

# **Dataset cleaning and other things to check**

# In this video

- Dataset cleaning
  - Constant features
  - Duplicated features
- Other things to check
  - Duplicated rows
  - Check if dataset is shuffled

# Duplicated and constant features

<i>is_train</i>	<b>f0</b>	<b>f1</b>	<b>f2</b>	<b>f3</b>	<b>f4</b>	<b>f5</b>
<i>True</i>	13	H	1.2	1.2	A	C
<i>True</i>	13	H	36.6	36.6	B	A
<i>False</i>	13	H	0	0	A	C
<i>False</i>	13	G	-14	-14	C	B

# Duplicated and constant features

<b><i>is_train</i></b>	<b><i>f0</i></b>	<b><i>f1</i></b>	<b><i>f2</i></b>	<b><i>f3</i></b>	<b><i>f4</i></b>	<b><i>f5</i></b>
<i>True</i>	<i>13</i>	H	1.2	1.2	A	C
<i>True</i>	<i>13</i>	H	36.6	36.6	B	A
<i>False</i>	<i>13</i>	H	0	0	A	C
<i>False</i>	<i>13</i>	G	-14	-14	C	B

```
| traintest.nunique(axis=1) == 1
```

# Duplicated and constant features

<i>is_train</i>	<b>f0</b>	<i>f1</i>	<b>f2</b>	<b>f3</b>	<b>f4</b>	<b>f5</b>
<i>True</i>	13	H	1.2	1.2	A	C
<i>True</i>	13	H	36.6	36.6	B	A
<i>False</i>	13	H	0	0	A	C
<i>False</i>	13	G	-14	-14	C	B

```
| train.nunique(axis=1) == 1
```

# Duplicated and constant features

<i>is_train</i>	<b>f0</b>	<b>f1</b>	<i>f2</i>	<i>f3</i>	<b>f4</b>	<b>f5</b>
<i>True</i>	13	H	1.2	1.2	A	C
<i>True</i>	13	H	36.6	36.6	B	A
<i>False</i>	13	H	0	0	A	C
<i>False</i>	13	G	-14	-14	C	B

```
| traintest.T.drop_duplicates()
```

# Duplicated and constant features

<i>is_train</i>	<b>f0</b>	<b>f1</b>	<b>f2</b>	<b>f3</b>	<i>f4</i>	<i>f5</i>
<i>True</i>	13	H	1.2	1.2	A	C
<i>True</i>	13	H	36.6	36.6	B	A
<i>False</i>	13	H	0	0	A	C
<i>False</i>	13	G	-14	-14	C	B

```
for f in categorical_feats:  
    traintest[f] = raintest[f].factorize()  
traintest.T.drop_duplicates()
```

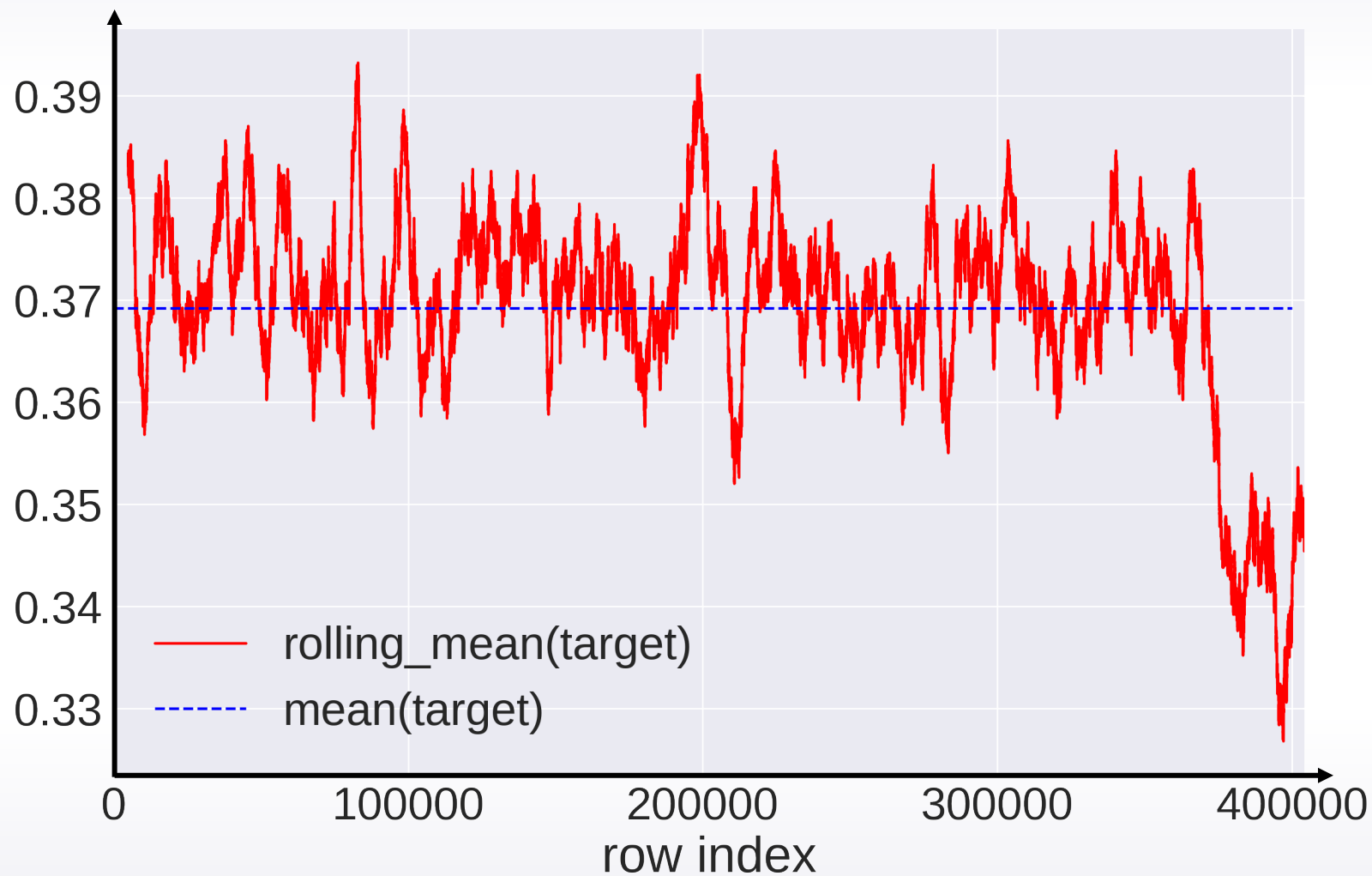
# Duplicated rows

<b>f1</b>	<b>f2</b>	<b>f3</b>	<b>y</b>
13	34r9	A	<b>0</b>
13	34r9	A	<b>1</b>
13	34r9	A	<b>1</b>

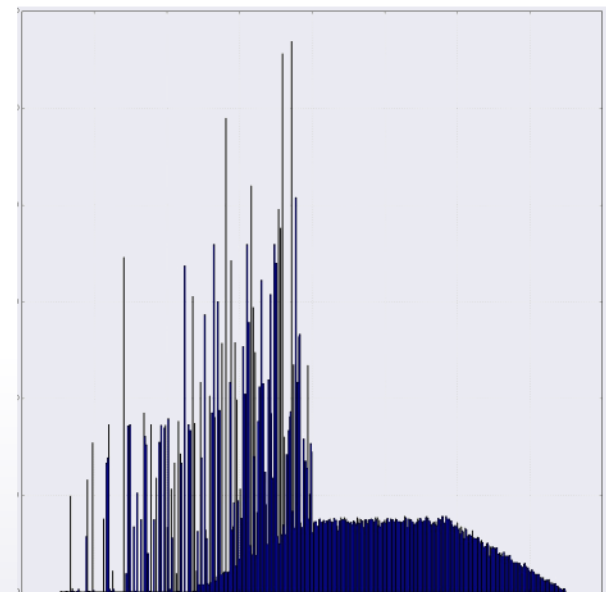
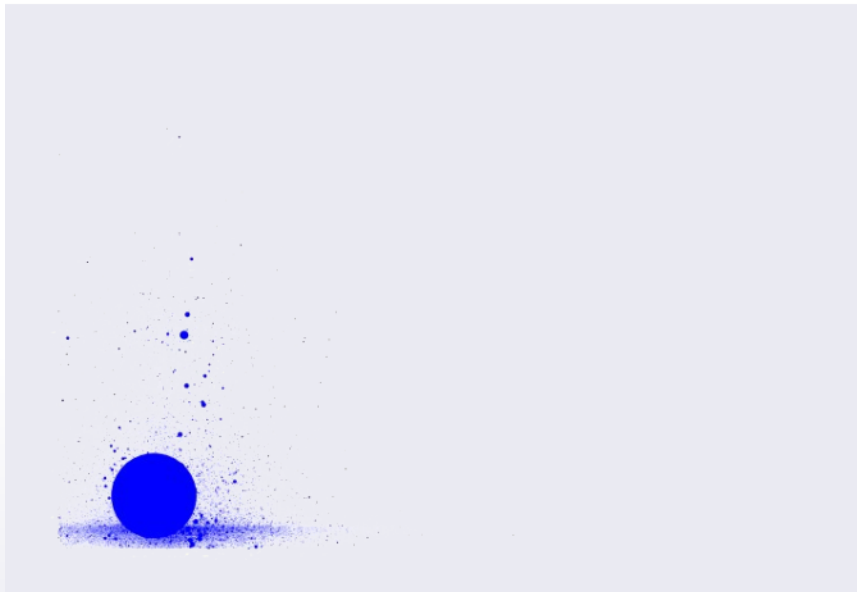
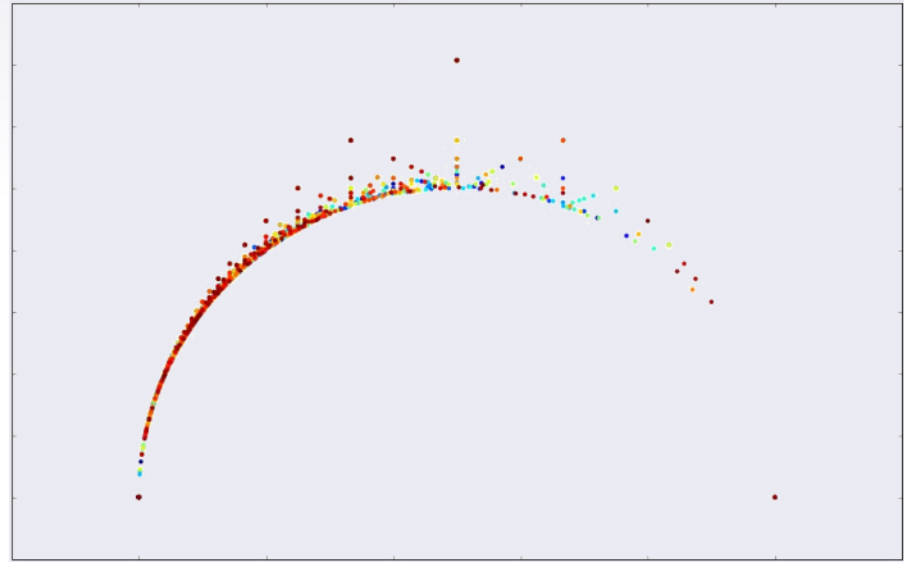
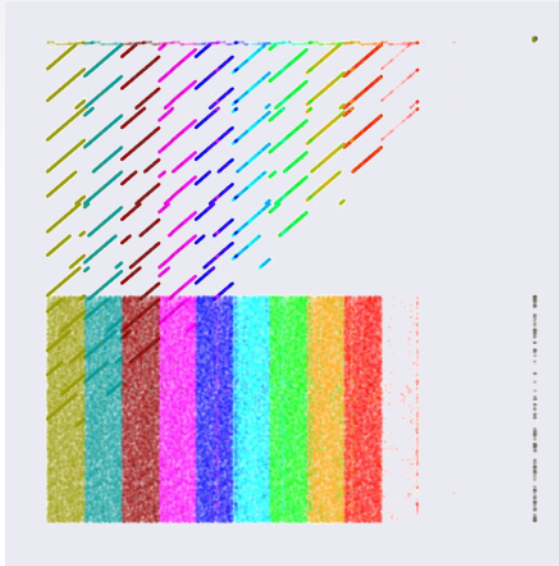
- Check if same rows have same label
- Find duplicated rows, understand why they are duplicated



# Check if dataset is shuffled



# Cool visualizations



# EDA check list

- Get domain knowledge
  - Check if the data is intuitive
  - Understand how the data was generated
- 

- Explore individual features
  - Explore pairs and groups
- 

- Clean features up
- 

- Check for leaks! (later in this course)