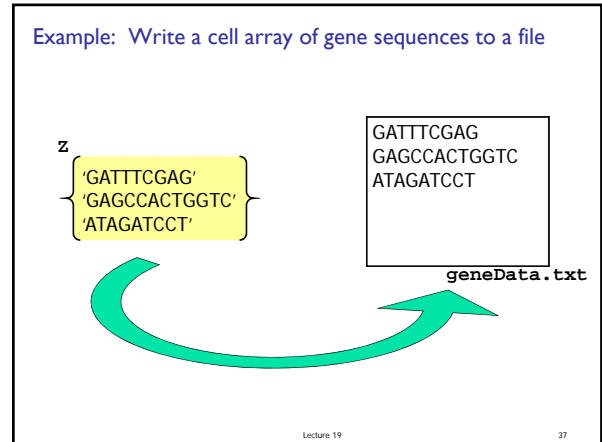


- Previous Lecture:
 - Cell arrays
- Today's Lecture:
 - File input/output
- Announcements:
 - Project 5 due Thurs Nov 6th at 11pm
 - Next week regular lecture—live, in person—resumes



- A 3-step process to read data from a file or write data to a file
1. (Create and) **open** a file
 2. **Read** data from or **write** data to the file
 3. **Close** the file

I. Open a file

```
fid = fopen('geneData.txt', 'w');
```

Annotations:

- `fid`: An open file has a file ID, here stored in variable `fid`
- `fopen`: Built-in function to open a file
- `'geneData.txt'`: Name of the file (created and) opened. `txt` and `dat` are common file name extensions for plain text files
- `'w'`: `'w'` indicates that the file is to be opened for writing. Use `'a'` for appending

2. Write (print) to the file

```
fid = fopen('geneData.txt', 'w');

for i=1:length(Z)
    fprintf(fid, '%s\n', Z{i});
end
```

Annotations:

- `fid`: Printing is to be done to the file with ID `fid`
- `%s\n`: Substitution sequence specifies the *string* format (followed by a new-line character)
- `Z{i}`: The i^{th} item in cell array `Z`

3. Close the file

```
fid = fopen('geneData.txt', 'w');

for i=1:length(Z)
    fprintf(fid, '%s\n', Z{i});
end

fclose(fid);
```

```
function cellArray2file(CA, fname)
% CA is a cell array of strings.
% Create a .txt file with the name
% specified by the string fname.
% The i-th line in the file is CA{i}

fid= fopen([fname '.txt'], 'w');
for i= 1:length(CA)
    fprintf(fid, '%s\n', CA{i});
end
fclose(fid);
```

Lecture 19 43

Reverse problem: Read the data in a file line-by-line and store the results in a cell array

GATTCGAG
 GAGCCACTGGTC
 ATAGATCCT

z
 { 'GATTCGAG'
 'GAGCCACTGGTC'
 'ATAGATCCT' }

geneData.txt

How are lines separated?
 How do we know when there are no more lines?

Lecture 19 44

In a file there are hidden “markers”

GATTCGAG ●
 GAGCCACTGGTC ●
 ATAGATCCT ●

- Carriage return marks the end of a line
- eof marks the end of a file

geneData.txt

Lecture 19 45

- Read data from a file
1. **Open** a file
 2. **Read** it line-by-line until eof
 3. **Close** the file
- Lecture 19 46

1. Open the file

```
fid = fopen('geneData.txt', 'r');
```

An open file has a file ID, here stored in variable **fid**

Name of the file opened. **txt** and **dat** are common file name extensions for plain text files

'r' indicates that the file has been opened for reading

Built-in function to open a file

Lecture 19 47

2. Read each line and store it in cell array

```
fid = fopen('geneData.txt', 'r');

k= 0;
while ~feof(fid)
    k= k+1;
    Z{k}= fgetl(fid);
end
```

False until end-of-file is reached (points to ~feof(fid))

Get the next line (points to fgetl(fid))

Lecture 19 48

3. Close the file

```
fid = fopen('geneData.txt', 'r');

k= 0;
while ~feof(fid)
    k= k+1;
    Z{k}= fgetl(fid);
end
fclose(fid);
```

```
function CA = file2cellArray(fname)
% fname is a string that names a .txt file
% in the current directory.
% CA is a cell array with CA{k} being the
% k-th line in the file.

fid= fopen([fname '.txt'], 'r');
k= 0;
while ~feof(fid)
    k= k+1;
    CA{k}= fgetl(fid);
end
fclose(fid);
```

A Detailed Read-File Example

From the protein database at

<http://www.rcsb.org>

we download the file **1b18.dat** which encodes the amino acid information for the protein with the same name. We want the xyz coordinates of the protein's "backbone."

The file has a long "header"

```
HEADER MEMBRANE PROTEIN 23-JUL-98 1BL8
TITLE POTASSIUM CHANNEL (KCSA) FROM STREPTOMYCES LIVIDANS
COMPND MOL_ID: 1;
COMPND 2 MOLECULE: POTASSIUM CHANNEL PROTEIN;
COMPND 3 CHAIN: A, B, C, D;
COMPND 4 ENGINEERED: YES;
COMPND 5 MUTATION: YES
SOURCE MOL_ID: 1;
SOURCE 2 ORGANISM_SCIENTIFIC: STREPTOMYCES LIVIDANS;
```

Need to read past hundreds of lines that are not relevant to us.

Eventually, the xyz data is reached...

```
MTRIX1 2 -0.736910 -0.010340 0.675910 112.17546 1
MTRIX2 2 0.004580 -0.999940 -0.010300 53.01701 1
MTRIX3 2 0.675980 -0.004490 0.736910 -43.35083 1
MTRIX1 3 0.137220 -0.931030 0.338160 80.28391 1
MTRIX2 3 0.929330 0.002860 -0.369240 -33.25713 1
MTRIX3 3 0.342800 0.364930 0.865630 -31.77395 1

ATOM 1 N ALA A 23 65.191 22.037 48.576 1.00181.62 N
ATOM 2 CA ALA A 23 66.434 22.838 48.377 1.00181.62 C
ATOM 3 C ALA A 23 66.148 24.075 47.534 1.00181.62 C
```

↑
Signal: Lines that begin with 'ATOM'

↑ x ↑ y ↑ z

Where exactly are the xyz data?

1-4 14-15 33-38 41-46 49-54 ← Column nos. of interest

```
ATOM 14 N HIS A 25 68.656 24.973 44.142 1.00128.26 N
ATOM 15 CA HIS A 25 69.416 24.678 42.939 1.00128.26 C
ATOM 16 C HIS A 25 68.843 23.458 42.227 1.00128.26 C
ATOM 17 O HIS A 25 68.911 23.354 41.007 1.00128.26 O
ATOM 18 CB HIS A 25 70.881 24.416 43.300 1.00154.92 C
ATOM 19 CG HIS A 25 71.188 22.977 43.573 1.00154.92 C
ATOM 20 ND1 HIS A 25 71.886 22.184 42.689 1.00154.92 N
ATOM 21 CD2 HIS A 25 70.877 22.182 44.625 1.00154.92 C
ATOM 22 CE1 HIS A 25 71.993 20.963 43.183 1.00154.92 C
ATOM 23 NR2 HIS A 25 71.388 20.935 44.356 1.00154.92 N
ATOM 24 N TRP A 26 68.271 22.546 43.005 1.00 87.09 N
ATOM 25 CA TRP A 26 67.702 21.311 42.475 1.00 87.09 C
ATOM 26 C TRP A 26 66.187 21.378 42.339 1.00 87.09 C
ATOM 27 O TRP A 26 65.577 20.508 41.718 1.00 87.09 O
```

↑ x ↑ y ↑ z

Just getting what you need from a data file

- Read past all the header information
- When you come to the lines of interest, collect the xyz data
 - Line starts with 'ATOM'
 - Cols 14-15 is 'CA'

Lecture 19

55

```

fid = fopen('1b18.dat', 'r');
x=[];y=[];z=[];
while ~feof(fid)
    s = fgetl(fid);
    if strcmp(s(1:4),'ATOM')
        if strcmp(s(14:15),'CA')
            x = [x; str2double(s(33:38))];
            y = [y; str2double(s(41:46))];
            z = [z; str2double(s(49:54))];
        end
    end
end
fclose(fid);
    
```

Iterate Until End of File

Lecture 19

58

A detailed sort-a-file example

Suppose each line in the file `statePop.txt` is structured as follows:

- Cols 1-14: State name
- Cols 16-24: Population (millions)

The states appear in alphabetical order.

Lecture 19

62

A detailed sort-a-file example

Create a new file

`statePopSm2Lg.txt`

that is structured the same as `statePop.txt` except that *the states are ordered from smallest to largest according to population.*

Alabama	4557808
Alaska	663661
Arizona	5939292
Arkansas	2779154
California	36132147
Colorado	4665177
:	:
:	:

- Need the pop as numbers for sorting.
- Can't just sort the pop—have to maintain association with the state names.

Lecture 19

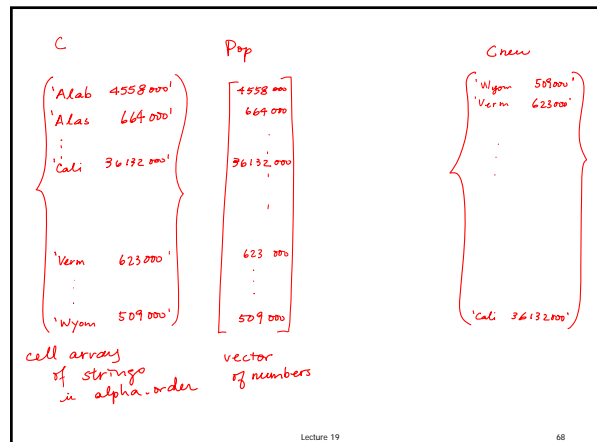
64

First, read the file and store each line in a cell of a cell array

```
C = file2cellArray('StatePop');
```

Lecture 19

65



Lecture 19

68

Next, get the populations into a **numeric vector**

```
C = file2cellArray('StatePop');
n = length(C);
pop = zeros(n,1);
for i=1:n
    s = C{i};
    pop(i) = str2double(s(16:24));
end
```

Converts a string representing a numeric value (digits, decimal point, spaces) to the numeric value → scalar of type double. E.g., `x=str2double(' 3.24 ')` assigns to variable `x` the numeric value 3.24

Lecture 19

67

Built-In function sort

Syntax: `[y,idx] = sort(x)`

x:

10	20	5	90	15
----	----	---	----	----

y:

5	10	15	20	90
---	----	----	----	----

idx:

3	1	5	2	4
---	---	---	---	---

`y(1) = x(3) = x(idx(1))`

Lecture 19

71

Built-In function sort

Syntax: `[y,idx] = sort(x)`

x:

10	20	5	90	15
----	----	---	----	----

y:

5	10	15	20	90
---	----	----	----	----

idx:

3	1	5	2	4
---	---	---	---	---

`y(k) = x(idx(k))`

Lecture 19

76

C	Pop	s	idx	Cnew																																														
<table border="1"> <tr><td>'Alab'</td><td>4558000</td></tr> <tr><td>'Alas'</td><td>664000</td></tr> <tr><td>...</td><td>...</td></tr> <tr><td>'Cali'</td><td>36132000</td></tr> <tr><td>...</td><td>...</td></tr> <tr><td>'Verm'</td><td>623000</td></tr> <tr><td>...</td><td>...</td></tr> <tr><td>'Wyom'</td><td>509000</td></tr> </table>	'Alab'	4558000	'Alas'	664000	'Cali'	36132000	'Verm'	623000	'Wyom'	509000	<table border="1"> <tr><td>4558000</td></tr> <tr><td>664000</td></tr> <tr><td>...</td></tr> <tr><td>36132000</td></tr> <tr><td>...</td></tr> <tr><td>623000</td></tr> <tr><td>...</td></tr> <tr><td>509000</td></tr> </table>	4558000	664000	...	36132000	...	623000	...	509000	<table border="1"> <tr><td>509000</td></tr> <tr><td>623000</td></tr> <tr><td>...</td></tr> <tr><td>36132000</td></tr> <tr><td>...</td></tr> <tr><td>623000</td></tr> <tr><td>...</td></tr> <tr><td>36132000</td></tr> </table>	509000	623000	...	36132000	...	623000	...	36132000	<table border="1"> <tr><td>50</td></tr> <tr><td>45</td></tr> <tr><td>...</td></tr> <tr><td>5</td></tr> <tr><td>...</td></tr> <tr><td>5</td></tr> </table>	50	45	...	5	...	5	<table border="1"> <tr><td>'Wyom'</td><td>509000</td></tr> <tr><td>'Verm'</td><td>623000</td></tr> <tr><td>...</td><td>...</td></tr> <tr><td>'Cali'</td><td>36132000</td></tr> </table>	'Wyom'	509000	'Verm'	623000	'Cali'	36132000
'Alab'	4558000																																																	
'Alas'	664000																																																	
...	...																																																	
'Cali'	36132000																																																	
...	...																																																	
'Verm'	623000																																																	
...	...																																																	
'Wyom'	509000																																																	
4558000																																																		
664000																																																		
...																																																		
36132000																																																		
...																																																		
623000																																																		
...																																																		
509000																																																		
509000																																																		
623000																																																		
...																																																		
36132000																																																		
...																																																		
623000																																																		
...																																																		
36132000																																																		
50																																																		
45																																																		
...																																																		
5																																																		
...																																																		
5																																																		
'Wyom'	509000																																																	
'Verm'	623000																																																	
...	...																																																	
'Cali'	36132000																																																	
cell array of strings in alpha-order	vector of numbers	vector of indices (rank i)																																																

Lecture 19

77

Sort from little to big

```
% C is cell array read from statePop.txt
% pop is vector of state pop (numbers)
[s,idx] = sort(pop);
Cnew = cell(n,1);
for i=1:length(C)
    ithSmallest = idx(i);
    Cnew{i} = C{ithSmallest};
end

cellArray2file(Cnew,'statePopSm2Lg')
```

Lecture 19

78

Wyoming	509294
Vermont	623050
North Dakota	636677
Alaska	663661
South Dakota	775933
Delaware	843524
Montana	935670
:	:
:	:
Illinois	12763371
Florida	17789864
New York	19254630
Texas	22859968
California	36132147

Lecture 19

79