

data pre-process:

1. text  $\leftarrow$  original text file
2. char2idx  $\leftarrow$  {char: index} dictionary
3. text\_as\_int  $\leftarrow$  [18 47 56 57 ...] text file represented by integer.
4. char\_dataset = tf.data.Dataset.from\_tensor\_slices(text\_as\_int)

char\_dataset  
↓  
tensor 格式  
text\_as\_int.

from\_tensor\_slices. example:  
dataset = tf.data.Dataset.from\_tensor\_slices([1, 2, 3])  
for element in dataset:  
 print(element)  
tf.Tensor(1, shape=(), dtype=int32)  
:  
2  
:  
3

5. sequences = char\_dataset.batch(seq\_length + 1, drop\_remainder=True)
- sequences. 相当于把  $[c_1, c_2, c_3 \dots c_n]$

↓  
 $\begin{bmatrix} [c_1, c_2, \dots, c_m] \\ [c_{m+1}, c_{m+2}, \dots, c_{2m}] \\ \vdots \\ [c_n] \end{bmatrix}$   $m = \text{batch}$   
here  $m = \text{seq\_length} + 1$   
 $\underline{m = 100}$

6. sequences is duplicated, processed and split to 2 parts,  
 $[\text{Input\_data}, \text{Target\_data}]$   
↓  
dataset.

7. dataset:  $[\text{Input\_data}, \text{Target\_data}]$ ,  
where

Input\_data = sequences,

Target\_data = sequences removed first char

Then:

8. dataset = dataset.shuffle(10000).batch(64)

shape(dataset):  $\underline{((64, 100), (64, 100)), \dots}$

Then dataset is used in training

model.fit(dataset, epochs, callbacks)

dataset 分为 training 和 validation 两部分