## Identifying Genes and Interactions through a Forest Approach

### Heping Zhang

And Xiang Chen, Ching-Ti Liu, Meizhuo Zhang

Presented at BIRS workshop, Canada

YALE UNIVERSITY
School of Medicine

---

## Outline

Introduction to trees and forests

Simulation studies

**Genetic variants associated with age-related macular degeneration**

Yale University
School of Medicine

2

---

## Complex Traits

Diseases that do not follow Mendelian Inheritance Pattern

Genetic factors, Environment factors, G-G and G-E interactions

Interactions: effects that deviate from the additive effects of single effects

**Genetic variants have been identified for Age-related Macular Degeneration, Diabetes, Inflammatory Bowel Disorders, etc.**

Yale University
School of Medicine

3

---

## Data Structure

| 1  |  | * | ~ |
| 2  |  | * | ~ |
| … |  | … | … |
| 25 |  | * | ~ |
| 26 |  | * | ~ |
| … |  | … | … |
| 72 |  | * | ~ |

Yale University
School of Medicine

4

---

## Recursive Partitioning

A technique to identify heterogeneity in the data and fit a simple model (such as constant or linear) locally, and this avoids pre-specifying a systematic component.

Yale University
School of Medicine

5

---

## Leukemia Data

Source: http://www-genome.wi.mit.edu/cancer
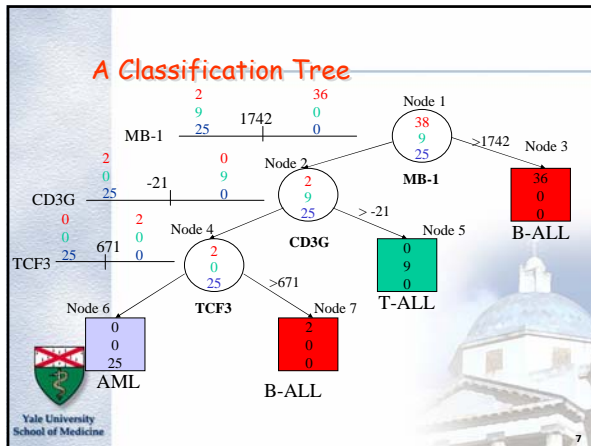
Contents:

• 25 mRNA - acute myeloid leukemia (AML)

• 38 - B-cell acute lymphoblastic leukemia (B-ALL)

• 9 - T-cell acute lymphoblastic leukemia (T-ALL)

• 7,129 genes

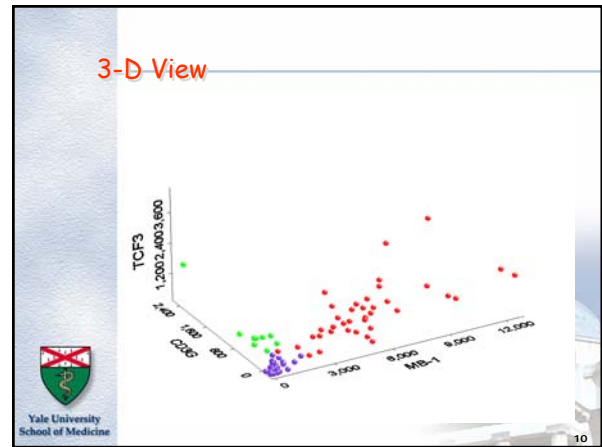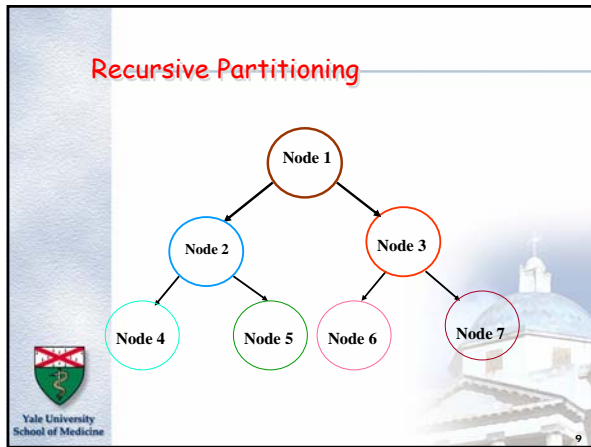Question: are the microarray data useful in classifying different types of leukemia?

Yale University
School of Medicine

6

---

**Slide 7: A Classification Tree**



A Classification Tree

Node 1: 38 / 9 / 25, MB-1, >1742 → Node 3: 36 / 0 / 0, B-ALL

Node 2: 2 / 9 / 25, CD3G, > -21 → Node 5: 0 / 9 / 0, T-ALL

Node 4: 2 / 0 / 25, TCF3, >671 → Node 7: 2 / 0 / 0, B-ALL

Node 6: 0 / 0 / 25, AML

MB-1: 2 / 9 / 25 — 1742 — 36 / 0 / 0
CD3G: 2 / 0 / 25 — -21 — 0 / 9 / 0
TCF3: 0 / 0 / 25 — 671 — 2 / 0 / 0

**Slide 8: Node Splitting**

Node Splitting

Click to see the diagram

**Slide 9: Recursive Partitioning**

Recursive Partitioning



Node 1 → Node 2, Node 3
Node 2 → Node 4, Node 5
Node 3 → Node 6, Node 7

**Slide 10: 3-D View**

3-D View



**Slide 11: Genomics**

Genomics

Adequate sample size for parameters of interest. Often, we have hundreds or thousands of observations for the inference on a few parameters. We can try to settle an "optimal" model.

In this information age, we have more and more variables but the access to the number of study subjects remains the same. One model can no longer provide an adequate summary of the information.
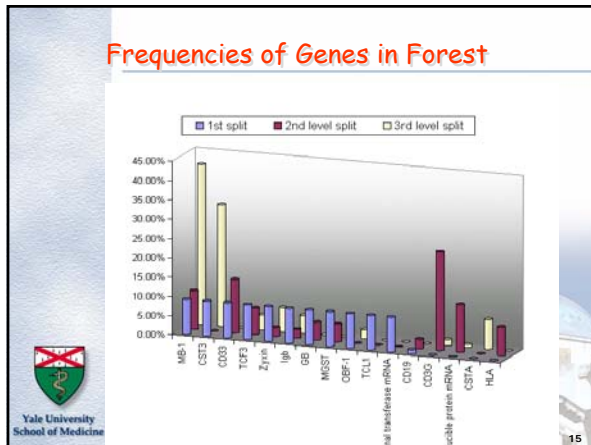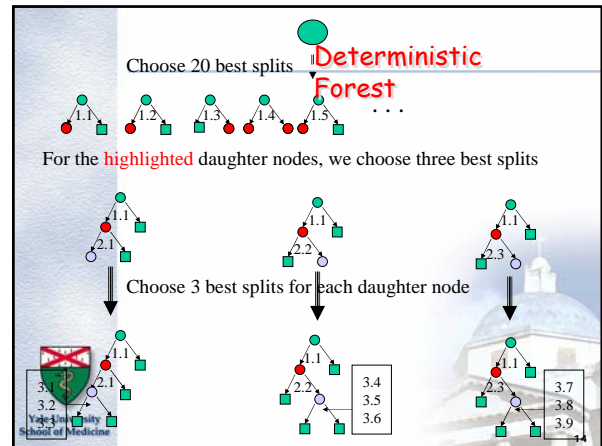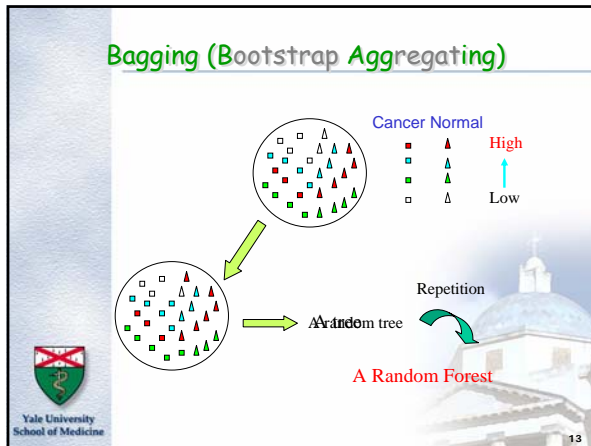
**Slide 12: Forests**

Forests

To identify a constellation of models that collectively help us understand the data.

For example, we can select and rank the genes whose expressions show a great promise of classifying tumor cells.
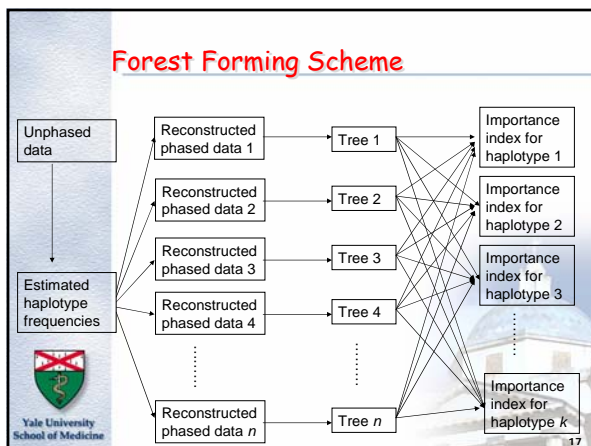
## Bagging (Bootstrap Aggregating)

Cancer Normal

High

Low

Repetition

A Random tree

A Random Forest

Yale University
School of Medicine

13

Choose 20 best splits

**Deterministic Forest**

1.1 1.2 1.3 1.4 1.5 . . .

For the highlighted daughter nodes, we choose three best splits

1.1 2.1   1.1 2.2   1.1 2.3

Choose 3 best splits for each daughter node

1.1 2.1   1.1 2.2   1.1 2.3

3.1
3.2   3.4
3.5
3.6   3.7
3.8
3.9

Yale University
School of Medicine

14

## Frequencies of Genes in Forest

■ 1st split   ■ 2nd level split   □ 3rd level split

45.00%
40.00%
35.00%
30.00%
25.00%
20.00%
15.00%
10.00%
5.00%
0.00%

MB-1 CST3 CD03 TGF3 Zyxin Igб GB MGST OBF-1 TCL1 CD19 CD3G nal transferase mRNA ucible protein mRNA CSTA HLA

Yale University
School of Medicine

15

## SNPs vs. Haplotypes

SNPs

✓ **Directly observed**
✓ **No uncertainty**
✗ **Less informative**
❖ **Tree approaches**

Haplotypes

✗ **Inferred from SNPs**
✗ **Uncertain**
✓ **More informative**
❖ **Forest approaches**

Yale University
School of Medicine

16

## Forest Forming Scheme

Unphased data

Reconstructed phased data 1 → Tree 1 → Importance index for haplotype 1

Reconstructed phased data 2 → Tree 2 → Importance index for haplotype 2

Reconstructed phased data 3 → Tree 3 → Importance index for haplotype 3

Reconstructed phased data 4 → Tree 4

Estimated haplotype frequencies

Reconstructed phased data $n$ → Tree $n$ → Importance index for haplotype $k$

Yale University
School of Medicine

17

## Haplotype Frequency Estimation

Existing haplotype frequency estimation software that output a set of haplotype pairs with corresponding frequencies for each subject in each region.
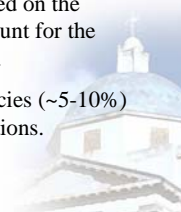
We used SNPHAP (Clayton 2006)

Yale University
School of Medicine

18

## Unphased to Phased Data

One unphased data expands to a large number of phased datasets.

In each region, an individual's haplotype pair is randomly selected based on the estimated frequencies to account for the uncertainty of the haplotypes.

Haplotypes with low frequencies (~5-10%) should have some representations.

Yale University
School of Medicine

19

## Trees Based on Phased Data

A tree is grown for each phased data set.

A random forest is formed for all phased data sets.

Yale University
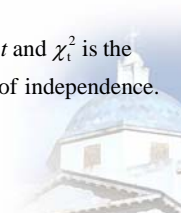School of Medicine

20

## Inference from the Forest

Importance of haplotype $h$ in tree $T$

$$V_h = \sum_{t \in T, t \text{ is split by } h} 2^{-L_t} \chi_t^2,$$

where $L_t$ is the depth of node $t$ and $\chi_t^2$ is the value of the $\chi^2$ - test statistic of independence.
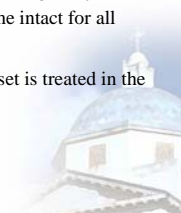
Yale University
School of Medicine

21

## Significance Level

Distribution of the maximum haplotype importance under null hypothesis is determined by permutation.

First, disease status is permuted among study subjects while keeping the genome intact for all individuals.

Then, each of the permuted data set is treated in the same way as the original data.

Yale University
School of Medicine

22

## Simulation Studies (2 loci)

- **300 cases and 300 controls**
- **Each region has 3 SNPs**
- **12 interaction models from Knapp *et. al.* (1994) and Becker *et. al.* (2005)**
- **2 additive models with background penetrance**
- **3 scenarios**
    - Neither region is in LD with the disease allele
    - One of the regions is in LD (D' = 0.5) with the disease allele
    - Both regions are in LD (D' = 0.5) with the disease allele

Yale University
School of Medicine

23

## Simulation Studies (2 loci)

Penetrance

|  |  | Region 2 | | |
|---|---|---|---|---|
|  |  | 0 | 1 | 2 |
| Region 1 | 0 | $f_{00}$ | $f_{01}$ | $f_{02}$ |
|  | 1 | $f_{10}$ | $f_{11}$ | $f_{12}$ |
|  | 2 | $f_{20}$ | $f_{21}$ | $f_{22}$ |

Yale University
School of Medicine

24

## Simulation Studies (2 loci)

| Model | $f_{22}$ | $f_{21}$ | $f_{20}$ | $f_{12}$ | $f_{11}$ | $f_{10}$ | $f_{02}$ | $f_{01}$ | $f_{00}$ | $f$ | $p_1$ | $p_2$ |
|---|---|---|---|---|---|---|---|---|---|---|---|---|
| Ep-1 | $f$ | $f$ | 0 | $f$ | $f$ | 0 | 0 | 0 | 0 | 0.707 | 0.210 | 0.210 |
| Ep-2 | $f$ | $f$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.778 | 0.600 | 0.199 |
| Ep-3 | $f$ | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0 | 0.900 | 0.577 | 0.577 |
| Ep-4 | $f$ | $f$ | 0 | $f$ | 0 | 0 | $f$ | 0 | 0 | 0.911 | 0.372 | 0.243 |
| Ep-5 | $f$ | $f$ | 0 | $f$ | 0 | 0 | 0 | 0 | 0 | 0.799 | 0.349 | 0.349 |
| Ep-6 | 0 | $f$ | $f$ | $f$ | 0 | 0 | $f$ | 0 | 0 | 1.000 | 0.190 | 0.190 |
| Het-1 | $g$ | $g$ | $f$ | $g$ | $g$ | $f$ | $f$ | $f$ | 0 | 0.495 | 0.053 | 0.053 |
| Het-2 | $g$ | $g$ | $f$ | $f$ | $f$ | 0 | $f$ | $f$ | 0 | 0.660 | 0.279 | 0.040 |
| Het-3 | $g$ | $f$ | $f$ | $f$ | 0 | 0 | $f$ | 0 | 0 | 1.000 | 0.194 | 0.194 |
| S-1 | $f$ | $f$ | $f$ | $f$ | $f$ | $f$ | $f$ | $f$ | 0 | 0.522 | 0.052 | 0.052 |
| S-2 | 1 | 1 | 1 | $f$ | $f$ | 0 | $f$ | $f$ | 0 | 0.574 | 0.228 | 0.045 |
| S-3 | 1 | 1 | $f$ | 1 | $f$ | 0 | $f$ | 0 | 0 | 0.512 | 0.194 | 0.194 |
| Ad-1 | $f$ | $f$ | 0.04 | $f$ | 0.304 | 0.02 | 0.01 | 0.01 | 0.01 | 0.799 | 0.349 | 0.349 |
| Ad-2 | $f$ | $f$ | 0.15 | $f$ | 0.324 | 0.10 | 0.05 | 0.05 | 0.05 | 0.799 | 0.349 | 0.349 |

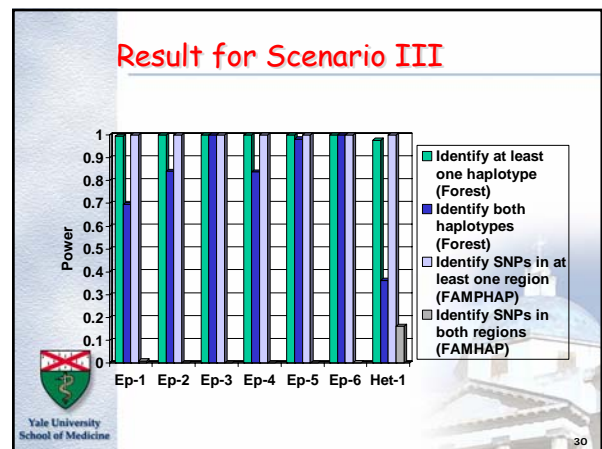$$g = 2f - f^2$$

25

## Simulation Studies (2 loci)
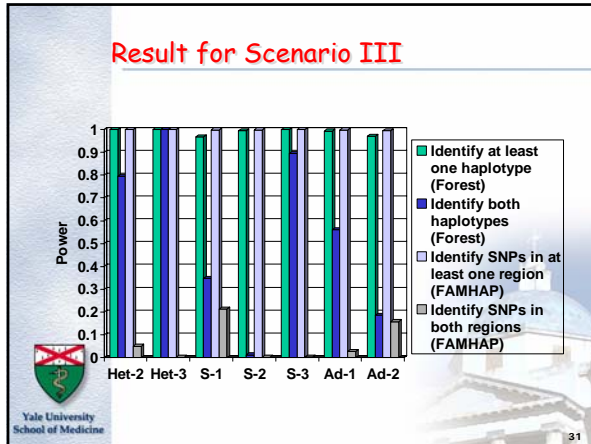
Benckmark:

**FAMHAP software from Becker et. al. (2005)**

26

## Result for Scenario I

False positive rate:

Our method: $< 1\%$

FAMHAP: $> 5\%$

27

## Result for Scenario II



Legend:
- Identify the correct haplotype (Forest)
- Identify an incorrect haplotype (Forest)
- Identify SNPs in the correct region (FAMHAP)
- Identify SNPs in the neutral region (FAMHAP)

Categories: Ep-1 Ep-2 Ep-3 Ep-4 Ep-5 Ep-6 Het-1

28

## Result for Scenario II



Legend:
- Identify the correct haplotype (Forest)
- Identify an incorrect haplotype (Forest)
- Identify SNPs in the correct region (FAMHAP)
- Identify SNPs in the neutral region (FAMHAP)

Categories: Het-2 Het-3 S-1 S-2 S-3 Ad-1 Ad-2

29

## Result for Scenario III



Legend:
- Identify at least one haplotype (Forest)
- Identify both haplotypes (Forest)
- Identify SNPs in at least one region (FAMHAP)
- Identify SNPs in both regions (FAMHAP)

Categories: Ep-1 Ep-2 Ep-3 Ep-4 Ep-5 Ep-6 Het-1

30

## Result for Scenario III



- Identify at least one haplotype (Forest)
- Identify both haplotypes (Forest)
- Identify SNPs in at least one region (FAMHAP)
- Identify SNPs in both regions (FAMHAP)

31

## Real Case Study

Age-related macular degeneration (AMD)
- **Leading cause of vision loss in elderly**
- **Affects more than 1.75 million individuals in the United States**
- **Projected to about 3 million by 2020**

Klein *et al.* (2005)
- **Case-control (96 AMD cases, 50 controls)**
- **~100,000 SNPs for each individual**
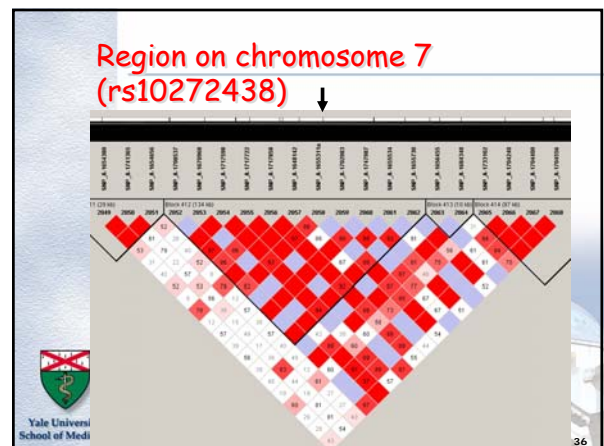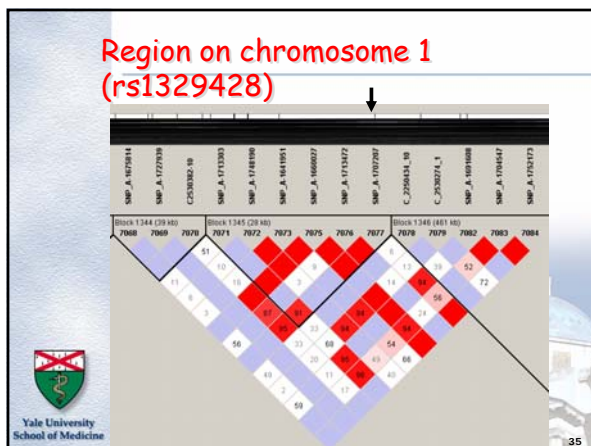- **CFH gene identified**

32

## Analysis Procedure

RTREE program
- **Each SNP is used as one covariate**
- **Two SNPs identified as potentially associated with AMD (rs1329428 on chromosome 1 and rs10272438 on chromosome 7)**

Hapview program: LD block construction
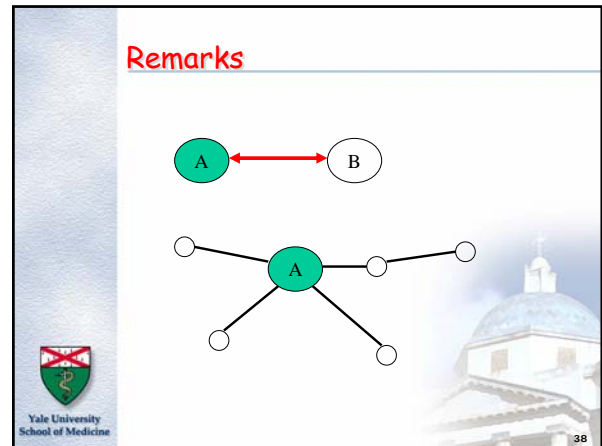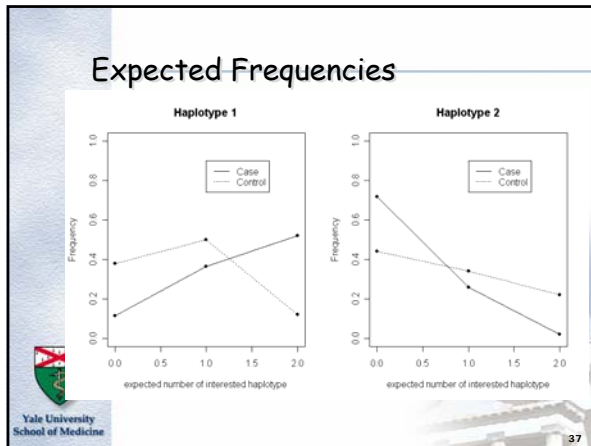- **6-SNP block for rs1329428**
- **11-SNP block for rs10272438**

Forest

33

## Result

Two haplotypes are identified
- **Most significant: ACTCCG in region 1 (p-value = 2e-6)**
  - Identical to Klein *et. al.* (2005)
  - Located in CFH gene
- **Another significant haplotype: TCTGGACGACA, in region 2 (p-value = 0.0024)**
  - Not reported before
  - Protective
  - Located in BBS9 gene

34

## Region on chromosome 1 (rs1329428)



35

## Region on chromosome 7 (rs10272438)



36

**Slide 37**

## Expected Frequencies



Yale University
School of Medicine

37

**Slide 38**

## Remarks



A — B

A

Yale University
School of Medicine

38

**Slide 39**

## Books

SEARCH INSIDE!™

CLASSIFICATION
AND
REGRESSION
TREES

Breiman Friedman
Olshen Stone

LOOK INSIDE!™

Statistics for Biology
and Health

Heping Zhang
Burton Singer

Recursive
Partitioning in
the Health
Sciences

Springer

Zhang HP and Singer B. *Recursive Partitioning in the Health Sciences*. Springer, 1999.

Yale University
School of Medicine

39

**Slide 40**

## Trees in Genetic Studies

**Zhang and Bonney (2000)**
**Nelson et al. (2001)**
**Bastone et al. (2004)**
**Cook, Zee and Ridker (2004)**
**Foulkes, De Gruttola and Hertogs (2004)**

Yale University
School of Medicine

40

**Slide 41**

## References on Forests

Breiman L. *Bagging predictors*. Machine Learning, 24(2):123-140, 1996.

Zhang HP. *Classification trees for multiple binary responses*. Journal of the American Statistical Association, 93: 180-193, 1998.

Zhang HP et al. *Cell and Tumor Classification using Gene Expression Data: Construction of Forests*. Proceedings of the National Academy of Sciences USA, 100: 4168-4172, 2003.

Yale University
School of Medicine

41