

## Buzz prediction on Twitter

This data-set contains examples of buzz events from Twitter about some topics.

### Objective

We want to predict if a topic could gain popularity in a social network.

### Predicted attribute :

“Annotation” (Buzz). This attribute is binary: 1 meaning “buzz observed”, 0 meaning “no buzz observed”. It is stored in the rightmost column.

### Relevant Information Paragraph:

Each instance covers seven days of observation for a specific topic (eg. overclocking...). Considering the couple day following this initial observation; if there is at least 500 additional active discussions by day (on average, with respect to the initial observation) then, the predicted attribute Buzz is True.

Observations are Independent and identically distributed.

Number of Instances: 126635

Number of Attributes: 77 + class attribute

Time representation: Each instance is described by 77 features, those describe the evolution of 11 “primary features” through time. Hence each feature name is postfixed with the relative time of observation. For instance, the value of the feature “Nb\_Active\_Discussion” at time t is given in “Nb\_Active\_Discussion\_t”.

## Attributes

### -- Number of Created Discussions (NCD) (columns [0,6])

-- Type : Numeric, integers only

-- Description : This feature measures the number of discussions created at time step t and involving the instance's topic.

-- Columns : From column 0 (NCD at relative time 0) to column 6 (NCD at relative time 6)

-- Abbreviations : NCD\_0, NCD\_1, NCD\_2, NCD\_3, NCD\_4, NCD\_5, NCD\_6

### -- Author Increase (AI) (columns [7,13])

-- Type : Numeric, integers only

-- Description : This feature measures the number of new authors interacting on

the instance's topic at time t (i.e. its popularity)

-- Columns : From column 7 (AI at relative time 0) to column 13 (AI at relative time 6)

-- Abbreviations : AI\_0, AI\_1, AI\_2, AI\_3, AI\_4, AI\_5, AI\_6

-- **Attention Level (measured with number of authors) (AS(NA))**  
(columns [14,20])

-- Type : Numeric, real in [0,1]

-- Description : This feature is a measure of the attention payed to a the instance's topic on a social media.

-- Columns : From column 14 (AS(NA) at relative time 0) to column 20 (AS(NA) at relative time 6)

-- Abbreviations : AS(NA)\_0, AS(NA)\_1, AS(NA)\_2, AS(NA)\_3, AS(NA)\_4, AS(NA)\_5, AS(NA)\_6

-- **Burstiness Level (BL) (columns [21,27])**

-- Type : Numeric, defined on [0,1]

-- Description : The burstiness level for a topic z at a time t is defined as the ratio of ncd and nad

-- Columns : From column 21 (BL at relative time 0) to column 27 (BL at relative time 6)

-- Abbreviations : BL\_0, BL\_1, BL\_2, BL\_3, BL\_4, BL\_5, BL\_6

-- **Number of Atomic Containers (NAC) (columns [28,34])**

-- Type : Numeric, integer

-- Description : This feature measures the total number of atomic containers generated through the whole social media on the instance's topic until time t.

-- Columns : From column 28 (NAC at relative time 0) to column 34 (NAC at relative time 6)

-- Abbreviations : NAC\_0, NAC\_1, NAC\_2, NAC\_3, NAC\_4, NAC\_5, NAC\_6

-- **Attention Level (measured with number of contributions) (AS(NAC))**  
(columns [35,41])

-- Type : Numeric, real in [0,1]

-- Description : This feature is a measure of the attention payed to a the instance's topic on a social media.

-- Columns : From column 35 (AS(NA) at relative time 0) to column 42 (AS(NAC) at relative time 6)

-- Abbreviations : AS(NAC)\_0, AS(NAC)\_1, AS(NAC)\_2, AS(NAC)\_3, AS(NAC)\_4, AS(NAC)\_5, AS(NAC)\_6

#### -- **Contribution Sparseness (CS) (columns [42,48])**

-- Type : Numeric, real in [0,1]

-- Description : This feature is a measure of spreading of contributions over discussion for the instance's topic at time t.

-- Columns : From column 42 (CS at relative time 0) to column 48 (CS at relative time 6)

-- Abbreviations : CS\_0, CS\_1, CS\_2, CS\_3, CS\_4, CS\_5, CS\_6

#### -- **Author Interaction (AT) (columns [49,55])**

-- Type : Numeric, integer.

-- Description : This feature measures the average number of authors interacting on the instance's topic within a discussion.

-- Columns : From column 49 (AT at relative time 0) to column 55 (AT at relative time 6)

-- Abbreviations : AT\_0, AT\_1, AT\_2, AT\_3, AT\_4, AT\_5, AT\_6

#### -- **Number of Authors (NA) (columns [56,62])**

-- Type : Numeric, integer.

-- Description : This feature measures the number of authors interacting on the instance's topic at time t.

-- Columns : From column 49 (NA at relative time 0) to column 55 (NA at relative time 6)

-- Abbreviations : NA\_0, NA\_1, NA\_2, NA\_3, NA\_4, NA\_5, NA\_6

#### -- **Average Discussions Length (ADL) (columns [63,69])**

-- Type : Numeric, real.

-- Description : This feature directly measures the average length of a

discussion belonging to the instance's topic.

-- Columns : From column 63 (ADL at relative time 0) to column 69 (ADL at relative time 6)

-- Abbreviations : ADL\_0, ADL\_1, ADL\_2, ADL\_3, ADL\_4, ADL\_5, ADL\_6

**-- Average Discussions Length (NAD) (columns [70,76])**

-- Type : Numeric, integer.

-- Description : This features measures the number of discussions involving the instance's topic until time t.

-- Columns : From column 70 (NAD at relative time 0) to column 76 (NAD at relative time 6)

-- Abbreviations : NAD\_0, NAD\_1, NAD\_2, NAD\_3, NAD\_4, NAD\_5, NAD\_6

**-- Annotation (column 77)**

-- Type : Numeric, integer: 0 or 1

-- Description : class attribute

-- Columns : 77

-- Buzz = 1

-- Non Buzz = 0