

Benchmark

Tramite l'analisi dell'esistente possiamo notare che esistono già alcuni siti che svolgono la funzione di "aggregatori", fornendo in un solo spazio una raccolta dei maggior corpus di italiano, fornendone una breve descrizione e il link alle piattaforme originali.

The screenshot shows the homepage of the Accademia della Crusca. At the top, there's a logo and a search bar labeled "Cerca nel sito...". Below the header, there are several menu items: L'Accademia, Attività, Pubblicazioni, Eventi, Scaffali digitali, Lingua Italiana, Archivio, and Biblioteca. A banner below the menu says "Banche dati, corpora e archivi testuali". Under this, there are sections for "Italiano scritto e parlato" and "ANALIAS_5_MT (Annotazione e ANAlisi Sintattica MuTtling-u)". There's also a section for "API (Archivio del Parlato Italiano)" and "BADIP (BAncA dati dell'Italiano Parlato)". On the right side, there's a "Agenda eventi" calendar for May, with various events listed. Below the calendar, there's a "Avvisi" section with a link to "La Commedia in cinese: in rete il ritmo moderno del Fronte Risorgimento".

banche dati, corpora e archivi testuali. In questo senso non può' essere utile ai fini dell'utilizzo dei corpora in classe perché comporterebbe un passaggio ulteriore da parte dell'insegnante/studente che è quello di selezionare i corpora e distinguerli da archivi e banche dati.

The screenshot shows the Treccani website under the heading "Principali corpora dell'italiano". It discusses the work of Roberto Busa and the Lexico di frequenza della lingua italiana contemporanea (LIF). It mentions the CoLFIS corpus and the Corpus e lessico di frequenza dell'italiano scritto (CoLFIS). Below this, there's a section titled "Esempi di uso" which includes a table titled "Occorrenze per milione di parole del nome 'opportunisto' secondo le concordanze". The table shows data for various categories like politica, cultura, sport, etc. The table is as follows:

CATEGORIA	OCCHERENZE
politica	2,7
cultura	1,1
sport	3,7
scuola	3,0
politologo	1,1
politologo	0,6
spazio	3,4

Primo fra tutti come risultato fornito da Google con stringa di ricerca

"corpus di italiano" troviamo:

Accademia della Crusca. Qui, tra i collegamenti utili, troviamo una bacheca di raccolta banche dati, corpora e archivi testuali. Seppur distinti fra italiano scritto e parlato, italiano antico, italiano di apprendenti, giornalistico, di internet, tecnico e trasmesso, nella bacheca troviamo assieme, come si evince dal titolo,

Secondo risultato proposto dal motore di ricerca Google è la definizione di **corpus di Treccani** che però non si limita alla stessa. Treccani dopo la definizione fornisce un elenco di alcuni corpus conosciuti, con link agli stessi e una breve descrizione.

Oltre a ciò questa voce è correlata da alcuni esempi di uso come possiamo osservare. Questo è certamente uno strumento utile per chi si sta approcciando per la prima volta ai corpora, ma non esaustivo per una classe di giovani studenti o studenti di italiano L2 che devono utilizzare i corpora come strumento didattico.

Corpora dell’italiano d’uso

Corpora dell’italiano scritto
CORIS/CODIS online
corpus di Italiano Scritto contemporaneo ricco di circa 100 milioni di parole.
Elaborato e prodotto dal Centro Interfacoltà di Linguistica Teorica e Applicata (CILTA) dell’Università di Bologna, prevede una licenza per l’accesso alla versione completa.

Corso e Lessico di Frequenza dell’Italiano Scritto (CILFIS)
il corpus di riferimento è costituito da testi tratti da quotidiani del periodo 1992 - 1994 (La Repubblica, La Stampa, Il Corriere della Sera), periodici e libri, compreso anche i libri letti dai madri e padri o professionisti. Conta di 3.150.078 record di testi. Al progetto partecipano i seguenti soggetti: Scuola Normale Superiore (Pisa), Istituto di Scienze e Tecnologie della Cognizione del CNR (Roma), Università di Salerno, Istituto di Linguistica Computazionale, Unità Staccata di Genova del CNR, Università di L'Aquila.

Dizionario italiano multimediale e multilingue d’Ortopografia e di Pronuncia
della RAI

versione online del vocabolario redatto a partire dal 1950 da Bruno Migliorini, Carlo Tagliavini e Piero Fiorelli (rv, agg e accr: da P. Fiorelli e T. Borri) per la sede fiorentina della RAI Televisione Italiana, ora disponibile in versione multilingue. Il corpus contiene oltre 92.000 voci di lessico della lingua italiana e oltre 37.000 di una sessantina di lingue diverse, presenta anche la registrazione fonetica dei voci.

la Repubblica Corpus
corpus molto ampio (circa 380mila parole) del lessico del quotidiano *la Repubblica*. Nel progetto curato dall’Università di Bologna, il corpus è stato immesso nel database e categorizzato per genere e tipi; gli articoli nel corpus sono strutturati nelle seguenti parti: titolo, sottotitolo, sommario, testo.

Corpora dell’italiano parlato

API Federico del Carbone Gallo
progetto coordinato dal prof. Federico Albano Leonini (CIRASS - Napoli) cui hanno partecipato la **Scuola Normale Sup. di Pisa**, il **CIRAS e il Dip. di Neuroscienze dell’Univ. Federico II** di Napoli, l’**Istituto Univ. Orientale di Napoli**, il **Politechnico di Bari**, l’**Uvr**, del Piemonte Orientale, l’**Univ. “Ce’ Foscari” di Venezia** e l’**Univ. di Pisa**.
Il corpus, raccolto in tre città italiane (Napoli, Bari, Pisa), consiste essenzialmente di dialoghi *in es* e comprende anche un campione di parlato infantile di bambini sordi e nonnourediti.

AVIP/API
progetto coordinato dal prof. Pier Marco Bertinetto (Scuola Normale Sup. - Pisa) cui hanno partecipato la **Scuola Normale Sup. di Pisa**, il **CIRAS e il Dip. di Neuroscienze dell’Univ. Federico II** di Napoli, l’**Istituto Univ. Orientale di Napoli** e il **Politechnico di Bari**. Sostanzialmente contiene lo stesso corpus del progetto API, che ne amplia lo spettro d’indagine, il campione e gli enti

Terzo risultato fornito da Google con query “corpus di italiano” è un ulteriore sito web **“Biblioteca Scuola Normale Superiore”** che riunisce in un unica bacheca dizionari e corpora dell’italiano. Distingue fra corpora dell’italiano scritto e parlato, settoriale e antico fornendone una breve descrizione e link alla risorsa originale.

Da questi esempi si può quindi comprendere che esistono vari siti che forniscono una bacheca da cui collegarsi a vari corpora **ma non forniscono delle istruzioni d’uso per gli stessi**. Nota positiva che va ovviamente mantenuta è la distinzione fra italiano d’uso e antico e italiano scritto e parlato, essenziale per poter aiutare insegnanti e studenti ad orientarsi nella scelta del corpus giusto per la lezione o gli esercizi.

Proseguendo nell’analisi del mercato a livello di contenuti esistenti, proponendo come stringa di ricerca **“esercizi corpus italiano”** sul motore Google otteniamo soprattutto dei risultati **PDF** di vari articoli circa l’utilità dell’utilizzo dei corpora a fini didattici, per la creazione di esercizi, ma non una raccolta di eventuali esercizi da proporre.

The screenshot shows the Valico platform interface. At the top, there are two tabs: "CORPORA UNITO" and "VALIDO". Below them is a map of Italy with several colored dots (green, yellow, red) scattered across the country, likely representing different exercise locations or types. To the left of the map is a sidebar with a list of links: Valico, Vinca, MorfoWeb, Vignette, Esercizi, Corpora Comparabili, and Pubblicazioni. At the bottom of the sidebar, there are links for "Esercizi on line" and "Scarica gli esercizi". The main area below the map has a "Mappa" and "Satellite" button, and a "CORPORA UNITO" logo.

Non si può ignorare in questo momento di analisi l’apporto fornito dalla piattaforma **Valico**, che ha una sezione esercizi, in cui è possibile osservare e scaricare alcuni esercizi proposti a studenti italiano l2, ed è inoltre possibile scaricare i risultati ottenuti in formato excel. Sebbene questi esercizi hanno un origine diversa, ovvero, vengono progettati dopo aver osservato i risultati ottenuti dal learner corpus e costruiti per valutare correttamente gli “errori” commessi dall’apprendente, è senz’altro interessante la scelta di mostrare i risultati ottenuti distinguendoli per tipo di esercizi

tramite legenda e ponendoli a seconda della provenienza degli apprendenti su una mappa geografia. Questo modo di esporre i risultati può senz’altro essere implementato su Corpora in Classe nella sezione learner corpora, per fornire agli studiosi uno strumento di visualizzazione ulteriore dei risultati.