

# Investment Data Standards Organization Best Practices

## **IDSO BEST PRACTICES** **Web Crawling**

Document # : IDSO-WC-BP-001

Version : 1.0

Effective Date : Draft

### **APPENDIX 2: CHECKLIST FOR ASSESSING RISK & COMPLIANCE OF WEB HARVESTING.**

#### **WEB HARVESTING REGULATIONS & COMPLIANCE**

#### **SECTION**

- |          |  |
|----------|--|
| <b>1</b> | All datasets should be obtained by legal means.  |
| <b>2</b> | Alternative Data purchased or owned by the organization should be assessed for web harvesting compliance.              |
| <b>3</b> | There should be a documented intent of use for each website that will be Web Harvested.                                |
| <b>4</b> | The organization should create a compliance working group to manage web harvesting within the organization.            |
| <b>5</b> | The compliance working group should insure that the compliance strategy conforms to the appropriate statutes and laws. |

#### **WEBSITE ASSESSMENT**

#### **SECTION**

- |           |  |
|-----------|--|
| <b>6</b>  | A data collector should assess a website according to the terms of its robots.txt.   |
| <b>7</b>  | A data collector should access websites in a way that the access does not interfere with or impose an undue burden on their operation.   |
| <b>8</b>  | A data collector should not access, download or transmit non-public website data.  |
| <b>9</b>  | A data collector should not circumvent logins or other access control restrictions such as captcha.  |
| <b>10</b> | A data collector should not utilize IP masking or rotation to avoid website restrictions.  |
| <b>11</b> | A data collector should respect valid cease and desist notices and the website's right to govern the terms of access to the website and data.                                      |
| <b>12</b> | A data collector should respect all copyright and trademark ownership and not act so as to obscure or delete copyright management information.                                     |
| <b>13</b> | A data collector should respect all contact terms with data provider formed by actual, human observed notice, and explicit agreement.  |
| <b>14</b> | A data collector should not act so as to significantly impair or compete with the business of the website.   |
| <b>15</b> | A data collector should make the best commercial efforts to preserve an image of the website, then current notification of terms and conditions, and robots.txt file with capture. |
| <b>16</b> | A data collector should not generate data in violation of the CAN-SPAM act.  |
| <b>17</b> | A data collector should seek to use a website's official API when possible.  |

*Confidential and Proprietary Information – Investment Data Standards Organization Best Practices  
Not to be Disclosed or Reproduced Without Prior Written Approval*

# Investment Data Standards Organization Best Practices

## **IDSO BEST PRACTICES** **Web Crawling**

Document # : IDSO-WC-BP-001

Version : 1.0

Effective Date : Draft

- |           |   |
|-----------|---|
| <b>18</b> | A data collector should not engage in any criminal or intentionally deceptive conduct or seek to exceed their access authorization. |
|-----------|---|

### **RISK MANAGEMENT**

### **SECTION**

- |           |  |
|-----------|--|
| <b>19</b> | The organization should have a risk assessment process for identifying the impact level of web harvesting. |
| <b>20</b> | A risk assessment should be performed for each website to identify web harvesting impact level.            |

### **WEB HARVESTING POLICIES & MANAGEMENT**

### **SECTION**

- |           |  |
|-----------|--|
| <b>21</b> | The organization should create at least one written procedure for web harvesting management.   |
| <b>22</b> | The definition of web harvesting should be defined in the procedure.   |
| <b>23</b> | A synopsis of all relevant privacy laws, regulations, and policies should be included in the procedure.  |
| <b>24</b> | The departments and individual roles and responsibilities for web harvesting should be defined.  |
| <b>25</b> | The roles, responsibilities, and response for web crawling incidents should be defined.  |
| <b>26</b> | The websites that will be crawled should be assessed to mitigate risk.   |
| <b>27</b> | Guidance should be provided on restrictions of data collection, disclosure, sharing, storage and use of web harvesting within the organization.                  |
| <b>28</b> | The procedure should specify an interval for review of Web Harvested websites and frequencies to determine whether the web harvesting is relevant and necessary. |
| <b>29</b> | The procedure should include instructions for handling a security or privacy breach as a result of web harvesting.   |

### **WEB HARVESTING REVIEW AT REGULAR INTERVALS AND AT LEAST ANNUALLY**

### **SECTION**

- |           |  |
|-----------|--|
| <b>30</b> | A review of Web Harvested websites and crawling frequencies should be conducted no less than once per year to determine whether the web harvesting of the particular website(s) is relevant and necessary. |
| <b>31</b> | The organization should update their Risk Assessment document at least annually.   |
| <b>32</b> | Organizations should regularly review audit records of inappropriate or unusual activity surrounding web harvesting processes.   |

*Confidential and Proprietary Information – Investment Data Standards Organization Best Practices  
Not to be Disclosed or Reproduced Without Prior Written Approval*

# Investment Data Standards Organization Best Practices

## **IDSO BEST PRACTICES** **Web Crawling**

Document # : IDSO-WC-BP-001

Version : 1.0

Effective Date : Draft

---

- 33** | The web harvesting procedure should be updated at least annually.

### **WEB HARVESTING INCIDENTS**

### **SECTION**

- 34** | A response plan to handle web harvesting incidents should be detailed in the web harvesting procedure.
- 35** | Information obtained through detection, analysis, containment, and recovery should be collected to help protect against future incidents.

### **TRAINING**

### **SECTION**

- 36** | The organization should provide training on web harvesting control policies including responsibilities, implementation, and access controls.
- 37** | The organization should provide incident response training for web harvesting.