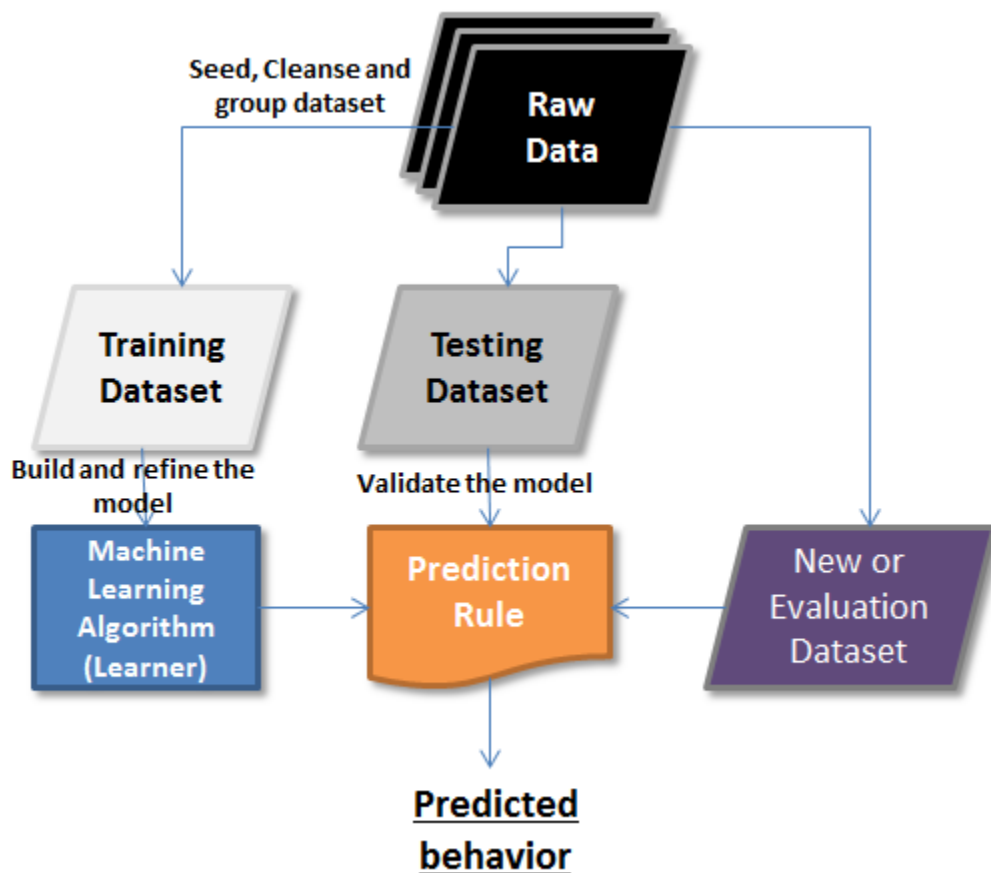
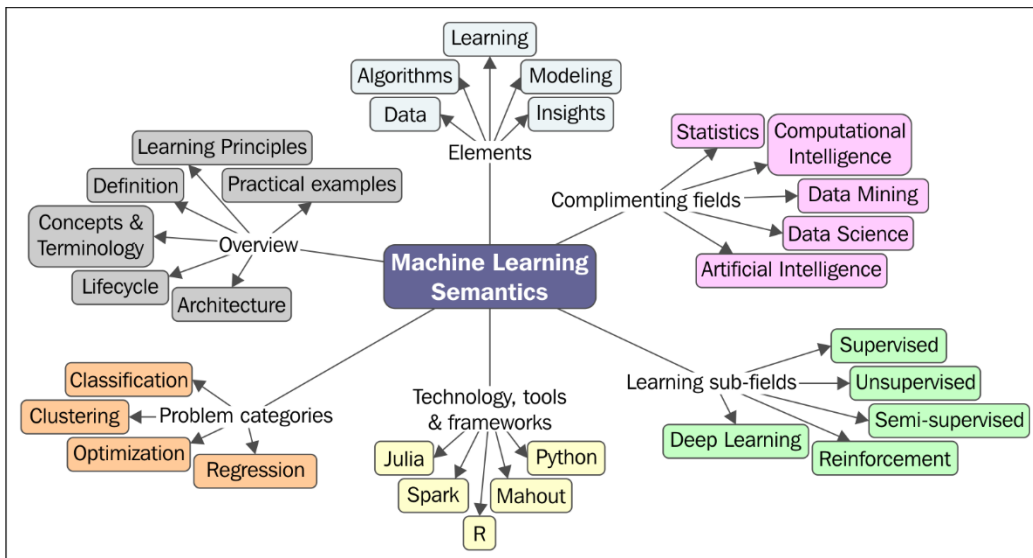
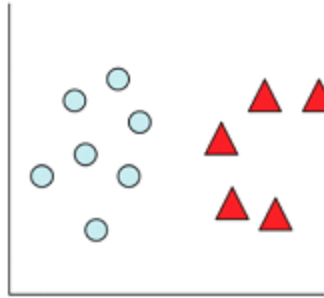
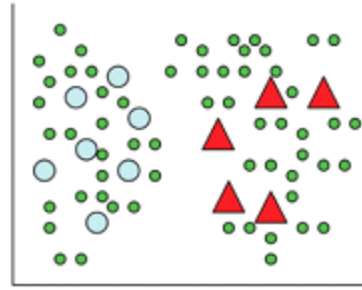


CHAPTER 1

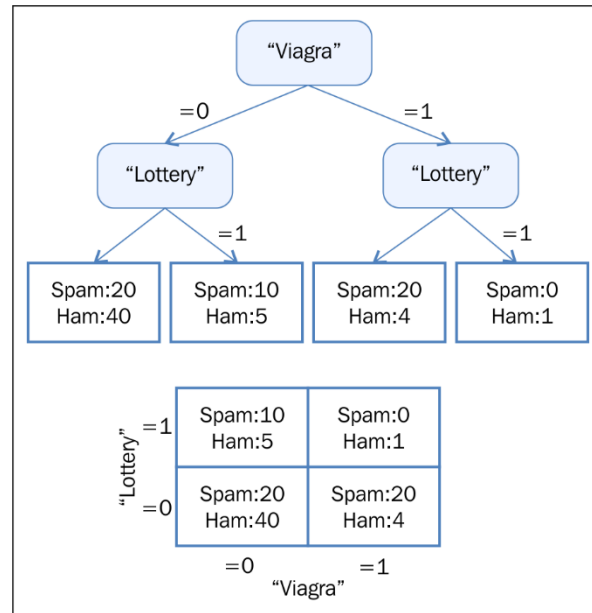


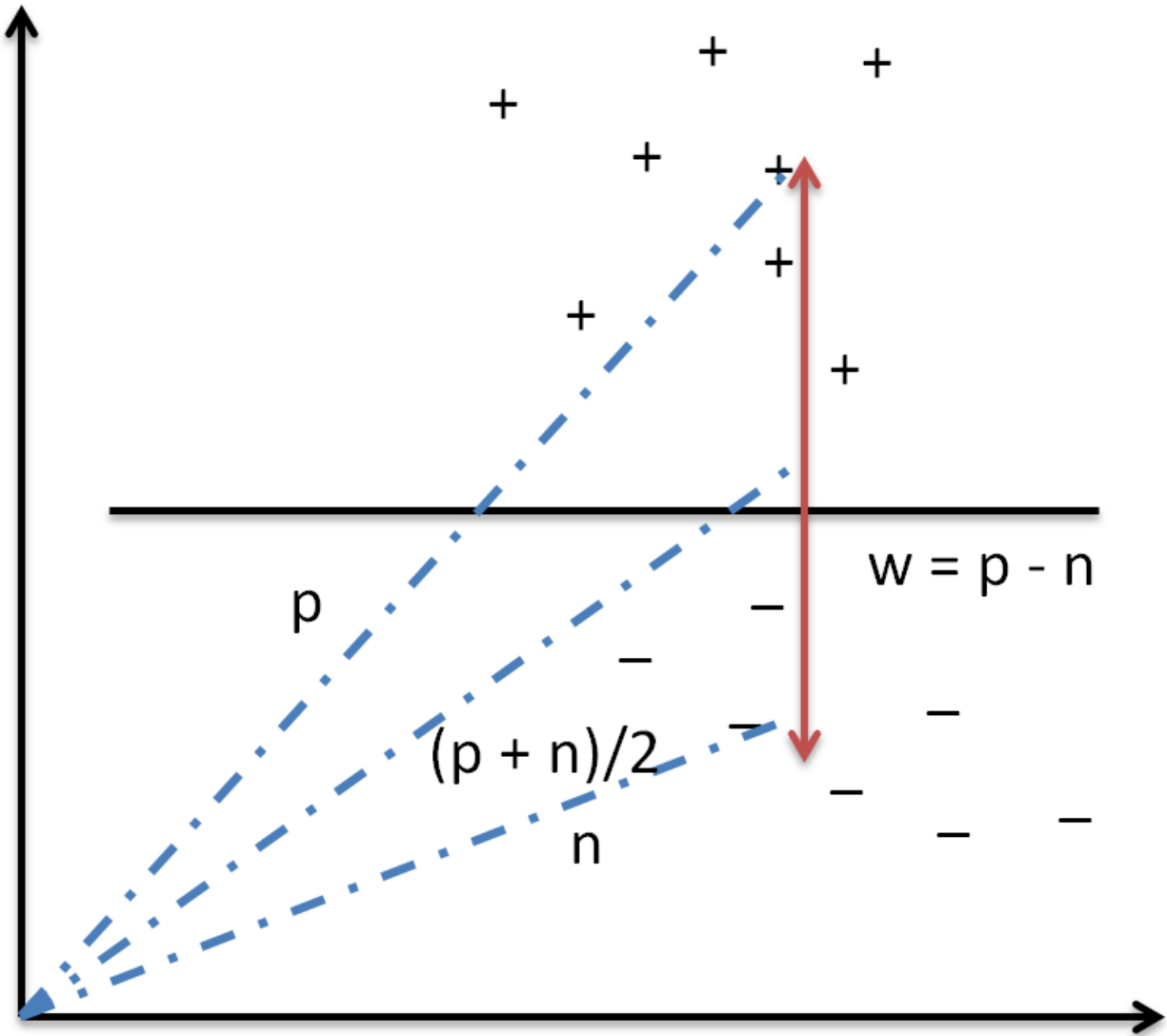


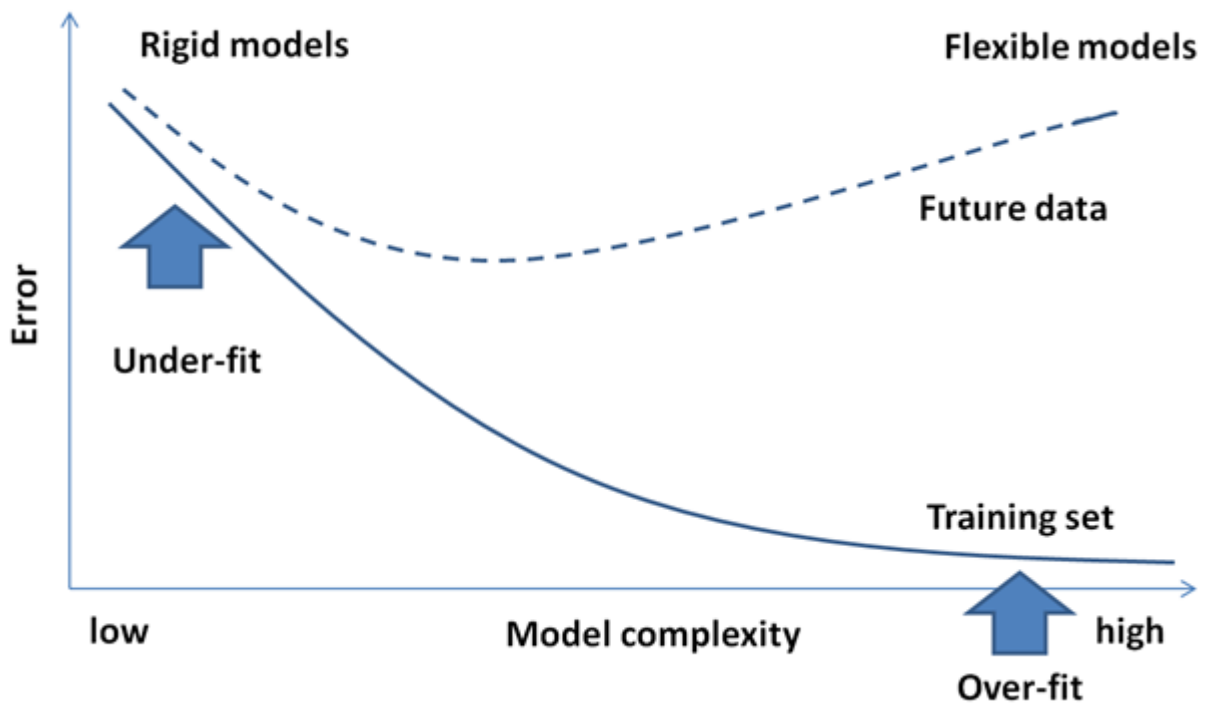
Labeled Data

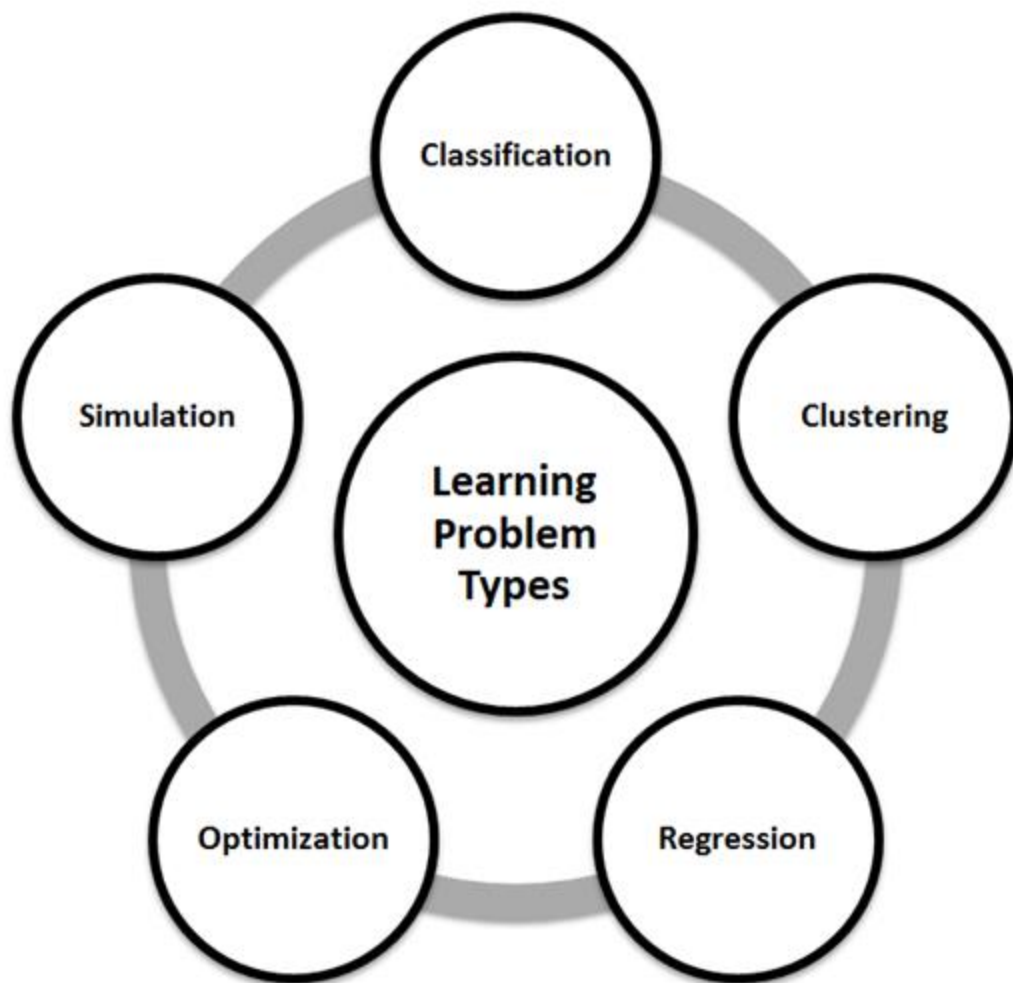


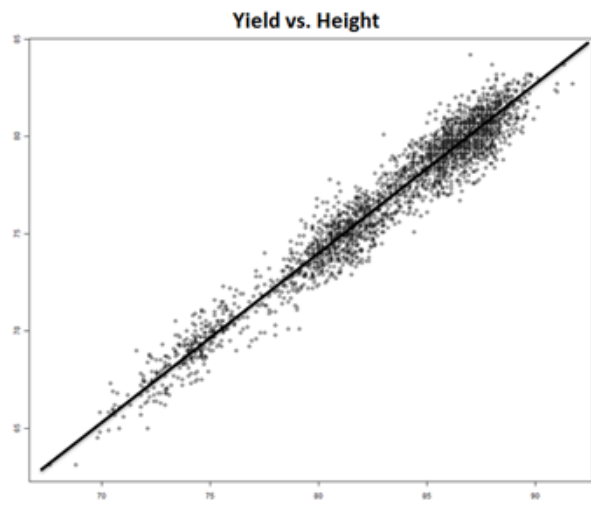
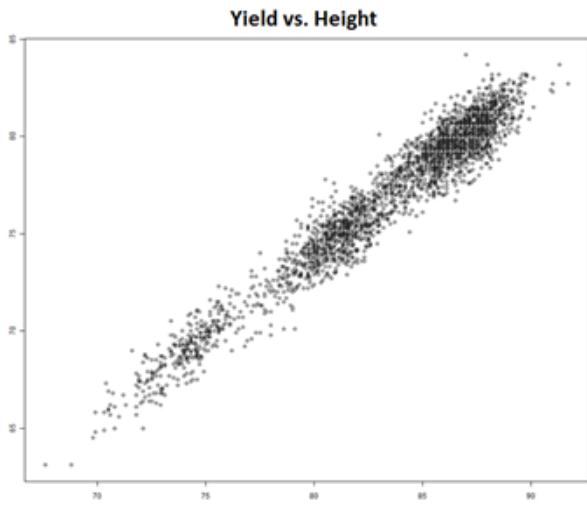
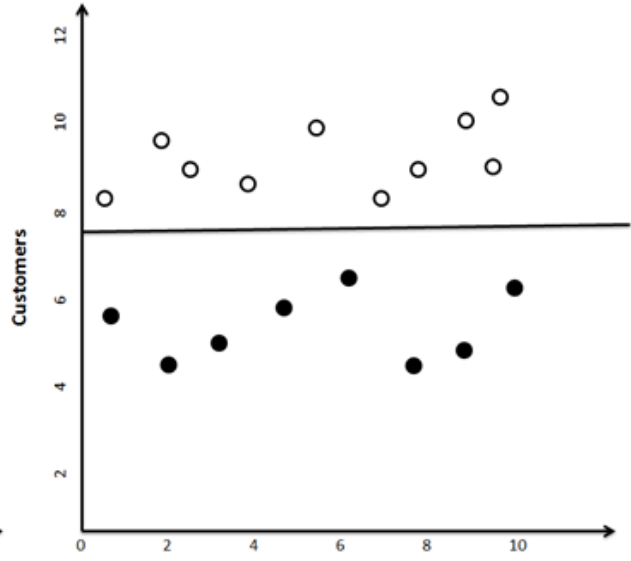
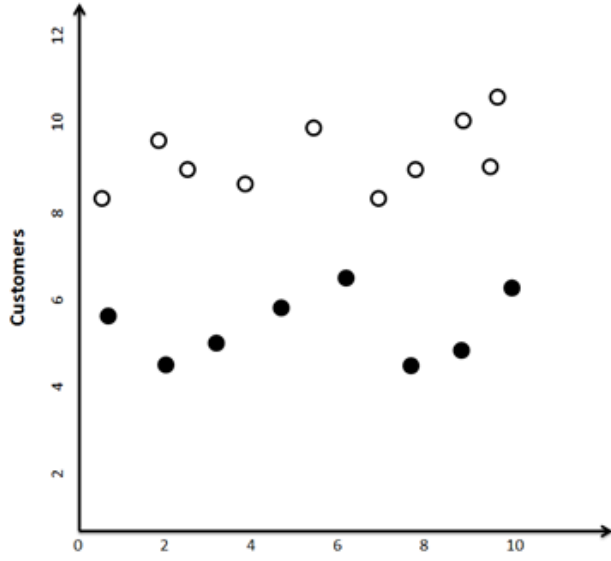
Labeled and Unlabeled Data

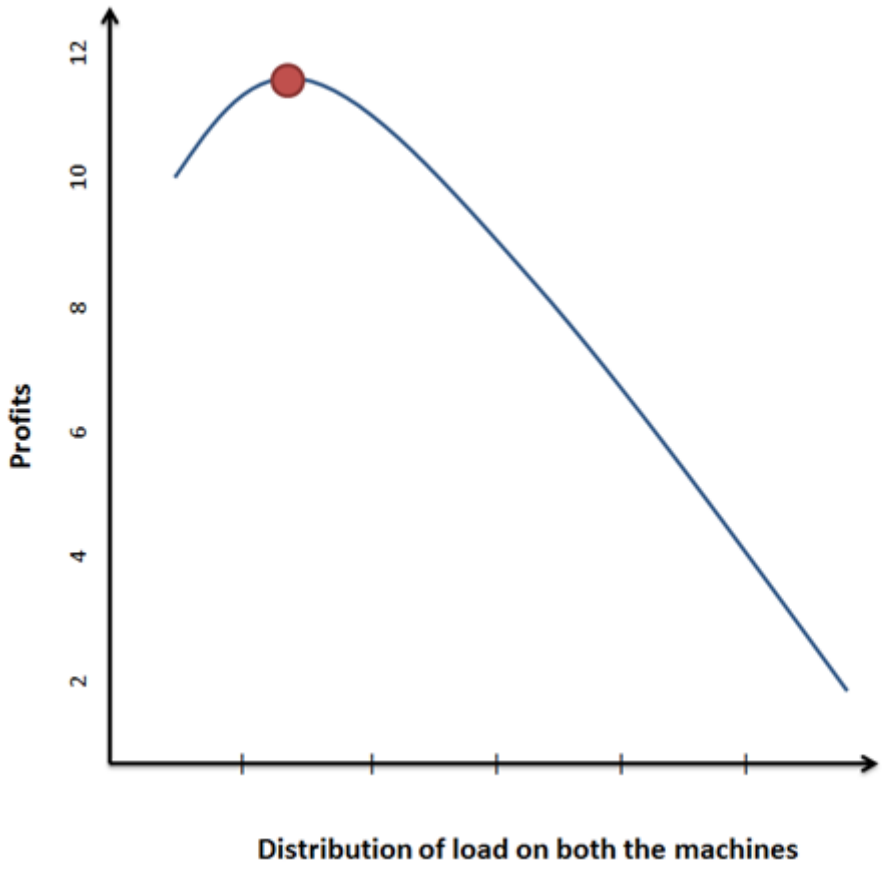


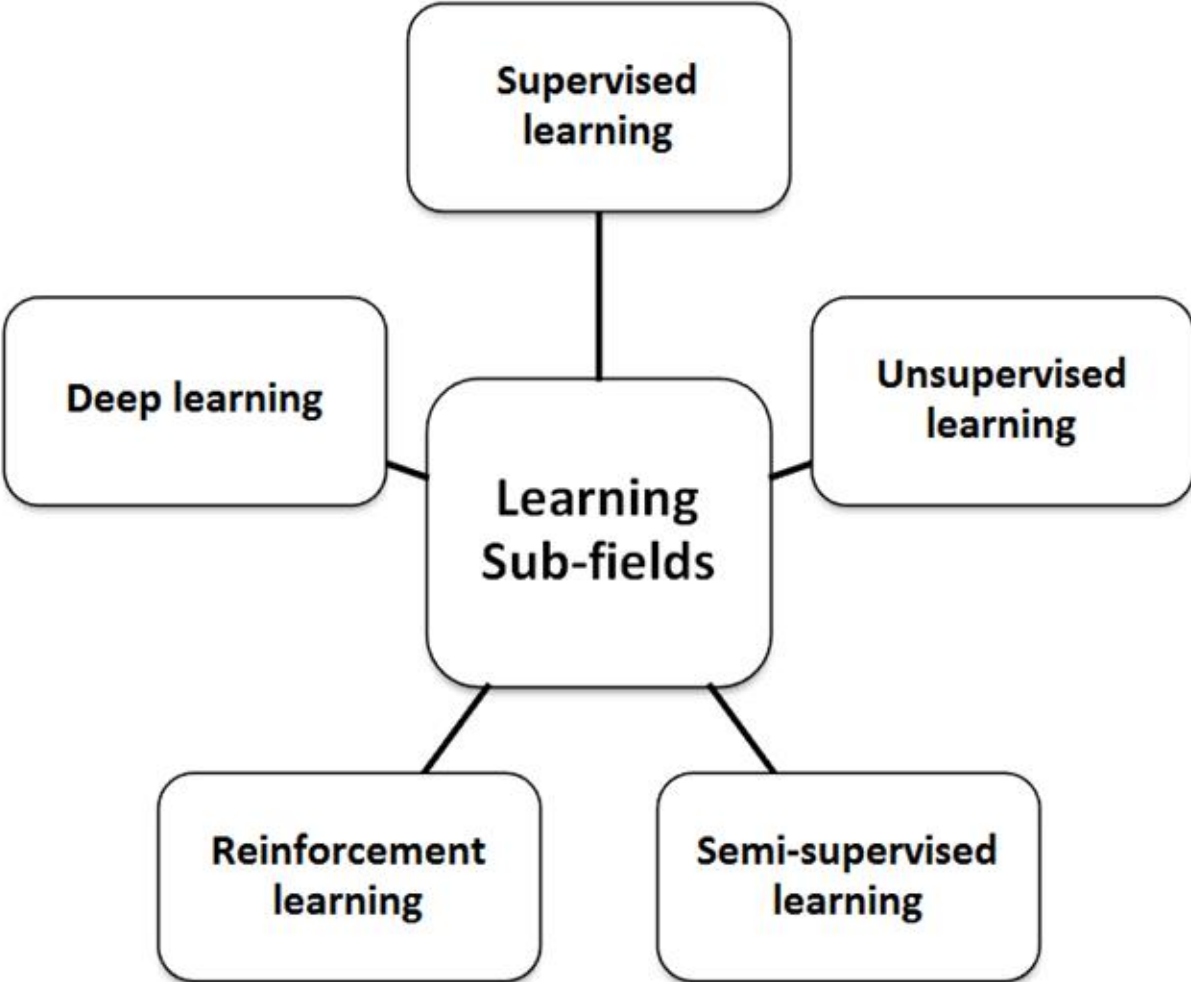


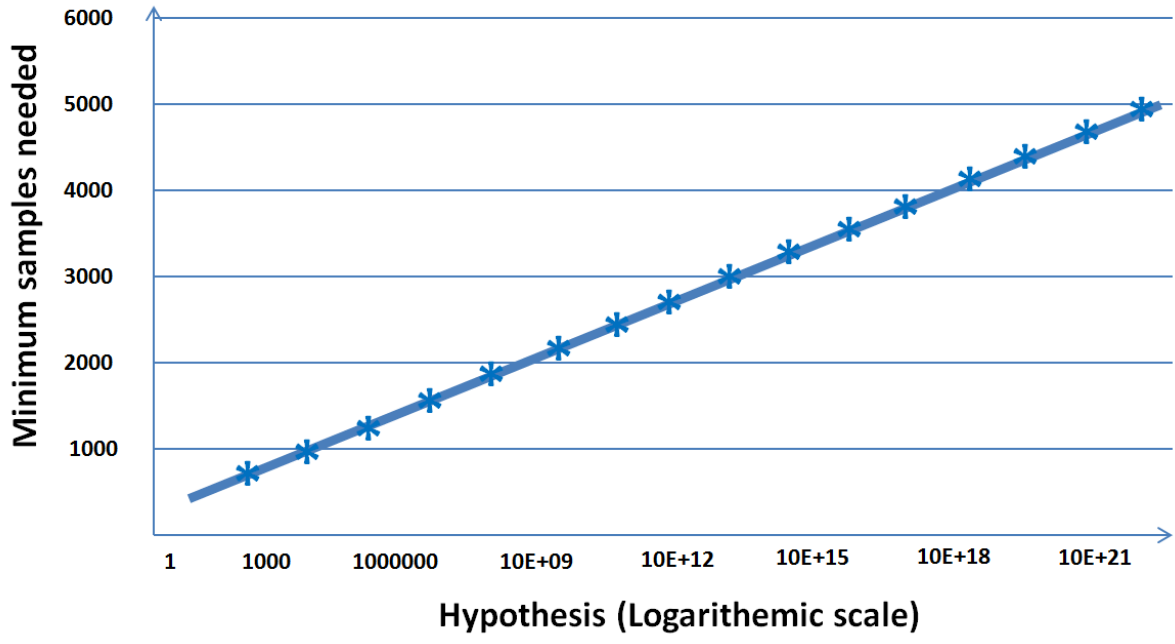








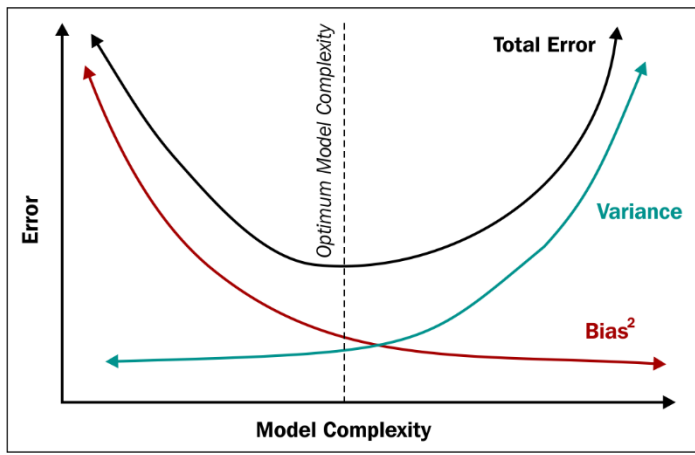


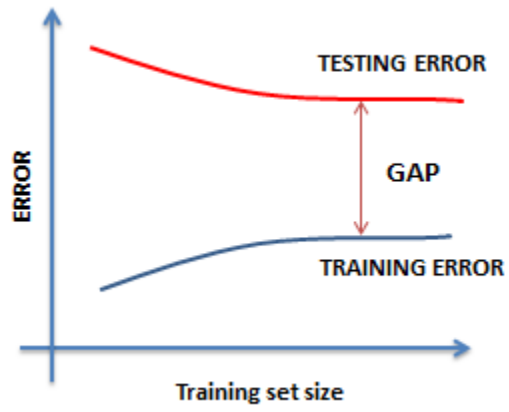
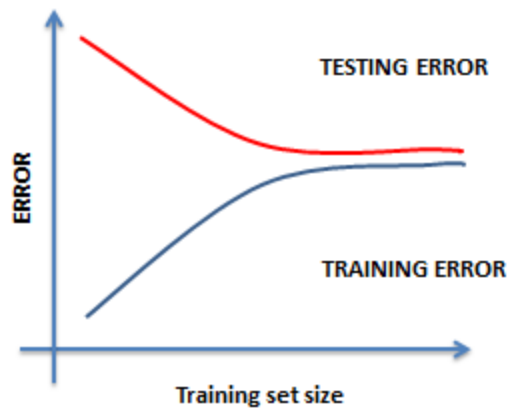


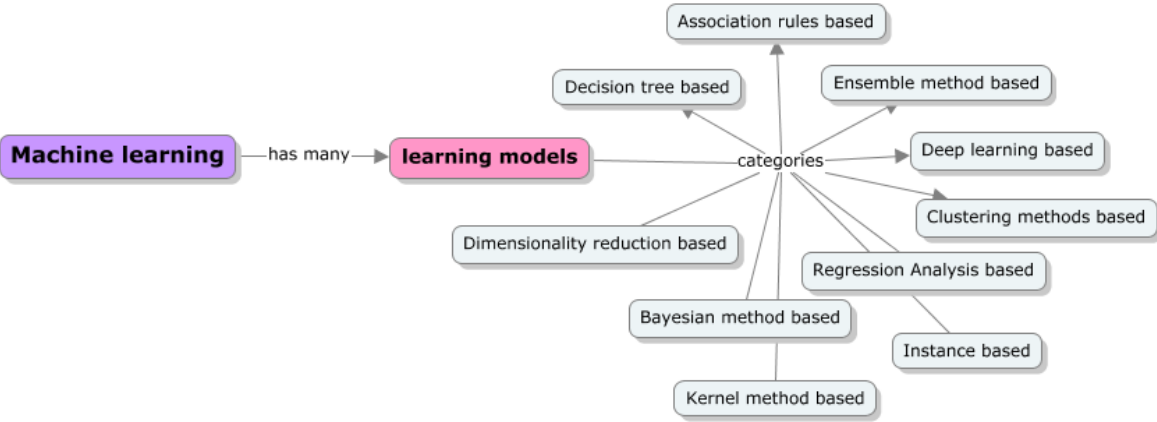
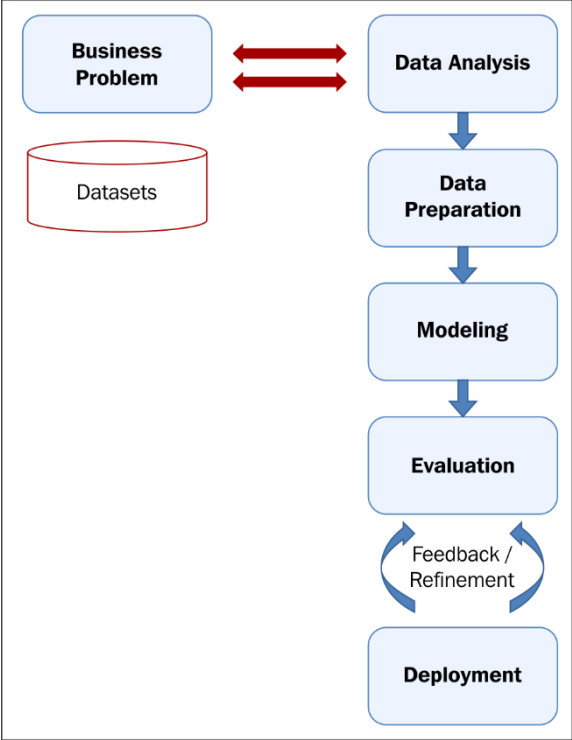
$$MSE = \frac{\sum_{i=1}^n (P_i - A_i)^2}{n}$$

$$MAE = \frac{\sum_{i=1}^n |P_i - A_i|}{n}$$

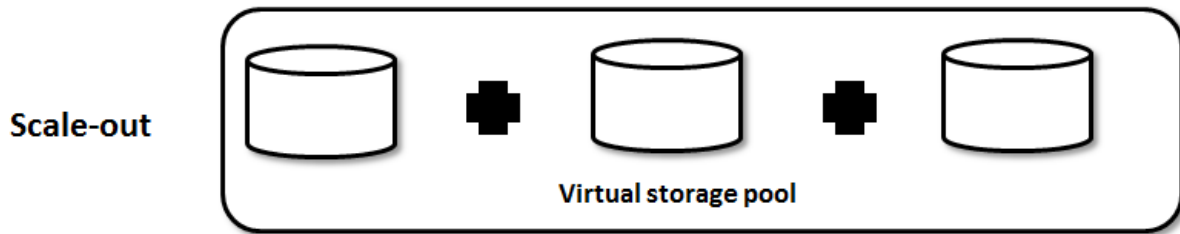
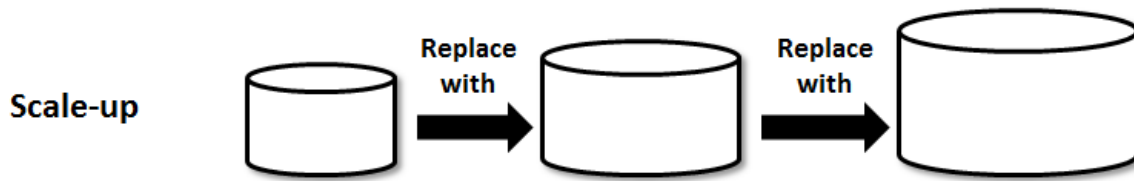
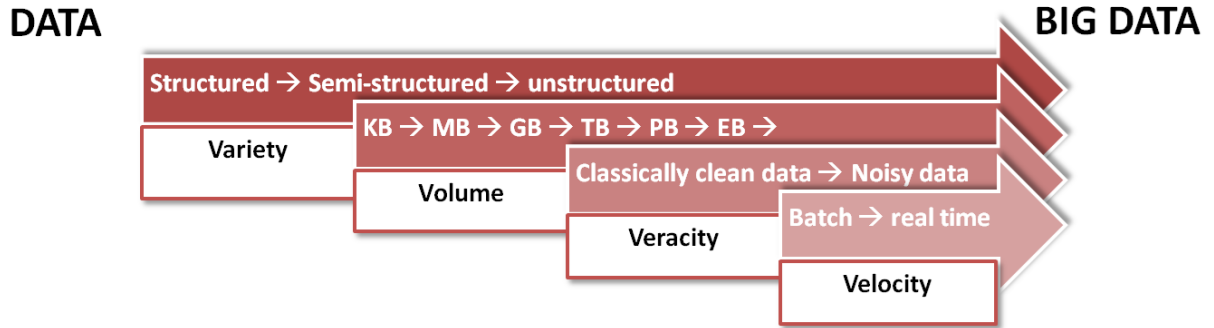
$$NMSE = \frac{MSE \text{ of developed model}}{MSE \text{ of naive model}}$$

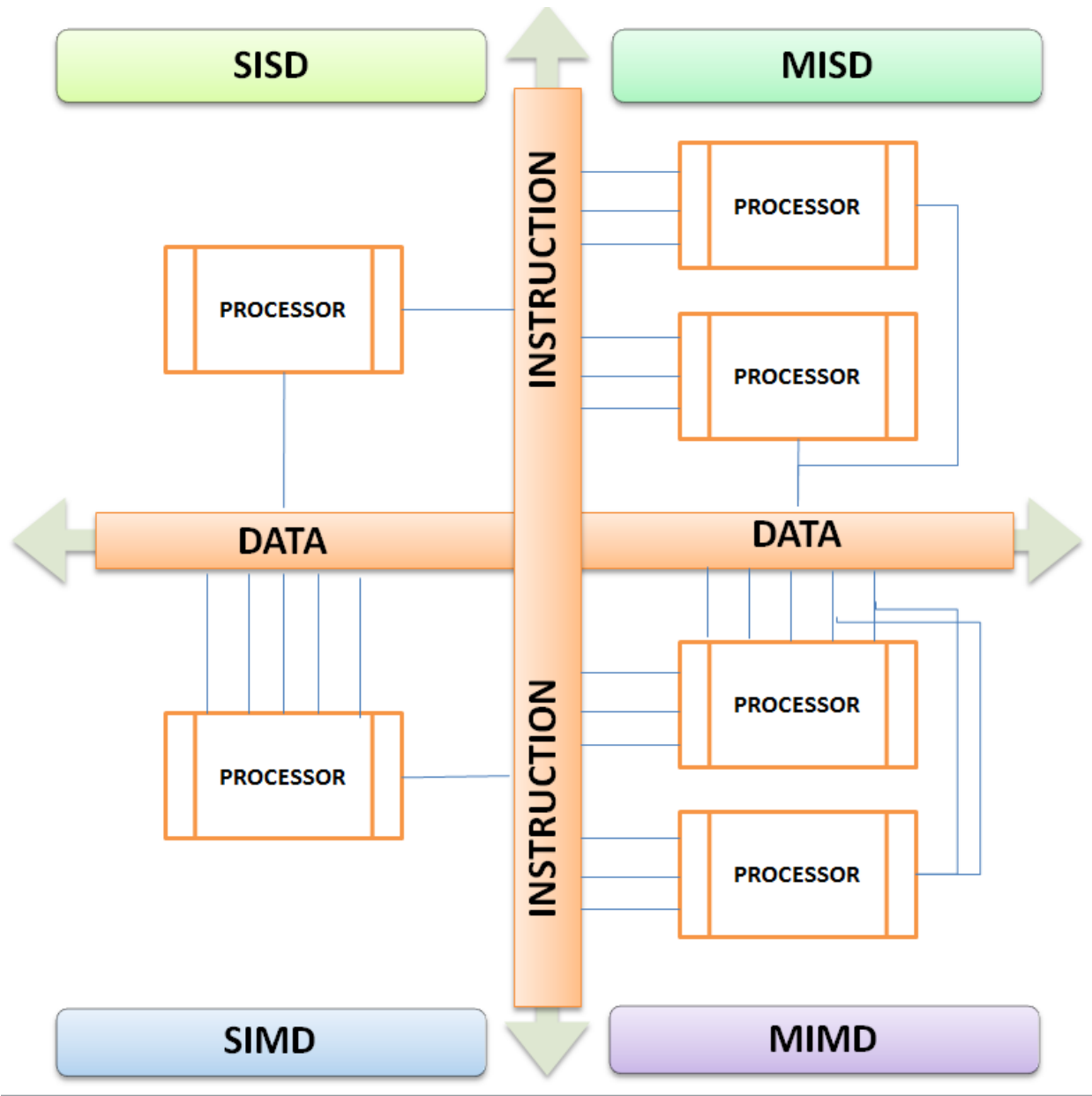


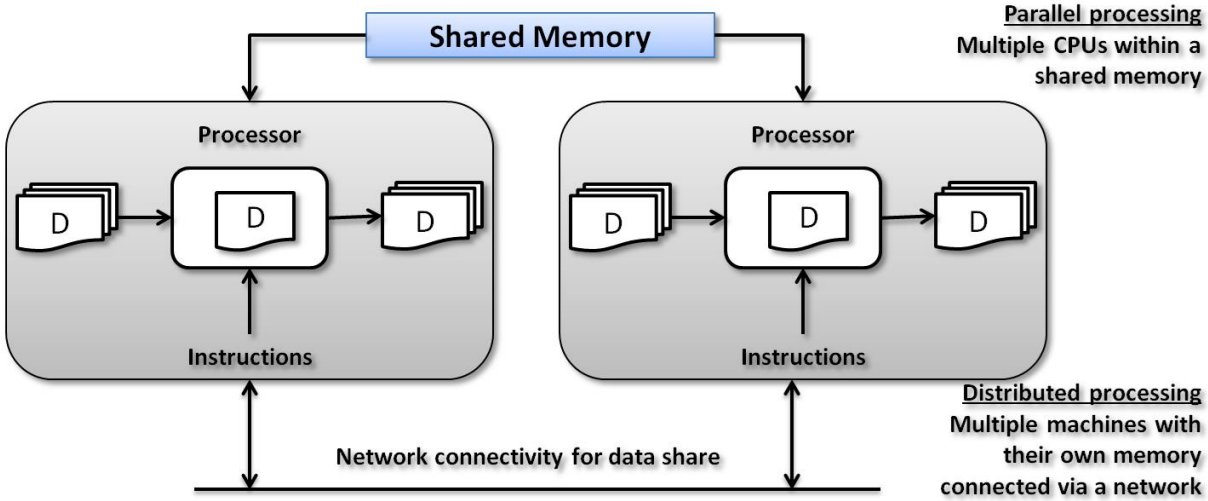
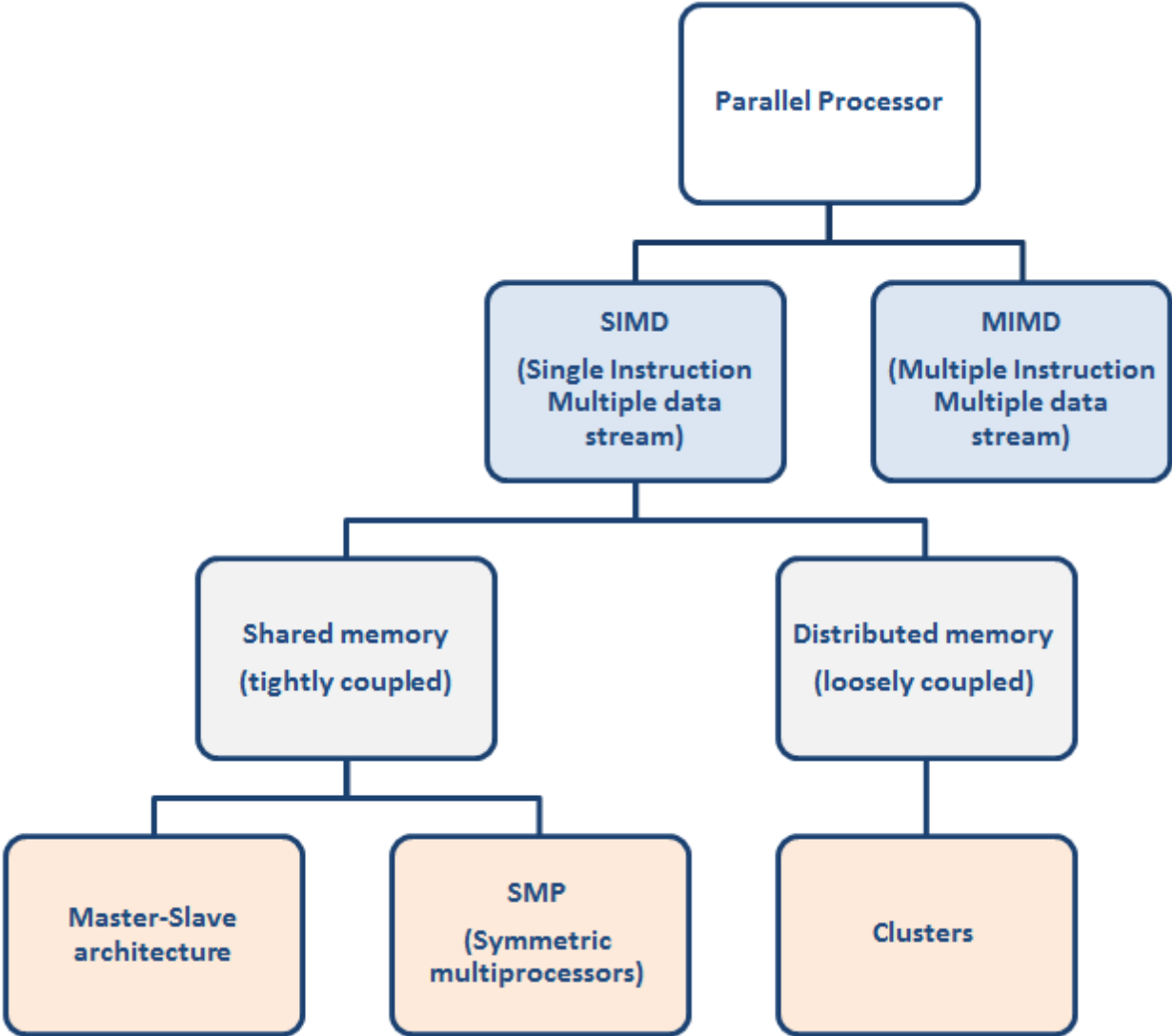




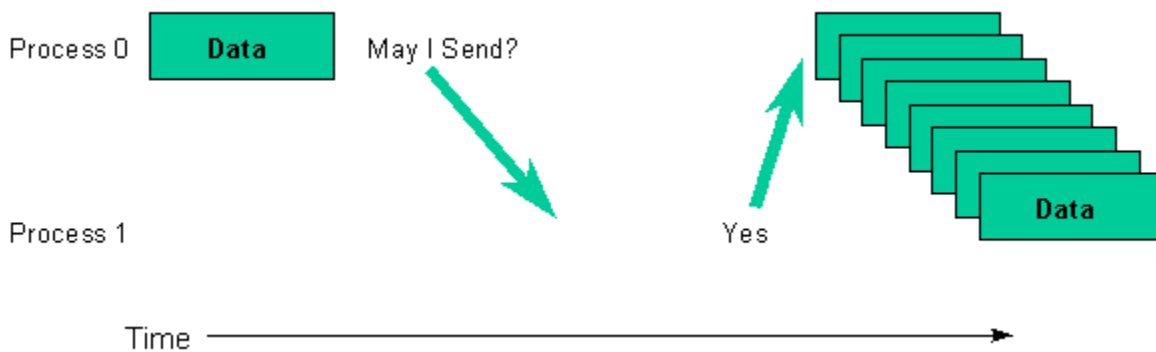
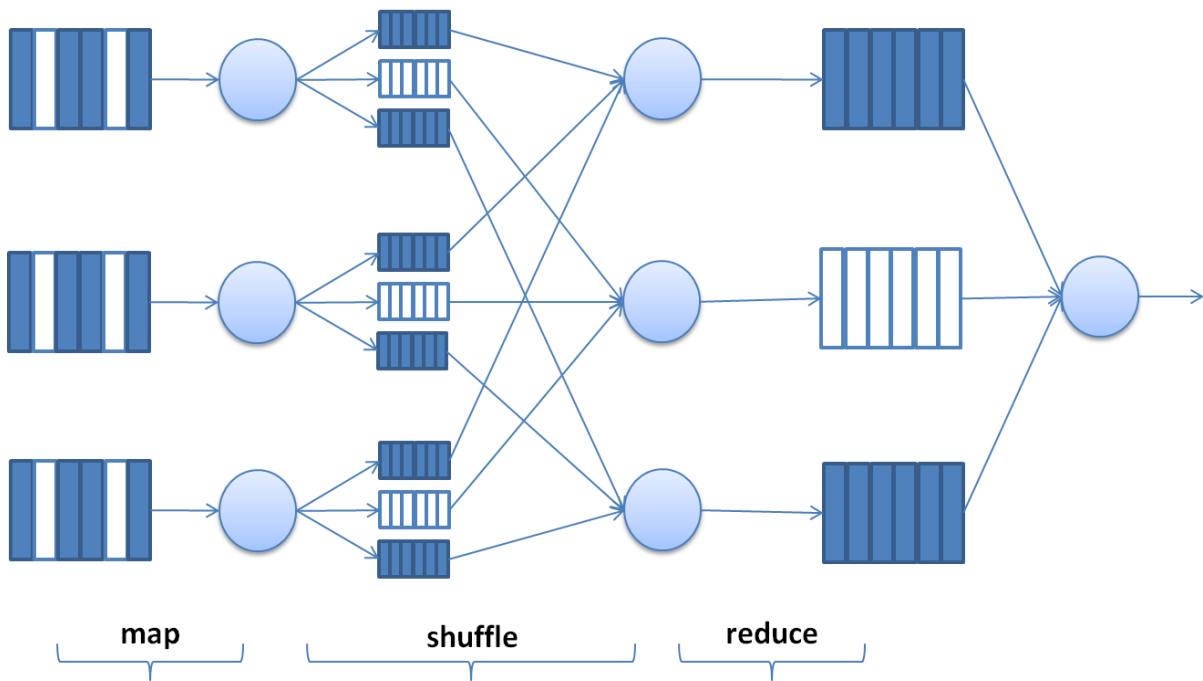
CHAPTER 2

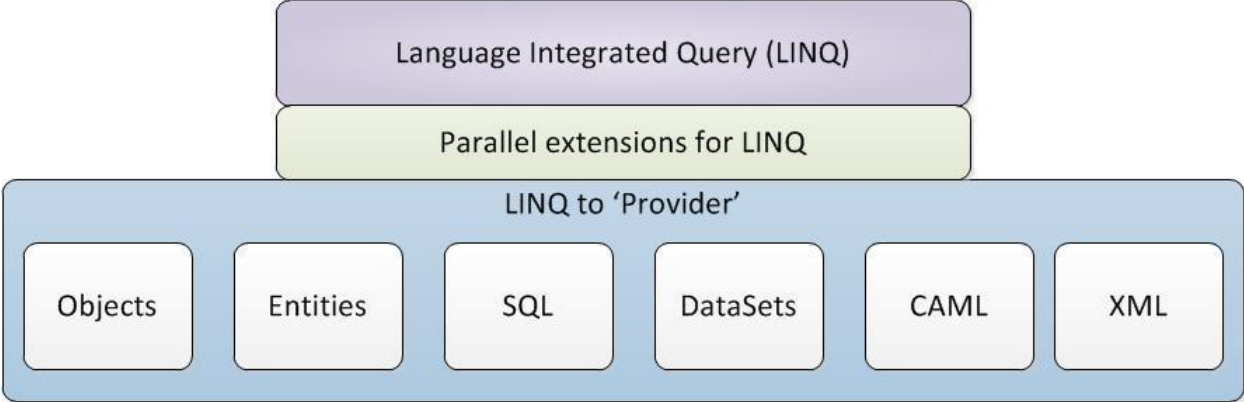






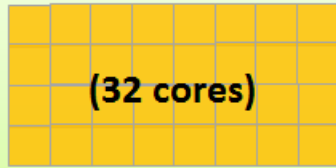
ID#	Model	Year	Color	Dealer	Price
4523	Civic	2002	Blue	MN	\$18,000
3476	Corolla	1999	White	IL	\$15,000
7623	Camry	2001	Green	NY	\$21,000
9834	Prius	2001	Green	CA	\$18,000
6734	Civic	2001	White	OR	\$17,000
5342	Altima	2001	Green	FL	\$19,000
3845	Maxima	2001	Blue	NY	\$22,000
8354	Accord	2000	Green	VT	\$18,000
4395	Civic	2001	Red	CA	\$17,000
7352	Civic	2002	Red	WA	\$18,000



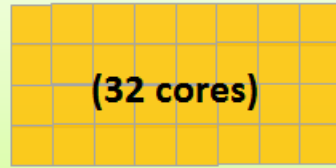


GPU

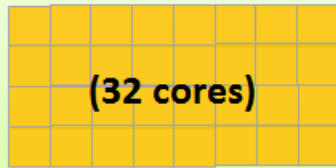
Multiprocessor 1



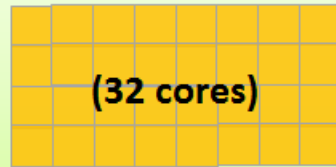
Multiprocessor 2



Multiprocessor 3



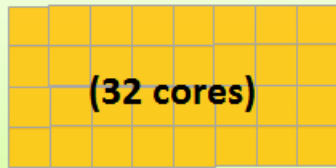
Multiprocessor 4



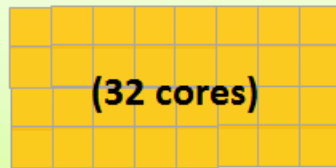
...

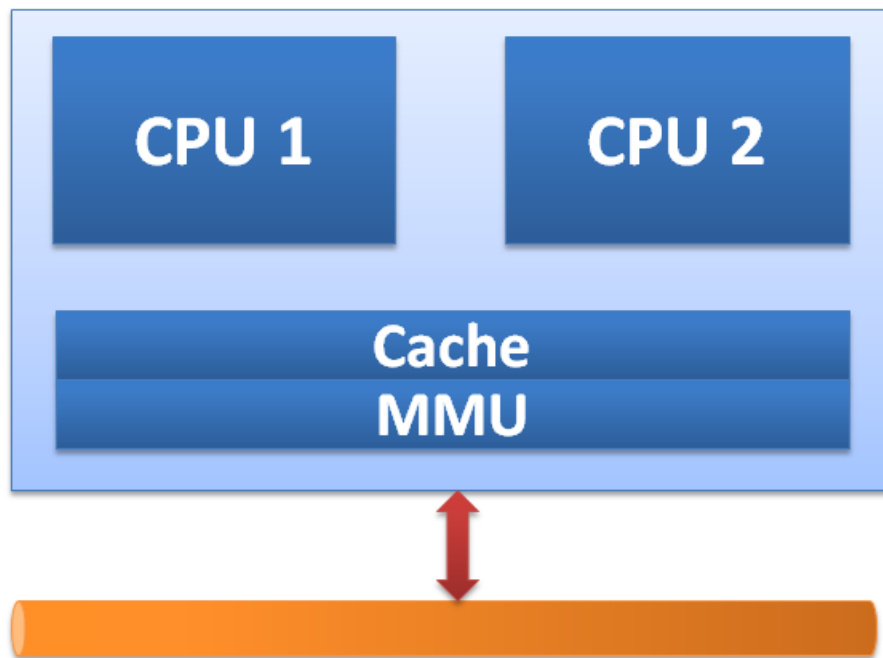
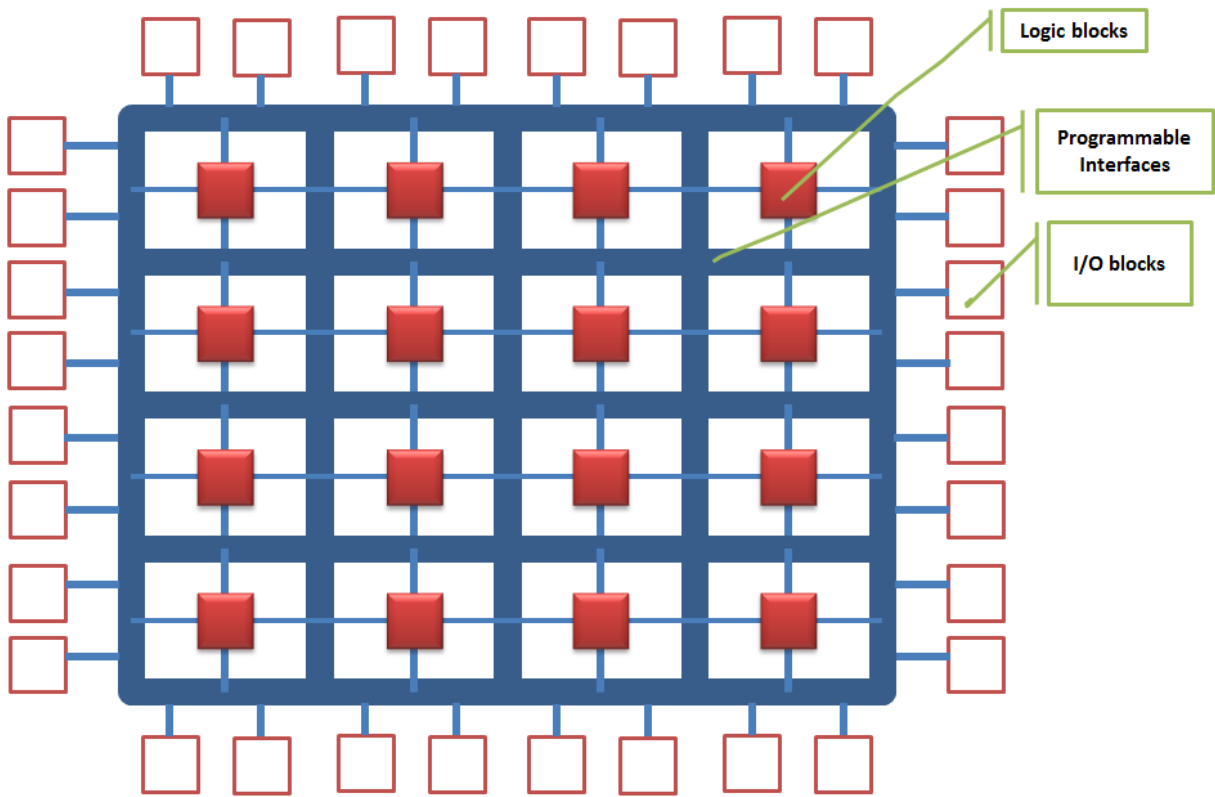
...

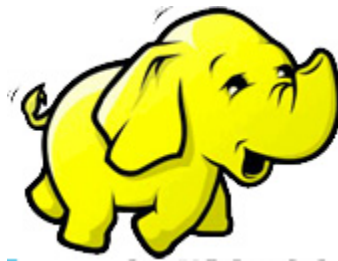
Multiprocessor 12



Multiprocessor 13







Year Evolution / Progress

2002-2003 Work on Nutch was started by Doug Cutting and Mike Cafarella

2003 - 2004 Google published work on GFS and MapReduce

2004 Doug Cutting added support for GFS and MapReduce to Nutch

2006 Hadoop spins out of Nutch when Yahoo hired Doug Cutting

2007 NY Times converts 4TB of image archives over 100 EC2s

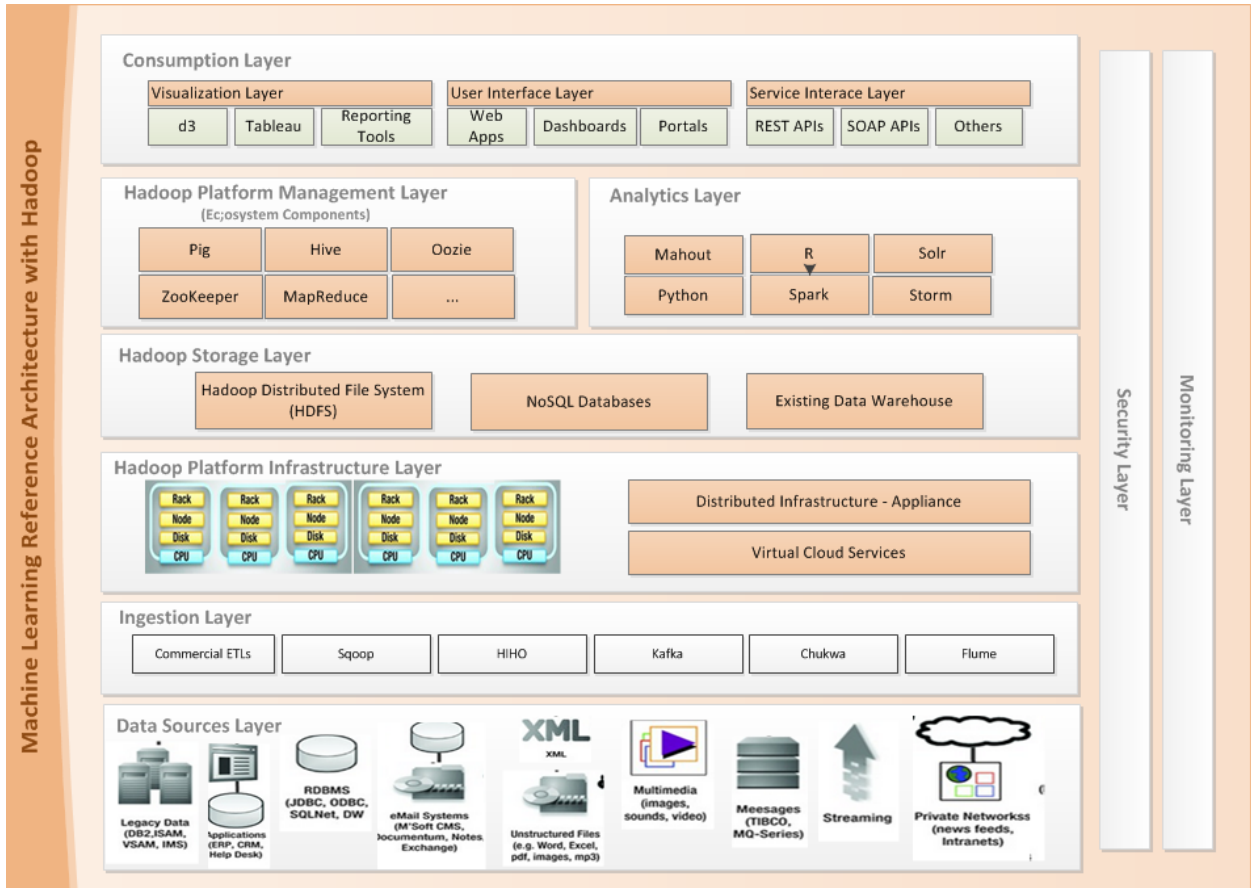
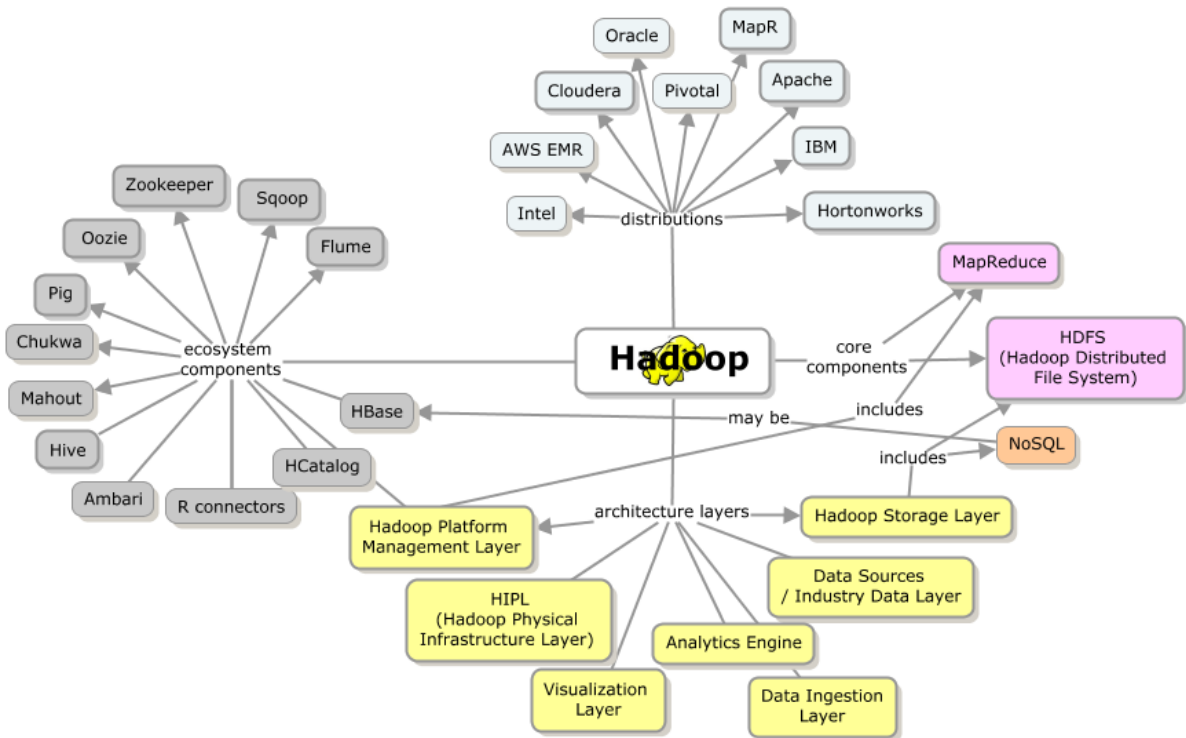
Facebook launched Hive, an SQL support for Hadoop

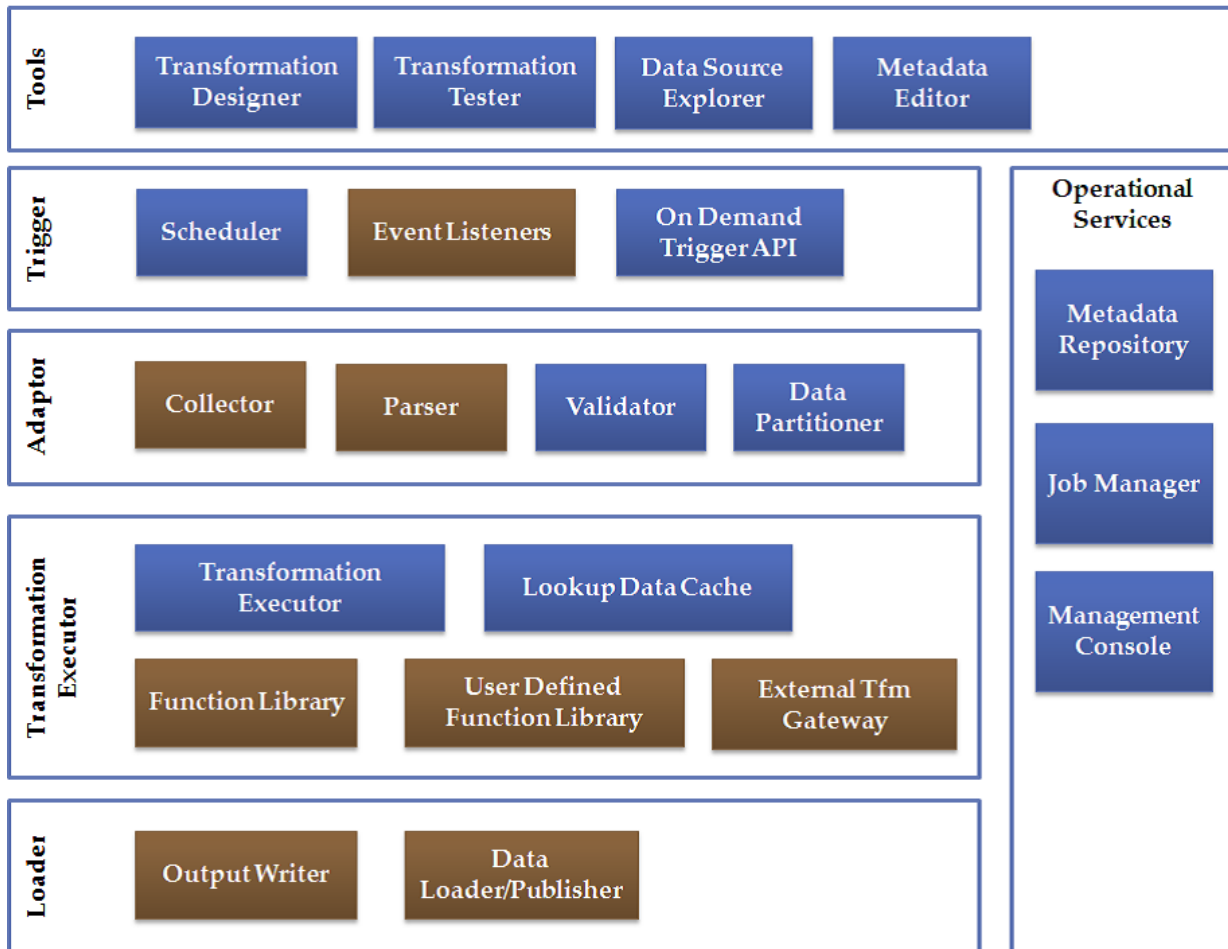
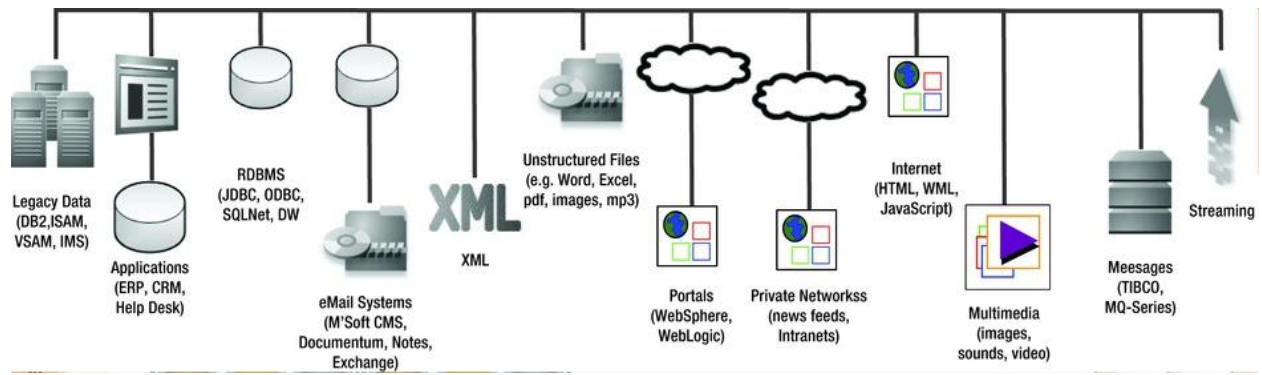
2008 Fastest sort over 910 nodes taking 3.5 mins

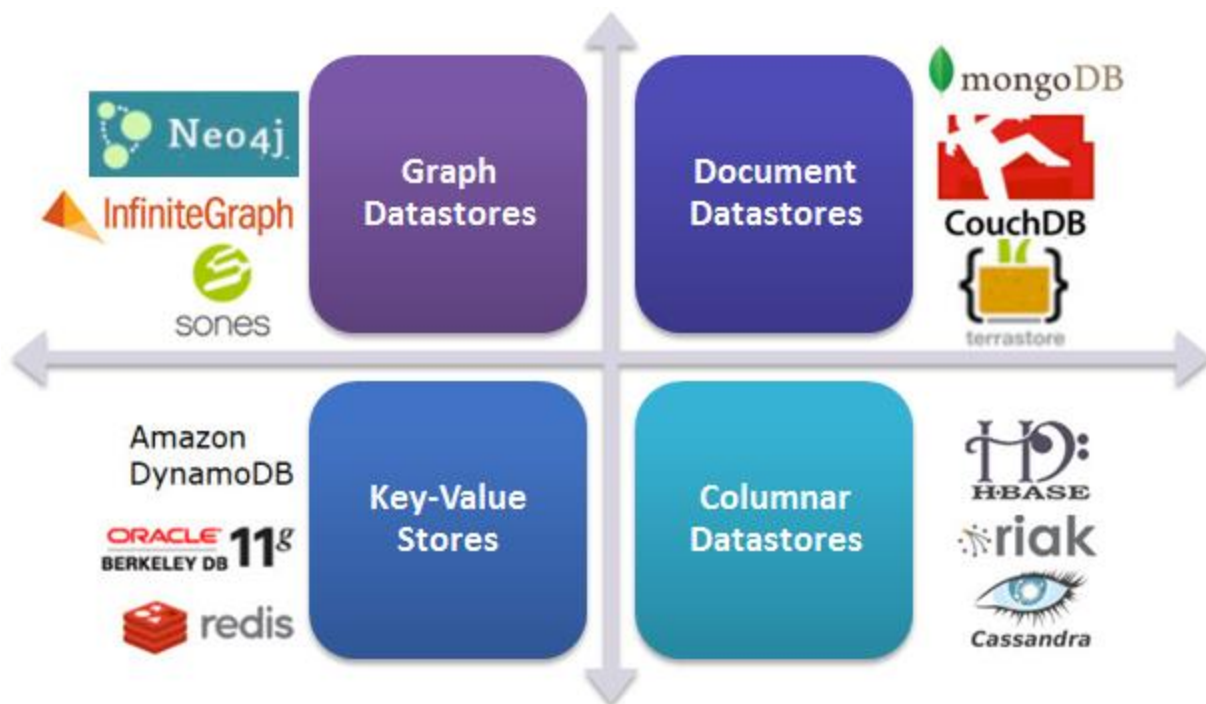
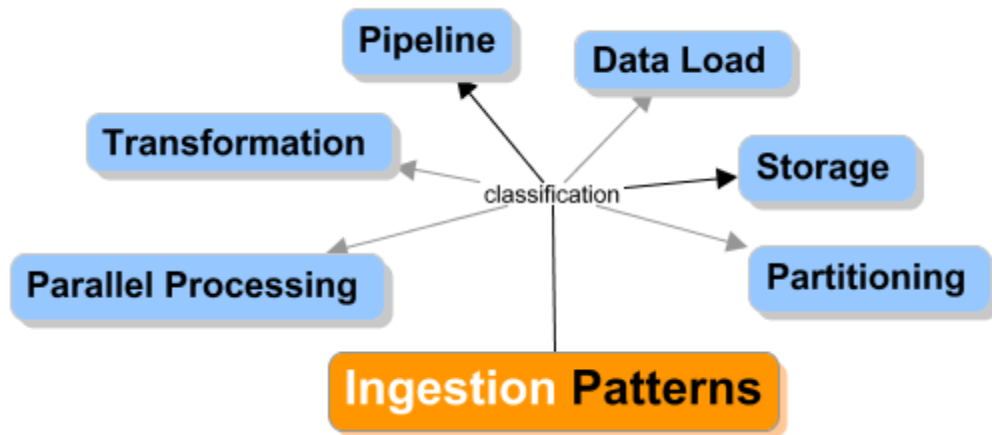
Cloudera founded

2009 First Hadoop Summit with 750 attendees

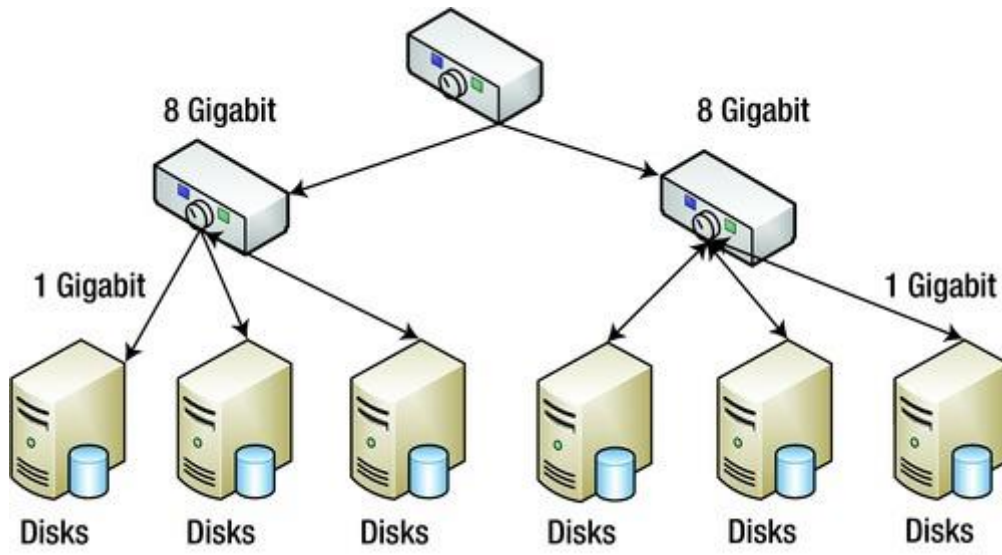
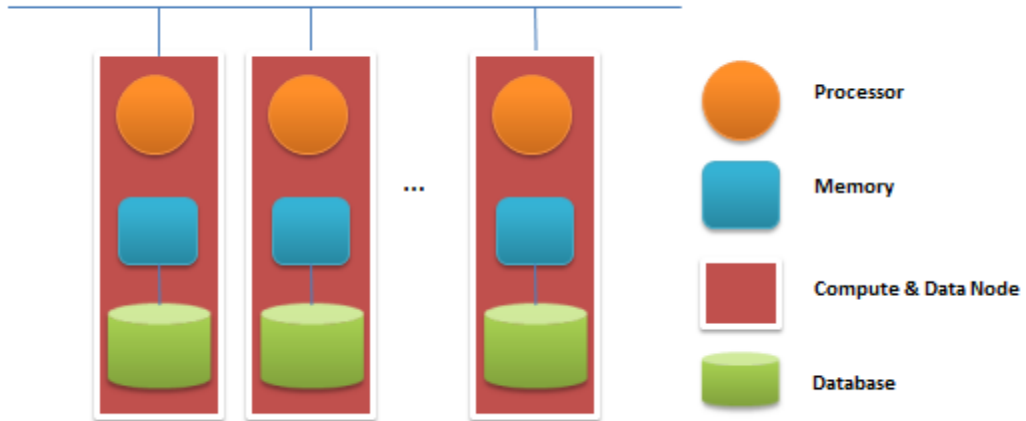
Doug Cutting joined Cloudera







Shared Nothing Data Architecture

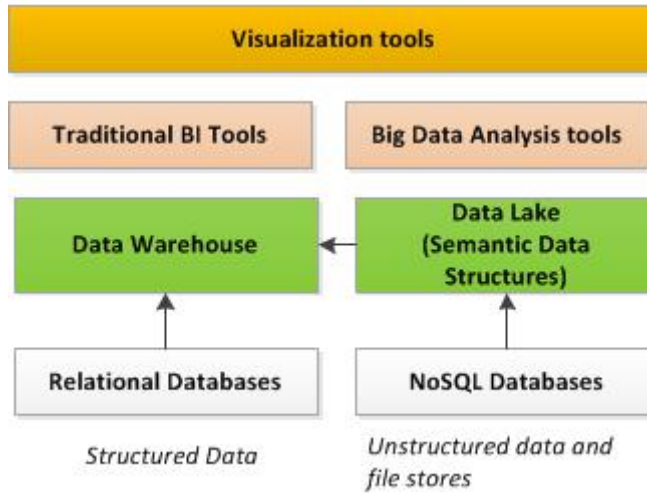


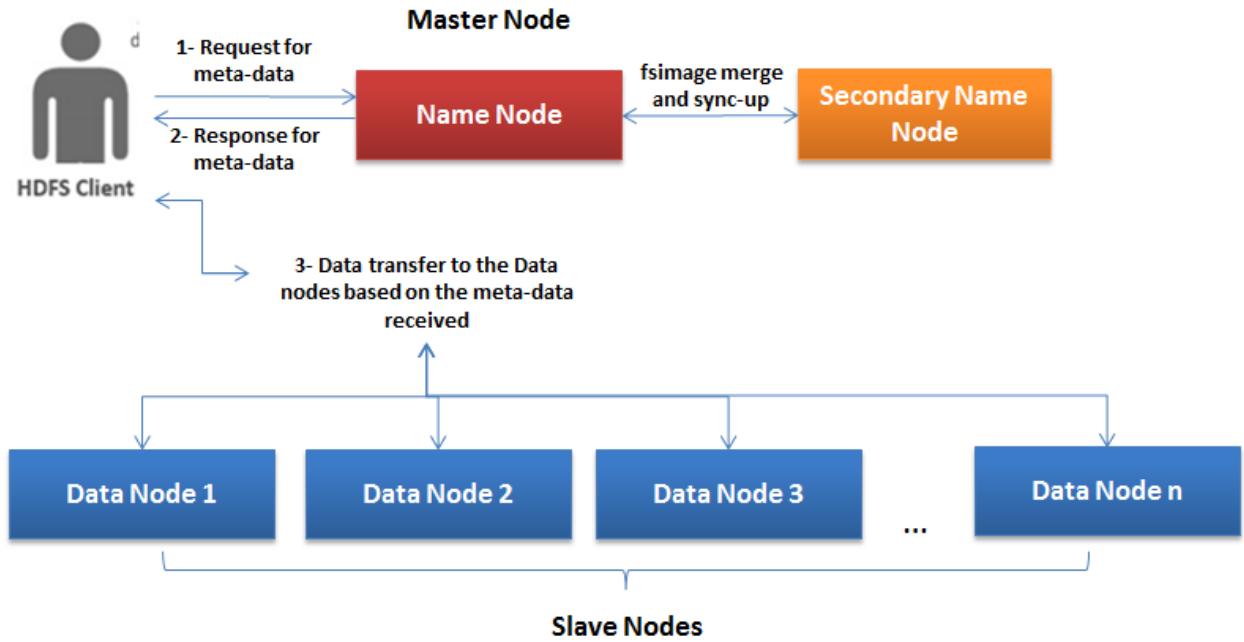
VISUALIZATIONS LAYER

Dashboards & Reports **Insight Analytics & Big Data Models** **Outlook Add-In & Excel Add-In** **BusinessObjects**

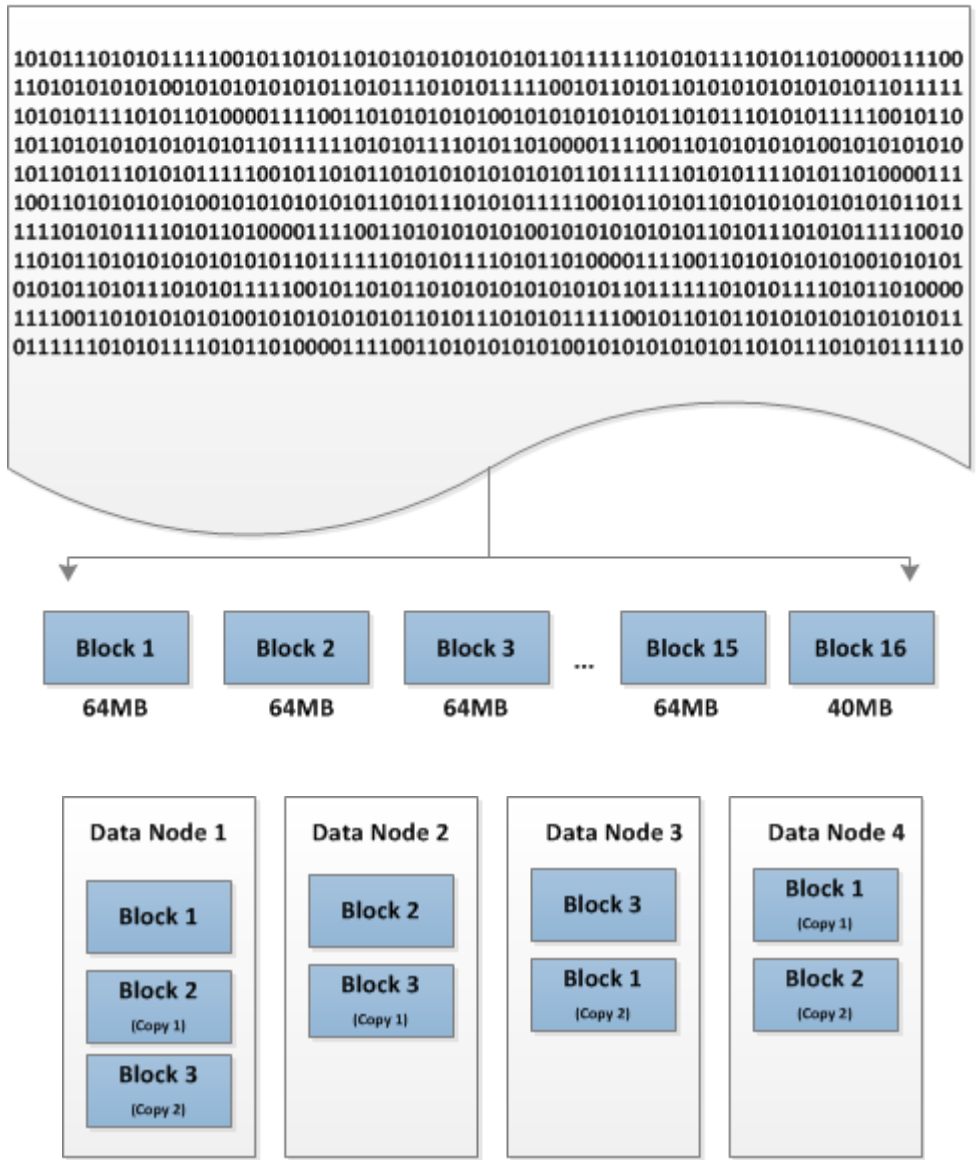
MOBILE ACCESS

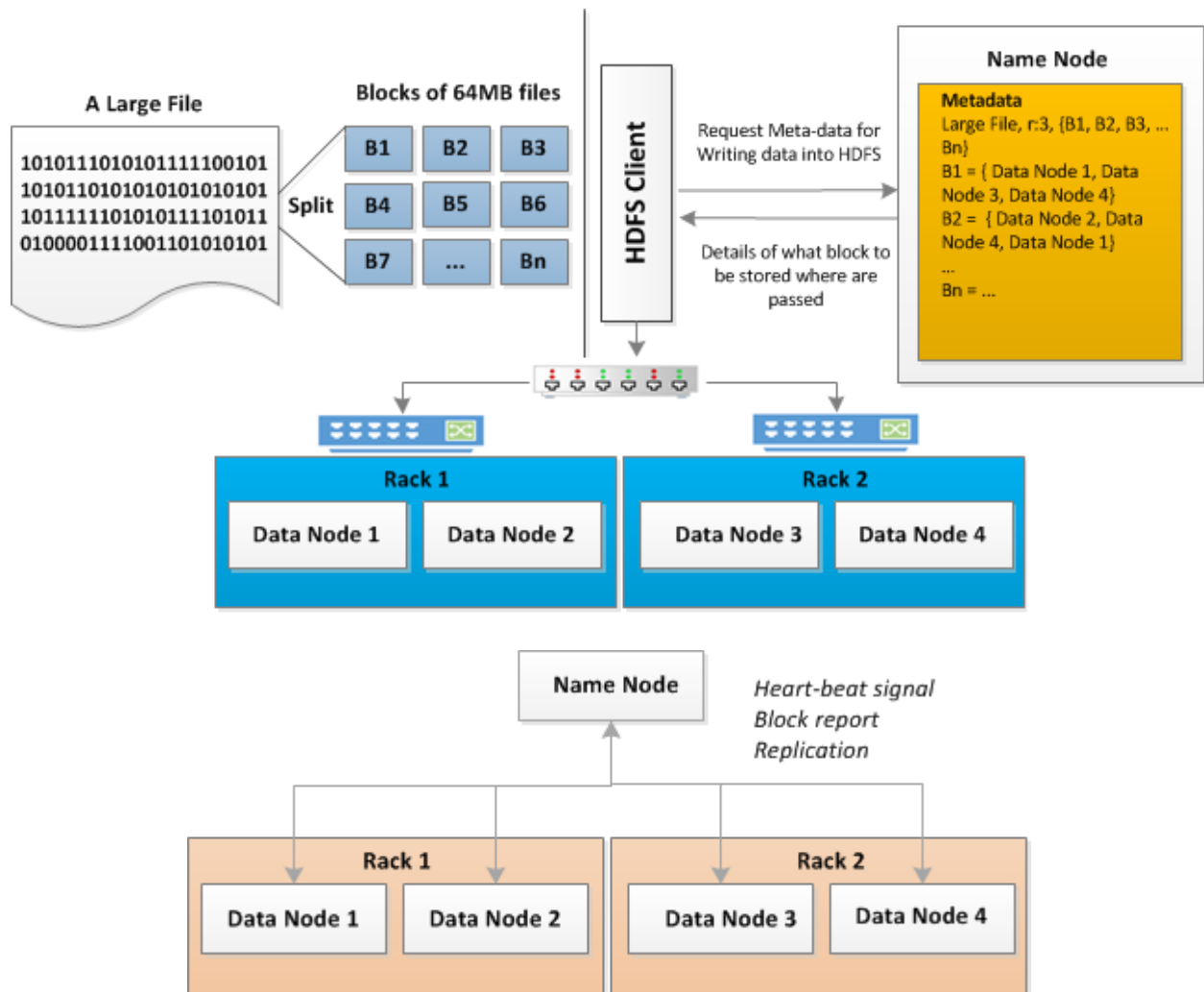
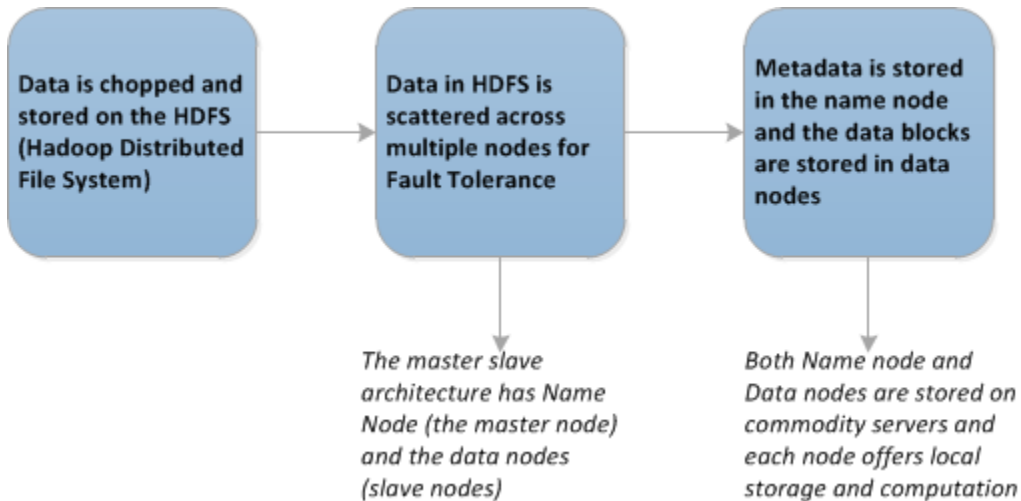
This section displays the visualization layer of the data architecture. It features four categories of tools: 'Dashboards & Reports' (showing a dashboard with bar and line charts), 'Insight Analytics & Big Data Models' (showing a complex analytics interface), 'Outlook Add-In & Excel Add-In' (showing data visualization within Microsoft Office applications), and 'BusinessObjects' (showing a traditional BI reporting interface). A green bar labeled 'MOBILE ACCESS' is positioned above the first two categories, indicating that these tools are accessible via mobile devices. Below each category are several screenshots of the respective software interfaces.





A Large File of size 1000 MB





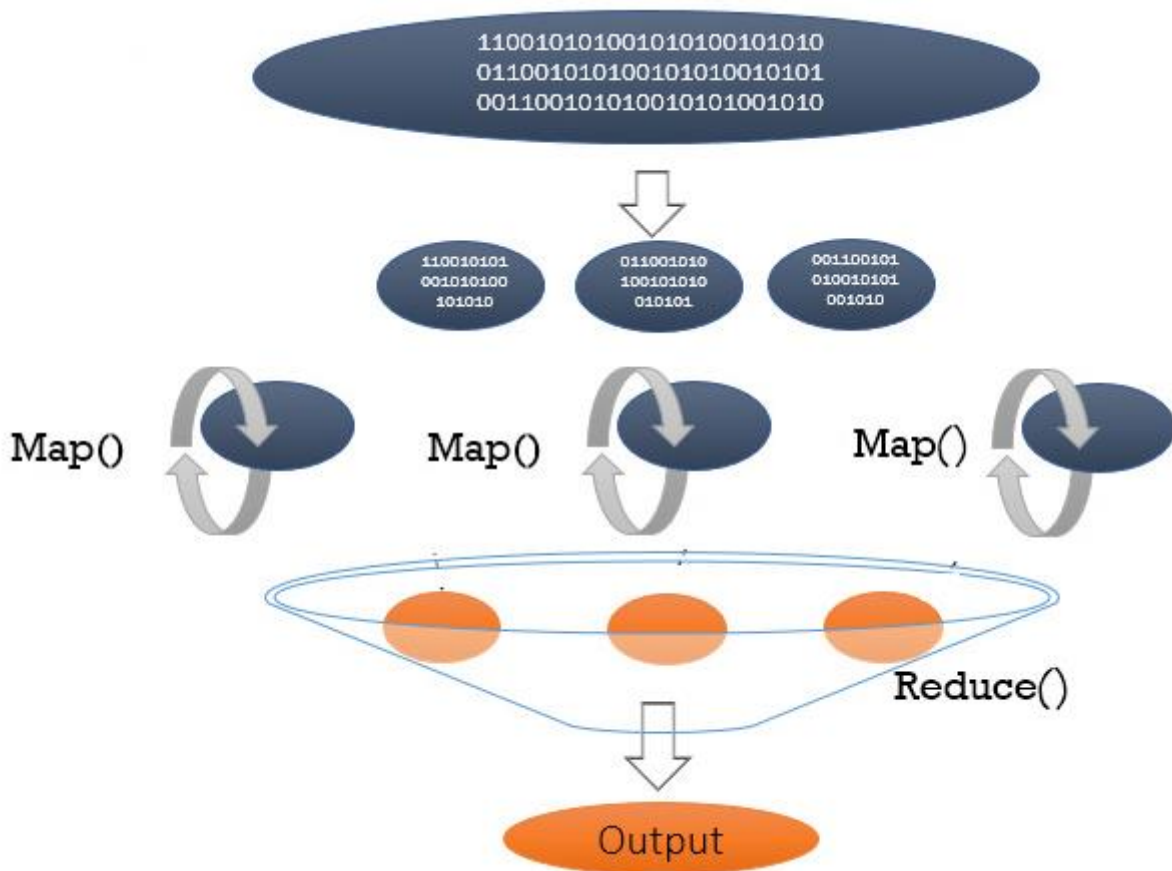
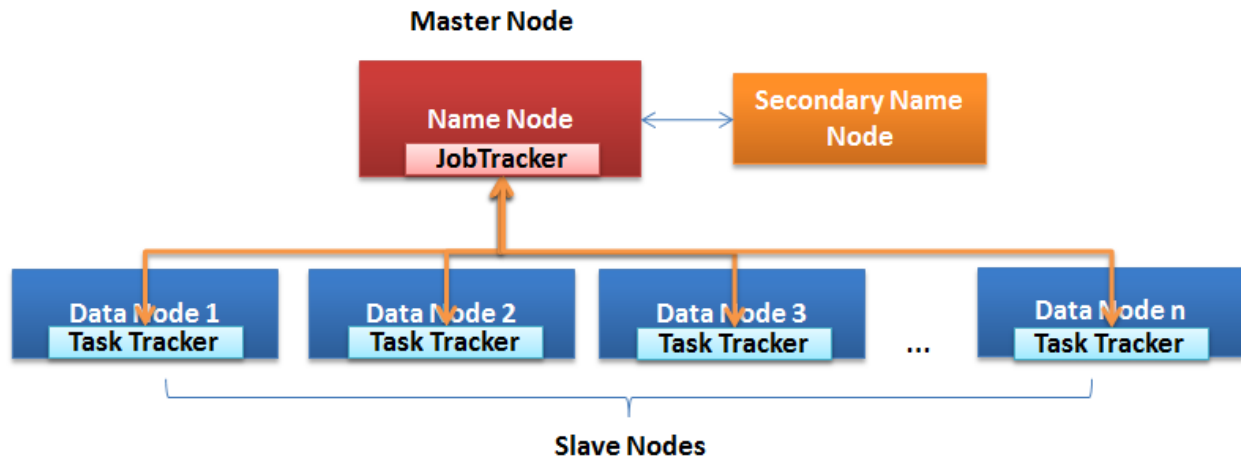
```
CA Hadoop Command Line
C:\apps\dist\hadoop-2.4.0.2.1.3.0-1948>hadoop fs
Usage: hadoop fs [generic options]
    [-appendToFile <localsrc> ... <dst>]
    [-cat [-ignoreCrc] <src> ...]
    [-checksum <src> ...]
    [-chgrp [-R] GROUP PATH...]
    [-chmod [-R] <MODE[,MODE]... : OCTALMODE> PATH...]
    [-chown [-R] [OWNER][:[GROUP]] PATH...]
    [-copyFromLocal [-f] [-p] <localsrc> ... <dst>]
    [-copyToLocal [-p] [-ignoreCrc] [-crc] <src> ... <localdst>]
    [-count [-q] <path> ...]
    [-cp [-f] [-p] <src> ... <dst>]
    [-createSnapshot <snapshotDir> [<snapshotName>]]
    [-deleteSnapshot <snapshotDir> <snapshotName>]
    [-df [-h] [<path> ...]]
    [-du [-s] [-h] <path> ...]
    [-expunge]
    [-get [-p] [-ignoreCrc] [-crc] <src> ... <localdst>]
    [-getfacl [-R] <path>]
    [-getmerge [-nl] <src> <localdst>]
    [-help [cmd ...]]
    [-ls [-d] [-h] [-R] [<path> ...]]
    [-mkdir [-p] <path> ...]
    [-moveFromLocal <localsrc> ... <dst>]
    [-moveToLocal <src> <localdst>]
    [-mv <src> ... <dst>]
    [-put [-f] [-p] <localsrc> ... <dst>]
    [-renameSnapshot <snapshotDir> <oldName> <newName>]
    [-rm [-f] [-r!-R] [-skipTrash] <src> ...]
    [-rmdir [--ignore-fail-on-non-empty] <dir> ...]
    [-setfacl [-R] [<-b!-k> <-n!-x <acl_spec>>] <path>][!--set <acl_spec> <pa
th>]]
    [-setrep [-R] [-w] <rep> <path> ...]
    [-stat [format] <path> ...]
    [-tail [-f] <file>]
    [-test [-d] [-d] [-s] <path>]
    [-text [-ignoreCrc] <src> ...]
    [-touchz <path> ...]
    [-usage [cmd ...]]

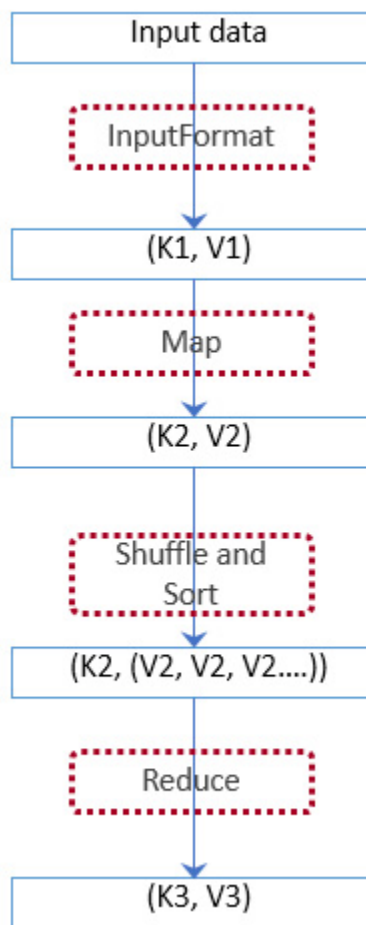
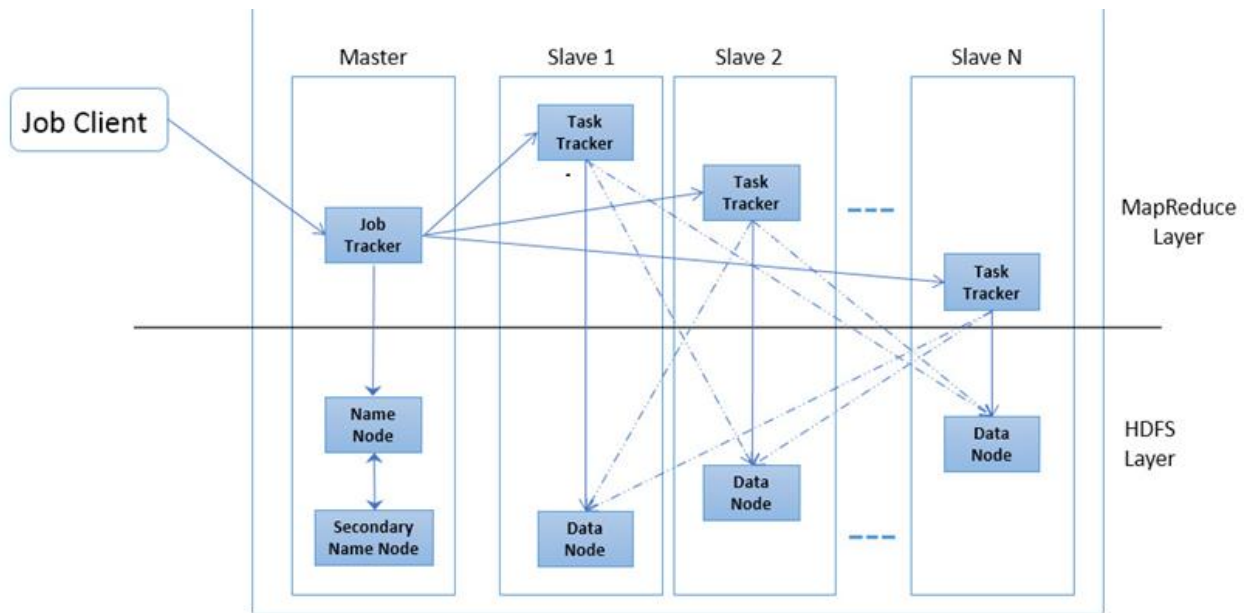
Generic options supported are
-conf <configuration file>      specify an application configuration file
-D <property=value>             use value for given property
-fs <localnamenode:port>        specify a namenode
-jt <localjobtracker:port>      specify a job tracker
-files <comma separated list of files> specify comma separated files to be co
pied to the map reduce cluster
-libjars <comma separated list of jars> specify comma separated jar files to
include in the classpath.
-archives <comma separated list of archives> specify comma separated archives
to be unarchived on the compute machines.

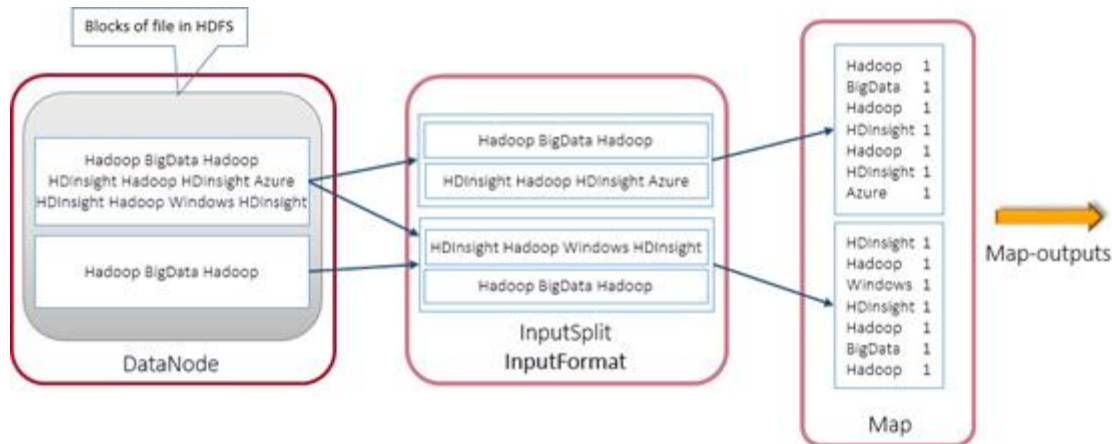
The general command line syntax is
bin/hadoop command [genericOptions] [commandOptions]

C:\apps\dist\hadoop-2.4.0.2.1.3.0-1948>_
```







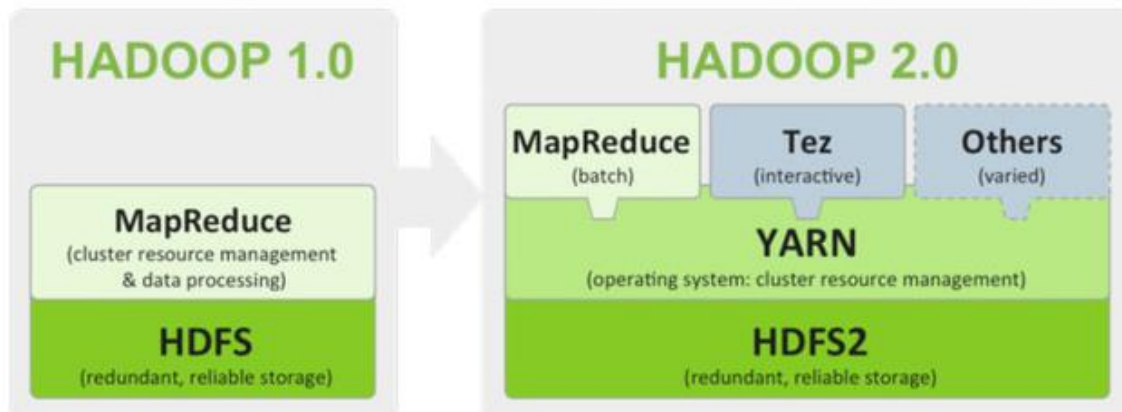


ARCHITECTURE COMPARISON

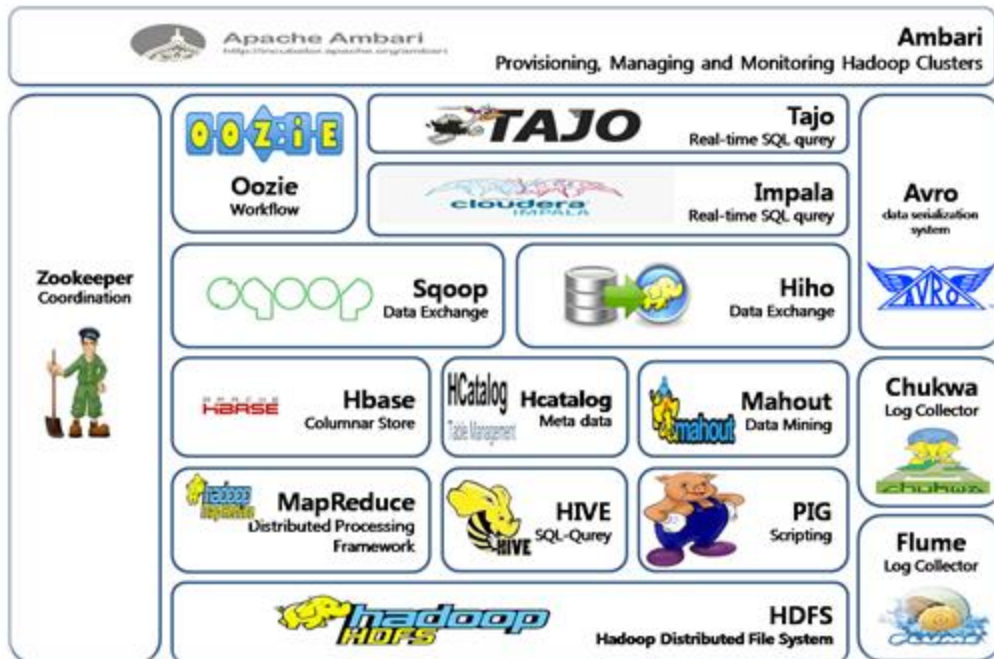
Hadoop 1.0 vs. Hadoop 2.0.

Single Use System
Batch Apps

Multi Use Data Platform
Batch, Interactive, Online, Streaming, ...



SOURCE: HORTONWORKS



Source: Internet

```

HTTP request sent, awaiting response... 200 OK
Length: 96316511 (92M) [application/x-gzip]
Saving to: `jdk-7u25-linux-x64.tar.gz'

100%[----->] 96,316,511  311K/s  in 5m 3s

2013-10-24 14:34:22 (311 KB/s) - `jdk-7u25-linux-x64.tar.gz' saved [96316511/96316511]

```

```

master@Hadoopupgrade:~$ java -version
java version "1.7.0_45"
Java(TM) SE Runtime Environment (build 1.7.0_45-b18)
Java HotSpot(TM) 64-Bit Server VM (build 24.45-b08, mixed mode)
master@Hadoopupgrade:~$ █

```



```
master@Hadoopupgrade:~$ sudo addgroup hadoop
Adding group 'hadoop' (GID 1001) ...
Done.
master@Hadoopupgrade:~$ sudo adduser --ingroup hadoop hduser
Adding user 'hduser' ...
Adding new user 'hduser' (1001) with group 'hadoop' ...
Creating home directory '/home/hduser' ...
Copying files from '/etc/skel' ...
Enter new UNIX password:
Retype new UNIX password:
passwd: password updated successfully
Changing the user information for hduser
Enter the new value, or press ENTER for the default
    Full Name []:
    Room Number []:
    Work Phone []:
    Home Phone []:
    Other []:
Is the information correct? [Y/n] Y
master@Hadoopupgrade:~$ █
```

```
<configuration>
<!-- Site specific YARN configuration properties -->
<property>
  <name>yarn.nodemanager.aux-services</name>
  <value>mapreduce_shuffle</value>
</property>
  <property>
    <name>yarn.nodemanager.aux-services.mapreduce.shuffle.class</name>
    <value>org.apache.hadoop.mapred.ShuffleHandler</value>
  </property>
</configuration>
```

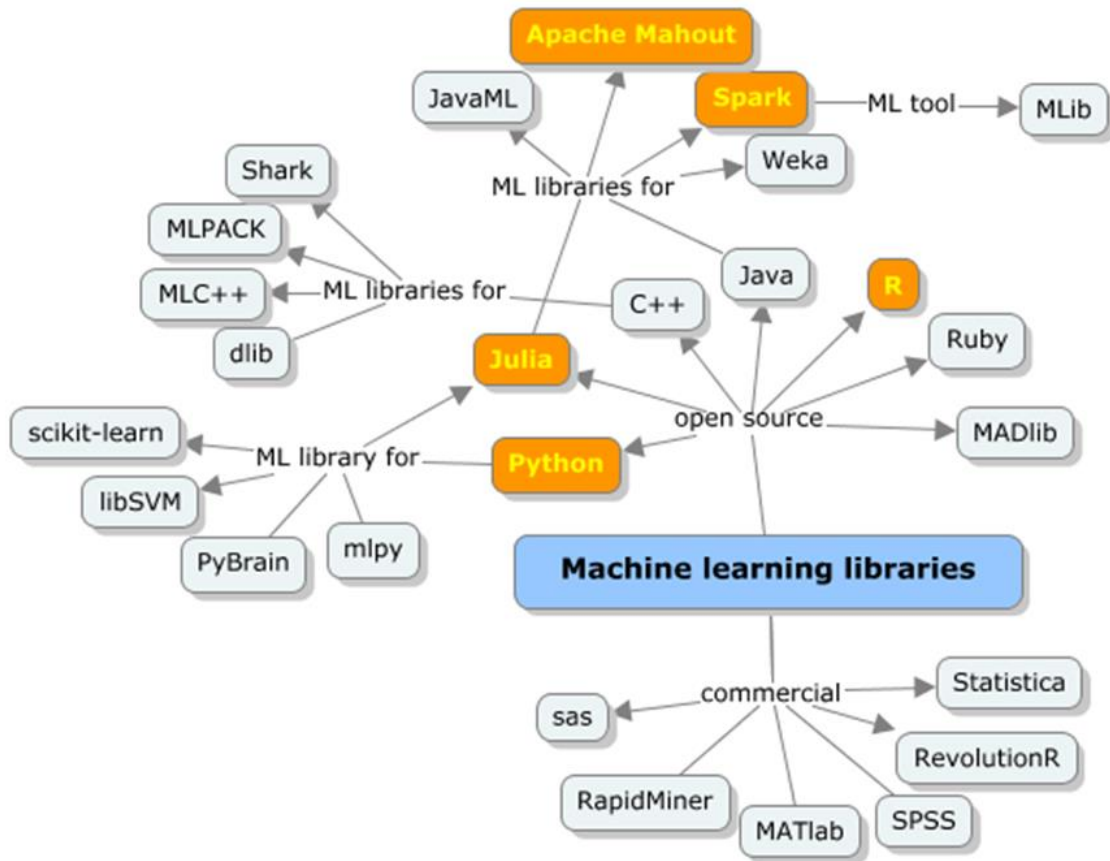
```
<configuration>
<property>
  <name>fs.default.name</name>
  <value>hdfs://localhost:9000</value>
</property>█
</configuration>
```

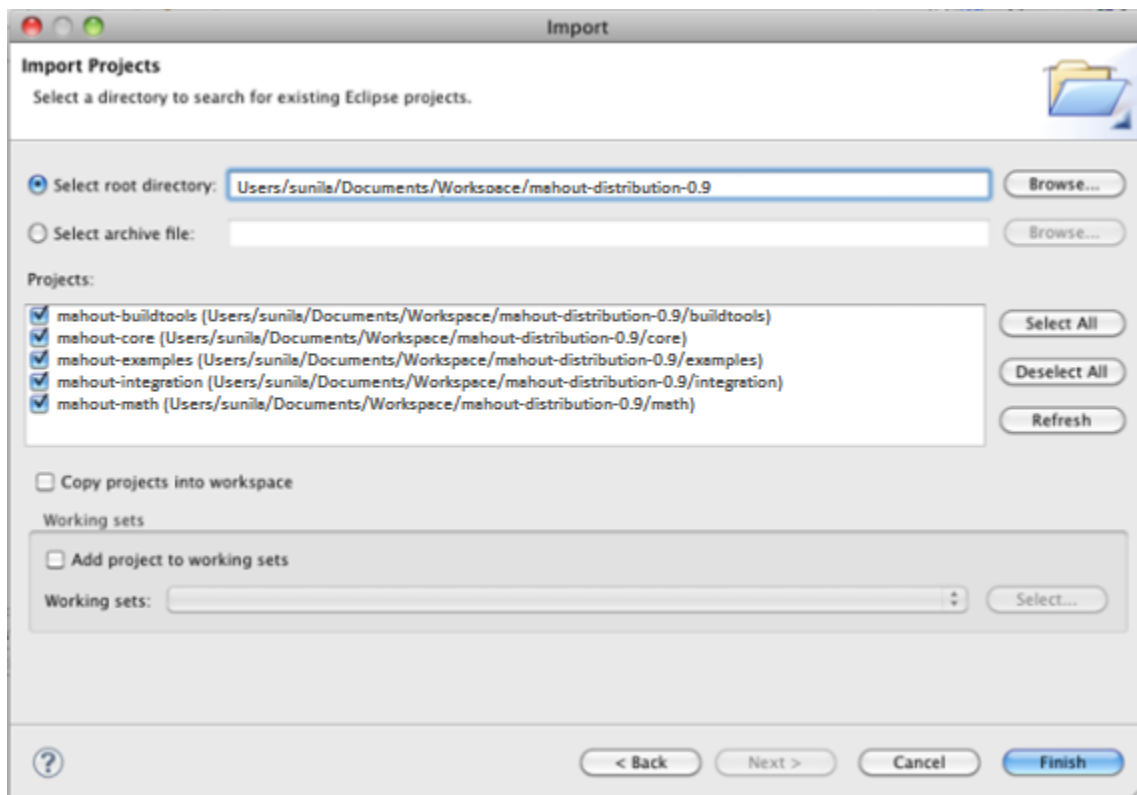
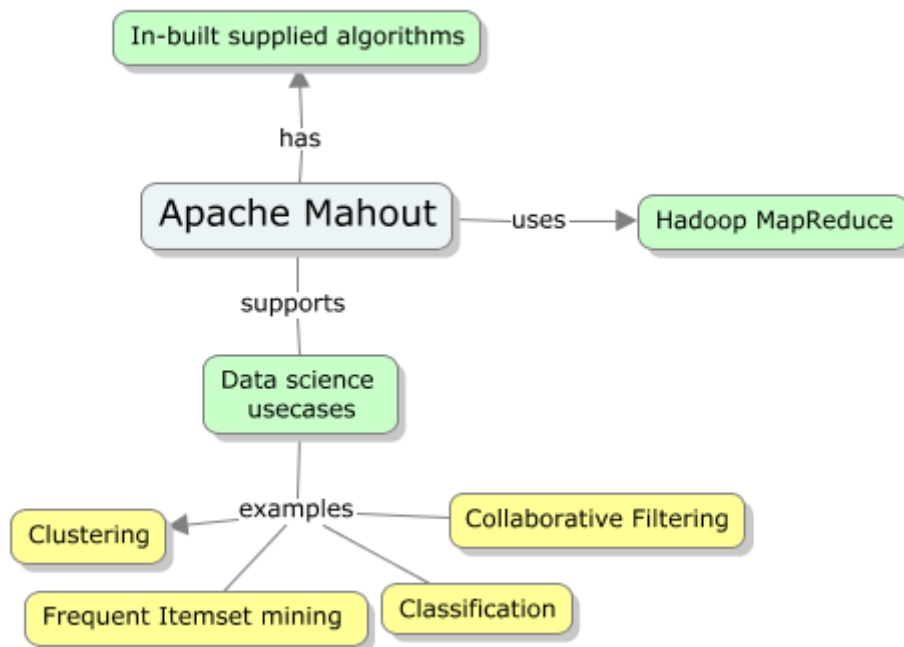
```
<configuration>
<property>
  <name>dfs.replication</name>
  <value>1</value>
</property>
  <property>
    <name>dfs.namenode.name.dir</name>
    <value>file:/usr/local/hadoop/yarn_data/hdfs/namenode</value>
  </property>
<property>
  <name>dfs.datanode.data.dir</name>
  <value>file:/usr/local/hadoop/yarn_data/hdfs/datanode</value>
</property>█
</configuration>
```

```
<configuration>
<property>
  <name>mapreduce.framework.name</name>
  <value>yarn</value>
</property>
</configuration>
```

```
! Set Hadoop-related environment variables
export HADOOP_PREFIX='/usr/local/hadoop'
export HADOOP_HOME='/usr/local/hadoop'
export HADOOP_MAPRED_HOME=${HADOOP_HOME}
export HADOOP_COMMON_HOME=${HADOOP_HOME}
export HADOOP_HDFS_HOME=${HADOOP_HOME}
export YARN_HOME=${HADOOP_HOME}
export HADOOP_CONF_DIR=${HADOOP_HOME}/etc/hadoop
! Native Path
export HADOOP_COMMON_LIB_NATIVE_DIR=${HADOOP_PREFIX}/lib/native
export HADOOP_OPTS="-Djava.library.path=${HADOOP_PREFIX}/lib"
! Java path
export JAVA_HOME='/usr/local/Java/jdk1.7.0_45'
! Add Hadoop bin/ directory to PATH
export PATH=${PATH}:${HADOOP_HOME}/bin:${JAVA_HOME}/bin:${HADOOP_HOME}/sbin
```

```
hduser@Hadoopupgrade:~$ hadoop-daemon.sh start namenode
starting namenode, logging to /usr/local/hadoop/logs/hadoop-hduser-namenode-Hadoopupgrade.out
hduser@Hadoopupgrade:~$ jps
1244 NameNode
1280 Jps
hduser@Hadoopupgrade:~$ hadoop-daemon.sh start datanode
starting datanode, logging to /usr/local/hadoop/logs/hadoop-hduser-datanode-Hadoopupgrade.out
hduser@Hadoopupgrade:~$ jps
1400 Jps
1244 NameNode
1332 DataNode
hduser@Hadoopupgrade:~$ yarn-daemon.sh start resourcemanager
starting resourcemanager, logging to /usr/local/hadoop/logs/yarn-hduser-resource-manager-Hadoopupgrade.out
hduser@Hadoopupgrade:~$ jps
1474 Jps
1244 NameNode
1433 ResourceManager
1332 DataNode
```





```

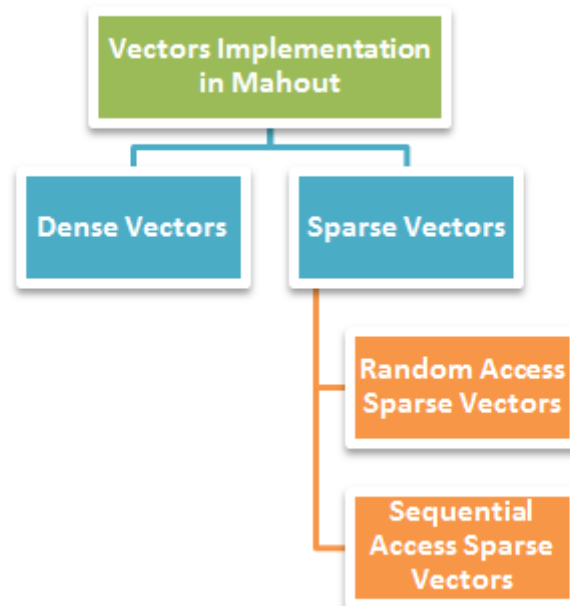
File Edit View Terminal Help
[INFO] --- maven-source-plugin:2.2.1:jar-no-fork (attach-sources) @ mahout-examples ---
[INFO]
[INFO] --- maven-install-plugin:2.5:install (default-install) @ mahout-examples ---
[INFO] Installing /home/abdulbasit/Softwares/Apache Mahout/examples/target/mahout-examples-0.9-SNAPSHOT.jar to /root/.m2/repository/org/apache/mahout/mahout-examples/0.9-SNAPSHOT/mahout-examples-0.9-SNAPSHOT.jar
[INFO] Installing /home/abdulbasit/Softwares/Apache Mahout/examples/pom.xml to /root/.m2/repository/org/apache/mahout/mahout-examples/0.9-SNAPSHOT/mahout-examples-0.9-SNAPSHOT.pom
[INFO] Installing /home/abdulbasit/Softwares/Apache Mahout/examples/target/mahout-examples-0.9-SNAPSHOT-job.jar to /root/.m2/repository/org/apache/mahout/mahout-examples/0.9-SNAPSHOT/mahout-examples-0.9-SNAPSHOT-job.jar
[INFO] Installing /home/abdulbasit/Softwares/Apache Mahout/examples/target/mahout-examples-0.9-SNAPSHOT-sources.jar to /root/.m2/repository/org/apache/mahout/mahout-examples/0.9-SNAPSHOT/mahout-examples-0.9-SNAPSHOT-sources.jar
[INFO]
[INFO] -----
[INFO] Building Mahout Release Package 0.9-SNAPSHOT
[INFO] -----
[INFO]
[INFO] --- maven-assembly-plugin:2.4:single (bin-assembly) @ mahout-distribution ---
[INFO] Assemblies have been skipped per configuration of the skipAssembly parameter.
[INFO]
[INFO] --- maven-assembly-plugin:2.4:single (src-assembly) @ mahout-distribution ---
[INFO] Assemblies have been skipped per configuration of the skipAssembly parameter.
[INFO]
[INFO] --- maven-install-plugin:2.5:install (default-install) @ mahout-distribution ---
[INFO] Installing /home/abdulbasit/Softwares/Apache Mahout/distribution/pom.xml to /root/.m2/repository/org/apache/mahout/mahout-distribution/0.9-SNAPSHOT/mahout-distribution-0.9-SNAPSHOT.pom
[INFO] -----
[INFO] Reactor Summary:
[INFO]
[INFO] Mahout Build Tools ..... SUCCESS [3.332s]
[INFO] Apache Mahout ..... SUCCESS [0.635s]
[INFO] Mahout Math ..... SUCCESS [1:31.616s]
[INFO] Mahout Core ..... SUCCESS [10:58.301s]
[INFO] Mahout Integration ..... SUCCESS [59.877s]
[INFO] Mahout Examples ..... SUCCESS [19.040s]
[INFO] Mahout Release Package ..... SUCCESS [0.060s]
[INFO] -----
[INFO] BUILD SUCCESS
[INFO] -----
[INFO] Total time: 13:53.340s
[INFO] Final Memory: 34M/335M
[INFO] -----
root@it:~/Softwares/Apache Mahout#

```

Applications



Evolutionary Algorithms	Classification	Clustering	Recommenders	Regression	FPM	Dimension Reduction
Utilities Lucene/Vectorizer		Math Vectors/ Matrices/SVD		Collections (Primitives)		Hadoop



```
RGui (32-bit)
File Edit View Misc Packages Windows Help

R Console
ISBN 3-900051-07-0
Platform: i386-pc-mingw32/i386 (32-bit)

R is free software and comes with ABSOLUTELY NO WARRANTY.
You are welcome to redistribute it under certain conditions.
Type 'license()' or 'licence()' for distribution details.

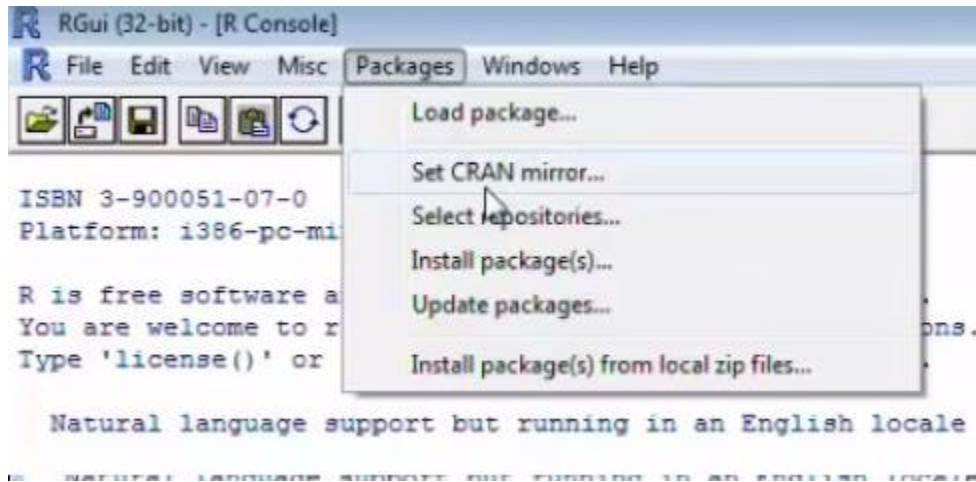
Natural language support but running in an English locale

R is a collaborative project with many contributors.
Type 'contributors()' for more information and
'citation()' on how to cite R or R packages in publications.

Type 'demo()' for some demos, 'help()' for on-line help, or
'help.start()' for an HTML browser interface to help.
Type 'q()' to quit R.

Loading required package: rcom
Loading required package: rscproxy
Warning messages:
1: package 'rcom' was built under R version 2.15.2
2: package 'rscproxy' was built under R version 2.15.2
[Previously saved workspace restored]

> |
```



R Console

```
Warning messages:
1: package 'rcom' was built under R version 2.15.2
2: package 'rscproxy' was built under R version 2.15.2
[Previously saved workspace restored]

> chooseCRANmirror()
> 1+2
[1] 3
> log(10)
[1] 2.302585
> 1+2
[1] 3
> a = 2
> a
[1] 2
> a <-2
> a
[1] 2
> a <-10
> a
[1] 10
> b <-5
> b
[1] 5
> |
```

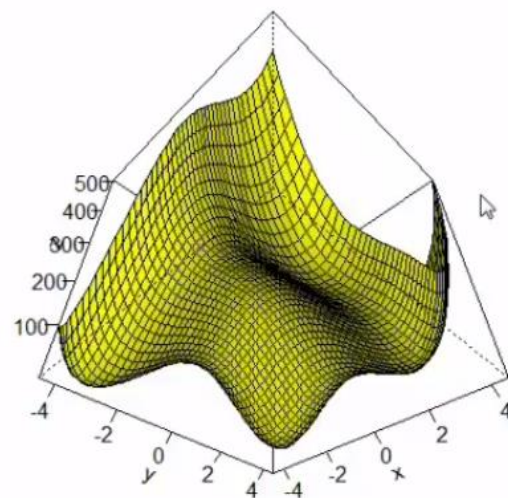
Untitled - R Editor

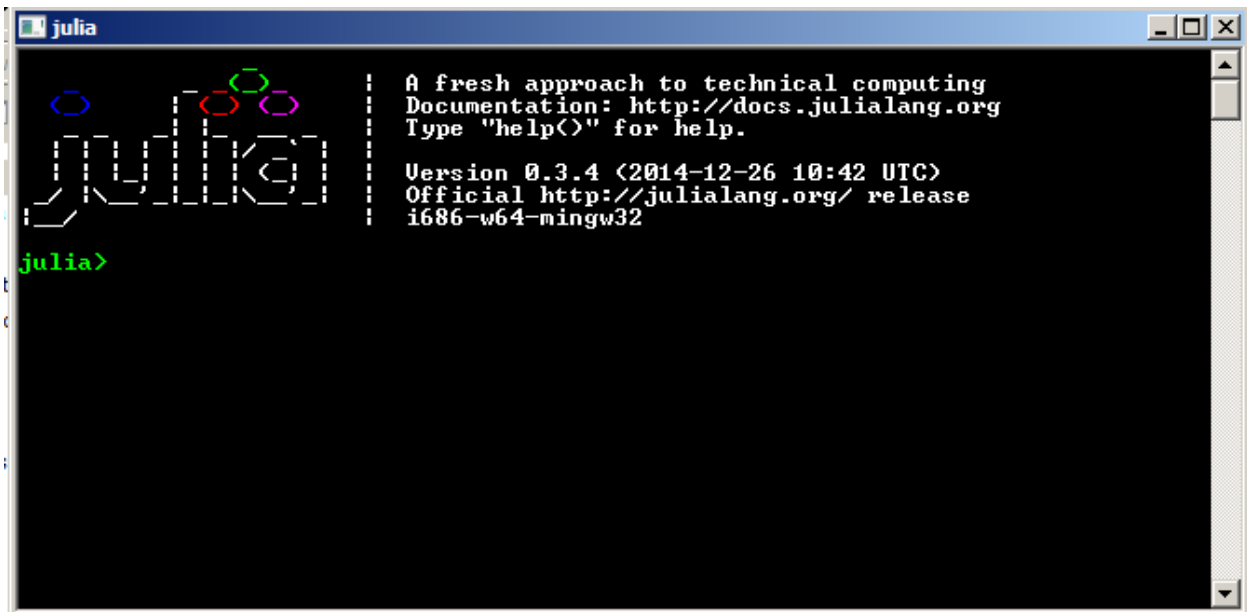
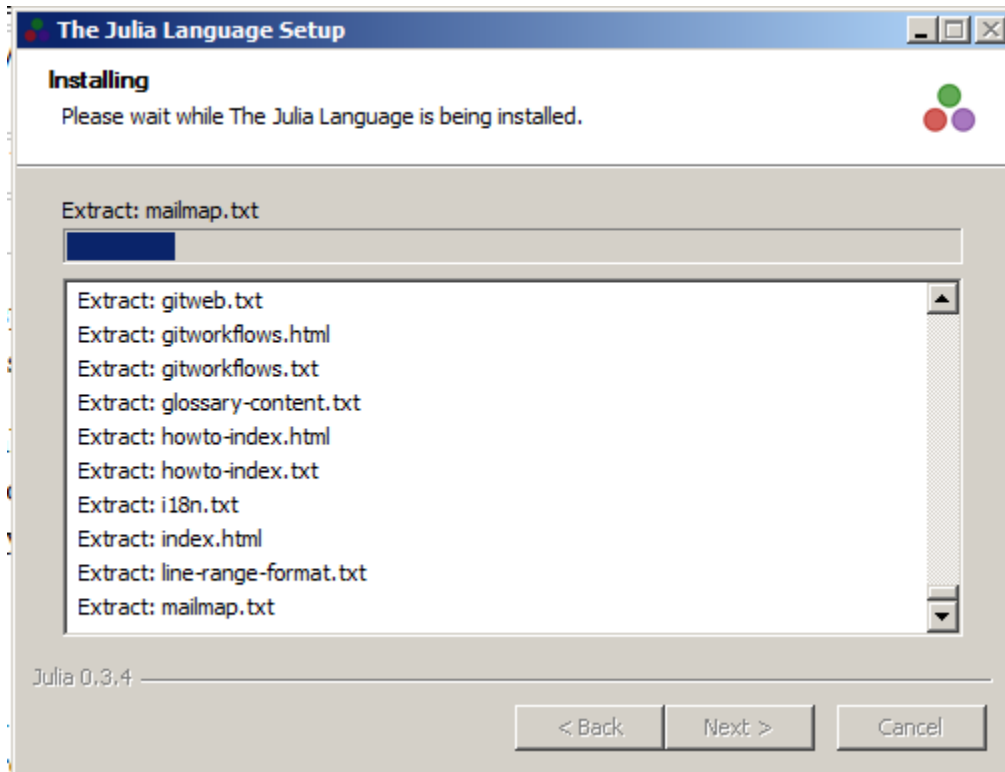
```
1+2
a = 2
a <-2
a <-10
b <-5
c<-a+|
```

```
f <- function(x1,y1) (1-x1)^2 + 100*(y1 - x1^2)^2
optim( c(0,0), f )$par
ror in y1 - x1^2 : 'y1' is missing

f <- function(x) (1-x[1])^2 + 100*(x[2]-x[1]^2)^2
optim( c(0,0), f )$par
] 0.9999564 0.9999085

f <- function(x1,y1) (x1^2 + y1 - 11)^2 + (x1^2 - y1)^2
x <- seq(-4.5,4.5,by=.2)
y <- seq(-4.5,4.5,by=.2)
z <- outer(x,y,f)
persp(x,y,z,phi=-45,theta=45,col="yellow",sha
```





mandelbrot.jl*

```
1 function mandel(z)
2     c = z
3     maxiter = 80
4     for n = 1:maxiter
5         abs(z) ≥ 2 && return n-1
6         z = z^2 + c
7     end
8     return max|
9 end
10
11 mandel(0) 80
12
13 mandel(x, y)
14
15 function mand
```

maxad
maxiter
maxabs!
maximum
maximum! + y*im)
maxintfloat
max_text_extents

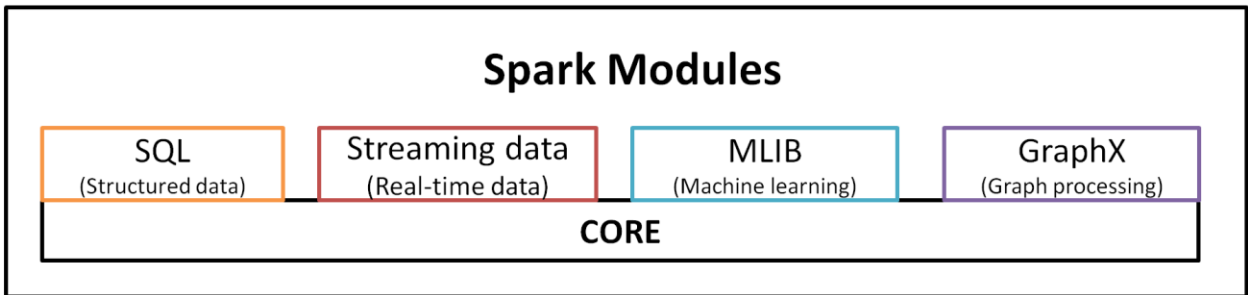
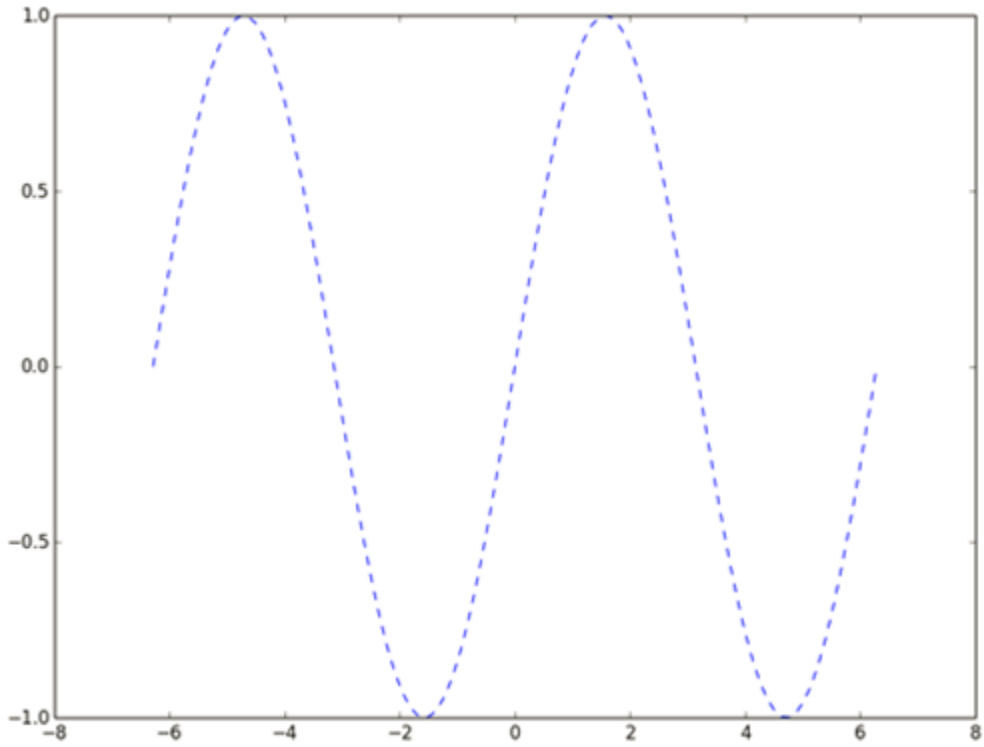
← → ↻ 🏠 <https://juliabox.org>



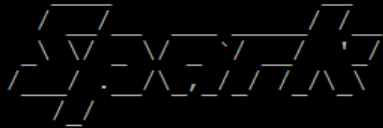
JuliaBox ^{beta}

Run Julia from the Browser. No setup.

The Julia community is doing amazing things.
We want you in on it!

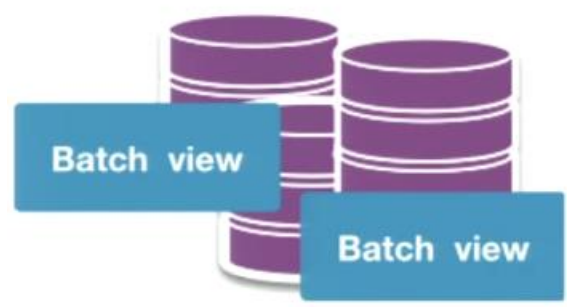
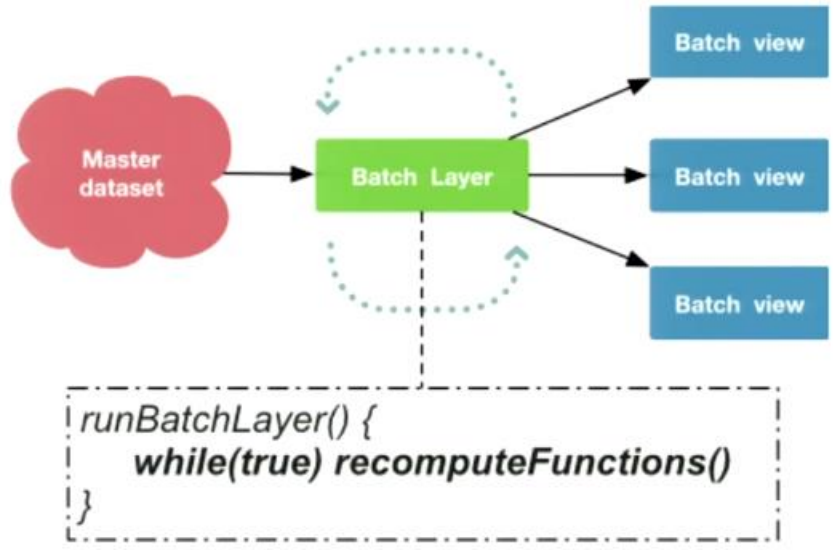


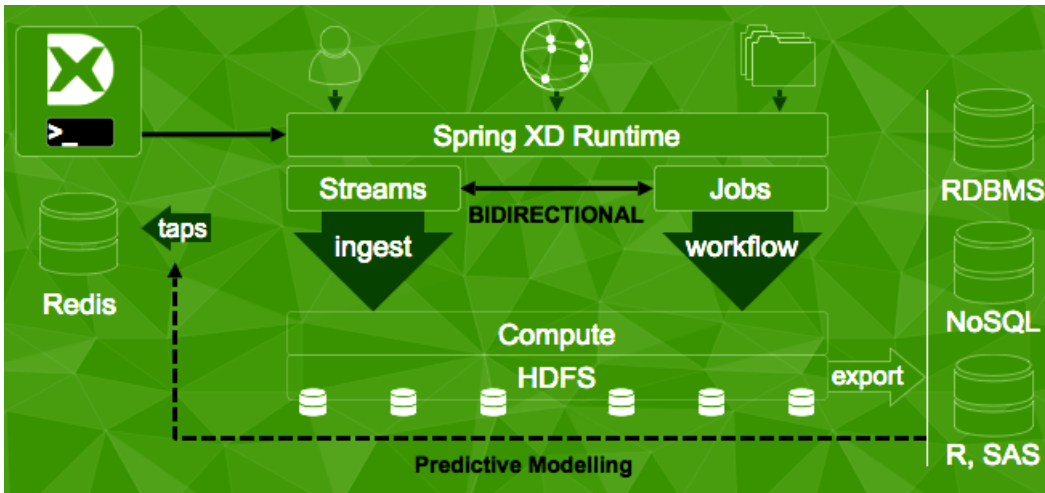
```
ubuntu@ip-172-31-21-139: ~/spark-1.2.0-bin-hadoop1
ubuntu@ip-172-31-21-139:~/spark-1.2.0-bin-hadoop1$ cd spark-1.2.0-bin-hadoop1/
ubuntu@ip-172-31-21-139:~/spark-1.2.0-bin-hadoop1$ ./bin/spark-shell
Spark assembly has been built with Hive, including Datanucleus jars on classpath
Using Spark's default log4j profile: org/apache/spark/log4j-defaults.properties
15/01/28 16:31:34 INFO SecurityManager: Changing view acls to: ubuntu
15/01/28 16:31:34 INFO SecurityManager: Changing modify acls to: ubuntu
15/01/28 16:31:34 INFO SecurityManager: SecurityManager: authentication disabled
; ui acls disabled; users with view permissions: Set(ubuntu); users with modify
permissions: Set(ubuntu)
15/01/28 16:31:34 INFO HttpServer: Starting HTTP Server
15/01/28 16:31:34 INFO Utils: Successfully started service 'HTTP class server' o
n port 59689.
Welcome to

 version 1.2.0

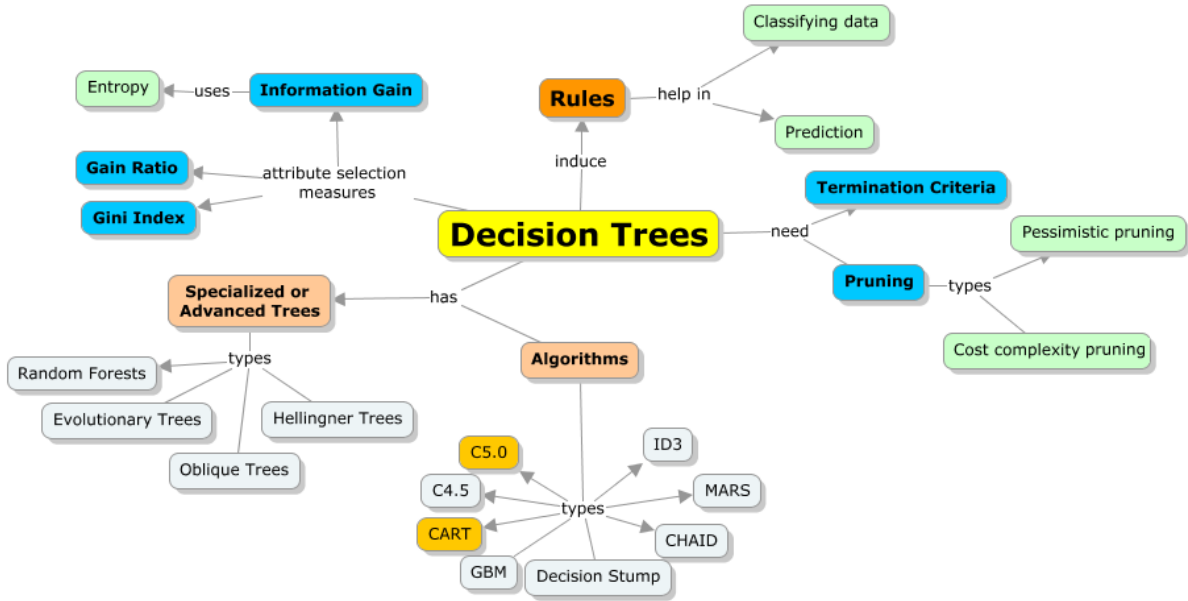
Using Scala version 2.10.4 (OpenJDK 64-Bit Server VM, Java 1.7.0_75)
Type in expressions to have them evaluated.
Type :help for more information.
15/01/28 16:31:41 INFO SecurityManager: Changing view acls to: ubuntu
15/01/28 16:31:41 INFO SecurityManager: Changing modify acls to: ubuntu
15/01/28 16:31:42 INFO MemoryStore: MemoryStore started with capacity 267.3 MB
15/01/28 16:31:42 INFO HttpFileServer: HTTP File server directory is /tmp/spark-
02c7a588-7925-444a-a5de-80513134293a
15/01/28 16:31:42 INFO HttpServer: Starting HTTP Server
15/01/28 16:31:42 INFO Utils: Successfully started service 'HTTP file server' on
port 60397.
15/01/28 16:31:42 INFO Utils: Successfully started service 'SparkUI' on port 404
0.
15/01/28 16:31:42 INFO SparkUI: Started SparkUI at http://ip-172-31-21-139.us-we
st-2.compute.internal:4040
15/01/28 16:31:43 INFO Executor: Using REPL class URI: http://172.31.21.139:5968
9
15/01/28 16:31:43 INFO AkkaUtils: Connecting to HeartbeatReceiver: akka.tcp://sp
arkDriver@ip-172-31-21-139.us-west-2.compute.internal:33030/user/HeartbeatReceiv
er
15/01/28 16:31:43 INFO NettyBlockTransferService: Server created on 44185
15/01/28 16:31:43 INFO BlockManagerMaster: Trying to register BlockManager
15/01/28 16:31:43 INFO BlockManagerMasterActor: Registering block manager localh
ost:44185 with 267.3 MB RAM, BlockManagerId(<driver>, localhost, 44185)
15/01/28 16:31:43 INFO BlockManagerMaster: Registered BlockManager
15/01/28 16:31:43 INFO SparkILoop: Created spark context..
Spark context available as sc.

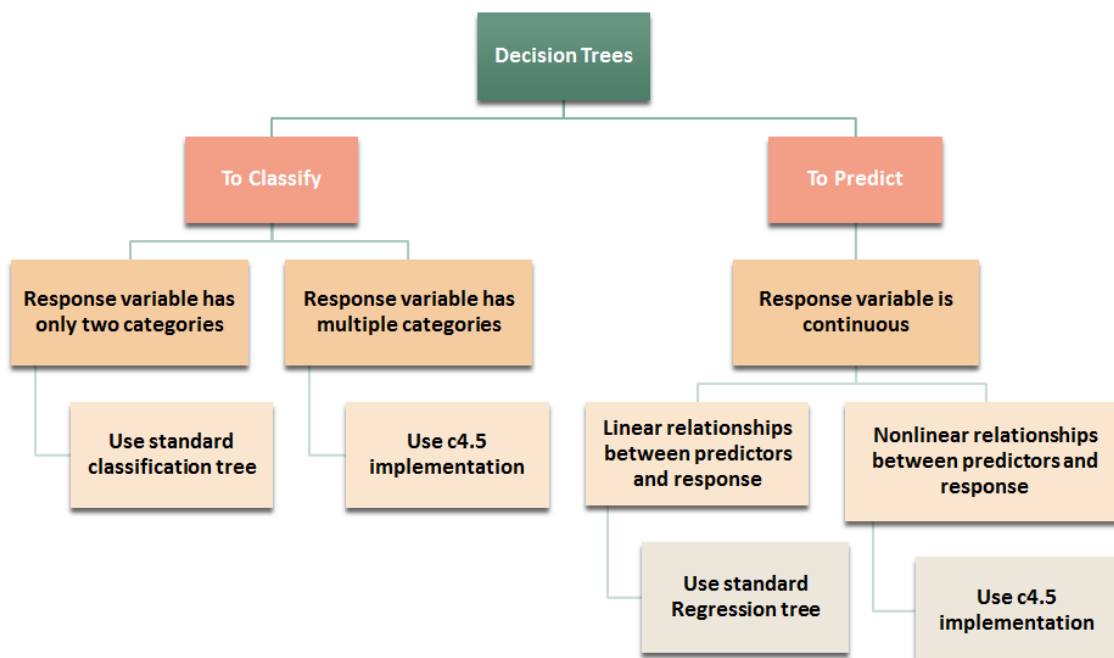
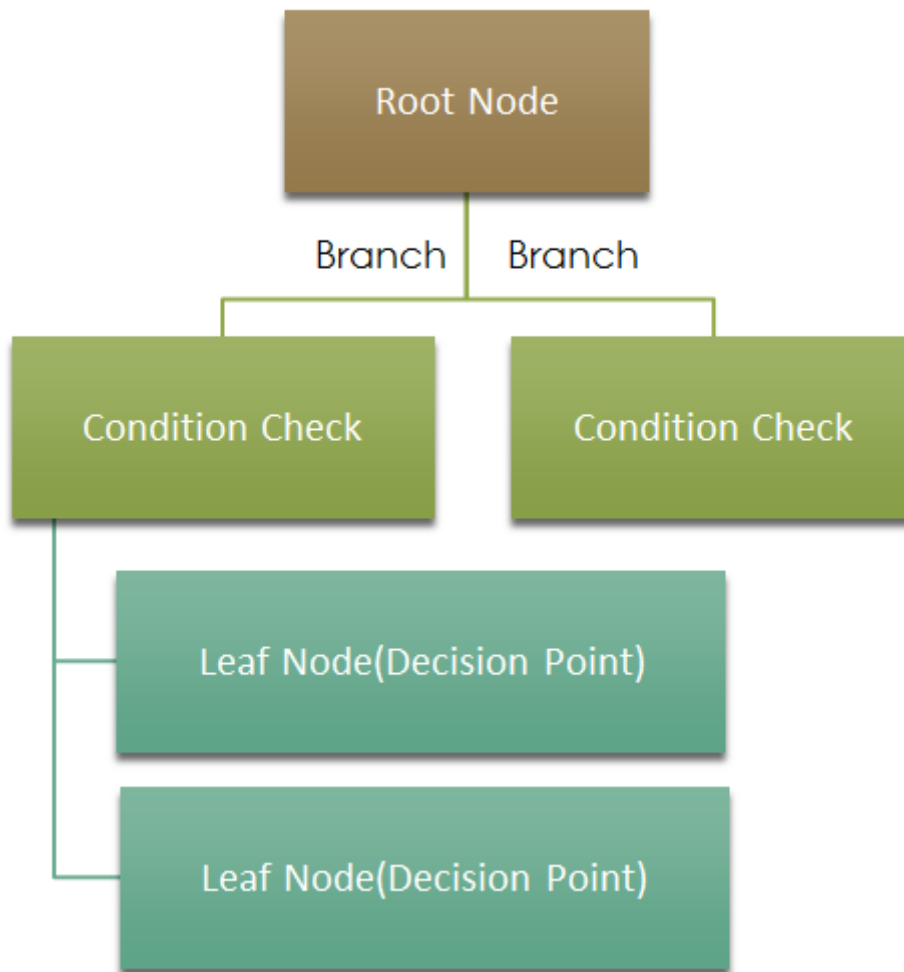
scala> █
```





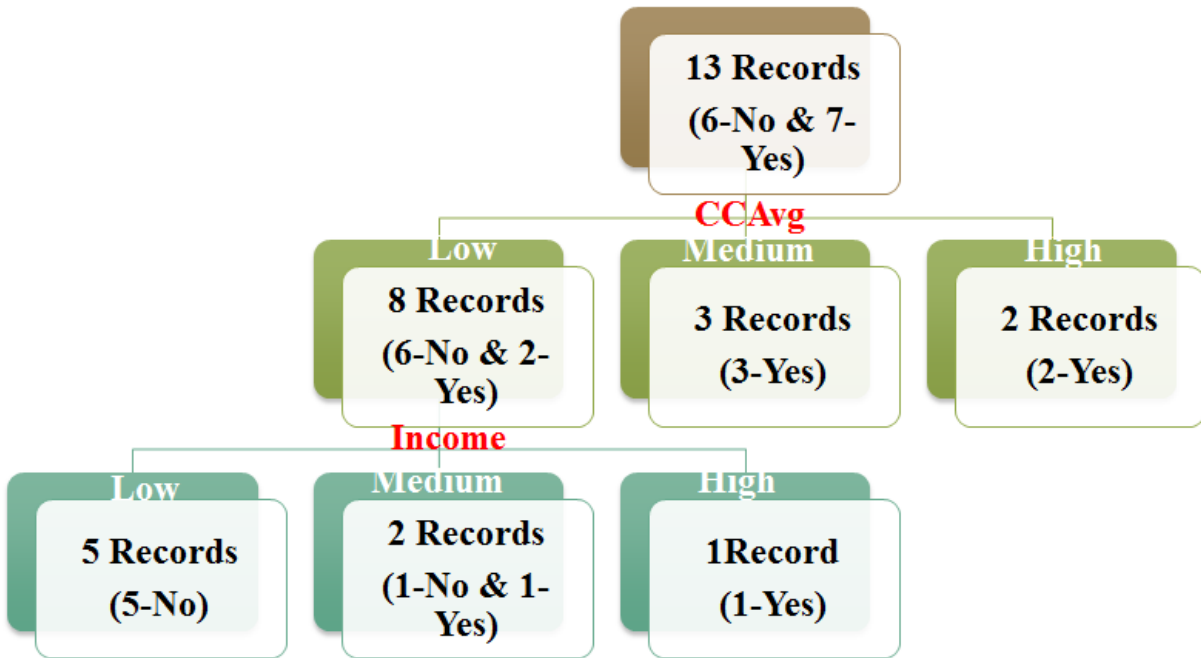
CHAPTER 5

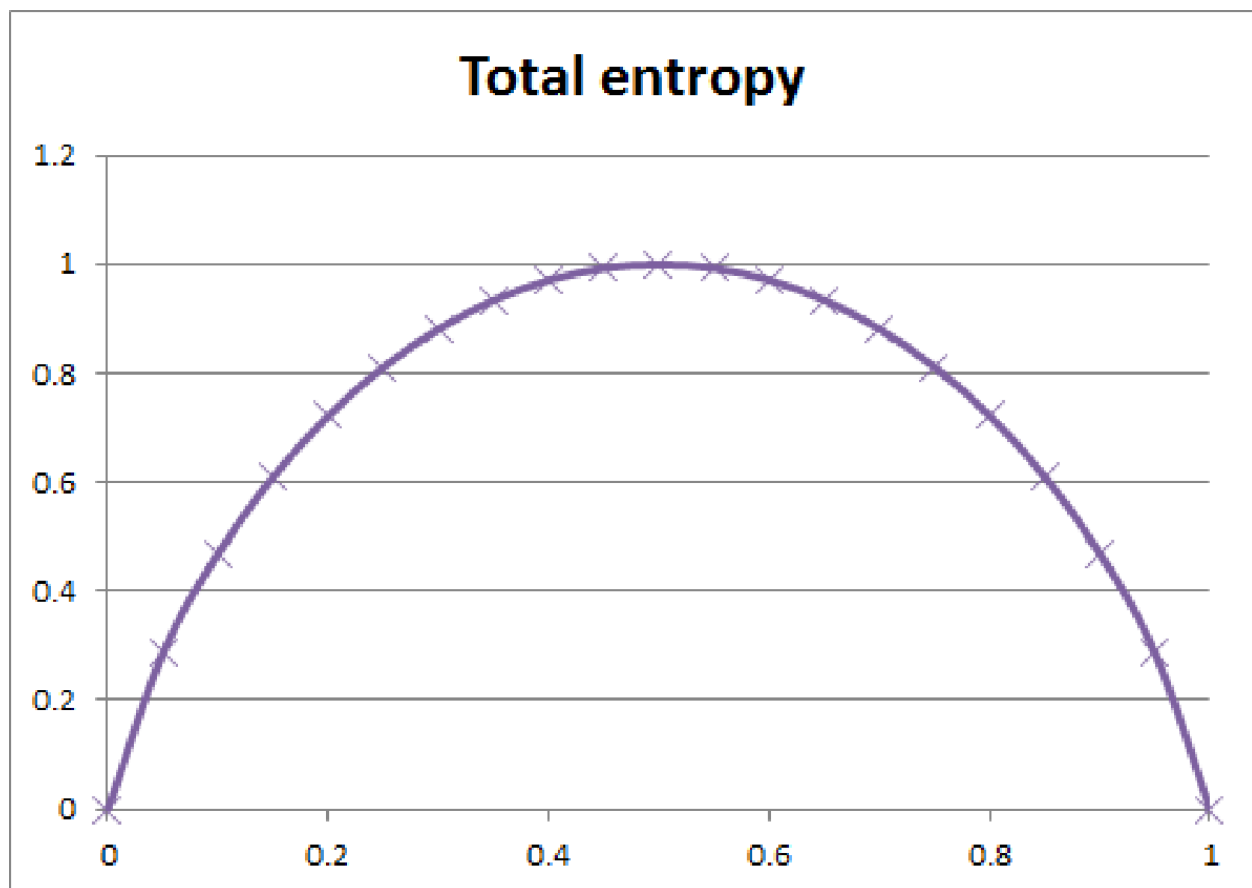
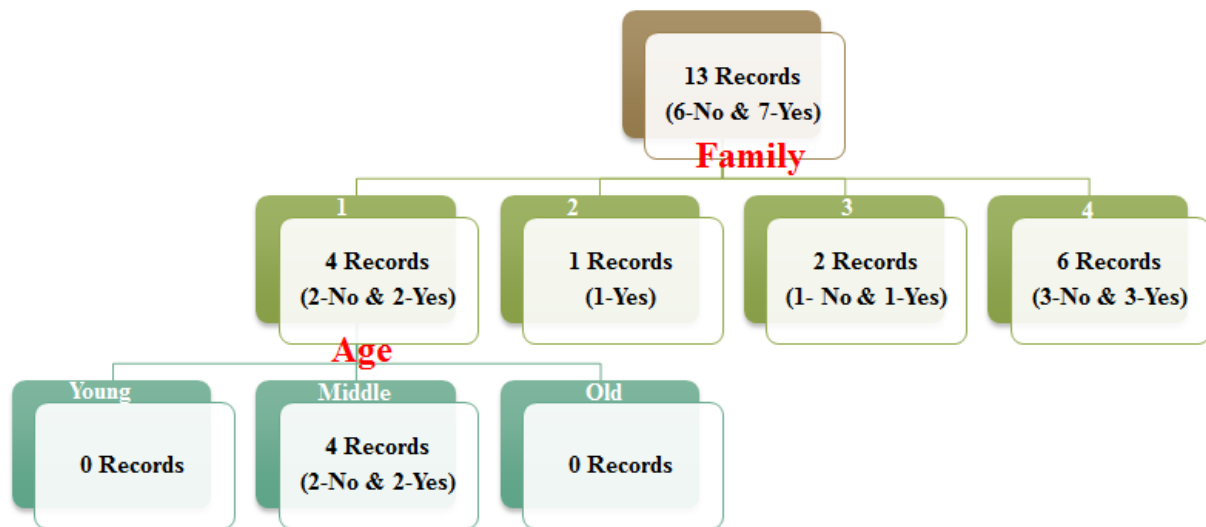




Source: <http://www.simafore.com/blog/bid/62482/2-main-differences-between-classification-and-regression-trees>

ID	Age	Experience	Income	Family	CCAvg	Personal Loan
1	25	1	49	4	1.60	0
2	45	19	34	3	1.50	0
3	39	15	11	1	1.00	0
4	35	9	100	1	2.70	0
5	35	8	45	4	1.00	0
6	37	13	29	4	0.40	0
10	34	9	180	1	8.90	1
17	38	14	130	4	4.70	1
19	46	21	193	2	8.10	1
30	38	13	119	1	3.30	1
39	42	18	141	3	5.00	1
43	32	7	132	4	1.10	1
48	37	12	194	4	0.20	1





$$E = - \sum_{i=1}^m p_i \log_2(p_i)$$

$$E_A = \sum_{i=1}^v \frac{D_i}{D} E(D_i)$$

ID	Age	Income	Family	CCAvg	Personal Loan
1	Young	Low	4	Low	0
2	Old	Low	3	Low	0
3	Middle	Low	1	Low	0
4	Middle	Medium	1	Low	0
5	Middle	Low	4	Low	0
6	Middle	Low	4	Low	0
10	Middle	High	1	High	1
17	Middle	Medium	4	Medium	1
19	Old	High	2	High	1
30	Middle	Medium	1	Medium	1
39	Old	Medium	3	Medium	1
43	Young	Medium	4	Low	1
48	Middle	High	4	Low	1

$$-\frac{6}{13} \log_2\left(\frac{6}{13}\right) - \frac{7}{13} \log_2\left(\frac{7}{13}\right) = 0.995727$$

		CCAvg		
	Fraction	Loan=No	Loan=Yes	Entropy
Low	0.615385	0.25	0.75	0.81
Medium	0.230769	0.00	1.00	0.00
High	0.153846	0.00	1.00	0.00
Entropy	0.499248			

$$Entropy_{CCAvg} = \frac{8}{13}E(6,2) + \frac{3}{13}E(0,3) + \frac{2}{13}E(0,2) = 0.499248$$

$$I_{CCAvg} = 0.995727 - 0.499248 = 0.496479$$

$$I_{Family} = 0.995727 - 0.923077 = 0.07265$$

$$\left(\frac{6}{13}\right)^2 + \left(\frac{7}{13}\right)^2$$

	CCAvg			
	Fraction	Loan=No	Loan=Yes	Gini
Low	0.615385	0.250000	0.750000	0.625
Medium	0.230769	0.000000	1.000000	1.000
High	0.153846	0.000000	1.000000	1.000
Gini Index			0.769231	

		Family		
	Fraction	Loan=No	Loan=Yes	Entropy
1	0.307692	0.5	0.5	0.5
2	0.076923	0	1	1
3	0.153846	0.50	0.50	0.5
4	0.461538	0.50	0.50	0.5
Gini Index			0.538462	

Men			Female		
Number	Age	Married	Number	Age	Married
Man 1	1	No	Woman 1	41	Yes
Man 2	2	No	Woman 2	42	Yes
...
Man 40	40	No	Woman 59	99	Yes
			Woman 60	99	No

Partition D into D_{train} (training / "growing"), $D_{\text{validation}}$ (validation / "pruning")

Build complete tree T on D_{train}

UNTIL accuracy on $D_{\text{validation}}$ decreases DO

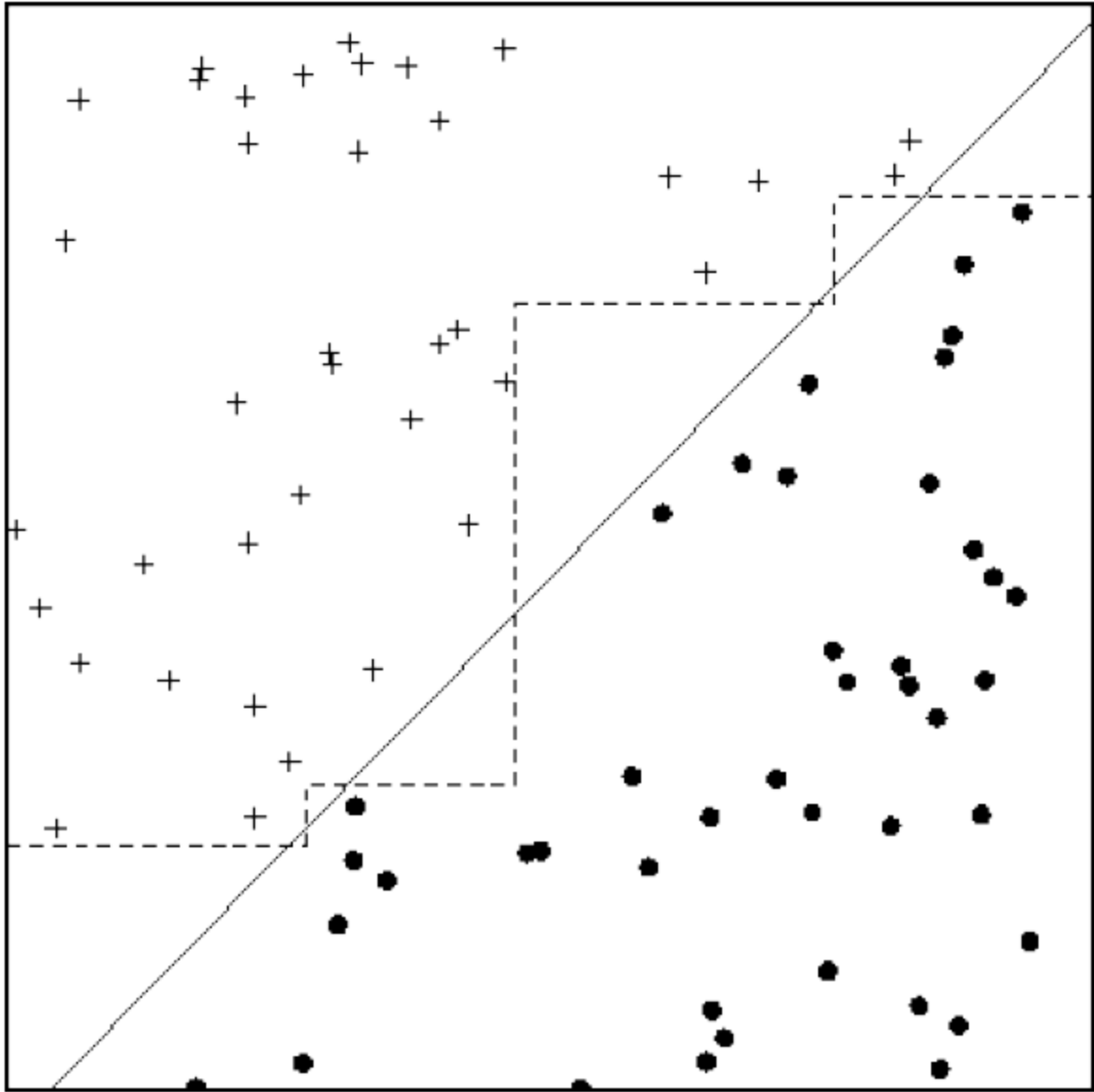
FOR each non-leaf node candidate in T

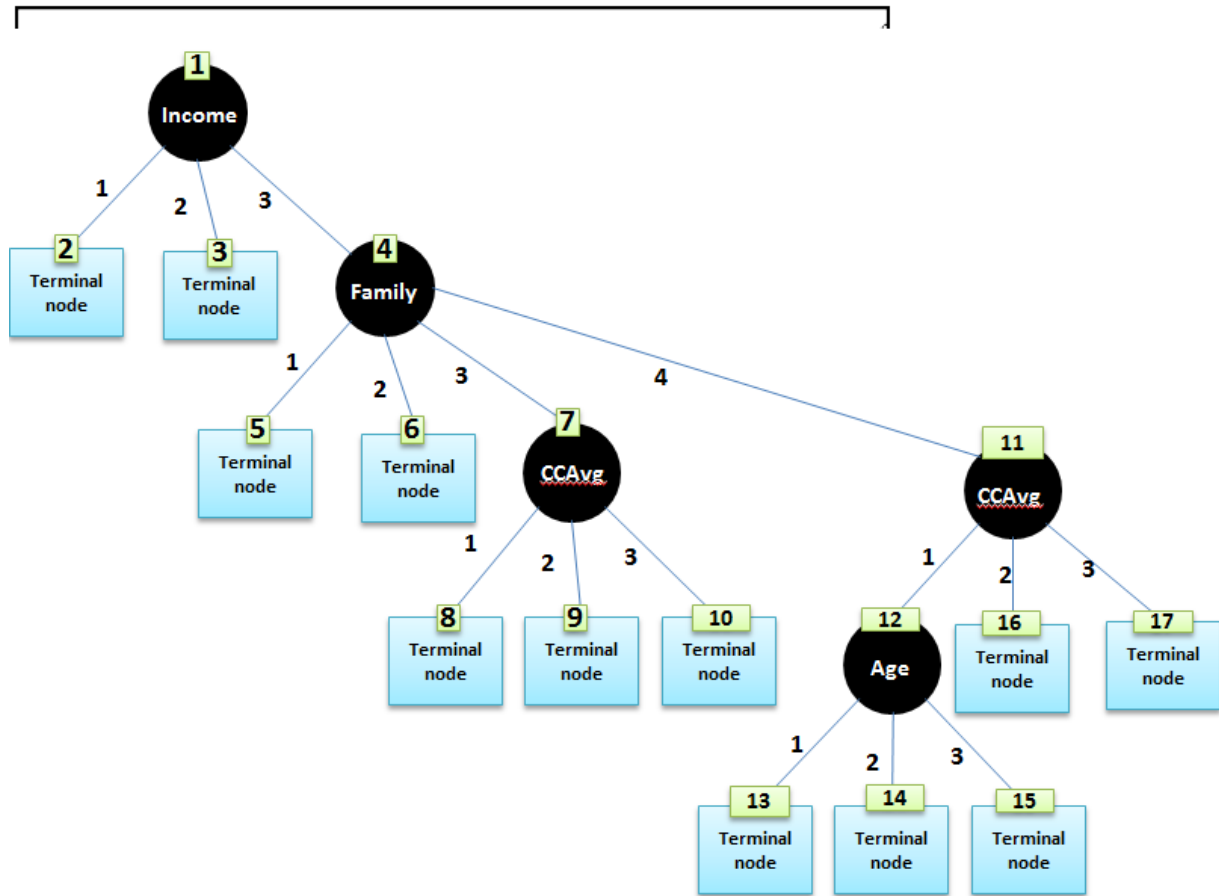
Temp[candidate] \leftarrow Prune (T , candidate)

Accuracy[candidate] \leftarrow Test (Temp[candidate], $D_{\text{validation}}$)

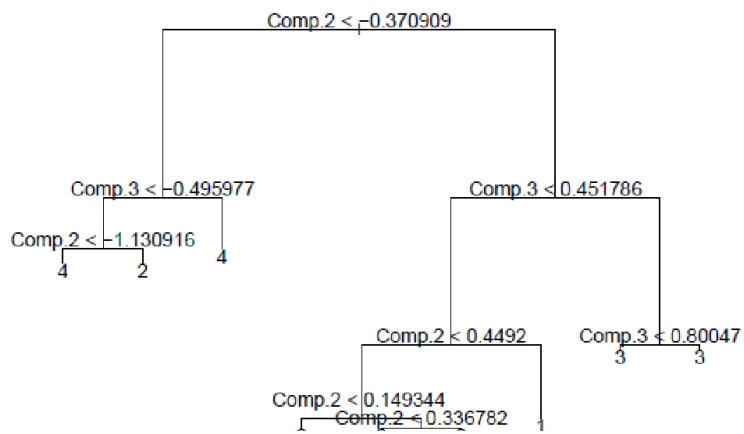
$T \leftarrow T' \in \text{Temp}$ with best value of Accuracy (best increase; greedy)

RETURN (pruned)

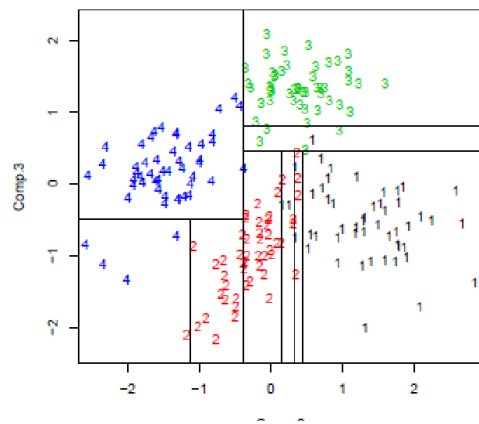




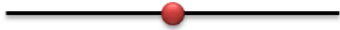
Axis-Parallel Tree



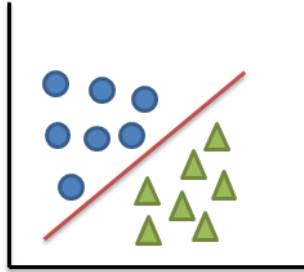
Associated Decision Boundaries



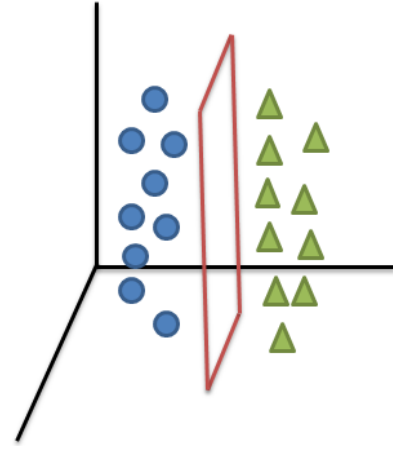
1 D



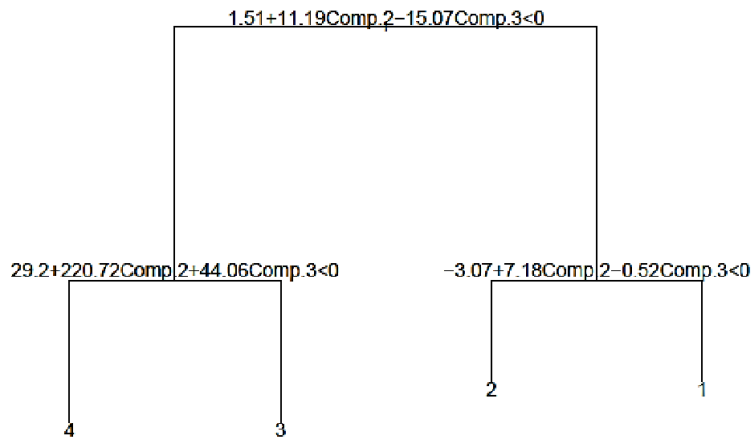
2 D



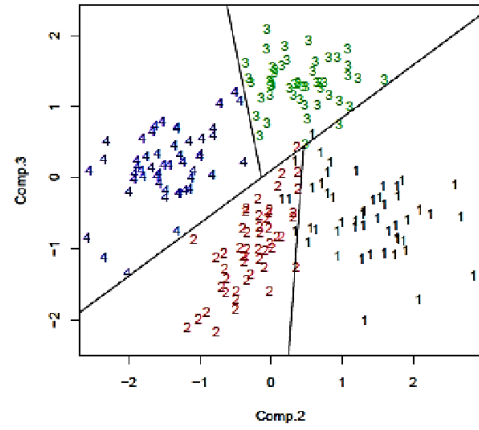
3 D

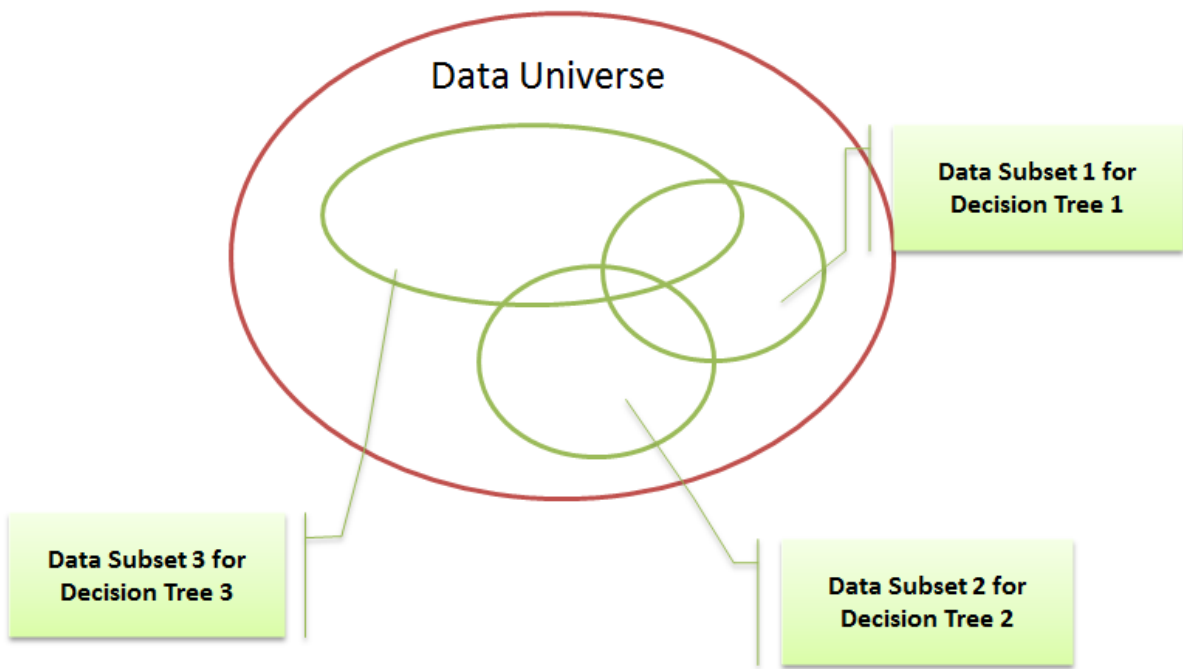


Oblique Tree



Associated Decision Boundaries





S No	Variable		
1	X1		X1, X2, X3 variable set for Decision Tree 1
2	X2		
3	X3		X3, X4, X5 variable set for Decision Tree 2
4	X4		
5	X5		
...	...		Xa, Xb, Xn variable set for Decision Tree 3
N	Xn		

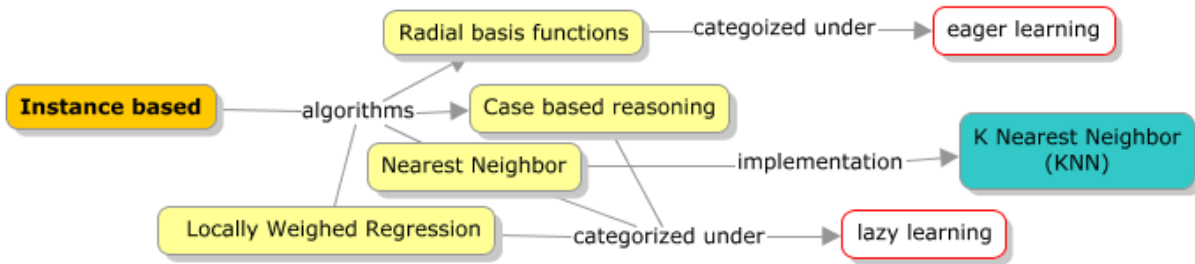
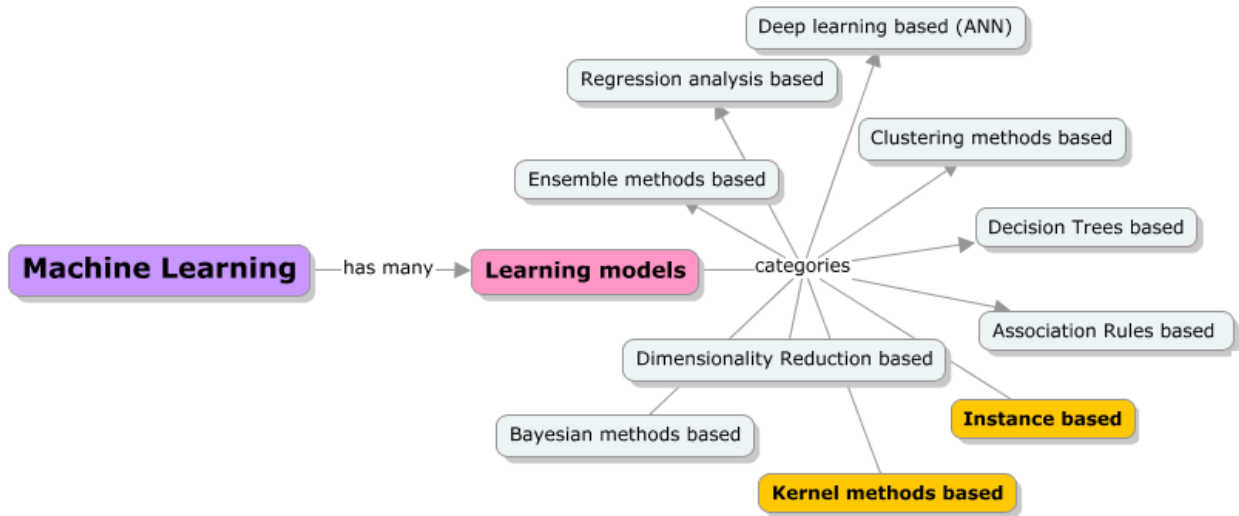
$$\text{MSE}[\hat{f}_m(p)] = O\left(\frac{1}{m^{4/(D+4)}}\right)$$

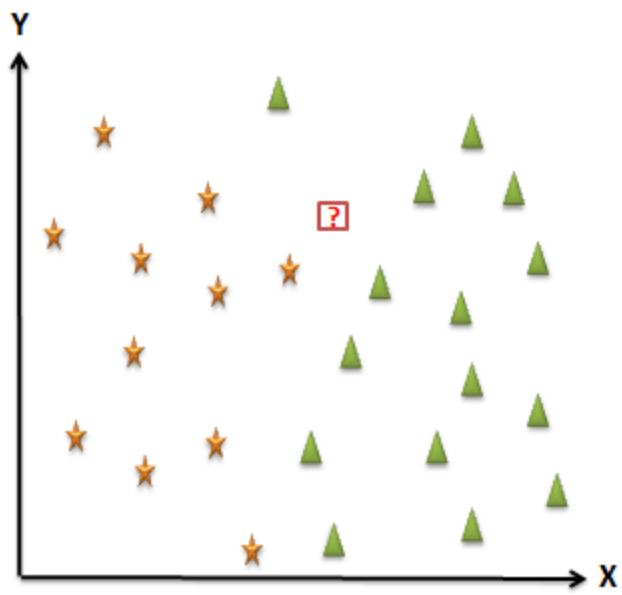
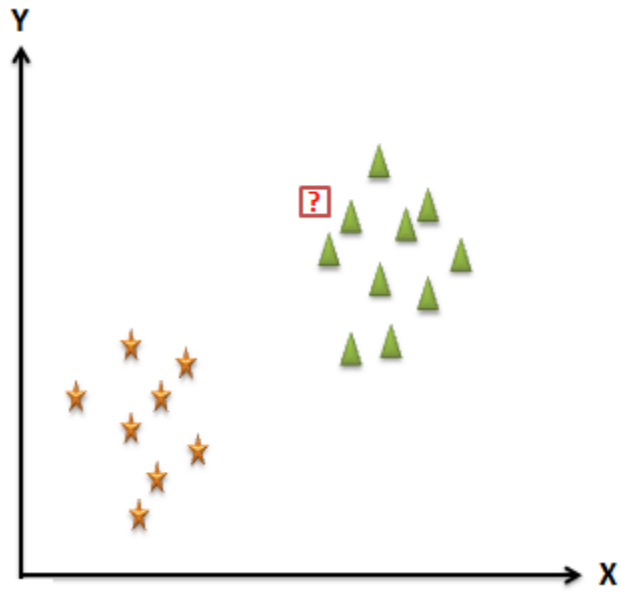
Dimensionality	Required Sample Size
1	4
2	19
5	786
7	10,700
10	842,000

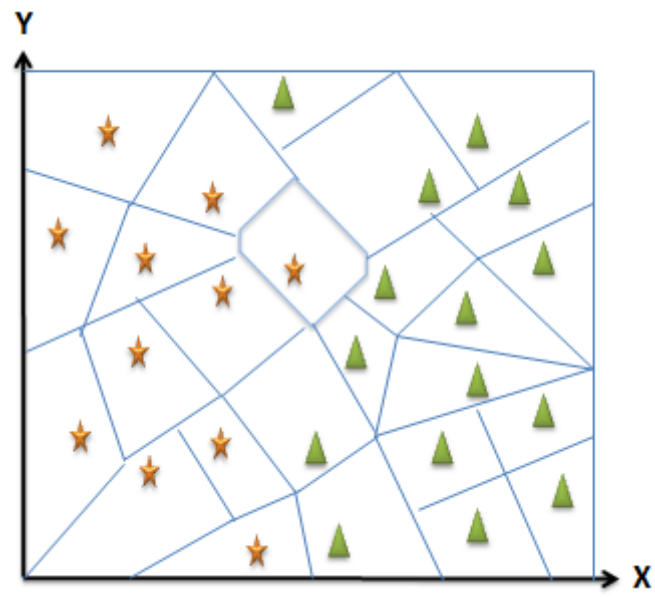
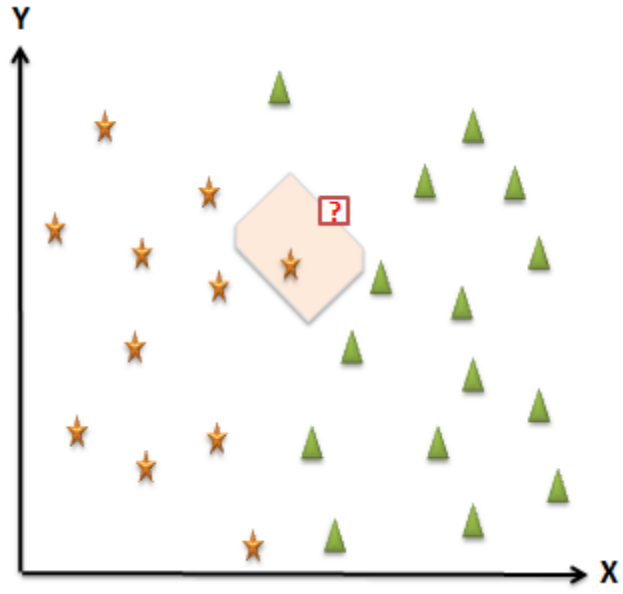
$$d_H(P(Y_+), P(Y_-)) = \sqrt{\sum_{i \in V} \left(\sqrt{P(Y_+ | X_i)} - \sqrt{P(Y_- | X_i)} \right)^2}$$

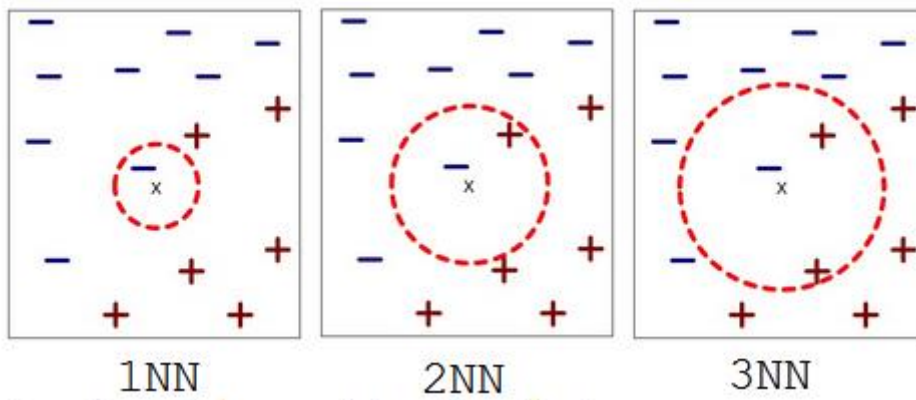
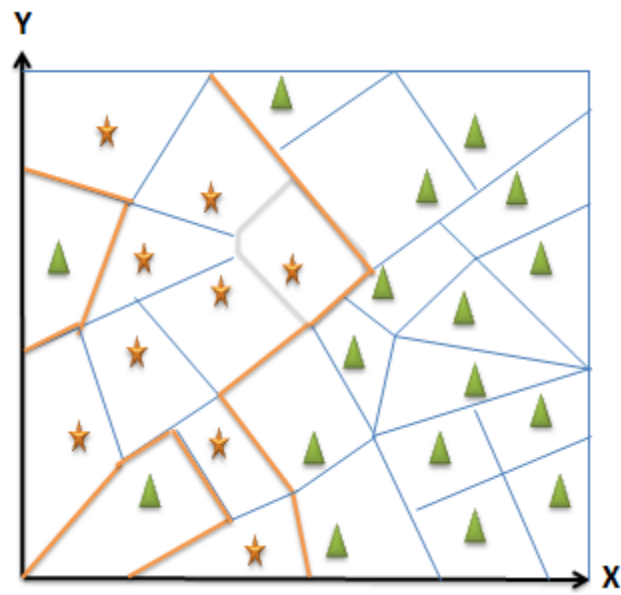
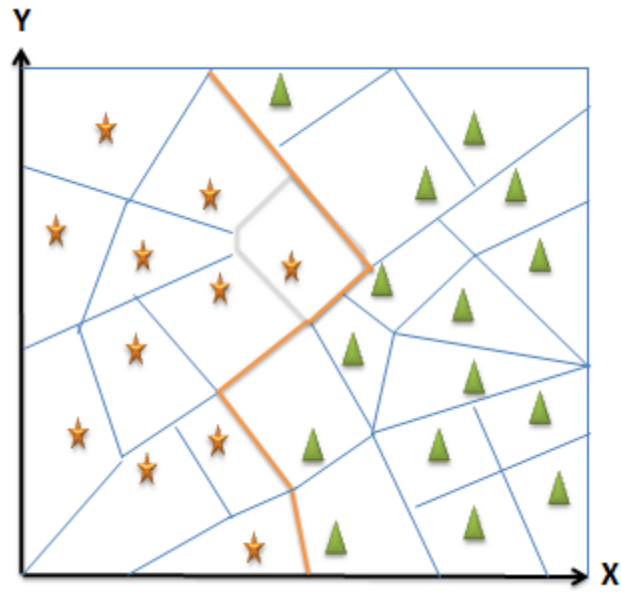
High	Low	0
High	High	1

CHAPTER 6





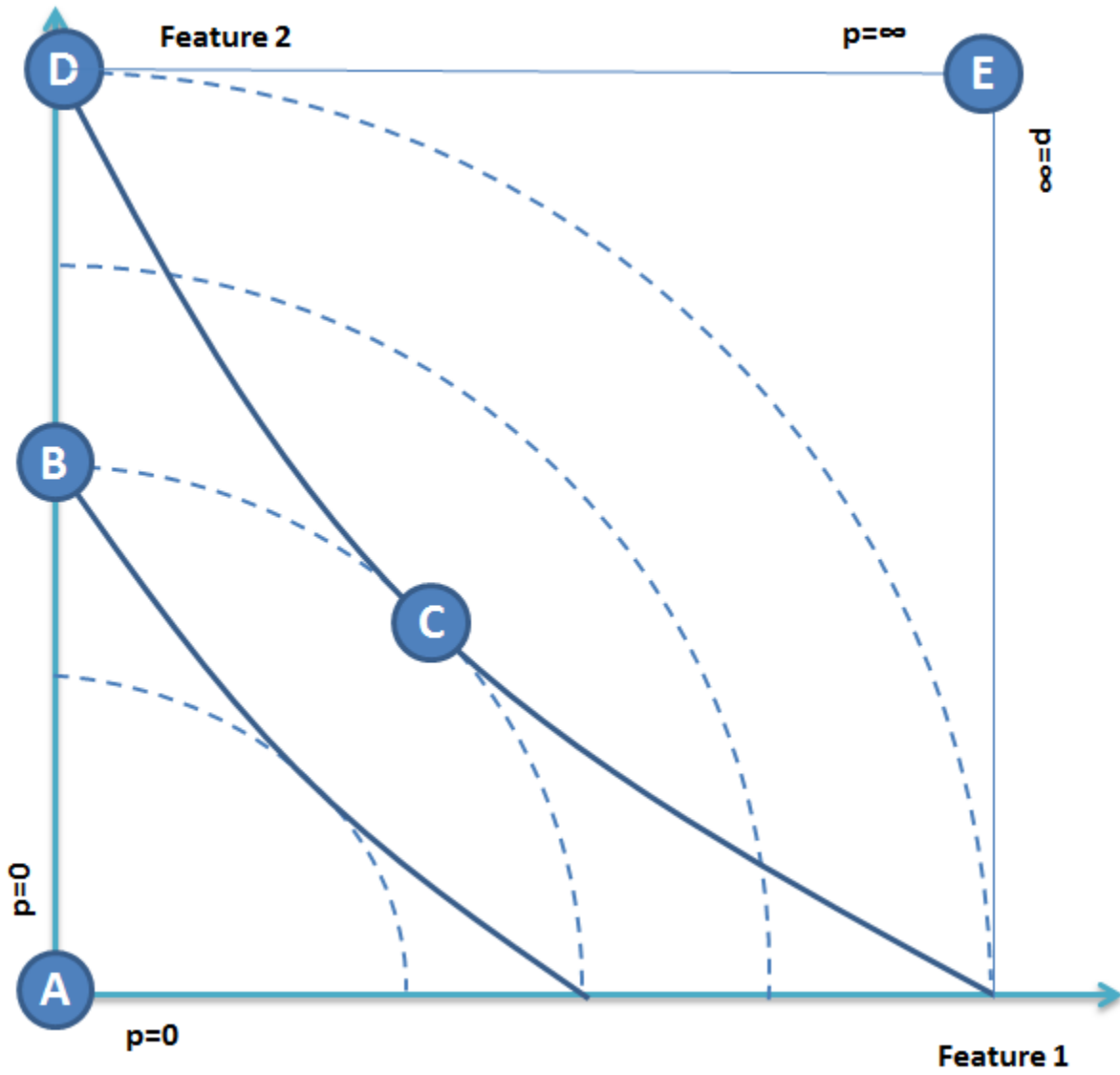


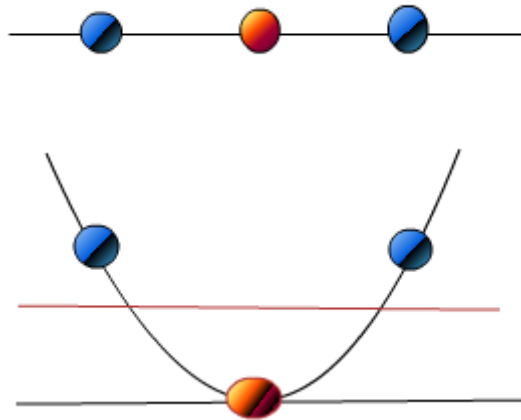
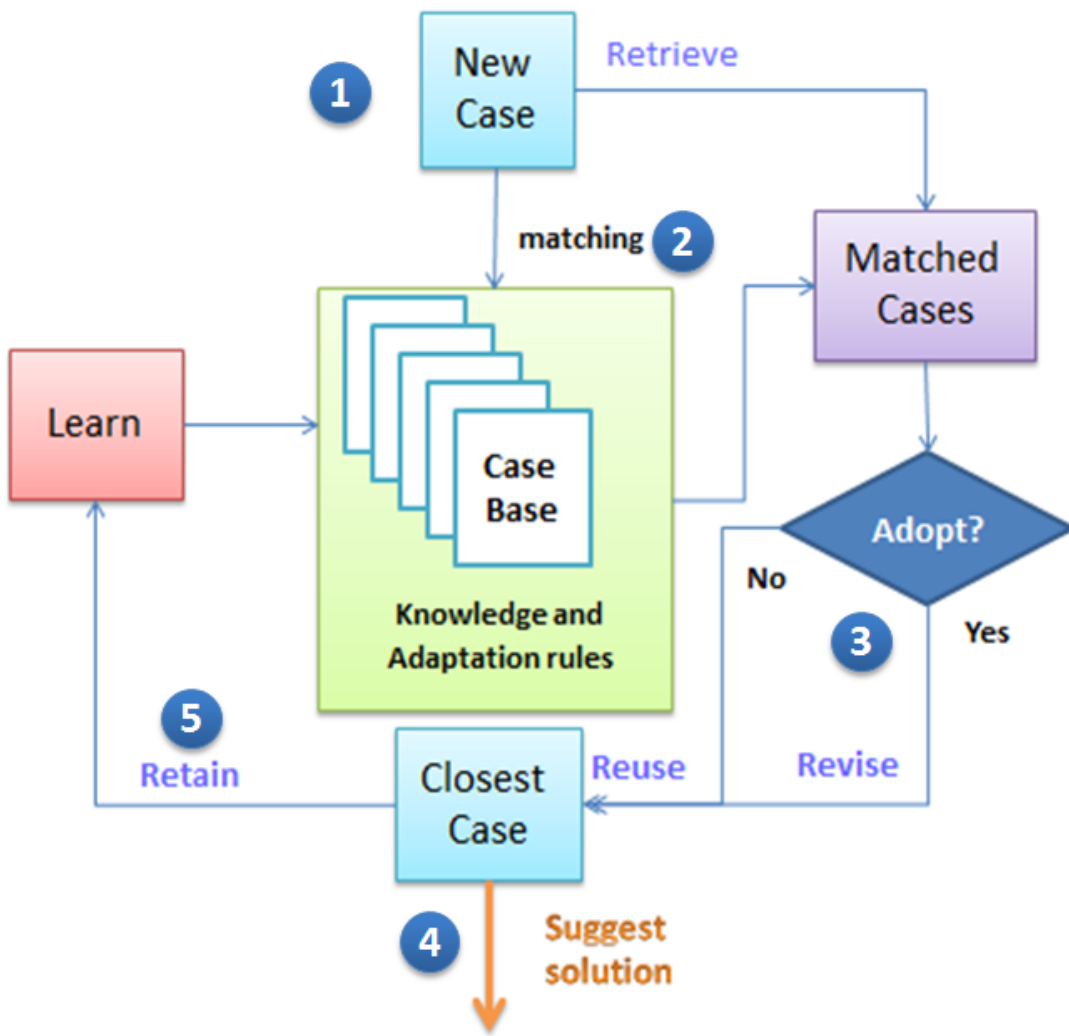


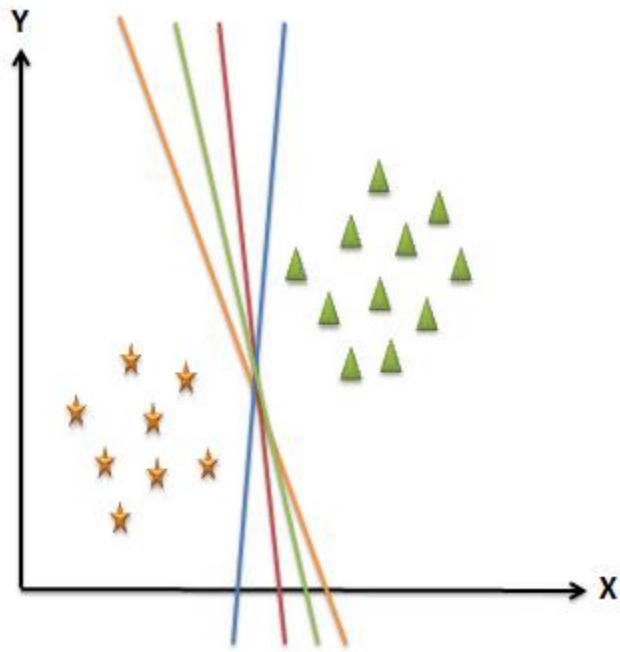
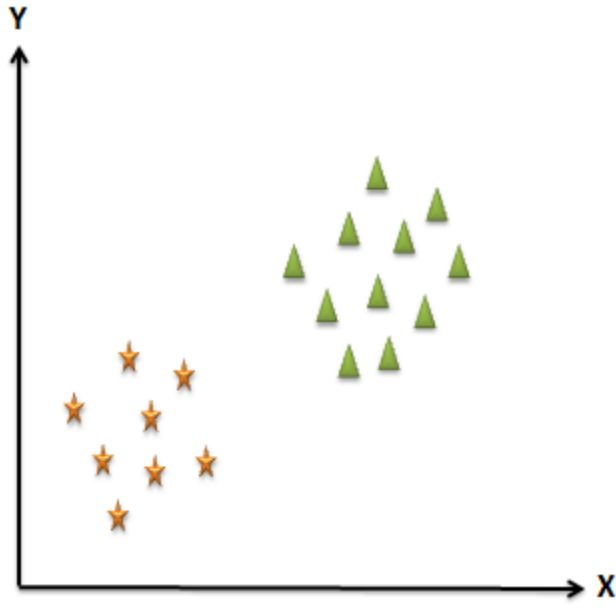
$$D(x, x') = \sqrt{\sum |x_d - x'_d|^2}$$

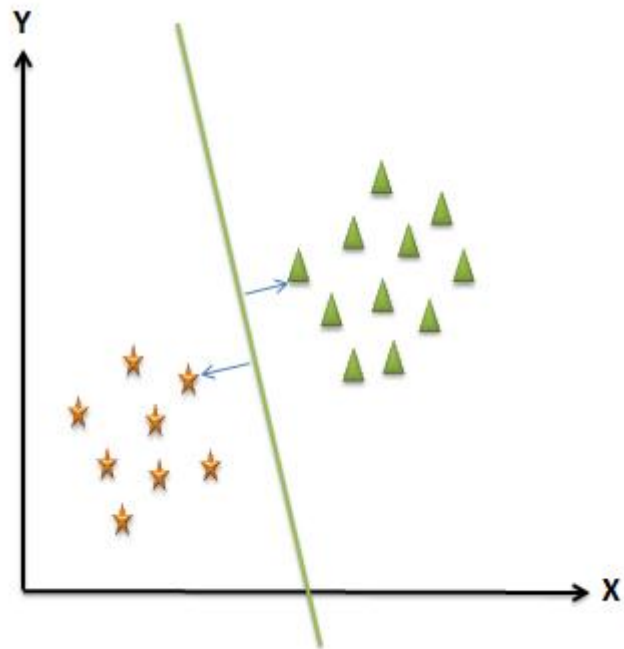
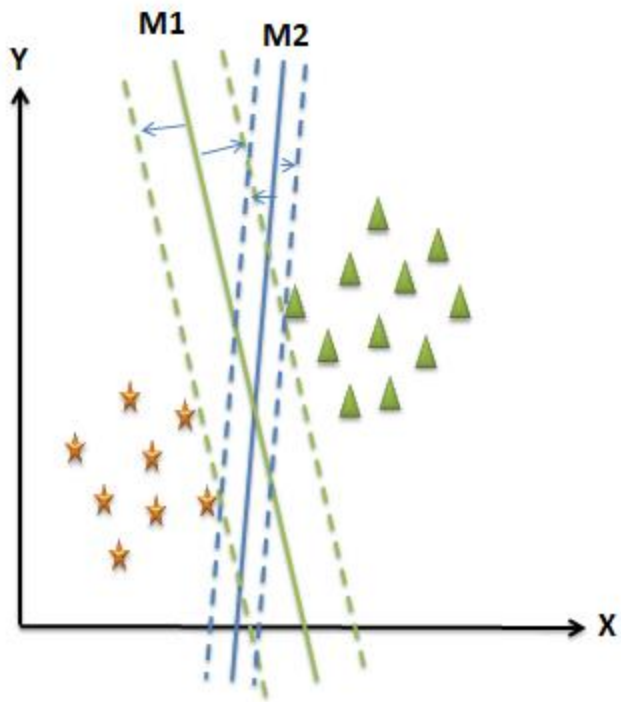
$$D(x, x') = \sum_d 1_{x_d \neq x'_d}$$

$$D(x, x') = \sqrt[p]{\sum |x_d - x'_d|^p}$$



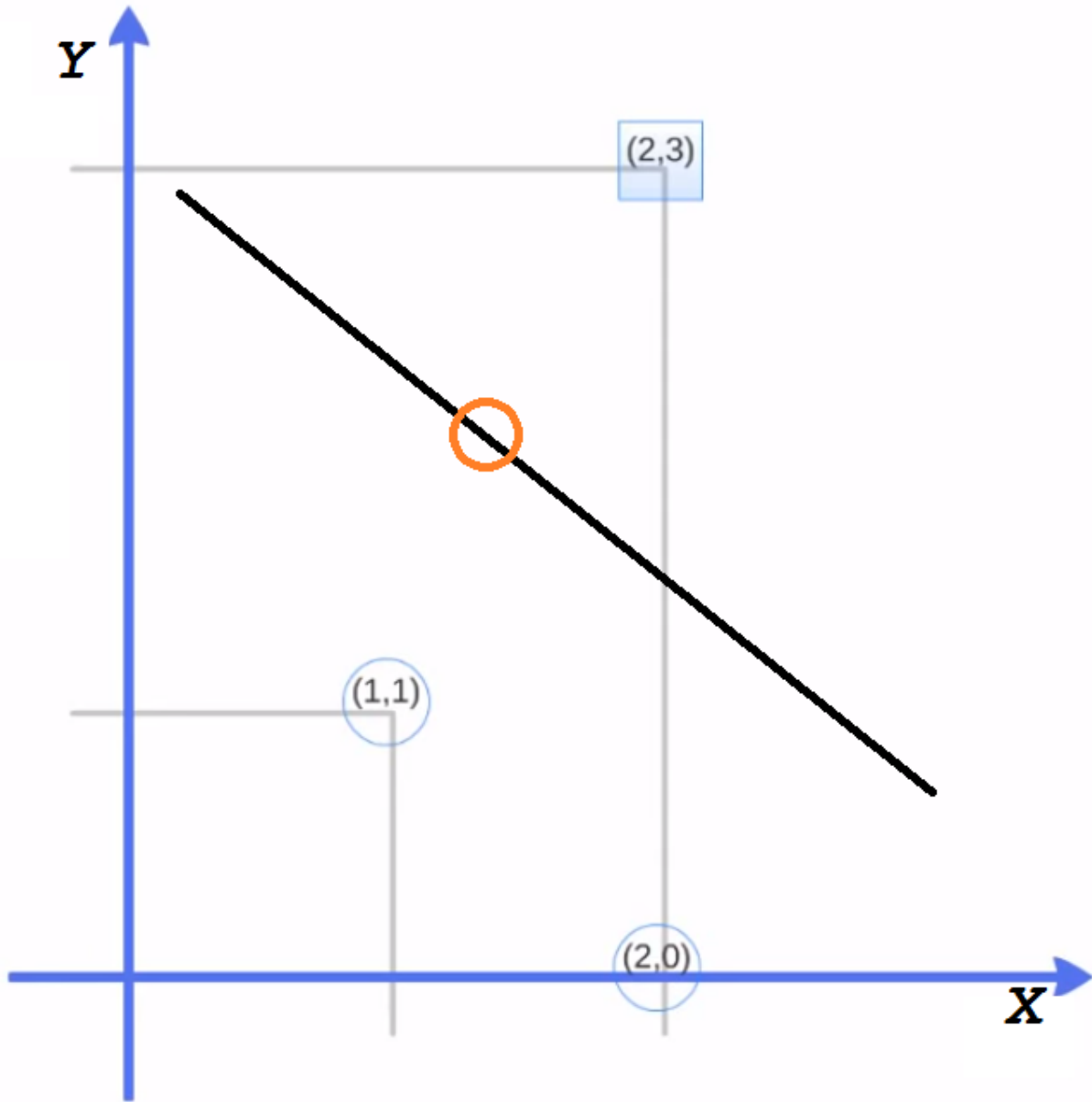






$$\mathbf{a}^{\text{ev}} = \sum_{i=0}^N \lambda_i y_i \vec{x}_i$$

$$\sum_{i=0}^N \lambda_i y_i = 0$$



$$\omega_0 = 1 - 8a \quad 3a + 1 - 8a = -1$$

$$\therefore 5a = 2$$

$$a = \frac{2}{5}$$

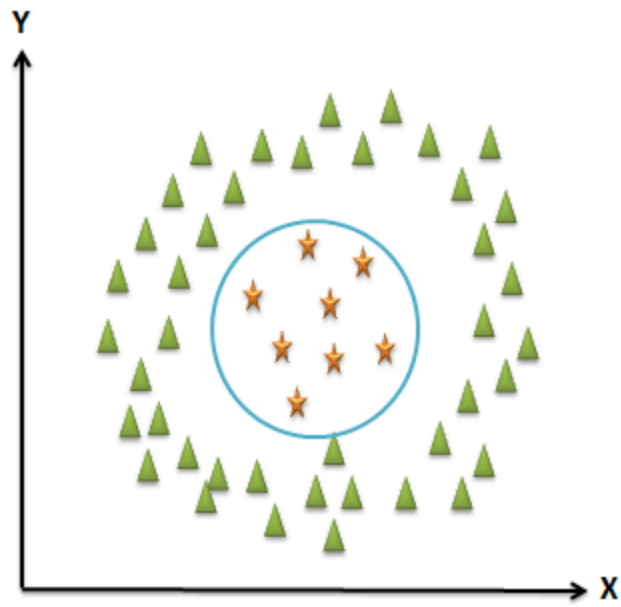
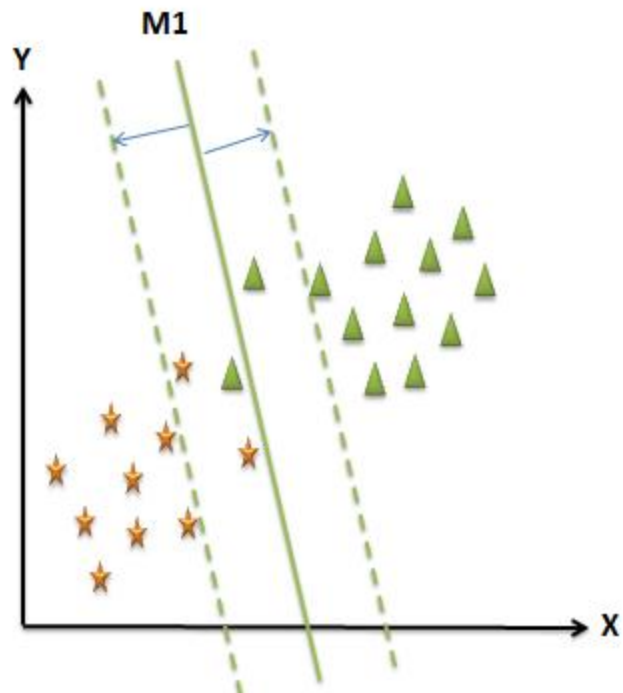
$$\omega_0 = 1 - 8 \cdot \frac{2}{5} = \frac{5 - 16}{5}$$

$$\omega_0 = -\frac{11}{5}$$

$$\bar{w} = \left(\frac{2}{5}, \frac{4}{5} \right)$$

$$g(\vec{x}) = \frac{2}{5}x_1 + \frac{4}{5}x_2 - \frac{11}{5}$$

$$g(\vec{x}) = x_1 + 2x_2 - 5.5$$



$$\text{Minimize : } \Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$\text{Subject to : } d_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 \quad \forall i$$



Introduce slack variables $\xi_i \geq 0$



$$\text{Minimize : } \Phi(\mathbf{w}) = \frac{1}{2} \mathbf{w}^T \mathbf{w}$$

$$\text{Subject to : } d_i(\mathbf{w}^T \mathbf{x}_i + b) \geq 1 - \xi_i \quad \forall i$$

Polynomial:

$$K_p(\mathbf{X}, \mathbf{Y}) = (1 + \mathbf{X} \cdot \mathbf{Y})^p$$

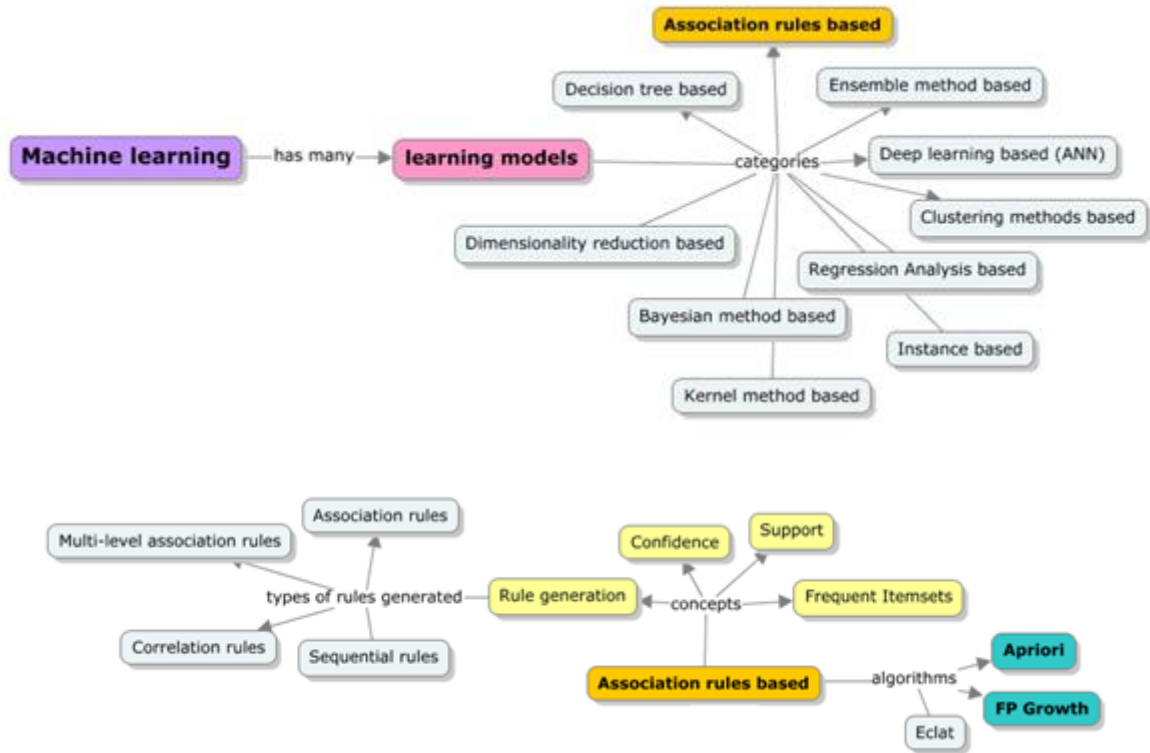
Radial Basis Function
(RBF) or Gaussian:

$$K_r(\mathbf{X}, \mathbf{Y}) = e^{-\frac{1}{2\sigma^2} \|\mathbf{X} - \mathbf{Y}\|_2^2}$$

Hyperbolic Tangent:

$$K_s(\mathbf{X}, \mathbf{Y}) = \tanh(\beta_0 \mathbf{X} \cdot \mathbf{Y} + \beta_1)$$

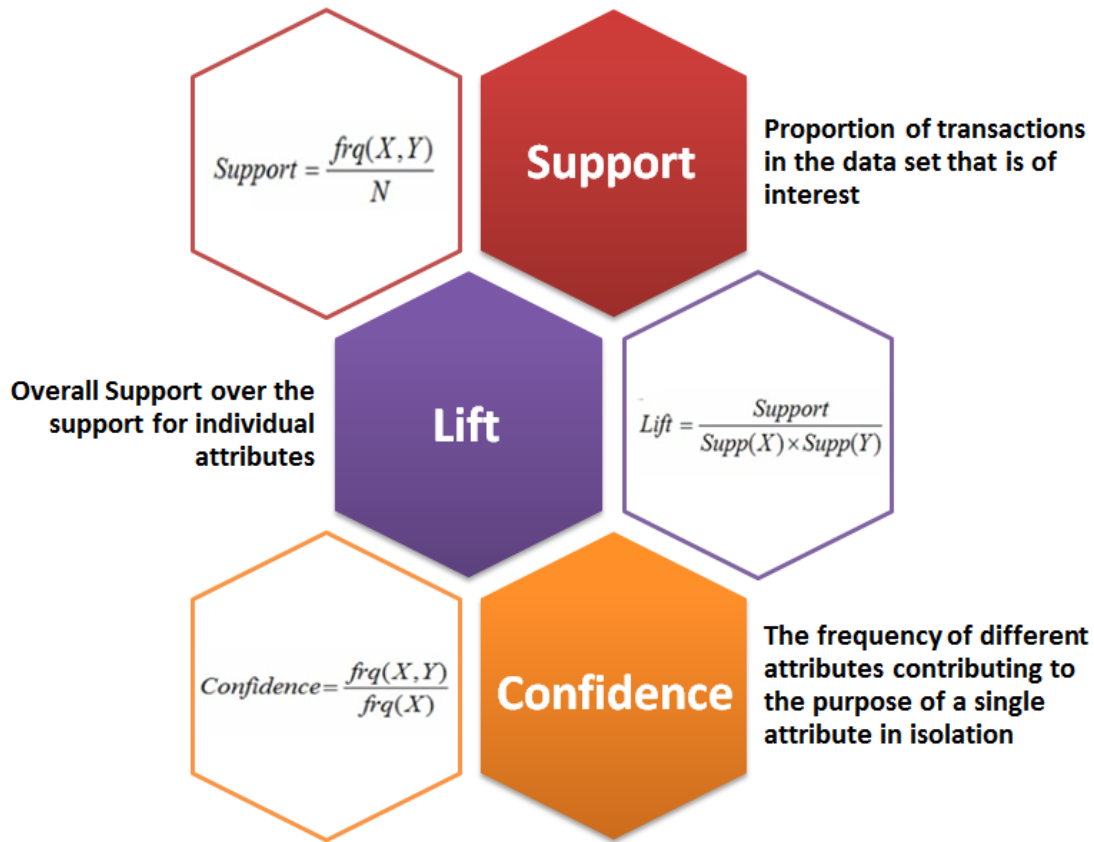
CHAPTER 7

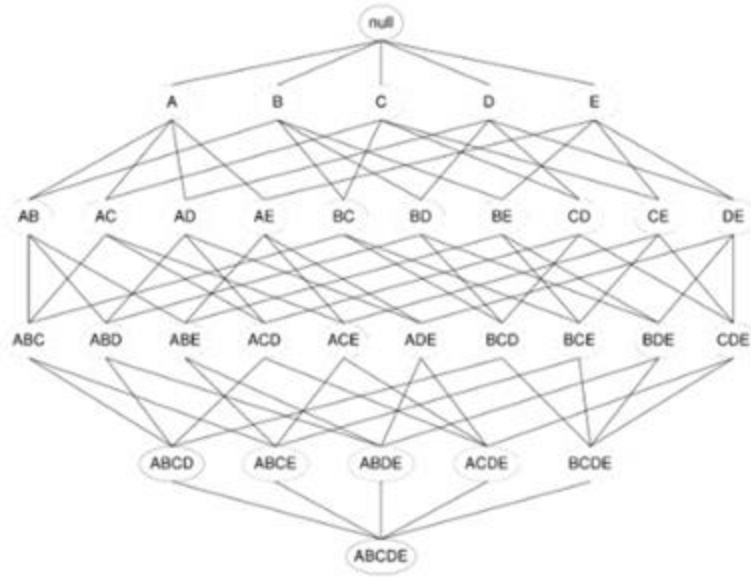


$$Support = \frac{freq(X, Y)}{N}$$

$$Confidence = \frac{freq(X, Y)}{freq(X)}$$

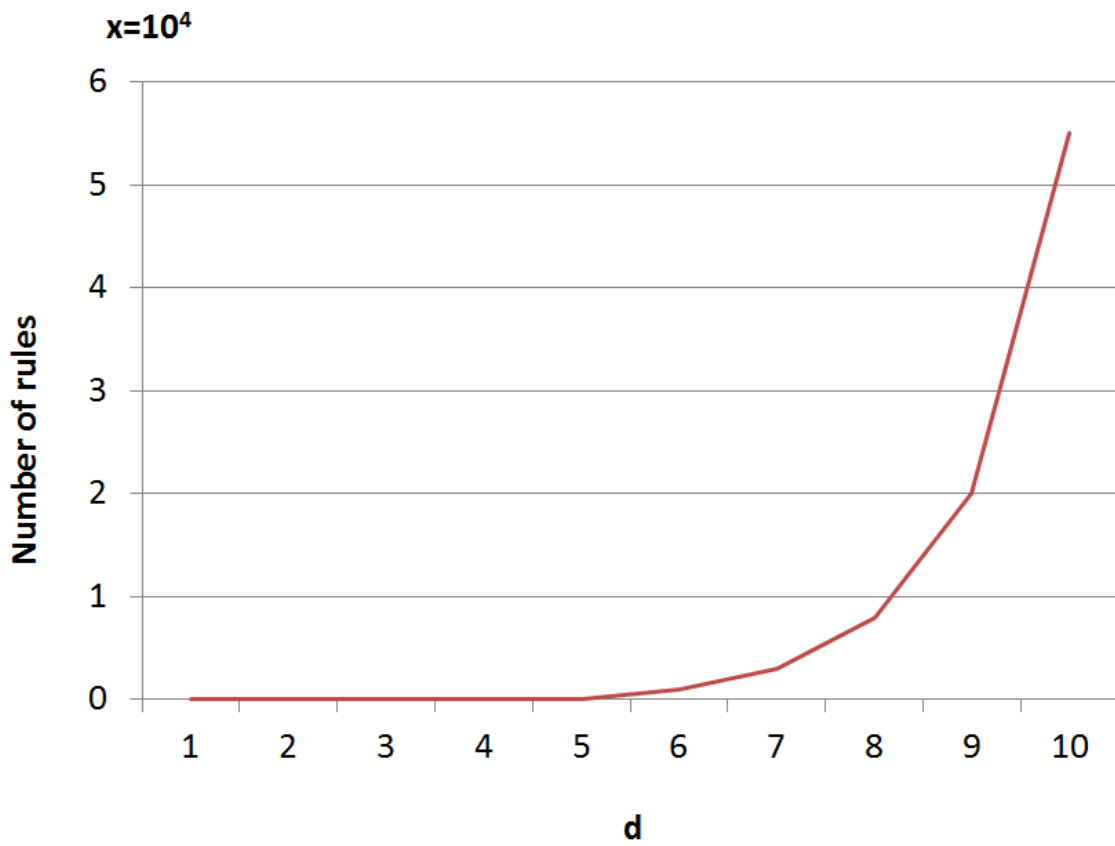
$$Lift = \frac{Support}{Supp(X) \times Supp(Y)}$$



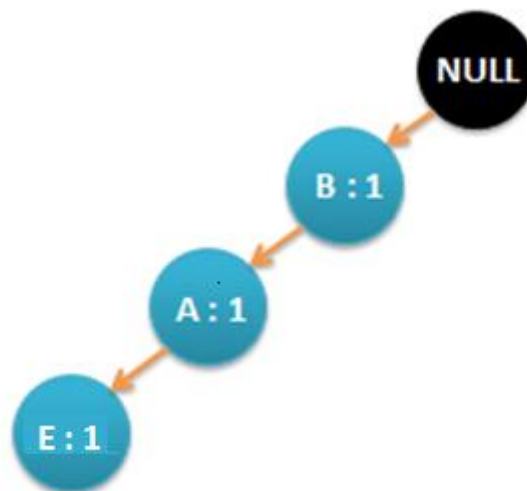
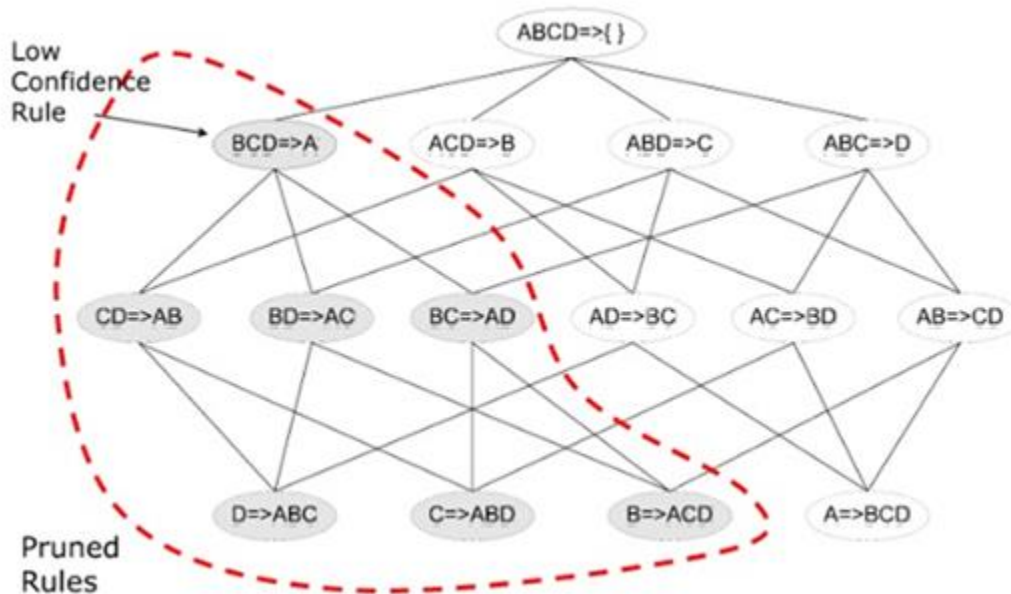


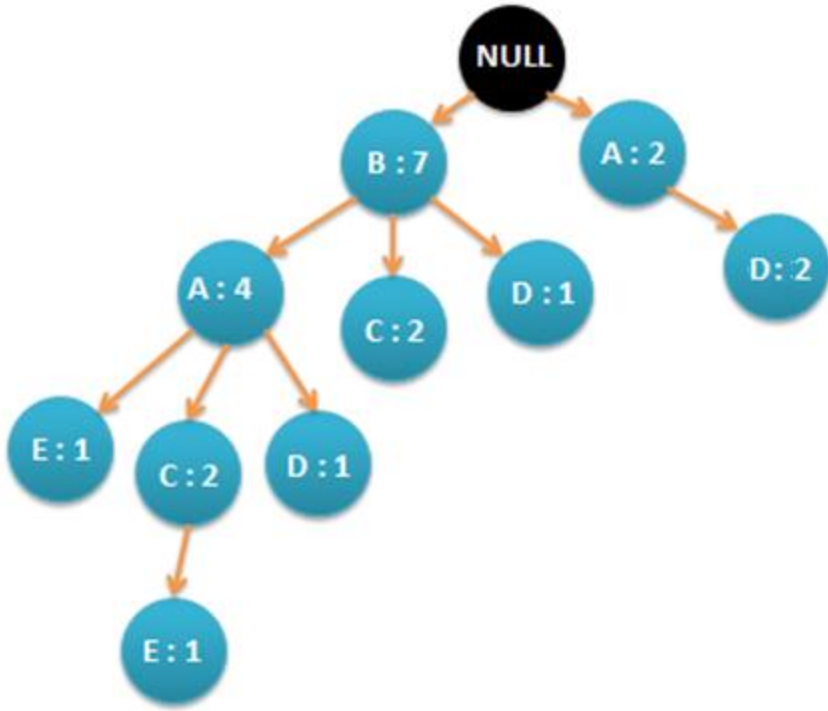
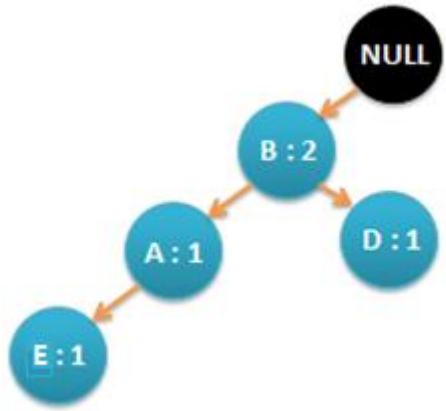
$$\sum_{k=1}^{d-1} \left[\binom{d}{k} \times \sum_{j=1}^{d-k} \binom{d-k}{j} \right]$$

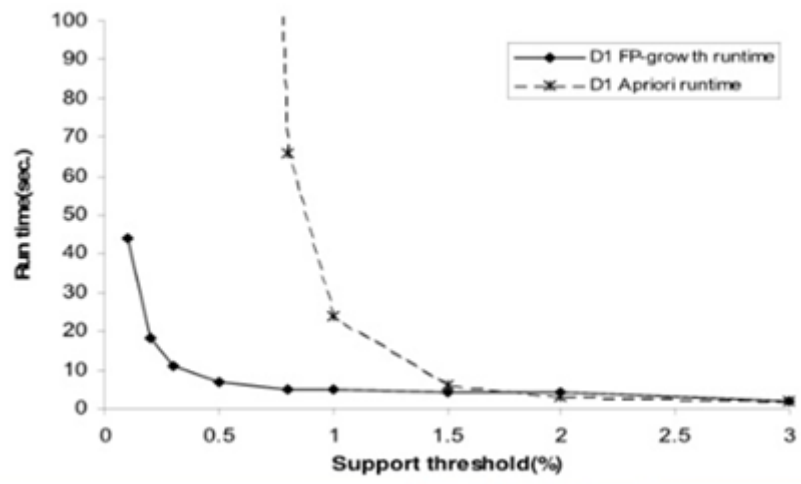
$$= 3^d - 2^{d+1} + 1$$



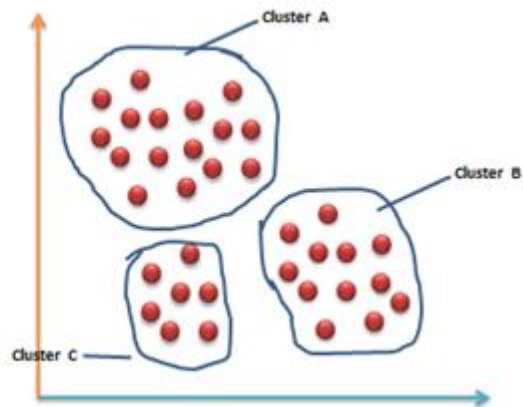
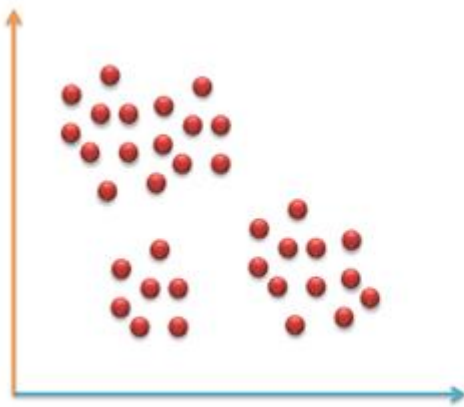
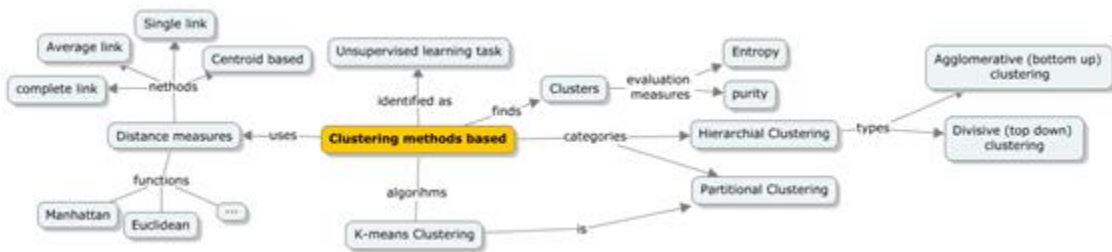
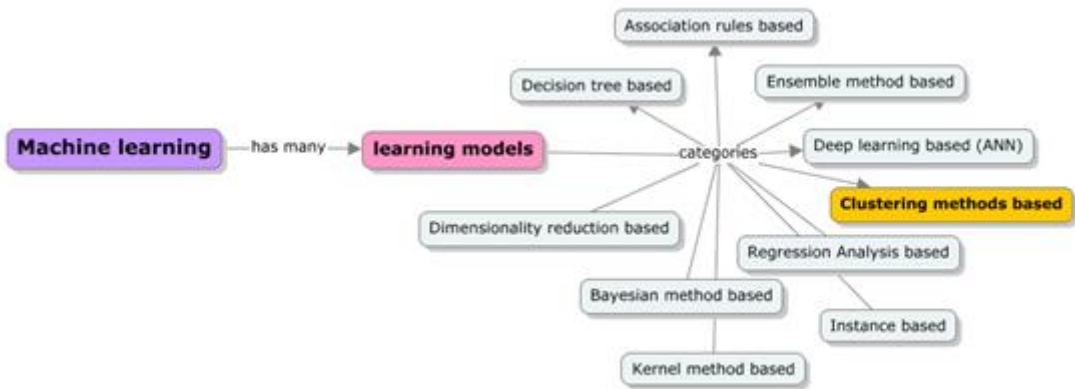
C_k : Candidate itemset of size k
 L_k : frequent itemset of size k
 $L_1 = \{\text{frequent items}\};$
for ($k = 1; L_k \neq \emptyset; k++$) **do begin**
 C_{k+1} = candidates generated from L_k ;
 for each transaction t in database **do**
 increment the count of all candidates in C_{k+1}
 that are contained in t
 L_{k+1} = candidates in C_{k+1} with min_support
 end
return $\cup_k L_k$;

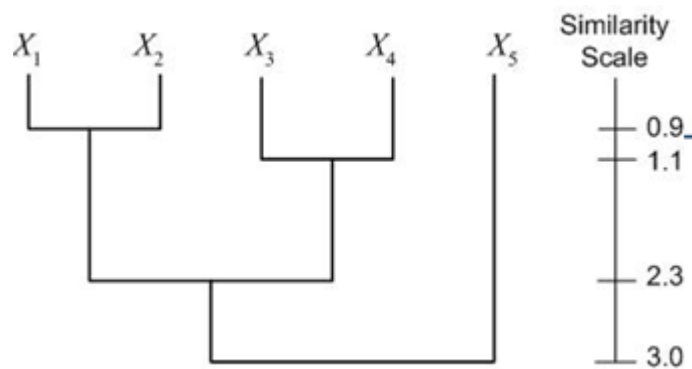
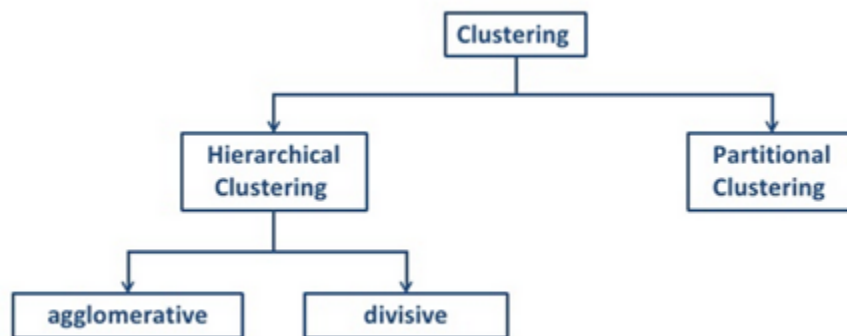
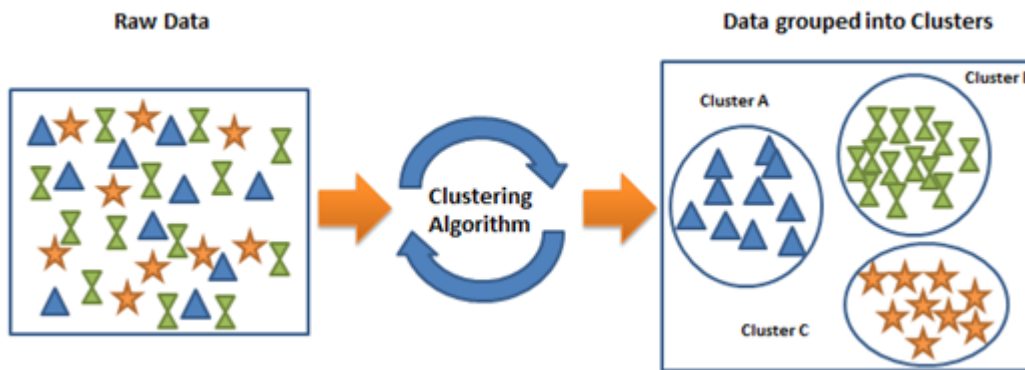


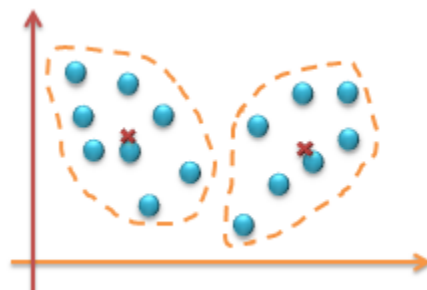
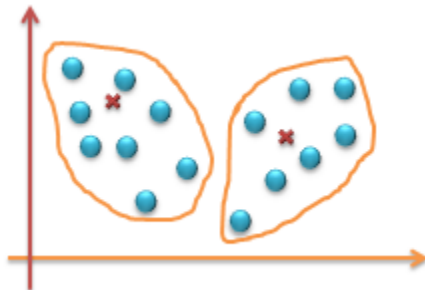
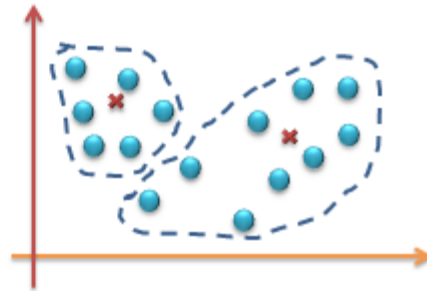
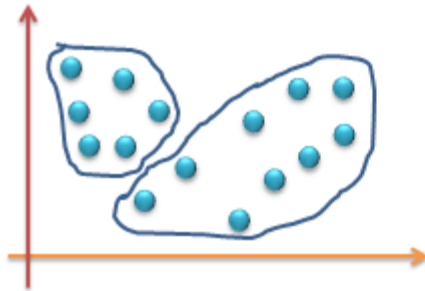
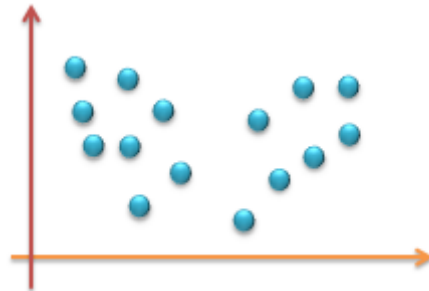
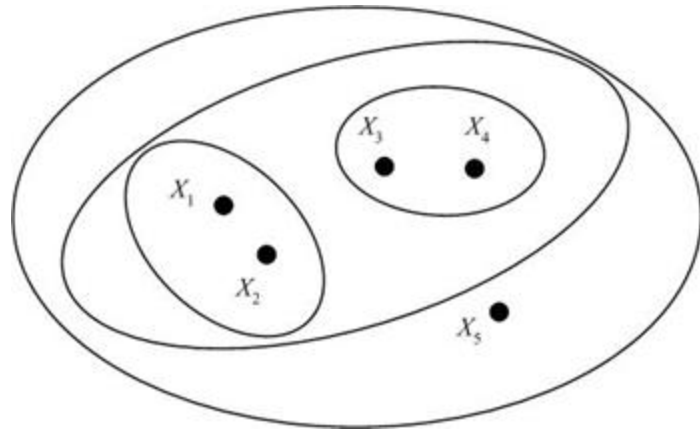


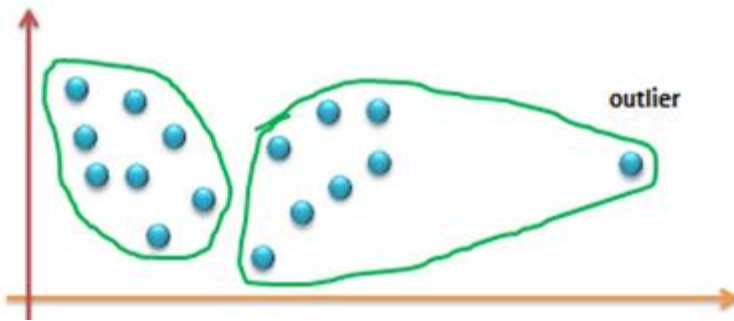
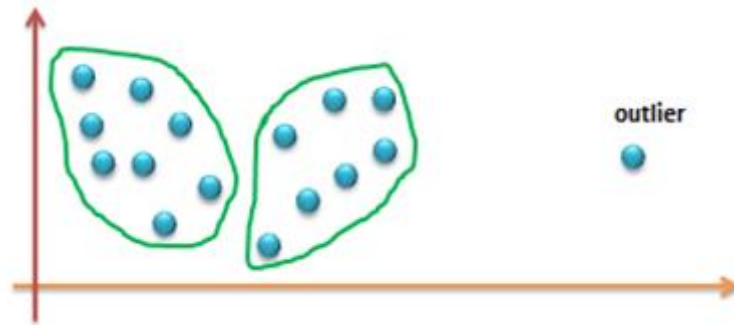
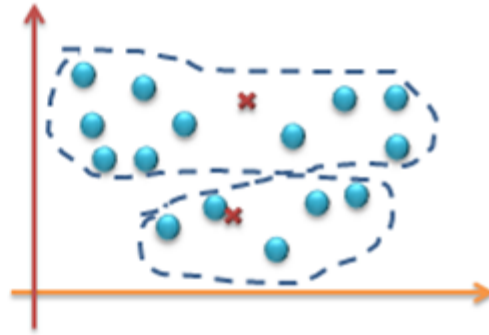
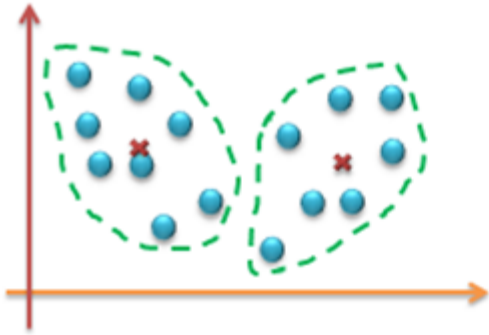
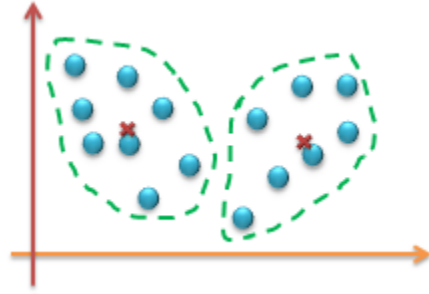
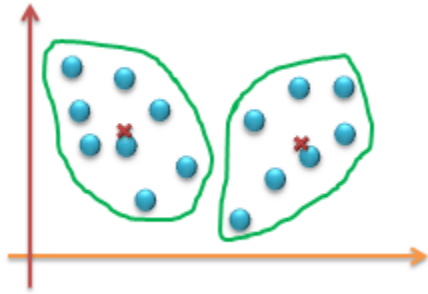


CHAPTER 8

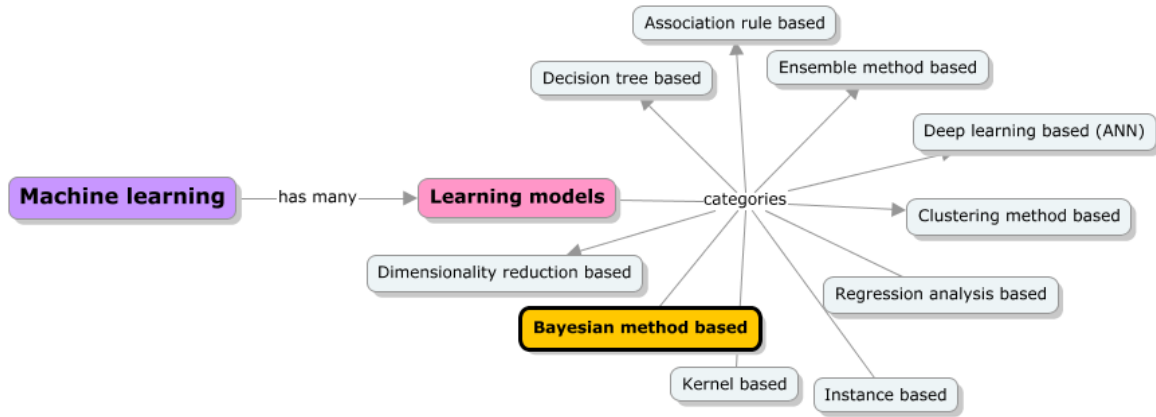




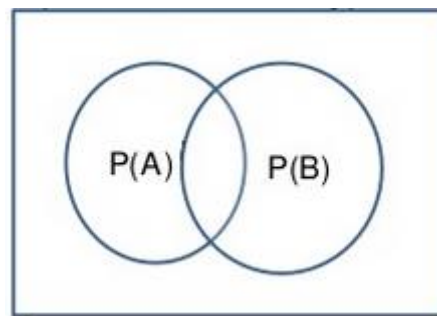
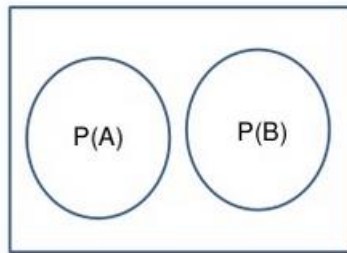




CHAPTER 9



$$\sigma = \sqrt{\frac{1}{n-1} \sum_{i=1}^{i=n} (x_i - \bar{x})^2}$$



$$\sum_x \sum_y P(X = x \text{ and } Y = y) = 1$$

$$P(X = x) = \sum_y P(X = x, Y = y) = \sum_y P(X = x | Y = y)P(Y = y)$$

X1	X2	Y (X1+X2)	Z (X1-X2)
1	1	2	0
1	2	3	-1
1	3	4	-2
1	4	5	-3
1	5	6	-4
1	6	7	-5
2	1	3	1
2	2	4	0
2	3	5	-1
2	4	6	-2
2	5	7	-3
2	6	8	-4
3	1	4	2
3	2	5	1
3	3	6	0
3	4	7	-1
3	5	8	-2
3	6	9	-3
4	1	5	3
4	2	6	2
4	3	7	1
4	4	8	0
4	5	9	-1
4	6	10	-2
5	1	6	4
5	2	7	3
5	3	8	2
5	4	9	1
5	5	10	0
5	6	11	-1
6	1	7	5
6	2	8	4
6	3	9	3
6	4	10	2
6	5	11	1
6	6	12	0

		Y											
		2	3	4	5	6	7	8	9	10	11	12	Marginal Z
Z	-5	0	0	0	0	0	1/36	0	0	0	0	0	1/36
	-4					1/36		1/36					1/18
	-3				1/36		1/36		1/36				1/12
	-2			1/36		1/36		1/36		1/36			1/9
	-1		1/36		1/36		1/36		1/36		1/36		5/36
	0	1/36		1/36		1/36		1/36		1/36		1/36	1/6
	1		1/36		1/36		1/36		1/36		1/36		5/36
	2			1/36		1/36		1/36		1/36			1/9
	3				1/36		1/36		1/36				1/12
	4					1/36		1/36					1/18
	5						1/36						1/36
		1/36	1/18	1/12	1/9	5/36	1/6	5/36	1/9	1/12	1/18	1/36	

Discrete	Continuous
Bernouli	Normal
Binomial	T distribution
Negative binomial	Gamma
Geometric	Chi Square
Poisson	Exponential
	Weibull
	F Distribution

Possible demand X	Number of days	Probability [P(X)]
3	3	0.06
4	7	0.14
5	12	0.24
6	14	0.28
7	10	0.2
8	4	0.08

Possible demand X	Probability [P(X)]	Weighted Value [XP(X)]
3	0.06	0.18
4	0.14	0.56
5	0.24	1.2
6	0.28	1.68
7	0.2	1.4
8	0.08	0.64
	1	E(X) = 5.66

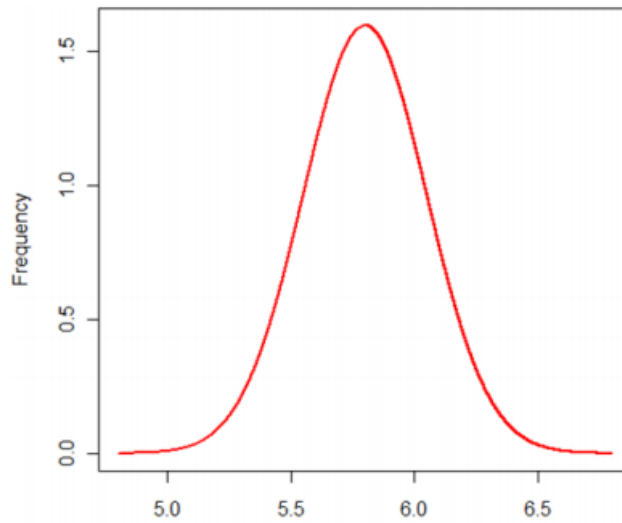
Possible demand X	Probability [P(X)]	Weighted Value [XP(X)]	Squared demand (X ²)	Weighted Square [X ² P(X)]
3	0.06	0.18	9	0.54
4	0.14	0.56	16	2.24
5	0.24	1.2	25	6
6	0.28	1.68	36	10.08
7	0.2	1.4	49	9.8
8	0.08	0.64	64	5.12
	1	E(X) = 5.66		E(X²) = 33.78

$$\lim_{n \rightarrow \infty} c_r^n \left(\frac{\lambda}{n}\right)^r \left(1 - \frac{\lambda}{n}\right)^{n-r}$$

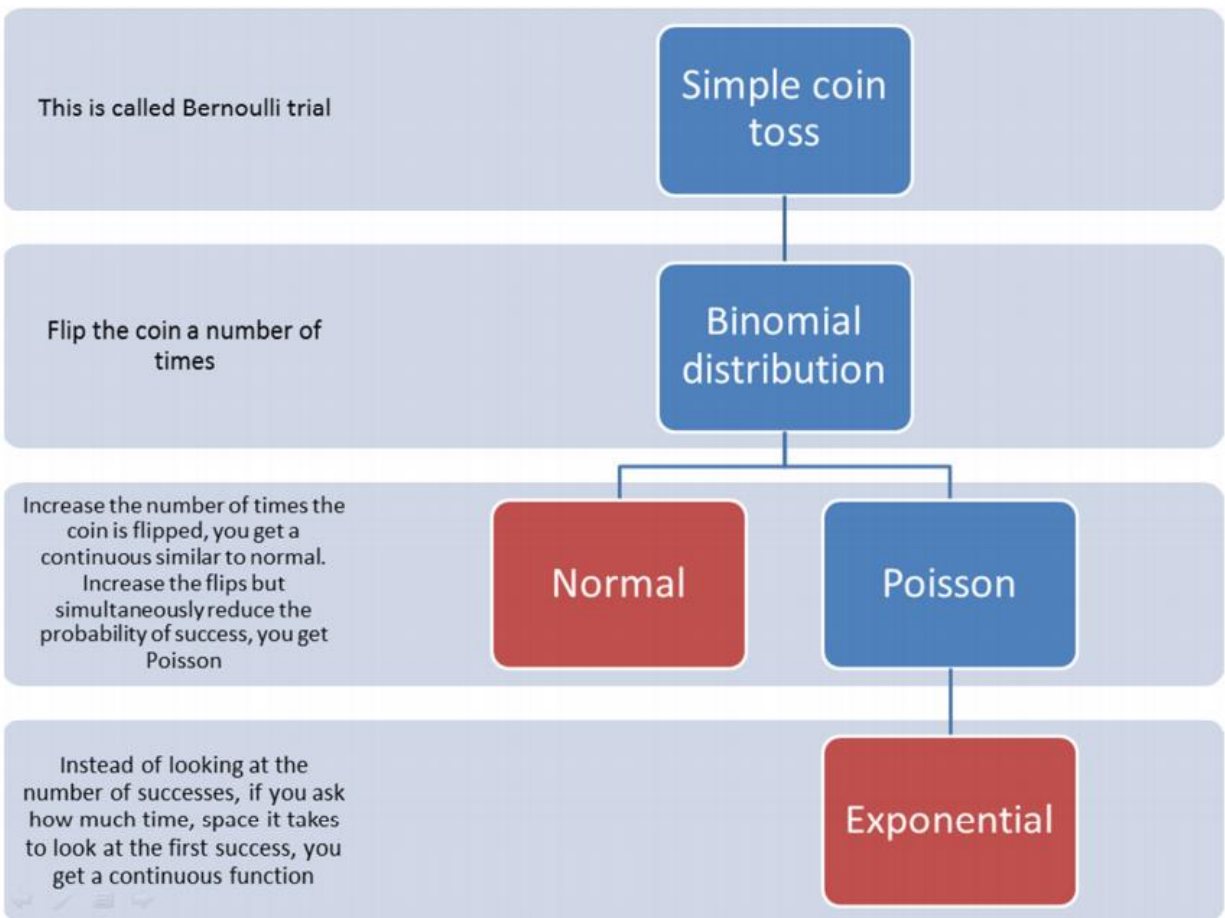
$$\lim_{n \rightarrow \infty} \frac{n!}{r!(n-r)!} \frac{\lambda^r}{n^r} \left(1 - \frac{\lambda}{n}\right)^{n-r}$$

$$P(X = r) = \frac{\lambda^r e^{-\lambda}}{r!}$$

$$P(X = 0) = \frac{e^{-\lambda} \lambda^0}{0!} = e^{-\lambda}$$



$$f(x) = \frac{1}{\sigma\sqrt{2\pi}} e^{-(x-\mu)^2/2\sigma^2}$$



$$p(A|B) = \frac{p(A)p(B|A)}{p(B)}$$

$$p(H|D) = \frac{p(H)p(D|H)}{p(D)}$$

$$p(H|D) = \frac{\overset{\text{Prior}}{p(H)} \overset{\text{Likelihood}}{p(D|H)}}{\underset{\text{Evidence}}{p(D)}} \underset{\text{Posterior Probability}}{\quad}$$

TRAINMULTINOMIALNB(C, ID)

```

1 V ← EXTRACTVOCABULARY(ID)
2 N ← COUNTDOCS(ID)
3 for each c ∈ C
4   do Nc ← COUNTDOCSINCLASS(ID, c)
5     prior[c] ← Nc/N
6     textc ← CONCATENATETEXTOFALLDOCSINCLASS(ID, c)
7     for each t ∈ V
8       do Tct ← COUNTTOKENSOFTERM(textc, t)
9       for each t ∈ V
10        do condprob[t][c] ←  $\frac{T_{ct}+1}{\sum_{t'} (T_{ct'}+1)}$ 
11 return V, prior, condprob
```

APPLYMULTINOMIALNB(C, V, prior, condprob, d)

```

1 W ← EXTRACTTOKENSFROMDOC(V, d)
2 for each c ∈ C
3   do score[c] ← log prior[c]
4     for each t ∈ W
5       do score[c] += log condprob[t][c]
6 return arg maxc ∈ C score[c]
```

```

TRAINBERNOULLINB(C, D)
1  $V \leftarrow \text{EXTRACTVOCABULARY}(\mathbf{D})$ 
2  $N \leftarrow \text{COUNTDOCS}(\mathbf{D})$ 
3 for each  $c \in \mathbf{C}$ 
4 do  $N_c \leftarrow \text{COUNTDOCSINCLASS}(\mathbf{D}, c)$ 
5    $\text{prior}[c] \leftarrow N_c/N$ 
6   for each  $t \in V$ 
7   do  $N_{ct} \leftarrow \text{COUNTDOCSINCLASSCONTAININGTERM}(\mathbf{D}, c, t)$ 
8      $\text{condprob}[t][c] \leftarrow (N_{ct} + 1)/(N_c + 2)$ 
9 return  $V, \text{prior}, \text{condprob}$ 

```

```

APPLYBERNOULLINB(C,  $V, \text{prior}, \text{condprob}, d$ )
1  $V_d \leftarrow \text{EXTRACTTERMSFROMDOC}(V, d)$ 
2 for each  $c \in \mathbf{C}$ 
3 do  $\text{score}[c] \leftarrow \log \text{prior}[c]$ 
4   for each  $t \in V$ 
5   do if  $t \in V_d$ 
6     then  $\text{score}[c] += \log \text{condprob}[t][c]$ 
7     else  $\text{score}[c] += \log(1 - \text{condprob}[t][c])$ 
8 return  $\arg \max_{c \in \mathbf{C}} \text{score}[c]$ 

```

$$d = \langle t_1, \dots, t_k, \dots, t_{n_d} \rangle, t_k \in V$$

$$d = \langle e_1, \dots, e_i, \dots, e_M \rangle,$$

$$e_i \in \{0, 1\}$$

$$\hat{P}(X = t | c)$$

$$\hat{P}(U_i = e | c)$$

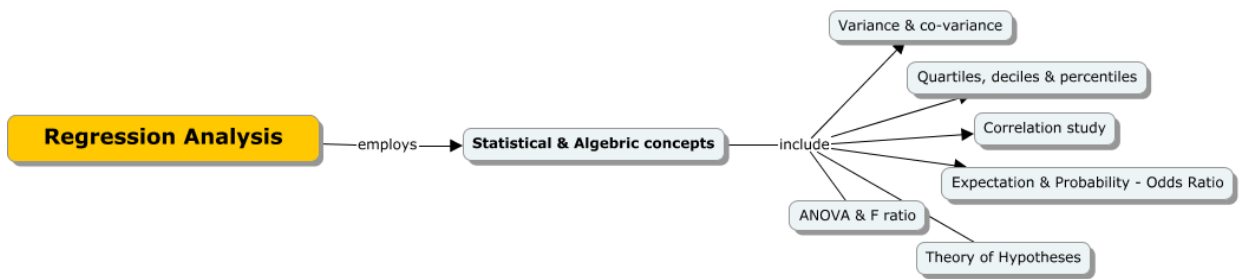
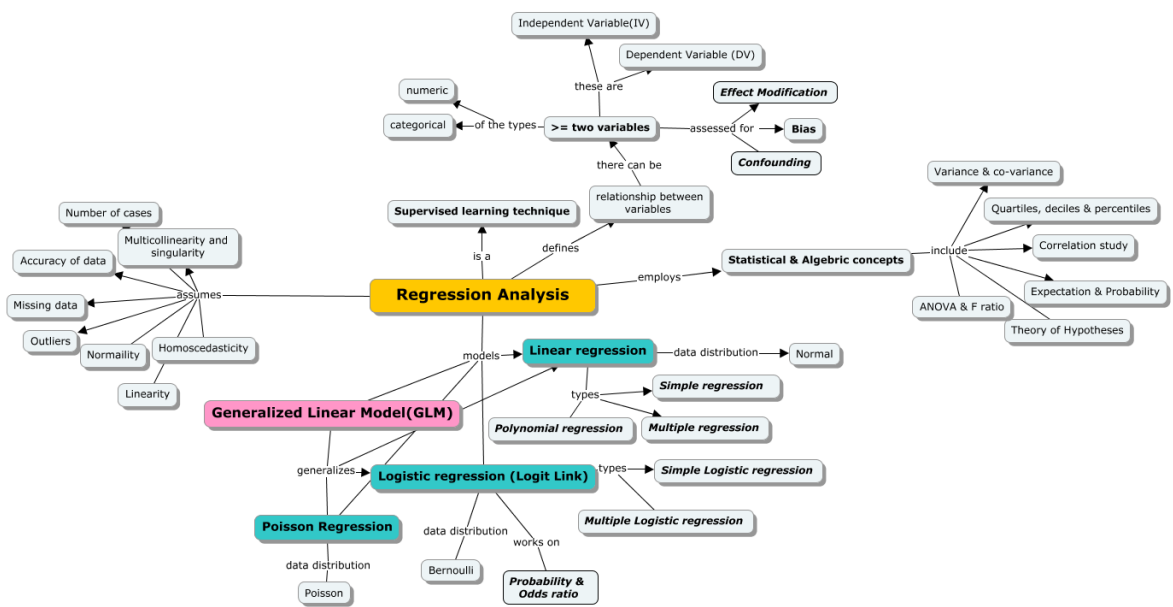
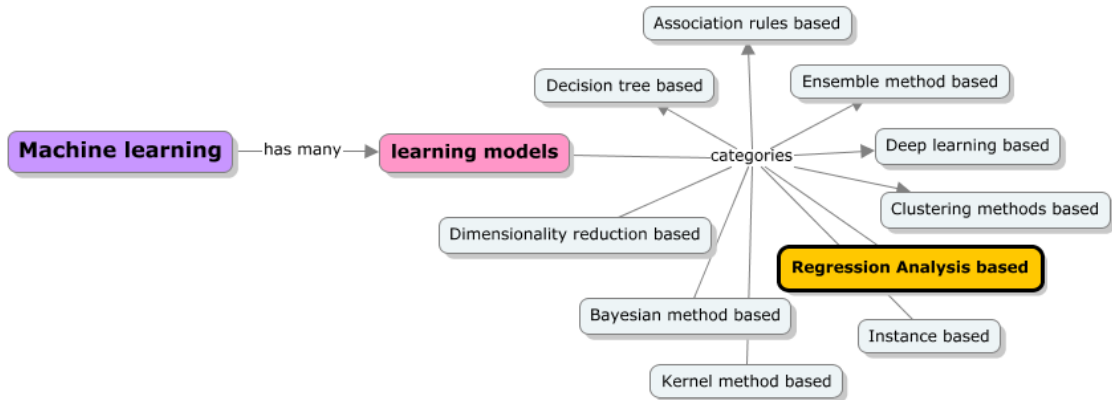
$$\hat{P}(c) \prod_{1 \leq k \leq n_d} \hat{P}(X = t_k | c)$$

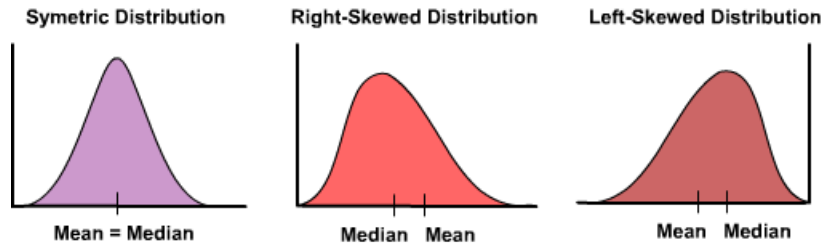
$$\hat{P}(c) \prod_{t_i \in V} \hat{P}(U_i = e_i | c)$$

$$\hat{P}(X = the | c) \approx 0.05$$

$$\hat{P}(U_{the} = 1 | c) \approx 1.0$$

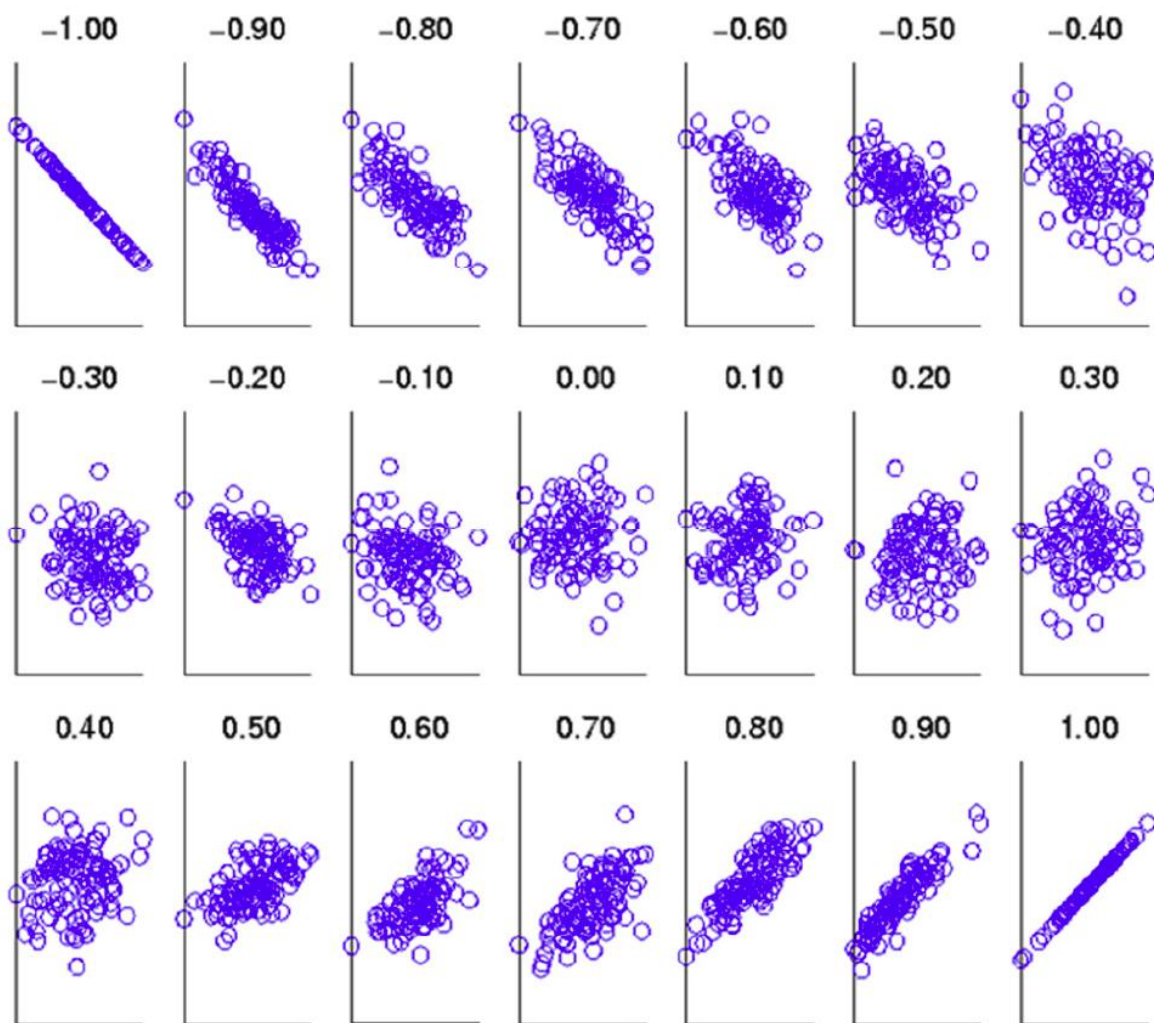
CHAPTER 10





$$\rho_{xy} = \frac{\sum_1^n (x_i - \bar{x})(y_i - \bar{y})}{n}$$

$$\text{Cor}(x, y) = \rho_{xy} / \sigma_x \sigma_y$$



$$\text{Mean} = \mu = \int_a^b xP(x)dx$$

$$\text{Variance} = \sigma^2 = \int_a^b x^2P(x)dx$$

Company A	Company B	Company C
498.3	243.7	250.15
515.25	245.75	250.25
506.4	242.7	250.25
504.8	244.65	253.55
536.95	250.95	236.8
512.55	227.4	219.1
525.65	240.3	206.5
538.95	243.8	216.45
510.45	235.25	217.9
503	238.35	215.3
500.75	231.4	218.15
496.45	228.15	217.55
492.3	219.7	215.15
496.6	218	205.4

Company A	Company B	Company C	A (returns)	B (returns)	C (returns)
498.3	243.7	250.15			
515.25	245.75	250.25	0.0340157	0.008412	0.0004
506.4	242.7	250.25	-0.017176	-0.012411	0
504.8	244.65	253.55	-0.00316	0.0080346	0.013187
536.95	250.95	236.8	0.0636886	0.0257511	-0.06606
512.55	227.4	219.1	-0.045442	-0.0938434	-0.07475
525.65	240.3	206.5	0.0255585	0.0567282	-0.05751
538.95	243.8	216.45	0.025302	0.0145651	0.048184
510.45	235.25	217.9	-0.052881	-0.0350697	0.006699
503	238.35	215.3	-0.014595	0.0131775	-0.01193
500.75	231.4	218.15	-0.004473	-0.0291588	0.013237
496.45	228.15	217.55	-0.008587	-0.0140449	-0.00275
492.3	219.7	215.15	-0.008359	-0.037037	-0.01103
496.6	218	205.4	0.0087345	-0.0077378	-0.04532

	A (returns)	B (returns)	C (returns)
	0.0340157	0.008412	0.0004
	-0.017176	-0.012411	0
	-0.00316	0.0080346	0.013187
	0.0636886	0.0257511	-0.06606
	-0.045442	-0.0938434	-0.07475
	0.0255585	0.0567282	-0.05751
	0.025302	0.0145651	0.048184
	-0.052881	-0.0350697	0.006699
	-0.014595	0.0131775	-0.01193
	-0.004473	-0.0291588	0.013237
	-0.008587	-0.0140449	-0.00275
	-0.008359	-0.037037	-0.01103
	0.0087345	-0.0077378	-0.04532
Sum	0.0026265	-0.1026342	-0.18764
Mean	0.000202	-0.0078949	-0.01443

	A (returns)	B (returns)	C (returns)
	-0.05288	-0.09384	-0.07475
	-0.04544	-0.03704	-0.06606
	-0.01718	-0.03507	-0.05751
	-0.01459	-0.02916	-0.04532
	-0.00859	-0.01404	-0.01193
	-0.00836	-0.01241	-0.01103
Median	-0.00447	-0.00774	-0.00275
	-0.00316	0.008035	0
	0.008735	0.008412	0.0004
	0.025302	0.013177	0.006699
	0.025558	0.014565	0.013187
	0.034016	0.025751	0.013237
	0.063689	0.056728	0.048184

A (returns)	B (returns)	C (returns)	Cov(A,B)	Cov(B,C)
-0.0528806	-0.093843	-0.07475	0.004562	0.005184
-0.0454419	-0.103097	0.00133	0.001505	
-0.0171761	-0.03507	-0.05751	0.000472	0.001171
-0.014595	-0.029159	-0.04532	0.000315	0.000657
-0.0085871	-0.014045	-0.01193	5.41E-05	-1.50E-05
-0.0083594	-0.012411	-0.01103	3.87E-05	-1.50E-05
A (returns)	B (returns)	C (returns)	Cov(A,B)	Cov(B,C)
-0.0044732	-0.007738	-0.00275	-7.30E-07	-1.84E-06
-0.0031596	0.0080346	0	-5.40E-05	0.00023
0.0087345	0.008412	0.0004	0.000139	0.000242
0.025302	0.0131775	0.006699	0.000529	0.000445
0.0255585	0.0145651	0.013187	0.00057	0.00062
0.0340157	0.0257511	0.013237	0.001138	0.000931
0.0636886	0.0567282	0.048184	0.004103	0.004047
		Covariance	0.001015	0.001154
		Correlation	0.939734	0.942151

$$E\left(\sum_{i=1}^n a_i x_i\right) = \sum_{i=1}^n a_i E(x_i)$$

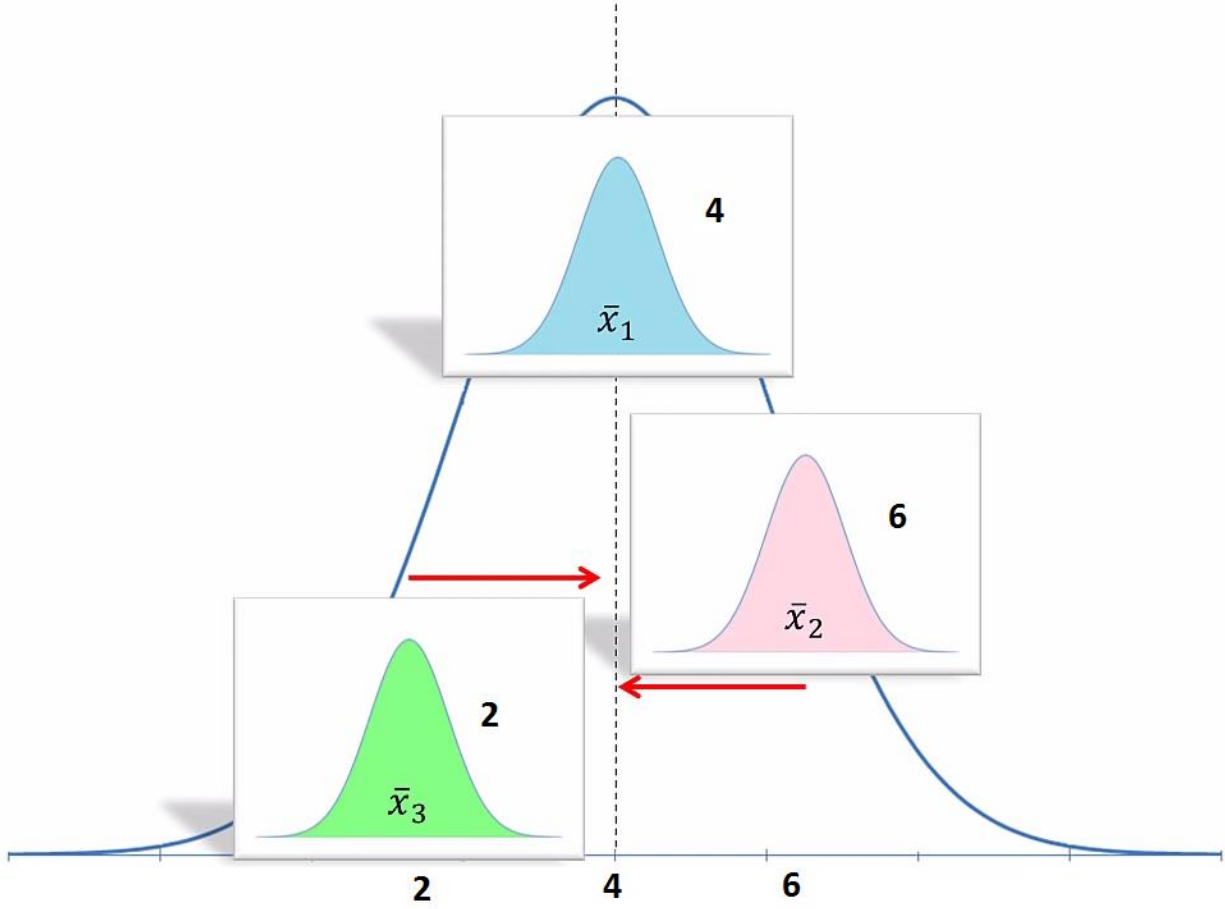
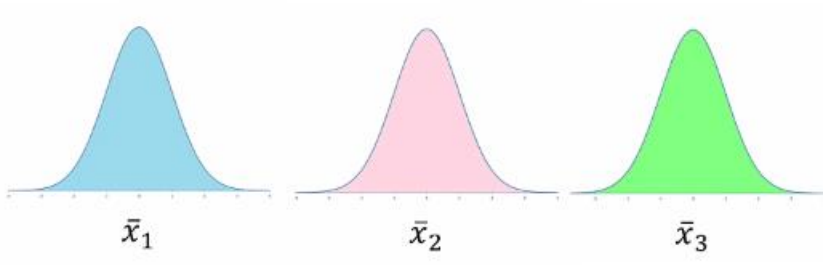
$$v \sum a_i X_i = \sum a_i^2 v(x_i) + 2 \sum \sum a_i a_j cov(x_i, x_j)$$

	Gold	IT	Bank
Returns	15	25	17
SD	5	15	10

	Gold	IT	Bank
Gold	1	-1	-0.5
IT	-1	1	0.5
Bank	-0.5	0.5	1

Portfolio1	0	1	0	25
Portfolio2	0	0.9	0.1	24.2
Portfolio3	0	0.8	0.2	23.4
Portfolio4	0	0.7	0.3	22.6
Portfolio5	0	0.6	0.4	21.8

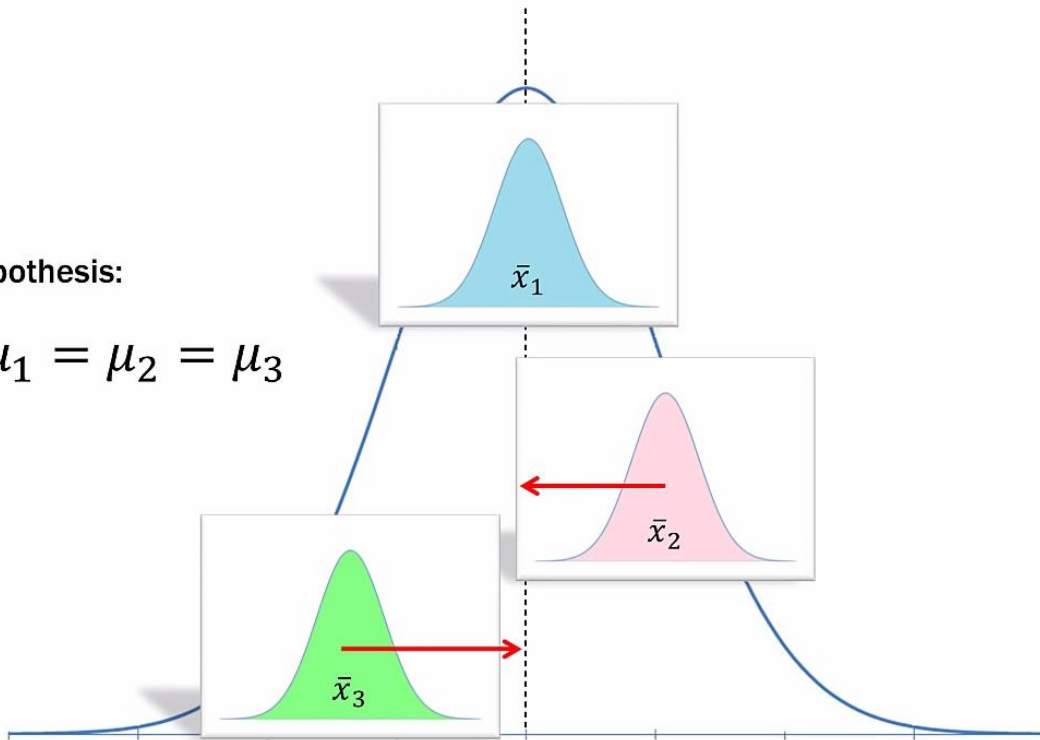
wg	wi	wb	SD
0	1	0	15
0	0.9	0.1	14.03
0	0.8	0.2	13.11
0	0.7	0.3	12.28
0	0.6	0.4	11.53



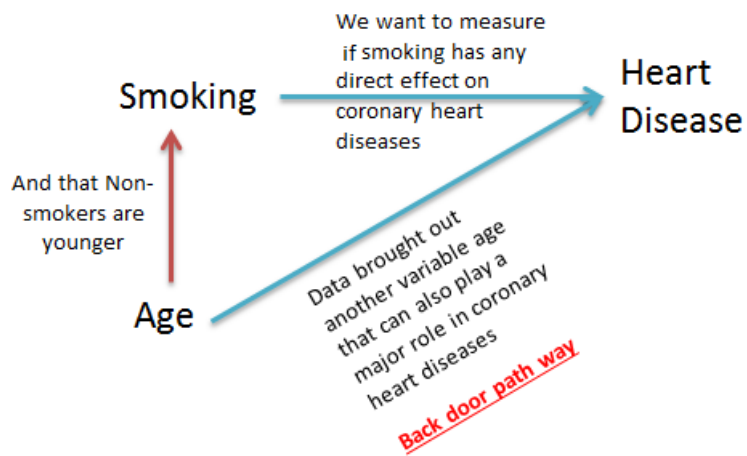
	Variation	DOF
Total	30	8(mn-1)
Within	6	6(m(n-1))
Between	24	2(m-1)

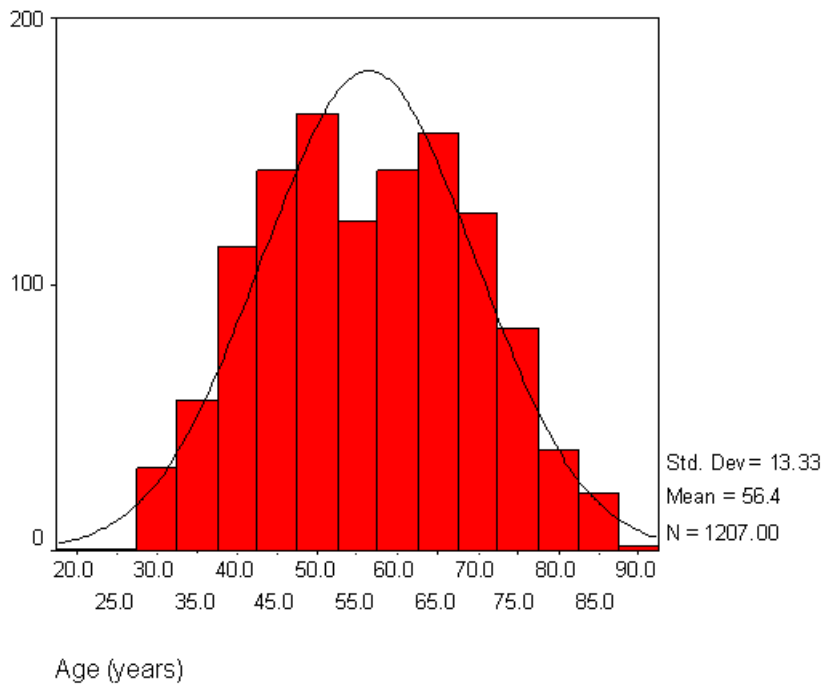
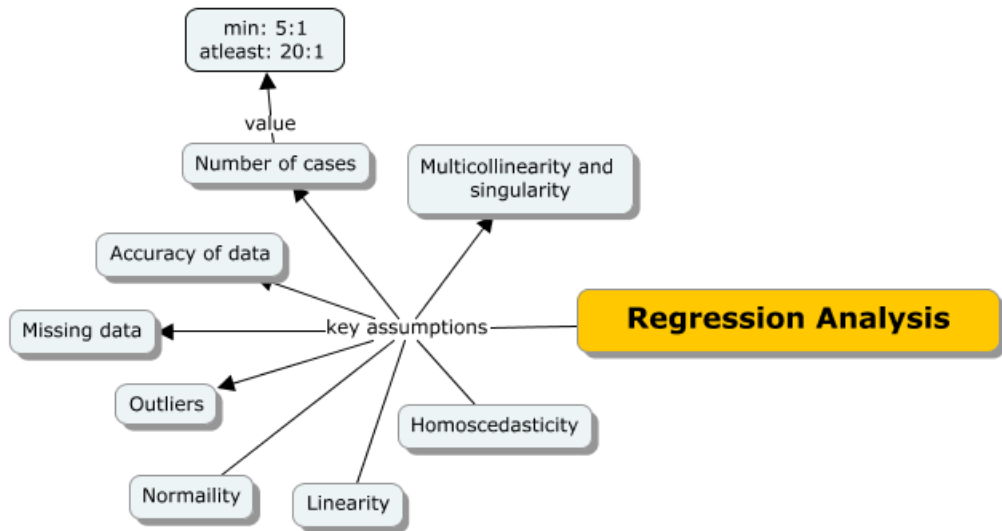
Null hypothesis:

$$H_0: \mu_1 = \mu_2 = \mu_3$$



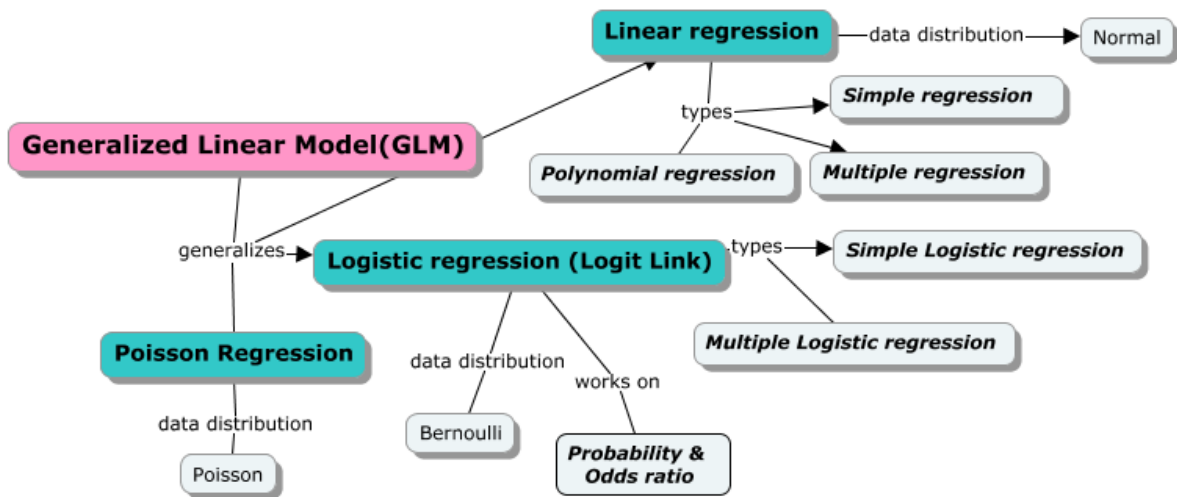
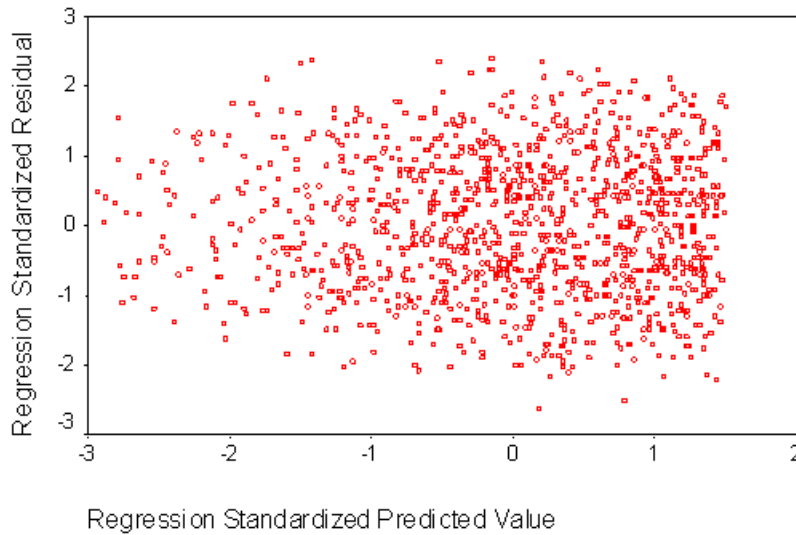
Direct path way



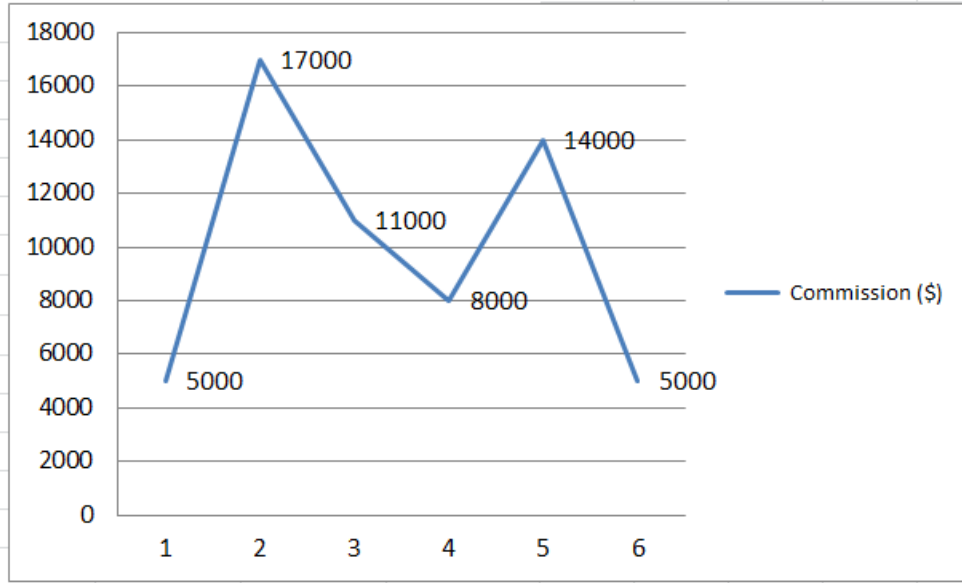


Scatterplot

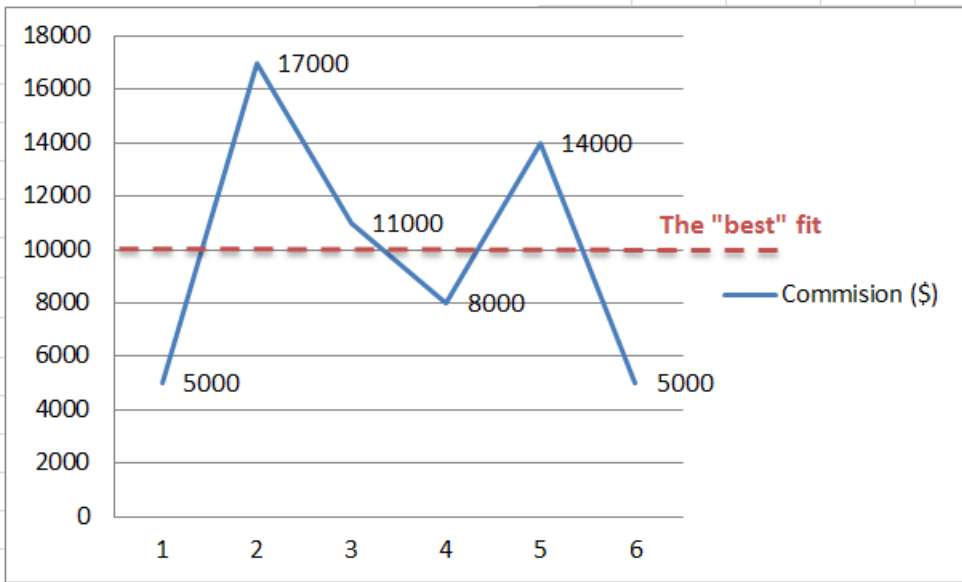
Dependent Variable: Age (years)



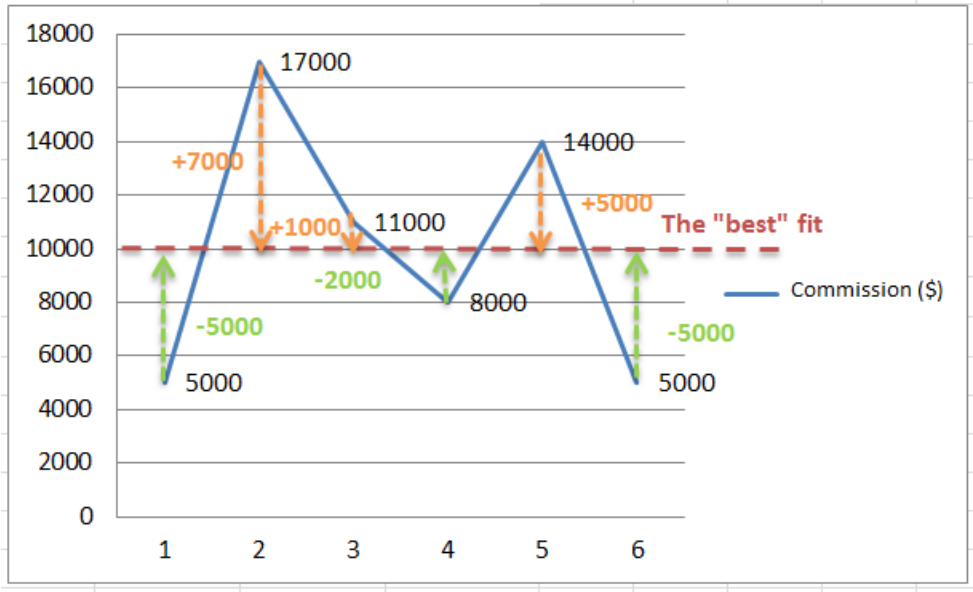
Trns #	Commision (\$)
1	5000
2	17000
3	11000
4	8000
5	14000
6	5000



Trns #	Commision (\$)
1	5000
2	17000
3	11000
4	8000
5	14000
6	5000
Mean	10000

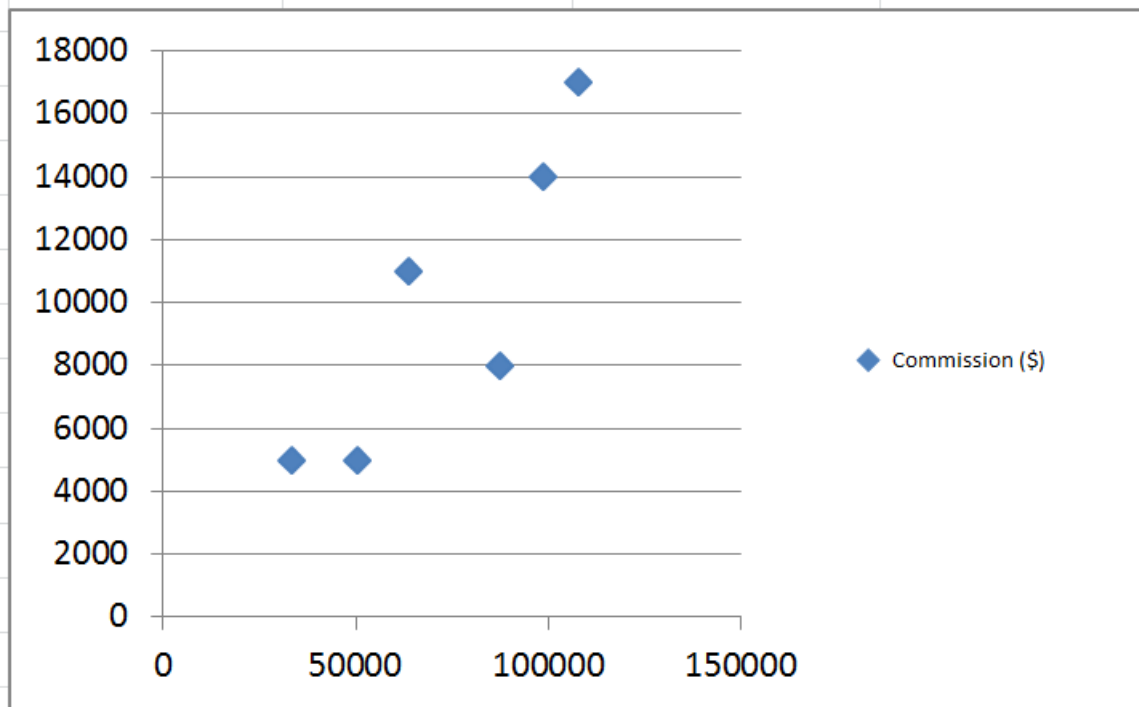


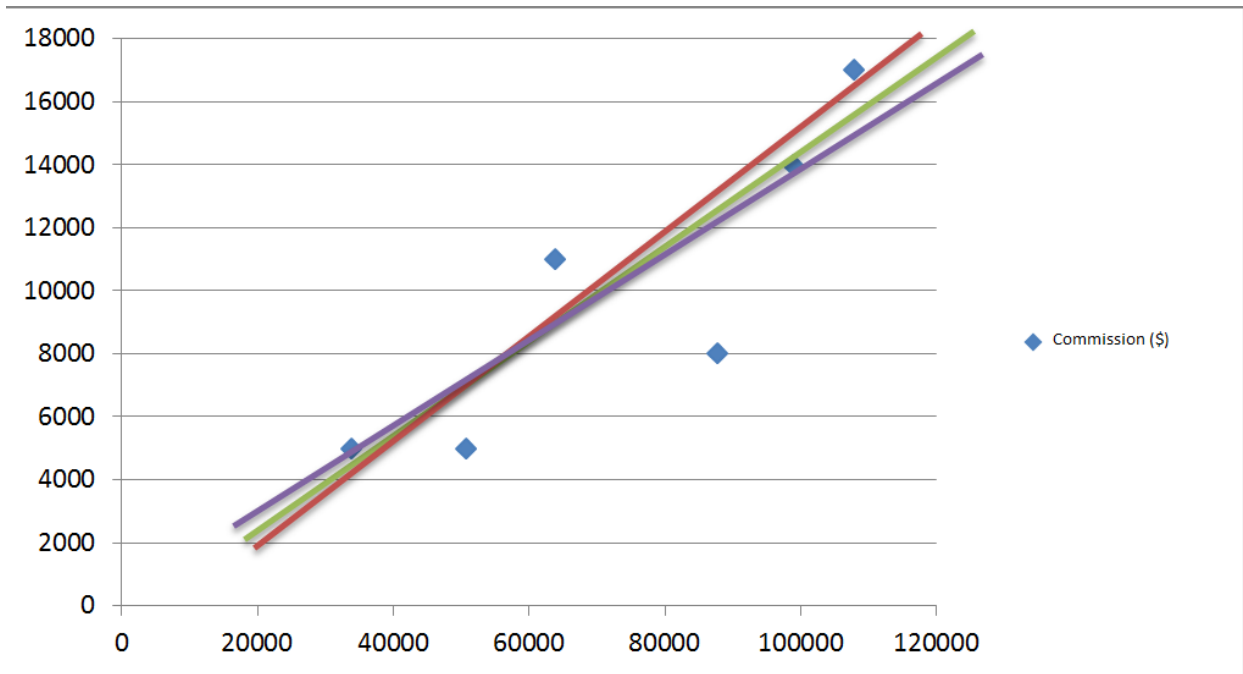
Trns #	Commision (\$)	error
1	5000	-5000
2	17000	7000
3	11000	1000
4	8000	-2000
5	14000	4000
6	5000	-5000
Total Error		0



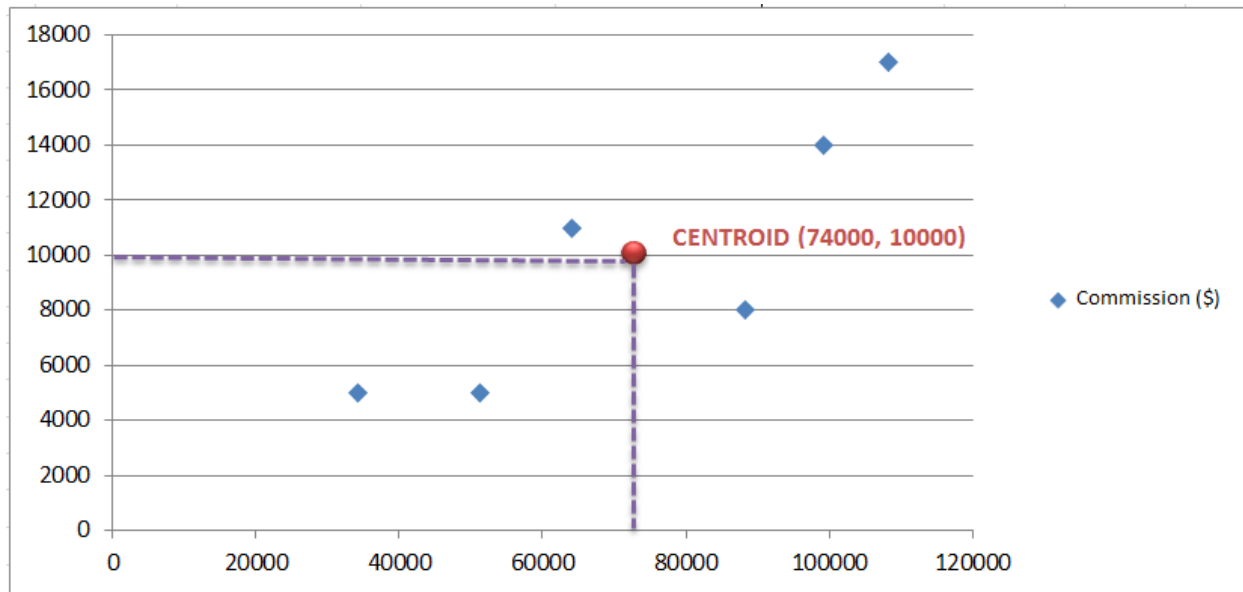
Trns #	Commision (\$)	error	error ^ 2
1	5000	-5000	25000000
2	17000	7000	49000000
3	11000	1000	1000000
4	8000	-2000	4000000
5	14000	4000	16000000
6	5000	-5000	25000000
		SSE	120000000

Transaction (\$)	Commision (\$)
34000	5000
108000	17000
64000	11000
88000	8000
99000	14000
51000	5000





	Transaction (\$)	Commision (\$)
	34000	5000
	108000	17000
	64000	11000
	88000	8000
	99000	14000
	51000	5000
Mean	74000	10000



$$\hat{y}_i = b_0 + b_1x_i$$

$$b_1 = \frac{\sum(x_i - \bar{x})(y_i - \bar{y})}{\sum(x_i - \bar{x})^2}$$

\bar{x} = mean of the independent variable x_i = value of independent variable

\bar{y} = mean of the dependent variable y_i = value of dependent variable

	Transaction (\$)	Commision (\$)	Txn Deviation	Comm Deviation	Dev Product	Square Txn Dev
	34000	5000	-40000	-5000	200000000	1600000000
	108000	17000	34000	7000	238000000	1156000000
	64000	11000	-10000	1000	-10000000	100000000
	88000	8000	14000	-2000	-28000000	196000000
	99000	14000	25000	4000	100000000	625000000
	51000	5000	-23000	-5000	115000000	529000000
Mean	74000	10000		Sum	615000000	4206000000

$$\hat{y}_i = b_0 + b_1x_i$$

$$b_0 = -0.8188$$

$$b_1 = 0.1462$$

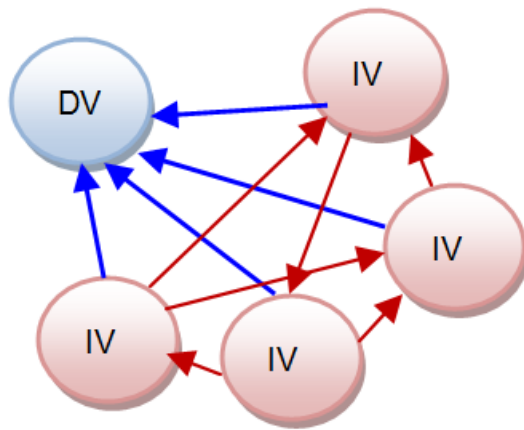
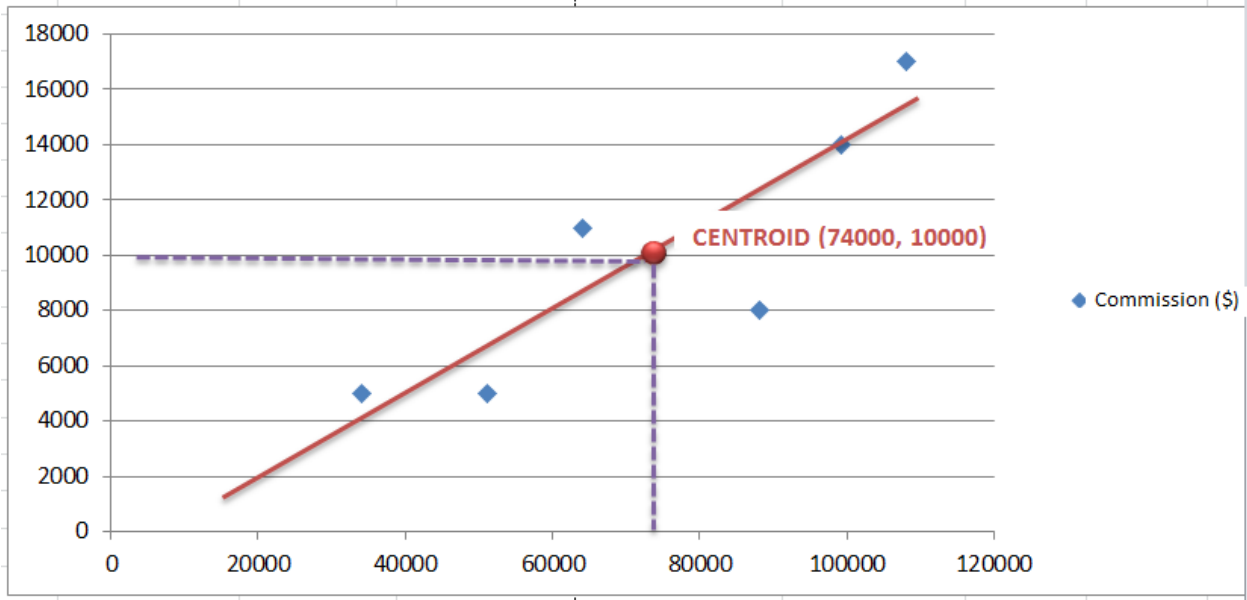
intercept

slope

$$\hat{y}_i = -0.8188 + 0.1462x$$

OR

$$\hat{y}_i = 0.1462x - 0.8188$$



$$y = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p + \epsilon$$

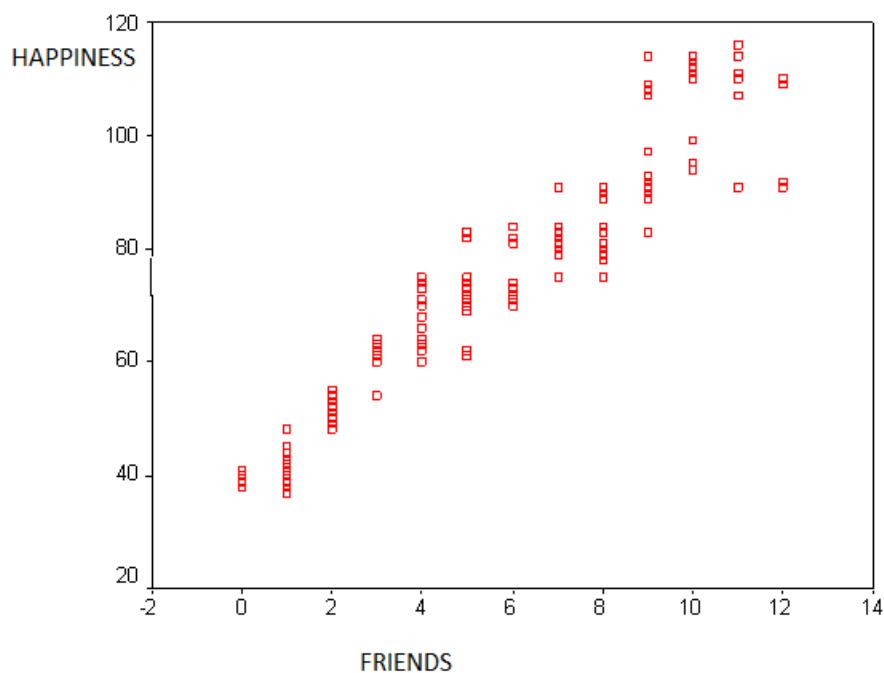
linear parameters error

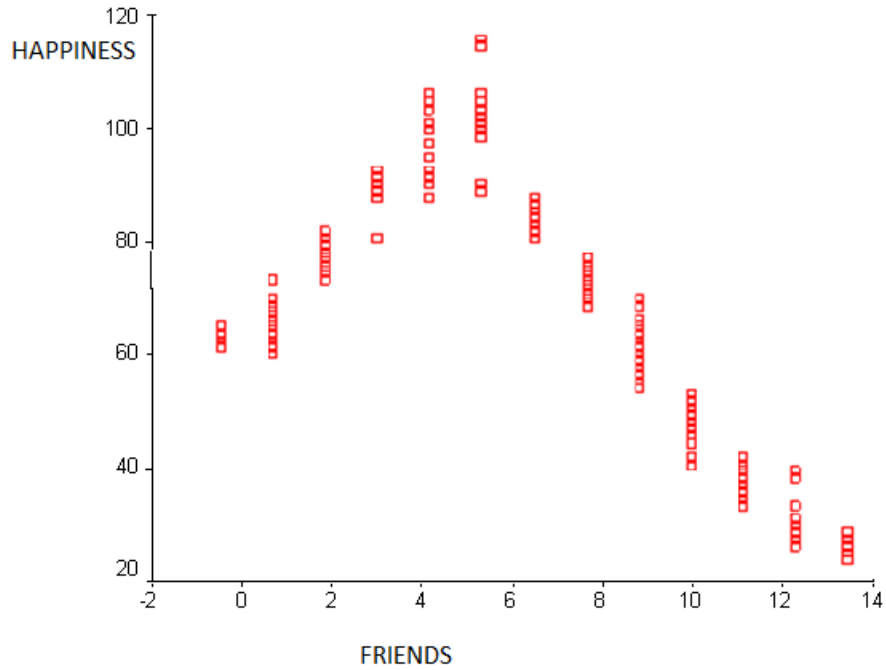
$$E(y) = \beta_0 + \beta_1 x_1 + \beta_2 x_2 + \cdots + \beta_p x_p$$

error term assumed to be zero

$$\hat{y} = b_0 + b_1 x_1 + b_2 x_2 + \cdots + b_p x_p$$

$b_0, b_1, b_2, \dots, b_p$ are the estimates of $\beta_0, \beta_1, \beta_2, \dots, \beta_p$
 \hat{y} = predicted value of the dependent variable



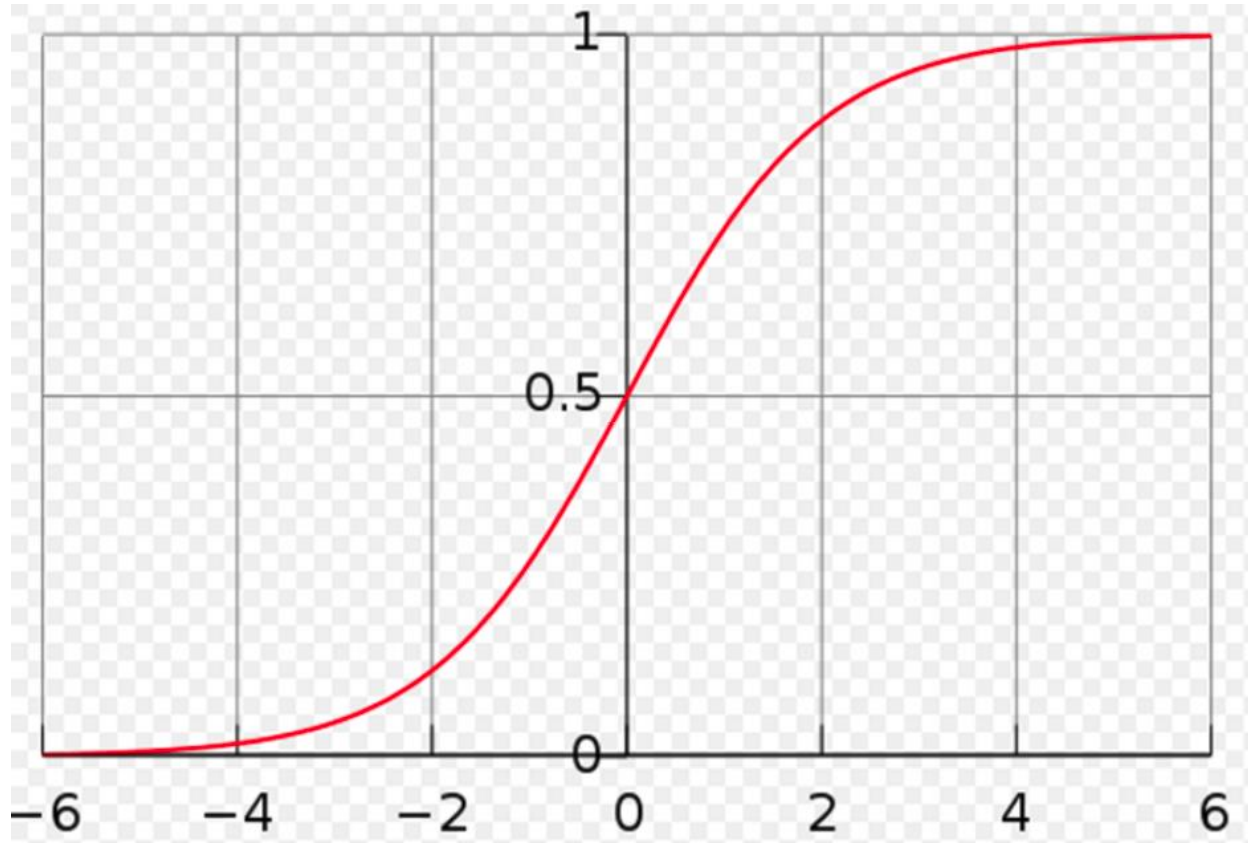


$$p_i = \frac{1}{1+e^{-(a+bx_i)}} \quad (\text{This is called the logistic response function})$$

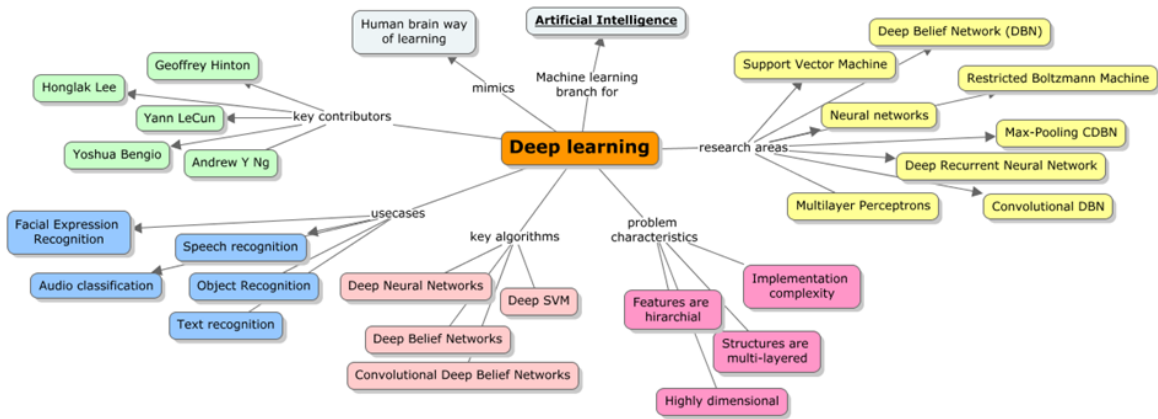
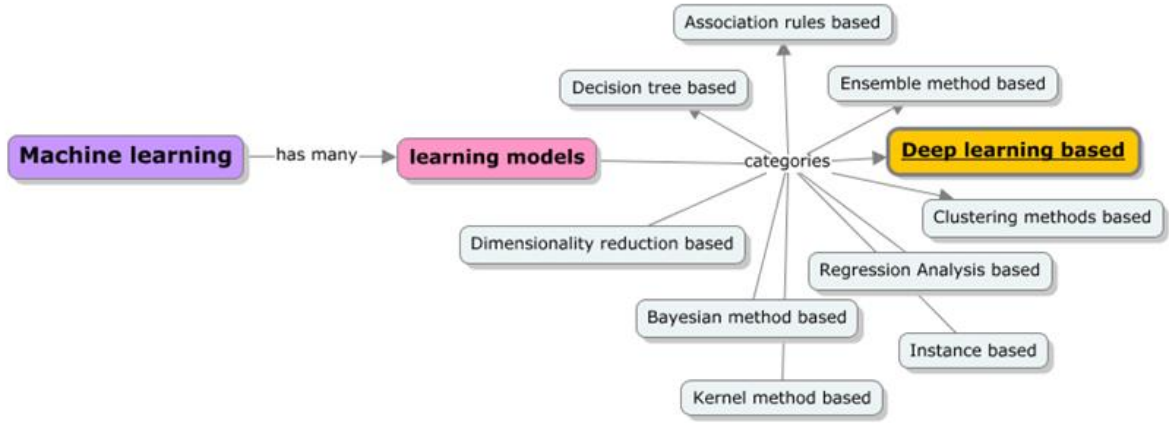
$$\log \frac{p_i}{1-p_i} = a+bx_i$$

$$\frac{p_i}{1-p_i} = e^{\hat{a}+\hat{b}x_i}$$

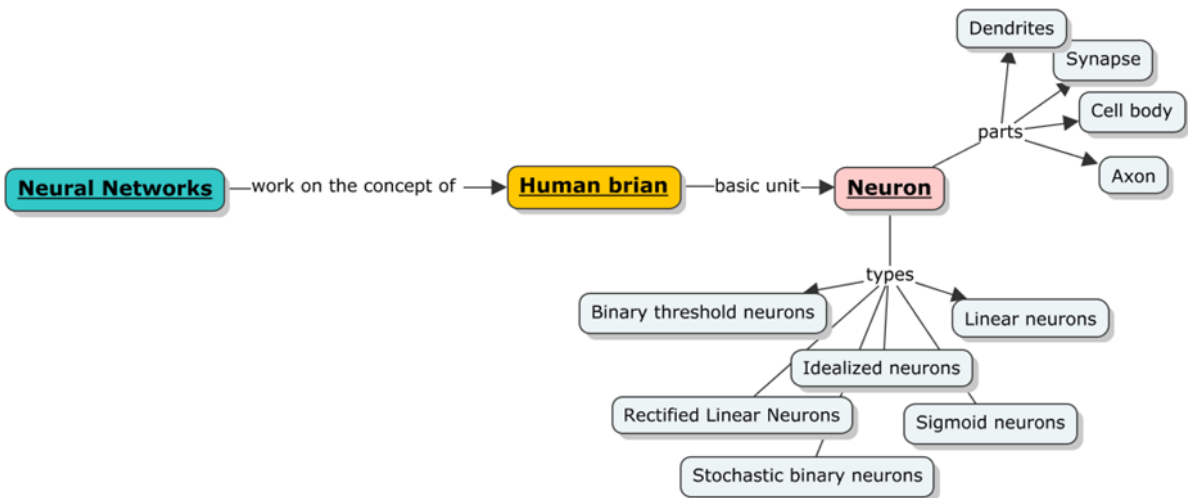
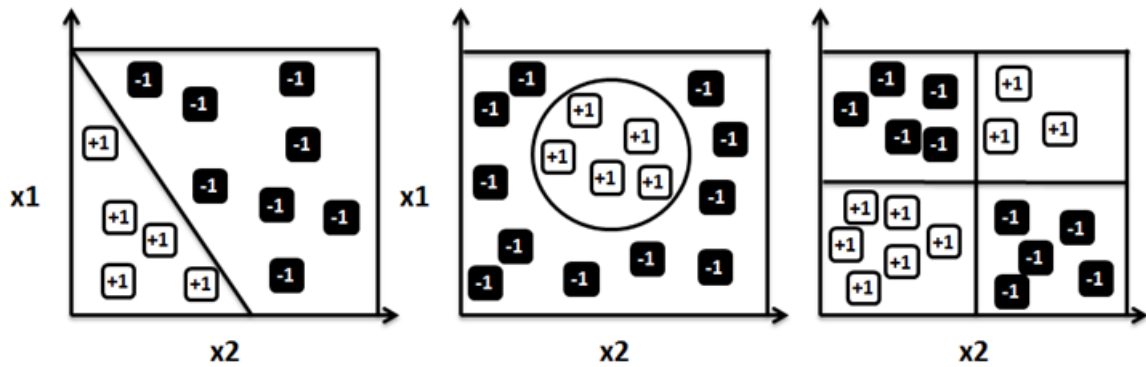
$$\hat{p}_i = \frac{e^{\hat{a}+\hat{b}x_i}}{1+e^{\hat{a}+\hat{b}x_i}}$$

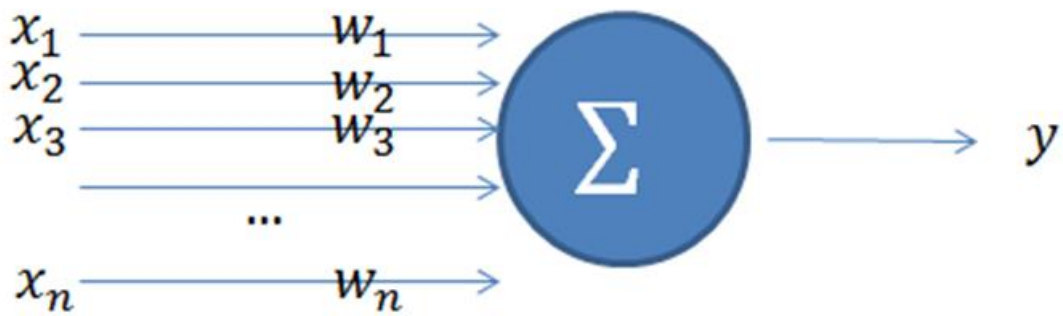
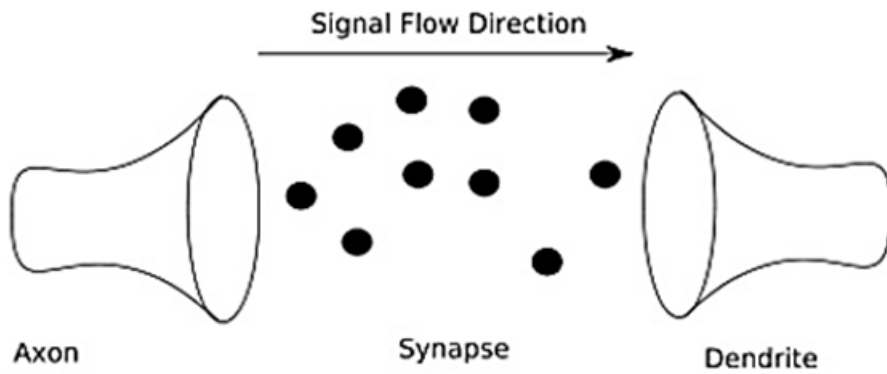
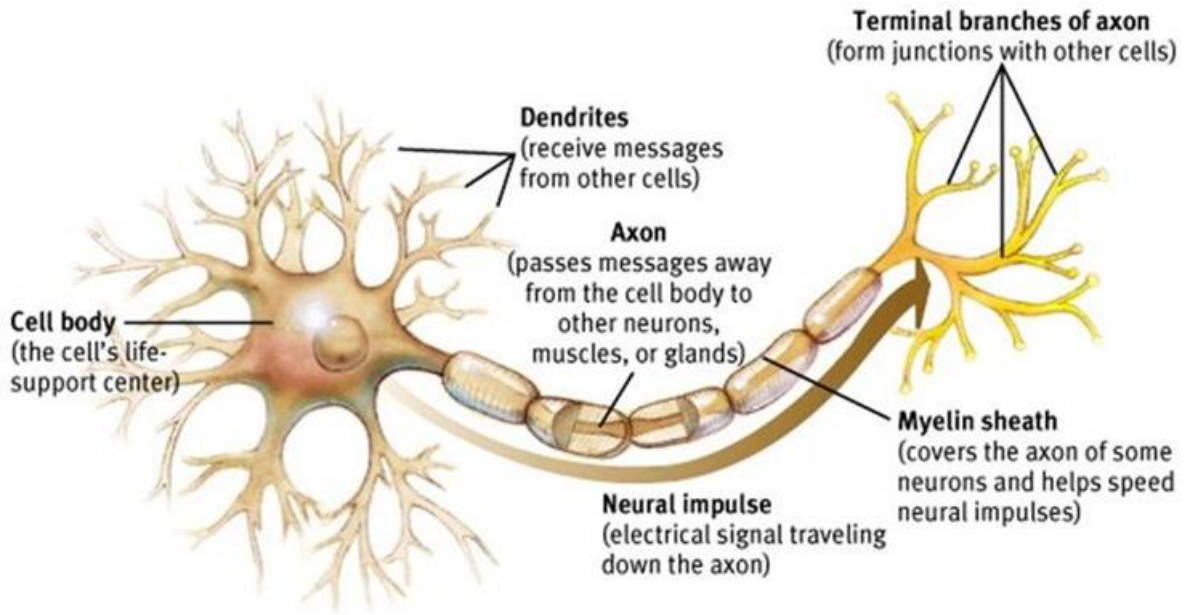


CHAPTER 11

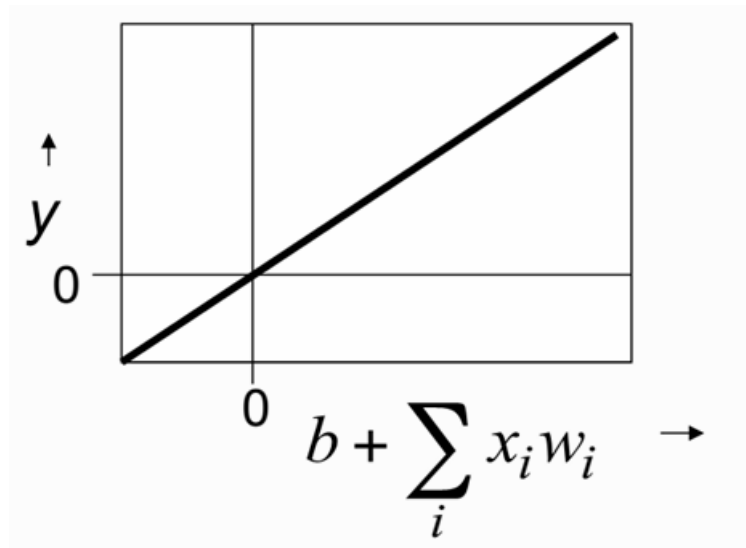


287635



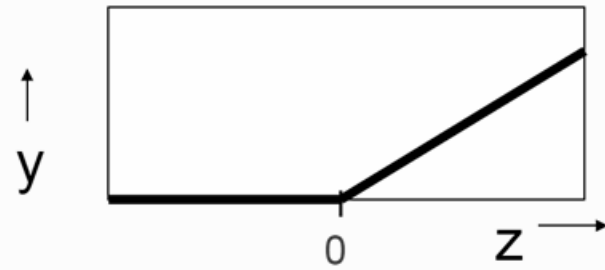


$$y = b + \sum_{i=1}^n (w_i x_i)$$



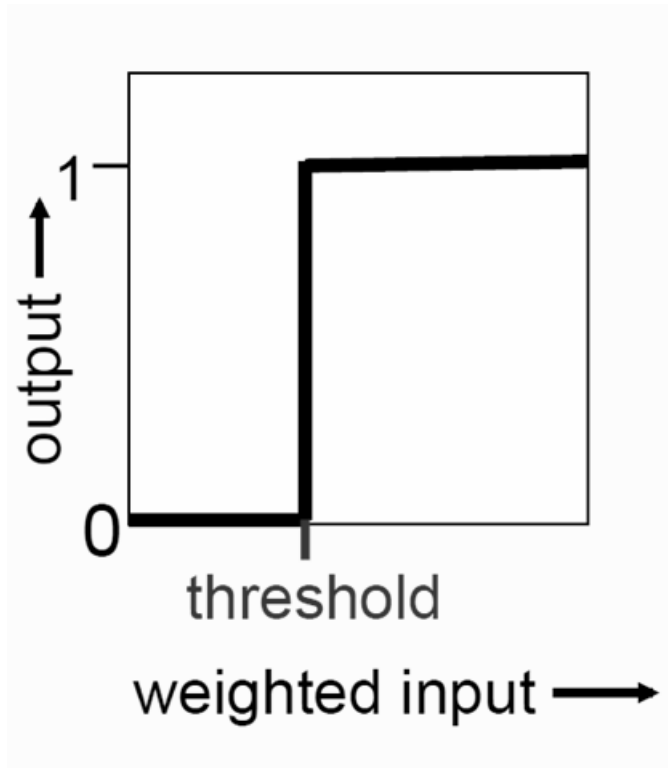
$$z = b + \sum_i x_i w_i$$

$$y = \begin{cases} z & \text{if } z > 0 \\ 0 & \text{otherwise} \end{cases}$$



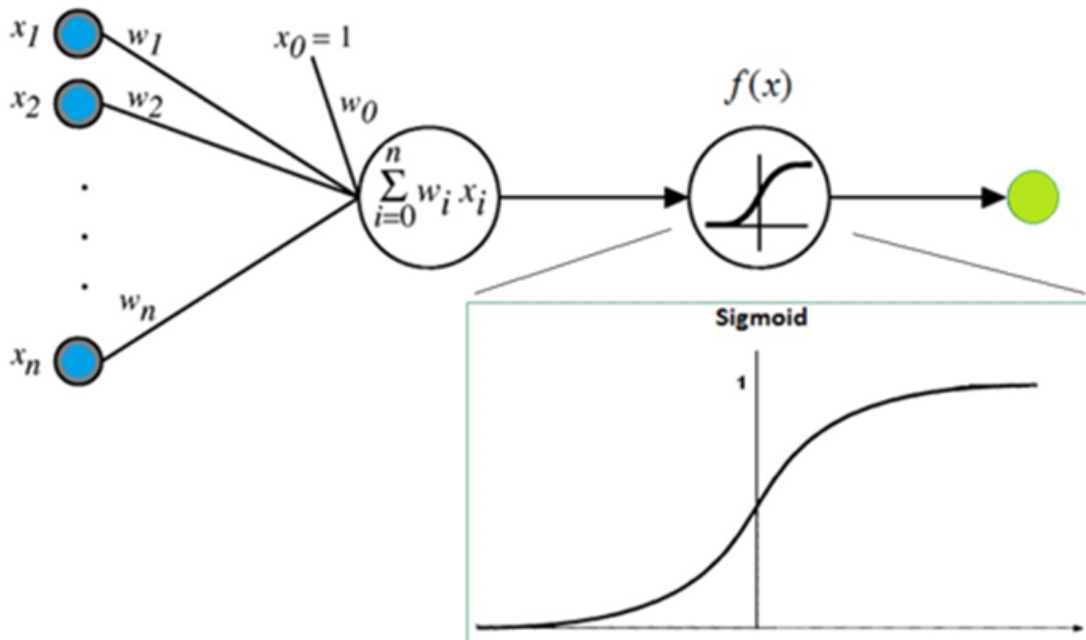
$$z = \sum_{i=1}^n (w_i x_i)$$

$$z = b + \sum_{i=1}^n (w_i x_i)$$



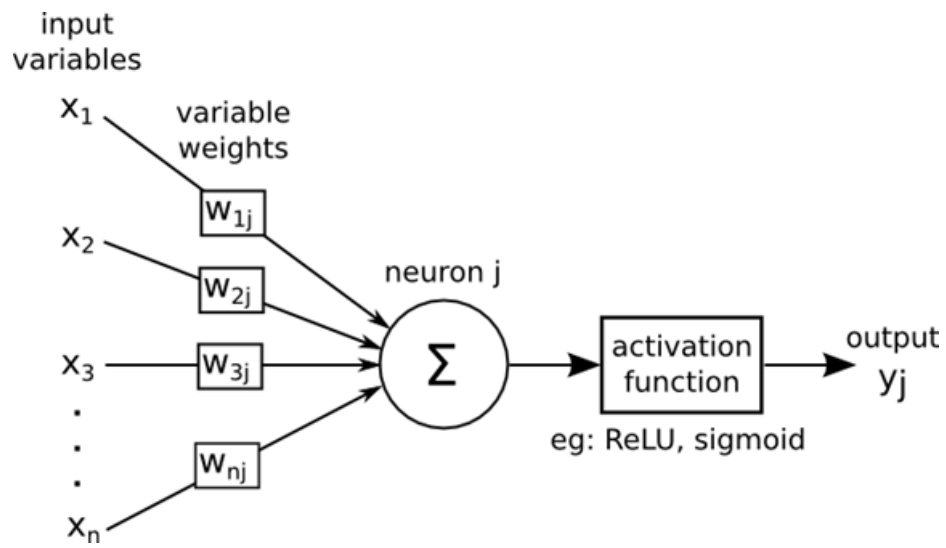
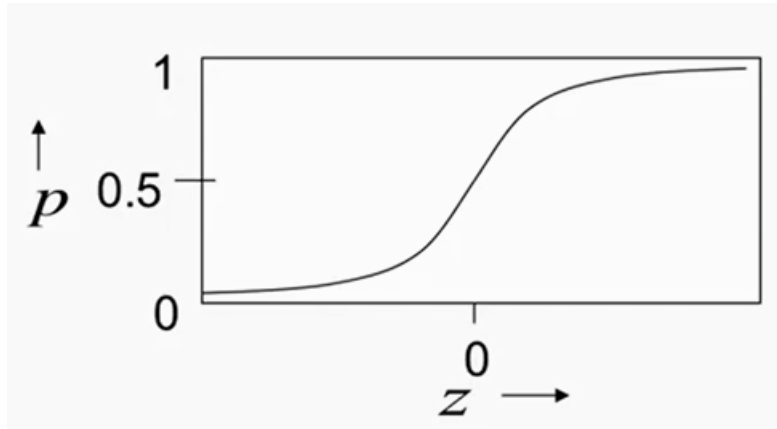
$$z = b + \sum_{i=1}^n (w_i x_i)$$

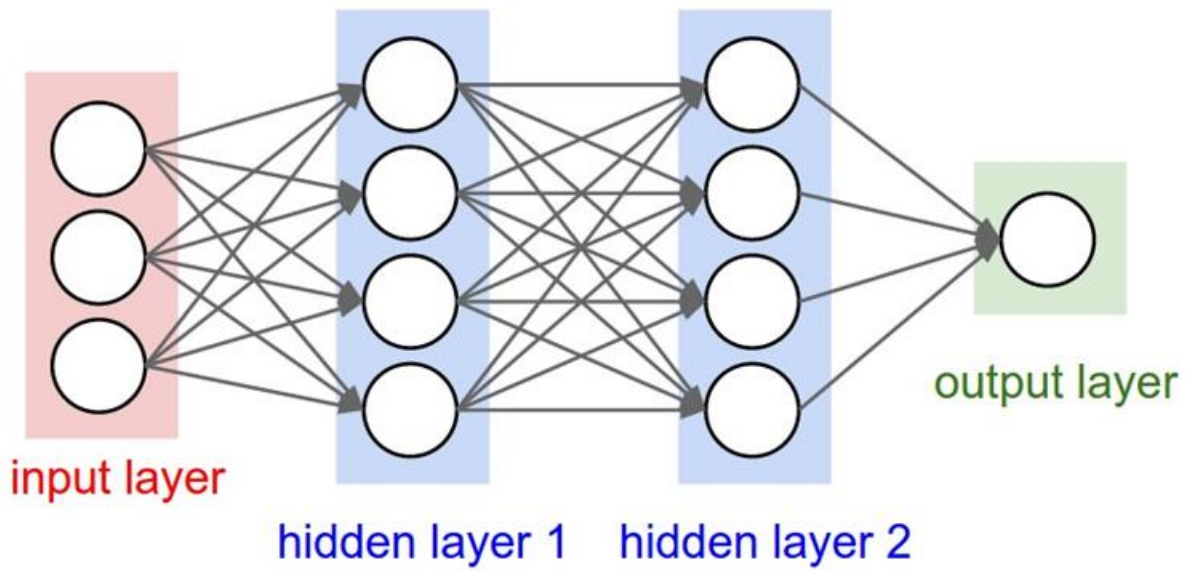
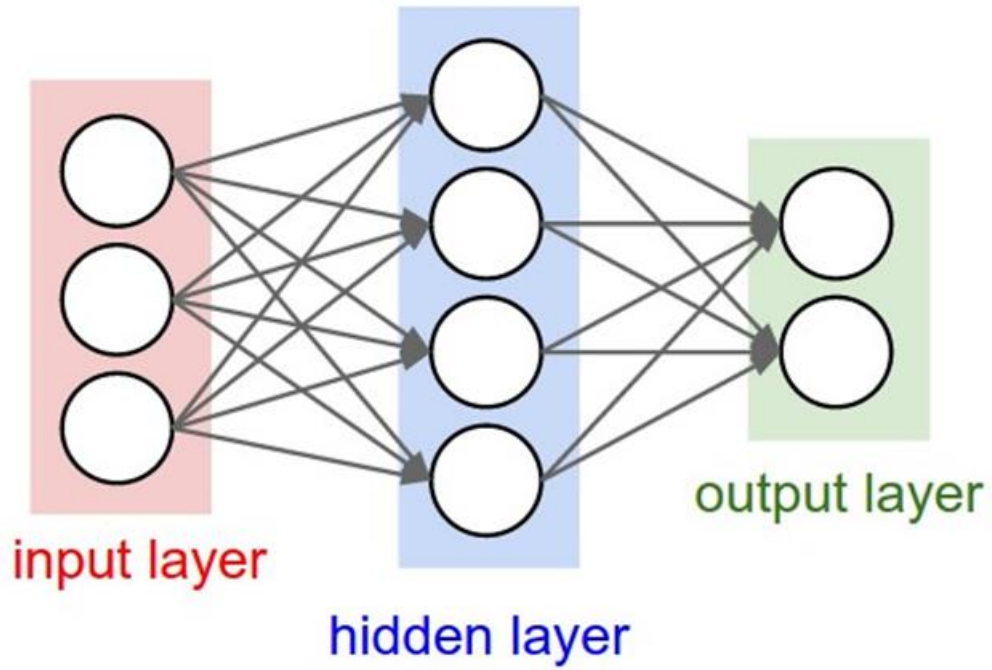
$$y = \frac{1}{1 + e^{-z}}$$

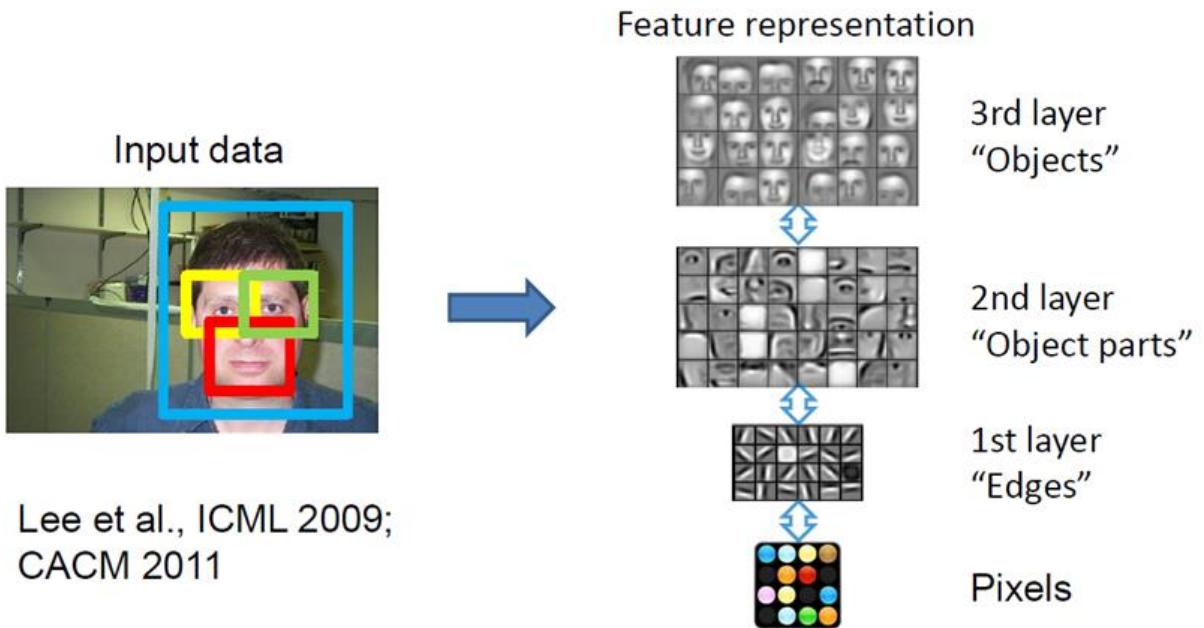
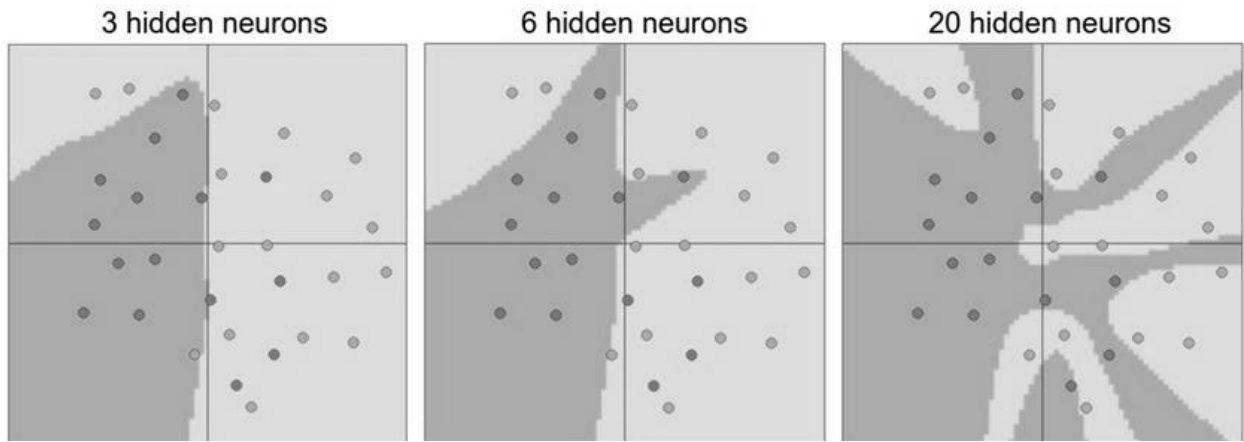


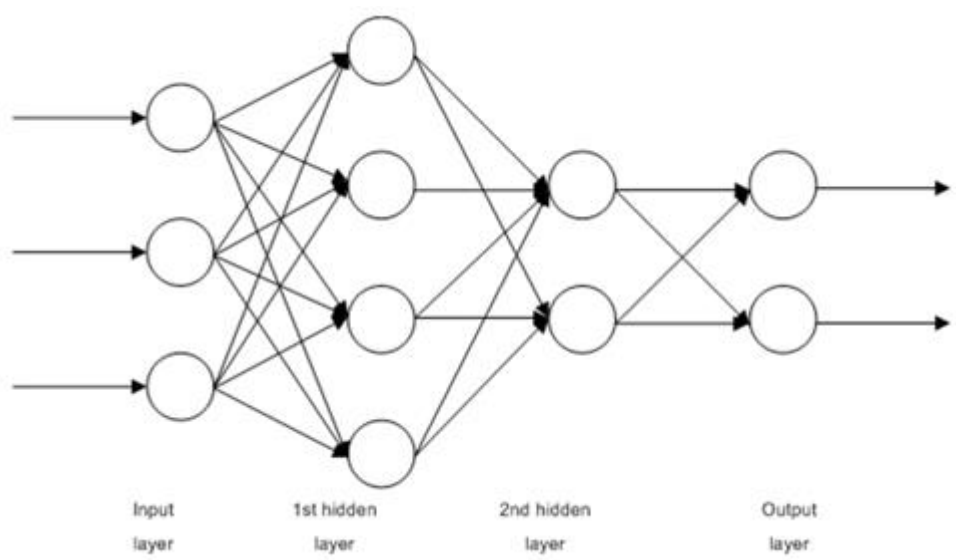
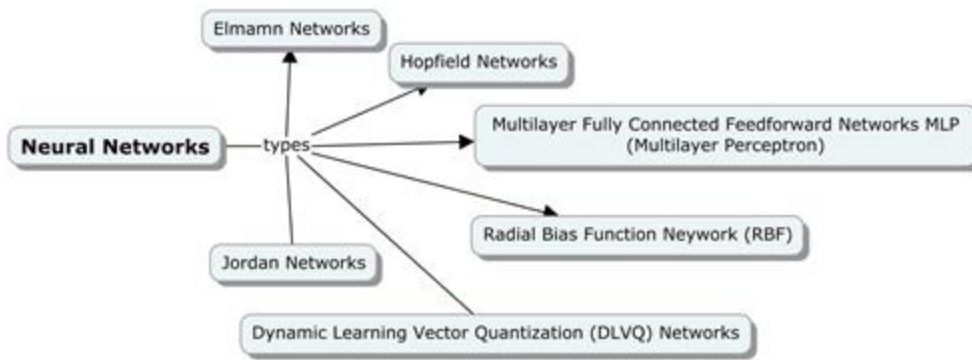
$$z = b + \sum_{i=1}^n (w_i x_i)$$

$$p(s=1) = \frac{1}{1 + e^{-z}}$$





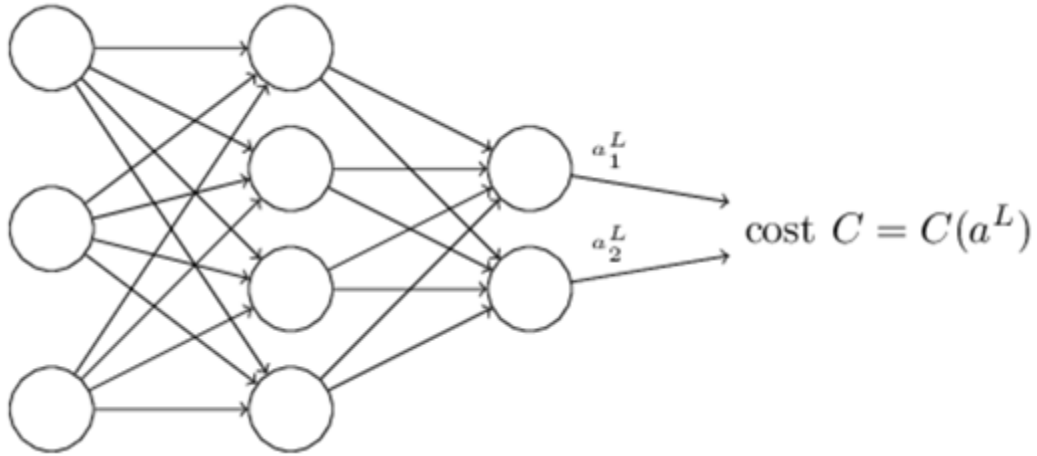




$$C = \frac{1}{2n} \sum_x \|y(x) - a^L(x)\|^2$$

$$C = \frac{1}{n} \sum_x C_x$$

$$C_x = \frac{1}{2} \|y - a^L\|^2$$



$$C = \frac{1}{2} \|y - a^L\|^2 = \frac{1}{2} \sum_j (y_j - a_j^L)^2$$

$$\partial C / \partial w_{jk}^l \text{ and } \partial C / \partial b_j^l$$

$$\delta_j^l \equiv \frac{\partial C}{\partial z_j^l}$$

Measures how the activation function changes at the current position in the network

$$\delta_j^L = \left(\frac{\partial C}{\partial a_j^L} \right) \sigma'(z_j^L)$$

Measures how the cost function changes based on the jth activation output

$$\delta^l = ((w^{l+1})^T \delta^{l+1}) \odot \sigma'(z^l)$$

Transpose of the weight matrix at l+1, this moves the error backwards through the network

Hadamard product, this moves the error backward through the activation layer l

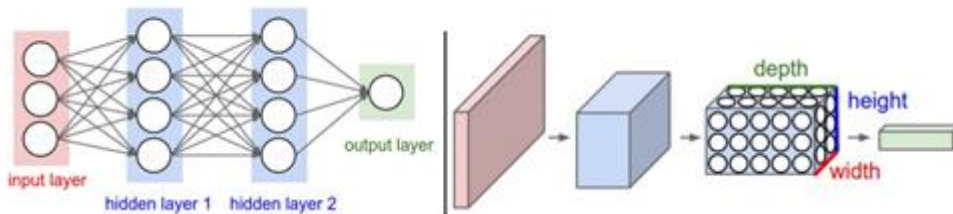
$$\begin{bmatrix} 1 \\ 2 \end{bmatrix} \odot \begin{bmatrix} 3 \\ 4 \end{bmatrix} = \begin{bmatrix} 1 * 3 \\ 2 * 4 \end{bmatrix} = \begin{bmatrix} 3 \\ 8 \end{bmatrix}$$

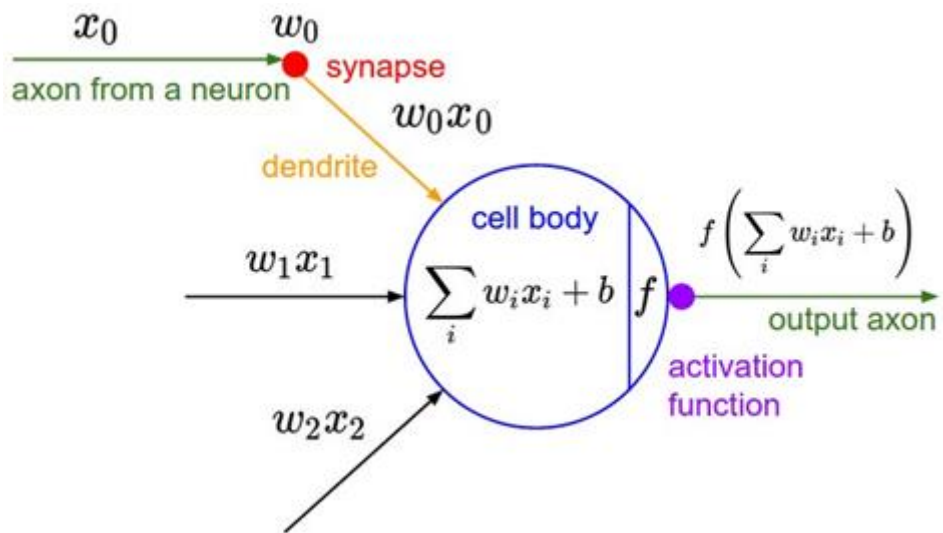
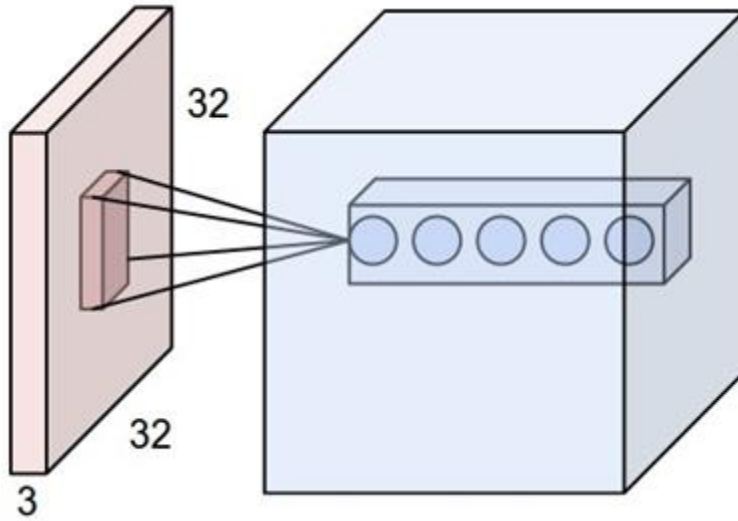
$$\frac{\partial C}{\partial b_j^l} = \delta_j^l$$

$$\frac{\partial C}{\partial b} = \delta$$

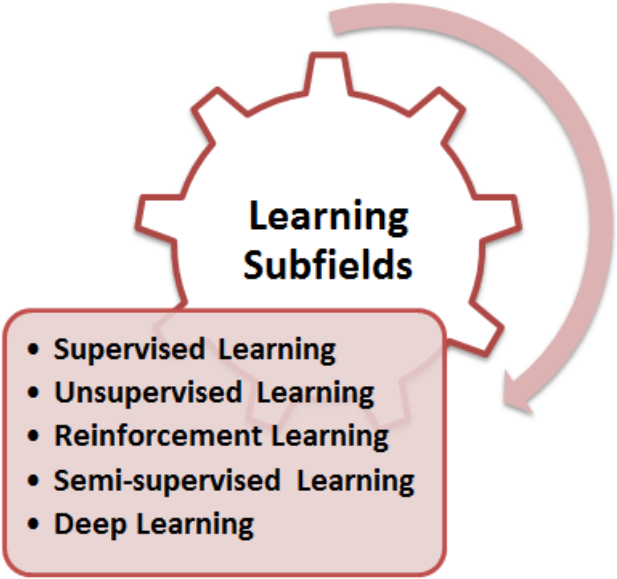
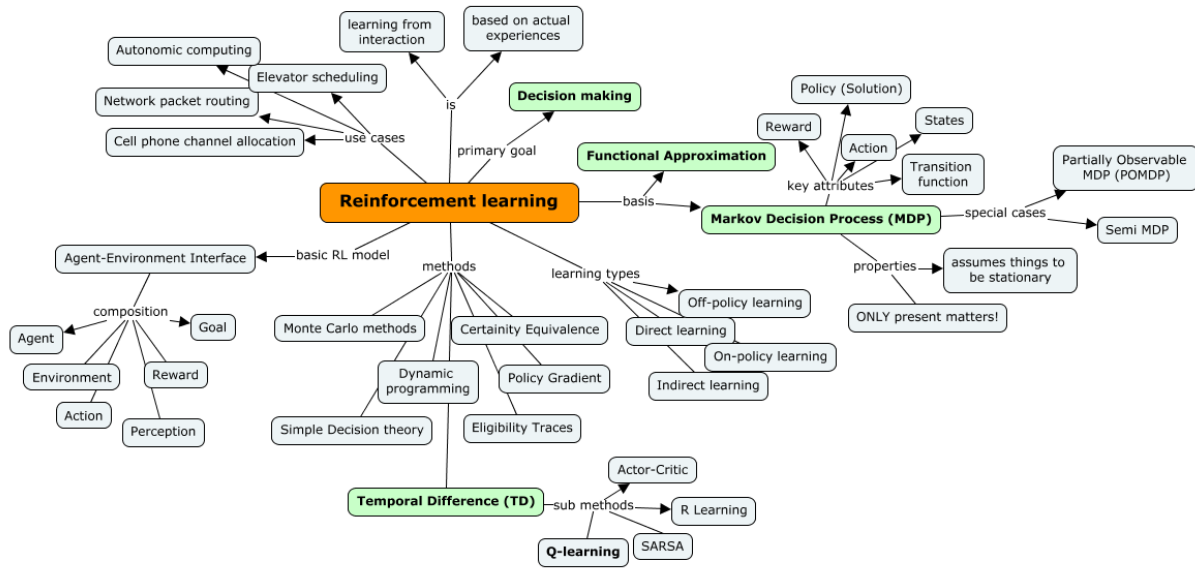
$$\frac{\partial C}{\partial w_{jk}^l} = a_k^{l-1} \delta_j^l$$

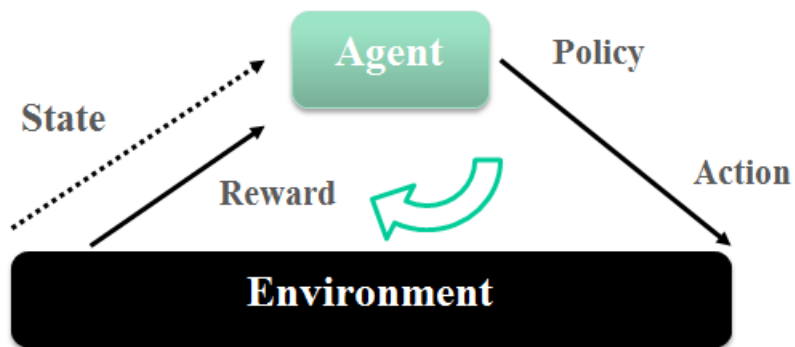
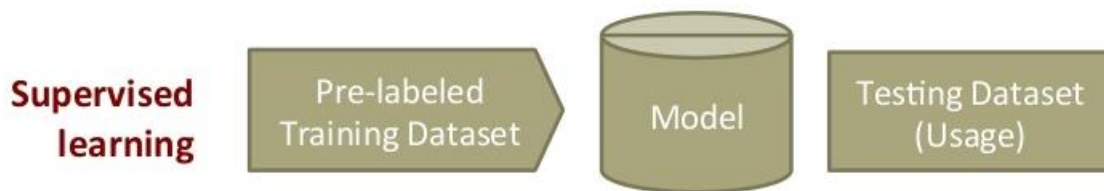
$$z^l = w^l a^{l-1} + b^l \text{ and } a^l = \sigma(z^l)$$



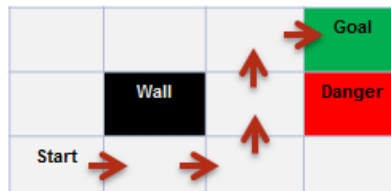
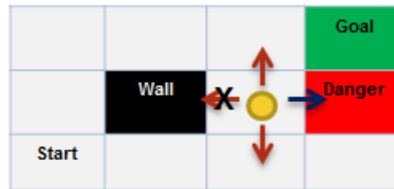
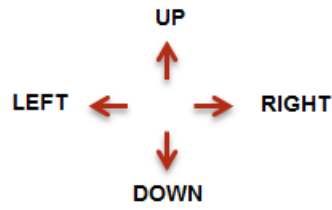
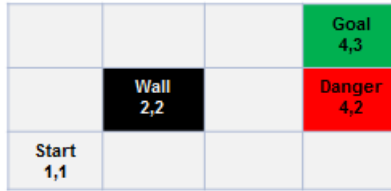


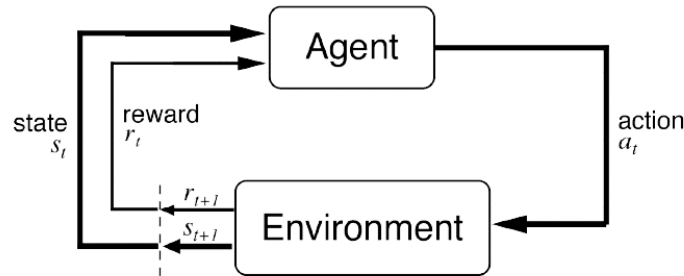
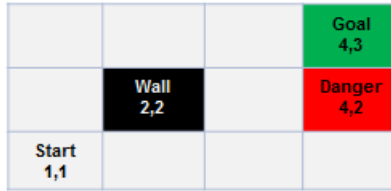
CHAPTER 12





1,3	2,3	3,3	4,3
1,2	2,2	3,2	4,2
1,1	2,1	3,1	4,1





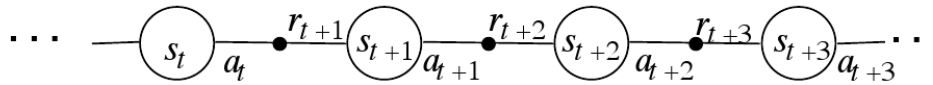
Agent and environment interact at discrete time steps : $t = 0, 1, 2, \dots$

Agent observes state at step t : $s_t \in S$

produces action at step t : $a_t \in A(s_t)$

gets resulting reward : $r_{t+1} \in \mathfrak{R}$

and resulting next state : s_{t+1}



$$\pi^* = \operatorname{argmax} \pi \left(E \left[\sum_{t=0}^{\infty} r^t R(s_t) / \pi \right] \right)$$

$$R(s) \neq U^\pi(s) = E \left[\sum_{t=0}^{\infty} r^t R(s_t) / \pi, s_0 = s \right]$$

$$\pi^* = \operatorname{argmax} \pi E \left[\sum_{s_1} T(s, a, s_1) U(s_1) \right]$$

$$U(s) = R(s) + \gamma \operatorname{max} \pi E \left[\sum_{s_1} T(s, a, s_1) U(s_1) \right]$$

$$\begin{aligned}
V^*(s) &= \max_{a \in \mathcal{A}(s)} Q^{\pi^*}(s, a) \\
&= \max_a E_{\pi^*} \left\{ R_t \mid s_t = s, a_t = a \right\} \\
&= \max_a E_{\pi^*} \left\{ \sum_{k=0}^{\infty} \gamma^k r_{t+k+1} \mid s_t = s, a_t = a \right\} \\
&= \max_a E_{\pi^*} \left\{ r_{t+1} + \gamma \sum_{k=0}^{\infty} \gamma^k r_{t+k+2} \mid s_t = s, a_t = a \right\} \\
&= \max_a E \{ r_{t+1} + \gamma V^*(s_{t+1}) \mid s_t = s, a_t = a \} \\
&= \max_{a \in \mathcal{A}(s)} \sum_{s'} \mathcal{P}_{ss'}^a \left[\mathcal{R}_{ss'}^a + \gamma V^*(s') \right].
\end{aligned}$$

1. Initialization

$V(s) \in \Re$ and $\pi(s) \in \mathcal{A}(s)$ arbitrarily for all $s \in \mathcal{S}$

2. Policy Evaluation

Repeat

$\Delta \leftarrow 0$

For each $s \in \mathcal{S}$:

$v \leftarrow V(s)$

$V(s) \leftarrow \sum_{s'} \mathcal{P}_{ss'}^{\pi(s)} \left[\mathcal{R}_{ss'}^{\pi(s)} + \gamma V(s') \right]$ =

$\Delta \leftarrow \max(\Delta, |v - V(s)|)$

until $\Delta < \theta$ (a small positive number)

3. Policy Improvement

policy-stable \leftarrow true

For each $s \in \mathcal{S}$:

$b \leftarrow \pi(s)$

$\pi(s) \leftarrow \arg \max_a \sum_{s'} \mathcal{P}_{ss'}^a \left[\mathcal{R}_{ss'}^a + \gamma V(s') \right]$

If $b \neq \pi(s)$, then *policy-stable* \leftarrow false

If *policy-stable*, then stop; else go to 2

Initialize V arbitrarily, e.g., $V(s) = 0$, for all $s \in \mathcal{S}^+$

Repeat

$$\Delta \leftarrow 0$$

For each $s \in \mathcal{S}$:

$$v \leftarrow V(s)$$

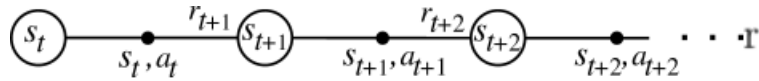
$$V(s) \leftarrow \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V(s')]$$

$$\Delta \leftarrow \max(\Delta, |v - V(s)|)$$

until $\Delta < \theta$ (a small positive number)

Output a deterministic policy, π , such that

$$\pi(s) = \arg \max_a \sum_{s'} \mathcal{P}_{ss'}^a [\mathcal{R}_{ss'}^a + \gamma V(s')]$$



$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_{t+1} + \gamma Q(s_{t+1}, a_{t+1}) - Q(s_t, a_t)].$$

Initialize $Q(s, a)$ arbitrarily

Repeat (for each episode):

Initialize s

Choose a from s using policy derived from Q (e.g., ϵ -greedy)

Repeat (for each step of episode):

Take action a , observe r, s'

Choose a' from s' using policy derived from Q (e.g., ϵ -greedy)

$$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma Q(s', a') - Q(s, a)]$$

$s \leftarrow s'; a \leftarrow a'$;

until s is terminal

$$Q(s_t, a_t) \leftarrow Q(s_t, a_t) + \alpha [r_{t+1} + \gamma \max_a Q(s_{t+1}, a) - Q(s_t, a_t)].$$

Initialize $Q(s, a)$ arbitrarily

Repeat (for each episode):

Initialize s

Repeat (for each step of episode):

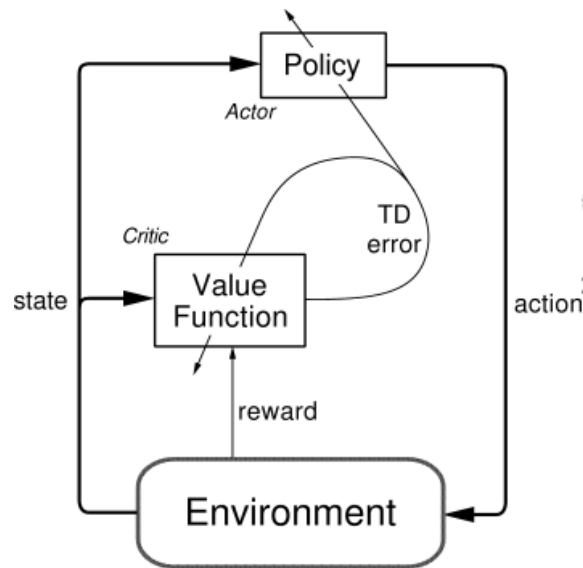
Choose a from s using policy derived from Q (e.g., ϵ -greedy)

Take action a , observe r, s'

$Q(s, a) \leftarrow Q(s, a) + \alpha [r + \gamma \max_{a'} Q(s', a') - Q(s, a)]$

$s \leftarrow s'$;

until s is terminal



Initialize ρ and $Q(s, a)$, for all s, a , arbitrarily

Repeat forever:

$s \leftarrow$ current state

Choose action a in s using behavior policy (e.g., ϵ -greedy)

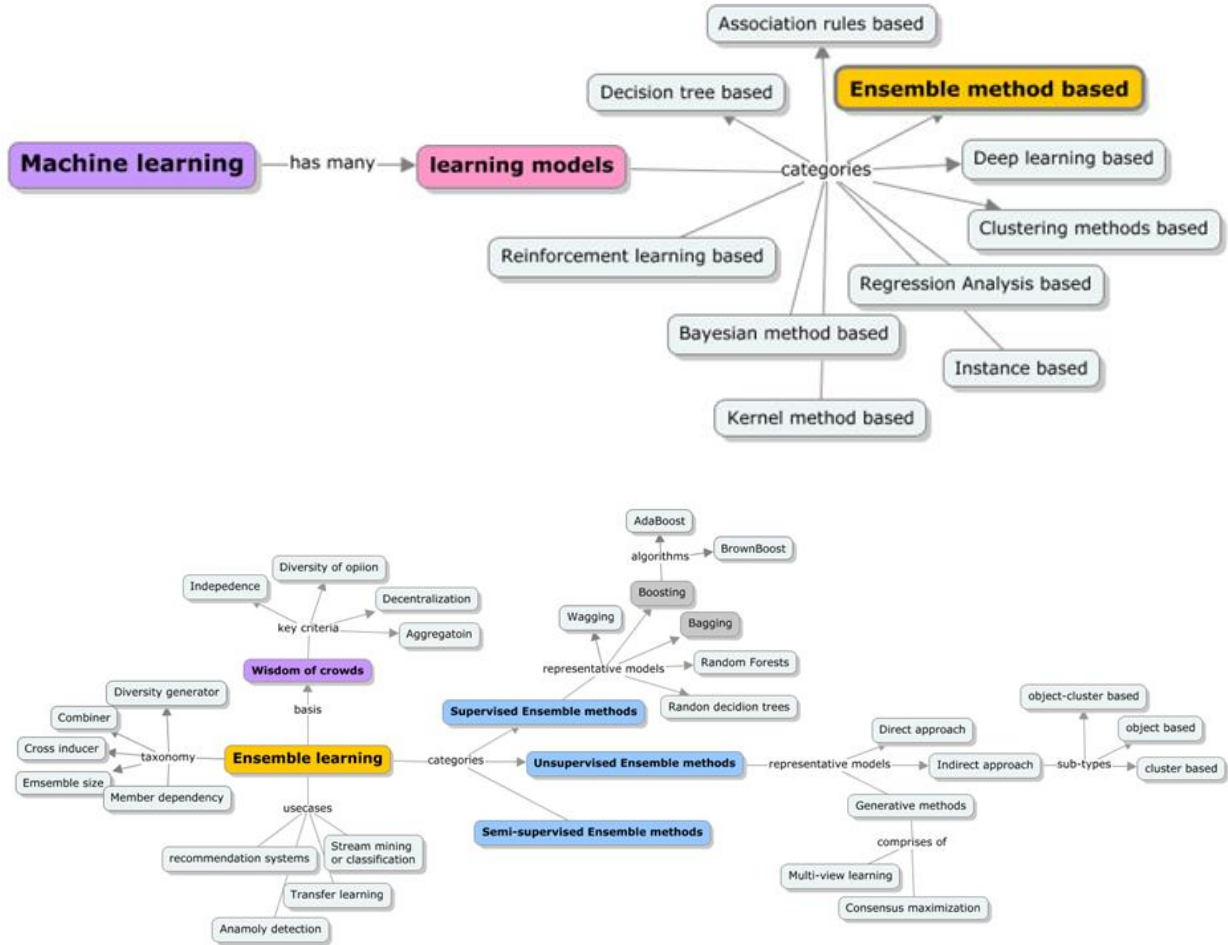
Take action a , observe r, s'

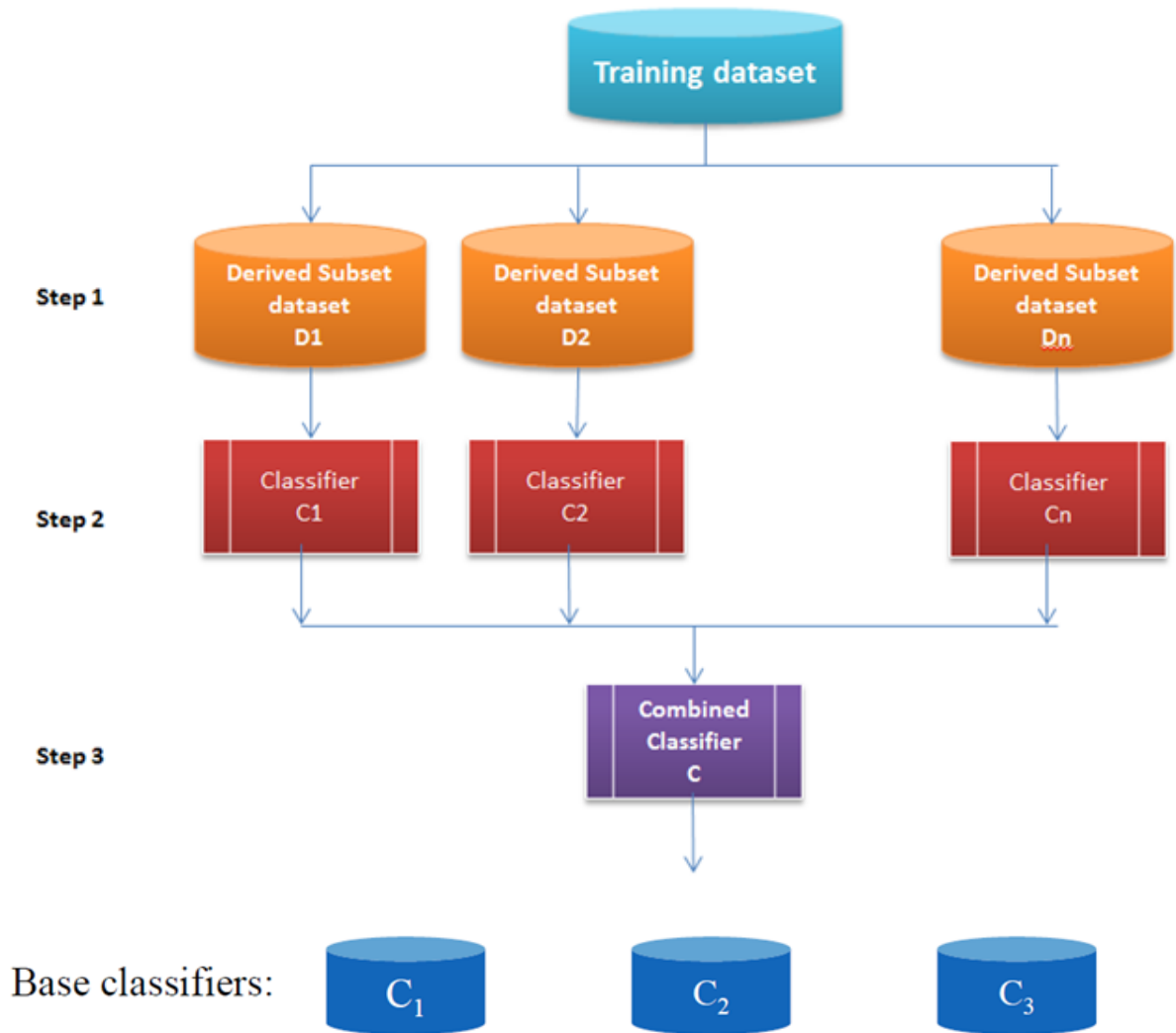
$Q(s, a) \leftarrow Q(s, a) + \alpha [r - \rho + \max_{a'} Q(s', a') - Q(s, a)]$

If $Q(s, a) = \max_a Q(s, a)$, then:

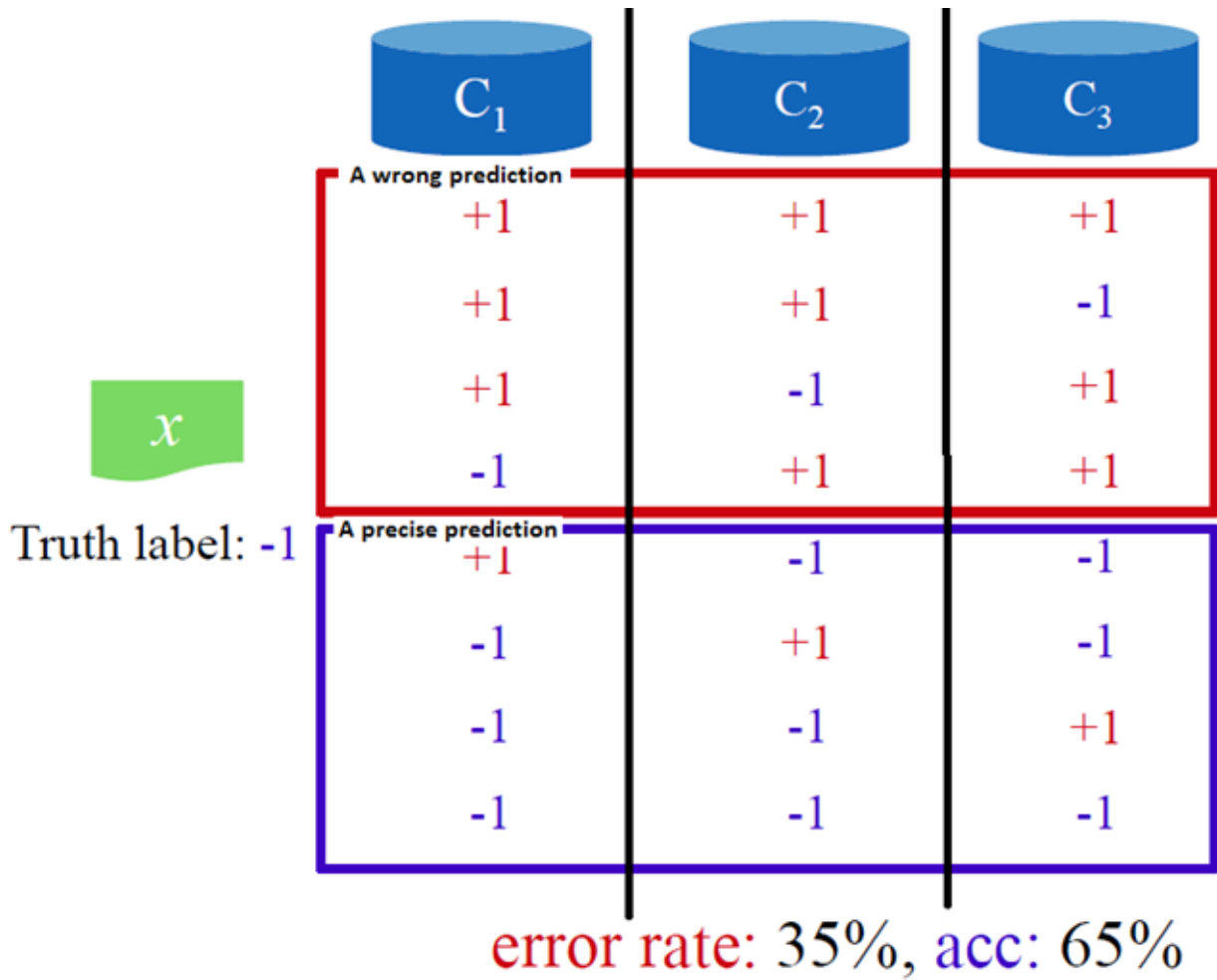
$\rho \leftarrow \rho + \beta [r - \rho + \max_{a'} Q(s', a') - \max_a Q(s, a)]$

CHAPTER 13



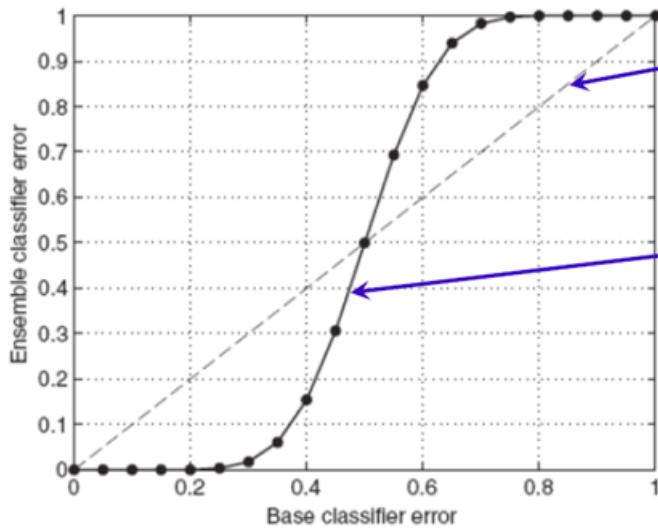


A test instance: x



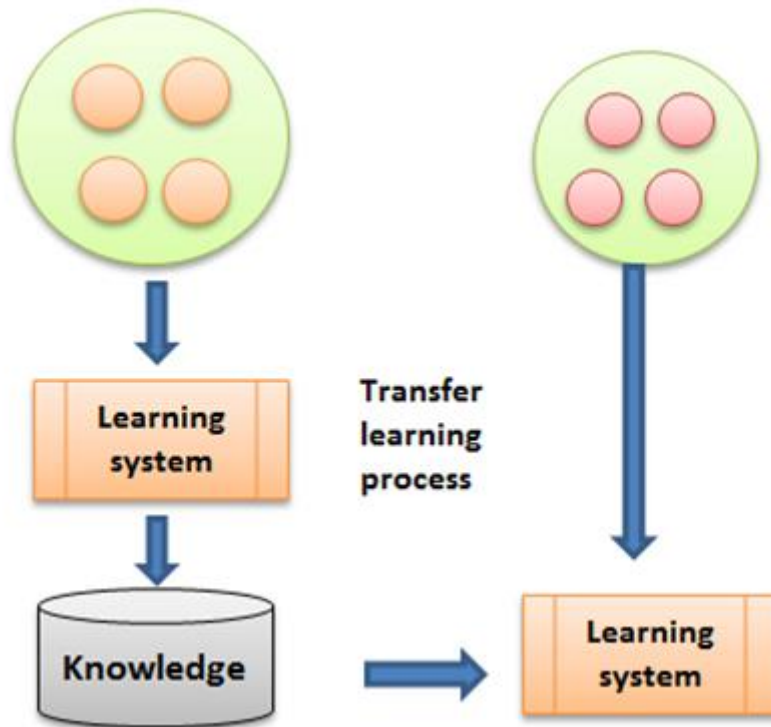
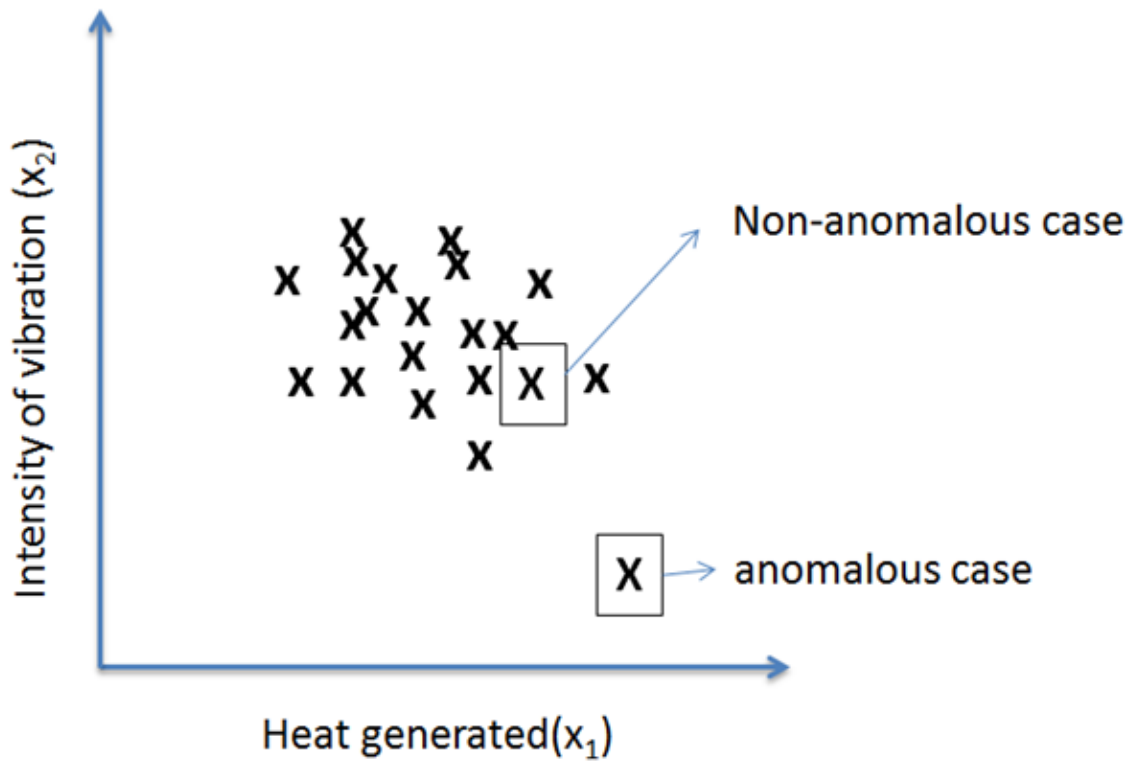
$$\sum_{i=2}^3 \binom{3}{i} \varepsilon^i (1-\varepsilon)^{3-i} = 3 \times 0.35^2 \times 0.65 + 1 \times 0.35^3 \times 1 = 0.2817$$

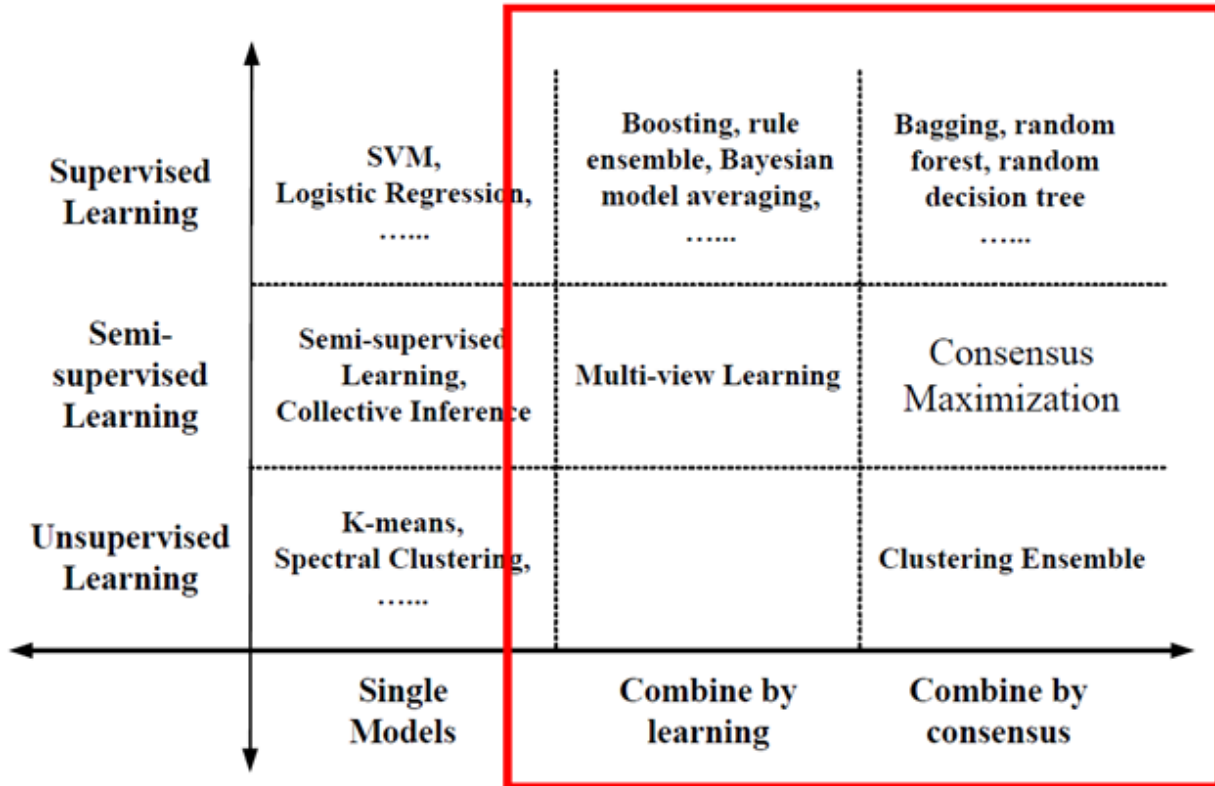
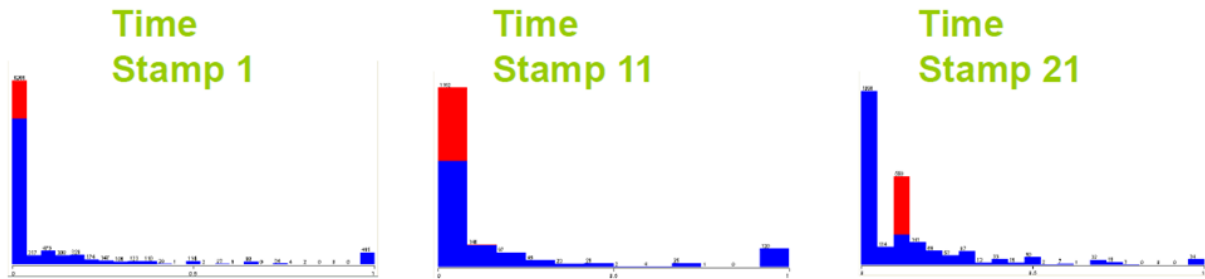
$$\sum_{i=13}^{25} \binom{25}{i} \varepsilon^i (1-\varepsilon)^{25-i} = 0.06$$



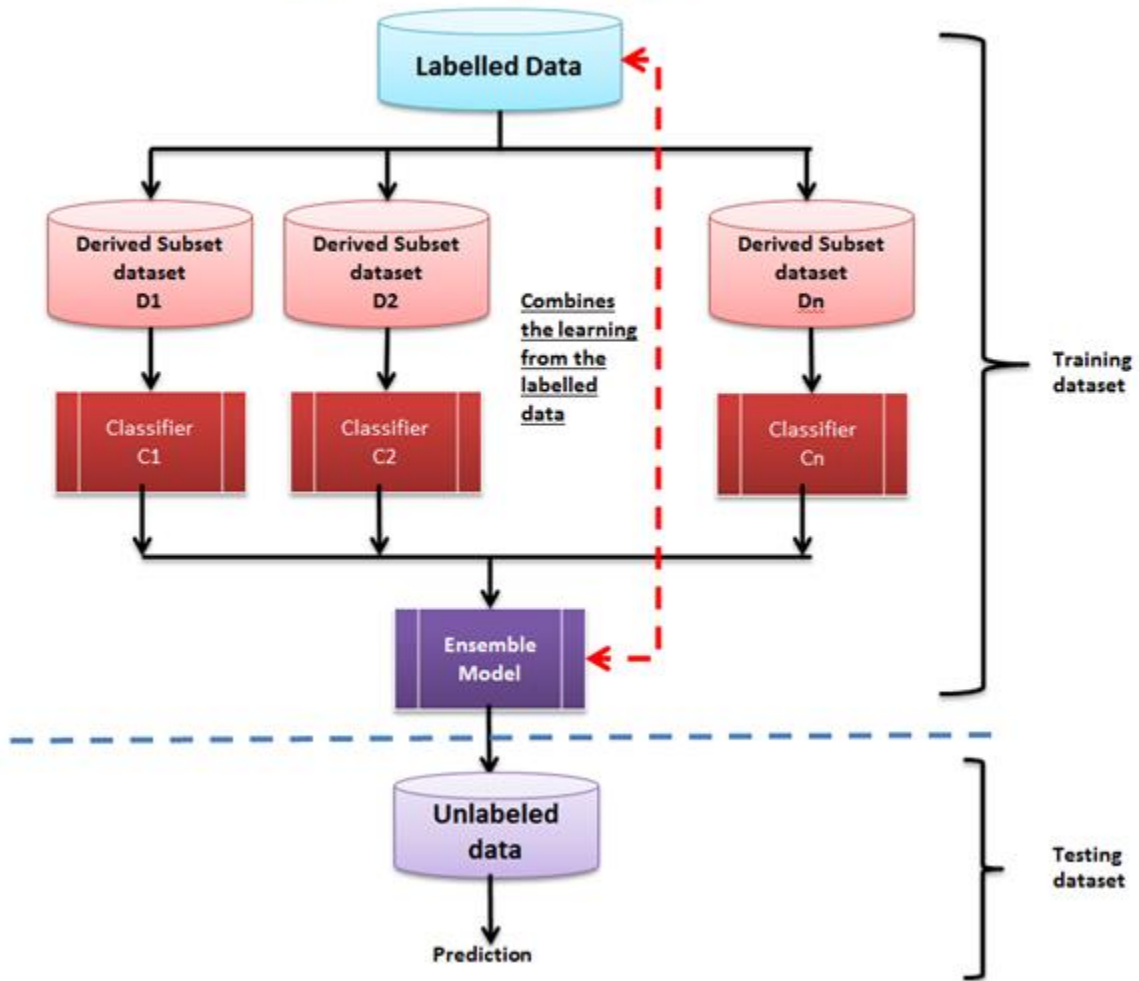
The base classifiers are identical (perfectly correlated)

The base classifiers are independent

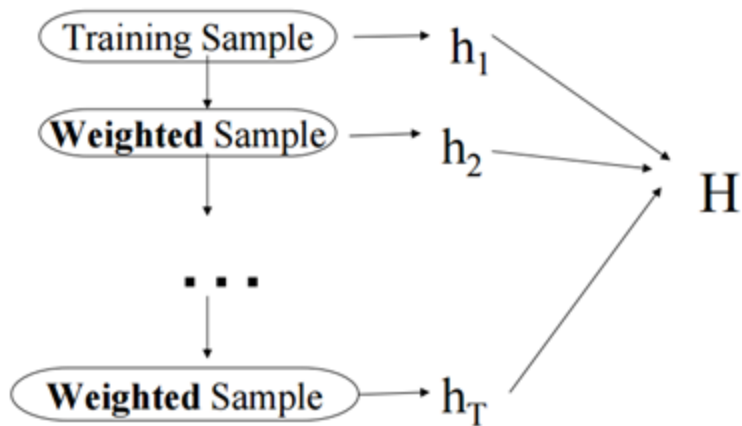
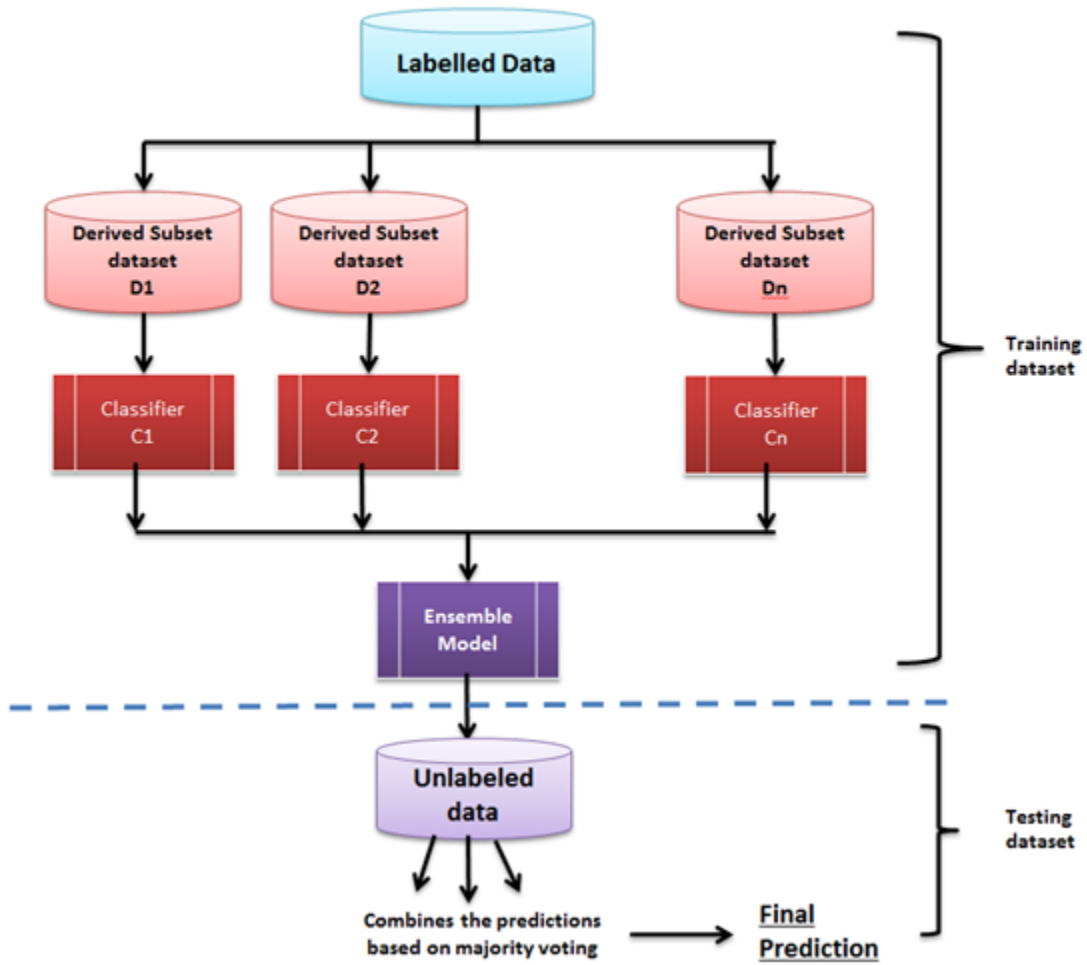




Ensemble method for Supervised learning
Combining the "learning" technique



Ensemble method for Supervised learning
Combining the "consensus" technique



$$H(x) = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(x)\right)$$

Given $(x_1, y_1), \dots, (x_m, y_m)$ where $x_i \in X$, $y_i \in \{-1, +1\}$

Initialise weights $D_1(i) = 1/m$

Iterate $t=1, \dots, T$:

- Train weak learner using distribution D_t
- Get weak classifier: $h_t: X \rightarrow \mathbb{R}$
- Choose $\alpha_t \in \mathbb{R}$
- Update: $D_{t+1}(i) = \frac{D_t(i) \exp(-\alpha_t y_i h_t(x_i))}{Z_t}$

- where Z_t is a normalization factor (chosen so that D_{t+1} will be a distribution), and α_t :

$$\alpha_t = \frac{1}{2} \ln\left(\frac{1 - \epsilon_t}{\epsilon_t}\right) > 0$$

Output – the final classifier

$$H(x) = \text{sign}\left(\sum_{t=1}^T \alpha_t h_t(x)\right)$$

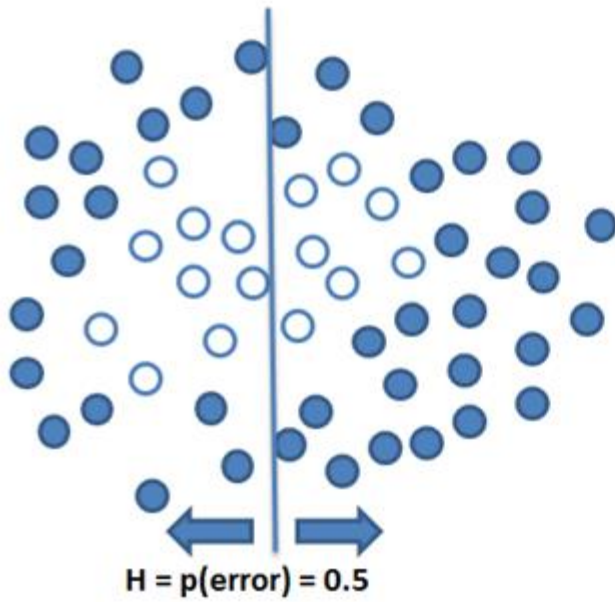


Data points class label

- $y_t = +1$
- $y_t = -1$

Data points weight value

$w_t = 1$

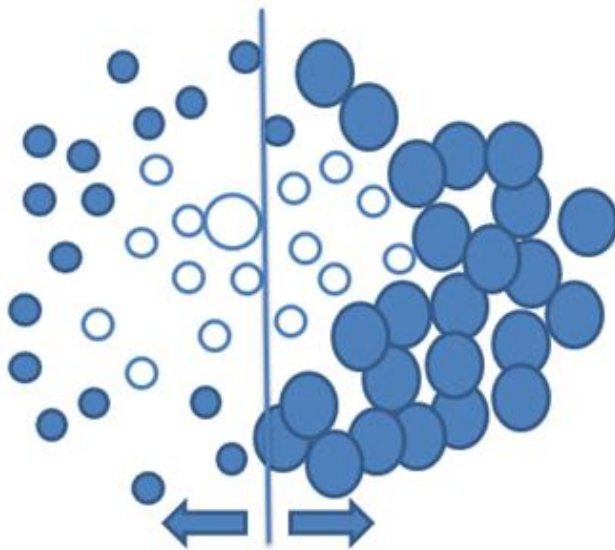


Data points class label

- $y_t = +1$
- $y_t = -1$

Data points weight value

$w_t = 1$

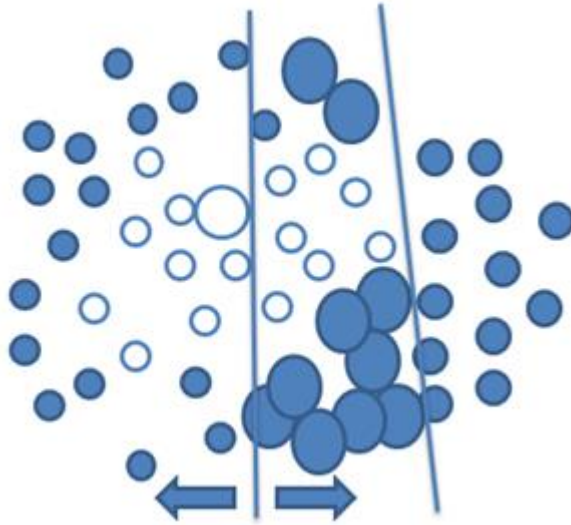


Data points class label

- $y_t = +1$
- $y_t = -1$

Weight value updated

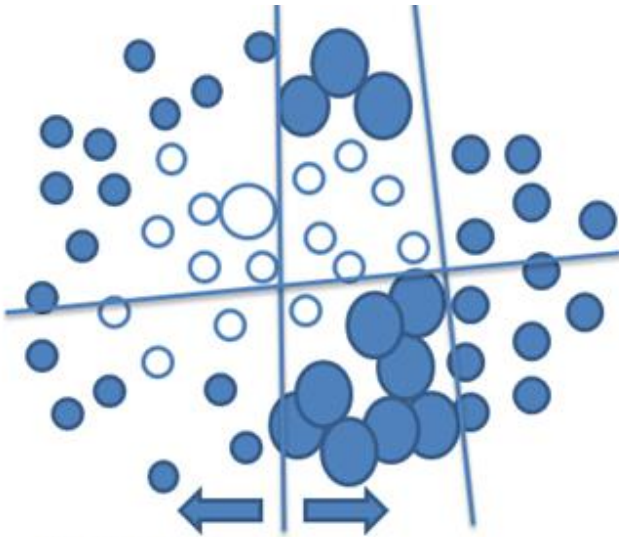
$w_t = w_t \exp\{-y_t, H\}$



Data points class label

- $y_t = +1$
- $y_t = -1$

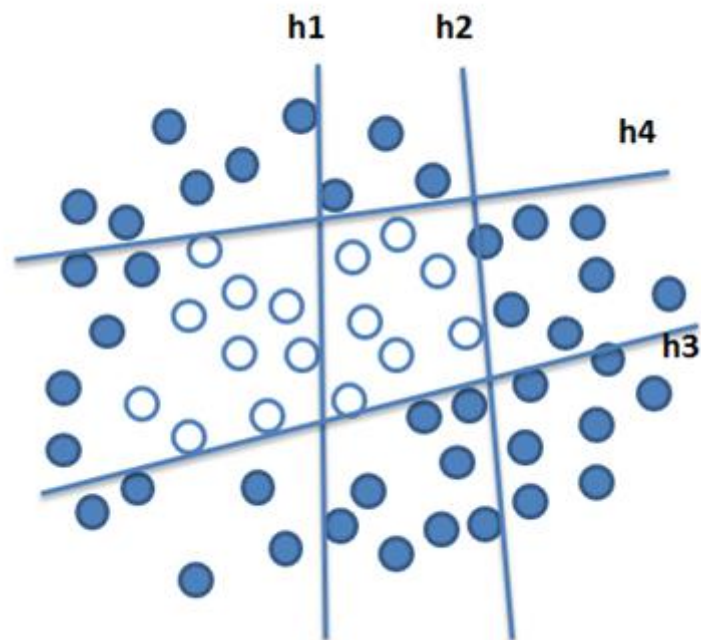
Weight value updated
 $w_t = w_t \exp\{-y_t, H\}$



Data points class label

- $y_t = +1$
- $y_t = -1$

Weight value updated
 $w_t = w_t \exp\{-y_t, H\}$



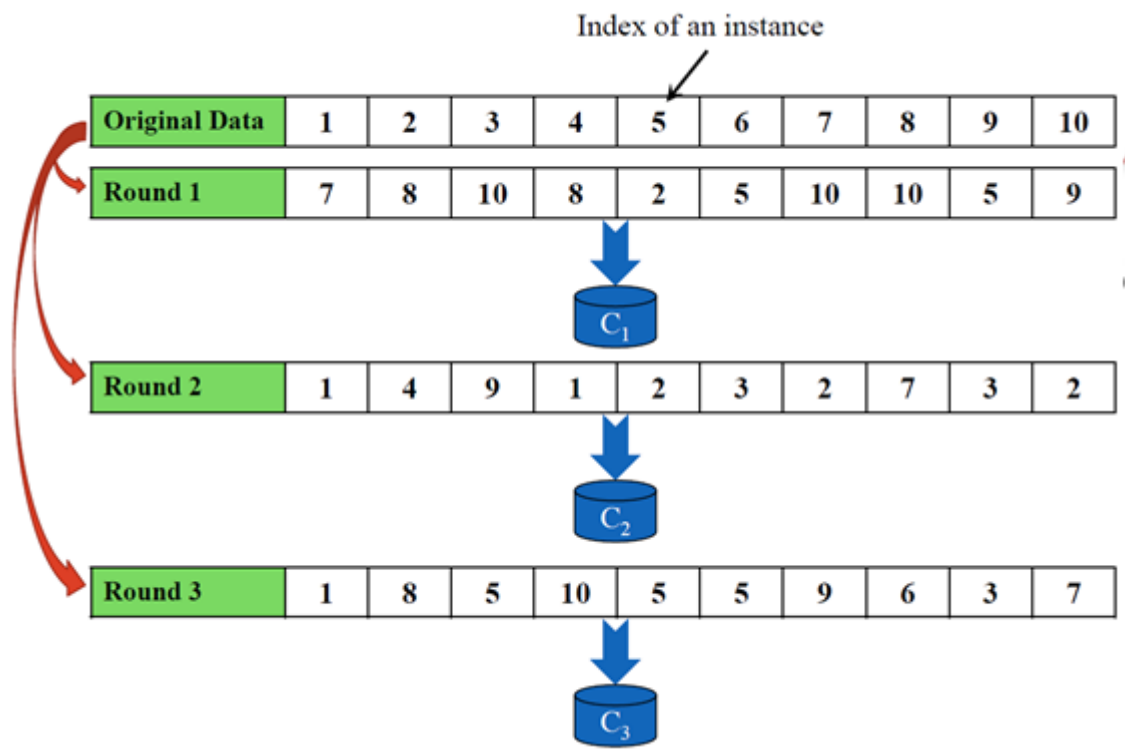
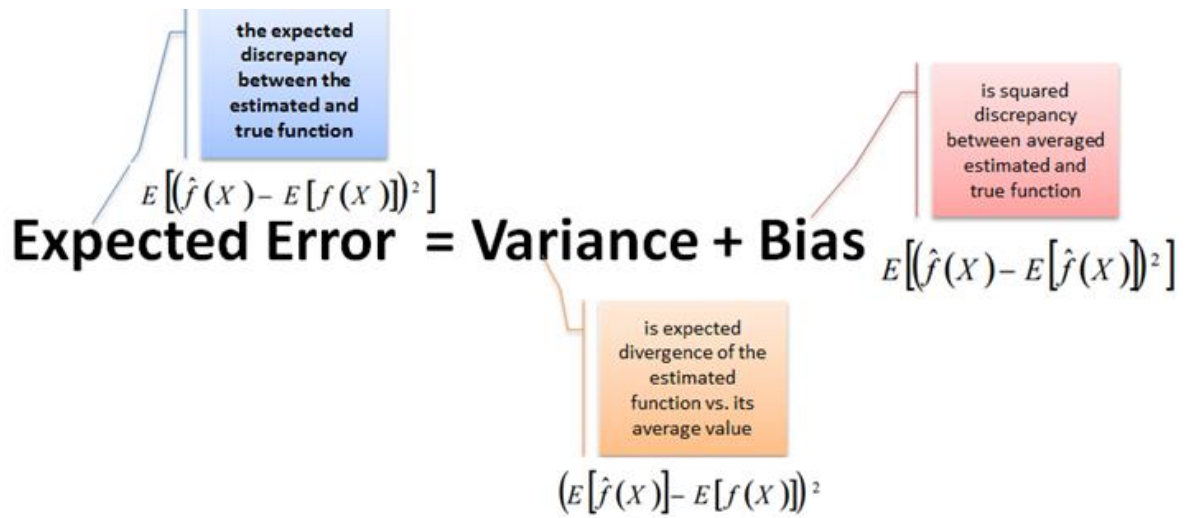
BAGGING

Training phase

1. Initialize the parameters
 - $\mathcal{D} = \emptyset$, the ensemble.
 - L , the number of classifiers to train.
2. For $k = 1, \dots, L$
 - Take a bootstrap sample S_k from \mathbf{Z} .
 - Build a classifier D_k using S_k as the training set.
 - Add the classifier to the current ensemble, $\mathcal{D} = \mathcal{D} \cup D_k$.
3. Return \mathcal{D} .

Classification phase

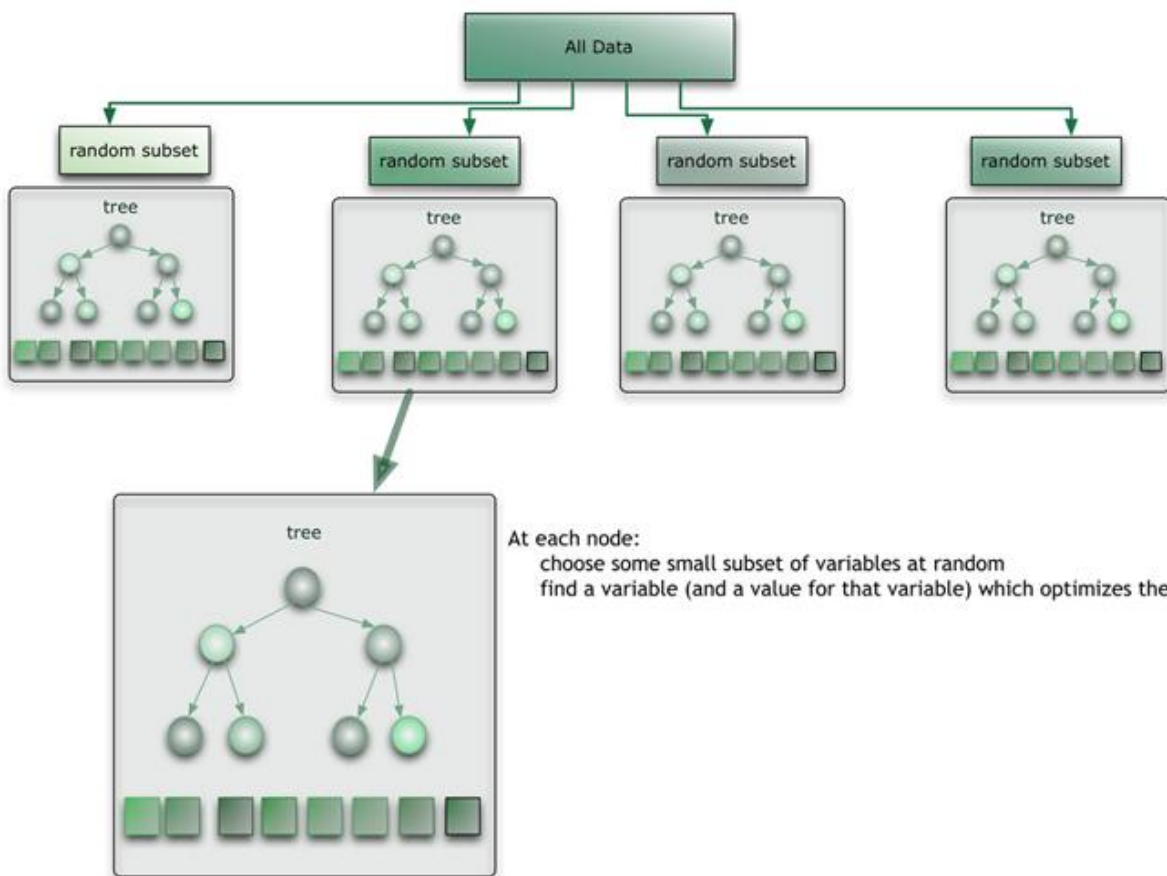
4. Run D_1, \dots, D_L on the input \mathbf{x} .
5. The class with the maximum number of votes is chosen as the label for \mathbf{x} .



Require: I (an inducer), T (the number of iterations), S (the training set), d (weighting distribution).

Ensure: $M_t; t = 1, \dots, T$

- 1: $t \leftarrow 1$
- 2: **repeat**
- 3: $S_t \leftarrow S$ with random weights drawn from d .
- 4: Build classifier M_t using I on S_t
- 5: $t++$
- 6: **until** $t > T$



Algorithm 1 Friedman's Gradient Boost algorithm

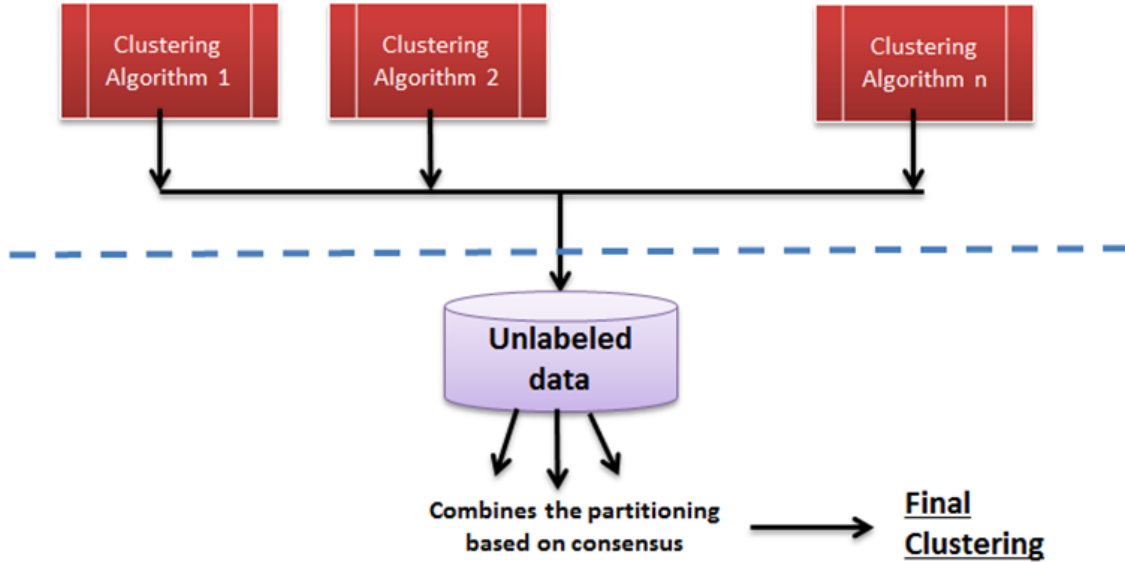
Inputs:

- input data $(x, y)_{i=1}^N$
- number of iterations M
- choice of the loss-function $\Psi(y, f)$
- choice of the base-learner model $h(x, \theta)$

Algorithm:

- 1: initialize \hat{f}_0 with a constant
 - 2: **for** $t = 1$ to M **do**
 - 3: compute the negative gradient $g_t(x)$
 - 4: fit a new base-learner function $h(x, \theta_t)$
 - 5: find the best gradient descent step-size ρ_t :
$$\rho_t = \arg \min_{\rho} \sum_{i=1}^N \Psi [y_i, \hat{f}_{t-1}(x_i) + \rho h(x_i, \theta_t)]$$
 - 6: update the function estimate:
$$\hat{f}_t \leftarrow \hat{f}_{t-1} + \rho_t h(x, \theta_t)$$
 - 7: **end for**
-

Ensemble method for Unsupervised learning
Combining the "consensus" technique



base clustering models



objects →

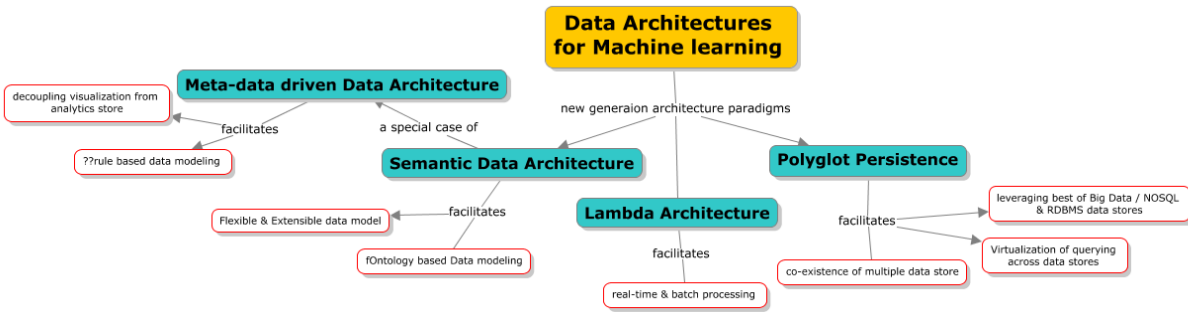
	C_1	C_2	C_3	C
v_1	1	1	1	1
v_2	1	2	2	2
v_3	2	1	1	1
v_4	2	2	2	2
v_5	3	3	3	3
v_6	3	4	3	3

they may not represent the same cluster!

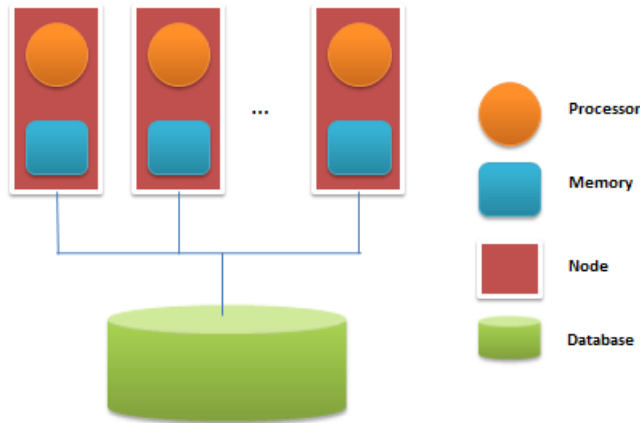
The goal: get the consensus clustering



CHAPTER 14



Shared Disk Data Architecture



Archives
Scanned documents, statements, medical records, e-mails etc..



Media
Images, video, audio etc.



Data Storages
RDBMS, NoSQL, Hadoop, file systems etc.



Docs
XLS, PDF, CSV, HTML, JSON etc.



Social Networks
Twitter, Facebook, Google+, LinkedIn etc.



Machine Log Data
Application logs, event logs, server data, CDRs, clickstream data etc.



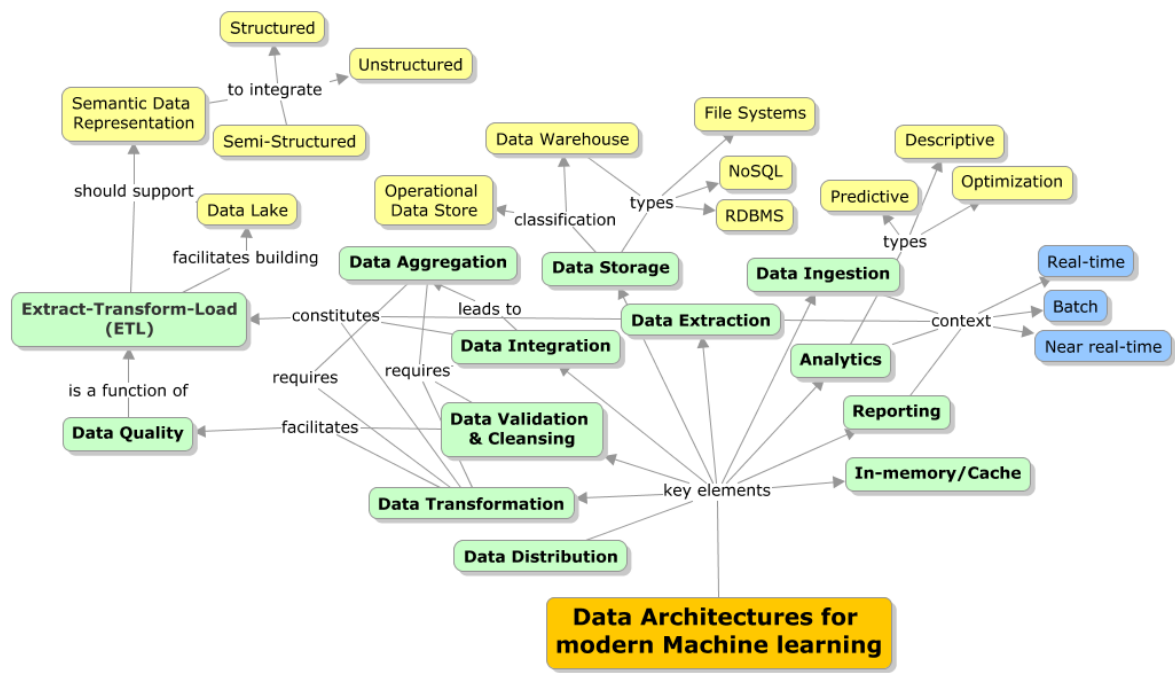
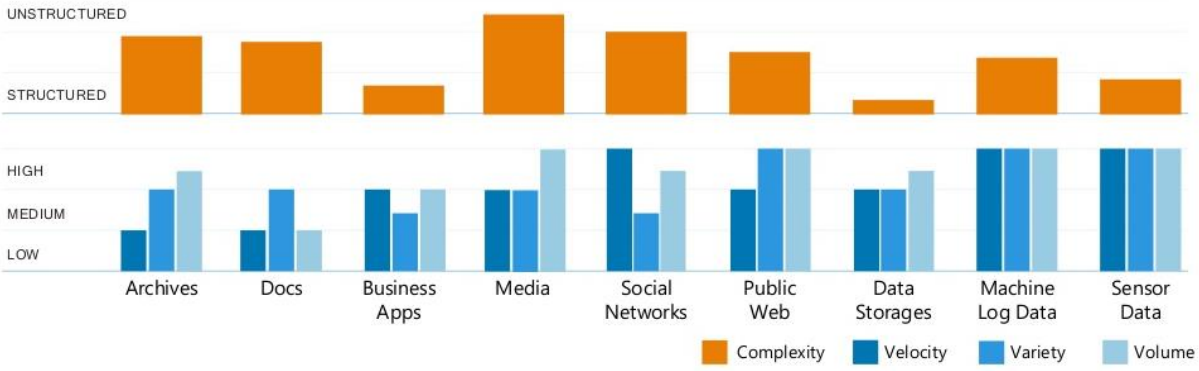
Business Apps
CRM, ERP systems, HR, project management etc.

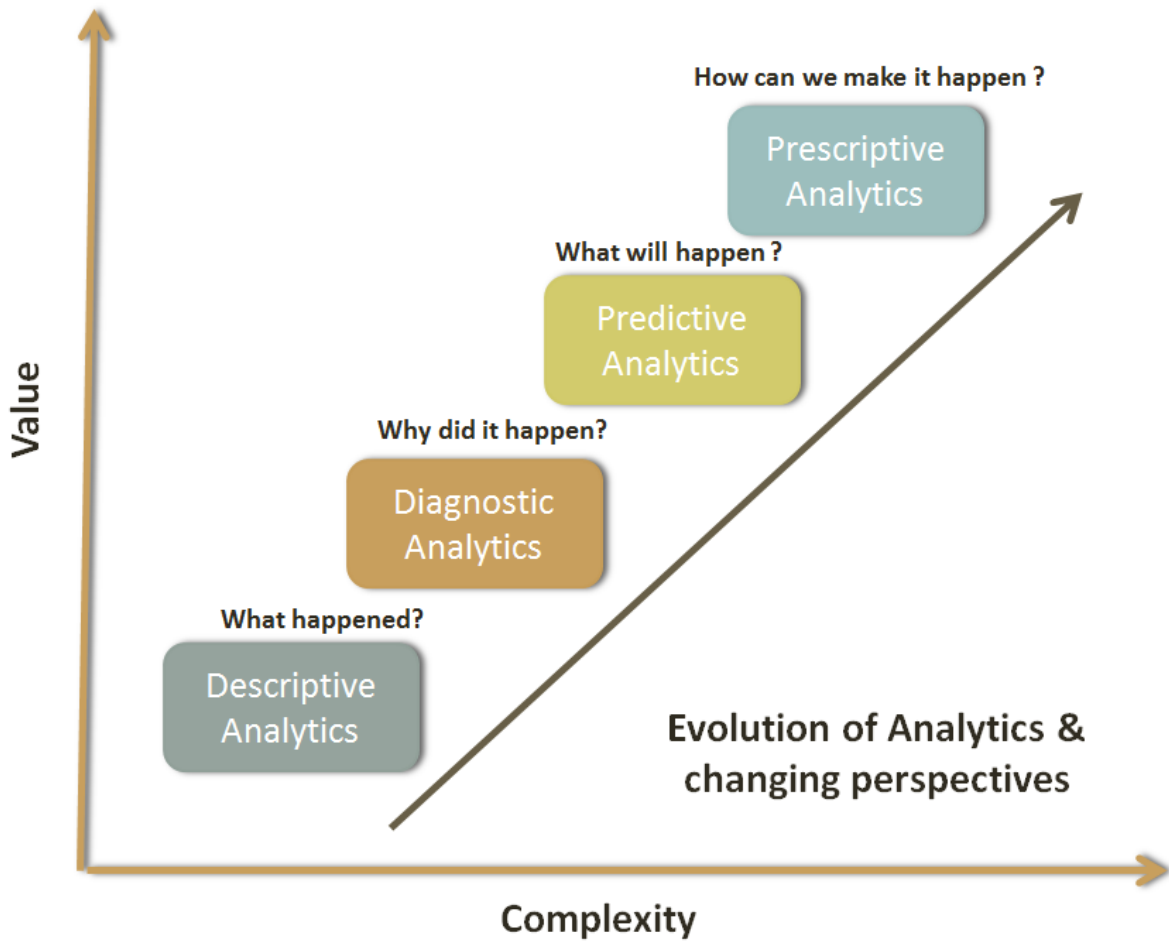


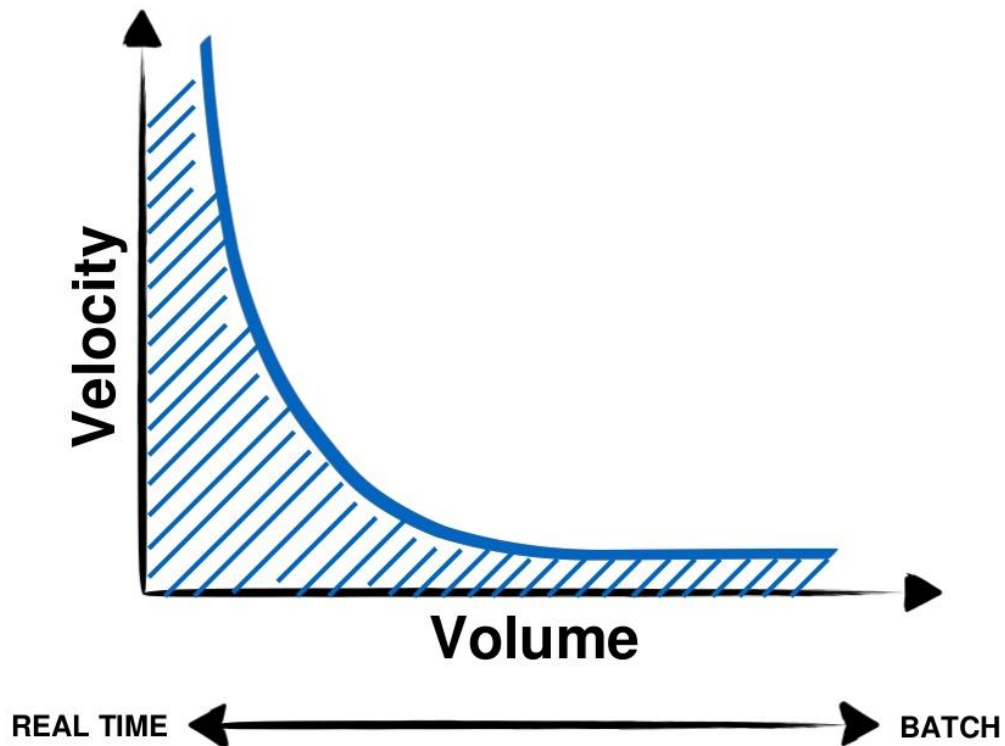
Public Web
Wikipedia, news, weather, public finance etc



Sensor Data
Smart electric meters, medical devices, car sensors, road cameras etc.







Semantified Common Data Repository

Data Ingestion

Data Management

Query Management

delivering

- Low Cost, High Performance Storage
- Flexible, Easy-to-Use Data Organization
- Performance-Optimized Analytics
- Automation of most manual Development and Query Activities
- Self-Service End-User Features
- Intelligent Processing

