

Data analysis was dealt at two levels, one considering player’s information given its performance in a single match (match level) and one taking into account player’s statistics over multiple games (player level). In order to distinguish between these cases, two distinct datasets were employed,  $\mathcal{M}$  and  $\mathcal{P}$  respectively. The first one was created with data gathered directly from the Open Dota platform and the second one by aggregating such data on the base of singular users. Moreover, for our analysis we also employed the  $\overline{\mathcal{M}}$  datasets by imposing a constraint on the entries of  $\mathcal{M}$  and by adding some features thanks to our domain knowledge. We now discuss the description of the features utilized in each of these datasets.

## 1 Match Features

The features illustrated in this Section are those used for the  $\mathcal{M}$  dataset.

### Panoramic

Panoramic features provide an overview of the match at its conclusion. It includes ids of the game, the player and the hero employed, some numbers that supply a shallow summary of the performance statistics, party-related information, and reports about the objects equipped. Instead of directly consider which items are present in player’s inventory and backpack, we categorized them on the base of their purposes within the game, highlighting characteristics that we identified as meaningful. Besides that, we also added a *build* feature, expressed as a number between 0 and 3, in order to better understand whether the items utilized are more oriented toward a defensive (value closer to 0) or an offensive (value closer to 3) gaming style. Features comprehended in this section therefor are:

- match\_id
- account\_id
- hero\_id
- player\_slot
- level

- kills
- deaths
- assists
- last\_hits
- denies
- net\_worth
- hero\_damage
- hero\_healing
- support\_cons
- build
- farm
- mobility
- incomplete
- party\_id
- party\_size

Note that attributes `party_id` and `player_slot` are not used in the predictions, since they are only needed to extract more meaningful information about party size and chat analysis respectively.

## Benchmarks

This kind of attributes take care of comparison with respect to recent performances on the hero played, considering not just the matches of the player under examination but every match of Dota 2. The “comparison” is achieved by storing two values for each feature, one that indicates the actual player’s score (raw value) and one that shows the percentage of recent games with that same hero with a score worse than that (pct value). In here are collected:

- gold\_per\_min.raw
- gold\_per\_min.pct
- xp\_per\_min.raw
- xp\_per\_min.pct
- kills\_per\_min.raw
- kills\_per\_min.pct
- last\_hits\_per\_min.raw
- last\_hits\_per\_min.pct
- hero\_damage\_per\_min.raw
- hero\_damage\_per\_min.pct
- hero\_healing\_per\_min.raw
- hero\_healing\_per\_min.pct
- stuns\_per\_min.raw
- stuns\_per\_min.pct
- tower\_damage.raw
- tower\_damage.pct
- lhten.raw
- lhten.pct

## Performance

Performance features are used to evaluate gaming skills by collecting:

- multi\_kills
- kill\_streaks

- stuns
- camps\_stacked
- life\_state\_dead
- buybacks
- pings

## Laning

These attributes are trivially associated with lane role, which can be roaming, safe, mid, or off. Since this aspect is particularly crucial in early phase of the match, lane efficiency (i.e., percentage of lane gold obtained) and denials are considered at ten minutes. Usually also last hits at ten are reported, however we already accounted for such value in the benchmarks section and therefore we avoid to consider that again to prevent duplicates. Thus, the features belonging to this section are simply:

- lane\_role
- lane\_efficiency
- dn@10

## Farm

Farm features provide more specific details about the kills and are respectively:

- hero\_kills
- creeps
- neutral\_kills
- ancient\_kills
- tower\_kills
- courier\_kills

- roshan\_kills
- observer\_kills
- necronomicon\_kills

## Actions

Such attributes lay out the amount of times a player executed a certain kind of action, providing a more elaborated distinction among them. Specifically, they are divided on the base of the type (move, attack or cast spell) and on the target they are aiming to:

- actions\_per\_minute
- action\_move\_to\_position
- action\_move\_to\_target
- action\_attack\_move
- action\_attack\_target
- action\_cast\_position
- action\_cast\_target
- action\_cast\_no\_target
- action\_hold\_position
- action\_radar

## Objectives

Illustrate the amount of damage dealt to towers, barracks, ancient, and Roshan. In the case of the first two, data is reported for each independent structure, with code names based on their position:

- t1
- m1

- b1
- t2
- m2
- b2
- t3
- m3
- b3
- raxt
- raxm
- raxb
- t4
- ant
- roshan

## Runes

Features about runes simply outline the number of times a certain type of rune was used by the player. Each rune has different effects and spawn on fixed positions of the game map. In here are gathered:

- double damage
- haste
- illusion
- invisibility
- regeneration
- bounty
- arcane

## Vision

Vision-related attributes are associated with the purchase and usage of two particular items that are the observer and the sentry; Both of them are used to make visible something that normally would not be, that are ground and invisible enemies respectively. Also features about their average duration are gathered, since, despite their time limited usage, they can be eliminated by enemies when discovered:

- purchase\_observer
- use\_observer
- duration\_observer
- purchase\_sentry
- use\_sentry
- duration\_sentry
- purchase\_dust
- use\_smoke\_of\_deceit
- pur\_smoke\_of\_deceit
- pur\_gem

## Cosmetics

*Cosmetics prices* is the only feature comprised into this section. Note that purchasable items are for customization only (they do not provide any in-game advantage).

## Chat

Features about chats assemble numeric information about player's chat usage, differentiating messages between radial (pre-written messages) and typed ones. A variety of values describing the amount times chat is utilized, messages lengths and chat ranking therefore are collected with the purpose to evaluate how much the player interacts with others through chat:

- msg\_tot
- msg\_radial
- msg\_chat
- msg\_tot\_ratio
- msg\_radial\_ratio
- msg\_chat\_ratio
- msg\_char\_sum
- len\_msg\_avg
- len\_msg\_mode
- len\_msg\_median
- rank\_msg\_tot
- rank\_chat
- rank\_radial

## Extra Features

In this section are gathered a range of different useful attributes that do not belong to other groups, some information derived by metadata and values slightly re-elaborated from those originally extrapolated by Open Dota, in particular about the hero employed for the match (e.g., its role, complexity, etc.):

- rank
- dota\_plus
- randomed
- weekday
- hour



- Carry
- Support
- Nuker
- Disabler
- Pusher
- Initiator
- Escape
- Durable
- Jungler
- hero\_primary\_attr
- hero\_attack\_type
- hero\_complexity
- hero\_gender
- hero\_species
- lobby\_type
- radiant\_win
- win
- teamfight\_participation
- perc\_gold\_spent

## 2 Match Features with Domain Knowledge

For the creation of dataset  $\overline{\mathcal{M}}$ , we derived some additional features from chat analysis. We adopted two strategies:

1. **Categorization of chat messages** (i.e., wrote by the players). Here we defined common words used in different contexts, i.e., laugh, slang, bad and good behaviour, and provocative messages at the end of the game (e.g., “gg ez”). We also searched for messages containing only “?” (which is highly provocative for the other team), number of question marks, exclamative marks, capital letters, early messages (sent before the game begins, usually sent to make noise or interact with the other team), and messages after a kill (which can be for complaining, provoking, or similar reasons). We identified such words by exploring websites (e.g., Dota2 forums, urban dictionary), chats of multiple games, and by our gaming experience. Then, we count the occurrences of such words in the chats typed by the player. The features obtained according to this procedure are:
  - counterLaugh
  - counterThank
  - counterDeny
  - ggez\_counter
  - bad\_counter
  - slang\_counter
  - question\_mark\_counter
  - counter\_capital
  - counter\_question
  - counter\_exclamation
  - earlySpeak
  - spam\_kill
2. **Categorization of chatweel messages** (i.e., pick from pre-defined messages, and automatically translated by the game to the receiver’s language). Here we categorized chatwheel messages that are pre-defined

by the game. These are useful since we do not incur in translation problems. We manually analyzed the file `chat_wheel.json` (containing, in english, all the chatwheel messages) and identified various types of messages, i.e., tactics, good behavior, and caster phrases. Moreover, we extracted which of them are reproduced as a sound (audible either in local or global chat), or are sprays left on the ground. Then, there is a second chatwheel, customised for each hero. Such messages contains general information, such as if they are laugh or deny messages, which we extracted as feature. Last, we extracted information such as whether the message was sent locally (same team) or globally (to both teams), and whether the message was an hero message. Therefore, we extracted the features:

- counterTattic
- counterGB
- counterCaster
- counterCasterIn
- counterSounds
- counterSoundsAllChat
- counterSoundsInChat
- counterAllChat
- hero\_msg

All the aforementioned features can be of extreme help for our task. For instance, the use of slang could be higher for younger players. Provocative messages can relate to both neurotic and extrovert people. Agreeableness could be related to the use of good behavior messages. Conscious people could use more tactics message, and openness could relate to early messages or laugh messages after a kill. The gender could be affected from several of such typing and hero chatwheel usage. Last, some hero messages are available only on purchase, which can be an indicator of the occupation and `buy_content` behavior. For this reason, the features reported in this Section were **added** to those presented for  $\mathcal{M}$  to create  $\overline{\mathcal{M}}$ .

### 3 Player Features

In the following we illustrate the manipulation process we used in order to build the  $\mathcal{P}$  dataset. It was created by aggregating data with the same *account\_id* within the  $\mathcal{M}$  database, described in the precedent Section, and using as new values the mean of such grouped data, therefore the numerical features are the same.

Some adjustments had to be done for **non numerical features**, since it's not possible to compute the average of Booleans nor categorical values:

- For the feature indicating *hero\_id* we kept the most frequently used parameter;
- Dropped *match\_id*, *party\_id*, *player\_slot*, and *radiant\_win*;
- Columns of extra features about hero's role were kept, but Boolean values were substituted with the percentage of matches in which the hero of choice adhered to the respective role in-game. The same was done for *dota\_plus*;
- Lane role, hour, and weekday were expanded in a sort of one-hot encoding fashion. That is, a column was created for each possible value of both features and the ratio between the number of games with such value and the total number of matches made by the player was used as new feature entry. The same transformation was applied to hero's attributes accounted in the extra features section of match level dataset. Therefore, the newly added columns are:

- *lane\_role\_mid*
- *lane\_role\_off*
- *lane\_role\_roaming*
- *lane\_role\_safe*
- *mon*
- *tue*
- *wed*
- *thu*
- *fri*

- sat
  - sun
  - night
  - morning
  - afternoon
  - evening
  - hero\_pr\_attr\_agi
  - hero\_pr\_attr\_int
  - hero\_pr\_attr\_str
  - hero\_melee\_atk
  - hero\_complexity\_1
  - hero\_complexity\_2
  - hero\_complexity\_3
  - hero\_gender\_A
  - hero\_gender\_F
  - hero\_gender\_M
  - hero\_species\_beast
  - hero\_species\_human
  - hero\_species\_humanoid
  - hero\_species\_spirit
- For each player we considered the three most used characters and report their mode for the values of primary attribute, attack type, complexity, gender, and species. We also added an attribute to represent the percentage of matches played with the top heroes over the total number of games. For limit cases we decided to simply perform the same procedure with the available number of heroes. Therefore, the added features are:
    - top\_hero\_pr\_attr
    - top\_hero\_atk\_type
    - top\_hero\_complexity

- top\_hero\_gender
  - top\_hero\_species
  - top\_usage
- we introduced a new feature to represent the *versatility* of the user, which highlights how much the player is prone to use different heroes and is expressed as a number between 0 and 1;
- information about lobby type (normal or ranked) was unified with the number of wins and, therefore, for every player are reported the amounts of:
  - normal\_games
  - normal\_win
  - ranked\_games
  - ranked\_win
  - win
  - games