

Lecture notes
Knowledge Discovery in Databases II
Winter Semester 2013/2014

Lecture 9: Stream classification

Lectures: PD Dr Matthias Schubert, Dr. Eirini Ntoutsis
Tutorials: Erich Schubert

Notes © 2012 Eirini Ntoutsis

[http://www.dbs.ifi.lmu.de/cms/Knowledge_Discovery_in_Databases_II_\(KDD_II\)](http://www.dbs.ifi.lmu.de/cms/Knowledge_Discovery_in_Databases_II_(KDD_II))

- Motivation
- Data streams
- Data stream clustering
- Data stream clustering algorithms
- Data stream classification
- Data stream classifiers
- Data stream classifiers evaluation

- So far: focus on batch learning using small datasets
- Batch learning:
 - Whole training data is available to the learning algorithm
 - Data instances are processed multiple times
 - e.g. k-Means clusterer (c.f., slide 4), ID3 classifier (c.f., slide 5)
- Assumption:
 - Instances are generated by some stationary probability distribution

- But, most interesting applications come from dynamic environments where data are collected over time
 - e.g., customer transactions, call records, customer click data.
- Batch learning is not sufficient anymore ...
- Algorithms should be able to incorporate new data
 - There are algorithms that are incremental by nature, e.g. kNN classifiers, Naïve Bayes classifiers
 - But the majority of the algorithms need changes to make incremental induction, e.g., ID3, k-Means, DBSCAN.
- Algorithms should be able to deal with non-stationary data generation processes
 - Incorporate concept drift
 - Forget outdated data and adapt to the most recent state of the data

- Cluster evolution

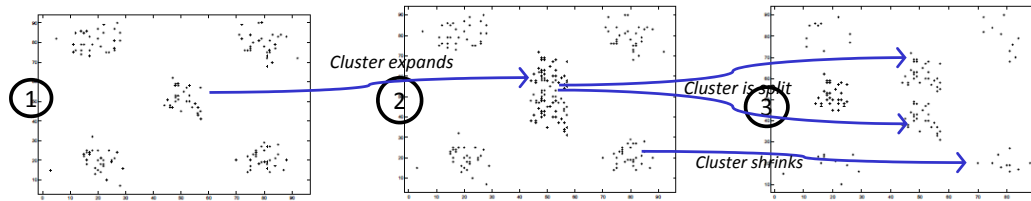


Figure: Data records at three consecutive time stamps, the clustering gradually changes (from: *MONIC - Modeling and Monitoring Cluster Transitions*, Spiliopoulou et al, KDD 2006)

- Concept drift

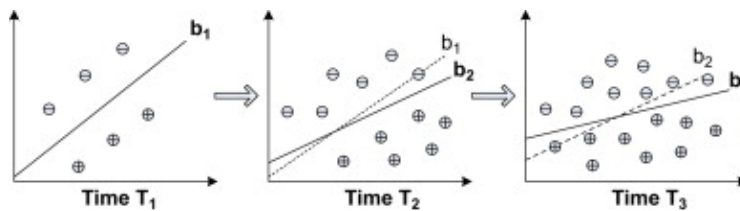


Fig. 1. An illustration of concept drifting in data streams. In the three consecutive time stamps T_1 , T_2 and T_3 , the classification boundary gradually drifts from b_1 to b_2 and finally to b_3 . (from: *A framework for application-driven classification of data streams*, Zhang et al, Journal Neurocomputing 2012)

- Banks (credit card transactions, loan applications,...)
- Telecommunication (call records, sms, www usage,...)
- Health care systems (customer records in a hospital,...)
- Retail industry (transactions in a supermarket,...)
- WWW (content, interactions, TCP/IP traffic, customer click data,...)
- Science (experiment results,...)
- ...

facebook

COMPANY

- Press Room
- Factsheet
- Statistics**
- Timeline
- Management Team
- Founder Bios
- Platform
- B-Roll
- Press Releases & Announce...

CONTACTS

- Speaker Requests
- Interview Requests
- Facebook Stories

Statistics

Blog · About · Press Releases RSS

People on Facebook More than 800 million active users
More than 50% of our active users log on to Facebook in any given day
Average user has 130 friends

Activity on Facebook More than 900 million objects that people interact with (pages, groups, events and community pages)
Average user is connected to 80 community pages, groups and events
On average, more than 250 million photos are uploaded per day

Global Reach More than 70 languages available on the site
Approximately 80% of users are outside of the United States
Over 300,000 users helped translate the site through the translations application

Platform On average, people on Facebook install apps more than 20 million times every day
Every month, more than 500 million people use an app on Facebook or experience Facebook Platform on other websites
More than 7 million apps and websites are integrated with Facebook

Mobile More than 350 million active users currently access Facebook through their mobile devices
More than 475 mobile operators globally work to deploy and promote Facebook mobile products

[Source: http://www.facebook.com/press/info.php?statistics](http://www.facebook.com/press/info.php?statistics)

#twitter#

by the numbers

The success of Twitter has been nothing short of amazing. What follows is a brief roundup of some of the latest growth and engagement metrics that Twitter recently published.

Quick Twitter Stats

- ✓ **1 billion** number of tweets posted per week
- ✓ **3 years, 2 months and 1 day** the time in between the first tweet and the billionth tweet
- ✓ **177 million** number of tweets sent on March 11, 2011
- ✓ **456** tweets per second (TPS) when Michael Jackson died on June 25, 2009 (a record at that time)
- ✓ **6,939** tweets per second, set a second after midnight in Japan on New Year's Day.
- ✓ **572,000** number of new accounts created on March 12, 2011
- ✓ **460,000** average number of new accounts per day during February, 2011.

Increase in the average number of tweets per day (TPD) for each month

Increase in the number of mobile users during 2010

Between 2008 and 2011 there was a 5000% increase in the number of employees at Twitter

Month	Number of Employees
January 2008	8
January 2009	29
January 2010	130
January 2011	330
March 2011	400

Source: twitter.com

Tweets per second

- 6,939** JAN 1 New Years
- 4,064** FEB 6 Superbowl
- 5,530** MAR 11 Japanese earthquake and tsunami
- 3,966** APR 28 UK Royal Wedding
- 5,106** MAY 2 Raid on Osama bin Laden
- 6,303** MAY 28 UEFA Champions League
- 5,531** JUN 3 NBA Finals
- 6,436** JUN 27 BET Awards
- 4,995** JUL 11 Home Run Derby
- 7,166** JUL 17 Brazil eliminated from the Cope America
- 7,196** JUL 17 End of FIFA Women's World Cup
- 5,449** AUG 23 East Coast earthquake
- 7,064** SEP 25 Steve Jobs resigns
- 8,868** SEP 28 MTV Video Music Awards
- 7,671** SEP 20 Troy Davis executed
- 6,049** OCT 6 Steve Jobs passes away



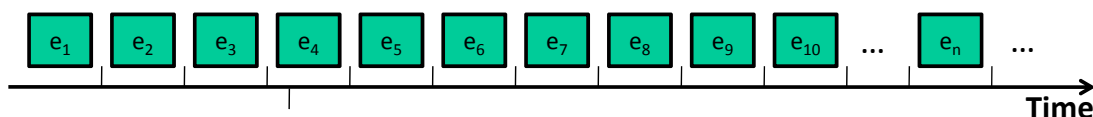
- Experiments at CERN are generating an entire petabyte (1PB=10⁶ GB) of data every second as particles fired around the Large Hadron Collider (LHC) at velocities approaching the speed of light are smashed together
- “We don’t store all the data as that would be impractical. Instead, from the collisions we run, we only keep the few pieces that are of interest, the rare events that occur, which our filters spot and send on over the network,” he said.
- This still means CERN is storing 25PB of data every year – the same as 1,000 years' worth of DVD quality video – which can then be analyzed and interrogated by scientists looking for clues to the structure and make-up of the universe.

Source: <http://public.web.cern.ch/public/en/LHC/Computing-en.html>

Source: <http://www.v3.co.uk/v3-uk/news/2081263/cern-experiments-generating-petabyte>

*“Τα πάντα ρεῖ καὶ οὐδὲν μένει”
 (“Ta panta rhei kai ouden menei”)
 “Everything flows, nothing stands still”
 Heraclitus (535-475 BC)*

- Data evolve over time as new data arrive and old data become obsolete
- We can distinguish between:
 - Dynamic data arriving at a low rate
 - Incremental methods might work for such cases
 - Data streams: possible infinite sequence of elements arriving at a rapid rate
 - new methods are required to deal with the amount and complexity of these data



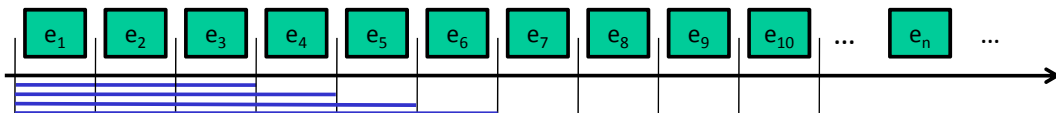
time	duration	protocol_type	service	flag	src_bytes	dst_bytes	...	class
t ₁	0	tcp	http	SF	181	5450	...	normal
t ₂	0	tcp	http	SF	239	486	...	normal
...
t ₇₈₃₈	0	icmp	ecr_i	SF	1032	0	...	smurf
t ₇₈₃₉	0	icmp	ecr_i	SF	1032	0	...	smurf
...
t ₇₀₅₃₁	0	tcp	private	S0	0	0	...	neptune
t ₇₀₅₃₂	0	tcp	private	S0	0	0	...	neptune
...
t ₄₉₂₃₁₀	0	tcp	http	SF	244	7161	...	normal
t ₄₉₂₃₁₁	0	tcp	http	SF	258	9517	...	normal
...

- The dataset consists of TCP connection records of LAN network traffic managed by Lincoln Labs.
- A connection is a sequence of TCP packets starting and ending at some well defined times, between which data flows to and from a source IP address to a target IP address under some well defined protocol.
- Connections are described in terms of 42 features like duration, protocol_type, service, flag, src_bytes, dst_bytes etc.,
- Each connection is labeled as either normal, or as an attack, with exactly one specific attack type. There are 4 main categories of attacks: DOS, R2L, U2R, PROBING and are further classified into attack types, like buffer-overflow, guess-passwd, neptune etc.
- Most of the connections in this dataset are normal, but occasionally there could be a burst of attacks at certain times.

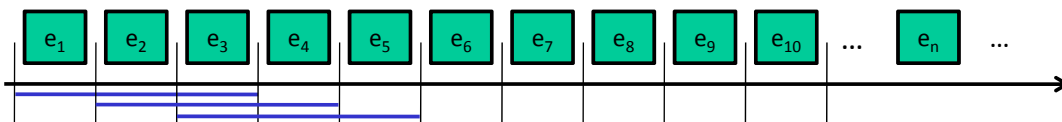
More on this dataset: <http://kdd.ics.uci.edu/databases/kddcup99/kddcup99.html>

- Data Mining over stream data is even more challenging:
 - Huge amounts of data → only a small amount can be stored in memory
 - Arrival at a rapid rate → no much time for processing
 - The generative distribution of the stream might change over time rather than being stationary → adapt and report on changes
- Requirements for stream mining algorithms:
 - Use limited computational resources:
 - Bounded memory
 - Small processing time
 - No random access to the data
 - Only 1 look at the data (upon their arrival)

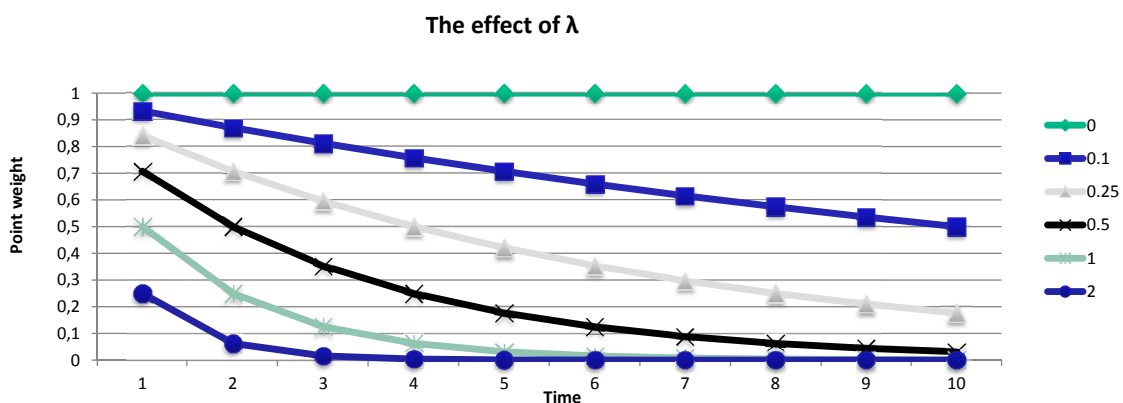
- Usually we are not interested in the whole history of the stream but only in the recent history.
- There are different *ageing/weighting mechanisms or window models* that reflect which part of the stream history is important
 - Landmark window:
 - Nothing is forgotten.
 - All points have a weight $w=1$.



- Sliding window:
 - Remember only the n more recent entries.
 - All points within the window have a weight $w=1$, for the rest: $w=0$.



- Damped window:
 - Data are subject to ageing according to a fading function $f(t)$, i.e., each point is assigned a weight that decreases with time t via $f(t)$.
 - A widely used fading function in temporal applications is the exponential fading function: $f(t)=2^{-\lambda t}$, $\lambda>0$.
 - The decay rate λ determines the importance of historical data
 - The higher the value of λ , the lower the importance of old data



- We focus on:
 - Data stream clustering
 - Data stream classification

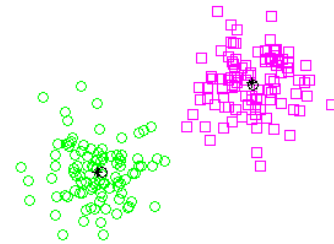
- One of the core tasks in DM
 - Used as either a standalone tool or as a preprocessing tool
- The (batch) clustering problem:
 - Given a set of measurements, observations, etc., the goal is to group the data into groups of similar data (clusters)
- The data stream clustering problem:
 - Maintain a continuously consistent good clustering of the sequence observed so far, using a small amount of memory and time.
- This implies:
 - Incremental clustering
 - Maintaining cluster structures that evolve over time
 - Working with summaries instead of raw data

- Traditional clustering methods require access upon the whole dataset
 - Online maintenance of clustering
- The underlying population distribution might change: drifts/ shifts of concepts
 - One clustering model might not be adequate to capture the evolution
- The role of outliers and clusters are often exchanged in a stream
 - Clear and fast identification of outliers

	Static clustering	Dynamic/Stream clustering
Partitioning methods	<ul style="list-style-type: none"> • k-Means • k-Medoids 	<ul style="list-style-type: none"> • Leader • Simple single pass k-Means • STREAM k-Means • CluStream
Density-based methods	<ul style="list-style-type: none"> • DBSCAN • OPTICS 	<ul style="list-style-type: none"> • DenStream • incDBSCAN * • incOPTICS *
Grid-based methods	<ul style="list-style-type: none"> • STING 	<ul style="list-style-type: none"> • DStream

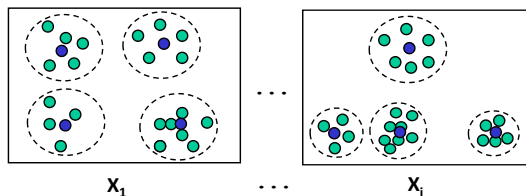
(*) These methods require access to the raw data (this access might be limited though)

- Goal: Construct a partition of a set of objects into k clusters
 - e.g. k-Means, k-Medoids
- Two types of methods:
 - Adaptive methods:
 - Leader (Spath 1980)
 - Simple single pass k-Means (Farnstrom et al, 2000)
 - STREAM k-Means (O'Callaghan et al, 2002)
 - Online summarization - offline clustering methods:
 - CluStream (Aggarwal et al, 2003)



- The simplest single-pass partitioning algorithm
- Whenever a new instance p arrives from the stream
 - Find its closest cluster (leader), c_{clos}
 - Assign p to c_{clos} if their distance is below the threshold d_{thresh}
 - Otherwise, create a new cluster (leader) with p
- + 1-pass and fast algorithm
- + No prior information on the number of clusters
- Unstable algorithm
- It depends on the order of the examples
- It depends on a correct guess of d_{thresh}

- An extension of k-Means for streams
 - The iterative process of static k-Means cannot be applied to streams
 - Use a buffer that fits in memory and apply k-Means locally in the buffer
- Stream is processed in chunks X_1, X_2, \dots , each fitting in memory
 - For each chunk X_i
 - Apply k-Means locally on X_i (retain only the k centers)
 - $X' \leftarrow i * k$ weighted centers obtained from chunks $X_1 \dots X_i$
 - Each center is treated as a point, weighted with the number of points it compresses
 - Apply k-Means on X' output the k centers



- The stream clustering process is separated into:
 - an online **micro-cluster** component, that summarizes the stream locally as new data arrive over time
 - Micro-clusters are stored in disk at snapshots in time that follow a pyramidal time frame.
 - an offline **macro-cluster** component, that clusters these summaries into global clusters
 - Clustering is performed upon summaries instead of raw data

- Assume that the data stream consists of a set of multi-dimensional records X_1, \dots, X_n, \dots , arriving at T_1, \dots, T_n, \dots : $X_i = (x_i^1, \dots, x_i^d)$
- The micro-cluster summary for a set of d-dimensional points (X_1, X_2, \dots, X_n) arriving at time points T_1, T_2, \dots, T_n is defined as:

$$\text{CFT} = (\text{CF2}^x, \text{CF1}^x, \text{CF2}^t, \text{CF1}^t, n)$$

$\sum_{i=1}^N X_i^2$

$\sum_{i=1}^N X_i$

$\sum_{i=1}^N T_i^2$

$\sum_{i=1}^N T_i$

- Easy calculation of basic measures to characterize a cluster:
 - Center: $\frac{\text{CF1}^x}{n}$
 - Radius: $\sqrt{\frac{\text{CF2}^x}{n} - \left(\frac{\text{CF1}^x}{n}\right)^2}$
- Important properties of micro-clusters:
 - Incrementality: $\text{CFT}(C_1 \cup p) = \text{CFT}(C_1) + p$
 - Additivity: $\text{CFT}(C_1 \cup C_2) = \text{CFT}(C_1) + \text{CFT}(C_2)$
 - Subtractivity: $\text{CFT}(C_1 - C_2) = \text{CFT}(C_1) - \text{CFT}(C_2)$, $C_1 \supseteq C_2$

- A fixed number of q **micro-clusters** is maintained over time
- **Initialize**: apply q -Means over *initPoints*, built a summary for each cluster
- **Online** micro-cluster maintenance as a new point p arrives from the stream
 - Find the closest micro-cluster clu for the new point p
 - If p is within the max-boundary of clu , p is absorbed by clu
 - o.w., a new cluster is created with p
 - The number of micro-clusters should not exceed q
 - Delete most obsolete micro-cluster or merge the two closest ones
- **Periodic** storage of micro-clusters snapshots into disk
 - At different levels of granularity depending upon their recency
- **Offline** macro-clustering
 - Input: A user defined time horizon h and number of macro-clusters k to be detected
 - Locate the valid micro-clusters during h
 - Apply k-Means upon these micro-clusters $\rightarrow k$ macro-clusters

- Initialization
 - Done using an offline process in the beginning
 - Wait for the first *InitNumber* points to arrive
 - Apply a standard *k*-Means algorithm to create *q* clusters
 - For each discovered cluster, assign it a unique ID and create its micro-cluster summary.
- Comments on the choice of *q*
 - much larger than the natural number of clusters
 - much smaller than the total number of points arrived

- A fixed number of *q* micro-clusters is maintained over time
- Whenever a new point *p* arrives from the stream
 - Compute distance between *p* and each of the *q* maintained micro-cluster centroids
 - *clu* ← the closest micro-cluster to *p*
 - Find the max boundary of *clu*
 - It is defined as a factor of *t* of *clu* radius
 - If *p* falls within the maximum boundary of *clu*
 - *p* is **absorbed** by *clu*
 - Update *clu* statistics (incrementality property)
 - Else, create a **new** micro-cluster with *p*, assign it a new cluster ID, initialize its statistics
 - To keep the total number of micro-clusters fixed (i.e., *q*):
 - **Delete** the most obsolete micro-cluster or
 - If its safe (based on how far in the past, the micro-cluster received new points)
 - **Merge** the two closest ones (Additivity property)
 - When two micro-clusters are merged, a list of ids is created. This way, we can identify the component micro-clusters that comprise a micro-cluster.

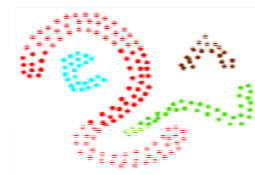
- Micro-clusters snapshots are stored at particular moments
- If current time is t_c and user wishes to find clusters based on a history of length h
 - Then we use the subtractive property of micro-clusters at snapshots t_c and t_c-h
 - In order to find the macro-clusters in a history or time horizon of length h
- How many snapshots should be stored?
 - It will be too expensive to store snapshots at every time stamp
 - They are stored in a pyramidal time frame
- It is an effective trade-off between the storage requirements and the ability to recall summary statistics from different time horizons.

- The offline step is applied on demand upon the q maintained micro-clusters instead of the raw data
- User input: time horizon h , # macro-clusters k to be detected
- Find the active micro-clusters during h :
 - We exploit the subtractivity property to find the active micro-clusters during h :
 - Suppose current time is t_c . Let $S(t_c)$ be the set of micro-clusters at t_c .
 - Find the stored snapshot which occurs just before time t_c-h . We can always find such a snapshot h' . Let $S(t_c-h')$ be the set of micro-clusters.
 - For each micro-cluster in the current set $S(t_c)$, we find the list of ids. For each of the list of ids, find the corresponding micro-clusters in $S(t_c-h')$.
 - Subtract the CF vectors for the corresponding micro-clusters in $S(t_c-h')$
 - This ensures that the micro-clusters created before the user-specified horizon do not dominate the result of clustering process
- Apply k-Means over the active micro-clusters in h to derive the k macro-clusters
 - Initialization: seeds are not picked up randomly, rather sampled with probability proportional to the number of points in a given micro-cluster
 - Distance is the centroid distance
 - New seed for a given partition is the weighted centroid of the micro-clusters in that partition

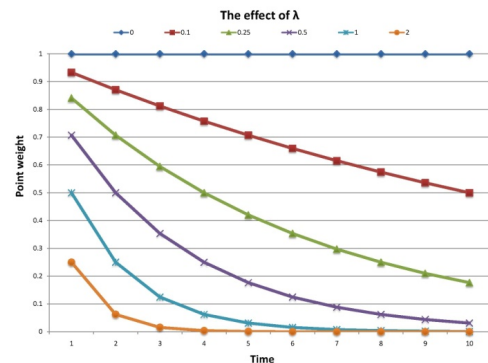
- + CluStream clusters large evolving data stream
- + Views the stream as a changing process over time, rather than clustering the whole stream at a time
- + Can characterize clusters over different time horizons in changing environment
- + Provides flexibility to an analyst in a real-time and changing environment

- Fixed number of micro-clusters maintained over time
- Sensitive to outliers/ noise

- Clusters as regions of high density surrounded by regions of low density (noise)
 - Density is measured locally, in the ϵ -neighborhood of each point
 - e.g. DBSCAN, OPTICS
- Very appealing for streams
 - No assumption on the number of clusters
 - Discovering clusters of arbitrary shapes
 - Ability to handle outliers and noise
- But, they miss a clustering model (or it is too complicated)
 - Clusters are represented by all their points
- Solution: Describe clusters as set of summaries
 - DenStream (Cao et al, 2006)



- The online-offline rationale is followed:
 - Online summarization as new data arrive over time
 - Core, potential core and outlier micro-clusters
 - Offline clustering over the summaries to derive the final clusters
 - A modified version of DBSCAN over the summaries
- Data are subject to ageing according to the exponential ageing function (damped window model)
 - $f(t)=2^{-\lambda t}, \lambda>0$
 - λ the decay rate
 - higher λ , less important the historical data comparing to more recent data



- The micro-cluster summary at time t for a set of d -dimensional points (p_1, p_2, \dots, p_n) arriving at time points T_1, T_2, \dots, T_n is:

$$MC = (CF^1, CF^2, w)$$

$\sum_{i=1}^n f(t-T_i)p_i$

$\sum_{i=1}^n f(t-T_i)p_i^2$

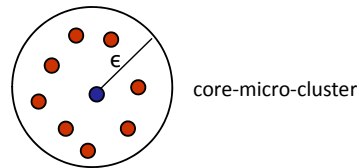
$\sum_{i=1}^n f(t-T_i)$

- Easy computation of basic measures:
 - Center: $c = \frac{CF^1}{w}$
 - Radius: $r = \sqrt{\frac{CF^2}{w} - \left(\frac{CF^1}{w}\right)^2}$
- A micro-cluster summary c_p can be maintained incrementally
 - If a new point p is added to c_p : $c_p = (CF^2+p^2, CF^1+p, w+1)$
 - If no point is added to c_p for time interval δt :
 - $c_p = (2^{-\lambda\delta t}*CF^2, 2^{-\lambda\delta t}*CF^1, 2^{-\lambda\delta t}*w)$

- Core (or dense) micro-clusters

- $(w \geq \mu) \ \& \ (r \leq \epsilon)$

- ϵ : the radius threshold
 - μ : the weight threshold

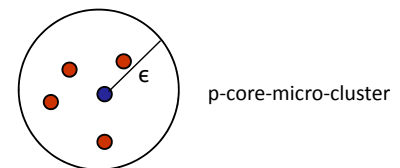


- But, in an evolving stream, the role of clusters and outliers often interchange:

- Should provide opportunity for the gradual growth of new clusters
 - Should promptly get rid of the outliers

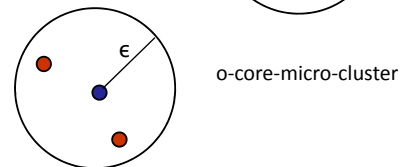
- Potential core micro-clusters

- $(w \geq \beta * \mu) \ \& \ (r \leq \epsilon), \ 0 < \beta \leq 1$



- Outlier micro-clusters

- $(w < \beta * \mu) \ \& \ (r \leq \epsilon), \ 0 < \beta \leq 1$



- Two lists of p-micro-clusters and o-micro-clusters are maintained over time
- **Initialize**: apply DBSCAN over *initPoints* → p-micro-clusters
- **Online** step as a new point *p* arrives from the stream
 - Try to assign it to its closest p-micro-cluster
 - If this is not possible, try to assign it to its closest o-micro-cluster
 - Check if this o-micro-cluster can be upgraded
 - If both are not possible, create a new o-micro-cluster with *p*
- **Periodic** micro-cluster maintenance based on data ageing
- **Offline** macro-clustering
 - On demand, upon user request apply a modified version of DBSCAN over p-micro-clusters to derive the final clusters

Two lists of p-micro-clusters and o-micro-clusters are maintained over time

- When a new point d arrives
 - Find its closest p-micro-cluster $pclu$
 - If the updated radius of $pclu \leq \epsilon$, merge d to $pclu$
 - o.w. find its closest o-micro-cluster $oclu$
 - If the updated radius of $oclu \leq \epsilon$, merge d to $oclu$
 - Check if $oclu$ can be upgraded to a p-micro-cluster (if now $w \geq \beta * \mu$)
 - o.w., create a new o-micro-cluster with d
- Periodic p- and o-micro-clusters update due to ageing
 - Delete a p-micro-cluster when $w < \beta * \mu$
 - Delete an o-micro-cluster when $w < \xi$ (expected weight based on its creation time)
 - The longer an o-micro-cluster exists, the higher its weight is expected to be
 - The number of p- and o-micro-clusters is bounded

• Upon request, apply a variant of DBSCAN over the set of online maintained p-micro-clusters

– Each p-micro-cluster c_p is treated as a virtual point located at the center of c_p with weight w .

• Core-micro-clusters (redefined)

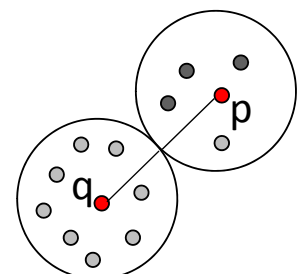
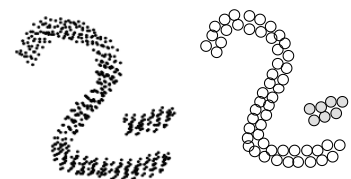
• Directly density reachable (redefined)

- c_p is directly density reachable from c_q if:
 - c_q is a c-micro-cluster and
 - $\text{dist}(c_p, c_q) \leq 2\epsilon$ (i.e. they are tangent or intersecting)

• Density reachable (redefined)

– A p-micro-cluster c_p is density reachable from a p-micro-cluster c_q if there is a chain of p-micro-clusters $c_{p1}=c_q, c_{p2}, \dots, c_{pn}=c_p$.

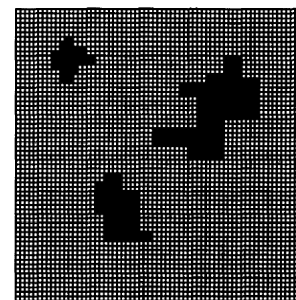
• Density connected (redefined)



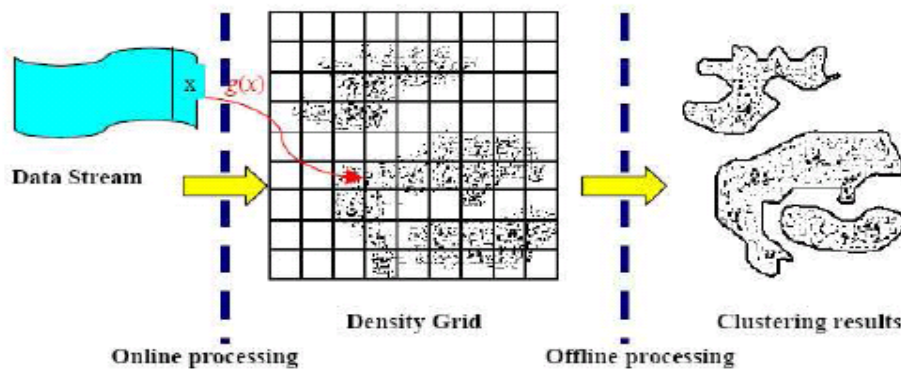
- + DenStream clusters large evolving data stream
- + Discover clusters of arbitrary shapes, following the density-based paradigm
- + No assumption on the number of clusters
- + Noise/ outlier handling

- The choice of the parameters ϵ , β , μ
- Constant parameters over time, what about clusters with different density

- A grid structure is used to capture the density of the dataset.
 - A cluster is a set of connected dense cells
 - e.g. STING
- Appealing features
 - No assumption on the number of clusters
 - Discovering clusters of arbitrary shapes
 - Ability to handle outliers
- In case of streams
 - The grid cells “constitute” the summary structure
 - Update the grid structure as the stream proceeds
 - DStream (Chen & Tu, 2007)



- The online-offline rationale is followed:
 - Online mapping of the new data into the grid
 - Offline computation of grid density and clustering of dense cells

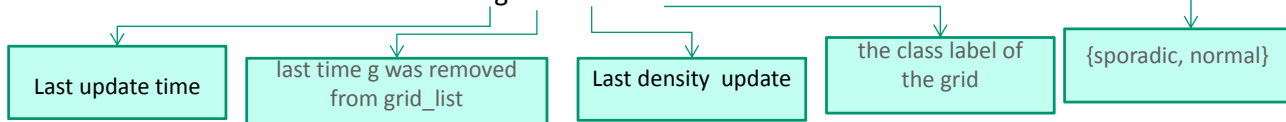


- Data ageing (damped window model):
 - $D(x,t) = \lambda^{t-t_c}$, t_c is the arrival time for point x , t is the current timepoint
 - λ in $(0,1)$ is the *decay factor*

- The density of a grid cell g at time t :
$$D(g,t) = \sum_{x \in E(g,t)} D(x,t)$$

- The characteristic vector of a grid cell g is defined as:

$(t_g, t_m, D, \text{label}, \text{status})$

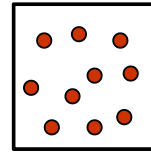


- The grid density can be updated incrementally

$$D(g, t_n) = \lambda^{t_n - t_l} D(g, t_l) + 1 \quad t_n: \text{the new record arrival time; } t_l: \text{the last record arrival}$$

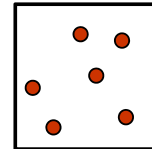
- The density of a grid is constantly changing over time.
- Dense grid cells

$$D(g, t) \geq \frac{C_m}{N(1-\lambda)} = D_m \quad C_m > 1$$



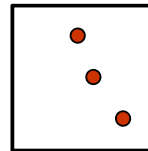
- Transitional grid cells

$$D(g, t) \leq \frac{C_l}{N(1-\lambda)} = D_l \quad 0 < C_l < 1$$



- Sparse grid cells

$$\frac{C_l}{N(1-\lambda)} \leq D(g, t) \leq \frac{C_m}{N(1-\lambda)}$$



N: #cells in the grid

- procedure D-Stream
- $t_c = 0;$
- initialize an empty hash table *grid_list*;
- while** data stream is active **do**
- read record $x = (x_1, x_2, \dots, x_d);$
- determine the density grid g that contains $x;$
- if (g not in *grid_list*) insert g to *grid_list*;
- update the characteristic vector of $g;$
- if** $t_c == gap$ **then**
- call *initial_clustering*(*grid_list*);
- end if**
- if** $t_c \bmod gap == 0$ **then**
- detect and remove sporadic grids from *grid_list*;
- call *adjust_clustering*(*grid_list*);
- end if**
- $t_c = t_c + 1;$
- end while**
- end_procedure**

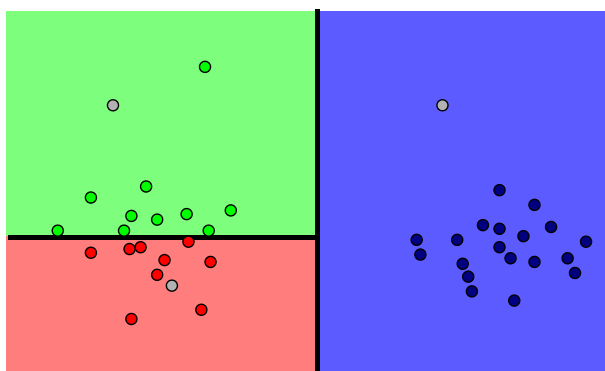
Grid update

Initialization

Clustering

- + DStream clusters large evolving data stream
- + It can discover clusters of arbitrary shapes
- + No assumption on the number of clusters
- + Distinguishes noise and outliers
- + The grid provides a level of abstraction over the data

- The choice of the grid parameters
- Fixed grid parameters

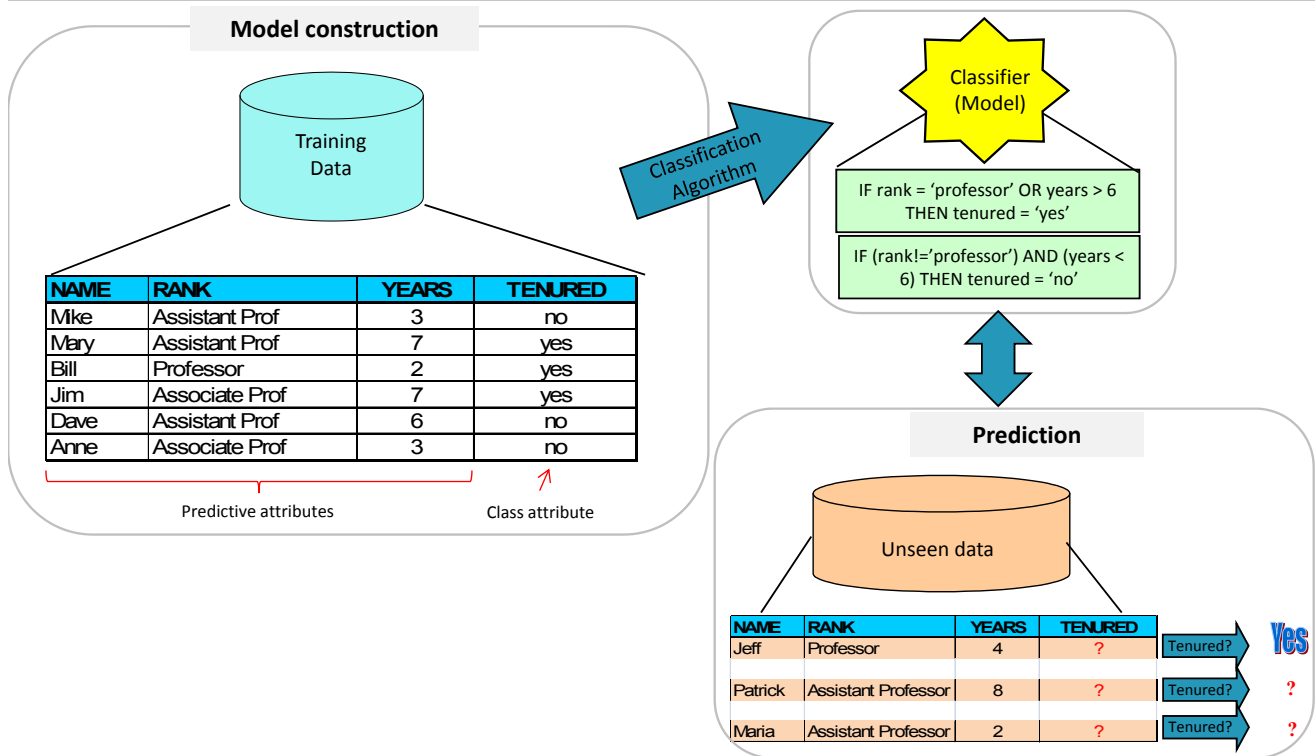


- Screw
 - Nails
 - Paper clips
- } Training data
- New object

Task:

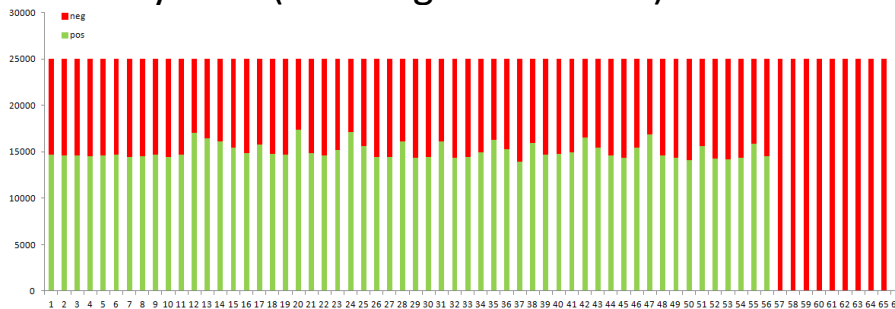
Learn from the already classified training data, the rules to classify new objects based on their characteristics.

The result attribute (class variable) is nominal (categorical)



- So far, classification as a batch/ static task
 - The whole training set is given as input to the algorithm for the generation of the classification model
 - When the performance of the model drops, a new model is generated from scratch over a new training set
- But, in a dynamic environment data change continuously
 - Batch model re-generation is not appropriate/sufficient anymore
- Need for new classification algorithms that
 - update existing models by incorporating new data
 - deal with non-stationary data generation processes (concept drift)
 - subject to:
 - Resource constraints (processing time, memory)
 - Single scan of the data (one look, no random access)

- Non-stationary data (evolving distribution)



Evolving data distribution (from the BA of Alina's Sinelnikova, "Sentiment analysis in the Twitter stream", LMU 2012.)

- Concept drift

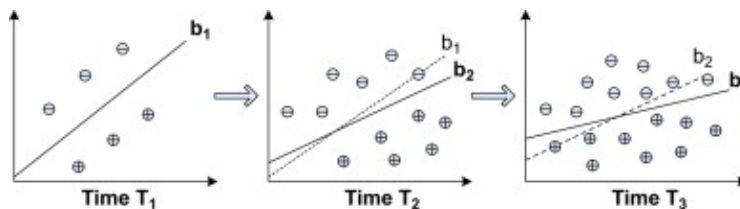


Fig. 1. An illustration of concept drifting in data streams. In the three consecutive time stamps T_1 , T_2 and T_3 , the classification boundary gradually drifts from b_1 to b_2 and finally to b_3 . (from: *A framework for application-driven classification of data streams*, Zhang et al, Journal Neurocomputing 2012)

- The batch classification problem:
 - Given a **finite** training set $D = \{(x, y)\}$, where $y = \{y_1, y_2, \dots, y_k\}$, $|D| = n$, find a function $y = f(x)$ that can predict the y value for an unseen instance x
- The data stream classification problem:
 - Given an **infinite** sequence of pairs of the form (x, y) where $y = \{y_1, y_2, \dots, y_k\}$, find a function $y = f(x)$ that can predict the y value for an unseen instance x
- Example applications:
 - Fraud detection in credit card transactions
 - Churn prediction in a telecommunication company
 - Sentiment classification in the Twitter stream
 - Topic classification in a news aggregation site, e.g. Google news
 - ...

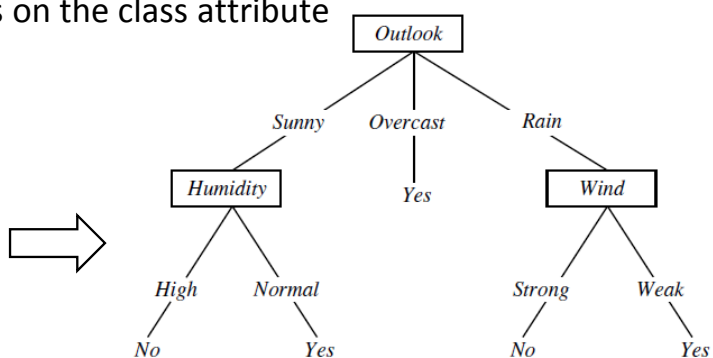
- Decision trees
- Ensemble methods

(Batch) Decision Trees (DTs)

- Training set: $D = \{(x,y)\}$
 - predictive attributes: $x = \langle x_1, x_2, \dots, x_d \rangle$
 - class attribute: $y = \{y_1, y_2, \dots, y_k\}$
- Goal: find $y = f(x)$
- Decision tree model
 - nodes contain tests on the predictive attributes
 - Leaves contain predictions on the class attribute

Training set

Day	Outlook	Temperature	Humidity	Wind	PlayTennis
D1	Sunny	Hot	High	Weak	No
D2	Sunny	Hot	High	Strong	No
D3	Overcast	Hot	High	Weak	Yes
D4	Rain	Mild	High	Weak	Yes
D5	Rain	Cool	Normal	Weak	Yes
D6	Rain	Cool	Normal	Strong	No
D7	Overcast	Cool	Normal	Strong	Yes
D8	Sunny	Mild	High	Weak	No
D9	Sunny	Cool	Normal	Weak	Yes
D10	Rain	Mild	Normal	Weak	Yes
D11	Sunny	Mild	Normal	Strong	Yes
D12	Overcast	Mild	High	Strong	Yes
D13	Overcast	Hot	Normal	Weak	Yes
D14	Rain	Mild	High	Strong	No



- Basic algorithm (ID3, Quinlan 1986)
 - Tree is constructed in a top-down recursive divide-and-conquer manner
 - At start, all the training examples are at the root node

Main loop:

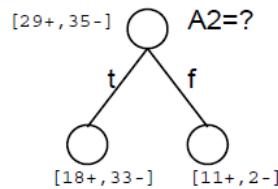
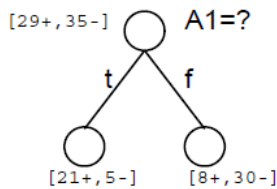
1. $A \leftarrow$ the “best” decision attribute for next *node*
2. Assign A as decision attribute for *node*
3. For each value of A , create new descendant of *node*
4. Sort training examples to leaf nodes
5. If training examples perfectly classified, Then STOP, Else iterate over new leaf nodes

Attribute selection measures:

- Information gain
- Gain ratio
- Gini index

(check Lecture 4, KDD I)

- But, which attribute is the best?



Goal: select the most useful attribute for classifying examples.

- useful \rightarrow the resulting partitioning is as pure as possible
- pure partition: all its instances belong to the same class.

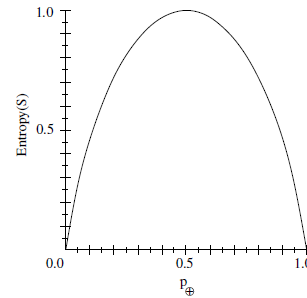
- Used in ID3
- It uses entropy, a measure of pureness of the data
- The information gain $Gain(S,A)$ of an attribute A relative to a collection of examples S measures the gain reduction in S due to splitting on A :

$$Gain(S, A) = \underbrace{Entropy(S)}_{\text{Before splitting}} - \sum_{v \in \text{values}(A)} \frac{|S_v|}{|S|} \underbrace{Entropy(S_v)}_{\text{After splitting on A}}$$

- Gain measures the expected reduction in entropy due to splitting on A
- The attribute with the higher entropy reduction is chosen

- Let S be a collection of positive and negative examples for a binary classification problem, $C=\{+, -\}$.
- p_+ : the percentage of positive examples in S
- p_- : the percentage of negative examples in S
- Entropy measures the impurity of S :

$$Entropy(S) = -p_+ \log_2(p_+) - p_- \log_2(p_-)$$



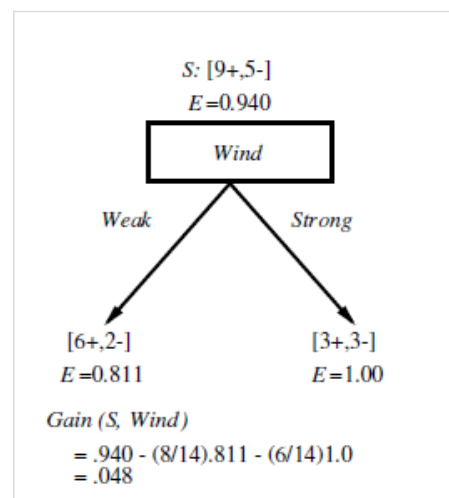
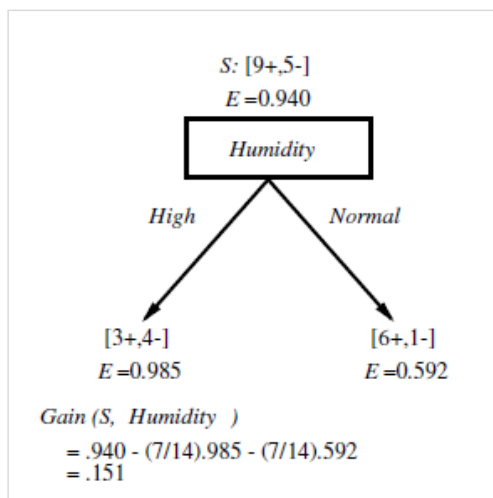
Examples :

- Let $S: [9+,5-]$ $Entropy(S) = -\frac{9}{14} \log_2(\frac{9}{14}) - \frac{5}{14} \log_2(\frac{5}{14}) = 0.940$
- Let $S: [7+,7-]$ $Entropy(S) = -\frac{7}{14} \log_2(\frac{7}{14}) - \frac{7}{14} \log_2(\frac{7}{14}) = 1$
- Let $S: [14+,0-]$ $Entropy(S) = -\frac{14}{14} \log_2(\frac{14}{14}) - \frac{0}{14} \log_2(\frac{0}{14}) = 0$

in the general case
(k -classification problem)
 $Entropy(S) = \sum_{i=1}^k -p_i \log_2(p_i)$

- Entropy = 0, when all members belong to the same class
- Entropy = 1, when there is an equal number of positive and negative examples

- Which attribute to choose next?



- Idea: In order to pick the best split attribute for a node, it may be sufficient to consider only a small subset of the training examples that pass through that node.
 - No need to look at the whole dataset (which is infinite in case of streams)
 - E.g., use the first few examples to choose the split at the root
- Problem: How many instances are necessary?
 - Hoeffding bound!
- Hoeffding tree variations
 - Very Fast DTs (VFDT), Domingos and Hulten, KDD 2000.
 - VFDTc, Gamma et al, KDD 2003
 - ...

- Consider a real-valued random variable r whose range is R
 - e.g., for a probability the range is one,
 - for an information gain the range is $\log_2(c)$, where c is the number of classes
- Suppose we have n independent observations of r and we compute its mean \bar{r}
- The Hoeffding bound states that with probability $1-\delta$ the true mean of the variable μ_r will not differ by more than ϵ from the estimated mean after n independent observations, i.e., $P(\mu_r \geq \bar{r} - \epsilon) = 1-\delta$, where:

$$\epsilon = \sqrt{\frac{R^2 \ln(1/\delta)}{2n}}$$

- This bound holds true regardless of the distribution generating the values, and depends only on the range of values, number of observations and desired confidence.
 - A disadvantage of being so general is that it is more conservative than a distribution-dependent bound

- Let $G()$ be the heuristic measure for choosing the split attribute at a node
- After seeing n instances at this node, let
 - X_a : be the attribute with the highest observed $G()$
 - X_b : be the attribute with the second-highest observed $G()$
- $\overline{\Delta G} = \overline{G(X_a)} - \overline{G(X_b)} \geq 0$ the difference between the 2 best attributes
- $\overline{\Delta G}$ is the random variable being estimated by the Hoeffding bound
- Given a desired δ , the Hoeffding bound guarantees that with probability $1-\delta$, X_a is the correct choice for this node, if n instances have been seen so far in this node and $\overline{\Delta G} > \epsilon$.
 - In this case, the sample size is enough to decide on X_a
- Otherwise, i.e., if $\overline{\Delta G} < \epsilon$, the sample size is not enough for a stable decision.
 - With R and δ fixed, the only variable left to change ϵ is n
 - We need to extend the sample by seeing more instances, until ϵ becomes smaller than $\overline{\Delta G}$

- **Input:** δ desired probability level.
- **Output:** \mathcal{T} A decision Tree
- **Init:** $\mathcal{T} \leftarrow$ Empty Leaf (Root)
- While (TRUE)
 - Read next Example
 - Propagate Example through the Tree from the Root till a leaf
 - Update Sufficient Statistics at leaf
 - If $leaf(\#examples) > N_{min}$
 - Evaluate the merit of each attribute
 - Let A_1 the best attribute and A_2 the second best
 - Let $\epsilon = \sqrt{R^2 \ln(1/\delta) / (2n)}$
 - If $G(A_1) - G(A_2) > \epsilon$
 - Install a splitting test based on A_1
 - Expand the tree with two descendant leaves

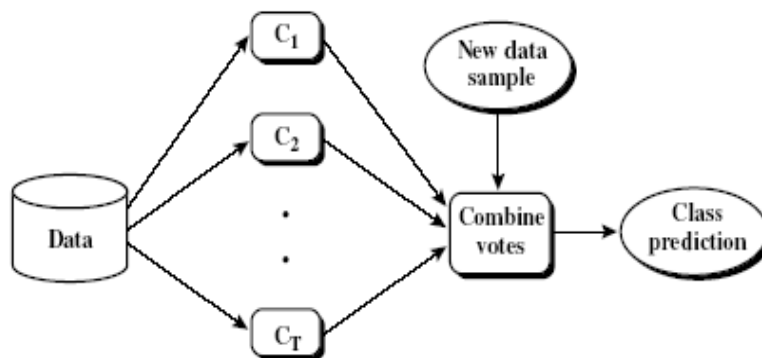
Those needed by the heuristic evaluation function $G()$

The evaluation of $G()$ after each instance is very expensive.
 → Evaluate $G()$ only after N_{min} instances have been observed since the last evaluation.

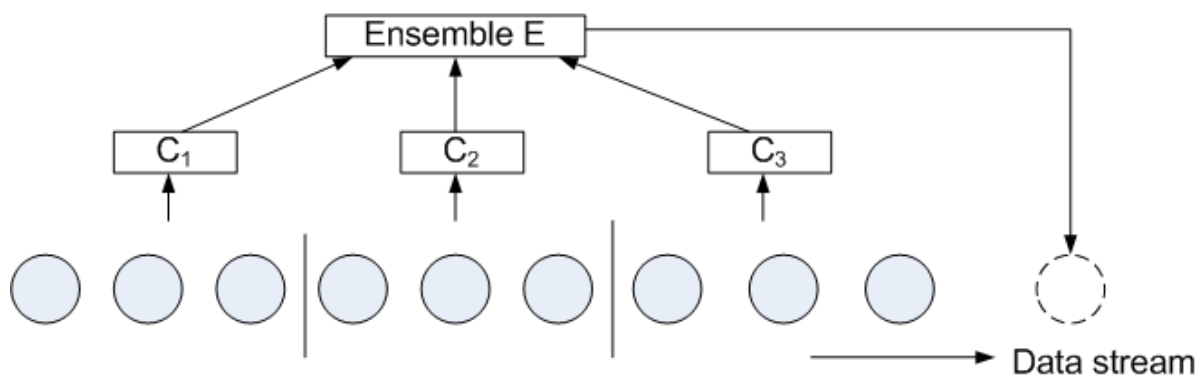
- Hoeffding tree
 - Leaf nodes are transformed into decision nodes by splitting on an attribute
 - Hoeffding trees exploit the fact that a small sample can often be enough to choose an optimal splitting attribute
 - Hoeffding bound
- Adaptive size Hoeffding tree (ASHT)
 - The tree has a maximum size (# of splitting nodes)
 - After one node splits, if the number of split nodes of the ASHT tree is higher than the maximum value, then it deletes some nodes to reduce its size

- Decision trees
- Ensemble methods

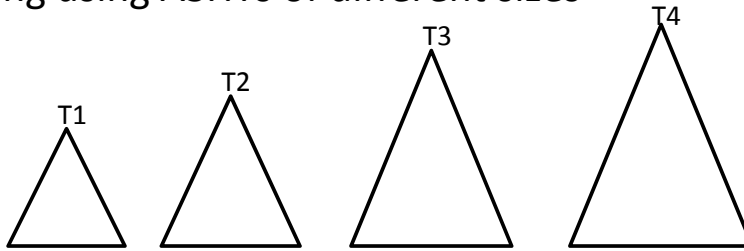
- Idea:
 - Instead of a single model, use a combination of models to increase accuracy
 - Combine a series of T learned models, M_1, M_2, \dots, M_T , with the aim of creating an improved model M^*
 - To predict the class of previously unseen records, aggregate the predictions of the ensemble



- Idea: Use a **divide-and-conquer** manner to build models from continuous data streams



- Bagging using ASHTs of different sizes



- Smaller trees adapt more quickly to changes
- Larger trees perform better during periods with no or little change
- The max allowed size for the n^{th} ASHT tree is twice the max allowed size for the $(n-1)^{\text{th}}$ tree.
- Each tree has a weight proportional to the inverse of the square of its error
- The goal is to increase bagging performance by tree diversity

- The quality of a classifier is evaluated over a *test set*, different from the training set
- For each instance in the test set, we know its true class label
- Compare the predicted class (by some classifier) with the true class of the test instances
- Terminology
 - Positive tuples: tuples of the main class of interest
 - Negative tuples: all other tuples
- A useful tool for analyzing how well a classifier performs is the *confusion matrix*
- For an m -class problem, the matrix is of size $m \times m$
- An example of a matrix for a 2-class problem:

		Predicted class		totals
		C_1	C_2	
Actual class	C_1	TP (true positive)	FN (false negative)	P
	C_2	FP (false positive)	TN (true negative)	N
Totals		P'	N'	

- Accuracy/ Recognition rate:
 - % of test set instances correctly classified

$$accuracy(M) = \frac{TP + TN}{P + N}$$

	C ₁	C ₂	totals
C ₁	TP (true positive)	FN (false negative)	P
C ₂	FP (false positive)	TN (true negative)	N
Totals	P'	N'	

classes	buy_computer = yes	buy_computer = no	total	recognition(%)
buy_computer = yes	6954	46	7000	
buy_computer = no	412	2588	3000	
total	7366	2634	10000	95.42

- Error rate/ Missclassification rate: $error_rate(M) = 1 - accuracy(M)$

$$accuracy(M) = \frac{FP + FN}{P + N}$$

- More effective when the class distribution is relatively balanced
 - Check Lecture 4, KDD I for more evaluation measures also if classes are imbalanced!

- Holdout method
 - Given data is randomly partitioned into two independent sets
 - Training set (~2/3) for model construction, Test set (~1/3) for evaluation
- Cross-validation (k-fold cross validation, k = 10 usually)
 - Randomly partition the data into k mutually exclusive subsets D₁, ..., D_k each approximately equal size
 - Training and testing is performed k times
 - At the i-th iteration, use D_i as test set and others as training set
 - Accuracy is the avg accuracy over all iterations
- Bootstrap: Samples the given training data uniformly with replacement
 - i.e., each time a tuple is selected, it is equally likely to be selected again and re-added to the training set
- Check Lecture 4, KDD I for more evaluation methods, their pros and cons!

- Holdout evaluation
 - 2 separate datasets for training (~70% - 80% of the dataset) and testing (~20%-30% of the dataset)
 - Train model on training set
 - Test model on test set
 - Static test set!!!

- Prequential evaluation (Interleaved test-then-train)
 - One dataset for training and testing
 - Models are first tested then trained in each instance
 - Test set is dynamic!!!

- Accuracy

- Kappa measure
 - normalizes the accuracy of a classifier p_0 by that of a chance predictor p_c

$$k = \frac{p_0 - p_c}{1 - p_c}$$

Kappa value	Classifier's performance
0%-20%	bad
21%-40%	fair
41%-60%	moderate
61%-80%	substantial
81%-100%	(almost) perfect

- Both measures are computed based on the most recent samples through some
 - sliding window
 - fading function

- J. Gama, Knowledge Discovery from Data Streams, Chapman and Hall/CRC, 2010.
- C. Aggarwal, Data Streams: Models and Algorithms, Springer, 2007.
- C. C. Aggarwal, J. Han, J. Wang, P. S. Yu: A framework for clustering evolving data streams. VLDB, 2003.
- F. Cao, M. Ester, W. Qian, A. Zhou: Density-Based Clustering over an Evolving Data Stream with Noise. SDM, 2006.
- Y. Chen, L. Tu: Density-Based Clustering for Real-Time Stream Data. KDD, 2007.
- F. Farnstrom, J. Lewis, C. Elkan: Scalability for clustering algorithms revisited. ACM SIGKDD Explorations Newsletter 2(1):51-57, 2000.
- S. Guha, A. Meyerson, N. Mishra, R. Motwani, L. O' Callaghan: Clustering data streams: Theory and practice. IEEE TKDE 15(3):515–528, 2003.
- L. O'Callaghan, N. Mishra, A. Meyerson, S. Guha, R. Motwani: Streaming-Data Algorithms for High-Quality Clustering. ICDE, 2002.
- www.utdallas.edu/~lkhan/Spring2008G/Charu03.ppt