**ArsDigitaUniversity**
**Month5:Algorithms    -ProfessorShaiSimonson**

**Lectures**

Thestudyofalgorithmsconcentratesonthehighleveldesignofdata structuresandmethodsforusingthemtosolveproblems.Thesubjectis highlymathematical,butthemat hematicscanbecompartmentalized,allowing astudenttoconcentrateon *what* ratherthan *why.*Theassumedprerequisiteis thatastudentcantakeadescriptionofanalgorithmandrelevantdata structures,anduseaprogrammingtooltoimplementthealgor ithm.Formost computerscientists,thisisexactlyhowtheymightinteractwithalgorithmsin theirfuturecareers.Keepingthisinmind,wheneverwe *write*thecodeforan algorithmwewilluseapseudoproceduralC -likelanguage,andthestudentis expectedtobeabletoimplementthedetailsusingOOPideasinJavaorC++, orfunctionalstylelikeScheme.

Acompleteunderstandingofalgorithmsismorethanjustlearninga fewparticularmethodsforafewparticularproblems.Thecoursefocusesnot justondetailsofparticularalgorithmsbutonstylesandpatternsthatcanbe usedinnewsituations.Thesecondfocusofthecourseisteachingthetoolsthat helpyoudistinguishbetweenproblemsthatareefficientlysolvableandones thatarenot.

Let'sgettosomeactualexamplesofthislatterpointbylistingpairs ofproblemsthatalthoughsuperficiallysimilar,haveoneproblemthatis efficientlysolvableandonethatisNP -complete.Fornow,NP -Complete meansthattheproblemasfarasanyreas onablepersonisconcernedhasno efficientsolution.Wewilllatergivearealdefinitionofthetermandmakethe intuitionmoreclear.

|   |   |   |
|---|---|---|
| 1.EulerCircuit | vs. | HamiltonianCircuit |
| 2.ShortestPath | vs. | LongestPath |
| 3.MinimumSpanningTree | vs. | DominatingSet |
| 4.EdgeCover | vs. | VertexCover |
| 5.MinCut | vs. | MaxCut |
| 6.2 -DimMatching | vs. | 3-DimMatching |
| 7.2 -Colorability | vs. | Colorability |
| 8.2 -Satisfiability | vs. | Satisfiability |

FormaldefinitionsoftheseproblemscanbefoundinGareyand Johnson'scomprehensivelistofNP -Completeproblemsfoundin:

ComputersandIntractibility:
AguidetothetheoryofNP -completeness
MichaelR.GareyandDavidS.Johnson
W.H.Freeman,1979.

Herearesomeinformaldescriptions:

**1.Eule rCircuit               vs.        HamiltonianCircuit**

EulerCircuitasks,givenanundirectedgraph,whetheryoucantrace theedgesstartingandendingatthesameplace,andtracingeachedgeexactly once.(Agivenvertexcanbetracedthroughmultipletimes).

Hamiltoniancircuitasks,givenanundirectedgraph,whetheryoucan tracethroughtheverticesofthegraph,startingandendingatthesameplace andtracingthrougheachvertexexactlyonce.

**2.ShortestPath               vs.        LongestPath**

ShortestPathasks,givenaw    eightedundirectedgraphandtwo verticesinthegraph,whatistheshortestpathbetweenthetwovertices.(The assumptionisthatthererearenonegativeweightcyclesinthegraph).

LongestPathisthenaturalanalogueofshortestpath.(The assumptionisthattheremaybepositiveweightcycles).

### 3.MinimumSpanningTree             vs.        DominatingSet

MinimumSpanningTree,givenaweightedundirectedgraph,asks fortheminimumweighttreethatcontainsalltheverticesinthegraph.Thetree mustbeasub  graphofthegivengraph.Sayyourtowngetsflooded,this problemasksfortheminimummilesofroadsthatneedtoberepavedto reconnectallthehousesintown.

DominatingSet,givenanundirectedgraph,asksfortheminimum sizesetofvertices,such  thateveryvertexiseitherinthissetorelseis connectedtoitdirectlybyanedge.Thisislikeyourtowngetsfloodedandyou wanttoputrescuestationsatintersectionssothateveryroadhasatleastone rescuestationateachend.

### 4.EdgeCo   ver                        vs.        VertexCover

EdgeCover,givenanundirectedgraph,asksforthesmallestsetof edgessuchthateveryvertexinthegraphisincidenttoatleastoneoftheedges.

VertexCover,givenanundirectedgraph,asksforthesmallestsetof verticessuchthateveryedgeinthegraphisincidenttoatleastoneofthe vertices.

### 5.MinCut                             vs.        MaxCut

MinCutasks,givenaweightedundirectedgraph,whatisthe minimumweightsetofedgeswhoseremovalseparatesthegraphintotwoor morediscon nectedcomponents.

MaxCutisthenaturalanalogueofMinCut.

### 6.2 -DimMatching                     vs.        3-DimMatching

2-DimMatching,isalsocalledthemarriageproblem.Givena certainnumberofmenandanequalnumberofwoman,andalistofpairsof men/womanwho  arewillingtobemarried,isthereawaytoarrangemarriages sothateveryonegetspairedup,andallpairsareinthepreferredlist.

3-DimMatchingisthenatural    *3-gender*analogueto2  -Dim Matching,whereeachmarriagemusthaveoneofeachoftheth        reegenders.

### 7.2 -Colorability                    vs.        Colorability

Colorabilityisthefamousproblemthatasksfortheminimum numbercolorsneededtocoloragraph,sothatnotwoconnectedverticeshave thesamecolor.Noteforaplanargraph,the4         -colortheoremim pliesthatthe numberisnolargerthanfour.

2-Colorabilityaskssimplywhetheragivengraphcanbecoloredwith atmosttwocolors.Thisequivalenttodeterminingwhetherornotagraphisbi          - partite.

### 8.2  -Satisfiability            vs.        Satisfiability

Satisfiabilityasks,givenawffinconjunctivenormalform,istherea T/Fassignmenttothevariables,suchthatthevalueoftheformulaendsup beingtrue.

2-Satistherestrictedversionof Satisfiabilitywhereeveryconjunct hasexactlytwovariables.

Trytoguesswhichoneofeachpairisthehardoneandwhichoneis theeasyone.Ishouldpointoutthatthe        *easy*onedoesnotnecessarilyhavean easilydesignedalgorithm.

### Algorithmscanbecategorizedbystyleandbyapplication.

Commonlyuse dstylesaredivideandconquer(recursion),dynamic programming(bottom -upormemoization),andgreedystrategy(dothebest thinglocallyandhopeforthebest).Commonapplicationcategoriesinclude mathematics,geometry,graphs,stringmatching,sort    ingandsearching. Combinatorialalgorithmsarealargercategoryincludinganyalgorithmthat mustconsiderthebestresultamongalargesamplespaceofpossibilities. ManycombinatorialproblemsareNP    -Complete.Donotconfusedynamic programmingwit h *linearprogramming.*   Thelatterreferstoaparticular probleminlinearalgebraandnumericalanalysiswithimportantapplicationin industryandoperationsresearch.Itisnotastyleortechnique,nordoesithave anythingtodowithprogramsthatru     ninlineartime.Itisaproblemthatcanbe usedtosolvemanycombinatorialproblems.

### CorrectnessandAnalysisofAlgorithms

Alargeconcernofthiscourseisnotjustthedesignofalgorithmsand thetechniquesusedindesigningthen,butthepro       ofsthatthealgorithmswork andtheanalysisofthetimeandspacerequirementsofthealgorithms.

Behindeveryalgorithmthereisaproofofcorrectnessoftenbasedon manytheoremsandlemmas.Theproofsareoftenbymathematicalinduction. Whenyou aredesigningyourownalgorithms,youmustbeabletoconvince yourselfthattheywork.Whenyoupublishorsharethem,youcan'treallyon yourownconfidenceandinstincts,butmustprovethattheywork.These proofscanoftenbequitetediousandtec       hnical,butunderstanding *why*an algorithmworkgivesyoubettertoolsforcreatingnewalgorithmsthanmerely knowinghowanalgorithmworks.

Therearemanywaystoanalyzethetimeandspacerequirementsof analgorithm.Wedescribethetimerequirem        entsofanalgorithmasafunction oftheinputsize,ratherthanalistofmeasuredtimesforparticularinputsona particularcomputer.Thelattermethodisusefulforengineeringissuesbutnot usefulforgeneralcomparisons.IntheearlydaysofCS            beforetheorywaswell developed,themethodofmeasuringactualtimesofprogramrunsgavenofair waytocomparealgorithmsduetotherandomengineeringdifferencesbetween computers,implementationsandlanguages.However,itshouldbeemphasized thattheoryisnotalwaysafairassessmenteither,andideallyonecalculatesthe timecomplexitytheoreticallywhileengineersfinetunetheconstantfactors,for practicalconcerns.Therearemanyexamplesofthis.Fibonacciheapsarethe fastestdatastr uctruesforcertainalgorithmsbutinpracticerequiretoomuch overheadtomakethemworthwhile.TheSimplexalgorithmofDantzigis worstcaseexponentialtimebutinpracticerunswellonreallifeproblems.

Thereisnoperfecttheoryformodelingall     aspectsofinputdistributionsand
timemeasurement.

Usually,thetimerequirementsarethemainconcern,sothespace
requirementsbecomesasecondaryissue.Onewaytoanalyzethetime
complexityofanalgorithmisworstcaseanalysis.Hereweimagin      ethatwe
nevergetlucky,andcalculatehowlongtheprogramwilltakeintheworstcase.
Theaveragecaseanalysismaybemoreuseful,whentheworstcasedoesnot
showupveryoften,andthismethodaveragesthetimecomplexityofthe
algorithmoveral lpossibleinputs,weightedbytheinputdistribution.Average
caseanalysisisnotascommonasworstcaseanalysis,andaveragecase
complexityisoftendifficulttocomputeduetotheneedforcarefulprobabilistic
analysis.

Athirdmethodofmeasur ingtimecomplexityiscalled *amortized*
analysis.Amortizedanalysisismotivatedbythefollowingscenario.Let'ssay
thatthereisanalgorithmwhoseaveragecasecomplexityislineartime,
howeverinpracticethisalgorithmisruninconjunctionwith        someother
algorithms,whoseworstcasecomplexitiesareconstanttime.Itturnsoutthat
whenyouusethesealgorithms,theslowoneisusedmuchlessfrequentlythan
thefastones.Somuchlessthatwecandistributethelinearcostoftheslow
oneove rthefastcostssothatthetotal      *amortized* timeovertheuseofallthe
algorithmsisconstanttimeperrun.Thiskindofanalysiscomesupwithwith
thefancierdatastructureslikeFibonacciheaps,whichsupportacollectionof
algorithms.Theamorti zedanalysisistheonlywayinwhichthefancierdata
structurecanbeprovedbetterthanthestandardbinaryheapdatastructure.

### LowerBoundsandNP -Completeness

Mostofthetimewewilldoworstcaseanalysisinthiscourse.This
givesusanupperb oundonthetimerequirementsofanalgorithm.Itisalso
useful,however,togetlowerboundsonourtimerequirementsforsolvinga
problem.Theselowerboundargumentsaremuchharder,becauseinsteadof
simplyanalyzingaparticularmethod,theyrequ irethatweanalyzeallpossible
methodstosolveaproblem.Onlywhenwecanquantifyoverallpossible
solutionscanweclaimthataparticularproblem.Afamouslowerboundresult
isthatsortingingeneralrequiresatleastO(nlogn)time.

Mosto fthetime,itistoohardtoproveanylowerboundtime
complexityonaproblem.Inthatcase,aweakerwaytoshowalowerboundon
solvingaproblem,istoprovethattheproblemisNP      -complete.Intuitively,
thismeansthatalthoughwecannotproveth      attheproblemrequiresmorethan
polynomialtime,wecanprovethatfindingapolynomialtimeforourproblem
wouldimplyapolynomialtimealgorithmforhundredsofotherproblems
(thoseinNP)forwhichnopolynomialtimealgorithmsarecurrentlyknown.
ThatisanNP -CompleteisthehardestofacollectionofproblemscalledNP.
Moredetailsonthislater.Itsufficestosaythatmanyhardproblemsthatdefy
apolynomialapproachareinthiscollectioncalledNP.

Determiningthefrontierbetweenwhic hversionsofaproblemare
polynomialtimesolvableandwhichareNP      -complete,isaskillhonedby
experience.Youwillgetplentyofpracticewiththissortofthinginthis
course.

### CopingwithNP -Completeness

Forthemostpart,provingthataproblem      isNP -completeisa
negativeresult.Allitreallygivesyouistheknowledgethatyoushouldnot
wasteyourtimetryingtocomeupwithafastalgorithm.However,thereare
manywaysofcopingwithanNP      -Completeproblem.

1. Lookforspecialcasesthat    maybepolynomialtimesolvable. Determinethefrontierfortheproblem.
2. Trytodesignan    *approximation* algorithm. Thisisanalgorithm thatwillnotgivetherightanswer,butwillgetananswerwithin acertainpercentageofthecorrectanswer.
3. Trya *probabilistic*algorithm.Thisisanalgorithmthatwillget therightanswersomepercentofthetime.
4. Experimentbyengineeringtheexponentialtimealgorithmsto seehowfaryoucangetwithcurrentspeedsandtechnologies.
5. Useadifferentcomputational    model(DNAforexample).

### MathematicalPreliminaries

Discretemathematicsisusedalotinstudyingalgorithms.Proofsby inductionaboundinprovingthecorrectnessofanalgorithm.Recurrence equationsandsumsfortheanalysisofalgorithmsarecruc        ial.Countingand probabilitycomeupeverywhere.Graphsandtreesarestandarddatastructures. LogicandBooleanalgebracomeupinmanyexamples.Finally,Big        -O notationfordescribingasymptoticcomplexityisastandardtool.

### MaxandMinAlgorit  hms –AWarmUp

Therearesimpleiterativeandrecursivealgorithmsforcalculating max(ormin)ofalistofnnumbers.Onemethodrecursivelyfindsthemaxof n-1ofthenumbers,andthencomparesthistothelastnumber.Therecurrence equationis  $T(n)=T(n$  -1)+1andT(1)=0.ThesolutiontothisisT(n)=n        -1, whichagreeswiththesimpleiterativemethodofinitializingaLargestSoFarto thefirstnumber,andthencomparingtherestofthenumbersonebyoneto LargestSoFar,andswappingwhen    necessary.

Anotherrecursivemethodtofindthemaxofalistofnnumbersisto splitthelistintotwoequalparts,recursivelyfindthemaxofeach,andcompare thetwovalues.TherecurrenceequationforthisisT(n)=2T(n/2)+1,T(1)=0, whichalso hasasolutionofn    -1.

Nowwhatifwewantedtocalculateboththemaxandmin.Wecould justdoanyoftheabovealgorithmstwice,giving2n        -2steps.Butcanwedo better?Wewillnotdobetterasymptotically.ThatiswewillstillneedO(n), butc anwemaketheconstantfactorsbetter?Normally,wedonotcareso muchaboutthis,butsometimesconstantfactorsareanimportantpractical matter,andhereitmotivatestheideaofdoingbetterthanwhatyougetwith bruteforce.

Wecancalculatethe  maxandminsimultaneouslybyusingthe followingrecursiveidea.Noteitispossibletodescribethisiteratively,butthe recursiveversionemphasizesthetechniquemoreelegantly.Theideaisthatwe canprocesstwonumbersatatime.Werecursively        computethemaxandmin ofn -2ofthenumbersinthelist.Thenwecomparethelargerofthetwo remainingnumberstothemaxandthelowerofthetworemainingnumbersto themin.ThismethodgivesarecurrenceequationofT(n)=T(n        -2)+3,T(2)= 1.Thesolutiontothisis(3n        -1)/2.Wecanalsosplitthelistintotwo,and recursivelycomputetheminandmaxofeachlist.Thenwecomparethemax ofonetothemaxoftheother,andtheminofonetotheminoftheother.This givesT(n)=2T(n/2)    +2,T(2)=1.Howdoesthismethodcompare?Knowing thesolutionstorecurrenceequationshelpsyouchoosewhichalgorithm variationtotry!

### TheMaxandSecondBiggest

Wecouldtrytodothesametrickforfindingthemaxandsecond biggestnumbersim ultaneously.Howexactly?Butthereisanadditionalthing wecanleverageinthiscase.Wecanruna     *tournament*tofindthemax,that takesn -1steps.Thistournamentisexactlythesameasthen/2recursionofthe previousexample.Wekeeptrackof     alltheelementsthatlostamatchwiththe eventualwinner.Therewillbelgnofthesemoreorless.Itisthelargestof thesewhichisthesecondbiggest.Thismethodtakesn     -1+lgn   –1 comparisons.Howdoesthiscomparewith(3n     -1)/2?Thedrawba ckhereis thatwemustdosomeextraworkinkeepingtrackofwhoplayedtheeventual winner.Nomatterhowyoudothis,itdoesaffecttheconstantfactorsofthe overallalgorithm.Therefore,comparingthemethodsatthislevelrequires experimentsan dmeasurementsmorethanjustanalysis.

### Sorting

Thebreadandbutterofalgorithmsaresortingalgorithms.Theydeal withthesimplestandmostcommonkindofproblem,andstillexhibitawide varietyoftechniques,datastructures,andmethodsofanal     ysisthatareusefulin otheralgorithms.Therearemanywaystocomparesortingalgorithms,butthe mainwayissimplybytimecomplexity.Otheraspectsofsortingrelatetothe practicalissuesofconstantfactors,recursiveoverhead,spacerequirement     s, performanceonsmalllists,suitabilityforparallelcomputationandwhetherthe sortpreservestheorderofequalsizedelements(stable).Therearehundredsof sortingalgorithms,ofwhichwestudyarepresentativesampling.

### O(n^2)TimeSortingAlgo   rithms

BubbleSortandInsertionSortaretwoO(n^2)sortingalgorithms. Despitetheslowtimecomplexity,theyhavetheirplaceinpractice.Theformer isexcellentforsortingthingsthatarealmostalreadysorted.Thelatteris excellentforsmall  lists.Boththesespecialtraitscanbeexplainedintuitively andcanbedemonstratedbyexperiments.Thisillustratesthatthetheorydoes notalwaysaddresstherealitysoexactly.Agoodcomputerscientistwill combineboththetheoryandengineeringto     explorethetruth.

### BubbleSort

```
SwitchMade=true;
for(i=1;(i<n)andSwitchMade;i++)
        SwitchMade=false;
        for(j=1;j<=n      -i;j++)
                if(a[j]>a[j+1])then{swap(a[j],a[j+1]);
SwitchMade=true;}
```

BubbleSortworksbydoingn    -1Bubbl eprocedures.EachBubble procedure(theinnerloop)comparesadjacentelementsanddecideswhetherto swapthem,proceedingdownwardthroughthearrayuptotheslotlookedatlast inthelastBubbleprocedure.

TheSwitchMadeflagisusedtoremember     whetheranyswapswere madeinthelastbubbleloop.Ifnot,thenthelistissortedandwearedone.An exampleisshownbelowwherethelistofnumbersisshownaftereachiteration oftheouterloop.

| 10 | 8 | 5 | 1 | 16 | 13 |
|----|---|---|---|----|----|
| 8 | 5 | 1 | 10 | 13 | 16 |
| 5 | 1 | 8 | 10 | 13 | 16 |
| 1 | 5 | 8 | 10 | 13 | 16 |

Notethatwiththeithiterationoftheouterloop,theinnerloopbubble procedurepushestheithlargestvaluetothebottomofthearraywhiletheother

lightervalues *bubble* upinalesspredictableway.Theworstcaseisthatthe
algorithmdoesn -1+n -2+n -3+…+1comparisons.Thesearethetriangle
numbersfromdiscretemathwhichareO(n^2).Youcanseethatifthearrayis
sortedtobeginwith,thenthealgorithmtakesO(n)time.Ifthearraygets
sortedearlierthanexpected,th  enthealgorithmstopsimmediately.

### InsertionSort

```
forj=2ton            {
          next=a[j];
          for(k=j  -1;((k>0)&&(a[k]>next));     --k)a[k+1]=a[k];
          a[k+1]=next;
}
```


Thisalgorithmsortsbykeepingasortedareafrom1toj,and
expandingitbyonewith   eachcompleteexecutionoftheinnerloop.Theinner
loopinsertsthenextelementintotheappropriateslotinthesortedareafrom1
toj -1.Itdoesthisbylookingattheareainreverseordershiftingelementsto
therightuntilitreachesanelement    whichitisnotlessthan.Itinsertsthe
elementinthatspot.Notethatlookingintheoppositeorderwouldtakeextra
time.ThisshouldremindyouofthecollapsingloopsinyourSameGame
programfromlastmonth.Anexampleisshownbelow:

| 10 | 8  | 5  | 1  | 16 | 13 |
|----|----|----|----|----|----|
| 8  | 10 | 5  | 1  | 16 | 13 |
| 5  | 8  | 10 | 1  | 16 | 13 |
| 1  | 5  | 8  | 10 | 16 | 13 |
| 1  | 5  | 8  | 10 | 13 | 16 |

TheworstcasetimecomplexityofthisalgorithmisalsoO(n^2)when
thelistwasoriginallyinreverseorder.Thesametrianglenumbersumas
BubbleSortshowsup,howeverunlikeBubbl      eSort,itdoesnotstopearlyifthe
listwasalmostsortedtobeginwith.ItdoeshoweverstilluseonlyO(n)total
operationsonasortedlist.Italsohasparticularlyfewoperationsintheinner
loop,makingtheconstantfactoroverheadverylow.I       tactuallyrunsfasterthan
manyO(nlogn)algorithmsforsmallvaluesofn,sincetheconstantfactors
and/orrecursiveoverheadfortheO(nlogn)algorithmsarerelativelyhigh.Of
courseforlargen,theO(n^2)algorithmsarenotpractical.

Thebook doesacarefulanalysisoftheaveragecaseforInsertion
Sortbysummingupallthetimerequirementsoftheequallylikelypossibilities,
dividingbythenumberofpossibilities,andshowingthatittooisO(n^2).
(pages8 -9).Itisoftenthecasetha      tworstcaseisthesameasaveragecase.
Thatiswhywerarelycalculateaveragecasetimecomplexity.

AnotableexceptionisQuicksortwhichisworstcaseO(n^2)and
averagecaseO(nlogn).WecalculateaveragecaseforQuicksortbecauseit
explains whythethingworkssowelldespiteitsbadworst       -casecomplexity.
Sometimespracticemotivatestheoryandsometimestheorymotivatespractice.

### Heapsort,MergesortandQuicksort  –O(nlogn)Time Algorithms

### Mergesort

MergesortisanO(nlogn)sortth       atworksrecursivelybysplittingthe
listintotwohalves,recursivelysortingeachhalf,andthen       *merging* theresults.
Themergealgorithmtakestwosortedlistsandcreatesamergedsortedlistby
examiningeachlistandcomparingpairsofvalues,putt       ingthesmallerofeach
pairofvaluesinanewarray.Therearetwopointers,oneforeacharrayof
numbers,andtheystartatthebeginningofthearray.Thepointertothelist

fromwhichthesmallervaluewastakenisincremented.Atsomepointone          of
thepointershitstheendofthearray,andthentheremainderoftheotherarray
issimplycopiedtothenewarray.

Thetimecomplexityformergingisthesumofthelengthsofthetwo
arraysbeingmerged,becauseaftereachcomparisonapointertoon          earrayis
moveddownandthealgorithmterminateswhenthepointersarebothatthe
endsoftheirarrays.Thisgivesarecurrenceequationof$T(n)=2T(n/2)+O(n)$,
$T(2)=1$,whosesolutionis$O(n\log n)$.

OneimportantfeatureofMergesortisthatisno          t *in-place*.Thatis,it
usesextraspaceproportionaltothesizeofthelistbeingsorted.Mostsortsare
*in-place,* includinginsertionsort,bubblesort,heapsortandquicksort.
Mergesorthasapositivefeatureaswellinthatthewholearraydoesnot          needto
beinRAMatthesametime.Itiseasytomergefilesoffdisksortapesin
chunks,soforthiskindofapplication,mergesortisappropriate.Youcanfind
aCversionofmergesortinassignment1ofHowComputersWork(month3).

### Heapsort

Heapsortisthefirstsortwediscusswhoseefficiencydepends
stronglyonanabstractdatatypecalleda          *heap.* Aheapisabinarytreethatis
ascompleteaspossible.Thatis,wefillitinonelevelatatimefromrightto
lefttoneachlevel.Ithasthe          propertythatthedatavalueateachnodeisless
thanorequaltothedatavalueatitsparent.(Notethatthereisanotherabstract
datatypecalledabinarysearchtreethatisnotthesameasaheap.Also,there
isanother *heap*usedinthecontexto          fdynamicallocationofstorageand
garbagecollectionforprogramminglanguagessuchasJavaorScheme.This
otherheaphasnothingtodowithourheap.Theheapfromdynamicmemory
allocationhasmoreofausualEnglishmeaningasinaheapoffreeof          memory,
andisactuallymorelikealinkedlist.)

Aheapsupportsanumberofusefuloperationsonacollectionofdata
valuesincludingGetMax(),Insert(x),andDeleteMax().Theeasiestwayto
implementaheapiswithasimplearray,whereA[1]isther          oot,andthe
successiveelementsfilleachlevelfromlefttoright.Thismakesthechildren
ofA[i]turnupatlocationsA[2i]andA[2i+1].Hencemovingfromaparentto
achildorviceversa,isasimplemultiplicationorintegerdivision.Heapsalso
allowchanginganydatavaluewhilemaintainingtheheapproperty,Modify(i,
x),whereiistheindexofthearrayandxisthenewvalue.Heapsareauseful
waytoimplementpriorityqueuesthatisacommonlyusedabstractdatatype
(ADT)likesstacksan  dqueues.

ToGetMax(),weneedonlypullthedatavaluefromtherootofthe
tree.TheotheroperationsInsert(x),DeleteMax()andModify(i,x)requiremore
carefulwork,becausethetreeitselfneedstobemodifiedtomaintaintheheap
property.Themod  ificationandmaintenanceofthetreeisdonebytwo
algorithmscalled *Heapify* (page143)  and *Heap-Insert* (page150).These
correspondtotheneedtopushavalueupthroughtheheap(Heap          -Insert)or
downthroughtheheap(Heapify).Ifavalueissmaller          orequaltoitsparentbut
smallerthanatleastoneofitschildren,wepushthevaluedownwards.Ifa
valueislargerorequalthanbothitschildren,butlargerthanitsparent,thenwe
pushitupwards.SometextscallthesetwomethodssimplyÂ          *PushUp* and
*PushDown.*Thedetailsofthesetwomethodsusingexampleswillbeshownin
class.ThetimecomplexityfortheseoperationsisO(h)wherehistheheightof
thetree,andinaheaphisO(lgn)becauseitissoclosetoperfectlybalanced.
Analter natewaytocalculatethecomplexityistherecurrence$T(n)=T(2n/3)+$
$O(1),T(1)=0$,whosesolutionis$O(\log n)$.Therecurrencecomesfromthe
factthattheworstcasesplittingofaheapis2/3and1/3(page144)onthetwo
children.

Heapsortworks  intwophases.Thefirstphaseistobuildaheapout ofanunstructuredarray.Thenextstepis:

```
forindex=lastto1     {
            swap(A[0],A[index]);
            Heapify( 0);
}
```

Wewilldiscussthebuildheapphaseinclass,andthereisaproblem onitinyourPset.     Itisalsodiscussedatlengthinthetext.Thenextphase worksassumingthearrayisaheap.Itcomputesthelargestvalueinthearray, andthenextetc.,byrepeatedlyremovingthetopoftheheapandswappingit withthenextavailableslotworking      backwardsfromtheendofthearray. Everyiterationneedstorestoretheheapproperty sinceapotentiallysmallvalue hasbeenplacedatthetopoftheheap.Aftern      -1iterations,theheapissorted. SinceeachPushUpandPushDowntakesO(lgn)andwed       oO(n)ofthese,that givesO(nlogn)totaltime.

Heapsorthassomenicegeneralizationsandapplicationsasyouwill seeinyourPset,anditshowstheuseofheaps,butitisnotthefastestpractical sortingalgorithms.

### Quicksort

Quicksortandits   variationsarethemostcommonlyusedsorting algorithms.Quicksortisarecursivealgorithmthatfirst     *partitions* thearrayin placeintotwopartswherealltheelementsofonepartarelessthanorequalto alltheelementsofthesecondpart.Aftert      hisstep,thetwopartsarerecursively sorted.

TheonlypartofQuicksortthatrequiresanydiscussionatallishow todothepartition.Onewayistotakethefirstelementa[0]andsplitthelist intopartsbasedonwhichelementsaresmalleror l      argerthanA[0].Therearea numberofwaystodothis,butitisimportanttotrytodoitwithoutintroducing O(n)extraspace,andinsteadaccomplishthepartitioninplace.Theissueis thatdependingonA[0],thesizeofthetwopartsmaybesimilar       orextremely unbalanced,intheworstcasebeing1andn      -1.TheworstcaseofQuicksort thereforegivesarecurrenceequationofT(n)=T(n      -1)+O(n),T(1)=0,whose solutionisO(n^2).

Thepartitionmethodwewillreviewinclasskeepspointerstotw       o endsofthearraymovingthemclosertoeachotherswappingelementsthatare inthewrongplaces.Itisdescribedonpages154      -155.Analternativepartition algorithmisdescribedinproblem8    -2onpage168.

ThequestioniswhyisQuicksortcalleda      nO(nlogn)algorithmeven thoughitisclearlyworstcaseO(n^2)?Ithappenstorunasfastorfasterthan O(nlogn)algorithmssowebetterfigureoutwherethethroryismessingup.It turnsoutthatifwecalculatetheaveragecasetimecomplexity       ofQuicksort,we getanO(nlogn)result.Thisisveryinterestinginthatagreeswithwhatwesee inpractice.Moreoveritrequiresthesolutionofacomplicatedrecurrence equation,T(n)=(2/n)(T(1)+T(2)+…+T(n      -1))+O(n),T(1)=0,whose solutionisobtainedbyguessingO(nlogn)andverifyingbymathematical induction,atechniquewithwhichyoumaybefamiliar.Thesolutionalso requirestheclosedformsummationofklogkfork=1ton,anothertechnique fromdiscretemathematicsthat    youhaveseenbefore.

### BucketSortandRadixSort    –LinearTimeSortsforSpecial Cases

### CountingSort

Theideabehindbucketsortisbasedonasimplerideaourtextcalls countingsort.Thismethodisequivalenttothefollowingwaytosortquizze        s whosegradescanbeanythingbetween0and10.Setup11placesonyourdesk andmarkthem0through10.Thengothrougheachquizoneatatimeand placeitinthepileaccordingtoitsgrade.Whenyouarefinished,gatherthe quizzestogethercollec  tingthepilesinorderfrom0to10.Thismethod generalizestosortinganarrayofintegerswherethenumberofdatavaluesis limitedtoarangebetween0andm.Thetimecomplexityofcountingsortis O(n+m)wherenisthenumberofdatavalues,and       misthelargestpossibledata value.NoteitisO(n),becausewecanplacethequizgradesintheir appropriateslotswithaloopofinstructionslike:fori=0ton       -1{B[A[i]]++}. B[j]holdsthenumberofquizzeswithgradej.Thealgorithmisanaddi        tional O(m)becausewemustinitializeandcollectthepilesattheend.Thisisdone by:forj=0tom      -1{fork=1toB[j]{print(j);}}.Thislooptakestimeequal tothemaximumofnandm.NoteitdoesnottakeO(nm).

### BucketSort

Countingsor  tisO(n)whenevermisO(n),butitcanbeveryslowif misO(2^n).BucketSortisawaytoextendCountingSortwhenthevaluesare notlimitedinsize.Instead,weartificiallydividethennumbersintondifferent groups.Soforexampleifwehave       100five -digitpositivenumberswherethe maximumis99999,thenwedividetherangeinto100differentintervals.We candothisbysimplyusingthefirsttwodigitsofthenumberasitsinterval value,sothat45678wouldbeplacedininterval45.(In          general,tofindthe correctintervalforadataentry,wewouldhavetodividethedatabym/n, wheremisthemaximumvalueandnisthenumberofintervals.)Thesort worksbylookingateachvalueandplacingitinitsappropriateinterval.Each intervalhasalinkedlistthatholdsallthevaluesinthatintervalsincetheremay ofcoursebemorethanoneinanyinterval.Afterallthevaluesareplacedin someinterval,eachintervalissortedandthentheyarecollectedtogetherin orderofinterv  alsize.

TheimplicitassumptionthatmakesthissortO(n)time,isthatthe distributionofvaluesintointervalsisuniform,henceBucketSortletsustrade theassumptionofuniformdistributionforCountingSort'sassumptionofa limitednumberofva  lues.Noteiftherewasoneintervalthatcontainedallthe numbers,thesortwouldtimeatbestO(nlogn).Theanalysisassuming uniformdistributionrequiresalittleprobabilityanddiscretemath(seepage 182).

### RadixSort

RadixSortisanother   generalizationofCountingSort,wherewe assumenothingaboutthelistofnumbers.RadixSortisanideaoriginallyseen inthepunchcardmachinesinventedbyHermannHollerithtododothe1890USA census.ThesametrickwasusedinIBMpunchcardsoft        he1960'sand1970's. AneatfeatureofRadixSortisthatitmustusea        *stable*sortasasubroutine. Recallthatastablesortisonethatpreservestheorderofequalvalued elements.

RadixSortworksbyrepeatedlysortingthenumbersbylookingat        the digits,fromrighttoleft.Agoodexamplewouldbesortingstringsoflength4 inalphabeticalorder.Saywewanttosort:

SHAI
FRAN
SHAM
FANA
FRAM

Wemake26boxeslabeledA   -Z,andwesortthestringsusing countingsortontherightmostcharac   ter.Toimplementthis,weuseanarrayof linkedlistswhoseindicesaretheAthroughZ.Afterthefirststep,allstringsin theoriginalinputarraythatendinA,areinthelinkedlistheadedbyAetc. Thenwecopyallthestringsfromthelinked       listsinorderbackintotheoriginal inputarray,overwritingtheoriginalarray.Notethatthisisa       *stable* process. Thisgivesthelistbelow:

FANA
SHAI
SHAM
FRAM
FRAN

Werepeatthisstepsortingonthecolumnthatisthesecondtothe right,makingsure(veryimportant)topreservetheorderofequalcharacters. (Theonlytimethispreservationdoesnotmatterisontheveryfirstiteration.) Thisgives:

SHAI
SHAM
FRAM
FRAN
FANA

Youshouldnotethatafterthisstepthestringsaresortedcorrectl       yif welookjustatthelasttwocharacters.Aftereachsubsequentstep,thissorting willbecorrectforonemorecolumntotheleft.Youcanprovenaturallyby induction,thatitworksingeneral.

Herearethelasttwostages:

FANA
SHAI
SHAM
FRAM
FRAN

FANA
FRAM
FRAN
SHAI
SHAM

Thesamealgorithmworksonintegersbysplittingeachintegerup intoitsdigits.Itisfinetousethebinarydigits.Thetimecomplexityis$O(n+k)$ foreachstep,wherenisthenumberofelementsandkisthenumberof differentdigits(2forbinary).Therearedstepswheredisthenumberofdigits ineachnumbergivingatotalof$O(d(n+k))$.Sincekisconstantanddisworst case$O(logn)$thenradixsortworksin$O(nlogn)$worstcase.Itcanbelinear timewhendh   appenstobe$O(n)$.

Onethingtonotewithradixsortisthatifwesortfromthemost significantbittotheleastsignificantbit,thenwedonotactuallysortthearray. Forexample:

| 356 | 189 | 338 | 185 | 266 | 325 | turns |

into:

| 189 | 185 | 266 | 356 | 338 | 325 | which |

turnsinto:

| 325 | 338 | 356 | 266 | 185 | 189 | |

andwearegettingnowherefast.

Ifwetrytofixthisbysortingsubsetsofthelistineachsubsequent iteration,weendtakingtoomuchspaceandtoomuchtime.HowmuchextraI leavetoyoutothinkabout.

Thisprocessingfromleastsignificanttomostsignificantseems unintuitiveatfirst,butisactuallythekeytothewholealgorithm.Itallowsus toreusethesameoriginalCountingSortarrayforeachiteration,anditkeeps thetimecomplexitydown.

### LowerBoundsonSorting

Thereisawellknownargumentthatanyalgorithmusing comparisonsrequiresatleastO(nlogn)comparisonstosortalistofn numbers.Theargumentturnsanyalgorithmintoadecisiontree,whichmust have *atleast* n!leaves.Each leafrepresentsaparticularpermutationofthe inputlist,andsincetheinputlistisarbitrary,theremustbeatleastn!leaves. Fromdiscretemath,werecallthatabinarytreewithmleaveshasdepthatleast lgm,andherethatgiveslgn!whichi        sO(nlogn).Hencewedonotexpecta generalsortingalgorithmusingstandardmethodsofcomputationtoeverdo betterthanwealreadyknowhowtodo.Thisideaofdecisiontreescanbeused togetprimitivelowerboundsoncertainotherproblems,but        lowerboundsin generalareelusiveformostproblems.

### MedianandtheKthLargestinLinearTime

WecancertainlyfindthemedianofalistbysortinginO(nlogn) time.Thequestioniswhetherwecandobetter.Heapsortcanbegeneralized inanobvi ouswaytogetthekthlargestvalueinO(klogn).Aprobleminyour Psetdiscussesawaytogetthisdowntoo(klogk).Unfortunately,whenk=n/2 (themedian),boththesearestillO(nlogn).Isthereanywaytofindthe medianofalistofnnumber    inO(n)time?

Thereisarecursiveschemethatsolvesthekthlargestproblemin lineartime,althoughtheconstantfactorandoverheadcanbequitehigh.Lets seehowitworks.

Wearrangethelistintoa2     -dimensionalarrayof5byn/5.Wethen findthemedianofeachcolumnandpartitionit.Wenowhaven/5columns, eachoneofwhichhasthehighertwovaluesontopandthelowertwoonthe bottom.Wethenlookatthemiddlerow,andrecursivelycalculateitsmedian, m,andthenpartition(alaQu    icksort)therowsothattheleftsidehasnumbers smallersthanthemedianandtherightsidehasnumberslargerorequal.The partitioningmovescolumnsalongwiththeirmiddleelements,sothatatthis pointwehavetheupperleftquadrantsmallerthan        m,andthelowerright quadrantlargerorequaltom.Theremaingtwoquadrantsmustbechecked elementbyelementtocompletelypartitionthearrayintotwopartsonesmaller thanm,andonelarger.CallthesetwopartsS_1andS_2respectively.If        S_1 hasmorethankelementsthenwerecursivelycallouralgorithmonS_1. OtherwisewerecursivelycallitonS_2.Wewilldoadetailedexamplein class.

Therecurrenceequationforthisprocessisamess.Tofindthe medianofacolumnof5elemen    tscanbedonein6comparisons,sothisstep takes6n/5.TorecursivelyfindmedianofmiddlerowtakesT(n/5).To partitionthearrayandmovethecolumnsaroundtakestimeaboutn/5+n.To constructS_1andS_2,includingcheckingeachelementinth        eupperrightand lowerleftquadrants,takesntime.TorecursivelycallthealgorithmonS_1or S_2takesworstcaseT(3n/4),becauseatleast¼ofthearrayisineachset. ThisgivesT(n)=T(n/5)+T(3n/4)+17n/5,T(5)=6.

Solvingthisexplicitl yisdifficult,butwecanguessthatthesolution islinear,andproveitisbyinduction.Theconstantthatworksfortheproofis constructed,andmaybefairlylarge.Thequestionofwhetherthisalgorithmis practicalneedstoconsidertheactualda taandsizeoftheliststhatweare processing.

ThekeypointaboutthisrecurrenceequationisthatitresemblesT(n) =T(n/2)+O(n).IfitresembledT(n)=2T(n/2)+O(n)thenthecomplexity wouldbeO(nlogn)ratherthanO(n).Thereasonforthis        isthat3n/4+n/5<n, andthisexplainswhyweusefiverowsinthismethod.Fiveisthesmallest numberofrowsthatallowsalineartimerecurrence.Withthreerowswe wouldhaveT(n/3)+T(3n/4)andthatwouldgiveanO(nlogn)recurrence.

### Data Structures

Thereareanumberofbasicdatastructuresthatcomputerscientists useindesigningalgorithms.Theseincludestacks,queues,linkedlists,priority queues,heaps,treesofvarioustypes,graphs,andhashing.Forthemostpart, thissmall collectionanditsvariationsareenoughtohandlealmostanyproblem youapproach.Yourarelyneedtodesignanewdatastructurefromscratch,but youmaywellneedtodesignavariationofoneofthese,oratleastknowwhich oneisappropriateforwhi chtask.Exceptforheaps,Red    -Blackbinarysearch treesandgraphs,wewillleavethediscussionsofbasicdatastructuresand implementationsforrecitationsandreviews.

### BinarySearchTrees

Theflipsideofsortingissearching.Searchingisperhaps       evenmore fundamentalthansorting,inthatoneofthelargestspecialtiesincomputer science,databases,isconcernedprimarilywithwaysoforganizingand managingdatathatmaintainintegrity,consistencyandallowforgeneral searching.Thetheory ofdatabasesrequiresaseparatecourse,andthedata structureusuallyusedfortheunderlyingphysicallayerofadatabaseisaB       - tree,orvariationsthereof.WewillleavethestudyofB        -treesforthedatabase course,anddiscusssearchingatitrelat     edtosimplersmallerscaleapplications.

Oneofthefirstalgorithmsthatchildrendiscoveriswhatwecall binarysearch.Itiswhenachildtriestoguessasecretnumber,andisgiven highorlowanswers,inordertodeterminehisnextguess.Ifth        erangeof numberis1 -16,thechildwouldguess8;ifthatistoohigh,thenthenextguess wouldbe4,etc.untilthecorrectvalueisfound.Thisprocessisnaturally recursiveandcutsthelistinhalfwitheachsubsequentguess.Therecurrence equationisT(n)=T(n/2)+1T(1)=0,andthesolutionisT(n)=O(logn).

Binarysearchcanbedoneonelementsstoredinanarray,but althoughthisallowssearchesthatareO(logn)time,theinsertionsneedO(n). Ifweuselinkedliststheinsertio nsareO(1)butthesearchisO(n).Thedata canalsobestoredina *binarysearchtree* .Abinarysearchtreeisadata structurethatsupportssearching,insertionsanddeletions.Eachnodestoresa datavalueonwhichwewillsearch.Allnumbersint heleftsubtreeofanode aresmallerthanthenumberinthenode,andallnumbersintherightsubtreeare largerthanorequal.Wekeeppointersfromeachnodetoitschildren,andwe canalsoincludeapointertoitsparent.Tofindavaluewecompare        ittothe rootandmoveleftorrightdependingontheanswer.Whenwegetanequal result,westopandreturntheinformationassociatedwiththatvalue.

Insertionsaredonebyaddinganewleaftothetreeintheappropriate placeafterwehitanu llpointer.Deletionsareabittrickier.Exampleswillbe showninclass.

**TreeTraversals**

Beforewediscussdeletionsfromabinarysearchtree,itis worthwhiletoreviewtreetraversals.Therearethreenaturalrecursivewaysto traverseatree, andoneothernon -recursiveway.Therecursivewaysarecalled inorder,preorderandpostorder.Inordertraversalmeansthatwefirst recursivelytraversetheleftsubtree,thentheroot,thenrecursivelytraversethe rightsubtree.Preordertraversal( whichisalsodefinableond -arytrees)means thatwevisittherootandthenrecursivelytraversethesubtrees.Thisislike depthfirstsearch.Postordermeanstheopposite,thatis,firstrecursively traversethesubtreesthenvisittheroot.Therea        reapplicationsforeachof these,andgoodexamplescomefromcompilerdesignandparsing.However, theinordertraversalisveryusefulforbinarysearchtrees.        *Aninordertraversal ofabinarysearchtree,printsoutthevaluesinsortedorder.*

Thenaturalnon -recursivetraversalofatreeiscalled        *levelorder.* Itis associatedwithbreadthfirstsearch,andjustasdepthfirstsearchandthe recursivetraversalsusestacksasafundamentaldatastructure,sodoesbreadth firstsearchandlevelo    rdertraversalsuseaqueue.

Backtodeletionofnodesinabinarysearchtree…Todeleteanode inabinarysearchtree,weneedtousetheinordertraversal.Theideaisthatwe donotwanttolosethepropertyofourorderedstructure.Whensearc        hingor addingthisisnoproblem,butdeletionwillmanglethetree.Thetrickistotry todeletealeafifpossiblebecausethisdoesnotaffectthetreeordered structure.Whenthenodewemustdeleteisnotaleaf,butithasonlyonechild, thenwe canjustdeletethenodeby        *splicing* itaway.Whenthenodehastwo childrenthenweusethefollowingtrick.Wefinditssuccessorintheinorder traversal,wespliceoutthesuccessor(whichwecanprovemusthaveatmost onechild),andreplacethev    alueofournodewiththatofthedeletedsuccessor. Intuitivelythisworks,becausewearebasicallyreplacingthenodetobedeleted withthenexthighestvalueinthetree.Thefactthatthisnexthighestvalue musthaveatmostonechildandtherefor        ecanbesplicedoutisveryhelpful. Therearemanywaystofindtheinordersuccessorofanodebutasimpleoneis justtodoalineartraversaltimeinordertraversalandstorethevaluesinan array.Thetextgivesfasterandmoreefficientwaystha        tdonothavetotraverse thewholetree(page249).Usingasuccessorguaranteesthatthebinarytree retainsitsorderedstructure.

**What'stheProblemwithBinarySearchTrees?**

Theproblemwithbinarysearchtreesisthattheycangetthinand scrawny,andtosupportfastinsertions,searchesanddeletionstheymustbefat andbushy.Alltheoperationsweperformonabinarysearchtreetaketime proportionaltotheheightofthetree.Butthetreeintheworstcasecanturn intoalongthinstraig    htline,sothetimecomplexitybecomesO(n)insteadof O(logn).

Wecanjustkeepthestuffinanarray,buttheninsertionstakeO(n)                – becausewehavetoshiftoverelementstomakeroomforanewleaf,like insertionsort.Thesolutionistocome        upwithsomesortofdatastructurethat hasthedynamicstructureofpointerslikeatree,butwhichisguaranteednever togrowtoothinandscrawny.Historically,therehavebeenanumberof candidatesincludingheightbalancedtreeslikeAVLtreesan        d2 -3trees,and weightbalancedtrees.Theheightbalancedtreeskeeptheleftandrightheights fromeachnodebalanced,whiletheweightbalancedtreeskeepthenumberof nodesineachrightandleftsubtreebalanced.Theyaresimilarintheorybut heightbalancedwonfavoroverweight    -balancedtrees.2    -3treesevolvedinto B-treesusedfordiskstorage,andAVLtreesgotreplacedbyRed        -Blacktrees, becausetheRed -Blacktreswereslightlymoreefficient.Anadvanceddata structurecalledaSplayt    reeaccomplishesthesamethingasRed    -Blacktrees, butusesamortizedanalysistodistributethecostoverthesetofalloperations

onthedatastructure.BinomialheapsandFibonacciheapsarealsoadvanced datastructuresthatdoforsimplebinaryhea     ps,whatsplaytreesdoforRed     - Blacktrees.

### Red-BlackTrees

WediscussonlyRed   -Blacktrees,leavingthesimplerandmore advanceddatastructuresforyourownpersonalstudyorrecitationsandreview. TheresultthatmakesoperationsonaRed     -Black treeefficient,isthattheheight ofaRed -Blacktreewithnnodesisatmost2lg(n+1),hencetheyare relativelybushytrees.

TheproofofthisresultdependsontheexactdefinitionofaRed          - Blacktreeandaproofbyinductionyoucanfindonpag        e264ofyourtext.A Red-BlackisabinarysearchtreewhereeachnodeiscoloredRedorBlack, everyRednodehasonlyBlackchildren,everyleafnode(nil)isBlack,andall thepathsfromafixednodetoanyleafcontainthesamenumberofBlack nodes.

WeshowhowtodosearchesandinsertionsonaBlack        -Redtree.The textshouldbeconsultedfordetailsondeletionsthatismildlymorecomplex (asitisingeneralbinarysearchtrees).Thedetailsofimplementationand pointerdetailsislefttoy     ou,withthetextprovidingplentyofhelponpages 266,268,273and274.

TosearchaRed   -Blacktreeyoujustdothenormalbinarytreesearch. SincetheheightisatmostO(logn),weareokay.Thehardpartistodo insertionsanddeletionsinO(log   n)whilemaintainingtheRed    -Blackproperty.

### InsertionsintoaRed   -BlackTree

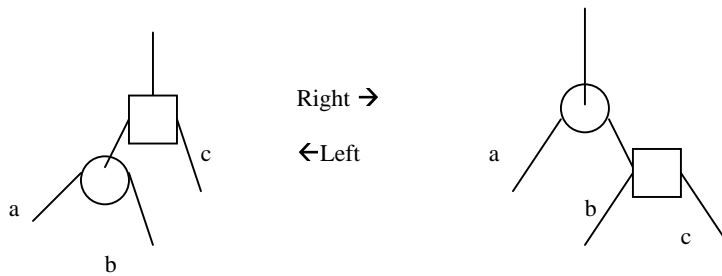Insertingavalueintoabinarysearchtreetakesplaceataleaf.What propertiesoftheRed   -Blacktreemightthisviolate.Ifwecolorthenewnode Red,andmakeitsn    ilchildrenBlack,thenthenumberofBlacknodesonany pathhasnotchanged,allnodesarestilleitherRedorBlack,allleavesare Black,andthechildrenofthenewRednodeareBlack.Theonlypropertyto worryaboutiswhethertheparentofthenew        leafisalsocoloredRed,which wouldviolatetheruleaboutRednodeshavingtohaveonlyBlackchildren.

Fixingthispropertyisnotsimple.Wewillneedtorecolornodes, pushinguptheRed   -Redclashuntilfinallygetridofit.Intheworstcasew          e needtorecolorvaluesallthewayupthetree.Inordertodothisrecoloring,we require1or2 *rotations*attheendwhichinvolveamutationofthesearchtreein additiontoarecoloring.

### Rotations

Arotationisawaytoreshapeabinarysearch         treethatmaintainsthe structureoftheinordertraversal.Itisanimportanttoolformanagingany height-balancedtrees.

Leftandrightrotationsareinversesofoneanotherandthebasic movementsareshownaspicturesbelow.Thedetailsoftheco       dethatactually movesthepointerstoaccomplishthis,canbefoundonpage266ofyourtext.

Right →

←Left

Youcancheckthattheinordertraversalofthesetwotreesisthe same.Theimprovementisthattheheightsofthesubtreeshavebecomemore balanced.Forexampleifthesubtreeat'a'wasabitlong,thenarightrotation willbalanceitupmore.Itiskindoflikepullingupdroopysocks,andmoving theslackovert hetop.

IwilldoadetaileddescriptionofinsertinganelementintoaRed                   - Blacktree,showingthecaseswheretheRed        -Redclashispushedupwards. WhentheUncleofthebottomRednodeisalsoRed,wecanjustpushitupand recolor.Theproblemoccu    rswhentheUncleofthebottomRednodeisBlack. Inthiscase1orpossiblytworotationsarenecessarytorestoretheBlack            -Red properties.Thegoodnewsisthatifanyrotationisnecessary,thenthatendsall futureRed -Redclashes,andwecanstop      pushingtheclashupthetree.See classexamplesinrealtimefordetails,orseeyourtextonpage271.

### GraphAlgorithms

Graphalgorithmsrepresentthemostdiversecollectionof applicationsandproblemsofanyalgorithmcategory.Graphsareused            to representgames,networks,processdependencies,scheduling,physicaland virtualconnectivity.ManygraphproblemsareNP        -complete.Thereafew basicgraphalgorithmsthataresolvableinpolynomialtimeincludingminimum spanningtree,shortestpath   ,maximumflow,andmaximummatching.There areanumberofgeneralalgorithmsongraphsthatareflexibleandcanbeused tosolveahostofotherproblems.Thesegeneraltechniquesaredepthfirst searchandbreadthfirstsearch.Theformerinparticu        larhasgreatflexibility anditsvariationsaremyriad.Itcanbeusedtofindcycles,connected components,stronglyconnectedcomponents,bi    -connectedcomponents, triangles,andtoplogicalorderings.DFScanalsobeusedasthebasisforbrute forcec ombinatorialalgorithmsthatareNP     -completeandcomeupinAIrelated applications.

Therearemanykindsofgraphs,andplentyoftheoremsaboutgraphs thatgiveusagoodfoundationonwhichtobuildouralgorithms.Mostofthese canbelookedupiny    ourdiscretemathtext,andreviewedifnecessaryin recitation.Wewillfocusfirstonhowwecanrepresentagraphinsidea computer;thatis,whatdoesagraphdatastructurelooklike?

Thethingthatmakesgraphalgorithmsalittlehardforbeginn       ersis thatweareusedtosolvinggraphalgorithmsbyusingoureyesandallthebuilt invisualprocessingthatgoeswiththat.Forexample,ifIshowyouapicture withasquareandatriangle,itiskindofobviouswhattheconnected components are,a ndwhetherornotthereisacycle.Thebeginnerdoesnot oftenknowjustwhatisnecessarytodescribeinagraphalgorithm.Thebest waytogetagoodsenseforgraphalgorithmsistolookatthegraphthewaya computerseesit,andthentrytodescri       beyourmethodoralgorithm.

### TheGraphDataStructure

Graphsarestoredasanarrayoflinkedlists.Eachnodexhasaslot inthearrayandeachhasalinkedlistthatholdsthenumbersofallthenodes thatadjacenttox.Forexample:

0:A →3 →4 →nil
1:B →2 →5 →7 →nil
2:C →1 →2 →7 →nil
3:D →0 →4 →nil
4:E →0 →3 →nil
5:F →1 →2 →7 →nil
6:G →2 →7 →nil
7:H →1 →2 →5 →nil

Quick!Whatdoesthisgraphlooklike?Whatareitsconnected components?Doesithaveacycl    e?Ifyoudrawthepicture,youwillbeableto answerthesequestionsimmediately.Butyoushoulddesignyouralgorithms withoutbeingabletoseethepicture,becausethatisthedatastructurethatyour programwilluse.Maybeonedaywewillhaveda           tastructureswhosemethods aretheanalogueofourvisualprocessing?Don'tholdyourbreath.

Therearemanywaystoaugmentthisdatastructure.Header informationliketheindegreeoroutdegreeofanodecanbeadded.Weightson eachedgecanbein   dicatedwithanotherfieldineachlinkedlistnode.Note thatinanundirectedgrapheachedgeappearsintwodistinctlinkedlistnodes. Ingeneral,anyalgorithmthatcanrunbyaconstantnumberoftraversalsofa graphdatastructurerunsintimeO(     n+e)wherenisthenumberofnodesinthe graphandeisthenumberofedges.

Ofcourseagraphcanalsobestoredinatwodimensionalarray, whichisusefulforcertainthings.Matrixmultiplicationandvariationsofinner product,helpcalculatet  henumberofwalksandpathsandrelatedstatsabouta graph,asyousawindiscretemath.However,formostalgorithmsthetwo dimensionalmethodsjustslowsdowntimecomplexityfromO(n+e)to O(n^2).Whenthegraphhasalotofedgesthisslowdownd          oesn'tmattermuch.

### AWarmUpforGraphAlgorithms    –TopologicalSorting

Adirectedgraphcanbeusedtorepresentdependenciesbetween nodes.Forexample,thenodesmayrepresentsectionsofalargesoftware projectandanedgefromnodextonodey         meansthatxmustbecomplete beforeycanbecompleted.Orthenodesrepresentcoursesandanedgefromx toymeansthatxisaprerequisiteofy.A     *topologicalordering*  or *topological sort*ofadirectedgraphisalistofthenodesinorderwherethe            edgesallpoint inalefttorightdirection.Inthecaseofcourses,atopologicalsortofthe coursesisanorderingofthecoursesguaranteeingcorrectprerequisites.We willdiscussanalgorithmtofindatopologicalsortofadigraph.Notethatlat             er onwhenwediscussdepthfirstsearch,therewillbeanothermoreelegant methodbasedonDFS.

Onestrategytotopologicallysortadigraphistorepeatedlydeletea nodeofin  -degree0fromthegraph.Howcanwedothis?

Whilethegraphisnotemp     tydo
a.    Findanodeofindegreezero.
b.    Deleteitfromthegraph.

Howcanweaccomplishthestepsaboveandhowmuchtimedoes eachtake?Tocheckifagraphisempty,wewouldhavetolookthroughthe wholearrayandcheckfornil'swhichtakesO(n)time.          Toaccomplishthisin constanttime,wecanjustkeepastaticvariablethatholdsthenumberofnodes inthegraphandcheckifthisiszero.ThismightrequireO(n)timeonceatthe

creationtimeofthedatastructureintheappropriateconstructor.Th        isvariable
wouldbemodifiedbydelete.

Findinganodeofin    -degreezerocanbedonebytraversingthegraph
datastructureandmarkingwhichnodesgetvisited.ThistakesO(n+e)time.
Deletinganodefromthegraphcanbedonebysettingthepointert        onilonthat
node,andtraversingtheotherlinkedlists,deletingthenodeanytimeitisfound.
ThisisalsoO(n+e)time.Theloopgetsexecutedatmostntimessowegeta
totaltimecomplexityofO(n(n+e)).

Thisiswaytooslow.Wecandobetterby      preprocessingthegraph,
calculatingtheindegreeofeachnodeinadvanceandstoringtheminheader
nodes.NowwecansolvethewholeprobleminO(e).Firstwefindanodeof
indegree0,whichtakesO(n)time.Thenwetraverseitslinkedlistandfor        each
nodeinthelist,wesubtract1fromtheappropriateheadernode.Afterwe
finishtraversingthelist,iftheindegreeofanynodeturnsto0weoutputit,and
thenwetraverseitslinkedlist.

Thenaturaldatastructureisaqueuethatholdsthe       nodeswhose
linkedlistsmustbetraversed.Weinitializethequeuetoallthenodesthat
initiallyhaveindegreezero.Whilethequeueisnotempty,wedeleteanode
fromthequeue,traverseitslinkedlist,subtractonefromtheappropriateheader
node,andiftheheadernodeturnedtozero,thenaddittothequeue.Thequeue
makessurethatwedeleteanodecompletelybeforedeletinganynodesto
whichitpointed.Inthisimplementationwedon'tactuallydeleteanynodes
fromthegraph,insteadwe    justmodifytheindegreevalues.Adetailed
examplewillbedoneinclass.

Notethatifthedigraphhasacycle,thenthisalgorithmwillnever
terminate.

### MinimumSpanningTree

Let'sconsiderafamousproblemongraphswhosesolutionwillhelp
use laboratemoreondatastructuresandtheirrelationshiptoalgorithms.
Prim'salgorithmwilluseapriorityqueuethatcanbeimplementedwithRed        -
Blacktreesorheaps,andKruskal'salgorithmcanbeimplementedusinganew
datastructurecalledtheUnion    -Finddatastructurewhichiscomposedoftrees
andlinkedlists.

Bothalgorithmsusea    *greedy* strategy.Thismeansthattheoverall
problemissolvedbyrepeatedlymakingthechoicethatisbestlocally,and
hopingthatthecombinationisbestglobally.    Itisnotreasonabletoexpect
greedystrategiestowork,yettheysometimesdo.Forexample,ifyouwanted
tofindthebestmoveinachessgame,itwouldbenaïvetothinkthattakingthe
opponentsqueenisalwaysbetterthanaslowdefensivepawnmov        e.Thequeen
capturemightbefollowedbyyoubeingcheckmatedthenextmove,whilethe
defensivemovemayresultinawinforyouthirtymovesdowntheline.
Nevertheless,therearecircumstanceswhenthegreedystrategyworksanda
generalmathematical discussionofwhatthesecircumstancesarebringsusto
someseriousmathematicsabout    *Matroids*andtheMatroidIntersection
Problem.TheinterestedreaderisreferredtoLawler'sbook,Combinatorial
Optimization –NetworksandMatroids.

### Prim'salgorith m

Thebasicideaistostartatsomearbitrarynodeandgrowthetreeone
edgeatatime,alwaysaddingthesmallestedgethatdoesnotcreateacycle.
WhatmakesPrim'salgorithmdistinctfromKruskal'sisthatthespanningtree
growsconnectedfromth   estartnode.Weneedtodothisn       -1timestomake
surethateverynodeinthegraphisspanned.Thealgorithmisimplemented

withapriorityqueue.Theoutputwillbeatreerepresentedbyaparentarray whoseindicesarenodes.

Wekeepapriorityque    uefilledwithallthenodesthathavenotyet beenspanned.The *value* ofeachofthesenodesisequaltothesmallestweight oftheedgesthatconnectittothepartialspanningtree.

1.   InitializethePqueuewithallthenodesandsettheirvaluestoa numberlargerthananyedge,setthevalueoftherootto0,and theparentoftheroottonil.
2.   WhilePqueueisnotemptydo{
        LetxbetheminimumvaluenodeinthePqueue;Deletex;
        Foreverynodeyinx'sadjacencylistdo{
                IfyisinPqueueandtheweig        htontheedge(x,y)
        islessthanvalue(y){
        Setvalue(y)toweightof(x,y);Setparentofytox.
                }
        }
}

Anexamplewillbedoneinclass,andyoucanfindoneinyourtext onpage508. Therear  esomedetailsintheimplement    ationthatare  not explicit herebutwhicharecrucial.    Forexample,whenwesetvalue(y)toweightof (x,y),the  heapmustbemod  ifiedandthiscanonlybedoneifyouhavean invertedindexfromthevaluestotheheaplocationswhereeachislocated. Thiscan  bestoredinan    aarraywhosevaluesmustbecontinuallyupdated duringanyheapmodi  fications.

TheanalysisofPrim'salgorithmwillshowanO(elogn)time algorithm,ifweuseaheaptoimplementthepriorityqueue.Step1runsin O(n)time,becauseitisjustbuildingaheap.Step2hasaloopthatrunsO(n) timesbecauseweaddanewnodetothespanningtreewitheachiteration. Eachiterationrequires   ustofindtheminimumvaluenodeinthePqueue.Ifwe useasimplebinaryheap,thistakesO(logn).Thetotalnumberoftimeswe executethelasttwolinesinstep2isO(e)becauseweneverlookatanedge morethantwice,oncefromeachend.Eachti       meweexecutethesetwolines, theremaybeaneedtochangeavalueintheheap,whichtakesO(lgn).Hence thetotaltimecomplexityisO(n)+O(nlogn)+O(elogn),andthelastterm dominates.

Afinalnoteisthatifweuse       an adjacencymatrix  insteadof adjacencylists,Prim 'sideacanbeimpl    ementedwithoutaheap,andthe minimumcalculationO(n)getsdonewhilewe         decidewhethertoupdatethe valuesof  thenodesineachiteration.Thisgivesniterations         ofO(n)and   an O(n^2)totalalgorithmwithmuchsimplerdatastructures.Thequestionof whichmethodofimpleme  ntationisbestdependsobvi  ouslyontherelationship betweeneandn.Thenumber       ofedges isofcourseboundedbetweenO(n)and O(n^2).Forsparsegraphs(fewedges)the       heapa ndadjacencylistmethodis better,andfor  densegraphsthe  two-dimadjacencymatrixi  sbest withouta heapisbetter.

### Kruskal'sAlgorithm

Kruskal'salgorithmalsoworksb  ygrowingthetreeoneedgeata time,addingthesmallestedgethatdoesnotcreateacycle.However,his algorithmdoesnotinsistontheedgesbeingconnecteduntilthefinalspanning treeiscomplete.Westartwithndistinctsinglenodetrees,andt        hespanning treeisempty.Ateachstepweaddthesmallestedgethatconnectstwonodesin differenttrees.

Inordertodothis,wesorttheedgesandaddedgesinascending orderunlessanedgeisalreadyinatree.Thecodeforthisisshownbelow:

```
Foreachedge(u,v)inthesortedlistinascendingorderdo{
     Ifuandvareindifferenttreesthenadd(u,v)tothespanningtree,
andunionthetreesthatcontainuandv.
```

Henceweneedsomedatastructuretostoresetsofedges,whereeach setrepr esentsatreeandthecollectionsofsetsrepresentsthecurrentspanning *forest.* Thedatastructuremustsupportthefollowingoperations:Union(s,t)      – whichmergestwotreesintoanewtree,andFind      -Set(x)whichreturnsthetree containingnodex.

ThetimecomplexityofKruskal'salgorithmwilldependcompletely ontheimplementationofUnionandFind     -Set.TheUnion   -Finddatastructure canbeimplementedinanumberofways,thebestoneshowingitsefficiency onlythrough *amortizedanalysis.*

### Union-FindDataStructure

TheUnion -Finddatastructureisusefulformanaging      *equivalence classes*,andisindispensableforKruskal'salgorithm.Itisadatastructurethat helpsusstoreandmanipulateequivalenceclasses.Anequivalenceclassis simplyasetofthingsthatareconsideredequivalent(satisfiesreflexive, symmetricandtransitiveproperties).Eachequivalenceclasshasa representativeelement.Wecanuniontwoequivalenceclassestogether,create anewequivalenceclass,orfindther     epresentativeoftheclasswhichcontainsa particularelement.Thedatastructurethereforesupportstheoperations, Makeset,UnionandFind.TheUnion    -Findcanalsobethoughtofasawayto maipulatedisjointsets,whichisjustamoregeneralviewof       equivalence classes.

Makeset(x)initializesasetwithelementx.Union(x,y)willunion twosetstogether.Find(x)returnsthesetcontainingx.Oneniceandsimple implementationofthisdatastructureusedatreedefinedbyaparentarray.A seti sstoredasatreewheretherootrepresentstheset,andalltheelementsin thesetaredescendentsoftheroot.Find(x)worksbyfollowingtheparent pointersbackuntilwereachnil(theroot'sparent).Makeset(x)justinitializes anarraywithparen tequaltonil,anddatavaluex.Union(x,y)isdoneby pointingtheparentofxtoy.MakesetandUnionareO(1)operationsbutFind isanO(n)operation,becausethetreecangetlongandthin,dependingonthe orderoftheparametersinthecallstot      heUnion.Inparticularitisbadtopoint thetallertreetotherootoftheshortertree.

WecanfixthisbychangingUnion.Union(x,y)willnotjustsetthe parentofxtoy.Insteaditwillfirstcalculatewhichtree,xory,hasthegreater numberofnodes.Thenitpointstheparentofthetreewiththefewernodesto therootofthetreewiththegreaternodes.Thissimpleideaguarantees(aproof byinductionisonpage453),thattheheightofatreeisatmostlgn.This meansthattheFind    operationhasbecomeO(logn).

### AnalysisofKruskal'sAlgorithm

WiththeUnion -Finddatastructureimplementedthisway,Kruskal's algorithmcanbeanalyzed.ThesortingoftheedgescanbedoneinO(eloge) whichisO(elogn)foranygraph(why?).      Foreachedge(u,v)wecheck whetheruandvareinthesametree,thisisdonewithtwocallstoFindwhich isO(logn),andweunionthetwoifnecessarywhichisO(1).Thereforethe loopisO(elogn).Hencethetotaltimecomplexityis O(elogn).

Itturnsthatifwethrowinanothercleverheuristictoour implementationoftheUnion   -Finddatastructure,wecanimprovethe

performanceofthealgorithm,buttheanalysisrequiresan        *amortizedanalysis.*
Theheuristiciseasytodescribeandthefin        alresulttoo,buttheanalysisisa
littleinvolved,andthereading(22.4)isoptional.Wewillhaveanoptional
recitationforanyoneinterestedinstudyingthedetails.

Thetrickheuristiciscalled    *pathcompression.*  Itisanotherwayto
makethet reesevenshorter.EverytimewetraversepointersbacktodoaFind,
wepointallnodesuptotherootofthetree.Thisisdoneintwophasesbyfirst
findingtherootasnormal,andthengoingbacktoreassigntheparentsofall
visitednodestother    oot.AlthoughtheworstcaseremainsO(logn)foraFind,
theextraconstanttimeworkforpathcompressionbalancestheoccasionallogn
searches.Thatis,everytimewehavealongsearchinaFind,itmakesmany
otherFindsearchesshort.Thedetails    arehairy,buttheupshotisthatp
operationsofUnionandFindusingweightedunionandpathcompressiontakes
timeO(plg*n).Thatis,eachoperationontheaverageistakingO(lg*n).
HenceKruskal'salgorithmrunsintimeO(elg*n).

### TheFunctionl g*n

Notethatlg*nisaveryslowgrowingfunction,muchslowerthanlg
n.Infactisslowerthanlglgn,oranyfinitecompositionoflgn.Itisthe
inverseofthefunction$f(n)=2^2 2^ \ldots ^2$,ntimes.For$n>=5$,$f(n)$isgreater
thanthenumberof    atomsintheuniverse.Henceforallintentsandpurposes,
theinverseof$f(n)$foranyreallifevalueofn,isconstant.Fromanengineer's
pointofview,Kruskal'salgorithmrunsinO(e).Noteofcoursethatfroma
theoretician'spointofview,atrue    resultofO(e)wouldstillbeasignificant
breakthrough.Thewholepictureisnotcompletebecausetheactualbestresult
showsthatlg*ncanbereplacedbytheinverseof$A(p,n)$whereAis
Ackermann'sfunction,afunctionthatgrowsexplosively.Thei    nverseof
Ackermann'sfunctionisrelatedtolg*n,andisanicerresult,buttheproofis
evenharder.

### AmortizedAnalysis

Amortizedanalysisiskindoflikeaveragecaseanalysisbutnotquite.
Inaveragecaseanalysis,wenoticethattheworstcasef        oranalgorithmisnota
goodmeasureofwhatturnsupinpractice,sowemeasurethecomplexityofall
possiblecasesofinputsanddividebythenumberofdifferentcases.In
Quicksort,forexample,thismethodofaveragecaseanalysisresultedinanO(    n
logn)timeanalysisratherthan$O(n^2)$whichisitsworstcaseperformance.

Inamortizedanalysis,wearenotworkingwithonealgorithm,rather
weareworkingwithacollectionofalgorithmsthatareusedtogether.
Typicallythisoccurswhenconside    ringoperationsformanipulatingdata
structures.Itispossiblethatoneofthealgorithmstakesalongtimeworstor
averagecase,butthatithappensrelativelyinfrequently.Sothatevenifthe
worstoraveragecaseforoneoperationisO(n),itisp        ossiblethatthisis
balancedwithenoughO(1)operationstomakeamixtureofpoperationshavea
timecomplexityofO(p),orO(1)peroperation.

Wewillnotgetintothedetailsoftheproofthatpoperationsof
Union-FindneedingatmostO(plg*n).H        owever,youshouldgetaflavorof
whatamortizedanalysisisusedfor,andthisproblemisaperfectmotivationfor
that.Let'sdiscussasimplerdatastructurewithacollectionofoperations
whoseefficiencycanbemoreeasilyanalyzedwithamortizeda        nalysis.

### StacksandaSimpleAmortizedAnalysis

Onewaytoimplementastackiswithalinkedlist,whereeachpush
andpoptakesO(1)time.Anotherwayiswithanarray.Theproblemwithan
arrayisthatweneedawaytomakeitlargerdynamically.        Onewaytodothis
iscalled *arraydoubling.*  Inthisscheme,apopisthesameasyouwouldexpect

andusesO(1),butapushcanbeO(1)orO(n)dependingonwhetherthearray
isfullandneedstobedynamicallyextended.Theideaisthatwewilldouble
thesizeofthearrayanytimeapushwilloverflowthecurrentspace.Ifapush
demandsadoublingofthearrayittakesO(n)todoit.

Henceinthisscheme,popsareO(1)butpushesareO(n)worstcase.
ThethingisthattheO(n)pushesdon'thappen         thatoften.Wecancalculatethis
explicitlywithanexample.Let'ssaywearepushingnineelementsintothe
array.Thearrayneeds    togetdoubledwhenweaddthe2 $^{nd}$,3 $^{rd}$,5 $^{th}$,and9 $^{th}$
elements.Thetimeforthesedoublingsis1,2,4,and8stepsre           spectively.The
timefortheactualpushesis9.Thisgivesatotaltimeof2(8)       $-1+(8+1)=$
3(8).Ingeneral,thetimetopushnelementsintothearrayis3n,(recallthethat
sum1+2+4+…+nequals2n    -1).Thismeansthatovernpushesweusean
averageof3stepsperpush,eventhoughtheworstcasepushisO(n).

Therearemanywaystothinkaboutamortizedanalysis,butIthink
theaboveideawillgiveyoutheflavorintheclearestway.Anotherwayto
thinkofit,isthatweaddtwostepsint        oasavingsaccount,everytimewedoa
fastpush,makingtheexpenditureforeachfastpushthreeinsteadofone.Then
wecashinonthissavingsonaslowpush,bywithdrawingn         -2stepsforthe
doubling.Thiswayeachlongpush(n+1steps)isaccompli        shedwiththree
stepsplusthen    -2wesavedupfromtheshortpushes.

Itisthisaccountingschemethatmustbedeterminedandspecifiedin
everyamortizedanalysis.Eachproblemrequiresitsowningenuityand
cleverness.

### DepthandBreadthFirstSear    ch

Withanydatastructurethefirstbasicalgorithmthatwewriteisone
thattraversesit.Inthecaseofgraphs,thetwobasicmethodsoftraversalare
breadthfirstanddepthfirst.Itturnsoutthateachoneofthese,butdepthfirst
searchinpar  ticular,canbemodifiedtosolveavarietyofproblemsongraphs
whileittraversesthegraph.Boththesealgorithmsrunintimeproportionalto
theedgesofthegraph.

### BreadthFirstSearch

Breadthfirstsearchtraversesagraphinwhatissometimec         alled *level
order.*Intuitivelyitstartsatthesourcenodeandvisitsallthenodesdirectly
connectedtothesource.Wecalltheselevel1nodes.Thenitvisitsallthe
unvisitednodesconnectedtolevel1nodes,andcallstheselevel2nodesetc.

Thesimplewaytoimplementbreadthfirstsearchisusingaqueue.
Infactwhenhearyou    *breadth*youshouldthink   *queue*,andwhenyouhear
*depth*youshouldthink    *stack*.Wehavethealgorithmoutputatreerepresenting
thebreadthfirstsearch,andstoreth     elevelofeachnode.Thecodecanbe
foundonpage470ofyourtext.Hereisaperhapsmorereadableversion.We
haveaparentarraytostorethesearchtree,aLevelarraytostorethelevel,and
avisitedarraytorememberwhohasalreadybeenplaced       onthequeue.The
bookusesathreevaluedcolorsystemwhite,greyandblackinsteadofthis
Booleanarray.Idon'tknowwhythisisnecessaryandIamnotsureitis..

```
Initialize:  QueueQ=source;level[source]=0,p[source]=nil;
             Foralln odesxdovisited[x]=false;
             visited[source]=true;

Loop:WhileQisnotemptydo{
                    x=deleteq(Q);
                    forallyadjacenttoxdo
                         ifvisited[y]=false{
```

```
                                        visited[y]=true;level[y]=level[x]+1;
                                        p[y]=x;addq(Q,y)}
                        }
```

Thetotaltimefor   initializingisO(n)andthetotaltimeforthe queuingoperationsisO(n)becauseeachnodeisputonthequeueexactlyonce. ThetotaltimeinthemainloopisO(e)becausewelookateachedgeatmost twice,oncefromeachdirection.ThisgivesatimecomplexityofO(n+e).

### DepthFirstSearch

Depthfirstsearch(DFS)traversesagraphbygoingasdeeplyas possiblebeforebacktracking.Itissurprisinglyrichwithpotentialforother algorithms.Italsoreturnsasearchtree.Itdoesnotreturnt         helevelofeach node,butcanreturnanumberingofthenodesintheorderthattheywere visited.Wefirstshowadepthfirstsearchskeletonanddefinethedifferent kindsclassesofedges.Thenweshowhowtoaugmenttheskeletontosolve twoveryba sicalgorithms:topologicalsorting,connectedcomponents.Eachof leveragesthepowerofDFSatadifferentlocationintheskeleton.We concludewithasophisticateduseofDFSthatfindsstronglyconnected componentsofadirectedgraph.Youmayreca        llthatinmonth0wediscussed amethodinlinearalgebrausingmatrixmultiplicationthatsolvedthis algorithminO(n^3).OurmethodwillworkinO(n+e).Thereareother sophisticatedusesofdepthfirstsearchincludinganalgorithmtofindbi       - connectedcomponentsinundirectedgraphs,andanalgorithmtodetermine whetheragraphisplanar.Neitheroneoftheseproblemshasanobviousbrute forcesolutionandcertainlynotanefficientone.

AsimilarDFSskeletoncanbefoundinyourtextonpage         478.

### DepthFirstSearchSkeleton

DFS(G,s)

Marksvisited;Dfsnum[s]=count;count++;
//countisaglobalcounterinitializedto1.
/*Processs  –previsitstage*/
RecursiveLoop:Foreveryyadjacenttosdo
            ifyisunvisitedthen{DFS(G,y);pa     rent[y]=x;}else…
                                /*processedges{s,y}*/;
/*Processs  –postvisitstage*/
Marksfinished;

Thepointsatwhichthisalgorithmcanbeaugmentedarethreefold:

1.  Afterthenodeismarked,beforelookingforwardonitsedges,
    (previsitstage).
2.  Whileweprocesseachedge(processedgestage).
3.  Afterallchildrenofanodehavebeensearched(postvisitstage).

Stagetwoprocessesedges.Edgescanbeclassifiedintofour categories(onlythefirsttwoexistforundirectedgraphs):treeedges,back edges,crossedgesanddescendantedges.Definitionscanbefoundonpage 482ofyourtextbutthebestwaytounderstandtheclassificationistogo throughaDFSexample,exhibitthesearchtreeandidentifythedifferentedges. Wewilldothisinclass.     Anexampleappearsinthebookonpage479.

### DFSinComparisonwithBFS

Itisstagethreethatgivesdepthfirstsearchallitspotentialpower.Atthis postordertime,manynodeshavealreadybeenexaminedandalotof

informationmaybeavailableto    thealgorithm.Thereisnoanaloguetothisin
breadthfirstsearchbecausethereisnorecursion.Thisiswhythereareso
manymoreapplicationsofDFSthanthereareforBFS.

### ConnectedComponents  –ASimpleApplication

Oureyescanpickouttheconn    ectedcomponentsofanundirected
graphbyjustlookingatapictureofthegraph,butitismuchhardertodoit
withaglanceattheadjacencylists.BothBFSandDFScanbeusedtosolve
thisproblem,becausethereisnopostordertimeprocessing.The          basictrickis
towrapthesearchinaloopthatlookslikethis:

Foreachnodexinthegraphdo
         Ifxisunvisited{markx;Search(x);}

DoingthiswithDFS,wecankeepaglobalcounterandassigna
numbertoeachconnectedcomponent.Thecounterki       sinitializedto1.During
the *processedge* stage,wethrowanyedgevisitedonastack.Whenwefinish
Search(x);thenbeforewegobackuptothetopoftheforloop,wepopall
edgesoffthestackuntilthestackisempty.Thenweoutputtheseedges          witha
headerthatsaysconnectedcomponentk,andweincrementk.

Fordirectedgraphs,itispossibleforasearchinasingleunderlying
connectedcomponenttofinishwithouttraversingallofthenodes.Thisis
becausethedirectionofthearrowsmigh       teffectivelydisconnectonepartofthe
graphoneanothereventhoughtheyareinthesameunderlyingundirected
connectedcomponent.Seethepicturebelow.



HenceindirectedgraphsacallDFSisalwayswrappedinalooplike
theon eabove.Inbothalgorithmsthatfollow,wemakeuseofthepostorder
processingavailableinDFS,andweassumethattheDFScalliswrappedina
loopthatchecksunvisitednodes.

### TopologicalSorting  –RevisitedwithDFS

Whenweintroducedthegraphda    tastructure,wedesignedan
algorithmfortopologicalsorting.Herewedesignanalternativealgorithm
usingDFS.Thisalgorithmdependsverystronglyonpostordertime
processing.Weuseaglobalcounterinapostorderprocessingstepthatassigns
finishingtimestoeachnode.Thenodesarepushedontoastackintheorder
theyfinish.Whenthesearchisover,poppingthestackliststhenodesin
topologicalorder.Thismakessensebecauseweonlyfinishanodeafterallits
descendantsarefinished .Hencethefinishingorderisthereverseofa
topologicalsorting.

### StronglyConnectedComponents  –ACoolApplicationofDFS

Thedescriptionofthisalgorithmlikemostvariationsofdepthsearch
isdeceptivelysimpletodescribebuttediousandcompl       extoprovethatit
works.

1.    CallDFSandcalculatefinishingtimesforeachnode.

2.  CallDFSonthetranspose(reversealledges),butconsiderthe
    nodesinreversefinishingtimeorder.
3.  EachconnectedtreeintheDFSforestoftreesisaseparate
    stronglyc onnectedcomponent.

Thisalgorithmseemstoworklikemagic,anditisindeedabit
amazing.Yourtextspendsfourpagesandthreetheoremsconvincingyouof
this,andwedonotreproducethatinformationhere.

### ShortestPathAlgorithms

Thesealgorit hmsareperhapsthemostwellknownalgorithmoutside
ofsorting.MostpeoplehaveheardofDijkstra'sshortestpathalgorithm.The
presentationherefollowsthespiritofTarjanin *DataStructureandNetwork
Algorithms.* Ourtext'spresentationisvery    similar.Bytheway,thepreviously
mentionedbookisanexcellentresourceforadvanceddatastructuresandbasic
graphalgorithms.

Theshortestpathalgorithmgivesagraphandastartnode,andasks
fortheshortestpathsfromthatnodetoeveryother    nodeinthegraph.Notethat
itdoesnotsaveanytimeingeneraltoconsideraparticulargoalnode,hencewe
mightaswellcalculatetheshortestpathstoalltheothernodes.Thereare
versionwhereyouwishtocalculatealltheshortestpathsbetwee    nanytwo
nodescanbesolvedbyrepeatingthesingesourcealgorithmntimes,orwitha
completelydifferenttechniqueusingdynamicprogramming.Wespeakabout
theallpairsshortestpathalgorithmnextweekwhendiscussingthetopicof
dynamicprogramm ing.

TheshortestpathproblemisNP    -completewhenthegraphhas
negativeweightcycles,hencethelongestpathproblemisNP    -completeinthe
presenceofpositiveweightcycles,whichisthecommonsituationwithagraph.

### SingleSourceShortestPath

Theoutputofthisalgorithmisashortestpathtreeandadistance
array.Thearraydist[x]storestheshortestpathdistancefromthesourcetox,
wheredist[s]=0.Thetreeisstoredviaaparentarray,(liketheminimum
spanningtree),whereparent[ source]isnil.Themaintoolusedinthealgorithm
iscalled *scanning* anode.Scanninganodelooksatallthenodesadjacenttoit
anddecideswhethertochangetheirparentanddistvalues.Yourtextbookcalls
anaffirmativedecision *relaxing* anode. Ideally,wewouldliketoscaneach
nodeexactlyonce,butthatisnotalwayspossible.

```
Scan(x)

Foreverynodeyadjacenttoxdo{
        Ifdist[x]+length(x,y)<dist[y]{
                dist[y]=dist[x]+length(x,y);parent[y]=x;
        }
}
```

Thecodelooksateac   hnodeyadjacenttoxandseeswhetherornot
thepathfromthesourcetoythroughxisshorterthantheshortestpath
currentlyknownfromthesourcetoy.Ifitisbetter,thenwechangethecurrent
distandparentvaluesfory.

Westartthealgorith   mbyinitializingthedistandparentarrays,and
scanningthesourcenode.

**ShortestPathSkeleton**

Initialize:                    forallnodesxdo{dist[x]=MAX;parent[x]=
nil;}

                              dist[s]=0;

Main:              Scan(s);
ScanningLoop:Scantheothernodesinsomesuitableord          er;

**Dijkstra'sAlgorithm**

Thepointnowistoconsidertheorderinwhichtoscanthenodes,
andthisdependsonthekindofgraphwearegiven.Foragraphwithonly
positiveedges,weuseagreedyapproachofscanningnodesinascendingorder
ofcurr entdistancevalues.ThisiscalledDijkstra'salgorithm.Onceanodeis
scanned,itisneverscannedagainbecausewecanprovethatscanninginthis
ordermeansthatthedistandparentvaluesforthatnodewillsubsequently
neverchange.Ifwenever    *relax* anodeafteritisscannedthenthereisnoneed
toscaniteveragain.

Anefficientwaytoimplementthisalgorithmistokeepaheapofthe
currentlyunscannednodesbytheirdistvalues,andmaintaintheheapthrough
possiblechangesindistvalu  es.GettingthenextnodetoscanisO(logn),and
wescaneachnodeexactlyonceduetothetheoremwementioned.Thisgives
O(nlogn).Wealsomustconsiderthemaintenanceoftheheap,whichtakes
O(logn)butcanhappenasmanyasO(e)times.Henc       ethetotaltime
complexityisO(elogn).Withoutusingaheap,thetimecomplexitycanbe
analyzedtobeO(n^2).Notethatifyouusead       -heap(aheapwithdchildren,
whered=2+e/n),oriftheedgeweightsarerestrictedtosmallintegers,then
wecanimprovethesetimecomplexities,butwewillnottalkaboutthese
advancedvariations.

ExamplesofDijkstra'salgorithmcanbefoundinyourtext(page
528)andwewilldooneinclass.Notethatwedidnotdiscussjustwhat
happenswithDijkstra' salgorithminthepresenceofnegativeweightedges.I
leavethisforyoutothinkabout.

**AcyclicDirectedGraphs  –TopologicalOrderScanning**

Anotherwaytoguaranteethatonceanodeisscannedthatitnever
needstobescannedagain,istoscanth       enodesintopologicalsortedorder.In
thiscase,noscannednodecaneverbe    *relaxed* later,becausetherearenoedges
comingbacktothisnodeinatopologicalordering.Nofancytheoremhere,
justsimplecommonsense.Ofcourseconstructingatopolo       gicalorderingis
onlypossibleinadirectedacyclicgraph.Notethisworkswhetherornotthe
graphhasnegativeweightedges.

**TheBellmanFordShortestPathAlgorithmforGraphswith
NegativeWeightEdgesbutNoNegativeWeightCycles      –BreadthFirst
Scanning**

Theshortestpathproblemis NP    -completeinthepresenceofnegative
weightcycles,butitcanbesolvedinpolynomialtimeforagraphwithnegative
weightedgesandcycles,aslongastherearenonegativeweightcycles.

Thealgorithmuses  abreadthfirstscanningorder.Themain
differencebetweenthisalgorithmandthepreviousalgorithms,isthatinthis
casewecannotguaranteethateverynodewillbescannedexactlyonce.We
mayhavetorescananodemultipletimes.Thekeyisthatw        emustscaneach

nodeatmostntimes.ThisresultsinanO(ne)timealgorithm.Thedetails
behindtheseclamsarenotatallobvious.

Tounderstandwhyeachnodeisscannedatmostntimesandwhythe
complexityisO(ne),ithelpstotalkaboutthede        tailsofimplementation.
Breadthfirstscanningcallsfortheuseofaqueue.Weinitializeaqueuetothe
sourcenode,andwhilethequeueisnotempty,weremoveanodefromthe
queue,scanit,andsettheparentanddistvaluesofitsadjacentnode
appropriately.Weaddanodetothequeuewheneveritisrelaxed(thatis,when
itsparentanddistvaluesarechanged).Noteifanode,isalreadyonthequeue
whenitisrelaxed,thenwedonotputitonagain,wesimplyleaveiton.This
impliesthata tanygiventimenonodeisonthequeuemorethanonce.

Let'sanalyzethealgorithmbylookingat    *phases.* The0thphaseof
thequeueiswhenitconsistsofjustthesourcenode.Theithphaseconsistsof
thenodesonthequeueafterthei        -1stphasen odeshavebeenremoved.There
isacrucialtheoremprovedbyinductionthatstatesthatifthereisashortest
pathfromthesourcetoanodexcontainingkedges,thenjustbeforethekth
phaseofthealgorithm,dist[x]willequalthelengthofthispath        .Sinceanypath
withoutcyclesfromthesourcetoanodecancontainatmostn        -1edges,this
meansthealgorithmneedsatmostO(n)phases.Moreover,sinceeach
individualphasehasnoduplicatenodesonthequeue,atmostnnodescanbe
onthequeuein  agivenphase.Processingaphasemeansdeletingandscanning
thesennodes.ThisprocessingtakesO(e)time,becausetheworstcaseisthat
welookateveryedgeadjacenttotheenodes.SinceweprocessatmostO(n)
phaseswithO(e)timeperphase,thi    sgivestheO(ne)timecomplexity.


### GeometricAlgorithms

Geometricalgorithmsarewonderfulexamplesforprogramming,
becausetheyaredeceptivelyeasytodowithyoureyes,yetmuchharderto
implementforamachine.

Wewillconcentrateonaparticula   rproblemcalledconvexhull,
whichtakesasetofpointsintheplaneasitsinputandoutputstheirconvex
hull.Wewillstayawayfromformaldefinitionsandproofshere,sincethe
intuitiveapproachwillbeclearerandwillnotleadyouastray.Toun        derstand
whata *convexhull* is,imaginethatanailishammeredinateachpointinthe
givenset,theconvexhullcontainsexactlythosepointsthatwouldbetouched
byarubberbandwhichwaspulledaroundallthenailsandletgo.The
algorithmisused  asawaytogetthenaturalborderofasetofpoints,andis
usefulinallsortsofotherproblems.

ConvexHullisthesortingofgeometricalgorithms.Itis
fundamental,andastherearemanymethodsforsorting,eachofwhich
illustratesanewtech   nique,soitisforconvexhull.

### GrahamScan

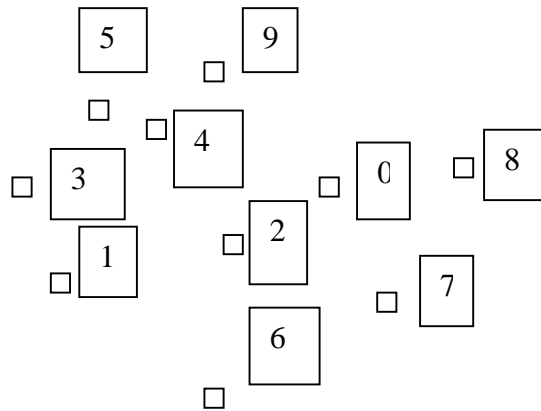TheparticularalgorithmwewillimplementforConvexHullisdueto
RonGrahamandwasdiscoveredin1972.GrahamScan,asitiscalled,works
bypickingthelowestpointp,i.e.theonewiththeminimump.        yvalue(note
thismustbeontheconvexhull),andthenscanningtherestofthepointsin
counterclockwiseorderwithrespecttop.Asthisscanningisdone,thepoints
thatshouldremainontheconvexhull,arekept,andtherestarediscarded
leaving onlythepointsintheconvexhullattheend.

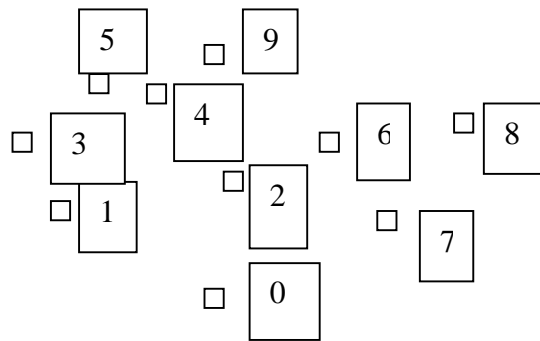Toseehowthisisdone,imaginefirstthat,byluck,allthepoints
scannedareactuallyontheconvexhull.Inthatcase,everytimewemovetoa
newpointwemakealeftturnwithrespecttothelinede        terminedbythelast
twopointsonthehull.Therefore,whatGrahamScandoes,istocheckifthe

nextpointisreallyaleftturn.IfitisNOTaleftturn,thenitbacktrackstothe
pairofpointsfromwhichtheturnwouldbealeftturn,anddiscards          allthe
pointsthatitbacksupover.Becauseofthebacktracking,weimplementthe
algorithmwithastackofpoints.Anexampleisworthathousandwords.The
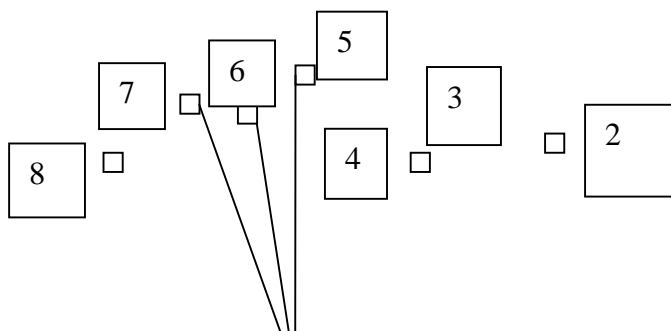inputlistofpointsis:

(A,0,0)(B,    -5, -2)(C,  -2, -1)(D,  -6,0)(E,   -3.5,1)(F,   -4.5,
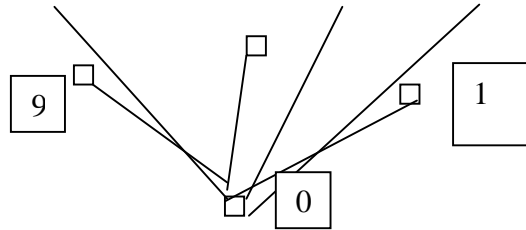1.5)

(G, -2.5, -5)(H,1,    -2.5)(I,2.5,.5)(J,     -2.2,2.2).

Thearrayofinputpointsisshownabovelabeledbyindexinthe
array(ratherthantheircharlabel).ThepointlabeledAisinindex0,Bisin
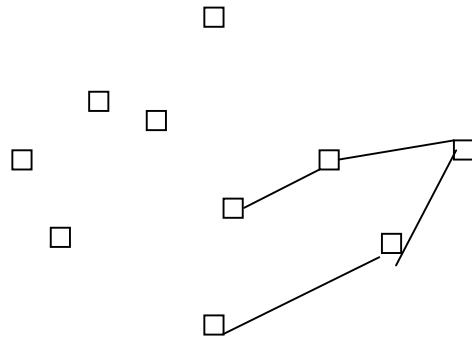index1,etc.Thelowestpointiscomputedandswappedwiththepointinindex
0ofthearray,asshownbelow.

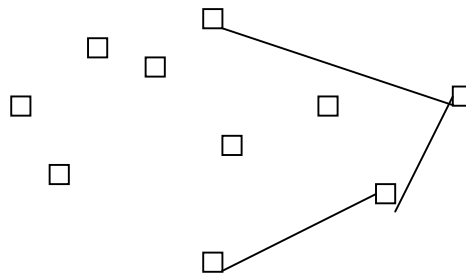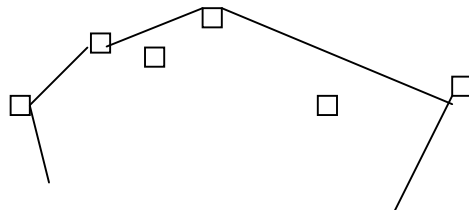Thepointsarethensortedbytheirpolarangleswithrespecttothe
lowestpoint.

   Thepointsaresortedandrearrangedinthearrayasshownabove.
Theturnfromline0   -1topoint2isleft,from1     -2to3isleft,from2     -3to4is
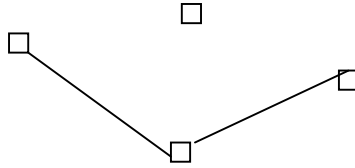left.Nowthestackcontainsthepoints01234.Thisrepresentsthepartialhull
int hefigurebelow.



   Theturnfromline3    -4topoint5isright,sowepopthestack.The
turnfrom2 -3to5isright,sowepopagain.Theturnfrom1           -2to5isleft,so
wepush5onthestack.Thestacknowhas0125,andthe            picturelookslike
this:



   Theturnfromline2    -5to6isleftso6ispushedonthestack.Then
theturnfrom5  -6to7isright,so6ispoppedand7ispushedbecausetheturn
fromline2 -5to7isleft.Therestoftheturns         areleft,so8and9arepushedon
thestack.Thefinalstackis0125789,andtheconvexhullisshownbelow:

**GrahamScanPseudo -code:**Thealgorithmtakesanarrayofpointsand returnsanarrayofpointsrepresenting    theconvexhull.

1.  Findthelowestpointp,(thepointwiththeminimumycoordinate).If thereismorethanonepointwiththeminimumycoordinate,thenusethe leftmostone.
2.  Sorttheremainingpointsincounterclockwiseorderaroundp.Ifany pointshavethesameanglewithrespecttop,thensortthembyincreasing distancefromp.
3.  Pushthefirst3pointsonthestack.
4.  Foreachremainingpointcinsortedorder,dothefollowing:
          b=thepointontopofthestack.
          a=thepointbelowthatonthes        tack.
          WhilealeftturnisNOTmadewhilemovingfromatobtocdo
                  popthestack.
                  b=thepointontopofthestack.
                  a=thepointbelowthatonthestack.
          Pushconthestack.
5.Returnthecontentsofthestack.

**ImplementationDetailsfor thePointClass** :

**PrivateData:**
        Westartbydefiningasimplegeometricclass      *point* anddecidingon theappropriateprivatedataandmemberfunctions.A       *point*shouldhavethree fields:twoarefloatforthe    *x*and *y*coordinates,andoneisacharforthe      name ofthepoint.

**Constructors:**
        Athreeparameterconstructorshouldbecreatedtosetuppoints.

**Methods:**
        Anoutputmethodtoprintoutapointbyprintingitsname(char) alongwithitscoordinates.

        Accessormethodstoextractthe    *x*or *y*co ordinatesofapoint.

        Astaticdistancemethodtodeterminethedistancefromonepointto another.

        A *turn-orientation* methodthattakestwopoints    *b*and *c*andreturns whetherthe *sweepingmovement* fromthelinea  -btothelinea  -cgoesclockwise (1), counterclockwise( -1)orneither(0).(Theresultisneither(0)whena,b andcareallonthesameline.)Thisfunctionisnecessaryfordecidingwhether alefttorrightturnismadewhenmovingfrom      *a*to *b*to *c*instep4ofthe pseudo-codeabove.It    isalsousefulforsortingpointsbytheirpolarangles.

        Itmaynotbeobvioushowtoimplementthisfunction.Onemethod isbasedontheideaofthecrossproductoftwovectors.Let        *a*, *b*and *c*be

points,where $x$ and $y$ areaccessormethodstoextractt    he $x$ and $y$ coordinates respectively.

if( $c.x - a.x)(b.y - a.y)>(c.y  - a.y)(b.x - a.x)$ thenthemovementfromline    $a$-$b$ toline  $a$-$c$ isclockwise.
if $(c.x - a.x)(b.y - a.y)<(c.y  - a.y)(b.x - a.x)$ thenthemovementfromline    $a$-$b$ toline  $a$-$c$ iscountercloc  kwise.
Otherwisethethreepointsareco    -linear.

Tounderstandthisintuitively,concentrateonthecasewherethelines a-banda  -cbothhavepositiveslope.Aclockwisemotioncorrespondstothe linea -bhavingasteeper(greater)slopethanlinea      -c. Thismeansthat  $(b.y - a.y)/(b.x - a.x)>(c.y  - a.y)/(c.x - a.x)$ .Multiplythisinequalityby  $(c.x - a.x)(b.x - a.x)$ andwegettheinequalitiesabove.

Thereasonfordoingthemultiplicationandtherebyusingthis      *cross product* is:
1.Toavoidhavi   ngtocheckfordivisionbyzero,and
2.Sothattheinequalityworksconsistentlyforthecaseswherebothslopesare notnecessarilypositive.(Youcancheckforyourselfthatthisistrue).

GrahamScanshouldbecodedusinganabstractSTACKclassof points.Thesortinginsteptwocanbedonebycomparingpairsofpointsvia theturn -orientationmethodwithrespecttothelowestpoint(object      $p$).An *interface*(ifyouuseJava)maybeconvenienttoallowthesortingofpoints.

**ANoteonComplexity**  :

ThecomplexityofGrahamScanis     $O(nlogn)$ .Wewilldiscuss informallywhatthismeansandwhyitistrue.Itmeansthatthenumberof stepsinthealgorithmisboundedasymptoticallybyaconstanttimesnlogn wherenisthenumberofpointsinth      einputset.Itistruebecausethemost costlystepisthesortinginstep2.Thisis      $O(nlogn)$ .Step1takestime   $O(n)$. Step3takes  $O(1)$.Step4istrickiertoanalyze.Itisimportanttonoticethat althougheachofthe   $O(n)$pointsareprocessed  ,andeachmightintheworst casehavetopopthestack    $O(n)$times,overallthisdoesNOTresultin    $O(n^2)$ . Thisisbecauseoverall,everypointisaddedtothestackexactlyonceandis removedatmostonce.Sothesumofallthestackoperationsis      $O(n)$.

Therearemany  $O(nlogn)$  and $O(n^2)$ algorithmsfortheconvexhull problem,justastherearebothforsorting.Fortheconvexhullthereisalsoan algorithmthatrunsin   $O(nh)$,where  $n$isthenumberofpointsintheset,andhis thenumberofpoi   ntsintheconvexhull.Forsmallconvexhulls(smallerthan *logn* )thisalgorithmisfasterthan    *nlogn* ,andforlargeconvexhullsitis slower.

### Jarvis'AlgorithmforConvexHull

Jarvis'algorithmusessomeofthesameideasaswesawinGraham Scan butitisalotsimpler.ItdoesnobacktrackingandthereforedoesNOT needtouseaSTACKclass,althoughitstillmakesuseoftheARRAYclass templatewithyourpointclass.

Asbefore,westartbyaddingthelowestpointtotheconvexhull. Thenwe  repeatedlyaddthepointwhosepolaranglefromthepreviouspointis theminimum.Thisminimumanglecomputationcanbedoneusingthe clockwise/counterclockwisememberfunction,similartohowthesortingstep (step2)ofGrahamScanusesthefunction.

Thecomplexityofthismethodis     $O(nh)$where  $h$isthenumberof pointsintheconvexhull,becauseintheworstcasewemustexamine      $O(n)$ pointstodeterminetheminimumpolarangleforeachpointinthehull.