

Tagging SNPs

Selection

2 Algorithms

A1 LD-Select

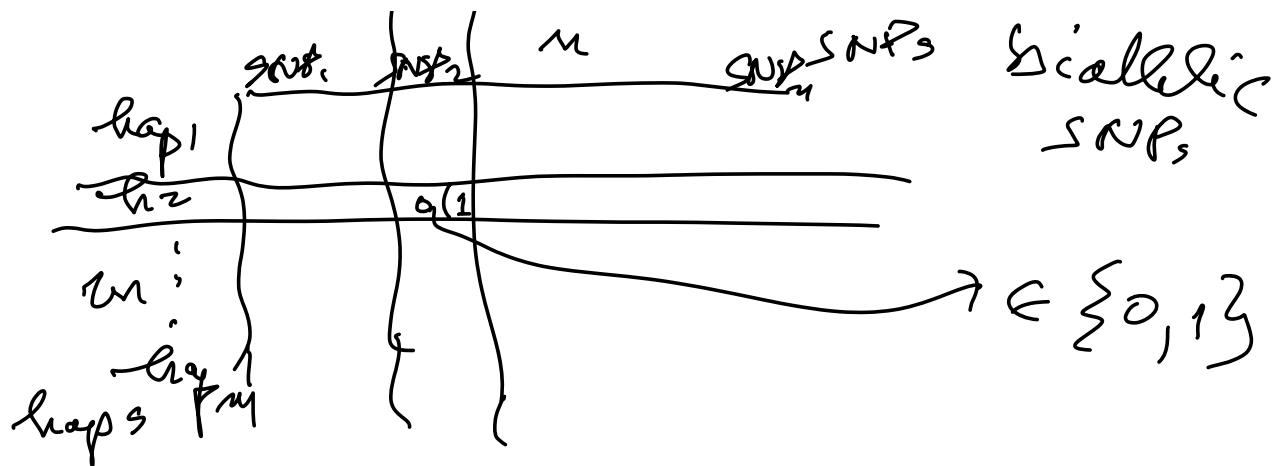
A2 Informativeness

A1 continuation

"most informative SNPs"

Input is a matrix

haplotypes \times SNPs

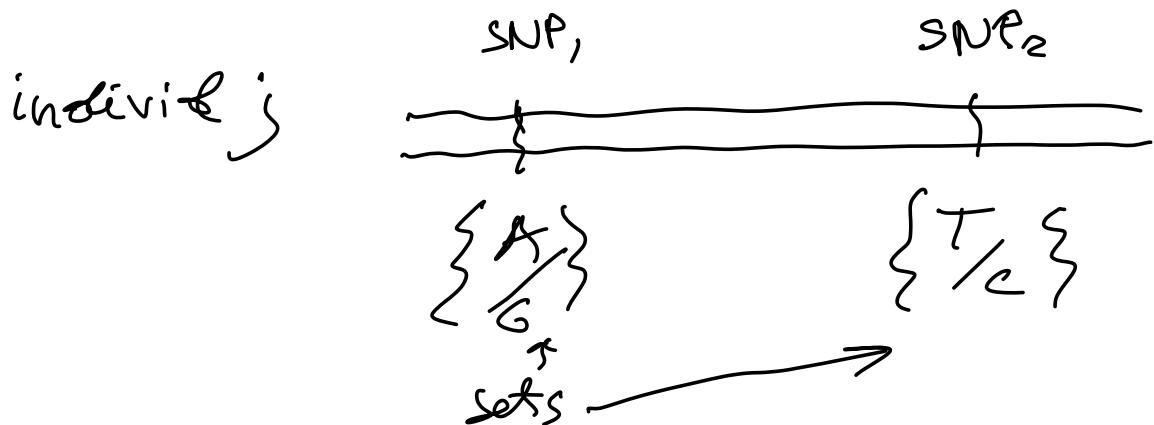


We have individuals genotyped at every SNP of a certain set of SNPs part of an array or assay of SNPs.

The result: for each SNP_i and individual j

output: $\{ \begin{matrix} A \\ T \end{matrix} \}$ the two

alleles of the individ. j at SNP_i.



$\{A, G\}, \{T, C\}, \dots \{A, A\}$

We need haplotypes

$\{A, G\}, \{T, C\}$

Possible haplotypes:

A	T	}
A	C	
G	T	4 but only
G	C	

2 occur
in individual;

With one individual only
we cannot infer which two
of the 4 haplotypes are
the true haplotypes of
individual j.

If we have a set of individuals
human genomes:

This problem is called PHASE
the haplotype phasing PB.

Maps: sequences over the alphabet $\{0, 1\}$
the two alleles are
• 0 and 1
0101110

Genotypes: sequences over the alphabet
 $\{0, 1, 2\}$

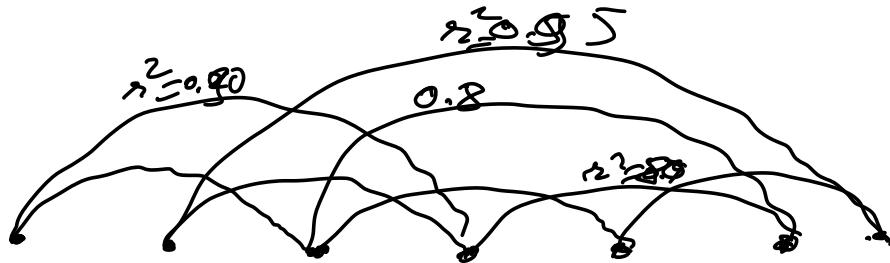
individual $\begin{cases} 0101110 \\ 0001001 \end{cases}$
 genotype (summing): $\frac{0201222}{\text{ex. hap over just 3 SNPs}}$
 $\overbrace{\text{genotype}}^{(1)}$
 $022 \rightarrow \begin{cases} 2010 \\ 2001 \end{cases} \rightarrow$ "explanations"
 $\rightarrow \begin{cases} 011 \\ 000 \end{cases}$

"2" = ambiguous
 { } determined

The phasing problem has a
 search space exponential
 in the number of '2's in the
 genotype sequence

INPUT: for the tagging
 SNPs Selection Φ_b , is
 haplotype $\begin{cases} 0101110 & 0111 \\ 0001010 & 101 \\ 1110100 & 01010 \end{cases}$ NP
 $\star_i \star_j \star_k$ $r^2(i, j)$

LD-select r^2



SNPs,

Complete undirected graph over the n SNPs/vertices and each edge is labelled by the r^2 value between the two vertices joined by the edge.

LD-select is a greedy algorithm. Two constraints as input to the algorithm:

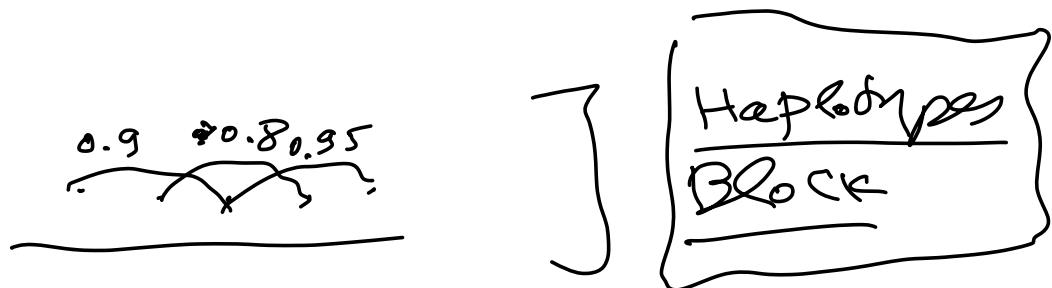
(1) only common SNPs are part of the input

$\geq 10\%$ of people in
the population have

(2) there are NPs
there is an ϵ_2^2 threshold
given: γ

Only edges with $\epsilon_2^2 \geq \gamma$
are considered

$$\text{e.g. } \epsilon_2^2 \geq 0.3$$



Max infoline : vertex that is connected to most other vertices with edges with $\ell_i \geq 7$.

Max infoline vertex $\xrightarrow{(MIV)}$ start a bin

$\{ MIV_1, \dots \} \rightarrow \text{Bin}_1$

$\Rightarrow MIV \quad \{ MIV_2, \dots \} \rightarrow \text{Bin}_2$

:

$\longleftarrow \text{Bin}_{805}$

The tagging SNPs :

$MIV_1, MIV_2, \dots, MIV_{805}$

Paper: "Selecting a Maximally infoline set of Single Nucleotide

polymorphisms for association analysis using linkage disequilibrium"

Am. J Hum. Gen (2004)

Cochran et al Leonid Kruglyak
Debbie Nickerson

2 populations

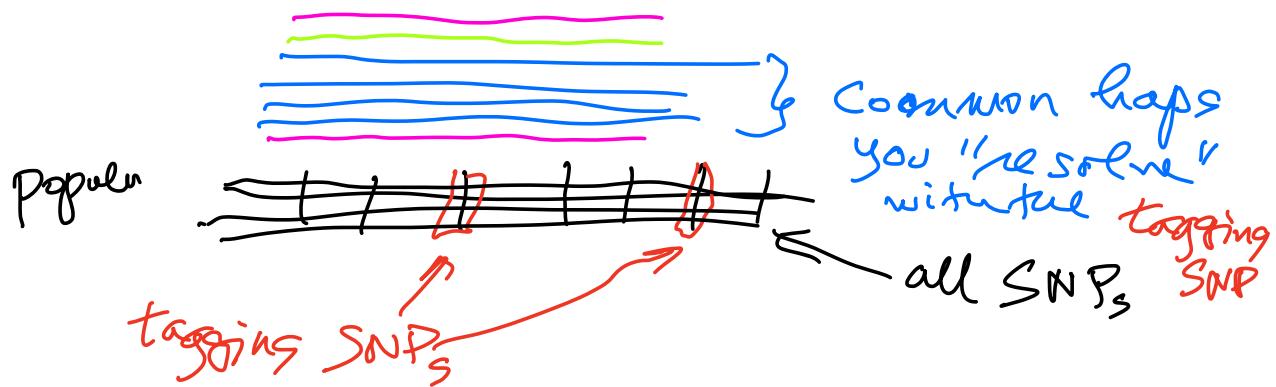
- African American AA
- European American EA

Inferred haplotypes (To Phase or
PHASE Not to Phase)

r^2 was calculated between pairs
of SNPs in each population

haplotypes $\begin{cases} \text{Common} & (\text{many people} \\ & \text{have these} \\ & \text{haplotypes}) \\ \text{rare} & \end{cases}$

One measure of success for selecting tagging SNPs is



Their LD-selected captured/recovered
 $\geq 80\%$ of the common haplotype

Capture the variation = haplotypes
 $\geq 80\%$ of common haplotype

Results

100 genes resequenced them
in 24 African Americans
23 European Americans

Gene average: 16.5 kb length
longest 45 kb length

8877 SNPs overall

7793 SNPs in AA pop

4620 SNPs in EA pop

Density 1 SNP in 200 bp

very small # of three-allelic SNPs

when $\text{MAF} > 10\%$ \Rightarrow 3178 SNPs
^{min allele freq} in AA
2375 in EA

Common SNPs 1 in 200 bp AA
1 in 700 bp EA

SNPs in
coding sequences (genes)
cSNPs 8% of
seq length
(135 kb)

The Second Algorithm

INFORMATIVE VALUES

LD-Select: the objective function
way: Minimum DOMINATING SET

Informativey: the objective function
is : Minimum SET COVER

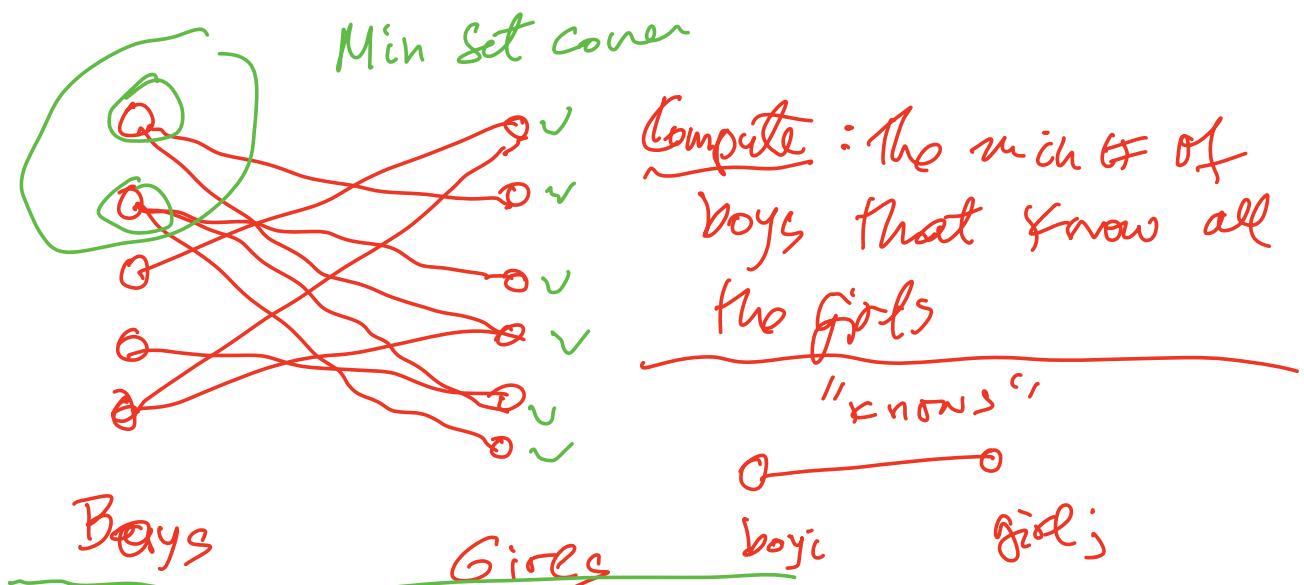
Dominating Set \leftrightarrow Set Cover

Both PBs are NP-Complete

reducible to each other

"Unify LD-Select & Informative?"

SET COVER



INFORMATIVENESS ALGORITHM

"information theory"

"bit": what a SNP does?

what classifier it gives

	SNP ₁	SNP ₂	SNP ₃	SNP ₄	SNP ₅	SNP ₆	SNP ₇	SNP ₈
indiv 1	0	0	0	0	1	0	1	1
indiv 2	1	0	1	0	1	1	1	1
indiv 3	0	0	0	1	1	0	0	0

$m = 8$

$n = 3$

column groups

SNP₃ graph

SNP₃ graph

indiv 1

indiv 2

indiv 3

D-graphs

D-edges

D-edges

A SNP provides one "bit" of information

indiv 1 and indiv 3 have different alleles at the SNP

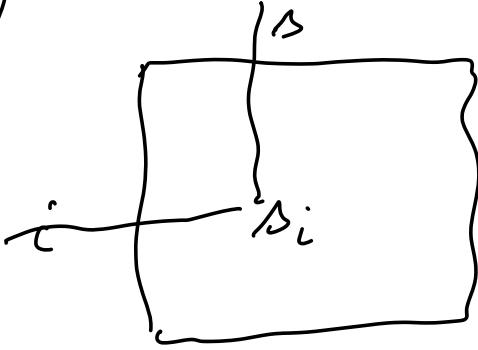
different alleles

D-edge "distinguishing" edge
 "different" edge

INPUT. A matrix M in haplotypes \times SNPs

$$S_i = M[i, s]$$

allele



Def For SNPs and two haplotypes i and j

let us consider

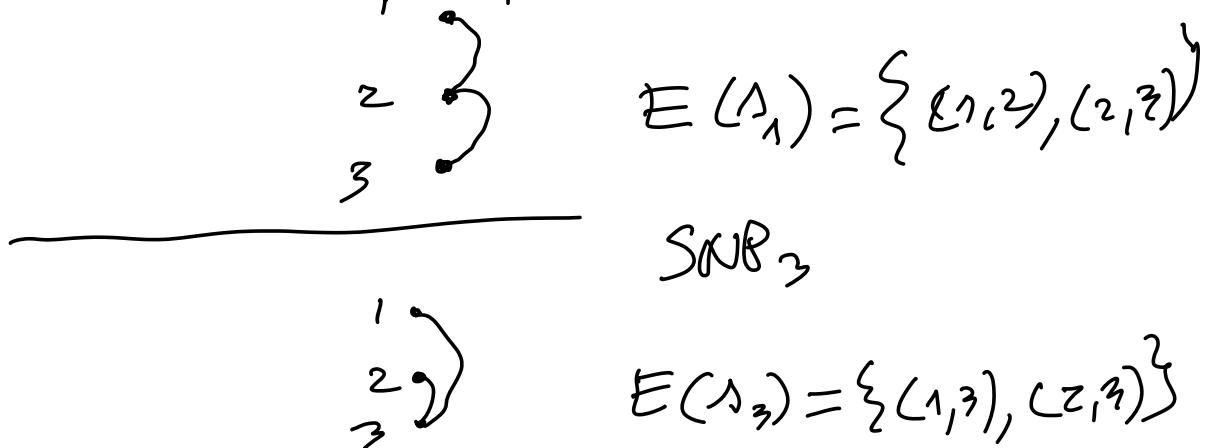
$$M[i, s] = M[j, s] \neq$$

How much information exists in one SNP about another SNP?
 Can i predict something about j ?

let $E(s) = \text{edge set of } s \text{ in } M$
 is the set of D-edges in the
 associated Column graph

In the above example :

Column graph of SNP₁



$E(s_1), \dots, E(s_8)$

The information in SNP is given
respect to SNP z

$$I(s, t) = \frac{|E(s) \cap E(t)|}{|E(t)|}$$

$|E(s)|$ = size of the edge set

$$E(\Delta_1) = \{(1,2), (2,3)\}$$

$$E(\Delta_2) = \emptyset$$

$$E(\Delta_3) = \{(1,3), (2,3)\}$$

$$E(\Delta_4) = \{(1,3), (2,3)\}$$

$$E(\Delta_5) = \emptyset$$

$$E(\Delta_6) = \{(1,2), (2,3)\}$$

$$E(\Delta_7) = \{(1,3), (2,3)\}$$

$$E(\Delta_8) = \{(1,3), (2,3)\}$$

$$I(\Delta_1, \Delta_8) = \frac{1}{2}$$

$$I(\Delta_1, \Delta_4) = \frac{1}{2}$$

Generalize to sets of SNPs

S, T two sets of SNPs

$$I(S, T) = \frac{|S \cap T|}{|T|} = \frac{\underbrace{\{(UE(s)) \cap (UE(t))\}}_{\substack{s \in S \\ t \in T}}}{\underbrace{|UE(s)|}_{s \in S}}$$

$S = \{s_1, s_3\}$

$$T = \{t_2, t_4, t_5, t_6, t_7, t_8\}$$

$$T(S, T) = \frac{3}{3} = 1$$

Look at $S = \{s_1, s_3\}$

Def the minimum informative

Subset of SNPs (more precisely)

the min set of SNPs of (MiS)
maximum information)

$S = \{s_1, s_3\}$ is the MiS for our example

Another mis is $S' = \{\Delta_4, \Delta_6\}$

Extensible to many SNPs not only pairwise.

A conservative extension of the pairwise case.

This was not a property of r^2 : Cannot be generalized conservatively to many SNPs.

"break the curse of the pairwise"

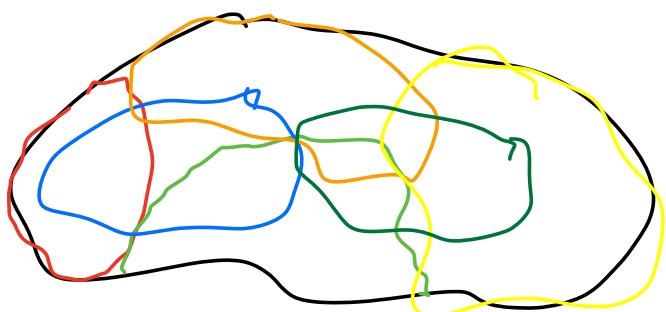
Information-theory based: 1 bit

The Set Cover Problem

Given: $U = \text{a universe (a set of elements)}$

$$R_i \subseteq U, \{R_i \mid i \in I\}$$

compute: A Set Cover is a collection of the R_i subsets whose union is the entire universe U .



An example:

$$U = \{a, b, c, d, e\}$$

$$R_1 = \{a, b, c\}, R_2 = \{a, b\}$$

$$R_3 = \{b, c, d\}, R_4 = \{c, d, e\}$$

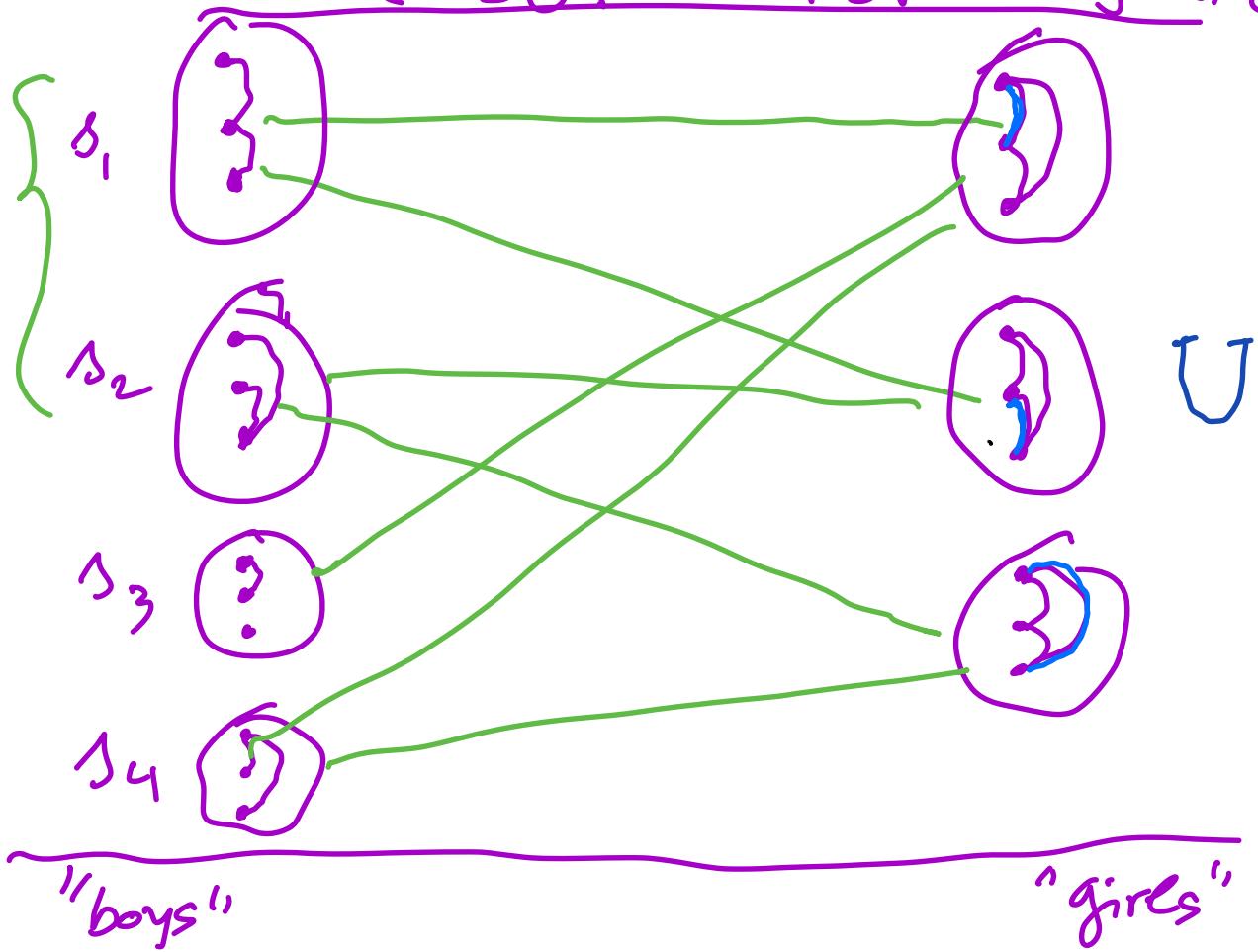
R₁, R₄ = a set cover

R₂, R₃ also a set cover

$\text{MiS} =$ is a set cover PB

i.e. the objective function

is the SET COVER objective



~~add a threshold~~
 $\text{MiS} \geq 80\%$

find the min no. of boys that know at least 80% of the girls

HW2 RB

LD-select with $\sigma_i^2 \geq 2$ e.g. 0.8
Informatics with $\sum \geq p$ e.g. 0.8

Find a correlation between the two :

In fact DOMINATING SET,
and SET COVER

are NP-complete \Rightarrow but

Reducible to each other:

- If you can solve one efficiently you can solve the other with same time complexity alg.
- The same for approximation algorithms.

Def dominating set

$G = (V, E)$ an undirected graph

V = vertex set

E = edge set

A subset $D \subseteq V$ is called a dominating set if every vertex outside D is connected by at least one edge to a vertex in D .

A reduction between

Dominating Set and Set Cover

① From Dominating SET to
 \nRightarrow SET COVER

Take a general graph $G = (V, E)$
with $V = \{1, 2, \dots, n\}$

We construct a SET COVER for G as follows.

$U = V$ universe of elements

consider a family of subsets of V defined as follows:

$$S = \{S_1, S_2, \dots, S_m\}$$

one for each vertex of G .

$$S_v = \{\text{the set of vertices adjacent to } v \text{ together with } v\}$$

Consider a dominating set of G .

let $C = \{S_v \mid v \in D\}$.

Want to show that C is a set cover. Indeed!

Moreover: $|C| = |D|$

Conversely: $W \subseteq V$

if $C = \{S_v : v \in W\}$ is
a set cover

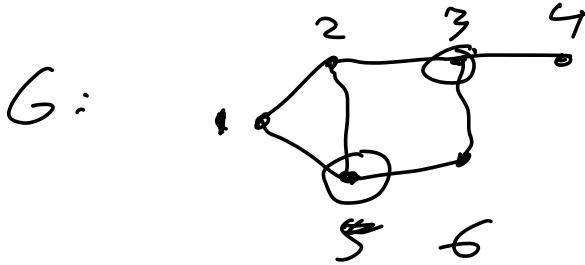
then want to show that
W is a dominating set for G.
Indeed!

One-to-one correspondence between
dominating sets in G and
set covers in (V, S) .

$\min - \max$

A very simple algorithm to construct
a set cover solution from
a dominating set in G.

An example : From Dominating Set to Set Cover



example of dominating set

$$D = \{3, 5\} \text{ min DS}$$

Construct the Set Cover Instance

$$S_1 = \{1, 2, 5\}, S_2 = \{1, 3, 5\}, S_3 = \{2, 4, 6\}$$
$$S_4 = \{1, 2, 6\}, S_5 = \{3, 5, 6\}$$

In G $D = \{3, 5\}$ is a dominating set

This corresponds to $C = \{S_3, S_5\}$
a set cover

(2) From SET COVER to
DOMINATING SET

A general Set Cover instance

$$(U, S)$$

$U = \text{universe}$

$S = \text{collection of}$
 $\text{subsets of } U$

$$S = \{S_i : i \in I\}$$

Assume U and I are disjoint

Construct a graph $G = (V, E)$
as follows :

- the set of vertices $V = U \cup I$
- the set of edges :

Type II $\{i, j\} \in E$

for every $i, j \in I$

Type II

$$\{i, u\} \subseteq E$$

for every $i \in I$
 $u \in S_i$

Note. I is a clique (completely connected graph)

Suppose that C is a set cover
solution to (U, S) pb.

$$C = \{S_i : i \in Z\}$$

want to show that Z is
a dominating set in G .

$$Z \subseteq I$$

then Z is a dominating set
for G .

- First for each $u \in U$ there

is $i \in Z$ such that $u \in S_i$

and by construction

u and i are adjacent
in G .

hence u is dominated by i

- Second, since Z must be nonempty,
each $i \in Z$ is adjacent to
a vertex in Z .

So Z is a dominating set.

Conversely: Let D be a dominating
set for G

Then it is possible to construct
another dominating set X

such that $|X| \leq |D|$

and $X \subseteq I$:

simply replace each $u \in D \cap V$
by a neighbor $i \in I$ of u .

Then $C = \{S_i : i \in X\}$ is a
set cover solution for (V, S)

and $|C| = |X| \leq |D|$.

Example: From set cover \Rightarrow Dominating set

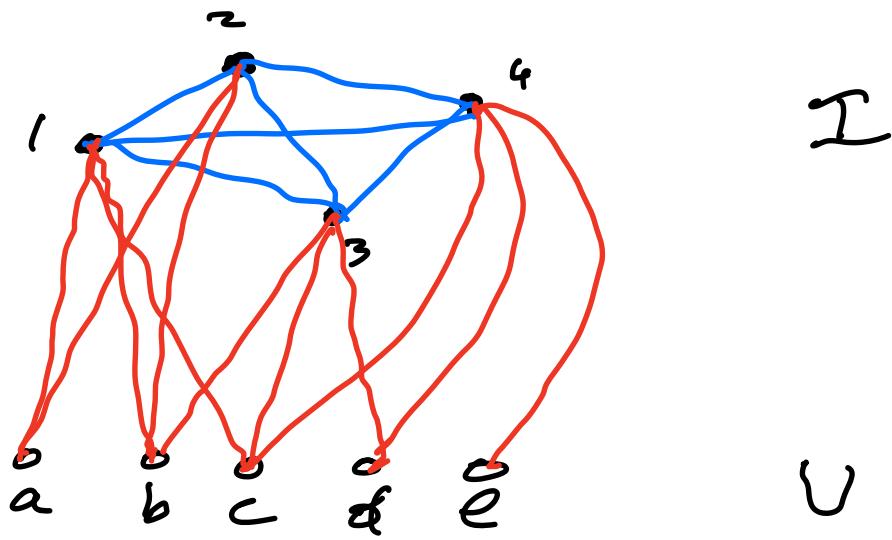
$$V = \{a, b, c, d, e\}$$

$$S_1 = \{a, b, c\}, S_2 = \{a, b\},$$

$$S_3 = \{b, c, d\}, S_4 = \{c, d, e\}$$

construct G : $V = V \cup I$

$$I = \{1, 2, 3, 4\}$$



In this example:

$C = \{S_1, S_4\}$ is a set cover
corresponding $D = \{1, 4\}$ it is a
dominating set

- $D = \{a, 3, 4\}$ is also a dominating set
 X constructed as above
 $X = \{1, 3, 4\}$ a dominating set
 $\subseteq I$.