जाता है।
7. **अनावश्यक कॉलम हटाना:** अनावश्यक कॉलम हटा दिए जाते हैं।
8. **डेटा निर्यात:** संसाधित डेटा को एक्सेल फ़ाइल में निर्यात किया जाता है।

## मुख्य अवधारणाएँ और तकनीकें

- \* **द्विटर API:** द्विटर उपयोगकर्ताओं के बारे में जानकारी प्राप्त करने के लिए।
- \* **Pandas:** डेटा हेरफेर और विश्लेषण के लिए।
- \* **Requests:** HTTP अनुरोध भेजने के लिए।
- \* **Tqdm:** प्रगति बार प्रदर्शित करने के लिए।
- \* **Itertools:** लेखकों के जोड़े बनाने के लिए।
- \* **OpenAlex API:** लेखकों के देश कोड प्राप्त करने के लिए।
- \* **CSV और Excel फ़ाइलें:** डेटा को संग्रहीत करने के लिए।

## त्रुटि प्रबंधन और प्रदर्शन

- \* **त्रुटि प्रबंधन:** कोड API कॉल के दौरान होने वाली त्रुटियों को संभालता है और उन्हें लॉग करता है। यदि कोई लेखक नहीं मिलता है, तो "अज्ञात" देश कोड का उपयोग किया जाता है।
- \* **प्रगति बार:** tqdm लाइब्रेरी का उपयोग करके प्रगति बार प्रदर्शित किया जाता है, जिससे उपयोगकर्ता को डेटा प्रोसेसिंग की प्रगति का पता चलता है।
- \* **दर सीमा प्रबंधन:** ट्विटर API की दर सीमाओं को ध्यान में रखा जाता है। यदि दर सीमा तक पहुँच जाती है, तो कोड एक संक्षिप्त विराम लेता है और फिर से प्रयास करता है।
- \* **सुरक्षा:** ट्विटर API की शर्तों का पालन किया जाता है।

## संभावित चुनौतियाँ और विचार

- \* **ट्विटर API की दर सीमाएं:** ट्विटर API की दर सीमाएं डेटा प्रोसेसिंग को धीमा कर सकती हैं। दर सीमाओं को संभालने के लिए रणनीति लागू की गई है।
- \* **डेटा की गुणवत्ता:** ट्विटर API से प्राप्त डेटा में त्रुटियाँ या अपूर्णता हो सकती है। कोड डेटा की गुणवत्ता की जांच करता है और त्रुटियों को संभालता है।
- \* **अपडेट किए गए API:** ट्विटर API में बदलाव आने से कोड को अपडेट करने की आवश्यकता हो सकती है।

## भविष्य के सुधार

- \* **डेटा सत्यापन:** डेटा की सटीकता सुनिश्चित करने के लिए डेटा सत्यापन तकनीकों को लागू किया जा सकता है।
- \* **अधिक उन्नत सहयोगी विश्लेषण:** अधिक उन्नत सहयोगी विश्लेषण तकनीकों को लागू किया जा सकता है, जैसे कि सह-लेखन नेटवर्क का निर्माण।
- \* **डेटाबेस में डेटा संग्रहीत करना:** बड़े डेटासेट के लिए, डेटा को डेटाबेस में संग्रहीत करने पर विचार किया जा सकता है।
- \* **विभिन्न प्रकार के ट्विटर डेटा का समर्थन:** कोड को अन्य प्रकार के ट्विटर डेटा, जैसे कि ट्वीट्स और ट्रेंडिंग विषयों का समर्थन करने के लिए बढ़ाया जा सकता है।

## सारांश

यह परियोजना ट्विटर से प्रोफाइल जानकारी निकालने के लिए एक शक्तिशाली और लचीला उपकरण है। यह ट्विटर API और pandas लाइब्रेरी का प्रभावी ढंग से उपयोग करता है ताकि ट्विटर उपयोगकर्ताओं के डेटा को संसाधित किया जा सके और एक संरचित प्रारूप में संग्रहीत किया जा सके। यह परियोजना ट्विटर डेटा विश्लेषण और अनुसंधान के लिए एक मूल्यवान संसाधन है। रखरखाव के लिए, Twitter API में होने वाले परिवर्तनों को ध्यान में रखते हुए नियमित अपडेट की आवश्यकता होगी। इस परियोजना को भविष्य में और अधिक उन्नत डेटा विश्लेषण तकनीकों को शामिल करने के लिए बढ़ाया जा सकता है।