

# Data Workshop #1

*dataworkshop.eu*

# DataWorkshop.eu

Data Workshop

[Intro](#)

[Goal](#)

[Approach](#)

[Prerequisite](#)

[Success metric](#)

[How to join?](#)

**Talk is cheap. Show me  
the data!**

Matters is only ready-made solution with actionable insights. The rest is secondary. Practice and learn.



# About me



**Vladimir Alekseichenko**

Love analyze data



Architect



slon1024



slon1024



vova@vova.me

# Agenda

- Prepare environment
- Build & evaluate a simple model
- while True:
  - Submit result to Kaggle
  - Improve model



# Environment

[github.com/dataworkshop/prerequisite](https://github.com/dataworkshop/prerequisite)

# Packages

github.com/**dataworkshop/prerequisite**

```
$ python run.py
seaborn-0.7.0 - OK
xgboost-0.4 - OK
matplotlib-1.5.1 - OK
IPython-4.1.2 - OK
numpy-1.11.0 - OK
pandas-0.18.0 - OK
sklearn-0.17.1 - OK
```

```
=====
All right, you are ready to go on Data Workshop!
```

```
$ python run.py
seaborn-0.6 should be upgraded to seaborn-0.7
xgboost-0.4 - OK
matplotlib-1.5.1 - OK
IPython-4.1.2 - OK
numpy-1.11.0 - OK
pandas-0.18.0 - OK
sklearn-0.17.1 - OK
```

```
=====
RECOMENDATION (without upgrade some needed features could be missing)
pip install --upgrade seaborn
```

```
$ python run.py
seaborn-0.7.0 - OK
xgboost - missing
matplotlib-1.5.1 - OK
IPython-4.1.2 - OK
numpy-1.11.0 - OK
pandas-0.18.0 - OK
sklearn-0.17.1 - OK
```

```
=====
REQUIRED
Please install those packages before Data Workshop: xgboost
pip install xgboost
More info how to install xgboost: http://xgboost.readthedocs.org/en/latest/build.html
```

# jupyter notebook



```
$ jupyter notebook
[I 22:17:17.650 NotebookApp] The port 8888 is already in use, trying another random port.
[I 22:17:17.650 NotebookApp] The port 8889 is already in use, trying another random port.
[I 22:17:17.651 NotebookApp] The port 8890 is already in use, trying another random port.
[I 22:17:17.651 NotebookApp] The port 8891 is already in use, trying another random port.
[I 22:17:17.657 NotebookApp] Serving notebooks from local directory: /Users/vova/src/github/dataworkshop/titanic/vladimir/tmp
[I 22:17:17.657 NotebookApp] 0 active kernels
[I 22:17:17.657 NotebookApp] The IPython Notebook is running at: http://localhost:8892/
[I 22:17:17.657 NotebookApp] Use Control-C to stop this server and shut down all kernels (twice to skip confirmation).
```



 jupyter

Files Running Clusters

Select items to perform actions on them.

Notebook list empty.

Upload New

Text File  
Folder  
Terminal

Notebooks  
Haskell  
Julia 0.3.8  
Python 2

2

# API review

sklearn

- `.fit(X_train, y_train)`
- `.transform(X)`
- `.predict(X_test)`



# Build a simple model

[github.com/dataworkshop/titanic](https://github.com/dataworkshop/titanic)

# Titanic

**15 April 1912**  
*104 years ago*



# Titanic - what happened?



- **PassengerId** -- A numerical id assigned to each passenger.
- **Survived** -- Whether the passenger survived (1), or didn't (0). Target variable.
- **Pclass** -- The class the passenger was in -- first, second and third.
- **Name** -- the name of the passenger.
- **Sex** -- The gender of the passenger -- male or female.
- **Age** -- The age of the passenger. Fractional.
- **SibSp** -- The number of siblings and spouses the passenger had on board.
- **Parch** -- The number of parents and children the passenger had on board.
- **Ticket** -- The ticket number of the passenger.
- **Fare** -- How much the passenger paid for the ticket.
- **Cabin** -- Which cabin the passenger was in.
- **Embarked** -- Where the passenger boarded the Titanic.

# Build more advanced model

[github.com/dataworkshop/titanic](https://github.com/dataworkshop/titanic)

# How to handle with...

- Missing values
- Categorical variables
- Pipeline
- Hyperparameter optimisation



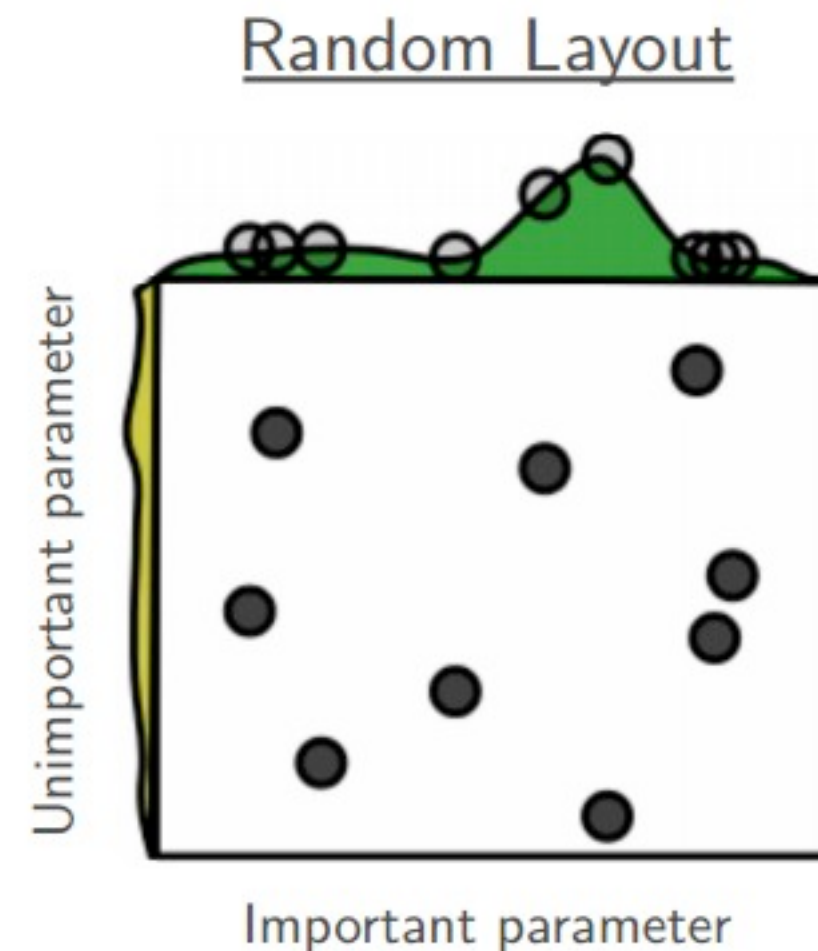
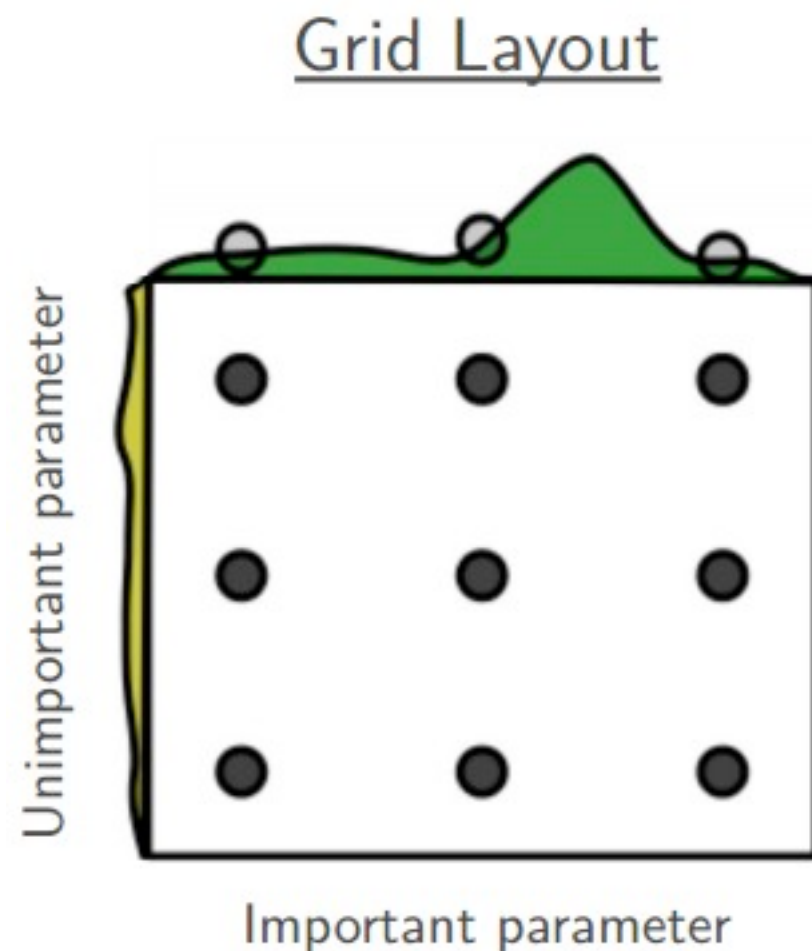
# Missing values

- Remove rows () or columns (features)
- Replace by value (-1, 0 and so on)
  - `pandas.fillna(value)`
- Use mean, median, most frequent
  - `sklearn.preprocessing.Imputer`

# Categorical Variables

- Convert into unique ID (integer)
  - “a”, “b”, “c” => 1, 2, 3  
sklearn.preprocessing.**LabelEncoder**
- One hot encoding
  - “a”, “b”, “c” => [1, 0, 0 ], [0, 1, 0], [0, 0, 1]  
pandas.**get\_dummies**
- Probability
  - ( #samples == {“a”, “b”, “c”} ) / #samples ALL

# Hyperparameter optimization



# Summary

Best results...  
but this **not** about ML :)

#	Δ1w	Team Name <small>* in the money</small>	Score <small>?</small>	Entries	Last Submission UTC (Best – Last Submission)
1	—	Prabhudatta Das *	<a href="#">1.00000</a>	1	<a href="#">Wed, 24 Feb 2016 15:14:43</a>
2	—	tryn2win	<a href="#">1.00000</a>	16	<a href="#">Wed, 02 Mar 2016 16:15:52 (-47.1h)</a>
3	—	norwayT	<a href="#">1.00000</a>	1	<a href="#">Thu, 03 Mar 2016 02:24:21</a>
4	—	determinedEurasians <small>👤</small>	<a href="#">1.00000</a>	1	<a href="#">Sat, 05 Mar 2016 15:25:43</a>
5	—	Formandos Puc Barreiro	<a href="#">1.00000</a>	2	<a href="#">Sun, 06 Mar 2016 17:13:36 (-22.7h)</a>
6	—	Nijida	<a href="#">1.00000</a>	4	<a href="#">Tue, 22 Mar 2016 08:43:34 (-12.9d)</a>
7	—	Eric_Chen	<a href="#">1.00000</a>	1	<a href="#">Thu, 17 Mar 2016 06:46:43</a>
8	—	Sajid Umair 2	<a href="#">1.00000</a>	1	<a href="#">Sun, 27 Mar 2016 13:10:32</a>
9	—	<small>👤</small> WooJungKim	<a href="#">1.00000</a>	3	<a href="#">Sat, 02 Apr 2016 16:58:46</a>
10	—	DaiHeping	<a href="#">1.00000</a>	13	<a href="#">Wed, 06 Apr 2016 09:52:51</a>
11	↑4	Prerit Ahuja	<a href="#">1.00000</a>	2	<a href="#">Mon, 18 Apr 2016 12:48:26</a>

# Three things

*if you can remember only three...*

- **Feature engineering** is important
- Manage **categorical variables** and **missing values** is important as well :)
- **Ensemble** is important



What Next?

## Active Competitions



### State Farm Distracted Driver Detection

Can computer vision spot distracted drivers?

**3 months**  
**410** teams  
**235** scripts  
**\$65,000**



### Santander Customer Satisfaction

Which customers are happy customers?

**9.8 days**  
**4662** teams  
**3129** scripts  
**\$60,000**



### Home Depot Product Search Relevance

Predict the relevance of search results on homedepot.com

**2.8 days**  
**2147** teams  
**1648** scripts  
**\$40,000**



### Expedia Hotel Recommendations

Which hotel type will an Expedia customer book?

**48 days**  
**313** teams  
**371** scripts  
**\$25,000**

Thank you