

IBM WebSphere Information Integrator
OmniFind Edition



文字分析整合

8.3 版

IBM WebSphere Information Integrator
OmniFind Edition



文字分析整合

8.3 版

使用此資訊及其支援的產品之前，請先閱讀「注意事項」中的一般資訊。

本文包含 IBM 的所有權資訊。乃依據授權合約提供並受著作權法保護。本書中的資訊不包括任何產品保證，且其陳述也不得延伸解釋。

您可以線上訂購 IBM 出版品，或可以透過當地的 IBM 業務代表來訂購：

- 若要線上訂購出版品，請造訪「IBM 出版品中心 (IBM Publication Center)」：www.ibm.com/shop/publications/order。
- 若要尋找當地的 IBM 業務代表，請造訪「IBM 全球聯絡站名錄 (IBM Directory of Worldwide Contacts)」：www.ibm.com/planetwide。

當您傳送資訊給 IBM 時，即授權予 IBM，IBM 得以其認為適當的方式來使用或分送資訊，而無需對您負任何責任。

© Copyright International Business Machines Corporation 2004, 2005. All rights reserved.

目錄

關於這些主題	v	建立停用字的 XML 檔案	51
這些主題的適用對象	v	建立停用字定義檔	52
語意搜尋的語言支援	1	自訂 Boost 字定義檔	55
 		建立 Boost 字的 XML 檔	55
自訂文字分析整合	3	建立 Boost 字定義檔	57
非結構化資訊管理架構 (UIMA) 概觀	3	企業搜尋中包括的文字分析	59
自訂分析整合的工作流程	4	語言識別	59
安裝及執行企業搜尋基本註解程式	6	非定義檔型斷詞法的語言支援	60
文字分析演算法	7	定義檔型斷詞法的語言支援	60
類型系統說明	8	日文的斷詞	62
分析及搜尋中的 XML 標記	11	日文的正體字變體	62
建立 XML 對 UIMA 類型對映配置檔	12	停用字移除	62
文字分析結果	16	字元正常化	63
特性路徑	17	 	
內建特性	18	企業搜尋文件	65
過濾器	20	 	
自訂分析結果的索引對映	21	WebSphere II OmniFind Edition 協助	
建立索引建置配置檔	22	工具	67
所選分析結果的資料庫對映	27	 	
在資料庫中儲存分析結果	28	企業搜尋的詞彙名詞解釋	69
建立 XML 對映配置檔	28	 	
儲存區類型對映	32	有關 WebSphere Information	
擷取文件中符合語意搜尋查詢的部份	36	Integration 的存取資訊	79
定義於企業搜尋的類型及特性	38	 	
定義於 UIMA 的類型及特性	41	提供文件的相關意見	81
語意搜尋應用程式	43	 	
語意搜尋查詢字詞	44	聯絡 IBM	83
搜尋應用程式中的同義字支援	47	商標	85
建立同義字的 XML 檔案	47	 	
建立同義字定義檔	48	索引	89
自訂停用字定義檔	51		

關於這些主題

請使用此資訊，在 IBM® WebSphere® Information Integrator OmniFind™ Edition 8.3 版系統中建立及部署語意搜尋解決方案。語意搜尋可讓您在搜尋查詢中搜尋較高層次的概念並指出關係，這些資訊可以利用文字分析來偵測。

WebSphere Information Integrator OmniFind Edition (WebSphere II OmniFind Edition) 提供一項技術，稱為企業搜尋。安裝 WebSphere II OmniFind Edition 產品時，會一併安裝企業搜尋的元件。在 WebSphere II OmniFind Edition 文件中會使用企業搜尋一詞，但在提及安裝路徑和產品包裝標籤時除外。

企業搜尋的文字分析文件包含下列主題：

- 企業中的語言支援簡介
- 如何在企業搜尋中整合自訂文字分析的相關指示
- 如何對映 XML 文件結構的相關指示
- 如何新增選取的分析結果至 JDBC 表格的相關指示
- 如何新增分析結果至企業搜尋索引以啓用語意搜尋的相關指示
- 如何在搜尋期間併入同義字、停用字及 Boost 字定義檔的相關指示
- 在文件處理程序期間自動執行哪些文字分析功能的概觀

這些主題的適用對象

此資訊主要適用於系統管理員及搜尋應用程式開發人員，因為他們負責在企業搜尋中建立及部署語意搜尋解決方案。

企業搜尋提供純文字文件的語意搜尋支援。文件會經過分析，然後會儲存結果，而在搜尋期間會存取這些結果。合併文字分析與可以檢索單字和文字跨距的企業搜尋，可啓用語意搜尋。搜尋時的語意支援目的在於增進文件搜尋結果，以便產生符合查詢的最佳可能文件集合。

使用此資訊可以瞭解如何在企業搜尋中整合自訂文字分析，以及如何使用同義字、停用字及 Boost 字定義檔來增進查詢結果。另請使用此資訊來瞭解在文件處理程序期間，企業搜尋全程提供的基本語言支援。

若要更有效地使用本資訊，您必須熟悉 Web 應用程式，並對要搜尋的資料來源有相關使用經驗。

語意搜尋的語言支援

企業搜尋針對大部份的印歐語言及亞洲語言 (如日文) 的純文字文件，提供語言學搜尋支援。

搜尋時的語言支援目的在於增進文件搜尋結果，以便產生符合查詢的最佳可能文件集合。

語言處理的執行分成兩個階段：處理要加入索引的文件時，以及當使用者在搜尋期間輸入查詢時。

企業搜尋只含有精細的語言功能，這是在判定輸入文件的語言及將文件輸入串流分段成單字或記號時的必要功能。

如果您知道搜尋主要限於基本關鍵字搜尋或使用文件結構的原生 XML 搜尋，則企業搜尋所含的語言處理會適當地涵蓋搜尋時需要的項目。

然而，如果您要在文件中搜尋純單字以外更特定的項目，則只有語言處理不一定能滿足需求，如下面範例所示：

- 在協同作業情況下，資訊不一定會明確地標示出來，如電子郵件中的位址或電話號碼。並未使用 **電話號碼** 一詞。而是在電子郵件中出現「您可以撥打 555-641-1805 來聯絡我」之類的詞組。
- 在競爭情報中，文件提到競爭者及其提供的貨物，或競爭者的網站在過去三個月從某一產品組移至另一產品組。
- 在客戶關係管理中，文件可能會提到舊金山區修車廠的汽車煞車問題。修車廠報告說明「由於油壓洩露而調整煞車皮」之類的狀況；此外，這些報告只提到修車廠的街名，而沒提到完整地址。
- 在研究中，文件說明特定蛋白質及它和同一段落中所提到的至少一種疾病之間的關係。在文獻中，該蛋白質至少有 20 種以上不同的名稱，而文件通常完全不會提到疾病這個字，只會提到疾病本身的名稱。

在這些範例中，在現今存在的廣大資源來源集中搜尋您需要的項目，這是一項新挑戰，需要比企業搜尋所提供的斷詞法層次及定義檔型分析更精確的分析。大部份需要的資訊都不會在原始文件中明確地標示或標記。而是必須分析資訊，才能辨識和尋只需要的概念，例如，人員、組織、位置、機能及產品之類的具名實體，以及這些實體之間的可能關係。

IBM Unstructured Information Management Architecture (UIMA) 是一種軟體架構，可協助您在企業搜尋中建置進階分析功能，以便在文件集中偵測及尋找需要的資訊。

相關概念

第 3 頁的『自訂文字分析整合』

「非結構化資訊管理架構 (UIMA)」是一種軟體架構，支援建立、探查、編製及部署文字分析功能。使用 UIMA，您可以建置自訂文字分析。

第 3 頁的『非結構化資訊管理架構 (UIMA) 概觀』

非結構化資訊管理架構 (UIMA) 是一種軟體架構，可協助您建置進階分析功能，以便在文件集中尋找特定資訊。

自訂文字分析整合

「非結構化資訊管理架構 (UIMA)」是一種軟體架構，支援建立、探查、編製及部署文字分析功能。使用 UIMA，您可以建置自訂文字分析。

UIMA 是一種開放式平台，可以針對每一種不同概念的分析功能來識別元件，並確保這些元件可以輕易地重覆使用並彼此結合。

UIMA 的中心概念是分析引擎，負責探查及代表純文字文件中的分析內容。分析邏輯元件稱為註解程式。註解程式的重點在於分析作業，且與任何其他處理程序無關。分析引擎可含有單一註解程式，或可以是許多引擎的組合，而每一個引擎又都含有註解程式。

由分析引擎所產生的推斷資訊稱為分析結果。理論上，分析結果會對應於您要搜尋的資訊。

進階語言分析包含許多不同分析作業的組合。分析是從 (舉例來說) 語言偵測和斷詞法開始，然後繼續詞性識別，接著執行深度文法剖析。最後一個步驟包括識別 (舉例來說) 特定化學物質和特定徵兆外觀之間的關係。分析程序中的每一個步驟都是後序步驟的必要作業。

UIMA 提供基本構成要素，供您建立、測試及部署自己的分析引擎。它不會以預先配置的分析引擎形式，提供可讓您在 UIMA 環境中部署的任何語言分析功能。

「UIMA 軟體開發套件」包含 UIMA 架構的 Java™ 實作，以供 UIMA 元件的實作、說明、組合及部署使用。它還提供了 Eclipse 型開發環境，其中包含一組用於 UIMA 的工具和公用程式。如需 Eclipse 的相關資訊，請參閱 www.eclipse.org。

若要使用 UIMA，您必須安裝「UIMA 軟體開發套件」。開發工具箱可在 IBM developerWorks® 中取得。請造訪 WebSphere Information Integrator 區域，以取得 <http://www.ibm.com/developerworks/db2/zones/db2ii/> 上的資訊。如需如何在「Eclipse 互動式開發環境」中安裝「UIMA 軟體開發套件」的相關指示，請參閱 UIMA 文件。

相關概念

第 1 頁的『語意搜尋的語言支援』

企業搜尋針對大部份的印歐語言及亞洲語言 (如日文) 的純文字文件，提供語言學搜尋支援。

『非結構化資訊管理架構 (UIMA) 概觀』

非結構化資訊管理架構 (UIMA) 是一種軟體架構，可協助您建置進階分析功能，以便在文件集中尋找特定資訊。

非結構化資訊管理架構 (UIMA) 概觀

非結構化資訊管理架構 (UIMA) 是一種軟體架構，可協助您建置進階分析功能，以便在文件集中尋找特定資訊。

特性結構是代表分析結果的基礎資料結構。特性結構是屬性值結構。每一個特性結構都屬於某一類型，而每一個類型都有一組指定的有效特性或屬性 (內容)，非常類似 Java 類別。特性含有範圍類型，指出特性必須具備的值類型，如 String。

大部分的分析演算法 (也稱為註解程式) 會以註解形式產生分析結果。註解是一種特殊類型的特性結構，主要用於語言分析處理。特性結構會跨越或包含輸入文字的片段，且是以輸入文字開頭及結束位置的詞彙定義。

例如，辨識貨幣表示式的註解程式針對文字 "100.55 US Dollars" 建立了類型 monetaryExpression 的註解，並將特性 currencySymbol 設為 "\$" 以替代該文字。

UIMA 中的所有註解程式都使用特性結構來儲存或讀取資訊，換句話說，所有資料的模型都設為特性結構。

類型系統以類型和特性來定義所有可能的特性結構，非常類似 Java 的類別階層。

所有特性結構都會顯示在稱為共用分析結構的中央資料結構中。您可以利用共用分析結構來處理所有資料交換。

共用分析結構包含下列物件：

- 純文字文件
- 類型系統說明，指出類型、次類型及其特性
- 分析結果，說明文件或文件的區域
- 索引儲存庫，支援分析結果的存取和疊代

相關概念

第 1 頁的『語意搜尋的語言支援』

企業搜尋針對大部份的印歐語言及亞洲語言 (如日文) 的純文字文件，提供語言學搜尋支援。

第 3 頁的『自訂文字分析整合』

「非結構化資訊管理架構 (UIMA)」是一種軟體架構，支援建立、探查、編製及部署文字分析功能。使用 UIMA，您可以建置自訂文字分析。

自訂分析整合的工作流程

您可以使用「UIMA 軟體開發套件」來建立和測試自訂文字分析演算法，然後在企業搜尋的文件集合上部署及執行。

若要開發分析演算法並在企業搜尋中執行：

1. 計畫及設計
 - a. 決定您要搜尋的資訊。您要擷取哪些文件？在特定搜尋作業中需要哪些概念及關係？例如，若要在製藥公司的內部網站上加強一般用途的搜尋，可能需要產品和員工名稱，而在研究和開發區的人員必須使用藥品名稱的變體，並瞭解藥品-原因-療法關係。
 - b. 指定在您要搜尋的文件中擷取資訊所需的文字分析類型。
 - c. 如果集合含有 XML 文件，請決定是否要在解決方案中利用 XML 標記。在企業搜尋中，您可以使用下列兩種方式來利用 XML 標記：

- 如果您可以在自訂分析中使用 XML 標記 (例如，文件含有有助於彙總或分類註解程式的 <summary> 或 <topic> 元素)，請定義 XML 對共用分析結構對映。
 - 如果要在查詢中依照文件的顯示方式來使用 XML 標記，請啓用原生 XML 對映。
- d. 決定您要利用語意搜尋存取哪些儲存在共用分析結構的文字分析結果資訊。定義共用分析結構對索引對映。
 - e. 決定是否要在關聯式資料庫儲存分析結果，例如，利用報告及資料採礦應用程式來探查趨勢及關聯。定義共用分析結構對 JDBC 表格對映。
 - f. 設計語意搜尋應用程式。決定搜尋使用者如何使用語意搜尋的其他功能。設計使用者介面。
2. 開發：「UIMA 軟體開發套件」活動
 - a. 定義個別的分析步驟。
 - b. 說明對映及分析演算法的類型系統。
 - c. 利用「UIMA 軟體開發套件」，開發每一個分析步驟的分析演算法 (註解程式) 並在分析引擎中嵌入註解程式。利用企業搜尋基本註解程式資料包中的基本功能 (語言識別及分段)，建置任何自訂分析。
 - d. 在 UIMA 中測試分析演算法之後，將分析引擎封裝為 PEAR 檔 (處理程序引擎保存檔)。保存檔只能含有您的分析演算法，而不是基本企業搜尋語言功能。
 3. 部署：企業搜尋活動
 - a. 將分析引擎保存檔 (.pear) 上傳到企業搜尋。提供分析元件的名稱，以便能在企業搜尋中參照它。
 - b. 關聯一或多個文件集合與您的分析引擎。
 - c. 如果適用，針對每一個集合，上傳及選取您為自訂分析所定義的 XML 元素對 UIMA 類型對映配置。
 - d. 如果適用，針對每一個集合，上傳及選取您為自訂分析所定義的資料庫對映配置。
 - e. 如果適用，針對每一個集合，上傳及選取您為語意搜尋所定義的索引對映配置。
 - f. 必要的話，請設定自訂語意搜尋應用程式，例如，將瀏覽器型搜尋使用者介面部署到應用程式伺服器。
 - g. 搜索、剖析及檢索語意搜尋集合中的文件，就像您在關鍵字型集合中所做的一樣。

相關工作

第 6 頁的『安裝及執行企業搜尋基本註解程式』

您可以使用企業搜尋基本註解程式資料包，開發以企業搜尋註解程式輸出為基礎的新註解程式，並在「UIMA 軟體開發套件 (SDK)」中測試自訂註解程式。

安裝及執行企業搜尋基本註解程式

您可以使用企業搜尋基本註解程式資料包，開發以企業搜尋註解程式輸出為基礎的新註解程式，並在「UIMA 軟體開發套件 (SDK)」中測試自訂註解程式。

一組基本註解程式包括：

- **語言 ID 註解程式**

偵測文件的語言。如需功能及配置參數，請參閱描述子檔案 `jlangid.xml`。

- **FROST 定義檔查閱註解程式**

依據 IBM LanguageWare 定義檔，提供分段及句子偵測。若為記號，則會產生其他語言資訊，例如，基礎詞形或詞形。如需功能及配置參數，請參閱描述子檔案 `jfrost.xml`。

- **空格記號器**

在所有歐洲語言文件上，執行以空格為基礎的分段或其他空格區隔 Script。此外，註解程式也可以對下列文字 Script 執行 n-gram 分段化：阿拉伯文、漢語、希伯來文、平假名、Katakana、寮文、蒙古文、泰文、YI 及 Hangul。如需功能及配置參數，請參閱描述子檔案 `jtok.xml`。

若要在 UIMA 中執行這些註解程式，您必須安裝「UIMA 軟體開發套件 (SDK)」。若要下載，請造訪 [IBM developerWorks](http://www-128.ibm.com/developerworks/db2/zones/db2ii/) 網站，網址為 <http://www-128.ibm.com/developerworks/db2/zones/db2ii/>。

企業搜尋基本註解程式資料包是一個壓縮檔，其中包含用於企業搜尋的文字分析註解程式。在企業搜尋中剖析文件時，這些註解程式一律任何自訂分析之前執行。

若要安裝註解程式資料包：

1. 在企業搜尋 (WebSphere Information Integrator OmniFind Edition) 安裝的 `ES_INSTALL_ROOT/packages/uima` 目錄中，註解程式資料包 `OF_base_annotators.zip`。
2. 將壓縮檔複製到 UIMA SDK 安裝的根目錄。
3. 解壓縮 zip 檔，將企業搜尋基本註解程式檔案加入 UIMA SDK 安裝的指定目錄結構。

安裝基本註解程式資料包後，註解程式描述子檔案會位於資料夾 `UIMA_SDK_INSTALL/docs/examples/descriptors/analysis_engine` 中。檔案 `of_tokenization.xml` 會依照企業搜尋的使用順序來列出基本註解程式。

描述子檔案所含的配置值與企業搜尋中使用的值相同。您可以在 UIMA SDK 中變更值以進行除錯。然而，請勿變更企業搜尋系統中的描述子檔案。變更這些檔案可能會造成系統不穩定或效能問題。

企業搜尋基本註解程式資料包只含有處理英文文件所需的定義檔。如果您要在開發環境中處理其他語言，請請遵循下列步驟：

1. 在企業搜尋安裝的 `ES_INSTALL_ROOT/configurations/parserservice/jediidata/frost/resources` 中，尋找企業搜尋定義檔。
2. 將定義檔的內容複製到本端 `UIMA SDK` 安裝的 `UIMA_SDK_INSTALL/data/frost/resources`。

若要驗證註解程式資料包是否已順利安裝：

1. 在下列目錄中開啓「共用分析結構 (CAS) 視覺除錯器 (CVD)」：
`UIMA_SDK_INSTALL/bin/cvd[.bat/.sh]`。
2. 按一下執行 → 載入 TAE。
3. 在 `UIMA_SDK_INSTALL/docs/examples/descriptors/analysis_engine` 目錄中選取文字分析引擎指定元檔案 `of_tokenization.xml`。
4. 載入範例文件，並執行文字分析引擎。您會在 CVD 中看到類型 `uima.tt.TokenAnnotation` 的註解。

若要使用企業搜尋註解程式以進行處理程序：

1. 如果自訂註解程式使用企業搜尋註解程式所定義的類型，請自訂註解程式指定元的 `typeSystem` 區段中併入檔案檔案 `of_typesystem.xml` 的參照。 `of_typesystem.xml` 檔案位於 `UIMA_SDK_INSTALL/docs/examples/descriptors/analysis_engine` 目錄中。請參閱 `analysis_engine` 目錄中的 `jtok.xml`，以取得如何併入描述子檔案參照的範例。

相關參考

2
2

第 38 頁的『定義於企業搜尋的類型及特性』

定義於企業搜尋的類型系統包含文件中間資料處理及基本語言分析。

文字分析演算法

「UIMA 軟體開發套件」包含 API 及工具，您可以用來建立註解程式 (分析演算法，包含類型系統說明) 並在分析引擎中嵌入這些註解程式。

UIMA 文件含有教學指導樣式手冊，可協助您建置這些元件。「軟體開發套件」提供測試及檢視結果的公用程式，以及檢索分析結果的小型語意搜尋引擎。您也可以對儲存在索引中的資訊執行更多進階搜尋。

「UIMA 軟體開發套件」不提供任何預先配置的分析引擎。然而，您可以在 UIMA 環境中使用企業搜尋提供的基本註解程式。在 UIMA 環境中進行開發時，如需如何在文字分析演算法之前併入語言偵測及分段功能的相關資訊，請參閱 UIMA 文件。

使用「UIMA 軟體開發套件」完成分析引擎的開發和測試之後，如果想要在企業搜尋中對文件集合執行這些演算法，則必須建立 PEAR (處理程序引擎保存) 檔。這個保存檔含有在企業搜尋中部署自訂分析功能作為分析引擎所需的全部資源。如需建立保存檔所需的所有處理程序步驟說明，請參閱「軟體開發套件」提供的 UIMA 文件。

保存檔只能含有您的自訂分析，即使它是以企業搜尋所提供的基本語言功能為基礎。基本企業搜尋分析步驟一律在任何自訂分析之前執行。

若要瞭解如何在企業搜尋中配置及部署語意搜尋解決方案，請執行 <http://www.ibm.com/developerworks/db2/zones/db2ii/> 所述的教學指導。教學指導會引導您完成在企業搜尋中部署自訂文字分析演算法的必要步驟，並顯示如何在查詢中使用分析結果來增進搜尋結果。

相關工作

第 6 頁的『安裝及執行企業搜尋基本註解程式』

您可以使用企業搜尋基本註解程式資料包，開發以企業搜尋註解程式輸出為基礎的新註解程式，並在「UIMA 軟體開發套件 (SDK)」中測試自訂註解程式。

類型系統說明

類型系統說明解釋在自訂分析中所使用的特性結構 (代表分析結果的基礎資料結構)。

定義於類型系統說明的相同類型必須由含有註解程式 (分析演算法) 的分析引擎使用，而在與自訂分析相關的所有對映檔中，XML 對映配置檔、索引建置配置檔或具 JDBC 功能的資料庫配置對映檔也是如此。

註解程式的類型系統說明可以是註解程式描述子的一部份，或可以內含於不同的類型系統描述子檔案中。有時，它是相同分析引擎所含的另一個註解程式的描述子。

類型系統說明必須是分析引擎保存檔 (.pear 檔) 的一部份，該檔案是從 UIMA 環境匯入企業搜尋。

接下來相同的範例類型系統說明會在所有主題中使用，這些主題討論使用自訂分析時所能選取的不同對映類型。

下列類型系統說明範例說明治安報告，其中包含嫌犯、犯案地點、犯案時間及日期等相關資訊：

```
<?xml version="1.0" encoding="UTF-8"?>
<typeSystemDescription>
  <name>Police Reports Type System</name>
  <description>Type system description for
    police reports</description>
  <version>1.0</version>
  <types>
    <typeDescription>
      <name>com.ibm.omnifind.types.PoliceReport</name>
      <description>Annotates a police report</description>
      <superTypeName>uima.tcas.Annotation</superTypeName>
      <features>
        <featureDescription>
          <name>time</name>
          <description>Time the crime was reported to have happened
            </description>
          <rangeTypeName>com.ibm.omnifind.types.Time</rangeTypeName>
        </featureDescription>
        <featureDescription>
          <name>date</name>
          <description>When the crime happened</description>
          <rangeTypeName>com.ibm.omnifind.types.Date</rangeTypeName>
        </featureDescription>
        <featureDescription>
          <name>location</name>
          <description>Where the crime took place</description>
          <rangeTypeName>com.ibm.omnifind.types.City</rangeTypeName>
        </featureDescription>
        <featureDescription>
          <name>knownSuspects</name>
          <description>Contains annotations of type Suspect</description>
          <rangeTypeName>uima.cas.FSArray</rangeTypeName>
        </featureDescription>
        <featureDescription>
          <name>crimeDescription</name>
          <description>Short description of the crime</description>
          <rangeTypeName>uima.cas.String</rangeTypeName>
        </featureDescription>
      </features>
    </typeDescription>
    <typeDescription>
      <name>com.ibm.omnifind.types.City</name>
      <description>The name of a city</description>
```



```

<superTypeName>uima.tcas.Annotation</superTypeName>
<features>
  <featureDescription>
    <name>cityName</name>
    <description>The name of the city</description>
    <rangeTypeName>uima.cas.String</rangeTypeName>
  </featureDescription>
  <featureDescription>
    <name>cityDistrict</name>
    <description>The name of the district</description>
    <rangeTypeName>uima.cas.String</rangeTypeName>
  </featureDescription>
</features>
</typeDescription>
<typeDescription>
<name>com.ibm.omnifind.types.Person</name>
<description>A person annotation</description>
<superTypeName>uima.tcas.Annotation</superTypeName>
<features>
  <featureDescription>
    <name>role</name>
    <description>For example, suspect or witness</description>
    <rangeTypeName>uima.cas.String</rangeTypeName>
  </featureDescription>
  <featureDescription>
    <name>firstName</name>
    <description>The first name of the person</description>
    <rangeTypeName>uima.cas.String</rangeTypeName>
  </featureDescription>
  <featureDescription>
    <name>surName</name>
    <description>The surname of the person</description>
    <rangeTypeName>uima.cas.String</rangeTypeName>
  </featureDescription>
  <featureDescription>
    <name>title</name>
    <description>For example, Mr. or Ms.</description>
    <rangeTypeName>uima.cas.String</rangeTypeName>
  </featureDescription>
  <featureDescription>
    <name>gender</name>
    <description>Male or female</description>
    <rangeTypeName>uima.cas.String</rangeTypeName>
  </featureDescription>
</features>
</typeDescription>
<typeDescription>
<name>com.ibm.omnifind.types.Suspect</name>
<description>A found suspect</description>
<superTypeName>com.ibm.omnifind.types.Person</superTypeName>
<features>
  <featureDescription>
    <name>description</name>
    <description>Suspect description,
    for example, bearded with dark glasses</description>
    <rangeTypeName>uima.cas.String</rangeTypeName>
  </featureDescription>
</features>
</typeDescription>
<typeDescription>
<name>com.ibm.omnifind.types.Date</name>
<description>A date</description>
<superTypeName>uima.tcas.Annotation</superTypeName>
<features>
  <featureDescription>
    <name>year</name>
    <description>The year, for example, 2005</description>

```

```

        <rangeTypeName>uima.cas.Integer</rangeTypeName>
    </featureDescription>
    <featureDescription>
        <name>month</name>
        <description>The month in digits, for example, 7</description>
        <rangeTypeName>uima.cas.Integer</rangeTypeName>
    </featureDescription>
    <featureDescription>
        <name>day</name>
        <description>The day in digits</description>
        <rangeTypeName>uima.cas.Integer</rangeTypeName>
    </featureDescription>
    <featureDescription>
        <name>dayOfWeek</name>
        <description>The day of the week, for example, Monday</description>
        <rangeTypeName>uima.cas.String</rangeTypeName>
    </featureDescription>
    <featureDescription>
        <name>quarter</name>
        <description>The quarter, for example, Q1-2005</description>
        <rangeTypeName>uima.cas.String</rangeTypeName>
    </featureDescription>
    <featureDescription>
        <name>englDate</name>
        <description>Date as mm/dd/yyyy</description>
        <rangeTypeName>uima.cas.String</rangeTypeName>
    </featureDescription>
</features>
</typeDescription>
<typeDescription>
    <name>com.ibm.omnifind.types.Time</name>
    <description>A time</description>
    <superTypeName>uima.tcas.Annotation</superTypeName>
    <features>
        <featureDescription>
            <name>hours</name>
            <description>Hours from 00-23</description>
            <rangeTypeName>uima.cas.Integer</rangeTypeName>
        </featureDescription>
        <featureDescription>
            <name>minutes</name>
            <description>Minutes in the hour</description>
            <rangeTypeName>uima.cas.Integer</rangeTypeName>
        </featureDescription>
        <featureDescription>
            <name>timeOfDay</name>
            <description>Time periods, such as morning, noon</description>
            <rangeTypeName>uima.cas.String</rangeTypeName>
        </featureDescription>
    </features>
</typeDescription>
</types>
</typeSystemDescription>

```

相關參考

- 2 第 38 頁的『定義於企業搜尋的類型及特性』
- 2 定義於企業搜尋的類型系統包含文件中間資料處理及基本語言分析。
- 2 第 41 頁的『定義於 UIMA 的類型及特性』
- 2 「UIMA 軟體開發套件」定義在文字分析時可以在文件中探查到的基本語言類型及特性。

分析及搜尋中的 XML 標記

您可以將文件中 XML 結構的資訊直接對映到共用分析結構，而不必撰寫 UIMA 註解程式。

如果集合中的文件是 XML，且您要在文字分析或語意搜尋時利用 XML 標記，則可以選擇下列選項：

原生 XML 搜尋

如果要在語意搜尋期間依照文件中的顯示原樣來使用所有 XML 標示及屬性，則請使用此選項。比方說，如果您的帳單文件含有 `<addressee>` 元素，則啓用原生 XML 搜尋可讓您在語意搜尋查詢中使用這個標示，以在這個元素中搜尋特定的客戶名稱。

使用此選項，則會在共用分析結構中使用 `com.ibm.es.tt MarkupTag` 類型來顯示文件的 XML 結構。針對每一個 XML 標示，都會建立此類型的註解。這個註解含有標示名稱、其屬性及屬性內容。此資訊一定會加以檢索，且可以存取它來執行語意搜尋。

原生 XML 搜尋不需要對映配置檔。您可以從企業搜尋的管理主控台啓用原生 XML 搜尋。

XML 元素對 UIMA 類型對映

在下列情況下，請使用此選項：

- 某些 XML 元素的語意明確，且可以用於進一步的文字分析步驟。這些分析步驟可以直接在從 XML 結構所建立的註解和特性上運作，且會和潛在的不同原始文件格式隔離。例如，帳單文件中的元素 `<addressee>` 通常含有客戶名稱。使用 XML 元素對類型對映，這個元素的內容可以直接對映到類型 `Customer` 的註解。然後，註解程式可以利用 `Customer` 註解周遭的資訊，推斷「客戶所在地」關係。
- 您想要將客戶註解程式的處理範圍限制於 XML 輸入中的特定指定區域。例如，您可能想要在偵測汽車問題的註解程式中，限制 `<technicianComment>` 標示的內容。
- 您想要將文字分析處理程序和後續的搜尋限制於 XML 文件的特定部份，然後過濾出無關或非文字內容。
- 您想要將含有不同名稱 (例如，`<mainHeading>` 或 `<doc>`) 的 XML 標示對映到要在語意搜尋中使用的一般跨距 (例如，標題)。

在這些情況下，您必須建立 XML 對 UIMA 類型對映配置檔，以定義特性結構。您在配置檔中定義的特性結構會在剖析文件時建立，並由自訂註解程式存取。

您可以在文件集中使用多個配置檔。哪一種配置用於哪一個 XML 文件，是由 `<identifier>` 元素決定。配置檔中的 `<identifier>` 元素必須符合 XML 文件中的根元素。比方說，如果文件的根元素是 `doc`，則配置檔中 `<identifier>` 元素的值必須也是 `"doc"`。

如果找不到相符的項目，則程式會搜尋 `<identifier>` 元素設為預設值的配置檔。如果找不到預設的配置，則文件的文字區段 (沒有標示資訊) 會對映到共用分析結構中的文件註解。

如果要擷取只在文件相關部份才有的資訊，並忽略無關的部份，只要指定文件中哪些 XML 元素含有相關資訊即可。這稱為內容擷取。例如，您可以擷取在標題和主體元素中指定的輸入，並忽略作者、日期、ID 和發佈者中的輸入。

內容擷取可以增進下列 XML 文件類型的分析處理程序：

- 含有大量內容且不適用於分析的文件，例如，二進位附件。使用內容擷取會大幅降低文件大小，加速處理程序並避免因不當資料而產生的分析錯誤。
- 文件中的文件文字散佈著無關的文字，例如，在 <note> 標示中含有編輯資訊的文件。分析文件內容時，忽略此資訊可以有更好的結果。

原生 XML 搜尋和 XML 元素對 UIMA 類型對映中的內容擷取選項彼此矛盾，因為要考慮所有內容或只能考慮指定的內容。如果指定內容擷取，則會忽略原生 XML 對映。沒有內容擷取，您可以同時使用 XML 元素對 UIMA 類型對映和原生 XML 搜尋。

在配置檔中使用的所有類型和特性，都必須在自訂分析步驟的類型系統說明中加以說明。您可以使用元件描述子編輯器 Eclipse 外掛程式，在 UIMA 環境中建立類型系統描述子。這個外掛程式可讓您建立描述子檔案，而無需瞭解必要的 XML 語法。

完成建置及測試自訂分析後，請使用 UIMA PEAR (處理程序引擎保存) 產生精靈，建立含有自訂分析檔案 (包含類型系統說明) 的保存檔。

然後，您可以使用企業搜尋的管理主控台，上傳自訂分析及 XML 元素對 UIMA 類型對映配置檔到企業搜尋。

相關工作

『建立 XML 對 UIMA 類型對映配置檔』

在 XML 對 UIMA 類型對映配置檔中，您可以使用所有配置選項，將 XML 對映到 UIMA 資料類型。

建立 XML 對 UIMA 類型對映配置檔

在 XML 對 UIMA 類型對映配置檔中，您可以使用所有配置選項，將 XML 對映到 UIMA 資料類型。

關於本作業

XML 對 UIMA 類型對映配置檔必須符合下面範例所顯示的綱目。

範例 XML 治安報告含有用於犯罪類型、犯罪日期、犯罪地點、報告員警、該員警任職的警方轄區、嫌犯說明及摘要的 XML 標示。後面接著主體區段。例如：

```
<report>
  <doc>
    <crimeType>Car theft</crimeType>
    <crimeDate>04/23/05 09:23 pm</crimeDate>
    <crimeLocation>27 Main Street, Brynston, Springfield, New Jersey</crimeLocation>
    <reportingOfficer rank="Lt">Jakob
      <lastname>Collins</lastname>
    </reportingOfficer>
    <policePrecinct>14th Precinct</policePrecinct>
    <suspectDescription>Male, dark haired, dark glasses,
      blue jeans with dark, probably black,
      jacket</suspectDescription>
    <abstract>A Mercedes CLK was stolen on 04/23/2005 from a parking
      lot in front of the Blue Lagoon restaurant on
      27 Main Street, Brynston.(serial number: 32 2761 50871)</abstract>
```

```
<body>A Mercedes CLK was stolen on 04/23/2004 from a parking
lot in front of the Blue Lagoon restaurant on 27 Main Street,
Brynston.(serial number: 32 2761 50871)
```

It has a black color and wide Michelin tires.

```
Eyewitnesses in front of the restaurant saw two darkly dressed
males drive away in the car at high speed. The car was
found abandoned on Aliway Ave in Brooklyn. The fuel tank was empty.
The seats were badly stained and the back seat was vandalized.
Nothing was stolen out of the car....</body>
</doc>
<image>
  </-- image of the crime scene as a base64-encoded string -->
</image>
</report>
```

依據範例報告，配置檔可以有如下結構。範例使用為治安報告實務範例所定義的類型系統。

```
<?xml version="1.0"?>
<xmlCasInitializerConfiguration
  xmlns="http://www.ibm.com/2005/uma/jedii_ci_xml">

  <identifier>Default</identifier>
  <description>Sample configuration</description>

  <contentElements>
    <element>/report/doc</element>
  </contentElements>

  <elementToTypeMappings>
    <elementToTypeMapping>
      <element>//doc//reportingOfficer</element>
      <type>com.ibm.omnifind.types.Person</type>
      <featureValueAssignment>
        <feature>role</feature>
        <basicValue default="Reporting officer">
          </basicValue>
        </featureValueAssignment>
        <featureValueAssignment>
          <feature>gender</feature>
          <basicValue default="male"
            useAttributeValue="sex">
          </featureValueAssignment>
        <featureValueAssignment>
          <feature>surName</feature>
          <values concatenate="true" delimiter="">
            <basicValue useAttributeValue="rank"
              default="Lt">
            <basicValue useElementContent="lastName">
          </values>
          </featureValueAssignment>
        </elementToTypeMapping>
      <elementToTypeMapping>
        <element>//doc</element>
        <type>com.ibm.omnifind.types.PoliceReport</type>
        <featureValueAssignment>
          <feature>crimeDescription</feature>
          <basicValue useElementContent="abstract"
            trim="true">
          </basicValue>
          </featureValueAssignment>
        </elementToTypeMapping>
      </elementToTypeMappings>

</xmlCasInitializerConfiguration>
```

限制

XML 對映配置檔分成兩個區段：

<contentElements> 元素

如果要擷取特定內容，請使用此元素。配置範例檔會在文件的 <doc> 區段中擷取內容，並忽略文件中的其他區段。在 XML 治安報告中，影像可能會很大，且對文字處理不是非常有幫助。指定 <doc> 作為內容元素且不指定 <image> 後，會在開始任何文字處理之前先過濾出影像。

<elementToTypeMappings>

使用此元素可以指定文件中的那些個別 XML 元素（指定於 <elementToTypeMapping> 元素）要對映到共用分析結構中的哪些特性結構。

如果您使用內容擷取選項，則 <contentElements> 區段指定的 XML 元素必須包含 <elementToTypeMappings> 區段指定的 XML 元素。

程序

若要建立 XML 對 UIMA 類型對映配置檔：

1. 建立 XML 檔案。若要避免 XML 語法錯誤，請使用 XML 編輯器或 XML 編寫工具以驗證 XML。配置檔的 XSD 綱目稱為 configuration.xsd，且內含於企業搜尋安裝的 `ES_INSTALL_ROOT/packages/uima/` 中。
2. 在 `<xmlCasInitializerConfiguration xmlns="http://www.ibm.com/2005/uima/jedii_ci_xml">` 元素中併入您的對映。名稱空間（在 `xmlns` 屬性中指定）必須如所示範例一模一樣。
3. 如果您要從文件區段中擷取特定內容，請新增 <contentElements> 元素及 <elementToTypeMappings> 元素，後者指定要將文件中哪些個別的 XML 元素對映到共用分析區域的哪些特性結構。
4. 新增 <identifier> 元素及 <description> 元素。ID 決定 XML 文件中使用的配置。ID 必須含有文件的根元素，如 doc。如果 ID 設為預設值，則文件的根元素就不適用，且配置對映會套用至任何 XML 文件。
5. 如果要擷取只在文件相關部份才有的資訊，請新增 <contentElements> 元素。它含有下列元件元素：
 - 一或多個 <element> 元素，其中含有文件中 XML 元素的路徑並遵循 XPath 語法，例如 `<element>/doc/crimeType</element>`。
6. 如果要指定文件中的哪些 XML 元素要對映到共用分析結構中的哪些特性結構，請新增 <elementToTypeMappings> 元素。有下列元件元素：
 - 一或多個 <elementToTypeMapping> 元素。這個元素必須含有下列巢狀元素：
 - <element> 元素是用來指定 XML 元素的路徑並遵循 XPath 語法：前導正斜線 (/) 表示已提供完整路徑。例如，根元素 doc 下的 abstract。兩個正斜線 (//) 表示任何路徑子集。例如，birthDate 必須發生在 reportingOfficer 內，雖然其他元素可以發生在這兩者之間。
 - <type> 元素，指定定義於類型系統說明的類型。它必須屬於類型 Annotation。
 - 零或多個 <featureValueAssignment> 元素。
7. 在 <featureValueAssignment> 元素中，命名 <feature> 元素中類型 String 的特性，並在 <basicValue> 元素中指定值。在 <values> 元素之間可以加入多個 <basicValue> 元素。

<basicValue> 元素可以有屬性。這些包括 useAttributeValue、useElementContent、default 及 trim。

如果要使用屬性值作為特性值，請使用 useAttributeValue。下面示範

```
<elementToTypeMapping>
  <element>/doc//reportingOfficer</element>
  <type>com.ibm.omnifind.types.Person</type>
  <featureValueAssignment>
    <feature>role</feature>
    <basicValue default="Reporting officer"/>
  </featureValueAssignment>
  <featureValueAssignment>
    <feature>gender</feature>
    <basicValue default="male" useAttributeValue="sex"/>
  </featureValueAssignment>
</elementToTypeMapping>
```

產生下列輸出：

- 針對文件中 <doc> XML 標示某處所發生的每一個 <reportingOfficer> XML 標示，都會建立一個類型 com.ibm.omnifind.types.Person 的特性結構。
- 如果 <reportingOfficer> 標示含有屬性 sex，則新建特性結構的特性 gender 會設為該屬性的值。

請使用屬性 useElementContent 來新增內容作為特性的值。例如，在下列配置片段中：

```
<elementToTypeMapping>
  <element>/doc</element>
  <type>com.ibm.omnifind.types.PoliceReport</type>
  <featureValueAssignment>
    <feature>crimeDescription</feature>
    <basicValue useElementContent="abstract" trim="true"/>
  </featureValueAssignment>
</elementToTypeMapping>
```

<doc> 中元素 <abstract> 所涵蓋的文字會變成特性結構 crimeDescription 的值。將會移除所有前導及尾端空白。

在下列情況下，可以在 <values> 元素之間指定多個值：

- 要設定的特性屬於 StringArray 類型。
- 利用區隔字元屬性，將許多字串連結成一個字串，並因而對映到類型 String 的特性。例如，職稱 Mr. 是常數、名字是屬性值，而 XML 元素會涵蓋姓氏：

```
<elementToTypeMapping>
  <element>/doc//reportingOfficer</element>
  <type>com.ibm.omnifind.types.Person</type>
  <featureValueAssignment>
    <feature>surName</feature>
    <values concatenate="true" delimiter=" ">
      <basicValue default="Mr."/>
      <basicValue useAttributeValue="rank"
        default="Lt."/>
      <basicValue useElementContent="lastName"/>
    </values>
  </featureValueAssignment>
</elementToTypeMapping>
```

字串特性值會如現狀從配置檔中擷取出來。值會保留所有前導或尾端空白。但會裁去類型及特性名稱中的空白。例如，<type> com.ibm.omnifind.types.Person </type> 會變成 <type>com.ibm.omnifind.types.Person</type>。

利用 `<condition>` 元素，設定屬性上的條件。例如，只有在屬性 `armed` 設為 `yes` 的文件中發生 `<suspectDescription>` 時，才能建立類型 `com.ibm.omnifind.types.Person` 的特性結構：

```
<elementToTypeMapping>
  <element>//suspectDescription</element>
  <type>com.ibm.omnifind.types.Person</type>
  <condition attribute="armed" value="yes"/>
</elementToTypeMapping>
```

依據範例治安報告及已定義的對映配置檔，建立下列特性結構：

com.ibm.omnifind.types.PoliceReport

- covered text: "Car theft 04/23/05 09:23 pm 27 Main Street, Brynston, Springfield, New Jersey Jakob Collins 14th Precinct Male, dark haired, dark glasses, blue jeans with dark, probably black, jacket A Mercedes CLK was ... Nothing was stolen out of the car.
- begin = 2
- end = 904
- knownSuspects = null
- crimeDescription = "A Mercedes CLK was stolen on 04/23/2005 from a parking lot in front of the Blue Lagoon restaurant on 27 Main Street, Brynston.(serial number: 32 2761 50871)"

com.ibm.omnifind.types.Person

- covered text = "Jakob Collins"
- begin = 112
- end = 127
- role = "Reporting officer"
- firstName = null
- surName = "Lt Collins"
- gender = "male"

建立 XML 檔後，您必須將它上傳至企業搜尋，然後利用企業搜尋管理主控台來選取其他自訂分析選擇的 XML 文件對映配置檔。

相關概念

第 11 頁的『分析及搜尋中的 XML 標記』

您可以將文件中 XML 結構的資訊直接對映到共用分析結構，而不必撰寫 UIMA 註解程式。

相關參考

第 8 頁的『類型系統說明』

類型系統說明解釋在自訂分析中所使用的特性結構 (代表分析結果的基礎資料結構)。

文字分析結果

所有文字分析結果都儲存在共用分析結構中。

註解程式通常會在共用分析結構中讀取和寫入。共用分析結構消費者 (CAS 消費者) 會對儲存在共用分析結構的分析結果執行最終的處理程序。企業搜尋含有兩種 CAS 消費者：

- 在搜尋引擎中檢索共用分析結構內容的消費者。這類消費者需要索引建置配置檔，您可以在企業搜尋管理主控台上使用自訂文字分析來選取該配置檔。
- 在關聯式資料庫中輸入特定分析結果的消費者。這類消費者也需要您在企業搜尋管理主控台上以自訂文字分析選項所選取的配置檔。

CAS 消費者只能讀取共用分析結構。

必要時，您可以在企業搜尋中部署自訂 CAS 消費者。如需如何撰寫消費者的相關資訊，請參閱 UIMA 文件。若要瞭解如何在企業搜尋中上傳及使用消費者，請參閱 IBM UIMA developerWorks 網站，網址為 <http://www.ibm.com/developerworks/db2/zones/db2ii/>。

相關概念

2

第 21 頁的『自訂分析結果的索引對映』

對文件集合執行自訂分析後，您可以使用企業搜尋中的搜尋引擎，從儲存在共用分析結構 (以自訂分析演算法建立) 中的資訊建置索引。

2

第 27 頁的『所選分析結果的資料庫對映』

在企業搜尋對文件集合執行自訂分析後，您可以將選取的文字分析結果儲存在具有 JDBC 功能的資料庫中。

特性路徑

特性路徑可讓您存取共用分析結構中的特性值，類似用來存取 XML 文件中 XML 元素的 XPath 陳述式。

如果您要存取結合複式特性 (例如，陣列值或指向另一個特性結構的特性) 的特性結構時，特性路徑是很有用的。利用特性路徑，您可以直接關聯特性值與特性結構，並將此值儲存在語意搜尋索引或資料庫中。

例如，考慮識別汽車及其樣式的註解程式。它會建立類型 `car` 且含有屬性 `make` 的註解。然而，`make` 不含實際公司 (例如，`Chevrolet`)，但含有類型 `Company` 的特性結構，而此特性結構本身含有字串值的屬性 `companyname`。若要啓用結合汽車名稱及公司名稱的語意查詢，可以使用特性路徑 `make/companyname` 將 `companyname` 的值連接至汽車註解所產生的汽車跨距。這會利用 `'car[@make="Chevrolet"]'` 來啓用查詢「給我含有 Chevrolet 所製汽車的文件」。

特性路徑是一連串具有下列內容的特性名稱 (`f1/.../fn`)：

- 特性路徑的值可以是 `String`、`Integer`、`Float` 或其中一種類型的陣列。
- 路徑中的所有特性 (從 `f1` 到 `fn-1`) 都必須具有複數類型，亦即，屬於類型 `uima.cas.TOP`、`uima.cas.FSArray`、`uima.cas.FSList` 或其中一個次類型。
- 路徑中的最後一個特性可以包含複數類型，此外，它也可以包含 `uima.cas.Float`、`uima.cas.Integer`、`uima.cas.String`、`uima.cas.FloatArray`、`uima.cas.IntegerArray`、`uima.cas.StringArray`、`uima.cas.FloatList`、`uima.cas.IntegerList` 或 `uima.cas.StringList` 的 (次) 類型。
- 您可以選擇性地鍵入特性。完整的類型名稱必須附加到特性名稱的前面，且必須以冒號區隔。例如，`f1/com.ibm.es.SomeType:f2/.../fn`。

您可以縮小特定特性的類型範圍。例如，考慮類型 `uima.cas.TOP` 的特性 `additionalInfo`。如果您知道特性 `additionalInfo` 的值實際上是屬於類型 `EmployeeInfo`，其中含有特性 `salary`，則可以利用 `additionalInfo/EmployeeInfo:salary` 來存取此特性。請注意，在本例中特性路徑 `additionalInfo/salary` 會造成錯誤，因為 `salary` 尚未在類型 `uima.cas.TOP` 中定義。

陣列或清單值的特性具有下列額外內容：

- 使用方括弧 (`[<number>]`) 來選取陣列或清單中的特定元素。陣列從零 (0) 開始。例如，若要選取公司 (`companies`) 陣列中的第一個元素，請使用 `companies[0]`。您可以使用特殊標記 `[last]` 來選取陣列中的最後一個項目，與它的大小無關，例如 `companies[last]`。
- 使用空的方括弧 (`[]`) 來表示所有元素。在一個特性路徑中，只容許使用一個空的方括弧 (`[]`)。比方說，如果有嫌犯陣列，則特性路徑 `knownSuspects[]/com.ibm.omnifind.types.Suspect:surName` 會將所有嫌犯的姓氏收集到 `String` 陣列。
- 在檢索期間使用傳回陣列的特性路徑時，則會連結 (以空格區隔) 陣列元素並寫入索引作為單一的多字詞屬性或欄位。
- 必須鍵入特性路徑中的下一個元素。類型名稱是陣列中的元素類型。例如，考慮類型 `Info` 的特性結構。此類型含有名稱為 `companies` 的特性，其範圍是 `FSArray`。陣列的元素屬於類型 `Company`。換言之，`Company` 具有名稱為 `profit` 的特性。若要取得第三家公司的利潤，請寫入 `companies[3]/Company:profit` (利用完整的類型名稱)。

內建特性

內建特性是預先定義的特性名稱，這些名稱具有特殊的語意。它們可以用來存取特性結構本身沒有的資訊，例如，特性結構的類型或註解的涵蓋文字。它們可以當成特性路徑中的最後一個或唯一的元素使用。

下列內建特性可以在兩個對映配置檔中使用：

- `fsId()` 傳回特性結構的 ID。傳回的 ID 是整數 (32 位元)。請使用這個內建特性來存取文件中完全符合查詢的部分。
- `typeName()` 以字串傳回共用分析結構物件類型。類型是包含任何名稱空間字首的完整類型名稱，例如 `uima.tcas.Annotation`。在資料庫環境定義中，如果您在相同的直欄中儲存類型及次類型，且想要知道註解或特性結構的實際類型，則 `typeName()` 特別有用。下列範例在角色直欄中儲存了人員類型，如嫌犯或證人。

```
<explicitMappingRule applyToSubTypes="false">
  <type>com.ibm.omnifind.types.Person</type>
  <table>sample.person</table>
  <featureMappings>
    <featureMapping>
      <feature>typeName()</feature>
      <column>role</column>
    </featureMapping>
  </featureMappings>
</explicitMappingRule>
```

- `coveredText()` 傳回共用分析物件所橫跨的文字。`coveredText()` 只適用於註解及其次類型。請勿在註解類型未納入的特性結構上使用這個內建特性。下列範例在 `suspectName` 直欄中儲存了嫌犯名稱。

```

<implicitMappingRule applyToSubTypes="false">
  <type>com.ibm.omnifind.types.Suspect</type>
  <relation>sample.person</relation>
  <featureMappings>
    <featureMapping>
      <feature>coveredText()</feature>
      <column>suspectName</column>
      <length>128</length>
    </featureMapping>
  </featureMappings>
</implicitMappingRule>

```

- [] 將控點傳回給現行儲存區項目 (陣列或清單)。特性暗示疊代，這表示資料庫表格中有項目，或陣列或清單的每一個元素有索引。下列範例取自容許內建函數 [:index] 的 JDBC 配置檔。

```

<implicitMappingRule applyToSubTypes="false">
  <type>uima.cas.FSArray</type>
  <table>sample.knownSuspects</table>
  <featureMappings>
    <featureMapping>
      <feature>uniqueId()</feature>
      <column>arrayId</column>
    </featureMapping>
    <featureMapping>
      <feature>[:index]</feature>
      <column>arrayIndex</column>
    </featureMapping>
    <featureMapping>
      <feature>[]/com.ibm.omnifind.types.Suspect:uniqueId()</feature>
      <column>suspectId</column>
    </featureMapping>
  </featureMappings>
</implicitMappingRule>

```

下列內建特性只能在 JDBC 對映配置檔中使用：

- uniqueId() 傳回特性結構的廣域唯一 ID。傳回的唯一 ID 是固定長度的字串 (27 個字元)，且是 fsId()、docId()、docTimestamp() 及現行片段號碼的連結，因為文件可以在企業搜尋中分成多個共用分析結構。

傳回的字串可以包含 a-z 及 A-Z 之間的任何字元、數字 0-9、分號 (;) 及冒號 (":")。

uniqueId() 的結果可以當成表格的主要索引鍵使用。

- objectId() 傳回註解或特性結構的 ID。objectId() 類似 uniqueId()，只是它不含 docTimestamp() 的結果。傳回的 ID 只在文件剖析過一次的集合中是唯一的。如果需要在所有文件及文件版本中都是唯一，則必須使用 uniqueId()。

內建特性 objectId() 的傳回字串是 16 個字元的固定長度，且可以含有 a-z 及 A-Z 之間的任何字元、數字 0-9、分號 (;) 及冒號 (":")。

如果 uniqueId() 或 objectId() 參照的特性結構是空的，則會採用定義於資料庫表格定義中的預設值，而不會儲存被參照類型的空物件。

- docId() 傳回文件 ID。傳回的值是屬於整數類型 (32 位元)。

下列範例顯示這些內建特性：

```

<explicitMappingRule applyToSubTypes="true">
  <type>com.ibm.omnifind.types.PoliceReport</type>
  <table>sample.PoliceReport</table>

```

```

<featureMappings>
  <featureMapping>
    <feature>uniqueId()</feature>
    <column>pol iceReportId</column>
  </featureMapping>
  <featureMapping>
    <feature>docId()</feature>
    <column>pol iceReportDocId</column>
  </featureMapping>
</featureMappings>
</explicitMappingRule>

```

- docUri() 傳回文件 URI。
- docTimestamp() 傳回處理文件時的時間 (毫秒)。若要追蹤文件版本 (例如，如果您要知道使用的文件版本是否為搜索器所傳送的最新版本)，這個內建特性是很有用的。

```

<explicitMappingRule applyToSubTypes="false">
  <type>com.ibm.omnifind.types.PoliceReport</type>
  <relation>sample.PoliceReport</relationcolumn>
</StoreFeature>
<featureMappings>
  <featureMapping>
    <feature>uniqueId()</feature>
    <column>pol iceReportId</column>
  </featureMapping>
  <featureMapping>
    <feature>docTimestamp()</feature>
    <column>reportVersion</column>
  </featureMapping>
</featureMappings>
</explicitMappingRule>

```

- parentId() 傳回含有儲存區對映的特性結構之 fsId()。parentId() 只有在儲存區對映的環境定義中有效。
- uniqueParentId() 傳回儲存區對映所含的註解或特性結構的 uniqueId()。這個內建特性也只適用於儲存區對映的環境定義中。
- [:index] 傳回現行儲存區項目 (陣列或清單) 的索引。

相關工作

第 36 頁的『擷取文件中符合語意搜尋查詢的部份』

您可以將相關的特性結構對映到索引或資料庫，並在語意搜尋查詢中指定跨距，就可以只擷取文件中完全符合查詢的部份。

過濾器

過濾器是用來限制索引和 JDBC 配置檔中的對映規則。只有在過濾條件為真時，分析結果才會加入索引或 JDBC 表格。

<filter> 元素是選用的，用來只將對映限制於有特定屬性值的特性。如果要將屬性當成檢索或新增至資料庫的開關時，這是非常有用的。例如，可在類型 EntityAnnotation 的註解中記錄人員及組織。設定稱為類型的特性為人員或組織。若只要擷取人員而不要組織，您可以將下列過濾條件加入對映規則：

```
<filter syntax="FeatureValue">type = "person"</filter>
```

每一個過濾表示式的格式如下：

```
<FeaturePath> <Operator> <Literal>
```

其中：

- FeaturePath 是共用分析結構中的特性路徑
- 運算子是 =、!=、<、<=、> 或 >=。請注意，< (且只有 <) 必須表示為 <。
- 文字是整數、浮點數 (不支援指數語法) 或以雙引號括住的字串文字，內含引號及使用反斜線作為跳出字元的反斜線。

<FeaturePath>、<Operator> 及 <Literal> 必須以空格區隔。

下列範例是有效的過濾器：

- <filter syntax="FeatureValue"> foo = "hello world" </filter>

特性 foo 含有字串 hello world。

- <filter syntax="FeatureValue"> foo < 42 </filter>

特性 foo 含有整數值 42。

- <filter syntax="FeatureValue"> make/company = "Chevrolet" </filter>

特性路徑 make/company，其中特性 make 所含的特性結構中具有隨附值 Chevrolet 的特性。

- <filter syntax="FeatureValue"> bar7 >= 0.5 </filter>

特性 bar7 含有浮點值 0.5。

自訂分析結果的索引對映

對文件集合執行自訂分析後，您可以使用企業搜尋中的搜尋引擎，從儲存在共用分析結構 (以自訂分析演算法建立) 中的資訊建置索引。

將分析結果對映到企業搜尋索引中的欄位、文字跨距及屬性，可讓您在查詢中使用此資訊。合併自訂分析與可以檢索單字和文字跨距的企業搜尋，可啟用語意搜尋。

使用索引建置配置檔，您可以決定要檢索共用分析結構中的哪些分析結果。

您可以使用不同的樣式，將共用分析結構中的特性結構對映到企業搜尋索引。

註解 如果利用註解樣式來檢索共用分析結構中的特性結構，則所有指定類型的註解都會在索引中儲存為可搜尋的跨距。

比方說，如果跨距特定文字區域的特性結構屬於 person 類型，並使用註解樣式加以檢索，則下列查詢是可行的：

表 1. 範例查詢

必要資訊	可能的查詢
給我至少含有一個人員名稱的所有文件	<person/>
給我人員註解中含有上司的所有文件	<person>boss</person>
給我在我的競爭對手之一的相同句子中提到 Lang 的所有文件	<sentence><person>Lang</person><competitor/></sentence>

特性結構的屬性也可以當成跨距的一部份來檢索。例如，考慮偵測汽車的註解程式，並將汽車樣式儲存為 car 註解的特性 make。這會啓用下列類型的查詢：「給我提到 Chevrolet 汽車樣式的文件」。

欄位 如果要在搜尋期間，利用企業搜尋的欄位搜尋功能來存取特性結構的內容，請使用此樣式。在此方式中，特性結構的內容可以顯示在搜尋結果中，或用於參數搜尋。

比方說，如果將藥品劑量對映到參數欄位，則可以使用下列查詢：「給我談到某些藥品劑量超過 100 毫克的所有文件」。

中斷 如果要將特定的特性結構解譯為清除區隔字元 (例如，小節或段落)，則請使用此樣式。企業搜尋預設為偵測句子及段落。只有在您的自訂分析在文件中偵測到額外的結構元素且您必須有不同的解譯方式時，才能使用此樣式。

也可以使用分析結果來影響企業搜尋中的文件相關性排序，即使是在簡式關鍵字查詢中。此程序分成兩步驟：

1. 使用「註解」或「欄位」對映樣式，將特性結構對映到可搜尋的跨距或欄位。
2. 使用企業搜尋管理主控台來定義 boost 類別，並將跨距或欄位名稱對映到這個 boost 類別。

如果使用者輸入的搜尋字詞內含於特性結構中，則文件的排序會較高。例如，考慮偵測人員和公司名稱的註解程式。將這些特性結構對映到跨距 (如「人員」及「公司」)，然後將這些跨距對映到 boost 類別，則 "gap" 的搜尋結果會讓談到 "Gap" 公司的文件排序高於只含有 "gap" 一詞的文件。

撰寫索引建置配置檔後，您就可以使用管理主控台將它上傳至企業搜尋。

相關工作

- 2 『建立索引建置配置檔』
使用索引建置配置檔，您可以決定要在共用分析結構中檢索哪些分析結果以啓用搜尋。
- 2

建立索引建置配置檔

使用索引建置配置檔，您可以決定要在共用分析結構中檢索哪些分析結果以啓用搜尋。

關於本作業

索引建置配置檔必須符合下面範例所顯示的綱目。範例配置檔是以治安報告實務範例中定義的類型系統為基礎。

```
<?xml version="1.0" encoding="UTF-8"?>
<indexBuildSpecification
xmlns="http://www.ibm.com/of/822/consumer/index/xml">
  <skipCondition>
    <type>com.ibm.uima.tt.DocumentAnnotation</type>
    <filter syntax="FeatureValue">toBeProcessed = 0</filter>
  </skipCondition>

  <indexBuildItem>
    <name>com.ibm.omnifind.types.Person</name>
    <indexRule>
      <style name="Annotation">
        <attributeMappings>
          <mapping>
```



```

        </mapping>
        <mapping>
          <feature>time/coveredText()</feature>
          <indexName>time</indexName>
        </mapping>
        <mapping>
          <feature>date/englDate</feature>
          <indexName>date</indexName>
        </mapping>
        <mapping>
          <feature>location/coveredText()</feature>
          <indexName>location</indexName>
        </mapping>
        <mapping>
          <feature>knownSuspects[]/com.ibm.omnifind.types.Suspect:surName</feature>
          <indexName>suspectsLastNames</indexName>
        </mapping>
      </attributeMappings>
    </style>
  </indexRule>
</indexBuildItem>
</indexBuildSpecification>

```

限制

索引對映配置檔必須含有所有您需要的分析結果，才能夠在查詢中搜尋。

程序

若要建立索引對映配置檔：

1. 建立 XML 檔案。若要避免 XML 語法錯誤，請選擇使用 XML 編輯器或 XML 編寫工具。配置檔的 XSD 綱目稱為 CasToIndexMapping.xsd，且內含於企業搜尋安裝的 `ES_INSTALL_ROOT/packages/uima/` 中。
2. 在 `<indexBuildSpecification xmlns="http://www.ibm.com/of/822/consumer/index/xml">` 元素中併入您的對映。名稱空間 (在 `xmlns` 屬性中指定) 必須如所示範例一模一樣。
3. 新增 `<skipCondition>` 元素以依據特定特性值，禁止檢索某些文件。這是選用的元素。在範例中，將不會檢索資料結構類型為 `com.ibm.uima.tt.DocumentAnnotation` 且特性 `toBeProcessed` 設為 0 的文件。
4. 新增一或多個 `<indexBuildItem>` 元素，其中含有共用分析結構中特定特性結構對索引中某一結構的對映。
5. 儲存並驗證 XML 檔。

<indexBuildItem> 元素

索引建置規格配置檔含有一或多個 `<indexBuildItem>` 元素。每一個說明均說明共用分析結構中特定特性結構對索引 (跨距或欄位) 中某一結構的對映。

`<name>` 元素含有特性結構類型。指定類型的方式有兩種：

- 完整類型名稱。例如，`com.ibm.omnifind.types.Suspect`
- 萬用字元。例如，`com.ibm.omnifind.types.*`。萬用字元只能在類型規格的尾端加入。

只使用 `uima.tcas.Annotation` 的次類型作為索引建置項目。如果特性結構是次類型 `uima.cas.TOP` (而不是屬於 `uima.tcas.Annotation`)，您可以利用從註解開始的特性路徑來存取這個特性結構。

如果類型 A 是類型 B 的次類型 (在範例中，`com.ibm.omnifind.types.Suspect` 是 `com.ibm.omnifind.types.Person` 的次類型)，且這兩個類型均定義了 `<indexBuildItem>` 元素 Ia 及 Ib，則處理程序如下：

- 定義於 Ib 的每一個索引規則會套用於類型 B 的特性結構和類型 A 的特性結構
- 定義於 Ia 的每一個索引規則只會套用於類型 A 的特性結構

在範例中，`com.ibm.omnifind.types.Person` 註解定義的 `<indexBuildItem>` 元素也會套用於 `com.ibm.omnifind.types.Suspect` 註解。建立嫌犯註解的兩個跨距：一個命名為「人員」，另一個則為「嫌犯」。

`<filter>` 元素是選用的，用來只將 `<indexBuildItem>` 對映限制於含有特定屬性值的特性結構。如果要將屬性當成檢索的開關時，這是非常有用的。例如，可在類型 `EntityAnnotation` 的註解中記錄人員及組織。設定稱為類型的特性為人員或組織。若只要擷取人員而不要組織，您可以新增下列過濾條件：

```
<filter syntax="FeatureValue">type = "person"</filter>
```

此外，您還可以選擇檢索不同跨距名稱下的人員和組織，例如：`person` 及 `organization`。若要執行此作業，請定義類型 `EntityAnnotation` 的兩個 `<indexBuildItem>` 元素，然後在 `type` 特性上使用兩個過濾器以觸發人員或組織。

<indexRule> 元素

每一個 `<indexBuildItem>` 元素都含有一個 `<indexRule>` 元素。每個 `<indexRule>` 元素都含有所有必要資訊，以用來將共用分析結構中的特性結構對映到索引作為欄位、註解或中斷樣式。註解及欄位樣式支援許多屬性。但您不能在企業搜尋中使用「UIMA 軟體開發套件」支援的詞彙樣式。

若為註解及欄位樣式，當您在索引中指定註解或欄位名稱時，可以使用下列選擇方案：

- 如果要在索引中以同一個名稱存取每一個特性結構，請使用 `fixedName`。在下列範例中，每一個類型 `Person` 的特性結構都會對映到索引中名稱為「人員」的跨距。

```
<indexBuildItem>
  <name>com.ibm.omnifind.types.Person</name>
  <indexRule>
    <style name="Annotation">
      <attribute name="fixedName" value="Person" />
    </style>
  </indexRule>
</indexBuildItem>
```

這會啟用「給我在人員名稱中含有上司的文件」之類的查詢。使用 XML 片段的查詢表示式如下：`@xmlf2::'<Person>Boss</Person>'`

- 如果註解儲存了不同的實體，而您想要依據註解的特定特性值來利用不同的跨距存取這些實體，則請使用 `nameFeature`。在下面範例中，`EntityAnnotation` 檢索為 `person` 或 `organization` 跨距，視名稱為 `type` 的特性值而定。特性也可以是特性路徑。

```

<indexBuildItem>
  <name>com.ibm.tt.EntityAnotation</name>
  <indexRule>
    <style name="Annotation">
      <attribute name="nameFeature" value="type" />
    </style>
  </indexRule>
</indexBuildItem>

```

這會啓用「給我在 WHO 組織的相關文件」(反對英文的 who 一詞) 之類的查詢。使用有限的 XPath 語法，查詢表示式如下：`@xslp::'/organization[ftcontains="WHO"]'`

- 如果未使用上述任何屬性，則會使用 `<indexBuildItem>` 元素中註解類型的短名稱。這是預設值。例如：

```

<indexBuildItem>
  <name>com.ibm.uima.tutorial.RoomNumber</name>
  <indexRule>
    <style name="Annotation" />
    <style name="Field" />
  </indexRule>
</indexBuildItem>

```

這個 `<indexBuildItem>` 元素會產生名稱爲 `RoomNumber` 的註解和欄位，並在其中輸入 `com.ibm.uima.tutorial.RoomNumber` 所涵蓋的文字。

`<style name="Annotation" />` 元素

`<style>` 元素中的註解指定如何在企業搜尋中存取跨距資訊。除了容許使用 `fixedName` 及 `nameFeature` 屬性外，此樣式還支援 `<attributemappings>` 元素。在這個元素中，可以將特性值對映到索引中結果跨距的屬性，以便後續在搜尋表示式中使用。

每一個對映都是在個別的 `<mapping>` 元素中執行。`<feature>` 元素含有特性路徑，而 `<indexName>` 元素含有在索引中用來儲存 `<feature>` 值的屬性名稱。例如，

```

  <mapping>
    <feature>make/companyname</feature>
    <indexName>company</indexName>
  </mapping>

```

這個 `<mapping>` 元素會將路徑 `make/companyname` 中的特性值直接儲存在索引屬性 `company` 中。

如果在文字分析期間使用的類型系統是複式的，其中包括許巢狀的特性結構，則將特性值對映到索引屬性特別有用。使用 `<mapping>` 元素時，會顯現相關屬性，以便讓您在查詢中使用它們，而不必清楚地瞭解原始類型系統結構。

`<style name="Field" />` 元素

`<style>` 元素中的欄位指定如何在企業搜尋中存取欄位資訊。除了 `fixedName` 及 `nameFeature` 屬性外，您還可以設定下列屬性。

參數 如果設爲 `true`，則可以利用參數搜尋來搜尋欄位值，例如，`#dosage:>100`

可搜尋的欄位

如果設爲 `true`，則可以在搜尋中使用欄位值，例如 `make:Bayer`

可傳回的

如果設爲 `true`，則會在搜尋結果中傳回欄位及其值。

欄位資訊一律是可以搜尋的內容，亦即，可以在一般關鍵字搜尋中存取欄位資訊。

選用屬性 `valueFeature` 定義哪些特性值要當成欄位值。如果特性結構是註解，且未設定屬性，則會使用註解的涵蓋文字作為欄位值。在範例中，

```
<indexBuildItem>
  <name>com.ibm.omnifind.types.Date</name>
  <indexRule>
    <style name="Field">
      <attribute name="fixedName" value="date"/>
      <attribute name="fieldSearchable"
        value="true"/>
      <attribute name="returnable" value="true"/>
    </style>
    <style name="Field">
      <attribute name="fixedName" value="hour"/>
      <attribute name="valueFeature" value="hour"/>
      <attribute name="parametric" value="true"/>
    </style>
  </indexRule>
  <filter syntax="FeatureValue">year="2005"</filter>
</indexBuildItem>
```

產生 `com.ibm.omnifind.types.Date` 的兩個欄位。一個欄位名稱為 `date` 且含有涵蓋的文字，例如，`5:15pm`。另一個欄位含有屬性 `hour` 的值。您可以在此使用 `'hour::<17'` 進行查詢。

<style name="Breaking" /> 元素

`<style>` 元素中的值 `Breaking` 不含任何進一步元素。

建立 XML 檔後，您必須將它上傳至企業搜尋，然後利用企業搜尋管理主控台來選取其他自訂分析選擇的索引對映配置檔。

相關概念

2

第 21 頁的『自訂分析結果的索引對映』

2

對文件集合執行自訂分析後，您可以使用企業搜尋中的搜尋引擎，從儲存在共用分析結構 (以自訂分析演算法建立) 中的資訊建置索引。

第 17 頁的『特性路徑』

特性路徑可讓您存取共用分析結構中的特性值，類似用來存取 XML 文件中 XML 元素的 XPath 陳述式。

相關參考

第 20 頁的『過濾器』

過濾器是用來限制索引和 JDBC 配置檔中的對映規則。只有在過濾條件為真時，分析結果才會加入索引或 JDBC 表格。

第 8 頁的『類型系統說明』

類型系統說明解釋在自訂分析中所使用的特性結構 (代表分析結果的基礎資料結構)。

所選分析結果的資料庫對映

在企業搜尋對文件集合執行自訂分析後，您可以將選取的文字分析結果儲存在具有 JDBC 功能的資料庫中。

這個版本只支援 DB2® Universal Database 8.2.2 版 (com.ibm.db2.jcc.DB2Driver 2.3 版) 及 Oracle 10g (oracle.jdbc.driver.OracleDriver 1.0 版)。

若為 DB2 Universal Database 及 Oracle，您可以選擇將分析結果直接插入資料庫，或產生相等的資料庫特定載入檔及對應的 Script 以執行載入命令。

將分析結果對映到資料庫中的表格，可讓您在後續的商業情報處理程序步驟中使用此資訊，或直接存取文件中符合語意搜尋查詢的相關部份。

XML 對映配置檔含有資料庫連線配置資訊，並說明哪些自訂分析結果要儲存在哪些表格和直欄中。配置檔中的表格和直欄名稱必須對應於在資料庫中建立的表格和直欄。

撰寫配置檔後，您就可以使用管理主控台將檔案上傳至企業搜尋。

相關工作

『建立 XML 對映配置檔』

若要將分析結果新增至資料庫，您必須建立配置檔，其中包含資料庫連線配置資訊的配置檔，以及自訂文字分析結果要儲存在哪些表格和直欄的說明。

在資料庫中儲存分析結果

若要在具有 JDBC 功能的資料庫中儲存選取的分析結果，則必須撰寫適用於企業搜尋的配置檔，且必要的 JDBC 驅動程式庫必須位於配置檔所定義的路徑。

若要在具有 JDBC 功能的資料庫中儲存分析結果：

1. 決定要在資料庫中儲存哪些分析結果。建立資料庫，其中的表格含有適當資料類型的必要直欄。

重要：建立您自己的 DB2 資料庫，以儲存選取的分析結果。請勿使用企業搜尋安裝所提供的 DB2 資料庫。

2. 在 XML 編輯器中，撰寫含資料庫配置資料的配置檔及您要儲存的分析結果。為了決定要在配置檔中併入哪些分析結果，您必須知道自訂分析所使用的基礎類型系統。
3. 將 JDBC 驅動程式庫放入可以從企業搜尋系統的索引節點存取的目錄中。
4. 使用企業搜尋管理主控台，上傳及選取含有自訂文字分析的配置檔。

建立 XML 對映配置檔

若要將分析結果新增至資料庫，您必須建立配置檔，其中包含資料庫連線配置資訊的配置檔，以及自訂文字分析結果要儲存在哪些表格和直欄的說明。

關於本作業

XML 對映配置檔必須符合下面範例所顯示的綱目。範例是以在治安報告實務範例中定義的類型系統為基礎。

在範例中，只有治安報告和這些治安犯罪報告中出現的城市會新增至資料庫。範例顯示內建特性及 <constant> 元素對映的用法。

```
<?xml version="1.0" encoding="UTF-8"?>
<cas2JdbcConfiguration xmlns="http://www.ibm.com/uima/consumer/jdbc/100/xml">
  <databaseConnection>
    <connectionUrl>db2://myMachine:myPort/myDatabase</connectionUrl>
    <driver type="jdbc">com.ibm.db2.jcc.DB2Driver</driver>

    <driverLibraries>
      <driverLibrary>C:\db2\db2jcc.jar</driverLibrary>
      <driverLibrary>C:\db2\db2jcc_license_cu.jar</driverLibrary>
    </driverLibraries>
  </databaseConnection>
</cas2JdbcConfiguration>
```

```

    <driverLibrary>C:\db2\db2jcc_license_cisuz.jar</driverLibrary>
</driverLibraries>

<authentication>
  <username>myUser</username>
  <password>myPassword</password>
</authentication>

<loadFile>
  <loadFileDirectory>/home/cas2jdbc/load/</loadFileDirectory>
  <loadScript>/home/cas2jdbc/load/load.sh</loadScript>
</loadFile>

</databaseConnection>

<jdbcMappingSpec>
<skipCondition>
  <name>com.ibm.uima.tt.DocumentAnnotation</name>
  <filter syntax="FeatureValue">toBeProcessed=0</filter>
</skipCondition>

<cas2JdbcMappings>
<explicitMappings>
  <explicitMappingRule applyToSubtypes="false">
<type>com.ibm.omnifind.types.PoliceReport</type>
<table>sample.policeReport</table>
<featureMappings>
  <featureMapping>
    <feature>uniqueId()</feature>
    <column>policeReportId</column>
  </featureMapping>
  <featureMapping>
    <feature>location/uniqueId()</feature>
    <column>crimeLocationId</column>
  </featureMapping>
</featureMappings>
  <filter syntax="FeatureValue">location/coveredText()="Los Angeles"</filter>
</explicitMappingRule>
</explicitMappings>

<implicitMappings>
<implicitMappingRule applyToSubtypes="false">
<type>com.ibm.omnifind.types.City</type>
<table>sample.City</table>
<featureMappings>
  <featureMapping>
    <feature>uniqueId()</feature>
    <column>crimeLocationId</column>
  </featureMapping>
  <featureMapping>
    <feature>coveredText()</feature>
    <column>cityName</column>
    <length>150</length>
  </featureMapping>
  <featureMapping>
    <constant>USA</constant>
    <column>country</column>
  </featureMapping>
</featureMappings>
</implicitMappingRule>
</implicitMappings>
</cas2JdbcMappings>
</jdbcMappingSpec>
</cas2JdbcConfiguration>

```

限制

建立您自己的 DB2 資料庫以儲存選取的分析結果。請勿使用企業搜尋安裝所提供的 DB2 資料庫。

程序

若要建立 XML 資料庫配置檔：

1. 建立 XML 檔案。若要避免 XML 語法錯誤，請選擇使用 XML 編輯器或 XML 編寫工具。配置檔的 XSD 綱目稱為 `CasToJDBCMapping.xsd`，且內含於企業搜尋安裝的 `ES_INSTALL_ROOT/packages/uima/` 中。
2. 在 `<cas2JdbcConfiguration xmlns="http://www.ibm.com/uima/consumer/jdbc/100/xml">` 元素中併入您的對映。名稱空間 (在 `xmlns` 屬性中指定) 必須如所示範例一模一樣。
3. 新增含有所有資料庫連線配置資訊的 `<databaseConnection>` 元素，以及說明儲存在資料庫或載入檔的分析結果對映規則的 `<jdbcMappingSpec>` 元素。

4. 將下列元件元素新增至 `<databaseConnection>` 元素：

- 必備元素：`<connectionUrl>` 元素。此元素含有資料庫連線 URL。您可以視 JDBC 驅動程式實作而定，以本端或遠端存取資料庫。
- 必備元素：`<driver>` 元素。此元素含有 JDBC 驅動程式類別的名稱，例如，適用於 DB2 的 `com.ibm.db2.jcc.DB2Driver` 或適用於 Oracle 的 `oracle.jdbc.driver.OracleDriver`。
- 必備元素：`<driverLibraries>` 元素。此元素列出驅動程式庫。每一個程式庫都列示在 `<driverLibrary>` 元素中。程式庫是在 DB2 或 Oracle 安裝目錄中。若為 DB2，程式庫是 `c:\your_db2_dir\db2jcc.jar`、`c:\your_db2_dir\db2jcc_license_cu.jar` 及 `c:\your_db2_dir\db2jcc_license_cisuz.jar`。若為 Oracle，包含的程式庫是 `c:\your_oracle_dir\classes12.zip`。
- 必備元素：`<authentication>` 元素。此元素含有資料庫的使用者名稱及密碼。
- 選用元素：`<loadFile>` 元素。此元素含有 `<loadFileDirectory>` 元素中的載入檔目錄，以及 `<loadScript>` 元素中的載入 Script 名稱。如果您沒有指定 `<loadFile>` 元素，則可利用 JDBC 將所有資料直接儲存在資料庫中。

當您使用資料庫特有的載入檔及 Script 時，您還必須新增所有資料庫配置參數。

5. 將下列元件元素新增至 `<jdbcMappingSpec>` 元素：

- 選用元素：`<skipCondition>` 元素。如果沒有定義略過條件，則會處理所有文件。

```
<skipCondition>
  <name>com.ibm.uima.tt.DocumentAnnotation</name>
  <filter syntax="FeatureValue">toBeProcessed=0</filter>
</skipCondition>
```

在範例中，將不考慮註解類型為 `com.ibm.uima.tt.DocumentAnnotation` 且特性 `toBeProcessed` 設為 0 的文件。

- `<cas2JdbcMappings>` 元素顯示哪些類型及特性對映至哪些資料庫表格及直欄。元素含有明確及隱含的對映區段。
6. 新增 `<explicitMappings>` 元素。這是必備元素。它必須具有一或多個定義明確對映的 `<explicitMappingRule>` 元素，且只能針對註解類型及其次類型定義。如果對映定義於明確對映區段，則符合對映定義的所有註解都會儲存在資料庫中。

7. 選用元素：新增 `<implicitMappings>` 元素。此元素支援所有特性結構類型。如果此元素存在，則必須至少含有一個 `<implicitMappingRule>` 元素。只有在另一個符合明確或隱含對映規則的註解參照相符註解類型時，才能將定義於隱含對映區段的對映加入資料庫。

隱含對映的目的是要讓您只儲存出現在特定環境定義的分析結果。比方說，如果註解類型 `com.ibm.omnifind.types.City` 的對映是隱含的，則只有明確對映區段中 `com.ibm.omnifind.types.PoliceReport` 對映定義所參照的城市會儲存在資料庫中。這表示只有治安報告所提到的城市會新增至資料庫。

如果有適用於「城市」註解的明確對映規則，則所有城市都會新增至資料庫。在這兩種情況下，如果某一城市有多個治安報告參照，也只會加入資料庫一次。

8. `<explicitMappingRule>` 及 `<implicitMappingRule>` 元素必須含有屬性 `applyToSubtypes`，如果該屬性設為 `true`，則不僅會儲存列示在 `<type>` 元素中的特性結構，還會儲存從該結構所衍生出來的所有特性結構。請將下列元件元素新增至 `<explicitMappingRule>` 及 `<implicitMappingRule>` 元素：

- 含有特性結構類型的 `<type>` 元素。
- `<table>` 元素包含資料庫綱目及表格名稱。語法遵循規則 `schema.table_name`，或如果未定義綱目，則只有 `table_name`。
- 含有一或多個 `<featureMapping>` 元素或一個 `<containerMapping>` 元素的 `<featureMappings>` 元素。
- 選用：`<filter>` 元素，含有每次對映規則相符時評估的條件。如果條件評估為 `true`，則特性結構的註解會儲存在資料庫中。在範例中，只有處理洛杉磯犯罪的治安報告會儲存在資料庫中。

9. `<featureMapping>` 元素元件結構會隨著您是否對映特性或常數而改變。

如果您要對映特性或特性路徑，則元件元素包括：

- 具有特性名稱的 `<feature>` 元素。必須定義類型元素中特性結構的特性。您也可以使用特性路徑建構或任何系統定義的內建特性。
- 選用：`<length>` 元素，其字串長度定義於指定的資料庫直欄。較長的字串會被截斷。
- `<column>` 元素含有要儲存特性值的直欄名稱。未在任何特性對映中使用的資料庫直欄會使用資料庫配置的預設值（通常是空值）。

請確定特性元素的值儲存在適當類型的直欄中。下表顯示哪些 UIMA 類型符合哪些資料庫類型。

表 2. UIMA 類型及相對應的資料庫類型之間的對映

UIMA 類型或內建特性	建議的 DB2 資料類型	建議的 Oracle 資料類型
Float	REAL	FLOAT
String	VARCHAR	VARCHAR2
Integer	INTEGER	INTEGER
uniqueId(), uniqueParentId()	CHAR(27)	CHAR(27)
objectId(), parentId()	CHAR(16)	CHAR(16)
docTimestamp()	BIGINT	LONG

若為常數，元件特性對映元素如下：

- `<constant>` 元素，含有常數值。
- `<column>` 元素，含有要新增常數值的直欄名稱。

10. `<containerMapping>` 元素含有儲存區類型特性 (陣列或清單) 的對映。此元素只適用於儲存區類型。有下列元件元素：

- 具有特性名稱的 `<feature>` 元素。您也可以使用特性路徑建構或任何系統定義的內建特性。
- `<table>` 元素包含資料庫綱目及表格名稱。語法遵循規則 `schema.table_name`，或如果未定義綱目，則只有 `table_name`。
- 一或多個 `<featureMapping>` 元素，其中含有特性結構名稱及要新增特性的直欄名稱。

11. 使用提供的綱目來儲存及驗證 XML 檔。

建立 XML 檔後，您必須將它上傳至企業搜尋，然後利用企業搜尋管理主控台來選取其他自訂分析選擇的資料庫對映配置檔。

相關概念

第 27 頁的『所選分析結果的資料庫對映』

在企業搜尋對文件集合執行自訂分析後，您可以將選取的文字分析結果儲存在具有 JDBC 功能的資料庫中。

第 17 頁的『特性路徑』

特性路徑可讓您存取共用分析結構中的特性值，類似用來存取 XML 文件中 XML 元素的 XPath 陳述式。

相關參考

第 20 頁的『過濾器』

過濾器是用來限制索引和 JDBC 配置檔中的對映規則。只有在過濾條件為真時，分析結果才會加入索引或 JDBC 表格。

第 18 頁的『內建特性』

內建特性是預先定義的特性名稱，這些名稱具有特殊的語意。它們可以用來存取特性結構本身沒有的資訊，例如，特性結構的類型或註解的涵蓋文字。它們可以當成特性路徑中的最後一個或唯一的元素使用。

第 8 頁的『類型系統說明』

類型系統說明解釋在自訂分析中所使用的特性結構 (代表分析結果的基礎資料結構)。

儲存區類型對映

儲存區類型是共用分析結構中的其中一個內建陣列或清單類型。儲存區類型對映是將陣列或清單值對映到關聯式資料庫的一種方式。

有兩種方法可以處理配置檔中的儲存區類型。方法之一是使用已定義的內建特性建構及一般鏈結表格，其中所含的陣列或清單是特性對映規則的值。因為不同的陣列或清單都儲存在相同的鏈結表格中，所以表格無法指出儲存資訊的關係。

在第二種方法中，以 `<containerMapping>` 元素所定義的鏈結表格定義明確地表示您要擁有的指定資訊之間的關係。

下面範例顯示一般鏈結表格對映可能的樣子。在治安報告和嫌犯之間有 n:m 關係，表示某一嫌犯會在一或多個治安報告中提及，且一份治安報告可以提及多個嫌犯。

範例中的一般 `sample.fsarray` 表格是治安報告和嫌犯之間的鏈結表格。如果除了特性類型 `com.ibm.omnifind.types.FSArray` 的 `com.ibm.omnifind.types.PoliceReport` 外，還有其他對映類型，也會對映到這個表格。您還是可以查詢表格以瞭解治安報告和嫌犯之間的正確關係，但是，您不能只是查看表格就斷定其中含有治安報告和可能嫌犯之間的關係或鏈結。

```

<cas2JdbcMappings>
  <explicitMappings>
    <explicitMappingRule applyToSubtypes="false">
      <type>com.ibm.omnifind.types.PoliceReport</type>
      <table>sample.policeReport</table>
      <featureMappings>
        <featureMapping>
          <feature>uniqueId()</feature>
          <column>policeReportId</column>
        </featureMapping>
        <featureMapping>
          <feature>knownSuspects/uniqueId()</feature>
          <column>suspectArrayId</column>
        </featureMapping>
        <featureMapping>
          <feature>location/cityName</feature>
          <column>city</column>
        </featureMapping>
      </featureMappings>
    </explicitMappingRule>
  </explicitMappings>

  <implicitMappings>
    <implicitMappingRule applyToSubtypes="false">
      <type>com.ibm.omnifind.types.Suspect</type>
      <table>sample.suspect</table>
      <featureMappings>
        <featureMapping>
          <feature>uniqueId()</feature>
          <column>suspectID</column>
        </featureMapping>
        <featureMapping>
          <feature>surName</feature>
          <column>lastName</column>
        </featureMapping>
        <featureMapping>
          <feature>說明</feature>
          <column>說明</column>
        </featureMapping>
      </implicitMappingRule>
    <implicitMappingRule applyToSubtypes="false">
      <type>uima.cas.FSArray</type>
      <table>sample.fsarray</table>
      <featureMappings>
        <featureMapping>
          <feature>uniqueId()</feature>
          <column>arrayId</column>
        </featureMapping>
        <featureMapping>
          <feature>[:index]</feature>
          <column>arrayIndex</column>
        </featureMapping>
        <featureMapping>
          <feature>[]/uniqueId()</feature>
          <column>suspectId</column>
        </featureMapping>
      </featureMappings>
    </implicitMappingRule>
  </implicitMappings>
</cas2JdbcMappings>

```

```

    </implicitMappingRule>
  </implicitMappings>
</cas2JdbcMappings>

```

下面顯示以上述一般對映規則為基礎的資料庫表格。

表 3. *sample.policeReport* 表格

policeReportId	suspectArrayId	city
aaa...1	bbb...1	Springfield
aaa...2	bbb...2	Ladysmith

表 4. *sample.fsarray* 表格

arrayId	arrayIndex	suspectId
bbb...1	1	ccc...1
bbb...1	2	ccc...2
bbb...2	1	ccc...3

表 5. *sample.suspect* 表格

suspectID	lastname	說明
ccc...1	Brown	Dark complexion
ccc...2	Smith	Wears glasses
...

範例顯示特性結構陣列的對映。您也可以將此類型的對映套用於 `StringArray`、`IntegerArray` 及 `FloatArray`。如果併入這些簡式值陣列的對映規則，請以 `[]` 取代 `[]/uniqueId()`。

相同的一般表格方法可以用於特性結構清單，以及簡式類型清單 (`StringList`、`IntegerList` 及 `FloatList`)。

更簡單的處理關係方法是使用明確的儲存區對映元素，它定義陣列或清單所含元素的疊代。

下面範例中的對映指出明確的鏈結表格。而在治安報告及嫌犯之間，再次存在 n:m 關係。但是，這次 `sample.reports_suspects` 表格是治安報告和嫌犯之間的鏈結表格。

在這種方式中，您不必考慮處理陣列 ID，或清單類型的頭尾項目對映。鏈結表格含有一個明確的關係。

```

<cas2JdbcMappings>
  <explicitMappings>
    <explicitMappingRule applyToSubtypes="false">
      <type>com.ibm.omnifind.types.PoliceReport</type>
      <table>sample.policeReport</table>
      <featureMappings>
        <featureMapping>
          <feature>uniqueId()</feature>
          <column>policeReportID</column>
        </featureMapping>
        <featureMapping>
          <feature>location/cityName</feature>
          <column>city</column>
        </featureMapping>
      </featureMappings>
    </explicitMappingRule>
  </explicitMappings>
</cas2JdbcMappings>

```

```

</featureMapping>
<featureMapping>
  <feature>knownSuspects</feature>
  <containerMapping>
    <table>sample.reports_suspects</table>
    <featureMapping>
      <feature>com.ibm.omnifind.types.PoliceReport
        /objectId()</feature>
      <column>policeReportId</column>
    </featureMapping>
    <featureMapping>
      <feature>knownSuspects/[]/objectId()</feature>
      <column>suspectId</column>
    </featureMapping>
  </containerMapping>
</featureMapping>
</featureMappings>
</explicitMappingRule>
</explicitMappings>

<implicitMappings>
  <implicitMappingRule applyToSubtypes="false">
    <type>com.ibm.omnifind.types.Suspect</type>
    <table>sample.suspect</table>
    <featureMappings>
      <featureMapping>
        <feature>objectId()</feature>
        <column>suspectID</column>
      </featureMapping>
      <featureMapping>
        <feature>surName</feature>
        <column>lastName</column>
      </featureMapping>
      <featureMapping>
        <feature>說明</feature>
        <column>說明</column>
      </featureMapping>
    </featureMappings>
  </implicitMappingRule>
</implicitMappings>
</cas2JdbcMappings>

```

<containerMapping> 元素是用來定義陣列所含元素的疊代。在範例中，sample.reports_suspects 鏈結表格含有 policeReportId 及 suspectId 直欄的鏈結。請勿將 <containerMapping> 元素巢狀化。

下面顯示以明確鏈結表格對映規則為基礎的資料庫表格。

表 6. sample.policeReport 表格

policeReportId	city
aaa...1	Springfield
aaa...2	Ladysmith

表 7. sample.reports_suspect 表格

policeReportId	suspectId
bbb...1	ccc...1
bbb...2	ccc...2
...	...

表 8. *sample.suspect* 表格

suspectID	lastname	說明
ccc...1	Brown	Dark complexion
ccc...2	Smith	Wears glasses
...

相關參考

第 18 頁的『內建特性』

內建特性是預先定義的特性名稱，這些名稱具有特殊的語意。它們可以用來存取特性結構本身沒有的資訊，例如，特性結構的類型或註解的涵蓋文字。它們可以當成特性路徑中的最後一個或唯一的元素使用。

擷取文件中符合語意搜尋查詢的部份

您可以將相關的特性結構對映到索引或資料庫，並在語意搜尋查詢中指定跨距，就可以只擷取文件中完全符合查詢的部份。

若要存取搜尋結果中特定註解類型的所有實例 (例如，若要取得所有人員)，請併入註解類型的欄位樣式對映，並在索引配置檔中將它標示為可傳回的。例如：

```
<indexBuildItem>
  <name>com.ibm.omnifind.types.Person</name>
  <indexRule>
    <style name="Annotation"/>
    <style name="Field">
      <attribute name="returnable" value="true"/>
    </style>
  </indexRule>
</indexBuildItem>
```

在此範例中，類型 `com.ibm.omnifind.types.Person` 的註解會對映到企業搜尋索引中名稱為「人員」的跨距，在語意搜尋時可以在其中存取它們。此外，註解的涵蓋文字 (例如，完整的人員名稱) 會儲存為可傳回的欄位。若要擷取這些註解值，請在搜尋查詢 (關鍵字或語意) 傳回的每一個結果物件上呼叫 `getFields("Person")`。此方法會傳回含有註解值的「字串」陣列，在此案例中，會傳回人員名稱。

然而，這種方式會傳回指定註解類型的所有實例，且如果您要將結果處理程序限制於完全符合查詢的文件時，則不適用。例如，文件可能會提到五個人。然而，在語意搜尋查詢 `<sentence><person/>IBM</sentence>` 中，使用者只想知道與 `IBM` 一詞出現在同一個句子中的人員。使用者並不想知道其他人。

若要存取並處理完全符合查詢的特性結構：

1. 使用註解對映樣式，將相關的特性結構類型對映到企業搜尋索引。例如：

```
<indexBuildItem>
  <name>com.ibm.omnifind.types.Person</name>
  <indexRule>
    <style name="Annotation"/>
  </indexRule>
</indexBuildItem>
```

2. 將相關的特性結構類型對映到 JDBC 表格。在對映的程序中，您必須併入文件 URI 和特性結構 ID 的兩個直欄。雖然可以將所有特性結構類型對映到相同的資料庫表格，但應將每一個類型對映到不同的表格。例如：

```

<explicitMappingRule applyToSubtypes="false">
  <type>com.ibm.omnifind.types.Person</type>
  <table>sample.person</table>
  <featureMappings>
    <featureMapping>
      <feature>objectId()</feature>
      <column>primaryId</column>
    </featureMapping>
    <!-- Contains the covered text of the annotation-->
    <featureMapping>
      <feature>coveredText()</feature>
      <column>personName</column>
    </featureMapping>
    <!-- Other mapping go in here-->
    <!-- To access the relevant person annotations in the query result-->
    <featureMapping>
      <feature>docUri()</feature>
      <column>docUri</column>
    </featureMapping>
    <featureMapping>
      <feature>fsId()</feature>
      <column>annotationId</column>
    </featureMapping>
  </featureMappings>
</explicitMappingRule>

```

3. 搜索、剖析及檢索文件。
4. 擷取符合查詢的實例 ID。在搜尋及索引 API (SI-API) 中，這些實例稱為目標元素。目標元素指定要傳回的輸入跨距。它的定義如下：
 - 在 XML 片段中，利用附加在前面的 # 記號來識別目標元素。在 XML 片段查詢中只容許使用一個 # 記號，該記號可以出現在任何位置。例如：
`$xmlf2::'<sentence><#person/>IBM</sentence>'`
 - 預設在 XPath 中，目標元素是 XPath 表示式中的最後一個欄位。
 - 利用方法 `Result.getProperty("TargetElement")` 來存取這些實例。傳回的內容是字串連接，其中含有所有出現 ID 並以空格區隔。內容中的每一個出現項目都可以轉換成整數值。
5. SI-API 不會自行傳回特性結構，只會傳回其出現 ID。這些 ID 對應於儲存在資料庫表格的 `fsId()` 值。若要擷取這些實例及相關資訊，您的應用程式必須：
 - a. 依據目標元素的跨距名稱，選取正確的資料庫表格。在範例中，應用程式含有從人員到 `sample.Person` 表格的對映。此資訊是從索引對映 (生產跨距名稱) 及 JDBC 對映 (產生表格名稱) 的配置檔所推斷出來的。
 - b. 針對搜尋結果中的每一個結果物件：
 - 1) 剖析 `Result.getProperty("TargetElement")` 傳回的字串以在出現 ID 上尋找。
 - 2) 使用結果 URI (可以使用 `Result.getDocumentId()` 來存取) 作為 `docUri` 直欄中的值，以及出現 ID 作為 `annotationId` 直欄中的值，以發出表格的 SELECT 陳述式。直欄名稱需視對映的檔案而定。直欄名稱取自前一個範例。

傳回的列含有特性結構的儲存資訊 (例如，涵蓋的文字)，或特性結構的特定屬性 (如「姓氏」或「出生地」)。

請確定資料庫的更新與企業搜尋中的索引更新同步。如果資料庫含有過期的資訊 (例如，因為您使用資料庫載入檔，且沒有更新資料庫，但已重新整理或重組了索引)，則部份出現 ID 可能不會存在於資料庫中。企業搜尋只會在索引中保留最新文件版本的記錄。因此，出現 ID 只適用於最新的文件。

如果在相同的資料庫表格中儲存了同一份文件的多個版本，則可能會有多列符合相同的出現 ID，每一列代表不同的文件版本。在此情況下，您必須定義文件版本直欄，並使用應用程式邏輯或內建特性（如 `docTimestamp()`）來輸入資料。這樣，您就可以過濾結果，只取得最新的文件版本。

相關概念

- 2 第 44 頁的『語意搜尋查詢字詞』
- 2 語意搜尋查詢字詞是以不透明字詞來傳達。

相關工作

第 22 頁的『建立索引建置配置檔』
使用索引建置配置檔，您可以決定要在共用分析結構中檢索哪些分析結果以啟用搜尋。

第 28 頁的『建立 XML 對映配置檔』
若要將分析結果新增至資料庫，您必須建立配置檔，其中包含資料庫連線配置資訊的配置檔，以及自訂文字分析結果要儲存在哪些表格和直欄的說明。

定義於企業搜尋的類型及特性

定義於企業搜尋的類型系統包含文件中間資料處理及基本語言分析。

檢索文件時，一定會發生文件語言識別及斷詞法形式的基本語言分析，與是否選取自訂分析無關。在基本文件分析期間，會在共用分析結構中加入下列資訊，您可以在自訂分析中使用此資訊：

- 類型 `com.ibm.es.tt.DocumentMetaData` 的文件中間資料
- 類型 `uima.tt.TokenAnnotation`、`uima.tt.SentenceAnnotation` 及 `uima.tt.ParagraphAnnotation` 的記號、句子及段落註解。記號註解包含特性 `lemma`。

定義於企業搜尋的類型系統不含任何文字分析專用的精確類型和特性。這些均內含於 UIMA 類型系統中，您可以在 UIMA 環境中定義自訂分析類型及特性時使用它們並加以延伸。您可能不需要延伸企業搜尋類型系統。

企業搜尋類型系統未定義於「UIMA 軟體開發套件 (SDK)」中。如果要在 UIMA 中撰寫註解程式時使用其中的任何類型，例如，如果您要存取文件機密保護資訊，或存取搜索器類型或文件類型，則必須在分析引擎的類型系統說明中重新定義類型。

下列類型及特性定義於企業搜尋中：

uima.tcas.Annotation

註解包含下列類型：

uima.tcas.DocumentAnnotation

文件註解含有下列特性：

esDocumentMetaData

含有類型 `com.ibm.es.tt.DocumentMetaData` 的文件中間資料

com.ibm.es.tt.ContentField

內容欄位註解含有下列特性：

參數 類型 `com.ibm.es.tt.CommonFieldParameters` 的內容欄位參數。

com.ibm.es.tt.Anchor

HTML 文件中錨點文字的錨點註解。它含有下列特性：

uri 錨點文字的目標 URI。特性值屬於類型 `uima.cas.String`。

com.ibm.es.tt.MarkupTag

標記資訊註解，例如，XML 標示的註解。標記資訊儲存於下列特性中：

名稱 標記標示的名稱。特性值屬於類型 `uima.cas.String`。

depth 巢狀深度。特性值屬於類型 `uima.cas.Integer`。

attributeName

特性屬性的名稱。特性值屬於類型 `uima.cas.StringArray`。

attributeValues

屬性值字串。特性值屬於類型 `uima.cas.StringArray`。

uima.CAS.TOP

類型系統的 root。有下列類型：

com.ibm.es.tt.DocumentMetaData

文件中間資料含有下列特性。這些特性連接至文件註解特性 `esDocumentMetaData`。

crawlerId

搜索器名稱。特性值屬於類型 `uima.cas.String`。

dataSource

下列其中一項資料來源類型：

- Web (代表從「Web 搜索器」產生的文件)
- NNTP (代表從「新聞群組搜索器」產生的文件)
- DB2 (代表從「DB2 搜索器」產生的文件)
- Notes® (代表從「Notes 搜索器」產生的文件)
- CM (代表從「內容管理搜索器」產生的文件)
- FS (代表從「UNIX® 檔案系統搜索器」產生的文件)
- WinFS (代表從「Windows® 檔案系統搜索器」產生的文件)
- Exchange (代表從「Exchange 搜索器」產生的文件)
- VBR (代表從「VeniceBridge 搜索器」產生的文件)

特性值屬於類型 `uima.cas.String`。

dataSourceName

搜索器的名稱 (資料來源)。特性值屬於類型 `uima.cas.String`。

docType

下列其中一項文件類型：

- text/html
- application/postscript
- application/pdf
- application/x-mspowerpoint
- application/msword

- application/x-msexcel
- application/rtf
- application/vnd.lotus-wordpro
- application/x-lotus-123
- application/vnd.lotus-freelance
- text/xml
- text/plain
- application/x-js-taro (Ichitaro)

特性值屬於類型 `uima.cas.String`。

securityTokens

文件機密保護記號。特性值屬於類型 `uima.cas.StringArray`。

date 文件日期。特性值屬於類型 `uima.cas.String`。

baseUri

網頁的基本 URI。特性值屬於類型 `uima.cas.String`。

metaDataFields

特性值屬於類型 `uima.cas.FSArray`。這個陣列中的每一個元素都是屬於類型 `com.ibm.es.tt.MetadataField`。

redirectUri

重新導向的 URL。特性值屬於類型 `uima.cas.String`。

mimeType

MIME 類型或文件類型，例如 XML。特性值屬於類型 `uima.cas.String`。

url 文件 URL。特性值屬於類型 `uima.cas.String`。

com.ibm.es.tt.CommonFieldParameters

共用欄位參數包括：

searchable

一種旗號，指出欄位是否為可搜尋的任意文字。

可搜尋的欄位

一種旗號，指出欄位是否可當成欄位搜尋。

參數 指示參數搜尋的旗號。

showInSearchResult

指示註解資料是否併入搜尋結果明細的旗號。

resolveConflict

解決 `MetadataPreferred`、`ContentPreferred` 及 `Coexist` 之間中間資料衝突的旗號。特性值屬於類型 `uima.cas.String`。

名稱 欄位的名稱。您可以使用欄位名稱來搜尋此欄位。特性值屬於類型 `uima.cas.String`。

com.ibm.es.tt.MetadataField

中間資料欄位資料不是文件內容的一部份，但卻儲存在 "text" 特性中：

參數 類型 `com.ibm.es.tt.CommonFieldParameters` 的中間資料欄位參數。

text 中間資料文字儲存在這個類型 `uima.cas.String` 的特性中。

相關參考

- 2 『定義於 UIMA 的類型及特性』
- 2 「UIMA 軟體開發套件」定義在文字分析時可以在文件中探查到的基本語言類型及特性。

定義於 UIMA 的類型及特性

「UIMA 軟體開發套件」定義在文字分析時可以在文件中探查到的基本語言類型及特性。

每一個分析引擎都有自己的類型系統說明，說明分析引擎中註解程式的輸入需求和輸出類型。類型系統說明是網域及應用程式特有的。

您可以延伸 UIMA 類型系統以併入您自己的類型及特性。在 UIMA 環境中，提供 Eclipse 外掛程式以協助您編輯註解程式的類型系統描述子。如需安裝及使用「元件描述子編輯器」外掛程式的詳細資訊，請參閱 UIMA 文件。

當您在 UIMA 環境中完成開發及測試分析引擎時，您所建立且含有分析引擎檔案的保存檔（.pear 檔）也會併入您的類型系統說明。

下列類型及特性定義於 UIMA：

uima.tcas.Annotation

註解包含下列類型：

uima.tcas.DocumentAnnotation

uima.tt.TTAnnotation

uima.tcas.DocumentAnnotation

文件註解包含下列特性：

種類 文件的種類名稱或標籤清單。特性值屬於類型 `uima.cas.FSList`。

languageCandidates

文件語言參照清單。特性值屬於類型 `uima.cas.FSList`。

id 文件識別的格式，如 URL。特性值屬於類型 `uima.cas.String`。

uima.tt.TTAnnotation

TT 註解包括下列類型：

uima.tt.DocStructureAnnotation

有關文件的結構資訊。文件結構註解包括下列類型：

uima.tt.SentenceAnnotation

包含開啓及關閉標點的句子。其中含有特性：

sentenceNumber

段落中句子的序號。在每一個段落的開頭重設為 1。特性值屬於類型 `uima.cas.Integer`。

uima.tt.ParagraphAnnotation

段落。其特性包括：

paragraphNumber

段落的序號。特性值屬於類型 `uima.cas.Integer`。

uima.tt.LexicalAnnotation

有關文件的內容資訊。詞彙註解包含下列類型：

uima.tt.CompPartAnnotation

複合字的一部份。許多日耳曼語系的複合字是寫在一起而沒有任何區隔空格的字。例如，德文單字 "Abteilungsleiter" (部門經理) 是由組件 "Abteilung" (部門) 及 "Leiter" (經理) 所組成。

uima.tt.TokenAnnotation

前後沒有空格的記號。其特性包括：

lemmaEntries

指定記號的所有可能詞形清單。每一個項目都是記號的可能定義檔項目。

詞形 lemmaEntries 中記號的所有可能詞形清單中的某一詞形。搜尋時會使用這個詞形。

tokenNumber

句子中記號的序號。在每一個句子的開頭重設為 1。特性值屬於類型 `uima.cas.Integer`。

tokenProperties

記號內容 (舉例來說) 是大寫字體或數值的記號。特性值屬於類型 `uima.cas.Integer`。

stopwordToken

表示為停用字的記號。特性值屬於類型 `uima.cas.Integer`。

synonymEntries

類型 `uima.tt.Synonym` 的項目參照清單。每一個項目都是記號的可能同義字項目。

normalizedCoveredText

註解所涵蓋文字的標準化表示法。特性值屬於類型 `uima.cas.String`。

uima.CAS.TOP

類型系統的 root。有下列類型：

uima.tt.KeyStringEntry

具有下列特性的字串：

key 字串。

uima.tt.Lemma

具有下列形態資訊的定義檔項目：

partOfSpeech

詞形的詞性完整編碼。

morphID

詞性資訊的完整編碼。

uima.tt.Synonym

類型為 `uima.tt.keyStringEntry` 的指定單字之同義字項目。

uima.tt.LanguageConfidencePair

具有下列特性的類型，說明文件語言選擇。

uima.tt.LanguageConfidencePair

languageConfidence

說明選擇的語言實際上如何符合文件語言的一種指示 (0 到 1 之間的浮點值)。

language

文件的語言 (ISO 值)。值屬於類型 `uima.cas.String`。

languageID

語言 ID。值是屬於類型 `uima.cas.Integer`。

uima.tt.CategoryConfidencePair

具有下列特性的類型，說明文件的種類選擇。

uima.tt.CategoryConfidencePair

具有下列特性的種類：

categoryString

種類名稱。值屬於類型 `uima.cas.String`。

categoryConfidence

種類如何符合文件的一種指示。值是屬於浮點類型。

mostSpecific

指出種類是否為文件專用的旗號 (屬於類型 `uima.cas.Integer`)。

分類法 (taxonomy)

種類所屬的分類架構名稱。文件可以使用不同分類架構的種類。值屬於類型 `uima.cas.String`。

相關參考

2
2

第 38 頁的『定義於企業搜尋的類型及特性』

定義於企業搜尋的類型系統包含文件中間資料處理及基本語言分析。

語意搜尋應用程式

四種類型的文件資訊儲存在企業搜尋索引中，您可以使用搜尋及索引 API (SI-API) 介面，在搜尋應用程式中查詢此資訊。

四種不同類型的資訊包括：

- 在文件中找到的文字字組，例如，電腦軟體之類的詞組。
- 跨距名稱，例如，含有 `<author>James</author>` 的 XML 文件，會產生跨距 `<author>`。
- 屬性名稱，例如，含有 `<author countryOfBirth=USA>James</author>` 的 XML 文件，會產生屬性 "countryOfBirth"。
- 屬性值，例如，USA 是屬性 "countryOfBirth" 的值。

SIAPI 查詢語言包含語意搜尋查詢字詞。該字詞指定細枝型樣。細枝是有葉子的小樹。每一片葉子都代表四種類型的資訊 (文字字組、跨距名稱等等)。樹的內部節點指定它們在文件中的出現如何彼此關聯。指定關係的內部節點有五種類型：

- and
- or
- not
- in_the_span_of
- attribute_in_the_span_of

如果文件中有出現葉子，則可用來滿足指定的語意搜尋字詞，並遵守內部節點指定的限制 (已定義的關係)。

語意搜尋查詢字詞有助於擷取更好品質的文件。您現在不只能利用單字和註解的 Boolean 組合來搜尋。還可以擷取 (舉例來說) *James* 出現在跨距具名作者，或 *ibm* 及搜尋詞彙出現在相同句子的文件。

語意搜尋查詢字詞

語意搜尋查詢字詞是以不透明字詞來傳達。

有兩種語法形式可用來在搜尋及索引 API (SIAPI) 中表示不透明字詞：

- XML 片段
- 限制的 XPath

XML 片段查詢字詞類似 XML 文件中平衡良好的片段。XML 片段查詢字詞的字首是不透明字詞符號 `@xmlf2::`，後面接著以單引號括住的 XML 片段表示式 ('...')。

然而，限制的 XPath 查詢字詞字首是 `@xmlxp::`，後面接著以單引號括住的 XPath 查詢 ('...')。

在搜尋及索引 API (SIAPI) 介面中使用一般查詢字詞時，每一個字詞都會有外觀修飾元：

加號 (+)

字詞必須出現。

字首 =

字詞必須完全相符。

字首 ~ 字元

考慮查詢字詞的同義字。

字尾 ~ 字元

考慮詞形與查詢字詞相同的單字。

下列範例顯示 XML 片段查詢。

@xmlf2::'<City>Springfield</City>'

尋找跨距 (註解) 城市含有字串 Springfield 的文件。

@xmlf2::'<Person gender="female">'

尋找其中註解女性人員的文件。

```
@xmlf2::'<Person><.or><@gender>female</@gender>  
<@title>Mrs</@title><@title>Ms</@title></or></Person>'
```

尋找依據性別或職稱指定人員為女性的文件。

```
@xmlf2::'<Person gender="male" role="suspect"/>  
<PoliceReport><@crimeDescription><.or>robbery theft</or>-accident  
</@crimeDescription></PoliceReport> <City>Springfield<.or>  
<@district>Brynston</@district><@district>Brooklyn</@district></or></City>'
```

尋找指定男性為嫌犯的文件，以及 policeReport 註解，其中的 crimeDescription 含有字串 robbery 及 theft，而不是字串 accident。文件也必須包含 Brynston 及 Brooklyn 區的城市註解。

相對應的 XPath 查詢含有下列結構：

```
@xmlxp::'//City[ftcontains(Springfield)]'
```

尋找跨距 (註解) 城市含有字串 Springfield 的文件。

```
@xmlxp::'//Person[@gender="female"]'
```

尋找其中註解女性人員的文件。

```
@xmlxp::'//Person[@gender="female" or @title ftcontains(Ms) or @title  
ftcontains(Mrs)]'
```

尋找依據性別或職稱指定人員為女性的文件。

```
@xmlxp::'//Person[@gender="male" and @role="suspect"] //PoliceReport  
[@crimeDescription ftcontains(robbery) or @crimeDescription ftcontains(theft)]  
//City [ (@district="Brynston" or @district="Brooklyn") and  
ftcontains(Springfield)]'
```

尋找指定男性為嫌犯的文件，以及 policeReport 註解，其中的 crimeDescription 含有字串 robbery 及 theft。文件也必須包含 Brynston 及 Brooklyn 區的城市註解。

搜尋應用程式中的同義字支援

使用者可以搜尋含有查詢詞彙同義字的文件，來擴充搜尋結果。

同義字通常包含多字詞彙，如 *WebSphere Information Integrator OmniFind* 之類的產品名稱。同義字定義檔所含的多字詞彙可在使用者查詢中正確地識別，而不需要以引號括住。

企業搜尋的「搜尋及索引 API (SI-API)」支援多種方法，以便使用者搜尋查詢字詞的同義字：

- SI-API 查詢語法支援以波浪符號 (~) 運算子，來延伸同義字。如果使用者在查詢詞彙開頭加上此運算子，則會延伸到該字詞的同義字。例如，查詢 ~WAS 會傳回討論 WebSphere Application Server，以及此縮寫的其他任何同義字的文件。
- 搜尋應用程式內可以使用 SI-API 同義字延伸介面，來啟用同義字延伸。查詢詞彙可以自動擴充為納入同義字，或者，搜尋應用程式可能包含選項，讓使用者指定搜尋結果是否傳回查詢詞彙的同義字。

在自動延伸同義字期間，會在所有查詢字詞和內容欄位上執行同義字查閱。搜尋結果會傳回含有查詢詞彙或其同義字的文件。搜尋結果也會顯示哪些詞彙已延伸到哪些同義字。

在以使用者主導的情況下，實際執行查詢之前，搜尋應用程式會向使用者顯示每一個查詢字詞已找到哪些同義字。然後，使用者可以選取哪些詞彙要併入搜尋，或重新定義搜尋來移除最初的查詢詞彙。在此實務範例中，使用者可以控制要在查詢中併入哪些詞彙：精確的同義語，或變化詞性及用法。

建立同義字的 XML 檔案

若要將企業搜尋中的查詢擴增為包括查詢字詞的同義字，必須在 XML 檔案中，指定可以彼此做為同義字的字組。

關於本作業

列出同義字的 XML 檔案，必須符合下列範例中顯示的綱目。

```
<?xmlversion="1.0" encoding="utf-8"?>
<synonymgroups xmlns="http://www.ibm.com/of/822/synonym/xml">
  <synonymgroup>
    <synonym>Think Pad</synonym>
    <synonym>Notebook</synonym>
    <synonym>Notebooks</synonym>
  </synonymgroup>
  <synonymgroup>
    <synonym>WebSphere Application Server</synonym>
    <synonym>WAS</synonym>
  </synonymgroup>
</synonymgroups>
```

限制

您必須將彼此是同義字的字組 (<synonym> 元素)，組織在 <synonymgroup> 元素中。同義字可以包括空格字元，但不可以包括標點字元，例如逗點 (,) 或垂直線 (|)，因為這些字元可能會影響企業搜尋查詢語法。

您必須列舉所有新增為同義字的可能詞彙變化，例如字組的單數和複數形式。您不需要列舉詞彙的正規化，例如移除重音或曲音 (企業搜尋會自動處理正規化)。比方說，如果要將詞彙 météo 加入為同義字，並不需要將詞彙 METEO 也加入。

程序

若要建立同義字清單，供進行企業搜尋：

1. 建立 XML 檔案。若要避免 XML 語法錯誤，請選擇使用 XML 編輯器或 XML 編寫工具。
2. 新增一個 <synonymgroup> 元素，然後對要視為同義字群組中其他字組之同義字的每個字組，插入 <synonym> 項目。

請確定在 <synonymgroups xmlns="http://www.ibm.com/of/822/synonym/xml"> 元素中包括對映。名稱空間 (在 xmlns 屬性中指定) 需要如所示範例一模一樣。

3. 重複之前的步驟，直到指定好要用來在企業搜尋集中搜尋文件的所有同義字。
4. 儲存並結束 XML 檔案。

建立好 XML 檔案之後，必須將該檔案轉換為同義字定義檔，才能新增到企業搜尋系統。

建立同義字定義檔

在 XML 檔案中建立或更新同義字清單之後，必須將 XML 檔案轉換為同義字定義檔。

關於本作業

若要建立同義字定義檔，請使用 WebSphere II OmniFind Edition 隨附的命令行工具 `essyndictbuilder`。這個工具是在 `ES_INSTALL_ROOT/bin` 目錄中。

工具的輸入是列出同義字的 XML 檔案，工具的輸出就是同義字定義檔。定義檔必須具有字尾 `.dic`。例如，`c:\mydictionaries\products.dic`。

這兩個檔案的預設位置，就是呼叫 Script 的目錄。如果已存在相同名稱的定義檔，Script 會產生錯誤。

程序

若要建立企業搜尋的同義字定義檔：

1. 以企業搜尋管理員的身份登入索引伺服器。此使用者 ID 是在安裝 WebSphere II OmniFind Edition 時指定。
2. 輸入下列命令，其中的 `XML_file` 是含有同義字清單之 XML 檔案的完整路徑，`DIC_file` 是同義字定義檔的完整路徑。

AIX、Linux 或 Solaris：`essyndictbuilder.sh XML_file DIC_file`

Windows：`essyndictbuilder.bat XML_file DIC_file`

建立好同義字定義檔後，使用企業搜尋管理主控台，將定義檔新增到企業搜尋系統，並將它與一或多個集合連結。

只有產生的 .dic 檔案，才會上載到企業搜尋系統。請確定將原始的 XML 檔案放在實施存取控制的環境中，並且確定您正常地備份檔案。您以後要更新同義字定義檔時，將會需要這個 XML 檔案。

自訂停用字定義檔

使用者可以定義要從查詢中移除的企業特有詞彙，以增進搜尋關聯。

企業搜尋提供兩種停用字支援：

- 語言特有的停用字識別，可以從多字查詢中移除所有常用的通用字，如 *a* 及 *the*。使用者無法修改每一個語言的停用字定義檔。這個停用字識別會自動在所有查詢中執行，以增進搜尋關聯性。
- 使用者定義或自訂停用字識別，可以從查詢中移除企業特有詞彙。這個停用字定義檔是由管理者定義，其中只含特殊詞彙。使用者定義的停用字定義檔並不會取代含有通用字的企業搜尋語言特有停用字定義檔。

使用者定義的停用字通常含有多字詞彙，如 *WebSphere Information Integrator OmniFind* 之類的產品名稱。停用字定義檔所含的多字詞彙可在使用者查詢中正確地識別，而不需要以引號括住。

德文中的複合字詞可以在查詢中正確識別。複合字詞是當作單一字使用之兩個或多個字的組合。詞彙化複合，如 *Reisebüro* (旅行社) 不會被視為複合。

會分解查詢中的複合字詞為組成複合的個別字詞。如果構成複合字的任何個別詞彙出現在停用字定義檔中，則不會從查詢中移除該複合字詞。

例如，查詢字詞 *Versicherungspolice* (保險原則) 會傳回含有複合字詞 *Lebensversicherungspolice* (人壽保險原則) 及 *Haftpflichtversicherungspolice* (第三方保險原則) 的文件。而在 *Haftpflicht* (第三方保險) 之類的查詢中，也會傳回第二個詞彙。即使單字 *Police* 列示在停用字定義檔中，也不會從查詢中移除複合查詢字詞 *Versicherungspolice*。

您必須在 XML 檔中列出企業特有詞彙，然後將該檔案轉換成停用字定義檔，才能將它新增至企業搜尋系統。

您可以在企業搜尋管理主控台上選取要使用哪一個停用字定義檔。您可以為每一個集合選取一個停用字定義檔。而一個停用字定義檔可由數個集合共用。

建立停用字的 XML 檔案

若要從查詢中移除企業特有的詞彙，您必須指定哪些單字在 XML 檔中定義為停用字。

關於本作業

列出停用字的 XML 檔必須符合下列範例所顯示的綱目。

```
<?xml version="1.0" encoding="UTF-8"?>
<stopWords xmlns="http://www.ibm.com/of/83/stopwordbuilder/xml">
  <stopWord>OmniFind Edition</stopWord>
  <stopWord>WAS</stopWord>
  <stopWord>...</stopWord>
</stopWords>
```

限制

停用字可以包含空格字元，但不能含有標點字元，如逗點 (,) 或垂直線 (|)，因為這些字元可能會影響企業搜尋查詢語法。

您不需要列舉詞彙的正規化，例如移除重音或曲音 (企業搜尋會自動處理正規化)。比方說，如果您要將詞彙 `météo` 併入為停用字，則不需要同時加入詞彙 `METEO`。

程序

若要建立企業搜尋的停用字清單：

1. 建立 XML 檔案。若要避免 XML 語法錯誤，請使用可以驗證 XML 的 XML 編輯器或 XML 編寫工具。
2. 針對每一個要當成停用字的單字，新增 `<stopWord>` 元素。

請務必在 `<stopWords xmlns="http://www.ibm.com/of/83/stopwordbuilder/xml">` 元素中併入對映。名稱空間 (在 `xmlns` 屬性中指定) 需要如所示範例一模一樣。

3. 重複之前的步驟，直到已指定使用者搜尋企業搜尋集合時要從查詢中移除的所有停用字為止。
4. 儲存並結束 XML 檔案。

建立 XML 檔後，您必須將它轉換成停用字定義檔，然後才能新增至企業搜尋系統。

建立停用字定義檔

在 XML 檔中建立或更新使用者定義的停用字清單後，您必須將 XML 檔轉換成停用字定義檔。

關於本作業

若要建立停用字定義檔，請使用 WebSphere II OmniFind Edition 隨附的命令行工具 `esstopworddictbuilder`。這個工具是在 `ES_INSTALL_ROOT/bin` 目錄中。

工具的輸入是列出停用字的 XML 檔，而工具的輸出就是停用字定義檔。定義檔必須具有字尾 `.dic`。例如，`c:\mydictionaries\productstopwords.dic`。

這兩個檔案的預設位置，就是呼叫 `Script` 的目錄。如果已存在相同名稱的定義檔，`Script` 會產生錯誤。

程序

若要建立企業搜尋的停用字定義檔：

1. 以企業搜尋管理員的身份登入索引伺服器。此使用者 ID 是在安裝 WebSphere II OmniFind Edition 時指定。
2. 輸入下列命令，其中 `XML_file` 是含有停用字清單的 XML 檔完整路徑，而 `DIC_file` 是停用字定義檔的完整路徑。

AIX、Linux 或 Solaris：

```
esstopworddictbuilder.sh XML_file DIC_file
```

```
Windows: esstopworddictbuilder.bat XML_file DIC_file
```

建立停用字定義檔後，請使用企業搜尋管理主控台將定義檔新增至企業搜尋系統，並建立它與一或多個集合的關聯性。

| 只有產生的 .dic 檔案，才會上載到企業搜尋系統。請確定將原始的 XML 檔案放在實
| 施存取控制的環境中，並且確定您正常地備份檔案。您需要此 XML 檔案以更新停用字
| 定義檔。

自訂 Boost 字定義檔

使用者可以定義特定的詞彙或多字詞彙，以提高或降低出現該詞彙的文件排序值。

Boost 定義檔中的每一個詞彙都與 Boost 因數相關聯，範圍從 -10 到 +10。您特別想在結果文件中看到的詞彙會配置較高的 Boost 因數，而那些完全不想讓它出現或與較高優先詞彙合併的那些詞彙，則會有較低的指定值。值 -1、0 及 1 沒有任何增值效果。

如果列示在 Boost 定義檔且具有特定 Boost 因數的查詢字詞出現在擷取的文件中，則文件排序值會視 Boost 值而升高或下降。指定給詞彙的 Boost 值是相對性的，因為它也會受其他因數影響。因此，如果詞彙 X 是以 B1 啟動且詞彙 Y 是以 B2 啟動，而 $B1 > B2$ ，則 $\text{boost}(X) \geq \text{boost}(Y)$ 。

Boost 字通常含有多字詞彙，如 *WebSphere Information Integrator OmniFind* 之類的產品名稱。Boost 字定義檔所含的多字詞彙可在使用者查詢中正確地識別，而不需要以引號括住。

德文中的複合字詞可以在查詢中正確識別。複合字詞是當作單一字使用之兩個或多個字的組合。詞彙化複合，如 *Reisebüro* (旅行社) 不會被視為複合。

會分解查詢中的複合字詞為組成複合的個別字詞。如果組成複合字詞的個別詞彙有 Boost 值，則即使在詞彙不是複合字詞的一部份時，指定的值會低於該詞彙的 Boost 值，還是會排序擷取的文件。這會擴大搜尋範圍，而在只找到少數文件含有完整複合字詞的情況下，這是非常有用的。

例如，查詢字詞 *Versicherungspolice* (保險原則) 會傳回含有複合字詞 *Lebensversicherungspolice* (人壽保險原則) 及 *Haftpflchtversicherungspolice* (第三方保險原則) 的文件。而在 *Haftpflcht* (第三方保險) 之類的查詢中，也會傳回後者。如果 Boost 字定義檔含有單字 *Police* (原則)，則含有複合查詢字詞 *Versicherungspolice* 的文件會有指定的 Boost 值。

您必須在 XML 檔中列出詞彙及其 Boost 值，然後將該檔案轉換成 Boost 字定義檔，才能將它新增至企業搜尋系統。

您可以在企業搜尋管理主控台上選取要使用哪一個 Boost 字定義檔。每一個集合可以選取一個 Boost 字定義檔。而一個 Boost 字定義檔可由數個集合共用。

建立 Boost 字的 XML 檔

若要提高或降低某些結果文件的重要性，您必須在 XML 檔中指定哪些單字會影響文件排序。

關於本作業

列出 Boost 字的 XML 檔必須符合下列範例所顯示的綱目。

```
<?xml version="1.0" encoding="UTF-8"?>
<boostTerms xmlns="http://www.ibm.com/of/83/boostbuilder/xml">
  <!-- group boost terms by boost value-->
  <boostTermList boost="5">
```



```

|         <!-- each term can specify the synonym expansion separatly-->
|         <term useVariants="true">OmniFind Edition</term>
|         <term useVariants="false">Edition</term>
|         <term useVariants>OmniFind</term>
|     </boostTermList>
|     <boostTermList boost="8">
|         <term useVariants="true">WAS</term>
|         <term>term9</term>
|     </boostTermList>
| </boostTerms>

```

限制

您必須將 `<boostTermList>` 元素中，共用相同 Boost 值的詞彙分組。Boost 值可以發生多次，例如，如果您要在 XML 檔中依字母順序排序 Boost 字。

Boost 字可以包含空格字元，但不能含有標點字元，如逗點 (,) 或垂直線 (|)，因為這些字元可能會影響企業搜尋查詢語法。

Boost 詞彙通常有變體，如字首語或縮寫。您可以在 Boost 字定義檔中列舉所有變體；但是，如果您計劃使用同義字定義檔及 Boost 字定義檔，且已在同義字定義檔中新增詞彙及其變體，則不需要將這些變體也加入 Boost 字清單。而只需要針對新增至 Boost 字定義檔的變體，將屬性 `useVariants` 設為 `true` 即可。如果這個詞彙列示在同義字定義檔中，則它在任何擷取文件中所發生的所有變體都會影響指定給這些文件的排序值。

您不需要列舉詞彙的正規化，例如移除重音或曲音 (企業搜尋會自動處理正規化)。比方說，如果您要將詞彙 `météo` 併入為 Boost 字，則不需要同時加入詞彙 `METEO`。

程序

若要建立企業搜尋的 Boost 字清單：

1. 建立 XML 檔案。若要避免 XML 語法錯誤，請選擇使用 XML 編輯器或 XML 編寫工具。
2. 在 `<boostTerms xmlns="http://www.ibm.com/of/83/boostbuilder/xml">` 元素中併入對映。名稱空間 (在 `xmlns` 屬性中指定) 需要如所示範例一模一樣。
3. 新增 `<boostTermList>` 元素，以分組所有共用指定 Boost 值的詞彙。

Boost 值的範圍可以從 -10 到 10。例如，`<boostTermList boost="-5">` 或 `<boostTermList boost="5">`。

含有指定詞彙的文件重要性會依指定的 Boost 值而提高或降低。

4. 針對使用指定 Boost 值的每一個詞彙，新增 `<term>` 元素。

如果要併入列示在同義字定義檔中 Boost 字的變體，請將 `<term>` 元素的 `useVariants` 屬性設為 `true`。預設值是 `false`。如果在同義字定義檔中找不到任何變體，則不會產生任何錯誤訊息。

5. 請重複之前的步驟，直到已指定使用者搜尋企業搜尋集合時要當成 Boost 字使用的所有詞彙為止。
6. 儲存並結束 XML 檔案。

建立 XML 檔後，您必須將它轉換成 Boost 字定義檔，然後才能新增至企業搜尋系統。

建立 Boost 字定義檔

在 XML 檔中建立或更新 Boost 字清單後，您必須將 XML 檔轉換成 Boost 字定義檔。

關於本作業

若要建立 Boost 字定義檔，請使用 WebSphere II OmniFind Edition 隨附的命令行工具 `esboostworddictbuilder`。這個工具是在 `ES_INSTALL_ROOT/bin` 目錄中。

工具的輸入是列出 Boost 字的 XML 檔，而工具的輸出就是 Boost 字定義檔。定義檔必須具有字尾 `.dic`。例如，`c:\mydictionaries\productboostwords.dic`。

這兩個檔案的預設位置，就是呼叫 `Script` 的目錄。如果已存在相同名稱的定義檔，`Script` 會產生錯誤。

程序

若要建立企業搜尋的 Boost 字定義檔：

1. 以企業搜尋管理員的身份登入索引伺服器。此使用者 ID 是在安裝 WebSphere II OmniFind Edition 時指定。
2. 輸入下列命令，其中 `XML_file` 是含有 Boost 字清單的 XML 檔完整路徑，而 `DIC_file` 是 Boost 字定義檔的完整路徑。如果您也想使用同義字定義檔，請在 Boost 定義檔名稱後加上同義字定義檔的完整路徑。同義字定義檔的命名是選用的。

UNIX：

```
esboostworddictbuilder.sh XML_file DIC_file SYNDIC_file
```

Windows：esboostworddictbuilder.bat XML_file DIC_file SYNDIC_file

建立 Boost 字定義檔後，請使用企業搜尋管理主控台將定義檔新增至企業搜尋系統，並建立它與一或多個集合的關聯性。

只有產生的 `.dic` 檔案，才會上載到企業搜尋系統。請確定將原始的 XML 檔案放在實施存取控制的環境中，並且具有適當的備份措施。您需要這個 XML 檔來更新 Boost 字定義檔。

相關工作

第 48 頁的『建立同義字定義檔』

在 XML 檔案中建立或更新同義字清單之後，必須將 XML 檔案轉換為同義字定義檔。

企業搜尋中包括的文字分析

企業搜尋所含的文字分析包括文件語言偵測及斷詞法。

處理文件時，企業搜尋會判定該文件的語言，然後將輸入文字的串流分成不同的記號單元。

在搜尋期間，使用者或應用程式必須手動選取查詢語言。查詢字串會在索引中分斷、分析及搜尋。

文件及查詢字串分析可以分成：

- 基本非定義檔型支援。這包括空格及 n-gram 斷詞法。
- 定義檔型語言支援。這包括單字和句子斷詞法，以及詞形還原。

語言處理涉及詞彙分析，也就是建立輸入文字替代表示法的程序，該輸入文字會將所有可用的定義檔資料關聯至輸入文字中所辨識的記號。使用進階語言處理程序，可以大幅加強搜尋品質。

相關概念

『語言識別』

企業搜尋必須先決定原始文件的語言，然後才能發生單字和句子斷詞法、字元正常化或詞形還原。

第 60 頁的『非定義檔型斷詞法的語言支援』

如果語言偵測及詞彙分析技術不支援文件的語言，則企業搜尋會提供 Unicode 型空格及 n-gram 斷詞法形式的基本支援。

語言識別

企業搜尋必須先決定原始文件的語言，然後才能發生單字和句子斷詞法、字元正常化或詞形還原。

企業搜尋可以自動偵測下列語言：

阿拉伯文	法文	韓文
中文 (繁體及簡體)	德文	波蘭文
捷克文	希臘文	葡萄牙文
丹麥文	希伯來文	俄文
荷蘭文	匈牙利文	西班牙文
英文	義大利文	瑞典文
芬蘭文	日文	土耳其文

企業搜尋中的語言處理程序會在檢索期間 (而非查詢處理程序期間) 偵測原始文件的語言。

在企業搜尋中，您可以指定自動偵測文件的語言或選取要使用的語言。

如果選取自動語言偵測，而剖析器無法判定文件的語言，則剖析器會使用您在企業搜尋管理主控台建立搜索器時指定的語言。

如果您沒有選取自動語言偵測，則一律使用您指定的語言。預設值是英文。

您可以使用基本語言獨立技術 (如空格斷詞法及 n-gram 斷詞法)，處理沒有語言專用定義檔的文件。

企業搜尋語言偵測技術最適合用於單語文件。如果文件使用多語言，則會嘗試判定文件中所使用的最主要語言。然而，分析結果不一定是讓人滿意的。

文件的語言可以用來將您的搜尋結果限制於特定語言的文件。比方說，如果搜尋有關 Jacques Chirac (席哈克) 的文件，您可以指定只在搜尋結果中併入以法文撰寫的文件。

相關概念

第 59 頁的『企業搜尋中包括的文字分析』

企業搜尋所含的文字分析包括文件語言偵測及斷詞法。

『非定義檔型斷詞法的語言支援』

如果語言偵測及詞彙分析技術不支援文件的語言，則企業搜尋會提供 Unicode 型空格及 n-gram 斷詞法形式的基本支援。

非定義檔型斷詞法的語言支援

如果語言偵測及詞彙分析技術不支援文件的語言，則企業搜尋會提供 Unicode 型空格及 n-gram 斷詞法形式的基本支援。

Unicode 型空格斷詞法

這個語言處理方法在單字之間使用空格 (或空白空間) 作為單字區隔字元。

N-gram 斷詞法

這個語言處理方法將 n 個字元的重疊順序視為單一個字。在許多擷取作業中，這個簡單的斷詞法方法就已經夠用了。

這些方法與任何語言定義檔無關，且不涉及更準確的語言處理技術，如基礎詞形還原。

N-gram 斷詞法適用於不使用空格作為區隔字元的語言，如泰文。相同的方法也適用於希伯萊文及阿拉伯文。雖然這兩種語言都使用空格區隔字元，但 n-gram 斷詞法傳回的結果會比 Unicode 型空格斷詞法的基本形更好。

相關概念

第 59 頁的『企業搜尋中包括的文字分析』

企業搜尋所含的文字分析包括文件語言偵測及斷詞法。

第 59 頁的『語言識別』

企業搜尋必須先決定原始文件的語言，然後才能發生單字和句子斷詞法、字元正常化或詞形還原。

定義檔型斷詞法的語言支援

如果正確地偵測到文件的語言且可以使用語言專用定義檔，則會套用適當的語言處理程序。

斷詞法是將輸入文件分成不同詞元的程序。此程序包括下列部份語言處理活動：

斷詞 (Word segmentation)

斷詞用於不在單字之間使用空格 (或區隔字元) 的語言，如日文和中文。

詞形還原 (Lemmatization)

詞形還原是一種語言處理形式，決定文字中每一個字形的詞形。單字的詞形包含了基礎詞形，並加上共用相同詞類部份的變化形。例如，*go* 的詞形包含 *go*、*goes*、*went*、*gone* 及 *going*。名詞的詞形分成單數和複數形 (如 *calf* 及 *calves*)。形容詞的詞形分成比較級和最高級形式 (如 *good*、*better* 及 *best*)。代名詞的詞形分成相同名詞的不同格 (如 *I*、*me*、*my* 及 *mine*)。

詞形還原需要檢索和搜尋的定義檔。

企業搜尋會檢索詞形和變化形，並還原查詢中所有變化形。詞形還原會尋找含有查詢中變化形變體的文件，以加強搜尋品質。例如，當查詢含有單字 *mouse* 時，會尋找含有單字 *mice* 的文件。

縮寫分割 (Contraction splitting)

識別縮詞並將它們分割成元件組件，以增進搜尋品質。例如：

wouldn't 分割成 *would* + *not*

Horse's 分割成 *Horse* + *is* 或 *'s* (以說明查詢的模糊性)

附著語素識別

附著語素是特殊的縮詞形式，可判定元件組件來增進搜尋品質。附著語素是作用類似詞綴及單字的一種元素。然而，附著語素很難識別，因為它們也是字形的一部份。與其他語素 (單字結構) 現象不同，附著語素發生在語法結構中，而它們與單字的連接並不是字形規則的一部份。例如：

reparti-lo-emos 分成元件 *repartir* + *lo* + *emos*

l'avenue 分成元件 *le* + *avenue*

dell'arte 分成元件 *dello* + *arte*。

非字母型字元識別

語言處理可辨識非字母型字元。依據內部語言相關邏輯，部份非字母型字元會當成不同類型的個別詞元傳回，而部份則會加以分組。

例如，撇號或連字號在附著語素中會被視為單字組件，而在不明縮寫的情況下，則會被視為句點 (句號)。語言處理會將部份特殊的字元順序辨識成記號，例如 URL、電子郵件位址及日期。

縮寫識別 (Abbreviation recognition)

語言處理會將定義檔中的縮寫辨識為一個詞元。如果縮寫不在定義檔中，則縮寫會被辨識為詞彙項目，但該縮寫將不會有任何相關的定義檔資訊。

正確地辨識縮寫是句子識別的重點。例如，縮寫結尾的句點不一定是句子的結尾。

句尾標記識別 (End-of-sentence marker recognition)

語言處理可正確地識別句子斷詞法的句尾標記。

定義檔型語言支援適用於下列語言：

阿拉伯文
中文 (簡體及繁體)
捷克文
丹麥文
荷蘭文
英文

義大利文
日文
韓文
挪威文 (巴克摩及耐諾斯克)
波蘭文
葡萄牙文 (國家及巴西)

芬蘭文
法文 (國家及加拿大)
德文 (國家及瑞士)
希臘文

俄文
西班牙文
瑞典文

相關概念

『日文的斷詞』

如果要以日文辨識純文字文件或查詢字串，則企業搜尋會利用最適用於日文的形態分析技術，執行相關的斷詞。

『日文的正體字變體』

日文使用許多正體字變體。Katakana 變體是最重要的，因為 Katakana 通常是用於外國字的拼寫和發音。日文經常使用許多 Katakana 變體。

日文的斷詞

如果要以日文辨識純文字文件或查詢字串，則企業搜尋會利用最適用於日文的形態分析技術，執行相關的斷詞。

這個最佳化的範例是單字分解。日文使用大量的複合字。這些字會分解成最佳大小的記號，以達到較好的搜尋結果。同時也會分解變化形和前置詞以增進搜尋效能。

相關概念

第 60 頁的『定義檔型斷詞法的語言支援』

如果正確地偵測到文件的語言且可以使用語言專用定義檔，則會套用適當的語言處理程序。

『日文的正體字變體』

日文使用許多正體字變體。Katakana 變體是最重要的，因為 Katakana 通常是用於外國字的拼寫和發音。日文經常使用許多 Katakana 變體。

日文的正體字變體

日文使用許多正體字變體。Katakana 變體是最重要的，因為 Katakana 通常是用於外國字的拼寫和發音。日文經常使用許多 Katakana 變體。

企業搜尋使用變體定義檔將一般的 Katakana 變體對映至其基本形式 (類似詞形)，以便能找到所有文件，包括含有查詢字串中 Katakana 單字之正體字變體的文件。

企業搜尋也支援一般的 Okurigana 變體，這是以平假名結尾的漢字單字。

相關概念

第 60 頁的『定義檔型斷詞法的語言支援』

如果正確地偵測到文件的語言且可以使用語言專用定義檔，則會套用適當的語言處理程序。

『日文的斷詞』

如果要以日文辨識純文字文件或查詢字串，則企業搜尋會利用最適用於日文的形態分析技術，執行相關的斷詞。

停用字移除

在企業搜尋中，會從多字查詢中移除所有停用字 (例如，*a* 及 *the* 之類的通用字)，以增加搜尋效能。

日文的停用字識別是以文法資訊為基礎。例如，企業搜尋會辨識單字是否為名詞或動詞，而其他語言則是使用特殊清單。

相關概念

『字元正常化』

字元正常化是可以增進檢索率的一種程序。使用字元正常化增進檢索率，表示即使文件不完全符合查詢，還是可以擷取更多文件。

字元正常化

字元正常化是可以增進檢索率的一種程序。使用字元正常化增進檢索率，表示即使文件不完全符合查詢，還是可以擷取更多文件。

企業搜尋使用 Unicode 相容性正常化，其中包含亞洲全形和半形字元的正常化。

例如，日文的全形英數字元可以正常化成半形字元，半形的 **Katakana** 字元可以正常化為全形字元等等。企業搜尋也可移除 **Katakana** 中間點（當成日文的複合字區隔字元使用）。

字元正常化的其他形式包括：

大小寫正常化

例如，搜尋 *usa* 時會尋找含有 *USA* 的文件。

曲音符號擴充

例如，搜尋 *schön* 時會尋找含有 *schoen* 的文件。

重音符號移除

例如，搜尋 *e* 時會尋找含有 *é* 的文件。

其他區別發音符號移除

例如，搜尋 *c* 時會尋找含有 *ç* 的文件。

連字擴充

例如，搜尋 *ae* 時會尋找含有 *Æ* 的文件。

所有正常化都適用於兩種方式。您可以在搜尋 *USA* 時尋找含有 *usa* 的文件，在搜尋 *é* 時尋找單字中有 *e* 的文件等等。這些正常化也可以合併。例如，搜尋 *METEO* 時尋找含有 *météo* 的文件。

正常化是以 Unicode 字元內容為基礎，和語言無關。例如，企業搜尋支援希伯萊文的區別發音符號移除，以及阿拉伯文的連字擴充。

相關概念

第 62 頁的『停用字移除』

在企業搜尋中，會從多字查詢中移除所有停用字（例如，*a* 及 *the* 之類的通用字），以增加搜尋效能。

企業搜尋文件

您可以閱讀 PDF 或 HTML 版的 WebSphere Information Integrator OmniFind Edition (企業搜尋) 文件。

WebSphere Information Integrator OmniFind Edition 安裝程式可自動安裝資訊中心。安裝程式會在搜尋伺服器上安裝資訊中心。如果沒有安裝資訊中心，則當您按一下說明時，會開啓 IBM 網站上的資訊中心。若要查看企業搜尋的 HTML 主題，請啓動資訊中心。

若要查看 PDF 文件，請移至 docs/locale/pdf。例如，若要尋找英文版文件，請移至 docs/en_US/pdf。您也可以 WebSphere Information Integrator OmniFind Edition 支援網站上檢視最新的 PDF 文件。

下表顯示可用的文件、檔名及位置。

表 9. 企業搜尋的 PDF 文件

標頭	標頭	標頭
<i>Installation Guide for Enterprise Search</i> (本書的主題也可在資訊中心中取得)	iiysi.pdf	docs/locale/pdf/
<i>Administering Enterprise Search</i> (本書的主題也可在資訊中心中取得。)	iiysa.pdf	docs/locale/pdf/
<i>Programming Guide and API Reference for Enterprise Search</i> (本書的主題也可在資訊中心中取得。)	iiysp.pdf	docs/locale/pdf/
<i>Messages for Enterprise Search</i> (本書的主題也可在資訊中心中取得。)	iiysm.pdf	docs/locale/pdf/
<i>Installation Requirements for Enterprise Search</i> (本書的主題也可在資訊中心中取得。)	iiysr.txt 或 iiysr.htm	docs/locale/ (此檔案也可以從「第一個步驟」程式啓動。)
版本注意事項	iiysn.pdf	只能在 IBM WebSphere Information Integrator OmniFind Edition 文件網站上取得。
WebSphere Information Integrator 資訊中心	不適用	

WebSphere II OmniFind Edition 協助工具

可以存取 IBM WebSphere Information Integrator OmniFind Edition 使用者介面及文件。

安裝程式

您可以使用快速鍵在 WebSphere II OmniFind Edition 安裝程式中移動。下表說明部份快速鍵。

表 10. 安裝程式的鍵盤快速鍵

動作	快速鍵
標示圓鈕	方向鍵
選取圓鈕	Tab 鍵
標示按鈕	Tab 鍵
選取按鈕	Enter 鍵
前往下一個或前一個視窗，或者取消	按下 Tab 鍵以標示出某一按鈕，然後按 Enter 鍵
使作用中視窗停用	Ctrl + Alt + Esc

企業搜尋管理主控台及資訊中心

管理主控台及資訊中心介面是瀏覽器型介面，您可以使用 Microsoft® Internet Explorer 或 Mozilla FireFox 來檢視。如需瀏覽器的快速鍵及其他協助工具特性清單，請參閱 Internet Explorer 或 FireFox 的線上說明。

PDF 文件

您可以檢視所有 PDF 格式的企業搜尋文件。只要使用 Adobe Acrobat 6.0 版，即可存取 PDF 文件。PDF 文件已結構化，應該可以使用大部份的螢幕閱讀器來讀取。

企業搜尋的詞彙名詞解釋

此名詞解釋定義用於企業搜尋介面及文件的詞彙。

存取控制清單 (access control list)

一種清單，用來識別哪些使用者可以存取相關物件，並指定使用者對該物件的存取權。

管理角色 (administrative role)

決定使用者在企業搜尋管理主控台中可使用功能的使用者分類。此角色也決定使用者可以管理哪些集合。

分析引擎 (analysis engine)

請參閱文字分析引擎。

分析結果 (analysis results)

資訊由註解程式所產生。和您要搜尋之資訊對應的分析結果，會寫入稱為「共用分析結構」的資料結構中。

註解 (annotation)

關於文字跨距的相關資訊。例如，註解可指出代表公司名稱的文字跨距。在 UIMA 中，註解是一種特殊的特性結構類型。

註解程式 (annotator)

執行特定語言分析作業，並產生和記錄註解的軟體元件。註解程式是分析引擎中的分析邏輯元件。

Boolean 搜尋 (boolean search)

使用運算子如 AND、NOT 及 OR 結合一或多個搜尋字詞的搜尋。

Boost 類別 (boost class)

可以影響文件在搜尋結果中相關性排序的一種規格。

Boost 字 (boost word)

可以影響文件在搜尋結果中相關性排序的單字。在查詢處理程序期間，含有 Boost 字的文件重要性會提高或降低，視該單字預先定義的評分而定。

種類 (category)

具相同內容的文件群組。

種類樹狀結構 (category tree)

顯示在企業搜尋管理主控台的種類階層結構。

認證 (certificate)

一種數位文件，可以將公開金鑰連結至憑證擁有人的身份，藉以鑑別憑證擁有人。憑證是由憑證管理中心發出。

憑證管理中心 (certificate authority)

發出憑證及鑑別電子交易中相關實體 (個人或組織) 的一種組織。憑證管理中心可保證交換資訊的雙方身份沒有問題。

字元正常化 (character normalization)

將字元的變體形式 (如大寫字體及區別發音符號) 簡化成一般形式的程序。

附著語素 (clitic)

語法上的功能完全不同，但在語音上連至另一個字的一種單字。附著語素的寫法可以和連結的單字相連接，或完全分開。附著語素的常見範例包括英文縮語的最後一部份 (*wouldn't* 或 *you're*)。

集合 (collection)

一組資料來源及搜索、剖析、建立索引及搜尋這些資料來源的選項。

共用分析結構 (common analysis structure)

一種結構，用來儲存文字分析引擎分析的文件。資訊會以註解及其他特性結構的形式，以共用分析結構儲存。

共用通訊層 (Common Communication Layer, CCL)

聯合 WebSphere Information Integrator OmniFind Edition 中各種元件 (控制器、剖析器、搜索器、索引器) 的通訊架構。

概念萃取 (concept extraction)

定義文字文件的重要字彙項目 (例如人名、地點及產品)，以及產生這些項目的清單的搜尋函數。另請參閱主題萃取。

搜索範圍 (crawl space)

與搜索器為建立索引從擷取項目讀取的指定型樣 (例如資料庫名稱、檔案系統路徑、網域名稱、IP 位址及 URL) 相符的一組來源。

搜索器 (crawler)

從資料來源擷取文件，以及收集可用於建立搜尋索引的資訊的軟體程式。

認證 (credential)

鑑別時獲得的詳細資訊，說明使用者、所有群組關聯及其他機密保護相關的身份屬性。認證可用來執行許多服務，如授權、審核及委任。

資料來源 (data source)

可擷取文件的任何資料儲存庫，例如 Web、關聯式及非關聯式資料庫、以及內容管理系統。

資料來源類型 (data source type)

依據存取資料通訊協定的資料來源分組。

移出佇列 (dequeue)

從佇列中移除項目。

區別發音符號 (diacritics)

加入字母以變更單字發音或區別類似單字的一種標記，如重音符號或德文曲音。

探索器 (discoverer)

搜索器的一項功能，決定哪些資料來源可讓搜索器擷取資訊。

識別名稱 (distinguished name)

唯一識別目錄項目的名稱。識別名稱是由逗點區隔的屬性:值組組成。另外，一組名稱值組 (如 CN=人名，及 C=國家或地區) 可唯一識別數位憑證的實體。

文件物件模型 (Document Object Model)

以物件樹方式檢視結構化文件 (如 XML 檔) 的系統，可以利用程式化的方式來存取及更新文件。

Domino® Document Manager 檔案櫃 (Domino Document Manager cabinet)

用來組織文件的 Domino Document Manager 資料庫。檔案櫃保留 Domino 資料庫。

Domino Document Manager 檔案庫 (Domino Document Manager library)

Domino Document Manager 資料庫，它是 Domino Document Manager 的進入點。

Domino Internet Inter-ORB Protocol (DIOP)

在伺服器上執行的一種伺服器作業，可以搭配 Domino Object Request Broker 一起使用，以便在使用 Notes Java 類別所建立的 Java Applet 與 Domino 伺服器之間通訊。瀏覽器使用者及 Domino 伺服器均使用 DIOP 來通訊及交換物件資料。

動態排序 (dynamic ranking)

查詢中詞彙的排序類型，這些詞彙是用來分析要搜尋的相關文件，以決定結果的排序。另請參閱文字計分。反義詞為靜態排序。

動態彙總 (dynamic summarization)

搜尋字詞強調顯示，且搜尋結果包含最能代表使用者搜尋文件概念的詞組的彙總類型。反義詞為靜態彙總。

加入佇列 (enqueue)

在佇列中加入項目。

企業搜尋管理員 (enterprise search administrator)

讓使用者可以管理整個企業搜尋系統的管理員。

跳出字元 (escape character)

為之後的一或多個字元抑制或選取特殊意義的字元。

外部資料來源 (external data source)

用於聯合，且未由 WebSphere Information Integrator OmniFind Edition 搜索、剖析或檢索的資料來源。外部資料來源的搜尋委託給那些資料來源的查詢應用程式程式設計介面。

特性路徑 (feature path)

用來存取 UIMA 特性結構中特性值的路徑。

特性結構 (feature structure)

代表文字分析結果的基礎資料結構。特性結構是屬性值結構。每一個特性結構都屬於某一類型，且每一類型都已指定一組有效的特性或屬性，非常類似 Java 類別。

聯合搜尋 (federated search)

可以在多個搜尋服務上執行搜尋，並傳回搜尋結果合併清單的一種搜尋功能。

聯合功能 (federation)

合併命名系統的程序，使聚集系統可以處理跨命名系統的複合式名稱。

欄位 (field)

記錄中最小的可識別部份。

限定欄位搜尋 (fielded search)

限於某個特定欄位的查詢。

任意文字搜尋 (free text search)

搜尋時，搜尋字組以沒有格式的文字表示。

全文檢索 (full text index)

參照資料項目以方便快捷尋找內含查詢字詞的文件的資料結構。

模糊搜尋 (fuzzy search)

傳回與搜尋字詞拼法相似的字詞的搜尋。

混合式搜尋 (hybrid search)

結合 Boolean 搜尋及任意文字搜尋。

身分識別管理 (identity management)

在機密保護儲存庫中加密使用者認證的能力。

索引 (index)

請參閱全文檢索。

索引佇列 (index queue)

要處理的索引重組要求或索引重新整理要求清單。

索引重新整理 (index refresh)

將新資訊新增至企業搜尋系統現存索引的程序。反義詞為索引重組。

索引重組 (index reorganization)

建置企業搜尋系統索引的程序。反義詞為索引重新整理。

資訊萃取 (information extraction)

一種概念萃取的類型，可自動辨識文字文件中重要的詞彙項目，如名稱、術語及表示式。

IP 位址 (IP address)

唯一的 32 位元位址，用來識別網路上的主機。

Java 資料庫連線功能 (Java Database Connectivity, JDBC)

Java 平台及大量資料庫之間資料庫獨立連線功能的業界標準。JDBC 介面提供 SQL 型資料庫存取的呼叫層次。

JavaScript™

在瀏覽器及 Web 伺服器中使用的一種 Web Scripting 語言。

JavaServer Pages (JSP)

一種伺服器 Scripting 技術，可以讓 Java 程式碼動態內嵌於網頁 (HTML 檔) 並在使用該網頁時執行，以便將動態內容傳回用戶端。

Java 虛擬機器 (Java virtual machine, JVM)

處理器的一種軟體實作，用來執行已編譯的 Java 程式碼 (Applet 及應用程式)。

Katakana

由兩種通用的日文語音字母之一所使用的符號組成的一種字集，主要用來依語音撰寫外文。

金鑰儲存庫檔案 (keystore file)

一種金鑰資料庫檔，其中含有儲存為簽章者憑證的公開金鑰，以及儲存在個人憑證中的私密金鑰。

語言識別 (language identification)

決定文件語言的一種企業搜尋功能。

詞形 (lemma)

字的標準格式。在高度變音的語言中，詞形是很重要的，如捷克。

詞形還原 (lemmatization)

查閱定義檔中指定單字詞形的程序。詞形還原不同於詞根檢索 (詞根檢索是一種演算法)，通常不使用列出語言單字的定義檔。

語彙聯繫關係 (lexical affinity)

搜尋單字在文件中彼此貼近的關係。語彙聯繫關係用於計算結果的關聯性。

檔案庫 (library)

一種系統物件，當成其他物件的目錄使用。另請參閱 Domino Document Manager 檔案庫。

連字 (ligature)

兩個以上連接的字元，使它們看起來像是一個字元，如結合 f 和 i 形成連字 fi。

輕量型目錄存取通訊協定 (Lightweight Directory Access Protocol, LDAP)

一種開放式通訊協定，使用 TCP/IP 存取支援 X.500 模型的目錄，且不需要更複雜的 X.500 Directory Access Protocol 的資源需求。

語言學搜尋 (linguistic search)

一種搜尋類型，利用還原成基礎詞形 (例如，*mice* 以 *mouse* 檢索) 或從基礎詞形擴充 (如複合字) 的詞彙來瀏覽、擷取及檢索文件。

鏈結分析 (link analysis)

以文件間超鏈結分析為基礎的一種方法，用來決定集合中的哪些頁面對使用者是重要的。

本端聯合器 (local federator)

聯合一組可搜尋物件的用戶端聯合器。

Lotus® QuickPlace® 工作區 (Lotus QuickPlace place)

由 Lotus QuickPlace 提供的一種 Web 集合點，可以讓散佈在不同地理環境的參與者在專案中分工合作，並在結構化的安全工作區中連線通訊。

Lotus QuickPlace 檔案室 (Lotus QuickPlace room)

Lotus QuickPlace 工作區的分割區，限於需要集體分擔工作的授權成員使用。

遮罩字元 (masking character)

用於代表搜尋字詞前面、中間及尾端選用字元的字元。遮罩字元通常用於尋找索引中詞彙的變體。另請參閱萬用字元。

MIME 類型 (MIME type)

一種網際網路標準，用來識別要在網際網路上傳送的物件類型。

模型種類 (model-based category)

用於決定文件主旨的預先定義詞彙分類法，使用具相似內容的文件，建立索引並搜尋文件。

監督 (monitor)

一種企業搜尋使用者，此使用者具有觀察集合層次程序的權限。

自然語言查詢 (natural language query)

分析書面詞句 (例如 "Who runs the finance department?") 而非簡單關鍵字集合的搜尋類型。

換行字元 (newline character)

造成列印或顯示位置向下移一行的一種控制字元。部份系統需要多個字元。

n-gram 斷詞法 (n-gram segmentation)

一種分析方法，會將指定字元數的重疊順序視為一個單字，而不是像 Unicode 型空格斷詞法使用空格來區隔單字。

不遵循指引 (no-follow directive)

網頁中的一種指引，指示網路蜘蛛 (如 Web 搜索器) 不遵循在那些網頁中找到的鏈結。

不檢索引 (no-index directive)

網頁中的一種指引，指示網路蜘蛛 (如 Web 搜索器) 不在索引中併入那些網頁的內容。

Notes 遠端程序呼叫 (Notes remote procedure call, NRPC)

Lotus Notes® 的架構層，用於所有 Notes 對 Notes 通訊。

操作員 (operator)

一種企業搜尋使用者，具有觀察、啟動及停止集合層次程序的權限。

參數搜尋 (parametric search)

一種搜尋類型，用來在指定範圍內尋找含有數值或屬性 (如日期、整數或其他數值資料類型) 的物件。

剖析器 (parser)

解譯新增至企業搜尋資料儲存區的文件的程式。剖析器會從文件中取出資訊，並準備好以供建立索引、搜尋及擷取之用。

工作區 (place)

可以在入口網站中看見的虛擬位置，個人和群組可以在該處協商分工合作。在入口網站中，每一個使用者都有私人工作用的個人工作區，而群組可以存取各種共用工作區 (可以是公用位置或限制工作區)。另請參閱 Lotus QuickPlace 工作區。

熱門等級 (popular ranking)

依據文件的熱門程度新增至文件現有等級的排序類型。

處理程序引擎保存檔 (processing engine archive)

一種 .pear zip 保存檔，其中含有 UIMA 分析引擎及所有在企業搜尋中用來它進行自訂分析所需的資源。

鄰近搜尋 (proximity search)

在相同句子、段落或文件中尋找特定單字的搜尋類型。

Proxy 伺服器 (proxy server)

一種伺服器，當成應用程式或 Web 伺服器所管理的 HTTP Web 要求媒介使用。Proxy 伺服器可以當成企業中內容伺服器的代用品。

快速鏈結 (quick link)

URI 及關鍵字與詞組之間的關聯。

排序 (ranking)

為查詢的搜尋結果中每一份文件指定一個整數值的程序。搜尋結果中的文件次序是依據查詢的相關性。較高的等級表示較為相符。另請參閱動態排序及靜態排序。

遠端聯合器 (remote federator)

聯合一組可搜尋物件的伺服器聯合器。

Robots Exclusion Protocol

一種通訊協定，可以讓網站管理者指示站台的哪些部份不得讓網路蜘蛛造訪。

檔案室 (room)

一種程式，可讓使用者建立文件以供他人讀取、回應他人意見，以及檢視專案狀態和截止時間。使用者也可以和其他在同一檔案室中的使用者聊天。另請參閱 Lotus QuickPlace 檔案室。

規則種類 (rule-based category)

依照規則建立的種類，指定哪些文件與哪些種類相關聯。例如，您可以定義規則，設定內含或不含特定單字或符合 URI 型樣的文件，與特定的種類相關。

範圍 (scope)

相關 URI 的群組，用來定義搜尋要求的範圍。

搜尋應用程式 (search application)

在企業搜尋系統中，為集合處理查詢、搜尋索引、傳回搜尋結果、以及擷取來源文件的程式。

搜尋快取 (search cache)

保留資料及先前搜尋要求的結果的緩衝區。

搜尋引擎 (search engine)

接受搜尋要求並將文件清單傳回給使用者的程式。

搜尋索引檔 (search index files)

一組檔案，其中的索引儲存於搜尋引擎。

搜尋結果 (search results)

符合搜尋要求的文件清單。

Secure Sockets Layer (SSL)

提供通訊私密性的一種安全通訊協定。

安全記號 (security token)

關於身分及機密保護的資訊，用於對集中文件的存取授權。不同的資料來源類型支援不同的安全記號類型。範例包括使用者角色、使用者 ID、群組 ID、以及可用於控制內容存取的其它資訊。

起點 URL (seed URL)

搜索的開始點。

斷詞法 (segmentation)

路徑控制將基本資訊單元分成較小單元 (稱為 BIU 區段) 的一種程序，以符合相鄰伺服器中較小的緩衝區大小。

servlet

在 Web 伺服器上執行的一種 Java 程式，可以產生回應 Web 用戶端要求的動態內容以延伸伺服器功能。Servlet 通常是用來連接資料庫與 Web。

軟錯誤頁面 (soft error page)

一種特殊頁面，可在 HTTP 伺服器無法傳回用戶端所要求的頁面時詳細說明問題，並配置 HTTP 伺服器傳回這些頁面，而不是只在回應中顯示標頭及指出發生什麼問題的回覆碼。

靜態排序 (static ranking)

一種排序類型，以有關要排序文件的因數 (如日期、指向文件的鏈結數等等) 來提高排序。反義詞為動態排序。

靜態彙總 (static summarization)

搜尋結果包含指定的已儲存文件摘要的彙總類型。反義詞為動態彙總。

詞幹 (stemming)

請參閱字幹。

停止單字 (stop word)

搜尋應用程式不處理的一種常用單字，如 *the*、*an* 或 *and*。

停止單字移除 (stop word removal)

從查詢中移除停用字的程序，以忽略常用字並傳回更多相關結果。

彙總 (summarization)

在搜尋結果中包括句子的程序，以簡短地描述文件的內容。請參閱動態彙總及靜態彙總。

同義字定義檔 (synonym dictionary)

一種定義檔，可以讓使用者在搜尋集合時搜尋查詢字詞的同義字。

分類法 (taxonomy)

依據相似性分出物件群組的分類。在企業搜尋中，分類法可組織資料入種類及子種類。另請參閱種類樹狀結構。

文字分析 (text analysis)

從文字中擷取語意及其他資訊以加強集中資料可擷取性的程序。

文字分析引擎 (text analysis engine)

一種軟體元件，負責尋找及表示文字中的環境定義及語意內容。

文字計分 (text-based scoring)

為文件指定整數值的程序，以表示文件與查詢字詞的相關性。較高的整數值表示與查詢較相符。另請參閱動態排序。

主題萃取 (theme extraction)

概念萃取的一種類型，自動辨識文字文件中重要的字彙項目，以取出文件的主題。另請參閱概念萃取。

記號 (token)

企業搜尋檢索的基本文字單元。記號可以是語言的單字或其他適用於檢索的文字單元。

記號器 (tokenizer)

一種文字斷詞法程式，可掃描文字並判斷文字系列是否可以識別為記號及其時機。

尾端字元 (trailing character)

A character that holds the last position in a word.

Unicode 型空格斷詞法 (Unicode-based white space segmentation)

分段的方法，使用 Unicode 字元內容來區分記號及分隔字元之間。

統一資源識別碼 (Uniform Resource Identifier, URI)

識別抽象或實體資源的精簡字串。

統一資源定位器 (Uniform Resource Locator, URL)

在電腦或網路如網際網路上，代表資訊資源的一連串字元。此一連串字元包括用於存取資訊資源的通訊協定縮寫名稱，此由通訊協定使用以尋找資訊資源。

通用資源名稱 (Universal Resource Name, URN)

一種網際網路通訊協定元素，由符合特定語法的短字串組成。字串包含可用來參照資源的名稱或位址。

非結構化資訊管理架構 (Unstructured Information Management Architecture, UIMA)

一種 IBM 架構，定義實作系統以進行非結構化資料分析的架構。

使用者代理程式 (user agent)

瀏覽 Web 並在其所造訪網站中留下自己的資訊的應用程式。在企業搜尋中，Web 搜索器是使用者代理程式。

Web 搜索器 (Web crawler)

網路蜘蛛軟體類別，可以擷取 Web 文件並遵循該文件中的鏈結以探索 Web。

加權字詞搜尋 (weighted term search)

某些字詞重要性較高的查詢。

萬用字元 (wildcard character)

用於代表搜尋字詞前面、中間或尾端選用字元的字元。

字幹 (word stemming)

語言正常化的程序，將單字的變化形簡化成常見形。例如，*connections*、*connective* 及 *connected* 之類的單字會還原成 *connect*。

XML 路徑語言 (XML Path Language, XPath)

唯一識別或定址來源 XML 文件組件的語言。XPath 也提供操作字串、數字和布林運算子的基本機能。

有關 WebSphere Information Integration 的存取資訊

您可以透過電話或從 Web 上取得 WebSphere Information Integration 產品的相關資訊。

這裡所提供的電話號碼只適用於美國：

- 若要訂購產品或取得一般資訊：1-800-IBM-CALL (1-800-426-2255)
- 若要訂購書籍：1-800-879-2755

如需 WebSphere Information Integration 的相關資訊，另請造訪 Web 的 www.ibm.com/software/data/integration/db2ii/。此站台含有下列項目的最新相關資訊：

- 產品文件
- 產品下載
- 修正套件
- 版本注意事項及其他支援文件
- WebSphere Information Integration 相關新聞
- Web 資源鏈結，如白皮書及 IBM Redbooks™
- 新聞群組及使用者群組的鏈結
- WebSphere Information Integration 產品線上資訊中心的鏈結
- 訂購書籍

若要存取產品文件：

1. 請造訪 Web 的 www.ibm.com/software/data/integration/db2ii/。
2. 從下拉清單選取產品，例如 WebSphere Information Integrator OmniFind Edition。
3. 按一下頁面左邊的 Support 鏈結。
4. 在 Learn 區段，選取您需要的鏈結。如果有適用於所選產品的資訊中心，您可以選取該資訊中心的鏈結。如需範例，請參閱第 80 頁的圖 1。

Learn

- **Product documentation and manuals** (2 items)
- **Redbooks** (1 item)
- **V8.2 Documentation and release notes**

Information Center

Provides fast, online centralized access to product information.

- [1.0](#)

圖 1. WebSphere Information Integration 支援網站上產品文件的鏈結範例

提供文件的相關意見

如果您對本資訊或其他 IBM WebSphere Information Integration 文件有任何意見，請將您的意見傳送給我們。

您的意見有助於 IBM 提供高品質的資訊。如果您對本資訊或其他 WebSphere Information Integration 文件有任何意見，請將您的意見傳送給我們。您可以使用下列任一方法來提供意見：

1. 您可以利用 www.ibm.com/software/awdtools/rcf/ 的線上讀者意見表，來傳送您的意見。
2. 請將您的意見以電子郵件寄至 comments@us.ibm.com。請併入產品名稱、產品的版本號碼，以及書籍資訊的名稱與產品編號 (如果有的話)。如果您對特定文字有意見，請說明該文字的位置 (例如，標題、表格號碼或頁碼)。

聯絡 IBM

若要聯絡美國或加拿大的 IBM 客戶服務中心，請撥打 1-800-IBM-SERV (1-800-426-7378)。

若要瞭解適用的服務選項，請撥打下列其中一個號碼：

- 美國：1-888-426-4343
- 加拿大：1-800-465-9600

若要尋找您所在國家或地區的 IBM 辦事處，請參閱 IBM Directory of Worldwide Contacts 網站，網址為 www.ibm.com/planetwide。

商標

本主題列出 IBM 商標及某些非 IBM 商標。

如需 IBM 商標的相關資訊，請參閱 <http://www.ibm.com/legal/copytrade.shtml>。

下列術語是其他公司的商標或註冊商標：

Java 及所有以 Java 為基礎的商標和標誌是 Sun Microsystems, Inc. 在美國及 (或) 其他國家的商標或註冊商標。

Microsoft、Windows、Windows NT 以及 Windows 標誌是 Microsoft Corporation 在美國及 (或) 其他國家的商標。

Intel、Intel Inside (標誌)、MMX 及 Pentium 是 Intel Corporation 在美國及 (或) 其他國家的商標。

UNIX 是 The Open Group 在美國及其他國家的註冊商標。

Linux 是 Linus Torvalds 在美國及 (或) 其他國家的商標。

其他公司、產品或服務名稱，可能是其他公司的商標或服務標誌。

注意事項

本資訊是針對 IBM 在美國所提供之產品與服務開發出來的。IBM 不見得會對所有國家或地區都提供本文件所提的各項產品、服務或功能。要知道在您所在地區是否可得到這些產品及服務時，請向當地的 IBM 服務代表查詢。而此處任何對於 IBM 產品、程式或服務的參考之處，並不表示或暗示只可以使用 IBM 的產品、程式或服務。任何未侵犯 IBM 的智慧財產權，任何功能相當的產品、程式或服務都可以取代 IBM 的產品、程式或服務。不過，使用者必須自行負責評估和驗證任何非 IBM 產品、程式或服務的作業。

在本文件中可能包含著 IBM 所擁有之專利或擱置專利申請的內容。本文件使用者並不享有前述專利之任何授權。您可以用書面方式來查詢授權，來函請寄到：IBM Director of Licensing IBM Corporation North Castle Drive Armonk, NY 10504-1785 U.S.A.

若要查詢二位元組 (DBCS) 資訊的授權事宜，請連絡您國家或地區的 IBM 智慧財產部門，或者用書面方式寄到：IBM World Trade Asia Corporation Licensing 2-31 Roppongi 3-chome, Minato-ku Tokyo 106-0032, Japan

下列段落不適用於英國或任何其他與當地法律相抵觸的國家或地區：IBM 公司係以『現狀』提供本出版品，且不作任何明示或默示的保證，包括但不僅限於非侵害、可售性或符合特定用途之暗示保證。有些地區不允許放棄在特定交易中的明示或默示保證，因此，這項聲明對您可能不適用。

本書中可能會有技術上的錯誤或排版印刷上的訛誤。因此，IBM 會定期修訂；並將修訂後的內容納入新版中。IBM 得隨時修改及/或變更本書中所說明的產品及/或程式，恕不另行通知。

本資訊中任何對非 IBM 網站的敘述僅供參考，為便利貴客戶之使用，而非為該網站背書。這些網站中的資料，並不包含在 IBM 產品的資料中，使用網站中的資料，須自行負擔風險。

在不造成您困擾或損及您個人權益的前提下，IBM 得以適切使用或散佈您以各種型式所提供的相關資訊。

本程式之獲授權者若希望取得本程式之相關資訊，以便達到下列目的：(i) 在獨立建立的程式與其他程式 (包括本程式) 之間交換資訊；以及 (ii) 相互使用已交換的資訊。則請與位於下列地址之人員連絡：

IBM Corporation J46A/G4
555 Bailey Avenue
San Jose, CA 95141-1003 U.S.A.

上述資料之取得有其條件，在某些情況下必須付費方得使用。

IBM 基於「IBM 客戶合約」、「IBM 國際程式授權合約」或雙方之間任何同等的合約等條款，提供本文件中所說的授權程式與其所有適用的授權資料。

任何此處涵蓋的執行效能資料都是在一個受控制的環境下決定出來的。因此，若在其他作業環境下，所得的結果可能會大大不同。有些測定已在開發階段系統上做過，不

過這並不保證在一般系統上會出現相同結果。再者，有些測定可能已透過推測方式評估過。但實際結果可能並非如此。本文件的使用者應依自己的特定環境，查證適用的資料。

非 IBM 產品的相關資訊，取自該產品供應商、發佈的聲明或其他公共來源。IBM 未測試這些產品，因此無法確認非 IBM 產品的效能、相容性或其他聲明。有關非 IBM 產品的功能問題，請洽該產品供應商。

有關 IBM 未來動向的任何陳述，僅代表 IBM 的目標而已，並可能於未事先聲明的情況下有所變動或撤回。

這個資訊中包含每日業務使用的報告和資料範例。為使說明盡可能完備，範例中包含個人、公司、品牌及產品的名稱。此等名稱皆屬虛構，凡有類似實際個人或企業所用之名稱及地址者，皆屬巧合。

著作權授權：

本資訊可包含原始語言的範例應用程式，用以說明各種作業平台上的程式設計技術。貴客戶得為開發、使用、行銷或散佈運用樣本程式之作業平台的應用程式程式介面所撰寫的應用程式之目的，免費複製、修改並散佈這些樣本程式。這些範例並未在所有情況下完整測試。故 IBM 不保證或默示保證這些樣本程式之可靠性、服務性或功能。貴客戶得為開發、使用、行銷或散佈符合 IBM 應用程式設計介面的應用程式之目的，免費複製、修改並散佈這些樣本程式。

這些範例程式的每個複本或任何部分，或任何衍生作品都必須包括以下版權聲明：

Outside In (®) Viewer Technology, ©1992-2005 Stellent, Chicago, IL., Inc. All Rights Reserved.

IBM XSLT Processor Licensed Materials - Property of IBM ©Copyright IBM Corp., 1999-2005. All Rights Reserved.

索引

索引順序以中文字，英文字，及特殊符號之次序排列。

〔四劃〕

- 支援語言
 - 定義檔型語言處理 60
 - 語言偵測 59
- 文件 65
- 日文的正體字變體 62

〔六劃〕

- 同義字定義檔
 - 建立 DIC 檔案 48
 - 建立 XML 檔案 47
 - 搜尋應用程式支援 47
- 字元正常化 63
- 存取文字分析結果
 - CAS 消費者的定義 17
- 存取自訂分析結果
 - 內建特性 18
 - 特性路徑的定義 17
 - 過濾器 20
- 自訂分析
 - 工作流程 4
 - 文字分析演算法 7
 - 在分析及搜尋中使用 XML 標記的方法 11
 - 具有 JDBC 能力之資料庫中的對映分析結果 27, 28, 32
 - 檢索自訂分析結果的方法 21
 - 類型系統說明範例 8

〔八劃〕

- 具有 JDBC 能力之資料庫中的對映自訂分析結果
 - 步驟 28
 - 說明 27
 - 儲存區類型 32
 - 儲存區類型對映 32
 - XML 對映配置檔 28
- 協助工具 67
- 定義檔型分析 60
- 定義檔型斷詞法 60
- 附著語素 60
- 非定義檔型分析 60
- 非定義檔型斷詞法 60

〔十一劃〕

- 停止單字移除 (stop word removal) 63
- 停用字 63
- 停用字定義檔
 - 建立 DIC 檔案 52
 - 建立 XML 檔案 51
 - 搜尋應用程式支援 51

〔十二劃〕

- 尋找企業搜尋文件 65
- 詞形 60
- 詞形還原 60

〔十三劃〕

- 搜尋伺服器
 - 同義字 XML 檔案 47
 - 建立 Boost 字定義檔 57
 - 建立同義字定義檔 48
 - 建立停用字定義檔 52
 - 停用字 XML 檔案 51
 - Boost 字 XML 檔 55
- 搜尋應用程式
 - 同義字支援 47
 - 停用字支援 51
 - Boost 字支援 55

〔十四劃〕

- 對映 XML 文件結構至 UIMA 類型
 - 建立配置檔 12
 - 說明 11
- 語言支援
 - 支援語言 60
 - 日文的正體字變體 62
 - 日文的斷詞 62
 - 字元正常化 63
 - 系統定義的類型及特性 38
 - 系統隨附的支援 59
 - 定義檔型斷詞法 60
 - 附著語素 60
 - 非定義檔型斷詞法 60
 - 停止單字移除 (stop word removal) 63
 - 詞形 60
 - 詞形還原 60
 - 語言偵測 59
 - 語意搜尋 43
 - 說明 1

- 語言支援 (繼續)
 - n-gram 斷詞法 60
 - Okurigana 變體 62
 - Unicode 正常化 63
 - Unicode 型空格斷詞法 60
- 語言偵測 59
- 語意搜尋
 - 語意搜尋查詢 44
 - 說明 43
 - 擷取文件中符合查詢的部份 36

〔十五劃〕

- 編製索引自訂分析結果
 - 建立配置檔 22
 - 說明 21

〔十八劃〕

- 斷詞法
 - 定義檔型 60
 - 非定義檔型 60
 - Unicode 型空格 60
- 斷詞，日文 62

B

- Boost 字定義檔
 - 建立 DIC 檔案 57
 - 建立 XML 檔案 55
 - 搜尋應用程式支援 55

D

- DIC 檔案
 - 同義字 48
 - 使用者定義的停用字 52
 - Boost 字 57

E

- esboostworddictbuilder.bat script 57
- esboostworddictbuilder.sh script 57
- esstopworddictbuilder.bat Script 52
- esstopworddictbuilder.sh Script 52
- essyndictbuilder.bat script 48
- essyndictbuilder.sh script 48

N

n-gram 斷詞法 60

O

Okurigana 變體 62

P

PDF 文件 65

S

Script

esboostworddictbuilder 57

esstopworddictbuilder 52

essyndictbuilder 48

U

UIMA

安裝基本企業搜尋註解程式 6

自訂文字分析支援 3

定義的類型及特性 41

基本概念 4

執行基本企業搜尋註解程式 6

說明 3

Unicode 正常化 63

Unicode 型空格斷詞法 60

W

WebSphere II OmniFind Edition 67

協助工具 67

讀者意見表

為使本書盡善盡美，本公司極需您寶貴的意見；懇請您閱讀後，撥冗填寫下表，惠予指教。

請於下表適當空格內，填入記號(√)；我們會在下一版中，作適當修訂，謝謝您的合作!

評估項目	評估意見	備註
正確性	內容說明與實際程序是否符合	<input type="checkbox"/> 是 <input type="checkbox"/> 否
	參考書目是否正確	<input type="checkbox"/> 是 <input type="checkbox"/> 否
一致性	文句用語及風格，前後是否一致	<input type="checkbox"/> 是 <input type="checkbox"/> 否
	實際產品介面訊息與本書中所提是否一致	<input type="checkbox"/> 是 <input type="checkbox"/> 否
完整性	是否遺漏您想知道的項目	<input type="checkbox"/> 是 <input type="checkbox"/> 否
	字句、章節是否有遺漏	<input type="checkbox"/> 是 <input type="checkbox"/> 否
術語使用	術語之使用是否恰當	<input type="checkbox"/> 是 <input type="checkbox"/> 否
	術語之使用，前後是否一致	<input type="checkbox"/> 是 <input type="checkbox"/> 否
可讀性	文句用語是否通順	<input type="checkbox"/> 是 <input type="checkbox"/> 否
	有否不知所云之處	<input type="checkbox"/> 是 <input type="checkbox"/> 否
內容說明	內容說明是否詳盡	<input type="checkbox"/> 是 <input type="checkbox"/> 否
	例題說明是否詳盡	<input type="checkbox"/> 是 <input type="checkbox"/> 否
排版方式	本書的形狀大小，版面安排是否方便閱讀	<input type="checkbox"/> 是 <input type="checkbox"/> 否
	字體大小，顏色編排，是否有助於閱讀	<input type="checkbox"/> 是 <input type="checkbox"/> 否
目錄索引	目錄內容之編排，是否便於查找	<input type="checkbox"/> 是 <input type="checkbox"/> 否
	索引語錄之排定，是否便於查找	<input type="checkbox"/> 是 <input type="checkbox"/> 否
※評估意見為"否"者，請於備註欄提供建議。		

其他：(篇幅不夠時，請另外附紙說明。)

上述改正意見，一經採用，本公司有合法之使用及發佈權利，特此聲明。
註：您也可將寶貴的意見以電子郵件寄至 NLSC01@tw.ibm.com，謝謝。

IBM WebSphere Information Integrator
OmniFind Edition
文字分析整合 8.3 版

SC40-2071-00

折疊線

110 台北市信義區松仁路 7 號 3 樓

臺灣國際商業機器股份有限公司
大中華研發中心 軟體國際部 啟



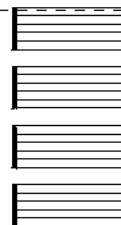
廣告回信
台灣北區郵政管理局 登記
北台字第 00176 號

(免貼郵票)

寄件人 姓名：
地址：

寄

折疊線



IBM



Java[™]
COMPATIBLE

SC40-2071-00



Spine information:



WebSphere II OmniFind
Edition

文字分析整合

8.3 版