

## #UNIX Assignment

### ##Data Inspection

#### ###Attributes of fang\_et\_al\_genotypes

here is my snippet of code used for data inspection:

```
ls -lh fang_et_al_genotypes.txt
```

```
file fang_et_al_genotypes.txt
```

```
head -n 10 fang_et_al_genotypes.txt
```

```
tail -n 10 fang_et_al_genotypes.txt
```

```
less fang_et_al_genotypes.txt
```

```
wc fang_et_al_genotypes.txt
```

```
wc -l fang_et_al_genotypes.txt
```

```
wc -w fang_et_al_genotypes.txt
```

```
wc -c fang_et_al_genotypes.txt
```

```
awk -F'\t' '{print NF; exit}' fang_et_al_genotypes.txt
```

```
awk -F'\t' " fang_et_al_genotypes.txt | sort -u
```

```
head -n 1 fang_et_al_genotypes.txt
```

```
cut -f1 fang_et_al_genotypes.txt | sort | uniq | head -n 20
```

```
cut -f1 fang_et_al_genotypes.txt | sort | uniq | wc -l
```

```
grep -E "\bNA\b|\b.\b|\b?\b" fang_et_al_genotypes.txt | head -n 10
```

```
cut -f1 fang_et_al_genotypes.txt | sort | uniq -d
```

By inspecting this file I learned that:

1. The file size is 11 MB
2. It is an ASCII text with very long lines
3. it has 2783 Lines
4. It has 2744038 words
5. It has 11051939 Bytes
6. It has 986 columns

#### ###Attributes of snp\_position.txt

here is my snippet of code used for data inspection:

```
ls -lh snp_position.txt

file snp_position.txt

head -n 10 snp_position.txt

tail -n 10 snp_position.txt

less snp_position.txt

wc snp_position.txt

wc -l snp_position.txt

wc -w snp_position.txt

wc -c snp_position.txt

awk -F'\t' '{print NF; exit}' snp_position.txt

awk -F'\t' " snp_position.txt | sort -u

head -n 1 snp_position.txt

cut -f1 snp_position.txt | sort | uniq | head -n 20

cut -f1 snp_position.txt | sort | uniq | wc -l

cut -f1 snp_position.txt | sort | uniq -d
```

By inspecting this file I learned that:

1. The file size is 81 KB
2. It is an ASCII text
3. it has 984 Lines
4. It has 13198 words
5. It has 82763 Bytes
6. It has 15 columns

##Data Processing

###Maize Data

Step 1: Extract ZMMIL, ZMMLR, and ZMMMR from the Group column `awk '$3 ~ /Group|ZMMIL|ZMMLR|ZMMMR/' fang_et_al_genotypes.txt > maize_data.txt`

Step 2: Transpose the extracted data `awk -f transpose.awk maize_data.txt > maize_transposed.txt`

Step 3: Sort the transposed file and add the header

```
head -n 1 maize_transposed.txt > header.txt tail -n +4 maize_transposed.txt | sort -k1,1 > sorted_maize.txt cat header.txt  
sorted_maize.txt > sorted_maize_with_header.txt
```

Step 4: Process SNP position data head -n 1 snp\_position.txt > snp\_header.txt tail -n +2 snp\_position.txt | sort -k1,1 > sorted\_snp.txt cat snp\_header.txt sorted\_snp.txt > sorted\_snp\_with\_header.txt cut -f 1,3,4 sorted\_snp\_with\_header.txt > snp\_trimmed.txt

Step 5: Join SNP data with maize data sed 's/Sample\_ID/SNP\_ID/' sorted\_maize\_with\_header.txt > maize\_final.txt join -1 1 -2 1 -t \$'\t' snp\_trimmed.txt maize\_final.txt > maize\_joined.txt

Step 6: Extract unknown and multiple positions grep -E "(Chromosome|unknown)" maize\_joined.txt > maize\_unknown.txt grep -E "(Chromosome|multiple)" maize\_joined.txt > maize\_multiple.txt

Step 7: Sort by increasing position and replace missing data head -n 1 maize\_joined.txt > maize\_header.txt tail -n +2 maize\_joined.txt | sort -k3,3n > maize\_sorted\_asc.txt cat maize\_header.txt maize\_sorted\_asc.txt | sed 's!?!?!?!g' > maize\_final\_asc.txt

Step 8: Sort by decreasing position and replace missing data tail -n +2 maize\_joined.txt | sort -k3,3nr > maize\_sorted\_desc.txt cat maize\_header.txt maize\_sorted\_desc.txt | sed 's!?!?!?!g' > maize\_final\_desc.txt

Step 9: Extract chromosomes (using a loop for efficiency) for i in {1..10}; do awk -v chr="\$i" '\(2 == chr' maize\_final\_asc.txt > maize\_chr\)\_asc.txt awk -v chr="\$i" '\(2 == chr' maize\_final\_desc.txt > maize\_chr\)\_desc.txt done

Here is my brief description of what this code does: I wrote brief description of each line in previous section

### ###Teosinte Data

here is my snippet of code used for data processing Step 1: Extract ZMPBA, ZMPIL, and ZMPJA from the Group column  
awk '\$3 ~ /Group|ZMPBA|ZMPIL|ZMPJA/' fang\_et\_al\_genotypes.txt > teosinte\_data.txt

Step 2: Transpose the extracted data awk -f transpose.awk teosinte\_data.txt > teosinte\_transposed.txt

Step 3: Sort the transposed file and add the header head -n 1 teosinte\_transposed.txt > header.txt tail -n +4 teosinte\_transposed.txt | sort -k1,1 > sorted\_teosinte.txt cat header.txt sorted\_teosinte.txt > sorted\_teosinte\_with\_header.txt

Step 4: Join SNP data with teosinte data sed 's/Sample\_ID/SNP\_ID/' sorted\_teosinte\_with\_header.txt > teosinte\_final.txt join -1 1 -2 1 -t \$'\t' snp\_trimmed.txt teosinte\_final.txt > teosinte\_joined.txt

Step 5: Extract unknown and multiple positions grep -E "(Chromosome|unknown)" teosinte\_joined.txt > teosinte\_unknown.txt grep -E "(Chromosome|multiple)" teosinte\_joined.txt > teosinte\_multiple.txt

Step 6: Sort by increasing position and replace missing data head -n 1 teosinte\_joined.txt > teosinte\_header.txt tail -n +2 teosinte\_joined.txt | sort -k3,3n > teosinte\_sorted\_asc.txt cat teosinte\_header.txt teosinte\_sorted\_asc.txt | sed 's!?!?!?!g' > teosinte\_final\_asc.txt

Step 7: Sort by decreasing position and replace missing data `tail -n +2 teosinte_joined.txt | sort -k3,3nr > teosinte_sorted_desc.txt` `cat teosinte_header.txt teosinte_sorted_desc.txt | sed 's!?!?!-!g' > teosinte_final_desc.txt`

Step 8: Extract chromosomes (using a loop for efficiency)

```
for i in {1..10}; do awk -v chr="$i" '\(2 == chr' teosinte_final_asc.txt > teosinte_chr\)_asc.txt awk -v chr="$i" '\(2 == chr'
teosinte_final_desc.txt > teosinte_chr\)_desc.txt done
```

Here is my brief description of what this code does: You can Find Breif Description of each line above that.

Following steps are for making folders for better undeerstanding:

Step 1: Create Folder Structure

Create Maize and Teosinte folders

```
mkdir -p Maize/increasing Maize/decreasing
```

```
mkdir -p Teosinte/increasing Teosinte/decreasing
```

Step 2: Move Maize Files

Move increasing position files for Maize

```
mv maize_chr*_asc.txt Maize/increasing/
```

Move decreasing position files for Maize

```
mv maize_chr*_desc.txt Maize/decreasing/
```

Move unknown and multiple files for Maize

```
mv maize_unknown.txt maize_multiple.txt Maize/
```

Step 3: Move Teosinte Files

Move increasing position files for Teosinte

```
mv teosinte_chr*_asc.txt Teosinte/increasing/
```

Move decreasing position files for Teosinte

```
mv teosinte_chr*_desc.txt Teosinte/decreasing/
```

Move unknown and multiple files for Teosinte

```
mv teosinte_unknown.txt teosinte_multiple.txt Teosinte/
```

Step 4: Create a Folder for Temporary Files

Create a folder named "temp\_files" in the current directory

```
mkdir -p temp_files
```

#### Step 5: Move Temporary Files into the Folder

Move all intermediate maize-related files

```
mv maize_data.txt maize_transposed.txt sorted_maize.txt header.txt maize_final.txt maize_joined.txt  
maize_sorted_asc.txt
```

```
maize_sorted_desc.txt maize_header.txt temp_files/
```

Move all intermediate teosinte-related files

```
mv teosinte_data.txt teosinte_transposed.txt sorted_teosinte.txt teosinte_final.txt teosinte_joined.txt  
teosinte_sorted_asc.txt
```

```
teosinte_sorted_desc.txt teosinte_header.txt temp_files/
```

Move SNP-related intermediate files

```
mv snp_header.txt sorted_snp.txt sorted_snp_with_header.txt snp_trimmed.txt temp_files/
```