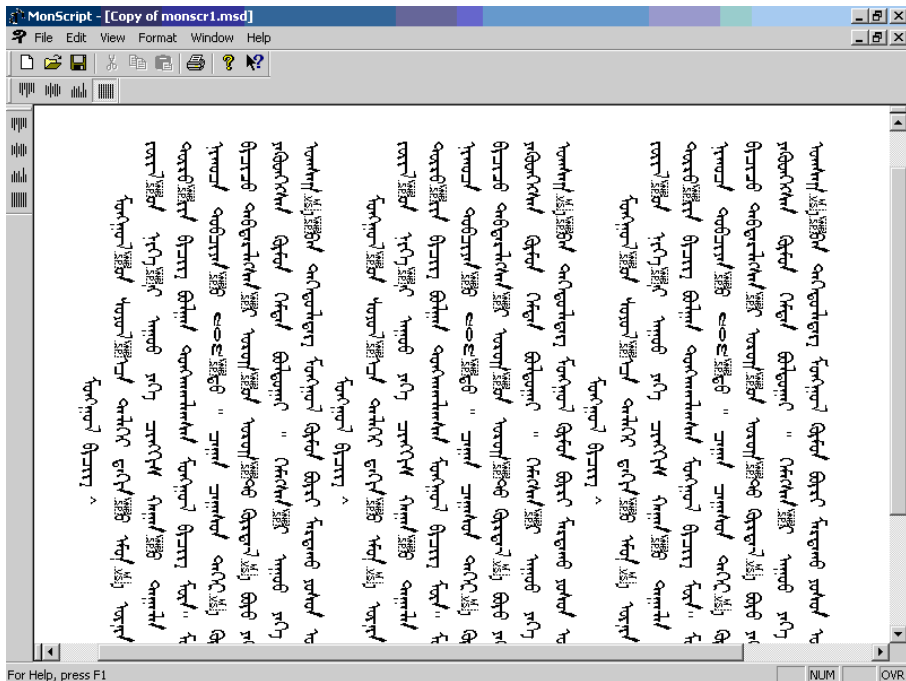


# Текст боловсруулалтын үндэс

Лекц №2

Регуляр  
илэрхийлэл

2020 он



# Регуляр илэрхийлэл

- Текстийн тэмдэгт мөрийг тодорхойлдог формал хэл
- Бид дараах үгсийн аль нэгийг хэрхэн олох вэ?
  - хэрэм
  - хэрэмнүүд
  - Хэрэм
  - Хэрэмнүүд



# Регуляр илэрхийлэл: Тусгаарлалт (Disjunctions)

- [] хаалт доторх үсэг

Паттерн	Нийцэл(Match)
[xX]эрэм	хэрэм, Хэрэм
[1234567890]	дурын тоо

- Муж [A-Z]

Паттерн	Нийцэл	
[A-Я]	Том үсэг	<u>Т</u> ом үсэг
[a-я]	Жижиг үсэг	<u>Ж</u> ижиг үсэг
[0-9]	Нэг оронтой тоо	Гарчиг <u>1</u> : Регуляр

# Регуляр илэрхийлэл: Тусгаарлалт дахь үгүйсгэл

- Үгүйсгэл `[^Ss]`
  - Энэ нь хаалтан дахь үсгүүдээс бусад нь гэсэн үг

Паттерн	Нийцэл
<code>[^А-Я]</code>	Том үсэг биш
<code>[^Хх]</code>	Х биш, мөн х биш
<code>[^e^]</code>	е болон ^ биш
<code>a^b</code>	<code>a^b</code>

< > ↺ 88 |  www.regexpal.com

RegEx Pal

From Dan's Tools

Regular Expression

`/[^e^]/g`

Test String

we looked  
then we saw him step

# Регуляр илэрхийлэл: Бусад тусгаарлалтууд

- Эсвэл холбоос |

Паттерн	Нийцэл
groundhog   woodchuck	
yours   mine	yours mine
a   b   c	= [abc]
[gG]roundhog   [Ww]oodchuck	



Photo D. Fletcher

# Регуляр илэрхийлэл: ? \* + .

Паттерн	Нийцэл	
<code>colou?r</code>	Өмнөх үсэг байхгүй байж болно	<u>color</u> <u>colour</u>
<code>oo*h!</code>	Өмнөх үсэг 0 эсвэл олон удаа давтагдаж болно	<u>oh!</u> <u>ooh!</u> <u>oooh!</u> <u>ooooh!</u>
<code>o+h!</code>	Өмнөх үсэг 1 эсвэл олон давтагдаж болно	<u>oh!</u> <u>ooh!</u> <u>oooh!</u> <u>ooooh!</u>
<code>baa+</code>		<u>baa</u> <u>baaa</u> <u>baaaa</u> <u>baaaaa</u>
<code>beg.n</code>	Ямар нэгэн	<u>begin</u> <u>begun</u> <u>begun</u> <u>beg3n</u>



Stephen C Kleene

Kleene \*, Kleene +

# Регуляр илэрхийлэл: **^** -эхлэх **\$** -төгсөх

Паттерн	Нийцэл
<b>^</b> [A-Z]	<u>P</u> alo Alto
<b>^</b> [^A-Za-z]	<u>1</u> <u>"Hello"</u>
\. <b>\$</b>	The end <u>.</u>
. <b>\$</b>	The end <u>?</u> The end <u>!</u>

# Жишээ

- Текст дэх “the” артиклийн бүх тохиолдол олох.

The

Буруу жишээ

[tT]he

other эсвэл theology буруу үр дүн буцаана

[^a-zA-Z][tT]he[^a-zA-Z]



# Алдаа

- Хоёр төрлийн алдааг засах дээр суурилах ёстой
  - Нийцэх ёсгүй нийцсэн тэмдэгт мөрүүд (there, then, other)
    - False positives (Type I)
  - Нийцэх ёстой нийцээгүй тэмдэгт мөрүүд (The)
    - False negatives (Type II)

# Алдаа

- ЭХ боловсруулалтад энэ төрлийн алдаануудтай байнга ажилладаг.
- Програмын алдааны түвшинг багасгахын тулд хоёр эсрэг оролдлогыг авч үздэг:
  - Зөв утгатай хэр ойр(accuracy) эсвэл хэр давтамжтай (precision) байгааг нэмэгдүүлэх нь (false positives хамгийн бага байлгахад туслана)
  - Хамрах хүрээ эсвэл recall -ыг нэмэгдүүлэх нь (false negatives хамгийн бага байлгахад тусална)

# Дүгнэлт

- Регуляр илэрхийллийн үүрэг их
  - регуляр илэрхийллүүд нь текст боловсруулах текстүүдийн хувь эхний модел байдаг
- Хүнд бодлогуудын даалгаварын хувьд, машин сургалтын ангилагчуудыг ашиглана
  - Гэвч регуляр илэрхийлэл ангилагч дотор өргөн хэрэглэгддэг
  - Ерөнхий тохиолдлыг илрүүлэхэд хэрэгтэй