

Unstructured Data Analysis

HAZIM FITRI

2025-04-13

Contents

Part 1: Extract Text from Files 1

Part 2: Web Scrapping 4

DIKW Pyramid (
hi

Part 1: Extract Text from Files

```
eg1<-read.table("Attachment 1/GC.txt",fill=T,header=F) #Data GC.txt
eg1[1,]
```

```
##           V1  V2  V3    V4   V5 V6           V7 V8 V9 V10 V11 V12 V13 V14 V15 V16
## 1 Gertrude Cox, The First Lady Of Statistics
##   V17 V18 V19 V20 V21 V22 V23 V24 V25
## 1
```

```
eg2<-read.csv("Attachment 1/GC.csv",header=F) #Data GC.csv
eg2[1,]
```

```
## [1] "Gertrude Cox, The First Lady Of Statistics"
```

```
#Using tm package
library(tm)# text mining
```

```
## Warning: package 'tm' was built under R version 4.4.3
```

```
## Loading required package: NLP
```

```
## Warning: package 'NLP' was built under R version 4.4.2
```

```
eg3<-c("Hi!", "Welcome to STQD6114", "Tuesday, 11-1pm")
mytext<-VectorSource(eg3)
mycorpus<-VCorpus(mytext)
inspect(mycorpus)
```

```
## <<VCorpus>>
## Metadata: corpus specific: 0, document level (indexed): 0
## Content: documents: 3
##
## [[1]]
## <<PlainTextDocument>>
```

```
## Metadata: 7
## Content: chars: 3
##
## [[2]]
## <<PlainTextDocument>>
## Metadata: 7
## Content: chars: 19
##
## [[3]]
## <<PlainTextDocument>>
## Metadata: 7
## Content: chars: 15

as.character(mycorpus[[1]])
```

```
## [1] "Hi!"
```

```
#Example using VectorSource
eg4<-t(eg1) #From example 1
a<-sapply(1:7,function(x) # 1:7 can be changed into "ncol(eg4)"
  trimws(paste(eg4[,x],collapse=" "),"right"))
mytext<-VectorSource(a)
mycorpus<-VCorpus(mytext)
inspect(mycorpus)
```

```
## <<VCorpus>>
## Metadata: corpus specific: 0, document level (indexed): 0
## Content: documents: 7
##
## [[1]]
## <<PlainTextDocument>>
## Metadata: 7
## Content: chars: 42
##
## [[2]]
## <<PlainTextDocument>>
## Metadata: 7
## Content: chars: 52
##
## [[3]]
## <<PlainTextDocument>>
## Metadata: 7
## Content: chars: 118
##
## [[4]]
## <<PlainTextDocument>>
## Metadata: 7
## Content: chars: 153
##
## [[5]]
## <<PlainTextDocument>>
## Metadata: 7
## Content: chars: 100
##
## [[6]]
## <<PlainTextDocument>>
## Metadata: 7
## Content: chars: 119
##
## [[7]]
## <<PlainTextDocument>>
## Metadata: 7
## Content: chars: 147
```

```
as.character(mycorpus[[1]])
```

```
## [1] "Gertrude Cox, The First Lady Of Statistics"
```

```
#Example using DirSource  
mytext<-DirSource("movies")  
mycorpus<-VCorpus(mytext)  
inspect(mycorpus)
```

```
## <<VCorpus>>  
## Metadata: corpus specific: 0, document level (indexed): 0  
## Content: documents: 5  
##  
## [[1]]  
## <<PlainTextDocument>>  
## Metadata: 7  
## Content: chars: 502  
##  
## [[2]]  
## <<PlainTextDocument>>  
## Metadata: 7  
## Content: chars: 400  
##  
## [[3]]  
## <<PlainTextDocument>>  
## Metadata: 7  
## Content: chars: 506  
##  
## [[4]]  
## <<PlainTextDocument>>  
## Metadata: 7  
## Content: chars: 387  
##  
## [[5]]  
## <<PlainTextDocument>>  
## Metadata: 7  
## Content: chars: 385
```

```
as.character(mycorpus[[1]])
```

```
## [1] "Spider-Man: Homecoming"  
## [2] ""  
## [3] "Thrilled by his experience with the Avengers, young Peter Parker returns home to live with his Aunt May. U
```

```
#Example using DataFrameSource  
eg5<-read.csv("Attachment 1/Doc6.csv",header=F) #Using doc6.csv  
docs<-data.frame(doc_id=c("doc_1","doc_2"),  
                 text=c(as.character(eg5[1,]),  
                       as.character(eg5[2,])),  
                 dmeta1=1:2,dmeta2=letters[1:2],  
                 stringsAsFactors=F)  
mytext<-DataframeSource(docs)  
mycorpus<-VCorpus(mytext)  
inspect(mycorpus)
```

```
## <<VCorpus>>  
## Metadata: corpus specific: 0, document level (indexed): 2  
## Content: documents: 2  
##  
## [[1]]
```

```
## <<PlainTextDocument>>
## Metadata: 7
## Content: chars: 15
##
## [[2]]
## <<PlainTextDocument>>
## Metadata: 7
## Content: chars: 413
```

```
as.character(mycorpus[[1]])
```

```
## [1] "Despicable Me 3"
```

```
eg5b = read.csv("Attachment 1/GC.csv", header=F)
```

Part 2: Web Scrapping

HTML Cheat Sheet

```
eg6<-readLines("https://en.wikipedia.org/wiki/Data_science")
```

```
## Warning in readLines("https://en.wikipedia.org/wiki/Data_science"): incomplete
## final line found on 'https://en.wikipedia.org/wiki/Data_science'
```

```
eg6[grep("\\h2",eg6)]
```

```
## [1] "\t<h2 class=\"vector-pinnable-header-label\">Contents</h2>"
## [2] "<div class=\"mw-heading mw-heading2\"><h2 id=\"Foundations\">Foundations</h2><span class=\"mw-editsection\"><sp
## [3] "<div class=\"mw-heading mw-heading2\"><h2 id=\"Etymology\">Etymology</h2><span class=\"mw-editsection\"><sp
## [4] "<div class=\"mw-heading mw-heading2\"><h2 id=\"Data_science_and_data_analysis\">Data science and data anal
## [5] "<div class=\"mw-heading mw-heading2\"><h2 id=\"Cloud_computing_for_data_science\">Cloud computing for data
## [6] "<div class=\"mw-heading mw-heading2\"><h2 id=\"Ethical_consideration_in_data_science\">Ethical considerati
## [7] "<div class=\"mw-heading mw-heading2\"><h2 id=\"See_also\">See also</h2><span class=\"mw-editsection\"><span
## [8] "<div class=\"mw-heading mw-heading2\"><h2 id=\"References\">References</h2><span class=\"mw-editsection\"><span"
```

```
eg6[grep("\\p",eg6)] #paragraph
```

```
## [1] "<html class=\"client-nojs vector-feature-language-in-header-enabled vector-feature-language-in-main-page
## [2] "<title>Data science - Wikipedia</title>"
## [3] "<script>(function(){var className=\"client-js vector-feature-language-in-header-enabled vector-feature-l
## [4] 'RLSTATE={\"ext.globalCssJs.user.styles\": \"ready\", \"site.styles\": \"ready\", \"user.styles\": \"ready\", \"
## [5] "<script>(RLQ=window.RLQ|| []).push(function(){mw.loader.impl(function(){return[\"user.options@12s5i\",fun
## [6] "});});});</script>"
## [7] "<link rel=\"stylesheet\" href=\"/w/load.php?lang=en&modules=ext.cite.styles%7Cext.uls.interlanguage%
## [8] "<script async=\"\" src=\"/w/load.php?lang=en&modules=startup&only=scripts&raw=1&skin=vec
## [9] "<link rel=\"stylesheet\" href=\"/w/load.php?lang=en&modules=site.styles&only=styles&skin=vec
## [10] "<meta name=\"robots\" content=\"max-image-preview:standard\">"
## [11] "<meta name=\"format-detection\" content=\"telephone=no\">"
## [12] "<meta property=\"og:image\" content=\"https://upload.wikimedia.org/wikipedia/commons/thumb/4/45/PIA23792
## [13] "<meta property=\"og:image:width\" content=\"1200\">"
## [14] "<meta property=\"og:image:height\" content=\"900\">"
## [15] "<meta property=\"og:image\" content=\"https://upload.wikimedia.org/wikipedia/commons/thumb/4/45/PIA23792
## [16] "<meta property=\"og:image:width\" content=\"800\">"
## [17] "<meta property=\"og:image:height\" content=\"600\">"
## [18] "<meta property=\"og:image:width\" content=\"640\">"
## [19] "<meta property=\"og:image:height\" content=\"480\">"
## [20] "<meta name=\"viewport\" content=\"width=1120\">"
## [21] "<meta property=\"og:title\" content=\"Data science - Wikipedia\">"
## [22] "<meta property=\"og:type\" content=\"website\">"
```

```
## [23] "<link rel=\"preconnect\" href=\"/upload.wikimedia.org\">"
## [24] "<link rel=\"alternate\" media=\"only screen and (max-width: 640px)\" href=\"/en.m.wikipedia.org/wiki/Data_science\">"
## [25] "<link rel=\"alternate\" type=\"application/x-wiki\" title=\"Edit this page\" href=\"/w/index.php?title=Data_science\">"
## [26] "<link rel=\"apple-touch-icon\" href=\"/static/apple-touch/wikipedia.png\">"
## [27] "<link rel=\"icon\" href=\"/static/favicon/wikipedia.ico\">"
## [28] "<link rel=\"search\" type=\"application/opensearchdescription+xml\" href=\"/w/rest.php/v1/search\" title=\"Search Wikipedia\">"
## [29] "<link rel=\"EditURI\" type=\"application/rsd+xml\" href=\"/en.wikipedia.org/w/api.php?action=rsd\">"
## [30] "<link rel=\"canonical\" href=\"https://en.wikipedia.org/wiki/Data_science\">"
## [31] "<link rel=\"license\" href=\"https://creativecommons.org/licenses/by-sa/4.0/deed.en\">"
## [32] "<link rel=\"alternate\" type=\"application/atom+xml\" title=\"Wikipedia Atom feed\" href=\"/w/index.php?title=Wikipedia:Atom_feed\">"
## [33] "<link rel=\"dns-prefetch\" href=\"/meta.wikimedia.org\" />"
## [34] "<link rel=\"dns-prefetch\" href=\"auth.wikimedia.org\">"
## [35] "<body class=\"skin--responsive skin-vector skin-vector-search-vue mediawiki ltr sitedir-ltr mw-hide-empty-ns ns\">"
## [36] "<div id=\"vector-main-menu-dropdown\" class=\"vector-dropdown vector-main-menu-dropdown vector-button-flip\">"
## [37] "<input type=\"checkbox\" id=\"vector-main-menu-dropdown-checkbox\" role=\"button\" aria-haspopup=\"true\" data-event-name=\"ui:dropdown\" />"
## [38] "<label id=\"vector-main-menu-dropdown-label\" for=\"vector-main-menu-dropdown-checkbox\" class=\"vector-dropdown-label\">"
## [39] "<span class=\"vector-dropdown-label-text\">Main menu</span>"
## [40] "<div class=\"vector-dropdown-content\">"
## [41] "<div id=\"vector-main-menu-unpinned-container\" class=\"vector-unpinned-container\">"
## [42] "<div id=\"vector-main-menu\" class=\"vector-main-menu vector-pinnable-element\">"
## [43] "<div class=\"vector-pinnable-header vector-main-menu-pinnable-header vector-pinnable-header-unpinned\">"
## [44] "<div data-feature-name=\"main-menu-pinned\">"
## [45] "<div data-pinnable-element-id=\"vector-main-menu\">"
## [46] "<div data-pinned-container-id=\"vector-main-menu-pinned-container\">"
## [47] "<div data-unpinned-container-id=\"vector-main-menu-unpinned-container\">"
## [48] "<div class=\"vector-pinnable-header-label\">Main menu</div>"
## [49] "<button class=\"vector-pinnable-header-toggle-button vector-pinnable-header-pin-button\" data-event-name=\"ui:toggle\">"
## [50] "<button class=\"vector-pinnable-header-toggle-button vector-pinnable-header-unpin-button\" data-event-name=\"ui:toggle\">"
## [51] "<div id=\"p-navigation\" class=\"vector-menu mw-portlet mw-portlet-navigation\">"
## [52] "<li id=\"n-mainpage-description\" class=\"mw-list-item\"><a href=\"/wiki/Main_Page\" title=\"Visit the main page\">Main page</a>"
## [53] "<div id=\"p-interaction\" class=\"vector-menu mw-portlet mw-portlet-interaction\">"
## [54] "<li id=\"n-help\" class=\"mw-list-item\"><a href=\"/wiki/Help:Contents\" title=\"Guidance on how to use Wikipedia\">Help</a>"
## [55] "<img class=\"mw-logo-icon\" src=\"/static/images/icons/wikipedia.png\" alt=\"Wikipedia logo\" aria-hidden=\"true\" />"
## [56] "<span class=\"mw-logo-container skin-invert\">"
## [57] "<img class=\"mw-logo-wordmark\" alt=\"Wikipedia wordmark\" src=\"/static/images/mobile/copyright/wikipedia-wordmark.png\" />"
## [58] "<img class=\"mw-logo-tagline\" alt=\"The Free Encyclopedia\" src=\"/static/images/mobile/copyright/wikipedia-tagline.png\" />"
## [59] "</span>"
## [60] "<div id=\"p-search\" role=\"search\" class=\"vector-search-box-vue vector-search-box-collapses vector-search-box-state\">"
## [61] "<a href=\"/wiki/Special:Search\" class=\"cdx-button cdx-button--fake-button cdx-button--fake-button--element\">"
## [62] "<span>Search</span>"
## [63] "<div class=\"vector-typeahead-search-container\">"
## [64] "<div class=\"cdx-typeahead-search cdx-typeahead-search--show-thumbnail cdx-typeahead-search--auto-expand\">"
## [65] "<form action=\"/w/index.php\" id=\"searchform\" class=\"cdx-search-input cdx-search-input--has-end\">"
## [66] "<div id=\"simpleSearch\" class=\"cdx-search-input__input-wrapper\" data-search-loc=\"header-move\">"
## [67] "<div class=\"cdx-text-input cdx-text-input--has-start-icon\">"
## [68] "<input type=\"text\" class=\"cdx-text-input__input\" value=\"\" />"
## [69] "<input type=\"text\" class=\"cdx-text-input__input\" value=\"\" />"
## [70] "<input type=\"text\" class=\"cdx-text-input__input\" value=\"\" />"
## [71] "<span class=\"cdx-text-input__icon cdx-text-input__start-icon\"></span>"
## [72] "<input type=\"hidden\" name=\"title\" value=\"Special:Search\" />"
## [73] "<button class=\"cdx-button cdx-search-input__end-button\">Search</button>"
## [74] "<div id=\"p-vector-user-menu-preferences\" class=\"vector-menu mw-portlet emptyPortlet\">"
## [75] "<div id=\"p-vector-user-menu-userpage\" class=\"vector-menu mw-portlet emptyPortlet\">"
## [76] "<nav class=\"vector-appearance-landmark\" aria-label=\"Appearance\">"
## [77] "<div id=\"vector-appearance-dropdown\" class=\"vector-dropdown\" title=\"Change the appearance of the page\">"
## [78] "<input type=\"checkbox\" id=\"vector-appearance-dropdown-checkbox\" role=\"button\" aria-haspopup=\"true\" data-event-name=\"ui:dropdown\" />"
## [79] "<label id=\"vector-appearance-dropdown-label\" for=\"vector-appearance-dropdown-checkbox\" class=\"vector-dropdown-label\">"
## [80] "<span class=\"vector-dropdown-label-text\">Appearance</span>"
## [81] "<div class=\"vector-dropdown-content\">"
## [82] "<div id=\"vector-appearance-unpinned-container\" class=\"vector-unpinned-container\">"
## [83] "<div id=\"p-vector-user-menu-notifications\" class=\"vector-menu mw-portlet emptyPortlet\">"
## [84] "<div id=\"p-vector-user-menu-overflow\" class=\"vector-menu mw-portlet\">"
## [85] "<li id=\"pt-sitesupport-2\" class=\"user-links-collapsible-item mw-list-item user-links-collapsible-item\">"
## [86] "<li id=\"pt-createaccount-2\" class=\"user-links-collapsible-item mw-list-item user-links-collapsible-item\">"
```

```
## [87] "<li id=\"pt-login-2\" class=\"user-links-collapsible-item mw-list-item user-links-collapsible-item\"><a c
## [88] "<div id=\"vector-user-links-dropdown\" class=\"vector-dropdown vector-user-menu vector-button-flush-right
## [89] \"<input type=\"checkbox\" id=\"vector-user-links-dropdown-checkbox\" role=\"button\" aria-haspopup=\"true
## [90] \"<label id=\"vector-user-links-dropdown-label\" for=\"vector-user-links-dropdown-checkbox\" class=\"vect
## [91] "<span class=\"vector-dropdown-label-text\">Personal tools</span>"
## [92] \"<div class=\"vector-dropdown-content\">"
## [93] "<div id=\"p-personal\" class=\"vector-menu mw-portlet mw-portlet-personal user-links-collapsible-item\"
## [94] \"<li id=\"pt-sitesupport\" class=\"user-links-collapsible-item mw-list-item\"><a href=\"https://don
## [95] "<div id=\"p-user-menu-anon-editor\" class=\"vector-menu mw-portlet mw-portlet-user-menu-anon-editor\" >
## [96] \"<Pages for logged out editors <a href=\"/wiki/Help:Introduction\" aria-label=\"Learn more about editin
## [97] \"<li id=\"pt-anoncontribs\" class=\"mw-list-item\"><a href=\"/wiki/Special:MyContributions\" title=
## [98] "<div class=\"mw-page-container\">"
## [99] \"<div class=\"mw-page-container-inner\">"
## [100] \"<nav id=\"mw-panel\" class=\"vector-main-menu-landmark\" aria-label=\"Site\">"
## [101] \"<div id=\"vector-main-menu-pinned-container\" class=\"vector-pinned-container\">"
## [102] \"<div class=\"vector-sticky-pinned-container\">"
## [103] \"<nav id=\"mw-panel-toc\" aria-label=\"Contents\" data-event-name=\"ui.sidebar-toc\" class=\"mw-t
## [104] \"<div id=\"vector-toc-pinned-container\" class=\"vector-pinned-container\">"
## [105] \"<div id=\"vector-toc\" class=\"vector-toc vector-pinnable-element\">"
## [106] \"<class=\"vector-pinnable-header vector-toc-pinnable-header vector-pinnable-header-pinned\""
## [107] \"<data-feature-name=\"toc-pinned\""
## [108] \"<data-pinnable-element-id=\"vector-toc\""
## [109] \"<h2 class=\"vector-pinnable-header-label\">Contents</h2>"
## [110] \"<button class=\"vector-pinnable-header-toggle-button vector-pinnable-header-pin-button\" data-event-nam
## [111] \"<button class=\"vector-pinnable-header-toggle-button vector-pinnable-header-unpin-button\" data-event-r
## [112] \"<ul class=\"vector-toc-contents\" id=\"mw-panel-toc-list\">"
## [113] \"<div class=\"vector-toc-text\">(Top)</div>"
## [114] \"<class=\"vector-toc-list-item vector-toc-level-1 vector-toc-list-item-expanded\""
## [115] \"<span class=\"vector-toc-numb\">1</span>"
## [116] \"<span>Foundations</span>"
## [117] \"<class=\"vector-toc-list-item vector-toc-level-1 vector-toc-list-item-expanded\""
## [118] \"<span class=\"vector-toc-numb\">2</span>"
## [119] \"<span>Etymology</span>"
## [120] \"<span class=\"vector-icon mw-ui-icon-wikimedia-expand\"></span>"
## [121] \"<span>Toggle Etymology subsection</span>"
## [122] \"<span class=\"vector-toc-numb\">2.1</span>"
## [123] \"<span>Early usage</span>"
## [124] \"<span class=\"vector-toc-numb\">2.2</span>"
## [125] \"<span>Modern usage</span>"
## [126] \"<class=\"vector-toc-list-item vector-toc-level-1 vector-toc-list-item-expanded\""
## [127] \"<span class=\"vector-toc-numb\">3</span>"
## [128] \"<span>Data science and data analysis</span>"
## [129] \"<li id=\"toc-Cloud_computing_for_data_science\""
## [130] \"<class=\"vector-toc-list-item vector-toc-level-1 vector-toc-list-item-expanded\""
## [131] \"<a class=\"vector-toc-link\" href=\"#Cloud_computing_for_data_science\">"
## [132] \"<span class=\"vector-toc-numb\">4</span>"
## [133] \"<span>Cloud computing for data science</span>"
## [134] \"<ul id=\"toc-Cloud_computing_for_data_science-sublist\" class=\"vector-toc-list\">"
## [135] \"<class=\"vector-toc-list-item vector-toc-level-1 vector-toc-list-item-expanded\""
## [136] \"<span class=\"vector-toc-numb\">5</span>"
## [137] \"<span>Ethical consideration in data science</span>"
## [138] \"<class=\"vector-toc-list-item vector-toc-level-1 vector-toc-list-item-expanded\""
## [139] \"<span class=\"vector-toc-numb\">6</span>"
## [140] \"<span>See also</span>"
## [141] \"<class=\"vector-toc-list-item vector-toc-level-1 vector-toc-list-item-expanded\""
## [142] \"<span class=\"vector-toc-numb\">7</span>"
## [143] \"<span>References</span>"
## [144] \"<header class=\"mw-body-header vector-page-titlebar\">"
## [145] "<div id=\"vector-page-titlebar-toc\" class=\"vector-dropdown vector-page-titlebar-toc vector-button-flush
## [146] \"<input type=\"checkbox\" id=\"vector-page-titlebar-toc-checkbox\" role=\"button\" aria-haspopup=\"true\"
## [147] \"<label id=\"vector-page-titlebar-toc-label\" for=\"vector-page-titlebar-toc-checkbox\" class=\"vector-c
## [148] "<span class=\"vector-dropdown-label-text\">Toggle the table of contents</span>"
## [149] \"<div class=\"vector-dropdown-content\">"
## [150] \"<div id=\"vector-page-titlebar-toc-unpinned-container\" class=\"vector-unpinned-container\">"
```

```
## [151] "\t\t\t\t\t<h1 id=\"firstHeading\" class=\"firstHeading mw-first-heading\"><span class=\"mw-page-title-ma
## [152] "<div id=\"p-lang-btn\" class=\"vector-dropdown mw-portlet mw-portlet-lang\" >"
## [153] "\t<input type=\"checkbox\" id=\"p-lang-btn-checkbox\" role=\"button\" aria-haspopup=\"true\" data-event-r
## [154] "\t<label id=\"p-lang-btn-label\" for=\"p-lang-btn-checkbox\" class=\"vector-dropdown-label cdx-button cd
## [155] "<span class=\"vector-dropdown-label-text\">50 languages</span>"
## [156] "\t<div class=\"vector-dropdown-content\">"
## [157] "\t\t\t\t\t<li class=\"interlanguage-link interwiki-ar mw-list-item\"><a href=\"https://ar.wikipedia.org/wi
Arabic\" lang=\"ar\" hreflang=\"ar\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-name
Azerbaijani\" lang=\"az\" hreflang=\"az\" data-title=\"Verilənlər elmi\" data-language-autonym=\"Azərbaycanca\" da
Bangla\" lang=\"bn\" hreflang=\"bn\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-name=
Bulgarian\" lang=\"bg\" hreflang=\"bg\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-loc
Catalan\" lang=\"ca\" hreflang=\"ca\" data-title=\"Ciència de les dades\" data-language-autonym=\"Català\" data-lan
Czech\" lang=\"cs\" hreflang=\"cs\" data-title=\"Data science\" data-language-autonym=\"Čeština\" data-language-lo
German\" lang=\"de\" hreflang=\"de\" data-title=\"Data Science\" data-language-autonym=\"Deutsch\" data-language-l
Estonian\" lang=\"et\" hreflang=\"et\" data-title=\"Andmeteadus\" data-language-autonym=\"Eesti\" data-language-lo
Greek\" lang=\"el\" hreflang=\"el\" data-title=\"E\" \\" data-language-autonym=\"E\" \\" data-language-local-
Spanish\" lang=\"es\" hreflang=\"es\" data-title=\"Ciencia de datos\" data-language-autonym=\"Español\" data-langua
Esperanto\" lang=\"eo\" hreflang=\"eo\" data-title=\"Datum-scienco\" data-language-autonym=\"Esperanto\" data-langu
Basque\" lang=\"eu\" hreflang=\"eu\" data-title=\"Datu zientzia\" data-language-autonym=\"Euskara\" data-language-
Persian\" lang=\"fa\" hreflang=\"fa\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-name=
French\" lang=\"fr\" hreflang=\"fr\" data-title=\"Science des données\" data-language-autonym=\"Français\" data-lan
Galician\" lang=\"gl\" hreflang=\"gl\" data-title=\"Ciencia de datos\" data-language-autonym=\"Galego\" data-langua
Korean\" lang=\"ko\" hreflang=\"ko\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-name=
Armenian\" lang=\"hy\" hreflang=\"hy\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-loc
Hindi\" lang=\"hi\" hreflang=\"hi\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-name=
Ido\" lang=\"io\" hreflang=\"io\" data-title=\"Cienco di datumi\" data-language-autonym=\"Ido\" data-language-local
Indonesian\" lang=\"id\" hreflang=\"id\" data-title=\"Ilmu data\" data-language-autonym=\"Bahasa Indonesia\" data-
Zulu\" lang=\"zu\" hreflang=\"zu\" data-title=\"INzululwazi yeMininingo\" data-language-autonym=\"IsiZulu\" data-l
Italian\" lang=\"it\" hreflang=\"it\" data-title=\"Scienza dei dati\" data-language-autonym=\"Italiano\" data-langu
Hebrew\" lang=\"he\" hreflang=\"he\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-name=
Kannada\" lang=\"kn\" hreflang=\"kn\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-name=
Kazakh\" lang=\"kk\" hreflang=\"kk\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-
Latvian\" lang=\"lv\" hreflang=\"lv\" data-title=\"Datu mācība\" data-language-autonym=\"Latviešu\" data-language-
Macedonian\" lang=\"mk\" hreflang=\"mk\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-
Malay\" lang=\"ms\" hreflang=\"ms\" data-title=\"Sains data\" data-language-autonym=\"Bahasa Melayu\" data-langua
Burmese\" lang=\"my\" hreflang=\"my\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-na
Dutch\" lang=\"nl\" hreflang=\"nl\" data-title=\"Datawetenschap\" data-language-autonym=\"Nederlands\" data-langua
Japanese\" lang=\"ja\" hreflang=\"ja\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-name=
Norwegian Bokmål\" lang=\"nb\" hreflang=\"nb\" data-title=\"Datavitenskap\" data-language-autonym=\"Norsk bokmål\"
Punjabi\" lang=\"pa\" hreflang=\"pa\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-name=
Polish\" lang=\"pl\" hreflang=\"pl\" data-title=\"Danologia\" data-language-autonym=\"Polski\" data-language-local
Portuguese\" lang=\"pt\" hreflang=\"pt\" data-title=\"Ciência de dados\" data-language-autonym=\"Português\" data-
Quechua\" lang=\"qu\" hreflang=\"qu\" data-title=\"Willakuy hamut&#039;ay\" data-language-autonym=\"Runa Simi\" da
Russian\" lang=\"ru\" hreflang=\"ru\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-n
Albanian\" lang=\"sq\" hreflang=\"sq\" data-title=\"Shkenca e të dhënave\" data-language-autonym=\"Shqip\" data-lan
Simple English\" lang=\"en-simple\" hreflang=\"en-simple\" data-title=\"Data science\" data-language-autonym=\"Simp
Sindhi\" lang=\"sd\" hreflang=\"sd\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-name=
Serbian\" lang=\"sr\" hreflang=\"sr\" data-title=\"\" \\" data-language-autonym=\"\" / srpski\" data-langua
Finnish\" lang=\"fi\" hreflang=\"fi\" data-title=\"Datatiede\" data-language-autonym=\"Suomi\" data-language-local
Tamil\" lang=\"ta\" hreflang=\"ta\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-name=
Thai\" lang=\"th\" hreflang=\"th\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-name=
Turkish\" lang=\"tr\" hreflang=\"tr\" data-title=\"Veri bilimi\" data-language-autonym=\"Türkçe\" data-language-lo
Ukrainian\" lang=\"uk\" hreflang=\"uk\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-loc
Urdu\" lang=\"ur\" hreflang=\"ur\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-name=
Vietnamese\" lang=\"vi\" hreflang=\"vi\" data-title=\"Khoa học dữ liệu\" data-language-autonym=\"Tiếng Việt\" data-
Cantonese\" lang=\"yue\" hreflang=\"yue\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-name=
Chinese\" lang=\"zh\" hreflang=\"zh\" data-title=\"\" \\" data-language-autonym=\"\" \\" data-language-local-name=
## [158] "\t\t\t\t<div class=\"after-portlet after-portlet-lang\"><span class=\"wb-langlinks-edit wb-langlinks-link\"
## [159] "\t\t\t\t\t<div class=\"vector-page-toolbar\">"
## [160] "\t\t\t\t\t\t\t<div class=\"vector-page-toolbar-container\">"
## [161] "\t\t\t\t\t\t\t\t\t<nav aria-label=\"Namespaces\">"
## [162] "<div id=\"p-associated-pages\" class=\"vector-menu vector-menu-tabs mw-portlet mw-portlet-associated-page
## [163] "\t\t\t\t<li id=\"ca-nstab-main\" class=\"selected vector-tab-noicon mw-list-item\"><a href=\"/wiki/Data_sc
## [164] "<div id=\"vector-variants-dropdown\" class=\"vector-dropdown emptyPortlet\" >"
```

```
## [165] "\t<input type=\"checkbox\" id=\"vector-variants-dropdown-checkbox\" role=\"button\" aria-haspopup=\"true
```

```
## [166] "\t<label id=\"vector-variants-dropdown-label\" for=\"vector-variants-dropdown-checkbox\" class=\"vector-
```

```
## [167] "\t<div class=\"vector-dropdown-content\">"
```

```
## [168] "<div id=\"p-variants\" class=\"vector-menu mw-portlet mw-portlet-variants emptyPortlet\" >"
```

```
## [169] "\t\t\t\t\t\t<div id=\"right-navigation\" class=\"vector-collapsible\">"
```

```
## [170] "<div id=\"p-views\" class=\"vector-menu vector-menu-tabs mw-portlet mw-portlet-views\" >"
```

```
## [171] "\t\t\t\t<li id=\"ca-view\" class=\"selected vector-tab-noicon mw-list-item\"><a href=\"/wiki/Data_science\
```

```
## [172] "\t\t\t\t\t\t\t\t<nav class=\"vector-page-tools-landmark\" aria-label=\"Page tools\">"
```

```
## [173] "<div id=\"vector-page-tools-dropdown\" class=\"vector-dropdown vector-page-tools-dropdown\" >"
```

```
## [174] "\t<input type=\"checkbox\" id=\"vector-page-tools-dropdown-checkbox\" role=\"button\" aria-haspopup=\"tru
```

```
## [175] "\t<label id=\"vector-page-tools-dropdown-label\" for=\"vector-page-tools-dropdown-checkbox\" class=\"vec
```

```
## [176] "\t<div class=\"vector-dropdown-content\">"
```

```
## [177] "\t\t\t\t\t\t\t\t\t\t\t\t<div id=\"vector-page-tools-unpinned-container\" class=\"vector-unpinned-container\">"
```

```
## [178] "<div id=\"vector-page-tools\" class=\"vector-page-tools vector-pinnable-element\">"
```

```
## [179] "\tclass=\"vector-pinnable-header vector-page-tools-pinnable-header vector-pinnable-header-unpinned\""
```

```
## [180] "\tdata-feature-name=\"page-tools-pinned\""
```

```
## [181] "\tdata-pinnable-element-id=\"vector-page-tools\""
```

```
## [182] "\tdata-pinned-container-id=\"vector-page-tools-pinned-container\""
```

```
## [183] "\tdata-unpinned-container-id=\"vector-page-tools-unpinned-container\""
```

```
## [184] "\t<div class=\"vector-pinnable-header-label\">Tools</div>"
```

```
## [185] "\t<button class=\"vector-pinnable-header-toggle-button vector-pinnable-header-pin-button\" data-event-na
```

```
## [186] "\t<button class=\"vector-pinnable-header-toggle-button vector-pinnable-header-unpin-button\" data-event-i
```

```
## [187] "<div id=\"p-cactions\" class=\"vector-menu mw-portlet mw-portlet-cactions emptyPortlet vector-has-collap
```

```
## [188] "\t\t\t\t<li id=\"ca-more-view\" class=\"selected vector-more-collapsible-item mw-list-item\"><a href=\"/wil
```

```
## [189] "<div id=\"p-tb\" class=\"vector-menu mw-portlet mw-portlet-tb\" >"
```

```
## [190] "\t\t\t\t<li id=\"t-whatlinkshere\" class=\"mw-list-item\"><a href=\"/wiki/Special:WhatLinksHere/Data_scienc
```

```
## [191] "<div id=\"p-coll-print_export\" class=\"vector-menu mw-portlet mw-portlet-coll-print_export\" >"
```

```
## [192] "\t\t\tPrint/export"
```

```
## [193] "\t\t\t\t<li id=\"coll-download-as-rl\" class=\"mw-list-item\"><a href=\"/w/index.php?title=Special:Downloa
```

```
## [194] "<div id=\"p-wikibase-otherprojects\" class=\"vector-menu mw-portlet mw-portlet-wikibase-otherprojects\">
```

```
## [195] "\t\t\tIn other projects"
```

```
## [196] "\t\t\t\t<li class=\"wb-otherproject-link wb-otherproject-commons mw-list-item\"><a href=\"https://commons.
```

```
## [197] "\t\t\t\t\t\t\t\t<div class=\"vector-sticky-pinned-container\">"
```

```
## [198] "\t\t\t\t\t\t\t\t\t\t<nav class=\"vector-page-tools-landmark\" aria-label=\"Page tools\">"
```

```
## [199] "\t\t\t\t\t\t\t\t\t\t<div id=\"vector-page-tools-pinned-container\" class=\"vector-pinned-container\">"
```

```
## [200] "\t\t\t\t\t\t\t\t\t\t<nav class=\"vector-appearance-landmark\" aria-label=\"Appearance\">"
```

```
## [201] "\t\t\t\t\t\t\t\t\t\t<div id=\"vector-appearance-pinned-container\" class=\"vector-pinned-container\">"
```

```
## [202] "\t\t\t\t\t\t\t\t<div id=\"vector-appearance\" class=\"vector-appearance vector-pinnable-element\">"
```

```
## [203] "\tclass=\"vector-pinnable-header vector-appearance-pinnable-header vector-pinnable-header-pinned\""
```

```
## [204] "\tdata-feature-name=\"appearance-pinned\""
```

```
## [205] "\tdata-pinnable-element-id=\"vector-appearance\""
```

```
## [206] "\tdata-pinned-container-id=\"vector-appearance-pinned-container\""
```

```
## [207] "\tdata-unpinned-container-id=\"vector-appearance-unpinned-container\""
```

```
## [208] "\t<div class=\"vector-pinnable-header-label\">Appearance</div>"
```

```
## [209] "\t<button class=\"vector-pinnable-header-toggle-button vector-pinnable-header-pin-button\" data-event-na
```

```
## [210] "\t<button class=\"vector-pinnable-header-toggle-button vector-pinnable-header-unpin-button\" data-event-i
```

```
## [211] "\t\t\t\t\t\t\t\t<div id=\"siteSub\" class=\"noprint\">From Wikipedia, the free encyclopedia</div>"
```

```
## [212] "\t\t\t\t\t\t\t\t<div id=\"mw-content-text\" class=\"mw-body-content\"><div class=\"mw-content-ltr mw-parser-out
```

```
## [213] "<style data-mw-deduplicate=\"TemplateStyles:r1236090951\">.mw-parser-output .hatnote{font-style:italic}.r
```

```
## [214] "<p class=\"mw-empty-elt\">"
```

```
## [215] "</p>"
```

```
## [216] "<figure class=\"mw-default-size\" typeof=\"mw:File/Thumb\"><a href=\"/wiki/File:PIA23792-1600x1200(1).jpg
```

```
## [217] "<p><b>Data science</b> is an <a href=\"/wiki/Interdisciplinary\" class=\"mw-redirect\" title=\"Interdisci
```

```
## [218] "</p><p>Data science also integrates domain knowledge from the underlying application domain (e.g., natura
```

```
## [219] "</p><p>Data science is \\a concept to unify <a href=\"/wiki/Statistics\" title=\"Statistics\">statistics
```

```
## [220] "</p><p>A <b>data scientist</b> is a professional who creates programming code and combines it with statis
```

```
## [221] "</p>"
```

```
## [222] "<meta property=\"mw:PageProp/toc\" />"
```

```
## [223] "<div class=\"mw-heading mw-heading2\"><h2 id=\"Foundations\">Foundations</h2><span class=\"mw-editsection
```

```
## [224] "<p>Data science is an <a href=\"/wiki/Interdisciplinarity\" title=\"Interdisciplinarity\">interdisciplina
```

```
## [225] "</p><p><a href=\"/wiki/Vasant_Dhar\" title=\"Vasant Dhar\">Vasant Dhar</a> writes that statistics emphas
```

```
## [226] "</p>"
```

```
## [227] "<div class=\"mw-heading mw-heading2\"><h2 id=\"Etymology\">Etymology</h2><span class=\"mw-editsection\">
```

```
## [228] "<div class=\"mw-heading mw-heading3\"><h3 id=\"Early usage\">Early usage</h3><span class=\"mw-editsectio
```



```
## [229] "<p>In 1962, <a href=\"/wiki/John_Tukey\" title=\"John Tukey\">John Tukey</a> described a field he called
## [230] "</p><p>The term \"data science\" has been traced back to 1974, when <a href=\"/wiki/Peter_Naur\" title=\"
## [231] "</p>"
## [232] "<div class=\"mw-heading mw-heading3\"><h3 id=\"Modern_usage\">Modern usage</h3><span class=\"mw-editsect
## [233] "<p>In 2012, technologists <a href=\"/wiki/Thomas_H._Davenport\" title=\"Thomas H. Davenport\">Thomas H.
## [234] "</p><p>The modern conception of data science as an independent discipline is sometimes attributed to <a
## [235] "</p><p>The professional title of \"data scientist\" has been attributed to <a href=\"/wiki/DJ_Patil\" ti
## [236] "</p>"
## [237] "<div class=\"mw-heading mw-heading2\"><h2 id=\"Data_science_and_data_analysis\">Data science and data an
## [238] "<figure class=\"mw-default-size\" typeof=\"mw:File/Thumb\"><a href=\"/wiki/File:EDA_example_-_Always_plo
## [239] "<figure class=\"mw-default-size\" typeof=\"mw:File/Thumb\"><a href=\"/wiki/File:Data_Science.png\" class=
## [240] "<p>Data analysis typically involves working with structured datasets to answer specific questions or sol
## [241] "</p><p>Data science involves working with larger datasets that often require advanced computational and
## [242] "</p>"
## [243] "<div class=\"mw-heading mw-heading2\"><h2 id=\"Cloud_computing_for_data_science\">Cloud computing for da
## [244] "<figure class=\"mw-default-size\" typeof=\"mw:File/Thumb\"><a href=\"/wiki/File:Cloud_computing_in_enabl
## [245] "<p><a href=\"/wiki/Cloud_computing\" title=\"Cloud computing\">Cloud computing</a> can offer access to la
## [246] "</p><p>Some distributed computing frameworks are designed to handle big data workloads. These frameworks
## [247] "</p>"
## [248] "<div class=\"mw-heading mw-heading2\"><h2 id=\"Ethical_consideration_in_data_science\">Ethical considerat
## [249] "<p>Data science involves collecting, processing, and analyzing data which often includes personal and ser
## [250] "</p><p>Machine learning models can amplify existing biases present in training data, leading to discrimin
## [251] "</p>"
## [252] "<div class=\"mw-heading mw-heading2\"><h2 id=\"See_also\">See also</h2><span class=\"mw-editsection\"><sp
## [253] "<ul><li><a href=\"/wiki/Python_(programming_language)\" title=\"Python (programming language)\">Python (p
## [254] "<li><a href=\"/wiki/R_(programming_language)\" title=\"R (programming language)\">R (programming languag
## [255] "<li><a href=\"/wiki/Topological_data_analysis\" title=\"Topological data analysis\">Topological data anal
## [256] "<li><a href=\"/wiki/List_of_free_and_open-source_software_packages#Data_science\" title=\"List of free an
## [257] "<div class=\"mw-heading mw-heading2\"><h2 id=\"References\">References</h2><span class=\"mw-editsection\"
## [258] "<style data-mw-deduplicate=\"TemplateStyles:r1239543626\">.mw-parser-output .reflist{margin-bottom:0.5em
## [259] "<div class=\"mw-references-wrap mw-references-columns\"><ol class=\"references\">"
## [260] "<li id=\"cite_note-1\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-1\"></a></b></span> <span
## [261] "<li id=\"cite_note-2\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-2\"></a></b></span> <span
## [262] "<li id=\"cite_note-3\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-3\"></a></b></span> <span
## [263] "<li id=\"cite_note-4\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-4\"></a></b></span> <span
## [264] "<li id=\"cite_note-5\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-5\"></a></b></span> <span
## [265] "<li id=\"cite_note-:2-6\"><span class=\"mw-cite-backlink\">^ <a href=\"#cite_ref-:2_6-0\"><sup><i><b>a</b></i></sup>
## [266] "<li id=\"cite_note-TansleyTolle2009-7\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-TansleyT
## [267] "<li id=\"cite_note-BellHey2009-8\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-BellHey2009-8
## [268] "<li id=\"cite_note-9\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-9\"></a></b></span> <span
## [269] "<li id=\"cite_note-10\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-10\"></a></b></span> <span
## [270] "<li id=\"cite_note-11\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-11\"></a></b></span> <span
## [271] "<li id=\"cite_note-12\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-12\"></a></b></span> <span
## [272] "<li id=\"cite_note-13\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-13\"></a></b></span> <span
## [273] "<li id=\"cite_note-:7-14\"><span class=\"mw-cite-backlink\">^ <a href=\"#cite_ref-:7_14-0\"><sup><i><b>a</b></i></sup>
## [274] "<li id=\"cite_note-15\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-15\"></a></b></span> <span
## [275] "<li id=\"cite_note-16\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-16\"></a></b></span> <span
## [276] "<li id=\"cite_note-Murtagh_2018_14-17\"><span class=\"mw-cite-backlink\">^ <a href=\"#cite_ref-Murtagh_20
## [277] "<li id=\"cite_note-18\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-18\"></a></b></span> <span
## [278] "<li id=\"cite_note-19\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-19\"></a></b></span> <span
## [279] "<li id=\"cite_note-20\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-20\"></a></b></span> <span
## [280] "<li id=\"cite_note-21\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-21\"></a></b></span> <span
## [281] "<li id=\"cite_note-22\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-22\"></a></b></span> <span
## [282] "<li id=\"cite_note-23\"><span class=\"mw-cite-backlink\"><b><a href=\"#cite_ref-23\"></a></b></span> <span
```

[283] "<li id=\"cite_note-24\">^
[284] "<li id=\"cite_note-25\">^
[285] "<li id=\"cite_note-26\">^
NSB-05-40, Long-Lived Digital Data Collections Enabling Research and Education in the 21st Century\". <i>www.n
[286] "<li id=\"cite_note-8-27\">^
[287] "<li id=\"cite_note-9-28\">^
[288] "<li id=\"cite_note-10-29\">^
[289] "<li id=\"cite_note-30\">^
115. doi:<a rel=\"n
[290] "<li id=\"cite_note-0-31\">^
[291] "<li id=\"cite_note-3-32\">^
1394. doi:<a rel=\"n
[292] "<li id=\"cite_note-4-33\">^
[293] "<li id=\"cite_note-11-34\">^
341. doi:<a rel=\"n
[294] "<li id=\"cite_note-35\">^
[295] "<li id=\"cite_note-36\">^
186. arXiv:
[296] "<div class=\"navbox-styles\"><style data-mw-deduplicate=\"TemplateStyles:r1129693374\">.mw-parser-output
[297] "Compression
[298] "Corruption
[299] "Deduplication
[300] "Exploration
[301] "Cooperatives
[302] "Philanthropy
[303] "Pre-proces
[304] "Preservation
[305] "Processing
[306] "Protection (privacy)
[307] "Publishing
[308] "Open data
[309] "Scraping
[310] "Stewardship
[311] "Topological data anal
[312] "Type
[313] "NewPP limit report"
[314] "Cache expiry: 2592000"
[315] "Reduced expiry: false"
[316] "Complications: [vary-revision-sha1, show-toc]"
[317] "Preprocessor visited node count: 2302/1000000"
[318] "Post-expand include size: 85027/2097152 bytes"
[319] "Template argument size: 780/2097152 bytes"
[320] "Highest expansion depth: 12/100"
[321] "Expensive parser function count: 4/500"
[322] "Unstrip recursion depth: 1/20"
[323] "Unstrip post-expand size: 142478/5000000 bytes"
[324] "Transclusion expansion time report (%,ms,calls,template)"
[325] " 65.45% 361.134 1 Template:Reflist"
[326] " 32.63% 180.049 16 Template:Cite_journal"
[327] " 13.97% 77.061 1 Template:Data"
[328] " 13.23% 73.001 1 Template:Navbar"
[329] " 11.32% 62.486 1 Template:Short_description"
[330] " 7.58% 41.816 6 Template:Cite_book"
[331] " 7.07% 39.038 7 Template:Cite_web"
[332] " 6.38% 35.213 4 Template:Cite_news"
[333] " 6.16% 34.012 2 Template:Pagetype"
[334] "<!-- Saved in parser cache with key enwiki:pcache:35458904:|#:idhash:canonical and timestamp 2025041216
[335] "</div><!--esi <esi:include src=\"/esitest-fa8a495983347898/content\" /> --><noscript><img src=\"https://
[336] "<div class=\"printfooter\" data-nosnippet=\"\">Retrieved from \"<a dir=\"ltr\" href=\"https://en.wikiped
[337] "–––<div id=\"catlinks\" class=\"catlinks\" data-mw=\"interface\"><div id=\"mw-normal-catlinks\" c
[338] "–––<li id=\"footer-info-lastmod\"> This page was last edited on 17 March 2025, at 07:41<span class=\"anony
[339] "–––<li id=\"footer-info-copyright\">Text is available under the <a href=\"/wiki/Wikipedia:Text_of_the_Crea
[340] "additional terms may apply. By using this site, you agree to the <a href=\"https://foundation.wikimedia.
[341] "–––<ul id=\"footer-places\">"

```

## [342] "\t<li id=\"footer-places-privacy\"><a href=\"https://foundation.wikimedia.org/wiki/Special:MyLanguage/Policy:Privacy_policy\">Privacy policy</a></li>\"
## [343] "\t<li id=\"footer-places-about\"><a href=\"/wiki/Wikipedia:About\">About Wikipedia</a></li>\"
## [344] "\t<li id=\"footer-places-disclaimers\"><a href=\"/wiki/Wikipedia:General_disclaimer\">Disclaimers</a></li>\"
## [345] "\t<li id=\"footer-places-contact\"><a href=\"//en.wikipedia.org/wiki/Wikipedia:Contact_us\">Contact Wikipedia</a></li>\"
## [346] "\t<li id=\"footer-places-wm-codeofconduct\"><a href=\"https://foundation.wikimedia.org/wiki/Special:MyLanguage/Policy:Wikimedia_foundation_code_of_conduct\">Wikimedia Foundation Code of Conduct</a></li>\"
## [347] "\t<li id=\"footer-places-developers\"><a href=\"https://developer.wikimedia.org\">Developers</a></li>\"
## [348] "\t<li id=\"footer-places-statslink\"><a href=\"https://stats.wikimedia.org/#/en.wikipedia.org\">Statistics</a></li>\"
## [349] "\t<li id=\"footer-places-cookiestatement\"><a href=\"https://foundation.wikimedia.org/wiki/Special:MyLanguage/Policy:Cookie_statement\">Cookie statement</a></li>\"
## [350] "\t<li id=\"footer-places-mobileview\"><a href=\"//en.m.wikipedia.org/w/index.php?title=Data_science&mobileaction=toggle_view_mobile\">Mobile view</a></li>\"
## [351] "\t<ul id=\"footer-icons\" class=\"noprint\">\"
## [352] "\t<li id=\"footer-copyrightico\"><a href=\"https://www.wikimedia.org/\" class=\"cdx-button cdx-button--fake-button cdx-button--fake-button--enabled cdx-button--weight-quiet cdx-button--icon-only\"><span><img alt=\"Wikimedia logo\" data-bbox=\"1000 17 1000 35\" data-cs=\"1\" data-kind=\"parent\" data-rs=\"2\"/></span></a></li>\"
## [353] "\t<li id=\"footer-poweredbyico\"><a href=\"https://www.mediawiki.org/\" class=\"cdx-button cdx-button--fake-button cdx-button--fake-button--enabled cdx-button--weight-quiet cdx-button--icon-only\"><span><img alt=\"MediaWiki logo\" data-bbox=\"1000 35 1000 53\" data-cs=\"1\" data-kind=\"parent\" data-rs=\"2\"/></span></a></li>\"
## [354] "\t\t\t\t<button class=\"cdx-button cdx-button--weight-quiet cdx-button--icon-only vector-sticky-header-search\" type=\"button\">\"
## [355] "<span>Search</span>\"
## [356] "\t\t\t\t<div class=\"vector-typeahead-search-container\">\"
## [357] "\t\t\t\t\t<div class=\"cdx-typeahead-search cdx-typeahead-search--show-thumbnail\">\"
## [358] "\t\t\t\t\t\t<form action=\"/w/index.php\" id=\"vector-sticky-search-form\" class=\"cdx-search-input cdx-search-input--has-start-icon\">\"
## [359] "\t\t\t\t\t\t\t<div class=\"cdx-search-input__input-wrapper\" data-search-loc=\"header-moved\">\"
## [360] "\t\t\t\t\t\t\t\t<div class=\"cdx-text-input cdx-text-input--has-start-icon\">\"
## [361] "\t\t\t\t\t\t\t\t\t<input type=\"text\">\"
## [362] "\t\t\t\t\t\t\t\t\t\tclass=\"cdx-text-input__input\">\"
## [363] "\t\t\t\t\t\t\t\t\t\t\ttype=\"search\" name=\"search\" placeholder=\"Search Wikipedia\">\"
## [364] "\t\t\t\t\t\t\t\t\t\t\t<span class=\"cdx-text-input__icon cdx-text-input__start-icon\"></span>\"
## [365] "\t\t\t\t\t\t\t\t\t\t\t<input type=\"hidden\" name=\"title\" value=\"Special:Search\">\"
## [366] "\t\t\t\t\t\t\t\t\t\t\t<button class=\"cdx-button cdx-search-input__end-button\">Search</button>\"
## [367] "\t\t\t\t\t\t\t\t<div id=\"vector-sticky-header-toc\" class=\"vector-dropdown mw-portlet mw-portlet-sticky-header\">\"
## [368] "\t\t\t\t\t\t\t\t\t<input type=\"checkbox\" id=\"vector-sticky-header-toc-checkbox\" role=\"button\" aria-haspopup=\"true\">\"
## [369] "\t\t\t\t\t\t\t\t\t<label id=\"vector-sticky-header-toc-label\" for=\"vector-sticky-header-toc-checkbox\" class=\"vector-dropdown-label-text\">\"
## [370] "<span class=\"vector-dropdown-label-text\">Toggle the table of contents</span>\"
## [371] "\t\t\t\t\t\t\t\t\t<div class=\"vector-dropdown-content\">\"
## [372] "\t\t\t\t\t\t\t\t\t\t<div id=\"vector-sticky-header-toc-unpinned-container\" class=\"vector-unpinned-container\">\"
## [373] "\t\t\t\t\t\t\t\t\t\t\t<div class=\"vector-sticky-header-context-bar-primary\" aria-hidden=\"true\"><span class=\"mw-potential\">\"
## [374] "\t\t\t\t\t\t\t\t\t\t\t\t<a href=\"#\" class=\"cdx-button cdx-button--fake-button cdx-button--fake-button--enabled cdx-button--weight-quiet\">\"
## [375] "<span></span>\"
## [376] "\t\t\t\t\t\t\t\t\t\t\t\t\t<a href=\"#\" class=\"cdx-button cdx-button--fake-button cdx-button--fake-button--enabled cdx-button--weight-quiet\">\"
## [377] "<span></span>\"
## [378] "\t\t\t\t\t\t\t\t\t\t\t\t\t<a href=\"#\" class=\"cdx-button cdx-button--fake-button cdx-button--fake-button--enabled cdx-button--weight-quiet\">\"
## [379] "<span></span>\"
## [380] "\t\t\t\t\t\t\t\t\t\t\t\t\t<a href=\"#\" class=\"cdx-button cdx-button--fake-button cdx-button--fake-button--enabled cdx-button--weight-quiet\">\"
## [381] "<span></span>\"
## [382] "\t\t\t\t\t\t\t\t\t\t\t\t\t<a href=\"#\" class=\"cdx-button cdx-button--fake-button cdx-button--fake-button--enabled cdx-button--weight-quiet\">\"
## [383] "<span></span>\"
## [384] "\t\t\t\t\t\t\t\t\t\t\t\t\t<a href=\"#\" class=\"cdx-button cdx-button--fake-button cdx-button--fake-button--enabled cdx-button--weight-quiet\">\"
## [385] "<span></span>\"
## [386] "\t\t\t\t\t\t\t\t\t\t\t\t\t<a href=\"#\" class=\"cdx-button cdx-button--fake-button cdx-button--fake-button--enabled cdx-button--weight-quiet\">\"
## [387] "<span></span>\"
## [388] "\t\t\t\t\t\t\t\t\t\t\t\t\t\t<button class=\"cdx-button cdx-button--weight-quiet mw-interlanguage-selector\" id=\"p-lang-btn\">\"
## [389] "<span>50 languages</span>\"
## [390] "\t\t\t\t\t\t\t\t\t\t\t\t\t\t\t<a href=\"#\" class=\"cdx-button cdx-button--fake-button cdx-button--fake-button--enabled cdx-button--weight-quiet\">\"
## [391] "<span>Add topic</span>\"
## [392] "<div class=\"mw-portlet mw-portlet-dock-bottom emptyPortlet\" id=\"p-dock-bottom\">\"
## [393] "<script>(RLQ=window.RLQ||[]).push(function(){mw.config.set({\"wgHostname\":\"mw-web.codfw.main-658b85fd8\"});\"</script>\"
## [394] "<script type=\"application/ld+json\">{\"@context\":\"https://schema.org\",\"@type\":\"Article\",\"name\":\"\"}</script>\"

```

```

#Using library XML
library(XML)

```

```

## Warning: package 'XML' was built under R version 4.4.3

```

```

doc<-htmlParse(eg6)
doc.text<-unlist(xpathApply(doc,'//p',xmlValue))
unlist(xpathApply(doc,'//h2',xmlValue))

```

```
## [1] "Contents"
## [2] "Foundations"
## [3] "Etymology"
## [4] "Data science and data analysis"
## [5] "Cloud computing for data science"
## [6] "Ethical consideration in data science"
## [7] "See also"
## [8] "References"
```

```
#Using library httr
library(httr)
```

```
## Warning: package 'httr' was built under R version 4.4.2
```

```
##
## Attaching package: 'httr'
```

```
## The following object is masked from 'package:NLP':
##
##      content
```

```
eg7<-GET("https://www.edureka.co/blog/what-is-data-science/")
doc<-htmlParse(eg7)
doc.text<-unlist(xpathApply(doc,'//p',xmlValue))
```

```
#Using library rvest
library(rvest)
```

```
## Warning: package 'rvest' was built under R version 4.4.2
```

```
eg8<-read_html("https://www.edureka.co/blog/what-is-data-science/")
nodesh<-html_nodes(eg8,'.col-lg-9 :nth-child(1)')
nodes<-html_nodes(eg8,'.color-4a div span , .btn-become-profesional-link+ p')
texts<-html_text(nodes)
#Selecting multiple pages
pages<-paste0('https://www.amazon.co.jp/s?k=skincare&crd=28HIW1TYLV9UM&sprefix=skincare%2Caps%2C268&ref=nb_sb_noss_2&page=')
pages<-paste0('https://www.amazon.com/s?k=condominium&crd=1ZNT0FW4JKRC3&sprefix=condominiu%2Caps%2C312&ref=nb_sb_noss_2&page=')
# Epage=
# read the first 3 pages
eg10<-read_html(pages[1])
nodes<-html_nodes(eg10,'.a-price-whole')
texts<-html_text(nodes)
Price<-function(page){
  url<-read_html(page)
  nodes<-html_nodes(url ,'.a-price-whole ')
  html_text(nodes)}
sapply(pages,Price)
```

```
## $`https://www.amazon.com/s?k=condominium&crd=1ZNT0FW4JKRC3&sprefix=condominiu%2Caps%2C312&ref=nb_sb_noss_2&page=1`
## [1] "0." "15." "12." "26." "23." "24." "23." "16." "15." "98."
## [11] "14." "9." "0." "5." "18." "16." "16." "9." "9." "21."
## [21] "69." "7." "0." "9." "29." "49." "7." "19." "0." "179."
## [31] "177." "0." "0." "0." "0."
##
## $`https://www.amazon.com/s?k=condominium&crd=1ZNT0FW4JKRC3&sprefix=condominiu%2Caps%2C312&ref=nb_sb_noss_2&page=2`
## [1] "0." "15." "12." "26." "23." "16." "15." "24." "23." "7." "0." "0."
## [13] "5." "9." "16." "9." "18." "16." "9." "21." "69." "98." "14." "9."
## [25] "7." "19." "0." "7." "0." "18." "10." "0." "0." "0." "0."
##
## $`https://www.amazon.com/s?k=condominium&crd=1ZNT0FW4JKRC3&sprefix=condominiu%2Caps%2C312&ref=nb_sb_noss_2&page=3`
## [1] "0." "0." "0." "9." "29." "49." "9." "0." "12."
```

```
do.call("c",lapply(pages,Price))
```

```
## [1] "0." "15." "12." "26." "23." "16." "15." "24." "23." "7."
## [11] "0." "0." "5." "9." "16." "9." "18." "16." "9." "21."
## [21] "69." "98." "14." "9." "7." "19." "0." "7." "0." "18."
## [31] "10." "0." "0." "0." "0." "0." "15." "12." "26." "23."
## [41] "24." "23." "16." "15." "98." "14." "9." "0." "5." "18."
## [51] "16." "16." "9." "9." "21." "69." "7." "0." "9." "29."
## [61] "49." "7." "19." "0." "179." "177." "0." "0." "0." "0."
## [71] "0." "0." "0." "18." "10." "29." "49." "9." "0." "9."
```