**2017 FAST CAMPUS**
**SCHOOL PRESENTATION**

**YOUR CAREER**
**TURNING POINT**

**SCHOOL**
**DATA SCIENCE**

Fast campus

2017 FAST CAMPUS SCHOOL

# DATA SCIENCE SCHOOL PROJECT ( 1 )

/

## REGRESSION ANALYSIS

DATA SCIENCE SCHOOL

**2017 FAST CAMPUS
SCHOOL PRESENTATION**

**YOUR CAREER
TURNING POINT**

**SCHOOL
DATA SCIENCE**

# 목차

DATA SCIENCE SCHOOL

**2017 FAST CAMPUS
SCHOOL PRESENTATION**

**YOUR CAREER
TURNING POINT**

**SCHOOL
DATA SCIENCE**

# 팀 공개 A반

1조
**박재근, 박재근, 박재근**

2조
**박재근, 박재근, 박재근**

3조
**박재근, 박재근, 박재근**

4조
**박재근, 박재근, 박재근**

5조
**박재근, 박재근, 박재근**

6조
**박재근, 박재근, 박재근**

7조
**박재근, 박재근, 박재근**

DATA SCIENCE SCHOOL

**2017 FAST CAMPUS
SCHOOL PRESENTATION**

**YOUR CAREER
TURNING POINT**

**SCHOOL
DATA SCIENCE**

## 팀 공개 B반

**1조**
윤병관, 백승민, 김인수

**2조**
배광빈, 송세현, 공명구

**3조**
최규형, 진미나, 염승식

**4조**
노범용, 유정오, 양영규, 김영법

**5조**
박상하, 김현규, 한상훈

**6조**
최윤철, 이기훈, 김성희, 안동순

**7조**
고정욱, 편설인, 이원재

DATA SCIENCE SCHOOL

**2017 FAST CAMPUS**
**SCHOOL PRESENTATION**

**YOUR CAREER**
**TURNING POINT**

**SCHOOL**
**DATA SCIENCE**

# 데이터셋 소개



Dataset

Training | Testing | Holdout Method

Cross Validation

Data Permitting:

Training | Validation | Testing | Training, Validation, Testing

Joseph Nelson @josephofiowa

DATA SCIENCE SCHOOL

**2017 FAST CAMPUS
SCHOOL PRESENTATION**

**YOUR CAREER
TURNING POINT**

**SCHOOL
DATA SCIENCE**

# 데이터셋 소개 : Train vs Test



**train.csv**

**test.csv**

Dataset

Training — **Testing** Holdout Method

**정답지(y값)가 주어진 데이터셋
모델링 작업을 위한 데이터셋**

**블라인드 테스트
정답지(y값) 이 없음
최종 테스트를 위한 데이터셋**

Cross Validation

Data Permitting:

Training | Validation | Testing | Training, Validation, Testing

Joseph Nelson @josephofiowa

DATA SCIENCE SCHOOL

**2017 FAST CAMPUS
SCHOOL PRESENTATION**

**YOUR CAREER
TURNING POINT**

**SCHOOL
DATA SCIENCE**

fast campus

# 데이터셋 소개 : Cross Validation



**모델완성**

**모의고사**

Joseph Nelson @josephofiowa

**2017 FAST CAMPUS
SCHOOL PRESENTATION**

**YOUR CAREER
TURNING POINT**

**SCHOOL
DATA SCIENCE**

# 데이터셋 소개 : Final Test

**2017 FAST CAMPUS**
**SCHOOL PRESENTATION**

**YOUR CAREER**
**TURNING POINT**

**SCHOOL**
**DATA SCIENCE**

# 주제 설명

1. ## House Prices: Advanced Regression Techniques
   Predict sales prices and practice feature engineering, RFs, and gradient boosting

2. ## New York City Taxi Trip Duration
   Share code and data to improve ride time predictions

3. ## Sberbank Russian Housing Market
   Can you predict realty price fluctuations in Russia's volatile economy?

4. ## Toyota Corolla Prices
   Predict used Toyota Corolla car prices

DATA SCIENCE SCHOOL

**2017 FAST CAMPUS**
**SCHOOL PRESENTATION**

**YOUR CAREER**
**TURNING POINT**

**SCHOOL**
**DATA SCIENCE**

# House Prices: Advanced Regression Techniques

## GOAL

It is your job to predict the sales price for each house.
For each Id in the test set, you must predict the value of the SalePrice variable.

## METRIC

Submissions are evaluated on Root-Mean-Squared-Error (RMSE) between the logarithm of the predicted value and the logarithm of the observed sales price. (Taking logs means that errors in predicting expensive houses and cheap houses will affect the result equally.)

$$\mathrm{RMSD}(\hat{\theta}) = \sqrt{\mathrm{MSE}(\hat{\theta})} = \sqrt{\mathrm{E}((\hat{\theta} - \theta)^2)}.$$

**2017 FAST CAMPUS
SCHOOL PRESENTATION**

**YOUR CAREER
TURNING POINT**

**SCHOOL
DATA SCIENCE**

House Prices: Advanced Regression Techniques

## SUBMISSION FILE FORMAT

The file should contain a header and have the following format:

```
Id,SalePrice
1461,169000.1
1462,187724.1233
1463,175221
...
```

## DATA SIZE

- train.csv : 1030rows, 81colums (@SalePrice)
- test.csv : 430rows, 80colums

DATA SCIENCE SCHOOL

**2017 FAST CAMPUS**
**SCHOOL PRESENTATION**

**YOUR CAREER**
**TURNING POINT**

**SCHOOL**
**DATA SCIENCE**

fast campus

# New York City Taxi Trip Duration

## GOAL

It is your job to predict the duration time for each id.
For each Id in the test set, you must predict the value of the trip_duration variable.

## METRIC

The evaluation metric for this competition is Root Mean Squared Logarithmic Error.
The RMSLE is calculated as

$$\epsilon = \sqrt{\frac{1}{n} \sum_{i=1}^{n} (\log(p_i + 1) - \log(a_i + 1))^2}$$

$\epsilon$ is the RMSLE value (score)
$n$ is the total number of observations in the (public/private) data set,
$p_i$ is your prediction of trip duration, and
$a_i$ is the actual trip duration for $i$.
$\log(x)$ is the natural logarithm of $x$

DATA SCIENCE SCHOOL

**2017 FAST CAMPUS
SCHOOL PRESENTATION**

**YOUR CAREER
TURNING POINT**

**SCHOOL
DATA SCIENCE**

fast
campus

## New York City Taxi Trip Duration

## SUBMISSION FILE FORMAT

The file should contain a header and have the following format:

```
id,trip_duration
id00001,978
id00002,978
id00003,978
id00004,978
...
```

## DATA SIZE

-   train.csv : 701778rows, 11colums (@dropoff_datetime, trip_duration)
-   test.csv : 346797rows, 9colums

DATA SCIENCE SCHOOL

**2017 FAST CAMPUS**
**SCHOOL PRESENTATION**

**YOUR CAREER**
**TURNING POINT**

**SCHOOL**
**DATA SCIENCE**

Fast campus

# Sberbank Russian Housing Market

## GOAL

It is your job to predict the sales price for each house.
For each Id in the test set, you must predict the value of the price_doc variable.

## METRIC

The evaluation metric for this competition is Root Mean Squared Logarithmic Error.
The RMSLE is calculated as

$$\epsilon = \sqrt{\frac{1}{n}\sum_{i=1}^{n}(\log(p_i + 1) - \log(a_i + 1))^2}$$

$\epsilon$ is the RMSLE value (score)
$n$ is the total number of observations in the (public/private) data set,
$p_i$ is your prediction of trip duration, and
$a_i$ is the actual trip duration for $i$.
$\log(x)$ is the natural logarithm of $x$

**2017 FAST CAMPUS
SCHOOL PRESENTATION**

**YOUR CAREER
TURNING POINT**

**SCHOOL
DATA SCIENCE**

Fast campus

# Sberbank Russian Housing Market

## SUBMISSION FILE FORMAT

The file should contain a header and have the following format:

```
id,price_doc
30474,7118500.44
30475,7118500.44
30476,7118500.44
...
```

## DATA SIZE

- train.csv : 21570rows, 292colums (@price_doc)
- test.csv : 8901rows, 291colums
- macro.csv : 2484rows, 100colums - data on Russia's macroeconomy and financial sector
          (could be joined to the train and test sets on the "timestamp" column)

DATA SCIENCE SCHOOL

**2017 FAST CAMPUS**
**SCHOOL PRESENTATION**

**YOUR CAREER**
**TURNING POINT**

**SCHOOL**
**DATA SCIENCE**

# Toyota Corolla Prices

## GOAL

It is your job to predict the sale price of a used automobile.
For each Id in the test set, you must predict the value of the Price variable.

## METRIC

Submissions are evaluated on Root-Mean-Squared-Error (RMSE) between the logarithm of the predicted value and the logarithm of the observed sales price. (Taking logs means that errors in predicting expensive houses and cheap houses will affect the result equally.)

$$\text{RMSD}(\hat{\theta}) = \sqrt{\text{MSE}(\hat{\theta})} = \sqrt{\text{E}((\hat{\theta} - \theta)^2)}.$$

**2017 FAST CAMPUS
SCHOOL PRESENTATION**

**YOUR CAREER
TURNING POINT**

**SCHOOL
DATA SCIENCE**

**fast
campus**

Toyota Corolla Prices

## SUBMISSION FILE FORMAT

The file should contain a header and have the following format:

```
Id,Price
3,13950
4,14950
6,12950
7,16900
...
```

## DATA SIZE

- train.csv : 1019rows, 39colums (@Price)
- test.csv : 417rows, 38colums

DATA SCIENCE SCHOOL

**2017 FAST CAMPUS
SCHOOL PRESENTATION**

**YOUR CAREER
TURNING POINT**

**SCHOOL
DATA SCIENCE**

# 주제 선정

**- House Prices
- NYC Taxi
- Sberbank House
- Toyota Corolla**

한 주제당 최대 **2조** 까지 선택 가능 **(main)**

한 주제 이상 분석 가능 **( optional )**

DATA SCIENCE SCHOOL

**2017 FAST CAMPUS**
**SCHOOL PRESENTATION**

**YOUR CAREER**
**TURNING POINT**

**SCHOOL**
**DATA SCIENCE**

fast campus

# 데이터 다운로드

# GITHUB

github.com/JKeun/dss-regression-datasets

DATA SCIENCE SCHOOL

**2017 FAST CAMPUS
SCHOOL PRESENTATION**

**YOUR CAREER
TURNING POINT**

**SCHOOL
DATA SCIENCE**

# 타임라인

~10/10
# EDA
**Exploratory data analysis**



DATA SCIENCE SCHOOL

**2017 FAST CAMPUS**
**SCHOOL PRESENTATION**

**YOUR CAREER**
**TURNING POINT**

**SCHOOL**
**DATA SCIENCE**

# 타임라인

~10/14

# FEATURE SELECTION

```python
var_real = ['Mfg_Month', 'Mfg_Year', 'HP', 'CC', 'Quarterly_Tax', 'Weight',
            'Guarantee_Period', 'KM','Age_08_04']
var_cat = ['Color','Fuel_Type', 'Met_Color', 'Automatic', 'ABS', 'Airbag_1',
           'Airbag_2', 'Airco', 'Automatic_airco', 'Boardcomputer', 'CD_Player', 'Central_Lock',
           'Powered_Windows', 'Power_Steering', 'Radio', 'Mistlamps', 'Sport_Model',
           'Backseat_Divider', 'Metallic_Rim', 'Radio_cassette', 'Parking_Assistant', 'Doors',
           'Tow_Bar','Cylinders','Gears','Mfr_Guarantee','BOVAG_Guarantee',]

print len(var_real), len(var_cat)
```

9 27

**2017 FAST CAMPUS**
**SCHOOL PRESENTATION**

**YOUR CAREER**
**TURNING POINT**

**SCHOOL**
**DATA SCIENCE**

fast campus

# 타임라인

~10/16
# MODELING
**OLS Regression**

**2017 FAST CAMPUS
SCHOOL PRESENTATION**

**YOUR CAREER
TURNING POINT**

**SCHOOL
DATA SCIENCE**

# 타임라인

~10/21
# MODEL SELECTION

**DATA SCIENCE SCHOOL**

**Milestone #12. 최종 회귀모델 결정**

**Price = 19,230 - 138.583(Age_08_04) - 0.0165(KM) + 546.796(InOut) - 396.891(Safe)**

**2017 FAST CAMPUS**
**SCHOOL PRESENTATION**

**YOUR CAREER**
**TURNING POINT**

**SCHOOL**
**DATA SCIENCE**

# 타임라인

~10/25
# PROJECT 발표

10/24 까지 프로젝트 **발표자료** 완성 ( Jupyter Notebook, PPT )

**10월 25일 수요일** 프로젝트 발표

DATA SCIENCE SCHOOL

**2017 FAST CAMPUS**
**SCHOOL PRESENTATION**

**YOUR CAREER**
**TURNING POINT**

**SCHOOL**
**DATA SCIENCE**

END

DATA SCIENCE SCHOOL

# THANK YOU ; - )