

Uniksowe filtry i wyrażenia regularne

Witold Paluszyński

Katedra Cybernetyki i Robotyki

Politechnika Wrocławska
<http://www.kcir.pwr.edu.pl/~witold/>

1995–2015

Ten utwór jest dostępny na licencji

 Creative Commons Uznanie autorstwa–
Na tych samych warunkach 3.0 Unported

Utwór udostępniany na licencji Creative Commons: uznanie autorstwa, na tych samych warunkach. Udziela się zezwolenia do kopiowania, rozpowszechniania i/lub modyfikacji treści utworu zgodnie z zasadami w/w licencji opublikowanej przez Creative Commons. Licencja wymaga podania oryginalnego autora utworu, a dystrybucja materiałów pochodnych może odbywać się tylko na tych samych warunkach (nie można zastrzec, w jakikolwiek sposób ograniczyć, ani rozszerzyć praw do nich).

W systemie Unix zaimplementowano szereg ciekawych programów przetwarzających na różne sposoby ciąg danych zorganizowanych w wiersze, to znaczy rekordy zakończone znakiem NEWLINE (ASCII 10).

Ponieważ programy te wykonują często bardzo proste operacje, często używa się więcej niż jednego na raz, korzystając z mechanizmu potoków:

```
cat plik | prog1 | prog2 | prog3
```

Takie przetwarzanie ma charakter filtrowania strumienia danych, stąd programy wykorzystywane w ten sposób nazywa się **filtrami tekstowymi**.

Typowa sytuacja wykorzystania potoków filtrów jest praca interakcyjna, gdzie użytkownik przeszukuje jakieś pliki lub dane dobierając właściwe filtry i ich parametry. Często na początku potoku pojawia się program `cat` czytający dane z pliku. Również typowo na końcu potoku bywa wywoływany program `more` lub `less` pozwalający przeglądać dane ekran po ekranie.

Warto przypomnieć, że **procesy w potoku pracują równolegle**, co daje istotne przyspieszenie w systemach wieloprocesorowych/wielordzeniowych.

Uniksowe filtry tekstowe — cat

Najprostszym z filtrów jest program `cat` czytający dane z wejścia i wypisujący je bez zmian na wyjście. `cat` posiada niewiele opcji, a ponadto jest niezwyczkle przydatnym programem, często wykorzystywanym w potokach filtrów.

Na przykład, wywołanie `cat` bez żadnych argumentów na początku potoku pozwala przetwarzać dane pisane z klawiatury. Dzięki temu użytkownik może szybko sam wpisać dane testujące dla jakiejś skomplikowanej filtracji:

```
cat | prog1 | prog2 | prog3
```

Dodatkowe możliwości uzyskujemy dzięki wykorzystaniu w potokach list poleceń, np.:

```
... | ( cat; cat plik ) | ...
... | cat - plik | ...
```

W pierwszym przykładzie do danych przesyłanych potokiem drugie wywołanie `cat` dopisuje zawartość pliku, i całość łącznie przesyana jest dalej. Drugi przykład ilustruje konwencję, na mocy której nazwa pliku w postaci minusa nie jest traktowana jako nazwa pliku dyskowego, tylko oznacza strumień `stdin`.

Uniksowe filtry tekstowe — tee

Program **tee** (ang. trójnik – odwołuje się do analogii hydraulicznej) przekazuje dane z wejścia na wyjście bez zmian, ale dodatkowo zapisuje cały strumień danych na pliku (z opcją **-a** dopisuje). Dzięki temu można łatwo zarejestrować postać pośrednią danych przetwarzanych przez jakiś potok filtracji.

Przykładowym zastosowaniem jest debugowanie skryptów filtrujących:

```
... | filter1 | tee dbout_f1 | filter2 | tee dbout_f2 | ...
...
... | \
filter1 | if [ -z "$DEBUG" ]; then cat; else tee dbout_f1; fi | \
filter2 | if [ -z "$DEBUG" ]; then cat; else tee dbout_f2; fi | \
...
...
```

Często przydatne jest dodanie „zdalnego sterowania” tym procesem za pomocą zmiennej eksportowanej do wywoływanego skryptu:

```
... | \
filter1 | if [ -z "$DEBUG" ]; then cat; else tee dbout_f1; fi | \
filter2 | if [ -z "$DEBUG" ]; then cat; else tee dbout_f2; fi | \
...
...
```

Przykład: wyświetl sekcję pliku od wiersza nr 1550 do 1750:

```
head -1750 /opt/csv/apache/logs/access_log | tail +1550
# lub alternatywnie (która wersja jest najlepsza?)
head -1750 /opt/csv/apache/logs/access_log | tail -201
tail +1550 /opt/csv/apache/logs/access_log | head -201
```

Head i **tail** są przydatnymi i prostymi w użyciu filtrami. **Head** wyświetla na wyjściu początkowy fragment pliku, domyślnie pierwszych 10 wierszy. Analogicznie, **tail** wyświetla końcowy fragment pliku, domyślnie ostatnie 10 wierszy.

Uniksowe filtry tekstowe — head i tail

Head i **tail** są przydatnymi i prostymi w użyciu filtrami. **Head** wyświetla na wyjściu początkowy fragment pliku, domyślnie pierwszych 10 wierszy. Analogicznie, **tail** wyświetla końcowy fragment pliku, domyślnie ostatnie 10 wierszy.

Tail ma jeszcze jedną przydatną opcję (**-f**). Oznacza ona, że po doczytaniu do końca pliku (i wyświetleniu na wyjściu zadanego fragmentu) należy czekać, i ponawiać próbę czytania. Gdy na końcu pliku pojawią się dalsze dane, należy je wyświetlać bez żadnych ograniczeń. Pozwala to śledzić pliki, do których dopisują coś na końcu pracujące programy. Przykład:

```
tail -f /opt/csv/apache/logs/access_log
```

Uniksowe filtry tekstowe — cut

Cut wycina fragment wiersza. Można określić wycinanie konkretnych znaków, lub pól (słów) wiersza. Niestety, przydatność **cut** jest znacznie ograniczona przez trudność w zdefiniowaniu separatora pól. W odróżnieniu od innych filtrów, separatorem może być tylko pojedynczy, sztywno określony znak, i domyślnie jest to tabulator.

```
# listing katalogu - wybierz pierwszy znak i nazwę pliku
ls -1 /etc | cut -c1,51-
```

```
# listing katalogu - prawa dostępu i nazwa (PORAZKA)
ls -1 /etc | cut -d' , -f 1,14
```

Jak widać, w wielu przypadkach trudno wyciąć programem **cut** właściwy fragment wiersza, i to ogranicza jego przydatność do wierszy o ścisłe zdefiniowanym formacie, np. plików CSV (*Comma Separated Values*).

Uniksowe filtry tekstowe — sort

Sort sortuje wiersze z wejścia. Domyślne jest sortowanie alfabetyczne całego wiersza, jednak można opcjami wybrać szereg alternatywny, jak sortowanie numeryczne. Można zdefiniować **dowolne pole** (słowo) wiersza jako **klucz sortowania**. Można również zdefiniować klucz drugiego rzędu (sortowanie wierszy z równym kluczem podstawowym), i dla niego oddziennie wybrać kryterium sortowania, i podobnie klucze dalszych rzędów.

Przykłady:

```
# sortowanie zawartosci katalogu po dlugosci plikow
ls -1 | sort -n -k5
```

```
# sort.katalogu po pierwszej literze nazwy, potem dlugosci
ls -lt /etc | sort -k8.1,8.2 -k5n
# sortowanie zawartosci archiwum tar po czasie utw pliku
tar tvf archiw.tar | sort -k7n,7 -k4M,4 -k5n,5 -k6d,6 -k8
```

Sort jest często używanym programem i warto nabrać wprawy w jego użyciu.

Uniksowe filtry tekstowe — `uniq`

`Uniq` służy do wykrywania i usuwania powtarzających się wierszy w ciągu wejściowym. Jego zastosowanie jest bardziej specjalistyczne, i często w połączeniu z innymi filtrami.

Przykłady

```
# lista nieortograficznych słów z plików
cat *.tex | ispell -l -t -d polish | sort | uniq

# lista imion użytkowników systemu
getent passwd | cut -d: -f5 | cut -d' ' -f1 | sort | uniq

# ile plików było tworzonych w poszczególnych dniach
ls -l | awk '{print $6}' | sort | uniq -c

Jak widać w powyższych przykładach, uniq jest często używany łącznie z sort.
```

WAŻNE: oba pliki wejściowe muszą być posortowane według pola klucza.

Uniksowe filtry tekstowe — `join`

Podobnie jak `diff`, `join` nie jest typowym filtrem, ponieważ pracuje na danych z dwóch różnych plików. Leczy on rekordy z dwóch plików w sposób podobny do bazodanowego operatora `join`. Rekordy z obu plików są do siebie dopasowywane według wartości określonego pola rekordów, stanowiącego klucz.

```
# lista użytkowników z symbolicznymi nazwami grup
sort -t: -k4 /etc/passwd > /tmp/passwd
sort -t: -k3 /etc/group > /tmp/group
join -j1 4 -j2 3 -o 2.1,0.1.5 -t: /tmp/passwd /tmp/group

# połączenie dwóch list numerów telefonów
cat /tmp/phone | cat /tmp/fax
!Name Phone Number
Don +1 123-456-7890
Hal +1 234-567-8901
Yasushi +2 345-678-9011
#
join -t"<tab>" -a 1 -a 2 -e 'unknown' -o 0.1.2,2.2 \
/tmp/phone /tmp/fax
```

Uniksowe filtry tekstowe — `diff`

`Diff` nie jest typowym filtrem, ponieważ jego rola nie jest filtrowanie danych w potoku. Jednak `diff` jest programem niezwykle przydatnym w wielu pracach. Porównuje on dwa pliki tekstowe, wiersz po wierszu, znajduje sekcje różniące te dwa pliki, i wyświetla je w postaci ciągu takich bloków różnic. Pozwala to na szybkie zorientowanie się czym różnią się dwa pliki, pod warunkiem, że stanowią one nieznacznie różniące się od siebie wersje tego samego dokumentu, programu, specyfikacji, czy innego typu pliku.

```
diff /etc/nsswitch.conf_old /etc/nsswitch.conf
11c11
< hosts: files dns
---
> hosts: files dns
```

Powyższe pliki różnią się tylko wierszem numer 11. Znaki `<` i `>` symbolizują zamianę jakiej należałoby dokonać, aby zrównać pierwszy plik z drugim.

`Diff` ma szereg opcji ułatwiających znajdowanie różnic w różnych sytuacjach, oraz zmieniających sposób ich prezentacji na wyjściu. Często przydatne są opcje ignorowania odstępów (spacji, tabulatorów) w porównywaniu wierszy, np. `-w`

Uniksowe filtry tekstowe — `tr`

Program `tr` realizuje transliterację, czyli podmianę jednych liter (znaków) innymi w stringach. Przykłady:

```
# zamiana małych liter na wielkie
echo Ala ma kota. | tr '[a-z]', '[A-Z]'

# zamiana wszystkich nie-liter na podkreslinki
echo 'ABCd1234' | tr -c '[alpha:]', '[_*]'

# zamiana polskich liter na łacinskie
echo Zażółć gęśla jązú |\
tr 'ąęłńóśćżż' 'aelcnoszzAEELCNOSZZ'

# konwersja polskich znaków ISO8859-2 na CP1250 \
tr '\261\352\346\263\361\363\266\274\277' \
'\245\251\206\210\344\242\230\253\276',
```

Wyrażenia regularne (1): podstawowe

Jednoznakowe wyrażenia regularne:

- — kropka pasuje do każdego znaku, dokładnie jednego
- — ciąg znaków w nawiasach kwadratowych pasuje do każdego znaku z wymienionych, albo należącego do przedziału strzałka na początku w nawiasie kwadratowym oznacza do pełnienie, tu znak niealfanumeryczny
- — dowolny znak niespecjalny — pasuje wyłącznie do samego siebie

Powtarzenia:

- — ciąg wyrażeń dopasowuje się do ciągu znaków jeśli kolejne wyrażenia dopasowują się do kolejnych podciągów znaków gwiazdka następująca za wyrażeniem regularnym ϵ oznacza wielokrotne (0 lub więcej razy) powtóżenie dopasowania do kolejnych podciągów ciągu znaków; każdy podciąg jest oddzielnie dopasowywany do wyrażenia ϵ

„Kotwice”:

- — pasuje do zerowego ciągu znaków, ale tylko na początku ciągu
- — analogicznie pasuje tylko na końcu łańcucha znaków

Wyrażenia regularne (3): język grep'a i egrep'a

Poza podanymi wyżej podstawowymi konstrukcjami wyrażeń regularnych, kilka dalszych konstrukcji również istnieje we wszystkich implementacjach. Jednak z pewnych względów historycznych, znalazły się one w dwóch oddzielnich podzbiorach, które będą temu nazwać, ze względu na ich implementację w dwóch podobnych programach, językami wyrażeń regularnych **grep'a** i **egrep'a**.

Język wyrażeń regularnych grep'a — wyrażenia zapamiętane:

- — dopasowanie jak do wyrażenia ϵ , z zapamiętaniem dopasowanego ciągu znaków, można się do niego odwołać konstrukcją $\backslash 1$ w dalszej części wyrażenia

Język wyrażeń egrep'a — alternatywy, nawiasy, i niezerowe powtóżenia:

- — alternatywa, dopasowanie do wyrażenia ϵ_1 lub do wyrażenia ϵ_2
- — dopasowanie jak dla wyrażenia ϵ
- — dopasowanie jak dla wyrażenia ϵ , lub pasuje do ciągu pustego
- — oznacza powtóżenie podobnie jak $*$, ale co najmniej jednokrotne dopasowanie musi wystąpić

Tak jak można się tego spodziewać, dodatkowe znaki interpretowane specjalnie w wyrażeniach **egrep'a** są normalnymi znakami w języku **grep'a**, i na odwrót.

Wyrażenia regularne (2): proste przykłady

[0-9] — pasuje do pojedynczej cyfry (jak również dowolnego stringa, który taka zawiera)

[0-9]* — pasuje do dowolnego ciągu cyfr, w tym również pustego stringa

[1-9][0-9]* — pasuje do dowolnego niepustego ciągu cyfr, bez wiodącego零

[A-Z][a-z]* — pasuje do napisu znakowego zaczynającego się od dużej litery nie pasuje do napisu bez dużych liter

ale pasuje do napisu z więcej niż jedną dużą literą na początku!

również pasuje do dowolnego zapisu z dużymi literami na końcu!

~[A-Z][a-z]*\$ — pasuje do dowolnego ciągu małych liter, również pustego, zaczynającego się od jednej dużej litery, i niczego! innego

~[A-Z]\$[a-z]* — niepoprawne wyrażenie regularne

[a-z]* [a-z]* — pasuje do dwóch dowolnych ciągów liter rozdzielonych spacją

w tym również pasuje do dowolnego ciągu dowolnych znaków zawierającego spację

```
fraza="A mnie jest szkoda lata."
expr "$fraza" : "A mnie" # sukces
expr "$fraza" : "szkoda" # porażka
expr "$fraza" : "" # sukces
expr "$fraza" : ",*" # sukces
expr "$fraza" : "~" # porażka(???) stdout: 0
expr "$fraza" : ".*\\"(lata*)" # sukces
expr "$fraza" : ".*\\"(.*)" # porażka(???) stdout: (nic)
```

Szybki test: do czego dopasowują się następujące wyrażenia?

```
[A-Z][a-z][a-z]*[ ]*[A-Z][a-z][a-z]*[_a-z][_a-z0-9]*0[1-9][0-9]-[1-9][0-9][0-9]-[0-9][0-9][0-9]
```

Jak widać, w przypadkach dopasowania pustego stringa **expr** sygnalizuje brak dopasowania. Szczególnie ostatni przypadek, gdzie oba podwyrażenia ***** konsumują zeroowe podstringi, wydaje się myły.

Wyszukiwanie wzorców: grep i egrep

Grep znajduje dopasowanie podanego wyrażenia regularnego we wszystkich wierszach strumienia wejściowego, lub zadanych plikach. Spróbuj rozszerzyć znaczenie poniższych przykładów:

```
grep money *
cat * | grep money
grep -n Count *.ch
grep -i kowalski spis.telef
ls -l | grep -v .cho$'
ls -l | grep '^.....w'
grep '^[:]*::' /etc/passwd
cat dictionary | grep '^...w.w..e.t$'
grep '^From, $MAIL | grep -v 'From szef',
cat text | grep '\([-A-Za-z][-A-Za-z]*\)[ ]*\^1 ,
egrep 'socket|pipe|msegget|semget|shmget', *.ch]
egrep 'socket|pipe|msegget|semget|shmget', *.ch]
```

Format polecenia: [adres1,adres2] operator [argumenty[modyfikatory]]

Adres w poleceniu seda może być liczba lub wzorcem (wyrażeniem regularnym). Operacja jest wykonywana tylko na wierszu, którego dotyczy adres, albo w przedziale wierszy określonym adresami (jeśli są dwa).

d wykasuje zawartość bufora (nic nie będzie wyświetlane na wyjście)

Operatory: p wyslij na wyjście zawartość bufora (oprócz wyśw. domyślnego)
q zakończ pracę (po przetworzeniu bieżącego wiersza)

```
sed 10q          # przepuszcza 10 pierwszych linii
sed /wzorzec/q   # wyświetla do linii z wzorcem
sed /wzorzec/d   # opuszcza linie z wzorcem (grep -v)
sed '/^$/d'      # opuszcza puste linie
sed -n /wzorzec/p # wyświetla tylko linie z wzor. (grep)
sed -n '/\begin{verbatim}/,\end{verbatim}/p'
```

Wyrażenia regularne (4): grep i egrep — zestawienie

Poniższe wyrażenia przedstawione są w kolejności malejącego priorytetu:

z dowolny znak niespecjalny pasuje do siebie samego
\z kasuje specjalne znaczenie znaku *z*
^ początek linii
\$ koniec linii
. dowolny pojedynczy znak
[abc...] dowolny znak spośród podanych, też przedziałów, np. a-zA-Z
[^abc...] dowolny znak spoza podanych, również mogą być przedziały
\n to do czego dopasowało się *n*-te wyrażenie \(\epsilon\)
*\epsilon** (tylko grep) zero lub więcej powtórzeń wyrażenia *\epsilon*
\epsilon+ jedno lub więcej powtórzeń wyrażenia *\epsilon* (tylko grep)
\epsilon? zero lub jedno wystąpienie wyrażenia *\epsilon* (tylko grep)
\epsilon_1\epsilon_2 *\epsilon_1* i następujące po nim *\epsilon_2*
\epsilon_1 \epsilon_2 *\epsilon_1* lub *\epsilon_2* (tylko grep)
\(\epsilon\)\(\epsilon\) zapamiętane wyrażenie regularne *\epsilon* (tylko grep)
(\epsilon) wyrażenie regularne *\epsilon* (tylko grep)

Sed: edytor strumieniowy

Edytor strumieniowy sed (stream editor) wczytuje dane z wejścia wierszu, na wczytanym wierszu wykonuje operacje zadane argumentem, i przetworzony wiersz wysyła na wyjście.

Oprócz przedstawionych wyżej operatorów seda: d, p, i q, najczęściej przydatnym jest operator podmiany stringów s. Wymaga on podania dwóch stringów jako argumentów po symbolu operatora. Pierwszym znakiem po s jest wybrany znak separatora, a potem dwa argumenty. Normalnie podmiana jest wykonywana jeden raz w wierszu, ale podanie modyfikatora g powoduje wykonanie podmiany dowolna liczby razy.

```
Sed 's/marzec/March/g'          # podmiana stringow (wiele razy)
sed 's/^~/I/'                   # indentacja (taby na pocz. linii)
sed '/./s/^~/I/'                # ulepszona indentacja
```

Pierwszy argument operatora s jest traktowany jako wyrażenie regularne typu grep'a, tzn. może zawierać operacje zapamiętywania \(\cdot\cdot\cdot\). Wtedy drugi argument może zawierać odwołania do zapamiętanych stringów \1, \2, itd. W przypadku wersji Gnu seda, możliwe jest również alternatywne stosowanie wyrażeń regularnych egrep'a. Operacja zapamiętywania jest wtedy niedostępna.

Sed posiada jeszcze kilka bardziej skomplikowanych operatorów, które wraz z sekwencjami pozwalają na pisanie złożonych wyrażeń, które są niekiedy bardzo trudne do zrozumienia i debugowania. Nie zmienia to faktu, że bardzo wiele przydatnych operacji można zrealizować czterema powyższymi operatorami.

Sed: przykład (1) komedia pomyłek

```
sierra-90> who
NAME      LINE        TIME      IDLE    PID   COMMENTS
witold   + vt04       Oct 21 04:46  2:45   238
witold   + ttyp0      Oct 21 04:46  2:43   292
witold   + ttyp1      Oct 21 04:46  .      291
witold   + ttyp2      Oct 21 04:46  .      290
sierra-91> who | sed 's/ .* / /,
NAME COMMENTS
witold
witold 292
witold 291
witold 290
sierra-92> who | sed 's/ .* [^ ] / /'
NAME COMMENTS
witold 38
witold 92
witold 91
witold 90
sierra-93> who | sed 's/ .* \([^\r\n ]\)/ \1 /'
```

Sed: przykład (2) — konwersja Latexa do HTMLa

```
# znaki specjalne HTML'a
s/\\&/\\&/g
s/</&lt;/;g
s/>/&gt;/g
;
# puste wiersze i komentarze
/^[\n]*$/i \
<p>
</p>
/^[\n]*%(\.*\)$/<!-- \1 -->/
];
# string cytowany \verb to prawdziwy problem
s#\\"verb\(.*)\(\[^\\n]*\)\1#\<tt\>\2</tt\>#g
];
# jednoznakowe roznosci
s/\n\n<br\/>/g
s/\n\n([#_\$\n])\n/\1/g
];
# te znaczniki mają swoje odpowiedniki
s#\\"underline{\([^\n]*\)}#\<u\>\1</u\>#g
s#\\"section{\([^\n]*\)}#\<h1\>\1</h1\>#
s#\\"subsection{\([^\n]*\)}#\<h2\>\1</h2\>#
s#\\"begin{enumerate}#\<ol\>#;
s#\\"begin{itemize}#\<ul\>#;
s#\\"begin{description}#\<dl\>#;
s#\\"item#\<li\>#;
```

Sed: kontynuacja przykładu

Jako wniosek z analizy powyższego przykładu, rozważmy zadanie napisania skryptu **sed**, który, filtrując ciąg wejściowy, wyświetli na wyjściu tylko pierwsze słowo (dla uproszczenia) z każdego wiersza. Rozważ poniższe rozwiązań tego zadania. Które z nich zawierają błędy, a które działają w pełni niezawodnie? Czym różni się działanie tych wersji „niezawodnych” między sobą?

```
sed 's/ .*$//,
sed 's/([^\r\n ] ) .*$/\1/',
sed 's/([^\r\n ] *) .*$/\1/',
sed 's/([a-zA-Z]* ) .*/\1/',
sed 's/([a-zA-Z][a-zA-Z]*) .*/\1/',
sed 's/([^\r\n ] ) .*/\1/',
sed 's/([^\r\n ] *) .*/\1/',
sed 's/([^\r\n ] * ) .*/\1/,'
```

Dla porównania rozważ możliwość wykorzystania następujących mechanizmów POSIXowych (patrz poniżej) do realizacji zadania: **<...>[:alpha:]** i **[space:]**. Spróbuj napisać dobrze rozwiązań problemu wykorzystując te mechanizmy. Które z nich stanowią istotne ulepszenie wersji nie-POSIXowej?

Sed: przykłady zaawansowane

Poniższy przykładowy skrypt **sed** skraca sekwencję pustych linii do pojedynczej pustej linii wykorzystując polecenie wczytywania kolejnych wierszy (**N**) i pełną zrealizowaną przez skok do etykiet (**b**):

```
# pierwszy pusty wiersz jawnie wypuszczamy na wyjście
/^$/p
:Empty
# dodajemy kolejne puste wiersze usuwając znaki NEWLINE
/^$/f N;s///;b Empty
}
# mamy wczytany niepusty wiersz, wypuszczamy go
{p;d;}
```

Skrypt w pełni kontroluje co jest wyświetlane na wyjściu i działa tak samo wywołany z opcją **-n** jak i bez niej.

Podobnie jak następujący, zaledwie dziesięciorzakowy skrypt który wyświetla plik wejściowy w odwrotnej kolejności wierszy: **!G;h;\$p;d**

Sed: podstawowe operatory

a\b etyk skończ na wyjście kolejne linie do linii nie zakończonej \\\
b\b etyk skok do etykiety
c\b zmień linie na następujący tekst, jak dla a
d\b skasuj linie
i\b wprowadź następujące linie przed innym wyjściem
1\b wyświetli linie, z wizualizacją znaków specjalnych
p\b wyświetl linie zakończ
q\b wczytaj plik, wypuść na wyjście
r\b plik zastap stary tekst *s₁* nowym *s₂*; jeden raz gdy brak modyfikatora z, wszystkie gdy z=g, wyświetlaj podstawienia gdy z=p, zapisz na pliku
gdy z=w plik skok do etykiety, gdy w bieżącej linii dokonane podstawienie
t\b etyk w\b plik zapisz linię na pliku
y/s₁/s₂/ zamień każdy znak z *s₁* na odpowiedni znak z *s₂*
= wyświetli bieżący numer linii
!polec wykonaj polecenie sada *polec* gdy bieżąca linia nie wybrana
:etyk etykieta dla polecień *b* i *t*
\{...\} grupowanie polecień

Wyrażenia regularne (5): POSIX — BRE i ERE

Specyfikacja POSIX porządkuje i rozszerza oryginalną koncepcję wyrażeń regularnych Uniksa. Uwzględnia ona, między innymi, specyfikację powtórzeń, klasy znaków, oraz lokalizacje, tzn. stosowany w danej lokalizacji zestaw znaków i konwencje równoważności i uporządkowania. Stanowi rozszerzenie wyrażeń regularnych grep i egrep, ale ze względu na ich wzajemną niekompatybilność, jej wynikiem jest definicja dwóch języków wyrażeń regularnych: BRE (Basic Regular Expressions) i ERE (Extended Regular Expressions).

W największym skrócie, warto zapamiętać:

BRE (zgodne z grepem) — wyrażenia regularne z operatorem zapamiętywania \(...\)\> i odwoływanie się do zapamiętanych stringów \1, \2, ...
ERE (zgodne z egrepem) — wyrażenia regularne z operatorem alternatywy \...|..., wyrażenia w nawiasach (...), wystąpienia opcjonalne ...?, oraz powtózenia co najmniej jeden raz ...+.

Oprócz powyższych, języki BRE i ERE różnią jeszcze szeregiem bardziej subtelnych drobiazgów, które nie będą tu szczegółowo omawiane.

Wyrażenia regularne (6): POSIX — inne konstrukcje

Standard POSIX wprowadził dodatkowo powtórzenia n-krotne:

$\epsilon^{\{n, m\}}$ powtóżenie: co najmniej *n*-razy, co najwyższej *m*-razy (grep)
 $\epsilon^{\{n, m\}}$ powtóżenie: co najmniej *n*-razy, co najwyższej *m*-razy (egrep)

Jednej z wartości *n* lub *m* można nie podać, co oznacza ograniczenie liczby powtórzeń tylko od dołu lub tylko od góry, o ile przecinek jest obecny. Podana jedna wartość, i brak przecinka, oznacza powtóżenie dopasowania ścisłe określona liczbę razy.

Inną, niezwykle przydatną konstrukcją, wprowadzoną w standardzie POSIX, są operatory \<...> wymuszające dopasowanie tylko na granicy słowa.

Standard POSIX rozszerzył też operator [] dopasowujący jeden znak o klasę znaków za pomocą wyrażenia [:klasa:], z następującymi klasami:

| | | | | | | | | | |
|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|------------|
| [:alnum:] | [:alpha:] | [:blank:] | [:cntrl:] | [:digit:] | [:graph:] | [:lower:] | [:print:] | [:space:] | [:xdigit:] |
|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|-----------|------------|

Wyrażenia regularne (7): przykłady wyrażeń ERE

Niektóre wyrażenia mają złożoną składnię i wymagania. Na przykład, adresy email:

`username@domain-spec`

Nazwa użytkownika musi być dowolnym niepustym ciągiem liter, cyfr, kropki, podkreśnika (podlogi), minusa i plusa.

Specyfikacja domeny musi składać się z niepustej liczby powtórzeń domen, rozdzielonych kropkami.

Domena musi być niepustym ciągiem liter, cyfr, minusa i podkreśnika (podlogi). Plusy i kropki są wykluczone (ale kropki występują między domenami).

Jako przypadek szczególny, domena główna (ostatni człon) musi składać się wyłącznie z liter, jednak nie mniej niż dwóch i nie więcej niż pięciu.

`~[a-zA-Z0-9_-\._\+\+)([a-zA-Z]\{2,5\})+[a-zA-Z]\{2,5\}$`

Zauważmy wygodę wielokrotnego użycia operatorów ERE powtóżenia 1 lub więcej razy (`+`) oraz operatora powtózenia od 2 do 5 razy (`{2,5}`).

Awk: uniwersalny filtr programowalny

`Awk` jest filtrem działającym, podobnie jak `sed`, na kolejnych wierszach. Jednak zamiast prostych operatorów o jednoznakowych nazwach, `awk` ma konstrukcje programowe przypominające język C. Dwukrokowy algorytm działania `awk`:

- 1.czyta wiersz z wejścia, dzieli na pola (słowa) dostępne jako: `$1, $2, ...`,
2. wykonuje cały swój program składający się z szeregu par: warunek-akcja.

Uwagi:

- par warunek-akcja może być wiele i w każdej może brakować warunku (domyślnie: prawda) albo akcji (domyślnie: wyświetlenie wiersza na wyjściu),
- w programie można używać zmiennych, które zachowują wartości pomiędzy wywołaniami programu dla kolejnych wierszy,
- zmiennych nie trzeba deklarować ani inicjalizować; w pierwszym użyciu są one inicjalizowane wartością 0 lub pustym stringiem, zależnie od operacji.

```
ls -1 ~student | awk '$5 > 100000 ,  
awk , {print $2, $1} , nazwa_pliku  
cat /etc/passwd | awk -F: '{print $4,$3}' | sort
```

Wyrażenia regularne (8): GNU grep

Wersja GNU programu `grep` implementuje całą funkcjonalność `grep` i `egrep`. Co więcej, wprowadza rozszerzenia pozwalające łączyć operacje tradycyjne dostępne tylko dla `grep` jak i `egrep`.

```
# operator dopasowania stringa do wyrażenia regularnego  
awk -F: , $7 ~ /bash$/ { print $1,$7 }' /etc/passwd  
  
# użycie zmiennych do zapamietania kontekstu między wierszami  
awk , $1 != prev { print; prev = $1 } ,  
  
# użycie zmiennych wbudowanych awk: NF i NR  
awk , NF > 5 { printf "Wiersz %d ma %d slow. "%NR,NF } ,  
  
# obliczanie długości stringa  
awk , { wd+=NF; ch+=length($0)+1 } END { print NR,wd,ch } ,  
  
# warunki specjalne do inicjalizacji i finalizacji  
awk , $1 < 0 { $1 = 0 } $1 > 100 { $1 = 100 } { print $0 } ,  
  
# połączenie z mechanizmami shella w wierszu wywołania  
awk , { s += $1 } END { print s } ,  
awk , { s += '$nr_pola' } END { print s } ,
```

Awk: użycie tablic

Awk pozwala na użycie tablic, jednak trochę innych niż typowe tablice w językach programowania. Tablice są indeksowane stringami, i nie deklaruje się ich rozmiaru. Z tego powodu nazywa się je **tablicami asocjacyjnymi**.

```
# sumowanie dowolnej liczby pozycji według nazwy
awk '{ sum[$1] += $2 } \
END { for (name in sum) print name, sum[name] } ,'

# zliczanie częstotliwości występowania słów w tekście
awk 'for (i=1; i<NF; i++) freq[$i++] \
END { for (word in freq) print word, freq[word] } ,'
```

Można też używać dwóch lub więcej indeksów tablicy. Warto jednak wiedzieć, że awk używa ich łącznie, jako jednego indeksu składającego się z obu stringów, plus oddzielającego je przecinką.

Awk: zmienne wbudowane

| | |
|----------|---|
| FILENAME | nazwa bieżącego pliku wejściowego |
| FS | znak podziela pół (domyślnie spacja i tab) |
| NF | liczba pól w bieżącym rekordzie |
| NR | numer kolejny bieżącego rekordu |
| OFMT | format wyświetlanego liczb (domyślnie %g) |
| OFS | napis rozdzielający pola na wyjściu (domyślnie spacja) |
| ORS | napis rozdzielający rekordy na wyjściu (domyślnie linefeed) |
| RS | napis rozdzielający rekordy na wejściu (domyślnie linefeed) |

Awk: uwagi o przenośności

Oryginalny uniwersalny awk był dość ograniczonym programem, i wkrótce po jego powstaniu pojawiła się wersja rozszerzona. Niestety, nie mogło się to dokonać w sposób całkowicie przenośny i nowa wersja zaczęła być instalowana pod nazwą **nawk** równolegle ze starej, do której instalowano link o nazwie **awk**. Jednak w duchu utrzymania kompatybilności z wcześniejszymi skryptami, które nie mają świadomości nowszych wersji, polecenie **awk** na wielu systemach uniksowych wywołuje bardzo ograniczonego oryginalnego **awk'a**.

Często dobrym sposobem jest **wywołanie awk jako nawk** — na wielu systemach istnieje taki program albo link. Jest to dobra forma przenośnego wywołania **awk** zapewniająca odcięcie się od wersji najstarszej.

Jednak istnieją jeszcze inne wersje **awk'a**, i na systemach linuksowych często instalowany jest program **mawk**, niestety różniący się drobnymi elementami. Z reguły główna robocka wersja **awk'a** na każdym systemie jest instalowana również pod nazwą **nawk** na potrzeby skryptów napisanych przenośnie godnie z powyższą zasadą.

Awk: operatory

| | |
|------------------------------------|---|
| w kolejności rosnącego priorytetu: | |
| = += -= *= /= %= | operatorzy przypisania podobne jak w C |
| | alternatywna logiczna typu „short-circuit” |
| && | koniunkcja logiczna typu „short-circuit” |
| ! | negacja wartości wyrażenia |
| > >= < <= != | operatorzy porównania |
| ~ ! ~ | (nie)dopasowanie wyrażeń regularnych do napisów |
| nic | konkatenacja napisów |
| + - | plus, minus |
| * / % | mnożenie, dzielenie, reszta |
| ++ -- | inkrement, dekrement (prefix lub postfix) |

Awk: funkcje wbudowane

```
kosinus, argument w radianach  
expr  
czyta następną linię z wejścia  
pozycja napisu s2 w s1; zwraca 0 jeśli nie ma  
części całkowita  
długość napisu znakowego  
logarytm naturalny  
sinus, argument w radianach  
podziel napis s względem c na części do tablicy a  
formatowanie napisu  
n-znakowy podciąg s począwszy od pozycji m  
substr(s,m,n)
```

Warto znać podstawowy zestaw filtrów tekstowych Uniksa, ponieważ realizują one bardzo proste algorytmy, które łatwo zapamiętać i ich używać. Jednocześnie łączenie tych filtrów pozwala czasem zaimplementować całkiem zaawansowane funkcje.

| | |
|-------|--|
| tac | wyswietlaj zawartość plików od końca |
| rev | wyswietlaj pliki odwracając kolejność znaków w wierszach |
| paste | łącz i wyświetl jako jeden wiersz kolejne wiersze z plików |

Inne przydatne filtry Uniksa

łączenie filtrów

Silą filtrów Uniksa leży w prostocie ich funkcjonalności, i łatwości łączenia w bardziej skomplikowane wyrażenia. Ilustracją tego może być poniższy przykład, który w pięciu operacjach znajduje 10 najczęstszych występujących słów w dowolnym zbiorze tekstów:

```
cat * | tr -cs "[:alpha:]" "\012*" \  
| sort \  
| uniq -c \  
| sort -nr \  
| head
```