# SIGNIFICANT VS PREDICTIVE
## VARIABLE SELECTION AND PREDICTION

Mengying Sun

Apr. 20th 2016

# Outline

- GWAS

- Missing Heritability

- Significant vs. Predictive Variables

- 3 Examples

- I Score

- Conclusions and comments

# GWAS

- **G**enome-**W**ide **A**ssociation **S**tudy

An examination of common genetic variants in different individuals to see if any variant is associated with a trait.

$$Trait\ Y \qquad SNPs\ [X_1, \dots X_j\ \dots X_m]$$

$$
\begin{matrix}
y_1 \\
y_2 \\
\vdots \\
y_n
\end{matrix}
\qquad
\begin{matrix}
x_{11} & \dots & x_{j1} & \dots & x_{m1} \\
\vdots & \dots & x_{jk} & \dots & \vdots \\
x_{1n} & \dots & x_{jn} & \dots & x_{mn}
\end{matrix}
\qquad
\begin{matrix}
j = 1,2, \dots m \\
n = 5k \\
m = 500k, 1000k
\end{matrix}
$$

$$n \times 1 \qquad\qquad n \times m$$

# GWAS

- **G**enome-**W**ide **A**ssociation **S**tudy

An examination of common genetic variants in different individuals to see if any variant is associated with a trait.

$$Y \quad \sim \quad X_j$$

$$\begin{matrix} y_1 \\ y_2 \\ \vdots \\ y_n \end{matrix} \qquad \begin{matrix} x_{11} & \dots & x_{j1} & \dots & x_{m1} \\ \vdots & \dots & x_{jk} & \dots & \vdots \\ x_{1n} & \dots & x_{jn} & \dots & x_{mn} \end{matrix}$$
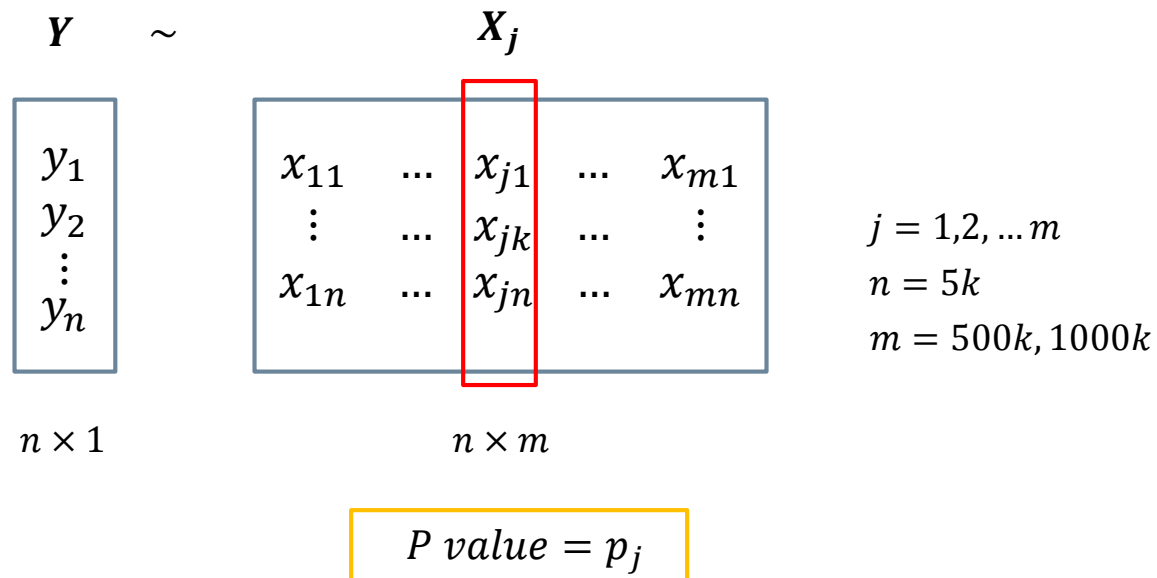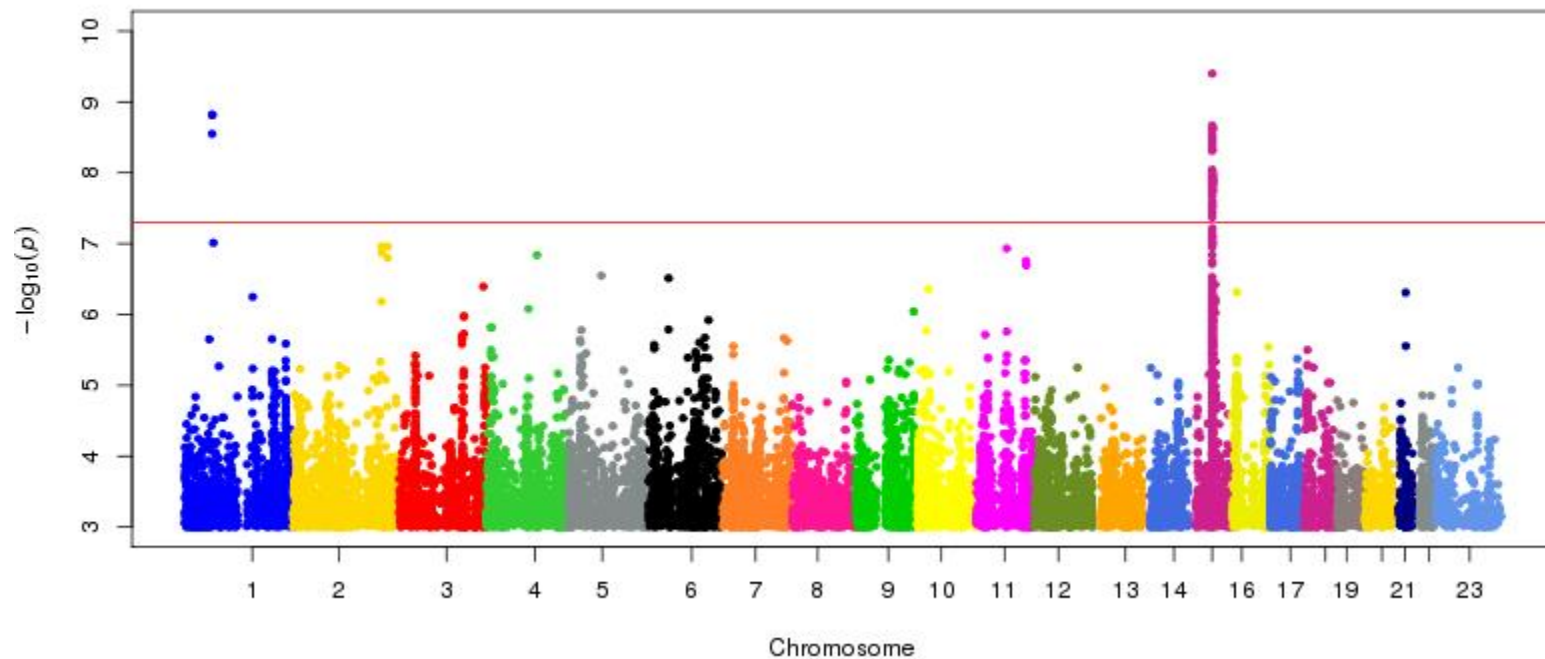
$$j = 1, 2, \dots m$$
$$n = 5k$$
$$m = 500k, 1000k$$

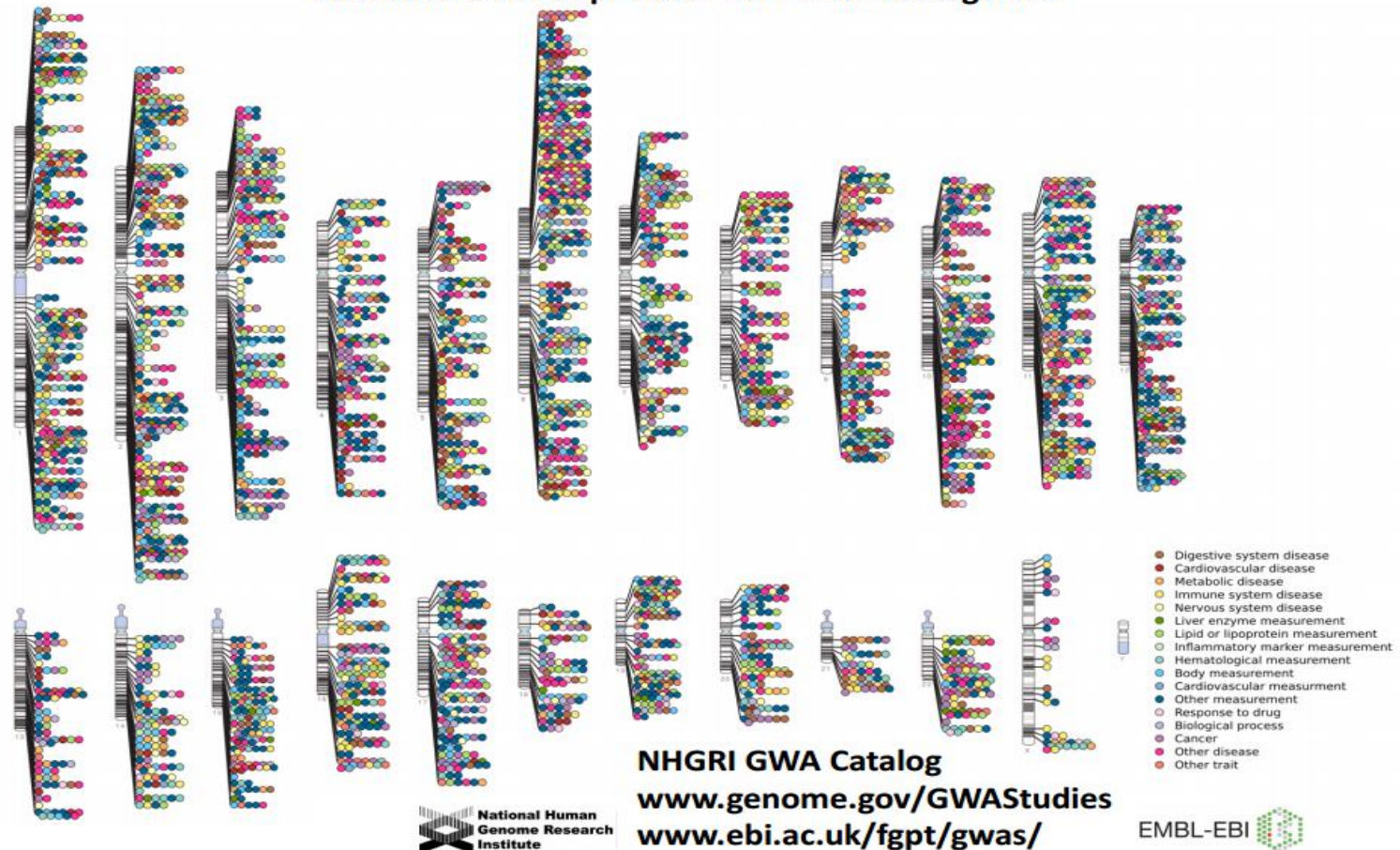$$n \times 1 \qquad\qquad n \times m$$

$$P\ value = p_j$$

# Manhattan Plot



Styrkarsdottir, Unnur, Gudmar Thorleifsson, et al. 2014. "Severe Osteoarthritis of the Hand Associates with Common Variants within the ALDH1A2 Gene and with Rare Variants at 1p31." *Nature Genetics* 46 (5): 498–502. doi:10.1038/ng.2957.

# Published GWA



**Published Genome-Wide Associations through 12/2013**
**Published GWA at p≤5X10⁻⁸ for 17 trait categories**

NHGRI GWA Catalog
www.genome.gov/GWAStudies
www.ebi.ac.uk/fgpt/gwas/

National Human Genome Research Institute

EMBL-EBI

Legend:
- Digestive system disease
- Cardiovascular disease
- Metabolic disease
- Immune system disease
- Nervous system disease
- Liver enzyme measurement
- Lipid or lipoprotein measurement
- Inflammatory marker measurement
- Hematological measurement
- Body measurement
- Cardiovascular measurment
- Other measurement
- Response to drug
- Biological process
- Cancer
- Other disease
- Other trait

Hindorff LA, MacArthur J, Morales J, Junkins HA, Hall PN, Klemm AK, and Manolio TA. A Catalog of Published Genome-Wide Association Studies.
www.genome.gov/gwastudies

# Validation

$Y - trait$

$X - significant\ variable\ set$

$$Y = X\beta + \epsilon$$

$$\hat{\beta} \implies \hat{y}$$

$$Y \sim \hat{y}$$

$r^2 - coefficient\ of\ determination$

$cor(y, \hat{y}) - correlation\ between\ y\ \&\ \hat{y}$

**Problem !!!**

# Missing Heritability

- Height is 80 – 90% heritable

- In GWAS, 40 genetic variants turned up associated with the height difference

- However, the variants accounted for only 5% of height's heritability

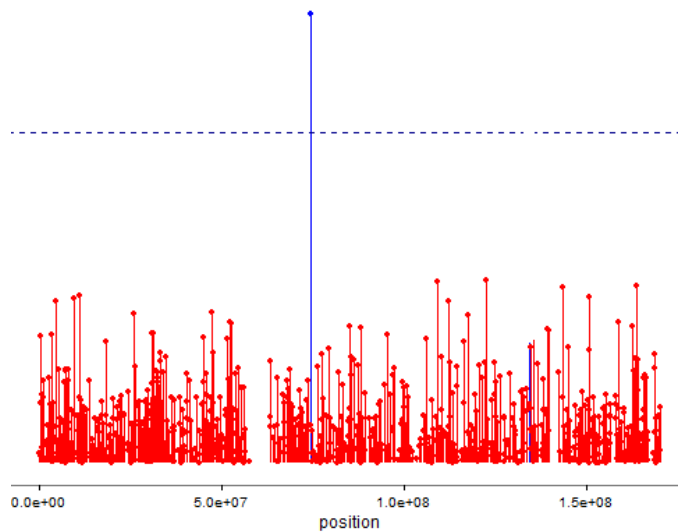| Disease | Number of loci | Proportion of heritability explained | Heritability measure |
|---|---|---|---|
| Age-related macular degeneration[72] | 5 | 50% | Sibling recurrence risk |
| Crohn's disease[21] | 32 | 20% | Genetic risk (liability) |
| Systemic lupus erythematosus[73] | 6 | 15% | Sibling recurrence risk |
| Type 2 diabetes[74] | 18 | 6% | Sibling recurrence risk |
| HDL cholesterol[75] | 7 | 5.2% | Residual[*] phenotypic variance |
| Height[15] | 40 | 5% | Phenotypic variance |
| Early onset myocardial infarction[76] | 9 | 2.8% | Phenotypic variance |
| Fasting glucose[77] | 4 | 1.5% | Phenotypic variance |

*Residual is after adjustment for age, gender, diabetes.

Manolio, Teri A., Francis S. Collins, Nancy J. Cox, David B. Goldstein, Lucia A. Hindorff, David J. Hunter, Mark I. McCarthy, et al. 2009. "Finding the Missing Heritability of Complex Diseases." *Nature* 461 (7265): 747–53. doi:10.1038/nature08494.
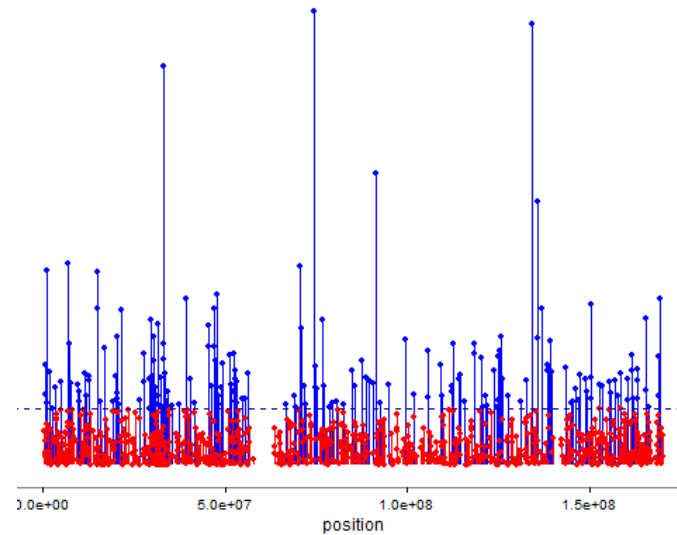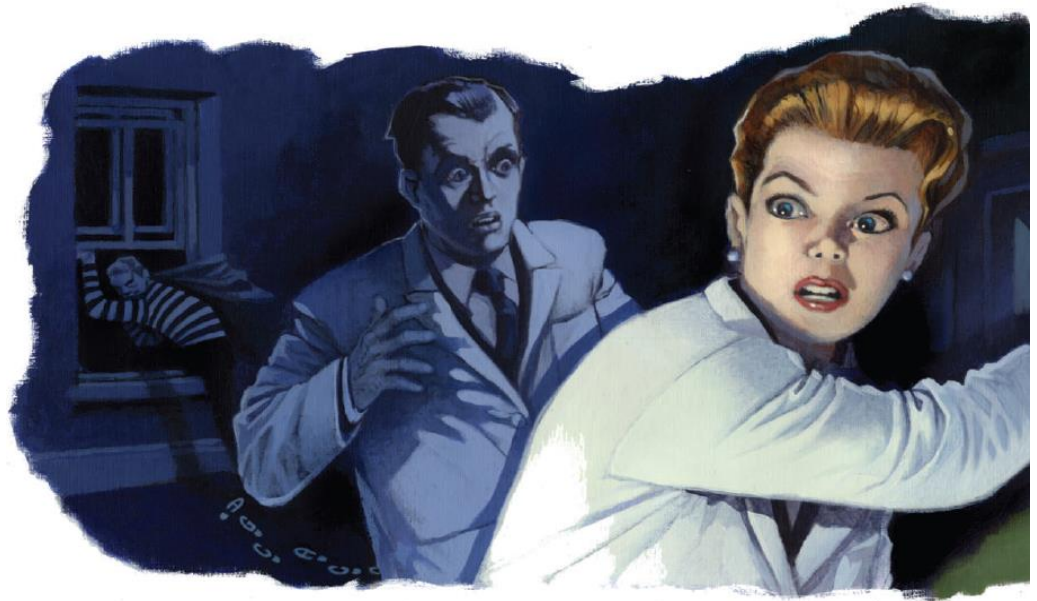
# Missing Heritability

- Single locus

- Complex traits

# Missing Heritability

- The way we select markers is WRONG ?

Maher, Brendan. 2008. "Personal Genomes: The Case of the Missing Heritability." *Nature News* 456 (7218): 18–21. doi:10.1038/456018a

# Paper

- Why significant variables aren't automatically good predictors

## Why significant variables aren't automatically good predictors

Adeline Lo[a], Herman Chernoff[b,1], Tian Zheng[c], and Shaw-Hwa Lo[c,1]

[a]Department of Political Science, University of California, San Diego, La Jolla, CA 92093; [b]Department of Statistics, Harvard University, Cambridge, MA 02138; and [c]Department of Statistics, Columbia University, New York, NY 10027

Thus far, genome-wide association studies (GWAS) have been disappointing in the inability of investigators to use the results of identified, statistically significant variants in complex diseases to

From the scientist's point of view there are two basic problems, complicated by the large size of the data set. These are variable selection and prediction. For variable selection, we wish to find a

# Significant vs. Predictive variables

- **Key Difference**

  **Different properties of their underlying distributions**

  $variable\ set\colon X$

  $f_0$: The distribution of X among control

  $f_1$: The distribution of X among cases

- **Variable selection (Significance)**

  $H_0\colon f_0 = f_1$   vs   $H_1\colon f_0 \neq f_1$

- **Prediction**

  $x, observed\ value\ of\ X\ for\ a\ single\ individual, belongs\ to\ f_0\ or\ f_1$

# Significant vs. Predictive variables

- **Variable selection (Significance)**

$H_0: f_0 = f_1$ vs $H_1: f_0 \neq f_1$

$$T_n$$

$$t_n = T_n(x)$$

$$P(T_n \geq t_n \mid H_0)$$

Reject $H_0$ if $P(T_n \geq t_n \mid H_0)$ is sufficiently small

# Significant vs. Predictive variables

- **Prediction**

$x$

$if\ f_0(x) < f_1(x), x\ is\ more\ likely\ from\ f_1$

$if\ f_1(x) < f_0(x), x\ is\ more\ likely\ from\ f_0$

Likelihood rate

$$prediction\ rate = 0.5 \sum_x \max(f_0(x), f_1(x))$$

# Significant vs. Predictive variables

- **Variable selection (Significance)**

$$H_0: f_0 = f_1 \quad \text{vs} \quad H_1: f_0 \neq f_1$$

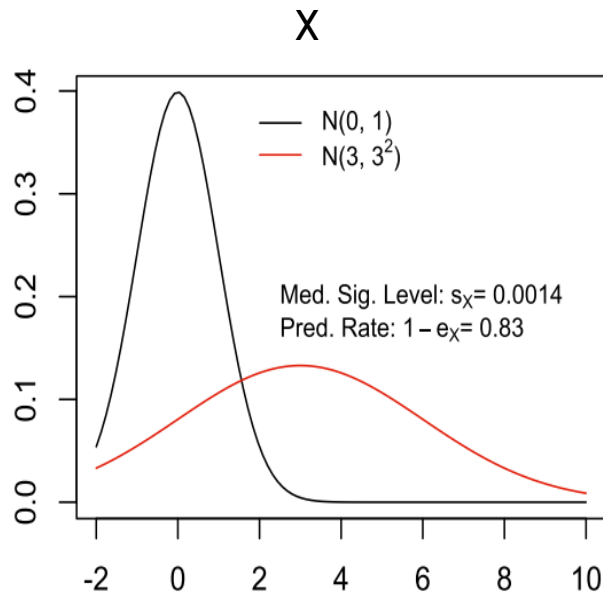Reject $H_0$ if $P(T_n \geq t_n \mid H_0)$ is sufficiently small

- **Prediction**

$$prediction\ rate = 0.5 \sum_x \max(f_0(x), f_1(x))$$

The former uses assumptions on, but no knowledge of, the EXACT distributions of the variables, whereas the latter, requires knowledge of both $f_0$ and $f_1$, which is the distribution of X

# Example 1

**Single Observation**

X



N(0, 1)
N(3, 3²)

Med. Sig. Level: $s_X = 0.0014$
Pred. Rate: $1 - e_X = 0.83$

$H_0: X \sim N(0, 1)$

$H_1: X \sim N(3, 3^2)$

- Classifying the state of an individual (H/D) who yields the the observation X

$e_x$

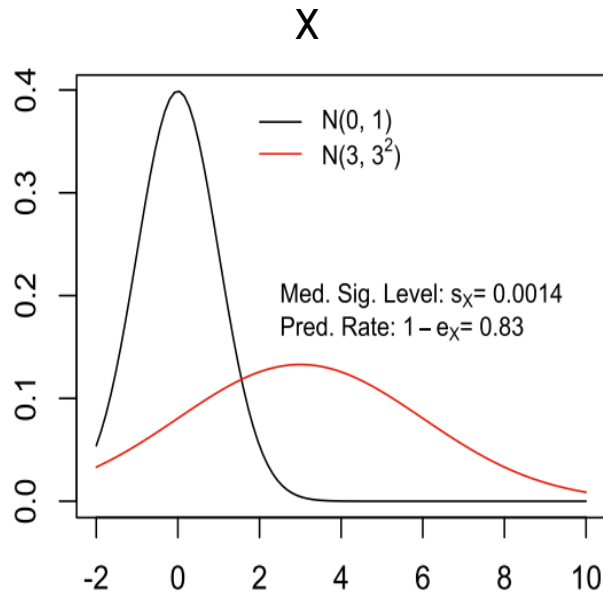|   | T | F |
|---|---|---|
| T | √ | e (c, H0) |
| F | e (c, H1) | √ |

- Reject H0 when likelihood ratio is large
- X is random, for every ratio c, 2*2
- choose c minimize the average of e(c,H0)+e(c,H1)

# Example 1

- **Single Observation**

X



Med. Sig. Level: $s_X = 0.0014$
Pred. Rate: $1 - e_X = 0.83$

Legend: N(0, 1), N(3, $3^2$)

$H_0: X \sim N(0, 1)$

$H_1: X \sim N(3, 3^2)$
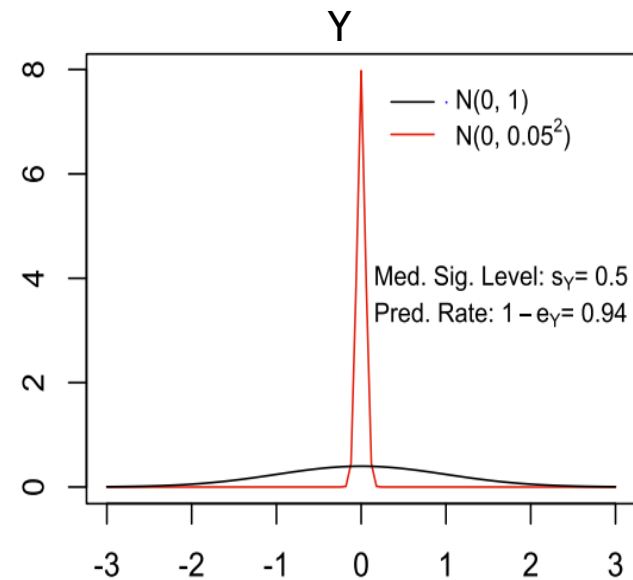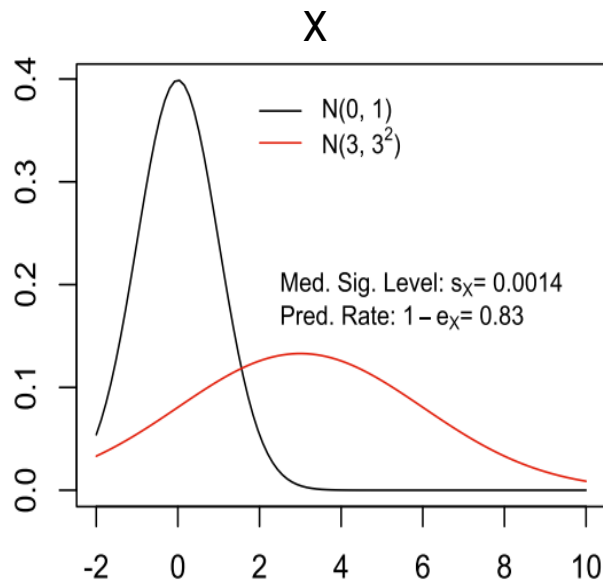
- Classifying the state of an individual (H/D) who yields the the observation X

$s_x$

- Large value of x will favor to reject H0
- For each observed x, P-value, choose the "median" value as significance level

# Example 1

- **Single Observation**



- Scientist's decision: whether observe X or Y
- Prefer X based on significance
- Prefer Y based on predictivity

# Example 1

- **Conclusion**

  - It happens that a variable that is more significant serves not as good as another less significant variable in prediction, based on the different distribution of those two variables.

- **Some doubts about Example 1**

  - Was expecting more extreme case (0.83 vs 0.94)
  - Choice of measures of significance and predictability (ex and sx) not commonly used which is biggest problem throughout the paper
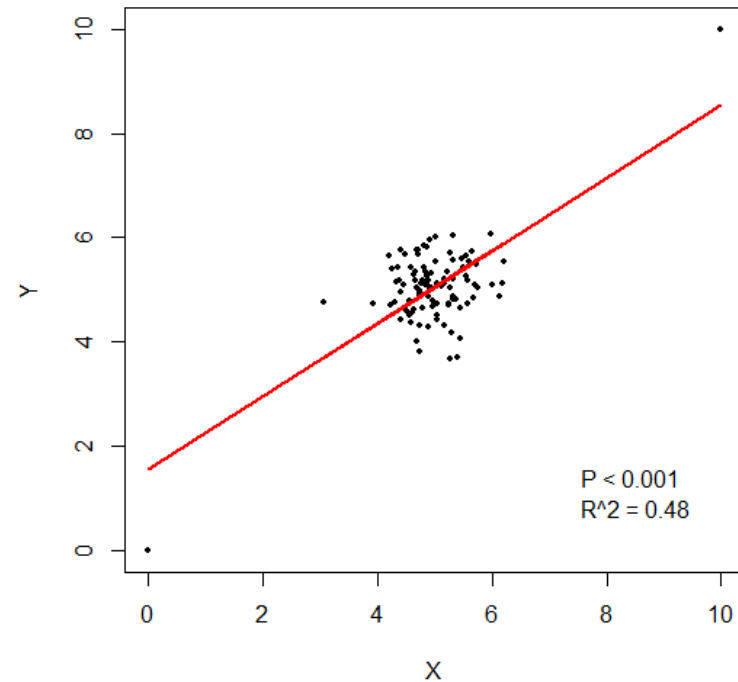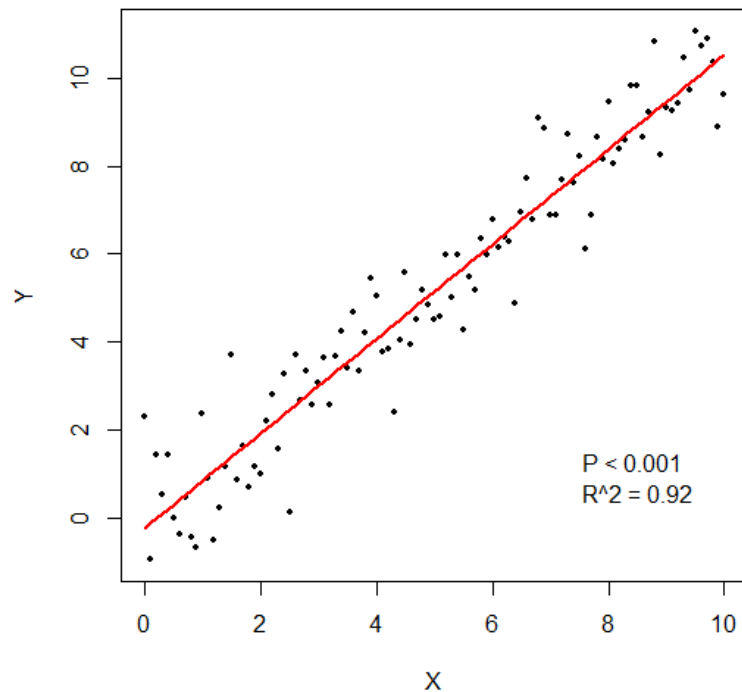  - Confusing and absurd

# What we can derive from the concept

- **Does distribution of X really affects the discrepancy between significance and predictivity?**

# What we can derive from the concept

- **An Extreme Case**  $y = x + \varepsilon, \qquad \varepsilon \sim N\,(0,1)$
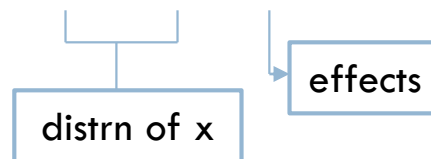
# What we can derive from the concept

- **Significance**

$$E(y \,|Xi) \sim N(X_i\beta, \sigma^2)$$

- **Prediction**

| X | 0 | 1 | 2 |
|---|---|---|---|
| MAF=p | $(1-p)^2$ | $2p(1-p)$ | $p^2$ |
| $EX = 2p$ | | $Var(X) = 2p(1-p)$ | |

$$Var(\hat{y}) = Var(X\beta + \varepsilon)$$

$$= \beta^2 Var(X)$$

$$= 2p(1-p)\beta^2$$

distrn of x

effects

# Example 2

- **1 SNP**

  - Outcome variable is the case or control

  - Explanatory variable is the reading on 1 SNP for each of 500 cases and 500 controls

  - Significance level using $X^2$ test

  - Prediction rate by calculating the proportion of correct classification among the 1000 individuals
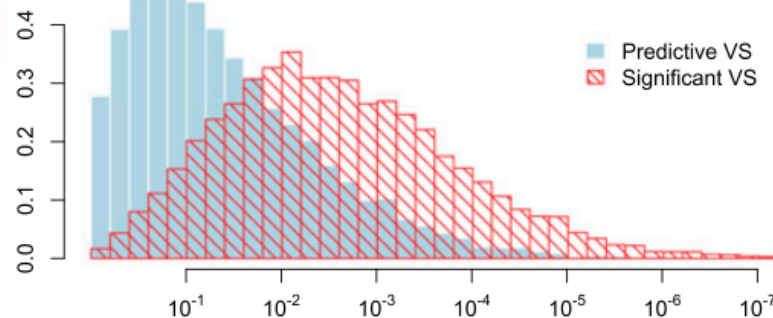
# Example 2

- **1 SNP**



Predictive variable sets | Significant variable sets

| genotype | 0 | 1 | 2 |
|----------|------|------|------|
| Cases | .274 | .500 | .226 |
| Controls | .226 | .500 | .274 |

Pred. rate = **0.524** Median P-value = $6 \times 10^{-2}$

Median I score = **1.544**

Distribution of P-values from Chi-squared test. Plotted on $-\log_{10}$ scale.

1 SNP

| genotype | 0 | 1 | 2 |
|----------|------|------|------|
| Cases | 0.01 | 0.19 | 0.80 |
| Controls | 0.04 | 0.16 | 0.80 |

Pred. rate = 0.515 **Median P-value = 3.5x10⁻³**

Median I score = 0.632

- Predictive VS
- Significant VS

Example 2

# Example 2

- **2 SNPs**



Predictive variable sets / Significant variable sets

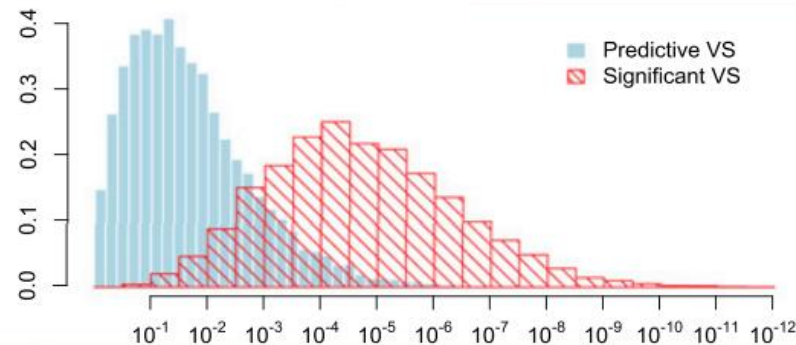| genotype | 00 | 01 | 02 | 10 | 11 | 12 | 20 | 21 | 22 |
|----------|------|------|------|------|------|------|------|------|------|
| Cases | .1 | .125 | .1 | .125 | .1 | .125 | .1 | .125 | .1 |
| Controls | .125 | .1 | .125 | .1 | .125 | .1 | .125 | .1 | .125 |

Pred. rate = **0.56**   Median P-value=3x10$^{-2}$

2 SNPs

| genotype | 00 | 01 | 02 | 10 | 11 | 12 | 20 | 21 | 22 |
|----------|------|------|------|------|------|------|------|------|------|
| Cases | .01 | .052 | .134 | .134 | .134 | .134 | .134 | .134 | .134 |
| Controls | .052 | .01 | .134 | .134 | .134 | .134 | .134 | .134 | .134 |

Pred. rate = 0.52  **Median P-value=2.8x10$^{-5}$**

Median I score = **1.854**

Distribution of P-values from Chi-squared test. Plotted on $-\log_{10}$ scale.

Median I score = 1.676

- Predictive VS
- Significant VS

Example 3

# Example 2

## Conclusion

- There is overlap between two variables set which means that to some extent significant variable serves good prediction ability and those highly predictive variable are also tend to be significant. However that is NOT necessary;
- In fact, large proportion of the predictive variables are not significant, vice versa.
- As the number of SNPs increase (complexity increase), the overlap become smaller and smaller
- But all results above based on their way of calculating the so-called "prediction rate" as measurement of prediction ability, which is a shortcoming of the paper.

# I score

- I score

$$I = \sum_{j=1}^{m1} \frac{n_j}{n} \frac{\left(\bar{Y}_j - \bar{Y}\right)^2}{s^2/n_j} = \frac{\sum_{j=1}^{m_1} n_j^2 \left(\bar{Y}_j - \bar{Y}\right)^2}{\sum_{i=1}^{n} (Y_i - \bar{Y})^2}$$

$Y_i$: $Y$ value of ith individual

$Y_j$: mean values of $Y$ in cell $j$
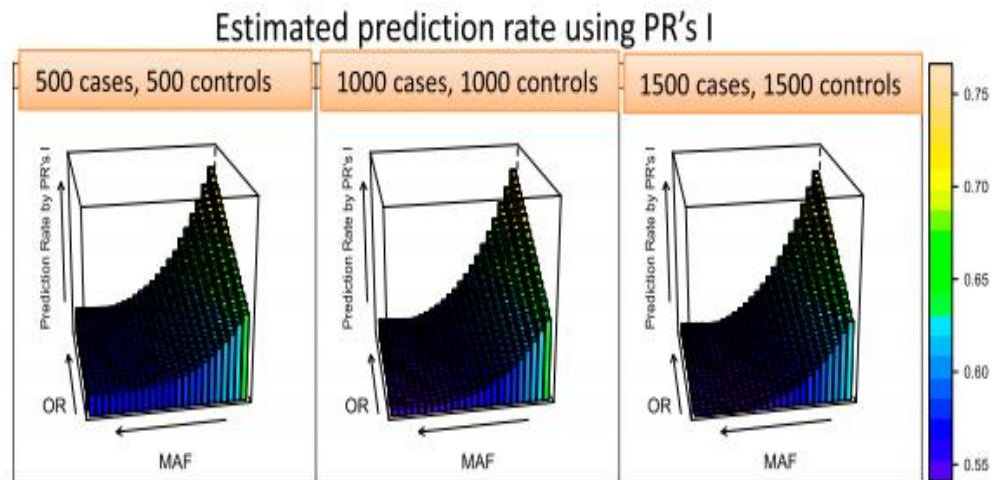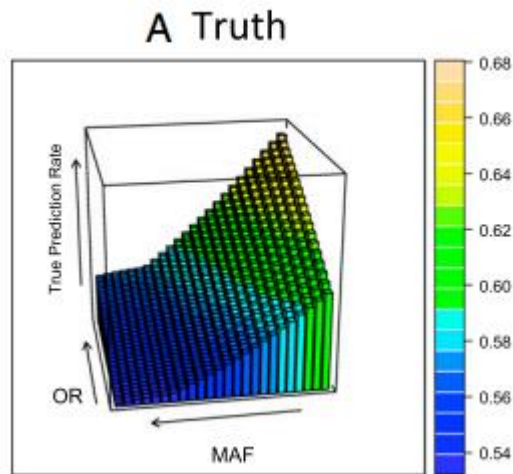
$\bar{Y}$: mean of all Ys

$m_1$: number of cells

- Weighted average of $X^2 s$
- Proportion of between-group variance divide by the total variance
- The bigger the I score, the better prediction ability
- Applied to small groups of

# Example 3

- 6 SNPs

# Conclusion and Comments

- Give us a sight into how the distribution of a variable will effect the significance and prediction ability.

- **Based on their way of measuring significance and prediction ability, which is not common(consistent) and powerful in real study**

- Establish a I score which to some extent useful in identifying whether a variable is predictive regardless of significance

- **Can only be used to small group of variables thus limited to be applied to general big genetic data when number of SNPs is huge**