



راهنمای استفاده از پیکره تک سندی پاسخ

تهیه شده توسط آزمایشگاه فناوری وب دانشگاه فردوسی مشهد

سفارش دهنده : سازمان فناوری اطلاعات و ارتباطات ایران

wtlab.um.ac.ir

ijaz.um.ac.ir

جمهوری اسلامی ایران

وزارت ارتباطات و فناوری اطلاعات

برای تولید پیکره تک سندی در ابتدا باید تعدادی سند از سایتهای خبرگزاری معروف انتخاب شود. در انتخاب خبرگزاری ها سعی شد تا آنهایی انتخاب شوند که بیشترین بازدید کننده را دارند و در ضمن دیدگاههای مختلف را پوشش دهند. بهمین منظور ۱۰۰ خبر با اندازه های مختلف از خبرگزاری های جدول ۱ دانلود و پالایش گردید:

جدول ۱- لیست خبرگزاری های انتخاب شده برای جمع آوری خبر

ردیف	نام خبرگزاری	آدرس خبرگزاری	شماره کد در نام گذاری
۱	خبرگزاری تابناک	http://www.tabnak.ir	TAB
۲	خبرگزاری پرس تی وی	http://www.presstv.ir	PTV
۳	خبرگزاری فارس	http://www.farsnews.com	FAR
۴	خبرگزاری ایرنا	http://irna.ir	IRN
۵	خبرگزاری همشهری	www.hamshahrionline.ir	HAM
۶	خبرگزاری الف	www.alef.ir	ALF
۷	خبرگزاری جام جم	www.jamejamonline.ir	JAM

این موضوعات در قالب دسته بندی کلی زیر جمع آوری شدند. در جمع آوری این پیکره سعی شده است تا همه انواع خبرها پوشش داده شود.

جدول ۲- دسته بندی موضوعات

ردیف	نام دسته	شماره کد در نام گذاری
۱	اقتصادی	EC
۲	فرهنگی	CU
۳	اجتماعی	SO
۴	سیاسی	PO
۵	ورزشی	SP
۶	علمی	SC

فرمت نام گذاری

جهت دسترسی سریع و مشخص به اسناد، فرمولی برای نام گذاری اسناد انتخاب شده است. اسناد در پیکره با فرمول زیر نام گذاری می شوند:

<شماره سند>.<تاریخ سند>.<کد دسته سند>.<کد خبرگزاری سند>

بنابراین سند با نام [PTV.PO.13901228.068.xml](#) یعنی سند با شماره کد ۰۶۸ که از نوع سیاسی بوده و از خبرگزاری پرس تی وی دانلود شده است و تاریخ آن هم مربوط به ۲۸ اسفند سال ۱۳۹۰ می باشد. این اسناد در فرمت xml و با تگ های زیر تولید شدند:

<?xml version="1.0" encoding="utf-8" ?>

<DOC>

<DOCNO>010</DOCNO>

<DATE-TIME>1391/02/02</DATE-TIME>

<CATEGORY>ورزشی</CATEGORY>

<TITLE>

در آکادمی ملی المپیک: دور جدید تمرینات تیم ملی والیبال بانوان از ۵ اردیبهشت می شود

</TITLE>

<SUMMARY>

خبرگزاری فارس: مرحله جدید اردوی تیم ملی والیبال بانوان بزرگسال برای حضور

در مسابقات هشت تیم برتر آسیا در آکادمی ملی المپیک برگزار می شود.

</SUMMARY>

<TEXT>

به گزارش خبرگزاری فارس، تیم ملی والیبال بانوان بزرگسال ایران که خود را برای حضور در مسابقات هشت تیم برتر آسیا آماده می کند به دلیل پر بودن سالن های آزادی در باشگاه پیکان تمرین می کند.

پس از بازگشایی هتل آکادمی قرار است ملی پوشان تمریناتشان را در این مکان پیگیری کنند.

تمرینات تیم ملی دختران نوجوان و جوان نیز بارها به دلیل نبود سالن و خوابگاه تعطیل شده است.

این اردو با حضور ۲۰ بازیکن از ۵ اردیبهشت آغاز می شود. پس از بازگشت تیم گیتی پسند از مسابقات جام باشگاه های آسیا تعدادی

فایل های xml به عنوان یک قالب استاندارد برای انتقال داده در محیط های مختلف نرم افزاری شناخته می شود. تقریباً در تمام زبان های برنامه سازی سطح بالا امکانات مناسبی برای کار با فایل های xml فراهم شده است. استفاده از ساختار سلسله مراتبی و بهره مندی از اصول ساخت یافتگی سبب شده است تا در دسره های احتمالی استفاده از پایگاه های داده در بسیاری از برنامه های کاربردی از بین برود. xml به عنوان یک پایه برای طراحی زبان های تحت وب نیز تلقی می شود. از این رو به اعتقاد بسیاری html در واقع گونه گسترش یافته xml در فضای وب محسوب می شود.

همانطور که در تصویر هم مشخص است مشخصات سند در تگهای xml قابل دسترسی است. مشخصات این تگ‌ها در جدول ۳ آمده است:

جدول ۳- لیست تگ های سندهای خبری	
توضیحات	تگ
حاوی شماره سند می باشد.	<DOCNO>
زمان انتشار خبر را مشخص می کند.	<DOC-TIME>
دسته خبر را مشخص می کند (سیاسی، اقتصادی، ورزشی و ...)	<CATEGORY>
عنوان خبر را مشخص می کند.	<TITLE>
خلاصه خبر که عموماً توسط خبرگزاریها تولید می شود در این تگ قرار می گیرد.	<SUMMARY>
متن خبر در این تگ قرار داده می شود.	<TEXT>

فرآیند تولید پیکره

برای تولید پیکره خلاصه انسانی، ۱۰ نفر از دانشجویان کارشناسی (زن - مرد) بکار گرفته شدند. با توجه به اینکه هر کدام از این افراد سلايق و دانش خاص خود را دارند، به منظور کاهش نظرات شخصی افراد، برای هر متن ۵ نفر خلاصه انسانی تولید کردند. بدین ترتیب در زمان ارزیابی، نتایج خلاصه های ماشینی با میانگین این ۵ خلاصه انسانی مقایسه می شود و تاثیر نظرات شخصی افراد خلاصه ساز کاهش می یابد.

هر کدام از این افراد، برای موضوعات تعیین شده، یک خلاصه چکیده ای و یک خلاصه ی استخراجی تولید کردند. در خلاصه سازی استخراجی، افراد شرکت کننده در فرآیند تولید پیکره، برای هر متن تک سندی، بین ۳ تا ۷ عدد از مهمترین جملات متن سند را انتخاب کردند. بنابراین مهمترین مساله در این مدل تعیین مهمترین جملات می باشد. در خلاصه سازی چکیده ای، افراد شرکت کننده در فرآیند تولید پیکره، برای هر متن تک سندی، حدود ۳ تا ۷ جمله که بیانگر محتوای اصلی متن می باشد، با قلم خود بازنویسی کردند.

بنابراین در مجموع، برای ۱۰۰ متن خبری، ۱۰۰۰ خلاصه تولید شد، که ۵۰۰ عدد چکیده ای و ۵۰۰ عدد/ستخر/اجی می باشد.

برای تولید پیکره توسط افراد، نرم افزاری به نام خلاصه یار تولید شد که از آن برای تولید خلاصه های مختلف می توان استفاده نمود. در بخش های بعدی به معرفی این سامانه پرداخته شده است.

لیست موضوعات تک سندی

همانطور که پیش تر اشاره شد، ۱۰۰ سند برای خلاصه سازی تک سندی جمع آوری شد که لیست عنوان موضوعات جمع آوری شده در جدول ذیل آورده شده است:

جدول ۴- لیست موضوعات اسناد برای تولید خلاصه های تک سندی

ردیف	موضوع	شماره سند
۱	جدایی نادر از... بالاخره خوب است یا نه	ALF.CU.13910117.019
۲	لزوم گسترش و بکارگیری تجارت الکترونیکی	ALF.EC.13901207.026
۳	قطر بازار جدید صادرات فرش ایرانی	ALF.EC.13910204.022
۴	کالاهایی که ۵ سال در گمرک خاک می خورد	ALF.EC.13910204.024
۵	گشایش بزرگترین نمایشگاه تجاری جهان در آلمان	ALF.EC.13910204.025
۶	تازه ترین دستاورد دانشمندان برای درمان کچلی	ALF.SC.13910131.018
۷	در گفتگو با یک جامعه شناس مطرح شد؛ کار برای ما ضد ارزش است	ALF.SO.13910126.023
۸	برخورد اماکن با خوابگاههای غیرمجاز	ALF.SO.13910130.017
۹	شمقدری، ایازی، اسلامی مهر، آقامحمدیان و...: معرفی اعضای شورای سیاستگذاری نخستین جشنواره فیلم ویدیویی تهران	FAR.CU.13910203.004
۱۰	با همکاری سینمای استرالیا: بازیگر انگلیسی فیلمساز شد	FAR.CU.13910203.005
۱۱	۸ تهدید برای بانک های ایرانی در سال ۹۱	FAR.EC.13910127.012
۱۲	نکته ای قابل توجه پس از نشست استانبول: نتایج زود هنگام مذاکرات	FAR.EC.13910129.014

	هسته ای	
FAR.EC.13910131.011	تأملی درباره گرانی های اخیر	۱۳
FAR.EC.13910203.008	نقش تبلیغات تجاری بر اشاعه‌ی مصرف‌زدگی	۱۴
FAR.SC.13901223.007	بدغذا خوردن افسردگی می‌آورد	۱۵
FAR.SC.13910122.006	۵ عقیده رایج اما اشتباه	۱۶
FAR.SP.13910202.010	در آکادمی ملی المپیک: دور جدید تمرینات تیم ملی والیبال بانوان از ۵ اردیبهشت می‌شود	۱۷
HAM.SC.13910204.038	هابل ۲۲ ساله شد؛ سومین تلسکوپ قدرتمند جهان از دید یک اخترشناس	۱۸
IRN.PO.13910202.036	وزیر خارجه انگلیس نیز بالاخره از خشونت های بحرین ابراز نگرانی کرد	۱۹
IRN.PO.13910203.028	انتخابات ریاست جمهوری فرانسه آغاز شد	۲۰
IRN.PO.13910203.037	نماینده دوما ی روسیه نتایج نشست استانبول را مثبت ارزیابی کرد	۲۱
IRN.PO.13910204.031	پایگاه خبری مالزیایی: دستیابی ایران به سیستم ساخت پهپاد آمریکایی یک موفقیت بزرگ محسوب می شود	۲۲
IRN.SO.13910204.029	فرمانده نیروی انتظامی: تلفات رانندگی با اجبار کاهش نمی یابد	۲۳
JAM.CU.13901101.003	فرهادی با چه کسانی رقابت کرد؟	۲۴
JAM.PO.13910203.034	کاهش مشارکت در انتخابات فرانسه	۲۵
PTV.PO.13910203.015	انتخابات ریاست جمهوری در فرانسه روز یک شنبه آغاز شد	۲۶
TAB.EC.13910125.030	حمایت از تولید ملی و افزایش صادرات	۲۷
TAB.EC.13910203.035	عباسعلی نورا : کاهش ارزش پول ملی عامل بازدارنده سرمایه‌گذاری	۲۸
TAB.PO.13910202.032	یک روز مانده به انتخابات فرانسه: «سارکوزی» در آستانه شکست!	۲۹
TAB.SP.13910203.033	مورینیو چگونه سد بارسا را شکست؟	۳۰
ALF.PO.13910203.021	"فرمول خون" اوضاع بحرین را نقش بر آب کرد	۳۱
ALF.PO.13910204.020	زنی که پایه های دیکتاتوری قذافی را لرزاند	۳۲
ALF.PO.13910204.027	هشدار ارتش کره شمالی به کره جنوبی	۳۳
FAR.SO.13910202.009	برای آموزش قوانین راهنمایی و رانندگی صورت گرفت: ابتکار معلم	۳۴

	سیرجانی در یادگیری خوب دانش‌آموزان	
FAR.SP.13910202.013	از نهم اردیبهشت در مجموعه ورزشی آزادی: بانوان دوچرخه سوار در اولین مرحله لیگ کورسی رقابت می‌کنند	۳۵
HAM.SC.13900713.047	برندگان نوبل فیزیک ۲۰۱۱ معرفی شدند	۳۶
HAM.SO.13900802.050	روایتی از شرکت رهبر انقلاب در سرشماری عمومی نفوس و مسکن	۳۷
HAM.SP.13910203.046	مرحله نیمه نهایی لیگ برتر بسکتبال از امروز کلید می‌خورد: آغاز بازی‌های دشوار برای فینالیست‌شدن	۳۸
IRN.CU.13910203.045	ارکستر آکادمیک تهران با عنوان خلیج فارس به صحنه می‌رود	۳۹
IRN.CU.13910204.055	مرکز دایره المعارف بزرگ اسلامی با ۴۰ اثر چاپ اولی در نمایشگاه کتاب حضور می‌یابد	۴۰
IRN.PO.13900801.051	رییس کمیته عالی مستقل انتخابات تونس: در برخی حوزه‌ها بیش از ۸۰ درصد واجدان شرایط رای دادند	۴۱
IRN.SC.13910131.052	گونه نادر سیاه‌گوش برای نخستین بار در استان مرکزی مشاهده شد	۴۲
IRN.SC.13910204.053	همزمان با روز زمین، لامپی با طول عمر بیش از ۲۰ سال عرضه می‌شود	۴۳
IRN.SC.13910204.054	محققان ایرانی نانوکامپوزیت زیستی برای پوشش دهی سبب تولید کردند	۴۴
IRN.SP.13910201.048	دوومیدانی قهرمانی کشور؛ رکورد پرش ارتفاع ایران شکسته شد	۴۵
IRN.SP.13910202.016	اسماعیل پور: رضا یزدانی می‌تواند در بازی‌های المپیک لندن جاودانه شود	۴۶
IRN.SP.13910204.049	تیم ملی شنای کشورمان با ۸ ورزشکار عازم مالزی می‌شود	۴۷
IRN.SP.13910204.057	میراسماعیلی: نتایج تیم ملی جودو در مسابقات آسیایی قابل پیش‌بینی نیست	۴۸
JAM.EC.13910203.056	هفته نزولی طلا	۴۹
JAM.SC.13910204.002	علم چیست؟	۵۰
JAM.SP.13910203.001	کیپسانگ فاتح ماراتن لندن شد	۵۱
PTV.PO.13910127.042	چند راکت منطقه دیپلماتیک پایتخت افغانستان را لرزاند	۵۲
PTV.PO.13910203.043	"مصر برای اسرائیل خطرناک‌تر از ایران است"	۵۳

PTV.PO.13910204.040	اتحادیه اروپا درباره تحریم سوریه به توافق رسیده است	۵۴
PTV.PO.13910204.044	ناظران اعزامی سازمان ملل به سوریه از شهرهای حمص، حما و رستن بازدید کردند	۵۵
PTV.PO.13910205.041	اتحادیه عرب خواهان راه حل سیاسی در سوریه	۵۶
PTV.PO.13910921.060	شکار هواپیمای امریکایی و نگرانی اسرائیل	۵۷
TAB.PO.13910202.039	بر پایه آمار انتخاباتی؛ نمایندگان که بیشترین درصد رأی را در انتخابات نهم به دست آوردند	۵۸
TAB.SO.13910120.059	قتل تازه در پرونده مردی با سابقه ۳ جنایت	۵۹
TAB.SO.13910202.058	مهمترین تفاوت‌های اینترنت ملی با اینترنت فعلی	۶۰
ALF.SO.13910203.075	تشدید طرح برخورد با موتورسواران متخلف	۶۱
HAM.EC.13900802.083	آیا مجلس مردم را قانع می‌کند؟ جزئیات گزارش اختلاس در کمیسیون اصل ۹۰	۶۲
HAM.SC.13901118.085	دانشمندان در جست‌وجوی داوطلب برای شبیه سازی سفر به مریخ	۶۳
HAM.SO.13900419.072	پژوهشگران در مجله لنست اعلام کردند: افزایش شمار خودکشی‌ها در اروپا به دنبال بحران اقتصادی ۲۰۰۸	۶۴
HAM.SO.13900524.087	شمار مظنونان پرونده اسیدپاشی مرگبار به ۵ نفر رسید	۶۵
HAM.SP.13900416.076	ایران قهرمان وزنه برداری جوانان جهان شد	۶۶
HAM.SP.13901125.070	انتظار نیم قرنی والیبال به پایان می‌رسد؟ آرزوی رسیدن به المپیک	۶۷
JAM.SC.13910204.084	چشم‌ها به کامپیوتر دروغ نمی‌گویند	۶۸
JAM.SO.13900307.066	روان‌شناسان معتقدند یکی از عوامل اضطراب دانش‌آموزان، فشار روانی والدین بر آنهاست: اضطراب امتحان	۶۹
JAM.SO.13901020.071	اغلب بازنشستگان زیر خط فقر هستند	۷۰
JAM.SO.13901203.081	طی دیمه سال جاری: تهران در صدر آمار طلاق است	۷۱
JAM.SO.13910204.062	بررسی ضرورت ارتباط والدین و مربیان مدرسه: مسیری هموار میان خانه اول و دوم	۷۲
JAM.SP.13900813.078	رییس کمیته فوتسال تشریح کرد: برنامه‌های تیم ملی فوتسال	۷۳
PTV.PO.13900507.097	تظاهرات ضد امریکایی و سعودی در یمن	۷۴

PTV.PO.13901228.068	ادامه سرکوب اعتراضات جنبش ۹۹ درصدی	۷۵
PTV.SO.13900809.096	افزایش تلفات زلزله ترکیه به ۶۰۰ نفر	۷۶
TAB.CU.13900209.074	آیه‌الله جوادی: تعیین‌کنندگان زمان ظهور دجالند	۷۷
TAB.CU.13900311.069	اثر جدید استاد فرشچیان رونمایی شد	۷۸
TAB.CU.13900424.073	اقامتگاه امام زمان (عج) کجاست؟	۷۹
TAB.CU.13900428.077	برگزاری هفته هنر یک استان کهن و تاریخی	۸۰
HAM.CU.13900115.094	نگرانی از خاموشی ارکستر سمفونیک تهران	۸۱
HAM.CU.13900218.093	نمایش مستند های جشنواره فیلم شهر در خانه هنرمندان	۸۲
HAM.CU.13901206.082	جذب ۵ هزار مربی قرآن در آموزش و پرورش در ۳ مرحله آزمون انجام شد	۸۳
HAM.CU.13910107.090	مهم‌ترین رویدادهای فرهنگی و ادبی مهر سال ۹۰	۸۴
HAM.CU.13910116.061	۹۰ درصد بودجه قرآنی سال گذشته معادل ۲۰۰ میلیارد تومان به حساب واریز شد	۸۵
HAM.EC.13900619.089	قیمت سکه در بورس، روی مرز ۷۰۰ هزار تومان	۸۶
JAM.CU.13900622.092	کتاب زندگی و آثار خوشنویسی استاد غلامحسین امیرخانی منتشر شد: یک عمر زندگی در خطوط	۸۷
JAM.CU.13900819.067	دري: كتاب «كلنل» دولت‌آبادی را مطالعه می‌کنم، اما ضرورتی برای انتشارش در داخل کشور وجود ندارد: شهرستان‌ها و روستاها؛ کانون توجه هفته کتاب	۸۸
JAM.CU.13901114.086	معاون صدا مطرح کرد: زبان سایبری بهترین روش برای معرفی فرهنگ عاشورا	۸۹
JAM.CU.13901214.063	جشنواره هنرهای تجسمی فجر امروز با معرفی ۱۳۰ برگزیده و اهدای ۹ نشان عالی به کار خود پایان می‌دهد: پایان طولانی‌ترین جشنواره فجر پس از ۳۴ روز	۹۰
JAM.SC.13900924.065	نگاهی به نقش روبات‌ها در پژوهش‌های دیرینه‌زیست‌شناسی: روبات‌های پرنده، کلید حل معمای پرواز می‌شوند	۹۱
JAM.SC.13901218.064	توهم و هذیان در روان‌های پریشان موج می‌زند: روان‌پریشی، فراگیرتر از دیابت	۹۲
PTV.EC.13900917.100	یک اقتصاددان یونانی: یورو فرو خواهد پاشید	۹۳

PTV.PO.13900120.099	عربستان و اردن در پس ناآرامی‌های سوریه	۹۴
PTV.PO.13900202.098	تجزیه لیبی در دستور کار ناتو	۹۵
PTV.PO.13900501.095	۹۱ کشته در بمب‌گذاری و تیراندازی نروژ	۹۶
TAB.CU.13900522.091	هما با اجازه مراجع صبحانه سرو می‌کند	۹۷
TAB.EC.13900331.079	روز شمار بحران ارزی خرداد ماه؛ پایان بحران ارزی، در نزدیکی مرز دلار ۱۱۰۰ تومانی	۹۸
TAB.EC.13900401.080	پنجاه درصد از جوانان یونانی بیکار هستند	۹۹
TAB.SO.13900412.088	صاحبخانه، مستاجر را مقابل سگ‌های گرسنه قرار داد	۱۰۰