

📁 ILLEGAL DRUG DETECTION DATASETS - HACKATHON PACKAGE

📁 Dataset Files Generated

1. drug_slang_dictionary.csv (200 entries)

Purpose: Comprehensive dictionary of illegal drug slang terms used in online communications

Columns:

- `slang_term`: The slang/street name (e.g., "nuggets", "snow", "molly")
- `drug_type`: Category of drug (marijuana, cocaine, heroin, methamphetamine, mdma, fentanyl, prescription, general)
- `definition`: Explanation of what the slang term means

Example entries:

- "nuggets" → marijuana → "Small marijuana buds"
- "snow" → cocaine → "Cocaine white powder"
- "molly" → mdma → "Pure MDMA"

2. drug_emoji_dictionary.csv (200 entries)

Purpose: Emoji codes used by drug dealers and buyers in messaging apps

Columns:

- `emoji`: The actual emoji symbol (🍁, ❄️, 🍊, 🍬)
- `meaning`: What the emoji represents in drug context
- `drug_category`: Type of drug or dealer activity it represents

Example entries:

- 🍁 → "All drugs (maple leaf)" → universal
- ❄️ → "Snow/cocaine" → cocaine
- 🍊 → "Drug supplier/plug" → dealer
- 🍬 → "Pills/medication" → prescription

3. synthetic_drug_conversations.csv (30,024 messages)

Purpose: Realistic synthetic drug dealing conversations across multiple platforms

Columns:

- `message_id`: Unique message identifier
- `conversation_id`: Groups messages into conversations
- `platform`: Communication platform (telegram, discord, whatsapp, reddit, snapchat)
- `sender_id`: Person sending the message (anonymized usernames)
- `recipient_id`: Person receiving the message
- `message_text`: The actual message content with slang/emojis
- `timestamp`: When message was sent (YYYY-MM-DD HH:MM:SS)
- `message_type`: Type of message (inquiry, availability, pricing, meetup, payment)
- `contains_slang`: 1 if message contains drug slang, 0 otherwise
- `contains_emoji`: 1 if message contains drug emojis, 0 otherwise
- `risk_level`: Assessed risk level (low, medium, high)
- `drug_category`: Primary drug type referenced in message

Example conversation flow:

1. "yo you around?" (inquiry)
2. "yeah I'm here, what you need?" (availability)
3. "looking for some fire weed 🍃" (inquiry with slang + emoji)
4. "got that good herb for \$50 per eighth" (pricing)
5. "usual spot in 20?" (meetup)
6. "cash only" (payment)

▮ Dataset Statistics

Drug Slang Distribution:

- Marijuana: 50 terms (25%)
- General terms: 40 terms (20%)
- Cocaine: 25 terms (12.5%)
- Heroin: 20 terms (10%)
- Prescription: 20 terms (10%)
- Methamphetamine: 15 terms (7.5%)
- Fentanyl: 15 terms (7.5%)
- MDMA: 15 terms (7.5%)

Emoji Categories:

- 20 emojis each for: marijuana, cocaine, heroin, prescription, methamphetamine, mdma, fentanyl, quality indicators, dealer signals
- 10 emojis each for: universal drug symbols, quantity indicators

Synthetic Conversations:

- 30,024 total messages across 5,000 conversations
- Average 6 messages per conversation
- 5 platforms: WhatsApp, Telegram, Discord, Reddit, Snapchat
- 25,526 messages contain slang terms (85%)
- 8,928 messages contain emojis (30%)
- Risk levels: High (33.5%), Medium (33.4%), Low (33.1%)

▮ Usage for Your Hackathon Project

Step 1: Drug Slang Dictionary

```
import pandas as pd
slang_df = pd.read_csv('drug_slang_dictionary.csv')
# Create lookup dictionary for real-time detection
slang_dict = dict(zip(slang_df['slang_term'], slang_df['drug_type']))
```

Step 2: Emoji Detection

```
emoji_df = pd.read_csv('drug_emoji_dictionary.csv')
# Create emoji to drug category mapping
emoji_dict = dict(zip(emoji_df['emoji'], emoji_df['drug_category']))
```

Step 3: Training Data

```
conversations_df = pd.read_csv('synthetic_drug_conversations.csv')
# Filter high-risk conversations for training
high_risk = conversations_df[conversations_df['risk_level'] == 'high']
# Train NLP model on message_text with risk_level as labels
```

Step 4: Platform Analysis

```
# Analyze patterns by platform
platform_stats = conversations_df.groupby('platform').agg({
    'contains_slang': 'mean',
    'contains_emoji': 'mean',
    'risk_level': lambda x: (x == 'high').mean()
})
```

Step 5: Network Graph Construction

```
import networkx as nx
# Build network from sender/recipient relationships
G = nx.from_pandas_edgelist(conversations_df,
                           source='sender_id',
                           target='recipient_id',
                           edge_attr=['timestamp', 'risk_level'])
```

⚠ Legal & Ethical Notice

These datasets are created for:

- ✓ Academic research and education
- ✓ Law enforcement training and detection systems
- ✓ Content moderation and safety tools
- ✓ Hackathon projects for public safety

NOT to be used for:

- ✗ Actual illegal drug transactions
- ✗ Facilitating illegal activities
- ✗ Harassment or targeting individuals

All data is synthetic and does not represent real individuals or actual illegal activities.

▢ Expected Model Performance

Based on dataset characteristics, your models should achieve:

- **Slang Detection:** 85-90% accuracy (comprehensive dictionary)
- **Emoji Detection:** 90-95% accuracy (clear emoji mappings)
- **Risk Classification:** 70-80% accuracy (realistic conversation patterns)
- **Platform Analysis:** 80-85% accuracy (distinct platform behaviors)

▢ Quick Start Commands

```
# Load and explore datasets
import pandas as pd

# Load all datasets
slang_df = pd.read_csv('drug_slang_dictionary.csv')
emoji_df = pd.read_csv('drug_emoji_dictionary.csv')
conversations_df = pd.read_csv('synthetic_drug_conversations.csv')

# Basic exploration
print(f"Slang terms: {len(slang_df)}")
print(f"Emojis: {len(emoji_df)}")
```

```
print(f"Messages: {len(conversations_df)}")  
print(f"Conversations: {conversations_df['conversation_id'].nunique()}")
```

Perfect for your 2-day hackathon timeline! 📅