

# *Beginner's guide to Machine Learning competitions*

EuroPython 2015

Christine Doig

Slides

**[bit.ly/ep2015-ml-tutorial](http://bit.ly/ep2015-ml-tutorial)**

Notebooks

**[bit.ly/ep2015-ml-tutorial-repo](http://bit.ly/ep2015-ml-tutorial-repo)**

# Christine Doig

Data Scientist, Continuum Analytics



**ch\_doig**



**chdoig**

**[chdoig.github.io](https://chdoig.github.io)**

## Concepts

Data Science

Machine Learning

Supervised learning

Classification

NLP

Sentiment analysis

## Setup

Anaconda

Kaggle

Competitions

Dataset

## Process

Feature preparation

Modeling

Optimization

Validation

# Concepts

Data Science

Machine Learning

Supervised learning

Classification

NLP

Sentiment analysis

# Setup

Anaconda

Kaggle

Competitions

Dataset

# Process

Feature preparation

Modeling

Optimization

Validation

45min

10 min

1h

5 min

1h

# Concepts

Data Science

Machine Learning

Supervised learning

Classification

NLP

Sentiment analysis

# Data Science

I IS DATA SCIENTIST







Machine Learning/Stats

Scientific Computing

Analytics

Distributed Systems

*data science*

Web

Machine Learning/Stats

Data Scientists/  
Modeler

Scientific Computing

Research/Computational  
Scientist

Analytics

Data/Business  
Analyst

Distributed Systems

Data Engineers/  
Architects

*data science*

Web

Developer

Machine Learning/Stats

Data Scientists/  
Modeler

Model

Scientific Computing

Research/Computational  
Scientist

Algorithm

Analytics

Data/Business  
Analyst

Report

Distributed Systems

Data Engineers/  
Architects

Pipeline/  
Architecture

Application

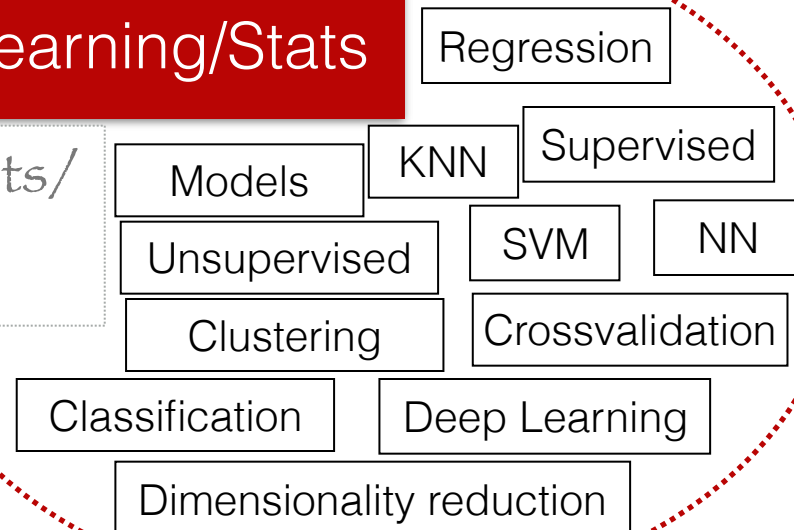
Web

Developer

*data science*

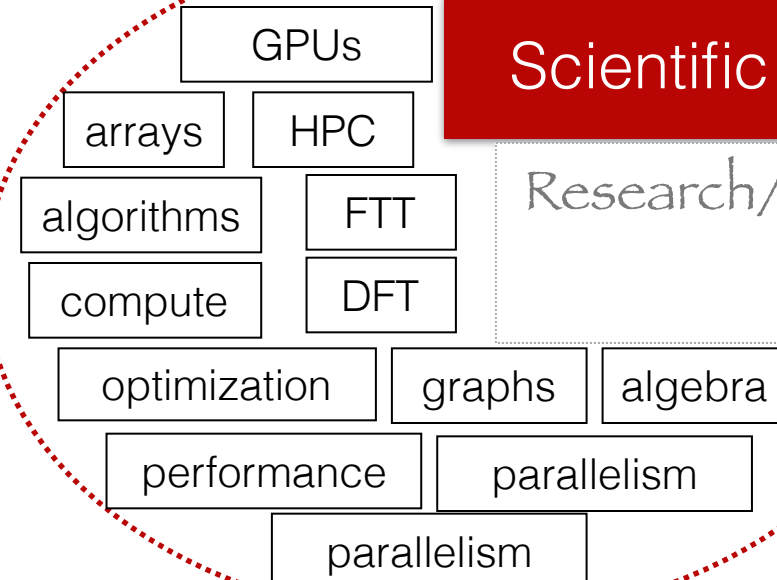
## Machine Learning/Stats

Data Scientists/  
Modeler



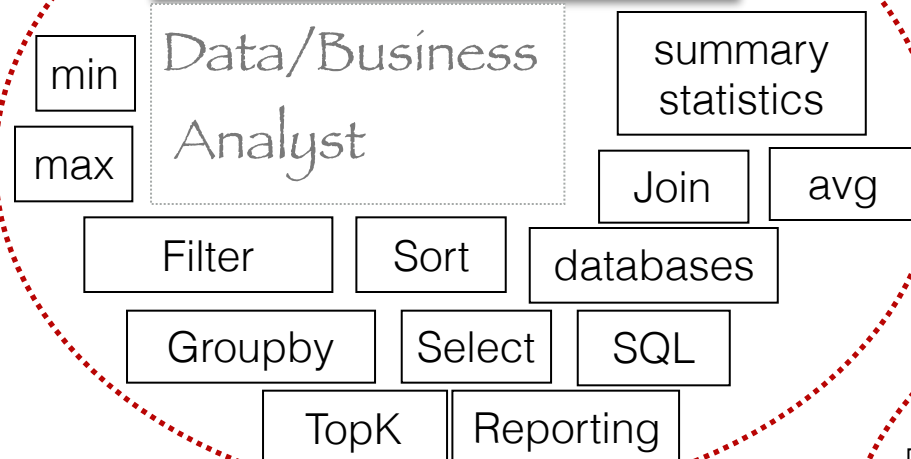
## Scientific Computing

Research/Computational  
Scientist



## Analytics

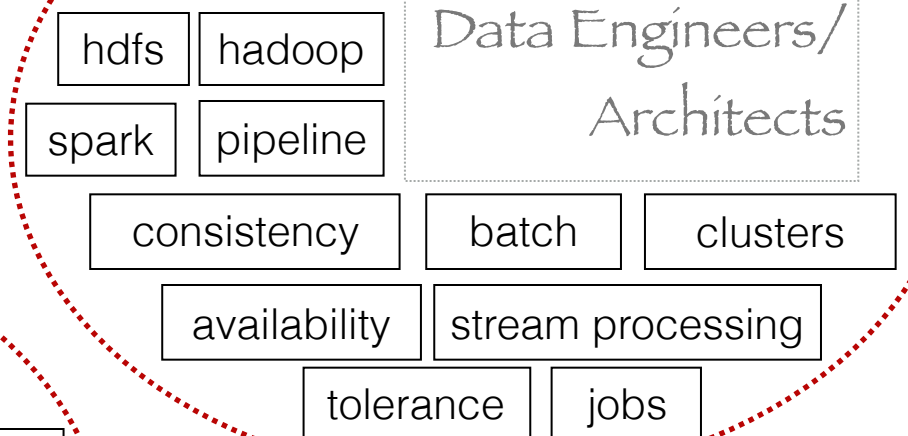
Data/Business  
Analyst



*data science*

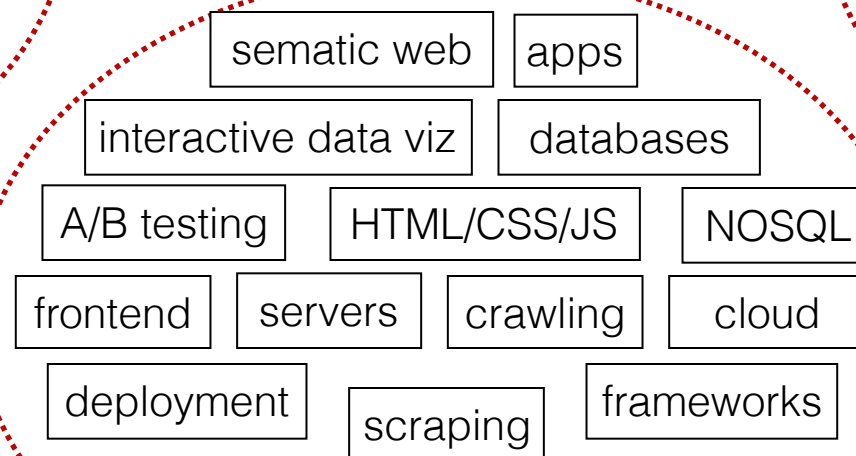
## Distributed Systems

Data Engineers/  
Architects



Web

Developer





## Machine Learning/Stats

Data Scientists/  
Modeler



## Numba

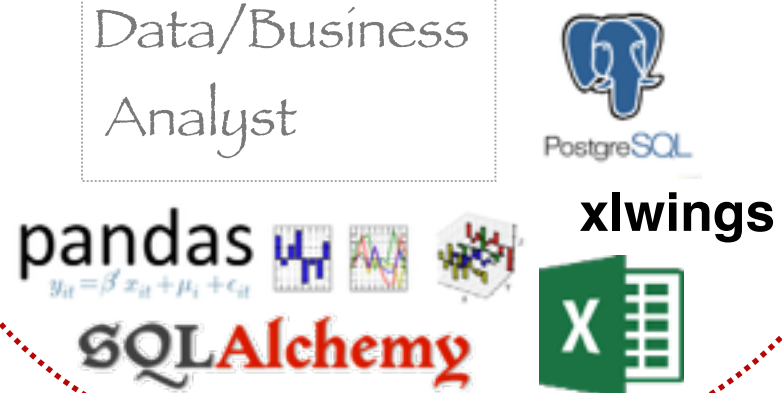
## Scientific Computing

Research/Computational  
Scientist



## Analytics

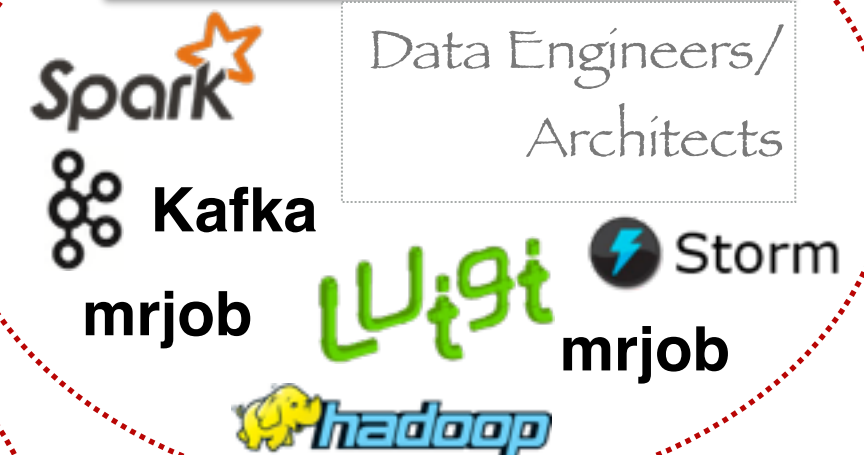
Data/Business  
Analyst



*data science*

## Distributed Systems

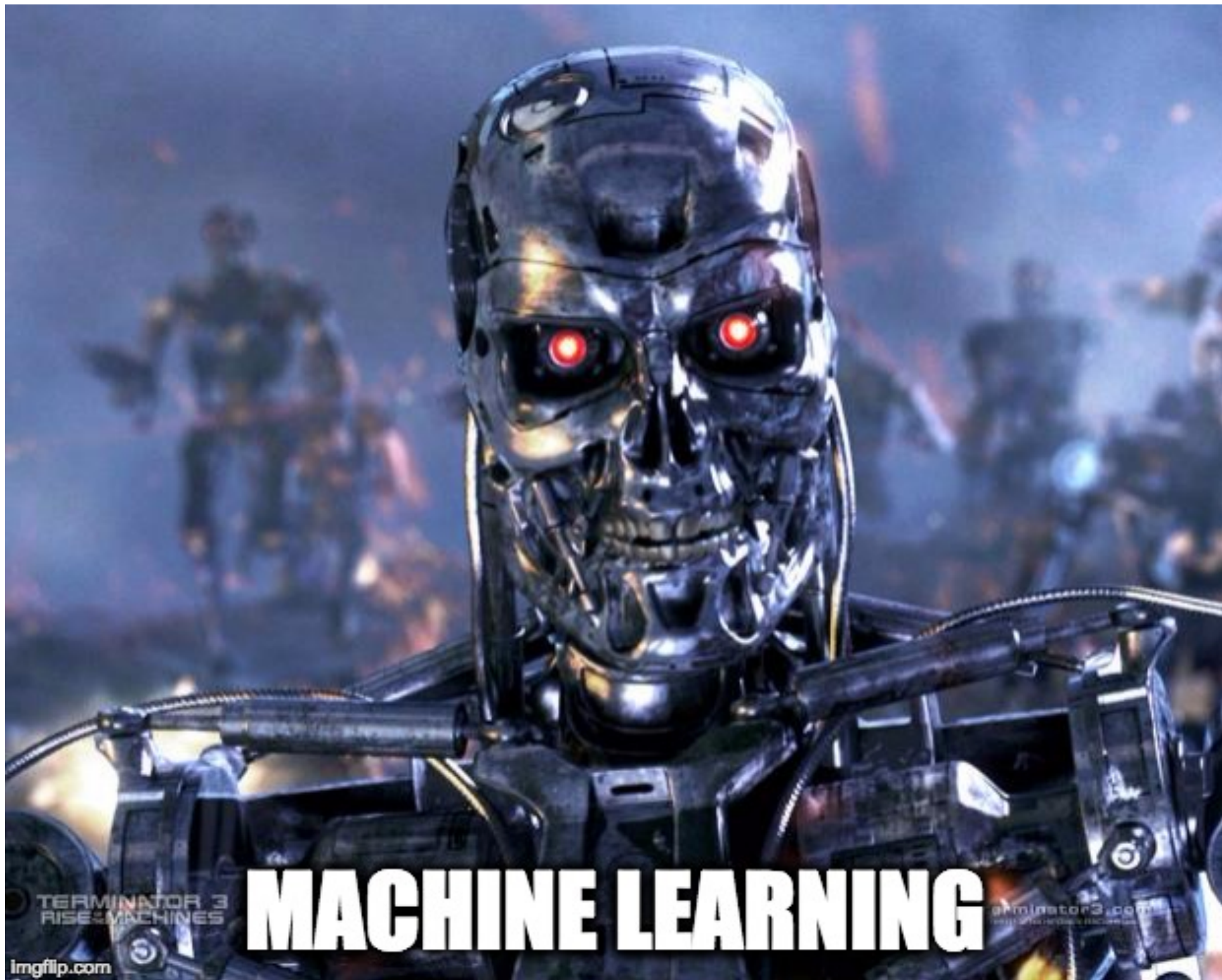
Data Engineers/  
Architects



Web

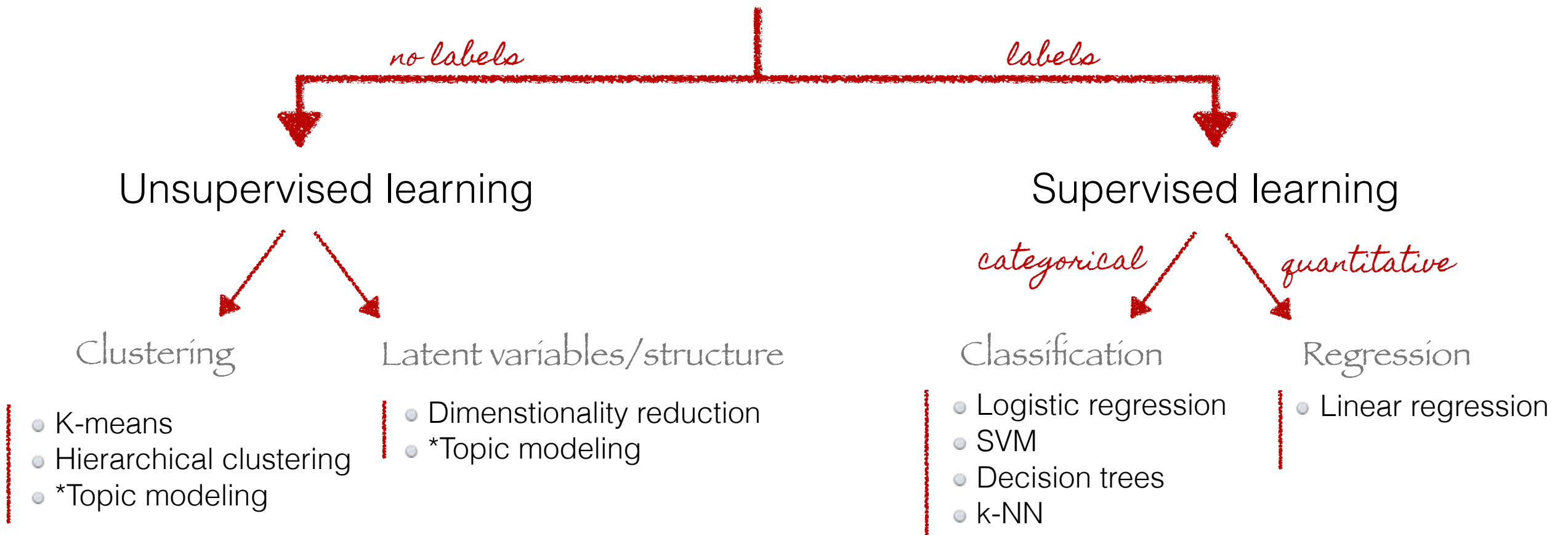
Developer

# Machine Learning

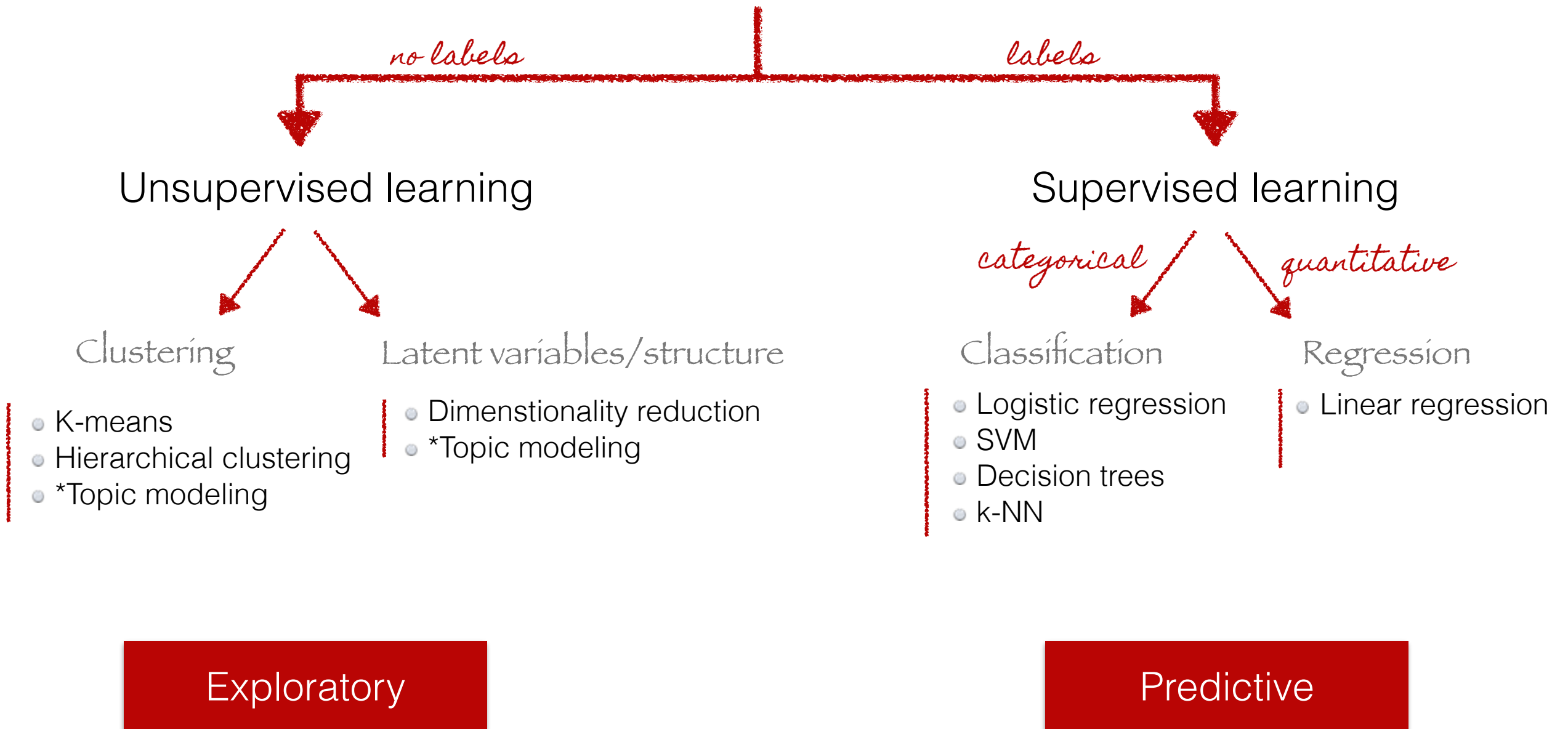




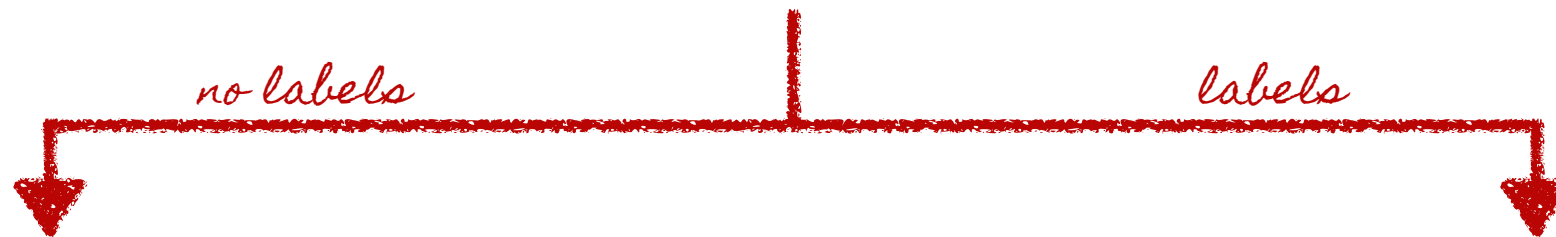
# Machine Learning



# Machine Learning



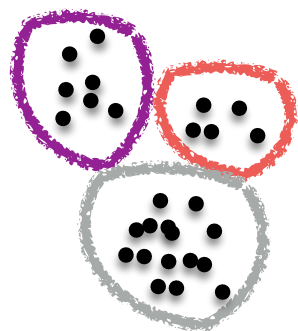
# Machine Learning



## Unsupervised learning

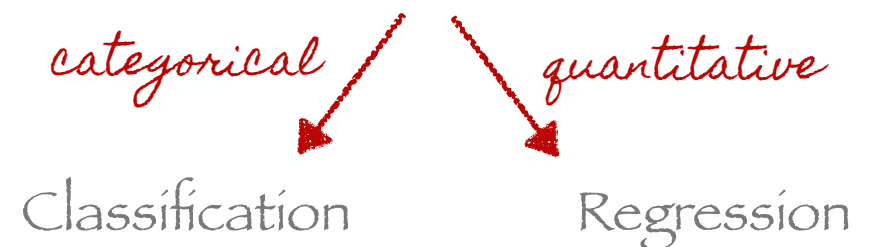
id	gender	age	job_id
1	F	67	1
2	M	32	2
3	M	45	1
4	F	18	2

Exploratory



group similar individuals together

## Supervised learning



id	gender	age	job_id	buy/click_ad	money_spent
1	F	67	1	Yes	\$1,000
2	M	32	2	No	-
3	M	45	1	No	-
4	F	18	2	Yes	\$300

Predictive

Classification

predict whether an individual is going to buy/click or not

Regression

predict how much is the individual going to spend

# Natural Language Processing

# Machine Learning

↳ Natural language processing

field concerned with the interactions between computers and human (natural) languages

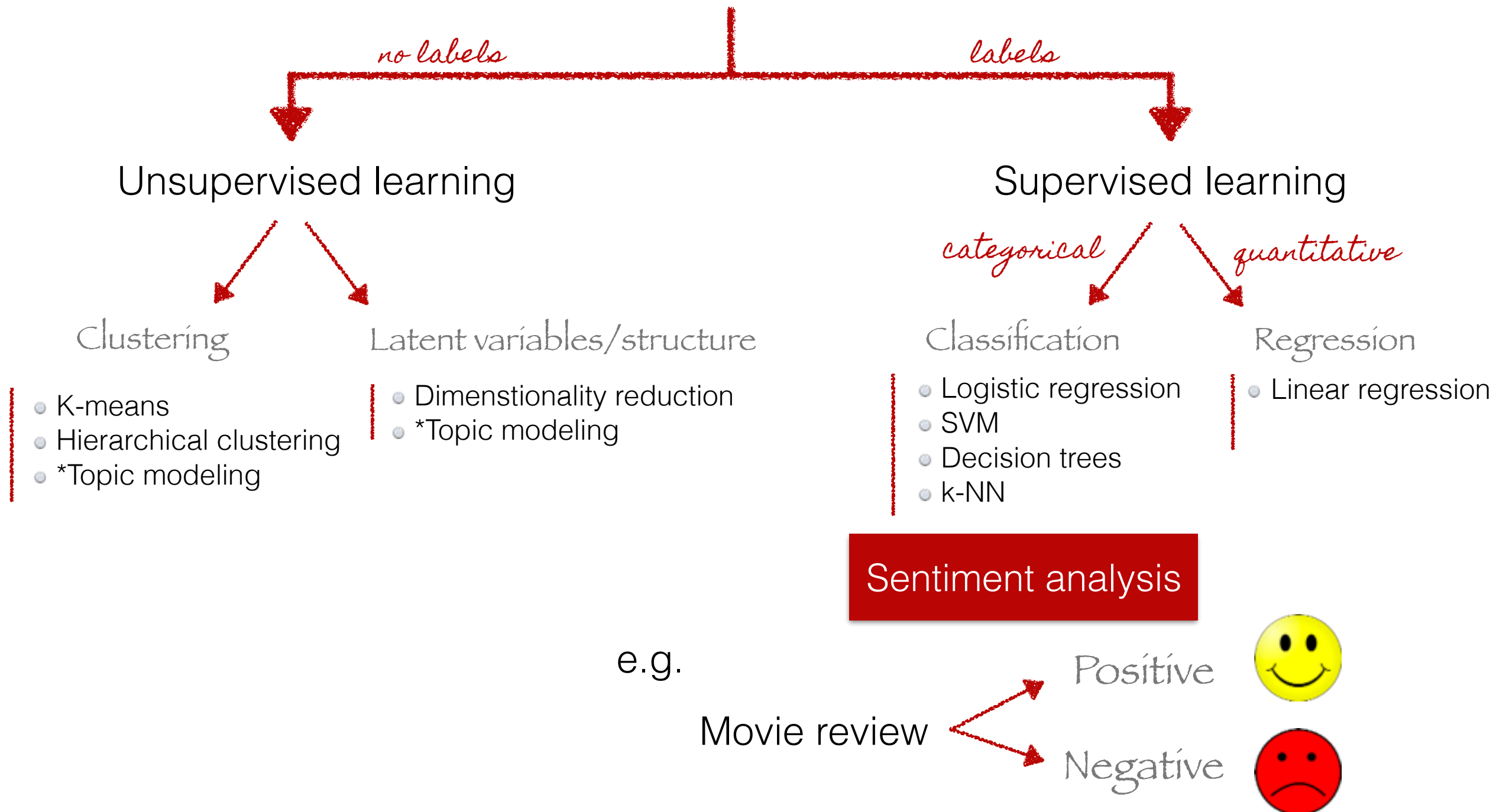
*tasks*

↳ Sentiment analysis

Extract subjective information on polarity (positive or negative) of a document (text, tweet, voice message...)

e.g online reviews to determine how people feel about a particular object or topic.

# Machine Learning





I love you! → Positive 😊

Setup



*Setup*

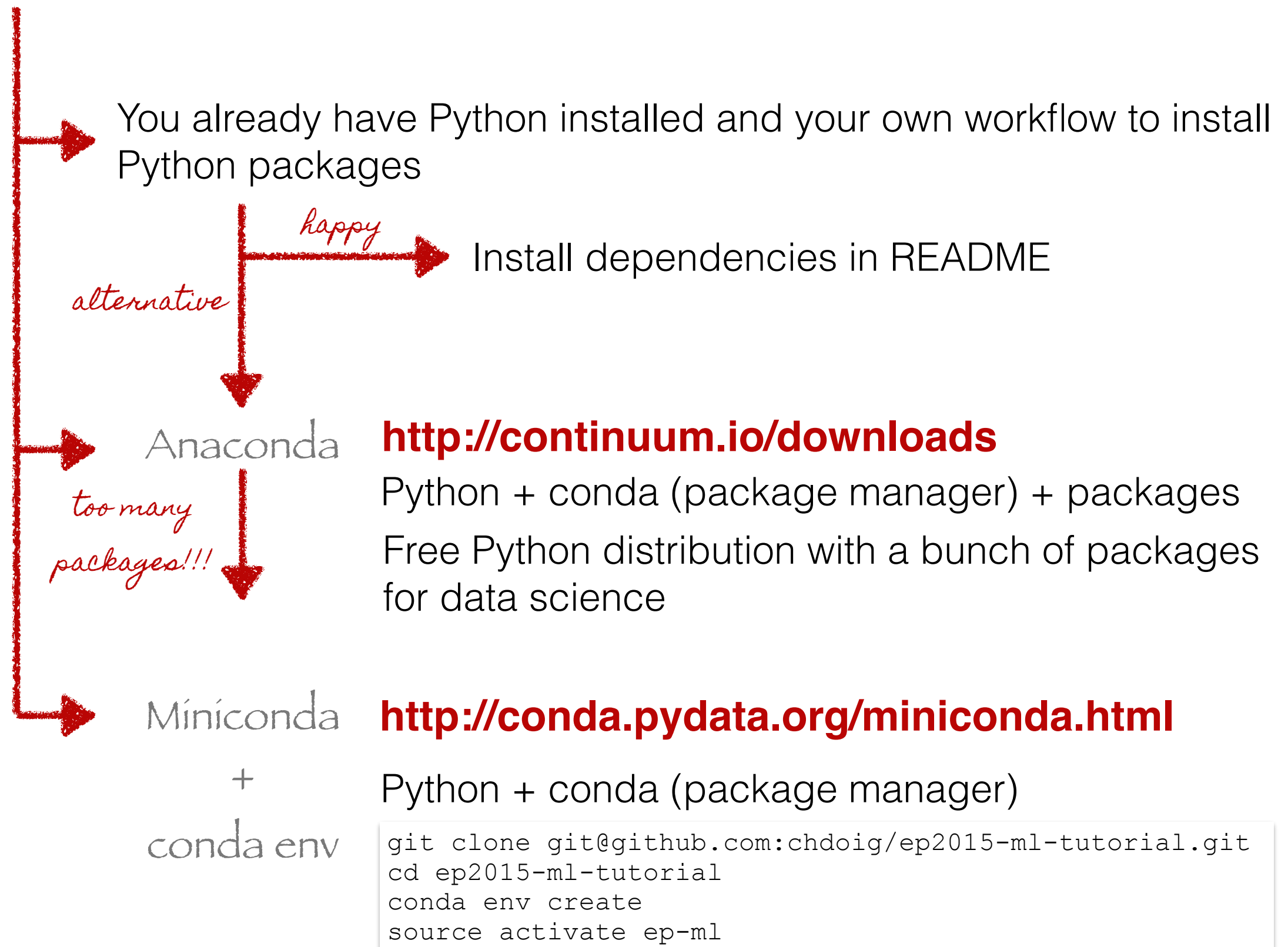
Anaconda

Kaggle

Competitions


Dataset

# Setup options



# Kaggle








<https://www.kaggle.com/>




HostCompetitionsScriptsJobsCommunity


ChristineLogout

Active Competitions

	 <div><b>Diabetic Retinopathy Detection</b> Identify signs of diabetic retinopathy in eye images</div>	8.4 days 598 teams \$100,000
	<b>Avito Context Ad Clicks</b> Predict if context ads will earn a user's click	9.4 days 342 teams \$20,000
	<b>ICDM 2015: Drawbridge Cross-Device Connections</b> Identify individual users across their digital devices	36 days 167 teams \$10,000
	<b>Liberty Mutual Group: Property Inspection Pred...</b> Quantify property hazards before time of inspection	40 days 949 teams \$25,000
	<b>Caterpillar Tube Pricing</b> Model quoted prices for industrial tube assemblies	43 days 814 teams \$30,000
	<b>Introducing Kaggle Scripts</b> Your code deserves better	2 months Swag



Christine  
[View /](#)  
[Edit Profile](#)



Is your company hiring?  
Are you on the job market?  
[Visit our jobs board >>](#)

Recent Jobs

State Farm - Research Statistician (Dunwoody, GA)  
Digital Catapult - Lead Technologist, Data (London)  
Microsoft - Software Developer/Data Wrangling Technologies (Redmond, WA)  
Imperva - Freelancer for Project (Redwood Shores, CA)  
Owlstone - Data Scientist (London or Cambridge, UK)  
Capital One - Hadoop Big Data Architect (Nottingham)

On the Forums

Regarding plotting data on map  
with latitude and longitude

hosts online machine learning competitions

# Kaggle Competition

<https://www.kaggle.com/c/word2vec-nlp-tutorial>



HostCompetitionsScriptsJobsCommunity ▾

ChristineLogout



Completed • Knowledge • 579 teams

## Bag of Words Meets Bags of Popcorn

Tue 9 Dec 2014 – Tue 30 Jun 2015 (18 days ago)

Dashboard

Home

Data

Make a submission

Information

Description

Evaluation

Rules

Part 1: For Beginners - Ba...

Part 2: Word Vectors

Part 3: More Fun With Wo...

Part 4: Comparing deep a...

Setting Up Your System

What is Deep Learning?

Forum

Leaderboard

My Team

GitHub

My Submissions

Competition Details » Get the Data » Make a submission

Data Files

File Name	Available Formats
sampleSubmission	<a href="#">.csv (276.17 kb)</a>
unlabeledTrainData.tsv	<a href="#">.zip (25.98 mb)</a>
testData.tsv	<a href="#">.zip (12.64 mb)</a>
labeledTrainData.tsv	<a href="#">.zip (12.96 mb)</a>

Data Set

The labeled data set consists of 50,000 IMDB movie reviews, specially selected for sentiment analysis. The sentiment of reviews is binary, meaning the IMDB rating < 5 results in a sentiment score of 0, and rating >=7 have a sentiment score of 1. No individual movie has more than 30 reviews. The 25,000 review labeled training set does

# Kaggle Competition

## Bag of Words Meets Bags of Popcorn

<https://www.kaggle.com/c/word2vec-nlp-tutorial>



Data

50,000 IMDB movie reviews

`labeledTrainData.tsv`

25,000 rows containing an id, sentiment, and text for each review.

`testData.tsv`

25,000 rows containing an id and text for each review

Task

predict the sentiment for each review  
in the test data set

*Process*

Feature preparation

Modeling

Optimization

Validation

Modeling

Neural Networks  
Logistic Regression

SVM

**Naive Bayes classifier**

Decision trees

Feature discretization

Feature scaling

Feature imputation

Feature preparation

Feature selection

**Feature extraction**

Hyperparameters

Optimization

Bagging

**Ensemble**

Boosting

Regularization **Random forest**

**Hold out method**

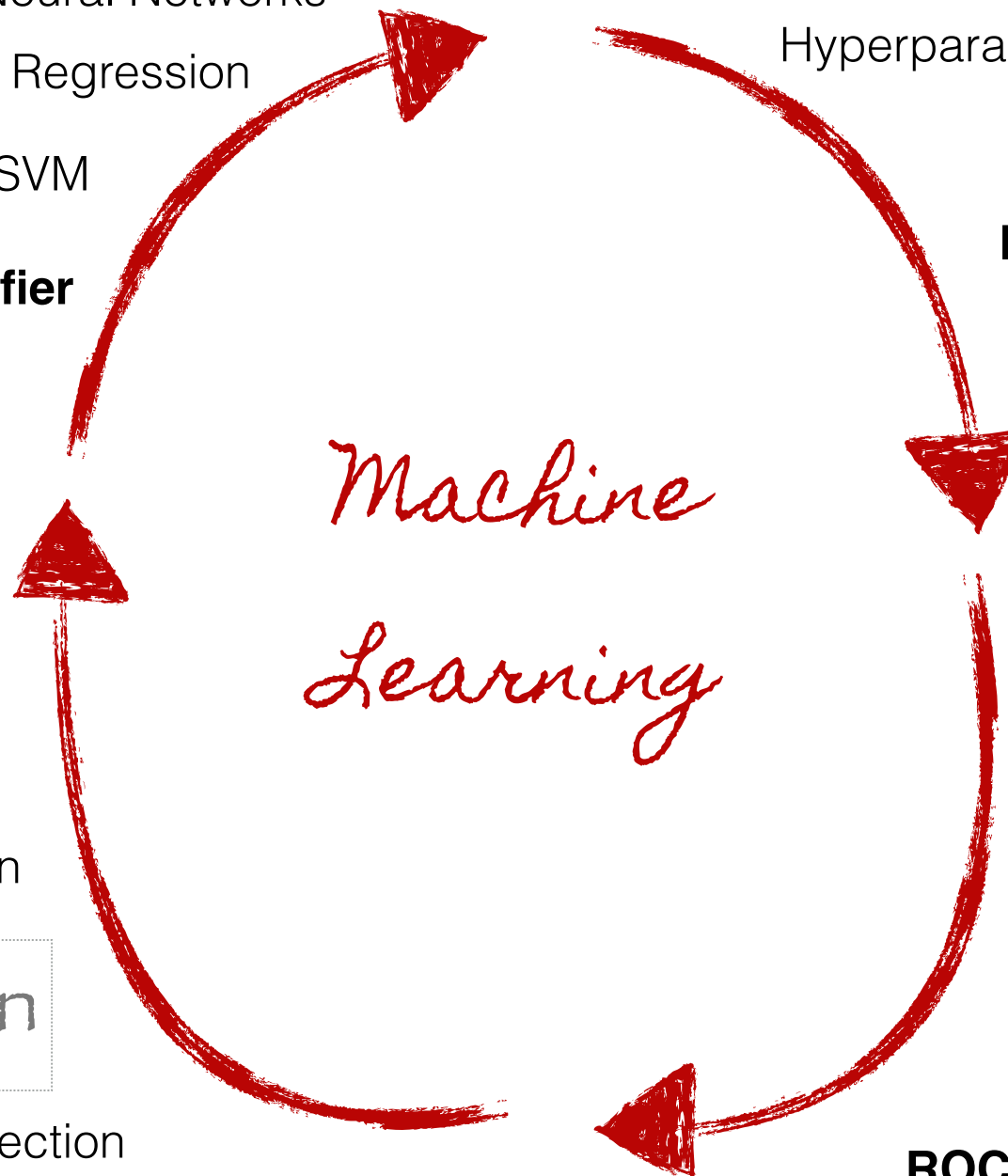
**Crossvalidation**

**Confusion matrix**

Validation

**ROC curve / AUC**

*Machine Learning*



# Feature preparation

## Feature extraction

the process of making features from available data to be used by the classification algorithms

M

Reviews



id	sentiment	review	count_words	terrible_word
1	0	the movie was terrible	4	1
2	1	I love it	3	0
3	1	Awesome! Love it!	3	0
4	0	I hated every minute	4	0



# Feature extraction Text

Tokenization

Stopwords

*Simple*

transition, metal, oxides,  
considered, generation,  
materials, field, electronics,  
advanced, catalysts, tantalum,  
v, oxide, reports, synthesis,  
material, nanometer, size,  
unusual, properties...

*Collocations*

transition\_metal\_oxides,  
considered, generation,  
materials, field, electronics,  
advanced, catalysts, tantalum,  
oxide, reports, synthesis,  
material, nanometer\_size,  
unusual, properties,  
sol\_gel\_method,  
biomedical\_applications...

*Entities*

transition, metal\_oxides,  
tantalum, oxide,  
nanometer\_size,  
unusual\_properties, dna,  
easy\_method,  
biomedical\_applications

*Combination*

transition, metal\_oxides,  
generation, tantalum, oxide,  
nanometer\_size,  
unusual\_properties, sol, dna,  
easy\_method,  
biomedical\_applications

*Lemmatization*

transition, metal, oxide,  
consider, generation, material,  
field, electronic, advance,  
catalyst, property...

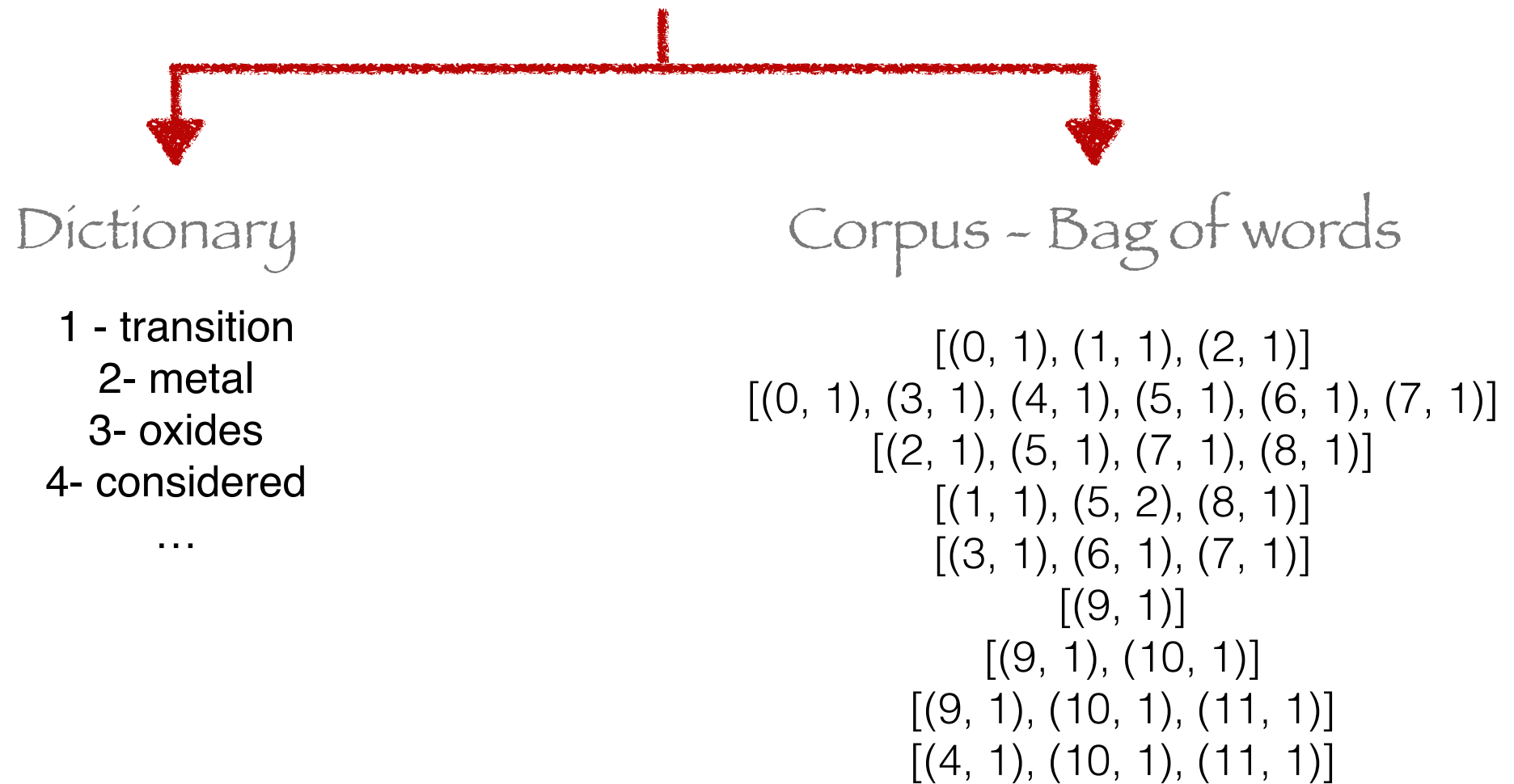
*language  
generic*

a  
above  
across  
after  
afterwards  
again  
against  
all  
...

*domain  
specific*

material  
temperature  
advance  
size  
....

# Vector Space



Feature\_extraction.ipynb

## Naive Bayes Classifier

$$P(A|B) = P(B|A) * P(A) / P(B)$$

id	sentiment	review	count_words	terrible_word
1	0	the movie was terrible	4	1
2	1	I love it	3	0
3	1	Awesome! Love it	1	0
4	0	I hated every minute	4	0

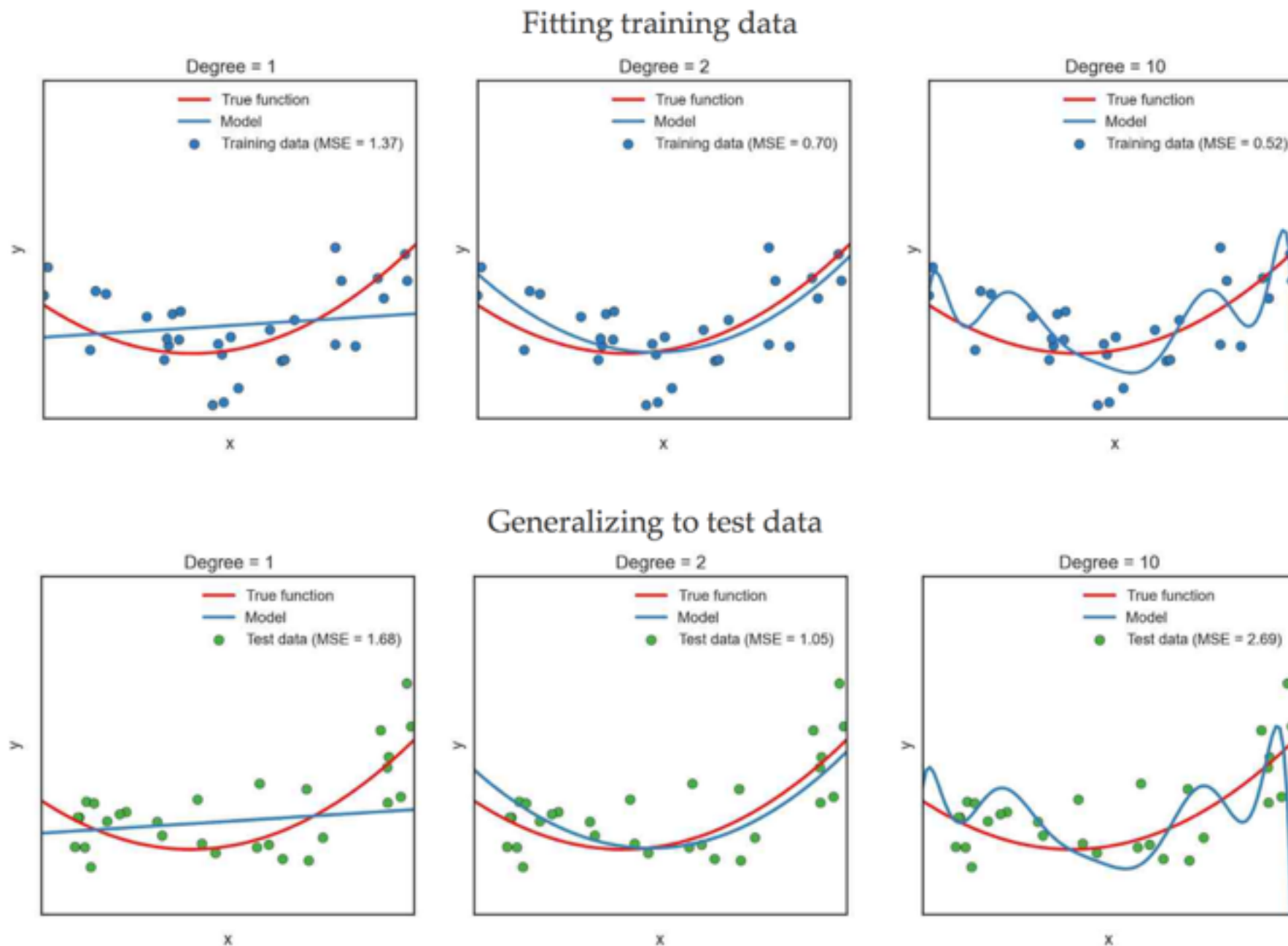
What's the probability of the review being positive if the word love appears in the review?

$$P(1 \mid \text{love}) = P(\text{love} \mid 1) * P(1) / P(\text{love}) = (2/2 * 2/4)/(2/4) = 100\%$$

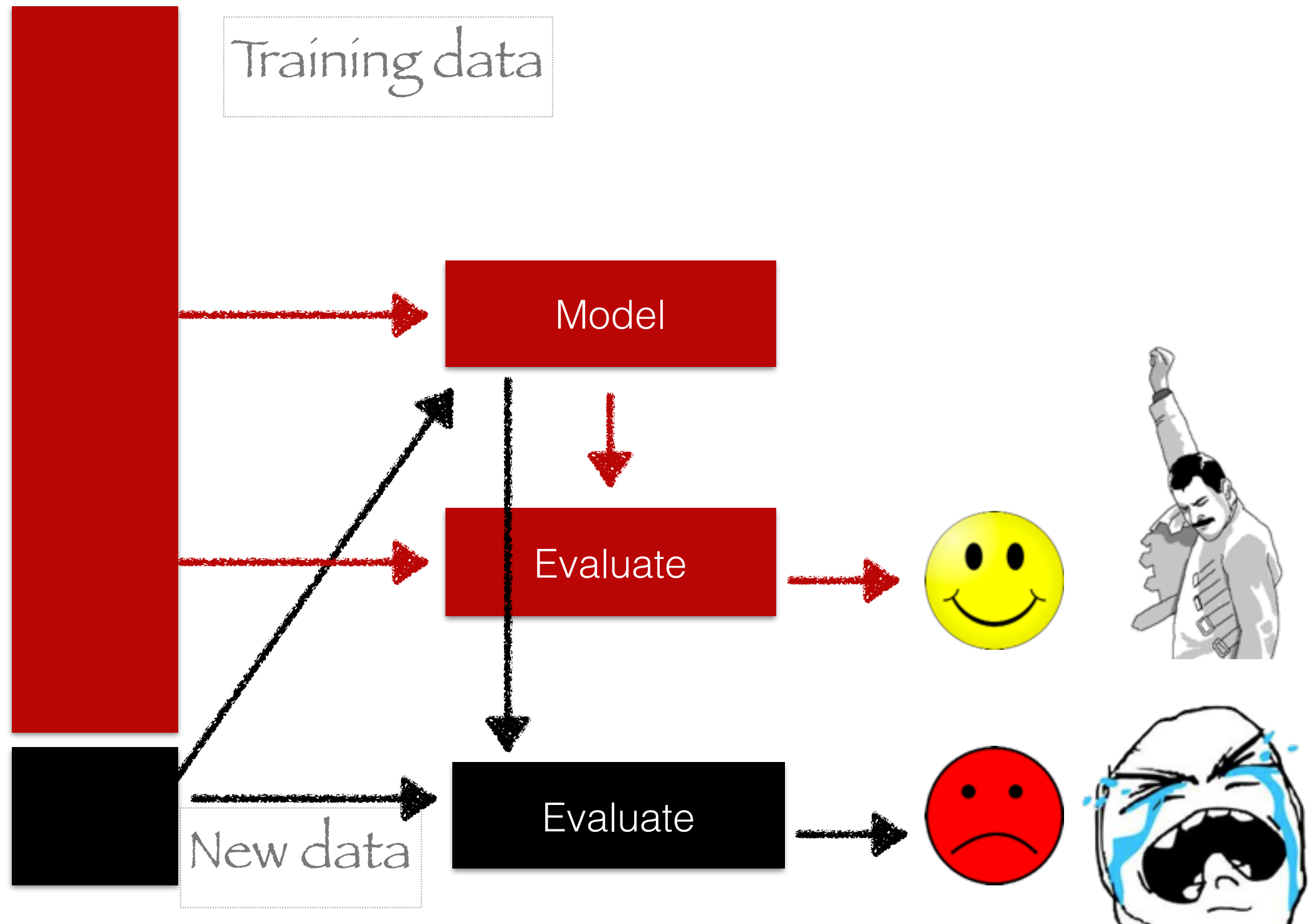
Modeling.ipynb

## Overfitting

*occurs whenever a model learns from patterns that are present in the training data but do not reflect the data-generating process. Seeing more than is actually there. A kind of data hallucination.*



# Validation



Validation.ipynb



Validation

## Hold out method



Training data

Test data

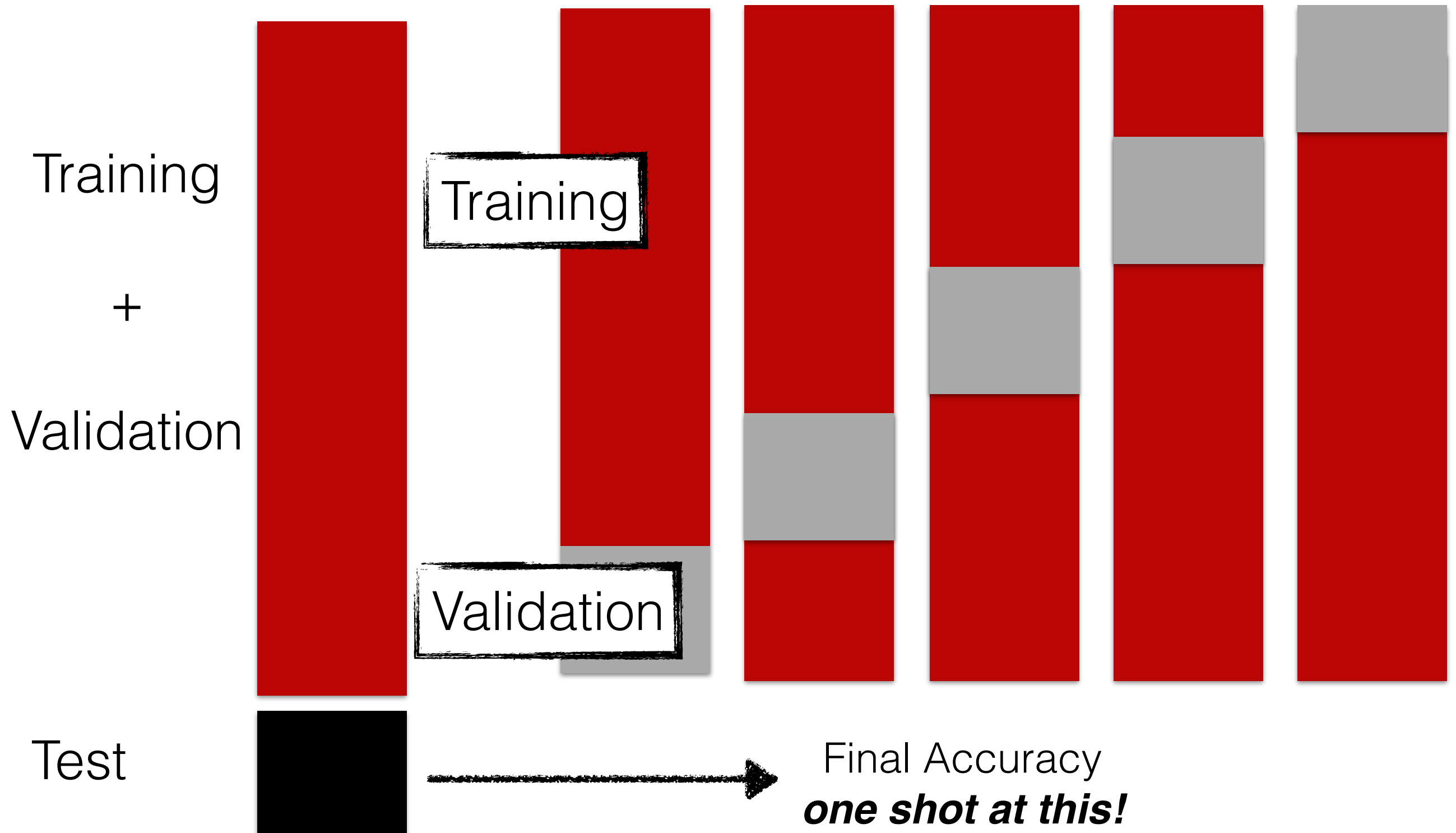


accuracy

# Crossvalidation

Accuracy = average(Round1, Round 2....)

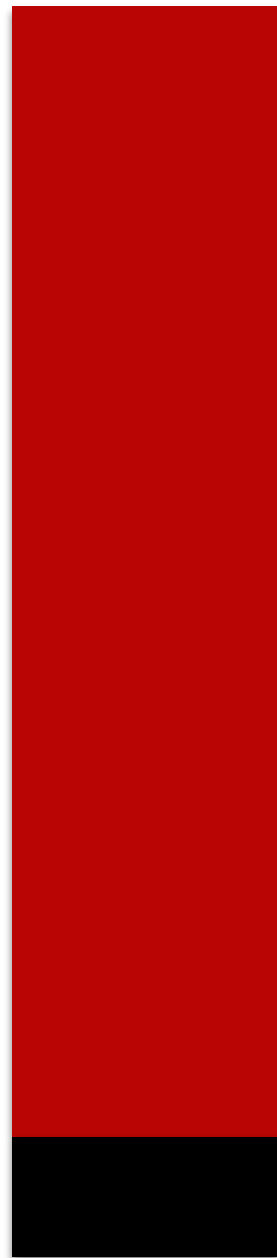
Accuracy in each round with validation set



Validation

## Confusion matrix

Real



Positive reviews  
95%

Negative reviews  
5%

Model prediction



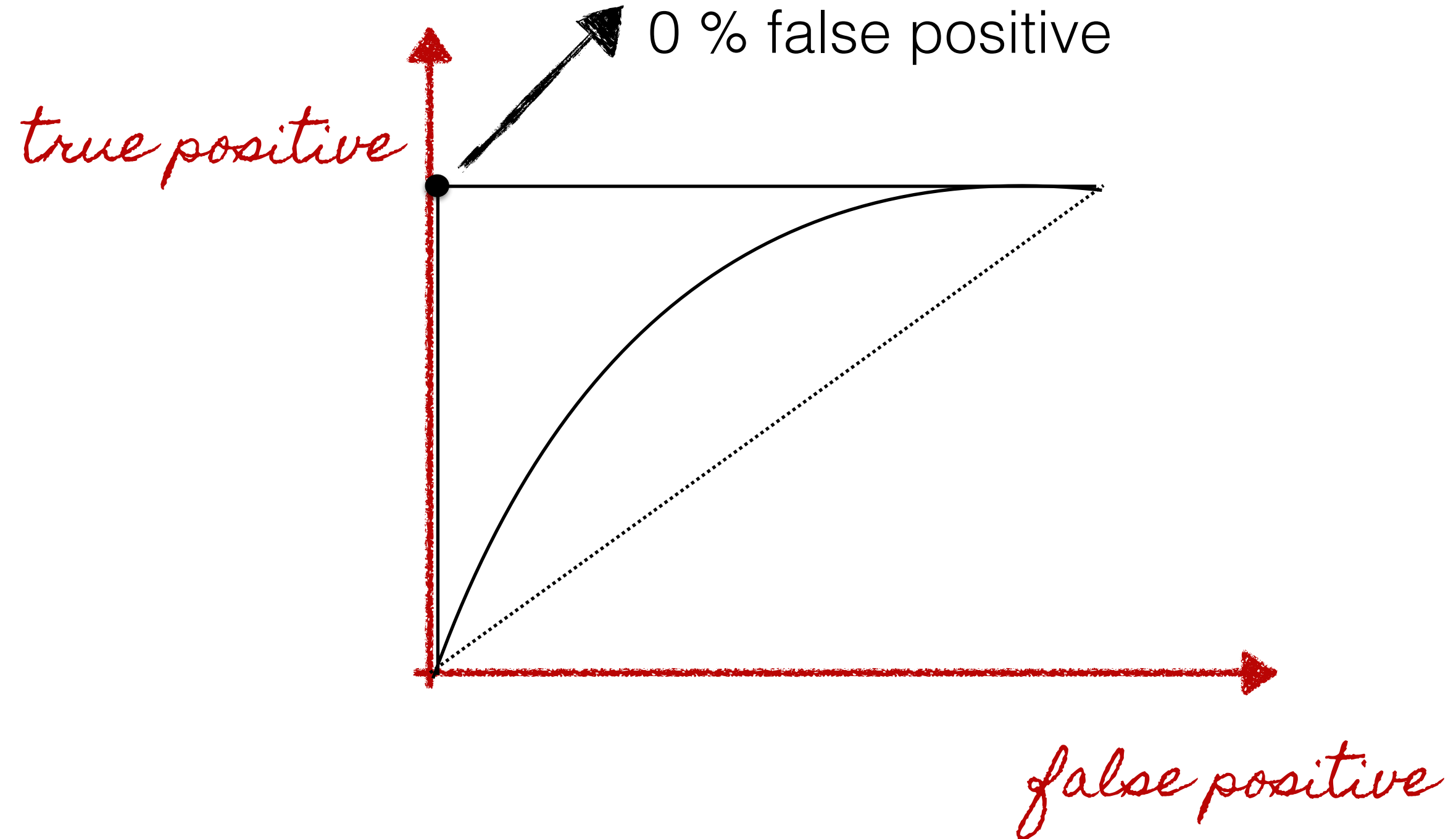
Accuracy  
95%

## Confusion matrix

model/real	positive	negative
positive	95	5
negative	0	0

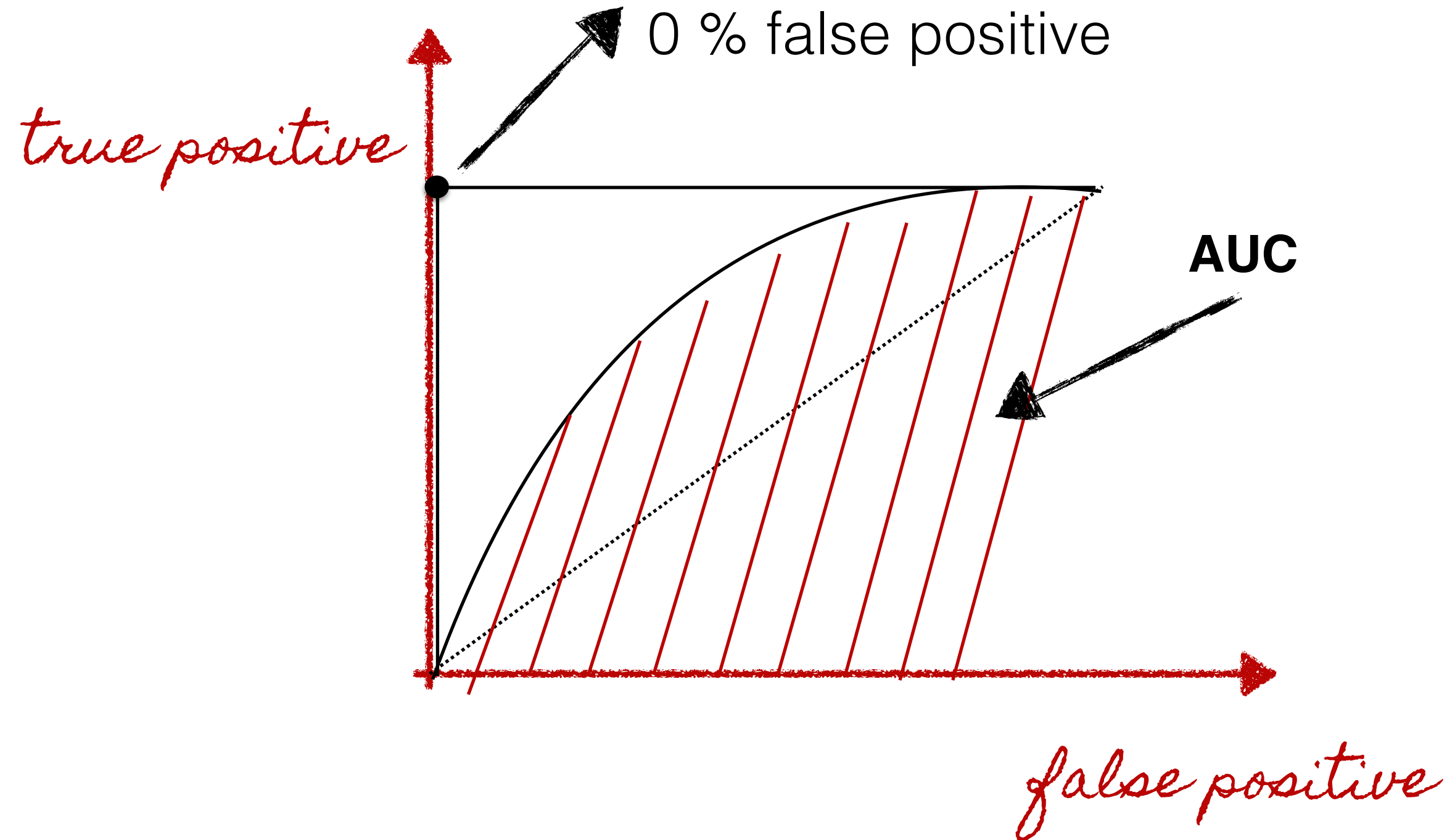
## ROC curve/ AUC

100% true positive  
0 % false positive



## ROC curve/ AUC

100% true positive  
0 % false positive



# Kaggle leaderboard

kaggle

Host

Competitions

Scripts

Jobs

Community ▾

Christine

Logout



Completed • Knowledge • 579 teams

## Bag of Words Meets Bags of Popcorn

Tue 9 Dec 2014 – Tue 30 Jun 2015 (19 days ago)


Dashboard ▾

### Public Leaderboard - Bag of Words Meets Bags of Popcorn

This leaderboard is calculated on all of the test data.

See someone using multiple accounts?

[Let us know.](#)

#	Δ1w	Team Name <small>* In the money</small>	Score <small>?</small>	Entries	Last Submission UTC (Best – Last Submission)
1	—	Zhang Yang *	1.00000	2	Mon, 08 Jun 2015 06:01:54
2	—	Puma Gigliotti	1.00000	14	Sat, 20 Jun 2015 02:34:35 (-2.1d)
3	↑214	Abner	1.00000	36	Thu, 25 Jun 2015 07:19:13
4	↑8	Cristian	0.99999	7	Tue, 30 Jun 2015 23:39:07
5	↓2	HeChen	0.99996	1	Tue, 05 May 2015 13:59:05
6	↓2	Alejandro Peláez	0.99259	25	Sun, 10 May 2015 01:10:37 (-0.7h)
7	↓2	Sebastian Raschka_	0.99156	5	Fri, 19 Jun 2015 18:13:32
8	↓2	 vgng	0.97663	22	Sun, 28 Jun 2015 23:24:02 (-7.8d)

## Ensemble methods



Classifier 1    Classifier 2    Classifier 3

id	cls_1	cls_2	cls_3	ensemble
1	0	0	0	0
2	0	1	1	1
3	1	1	1	1
4	0	0	1	0

e.g. majority voting

$w_1$

$w_2$

$w_3$

e.g. weighted voting



Ensemble.ipynb



Completed • Knowledge • 579 teams

## Bag of Words Meets Bags of Popcorn

Tue 9 Dec 2014 – Tue 30 Jun 2015 (19 days ago)

Dashboard ▾

Competition Forum

All Forums » Bag of Words Meets Bags of Popcorn

Search

« Prev  
Topic

## Post competition solutions

Next  
Topic »

Start Watching

▲  
6

Hi all,

Since the competition is over, it would be great to have a topic to discuss the solutions.

My submission has the LB score of 0.97663. It is based on the ensemble of NBSVM, Paragraph Vector and Gated Recurrent Neural Network. The code is at: <https://github.com/vinhkhuc/kaggle-sentiment-popcorn.git>

I also would like to know about other approaches.

## Concepts

Data Science

Machine Learning

Supervised learning

Classification

NLP

Sentiment analysis

## Setup

Anaconda

Kaggle

Competitions

Dataset

## Process

Feature preparation

Modeling

Optimization

Validation

Modeling

Neural Networks  
Logistic Regression

SVM

**Naive Bayes classifier**

**Decision trees**

Feature discretization

Feature scaling

Feature imputation

Feature preparation

Feature selection

**Feature extraction**

Hyperparameters

Optimization

Bagging

**Ensemble**

Boosting

Regularization **Random forest**

**Hold out method**

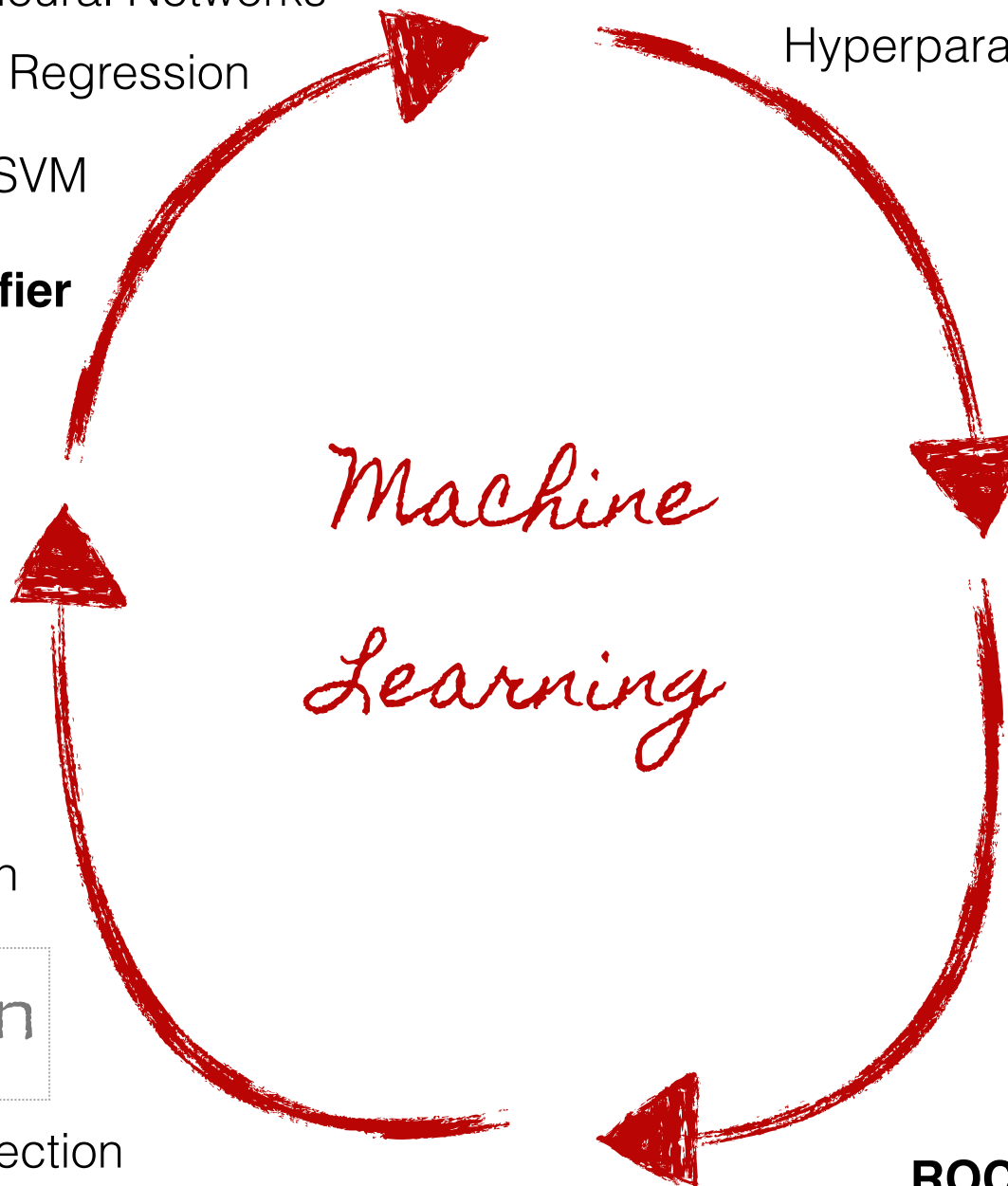
**Crossvalidation**

**Confusion matrix**

Validation

**ROC curve / AUC**

*Machine Learning*



Q&A