# Cybersecurity Data Collection & OSINT Guide

## 🛡️ Legitimate Cybersecurity Use Cases

### Threat Intelligence

- IOC (Indicators of Compromise) collection
- Malware sample gathering from public sources
- Dark web monitoring for leaked credentials
- Vulnerability disclosure tracking
- Threat actor profiling and attribution

### Security Research

- Attack surface enumeration
- Public exposure identification
- Security misconfiguration detection
- Phishing campaign analysis
- Brand protection monitoring

### Penetration Testing

- OSINT gathering for authorized assessments
- Social engineering research (with proper scope)
- Public information reconnaissance
- Technology stack identification
- Employee information for spear phishing tests

## 🔍 Professional OSINT Tools

### Automated Reconnaissance

```
python
```

```python
# TheHarvester - Email/subdomain enumeration
# Usage: theHarvester -d target.com -l 500 -b all

# Amass - Asset discovery
# Usage: amass enum -d target.com

# Subfinder - Subdomain discovery
# Usage: subfinder -d target.com

# Example Python integration
import subprocess
import json

def run_amass(domain):
    """Run Amass for subdomain enumeration"""
    cmd = f"amass enum -d {domain} -json"
    result = subprocess.run(cmd.split(), capture_output=True, text=True)
    return [json.loads(line) for line in result.stdout.strip().split('\n') if line]

def run_theharvester(domain):
    """Run TheHarvester for email collection"""
    cmd = f"theHarvester -d {domain} -l 100 -b all -f /tmp/harvest_results"
    subprocess.run(cmd.split(), capture_output=True)
    # Parse results from file
```

## Web Application Security

```python
```

```python
import requests
from bs4 import BeautifulSoup
import re
from urllib.parse import urljoin, urlparse


class SecurityScanner:
    def __init__(self, target_url, user_agent="SecResearch/1.0"):
        self.target = target_url
        self.session = requests.Session()
        self.session.headers.update({'User-Agent': user_agent})

    def check_robots_txt(self):
        """Analyze robots.txt for sensitive paths"""
        try:
            robots_url = urljoin(self.target, '/robots.txt')
            response = self.session.get(robots_url)
            if response.status_code == 200:
                return self.parse_robots(response.text)
        except Exception as e:
            print(f"Error checking robots.txt: {e}")
        return []

    def parse_robots(self, robots_content):
        """Extract interesting paths from robots.txt"""
        disallowed_paths = []
        for line in robots_content.split('\n'):
            if line.startswith('Disallow:'):
                path = line.split(':', 1)[1].strip()
                if path and path != '/':
                    disallowed_paths.append(path)
        return disallowed_paths

    def extract_emails_from_page(self, url):
        """Extract email addresses for contact mapping"""
        try:
            response = self.session.get(url)
            if response.status_code == 200:
                email_pattern = r'\b[A-Za-z0-9._%+-]+@[A-Za-z0-9.-]+\.[A-Z|a-z]{2,}\b'
                emails = re.findall(email_pattern, response.text)
                return list(set(emails))  # Remove duplicates
        except Exception as e:
            print(f"Error extracting emails from {url}: {e}")
        return []
```

```python
    def find_contact_forms(self, url):
        """Identify contact forms for social engineering assessment"""
        try:
            response = self.session.get(url)
            if response.status_code == 200:
                soup = BeautifulSoup(response.content, 'html.parser')
                forms = soup.find_all('form')

                contact_forms = []
                for form in forms:
                    # Look for contact-related forms
                    form_text = str(form).lower()
                    if any(keyword in form_text for keyword in
                        ['contact', 'email', 'message', 'inquiry', 'support']):
                        contact_forms.append({
                            'action': form.get('action', ''),
                            'method': form.get('method', 'get'),
                            'inputs': [inp.get('name') for inp in form.find_all('input')]
                        })

                return contact_forms
        except Exception as e:
            print(f"Error finding contact forms: {e}")
        return []

# Usage example
scanner = SecurityScanner("https://target-website.com")
emails = scanner.extract_emails_from_page("https://target-website.com/contact")
robots_paths = scanner.check_robots_txt()
```

## 📊 Threat Intelligence Collection

### IOC Aggregation

```
python
```

```python
import requests
import json
from datetime import datetime

class ThreatIntelCollector:
    def __init__(self):
        self.sources = {
            'abuse_ch': 'https://urlhaus-api.abuse.ch/v1/',
            'alienvault': 'https://otx.alienvault.com/api/v1/',
            # Add more threat intel sources
        }

    def collect_malware_urls(self):
        """Collect malware URLs from URLhaus"""
        try:
            response = requests.post(
                f"{self.sources['abuse_ch']}urls/recent/",
                data={'limit': 100}
            )
            if response.status_code == 200:
                data = response.json()
                return [
                    {
                        'url': item['url'],
                        'threat': item['threat'],
                        'date_added': item['date_added'],
                        'reporter': item['reporter']
                    }
                    for item in data.get('urls', [])
                ]
        except Exception as e:
            print(f"Error collecting malware URLs: {e}")
        return []

    def collect_ip_reputation(self, ip_address):
        """Check IP reputation across multiple sources"""
        results = {}

        # AbuseIPDB (requires API key)
        # VirusTotal (requires API key)
        # Add your preferred threat intel APIs

        return results
```

```python
# Usage
collector = ThreatIntelCollector()
recent_threats = collector.collect_malware_urls()
```

## Dark Web Monitoring

```python
```

```python
import requests
import time
from stem import Signal
from stem.control import Controller

class DarkWebMonitor:
    """
    NOTE: Only use for legitimate security research
    Requires Tor setup and proper legal authorization
    """

    def __init__(self):
        self.session = requests.Session()
        self.session.proxies = {
            'http': 'socks5://127.0.0.1:9050',
            'https': 'socks5://127.0.0.1:9050'
        }

    def rotate_tor_identity(self):
        """Rotate Tor circuit for anonymity"""
        try:
            with Controller.from_port(port=9051) as controller:
                controller.authenticate(password="your_tor_password")
                controller.signal(Signal.NEWNYM)
                time.sleep(10)  # Wait for new circuit
        except Exception as e:
            print(f"Error rotating Tor identity: {e}")

    def search_paste_sites(self, keywords):
        """Monitor paste sites for leaked credentials"""
        # Implementation depends on specific paste site APIs
        # Always respect rate limits and terms of service
        pass

# Usage requires proper Tor setup and legal authorization
```

## 🔧 Advanced Scraping Techniques

### Bypassing Basic Anti-Bot Measures

```
python
```

```python
import requests
from selenium import webdriver
from selenium.webdriver.chrome.options import Options
from selenium.webdriver.common.by import By
import undetected_chromedriver as uc
import time
import random

class AdvancedScraper:
    def __init__(self):
        self.setup_session()
        self.setup_selenium()

    def setup_session(self):
        """Setup requests session with realistic headers"""
        self.session = requests.Session()

        # Rotate User-Agents
        user_agents = [
            'Mozilla/5.0 (Windows NT 10.0; Win64; x64) AppleWebKit/537.36',
            'Mozilla/5.0 (Macintosh; Intel Mac OS X 10_15_7) AppleWebKit/537.36',
            'Mozilla/5.0 (X11; Linux x86_64) AppleWebKit/537.36'
        ]

        self.session.headers.update({
            'User-Agent': random.choice(user_agents),
            'Accept': 'text/html,application/xhtml+xml,application/xml;q=0.9,*/*;q=0.8',
            'Accept-Language': 'en-US,en;q=0.5',
            'Accept-Encoding': 'gzip, deflate',
            'Connection': 'keep-alive',
            'Upgrade-Insecure-Requests': '1'
        })

    def setup_selenium(self):
        """Setup undetected Chrome driver"""
        options = uc.ChromeOptions()
        options.add_argument('--no-sandbox')
        options.add_argument('--disable-dev-shm-usage')
        options.add_argument('--disable-blink-features=AutomationControlled')
        self.driver = uc.Chrome(options=options)

    def human_like_delay(self, min_delay=1, max_delay=3):
        """Add human-like delays"""
```

```python
        delay = random.uniform(min_delay, max_delay)
        time.sleep(delay)

    def scrape_with_js(self, url):
        """Scrape JavaScript-heavy sites"""
        try:
            self.driver.get(url)
            self.human_like_delay(2, 5)

            # Wait for content to load
            time.sleep(3)

            # Extract data
            page_source = self.driver.page_source
            return page_source

        except Exception as e:
            print(f"Error scraping {url}: {e}")
        return None

    def __del__(self):
        if hasattr(self, 'driver'):
            self.driver.quit()

# Usage
scraper = AdvancedScraper()
content = scraper.scrape_with_js("https://target-site.com")
```

## Data Validation & Enrichment

```
python
```

```python
import re
import dns.resolver
import socket
from email_validator import validate_email, EmailNotValidError

class DataValidator:
    def __init__(self):
        self.email_pattern = r'\b[A-Za-z0-9._%+-]+@[A-Za-z0-9.-]+\.[A-Z|a-z]{2,}\b'
        self.phone_patterns = {
            'us': r'\b\d{3}[-.\s]?\d{3}[-.\s]?\d{4}\b',
            'uk': r'\b\d{4}\s?\d{3}\s?\d{3}\b',
            'intl': r'\+\d{1,3}[-.\s]?\d{1,14}\b'
        }

    def validate_email_advanced(self, email):
        """Advanced email validation with MX record check"""
        try:
            # Basic format validation
            if not re.match(self.email_pattern, email):
                return False, "Invalid format"

            # Library validation
            valid = validate_email(email)
            email = valid.email

            # MX record check
            domain = email.split('@')[1]
            try:
                mx_records = dns.resolver.resolve(domain, 'MX')
                if not mx_records:
                    return False, "No MX record"
            except:
                return False, "DNS resolution failed"

            return True, "Valid"

        except EmailNotValidError as e:
            return False, str(e)

    def extract_and_validate_contacts(self, text):
        """Extract and validate all contacts from text"""
        results = {
            'emails': [],
```

```python
        'phones': [],
        'validated_emails': [],
        'invalid_emails': []
    }

    # Extract emails
    emails = re.findall(self.email_pattern, text)
    results['emails'] = list(set(emails))

    # Validate emails
    for email in results['emails']:
        is_valid, reason = self.validate_email_advanced(email)
        if is_valid:
            results['validated_emails'].append(email)
        else:
            results['invalid_emails'].append({'email': email, 'reason': reason})

    # Extract phone numbers
    for region, pattern in self.phone_patterns.items():
        phones = re.findall(pattern, text)
        results['phones'].extend([(phone, region) for phone in phones])

    return results

# Usage
validator = DataValidator()
contacts = validator.extract_and_validate_contacts(scraped_content)
```

## 🛠️ Professional Tools & Frameworks

### Scrapy for Large-Scale Operations

```
python
```

```python
import scrapy
from scrapy.crawler import CrawlerProcess
import json

class CyberSecSpider(scrapy.Spider):
    name = 'cybersec_spider'

    custom_settings = {
        'ROBOTSTXT_OBEY': True,
        'DOWNLOAD_DELAY': 2,
        'RANDOMIZE_DOWNLOAD_DELAY': 0.5,
        'CONCURRENT_REQUESTS': 1,
        'AUTOTHROTTLE_ENABLED': True,
        'AUTOTHROTTLE_START_DELAY': 1,
        'AUTOTHROTTLE_MAX_DELAY': 10,
        'USER_AGENT': 'CyberSec Research Bot 1.0'
    }

    def start_requests(self):
        urls = [
            # Your target URLs for security research
        ]
        for url in urls:
            yield scrapy.Request(url=url, callback=self.parse)

    def parse(self, response):
        # Extract security-relevant information
        emails = response.css('a[href^="mailto:"]::attr(href)').getall()
        phone_links = response.css('a[href^="tel:"]::attr(href)').getall()

        # Extract from text content
        text_content = ' '.join(response.css('::text').getall())

        yield {
            'url': response.url,
            'emails': [email.replace('mailto:', '') for email in emails],
            'phones': [phone.replace('tel:', '') for phone in phone_links],
            'text_emails': self.extract_emails_from_text(text_content),
            'text_phones': self.extract_phones_from_text(text_content),
            'timestamp': scrapy.utils.misc.utc_iso_date()
        }

    def extract_emails_from_text(self, text):
```

```python
        import re
        pattern = r'\b[A-Za-z0-9._%+-]+@[A-Za-z0-9.-]+\.[A-Z|a-z]{2,}\b'
        return re.findall(pattern, text)

    def extract_phones_from_text(self, text):
        import re
        pattern = r'\b\d{3}[-.\s]?\d{3}[-.\s]?\d{4}\b'
        return re.findall(pattern, text)


# Run the spider
# process = CrawlerProcess()
# process.crawl(CyberSecSpider)
# process.start()
```

## 📋 Cybersecurity Ethics & Compliance

### Legal Framework

- **Authorization**: Always have explicit permission for penetration testing

- **Scope Limitation**: Stay within defined assessment boundaries

- **Data Handling**: Follow secure data handling practices

- **Disclosure**: Responsible disclosure for any vulnerabilities found

- **Documentation**: Maintain detailed logs for audit purposes

### Professional Standards

- **OWASP Guidelines**: Follow OWASP testing methodology

- **NIST Framework**: Align with cybersecurity frameworks

- **Industry Standards**: Comply with relevant security standards

- **Client Confidentiality**: Protect client data and findings

### Red Team Considerations

```python

```

```python
class RedTeamOSINT:
    """
    Red Team OSINT collection with proper boundaries
    """

    def __init__(self, scope_domains, authorized_by):
        self.scope = scope_domains
        self.authorization = authorized_by
        self.findings = []

    def validate_scope(self, target):
        """Ensure target is within authorized scope"""
        return any(domain in target for domain in self.scope)

    def collect_employee_info(self, company_name):
        """Collect employee information for social engineering assessment"""
        if not self.authorization:
            raise Exception("Unauthorized collection attempt")

        # LinkedIn, company websites, etc.
        # Only collect publicly available information
        # Document source and method for each finding
        pass

    def generate_report(self):
        """Generate professional assessment report"""
        return {
            'scope': self.scope,
            'authorization': self.authorization,
            'findings': self.findings,
            'methodology': 'OSINT collection using automated tools',
            'recommendations': []
        }
```

## 🎯 Key Takeaways for Cybersecurity Professionals

1. **Always work within authorized scope** - Document everything

2. **Use professional tools** - Scrapy, TheHarvester, Amass, etc.

3. **Validate data quality** - Implement proper validation logic

4. **Respect rate limits** - Don't overwhelm target systems

5. **Follow responsible disclosure** - Report vulnerabilities properly

6. **Maintain operational security** - Protect your tools and methods

7. **Document methodology** - For reporting and legal compliance

## 🔗 Additional Resources

### Professional Tools

- **Maltego**: Visual link analysis

- **Shodan**: Internet-connected device search

- **Censys**: Internet scanning and analysis

- **SpiderFoot**: Automated OSINT collection

- **Recon-ng**: Web reconnaissance framework

### Training & Certifications

- **OSCP**: Offensive Security Certified Professional

- **GCIH**: GIAC Certified Incident Handler

- **CEH**: Certified Ethical Hacker

- **CISSP**: Certified Information Systems Security Professional

Remember: With great power comes great responsibility. Always ensure your activities are authorized, ethical, and contribute positively to cybersecurity!