# Session 01 - Exercises

**Deepika**

📋 workflowr ❗

Before you begin:

- Make sure that R is installed on your computer
- For this lab, we will use a few R libraries:

Set your working directory to your home directory using in R*

The data files are in the folder `/data/SISG2022M15/data/`.

## Case-Control Association Testing

### Introduction

We will be using the LHON dataset (https://raw.githubusercontent.com/joellembatchou/SISG2022_Association_Mapping/master/data/LHON.txt) covered in the lecture notes for this portion of the exercises. The LHON dataset is from a case-control study and includes both phenotype and genotype data for a candidate gene.

Let's first load the LHON data file into the R session.

```
       vars   n    mean     sd median trimmed     mad min max range   skew kurtosis
IID*      1 328 164.50  94.83  164.5  164.50  121.57   1 328   327   0.00    -1.21
GENO*     2 328   2.68   0.56    3.0    2.78    0.00   1   3     2  -1.54     1.38
PHENO*    3 328   1.73   0.45    2.0    1.78    0.00   1   2     1  -1.02    -0.95
         se
IID*   5.24
GENO*  0.03
PHENO* 0.02
```

### Exercises

Here are some things to look at:

1. Examine the variables in the dataset:

- How many observations?

```
[1] 328
```

- How many cases/controls?

```
     CASE CONTROL
       89     239
```

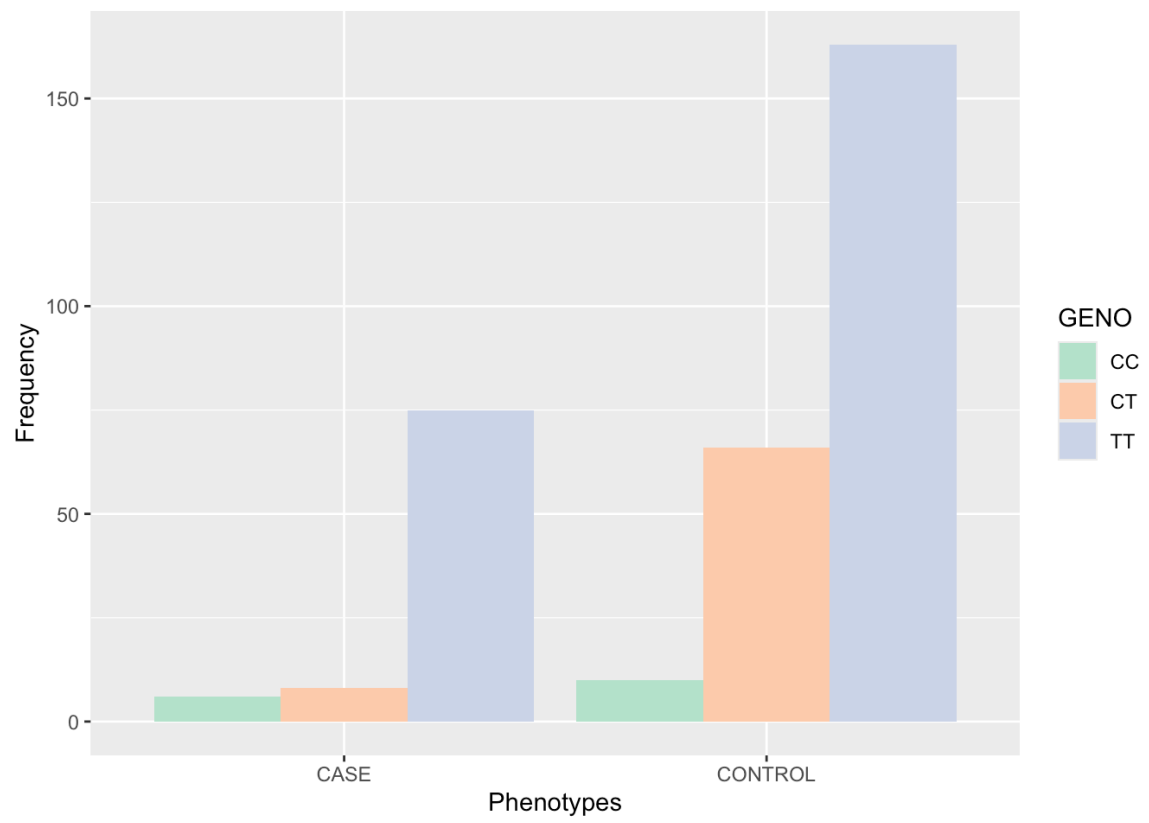- What is the distribution of the genotypes across cases/controls?

```
   Cell Contents
|-----------------------|
|                     N |
| Chi-square contribution |
|          N / Row Total |
|          N / Col Total |
|        N / Table Total |
|-----------------------|


Total Observations in Table:  328


          | GENO
    PHENO |        CC |        CT |        TT | Row Total |
----------|-----------|-----------|-----------|-----------|
     CASE |         6 |         8 |        75 |        89 |
          |     0.634 |     7.267 |     1.682 |           |
          |     0.067 |     0.090 |     0.843 |     0.271 |
          |     0.375 |     0.108 |     0.315 |           |
          |     0.018 |     0.024 |     0.229 |           |
----------|-----------|-----------|-----------|-----------|
  CONTROL |        10 |        66 |       163 |       239 |
          |     0.236 |     2.706 |     0.626 |           |
          |     0.042 |     0.276 |     0.682 |     0.729 |
          |     0.625 |     0.892 |     0.685 |           |
          |     0.030 |     0.201 |     0.497 |           |
----------|-----------|-----------|-----------|-----------|
Column Total |     16 |        74 |       238 |       328 |
          |     0.049 |     0.226 |     0.726 |           |
----------|-----------|-----------|-----------|-----------|
```

```
Warning: The dot-dot notation (`..count..`) was deprecated in ggplot2 3.4.0.
ℹ Please use `after_stat(count)` instead.
This warning is displayed once every 8 hours.
Call `lifecycle::last_lifecycle_warnings()` to see where this warning was
generated.
```
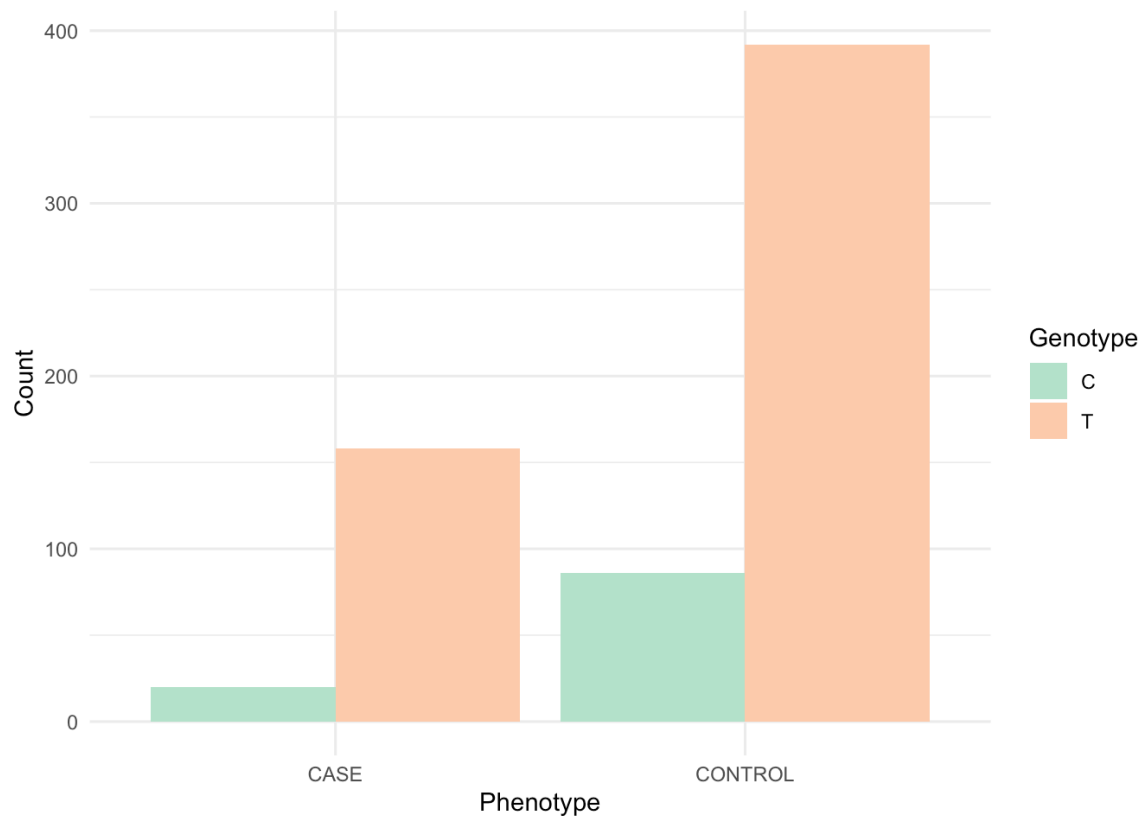
- What about for allele types?

```
   Cell Contents
|-----------------------|
|                     N |
| Chi-square contribution |
|         N / Row Total |
|         N / Col Total |
|       N / Table Total |
|-----------------------|


Total Observations in Table:  656


              |
      PHENO  |         C |         T | Row Total |
------------|-----------|-----------|-----------|
      CASE  |        20 |       158 |       178 |
            |     2.669 |     0.514 |           |
            |     0.112 |     0.888 |     0.271 |
            |     0.189 |     0.287 |           |
            |     0.030 |     0.241 |           |
------------|-----------|-----------|-----------|
   CONTROL  |        86 |       392 |       478 |
            |     0.994 |     0.192 |           |
            |     0.180 |     0.820 |     0.729 |
            |     0.811 |     0.713 |           |
            |     0.131 |     0.598 |           |
------------|-----------|-----------|-----------|
Column Total |       106 |       550 |       656 |
            |     0.162 |     0.838 |           |
------------|-----------|-----------|-----------|
```

2. Perform a logistic regression analysis for this data with `CC` as the reference genotype using the
   `glm()` function. (Hint: make sure to convert the phenotype to a binary 0/1 variable and specify
   `family = binomial(link = "logit")` in the `glm` call)

```
    CASE CONTROL
 0    0     239
 1   89       0
```

```
Call:
glm(formula = BinPheno ~ GENO, family = binomial(link = "logit"),
    data = LHON.df)

Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept)  -0.5108     0.5164  -0.989   0.3226
GENOCT       -1.5994     0.6378  -2.508   0.0122 *
GENOTT       -0.2654     0.5349  -0.496   0.6197
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 383.49  on 327  degrees of freedom
Residual deviance: 368.48  on 325  degrees of freedom
AIC: 374.48

Number of Fisher Scoring iterations: 4
```

3. Obtain odds ratios and confidence intervals for the `CT` and `TT` genotypes relative to the `CC` reference
   genotype. Interpret.

```
(Intercept)      GENOCT      GENOTT
  0.6000000   0.2020202   0.7668712
```

```
               2.5 %     97.5 %
(Intercept) 0.20413356 1.6156811
GENOCT      0.05710635 0.7223515
GENOTT      0.27431485 2.3258908
```

4. Is there evidence of differences in odds of being a case for the `CT` and `TT` genotypes (compared to `CC`)?

```
                  OR      2.5 %     97.5 %    p_value
(Intercept) 0.6000000 0.20413356 1.6156811 0.32256061
GENOCT      0.2020202 0.05710635 0.7223515 0.01215534
GENOTT      0.7668712 0.27431485 2.3258908 0.61973850
```

*Extra*: 5. Perform the logistic regression analysis with the additive genotype coding. Obtain odds ratios and confidence intervals. Is there evidence of an association? How does it compare with the 2-parameter model?

```
Call:
glm(formula = BinPheno ~ Dosage, family = binomial(link = "logit"),
    data = LHON.df)

Coefficients:
            Estimate Std. Error z value Pr(>|z|)
(Intercept)  -1.8077     0.4554  -3.970  7.2e-05 ***
Dosage        0.4787     0.2505   1.911   0.0559 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

(Dispersion parameter for binomial family taken to be 1)

    Null deviance: 383.49  on 327  degrees of freedom
Residual deviance: 379.47  on 326  degrees of freedom
AIC: 383.47

Number of Fisher Scoring iterations: 4
```

```
(Intercept)     Dosage
  0.1640322   1.6140439
```

```
Waiting for profiling to be done...
```

```
               2.5 %     97.5 %
(Intercept) 0.06293774 0.3801326
Dosage      1.01029266 2.7133859
```

```
Waiting for profiling to be done...
```

```
                    OR      2.5 %     97.5 %      p_value
(Intercept) 0.1640322 0.06293774 0.3801326 7.201892e-05
Dosage      1.6140439 1.01029266 2.7133859 5.594566e-02
```

# Association Testing with Quantitative Traits

## Introduction

We will be using the Blood Pressure dataset
(https://raw.githubusercontent.com/joellembatchou/SISG2022_Association_Mapping/master/data/bpdata.csv)
for this portion of the exercises. This dataset contains diastolic and systolic blood pressure measurements
for 1000 individuals, and genotype data at 11 SNPs in a candidate gene for blood pressure. Covariates such
as gender (sex) and body mass index (bmi) are included as well.

Let's first load the file into R.

```
        vars    n   mean     sd median trimmed    mad min  max range  skew
V1         1 1000 500.50 288.82  500.5  500.50 370.65   1 1000   999  0.00
sex*       2 1000   1.53   0.50    2.0    1.54   0.00   1    2     1 -0.13
sbp        3 1000 141.42  18.47  140.0  140.72  17.79  87  202   115  0.35
dbp        4 1000  82.61  10.84   82.0   82.60  10.38  47  117    70 -0.03
snp1*      5  988   2.31   0.67    2.0    2.39   1.48   1    3     2 -0.46
snp2*      6  978   2.65   0.54    3.0    2.72   0.00   1    3     2 -1.19
snp3*      7  960   1.39   0.56    1.0    1.32   0.00   1    3     2  1.07
snp4*      8  928   1.75   0.70    2.0    1.69   1.48   1    3     2  0.38
snp5*      9  916   1.20   0.43    1.0    1.11   0.00   1    3     2  1.94
snp6*     10  929   1.52   0.64    1.0    1.42   0.00   1    3     2  0.84
snp7*     11  986   2.44   0.65    3.0    2.53   0.00   1    3     2 -0.73
snp8*     12  984   1.35   0.55    1.0    1.26   0.00   1    3     2  1.29
snp9*     13  966   2.69   0.53    3.0    2.78   0.00   1    3     2 -1.45
snp10*    14  978   1.91   0.71    2.0    1.89   1.48   1    3     2  0.13
snp11*    15  979   2.69   0.51    3.0    2.76   0.00   1    3     2 -1.28
bmi       16  999  30.26   6.23   29.0   29.74   5.93  16   51    35  0.83
       kurtosis   se
V1        -1.20 9.13
sex*      -1.99 0.02
sbp        0.15 0.58
dbp        0.32 0.34
snp1*     -0.80 0.02
snp2*      0.42 0.02
snp3*      0.14 0.02
snp4*     -0.92 0.02
snp5*      2.91 0.01
snp6*     -0.35 0.02
snp7*     -0.51 0.02
snp8*      0.67 0.02
snp9*      1.15 0.02
snp10*    -1.03 0.02
snp11*     0.59 0.02
bmi        0.70 0.20
```

## Exercises

Here are some things to try:

1. Perform a linear regression of systolic blood pressure ( sbp ) on SNP3 using the lm() function. Compare the estimates, confidence intervals and p-values you get using:

- additive (linear) model

```
 CC  TC  TT
621 304  35
```

```
Call:
lm(formula = sbp ~ snp3Dosage, data = BP.df)

Residuals:
    Min      1Q  Median      3Q     Max
-55.974 -12.418  -0.974  10.582  60.582

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 140.4179     0.7219 194.506   <2e-16 ***
snp3Dosage    2.5556     1.0615   2.407   0.0163 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.33 on 958 degrees of freedom
  (40 observations deleted due to missingness)
Multiple R-squared:  0.006014,  Adjusted R-squared:  0.004976
F-statistic: 5.796 on 1 and 958 DF,  p-value: 0.01625
```

```
(Intercept)  snp3Dosage
 140.417909    2.555635
```

```
                  2.5 %      97.5 %
(Intercept) 139.0011786 141.834639
snp3Dosage    0.4724342   4.638837
```

```
             Estimate       2.5 %      97.5 %    p_value
(Intercept) 140.417909 139.0011786 141.834639 0.00000000
snp3Dosage    2.555635   0.4724342   4.638837 0.01625073
```

- dominant model

```
    CC  TC  TT
0 621   0   0
1   0 304  35
```

```
Call:
lm(formula = sbp ~ snp3Dom, data = BP.df)

Residuals:
    Min      1Q  Median      3Q     Max
-56.218 -12.428  -0.823  10.572  60.572

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  140.428      0.736 190.801   <2e-16 ***
snp3Dom        2.790      1.238   2.253   0.0245 *
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.34 on 958 degrees of freedom
  (40 observations deleted due to missingness)
Multiple R-squared:  0.005269,  Adjusted R-squared:  0.00423
F-statistic: 5.074 on 1 and 958 DF,  p-value: 0.02451
```

```
(Intercept)      snp3Dom
 140.428341     2.789948
```

```
               2.5 %      97.5 %
(Intercept) 138.9839938 141.872689
snp3Dom       0.3593814   5.220514
```

```
            Estimate       2.5 %      97.5 %     p_value
(Intercept) 140.428341 138.9839938 141.872689 0.00000000
snp3Dom       2.789948   0.3593814   5.220514 0.02450948
```

- recessive model

```
    CC  TC  TT
 0 621 304   0
 1   0   0  35
```

```
Call:
lm(formula = sbp ~ snp3Rec, data = BP.df)

Residuals:
    Min      1Q  Median      3Q     Max
-54.251 -12.501  -1.251  10.749  59.749

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept)  141.251      0.604 233.854   <2e-16 ***
snp3Rec        4.463      3.163   1.411    0.159
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.37 on 958 degrees of freedom
  (40 observations deleted due to missingness)
Multiple R-squared:  0.002074,  Adjusted R-squared:  0.001032
F-statistic: 1.991 on 1 and 958 DF,  p-value: 0.1586
```

```
(Intercept)      snp3Rec
 141.250811     4.463475
```

```
               2.5 %     97.5 %
(Intercept) 140.065471 142.43615
snp3Rec      -1.744423  10.67137
```

```
            Estimate      2.5 %    97.5 %    p_value
(Intercept) 141.250811 140.065471 142.43615 0.0000000
snp3Rec       4.463475  -1.744423  10.67137 0.1585706
```

- 2 parameter model

```
 CC  TC  TT
621 304  35
```

```
Call:
lm(formula = sbp ~ snp3, data = BP.df)

Residuals:
    Min      1Q  Median      3Q     Max
-55.931 -12.428  -0.931  10.572  60.572

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 140.4283     0.7361 190.773   <2e-16 ***
snp3TC        2.5026     1.2840   1.949   0.0516 .
snp3TT        5.2859     3.1868   1.659   0.0975 .
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.34 on 957 degrees of freedom
  (40 observations deleted due to missingness)
Multiple R-squared:  0.006019,  Adjusted R-squared:  0.003942
F-statistic: 2.898 on 2 and 957 DF,  p-value: 0.05563
```
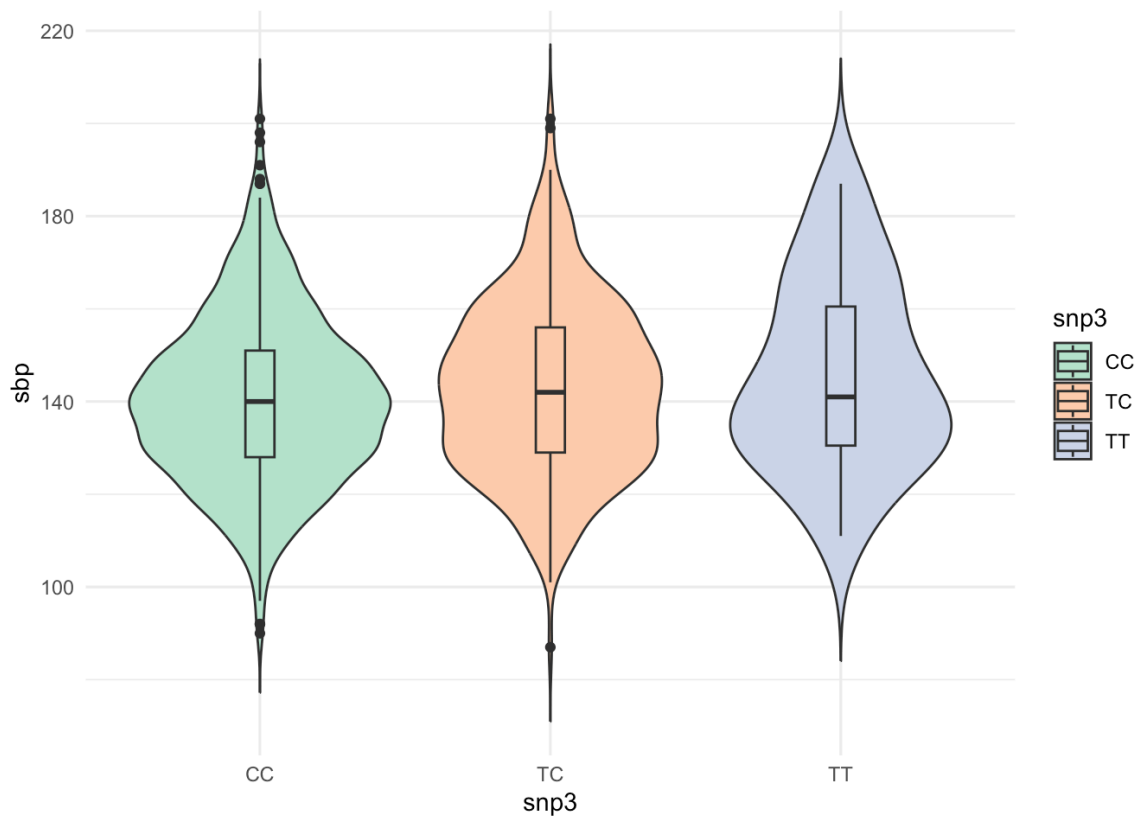
```
(Intercept)       snp3TC       snp3TT
 140.428341     2.502580     5.285944
```

```
                 2.5 %      97.5 %
(Intercept) 138.98378278 141.872900
snp3TC       -0.01723871   5.022398
snp3TT       -0.96798634  11.539875
```

```
            Estimate        2.5 %      97.5 %    p_value
(Intercept) 140.428341 138.98378278 141.872900 0.00000000
snp3TC        2.502580  -0.01723871   5.022398 0.05158487
snp3TT        5.285944  -0.96798634  11.539875 0.09750441
```

(Hint: for each case, first add a new column to the data frame, containing the 'predictor' variable you need.
Then do the regression using lm())

2. Provide a plot illustrating the relationship between sbp and the three genotypes at SNP3.

For question 3 and 4 below, R also has a 'formula' syntax, frequently used when specifying regression models with many predictors. To regress an outcome `y` on several covariates, the syntax is:

```
outcome ~ covariate1 + covariate2 + covariate3
```

3. Now redo the linear regression analysis of `sbp` from question 1 for the additive model, but this time adjust for `sex` and `bmi`. Do the results change?

```
Call:
lm(formula = sbp ~ snp3Dosage + sex + bmi, data = BP.df)

Residuals:
   Min     1Q Median     3Q    Max
-58.83 -12.81  -0.82  11.58  57.80

Coefficients:
             Estimate Std. Error t value Pr(>|t|)
(Intercept) 145.85380    3.00271  48.574  < 2e-16 ***
snp3Dosage    2.63566    1.05434   2.500   0.0126 *
sexMALE      -4.77580    1.17642  -4.060 5.32e-05 ***
bmi          -0.09837    0.09481  -1.038   0.2997
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.19 on 955 degrees of freedom
  (41 observations deleted due to missingness)
Multiple R-squared:  0.02402,   Adjusted R-squared:  0.02096
F-statistic: 7.836 on 3 and 955 DF,  p-value: 3.608e-05
```

```
 (Intercept)    snp3Dosage        sexMALE            bmi
145.85379689    2.63566473    -4.77580021    -0.09837466
```

```
                 2.5 %        97.5 %
(Intercept) 139.9611289 151.74646484
snp3Dosage    0.5665706   4.70475889
sexMALE      -7.0844766  -2.46712378
bmi          -0.2844395   0.08769022
```

Results with Covariates

```
               Estimate        2.5 %        97.5 %       p_value
(Intercept) 145.85379689 139.9611289 151.74646484 2.761186e-260
snp3Dosage    2.63566473   0.5665706   4.70475889 1.259244e-02
sexMALE      -4.77580021  -7.0844766  -2.46712378 5.318310e-05
bmi          -0.09837466  -0.2844395   0.08769022 2.997326e-01
```

Results without Covariates

```
              Estimate       2.5 %      97.5 %    p_value
(Intercept) 140.417909 139.0011786 141.834639 0.00000000
snp3Dosage    2.555635   0.4724342   4.638837 0.01625073
```

4. What proportion of the heritability of sbp is explained by all 11 SNPs combined? (contrast categorical coding vs additive coding for the genotypes)

Categorical Coding

```
Call:
lm(formula = sbp ~ snp1 + snp2 + snp3 + snp4 + snp5 + snp6 +
    snp7 + snp8 + snp9 + snp10 + snp11, data = BP.df)

Residuals:
    Min      1Q  Median      3Q     Max
-50.722 -11.967  -0.703  11.021  61.704

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 133.1726    12.4033  10.737   <2e-16 ***
snp1CT       -1.7048     4.5991  -0.371    0.711
snp1TT        1.9319     8.2839   0.233    0.816
snp2AT        0.7347     5.5923   0.131    0.896
snp2TT       -0.5118     6.9317  -0.074    0.941
snp3TC        4.7672     5.0211   0.949    0.343
snp3TT        6.6913     9.7904   0.683    0.495
snp4CT       -0.4778     3.5501  -0.135    0.893
snp4TT        2.3431     6.4874   0.361    0.718
snp5CT        1.1896     3.0462   0.391    0.696
snp5TT       -2.2787     7.5490  -0.302    0.763
snp6AG       -3.0266     2.0697  -1.462    0.144
snp6GG        2.1230     4.6650   0.455    0.649
snp7AT       -3.0873     3.9148  -0.789    0.431
snp7TT       -2.6319     4.3146  -0.610    0.542
snp8CT       -1.5509     3.6318  -0.427    0.669
snp8TT       -2.5507     7.3228  -0.348    0.728
snp9CT        6.0693     7.6170   0.797    0.426
snp9TT        4.7385     7.4517   0.636    0.525
snp10CT       1.4330     1.6466   0.870    0.384
snp10TT       1.9810     2.0699   0.957    0.339
snp11CT       4.8005     6.5175   0.737    0.462
snp11TT       4.0226     9.2775   0.434    0.665
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.2 on 707 degrees of freedom
  (270 observations deleted due to missingness)
Multiple R-squared:  0.02633,   Adjusted R-squared:  -0.003965
F-statistic: 0.8691 on 22 and 707 DF,  p-value: 0.6372
```

```
             Estimate       2.5 %      97.5 %      p_value
(Intercept) 133.1725671 108.820890 157.524245 5.131083e-25
snp1CT       -1.7048318 -10.734443   7.324779 7.109834e-01
snp1TT        1.9318980 -14.332013  18.195809 8.156641e-01
snp2AT        0.7346616 -10.244868  11.714191 8.955201e-01
snp2TT       -0.5118275 -14.121058  13.097403 9.411599e-01
snp3TC        4.7671827  -5.090918  14.625283 3.427288e-01
snp3TT        6.6912706 -12.530460  25.913001 4.945449e-01
snp4CT       -0.4777530  -7.447752   6.492246 8.929866e-01
snp4TT        2.3430912 -10.393724  15.079906 7.180747e-01
snp5CT        1.1896284  -4.791101   7.170358 6.962657e-01
snp5TT       -2.2787219 -17.099774  12.542330 7.628481e-01
snp6AG       -3.0265667  -7.090080   1.036946 1.440994e-01
snp6GG        2.1230367  -7.035788  11.281861 6.491746e-01
snp7AT       -3.0872887 -10.773220   4.598642 4.305930e-01
snp7TT       -2.6319274 -11.102829   5.838974 5.420516e-01
snp8CT       -1.5509162  -8.681351   5.579519 6.694831e-01
snp8TT       -2.5507399 -16.927782  11.826302 7.276973e-01
snp9CT        6.0693019  -8.885403  21.024006 4.258307e-01
snp9TT        4.7385208  -9.891650  19.368692 5.250505e-01
snp10CT       1.4329754  -1.799872   4.665823 3.844573e-01
snp10TT       1.9809712  -2.082851   6.044793 3.388680e-01
snp11CT       4.8005159  -7.995478  17.596510 4.616364e-01
snp11TT       4.0225845 -14.192153  22.237323 6.647219e-01
```

Let's check the model if we had used additive coding for all SNPs.

SNP 1

```
       0    1    2
 CC  119    0    0
 CT    0  444    0
 TT    0    0  425
```
SNP 2

```
       0    1    2
 AA   30    0    0
 AT    0  285    0
 TT    0    0  663
```
SNP 3

```
       0    1    2
 CC  621    0    0
 TC    0  304    0
 TT    0    0   35
```
SNP 4

```
       0    1    2
 CC  368    0    0
 CT    0  421    0
 TT    0    0  139
```
SNP 5

```
       0    1    2
 CC  742    0    0
 CT    0  162    0
 TT    0    0   12
```
SNP 6

```
       0    1    2
 AA  521    0    0
 AG    0  335    0
 GG    0    0   73
```
SNP 7

```
       0    1    2
 AA   85    0    0
 AT    0  381    0
 TT    0    0  520
```
SNP 8

```
       0    1    2
 CC  678    0    0
 CT    0  271    0
 TT    0    0   35
```
SNP 9

```
       0    1    2
 CC   30    0    0
 CT    0  239    0
 TT    0    0  697
```
SNP 10

```
        0   1   2
  CC 296   0   0
  CT   0 475   0
  TT   0   0 207
SNP 11

        0   1   2
  CC  20   0   0
  CT   0 264   0
  TT   0   0 695
```

```
Call:
lm(formula = sbp ~ snp1Dosage + snp2Dosage + snp3Dosage + snp4Dosage +
    snp5Dosage + snp6Dosage + snp7Dosage + snp8Dosage + snp9Dosage +
    snp10Dosage + snp11Dosage, data = BP.df.dos)

Residuals:
    Min      1Q  Median      3Q     Max
-53.638 -12.849  -0.522  11.032  61.683

Coefficients:
            Estimate Std. Error t value Pr(>|t|)
(Intercept) 137.36286    9.09141  15.109   <2e-16 ***
snp1Dosage    1.88456    4.03838   0.467    0.641
snp2Dosage   -1.95639    2.96674  -0.659    0.510
snp3Dosage    4.60730    4.65652   0.989    0.323
snp4Dosage    0.05946    3.11138   0.019    0.985
snp5Dosage   -0.26494    2.58719  -0.102    0.918
snp6Dosage   -1.17284    1.80185  -0.651    0.515
snp7Dosage   -0.28939    1.78362  -0.162    0.871
snp8Dosage    0.70702    2.78030   0.254    0.799
snp9Dosage    2.17197    2.54774   0.853    0.394
snp10Dosage   0.60685    1.01229   0.599    0.549
snp11Dosage  -0.39009    4.15347  -0.094    0.925
---
Signif. codes:  0 '***' 0.001 '**' 0.01 '*' 0.05 '.' 0.1 ' ' 1

Residual standard error: 18.17 on 718 degrees of freedom
  (270 observations deleted due to missingness)
Multiple R-squared:  0.01418,   Adjusted R-squared:  -0.0009268
F-statistic: 0.9386 on 11 and 718 DF,  p-value: 0.5022
```

```
              Estimate       2.5 %      97.5 %      p_value
(Intercept) 137.36286195 119.513941 155.211783 5.463978e-45
snp1Dosage    1.88455556  -6.043884   9.812996 6.408836e-01
snp2Dosage   -1.95638699  -7.780906   3.868132 5.098245e-01
snp3Dosage    4.60729872  -4.534719  13.749316 3.227861e-01
snp4Dosage    0.05945758  -6.049037   6.167953 9.847589e-01
snp5Dosage   -0.26493620  -5.344303   4.814431 9.184654e-01
snp6Dosage   -1.17284174  -4.710367   2.364684 5.153132e-01
snp7Dosage   -0.28939265  -3.791133   3.212348 8.711547e-01
snp8Dosage    0.70701756  -4.751477   6.165512 7.993403e-01
snp9Dosage    2.17196761  -2.829934   7.173869 3.942159e-01
snp10Dosage   0.60685115  -1.380549   2.594251 5.490393e-01
snp11Dosage  -0.39009298  -8.544499   7.764313 9.251992e-01
```

🔧 Session information