

# GenomeDataAnalysis6

大阪大学大学院医学系研究科 遺伝統計学  
東京大学大学院医学系研究科 遺伝情報学  
理化学研究所生命医科学研究センター システム遺伝学チーム

<http://www.sg.med.osaka-u.ac.jp/index.html>

## GenomeDataAnalysis6

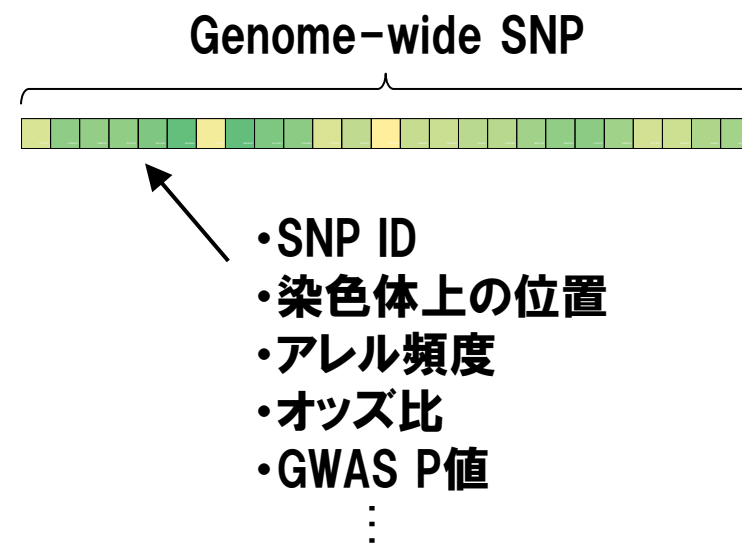
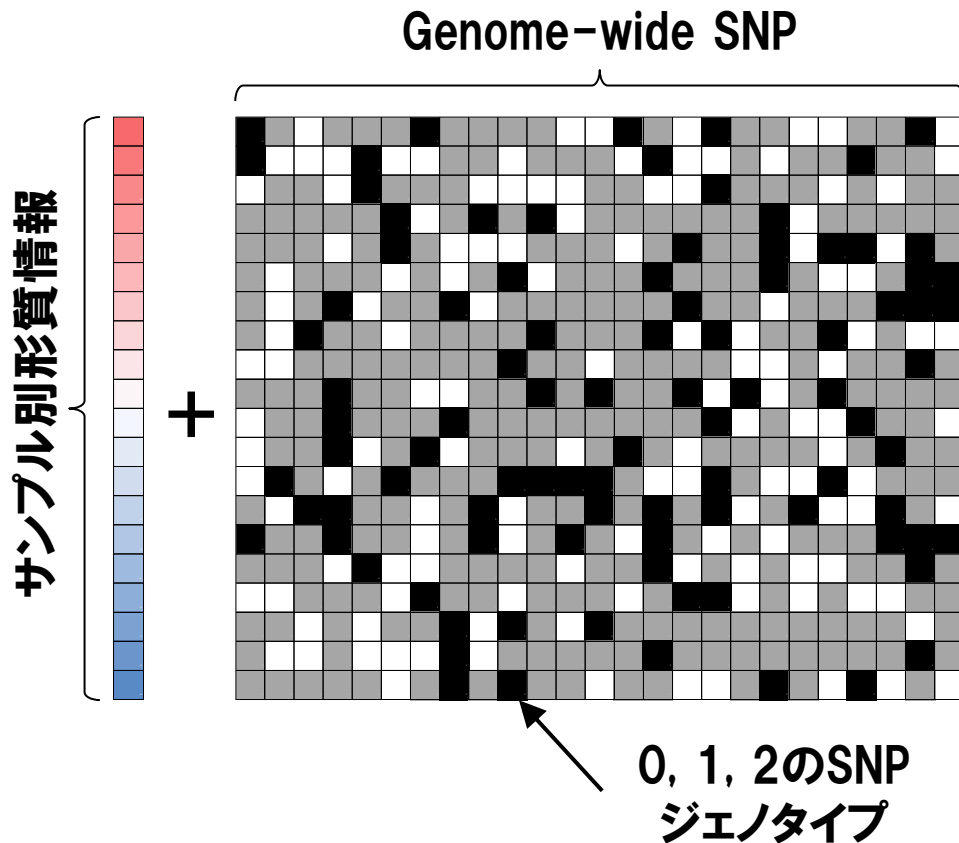
- ① **GWAS統計量を用いた解析手法とLDSC**
- ② **Anacondaを使ったLDSCのインストール**
- ③ **LDSC解析実習**

本講義資料は、Windows PC上で  
C:¥SummerSchoolにフォルダを配置すること  
を想定しています。

# ① GWAS統計量を用いた解析手法とLDSC

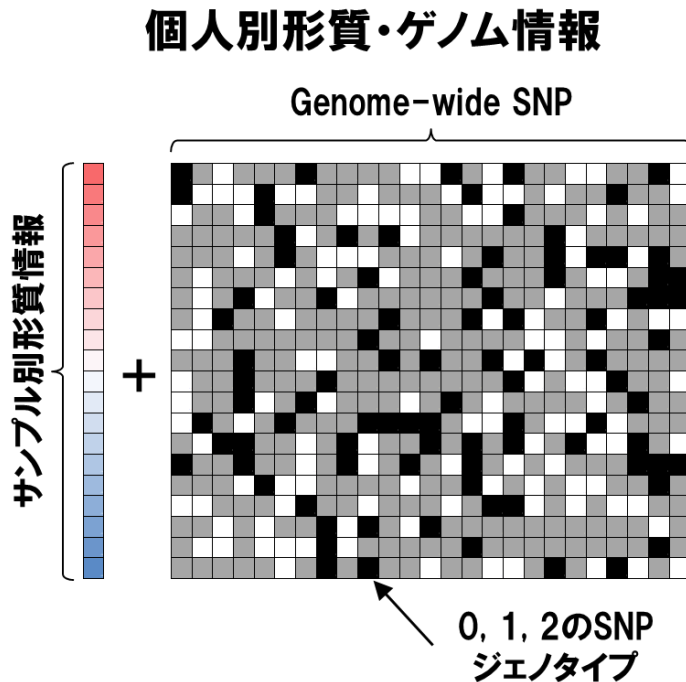
## 個人別形質・ゲノム情報

## GWAS統計量



- GWASゲノムデータ解析に際して使用する入力データは、大きく分けて2種類存在し、**個人別の形質情報とゲノム情報**(GWAS individual data)と、**GWAS統計量**(GWAS summary statistics)になります。

# ① GWAS統計量を用いた解析手法とLDSC

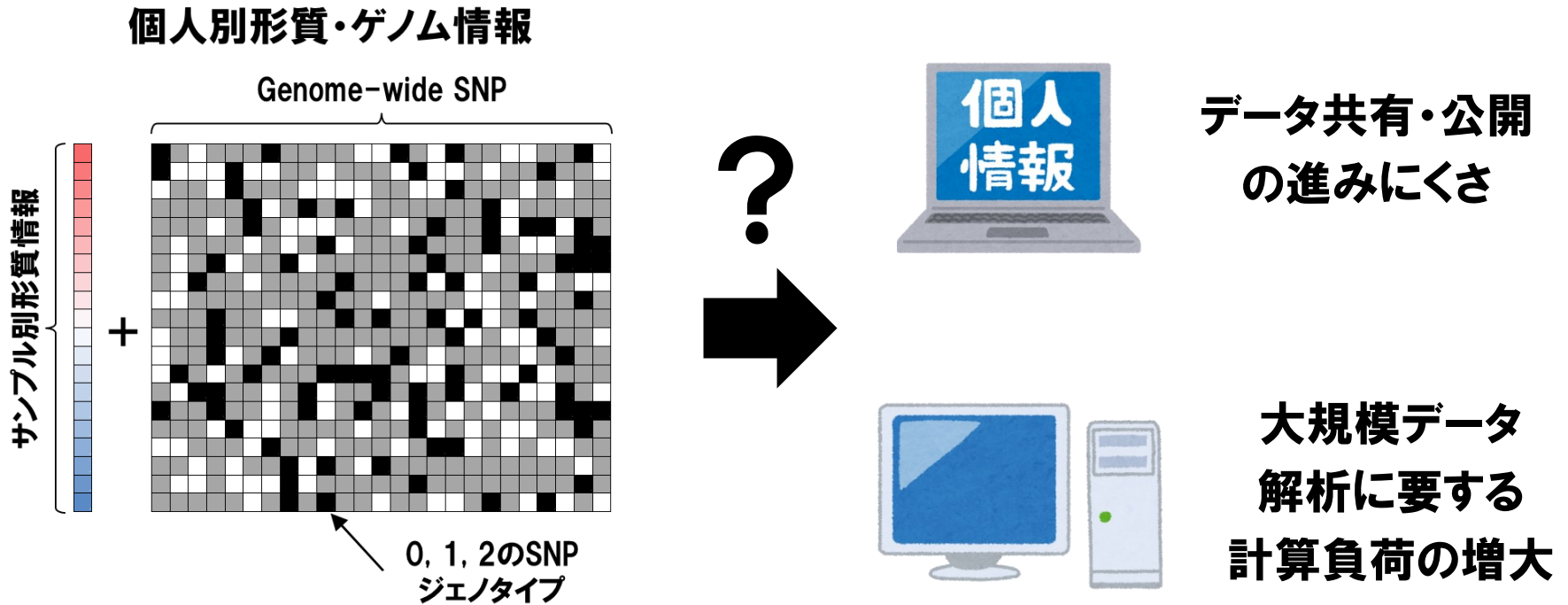


## 実施可能な解析手法

- GWAS
- eQTL解析
- Genotype imputation
- 選択圧解析
- Polygenic risk score推定
- Heritability推定
- Conditional解析
- Haplotype解析
- 

- 個人別の形質・ゲノム情報には、GWASゲノムデータ解析に必要な情報が(ほぼ)全て含まれています。
- GWASそのものも含め、**多彩なゲノムデータ解析が実施可能になります。**
- GenomeDataAnalysis 1～5で実施した演習内容は、いずれも個人別の形質・ゲノム情報を用いた解析手法が対象でした。

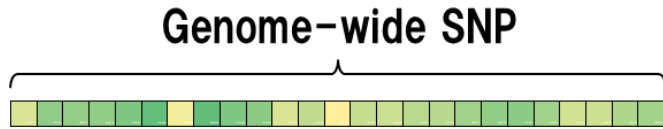
# ① GWAS統計量を用いた解析手法とLDSC



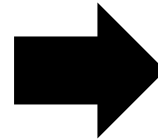
- 一方、個人別形質・ゲノム情報は個人情報に相当するため、**データの共有や公開が進みにくい**、というデメリットがあります。
- 特に、複数のGWASを統合するGWASメタアナリシスでは、全てのGWASの個人別情報へのアクセスが、**困難になります**。
- 数十万人規模の個人別情報解析は、**計算機資源の負荷も大きい**です。

# ① GWAS統計量を用いた解析手法とLDSC

## GWAS統計量



- SNP ID
- 場所
- アレル頻度
- オッズ比
- GWAS P値
- ⋮

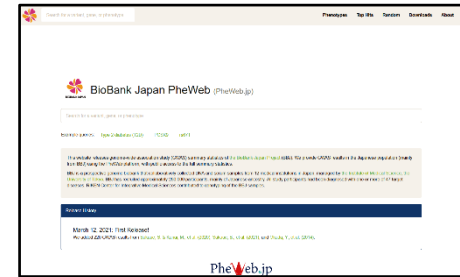


データの公開・  
入手が容易

## 公開データベース



<https://pan.ukbb.broadinstitute.org/>



<https://pheweb.jp/>

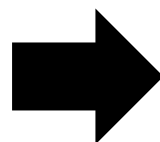
- 一方、GWAS統計量は個人情報に相当しないため、データの公開や入手が容易という利点があります。
- 現在、GWAS論文の出版時にはGWAS統計量を一般公開(=特定の手続きを行わずに自由に入手可能な状態)することが推奨されています。
- GWAS統計量を一般公開するデータベースの整備も進んでいます。

# ① GWAS統計量を用いた解析手法とLDSC

## GWAS統計量



- SNP ID
- 場所
- アレル頻度
- オッズ比
- GWAS P値
- …



データの公開・  
入手が容易

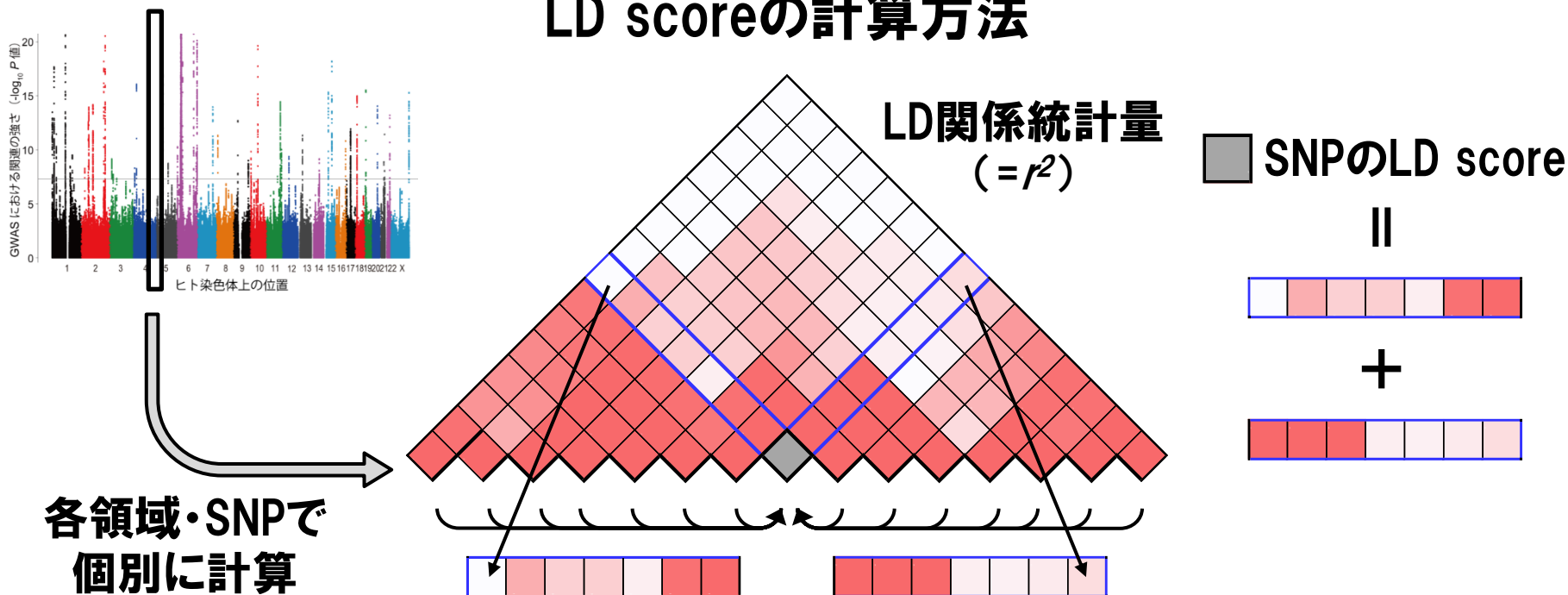
## 実施可能な解析手法

- GWASメタアナリシス
- Fine-mapping
- Functional annotation
- Gene-based関連解析
- 統計量imputation
- Conditional解析
- TWAS
- PRS予測モデル構築
- Mendelian randomization

- GWAS統計量を対象とした解析手法は、実施に際して要求される**計算機資源が比較的小規模**であり、ノートPCで実施可能な例が多いです。
- 解析手法の発達により、以前は個人別形質・ゲノム情報が必要と考えられていた解析も、GWAS統計量に基づき実施可能になっています。
- その一つに、**LDSCによるheritability推定**があります。

# ① GWAS統計量を用いた解析手法とLDSC

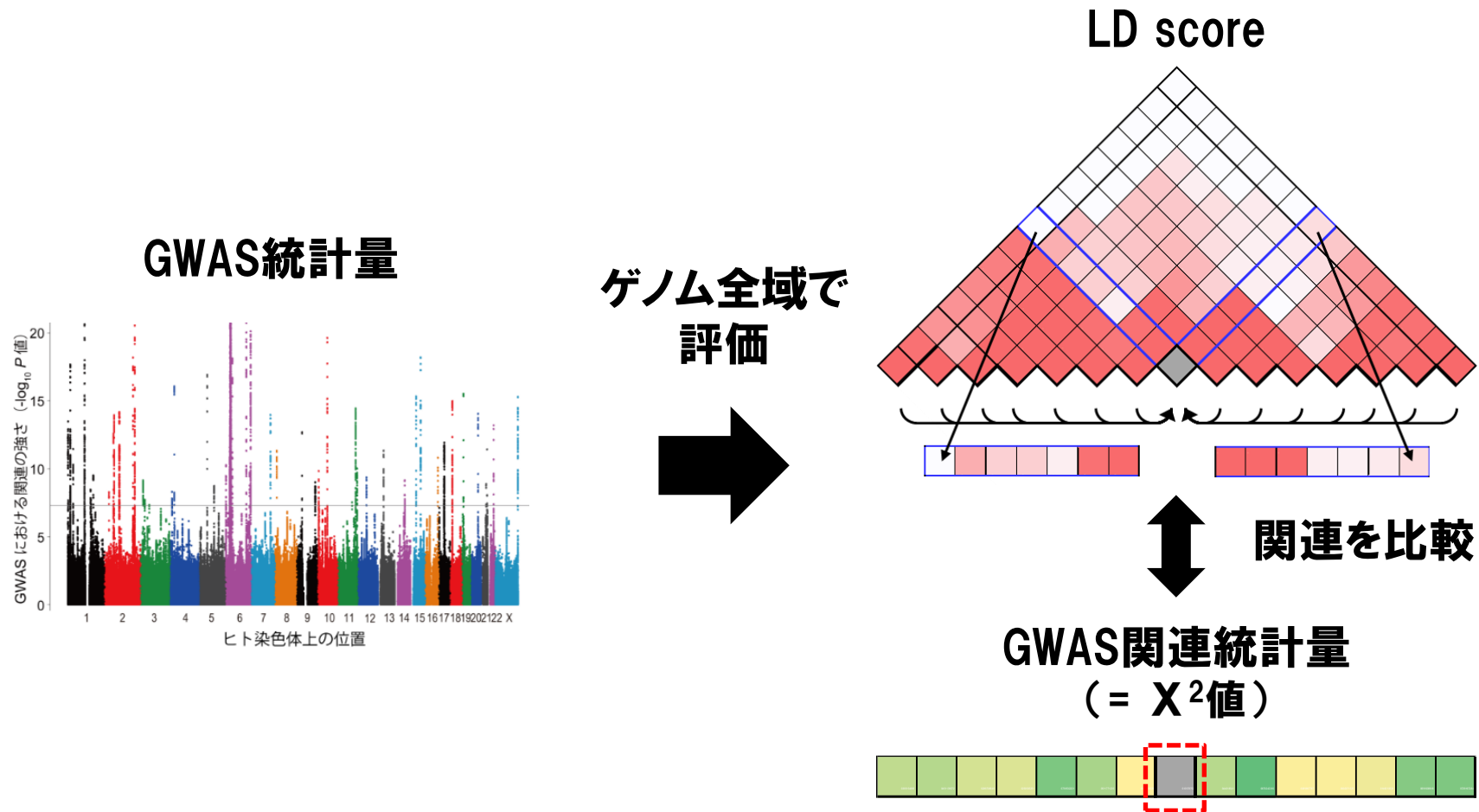
## LD scoreの計算方法



- **L**inkage **D**isequilibrium **S**core regression (LDSC) は、ゲノム全域に存在するコモンバリエーションにおけるLD scoreに基づく解析手法です。
- LD scoreは、予め各SNPにおいて、近傍のLD関係にあるSNPとのLD統計量を足し合わせて計算されています。
- LD scoreが高いSNP程、多くのSNPとLD関係にある、と解釈できます。  
(※LD scoreの計算自体には個人別ゲノムデータが必要になります。)



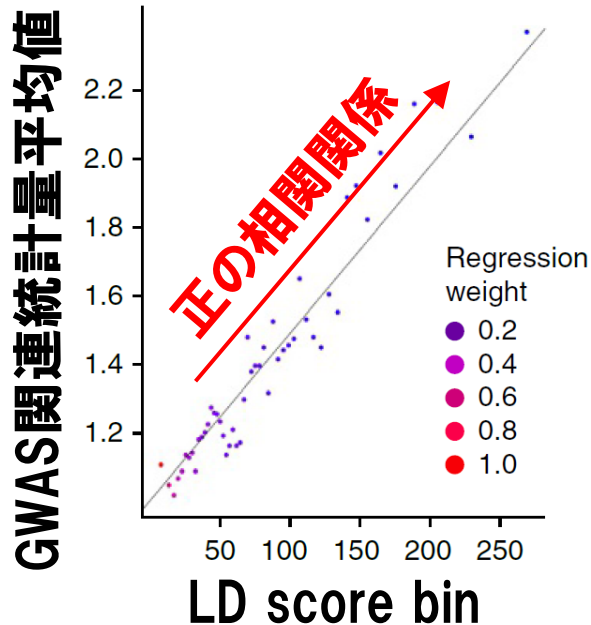
# ① GWAS統計量を用いた解析手法とLDSC



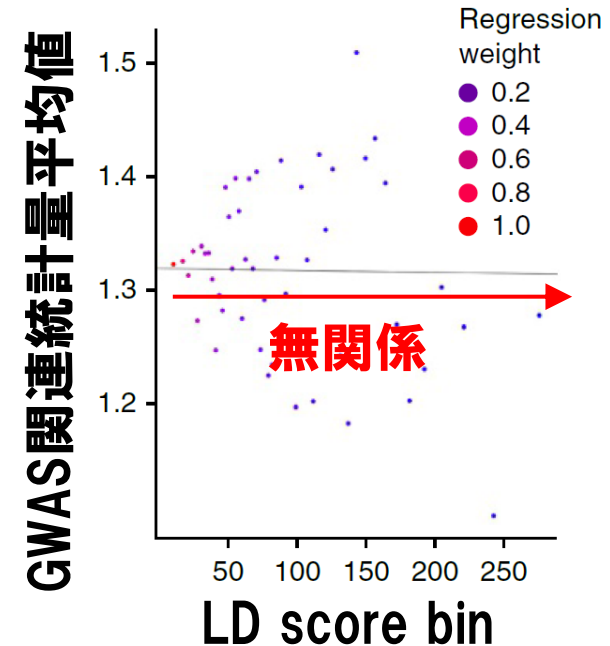
- LDSCでは、各SNPのLD scoreとGWAS関連統計量(=  $X^2$ 値)の関連をゲノム全域に渡って比較(=regression)することで、疾患の遺伝的背景(=heritability)を定量化します。

# ① GWAS統計量を用いた解析手法とLDSC

## Polygenic形質のGWAS



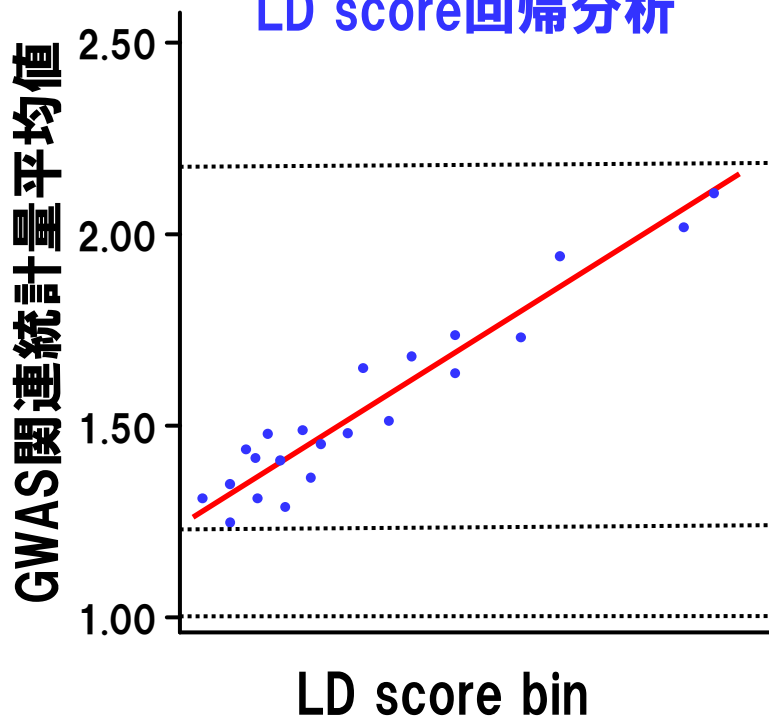
## 集団構造化のみのGWAS



- **Polygenic仮説**(=ゲノム全域の無数のSNPが弱い疾患リスクを有する)が成立する形質のGWASでは、**LD scoreが高いSNP程、GWAS関連統計量平均値も高くなる傾向**(正の相関関係)が認められます。
- GWAS統計量に疾患感受性が反映されず、**集団構造化等で見せかけの関連が存在する場合、両者は無関係になります。**

# ① GWAS統計量を用いた解析手法とLDSC

## LD score回帰分析



$$E[X_j^2] = 1 + Na + \frac{h_g^2 N}{M} l_j$$

Polygenicityに由来  
(= slope、傾き)

集団構造化に由来  
(= intercept、切片)

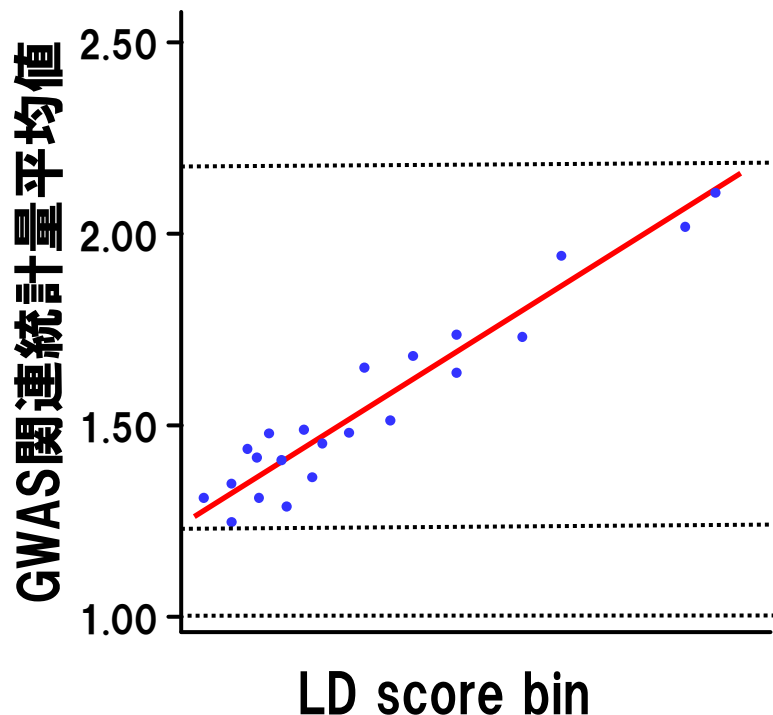
- $X^2$ : GWAS関連統計量
- $N$ : サンプル数
- $a$ : 集団構造化
- $h_g^2$ : heritability
- $M$ : SNP数
- $l$ : LD score
- $j$ : SNP

- LDSCにより、GWAS統計量における統計量のinflation(=  $\lambda_{GC}$ の1.00からの乖離)を、集団構造化からの由来(= intercept、切片)と、polygenicityからの由来(= slope、傾き)に分解することができる、とも解釈できます。
- 後者に基づき、heritabilityを推定することができます。

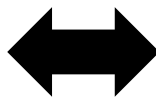
(Heritability=ゲノムワイドな遺伝子多型が説明可能な疾患の遺伝的背景の割合)

# ① GWAS統計量を用いた解析手法とLDSC

形質AのLD score回帰

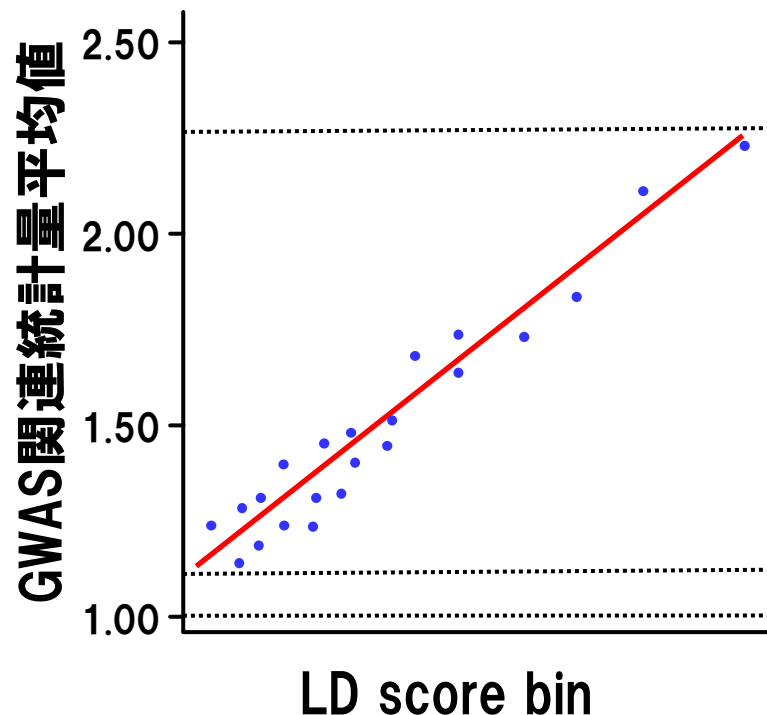


ゲノム全域  
で比較



形質Aと  
形質Bの  
遺伝的相関  
の定量化

形質BのLD score回帰

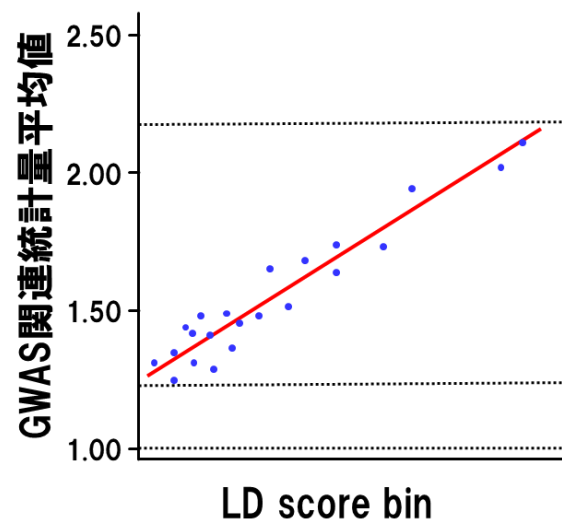


- 二つのpolygenicな形質に対するLD score回帰の結果をゲノム全域で比較することで、**遺伝的背景の相関関係(=genetic correlation)**を定量化することができます。

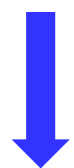
(正の相関と負の相関を区別する都合、GWAS関連統計量は正の値のみをとる $\chi^2$ 値(=Z値 $\times$ Z値)でなくZ値を用いて計算されます。) (Bulik-Sullivan BK et al. *Nat Genet* 2015)

# ① GWAS統計量を用いた解析手法とLDSC

## LD score回帰分析

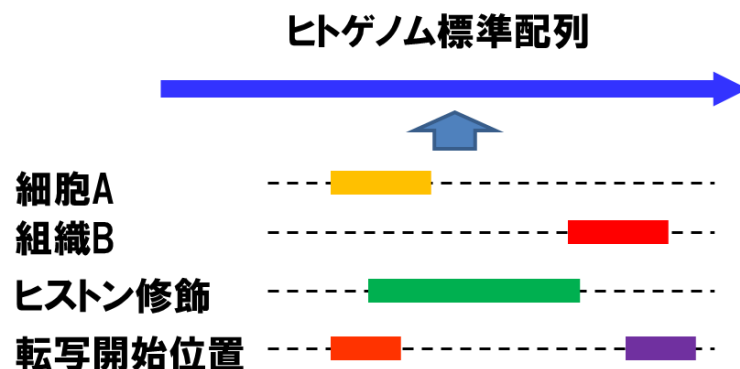


各カテゴリーに属する  
ゲノム領域に層別化  
した解析



各カテゴリーにおける  
heritabilityのenrichment  
を定量化

## 層別化カテゴリー情報



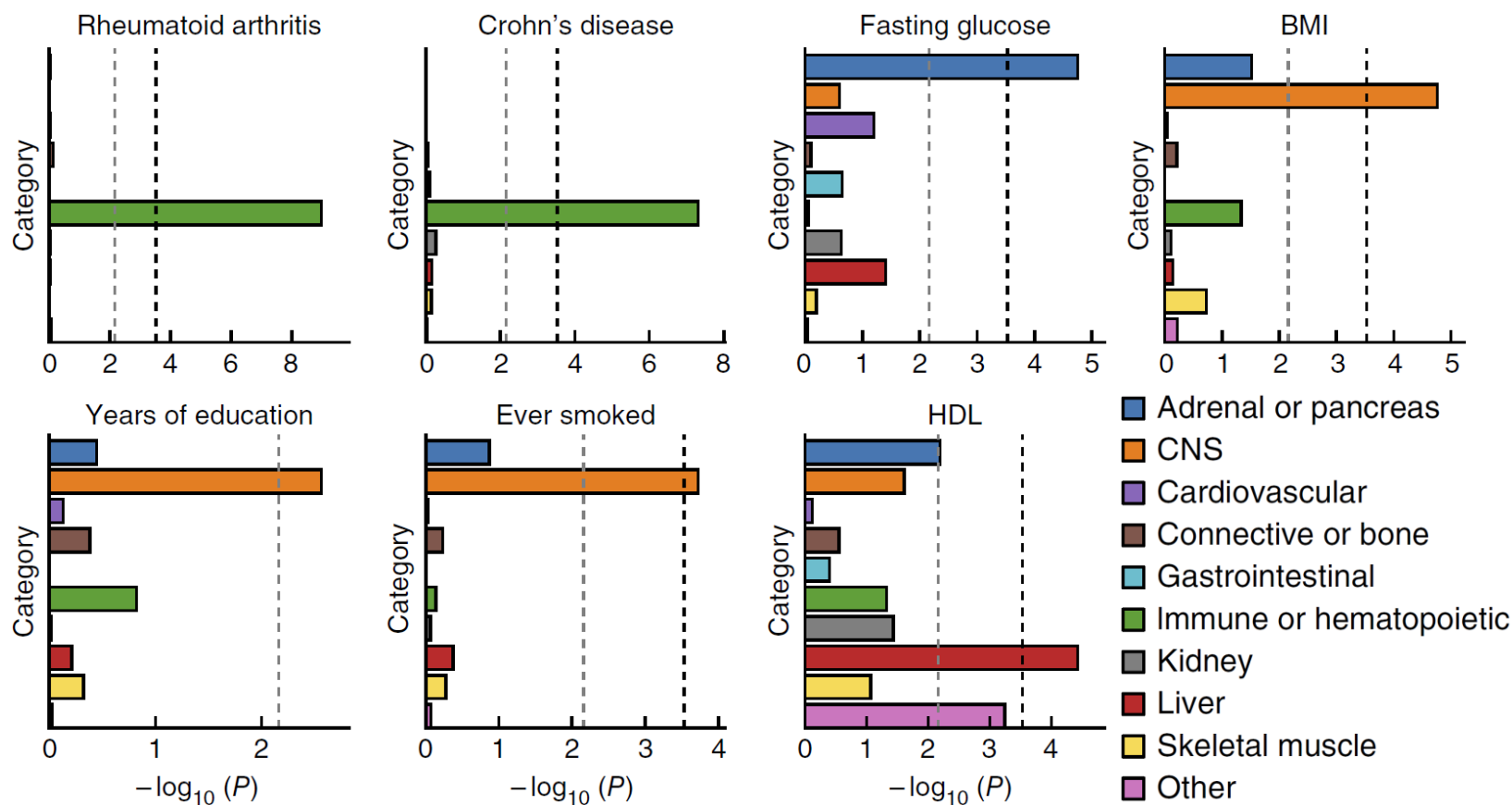
$$E[X_j^2] = 1 + Na + N \sum_C \tau_{C^l(j,C)}$$

- C: 層別化カテゴリー
- $\tau$ : SNP毎のheritabilityへの寄与

- ヒトゲノム領域を、細胞組織に特異的なエピゲノム修飾等のカテゴリーに層別化した上でLDSCを実施することで、カテゴリー別のheritability (=partitioned heritability)のenrichmentを定量化できます (=stratified LDSC)。
- どの細胞組織・エピゲノム修飾が病態に重要か、知ることができます。

# ① GWAS統計量を用いた解析手法とLDSC

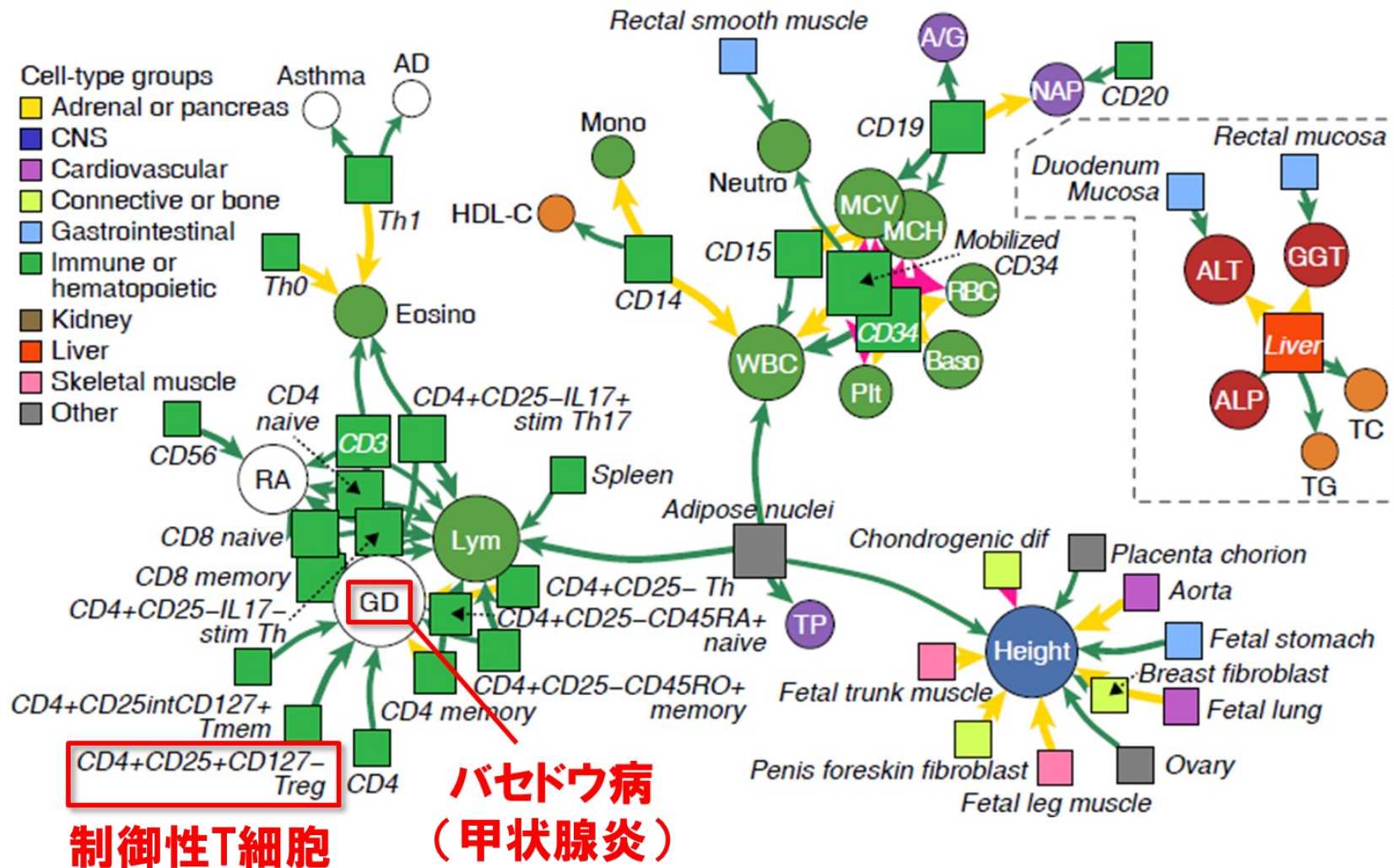
## GWAS統計量に対するStratified LDSC解析結果



• Stratified LDSC解析をGWAS統計量に適用した結果、**関節リウマチ・クローン病における免疫細胞、肥満・喫煙歴における中枢神経系の関与**など、疾患の細胞組織特異性が明らかになりました。



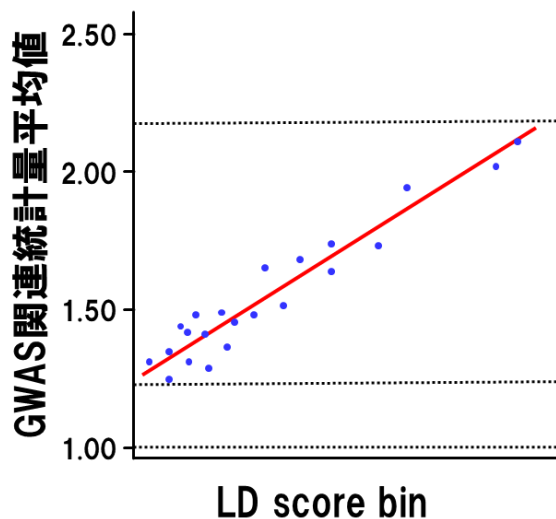
# ① GWAS統計量を用いた解析手法とLDSC



・日本人集団100形質、200細胞組織を対象にstratified LDSC解析を実施したところ、疾患と細胞組織のネットワークが再構築されました。

# ① GWAS統計量を用いた解析手法とLDSC

## LD score回帰分析



細胞組織別の  
遺伝子発現を  
考慮したLDSC解析



各細胞組織における  
heritabilityのenrichment  
を細分化して定量化

## 細胞組織特異的 遺伝子発現情報

	細胞組織A	細胞組織B	細胞組織C	細胞組織D	細胞組織E	細胞組織F
遺伝子A	白	赤	青	白	青	白
遺伝子B	赤	白	赤	青	赤	白
遺伝子C	赤	白	白	青	白	赤
遺伝子D	赤	白	青	青	白	青
遺伝子E	青	青	白	青	赤	白
遺伝子F	青	赤	赤	青	赤	青

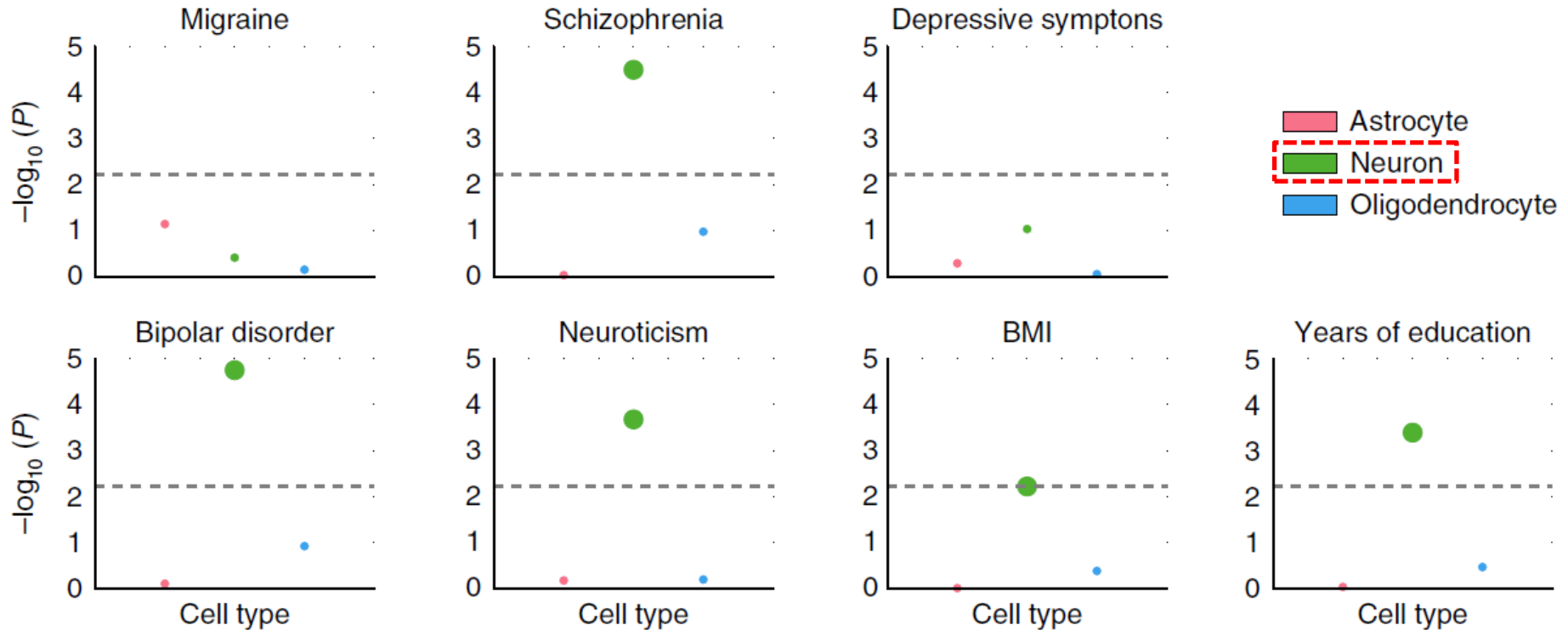
SEG = Specifically  
Expressed Genes

- 細分化された細胞組織特異性を検討する手法として、**細胞組織特異的遺伝子発現情報を考慮したLDSC(=LDSC-SEG)**も開発されています。
- どの細胞組織が疾患病態に重要か、より細かい細胞組織分類に基づき検討することができます。



# ① GWAS統計量を用いた解析手法とLDSC

## GWAS統計量に対するLDSC-SEG解析結果



- 中枢神経系に関わる形質のGWAS統計量にLDSC-SEG解析を適用した結果、Astrocyte、Neuron、Oligodendrocyteの3種類の神経細胞の中で、**Neuronの寄与が相対的に大きい事**が明らかになりました。

# ① GWAS統計量を用いた解析手法とLDSC

- LDSCは、heritabilityの推定や細胞組織特異性の解明など、疾患ゲノムデータ解析の根幹となる情報を、GWAS統計量に基づき小さい計算負荷で推定可能にした点が、画期的な解析手法でした。
- LDSCを実施する際の注意点として、「GWASサンプル数が数千名以上であること」が挙げられます。
- サンプル数が少ないGWASについては、個人別形質・ジェノタイプ情報を用いた手法など、異なるアプローチが必要になります。
- 「GWASのサンプル集団と同一のサンプル集団でLD scoreが計算されていること」、も注意点になります。
- 異なる人種集団間の遺伝的相関の検定にはLDSCは使用できず、異なる解析手法が必要となることに、留意してください。

# ① GWAS統計量を用いた解析手法とLDSC

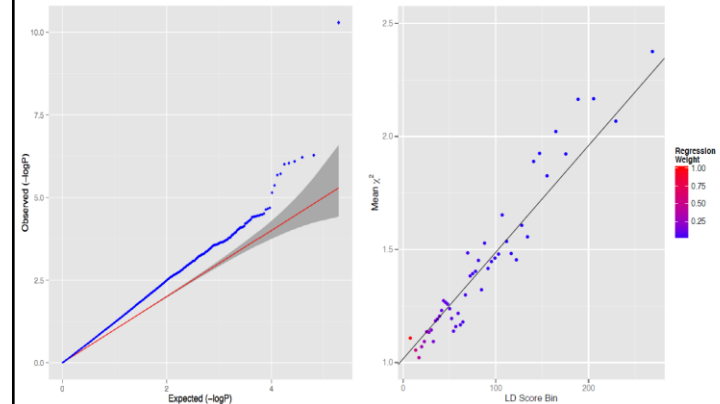
## Use and Interpretation of LD Score Regression

Brendan Bulik-Sullivan  
[bulik@broadinstitute.org](mailto:bulik@broadinstitute.org)  
PGC Stat Analysis Call



## Simulated Polygenicity

- $\lambda_{GC} = 1.30$ ; LD Score Regression intercept = 1.02



## LD Score Regression

- Regress  $\chi^2$  statistics against LD Score

$$E[\chi^2 | l_j] = Nh^2 l_j / M + Na + 1$$

- LD Score ( $l_j$ ) is a property of SNP  $j$ , defined as sum  $r^2$ , estimated as sum  $r^2 w$  / all other SNPs a 1cM window.
- $N$  is sample size.
- $M$  is # SNPs.
- $h^2$  is SNP-heritability.
- $a$  is inflation from pop strat/cryptic relatedness.

Bulik-Sullivan et al., Nat Genet, 2015

## URLs

- ldsc
  - [github.com/bulik/ldsc](https://github.com/bulik/ldsc)
  - [Installation instructions](#)
  - [FAQ](#)
- Tutorials / wiki
  - [github.com/bulik/ldsc/wiki](https://github.com/bulik/ldsc/wiki)
- Pre-computed European LD Scores
  - [broadinstitute.org/~bulik/eur\\_ldscores/](https://broadinstitute.org/~bulik/eur_ldscores/)
- ldsc\_users google group:
  - [groups.google.com/forum/?hl=en#forum/ldsc\\_users](https://groups.google.com/forum/?hl=en#forum/ldsc_users)

[https://www.med.unc.edu/pgc/wp-content/uploads/sites/959/2019/01/pgc\\_stat\\_bulik\\_2015.pdf](https://www.med.unc.edu/pgc/wp-content/uploads/sites/959/2019/01/pgc_stat_bulik_2015.pdf)

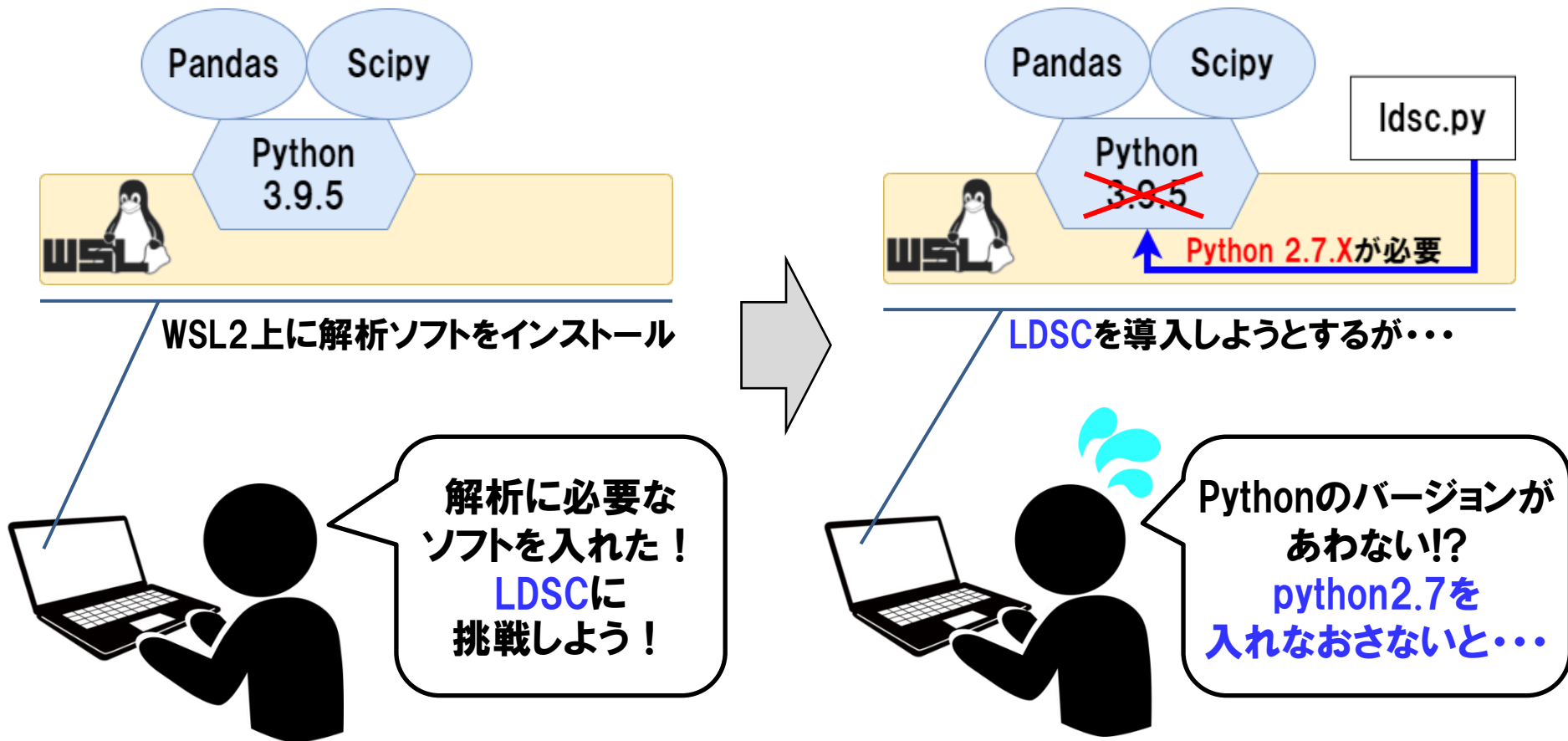
• LDSCの作者による説明スライドが公開されており、参考になります。<sup>19</sup>

## GenomeDataAnalysis6

- ① GWAS統計量を用いた解析手法とLDSC
- ② Anacondaを使ったLDSCのインストール
- ③ LDSC解析実習

本講義資料は、Windows PC上で  
C:¥SummerSchoolにフォルダを配置することを  
想定しています。

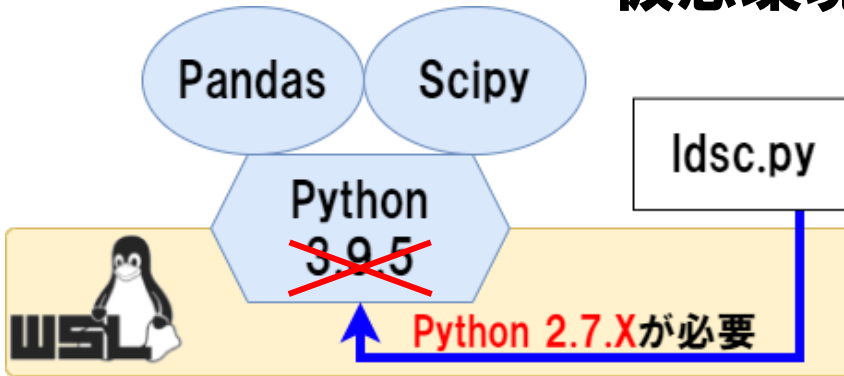
## ② Anacondaを使ったLDSCのインストール



- ソフトウェアは他のプログラムやパッケージに依存している事があります。
- 実際の解析では多くのソフトウェアを使用しますが、ソフトウェア同士でプログラムのバージョンが合わずに困ることがよくあります。

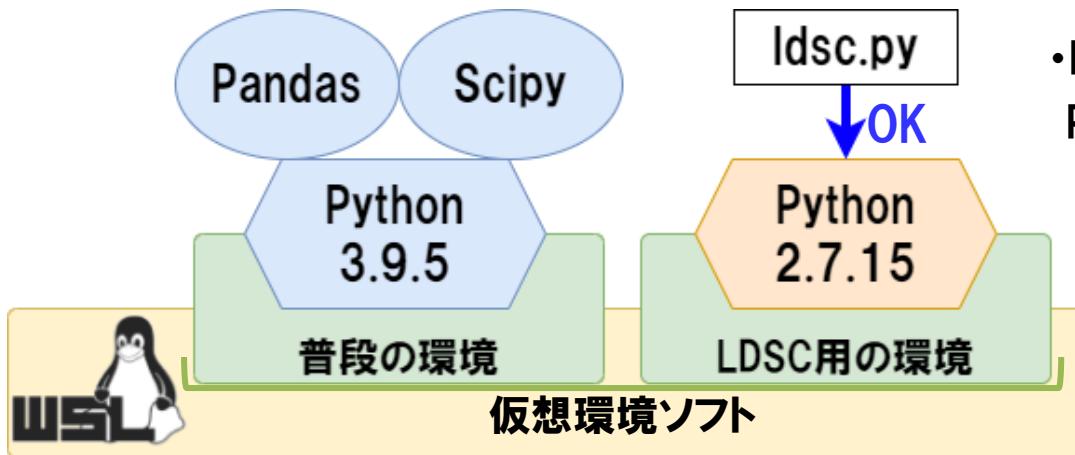
## ② Anacondaを使ったLDSCのインストール

### 仮想環境ソフトなし



- Pythonを一度すべて消して入れなおす。  
(→今までに構築したモジュールも消える)
- Python2.7を追加で入れる。  
(→**コマンドが重複する**など、混乱が生じる)

### 仮想環境ソフトあり



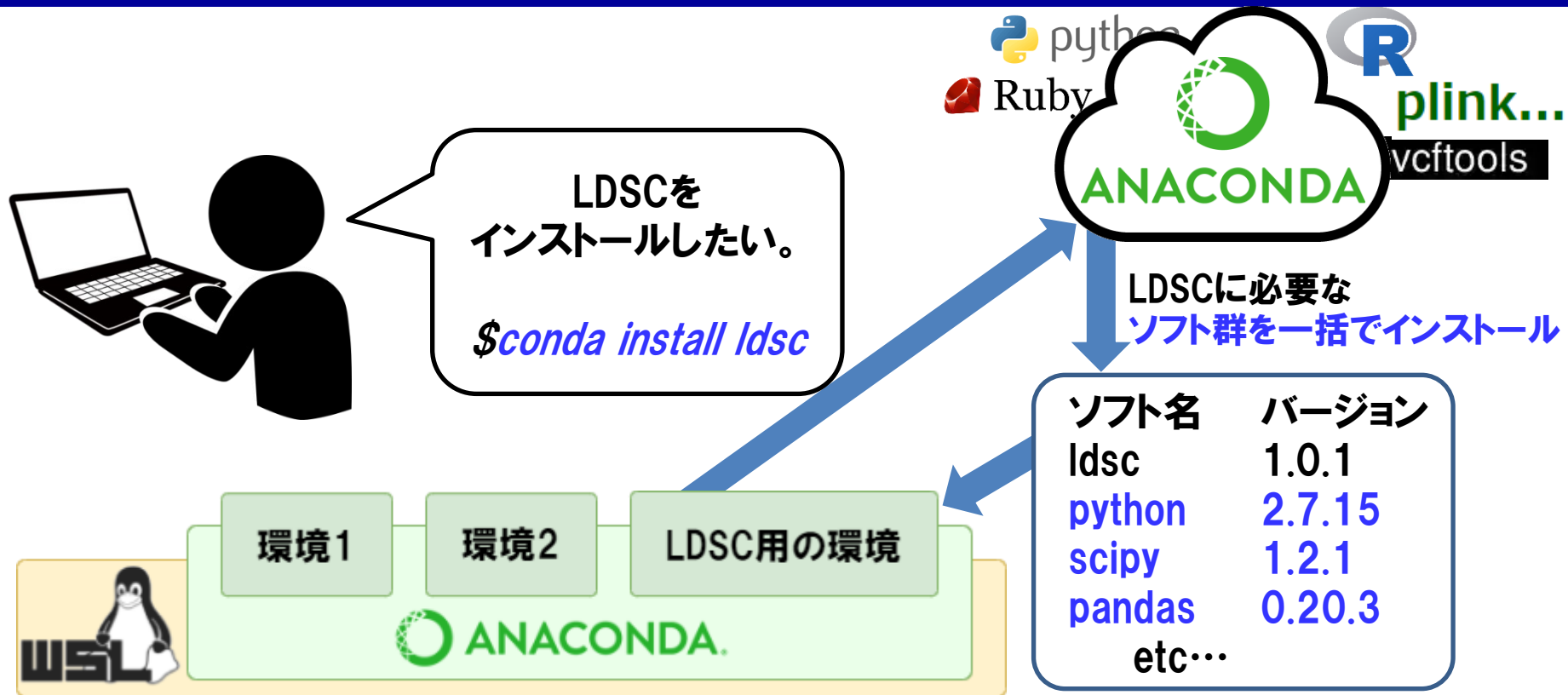
- LDSC用の環境を新規に作成してそこに Python2.7をインストールする。

環境の切り替えも  
簡単！



- ソフトの依存関係問題を解決する方法の一つに、**仮想環境**があります。  
**解析ごとに環境を分けて**ソフトをインストールできます。
- 様々な仮想環境ソフトがありますが、今回は**Anaconda**を使用します。

## ② Anacondaを使ったLDSCのインストール



- AnacondaはPythonとそのモジュールを管理するソフトですが、現在はPythonに限らずデータサイエンス用のソフトに広く対応しています。
- 主な機能は、仮想環境を構築すること、Anacondaのリポジトリ(ソフトが揃っているクラウドリソース)からソフトをインストールすることです。
- ソフト名を指定するだけで必要なソフトを一括インストールしてくれます。

## ② Anacondaを使ったLDSCのインストール

statgen@statgen-PC: ~

```
$ conda env create -f /mnt/c/SummerSchool/___for_install___/ldsc_setting.yml
```

### ldsc\_setting.yml

```
name: ldsc_env
channels:
- bioconda
- conda-forge
dependencies:
- r-base=3.6.1
- r-corrplot=0.84
- r-ggplot2=3.1.1
- python=2.7.18
- bitarray=0.8.3
  ⋮
```

→ 作成する環境名

→ ソフトを探すAnacondaのリポジトリ名

→ ダウンロード・インストールするソフト名

- LDSCの実習で必要なAnaconda環境は手順書で説明しています。
- 手順書で実行してもらった上記コマンドは、ymlファイルに記載しているソフトウェアをインストールしつつ、仮想環境を構築するコマンドでした。<sup>24</sup>



## ② Anacondaを使ったLDSCのインストール

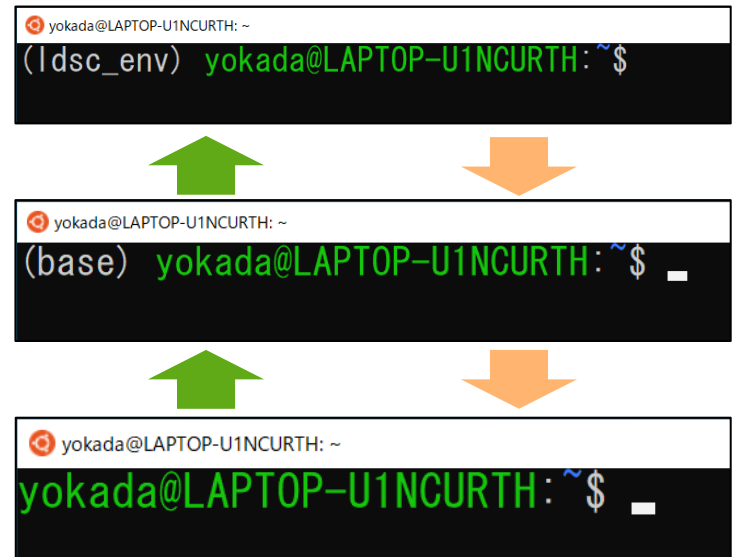
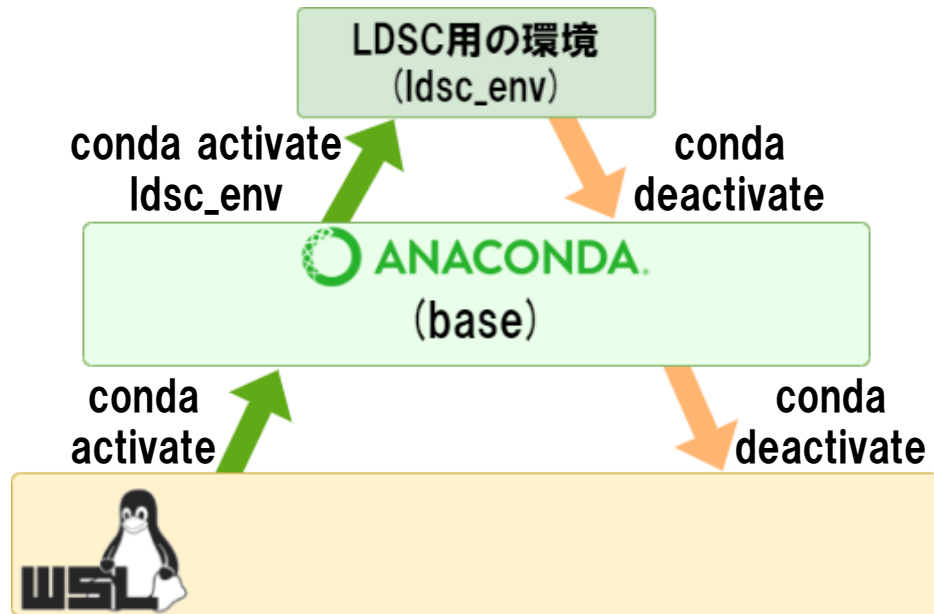
statgen@statgen-PC: ~

```
$ conda activate
```

```
$ conda activate ldsc_env
```

```
$ conda deactivate
```

```
$ conda deactivate
```



- 皆さんのPCには、事前に配布した手順書通りに行っていれば、既に”ldsc\_env”という環境ができています。
- コマンドを入力して、環境を切り替えてみましょう。

## ② Anacondaを使ったLDSCのインストール

### M1/M2 MacでAnaconda3がうまく使えない理由？

		x86_x64系		Arm v〇系	
命令セット アーキテクチャ		Intel64	AMD64	Arm v8	
使用ハード		PC・サーバー			スマートフォン
プロセッサ シリーズ	コンシューマPC	IntelCore	AMD Ryzen	Apple M	Apple A Qualcomm Snapdragon
	サーバー用PC	IntelXeon	AMD Ryzen Threadripper AMD EPYC	FUJITSU Processor A64FX	
特徴		高パフォーマンス・設計の高い汎用性			小型・省電力・省スペース・設計の独自性

殆どのWindows PC  
M1以前のMac

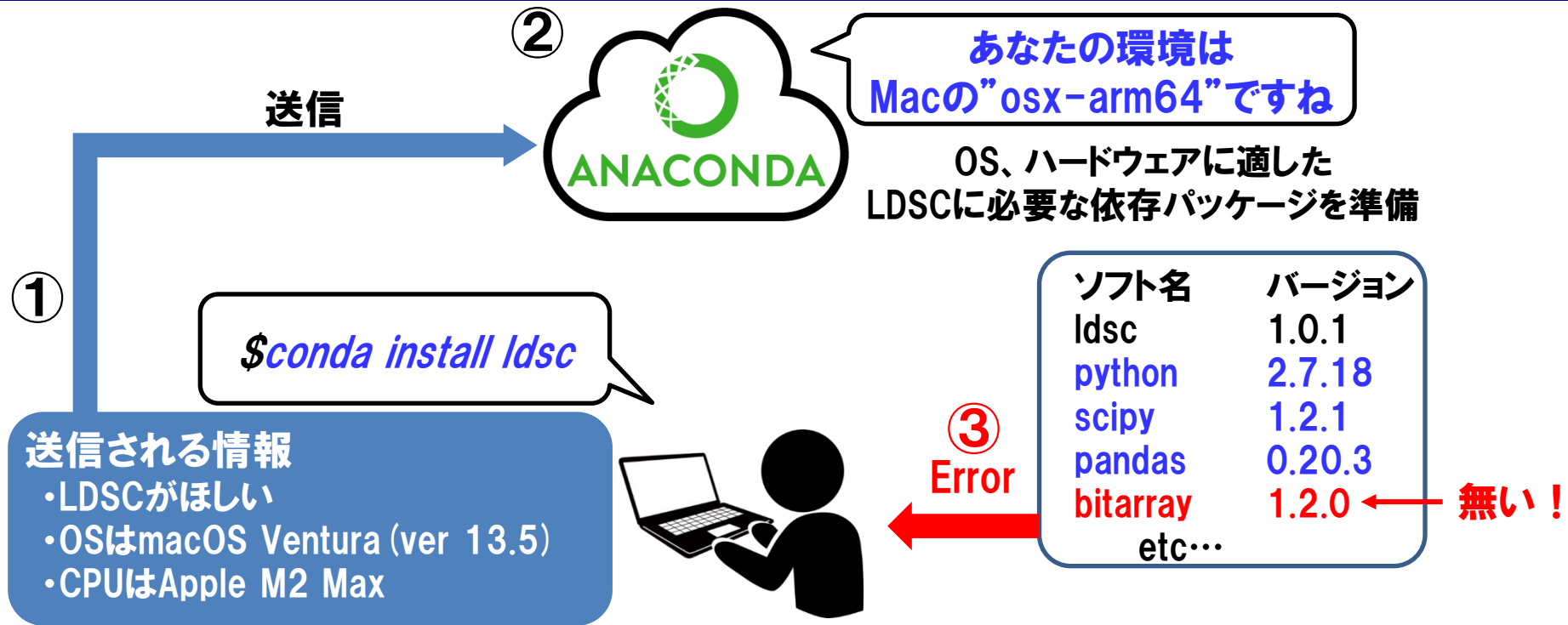


M1/M2 Mac



- M1/M2 MacでAnaconda3の実装の不具合が報告されています。
- ⇒CPU命令セットアーキテクチャ(Instruction Set Architecture : ISA)の違い。
- x84\_x64系統とArm系統のCPUにはソフトウェアの互換性はありません。
- Arm版CPUへのソフトウェア実装時には個別CPUにあわせたコンパイルが必要になります。

## ② Anacondaを使ったLDSCのインストール



- Anacondaは各CPU用のコンパイル済みのパッケージを配布している。
- まだM1/M2 CPUに用意されていないパッケージが相当数ある。
- 根本的な解決策は、必要なパッケージやライブラリを自分で把握してM1/M2 Mac上でソースコードからコンパイルするのが対策方法(エンドユーザーにはまだ敷居が高い?)。
- 今後対応が進んでいく可能性はあります。

## GenomeDataAnalysis6

- ① GWAS統計量を用いた解析手法とLDSC
- ② Anacondaを使ったLDSCのインストール
- ③ LDSC解析実習

本講義資料は、Windows PC上で  
C:¥SummerSchoolにフォルダを配置することを  
想定しています。

# ③ LDSC解析実習

The screenshot shows the GitHub repository page for `bulik/ldsc`. The repository is in the `master` branch and has 3 branches and 2 tags. The file list includes:

File/Folder	Description	Last Update
ContinuousAnnotations	Update quantile_h2g.r	3 years ago
ldscore	Fix globbing bug in splitp (#221)	13 months ago
test	Improve consistency of unit test Hsq_2D (#163)	2 years ago
.gitignore	Improve consistency of unit test Hsq_2D (#163)	2 years ago
CHANGELOG	Fix key error in allele_merge (#185)	2 years ago
LICENSE	update license + masthead	7 years ago
README.md	Update README.md	2 years ago
environment.yml	Update numpy to 1.16, make dependencies consistent (#212)	15 months ago
ldsc.py	v1.0.1 (#164)	2 years ago
make_annot.py	fix make_annot bug (#199)	15 months ago
munge_sumstats.py	Fix key error in allele_merge (#185)	2 years ago
requirements.txt	Update numpy to 1.16, make dependencies consistent (#212)	15 months ago
setup.py	Update numpy to 1.16, make dependencies consistent (#212)	15 months ago

The right sidebar shows repository details: **About** (LD Score Regression (LDSC)), **Releases** (2), and **Contributors** (13).

<https://github.com/bulik/ldsc>

[https://www.med.unc.edu/pgc/wp-content/uploads/sites/959/2019/01/pgc\\_stat\\_bulik\\_2015.pdf](https://www.med.unc.edu/pgc/wp-content/uploads/sites/959/2019/01/pgc_stat_bulik_2015.pdf)

- LDSCのソースコードはgithub上で公開されています。
- 実施方法の説明や解析に必要なデータのダウンロードリンク、FAQなど内容が充実しており、手順に沿った一通りの解析も実施可能です。<sup>29</sup>

### ③ LDSC解析実習

./GenomeDataAnalysis6/Analysis

./Input/

1000G\_EUR\_Phase3\_baseline/  
1000G\_frq/  
1kg\_eur/  
baseline/  
Cahoy\_1000Gv3\_ldscores/  
cell\_type\_groups/  
eas\_ldscores/  
GWASstats/  
weights\_hm3\_no\_hla/

./Output/

CellTypeSpecificity/  
Heritability/  
LDscore/

./ldsc/

.git/  
ContinuousAnnotations/  
ldscore/  
Test/  
.gitignore  
CHANGELOG  
environment.yml  
ldsc.log  
ldsc.py  
ldsc.pyc  
LICENSE  
make\_annot.py  
mungesumstats.py  
README.md  
requirements.txt  
setup.py

#### 解析実習内容

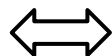
1. LD scoreの計算
2. Heritability推定
3. Genetic correlation解析
4. Stratified LDSC解析
5. LDSC-SEG解析

- LDSCには複数の解析方法があり、また多くのファイル群を扱います。
- 整理の観点から、入力ファイル群(./Input/)、出力ファイル群(./Output/)、実行ファイル群(./ldsc/)、に分類して配置しています。
- 今回の実習では、5種類のLDSC解析を実施します。

### ③ LDSC解析実習

./Input/

```
1000G_EUR_Phase3_baseline/  
1000G_frq/  
1kg_eur/  
baseline/  
Cahoy_1000Gv3_ldscores/  
cell_type_groups/  
eas_ldscores/  
GWASstats/  
weights_hm3_no_hla/
```



./Output/

```
CellTypeSpecificity/  
Heritability/  
LDscore/
```



./ldsc/

```
.git/  
ContinuousAnnotations/  
ldscore/  
Test/  
.gitignore  
CHANGELOG  
environment.yml  
ldsc.log  
ldsc.py  
ldsc.pyc  
LICENSE  
make_annot.py  
mungesumstats.py  
README.md  
requirements.txt  
setup.py
```

#### 解析実習内容

1. LD scoreの計算
2. Heritability推定
3. Genetic correlation解析
4. Stratified LDSC解析
5. LDSC-SEG解析

- 「LD scoreの計算」を行います。
- 主要な人種集団におけるLD scoreは、既に計算され、githubホームページ等で一般公開されており、必ずしも必要な作業ではありません。
- 独自のLD referenceパネルを持っている場合、LD scoreを改めて計算することで、より高解像度の解析が可能になります。

### ③ LDSC解析実習

```
(base) statgen@statgen-PC: ~
```

```
$ conda activate ldsc_env
```

※ファイル”LDSC\_Command.txt”を開いて、内容をLinux  
コマンドにコピー&ペーストして下さい。

Anaconda環境の起動



```
(ldsc_env) statgen@statgen-PC: ~
```

```
$ cd /mnt/c/SummerSchool/GenomeDataAnalysis6/Analysis/
```

```
(ldsc_env) statgen@statgen-PC: /mnt/c/SummerSchool/GenomeDataAnalysis6/Analysis
```

```
$ python ./ldsc/ldsc.py --bfile ./Input/1kg_eur/22 --l2 --ld-window-cm 1 --out  
./Output/LDscore/EUR_22
```

- AnacondaでLDSCの実行に必要な仮想環境を起動します。  
(LDSC解析演習は、Cygwin環境やM1 M2 macには対応していません。)
- 1000 Genomes Project EURサンプルの22番染色体のSNPジェノタイプファイル(PLINKフォーマット・379名・19,156SNP)を使って、LD scoreを計算します。



### ③ LDSC解析実習

”EUR\_22.LDscore.gz”  
各SNP毎のLD scoreファイル

```
1 | CHR SNP BP L2↓
2 | 22 rs9617528 16061016 1.271↓
3 | 22 rs4911642 16504399 1.805↓
4 | 22 rs140378 16877135 3.849↓
5 | 22 rs131560 16877230 3.769↓
6 | 22 rs7287144 16886873 7.226↓
7 | 22 rs5748616 16888900 7.379↓
8 | 22 rs5748662 16892858 7.195↓
9 | 22 rs5994034 16894090 2.898↓
10 | 22 rs4010554 16894264 6.975↓
11 | 22 rs4010558 16896762 7.379↓
12 | 22 rs3954571 16953560 5.242↓
13 | 22 rs11089179 16953727 4.606↓
14 | 22 rs9604821 17012935 5.327↓
15 | 22 rs2379965 17023514 5.290↓
16 | 22 rs2379981 17030792 4.194↓
17 | 22 rs4535153 17031072 4.201↓
```

LD score

MAFとLD score  
の相関係数

”EUR\_22.log”  
ログファイル

```
21 ↓
22 Summary of LD Scores in ./Output/LDscore/EUR_22.LDscore.gz↓
23 ..... MAF ..... L2↓
24 mean 0.2323 18.5353↓
25 std 0.1453 16.1039↓
26 min 0.0013 0.0657↓
27 25% 0.1042 7.8392↓
28 50% 0.2243 13.4837↓
29 75% 0.3549 22.9722↓
30 max 0.5000 109.7163↓
31 ↓
32 MAF/LD Score Correlation Matrix↓
33 ..... MAF ..... L2↓
34 MAF 1.0000 0.2749↓
35 L2 0.2749 1.0000↓
36 Analysis finished at Sat Aug 7 15:28:37 2021↓
37 Total time elapsed: 3.8s↓
```

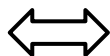
MAFとLD score  
の対応関係

- 各SNPに対して、LD scoreが計算されます。
- マイナーアレル頻度(=minor allele frequency; MAF)とLD scoreは**正の相関**があり、MAFが高い程、LD scoreが大きくなる傾向が認められます。  
(MAFが大きいSNP間において、LD指標 $r^2$ が高い値をとりやすいことに由来します<sup>33</sup>。)

### ③ LDSC解析実習

./Input/

```
1000G_EUR_Phase3_baseline/  
1000G_frq/  
1kg_eur/  
baseline/  
Cahoy_1000Gv3_ldscores/  
cell_type_groups/  
eas_ldscores/  
GWASstats/  
weights_hm3_no_hla/
```



./Output/

```
CellTypeSpecificity/  
Heritability/  
LDscore/
```



./ldsc/

```
.git/  
ContinuousAnnotations/  
ldscore/  
Test/  
.gitignore  
CHANGELOG  
environment.yml  
ldsc.log  
ldsc.py  
ldsc.pyc  
LICENSE  
make_annot.py  
munge_sumstats.py  
README.md  
requirements.txt  
setup.py
```

#### 解析実習内容

1. LD scoreの計算
2. Heritability推定
3. Genetic correlation解析
4. Stratified LDSC解析
5. LDSC-SEG解析

- 「Heritability推定」を行います。
- LD scoreには、東アジア人集団を対象に計算済みの値を使用します。
- ./Input/GWASstats/ に用意した、日本人集団の身長(Height)・肥満 (body mass index; BMI)・2型糖尿病(type 2 diabetes; T2D)のGWAS統計量を使用します。

### ③ LDSC解析実習

```
(ldsc_env) statgen@statgen-PC: /mnt/c/SummerSchool/GenomeDataAnalysis6/Analysis
```

```
$ python ./ldsc/ldsc.py --h2 ./Input/GWASstats/BMI_JPT.sumstats.gz --ref-ld-chr  
./Input/eas_ldscores/ --w-ld-chr ./Input/eas_ldscores/ --out ./Output/Heritability/Herit_BMI_JPT
```

```
$ python ./ldsc/ldsc.py --h2 ./Input/GWASstats/Height_JPT.sumstats.gz --ref-ld-chr  
./Input/eas_ldscores/ --w-ld-chr ./Input/eas_ldscores/ --out ./Output/Heritability/Herit_Height_JPT
```

```
$ python ./ldsc/ldsc.py --h2 ./Input/GWASstats/T2D_JPT.sumstats.gz --ref-ld-chr  
./Input/eas_ldscores/ --w-ld-chr ./Input/eas_ldscores/ --out ./Output/Heritability/Herit_T2D_JPT
```

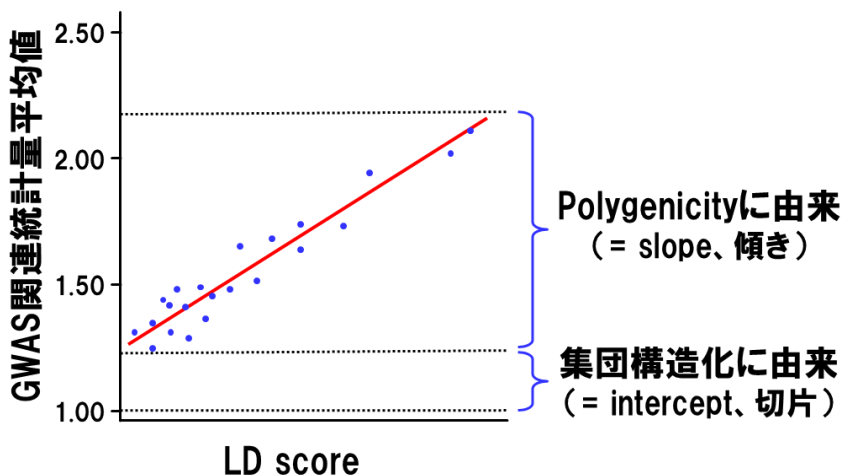
※ファイル”LDSC\_Command.txt”を開いて、内容をLinux  
コマンドにコピー&ペーストして下さい。

- GWAS統計量は、データベース(<https://pheweb.jp/>)からダウンロード後、`munge_sumstats.py`を用いてLDSC解析用にフォーマット変更しています。
- LDSC解析は、ゲノム全域の(全てではなく)主要なコモンバリエーションを対象とし、HapMap3プロジェクトのSNP(約100万箇所)を使用しています。  
(LDSC解析用フォーマットやSNPの内訳の詳細はgithubを参照してください。)

### ③ LDSC解析実習

”Herit\_XXX\_JPT.log”  
各形質のLDSC解析結果

LD score回帰分析



GWAS統計量  
全体のinflation

集団構造化に  
由来するinflation

	Lambda <sub>GC</sub>	Mean X <sup>2</sup>	Intercept	h <sup>2</sup>
BMI	1.49	1.68	1.09	0.169
Height	1.69	2.42	1.21	0.326
T2D	1.13	1.31	0.96	0.097

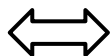
Heritability  
推定値

- GWAS統計量の帰無仮説からの乖離をinflation(=  $\lambda_{GC}$ の1.00からの乖離)と呼びます。LDSC解析により、集団構造化とpolygenicityに由来するinflationを分離することが可能になりました。
- いずれの形質も、inflationに占める集団構造化の影響は小さそうです。
- GWASデータで説明可能なheritabilityは、身長 > 肥満 > 2型糖尿病の順に高いようです。

### ③ LDSC解析実習

./Input/

```
1000G_EUR_Phase3_baseline/  
1000G_frq/  
1kg_eur/  
baseline/  
Cahoy_1000Gv3_ldscores/  
cell_type_groups/  
eas_ldscores/  
GWASstats/  
weights_hm3_no_hla/
```



./Output/

```
CellTypeSpecificity/  
Heritability/  
LDscore/
```



./ldsc/

```
.git/  
ContinuousAnnotations/  
ldscore/  
Test/  
.gitignore  
CHANGELOG  
environment.yml  
ldsc.log  
ldsc.py  
ldsc.pyc  
LICENSE  
make_annot.py  
mungesumstats.py  
README.md  
requirements.txt  
setup.py
```

#### 解析実習内容

1. LD scoreの計算
2. Heritability推定
3. Genetic correlation解析
4. Stratified LDSC解析
5. LDSC-SEG解析

- 「Genetic correlation解析」を行います。
- 引き続き、東アジア集団のLD scoreを使用します。
- 3形質の日本人集団のGWAS統計量(身長・肥満・2型糖尿病)の遺伝的相関をペアワイズで推定します。

### ③ LDSC解析実習

```
(ldsc_env) statgen@statgen-PC: /mnt/c/SummerSchool/GenomeDataAnalysis6/Analysis
```

```
$ python ./ldsc/ldsc.py --rg
```

```
./Input/GWASstats/BMI_JPT.sumstats.gz,./Input/GWASstats/Height_JPT.sumstats.gz --ref-ld-chr  
./Input/eas_ldscores/ --w-ld-chr ./Input/eas_ldscores/ --out ./Output/Heritability/Cor_BMI_Height_JPT
```

```
$ python ./ldsc/ldsc.py --rg
```

```
./Input/GWASstats/BMI_JPT.sumstats.gz,./Input/GWASstats/T2D_JPT.sumstats.gz --ref-ld-chr  
./Input/eas_ldscores/ --w-ld-chr ./Input/eas_ldscores/ --out ./Output/Heritability/Cor_BMI_T2D_JPT
```

```
$ python ./ldsc/ldsc.py --rg
```

```
./Input/GWASstats/Height_JPT.sumstats.gz,./Input/GWASstats/T2D_JPT.sumstats.gz --ref-ld-chr  
./Input/eas_ldscores/ --w-ld-chr ./Input/eas_ldscores/ --out ./Output/Heritability/Cor_Height_T2D_JPT
```

※ファイル”LDSC\_Command.txt”を開いて、内容をLinux  
コマンドにコピー&ペーストして下さい。

• Genetic correlation解析を実施する場合、2形質のGWAS統計量を一つのコマンドライン上で指定します。

### ③ LDSC解析実習

”Cor\_XXX\_XXX\_JPT.log”  
2形質間のLDSC解析結果

”Cor\_BMI\_T2D\_JPT.log”

```
47 Genetic_Covariance↓
48 -----↓
49 Total_Observed_scale_gencov: 0.0339 (0.0063)↓
50 Mean_z1*z2: 0.144↓
51 Intercept: 0.0215 (0.0122)↓
52 ↓
53 Genetic_Correlation↓
54 -----↓
55 Genetic_Correlation: 0.2524 (0.0538)↓
56 Z-score: 4.6937↓
57 P: 2.6834e-06↓
58 ↓
```

P値

遺伝的相関

3形質間の  
遺伝的相関関係

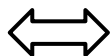
	BMI	Height	T2D
BMI	-	-0.070	0.252
Height	-0.070	-	-0.008
T2D	0.252	-0.008	-

- 各形質ペアごとに、遺伝的相関(genetic correlation)が計算されます。
- 正の相関か、負の相関か、相関は有意かどうか、を確認します。
- 肥満と2型糖尿病の場合、正の遺伝的相関( $r_g=0.252$ )が、有意に確認されました( $P = 2.6 \times 10^{-6}$ )。
- 一方、肥満と身長、2型糖尿病と身長の遺伝的相関は小さいようです。

### ③ LDSC解析実習

./Input/

1000G\_EUR\_Phase3\_baseline/  
1000G\_frq/  
1kg\_eur/  
baseline/  
Cahoy\_1000Gv3\_ldscores/  
cell\_type\_groups/  
eas\_ldscores/  
GWASstats/  
weights\_hm3\_no\_hla/



./Output/

CellTypeSpecificity/  
Heritability/  
LDscore/



./ldsc/

.git/  
ContinuousAnnotations/  
ldscore/  
Test/  
.gitignore  
CHANGELOG  
environment.yml  
ldsc.log  
ldsc.py  
ldsc.pyc  
LICENSE  
make\_annot.py  
munge\_sumstats.py  
README.md  
requirements.txt  
setup.py

#### 解析実習内容

1. LD scoreの計算
2. Heritability推定
3. Genetic correlation解析
4. Stratified LDSC解析
5. LDSC-SEG解析

- 「Stratified LDSC解析」を行います。
- UKバイオバンクの肥満GWAS統計量(UKBB BMI)と、欧米人集団を対象に計算されたLD score、細胞組織別のエピゲノム修飾状況ファイルを使用します。



### ③ LDSC解析実習

```
(ldsc_env) statgen@statgen-PC: /mnt/c/SummerSchool/GenomeDataAnalysis6/Analysis
```

```
$ python ./ldsc/ldsc.py --h2 ./Input/GWASstats/BMI_UKBB.sumstats.gz --w-ld-chr  
./Input/weights_hm3_no_hla/weights. --ref-ld-chr ./Input/cell_type_groups/CNS../Input/baseline/baseline.  
--overlap-annot --frqfile-chr ./Input/1000G_frq/1000G.mac5eur.  
--out ./Output/CellTypeSpecificity/BMI_UKBB_CNS --print-coefficients
```

```
$ python ./ldsc/ldsc.py --h2 ./Input/GWASstats/BMI_UKBB.sumstats.gz --w-ld-chr  
./Input/weights_hm3_no_hla/weights. --ref-ld-chr ./Input/cell_type_groups/Other../Input/baseline/baseline.  
--overlap-annot --frqfile-chr ./Input/1000G_frq/1000G.mac5eur.  
--out./Output/CellTypeSpecificity/BMI_UKBB_Other --print-coefficients
```

※ファイル”LDSC\_Command.txt”を開いて、内容をLinux  
コマンドにコピー&ペーストして下さい。

- **中枢神経系(”central nerve system; CNS”)とその他の細胞組織(”Other”)**  
の二つの細胞組織を対象に、stratified LDSC解析を実施します。
- **実際には、より多くの細胞組織に対して並行してstratified LDSC解析**  
を実施します。

### ③ LDSC解析実習

#### ”BMI\_UKIBB\_XXX.results” Stratified LDSC解析結果

Category	Proportion of SNPs	Proportion of $h^2$ (SE)	Enrichment (SE)	P値
CNS	0.149	0.399 (0.021)	2.68 (0.14)	6.2E-28
Other	0.203	0.323 (0.027)	1.59 (0.13)	7.5E-06

- **CNS特異的なエピゲノム修飾**がゲノム全体のSNPの14.9%を占める一方、それらのSNPはheritabilityの39.9%を説明し、**約2.68倍** (=0.399/0.149)、heritabilityが**enrich**していると考えられました。
- 一方、Other特異的なエピゲノム修飾は、そこまで強くenrichしていないようです。
- 相対的に、CNSの方が疾患病態への関与が強いと考えられます。

### ③ LDSC解析実習

./Input/

1000G\_EUR\_Phase3\_baseline/  
1000G\_frq/  
1kg\_eur/  
baseline/  
Cahoy\_1000Gv3\_ldscores/  
cell\_type\_groups/  
eas\_ldscores/  
GWASstats/  
weights\_hm3\_no\_hla/

./Output/

CellTypeSpecificity/  
Heritability/  
LDscore/

./ldsc/

.git/  
ContinuousAnnotations/  
ldscore/  
Test/  
.gitignore  
CHANGELOG  
environment.yml  
ldsc.log  
ldsc.py  
ldsc.pyc  
LICENSE  
make\_annot.py  
mungesumstats.py  
README.md  
requirements.txt  
setup.py

#### 解析実習内容

1. LD scoreの計算
2. Heritability推定
3. Genetic correlation解析
4. Stratified LDSC解析
5. LDSC-SEG解析

- 「LDSC-SEG解析」を行います。
- UKバイオバンクの肥満GWAS統計量(UKBB BMI)と、欧米人集団を対象に計算されたLD score、細胞組織特異的遺伝子発現情報を使用します。

### ③ LDSC解析実習

```
(ldsc_env) statgen@statgen-PC: /mnt/c/SummerSchool/GenomeDataAnalysis6/Analysis
$ python ./ldsc/ldsc.py --h2-cts ./Input/GWASstats/BMI_UKBB.sumstats.gz --ref-ld-chr
./Input/1000G_EUR_Phase3_baseline/baseline. --out
./Output/CellTypeSpecificity/BMI_UKBB_CellTypeSpecificity --ref-ld-chr-cts
./Input/Cahoy_1000Gv3_ldscores/Cahoy.v2.ldcts --w-ld-chr ./Input/weights_hm3_no_hla/weights.
```

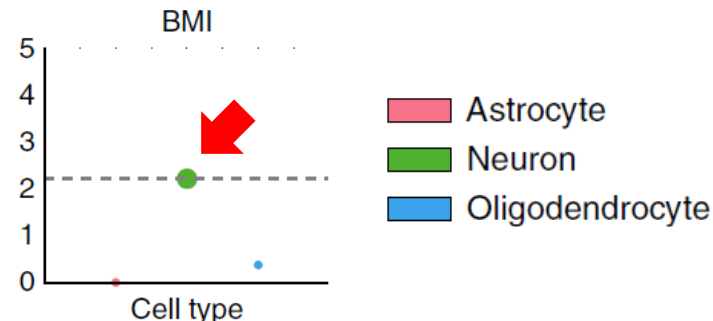
※ファイル”LDSC\_Command.txt”を開いて、内容をLinux  
コマンドにコピー&ペーストして下さい。

- **今回は、マウスの3種類の中枢神経細胞(Astrocyte、Neuron、Oligodendrocyte)における遺伝子発現情報を使用します。**

### ③ LDSC解析実習

” BMI\_UKBB\_CellTypeSpecificity.cell\_type\_results.txt”  
LDSC-SEG解析結果

Cell type	Coefficient (SE)	P値
Neuron	7.9E-09 (3.0E-09)	0.0044
Oligodendrocyte	7.3E-10 (3.5E-09)	0.42
Astrocyte	5.8E-09 (2.6E-09)	0.99



- 肥満のGWAS統計量にLDSC-SEG解析を適用した結果、Astrocyte、Neuron、Oligodendrocyteの3種類の神経細胞の中で、**Neuronの寄与が相対的に大きい事**が確認できました。

# 終わりに

- GWASゲノムデータ解析に際して使用する入力データのうち、GWAS統計量を用いた解析手法として、LDSCを取り上げてみました。
- LDSCは、heritabilityの推定や細胞組織特異性の解明など、**GWAS統計量を用いたデータ解析の応用可能性を広げた**点が画期的でした。
- 現在、GWAS統計量を用いた多くのデータ解析手法が開発されています。
- どのような手法で、どのような解析が可能か、総説等で最新の状況をアップデートしてみると、いいかもしれません。

