

# INTRODUCTION TO DATA SCIENCE

JMCT

SLIDES BY JPD

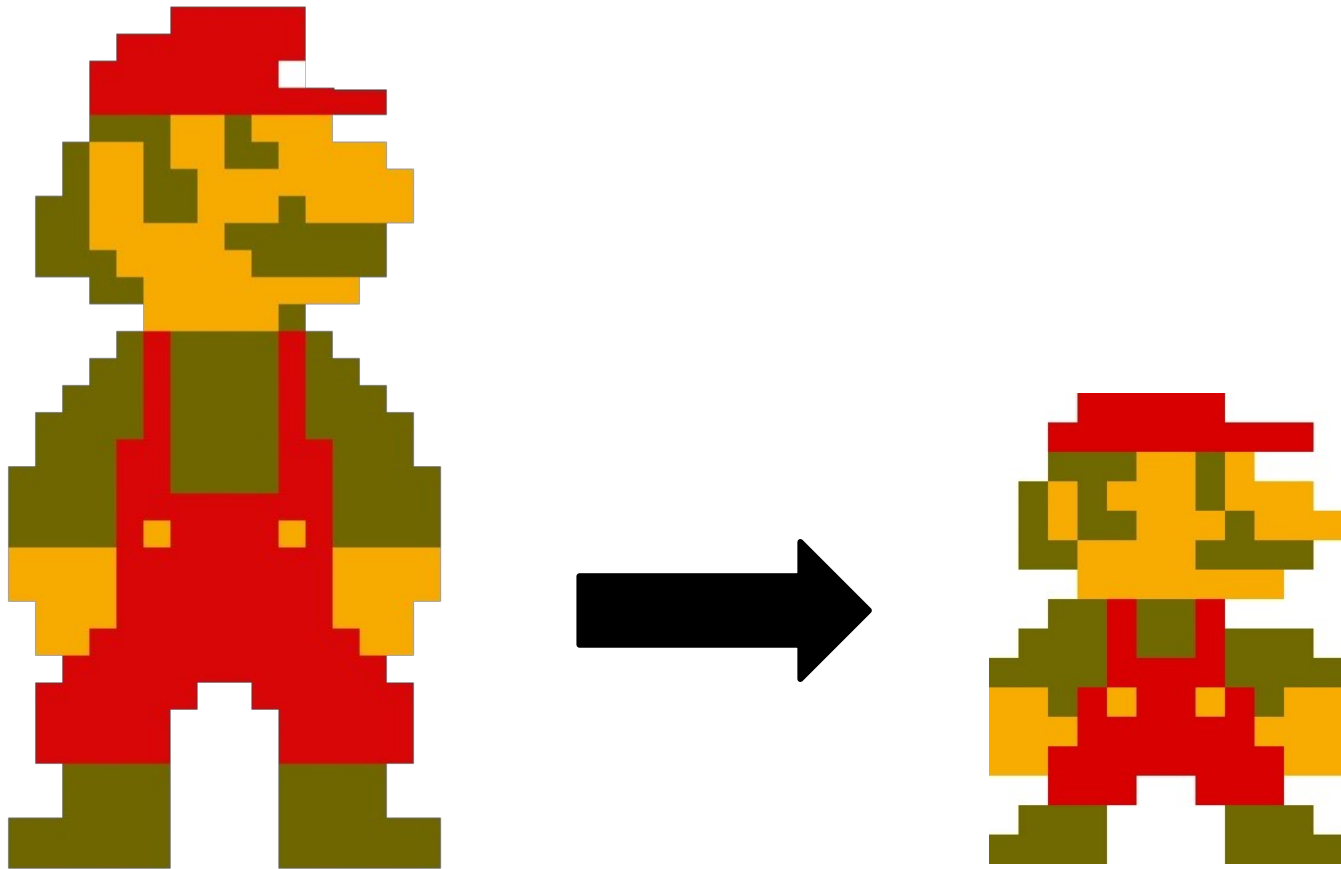
Lecture #24 – 04/28/2021

CMSC320



**COMPUTER SCIENCE**  
UNIVERSITY OF MARYLAND

# DIMENSIONALITY REDUCTION



Thanks to: Zico Kolter

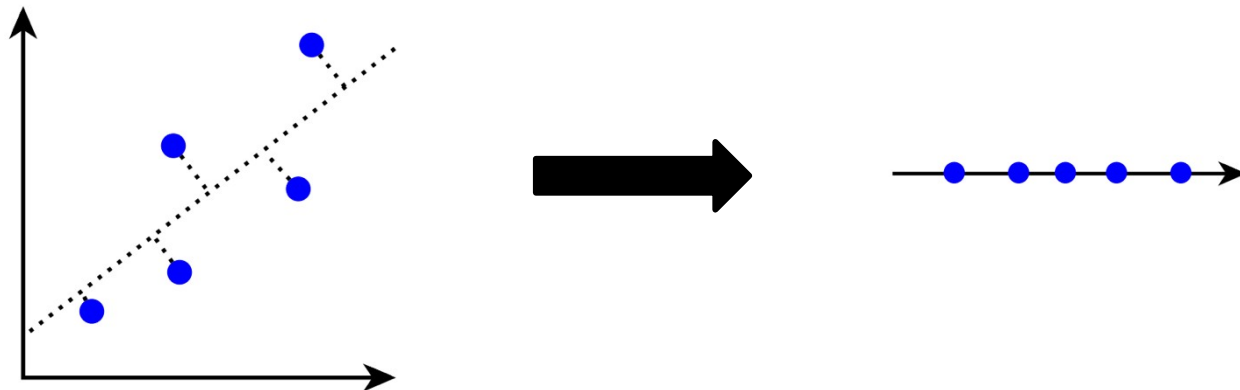
# PRINCIPAL COMPONENT ANALYSIS (PCA)

So you've measured lots of features ...

- Overfitting, interpretability issues, visualization, computation

Can we combine raw features into new features that yield a simpler description of the same system?

Principal component analysis (PCA) does this by preserving the axis of major variation in the data:



# PRINCIPAL COMPONENT ANALYSIS (PCA)

**Assume: data is normalized** ????????????

- Zero mean, unit (= 1) variance

**Hypothesis function:**

$$h_{\theta}(x) = UWx, \theta = \{U \in \mathbb{R}^{n \times k}, W \in \mathbb{R}^{k \times n}\}$$

- First multiply input by low rank matrix  $W$  (“compress” it), then map back into the initial space using  $U$

**Loss function: squared distance (like k-means)**

$$\ell(h_{\theta}(x), x) = \|h_{\theta}(x) - x\|_2^2$$

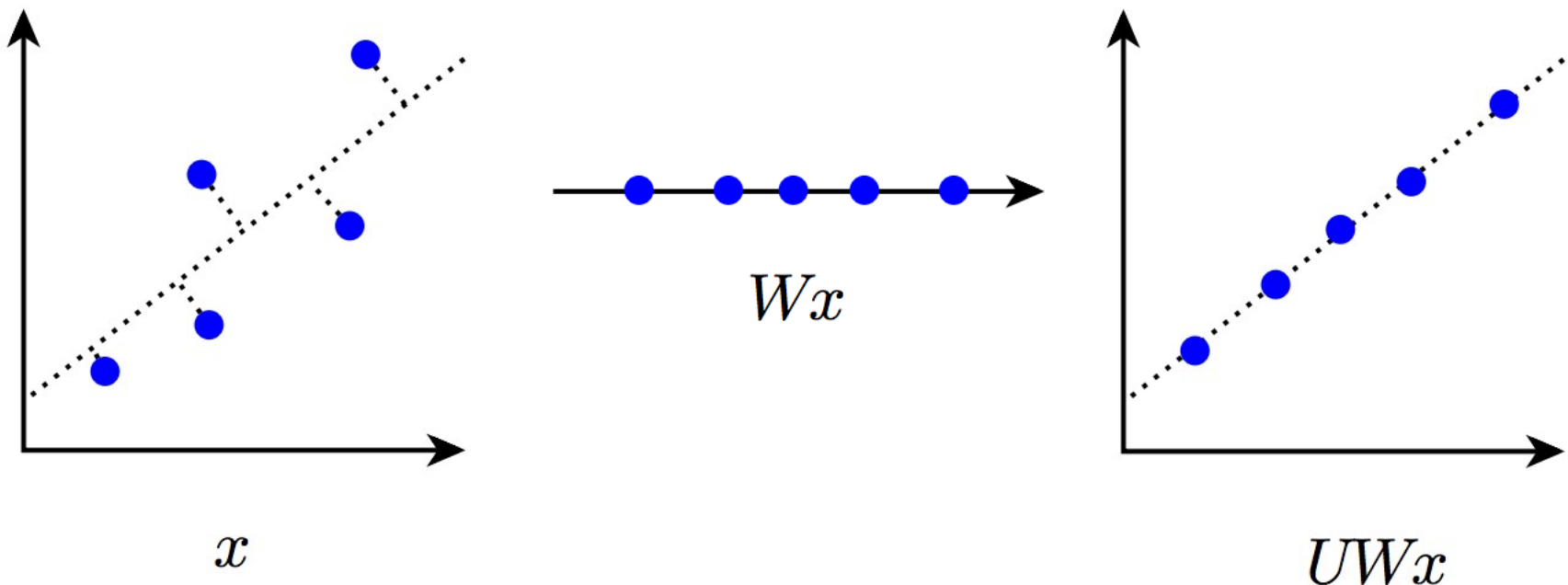
**Optimization problem:**

$$\underset{U, W}{\text{minimize}} \sum_{i=1}^m \|UWx^{(i)} - x^{(i)}\|_2^2$$

# PRINCIPAL COMPONENT ANALYSIS (PCA)

**Dimensionality reduction:** main use of PCA for data science applications

If  $h_{\theta}(x) = UWx$ , then  $Wx \in \mathbb{R}^k$  is a reduced (probably with some loss) representation of input features  $x$



# PRINCIPAL COMPONENT ANALYSIS (PCA)

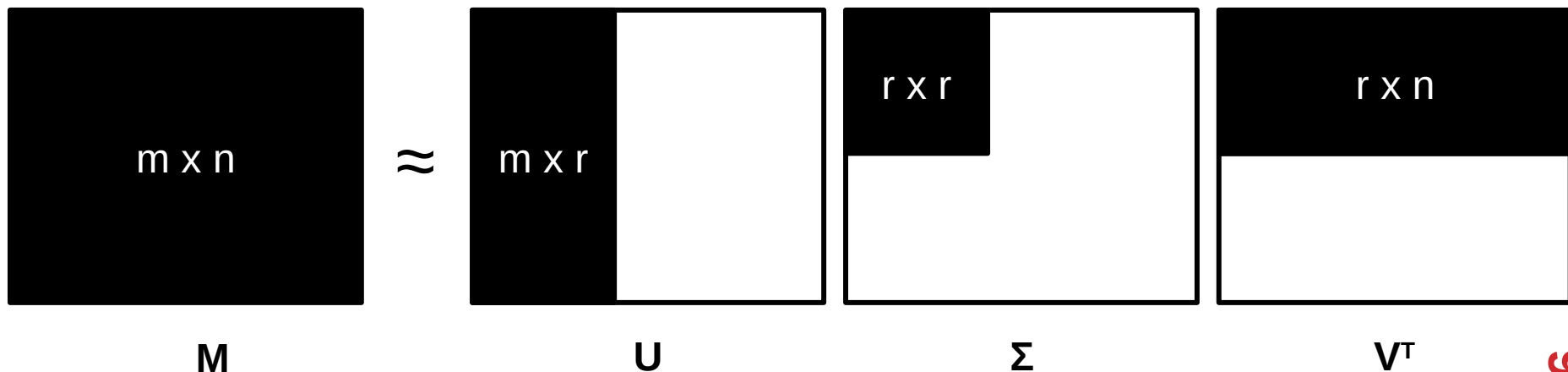
CMSC422  
MATH240

$$\underset{U, W}{\text{minimize}} \sum_{i=1}^m \|UW x^{(i)} - x^{(i)}\|_2^2$$

PCA optimization problem is non-convex ??????????????

We can solve the problem exactly using the singular value decomposition (SVD, from linear algebra):

- Factorize matrix  $M = U \Sigma V^T$  (also used to approximate)



# PRINCIPAL COMPONENT ANALYSIS (PCA)

Solving PCA exactly using the SVD:

1. Normalize input data, pick #components  $k$

2. Compute (exact) SVD of  $X = U \Sigma V^T$

3. Return:  $h_{\theta}(x) = UWx$

- $U = V_{:,1:k} \Sigma^{-1}_{1:k,1:k}$

- $W = V^T_{:,1:k}$

Loss is  $\sum_{i=k+1}^n \Sigma_{ii}^2$



# PCA IN PYTHON

Can roll your own PCA easily (assuming a call to SVD via SciPy or similar) ...

... or just use Scikit-Learn:

```
from sklearn.decomposition import PCA

X=np.array([[ -1, -1], [-2, -1], [-3, -2], [1, 1], [2, 1], [3, 2]])

# Fit PCA with 2 components (i.e., two final features)
pca = PCA(n_components=2)
pca.fit(X)
print(pca.explained_variance_ratio_)
```

```
[ 0.99244... 0.00755...]
```

Looks like our data basically sit on a line



# HOW TO USE PCA & FRIENDS IN PRACTICE

**Unsupervised learning methods are useful for EDA**

- Cluster or reduce to a few dimensions and visualize!

**Also useful as data prep before supervised learning!**

1. Run PCA, get  $W$  matrix
2. Transform  $\tilde{x}^{(i)} = Wx^{(i)}$  – (reduce colinearity, dimension)
3. Train and test your favorite supervised classifier

**Or use k-means to set up radial basis functions (RBFs):**

4. Get  $k$  centers  $\mu^{(1)}, \dots, \mu^{(k)}$
5. Create RBF features  $\phi_j^{(i)} = \exp\left(-\frac{\|x^{(i)} - \mu^{(j)}\|_2^2}{2\sigma^2}\right)$



# RECOMMENDER SYSTEMS & COLLABORATIVE FILTERING

# NETFLIX PRIZE

**Recommender systems:** predict a user's rating of an item

	Twilight	Wall-E	Twilight II	TFotF
User 1	+1	-1	+1	?
User 2	+1	-1	?	?
	-1	+1	-1	+1

**Netflix Prize: \$1MM to the first team that beats our in-house engine by 10%**

- Happened after about three years
- Model was **never used** by Netflix for a variety of reasons
  - Out of date (DVDs vs streaming)
  - Too complicated / not interpretable

# RECOMMENDER SYSTEMS

## Recommender systems feel like:

- Supervised learning (we know the user watched some movies, so these are like labels)
- Unsupervised learning (we want to find latent structure, e.g., genres of movies)

**They fall somewhere in between, in “Information Filtering” or Information Retrieval” ...**

- ... but we can still just phrase the problem in terms of hypothesis classes, loss functions, and optimization problems

# PREDICTION

## **Pure user information:**

- Age
- Location
- Profession/Salary

## **Pure item information:**

- Movie budget
- Main actors
- Is it a Netflix release?

## **User-item information:**

- Which items are most similar to those I've watched before?
- Which users are most similar to me, and what did they watch?

# COLLABORATIVE FILTERING

**Collaborative filtering (CF):** recommender systems that predict based only on the expressed preferences of other users for an item

$X =$

	$i_1$	$i_2$	$i_3$	$i_4$
$u_1$	1			3
$u_2$		2	5	
$u_3$		3		5
	4		4	

Rows are users

Cols are items

# MATRIX VIEW

Goal: “fill in” the matrix

	$i_1$	$i_2$	$i_3$	$i_4$
$u_1$	1	?	?	3
$u_2$	?	2	5	?
$u_3$	?	3	?	5
	4	?	4	?

The matrix is **sparse**, but the empty cells are not (necessarily) zero!

# APPROACHES TO CF

## User-user:

- Find users who look like me – based on items that we've both rated
- Predict scores for my unrated items as average of those users

## Item-item:

- Find similar items (based on scores from all users who have rated), predict scores for other users based off this

## Matrix factorization:

- Find a low-rank decomposition of  $X$  that agrees (exactly, approximately) at the observed values



# APPROACH #1: ITEM-BASED CF EX: INFER (USER 1, ITEM 3)

	Item 1	Item 2	Item 3	Item 4	Item 5
User 1	8	1	?	2	7
User 2	2	?	5	7	5
User 3	5	4	7	4	7
User 4	7	1	7	3	8
User 5	1	7	4	6	?
User 6	8	3	8	3	7

# HOW TO CALCULATE SIMILARITY (ITEMS 3 AND 5)?

	Item 1	Item 2	Item 3	Item 4	Item 5
User 1	8	1	?	2	7
User 2	2	?	5	7	5
User 3	5	4	7	4	7
User 4	7	1	7	3	8
User 5	1	7	4	6	?
User 6	8	3	8	3	7

# SIMILARITY BETWEEN ITEMS

Item 3	Item 4	Item 5
?	2	7
5	7	5
7	4	7
7	3	8
4	6	?
8	3	7

How should we calculate the similarity between two items (e.g., items 3 and 5)?

We've done this before in a different context!

# SIMILARITY BETWEEN ITEMS

Item 3	Item 5
?	7
5	5
7	7
7	8
4	?
8	7

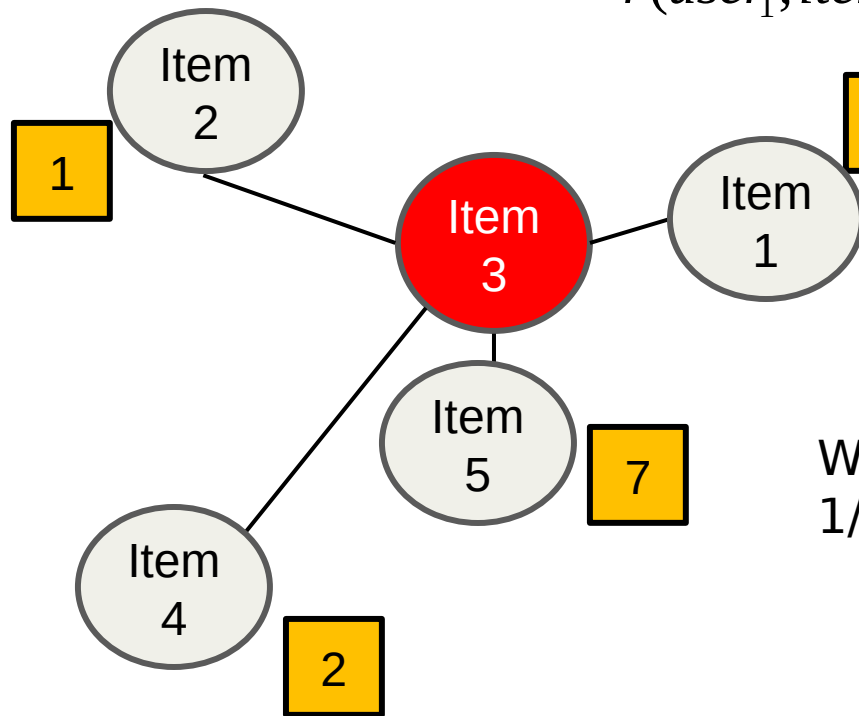
Only consider users (i.e., rows) who have rated both items (i.e., non-empty)

One approach: For each user:  
 Calculate difference in ratings for the two items

Take the average of this difference  
 $sim(item 3, item 5) = cosine( (5, 7, 7, 8), (5, 7, 8, 7) )$   
 over the users

Another approach: cosine similarity!  
 $= (5*5 + 7*7 + 7*8 + 8*7) / (\sqrt{5^2 + 7^2 + 7^2 + 8^2} * \sqrt{5^2 + 7^2 + 8^2 + 7^2})$

# PREDICTION: CALCULATING RANKING $R(\text{USER1}, \text{ITEM3})$



$$r(\text{user}_1, \text{item}_3) = \alpha * \{ r(\text{user}_1, \text{item}_1) \text{sim}(\text{item}_1, \text{item}_3) \\ + r(\text{user}_1, \text{item}_2) \text{sim}(\text{item}_2, \text{item}_3) \\ + r(\text{user}_1, \text{item}_4) \text{sim}(\text{item}_4, \text{item}_3) \\ + r(\text{user}_1, \text{item}_5) \text{sim}(\text{item}_5, \text{item}_3) \}$$

Where  $\alpha$  is a normalization factor, which is  $1/[\text{the sum of all } \text{sim}(\text{item}_i, \text{item}_3)]$ .

# CF IN PYTHON

- SKLearn *does not* have CF 'built in' :(
- Some alternatives
  - “Surprise” package: <http://surpriselib.com/>
  - ‘fast.ai’ library: <https://docs.fast.ai/collab.html>



(SOME MORE)

RECOMMENDER SYSTEMS (ISH)

# ASSOCIATION RULES

Last time: CF systems give predictions based on other users' scores of the same item

Complementary idea: Find rules that **associate** the presence of one set of items with that of another set of items

Customers who bought this item also bought



ThinkGeek Plush Unicorn Slippers, One Size, White  
★★★★★ 395  
\$7.77



Adult New Purple Unicorn Onesie Pajamas Kigurumi Cosplay Costumes Animal Outfit  
★★★★★ 168  
\$23.99 - \$28.99



EOS ~ Holiday 2015 Limited Edition Decorative Lip Balm Collection  
★★★★★ 156  
\$5.24 - \$22.99



# FORMAT OF ASSOCIATION RULES

## Typical Rule form:

- **Body**  $\Rightarrow$  **Head**
- **Body and Head can be represented as sets of items (in transaction data) or as conjunction of predicates (in relational data)**
- **Support and Confidence**
  - Usually reported along with the rules
  - Metrics that indicate the strength of the item associations

## Examples:

- $\{\text{diaper, milk}\} \Rightarrow \{\text{beer}\}$  [support: 0.5%, confidence: 78%]
- $\text{buys}(x, \text{"bread"}) \wedge \text{buys}(x, \text{"eggs"}) \Rightarrow \text{buys}(x, \text{"milk"})$  [sup: 0.6%, conf: 65%]
- $\text{major}(x, \text{"CS"}) \wedge \text{takes}(x, \text{"DB"}) \Rightarrow \text{grade}(x, \text{"A"})$  [1%, 75%]
- $\text{age}(x, 30-45) \wedge \text{income}(x, 50K-75K) \Rightarrow \text{owns}(x, \text{SUV})$
- $\text{age}=\text{"30-45"}, \text{income}=\text{"50K-75K"} \Rightarrow \text{car}=\text{"SUV"}$

# ASSOCIATION RULES: BASIC CONCEPTS

Let  $D$  be database of transactions

Transaction ID	Items
1000	A, B, C
2000	A, B
3000	A, D
4000	B, E, F

Let  $I$  be the set of items that appear in the database:

- e.g.,  $I = \{A, B, C, D, E, F\}$

Each transaction  $t$  is a subset of  $I$

A rule is an implication among itemsets  $X$  and  $Y$ , of the form by  $X \Rightarrow Y$ , where  $X \subset I$ ,  $Y \subset I$ , and  $X \cap Y = \emptyset$

- e.g.:  $\{B, C\} \Rightarrow \{A\}$

# ASSOCIATION RULES: BASIC CONCEPTS

## Itemset

- A set of one or more items
  - E.g.: {Milk, Bread, Diaper}
- k-itemset
  - An itemset that contains k items

## Support count ( $\sigma$ )

- Frequency of occurrence of an itemset (number of transactions in which it appears)
- E.g.  $\sigma(\{\text{Milk, Bread, Diaper}\}) = 2$

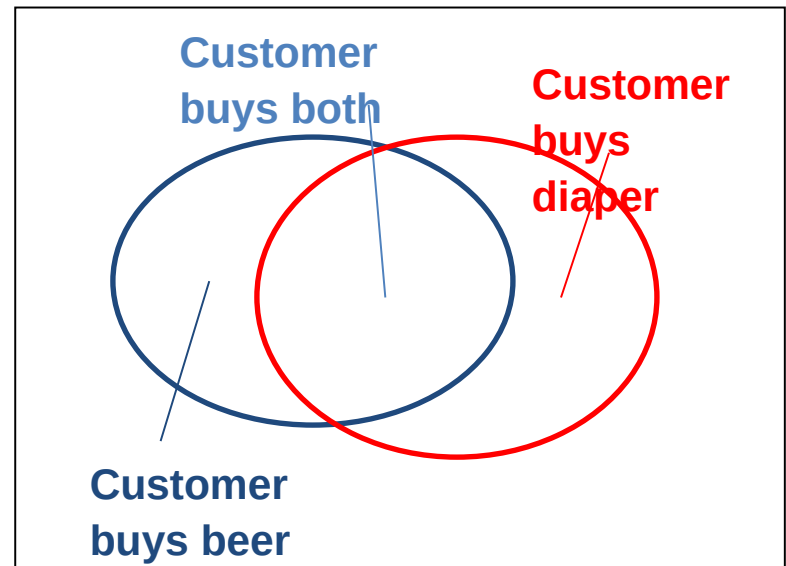
## Support

- Fraction of the transactions in which an itemset appears
- E.g.  $s(\{\text{Milk, Bread, Diaper}\}) = 2/5$

## Frequent Itemset

- An itemset whose support is greater than or equal to a *minsup* threshold

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke



# ASSOCIATION RULES: BASIC CONCEPTS

## Association Rule

- $X \Rightarrow Y$ , where X and Y are non-overlapping itemsets
- $\{\text{Milk, Diaper}\} \Rightarrow \{\text{Beer}\}$

## Rule Evaluation Metrics

- **Support (s)**
  - Fraction of transactions that contain both X and Y
  - i.e., support of the itemset  $X \cup Y$
- **Confidence (c)**
  - Measures how often items in Y appear in transactions that contain X

TID	Items
1	Bread, Milk
2	Bread, Diaper, Beer, Eggs
3	Milk, Diaper, Beer, Coke
4	Bread, Milk, Diaper, Beer
5	Bread, Milk, Diaper, Coke

## Example:

$\{\text{Milk, Diaper}\} \rightarrow \text{Beer}$

$$s = \frac{\sigma(\text{Milk, Diaper, Beer})}{|D|} = \frac{2}{5} = 0.4$$

$$c = \frac{\sigma(\text{Milk, Diaper, Beer})}{\sigma(\text{Milk, Diaper})} = \frac{2}{3} = 0.67$$

# ASSOCIATION RULES IN PRACTICE

**Orange3 is a {GUI, Python API, ...} that:**

- Enumerates frequent itemsets
- Performs association rule mining
- (Wrapper calls to, shared functionality with, Scikit-Learn)

```
conda install -c ales-erjavec orange3
```

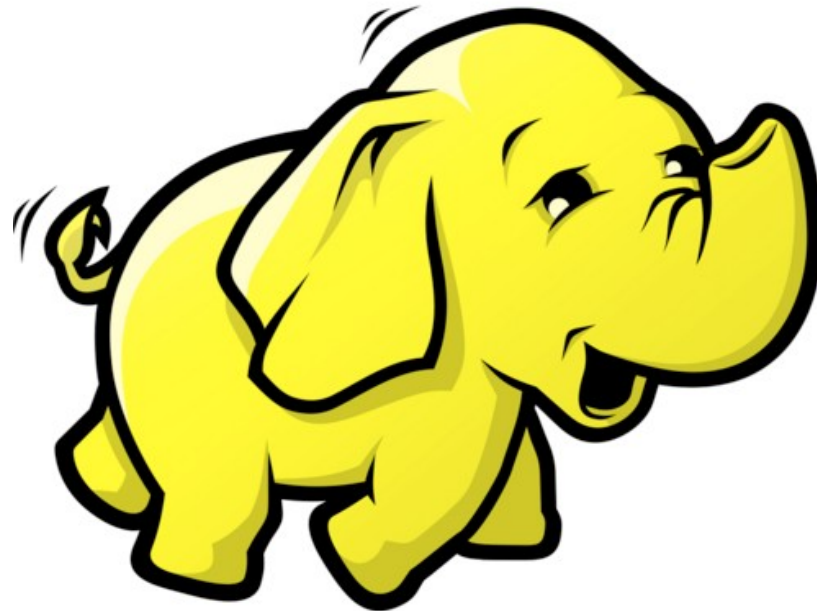
**More information:**

<https://blog.biolab.si/2016/04/25/association-rules-in-orange/>

**In general:**

- **Can be useful for interpretable, fast data mining**
- **Typically doesn't consider order, scalability issues ...**

# SCALING IT UP: BIG DATA & MAPREDUCE



*Thanks to: Jeff Dean, Sanjay Ghemawa, Zico Kolter*

# “Big data”



My laptop  
8GB RAM  
500GB Disk

**Big data?**  
No



Google Data Center  
??? RAM/Disk  
( $\gg$  PBs)

**Big data?**  
Yes



# Some notable inflection points

1. Your data fits in RAM on a single machine
2. Your data fits on disk on a single machine
3. Your data fits in RAM/disk on a “small” cluster of machines (you don’t need to worry about machines dying)
4. Your data fits in RAM/disk on a “large” cluster of machine (you need to worry about machines dying)

It’s probably reasonable to refer to 3+ as “big data”, but many would only consider 4



# Do you have big data?

If your data fits on a single machine (even on disk), then it's almost always better to think about how you can design an efficient single-machine solution, unless you have extremely good reasons for doing otherwise

scalable system	cores	twitter	uk-2007-05
GraphChi [10]	2	3160s	6972s
Stratosphere [6]	16	2250s	-
X-Stream [17]	16	1488s	-
Spark [8]	128	857s	1759s
Giraph [8]	128	596s	1235s
GraphLab [8]	128	249s	833s
GraphX [8]	128	419s	462s
Single thread (SSD)	1	300s	651s
Single thread (RAM)	1	275s	-

**Table 2: Reported elapsed times for 20 PageRank iterations, compared with measured times for single-threaded implementations from SSD and from RAM. GraphChi and X-Stream report times for 5 PageRank iterations, which we multiplied by four.**

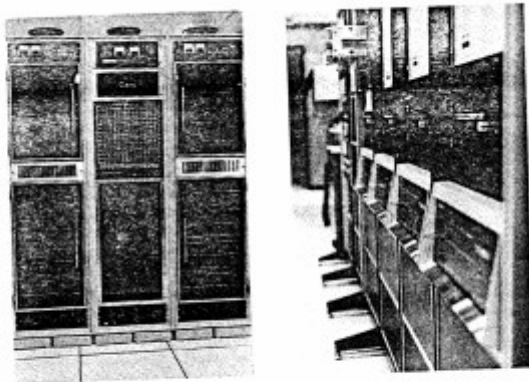
scalable system	cores	twitter	uk-2007-05
Stratosphere [6]	16	950s	-
X-Stream [17]	16	1159s	-
Spark [8]	128	1784s	$\geq 8000s$
Giraph [8]	128	200s	$\geq 8000s$
GraphLab [8]	128	242s	714s
GraphX [8]	128	251s	800s
Single thread (SSD)	1	153s	417s

**Table 3: Reported elapsed times for label propagation, compared with measured times for single-threaded label propagation from SSD.**

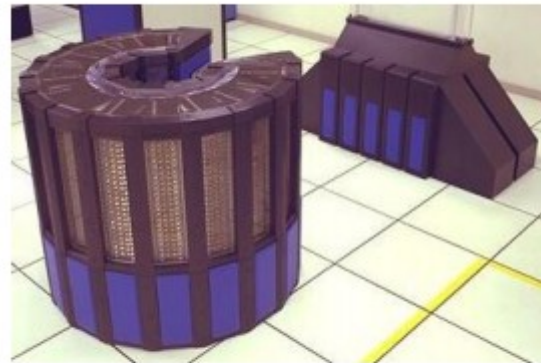
Tables from [McSherry et al., 2015 “Scalability! But at what COST”]

# Distributed computing

Distributed computing rose to prominence in the 70s/80s, often built around “supercomputing,” for scientific computing applications



1971 – CMU C.mmp  
(16 PDP-11 processors)



1984 – Cray-2  
(4 vector processors)

# Message passing interface



In mid-90s, researchers built a common interface for distributed computing called the message passing interface (MPI)

MPI provided a set of tools to run multiple processes (on a single machine or across many machines), that could communicate, send data between each other (all of “scattering”, “gathering”, “broadcasting”), and synchronize execution

Still common in scientific computing applications and HPC (high performance computing)

# Downsides to MPI

MPI is extremely powerful but has some notable limitations

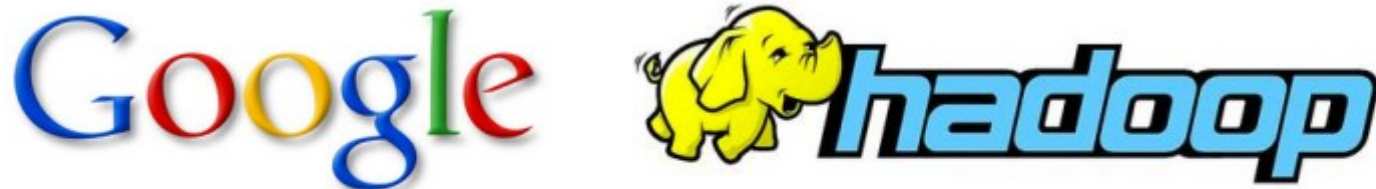
1. MPI is complicated: programs need to explicitly manage data, synchronize threads, etc
2. MPI is brittle: if machines die suddenly, can be difficult to recover (unless explicitly handled by the program, making them more complicated)

# A new paradigm for data processing

When Google was building their first data centers, they used clusters of off-the-shelf commodity hardware; machines had different speeds and failures were common given cluster sizes

Data itself was distributed (redundantly) over many machines, as much as possible wanted to do the computation on the machine where the data is stored

Led to the development of the MapReduce framework at Google [Ghemawat, 2004], later made extremely popular through the Apache Hadoop open source implementation



# AN EXAMPLE PROGRAM

**Present the concepts of MapReduce using the “typical example” of MR, Word Count**

- Input: a volume of raw text, of unspecified size (could be KB, MB, TB, **it doesn't matter!**)
- Output: a list of words, and their occurrence count.

**(Assume that words are split correctly; ignore capitalization and punctuation.)**

**Example:**

- **The doctor went to the store. =>**
  - The, 2
  - Doctor, 1
  - Went, 1
  - To, 1
  - Store, 1



# MAP? REDUCE?

**Mappers** read in data from the filesystem, and output (typically) modified data

**Reducers** collect all of the mappers output on the keys, and output (typically) reduced data

The outputted data is written to disk

All data is in terms of key-value pairs (“The” 📧 2)

# MAPREDUCE VS HADOOP

The paper is written by two researchers at Google, and describes their programming paradigm

Unless you work at Google, or use Google App Engine, you won't use it! (And even then, you might not.)

Open Source implementation is Hadoop MapReduce

- Not developed by Google
- Started by Yahoo!; now part of Apache

Google's implementation (at least the one described) is written in C++

Hadoop is written in Java



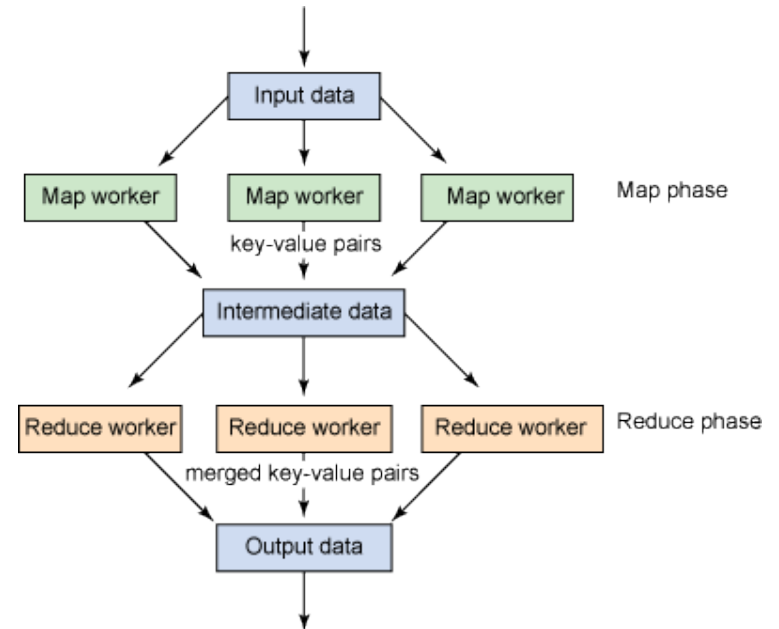
# MAJOR COMPONENTS

## User Components:

- Mapper
- Reducer
- Combiner (Optional)
- Partitioner (Optional) (Shuffle)
- Writable(s) (Optional)

## System Components:

- Master
- Input Splitter\*
- Output Committer\*
- \* You can use your own if you really want!



# KEY NOTES

**Mappers and Reducers are typically single threaded and deterministic**

- Determinism allows for **restarting of failed jobs**, or speculative execution

**Need to handle more data? Just add more Mappers/Reducers!**

- No need to handle multithreaded code
- Since they're all independent of each other, you can run (almost) arbitrary number of nodes

**Mappers/Reducers run on **arbitrary** machines. A machine typically multiple map and reduce slots available to it, typically one per processor core**

**Mappers/Reducers run entirely independent of each other**

- In Hadoop, they run in separate JVMs

# BASIC CONCEPTS

**All data is represented in key-value pairs of an arbitrary type**

**Data is read in from a file or list of files, from distributed FS**

**Data is chunked based on an input split**

- A typical chunk is 64MB (more or less can be configured depending on your use case)

**Mappers read in a chunk of data**

**Mappers emit (write out) a set of data, typically derived from its input**

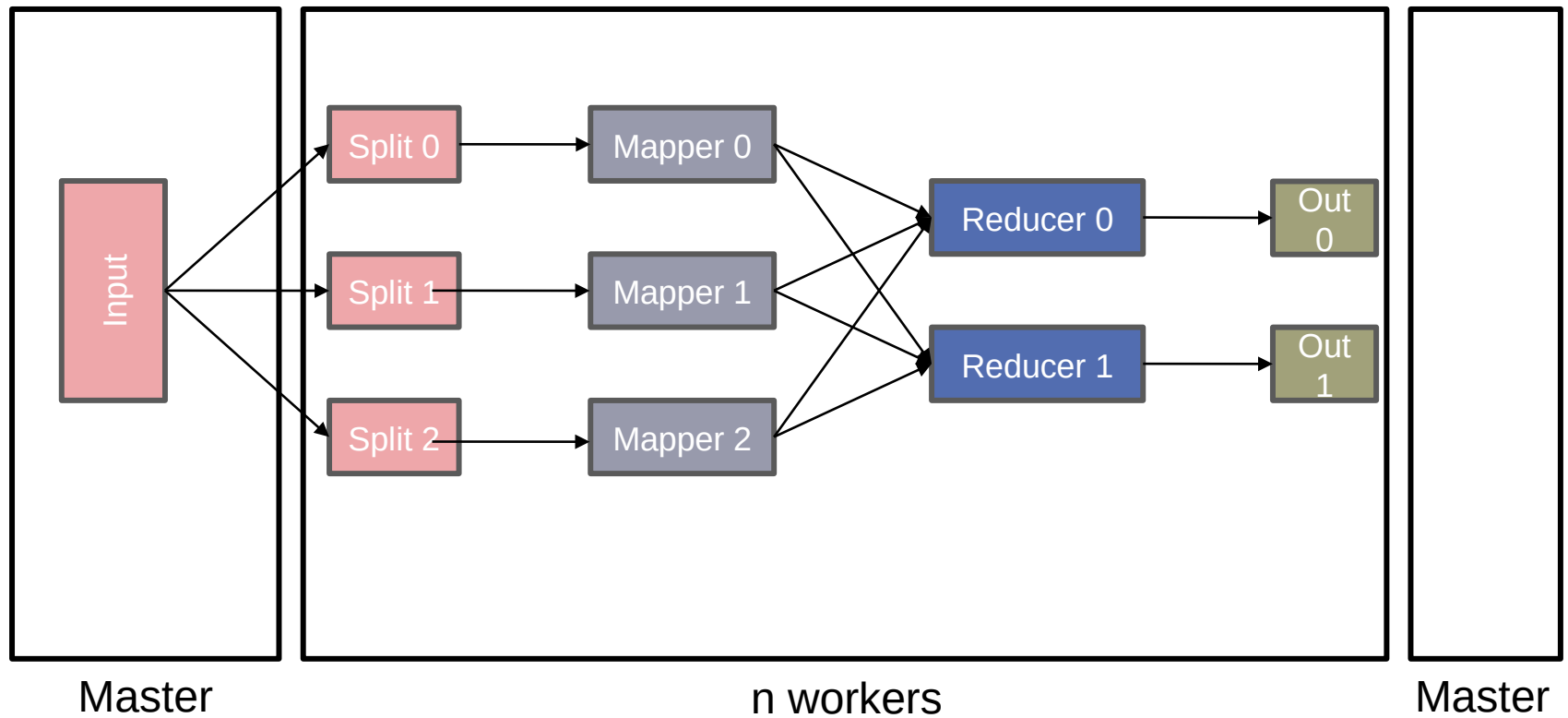
**Intermediate data (the output of the mappers) is split to a number of reducers**

**Reducers receive each key of data, along with ALL of the values associated with it (this means each key must always be sent to the same reducer)**

- Essentially, <key, set<value>>

**Reducers emit a set of data, typically reduced from its input which is written to disk**

# DATA FLOW



# INPUT SPLITTER

**Is responsible for splitting your input into multiple chunks**

**These chunks are then used as input for your mappers**

**Splits on logical boundaries. The default is 64MB per chunk**

- Depending on what you're doing, 64MB might be a LOT of data!  
You can change it

**Typically, you can just use one of the built in splitters, unless you are reading in a specially formatted file**

# MAPPER

**Reads in input pair <K,V> (a section as split by the input splitter)**

**Outputs a pair <K', V'>**

**Ex. For our Word Count example, with the following input: “The teacher went to the store. The store was closed; the store opens in the morning. The store opens at 9am.”**

**The output would be:**

- <The, 1> <teacher, 1> <went, 1> <to, 1> <the, 1> <store, 1>  
<the, 1> <store, 1> <was, 1> <closed, 1> <the, 1> <store, 1>  
<opens, 1> <in, 1> <the, 1> <morning, 1> <the 1> <store, 1>  
<opens, 1> <at, 1> <9am, 1>

# REDUCER

**Accepts the Mapper output, and collects values on the key**

- All inputs with the same key must go to the same reducer!

**Input is typically sorted, output is output exactly as is**

**For our example, the reducer input would be:**

- <The, 1> <teacher, 1> <went, 1> <to, 1> <the, 1> <store, 1>  
<the, 1> <store, 1> <was, 1> <closed, 1> <the, 1> <store, 1>  
<opens, 1> <in, 1> <the, 1> <morning, 1> <the 1> <store, 1>  
<opens, 1> <at, 1> <9am, 1>

**The output would be:**

- <The, 6> <teacher, 1> <went, 1> <to, 1> <store, 3> <was, 1>  
<closed, 1> <opens, 1> <morning, 1> <at, 1> <9am, 1>

# COMBINER

**Essentially an intermediate reducer**

- Is optional

**Reduces output from each mapper, reducing bandwidth and sorting**

**Cannot change the type of its input**

- Input types must be the same as output types



# OUTPUT COMMITTER

**Is responsible for taking the reduce output, and committing it to a file**

**Typically, this committer needs a corresponding input splitter (so that another job can read the input)**

**Again, usually built in splitters are good enough, unless you need to output a special kind of file**

# PARTITIONER (SHUFFLER)

**Decides which pairs are sent to which reducer**

**Default is simply:**

- `Key.hashCode() % numOfReducers`

**User can override to:**

- Provide (more) uniform distribution of load between reducers
- Some values might need to be sent to the same reducer
  - Ex. To compute the relative frequency of a pair of words  $\langle W1, W2 \rangle$  you would need to make sure all of word  $W1$  are sent to the same reducer
- Binning of results

# MASTER

**Responsible for scheduling & managing jobs**

**Scheduled computation should be close to the data if possible**

- Bandwidth is expensive! (and slow)
- This relies on a Distributed File System (e.g. GFS)!

**If a task fails to report progress (such as reading input, writing output, etc), crashes, the machine goes down, etc, it is assumed to be stuck, and is killed, and the step is re-launched (with the same input)**

**The Master is handled by the framework, no user code is necessary**

# MAPREDUCE IN PYTHON

```
def mapreduce_execute(data, mapper, reducer):
    values = map(mapper, data)

    groups = {}
    for items in values:
        for k,v in items:
            if k not in groups:
                groups[k] = [v]
            else:
                groups[k].append(v)

    output = [reducer(k,v) for k,v in groups.items()]
    return output
```

# MAPREDUCE IN PYTHON

Don't do the last slide ...

Python's `mrjob` library:

- write mappers and reducers in Python
- Deploy on Hadoop systems, Amazon Elastic MR, Google Cloud

```
from mrjob.job import MRJob

class WordOccurrenceCount(MRJob):
    def mapper(self, _, line):
        for word in line.split(" "):
            yield word, 1

    def reducer(self, key, values):
        yield key, sum(values)
```

# MAPREDUCE?

## **Good:**

- All you need to do is write a mapper and a reducer
- Can get away with not exposing any of the internals (data splitting, locality issues, redundancy, etc) if you're using a ready-made engine

## **Bad:**

- Lots of reading/writing from disk (in part because this helps with redundancy)
- Sometimes communication between processes is necessary
- Talk about later: parameter servers, GraphLab aka Dato, etc

*NEXT UP:*

REVIEW OF HYPOTHESIS TESTING  
(AND THEN A BUNCH OF STUFF LIKE PRIVACY,  
ETHICS, DEBUGGING DATA SCIENCE, ETC!)