



Buscando la paz interior...

twiter: @richardyantas5

TITLE:

Introduction to Data Science Competition

fuelle: www.prometec.net

Versión en L^AT_EX:
Msc(c). Richard Valentín Yantas Alcantaraaa

Contents

1	Introduction	1
1.1	Introductionnnn	1
1.1.1	Feature extractionnnñ from texts and image	1
1.1.2	Categorical and ordinal features	2
1.1.3	Competition Mechanics	2
1.1.4	Datetime and Coordinates	3
1.1.5	Handling missing	3
1.1.6	Handling missing	3
1.1.7	Kaggle overview	3
1.1.8	Numeric features	3
1.1.9	Recap of main algorithms	3
2	Exploratory Data Analysis	5
2.1	Exploratory Data Analysis	5
2.1.1	Building Intuition about the data	5
2.1.2	Data cleaning	5
2.1.3	Data splitting strategy	5
2.1.4	Exploratory data Analysis	6
2.1.5	Exploratory anonymized data	6
2.1.6	Validation and overfitting	6
2.1.7	Validation strategies	6
2.1.8	Visualizations	6

Chapter 1

Introduction

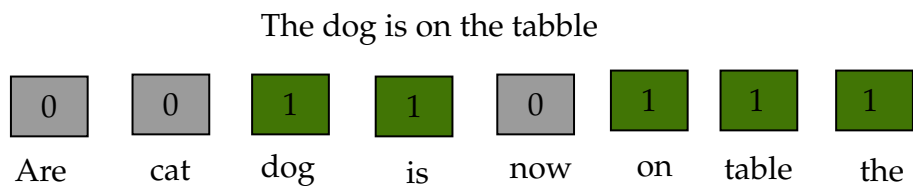
1.1 Introductionnnn

1.1.1 Feature extractionnnñ from texts and image

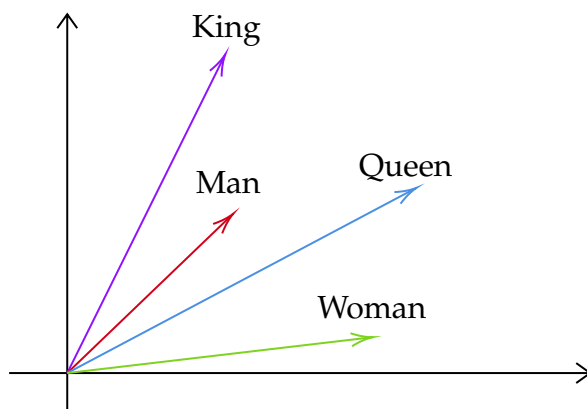
This is the titanic data sets

Text to vector

1. Bag of words:



2. Embeggins:



Text preprocessing

1. Lowercase
2. Lemmatization
3. Stemming

4. StopWords

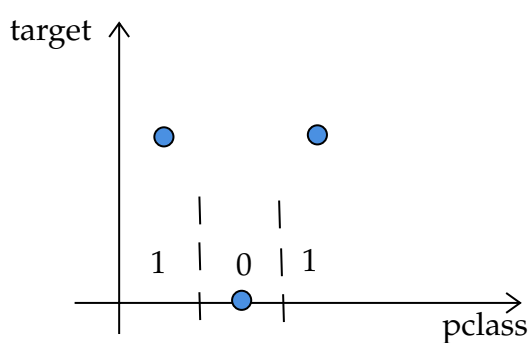
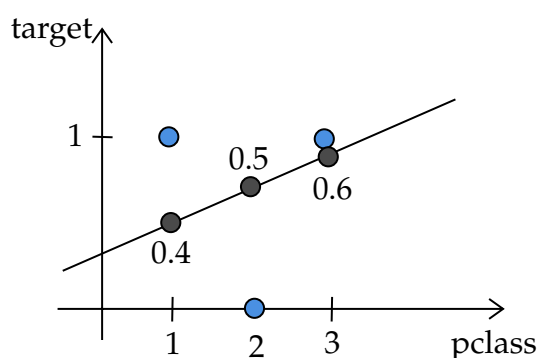
1.1.2 Categorical and ordinal features

Categorical

Ordinal Features

Label encoding

pclass	1	2	3
target	1	0	1



Frequency encoding

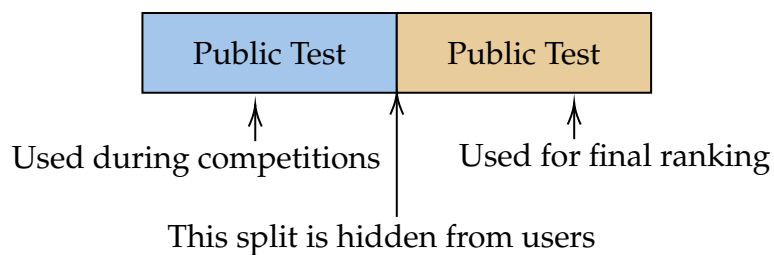
1.1.3 Competition Mechanics

1. Evaluation function Exists some evaluation functions like:

- Accuracy
- Logistic loss
- AUC
- RMSE
- MAE

2. Public/Private tests

You should submit predictions
for a whole test set



3. Sites:

- Kaggle
- Driven data
- CrowdAaqnalityx
- CodaLab
- DataScienceChallenge.net
- DataScience.net
- Single-Competition sites like KDD,VizDooM

4. Conclusions:

- Main Concepts(Data, Model, Submission, Evaluation, Leaderboard)
- Competition platforms
- Reasons for participating

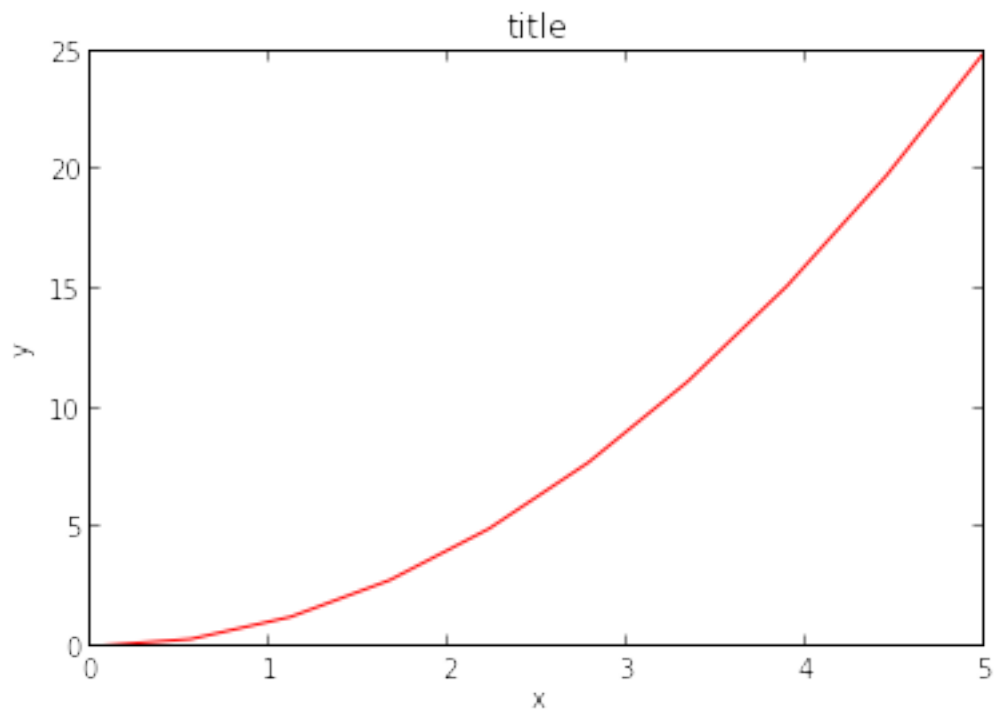
1.1.4 Datetime and Coordinates**1.1.5 Handling missing****1.1.6 Handling missing****1.1.7 Kaglle overview****1.1.8 Numeric features****1.1.9 Recap of main algorithms**

```
In [3]: %matplotlib inline
```

```
In [4]: from pylab import *
```

```
In [5]: x = linspace(0, 5, 10)
        y = x ** 2
```

```
In [6]: figure()
        plot(x, y, 'r')
        xlabel('x')
        ylabel('y')
        title('title')
        show()
```



Chapter 2

Exploratory Data Analysis

2.1 Exploratory Data Analysis

Nowadays, most ecological research is done with hypothesis testing and modelling in mind. However, Exploratory Data Analysis (EDA), which uses visualization tools and computer synthetic descriptors, is still required at the beginning of the statistical analysis of multidimensional data, in order to:

- Get an overview of the data
- Transform or recode some variables
- Orient further analyses

2.1.1 Building Intuition about the data

1. Getting domain knowledge It helps to deeper understand the problem
2. Checking if the data is intuitive To be agree with domain knowledge
3. Understanding how the data was generated A it is crucial to set up a proper Validation

2.1.2 Data cleaning

1. Constant features
2. Duplicated features
3. Duplicated rows
4. Check if datasets is shuffled

here graphics!! :D

2.1.3 Data splitting strategy

- 1.

- 2.1.4 Exploratory data Analysis**
- 2.1.5 Exploratory anonymized data**
- 2.1.6 Validation and overfitting**
- 2.1.7 Validation strategies**
- 2.1.8 Visualizations**