

CS 345: Homework 2 (Due: Sep 30, 2013 before 10:01am)

Rules. This is to be completed no later than the start of class on the day it is due. Earlier submission is also possible by email (read Handout 1).

Problem 1. (30 POINTS) Estimating the corpus size of Google and Bing

We are asking you to search in queries 1-6 certain uncorrelated terms, and in queries 7-9 to improvise and determine yourselves something better than 1-6. Choose two words for w_1 and w_2 that are not correlated to obtain as high and accurate estimate values for the corpus size of Bing and Google as possible. For a counterexample a poor choice for the pair (w_1, w_2) would be (tropical, fish) or (tropical, storm). The textbook suggests (tropical, Lincoln). Use something else (eg. your last name for w_1 or w_2) and NOT along the lines of those three examples. Because this is not a verbatim query for words w_1 , w_2 we use the $[]$ symbols in queries 7-9 instead of $\langle \rangle$.

No	Query Term(s)	GOOGLE No of hits	BING No of hits
1	$\langle tiger \rangle$		
2	$\langle algorithm \rangle$		
3	$\langle tiger algorithm \rangle$		
E1	What is the corpus size based on 1-3?		
4	$\langle cs345 \rangle$		
5	$\langle delicious \rangle$		
6	$\langle delicious cs345 \rangle$		
E2	What is the corpus size based on 4-6?		
7	$[w_1]$		
8	$[w_2]$		
9	$[w_1 w_2]$		
E3	What is the corpus size based on 7-9?		

Table 1: Table for Q3.7

Based on the hits reported for 1-3, 4-6 and 7-9 estimate through probabilistic reasoning as specified in class, textbook, and the notes what your estimate of the corpus size of Google and Bing would be by filling in lines E1, E2, E3. Are these estimates consistent with those reported elsewhere say in Homework 1? Explain. Can you improve upon them some other way? Explain.

Problem 2. (40 POINTS) Crawling...

(a) **Unix command line wget.** Fetch a copy of the course web-page.

(a.1) Login on an AFS account say on `afsconnect1.njit.edu`. You are going to use the `wget` command to get a copy of the course web page (a subset of all the files that are accessible to you). For help about the options of `wget` type `man wget` under AFS or `wget --help`. Then work as instructed below.

(a.2) Although `wget` by default can be used to copy a remote file locally, in order to grab information from a remote site that consists of more than one file (eg directories) you need to use some recursive options of `wget`. Read the manual for details on options of `wget`. Note also than on a UNIX system by typing `du -s dirname` you can get info about the file size of all files in the directory (and subdirectories) named `dirname`. The size may be in kilobytes or multiples of 512B (read the manual). Certain details are missing because you are expected to read the manual pages and also to improvise. Handout 0 should help as well.

(a.3) After you get familiarized with `wget` and its options do the following.

a.3.1 Grab the course web-page located `... alexg/courses/cs345/index.html`. Don't ask why part of the URL is missing. If you do it properly, you will collect close to 3MiB of information (note that i typed MiB not MB). Use proper options of `wget` to achieve the objective. What is the command (all arguments to `wget`) that you used to collect this information? **(5 points)**

a.3.2 `Wget` generates a log file while transferring files. You can capture into a file by redirecting output to a file as partially shown next

```
wget .... alexg/courses/cs345/index.html > & capturefilename
```

a.3.3 How many (remote) files did you capture locally? **(5 points)**

a.3.4 What was the total size (volume) of those files ? **(5 points)**

a.3.5 How many URLs were not found to contain information (and thus no file was transferred locally) ? What error did they generate? **(5 points)**

a.3.6 How many URLs were protected behind a password authorization? What error did they generate? **(5 points)**

(b) **webcrawler.** Download the zip file `webcrawler_2.1.zip`. It is available in the handouts section of the course web-page. Unzip it by typing `unzip webcrawler_2.1.zip` and it will unfold 7 Python files. You will run Python on the same machine you used in part (a). There are several versions of Python on AFS but the one that works correctly is version 2.6. Thus your invocation will involve something along the lines

```
/usr/bin/python2.6 modular_crawler.py
```

There are some additional options (2 to 3) required but definitely use the last option which is a numeric value that gives an upper bound for the number of files to be downloaded. The tricky part of this assignment is to formulate precisely the first and second options. The previous problem's issues (eg a.3.1 might offer some insight). (You don't want risk exploring all of the web.)

b.1 Run `webcrawler` on the course web-page and explore all options for a URL as specified in Handout 0. For the more productive and complete of these options, and the last option held at 100 what is the URL of the 100th page crawled? (If you do it properly the log file would grow to a size of at least 150KiB, if you do it incompletely it would be around 50KiB.) **(5 points)**

b.2 How many times is there a URL that contains the digit sequence 71854? What does it relate to? **(5 points)**

b.3 You can redirect the log into a file the way you did it before in a.3.2. A command `sort filename` on a UNIX system will sort the lines of `filename`. A `sort -u filename` will also eliminate duplicates. What is the total number of unique URLs listed? (Among the lines of the log the URLs start from character `http` as the first 4 characters of a line.) **(5 points)**

Problem 3. (20 POINTS)

Perform the following study for Google and Bing. We execute two queries in Google and Bing and the collected results can be found in the last 4 pages of this homework as well as in the Handouts section of the course web-page. One query(Q1) is CS 485 NJIT and the other(Q2) is CS 485 NJIT Web-search. The CS 485 course number is a generic special topics one and thus other NJIT CS 485 courses that ARE NOT RELEVANT TO the CS 485 NJIT Web-Search offered in Fall 2010 and Fall 2011 are output. Thus for Q1/Q2 relevant documents are only those of to a CS 485 course offered at NJIT with subtitle Web-search as can be recognized either by having the name or web-page of the instructor listed or the title of the course in the 3-5 lines of each result. (Each one of the 4 pages contains 10 results.)

The first query is not very well thought-out. The second is more refined. You are asked to evaluate those two queries and also Google's and Bing's capabilities in generating relevant documents. We compared both engines in HW1. We evaluate the engines' search capabilities now. **When you examine the output results tabulate by starting with Google Q1, then Q2, and then move to Bing and do not deviate from this order! Do not include ad links in considering relevant documents!**

(a) Tabulate for each engine and query and engine, the number of hits reported in the result page attached (easy always the same), the number of relevant documents (be cautious when you read the entries), and the precision for the single result page reported (do not get confused with the huge number reported by Google or Bing!).

(b) Fill the following table with data (columns 2 and 3). If Google is better give Google 1pt and Bing 0 in columns 4 and 5. If Bing is better, give Bing 1pt and Google 0. If it is a tie, give each one 1pt each. Add up the points of columns 4 and 5 and provide sum in the last row. Who wins?

	Values		Points	
	Google	Bing	Google	Bing
# Retrieved docs for Q1				
# Relevant docs for Q1				
Precision for Q1				
# Retrieved docs for Q2				
# Relevant docs for Q2				
Precision for Q2				
=====				
Number of point wins (sum)	-	-		

Date Posted: 9/23/2013

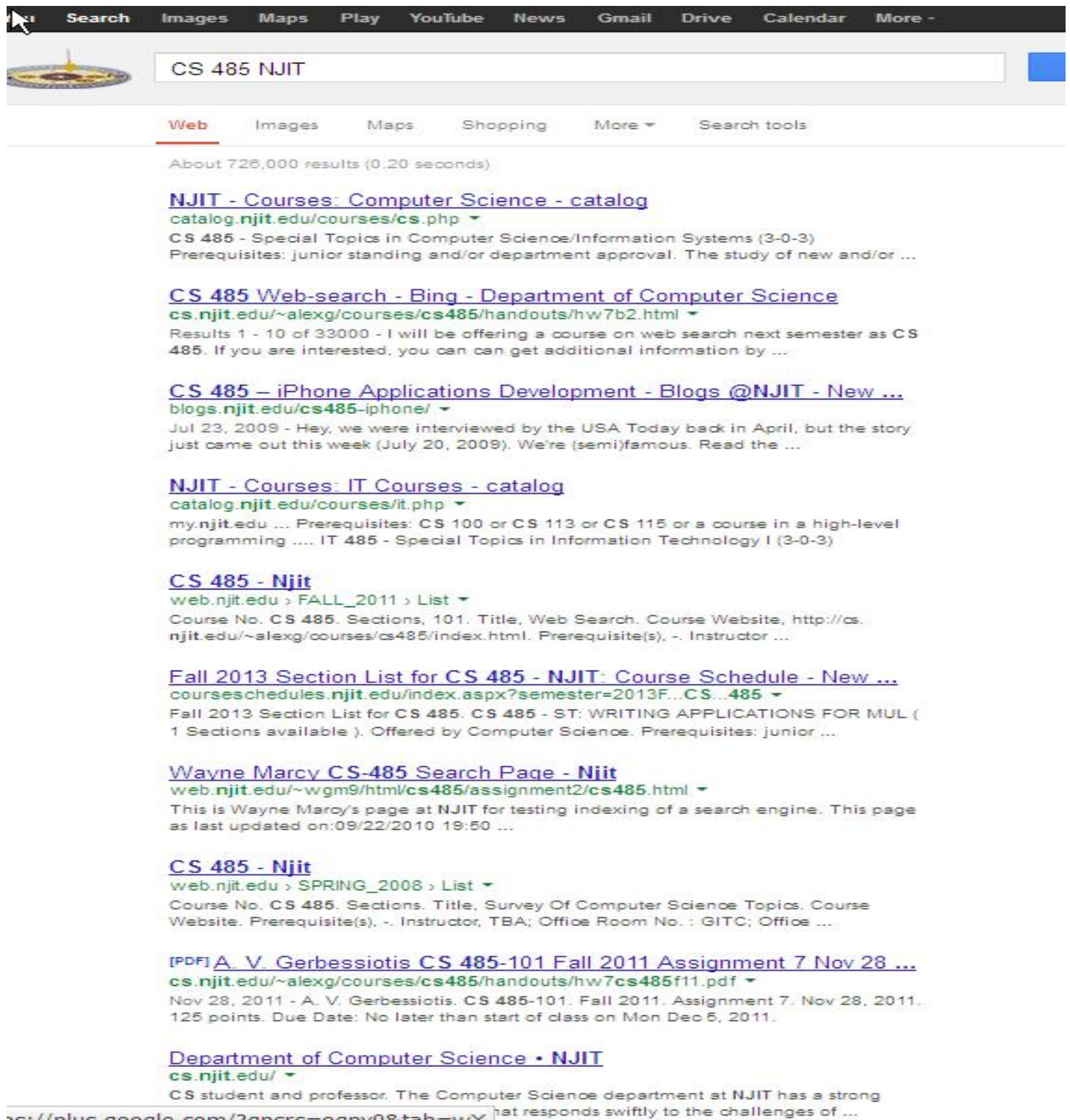


Figure 1: Google for CS 485 NJIT

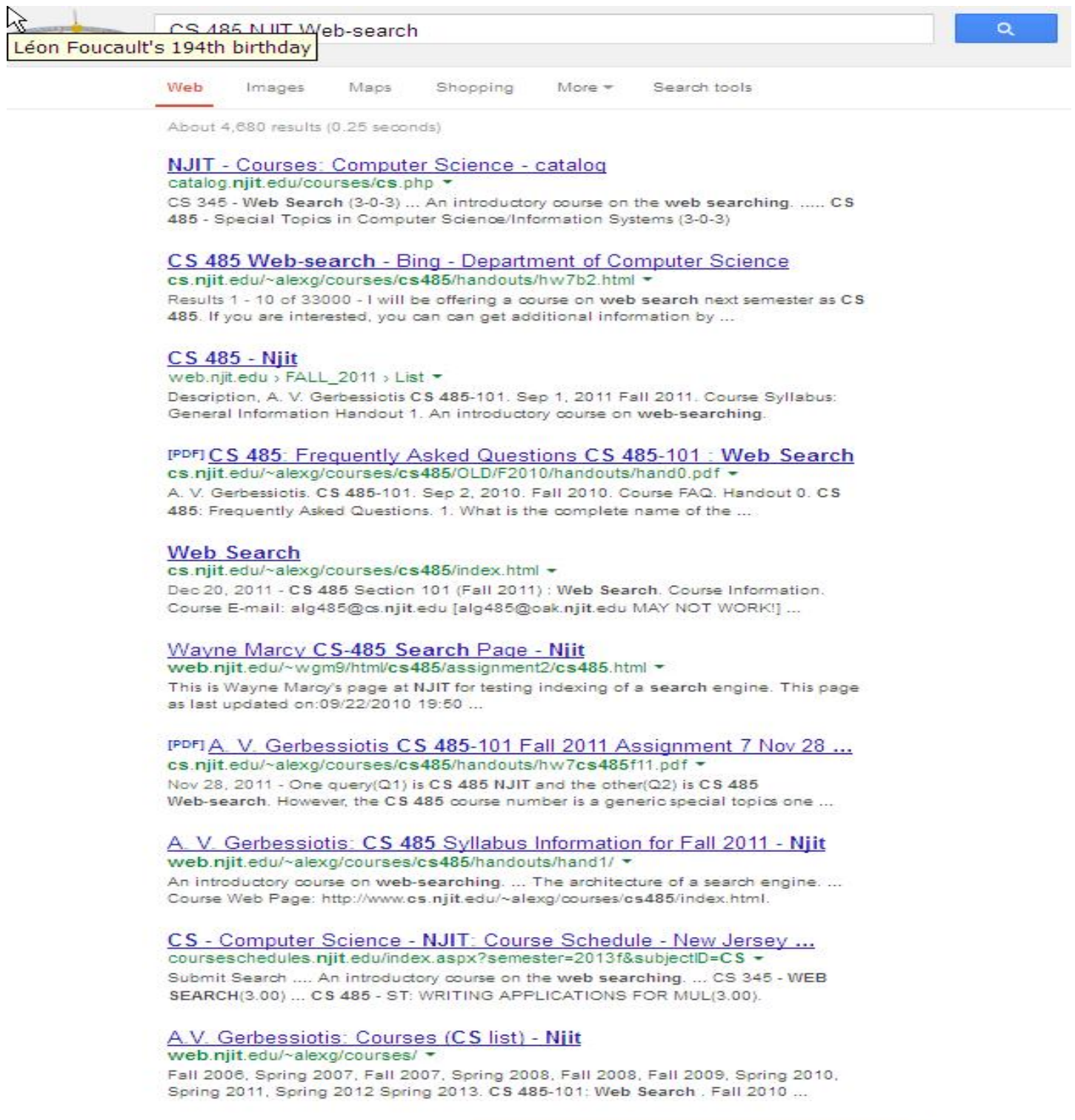


Figure 2: Google for CS 485 NJIT Web-search



WEB IMAGES VIDEOS MAPS NEWS MORE

CS 485 NJIT



1,260,000 RESULTS Any time ▾

[A. V. Gerbessiotis: Course CS485 Section 101 \(Fall 2011\): Web ...](#)

cs.njit.edu/~alexg/courses/cs485/index.html ↗

Course E-mail: alg485@cs.njit.edu [alg485@oak.njit.edu MAY NOT WORK!]
Instructor: A. Gerbessiotis; Office Hours: [Office hour schedule](#). Course Hours: Mon 6-9pm

[NJIT - Courses: Computer Science - NJIT - catalog: Home Page](#)

catalog.njit.edu/courses/cs.php ↗

CS 103*** - Computer Science with Business Problems ... and getting acclimated to NJIT. ... CS 485 - Special Topics in Computer Science/Information Systems ...

[CS 485 – iPhone Applications Development](#)

<https://blogs.njit.edu/cs485-iphone> ↗

Hey, we were interviewed by the USA Today back in April, but the story just came out this week (July 20, 2009). We're (semi)famous. Read the USA Today article and ...

[NJIT: Course Schedule](#)

courseschedules.njit.edu/index.aspx?semester=2009S&subjectID=CS&... ↗

My NJIT | CALENDAR | DIRECTORY | A-Z LINKS | CONTACT US | Search.
About Admissions Academics Student Life Research Working ...

[CS 485: Frequently Asked Questions - Department of Computer ...](#)

cs.njit.edu/~alexg/courses/cs485/handouts/hand0.pdf · PDF file

A. V. Gerbessiotis CS 485-101 Sep 1, 2011 Fall 2011 Course FAQ Handout 0 CS 485: Frequently Asked Questions 1. What is the complete name of the course?

[CS 485 – iPhone Applications Development» Blog Archive » USA ...](#)

<https://blogs.njit.edu/cs485-iphone/2009/07/23/usa-today-story...> ↗

Hey, we were interviewed by the USA Today back in April, but the story just came out this week (July 20, 2009). We're (semi)famous. Read the USA Today article and ...

[NJIT: Course Schedule](#)

courseschedules.njit.edu/index.aspx?semester=2012F&subjectID=CS ↗

An introductory course in computer science and programming ... and getting acclimated to NJIT. ... CS 485 - ST: ANDROID APPLICATION ...

[wgm9 NJIT Home Page](#)

web.njit.edu/~wgm9 ↗

Wayne Marcy aka wgm9@njit.edu Home Page This page as last updated on: 12/02/2010 5:37 PM. ... CS 485 Assignments Web Search: CS-485 Assignment 2 IS 270 Class ...

[CS](#)

www.njit.edu/registrar/schedules/courses/fall/2012F.CS.html ↗

<http://www.njit.edu/registrar/exams/index.php>: 03.00 ; 003: ... 485 ST: ANDROID APPLICATION DEVELOPM. Sect ... 491 COMPUTER SCIENCE PROJECT - ...

[NJIT - Courses: IT Courses - NJIT - catalog: Home Page](#)

catalog.njit.edu/courses/it.php ↗

... REGISTRAR | NJIT | CATALOG ARCHIVE: Information ... CS 100 or CS 113 or CS 115 or a course in a high-level programming language as ... IT 485 - Special Topics in ...



Connect to Facebook to your friends know.

[Learn more](#)

[Connect with Facebook](#)

Figure 3: Bing for CS 485 NJIT



WEB IMAGES VIDEOS MAPS NEWS MORE

CS 485 NJIT Web-search



ALSO TRY: NJIT CS 241 - CS 434 NJIT - CS 435 NJIT

38,600,000 RESULTS Any time ▾

[A. V. Gerbessiotis: Course CS485 Section 101 \(Fall 2011\): Web Search](#)

cs.njit.edu/~alexg/courses/cs485/index.html ↗

Course E-mail: alg485@cs.njit.edu [alg485@oak.njit.edu MAY NOT WORK!]
Instructor: A. Gerbessiotis; Office Hours: [Office hour schedule](#); Course Hours: Mon 6-9pm

[CS 485: Frequently Asked Questions - Department of Computer ...](#)

cs.njit.edu/~alexg/courses/cs485/handouts/hand0.pdf - PDF file

CS 485-101 : Web Search ... Use the alg485@cs.njit.edu address. The instructor or the assistant will respond (if there is going to be an assistant).

RELATED SEA

[NJIT CS 24](#)

[CS 434 NJI](#)

[CS 435 NJI](#)

[New Jersey](#)

[NYIT Comp](#)

[NJIT Step I](#)

[CS 667 Alg](#)

[NJIT Acade](#)

Related searches for CS 485 NJIT Web-search

[NJIT CS 241](#)

[New Jersey Inst of Technology](#)

[CS 434 NJIT](#)

[NYIT Computer Science](#)

[CS 435 NJIT](#)

[NJIT Step Program](#)

[NJIT - Courses: Computer Science - NJIT - catalog: Home Page](#)

catalog.njit.edu/courses/cs.php ↗

Computer Science: UNDERGRADUATE COURSES: CS 100 - Roadmap to Computing (3-0-3) An introduction to programming and problem solving skills using Python or ...

[wgm9 NJIT Home Page](#)

web.njit.edu/~wgm9 ↗

Wayne Marcy aka wgm9@njit.edu Home Page This page as last updated on: 12/02/2010 5:37 PM. ... CS 485 Assignments Web Search: CS-485 Assignment 2 IS 270 Class ...

[New Jersey Institute of Technology](#)

www.njit.edu ↗

NJIT is a public research university committed to educating a wide range of students to achieve their full potential, preparing them for entry into professional ...

[NJIT: Course Schedule](#)

courseschedules.njit.edu/index.aspx?semester=2012f&subjectID=CS ↗

An introductory course in computer science and programming ... CS 345 - WEB SEARCH ... CS 485 - ST: ANDROID APPLICATION ...

[Computer Science Course Information - UCS Home](#)

web.njit.edu/cs/cs_courses/index.php?d=lg&sem=FALL_2011 ↗

Computer Science Course Information: Level >Graduate >FALL_2011 >List No Information as of yet Partial Information Off Campus Course E-Learning Section

[CS](#)

www.njit.edu/registrar/schedules/courses/fall/2012F.CS.html ↗

<http://www.njit.edu/registrar/exams/index.php>: 03.00 : 003: ... 345 WEB SEARCH . Sect Call # Days Times Room ... 491 COMPUTER SCIENCE PROJECT - HONOR. Sect ...

[CS 485 – iPhone Applications Development](#)

<https://blogs.njit.edu/cs485-iphone> ↗

CS 485 – iPhone ... David is Web Architect for NJIT and leads much of our user interface and design efforts. ... New Jersey Institute of Technology, University ...

[A.V. Gerbessiotis: Courses \(CS list\) - UCS Home](#)

web.njit.edu/~alexg/courses/index.html ↗

Courses taught and being offered at the CS Department at NJIT [Last Update: Aug 2013] Graduate level courses. ... CS 485-101: Web Search . Fall 2010, Fall 2011.

Figure 4: Bing for CS 485 NJIT Web-search