

NAME: ALAMELU
ROLL NUMBER:25210013
LAB ASSIGNMENT – 2
Linux & Shell Scripting with Biological Data Files

PART-1

1. Open a new file called notes.txt in vi.

- Insert exactly one line of text:

Have a nice day

(Make sure there is no trailing space at the end.)

- Save and exit.

- Verify that the file contains exactly one line and 15 characters.

```
user@DESKTOP-A5PNEA5:/mnt/c/Users/user$ vi notes.txt
user@DESKTOP-A5PNEA5:/mnt/c/Users/user$ wc -l notes.txt
1 notes.txt
user@DESKTOP-A5PNEA5:/mnt/c/Users/user$ wc -m notes.txt
16 notes.txt
user@DESKTOP-A5PNEA5:/mnt/c/Users/user$ vi notes.txt
user@DESKTOP-A5PNEA5:/mnt/c/Users/user$ wc -m notes.txt
16 notes.txt
user@DESKTOP-A5PNEA5:/mnt/c/Users/user$ |
```

PART-2

2. Display the last four lines of sequence.fasta without opening the file in an editor.

```
user@DESKTOP-A5PNEA5:/mnt/c/Users/user$ cd /mnt/d/MTECH/SEM1/Biocomputing
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing$ cd Lab1
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab1$ tail -n 4 sequence.fasta
TAACTACTGATAAGTTACAAAAGTGTCTATCCTAAAGGGCAATACAGCCCTAGACTCTCCCAGGTAT
TTGACTCCTGCAGCAAAAAGGAAATTGAGGAAATAGAGCAAGCTATTTCTCAGAGGCAACTATATCACA
TAGACACCCCG
```

3. In sequence5.fasta, print all header lines (lines starting with >).

```
user@DESKTOP-A5PNEA5:/mnt/c/Users/user$ cd /mnt/d/MTECH/SEM1/Biocomputing/Lab2
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab2$ grep -c ">" sequence5.f
asta
13
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab2$ grep ">" sequence5.fast
a
>ahr
>clock
>hif1a
>hif2a
>hif3a
>npas1
>npas2
>npas3
>npas4
>sim1
>sim2
>arnt1
>bmal1
```

NAME: ALAMELU
ROLL NUMBER:25210013
LAB ASSIGNMENT – 2
Linux & Shell Scripting with Biological Data Files

4. Find all matches in sequence5.fasta where A is followed by any single character and then G.

```
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab2$ grep "A.G" sequence5.fasta
IFRTKHKLDFTPIGCDAKGRIVLGYTEAELCTRGSGYQFIHAADMLYCAESHIRMIKTGESGMIVFRLLT
DAARSRSSQETEVLYQLAHTLPFARGVSAHLDKASIMRLTISYLRMHRLCAAGEWNOVGAGGEPLDACYL
KALEGFVMVLTAEGD MAYLSENVSKHLGLSQLEIGHSIFDFIHPCDQEELQDALTPPTERCFSLRMKST
KEKSRNAARSRRGKENLEFFELAKLLPLPGAISSQLDKASIVRLSVTYLRLRRFAALGAPPWGLRAAGPP
AGLAPGRRGPAALVSEVFEQHLGGHILQSLDGFVFALNQEKGFLYISETVSIYGLSQVEMTGSSVFDYI
HPGDHSEVLEQLGLVQERSFFVRMKSTLTRGLHVKASGKYKVIHVTGRLRALGLVALGHTLPPAPLAELP
WLQRAGGFWVLQSVATVAGSGKSPGEHHVLWVSHVLSQAEGGQT
GASKARRDQINAEIRNLKELLPLAEADKVRLSYLHIMSLACIYTRKGVFAGGTPLAGPTGLLSAQELED
IVAALPGFLLVFTAEGKLLYLSESVSEHLGHSMDLVAGQDSIYDIIDPADHLTVRQQLTLDRLFRCRF
EKSKNAARTREKENSEFELAKLLPLSAITSQLDKASIIRLTTSYLMRVVFPEGLGEAWGHSRTSP
ETERSFFLRMKCVLAKRNAGLTCSGYKVIHCSGYLKIRNVGLVAVGHSLLPPSAVTEIKLHSNMFMRASL
EKSKNAARTREKENSEFELAKLLPLSAITSQLDKASIIRLTTSYLMRAVFPPEGLGDAWGQPSRAGP
ETERSFFLRMKCVLAKRNAGLTCSGYKVIHCSGYLKIRIVGLVAVGQSLPPSAITEIKLYSNMFMFRASL
ELKHLILEAADGFLFIVSCETGRVVYVSDSVTPVLNQPQSEWFGSTLYDQVHPDDVDKLRQLSTSRMCM
GSRRSFICRMRCGSSEPHFVVVHCTGYIKAKFCLVAIGRLQVTSSPNCIDMSNVCQPTFISRHNIEGIF
DELKHLILRAADGFLFVVGCDRGKILFVSES VFKILNYSQNDLIGQSLFDYLHPKDIKVKELSSSRLC
SGARRSFFCRMKNRPRKSFCTIHSTGYLKSNSCLVAIGRLHSHVVPQPVNGEIRVKSMEYVSRHAIDG
```

5. Find all matches in sequence5.fasta where P is followed by any character except A, then L.

```
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab2$ grep "P[^A]L" sequence5.fasta
QLHWQIPPENSPLMERCFCRLRCLLDNSSGFLAMNFQGLKYLPLPQLALFAIATPLQPPSILEIRTKNF
MRMKCTVTNRGRTVNLKSATWKVLHCTGQVKVYEPLLSCLIMCEPIQHPSHMDIPLDSKTFLSRHSMDM
LTSRGRTLNLKAATWKVLNCSGHMRAYEPPLQCLVLICEAIPHPGSLEPPLGRGAFLSRHSLDMKFTYCD
FTQLMLEALDGFIIAVTDDGSIYVSDSITPLLGHLPSDVMQNLNLFPEQEHSEVYKLSSEYKSDS
ELKHLILEAADGFLFIVSCETGRVVYVSDSVTPVLNQPQSEWFGSTLYDQVHPDDVDKLRQLSTSRMCM
```

6. Print all lines in sequence5.fasta that have exactly 2 consecutive Vs anywhere in the line.

```
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab2$ grep "VV" sequence5.fasta
AANFREGNLQEGEFLLQALNGFVLVVTTDALVFYASSTIQDYLGFGQSDVIHQSVYELIHTEDRAEFQR
IWLQTHYYITYHQWNSRPEFIVCTHTVVSYAEVRAE
TVIYNTKNSQPQCIVCNVVVSGIIQHDL
QMDNLYLKALEGFIAVVTDGDMIFLSENISKFMGLTQVELTGHSIFDFTHPCDHEEIRENLSSTERDFF
KFTYCDDRITELIGYHPEELLGRSAYEFYHALDSENMTKSHQNLCTKGQVVSGQYRMLAKHGGYVWLETQ
DRIAEVAGYSPDDLIGCSAYEVIHALDSDAVSKSIHTLLSKGQAVTGQYRFLARSGGYLWTQTQATVVSG
QTHYYITYHQWNSKPEFIVCTHSVVSYADVRVE
DYVHPGDHVEMAEQLGMTLERSFFIRMKSTLTRGVHIKSSGYKVIHITGRLRLRMGLVVVAHALPPPTI
ISESVLYLGFERSSELLCKSWYGLLHPEDLAHASAQHYRLAESGDIAEMVVRLQAKTGGWAWIYCLLY
EKSKNAARTREKENSEFELAKLLPLSAITSQLDKASIIRLTTSYLMRVVFPPEGLGEAWGHSRTSP
LDNVGRELGSLLQTLDFGFIVVAPDGKIMYISETASVHLGLSQVELTGNSIYEIHPADHDEMTAVLTA
LDGVAKELGSHLLQTLDFGVVVASDGKIMYISETASVHLGLSQVELTGNSIYEIHPADHDEMTAVLTA
SYATVVHNSRSSRPHCIVSVNYVLTIEYKEL
ELKHLILEAADGFLFIVSCETGRVVYVSDSVTPVLNQPQSEWFGSTLYDQVHPDDVDKLRQLSTSRMCM
GSRRSFICRMRCGSSEPHFVVVGTGYIKAKFCLVAIGRLQVTSSPNCIDMSNVCQPTFISRHNIEGIF
TFVDHRCVATVGYQPQELLGKNIVEFCHPEDQQLLRDSFQVVKLKGQVLSVHFRFRSKNQEWLWMRTSS
DELKHLILRAADGFLVVGCDRGKILFVSES VFKILNYSQNDLIGQSLFDYLHPKDIKVKELSSSRLC
SGARRSFFCRMKNRPRKSFCTIHSTGYLKSNSCLVAIGRLHSHVVVPQPVNGEIRVKSMEYVSRHAIDG
RWFSEFMPNPTKEVEYIVSTNTVVVL
```

NAME: ALAMELU
ROLL NUMBER:25210013
LAB ASSIGNMENT – 2
Linux & Shell Scripting with Biological Data Files

7. Print all lines in sequence5.fasta that contain either AA or DD.

```
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab2$ grep -E "AA|DD" sequence5.fasta
AANFREGLNLQEGFLLQALNGFVLVTTDALVFYASSTIQDYLGFGQSDVIHQSVYELIHTEDRAEFQR
IFRTKHKLDFTPIGCDAGRIVLGYTEAELCTRGSGYQFIHAADMLYCAESHIRMIKTGESGMIVFRLLT
RHSLEWKFLFLDHRAPPIIGYLPFEVLGTSGYDYYHVDDLENLAKCHEHLMQYKGKSCYYRFLTKGQQW
KEKSRDAARSRRSKESEVFYELAHQLPLPHNVSSHLDKASVMRLTISYLRVRKLLDAGDLIEDDMKAQM
NCFYLKALDGFVMVLTDGDMIIYISDNVNKYMGLTQFELTGHVSFDFTHPCDHEEMREMLTHNTQRSFFL
KEKSRDAARCRRSKETEVYELAHQLPLPHSVSSHLDKASIMRLAISFLRTHKLLSSVCSENESEAEADQ
KFTYCDDRITELIGYHPEELLGRSAYEFYHALDSENMTKSHQNLCTKGQVVSQYRMLAKHGGYVWLETQ
DAARSRRSQETEVLYQLAHTLPFARGVSAHLDKASIMRLTISYLRMHRLCAAGEWNQVGAGGEPLDACYL
LTSRGTLLNLKAATWVKVNLCSGHRMAYEPPLQCLVLICEAIPHPGSLEPPLGRGAFLSRHSLDMKFTYCD
DRIAEGVAGYSPDDLIGCSAYEYIHALDSDAVSKSIHTLLSKGQAVTGQYRFLARSGGYLWTQTQATVVSG
KEKSRNAARSRRGKENLEFFELAKLLPLPGAISSQLDKASIVRLSVTYLRLRRFAALGAPPWGLRAAGPP
AGLAPGRRGPAAALVSEVFEQHLGGHILQSLDGFVFALNQEGKFLYISETVSIYGLSQVEMTGSSVFDYI
LEWKFLFLDHRAPPIIGYLPFEVLGTSGYDYYHIDDELLARCHQHLMQFGKGKSCCYRFLTKGQQWIWL
SRDAARSRRGKENFEFELAKLLPLAAITSQLDKASIIRLTISYLMRDFANQGDPPWNLMEGPPPPNT
IVAALPGFLLVFTAEGKLLYLSVSEHLGHSMDLVAQGDSIYDIIDPADHLTVRQQLTLDRLFRCRF
EKSNAARTREKENSEFELAKLLPLPSAITSQLDKASIIRLTISYLMRAVFPFEGLDGAWGHSRTSP
EKSNAARTREKENSEFELAKLLPLPSAITSQLDKASIIRLTISYLMRAVFPFEGLDGAWGQPSRAGP
ELKHLILEAADGFLFIVSCETGRVVYVSDSVTPVLNQPSSEWFGSTLYDQVHPDDVDKLRQLSTSRMCM
DELKHLILRAADGFLFVVGCDRGKILFVSESVFKILNYSQNDLIGQSLFDYLHPKDIKVKQLSSSRLC
KFVFDORATAILAYLPQELLGTSCYEFHDDIGHLAECHROVLQTREKITTNCYKFKIKDGSFITLRS
```

8. Print only the sequence lines (ignore headers) from sequence5.fasta that contain the letter P.

```
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab2$ grep -v ">" sequence5.fasta | grep "P"
SNPSKRHRDRNLTELDRLASLLPFPQDVINKDLKSVLRLSVSYLRAKSFDDVALKSSPTERNGGQDNCR
QLHWQIPPENSPLMERCFCIRLCRLDNSSGFLAMNFQKGLKYLPPQLALFAIATPLQPPSILEIRTKNF
IFRTKHKLDFTPIGCDAGRIVLGYTEAELCTRGSGYQFIHAADMLYCAESHIRMIKTGESGMIVFRLLT
KNNRWTVVQSNARLLYKNGRPDYIIVTQRPLTDEEGTEHLR
VSRNKSEKKRRDQNFVLIKELGSMPLGNARKMDKSTVLQKSIDFLRKHKEITAQSDASEIRQDWKPTFLS
NEEFTQLMLEALDGGFLAINTDGSIIYVSESVTSLLEHLPSDLVDQSIFNFIPEGEHSEVYKILSTEYK
SKNQLEFCCHMLRGTDIDKEPSTYEVVKFIGNFKSLYEDRVCFVATVRLATPQFIKEMCTVEEPNEEFTS
RHSLEWKFLFLDHRAPPPIIGYLPFEVLGTSGYDYYHVDLENLAKCHEHLMQYKGKSCYYRFLTKGQQW
IWLQTHYYITYHQWNSRPEFIVCTHTVVSVAEVRAE
KEKSRDAARSRRSKESEVFYELAHQLPLPHNVSSHLDKASVMRLTISYLRVRKLLDAGDLIEDDMKAQM
NCFYLKALDGFVMVLTDGDMIIYISDNVNKYMGLTQFELTGHVSFDFTHPCDHEEMREMLTHNTQRSFFL
RMKCTLTSRGTMMIKSATWVKVLHCTGHIHVYKPPMTCLVLICEPIPHPSNIEIPLDSKTFLSRHSLDMK
FSYCDERITELMGYPEELLGRSIEYHYHALSDHLTKTHHDMFTKGQVTTGQYRMLAKHGGYVWVETQA
TVIYNTKNSQPQCIVCVNYVVGIIQHDL
KEKSRDAARCRRSKETEVYELAHQLPLPHSVSSHLDKASIMRLAISFLRTHKLLSSVCSENESEAEADQ
QMDNLYLKALEGFIAVVTQDGMIFLSENISKFMGLTQVELTGHISIFDFTHPCDHEEIRENLSSTERDFF
MRMKCTVTNRGTVNLKSATWVKVLHCTGQVKVYEPILLSCLIIMCEPIQHPSHMDIPLDSKTFLSRHSMDM
KFTYCDRITELIGYHPEELLGRSAYEFYHALDSENMTKSHQNLCTKGQVVSQYRMLAKHGGYVWLETQ
GTVIYNPRNLQPQCIMCVNYVLSEIEKNDV
DAARSRRSQETEVLYQLAHTLPFARGVSAHLDKASIMRLTISYLRMHRLCAAGEWNQVGAGGEPLDACYL
KALEGFVMVLTAEGDMAYLSENVSKHLGLSQLEIGHISIFDFIHPCDQEELQDALTPPTERCFSLRMKST
LTSRGTLLNLKAATWVKVNLCSGHRMAYEPPLQCLVLICEAIPHPGSLEPPLGRGAFLSRHSLDMKFTYCD
DRIAEGVAGYSPDDLIGCSAYEYIHALDSDAVSKSIHTLLSKGQAVTGQYRFLARSGGYLWTQTQATVVSG
GRGPQSESIVCVHFLISQVEETGV
KEKSRNAARSRRGKENLEFFELAKLLPLPGAISSQLDKASIVRLSVTYLRLRRFAALGAPPWGLRAAGPP
AGLAPGRRGPAAALVSEVFEQHLGGHILQSLDGFVFALNQEGKFLYISETVSIYGLSQVEMTGSSVFDYI
```


NAME: ALAMELU
ROLL NUMBER:25210013
LAB ASSIGNMENT – 2
Linux & Shell Scripting with Biological Data Files

```
GRGPQSESIVCVHFLISQVEETGV
KEKSRNAARRRKENLEFFELAKLLPLPGAISSQLDKASIVRLSVTYLRLRRFAALGAPPWGLRAAGPP
AGLAPGRRGPAAALVSEVFEQHLGGHILQSLDGFVFALNQEKGFLYISETVSIYLGLSQVEMTGSSVFDYI
HPGDHSEVLEQLGLVQERSFFVRMKSTLTKRGLHVKASGYKVIHVTGRLRALGLVALGHTLPPAPLAELP
LHGMMIVFRLSLGLTILACESRVSDHMDLGPSELVGRSCYQFVHGQDATRIRQSHVDLLDKGQVMTGYR
WLQRAGGFVWLQSVATVAGSGKSPGEHHVLWVSHVLSQAEGGQT
NKSEKKRRDQFNVLIKELSSMLPGNTRKMDKTTVLEKVIQFLQKHNEVSAQTEICDIQQDWKPSFLSNEE
FTQLMLEALDGFIIAVTTDGSIIYVSDSITPLLGHLPSDVMQDNLNLFPEQEHSEVYKILSSEYLSKSDS
DLEFYCHLLRGSNPKFPTYEYIKFVGNFRSYLGKEVCFIATVRLATPQFLKEMCIVDEPLEEFTSRHS
LEWKFLFLDHRAPPPIIGYLPFEVLGTSGYDYHIDDLELLARCHQHLMQFGKGKSCCYRFLTkgQQWIWL
QTHYYITYHQWNSKPEFIVCTHSVVSADVRVE
SRDAARRRKENFEFELAKLLPLPAAITSQLDKASIIRLTISYLMRDFANQGDPPWNLRMGPPPPNT
SVKVIGAQRRRSPSALAIEVFEAHLGSHILQSLDGFVFALNQEKGFLYISETVSIYLGLSQVELTGSSVF
DYVHPGDHVEAEQLGMTLERSFFIRMKSTLTKRGVHIKSSGYKVIHITGRLRLRMGLVVVAHALPPPTI
NEVRIDCHMFVTRVNMDLNIYCNRISDYMDLTPVDIVGKRCYHFIHAEDVEGIRHSHLDLLNKGQCVT
KYRWMQKNGGYIWIQSSATIAINAKNANEKNIWVNYLLSNPEYKDT
GASKARRDQINAEIRNLKELLPLAEADKVRLSYLHMSLACIYTRKGVFFAGGTPLAGPTGLLSAQELED
IVAALPGFLLVFTAEGKLLYLSSEVSEHLGHSMDLVAQGDSIYDIIDPADHLTVRQQLTLDRLFRCRF
NTSKSLRRQSAGNKLVLIRGRFAHNPVFTAFCAPLEPRPRPGPGPGPASLFLAMFQSRHAKDLALLD
ISESVLIYLGFERSELLCKSWYGLLHPEDLAHASAQHYRLLAESGDIQAEVVRLQAKTGGWAWIYCLLY
SEGPEGITANNYPISDMEAWSLRQQL
EKSNAARTREKENSEFELAKLLPLPSAITSQLDKASIIRLTTSYLMRNVFPEGLGEAWGHSSRTSP
LDNVGRELGSLLQTLDDGFIFVVPADGKIMYISETASVHLGLSQVELTGNSIYEYIHPADHDEMTAVLTA
EIERSFFLRMKCVLAKRNAGLTCGGYKVIHCSGYLKIRNVGLVAVGHSLPPSAVTEIKLHNSMFMFRASL
DMKLIFLDSRVAELTGYPEQDLIEKTLYHHVHGCDTFHLRCAHLLLVKGQVTTKYRFLAKHGGVWVWQ
SYATIVHNSRSSRPHCIVSVNYVLTDEYKGL
EKSNAAKTRREKENGFEYELAKLLPLPSAITSQLDKASIIRLTTSYLMRAVFPPEGLGDAWGQPSRAGP
LDGVAKELGSHLLQTLDDGFVAVSDGKIMYISETASVHLGLSQVELTGNSIYEYIHPSDHDEMTAVLTA
EIERSFFLRMKCVLAKRNAGLTCGYKVIHCSGYLKIRIVGLVAVGQSLPPSAITEIKLYNSMFMFRASL
DLKLIFLDSRVTEVTGYEQDLIEKTLYHHVHGCDVFHLRYAHLLLVKGQVTTKYRLLSKRGGVWVWQ
SYATVVHNSRSSRPHCIVSVNYVLTIEYKEL
NHSEIERRRRNKMTAYITELSDMVPTCSALARKPDKLTILRMVSHMKSRLGTGNTSTDGSYKPSFLTQDQ
ELKHLILEAADGFLFIVSCETGRVVYVSDSVTPVLNQPQSEWFGSTLYDQVHPDDVDKLRQQLSTSRMCM
GSRRSFICMRGCSSEPHFVVVHCTGYIAKAFCLVAIGRLQVTSSPNCTDMSNVCQPTTEFISRHNIEGIF
TFVDHRCVATVGYYQPQELLGKNIVEFCHPEDQQLRDSFQQVVKLGQVLSVMFRFRSKNQEWLWMRTSS
FTFQNPYSDEIEYIICTNTNVK
EAHSQIEKRRRDKMNSFIDELASLVPTCNAMSRKLDKLTVLRMAVQHMKTLRGATNPYTEANYKPTFLSD
DELKHLILRAADGFLFVVGCDRGKILFVSESFVKILNYSQNDLIGQSLFDYLPKDKIAKVKEQLSSSRCLC
SGARRSFFCRMKNRPRKSFCITIHSTGYLKSNSCLVAIGRLHSHVVPQPVNGEIRVKSMEYVSRHAIDG
KFVVDQRATAILAYLPQELLGTSCYEFHQDDIGHLAECHRQVLQTREKITTNCYKFKIKDGSFITLRS
RWFSSMNPWKEVEYIVSTNTVVL
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab2$
```

NAME: ALAMELU
ROLL NUMBER:25210013
LAB ASSIGNMENT – 2
Linux & Shell Scripting with Biological Data Files

PART 3

9. Store the filename sequence5.fasta in a variable called seq and print the number of sequences in it (headers count as sequences).

```
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab2$ seq="sequence5.fasta"
echo "Number of sequences in $seq:"
grep -c ">" $seq
Number of sequences in sequence5.fasta:
13
```

10. Store the pattern G{2,\} in a variable and search protein.fasta for sequence lines (ignore headers) with 2 or more consecutive Gs.

```
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab2$ var="G{2,\}"
grep -v ">" protein.fasta|grep "$var"
KPVKKKKIKREIKILENLRGGPNITLADIVKDPVSRTPALVFEHVNNTDFKQLYQTLTDYDIRFYMYEI
WERFVHSENQHLVSPEALDFLDKLLRYDHQSRLTAREAMEHPYFYTVVKDQARMGSSSMPGGSTPVSSAN
```

11. Store "Biocomputing" in a variable, export it, and verify that it is available inside a new shell started using:

```
bash -c 'echo $VARIABLE_NAME'
```

```
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab1$ var="Biocomputing"
export var
bash -c 'echo $var'
Biocomputing
```

NAME: ALAMELU
ROLL NUMBER:25210013
LAB ASSIGNMENT – 2
Linux & Shell Scripting with Biological Data Files

PART-4

12. Write a shell script that checks if sequence3.fasta exists in the current folder. If yes, print the number of lines. If no, print "Missing file".

```
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab2$ if [ -f "sequence3.fasta" ]; then
wc -l sequence3.fasta
else
echo "Missing file"
fi
19 sequence3.fasta
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab2$ cd /mnt/d/MTECH/SEM1/Biocomputing/Lab1
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab1$ if [ -f "sequence3.fasta" ]; then wc -l sequence3.fasta; else echo "Missing file"; fi
Missing file
```

13. Using a for loop, go through all .fasta files in the current directory and print: filename, number of sequences, and file size in characters.

```
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab1$ for file in *.fasta; do
echo "File name: $file"
echo "Sequences: $(grep -c ">" "$file")"
echo "File size (characters): $(wc -c < "$file")"
done
File name: protein.fasta
Sequences: 1
File size (characters): 467
File name: sequence.fasta
Sequences: 1
File size (characters): 79551
```

NAME: ALAMELU
ROLL NUMBER:25210013
LAB ASSIGNMENT – 2
Linux & Shell Scripting with Biological Data Files

```
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab1$ cd /mnt/d/MTECH/SEM1/Biocomputing/Lab2
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab2$ for file in *.fasta; do
echo "FILE NAME: $file"
echo "NO OF SEQUENCES: $(grep -c ">" "$file")"
echo "FILE SIZE(IN CHARACTERS): $(wc -c < "$file")"
done
FILE NAME: protein.fasta
NO OF SEQUENCES: 1
FILE SIZE(IN CHARACTERS): 467
FILE NAME: sequence.fasta
NO OF SEQUENCES: 1
FILE SIZE(IN CHARACTERS): 79551
FILE NAME: sequence1.fasta
NO OF SEQUENCES: 1
FILE SIZE(IN CHARACTERS): 974
FILE NAME: sequence2.fasta
NO OF SEQUENCES: 4
FILE SIZE(IN CHARACTERS): 1710
FILE NAME: sequence3.fasta
NO OF SEQUENCES: 2
FILE SIZE(IN CHARACTERS): 1000
FILE NAME: sequence4.fasta
NO OF SEQUENCES: 4
FILE SIZE(IN CHARACTERS): 2374
FILE NAME: sequence5.fasta
NO OF SEQUENCES: 13
FILE SIZE(IN CHARACTERS): 4229
```

14. Modify the above loop so that it only prints files with more than 3 sequences.

```
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab2$ for file in *.fasta
do
count=$(grep -c "^>" "$file")
if [ "$count" -gt 3 ]; then
echo "FILE NAME: $file"
echo "NO OF SEQUENCES: $(grep -c ">" "$file")"
echo "FILE SIZE(IN CHARACTERS): $(wc -c < "$file")"
fi
done
FILE NAME: sequence2.fasta
NO OF SEQUENCES: 4
FILE SIZE(IN CHARACTERS): 1710
FILE NAME: sequence4.fasta
NO OF SEQUENCES: 4
FILE SIZE(IN CHARACTERS): 2374
FILE NAME: sequence5.fasta
NO OF SEQUENCES: 13
FILE SIZE(IN CHARACTERS): 4229
```


NAME: ALAMELU
ROLL NUMBER:25210013
LAB ASSIGNMENT – 2
Linux & Shell Scripting with Biological Data Files

PART-5

15. From sequence5.fasta, extract only the sequence lines (no headers) that contain 3 or more cysteines (C). Save the output to a file named cys_rich.txt. Ensure the output file contains no empty lines.

```
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab2$ grep -v ">" sequence5.fasta | grep "C
.*C.*C" > cys_rich.txt
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab2$ less cys_rich.txt
user@DESKTOP-A5PNEA5:/mnt/d/MTECH/SEM1/Biocomputing/Lab2$ cat cys_rich.txt
SNPSKRHRDRNLNTELDRLASLLPFPQDVINKLDKLSVLRSLVSVYRAKSFDFVALKSSPTERNGGQDNCR
QLHWQIPPENSPLMERCFICRLRCLLDNSSGFLAMNFQGKLYLPPQLALFAIATPLQPPSILEIRTKNF
IFRTKHKLDFTPIGCDAGRIVLGYTEAELCTRGSGYQFIHAADMLYCAESHIRMIKTGESGMIVFRLLT
SKNQLEFCCHMLRGITDPKEPSTYEVVKFIGNFKSLYEDRVCFVATVRLATPQFIKEMCTVEEPNEEFTS
RHSLEWKFLFLDHRAPPIIGYLPFEVLGTSGYDYHVDLENLAKCHEHLMQYVGKGSVCYRFLTKGQQW
IWLQTHYYITYHQWNSRPEFIVCTHTVVSVAEVRAE
NCFYLKALDGFVMVLTDGDMYISDNVNKYMGLTQFELTGHSVDFDTHPCDHEEMREMLTHTQRSFFL
RMKCTLTSRGRMTMKSATWKVLHCTGHIHVYKPPMTCLVLICEPIPHPSNIEIPLDSKTFLSRHSMDK
FSYCDERITELMGYEPEELLGRSIEYEHALDSHDLTKTHDMFTKGQVTTGQYRMLAKRGGVVWVETQA
TVIYNTKNSQPQCIVCVNYVVSIGIQHDL
KEKSRDAARCRRSKETEVFYELAHPLPHSVSSHLDKASIMRLAISFLRTHKLLSSVCSENESEAEADQ
QMDNLYLKALEGFIADVTDGDMIFLSENISKFMGLTQVELTGHSIFDFTHPCDHEEIRENLSSSTERDFF
MRMKCTVTNRRGRTVNLKSATWKVLHCTGQVKVYEPLLSCLIMCEPIQHPSHMDIPLDSKTFLSRHSMDM
KFTYCDRITELIGYHPEELLGRSAYEFYHALDSENMTKSHQNLCTKGQVVSQYRMLAKHGGYVWLETQ
GTVIYNPRNLQPQCIMCVNYVLSIEKNDV
DAARSRSQETEVLYQLAHTLPFARGVSAHLDKASIMRLTISYLRMHRLCAAGEWNVQVAGGEPDACYL
KALEGFMVVLTAEGDMAYLSENVSKHLGLSQLELIGHSIFDIHPCDQEELQDALTPPTERCFSLRMKST
LTSRGRTLNLKAATWKVLNCSGHRMAYEPPLQCLVLICEAIPHPGSLPEPLGRGAFLSRHSMDMKFTYCD
DRIAEVAGVSPDDLIGCSAYEYIHALDSDAVSKSIHTLLSKGQAVTGQYRFLARSGGYLWTQTQATVVSG
GRGPQSEISIVCVHFLISQVEETGV
LHGHIIVFRLSLGLTILACESRVSDHMDLGPSELVGRSCYQFVHGQDATRIRQSHVDLLDKGQVMTGYR
NKSEKRRDQFNVLIKELSSMLPGNTRKMDKTTVLEKVIQGLQKHNEVSAQTEICDIQDQWPKPSFSLNEE
DLEFYCHLLRGLSNPKFPTVEYIKFVGNFRSYLGEKVCFIATVRLATPQFLKEMCIVDEPLEEFTSRHS
LEWKFLFLDHRAPPIIGYLPFEVLGTSGYDYHIDDLELLARCHQHLMQFGKGKSCCYRFLTKGQQWIWL
QTHYYITYHQWNSKPEFIVCTHSVSVYADVRVE
NEVRIDCHMFVTRVMDLNIYICENRISDYMDLTPVDIVGKRCYHFIHAEDVEGIRHSHDLLNKGQCVT
GASKARRDQINAEIRNLKELLPLAEADKVRLSYLHIMSLACIYTRKGVFVAGGTPLAGPTGLLSAQELED
IVAALPGFLLVFTAEGKLLYLSVSSEHLGHSMDVLAQGDSDIYDIIDPADHLTVRQQLTLDRLFRCRF
NTSKSLRRQSAGNKLVLIRGRFHAHNPVFTAFCAPLEPRPRPGPGPGPASFLAMFQSRHAKDLALLD
ISESVLIYLGFERSELLCKSWYGLLHPEDLAHASAQHYRLLAESGDIQAEVVRQLAKTGGWAWIYCLLY
EIERSSFFLRMKCVLAKRNAGLTCGGYKVIHCSGYLKIRNVGLVAVGHSLPPSAVTEIKLHSNMFMFRASL
DMKLIFLDSRVAELTGYPQDLIEKTLYHHVHGCDTFHLRCAHLLLVKGQVTTKYRFLAKHGGVWVWQ
SYATIVHNSRSSRPHCIVSVNYVLTDEYKGL
EIERSSFFLRMKCVLAKRNAGLTCGGYKVIHCSGYLKIRIVGLVAVGQSLPPSAITEIKLYSNMFMFRASL
DLKLIFLDSRVTEVTGYEPQDLIEKTLYHHVHGCDVFHLRYAHLLLVKGQVTTKYRLLSKRGGVWVWQ
SYATVVHNSRSSRPHCIVSVNYVLTIEYKEL
NHSEIERRRRNKMTAYITLSDMVPTCSALARKPKDLTILRMAVSHMKSLRGNTSTDGSYKPSFLTDQ
ELKHLILEAADGFLFIVSCETGRVYVSDSVTPVLNQPQSEWFGSTLYDQVHPDDVDKREQLSTRMCM
GSRRSFCRMRCGSEPHFVVHCTGYIAKAFCLVAIGRLQVTSNPCTDMSNVCQPTFISRHNIEGIF
TFVDHRCVATVGYQPQELLGKNIVEFCHPEDQQLLRDSFQQVVKLKGQVLSVMFRFSKNQEWLWMTSS
FTFQNPYSDEIEYIICNTNVK
EAHSQIEKRRRDKMNSFIDELASLVPTCNAMSRKLDKLTVLRLMAVQHMKTLRGATNPYTEANYKPTFLSD
DELKHLILRAADGFLFVVGCDRGKILFVSESVFKILNYSQNDLIGQSLFDYLHPKDIKVKELSSSRLC
SGARRSFFCRMKNRPRKSFTIHSYGKLSNLSCLVAIGRLSHSVVPQVNGEIRVKSMEYVSRHAIDG
KFVFDQRAATAILAYLPQELLGTSCYEYFHQDDIGHLAECHRQVLQTRKITTNCYKFKIKDGSFITLRS
```