

# DATA MINING AND TEXT MINING

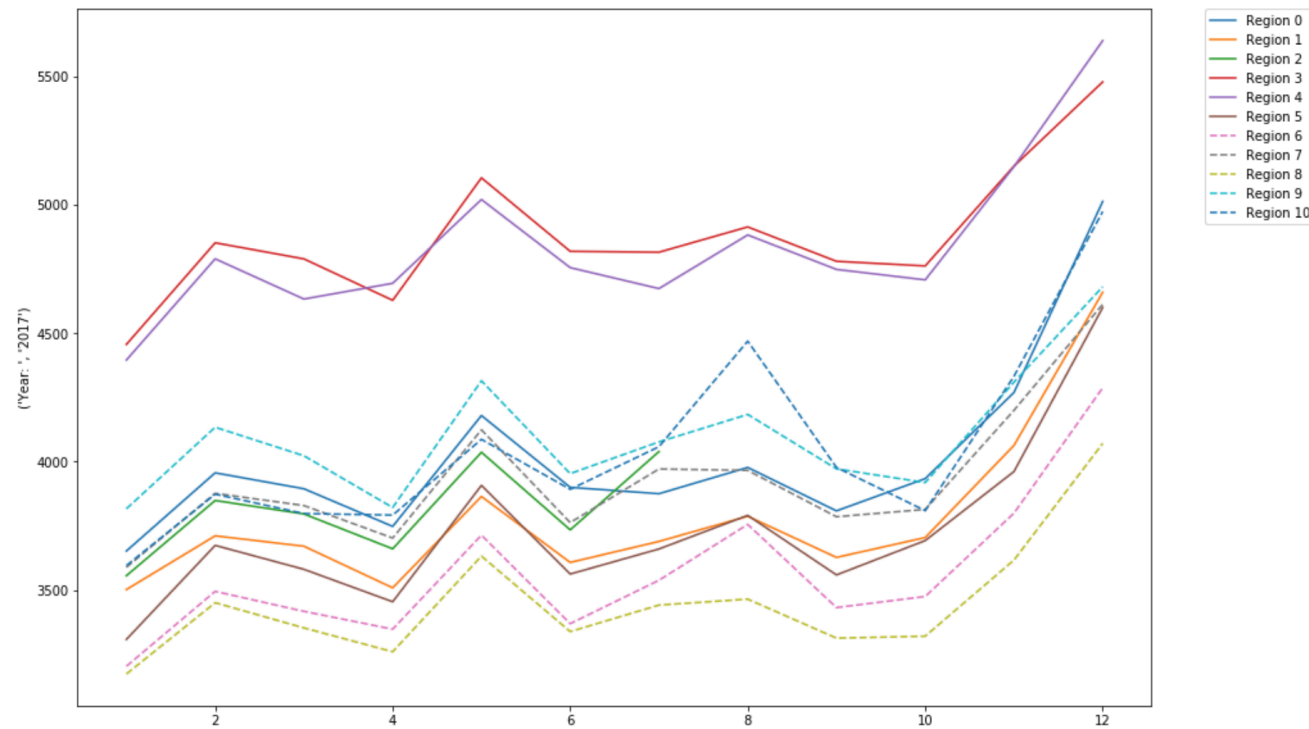


## CHALLENGE

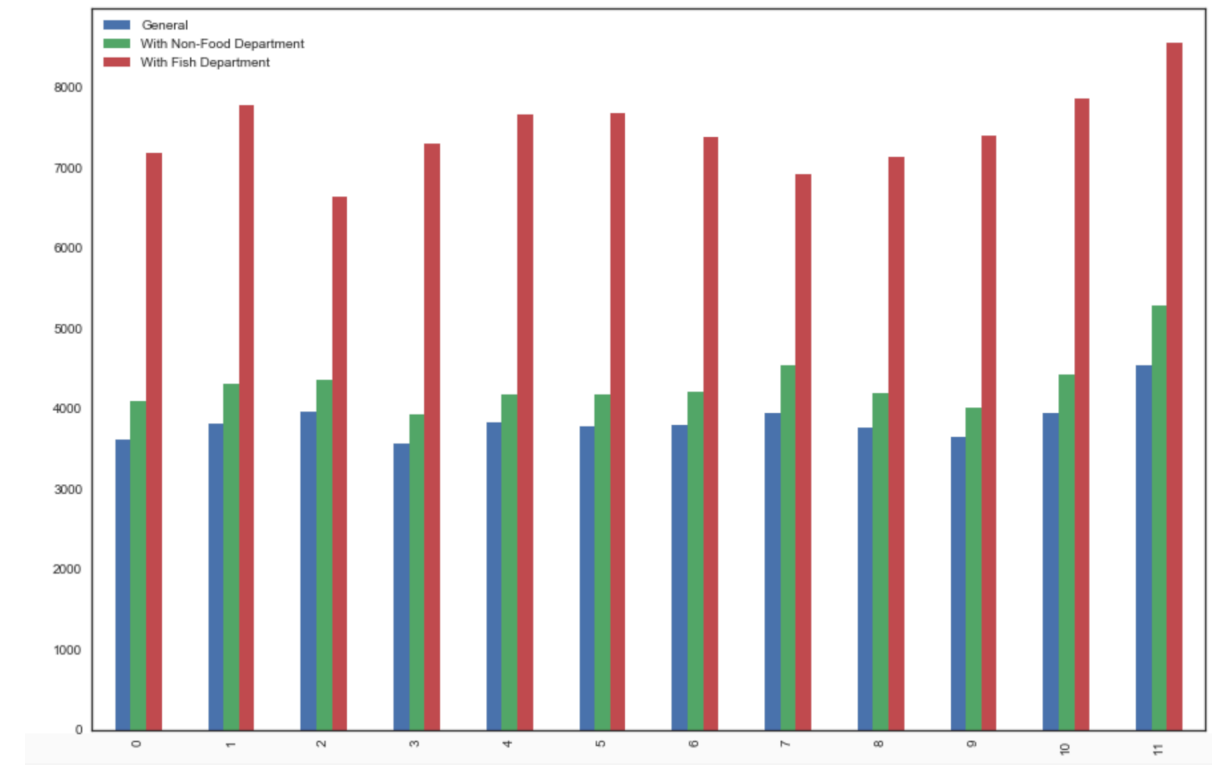
LORENZO NORCINI  
GUGLIELMO MENCHETTI

# DATA ANALYSIS

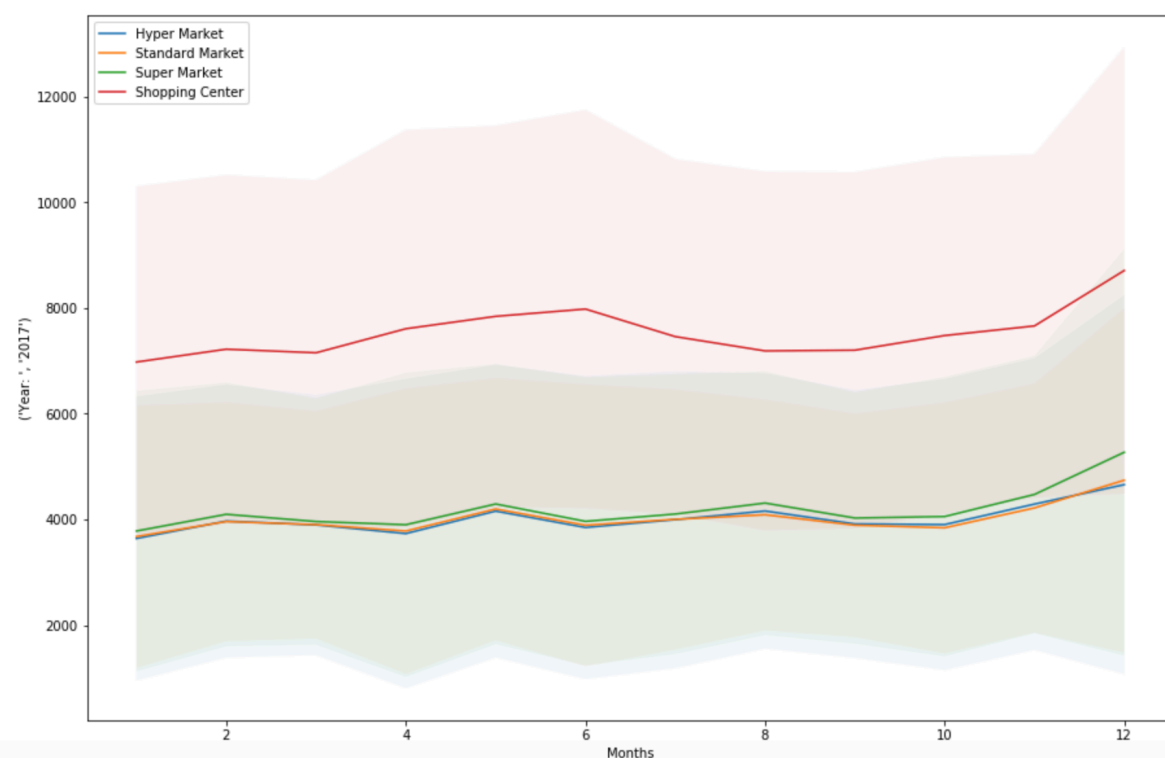
## Monthly Sales per Region



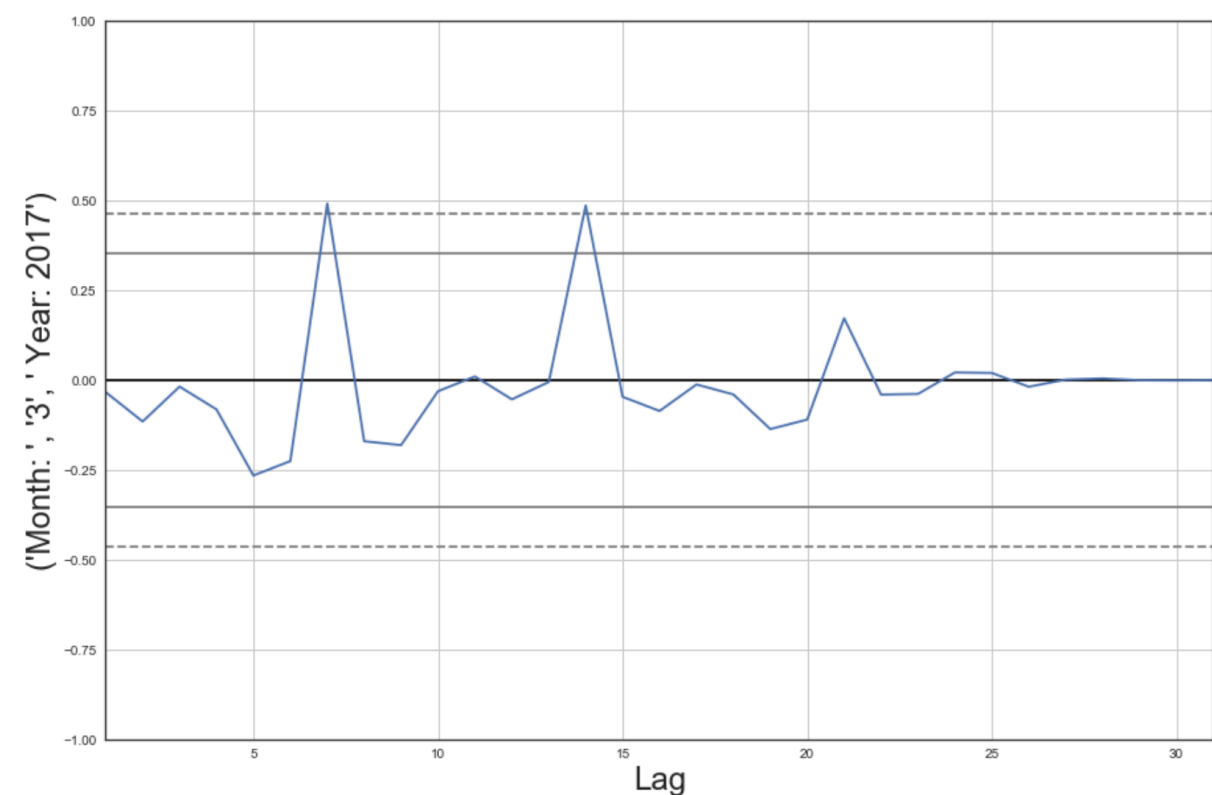
## Monthly Sales per Assortment Type



## Monthly Sales per Store Type



## Autocorrelation

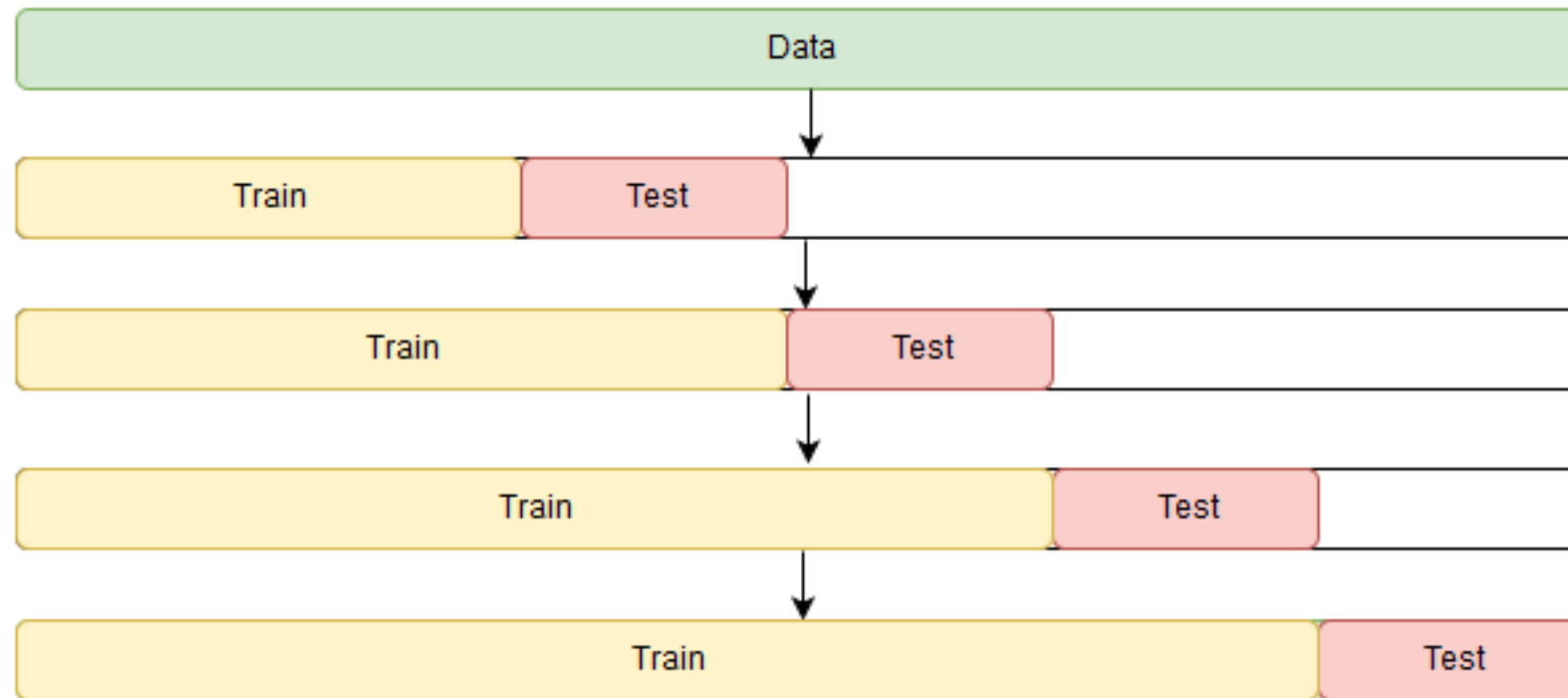


# FEATURES AND MODEL PARAMETERS

Store Features	Time Features	Meteo Features
<ul style="list-style-type: none"><li>• IsHoliday</li><li>• HasPromotions</li><li>• StoreType</li><li>• NearestCompetitor</li><li>• Region</li><li>• AssortmentType</li></ul>	<ul style="list-style-type: none"><li>• Rolling Mean (14, 30, 60, 90)</li><li>• Lag (1, 7, 14, 21, 28)</li><li>• Day of the week</li><li>• Month</li></ul>	<ul style="list-style-type: none"><li>• Events</li><li>• Cloud Cover</li><li>• Precipitation</li></ul>

Model	Parameters
SVM	<ul style="list-style-type: none"><li>• Error Penalty (0.1, 1, 10, 100)</li><li>• Kernel (rbf, linear)</li></ul>
Random Forest	<ul style="list-style-type: none"><li>• Number of Estimators (10, 20, 50, 100)</li><li>• Bootstrap (True, False)</li><li>• Max Features (all, log2, sqrt)</li><li>• Min samples for leaf (1, 2, 3, 4, 5, 6, 7, 8, 9, 10)</li></ul>
K Nearest Neighbours	<ul style="list-style-type: none"><li>• Number of Neighbours (5, 10, 20, 50, 75, 100)</li><li>• Distance Metric (Euclidian, Chebyshev, Manhattan)</li><li>• Weights (Uniform, Distance)</li></ul>
Adaboost + KNN	<ul style="list-style-type: none"><li>• Number of Estimators (10, 20, 50, 100)</li><li>• Number of Neighbours (5, 10, 20, 50, 75, 100)</li><li>• Distance Metric (Euclidian, Chebyshev, Manhattan)</li><li>• Weights (Uniform, Distance)</li></ul>

# FEATURES AND MODEL SELECTION



Time Series  
Cross Validation

2 Months Overlapped  
Test set size

Performance Metric

$$E_r = \frac{\sum_{i \in S_r} \sum_{j \in \{3,4\}} |a_{ij} - p_{ij}|}{\sum_{i \in S_r} \sum_{j \in \{3,4\}}^j a_{ij}}$$

**Region Error**

$$E = \frac{\sum_{r \in R} E_r}{|R|}$$

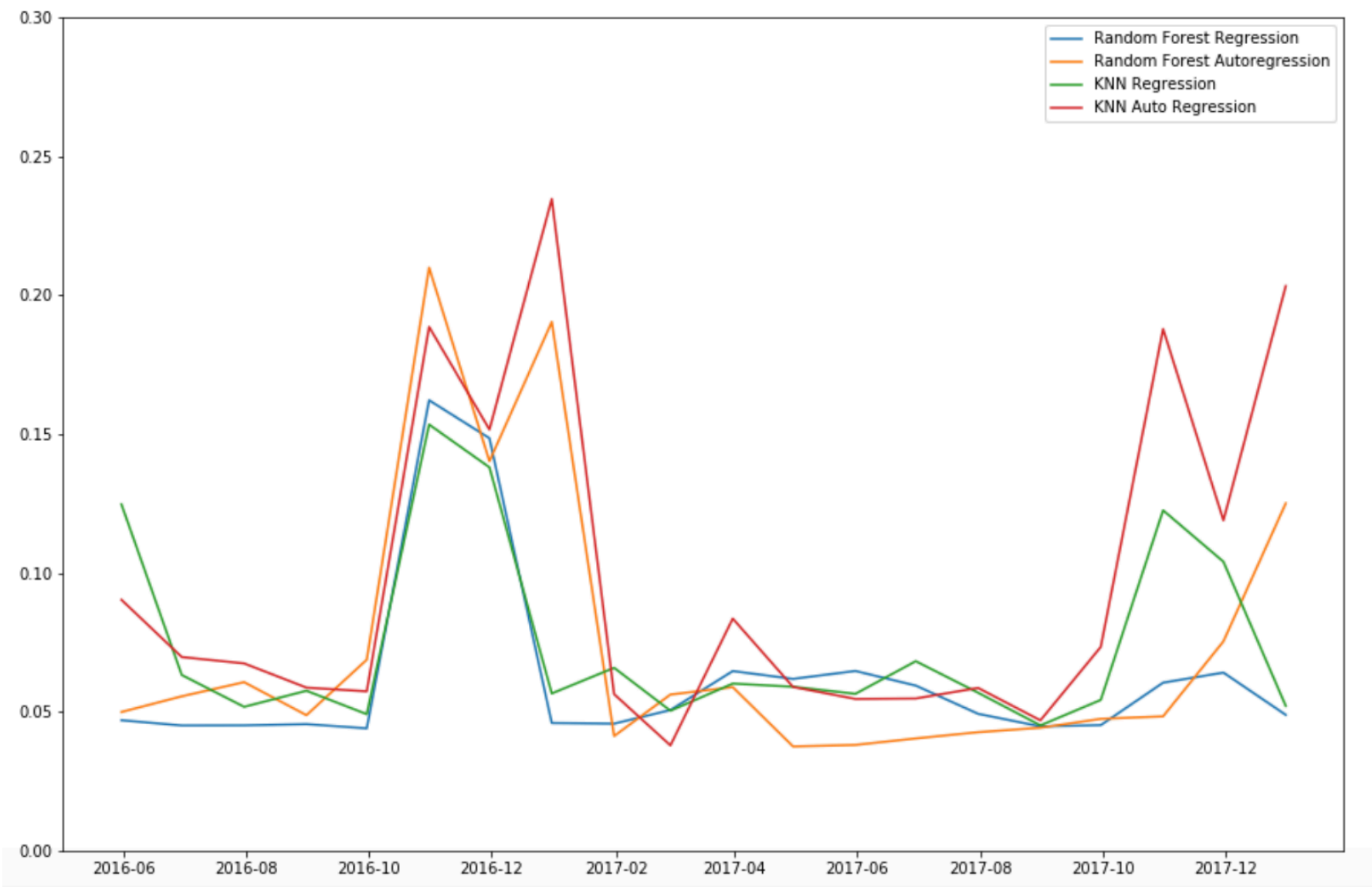
**Total Error**

$R$  Set of regions

$a_{ij}$  Actual Value of Store  $i$  and month  $j$      $p_{ij}$  Predicted Value of Store  $i$  and month  $j$

# PERFORMANCES

Type	Learner	$\mu$	$\sigma$
Regression	K-NN	0.087	0.05
Regression	Random Forest	0.062	0.03
Autoregression	K-NN	0.097	0.06
Autoregression	Random Forest	0.074	0.05



Performances on  
Test folds

Learning Curves

# OBTAINED MODELS

Standard Regression		Auto Regression
Features	<ul style="list-style-type: none"><li>• IsHoliday</li><li>• HasPromotions</li><li>• StoreType</li><li>• NearestCompetitor</li><li>• Day of the week</li><li>• Region</li><li>• Month</li><li>• AssortmentType</li></ul>	<ul style="list-style-type: none"><li>• IsHoliday</li><li>• HasPromotions</li><li>• StoreType</li><li>• NearestCompetitor</li><li>• WeekDay 'Region,</li><li>• Month</li><li>• AssortmentType</li><li>• Lag 7 Days</li><li>• Lag 14 Days</li><li>• Rolling mean 60 Days</li></ul>
Model	Random Forest	Random Forest
Parameters	<ul style="list-style-type: none"><li>• Number of Estimator : 100</li><li>• Bootstrap</li><li>• Max Features : all</li><li>• Min samples for leaf :5</li></ul>	<ul style="list-style-type: none"><li>• Number of Estimator : 50</li><li>• Bootstrap</li><li>• Max Features : all</li><li>• Min samples for leaf :1</li></ul>
$\mu$	0.062	0.074
$\sigma$	0.03	0.05

# DATA MINING AND TEXT MINING



**THANK YOU FOR YOUR  
ATTENTION**

**LORENZO NORCINI  
GUGLIELMO MENCHETTI**