| Lext Processing |
|---|
| Starts at 9:05 p.m. |
| |
| ## Agenda |
| 1. Arrays |
| 2. grep |
| 3. sed |
| 4. awk |
| 5. regex |
| 6. Program |
| grays |
| (i) Delaring an array. |
| |
| my_array = ("apple" "banano" "orange") |
| my_array: ("apple" "banano" "orange") my_array: (12345) |
| 2) Accessing an dray element |
| |

euro \$ 4 my-array [0] $\rightarrow apple$ euro \$ 4 my-array [1] → banano

echo \$ & my - array [@]}

Ly access all elements of an array.

Usina **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

— **

—

- appl banana orange

→ \$# size of away.

" & s numbers (@7 § "

anays 1 "1" "2" "3" 4" "5"

\$ (a) → treats each element as argument
"1" "2" "3"

- trads entire array as argument.

" 1 2 3 4 5 "

| → geup. |
|--|
| |
| -> grep-flag "patten" filenam. |
| |
| i -> case insesnsitive Appl -i |
| n -> prints line number along with the text |
| c -> counts the matches |
| w -> match whole word |
| B -> before |
| A -> after |
| C -> around |
| v-> Invert match |
| o -> only the matched part, instead of the whole line. |
| |
| - E → extended sugular expression. |
| |
| grep -E 'lemon tart' text |
| |
| SED Stram editor. |
| |
| Soarch |
| replan |
| |

Pount O substitution (S) delete 2 delution (d) insert (3) Prunting (P) append. () substitution (s) sed 's/original/replace/' file replaces first occurrence on each line. Sed 'S/original Suplare/g' file -i → in plane update.

(modify the contents of file) sed '2 s/appl/biryani/ file sed 's/apple/birgani/I, file

APPLE - biryani case Insensitive. 2 Seleting. sed'/pattun/d' file Sed '2d' file -, for deleting a line 3 Inserting & Appending Sed '3i \ Inserted line 'file sed '3a Appended line 'file Sed for printing. (p) Sed -n '2p' file

Suppressing all other lines

| AWK. | | | | | | |
|----------------|-----|-------|----------|----|--|--|
| | | , fie | elcls | | | |
| | 10, | C 2 | ر ا ر | Cy | | |
| R | | | | • | | |
| Rorards. R2 | | 0 | | | | |
| R ₃ | | G. | | | | |
| Ry | | | | | | |
| | ' | | · | | | |

\$ - fields

default separator → spare

-F',

awk 'NR == 2; ** perint \$ 49'

- F',

- pattun matching. print \$0

awk '/pattun/ Kprint 4' text

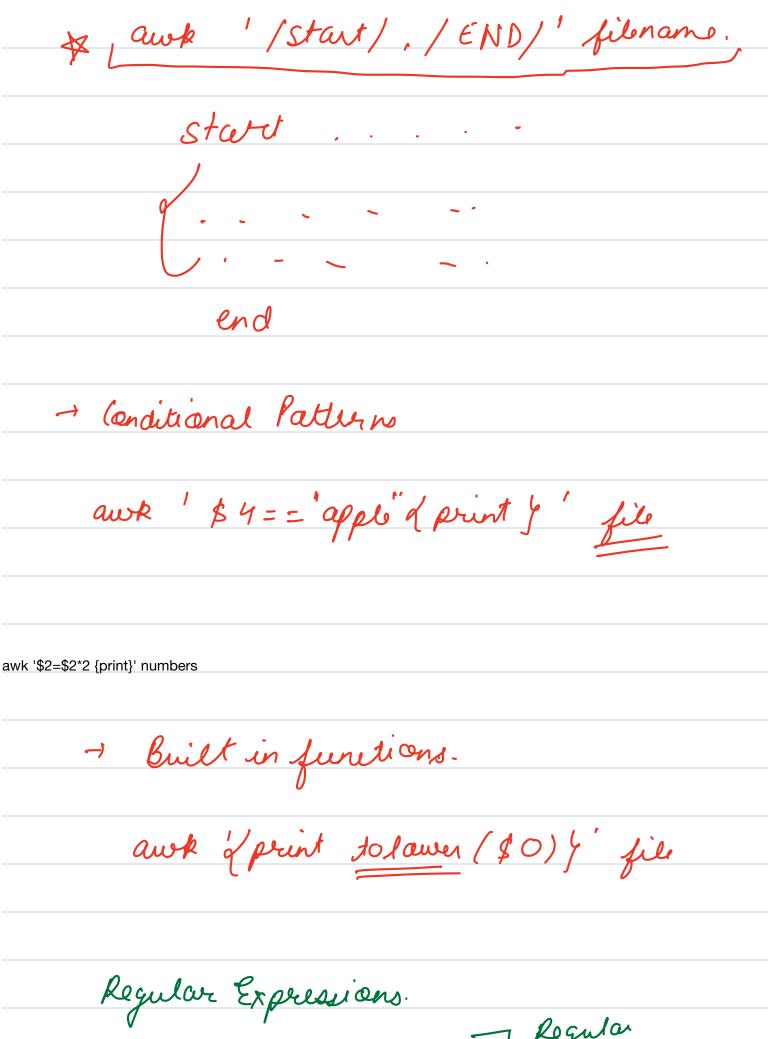
- changing field separator.

awk -F ',' '{print \$2}' awk_delim

- Lange Patterns.

fetch logs from pattern!

to pattern?



Rogular

· dot - any character in place of dot. grep "a.p" text matches -> aip 1 → position at start of line

\$\frac{1}{2}\$ → position at end of line grep "1 apple" text
grep "LE\$" text

3 Character class []

grep [AcJpple text [a-z] -> denotes all words SA-ZJ → " capital SO-97 → all numbers [0-9] - all numbers [] I risside brackets is used for negation. (4) Quartifiers * > zero or more occurrences + + 1 or more occurrences. grop ap*" file reeded px - 0 will also be fine. grep apt file needed.

| ang - match exact number of times |
|---|
| of n, & -1 atleast n times |
| dn, m g → Between n 2 m time |
| |
| string literals. Escape character. |
| Con on the same that |
| Escapi Maran. |
| |
| grep "\•" file |
| |
| W. Fetch emails Iroma file using |
| Y Fetch emails from a file using regular expressions. |
| |
| |
| |
| a-z -) all smallrage |
| A-7 → all upperease. |
| A-7 → all upperease. 0-9 → all numbers. |
| • |
| √ . |
| + |

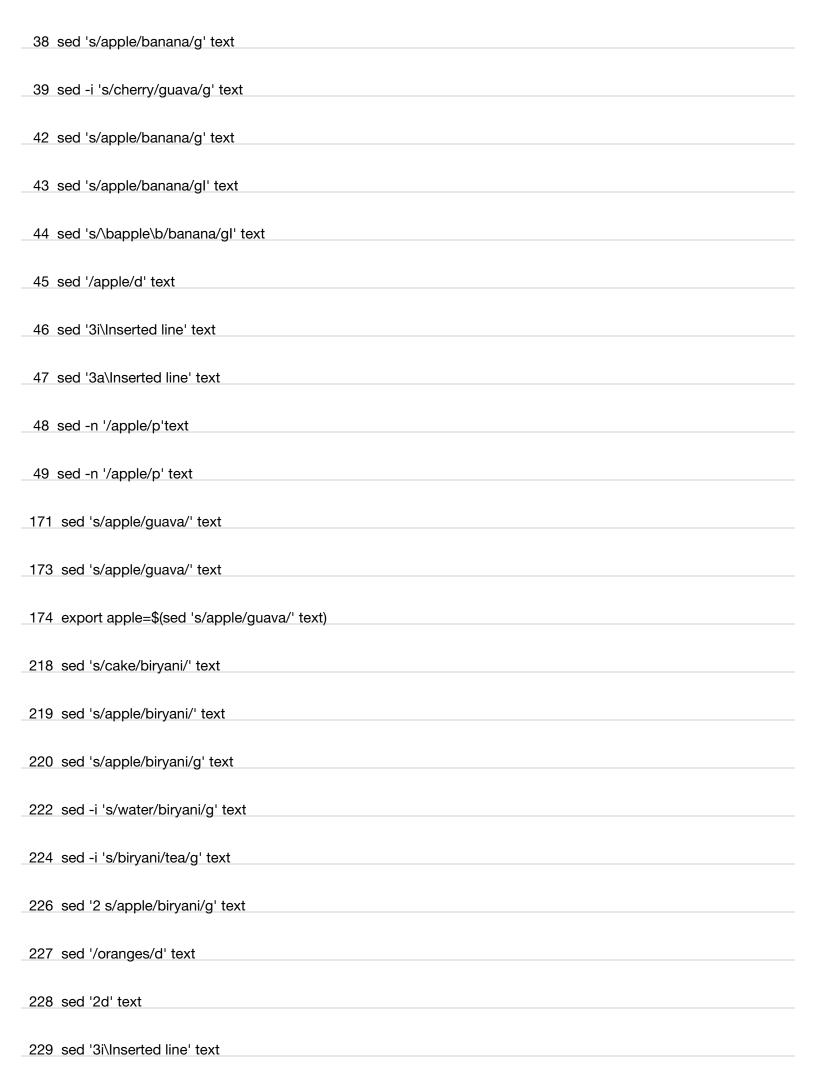
^[a-zA-z0-9._%+-]+@ [a-zA-z0-9._]+.[a-z1-2]~2,4\$

- **`[a-zA-Z0-9._%+-]+`**: - **`[]`**: This denotes a character class. - **`a-zA-Z`**: Matches any uppercase or lowercase letter. - **`0-9`**: Matches any digit. - **`.`**: Matches a literal dot. - **`_`**: Matches an underscore. - **`%`**, **`+`**, **`-`**: Matches these special characters. - **`+`**: This quantifier means "one or more" of the preceding characters. So, this part matches the local part of the

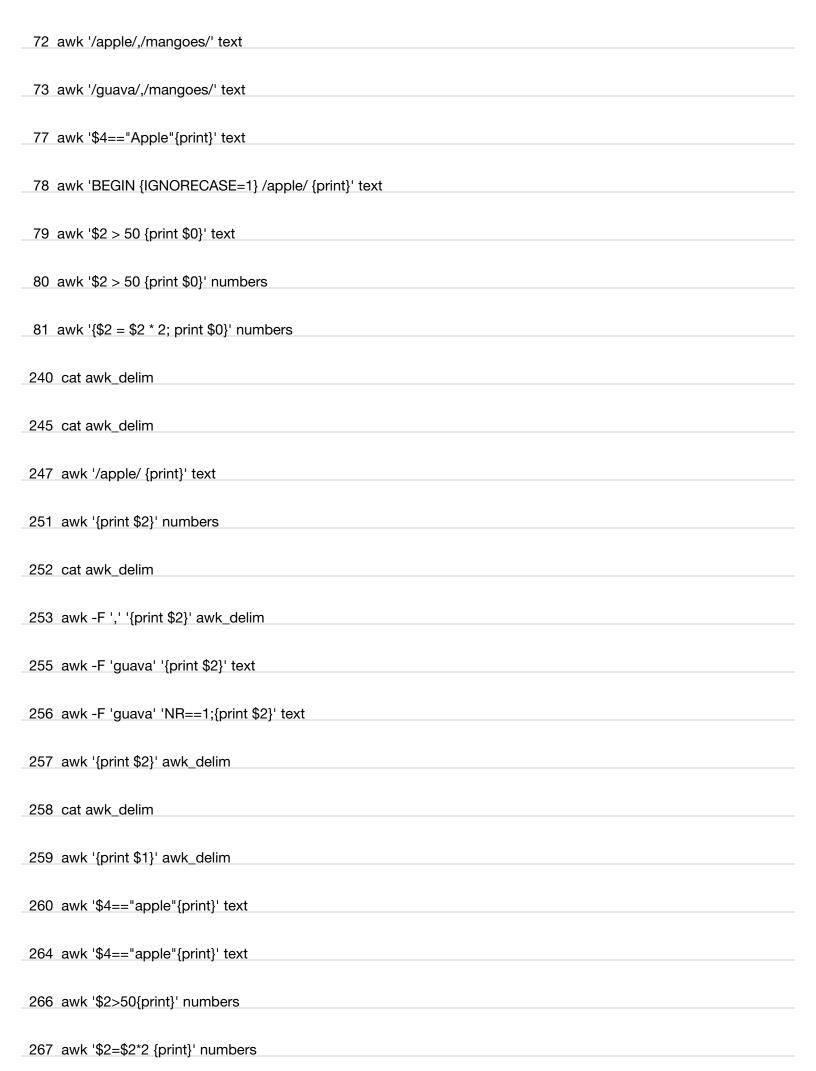
email (like 'user.name', 'user-name', 'user123', etc.).

1. **'@'**: This matches the '@' symbol, which separates the local part from the domain.

| 2. **`[a-zA-Z0-9]+``*: |
|---|
| |
| Similar to the legal part, this part metabos the domain name, allowing letters, digits, data, and hyphone |
| - Similar to the local part, this part matches the domain name, allowing letters, digits, dots, and hyphens. |
| - **`+`**: Again, means "one or more" characters. |
| 3. **`\.`**: Matches a literal dot (`.`). Since `.` is a metacharacter in regex (matching any character), we escape it with a |
| backslash to treat it as a normal character. |
| |
| 4 **`[a ¬A ¬][Q]`**· |
| 4. **`[a-zA-Z]{2,}`**: |
| |
| - This matches the top-level domain (like `com`, `net`, `org`, etc.). |
| - **`{2,}`**: Means "at least 2 characters" (to allow for domains like `.com` and `.org`). |
| |
| 5. **`\$`**: Asserts the end of the string, ensuring that nothing follows the email address. |
| |
| -> History from system. |
| SED |
| |
| sed '2d' text |
| 32 sed '2d' text |
| 33 sed '2p' text |
| 34 sed -n '2p' text |
| _ 0 1 364 11 2ρ t6λt |
| _35_sed 's/apple/banana/' text |
| 37 sed 's/apple/banana/' text |



| 230 sed '3a\Inserted line' text |
|--|
| 231 sed '2p' text |
| |
| 232 sed -n '2p' text |
| 233 sed -n '2,4p' text |
| _235_sed -n '2,4p' text |
| 236 sed -n '/apple/p' text |
| 249 sed -n "/apple/p" text |
| 318 history grep sed |
| |
| AWK. |
| 51 awk '\$2 > 50 {print}' numbers |
| 52 awk '\$2 > 50 {print \$0}' numbers |
| 53 awk '{print \$0}' numbers |
| 54 awk '{print}' numbers |
| 58 awk '{print \$3}' numbers |
| 59 awk '{NR==2;print \$3}' numbers |
| 60 awk 'NR==2{print \$3}' numbers |
| 62 vi awk_delim |
| 63 awk -F ',' 'NR==2{print \$3}' awk_delim |
| |
| 64 awk -F ',' '{print \$3}' awk_delim |
| _65_cat_awk_delim |
| 71 awk '/apple/,/oranges/' text |



| 314 grep -E "ERROR" data awk '{print \$1 \$2 \$3 }' |
|--|
| |
| 315 grep -E "ERROR" data awk '{print \$1" " \$2" " \$3 }' |
| 316 grep -E "ERROR" data awk '{print \$1" " \$2" " \$3 " "for (i=5; i<=NF; i++) printf \$i }' |
| 317 grep -E "ERROR" data awk '{print \$1" " \$2" " \$3 " " for (i=5; i<=NF; i++) printf \$i }' |
| 317 grep -L Linnon data awk (print φτ φε φο τοι (i=0, i<=ivi , i++) printi φι j |
| 319 history grep awk |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |
| |