

# Lexical analysis and parsing

## A. Reserved words

Reserved words are not reserved when used as fields. So `return=1` is illegal, but `foo.return=1` is fine.

## B. Round, square, and curly brackets

In Python and in C it's simple -- function are called using parentheses, and array subscripted with square brackets. In Matlab and in Fortran both are done with parentheses, so there is no visual difference between function call `foo(x)` and indexing array `foo`.

## C. Handling the white space.

There are four kinds of ws, each recognized by a dedicated rule: NEWLINE, COMMENT, ELLIPSIS, and SPACES. Only NEWLINE returns a token, the three others silently discard their input. NEWLINE collects adjacent `\n` characters and returns a single SEMI token:

```
def t_NEWLINE(t):
    r'\n+'
    t.lexer.lineno += len(t.value)
    if not t.lexer.parens and not t.lexer.braces:
        t.value = ";"
        t.type = "SEMI"
    return t
```

## Comments come in two flavors -- regular and MULTILINE.

The regular comments are discarded, while MULTILINE become doc strings, and are handled as expr statements. regular consume everything from % or # up to but not including n:

```
(%|\#).*
```

Multiline COMMENT rule drops leading blanks, then eats everything from % or # up to the end of line, *including* newline character. COMMENT leaves one newline character, or else funny things happen:

```
@TOKEN(r"(%|\#).*")
def t_COMMENT(t):
    if not options.do_magic or t.value[-1] != "!":
        t.lexer.lexpos = t.lexer.lexdata.find("\n", t.lexer.lexpos)
```

Multiline comments:

```
(^[ \t](%|\#).*\n)+
```

The pattern for multi-line comments works only if `re.MULTILINE` flag is passed to ctor.

Comments starting with `%!`  have a special meaning TBD

ELLIPSIS is the matlab way for line continuation. It discards everything between `...` and the newline character, including the trailing `\n`:

```
def t_ELLIPSIS(t):
    r"\.\.\.\.*\n"
    t.lexer.lineno += 1
    pass
```



- (1) a name or a number
- (2) literal string using single or double quote
- (3) (T) or [T] or {T} or T' or +T or -T

Terms end with:

- (1) an alphanumeric character \w
- (2) single quote (in octave also double-quote)
- (3) right parenthesis, bracket, or brace
- (4) a dot (after a number, such as 3.

The pattern for whitespace accounts for ellipsis as a whitespace, and for the trailing junk.

Terms start with:

- (1) an alphanumeric character
- (2) a single or double quote,
- (3) left paren, bracket, or brace and finally
- (4) a dot before a digit, such as .3 .

TODO: what about curly brackets ???

TODO: what about dot followed by a letter, as in field  
[foo .bar]

```
t.lexer.lineno += t.value.count("\n")
t.type = "COMMA"
return t
```

## Class matlabarray

Matlab arrays differ from numpy arrays in many ways, and class matlabarray captures the following differences:

### A. Base-one indexing

Following Fortran tradition, matlab starts array indexing with one, not with zero. Correspondingly, the last element of a N-element array is N, not N-1.

### B. C\_CONTIGUOUS and F\_CONTIGUOUS data layout

Matlab matrix elements are ordered in columns-first, aka F\_CONTIGUOUS order. By default, numpy arrays are C\_CONTIGUOUS. Instances of matlabarray are F\_CONTIGUOUS, except if created empty, in which case they are C\_CONTIGUOUS.

matlab	numpy
<pre>&gt; reshape(1:4,[2 2]) 1 3 2 4</pre>	<pre>&gt;&gt;&gt; a=matlabarray([1,2,3,4]) &gt;&gt;&gt; a.reshape(2,2,order="F") 1 3 2 4  &gt;&gt;&gt; a.reshape(2,2,order="C") 1 2 3 4</pre>

```
>>> a=matlabarray([1,2,3,4])
>>> a.flags.f_contiguous
True
>>> a.flags.c_contiguous
False
```

```
>>> a=matlabarray()
>>> a.flags.c_contiguous
True
>>> a.flags.f_contiguous
False
```

### C. Auto-expanding arrays

Arrays are auto-expanded on out-of-bound assignment. Deprecated, this feature is widely used in legacy code. In smop, out-of-bound assignment is fully supported for row and column vectors, and for their generalizations having shape

[1 1 ... N ... 1 1 1]

These arrays may be resized along their only non-singular dimension. For other arrays, new columns can be added to F\_CONTIGUOUS arrays, and new rows can be added to C\_CONTIGUOUS arrays.

matlab	numpy
<pre>&gt; a=[] &gt; a(1)=123 &gt; a 123</pre>	<pre>&gt;&gt;&gt; a=matlabarray() &gt;&gt;&gt; a[1]=123 &gt;&gt;&gt; a 123</pre>

### D. Create by update

In matlab, arrays can be created by updating a non-existent array, as in the following example:

```
>>> clear a
>>> a(17) = 42
```

This unique feature is not yet supported by smop, but can be worked around by inserting assignments into the original matlab code:

```
>>> a = []
>>> a(17) = 42
```

### E. Assignment as copy

Array data is not shared by copying or slice indexing. Instead there is copy-on-write.

### F. Everything is a matrix

There are no zero or one-dimensional arrays. Scalars are two-dimensional rather than zero-dimensional as in numpy.

### G. Single subscript implies ravel.

TBD

### H. Broadcasting

Boadcasting rules are different

## I. Boolean indexing

TBD

## J. Character string constants and escape sequences [ffd52d5fc5]

In Matlab, character strings are enclosed in single quotes, like 'this', and escape sequences are not recognized:

```
matlab> size('hello\n')
1    7
```

There are seven (!) characters in 'hello\n', the last two being the backslash and the letter n.

Two consecutive quotes are used to put a quote into a string:

```
matlab> 'hello''world'
hello'world
```

In Octave, there are two kinds of strings: octave-style (enclosed in double quotes), and matlab-style (enclosed in single quotes). Octave-style strings do understand escape sequences:

```
matlab> size("hello\n")
1    6
```

There are six characters in "hello\n", the last one being the newline character.

Octave recognizes the same escape sequences as C:

```
\ "  \a  \b  \f  \r  \t  \0  \v  \n  \\  \nnn  \xhh
```

where n is an octal digit and h is a hexadecimal digit.

Finally, two consecutive double-quote characters become a single one, like here:

```
octave> "hello""world"
hello"world
```

---

# Data structures

## A. Empty vector [], empty string "", and empty cellarray {}

matlab	numpy
<pre>&gt; size([]) 0 0  &gt; size('') 0 0  &gt; size({}) 0 0</pre>	<pre>&gt;&gt;&gt; matlabarray().shape (0, 0)  &gt;&gt;&gt; char().shape (0, 0)  &gt;&gt;&gt; cellarray().shape (0, 0)</pre>

## B. Scalars are 1x1 matrices

matlab	numpy
<pre>&gt; a=17 &gt; size(a) 1 1</pre>	<pre>&gt;&gt;&gt; a=matlabarray(17) &gt;&gt;&gt; a.shape 1 1</pre>

## C. Rectangular char arrays

Class char inherits from class matlabarray the usual matlab array behaviour -- base-1 indexing, Fortran data order, auto-expand on out-of-bound assignment, etc.

matlab	numpy
<pre>&gt; s='helloworld' &gt; size(s) 1 10 &gt; s(1:5)='HELLO' &gt; s HELLOworld &gt; resize(s,[2 5]) HELLO world</pre>	<pre>&gt;&gt;&gt; s=char('helloworld') &gt;&gt;&gt; print size_(s) (1,10) &gt;&gt;&gt; s[1:5]='HELLO' &gt;&gt;&gt; print s HELLOworld &gt;&gt;&gt; print resize_(s,[2,5]) HELLO world</pre>

## D. Row vector

matlab	numpy
<pre>&gt; s=[1 2 3]</pre>	<pre>&gt;&gt;&gt; s=matlabarray([1,2,3])</pre>

## E. Column vector

matlab	numpy
<pre>&gt; a=[1;2;3]  &gt; size(a) 3 1</pre>	<pre>&gt;&gt;&gt; a=matlabarray([[1],                   [2],                   [2]])  &gt;&gt;&gt; a.shape (3, 1)</pre>

## F. Cell arrays

Cell arrays subclass matlabarray and inherit the usual matlab array behaviour -- base-1 indexing, Fortran data order, expand on out-of-bound assignment, etc. Unlike matlabarray, each element of cellarray holds a python object.

matlab	numpy
<pre>&gt; a = { 'abc', 123 } &gt; a{1} abc</pre>	<pre>&gt;&gt;&gt; a=cellarray(['abc',123]) &gt;&gt;&gt; a[1] abc</pre>

## G. Cell arrays of strings

In matlab, cellstrings are cell arrays, where each cell contains a char object. In numpy, class cellstring derives from matlabarray, and each cell contains a native python string (not a char instance).

matlab	numpy
<pre>&gt; a = { 'abc', 'hello' }  &gt; a{1} abc</pre>	<pre>&gt;&gt;&gt; a=cellstring(['abc',                   'hello'])  &gt;&gt;&gt; a[1] abc</pre>

---

### Data structures

All matlab data structures subclass from matlabarray

### Structs

TBD

### Function pointers

Handles @

### String concatenation

Array concatenation not implemented

```
>>> ['hello' 'world']  
helloworld
```