

Polls Syllabus Assignments Sorting Implementation Summary

BIOSTAT615/815 - Objectives

- 1 Equip the ability to **implement** computational and/or statistical **ideas** into working **software programs**.
- 2 Learn computational cost management in developing statistical methods.
 - ✓ Understand the practical importance of computation cost in many statistical inference applications.
 - ✓ Develop the ability to estimate computational time and memory required for an algorithm given data size.
 - ✓ Develop the ability to improve computation efficiency of existing algorithms and to optimize the cost/accuracy trade-off.

Hyun Min Kang Biostatistics 615/815 - Lecture 1 September 4th, 2012 9 / 46

Why Is Statistical Computing Important?

- Now is "Big Data" era
 - \checkmark Next-generation sequencing studies often become > 100 TB in size
 - ✓ Storing MRI sessions of 1,500 subjects requires 5.4TB of data.
 - ✓ Google Maps has over 20 petabytes of data
 - Computation of simple statistics (e.g. mean) will take a long time.
- 2 Efficiency affects the feasibility of statistical methods
 - ✓ In statistical inference from genome-wide data, it is not common that more accurate methods takes much longer time than less accurate approximation (e.g. 40 years vs 1 day).
 - ✓ Implementation with low-level languages such as C++ has more room for speed improvements than higher level languages such as R or SAS.
 - ✓ Even with the same language, different algorithms can result in substantial gain in speed without losing accuracy.



BIOSTAT615/815 - Objectives

- 1 Equip the ability to **implement** computational and/or statistical **ideas** into working **software programs**.
- 2 Learn computational cost management in developing statistical methods.
- 3 Understand numerical methods useful for statistical inference
 - ✓ Learn numerical optimization methods for solving analytically intractable problems computationally
 - ✓ Understand a variety of randomized algorithms for robust and efficient estimation of computationally intractable problems to obtain deterministic solution.

Hyun Min Kang Biostatistics 615/815 - Lecture 1 September 4th, 2012 10 / 46



- **1** C++ Basics and Introductory Algorithms
 - Computational Time Complexity
 - Sorting
 - Divide and Conquer Algorithms
 - Searching
 - Key Data Structure and Standard Template Libraries
 - Dynamic Programming
 - Hidden Markov Models
 - Interfacing between C++ and R

15/815 - Lecture 1 September 4th, 2012 11 / 46 Hyun Min Kang Biostatistics 615/815 - Lecture 1 September 4th, 2012 12 / 46



- **1** C++ Basics and Introductory Algorithms
- 2 Numerical Methods and Randomized Algorithms
 - Random Numbers
 - Matrix Operations and Least Square Methods
 - Importance Sampling
 - Expectation-Maximization
 - Markov-Chain Monte Carlo (MCMC) Methods
 - Simulated Annealing
 - Gibbs Sampling

Hyun Min Kang Biostatistics 615/815 - Lecture 1 September 4th, 2012 13 / 4

Polls Syllabus Assignments Sorting Implementation Summary

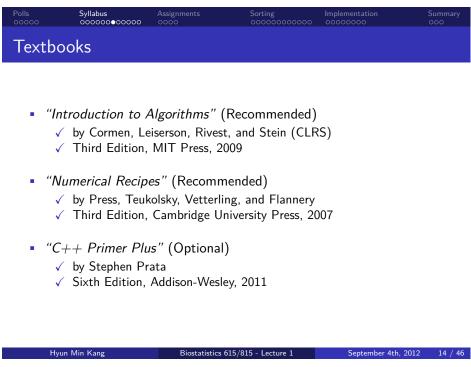
Assignments and Grading

BIOSTAT615

- Biweekly Assignments 50%
- Midterm Exam 25%
- Final Exam 25%

BIOSTAT815

- Biweekly Assignments 33%
 - Expected to solve extra problems on top of 615 assignments
- Midterm Exam 17%
- Final Exam 17%
- Projects, to be completed in pairs 33%





- Students may have little experience in C++ programming
 - But students must be strongly motivated to learn C++ programming
 - Students with no/little experience in C++ are expected to spend additional hours than other students to accomplish homework
- Students should be familiar with basic concept of probability distribution, hypothesis testing, and simple regression
 - But BIOSTAT601 is not strictly required
- Students are expected to know or willing to learn basics of R language.

Hyun Min Kang Biostatistics 615/815 - Lecture 1 September 4th, 2012 16 / 46

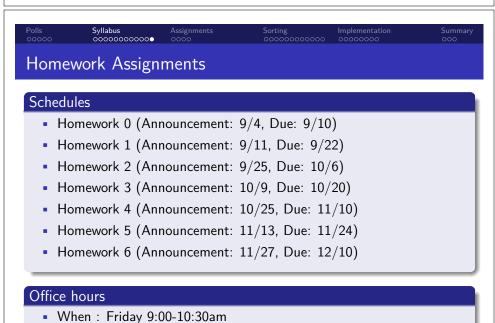


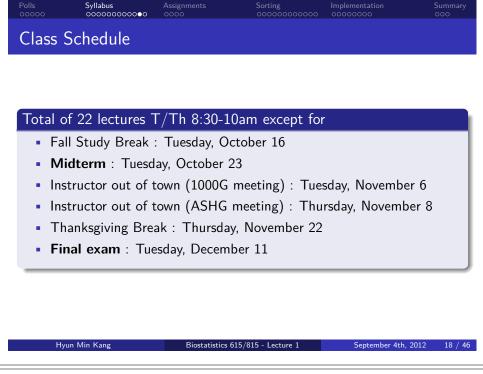
- Students are expected to have decent experience in C/C++/Java programming
 - And expected to be fluent in C++ enough to be able to accomplish term projects.
 - List of suggested projects will be announced in the next week.
 - Students must be strongly motivated to learn C++ programming
- Students should be familiar with basic concept of probability distribution, hypothesis testing, and simple regression
 - But BIOSTAT601 is not strictly required

Where: M4531 SPH II

 Students are expected to know or willing to learn basics of R language.

Hyun Min Kang Biostatistics 615/815 - Lecture 1 September 4th, 2012 17 / 46







- Programming language must be in C++
- Assignments must be submitted in both of two formats
 - Google Docs document shared only to instructor and grader
 - Compressed source codes ready to run in Linux environment.
- The grade of late submission of homework will be multiplied by a factor of $e^{-\Delta m/14400}$, where Δm is the difference between the due and submission time in minutes

Hyun Min Kang Biostatistics 615/815 - Lecture 1 September 4th, 2012 20 / 46



Steps to share a google document

- 1 Use your umich google account than other google accounts.
- Make sure to rename the document title into the following format : "[BIOSTAT615] (or [BIOSTAT815]) Homework 1 John Doe"
- 3 Click "Share" and type hmkang@umich.edu in "Add people", and give "Can edit" permission and do NOT modify after submission due.
- 4 When grading is finished, you will be notified by email.
- **6** Grader edit documents with comments and corrections if necessary. You can see comments and revision history.
- Make sure not to modify the the submitted homework until the grading finishes. The late submission penalty will rely on the last modified time.

Hyun Min Kang Biostatistics 615/815 - Lecture 1 September 4th, 2012 21 / 46

Honor code

- Honor code is STRONGLY enforced throughout the course.
 - The key principle is that all the code you produce must be on your own.
 - See http://www.sph.umich.edu/academics/policies/conduct.html for details.
- You are NOT allowed to share any piece of your homework with your colleagues electronically (e.g. via E-mail or IM), or by a hard copy.
- Discussion between students are generally encouraged
 - You may help your colleague setting up the programming environment necessary for the homework.
 - You may help debugging your colleagues' homework by sharing your trial and errors only up to non-significant fraction of your homework
 - Significant fraction of help can be granted if notified to the instructor, so that the contribution can be reflected in the assessment.
- If a break of honor code is identified, your entire homework (or exam) will be graded as zero, while incomplete submission of homework assignment will be considered for partial credit.

Hyun Min Kang Biostatistics 615/815 - Lecture 1 September 4th, 2012 23 / 46

Assignments are Very Important

- The main objective of the course is to develop the ability to implement software on your own.
- You will meet the instructor's teaching goal for BIOSTAT615 if you can accomplish your homework on your own.
 - If you got a good grade from BIOSTAT615, it should suggest that you are likely capable of use C++ and numerical methods for your research.
- You will meet the instructor's teaching goal for BIOSTAT815 if you can implement sophisticated statistical methods that are useful and convenient for others' use.
 - If you got a good grade from BIOSTAT815, it should suggest that you are likely capable of develop and release a software for your research community.

Hyun Min Kang Biostatistics 615/815 - Lecture 1 September 4th, 2012 22 / 46

Assignments

Sorting Imp

Implementation

Algorithms

An Informal Definition

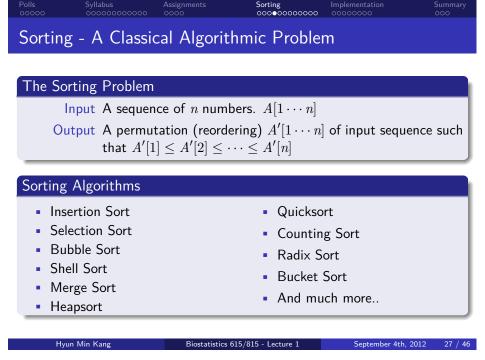
- An **algorithm** is a sequence of well-defined computational steps
- that takes a set of values as input
- and produces a set of values as output

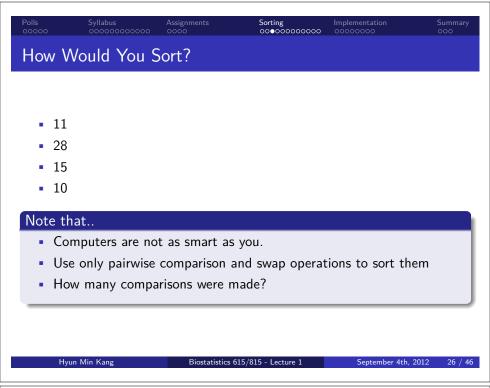
Key Features of Good Algorithms

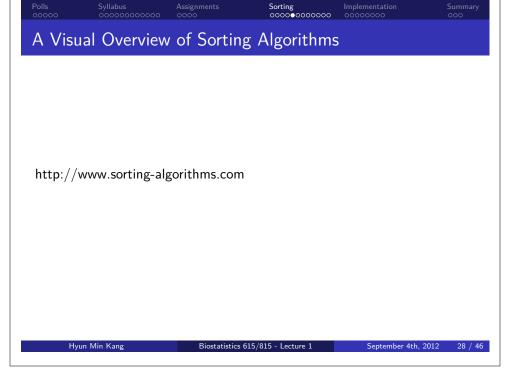
- Correctness
 - $\checkmark\,$ Algorithms must produce correct outputs across all legitimate inputs
- Efficiency
 - \checkmark Time efficiency : Consume as small computational time as possible.
 - \checkmark Space efficiency : Consume as small memory / storage as possible
- Simplicity
 - ✓ Concise to write down & Easy to interpret.

Hyun Min Kang Biostatistics 615/815 - Lecture 1 September 4th, 2012 24 / 46









Today - Insertion Sort

http://www.sorting-algorithms.com/insertion-sort

Algorithm InsertionSort

```
Data: An unsorted list A[1 \cdots n]
Result: The list A[1 \cdots n] is sorted
for j=2 to n do
   key = A[j];
   i = i - 1:
   while i > 0 and A[i] > key do
      A[i+1] = A[i];
      i = i - 1;
   end
   A[i+1] = key;
end
```

Key Idea of Insertion Sort

- For k-th step, assume that elements $a[1], \dots, a[k-1]$ are already sorted in order.
- Locate a[k] between index $1, \dots, k$ so that $a[1], \dots, a[k]$ are in order
- Move the focus to k+1-th element and repeat the same step

Correctness of InsertionSort

Loop Invariant

At the start of each iteration, $A[1 \cdots j-1]$ is loop invariant iff:

- $A[1 \cdots j-1]$ consist of elements originally in $A[1 \cdots j-1]$.
- $A[1 \cdots j-1]$ is in sorted order.

A Strategy to Prove Correctness

Initialization Loop invariant is true prior to the first iteration

Maintenance If the loop invariant is true at the start of an iteration, it remains true at the start of next iteration

Termination When the loop terminates, the loop invariant gives us a useful property to show the correctness of the algorithm

Correctness Proof (Informal) of INSERTIONSORT

Initialization

• When j=2, $A[1\cdots j-1]=A[1]$ is trivially loop invariant.

Maintenance

If $A[1 \cdots j-1]$ maintains loop invariant at iteration j, at iteration j+1:

- $A[j+1\cdots n]$ is unmodified, so $A[1\cdots j]$ consists of original elements.
- $A[1\cdots i]$ remains sorted because it has not modified.
- $A[i+2\cdots j]$ remains sorted because it shifted from $A[i+1\cdots j-1]$
- $A[i] \leq A[i+1] \leq A[i+2]$, thus $A[1 \cdots j]$ is sorted and loop invariant

Termination

• When the loop terminates (j=n+1), $A[1\cdots j-1]=A[1\cdots n]$ maintains loop invariant, thus sorted.

Hyun Min Kang

Biostatistics 615/815 - Lecture

September 4th, 2012

•

Summary - Insertion Sort

- One of the most intuitive sorting algorithm
- Correctness can be proved by induction using loop invariant property
- Time complexity is $O(n^2)$.

Polls Syllabus Assignments Sorting Implementation Summary 00000 0000000000 000000000 00000000 00000000

Time Complexity of INSERTIONSORT

$\begin{aligned} & \text{for } j = 2 \text{ to } n \\ & \text{do} \\ & key = A[j]; & c_2(n-1) \\ & i = j-1; & c_3(n-1) \\ & \text{while } i > 0 \text{ and } A[i] > key & c_4 \sum_{j=2}^n j \\ & \text{do} \\ & & A[i+1] = A[i]; & c_5 \sum_{j=2}^n (j-1) \\ & i = i-1; & c_6 \sum_{j=2}^n (j-1) \\ & \text{end} \\ & A[i+1] = key; & c_7(n-1) \end{aligned}$

$$T(n) = \frac{c_4 + c_5 + c_6}{2}n^2 + \frac{2(c_1 + c_2 + c_3 + c_7) + c_4 - c_5 - c_6}{2}n - (c_2 + c_3 + c_4 + c_7)$$

= $\Theta(n^2)$

Hyun Min Kang

end

Worst Case Analysis

Biostatistics 615/815 - Lecture 1

September 4th, 2012

Environment for homework: Connect to scs.itd.umich.edu

See http://www.itcs.umich.edu/scs/ and http://www.itcs.umich.edu/ssh/ for details

- Windows environment : Use PuTTY for command line and WinSCP for file transfer
- MacOS X or UNIX

Command line ssh uniqname@scs.itd.umich.edu File transfer scp [your-file]

uniqname@scs.itd.umich.edu:[path]/[to]/[destination]

Tip: Add the following line in ~/.cshrc file for more convenient command line (which shows current working directory).

set prompt="`whoami`@`hostname -s`:%~\$ "

Hyun Min Kang Biostatistics 615/815 - Lecture 1 September 4th, 2012 36 / 46



Steps for Homework 0

- ① Create a directory Private/biostat615/hw0/ under your home directory
 - Create a directory Private/biostat615/hw0/ using WinSCP
 - Or type mkdir --p ~/Private/biostat/hw0/ in the command line
 - Make sure your homework is in private space. If it is accessible by someone else, your homework will be discarded.
- 2 Create (or copy) a file (e.g. Private/biostat/hw0/helloWorld.cpp)
 - Directly type vi Private/biostat615/hw0/helloWorld.cpp
 - Or copy a file remotely using WinSCP or scp
- 3 Use basic Unix commands (e.g. cd ~/my/path/, pwd) to navigate between directories.
- 4 Compile and run the program (see the next slide)
- **5** Before submission, remove all the executable and object (.o) files, and compress your code (under Private/biostat615)

```
cd ~/Private/biostat615/
tar czvf uniqname_hw0.tar.gz hw0/
```

Hyun Min Kang

Biostatistics 615/815 - Lecture 1

September 4th, 2012

2 37 / 46

Implementing InsertionSort Algorithm

```
insertionSort.cpp - main() function
#include <iostream>
#include <vector>
void printArray(std::vector<int>& A);  // declared here, defined later
void insertionSort(std::vector<int>& A); // declared here, defined later
int main(int argc, char** argv) {
  std::vector<int> v; // contains array of unsorted/sorted values
                     // temporary value to take integer input
  while ( std::cin >> tok ) // read an integer from standard input
   v.push back(tok);
                            // and add to the array
  std::cout << "Before sorting:";</pre>
  printArray(v); // print the unsorted values
  insertionSort(v); // perform insertion sort
  std::cout << "After sorting:";</pre>
  printArray(v); // print the sorted values
  return 0;
```

Biostatistics 615/815 - Lecture 1

```
    Polls
    Syllabus
    Assignments
    Sorting
    Implementation
    Summary

    00000
    00000000000
    00000000000
    000000000
    0000000000
```

Getting Started with C++

```
Writing helloWorld.cpp

#include <iostream> // import input/output handling library
int main(int argc, char** argv) {
  std::cout << "Hello, World" << std::endl;
  return 0; // program exits normally</pre>
```

Compiling helloWorld.cpp

```
user@host:~$ cd ~/Private/biostat615/hw0/
user@host:~/Private/biostat615/hw0$ g++ -O -o helloWorld helloWorld.cpp
```

-O option will increase the computational speed. Use -g when you need debugging (e.g. with gdb)

Running helloWorld

```
user@host:~/Private/biostat615/hw0$ ./helloWorld

Hello, World

Hyun Min Kang Biostatistics 615/815 - Lecture 1 September 4th, 2012 38 /-
```

Implementing INSERTIONSORT Algorithm

insertionSort.cpp - printArray() function

```
// print each element of array to the standard output
void printArray(std::vector<int>& A) { // call-by-reference : will explain later
for(int i=0; i < A.size(); ++i) {
   std::cout << " " << A[i];
  }
  std::cout << std::endl;
}</pre>
```

Hyun Min Kang

Biostatistics 615/815 - Lecture

eptember 4th, 2012

40 /

Implementing InsertionSort Algorithm

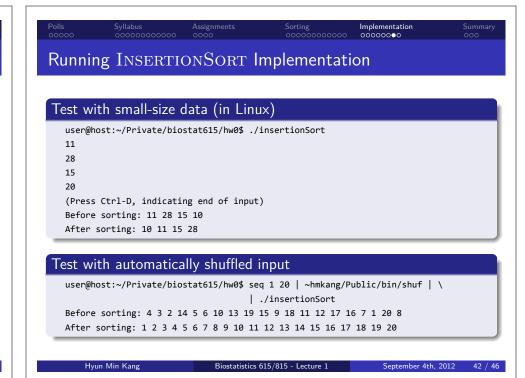
```
insertionSort.cpp - insertionSort() function
// perform insertion sort on A
void insertionSort(std::vector<int>& A) { // call-by-reference
  for(int j=1; j < A.size(); ++j) { // 0-based index</pre>
   int key = A[j]; // key element to relocate
   int i = j-1;
                 // index to be relocated
    while (i \ge 0) && (A[i] > key) ) { // find position to relocate
      A[i+1] = A[i]; // shift elements
                    // update index to be relocated
                   // relocate the key element
    A[i+1] = key;
```

How fast is INSERTIONSORT?

```
user@host::~/Private/biostat615/hw0$ time sh -c 'seq 1 100000 | \
                   ~hmkang/Public/bin/shuf | ./insertionSort > /dev/null'
0:03.33 elapsed, 3.334 u, 0.012 s, cpu 100.3%, 0 swaps, 0 rds, 0 wrts, \
pgs: 0 avg., 0 max.
user@host::~/Private/biostat615/hw0$ time sh -c 'seq 1 100000 | \
                   ~hmkang/Public/bin/shuf | sort -n > /dev/null'
0:00.16 elapsed, 0.157 u, 0.010 s, cpu 100.0%, 0 swaps, 0 rds, 0 wrts, \
pgs: 0 avg., 0 max.
```

default sort application is orders of magnitude faster than insertionSort

Biostatistics 615/815 - Lecture 1



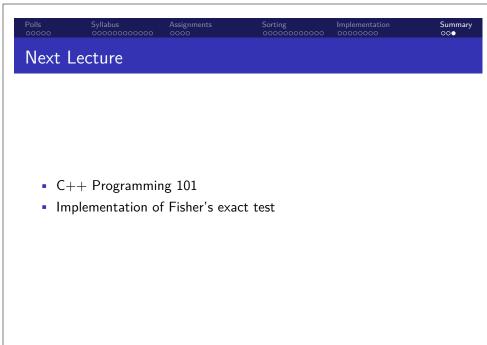


- Algorithms are sequences of computational steps transforming inputs into outputs
- Insertion Sort
 - √ An intuitive sorting algorithm
 - √ Loop invariant property
 - \checkmark $\Theta(n^2)$ time complexity
 - ✓ Slower than default sort application in Linux.
- A recursive algorithm for the Tower of Hanoi problem
 - √ Recursion makes the algorithm simple
 - √ Exponential time complexity
- C++ Implementation of the above algorithms.



- Implement the following two programs and send the output screenshots to the instructor (hmkang at umich dot edu) by E-mail
 - 1 HelloWorld.cpp
 - 2 TowerOfHanoi.cpp
- Briefly describe your operating system and C++ development environment with your submission
- This homework will not be graded, but mandatory to submit for everyone who wants to take the class for credit
- No due date, but homework 0 must be submitted prior to submitting any other homework.

Hyun Min Kang Biostatistics 615/815 - Lecture 1 September 4th, 2012 45 / 46 Hyun Min Kang



Biostatistics 615/815 - Lecture 1