

# Data Workshop #3

## Evaluating ML Models

[dataworkshop.eu](http://dataworkshop.eu)

Alekseichenko Vladimir

# 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

# Disclaimer

**Data Workshop** [*all time*] focuses on the **intuition** and **practical** tips.

*For a formal treatment, see something else<sup>\*</sup>.*

<sup>\*</sup> papers or classical machine learning books

# Environment

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

[github.com/dataworkshop/model\\_evaluation](https://github.com/dataworkshop/model_evaluation)

# 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).
```



Files    Running    Clusters

Select items to perform actions on them.



Notebook list empty.



Text File

Folder

Terminal

Notebooks

Haskell

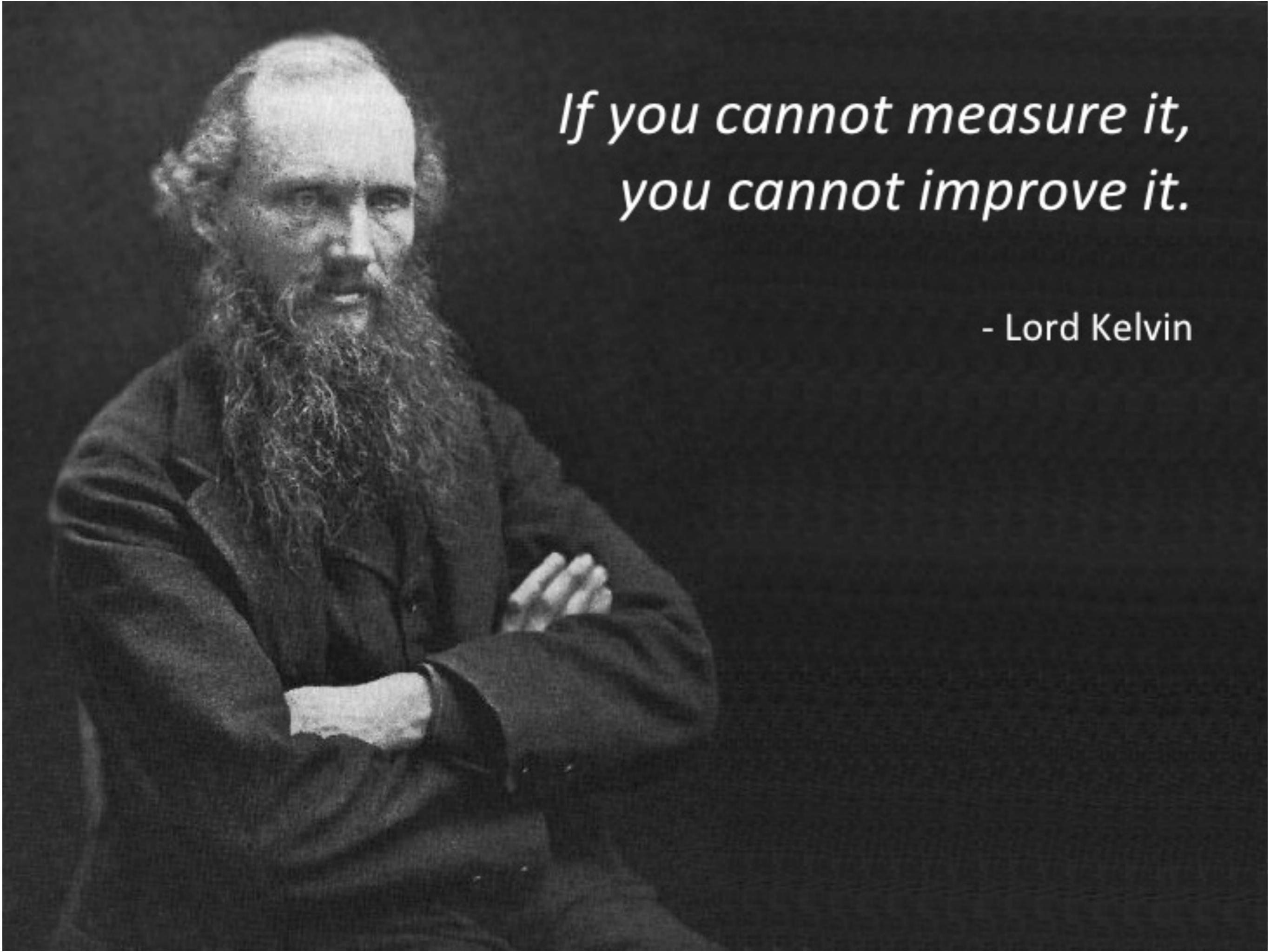
Julia 0.3.8

Python 2



# Motivation

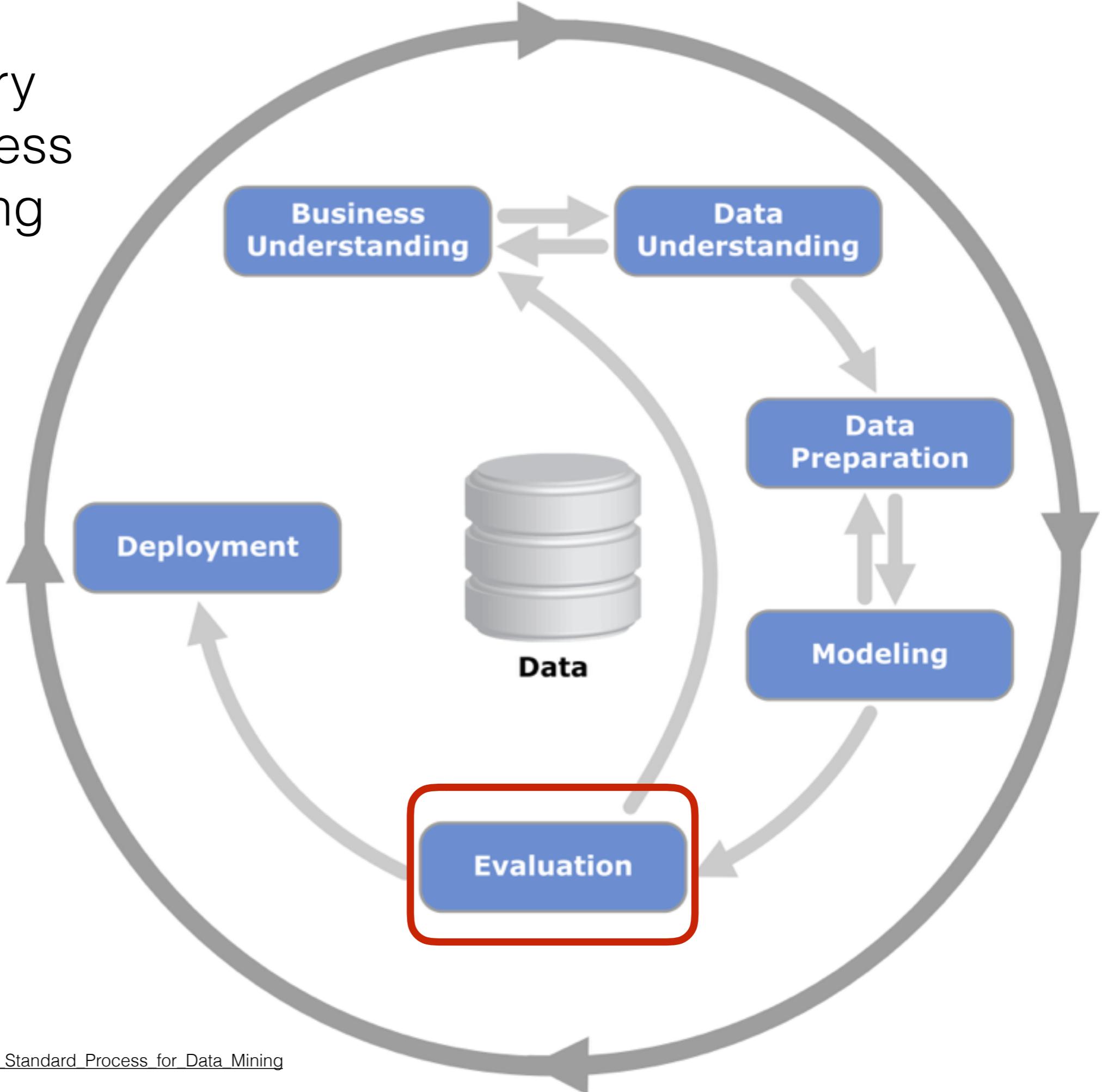
metrics & validation



*If you cannot measure it,  
you cannot improve it.*

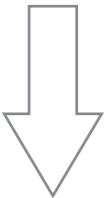
- Lord Kelvin

# Cross Industry Standard Process for Data Mining



## Understand Business & Data

Read and explore data

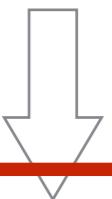


Chanel_ID	Client_ID	Product_ID	Demand
1	3	2	10
1	3	5	15
3	3	6	12

Evaluation is needed

## Feature Engineering

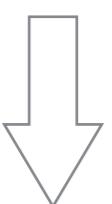
Create a new ones based on already exists



Chanel_ID	Client_ID	Product_ID	...	Demand	DemandLog
1	3	2	...	10	2.303
1	3	5	...	15	2.708
3	3	6	...	12	2.485

## Feature Selection

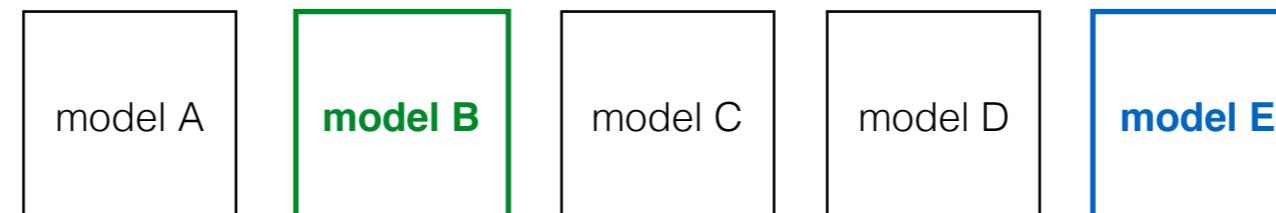
Select only useful features



Client_ID	Product_ID	...	DemandLog
3	2	...	2.303
3	5	...	2.708
3	6	...	2.485

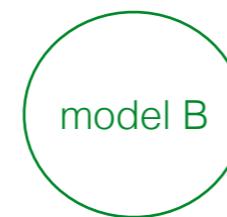
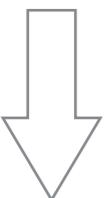
## Model Selection

Find the best model(s)

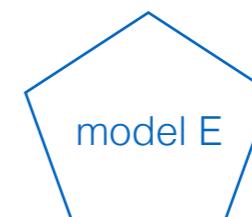


## Tuning Hyperparameters

Find the best hyperparameters for given model



model B **x0.6**



+ model E **x0.4**

## Ensemble Modeling

Combine few models into one more better

# Simple questions

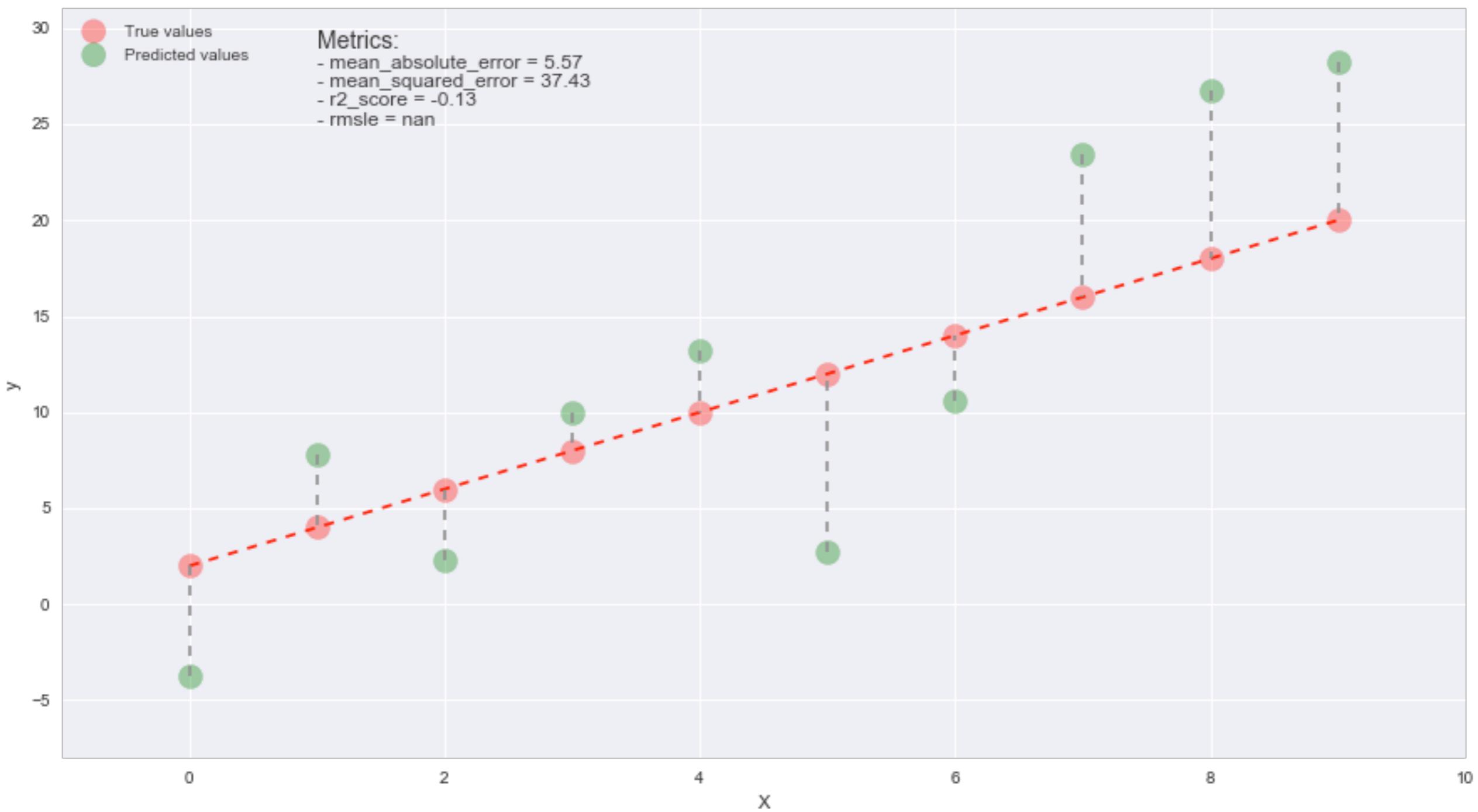
- How can i measure success for this project?
- How would i know when i've succeeded?
- Is it possible succeed for this project?

# Regression

mse, rmse, rmsle, r2

`sklearn.metrics.*`

# Regression



# Simple points about metrics

- Range
- Outliers
- Negative & Positive

# Classification

accuracy, confusion-matrix, f1, log-loss, auc

```
sklearn.metrics.*
```

# Classification

**Predicting class labels given input data**

- Binary - only two labels (e.g spam or ham)
- Multi-class - more than two possible classes (e.g. good, neutral, bad)

# Classification - binary



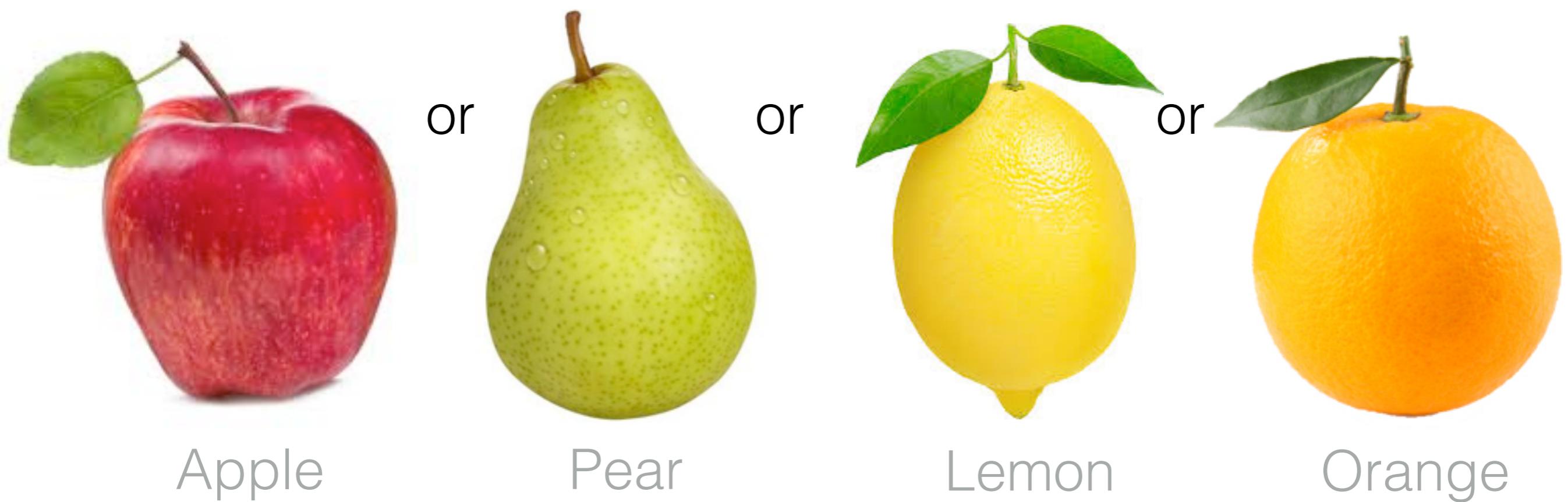
Ferrari

or



Lamborghini

# Classification - multi-class



# Accuracy

$$\frac{2 \text{ smiley faces}}{5 \text{ total faces}} = 40\%$$

#	actual	predicted	status
1			
2			
3			
4			
5			

# Accuracy Paradox

		prediction	
		yes	no
true	yes	1	5
	no	100	10000

		prediction	
		yes	no
true	yes	0	6
	no	0	10100

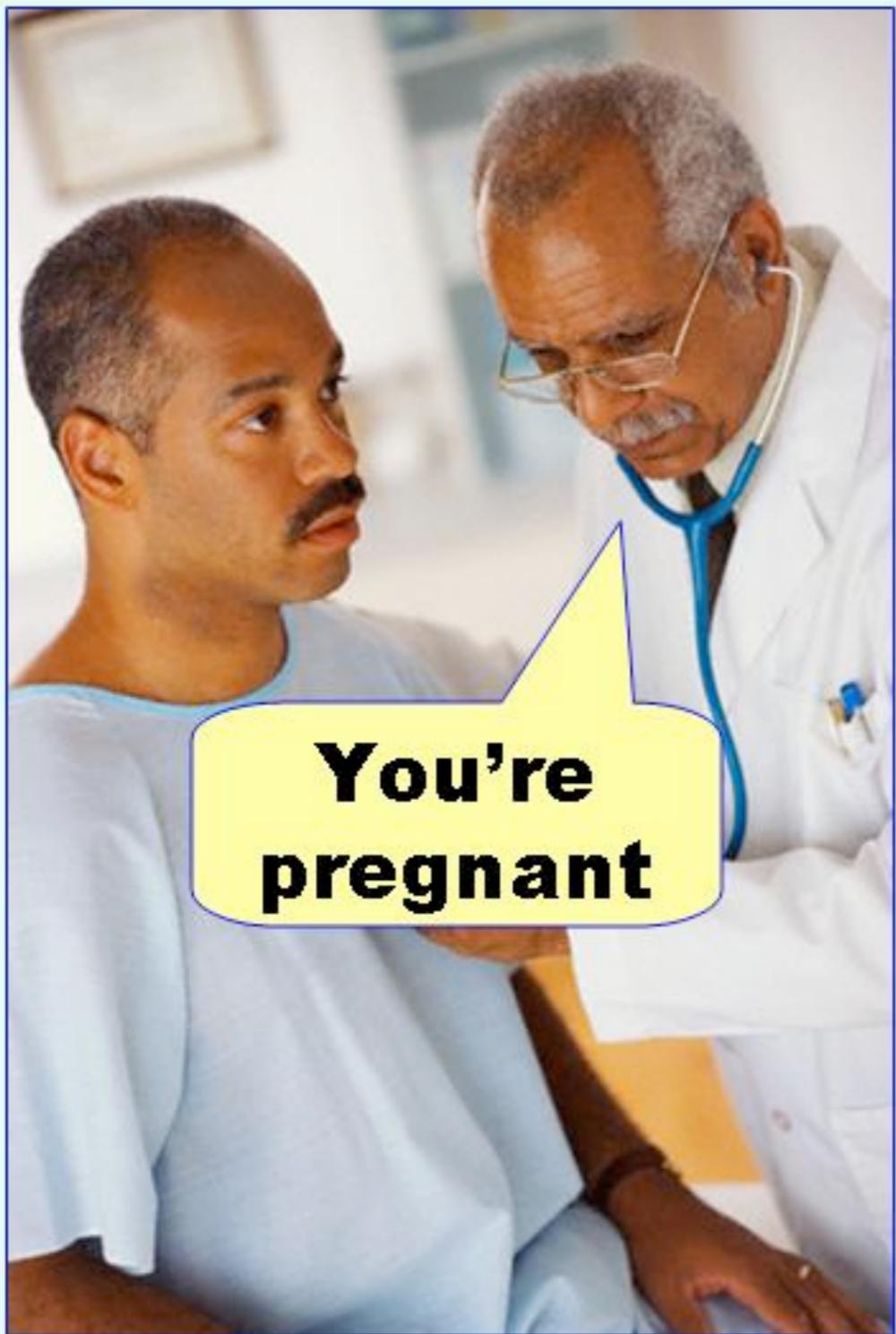
accuracy

$$\frac{1 + 10000}{1 + 5 + 100 + 10000} = 98.96\%$$

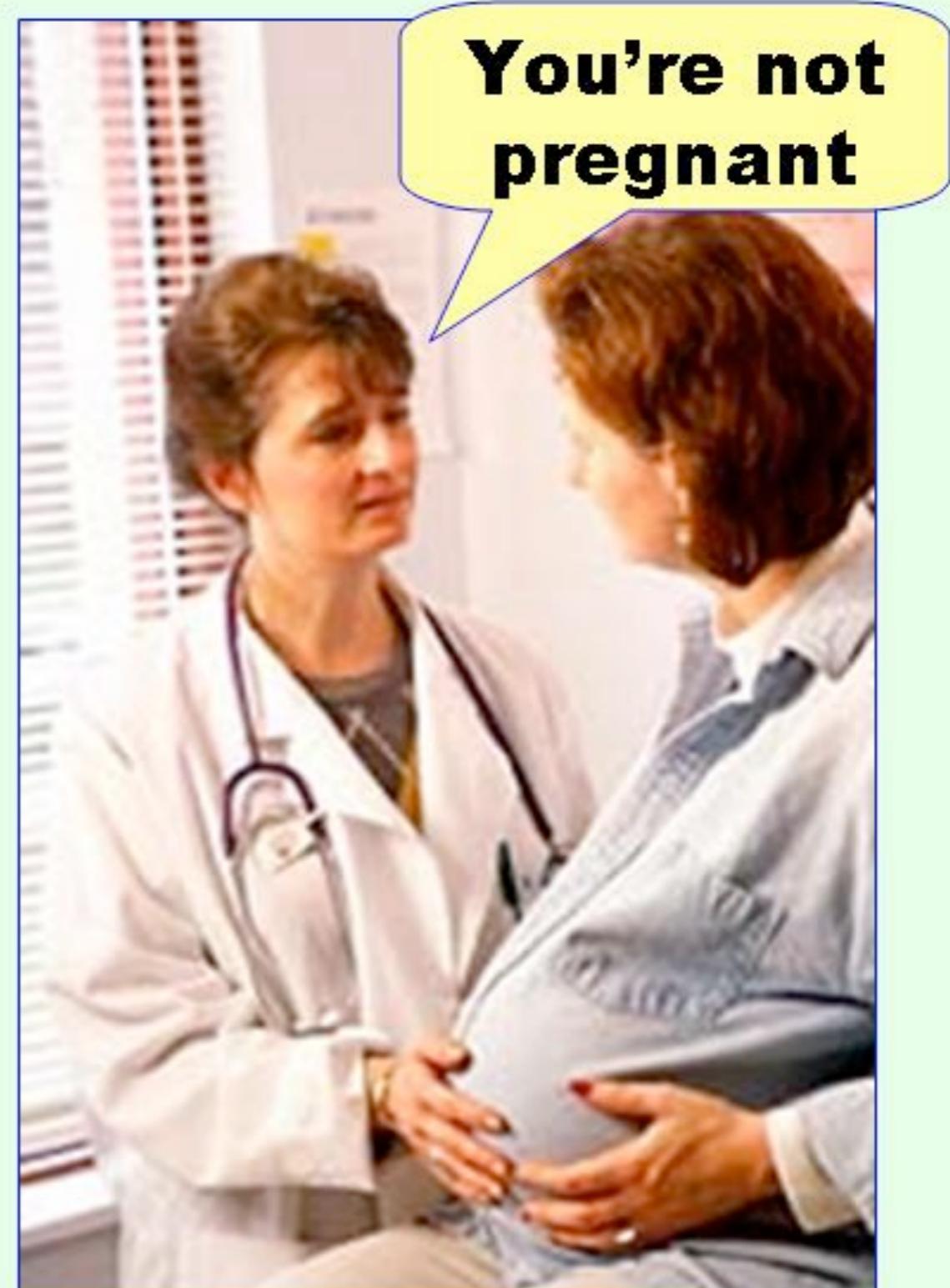
accuracy

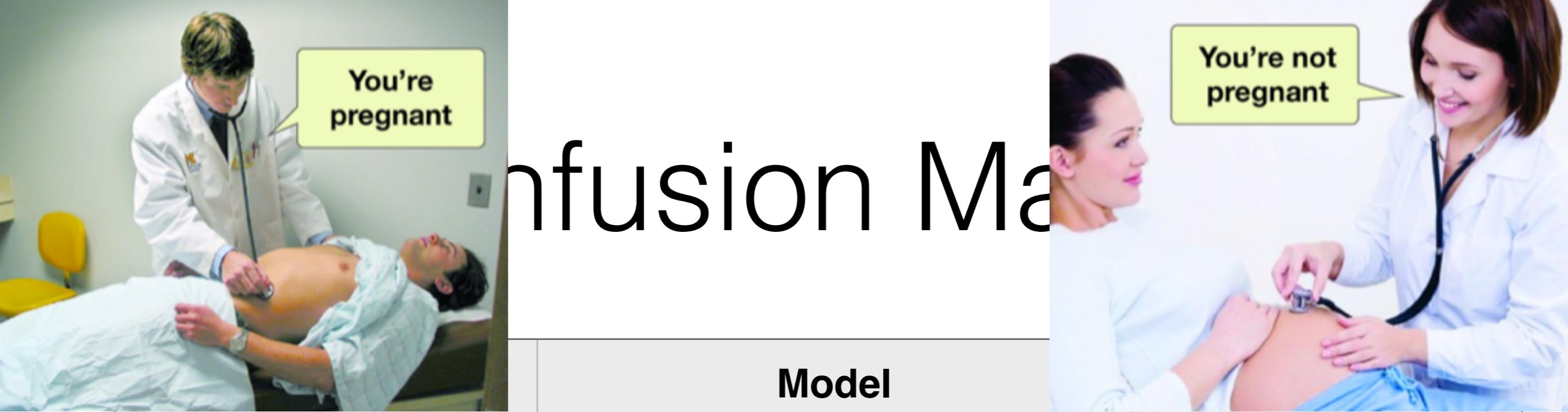
$$\frac{10100}{10100 + 6} = 99.94\%$$

## Type I error (false positive)



## Type II error (false negative)





# Confusion Matrix

		Model			
		Positive	Negative		
<b>TRUE</b>	Positive	a	b	Sensitivity, (recall, TPR)	b/(a+b)
	Negative	c	d	FPR	c/(c+d)
		Precision	False omission rate	Accuracy $= (a+d) / (a+b+c+d)$	
		$a/(a+c)$	$b/(b+d)$		

# ROC Curves

Receiver Operating Characteristic

```
sklearn.metrics.*
```

# Royal Observer Corps

**Post Instrument** - prior to the introduction of radar, were commonly used to spot and report aircraft positions.



# Receiver Operating Characteristic

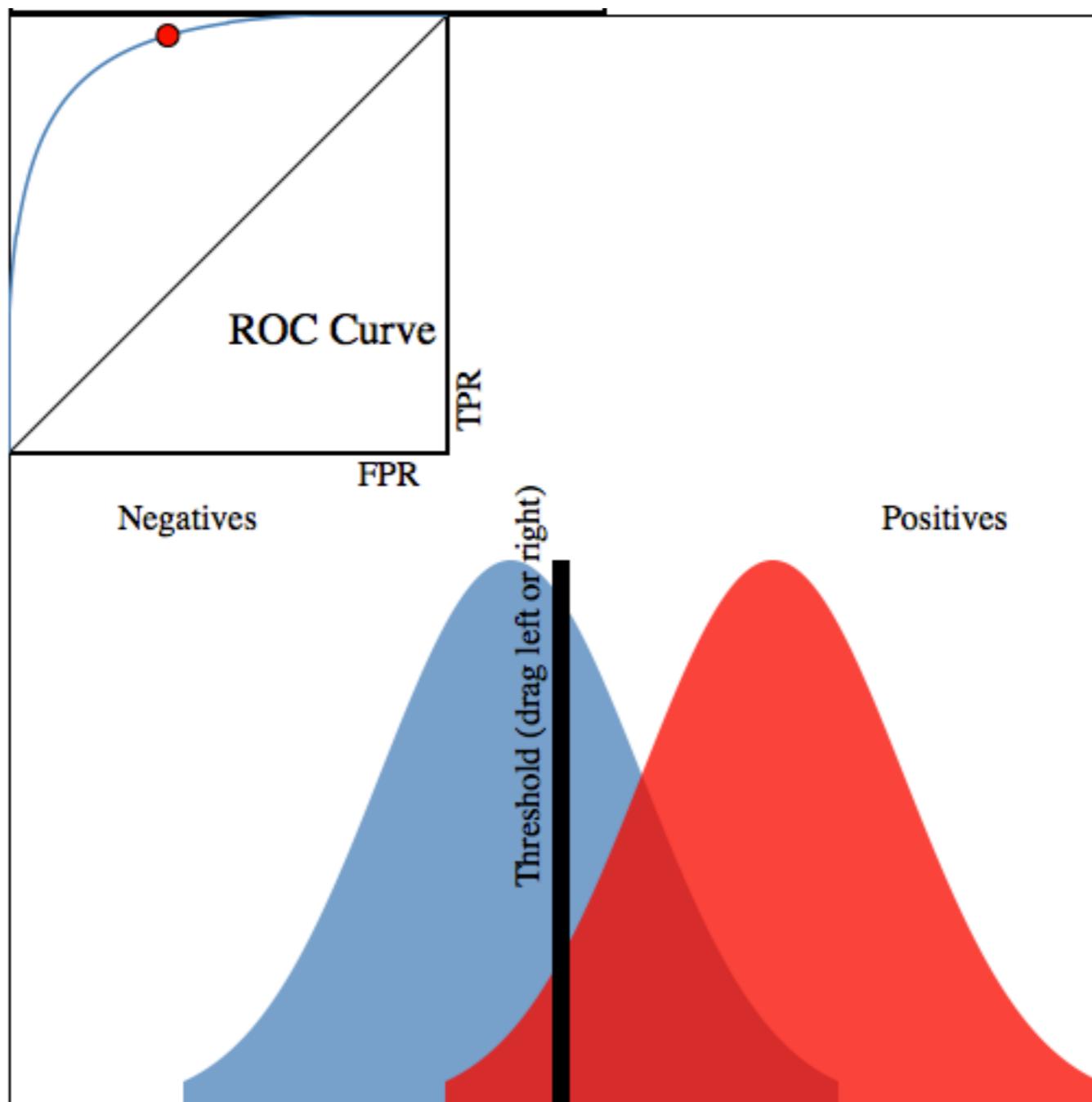


For a first radar, was a challenge to distinguish **bird and plane**

First\* Great Britain Radar

\* at least one of them

# ROC curves



# The Boy Who Cried Wolf



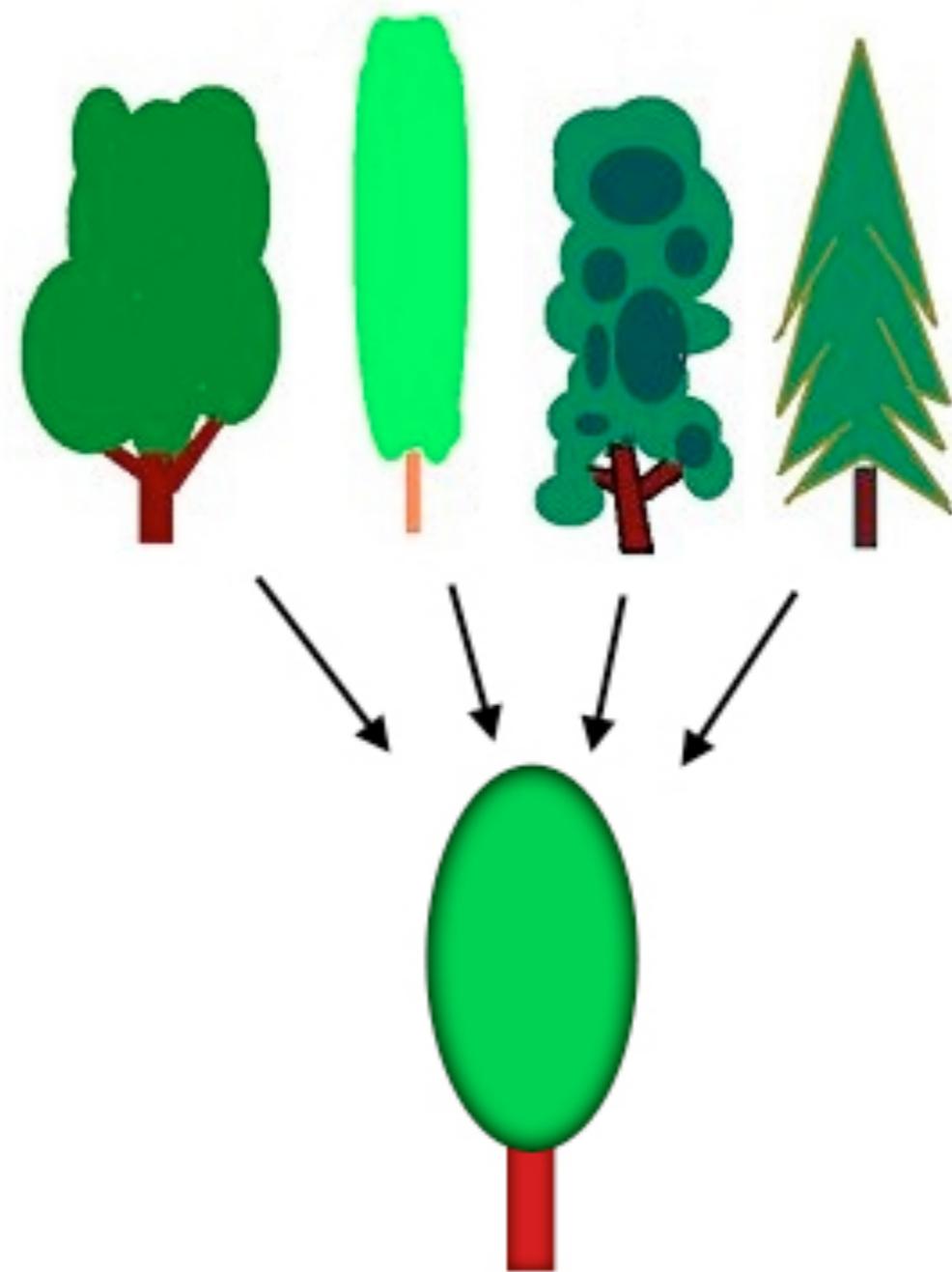
Ian Healey

# Validation

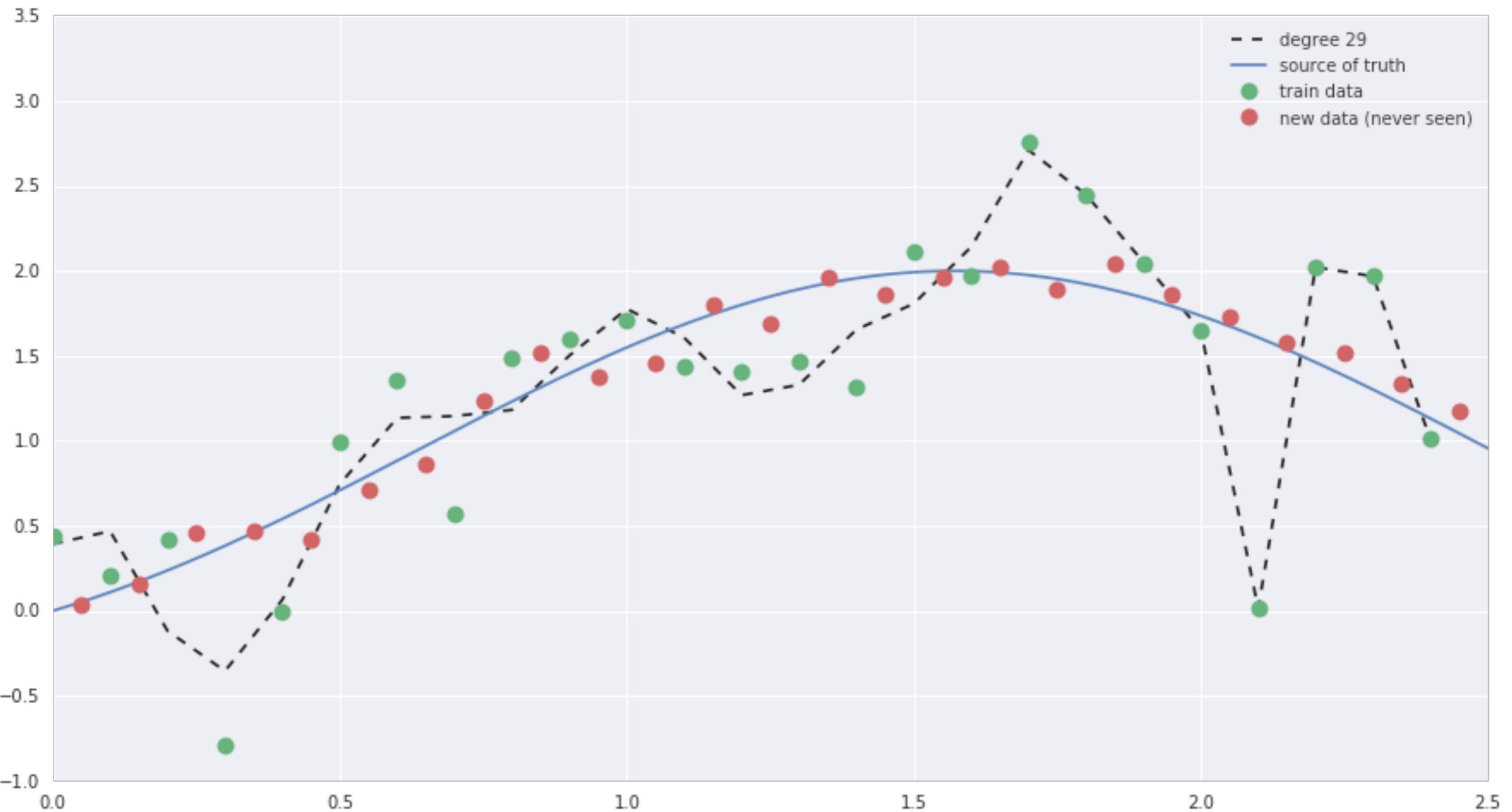
cross-validation (k-fold, nested), bootstrapping

```
sklearn.cross_validation.*
```

# Generalization



# Overfitting



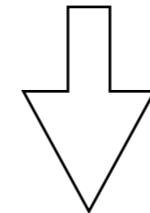
# Hold Out Validation

Split for two datasets

```
sklearn.cross_validation.train_test_split
```

# Bootstrapping

Age	Gender	Height
20	male	170
30	female	163
24	male	198
28	female	169



Age	Gender	Height
20	male	170
30	female	163
24	male	198

Age	Gender	Height
28	female	169

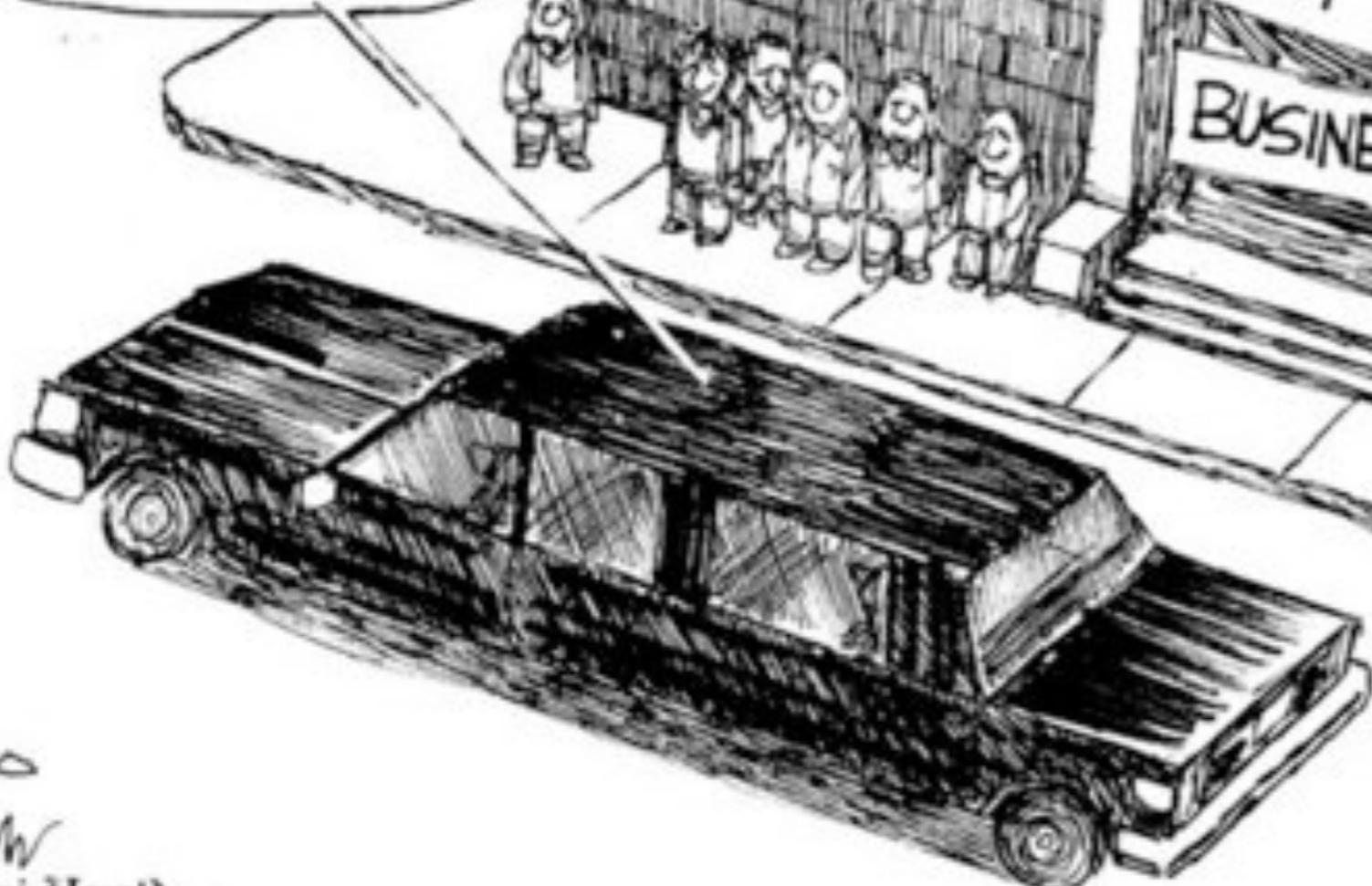
# Bootstrapping resampling

```
sklearn.utils.resample
```

# ACME BOOTSTRAPS, inc

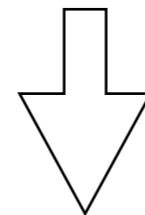
PULL YOURSELF  
UP BY YOUR  
BOOTSTRAPS!

OUT OF  
BUSINESS



# Bootstrapping

Age	Gender	Height
20	male	170
30	female	163
24	male	198
28	female	169



Age	Gender	Height
20	male	170
24	male	198
28	female	169

Age	Gender	Height
20	male	170
30	female	163
30	female	163

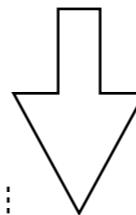
Age	Gender	Height
30	female	163
24	male	198
28	female	169

# Cross validation

k-fold, leave-one-out, stratified

# Cross-Validation: 3-fold

Age	Gender	Height
20	male	170
30	female	163
24	male	198
28	female	169



Age	Gender	Height
20	male	170
30	female	163
24	male	198

Age	Gender	Height
28	female	169
20	male	170
30	female	163

Age	Gender	Height
24	male	198
28	female	169
20	male	170

Age	Gender	Height
28	female	169

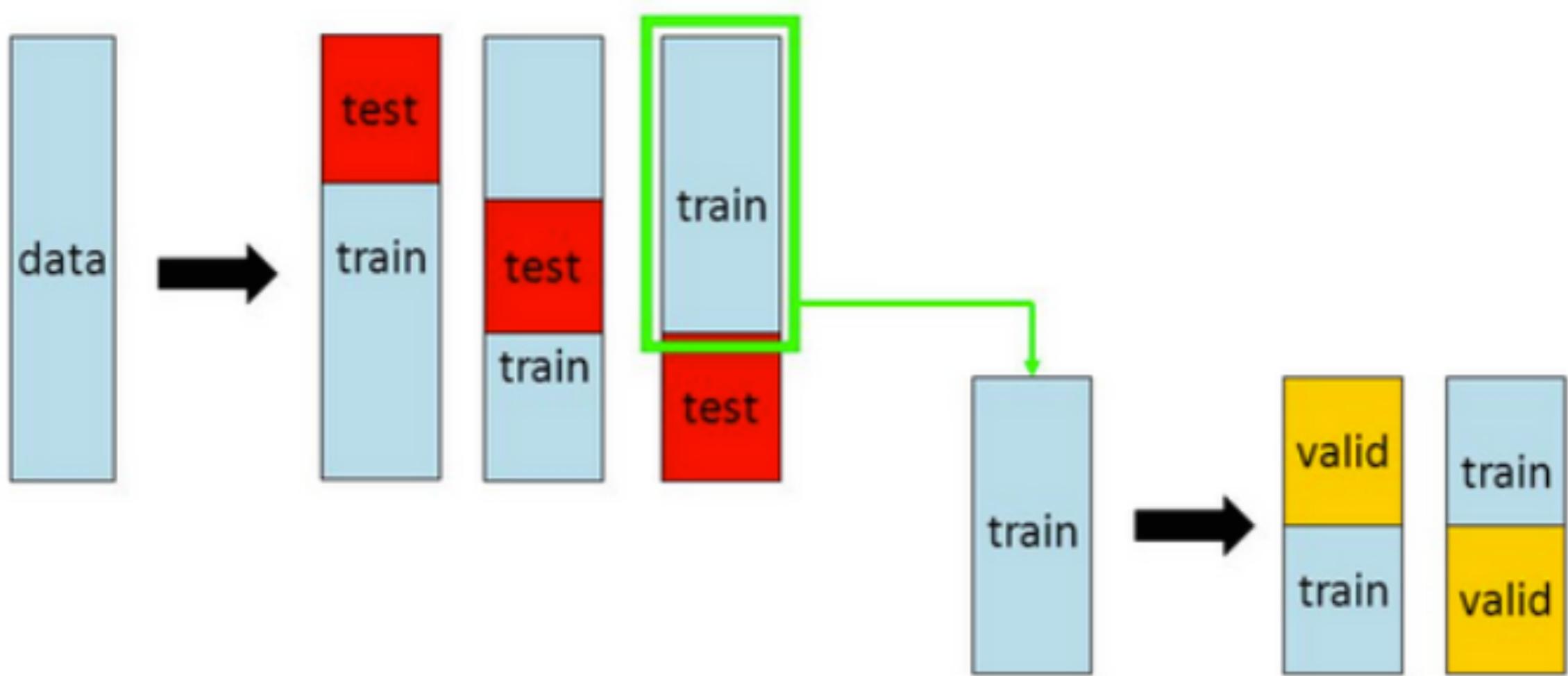
Age	Gender	Height
24	male	198

Age	Gender	Height
30	female	163

# Nested Cross Validation

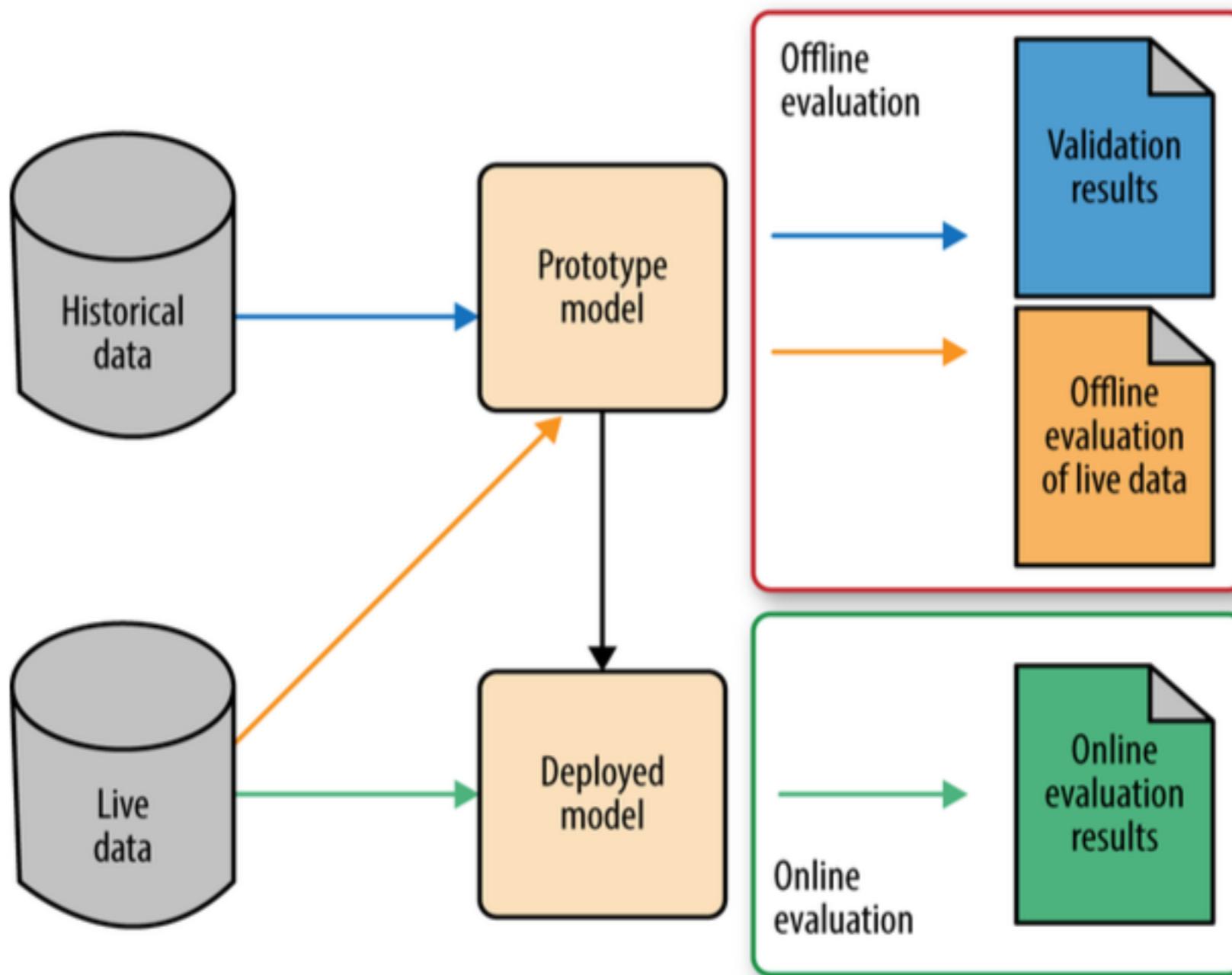
More advanced approach

# Nested Cross Validation



# Summary

# Development & evaluation model



# Three things

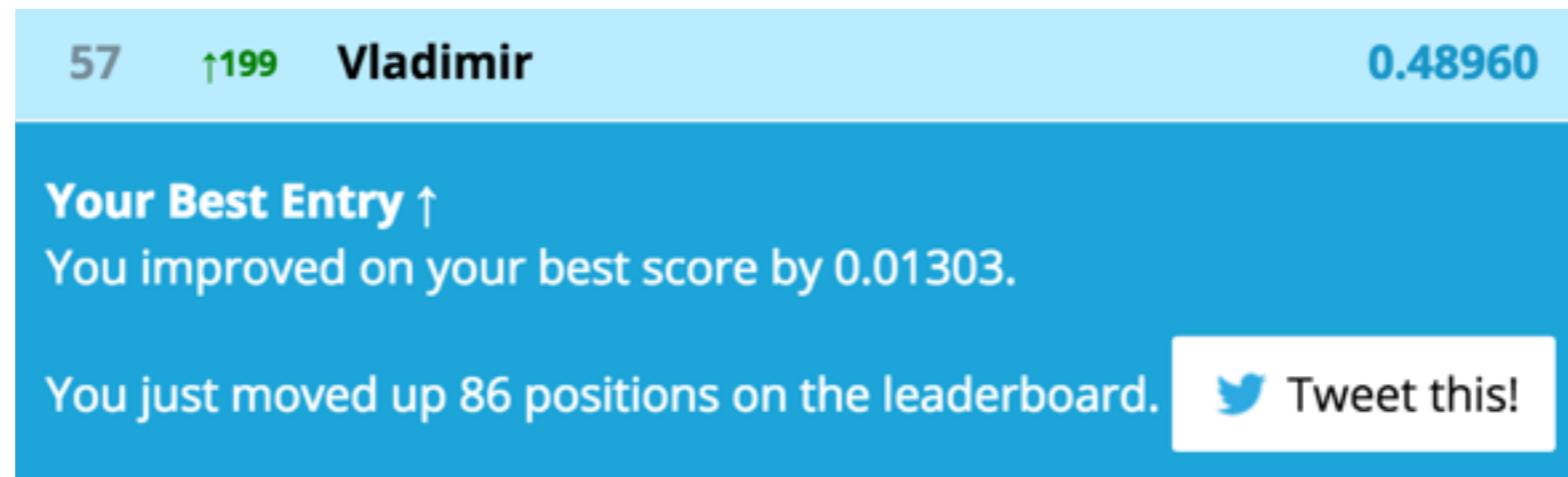
*if you can remember only three...*

- Understand expected success by metric[s]
- Never mix training data and evaluation data
- The distribution of data changes over time

New challenge:



Maximize sales and minimize returns of  
bakery goods



Are you ready to join?

Thank you