

# Web Data Mining

## Lecture 11: Mining Data Streams

**Jaroslav Kuchař & Milan Dojčinovski**

jaroslav.kuchar@fit.cvut.cz, milan.dojchinovski@fit.cvut.cz



Czech Technical University in Prague - Faculty of Information Technologies - Software and Web Engineering



Summer semester 2019/2020  
Humla v0.3

## Overview

- Introduction
- Stream Data Processing
- Approaches
- Technologies

## Introduction

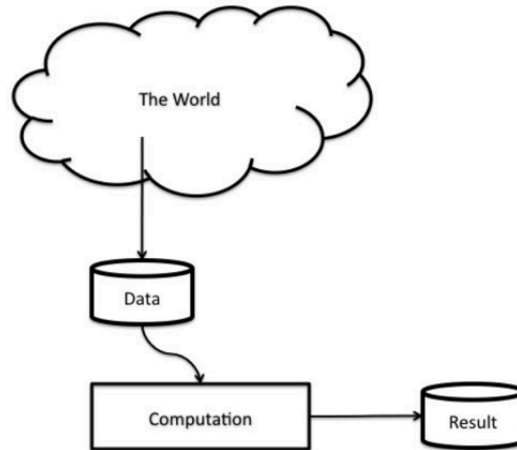
- **Motivation**
  - Traditional approaches assume to have all data available in the database
  - Data is not static
    - Data is growing all the time
    - Every second new data is generated
    - Data arrives in a stream and if it is not processed immediately or stored, then it is lost forever.
- **Application context**
  - When your data is too large to fit in the memory
  - When new data is constantly being generated, and/or is dependent upon time
- **Examples**
  - Web data
  - Sensor data
  - Telcos, activity data, Social networks, Real-time recommendations, Fraud or Spam detection, ....

## Internet Minute



## Data Mining Overview

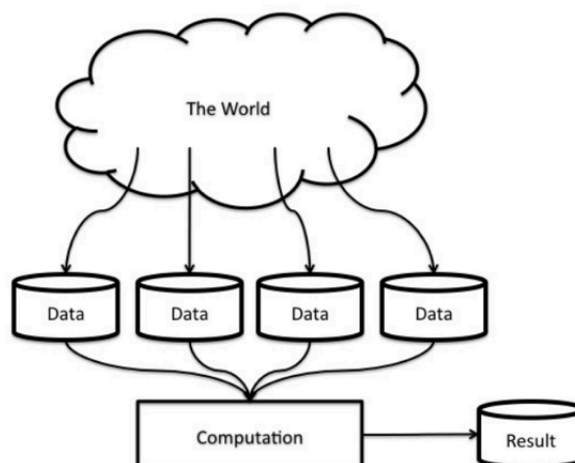
- Single Machine Data Mining



\* Edo Liberty, Jelani Nelson: Streaming Data Mining - KDD 2012 tutorial on practical algorithms in mining streaming data.

## Distributed Storage

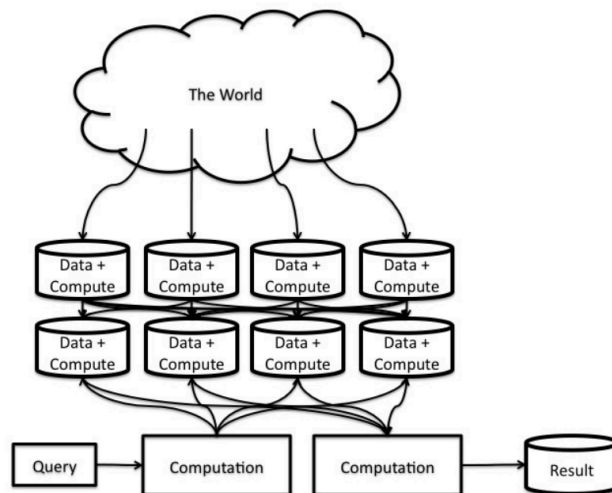
- Distributed Storage



\* Edo Liberty, Jelani Nelson: Streaming Data Mining - KDD 2012 tutorial on practical algorithms in mining streaming data.

## Distributed Computation

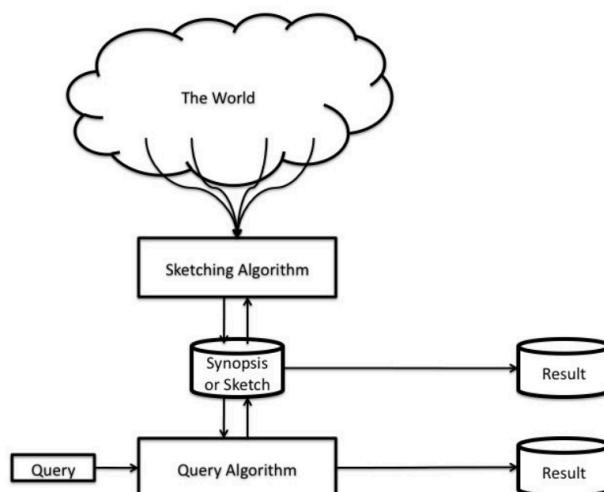
- Distributed Computation



\* Edo Liberty, Jelani Nelson: Streaming Data Mining - KDD 2012 tutorial on practical algorithms in mining streaming data.

## Streaming Model

- Streaming Model



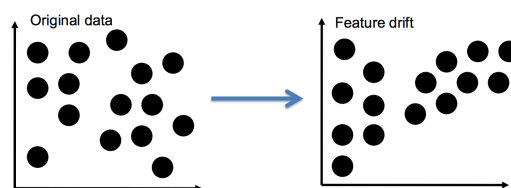
\* Edo Liberty, Jelani Nelson: Streaming Data Mining - KDD 2012 tutorial on practical algorithms in mining streaming data.

## Data Streams

- Characteristics
  - High volume (possibly infinite) of continuous data
  - Data arrive at a rapid rate
  - Data distribution changes on the fly
  - Limited access to historical data
  - Can't store them all
- Static vs Streaming data
  - "If the data distribution is stable, mining a data stream is largely the same as mining a large data set, since statistically we can draw and mine a sufficient sample"
- Challenges
  - The data is evolving
  - Finding and understanding changes
  - Maintaining an updated model

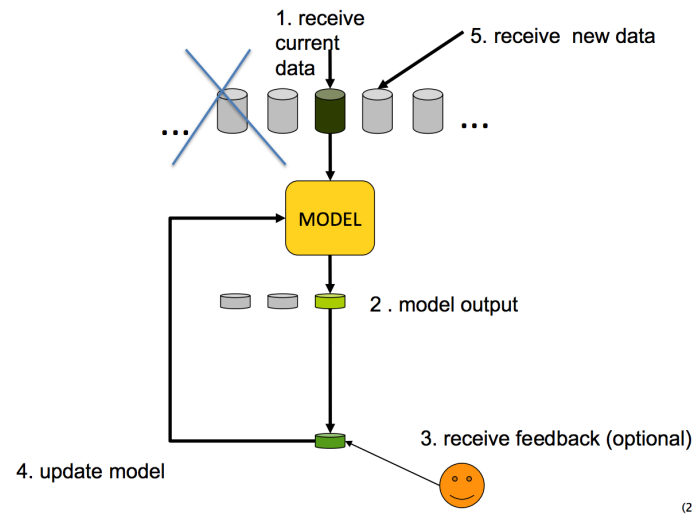
## Types of changes in Data

- Evolving/drifting data = data distribution changes over time
  - Feature drift
    - distribution of input data  $X$  changes
  - Real concept drift
    - relation between input  $X$  and target  $y$  changes
  - Arrival of new information
- Examples
  - Feature drift



# Mining Streaming Data

- Mining streaming data

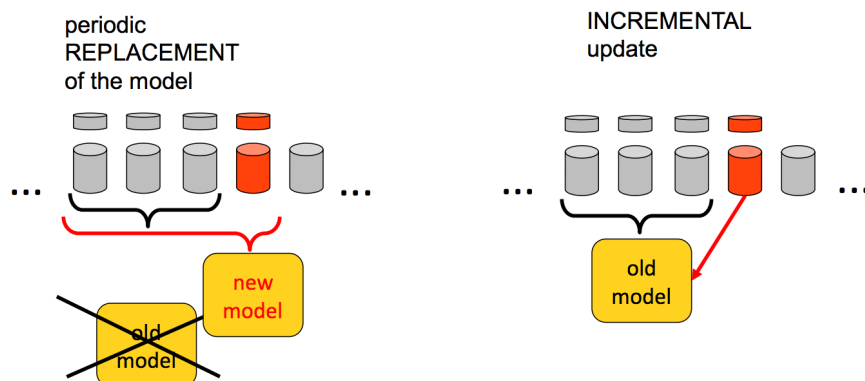


\* A. Bifet, J. Gama, R. Gavalda, G. Kreml, M. Pechenizkiy, B. Pfahringer, M. Spiliopoulou, I. Zliobaite: Advanced Topics on Data Stream Mining -

# Model Adaptation

- Model Adaptation

two main approaches



\* A. Bifet, J. Gama, R. Gavalda, G. Kreml, M. Pechenizkiy, B. Pfahringer, M. Spiliopoulou, I. Zliobaite: Advanced Topics on Data Stream Mining -

## Computing in Data Streams

- Axioms
  - *One pass*
  - *Low time per item - read, process, discard*
  - *Sublinear memory - only summaries or sketches*
  - *Anytime, real-time answers*
  - *The stream evolves over time*
- Characteristics
  - *Approximate answers are often OK*
  - *Algorithms use a source of independent random bits*
    - *So different runs give different outputs*
  - *But "most runs" are "approximately correct"*

## Overview

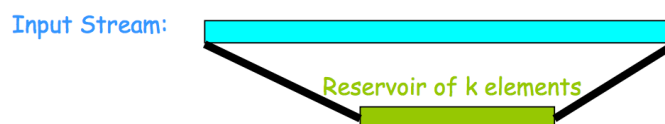
- Introduction
- **Stream Data Processing**
- Approaches
- Technologies

## Sampling

- Sampling
  - Find uniform random samples of an infinite data stream
- Input
  - Stream of data that arrive online
  - Sample size  $k$
  - Sample range
    - entire stream
    - most recent window (count-based or time-based)
- Output
  - $k$  elements chosen uniformly at random within the sample range

## Reservoir Sampling

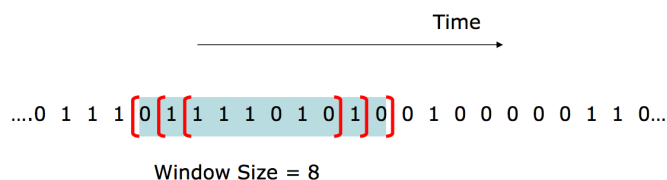
- Classical algorithm
  - Goal: maintains a fixed-size uniform random sample
  - Approach
    - Put the first  $k$  elements from the stream into the repository
    - When the  $i$ -th element arrives
      - Add it to reservoir  $S$  with probability  $p$
      - If added, randomly remove an element from  $S$
- Duplicates
  - Stream contains duplicate elements
  - Any value occurring frequently in the sample is a wasteful use of the available space





## Moving Window

- Moving Window
  - *Timeliness*
    - *old data are not useful*
  - *Restrict samples to a window of recent data*
    - *As new data arrives, old data "expires"*
  - *Reservoir sampling cannot handle data expiration*
- Naive algorithm
  - *Place a moving window of size  $N$  on the stream*
  - *an old element  $y$  expires when a new element  $x$  arrives*



## Other Processing

- Sketches
  - *Sampling techniques and sliding window models focus on a small part of the data*
  - *Sketches try to summarize the entire data, often at multiple levels of detail*
    - *build a small-space summary for a distribution vector (e.g., histogram)*
- Randomized Algorithms
  - *mainly consider Monte Carlo algorithms*
    - *has bounds on the running time but may not return the correct result*
  - *may err with some small, but controllable, probability.*

## Overview

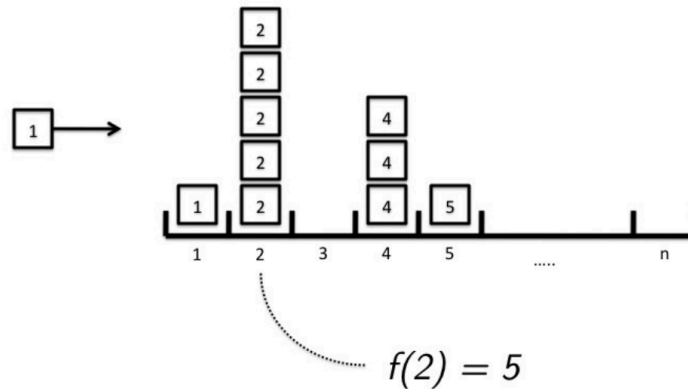
- Introduction
- Stream Data Processing
- **Approaches**
- Technologies

## Approaches

- Tasks
  - *Trivial*
    - *count items, sum values, find min/max, sample*
  - *Nontrivial*
    - *alert on new item, most frequent item, finding median*
- Examples of
  - *Item frequencies*
  - *Count distinct items*
  - ...

## Item frequencies - Naive approach

- Item frequencies
  - $n$  counters
  - Computing  $f(i)$  for all  $i$  is easy in  $O(n)$  space.



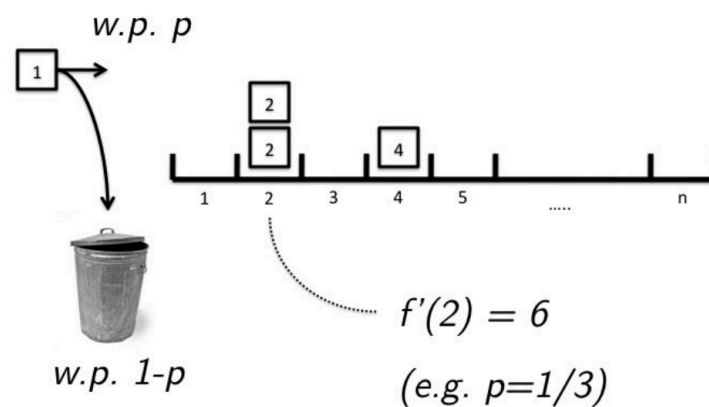
\* Edo Liberty, Jelani Nelson: Streaming Data Mining - KDD 2012 tutorial on practical algorithms in mining streaming data.

Lecture 11: Mining Data Streams - Jaroslav Kuchař & Milan Dojčinovski

– 21 –

## Item frequencies - sampling

- Item frequencies - sampling
  - We sample with probability  $p$  and estimate  $f(i) = 1/p A(i)$
  - $O(pN)$  space



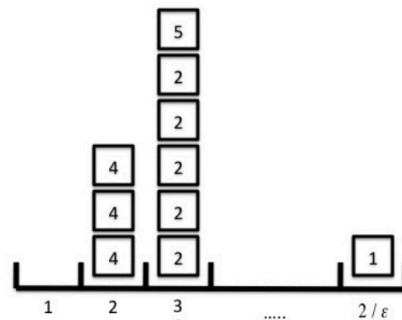
\* Edo Liberty, Jelani Nelson: Streaming Data Mining - KDD 2012 tutorial on practical algorithms in mining streaming data.

Lecture 11: Mining Data Streams - Jaroslav Kuchař & Milan Dojčinovski

– 22 –

## Item frequencies - Count-Min Sketch

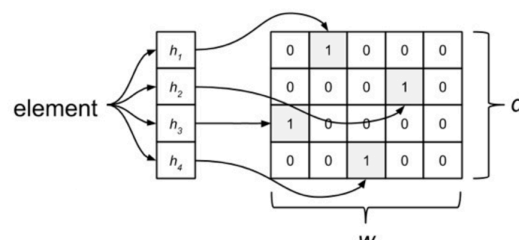
- Count-min sketch
  - limited amount of counters
  - Items are counted in buckets according to a hash function



$$f'(2) = A(h(2)) = A(3) = 6$$

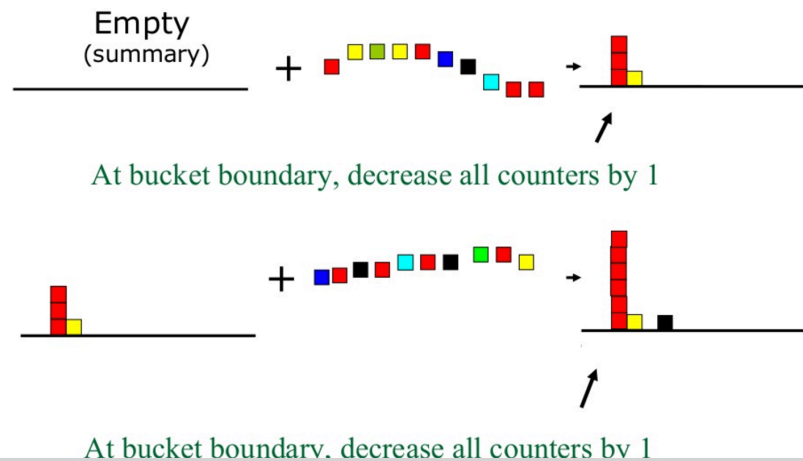
## Item frequencies - Count-Min Sketch

- Count-min sketch
  - approximate frequencies in sublinear space
  - matrix with  $w$  columns and  $d$  rows - initialized to zeros
  - each row has a hash function
  - When elements arrive
    - hash for each row
    - increment each counter by 1
  - $\text{freq}(\text{element}) = \min \text{ counter value}$
- Minimum value
  - possibility for collision for elements
  - may overestimate



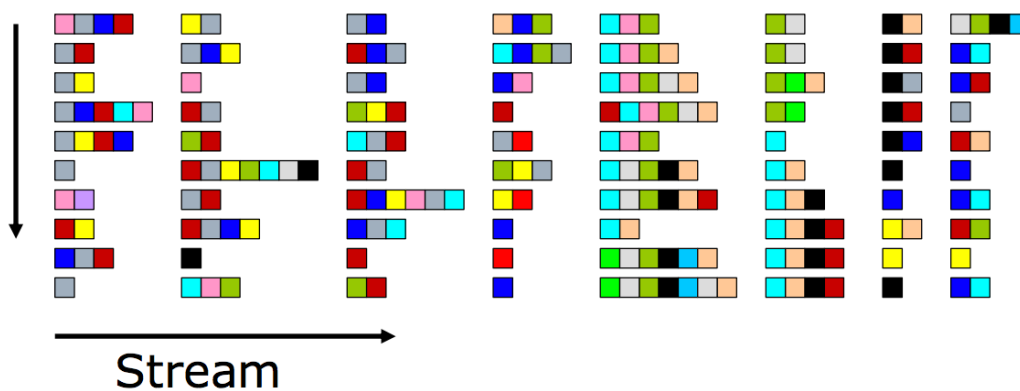
## Identifying frequent items - Lossy Counting

- Lossy Counting
  - Application in the Frequent Pattern Mining
    - Association Rules
  - Approach
    - Incoming stream is divided into buckets of size  $w$
    - Compute frequencies according to the bucket
    - At bucket boundary, decrease all counters by 1



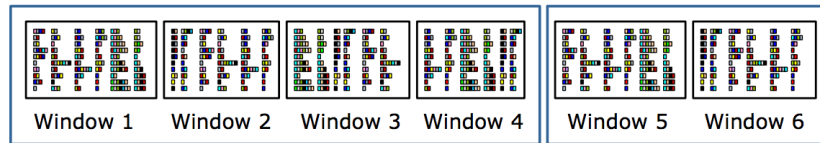
## Finding Frequent Itemsets

- Approach
  - load as many buckets as possible to the main memory
  - delete entry if the updated frequency is less than bucket number
  - decrease the frequency by number of buckets

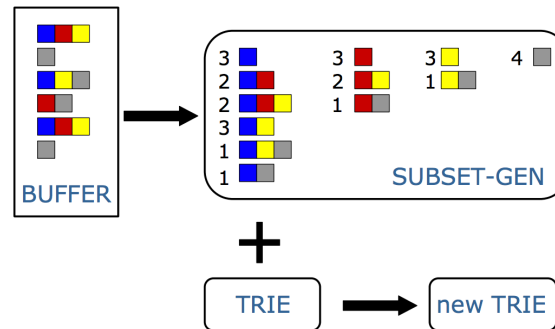


## Finding Frequent Itemsets (cont.)

- Load as many buckets as possible to the main memory



- Algorithm



\* Frequency Counts over Data Streams: <https://www.cse.ust.hk/vldb2002/VLDB2002-proceedings/slides/S10P03slides.pdf>.

## Count Distinct Items

- Flajolet–Martin algorithm
  - probabilistic computing
  - refined in *LogLog* and *HyperLogLog*
    - Billions of distinct values in 1.5KB of RAM with 2% relative error
- The main idea - Cardinality estimation
  - Algorithm
    - $n = 0$
    - For each input item
      - hash item into  $i$  bit string
      - count trailing zeros in the bit string
      - if the count  $> n$ 
        - $n = \text{count}$
    - Cardinality =  $2^n$
  - Assumptions
    - a good hash function
    - in random data, a sequence of  $n$  zero bits will occur once in every  $2^n$  elements, on average

## Count Distinct Items (cont.)

- Flajolet–Martin example

- ... 0 1 0 1 0 1 0 1 0 0 1 0 0 0
- $P(n \text{ trailing zeros}) = (1/2)^n$
- Number of seen hashes =  $2^n$ 
  - ... {000, 001, 010, 011, 100, 101, 110, 111}

id	H(id)	trailing zeros	max zeros
A	0111001001	0	0
B	1000110100	2	2
C	1101111000	3	3
D	1001010110	1	3
E	1101011100	2	3

- Assumptions

- if we had seen 2 distinct items, we would expect 1 trailing zero
- if we had seen 4 distinct items, we would expect 2 trailing zero
- reasonable estimation is  $2^n$

- Improvements

- more independent hashes or prefix and subhashes
- compute mean values

## Count Distinct Items (cont.)

```
1 # https://jeremykun.com/2016/01/04/hashing-to-estimate-the-size-of-a-stream/
2 import random
3 def randomHash(modulus):
4     a, b = random.randint(0, modulus-1), random.randint(0, modulus-1)
5     def f(x):
6         return (a*x + b) % modulus
7     return f
8
9 def average(L):
10     return sum(L) / len(L)
11
12 def numDistinctElements(stream, numParallelHashes=10):
13     modulus = 2**20
14     hashes = [randomHash(modulus) for _ in range(numParallelHashes)]
15     minima = [modulus] * numParallelHashes
16     currentEstimate = 0
17     for i in stream:
18         hashValues = [h(i) for h in hashes]
19         for i, newValue in enumerate(hashValues):
20             if newValue < minima[i]:
21                 minima[i] = newValue
22
23     currentEstimate = modulus / average(minima)
24     return currentEstimate
25
26 S = [random.randint(1, 2**20) for i in range(10000)]
27 for k in (10, 50, 100):
28     print(numDistinctElements(S, k))
29
30 # 9110.130321459601
31 # 9463.682310469314
32 # 10447.105708877154
```

# Computing Moving Average

- Sliding Window Approaches

```
1 import random
2 import numpy as np
3
4 full = []
5 queue1 = []
6 size = 3
7
8 def add1(value):
9     if len(queue1) < size:
10         queue1.append(value)
11     else:
12         del queue1[0]
13         queue1.append(value)
14     return sum(queue1)/len(queue1)
15
16 for i in range(10):
17     a = random.randint(0,100)
18     full.append(a)
19     print("add({}) - moving={} - total={}".format(a, add1(a), np.mean(full)))
```

```
1 add(83) - moving=83.0 - total=83.0
2 add(53) - moving=68.0 - total=68.0
3 add(28) - moving=54.666666666666664 - total=54.666666666666664
4 add(18) - moving=33.0 - total=45.5
5 add(42) - moving=29.333333333333332 - total=44.8
6 add(60) - moving=40.0 - total=47.333333333333336
7 add(39) - moving=47.0 - total=46.142857142857146
8 add(37) - moving=45.333333333333336 - total=45.0
9 add(11) - moving=29.0 - total=41.222222222222222
10 add(12) - moving=20.0 - total=38.3
```

## Overview

- Introduction
- Stream Data Processing
- Approaches
- Technologies



## Technologies

- Stream Processing Systems
  - can handle a nearly unlimited amount of data, but they only process one (true stream processing) or very few (micro-batch processing) items at a time
- Many existing technologies
  - Apache Storm
    - stream processing framework, extremely low latency, near real-time processing
  - Apache Samza
    - tightly tied to the Apache Kafka messaging system
  - Apache Spark
    - batch processing framework with stream processing capabilities
  - Apache Flink
    - stream processing framework that can also handle batch tasks
  - ...

## Apache Spark Example

- Docker Jupyter Notebook Python, Spark, Mesos Stack

```
1 | docker pull jupyter/pyspark-notebook
2 | docker run -it --rm --add-host=parent-host:10.0.0.10 -p 8880:8880 -v $PWD:/home/jovyan/work ju
```

- Processing stream

```
1 | from pyspark import SparkContext
2 | from pyspark.streaming import StreamingContext
3 |
4 | sc = SparkContext("local[2]", "NetworkWordCount")
5 | ssc = StreamingContext(sc, 1)
6 | lines = ssc.socketTextStream("parent-host", 9999)
7 | words = lines.flatMap(lambda line: line.split(" "))
8 | pairs = words.map(lambda word: (word, 1))
9 | wordCounts = pairs.reduceByKey(lambda x, y: x + y)
10 | wordCounts.pprint()
11 | ssc.start() # Start the computation
12 | ssc.awaitTermination() # Wait for the computation to terminate
```

- Generate stream

```
1 | nc -lk 9999
```

- Outputs

```
1 | (hello,1)
2 | (world,1)
```