

## Greg Lee-Podgorski - Head of Machine Learning & Technology Platforms @ RS

# CLASS PREDICTION USING CLASSIFICATION - TOP SELLING PRODUCTS

## I. Define & Scope

### The Brief

Based on internal RS corporate data on best-selling products we carry, use classification method(s) to analyze product features available in the dataset and find most suitable ones for class prediction (conversion rate below/above mean). Then create a visualization of the resulting analysis. Draw conclusions, if findings allow it. ¶

#### i) Intro & my role in the organization

In this exercise I'm going to apply classification method(s) of Machine Learning to our best-selling products in RS Components product catalogue. While classification algorithms consist of a number of algorithms available to be used, from naive Bayes, through Decision Tree(s) to Random Forest, none of those is perfect. I'm going to utilize couple of those aforementioned to mine for most correlated features that affect class prediction among those top selling products, trying to determine which class each observation belongs to and hoping to find some hidden patterns or shared characteristics in the process.

I happen to hold a position of Head of Machine Learning, Automation & Technology Platform within the Electrocomponents PLC Group (RS Components is a trading name in the UK and EMEA). Thus I am an outspoken advocate for increased use of ML in our organization and also a passionate believer in value that deploying ML algorithms into our operations can unlock for us. I'm hoping that by showing use cases of certain algorithms, I will lead by example, showing how ML algos can be used to find hidden patterns in data and provide business colleagues with valuable insights.

#### ii) Background

RS Components is a publicly listed corporation and trades on LSE (London Stock Exchange) under Electrocomponents PLC trade name. We have significant trading presence in each of the top economically most active regions - EMEA, APAC and Americas. Due to wildly varying standards across any number of electronic products specifications that have local requirements, whether voltage, number of prongs/holes in a electric socket, certification standards or just local engineer's preferences, standarization is a challenge for RS and thus we need to stock very large range of products. As it stands RS currently has over 500,000 SKU's (over 0.5 million product in stock) and over 2.5 million non-stock products which we can drop-ship for our clients from our global supplier base. That means that while most of online stores would be content with coming up with a Top 100 best selling products, out Best Sellers list features slightly over 10,000 products. Some of the products on that best-sellers' list vary only slightly from their "cousins", e.g. by different AC/DC output or electrical resistance (in Ohms).

To make matters even more complicated we carry a number of distinct categories of products, e.g. BLE (Board Level Electronics) - something a computer maker like Asus or Dell would need, or IA&C (Undustrial Automation & Control) - e.g. assembly line sensors monitoring health of a conveyor belt, or PPE (Personal Protection Equipment) - e.g. gloves or protective masks engineers or doctors can cover their faces with.

### **iii) Problem statement and potential business application**

With such huge range of products covering a number of distinct categories and sourced from over 2500 unique suppliers across different geographies, it's little wonder that our inventory management function is 3-digit strong when it comes to headcount and inventory planners are always playing a catch-up game. Whether it comes to updating ERP system with up-to-date stock counts from a number of out trade counters (brick and mortar store locations), keeping track what's about to run out in our warehouses located in key countries globally or even trying to assign priorities correctly across several different team within Inventory Management, while RS as a company noticeably improved our approach to managing stock, we are surely in need of further improvements. Those will be difficult to come by and even more so - to implement - unless we manage to identify significant patterns in underlying data.

Being able to predict if a given product, product family or newly sourced parts' range will have above or below average conversion rate will enable our marketing team to customize our sales offers and highlight promotions to online shoppers when they are interacting with those products that based on their features we expect to fight uphill battle to reach average conversion rate and direct inventory teams to stock less of those where we expect conversion rates to struggle. In more clear-cut cases they might even be able to terminate orders from our suppliers or renegotiate better deals so that we in turn can offer better pricing and thus provide a lift to those struggling conversion rate products.

[https://uk.rs-online.com/web/op/allproducts/?intcmp=UK-WEB\\_-\\_HP-TC0\\_-\\_201711\\_-\\_ALLCATEGORIES](https://uk.rs-online.com/web/op/allproducts/?intcmp=UK-WEB_-_HP-TC0_-_201711_-_ALLCATEGORIES)

All Products	Our Brands	New Products	My Account	Our Services
Home > All Products				
<b>All Products</b>				
<b>Abrasives &amp; Engineering Materials</b>	<b>Facilities Cleaning &amp; Maintenance</b>	<b>Pneumatics, Hydraulics &amp; Power Transmission</b>		
Aluminum Tubes, Sheets & Angles	Brooms, Mops, Buckets & Dust Pans	Electric Actuators		
Brass Tubes, Sheets & Rods	Car Care & Polishes	Hydraulic & Pneumatic Tools		
Bronze & Gunmetal Tubes & Rods	Cleaners, Degreasers & Removers	Hydraulic Adaptors, Fittings & Couplings		
Carbon Fibre & Felt Sheets	Cleaning Brushes	Hydraulic Cylinders, Pumps & Power Units		
Ceramic Rods, Sheets & Beads	Cleaning Wipes & Cloths	Hydraulic Fluids & Filtration		
Copper Rods, Sheets & Bars	Desktop Cleaning Equipment	Hydraulic Instrumentation & Switches		
Insulation Materials	Electronics Cleaners & Protective Coatings	Hydraulic Tubing & Hose		
Manual Sanding & Sharpening	Facility Maintenance Equipment	Hydraulic Valves & Manifolds		
Mild Steel Tube, Sheets & Angles	Floor Scrubbers & Accessories	Pneumatic & Hydraulic Pressure Gauges		
Plastic Rods, Sheets & Tubes	Greases, Oils & Lubricants	Pneumatic Adaptors, Fittings & Couplings		
Polishing & Finishing	Paint & Painting Supplies	Pneumatic Air Compressors, Boosters & Vacuum Pumps		
Rubber Sheets	Vacuum Cleaners & Carpet Cleaners & Accessories	Pneumatic Air Preparation		
Shim Stock	Washroom Equipment & Supplies	Pneumatic Counters, Logic Controllers & Timers		
Stainless Steel Tubes, Sheets & Angles	Waste Removal Bins & Accessories	Pneumatic Cylinders & Actuators		
<b>Adhesives, Sealants &amp; Tapes</b>	<b>Fasteners &amp; Fixings</b>	Pneumatic Instrumentation & Switches		
Adhesive, Sealant & Tape Dispensers	Channel Support Systems	Pneumatic Tubing & Air Hose		
Adhesives & Glues	Clips	Pneumatic Valves & Manifolds		
Hot Melt Glue Guns & Accessories	Fastener Kits	Power Transmission - Linear Bearings, Housings & Blocks		
Sealants	Girder Fixings	Power Transmission - Bearing Tools		
Tapes	Hooks & Eyes	Power Transmission - Belts		
<b>Automation &amp; Control Gear</b>	Levelling & Vibration Control	Power Transmission - Bushes & Collars		
Circuit Protection & Circuit Breakers	Magnets & Magnetic Strips	Power Transmission - Clutches & Brakes		
Contactors	Nails	Power Transmission - Couplings		
Control Relays	Nuts & Washers	Power Transmission - Gaskets, Seals & Packings		
Electric Motors, Motor Controllers & Peripherals	Pins, Keys & Retaining	Power Transmission - Gears		
Engineering Services	Rivets & Riveting Tools	Power Transmission - Linear Shafts, Rails, Ball Screws & Lead Screws		
Fluid Control Systems	Screws & Bolts	Power Transmission - Linear Slides, Guides & Positioning Tables		
Industrial Push Buttons, Pilot Lights & Control Stations	Sheet Metal & Panel Fasteners	Power Transmission - Pulleys		
Industrial Switches	Spacers & Standoffs	Power Transmission - Rod Ends & Spherical Bearings		
Machine Guarding & Safety	Threaded Rods & Studs	Power Transmission - Roller Chains & Accessories		
Panel Meters & Components	Wall Plugs, Anchors, Fixings & Kits	Power Transmission - Rotary Bearings		
PLCs, HMI & Data Acquisition	Wire Rope Suspension Systems	Power Transmission - Sprockets		
Sensors & Transducers	<b>Fuses, Sockets &amp; Circuit Breakers</b>	Vacuum Components		
Solenoids	Circuit Breakers	<b>Power Supplies &amp; Transformers</b>		
Sounders & Beacons	Consumer Units & Accessories	DC-AC Power Inverters		
Temperature Control & Process Heating	Electrical Installation Accessories	DC-DC Converters		
Timers & Counters	Fuse Holders	Portable Generators & Accessories		
<b>Batteries &amp; Chargers</b>	Fuse Tools & Testing	Power Conditioners		
Battery & Charger Accessories	Fuses	Power over Ethernet - PoE		
Battery Chargers & Power Banks	Fuses - PCB	Power Supplies - PSUs		
Non-Rechargeable Batteries	Lightning Protection	Renewable Energy		
Rechargeable Batteries	<b>HVAC, Fans &amp; Thermal Management</b>	Transformers		
<b>Cables &amp; Wires</b>	Air Conditioning & Climate Control Units	Uninterruptible Power Supplies - UPS		
Audio Video Cable	Air Filters & Accessories	<b>Power Tools, Soldering &amp; Welding</b>		
Cable Accessories, Ties & Tools	Air Management Accessories	Air Tools		
Cable Conduit, Trunking & Routing	Electronics Heating & Cooling	Drill Bits & Parts		
Cable Glands, Strain Relief & Grommets	Electronics Humidity & Pressure Control Devices	Milling & Lathing		
Coaxial Cable	Fan Parts & Accessories	Power Tool Accessories		
Computer Cable Assemblies	Fans	Power Tools		
Electrical Power & Industrial Cable	HVAC Ducting	Sanding Belts, Discs & Wheels		
Network & Communication Cable	HVAC Sensors & Controllers	Soldering		
Ribbon & Flat Cable	<b>Lighting</b>	Welding & Brazing		
Wire & Single Core Cable	Emergency Lighting & Safety Lighting	Workshop Tools		
Wire to Board Cable Assemblies	Fluorescent Lamps & Tubes			
<b>Computing &amp; Peripherals</b>	Halogen Lamps			
3D Printing & Scanning	HID Lamps			
Audio & Video	Incandescent Light Bulbs			
Barcode Readers & Accessories	Indicators & Indicator Components			
	Infrared Lamps			
		<b>Relays</b>		
		General Purpose Relay Accessories		
		General Purpose Relays		
		Interface Relay Modules & Accessories		
		Direct Dial-in		

#### iv) Data

The dataset has been sourced internally, with a help of a product analyst, who is a part of inventory management function and looks after best-selling products, making sure they are in stock with maximum availability. The Excel files utilized here has over 11,000 rows with multiple columns providing information such as Category, Technology type, number of orders, etc.

A	B	C	D	E	F	G	H	I	J	K	L	M
Article	Description	FSTD Global	Category	Cell	Technology	Visit	Product Views	In Stock Views	Cart Additions	Order %	Conversion Rate	Revenue
26	9092061 PCM INTERFACE PAD FOR HS300, 128X72.5MM	12/31/9999 I&T	EMECH	Thermal Management	98	101	0.713	6	2	0.0204	135.68	
27	291717 NYLON 6.6 HOSE CLIP,26.5-29.5MM	12/31/9999 TCFM	Consumables & Safety	Fasteners	281	256	0.835	43	17	0.0605	182.56	
28	5217528 DELRIN SPUR GEAR - 1.0 MODULE 25 TEETH	12/31/9999 IAAC	Mech & Flu Pwr	Pow Tran & Struc	226	247	0.781	36	15	0.0664	205.26	
29	9225141 SOLID STATE RELAY	12/31/9999 IACC	Sensors & Switches	Ctrl & Ind Relay	274	312	0.254	16	1	0.0036	211.08	
30	2664268 MEDIUM DISPOSABLE CO-POLYMER GLOVES	12/31/9999 TCFM	Consumables & Safety	PPE	274	264	0.989	41	6	0.0219	213.6	

### v) Data protection

The dataset used doesn't contain any customer-related info (other than number of orders and views in our online store) so there is no personally-identifiable information available. Thus there is no risk of any GDPR violation. As far as corporate data is concerned, there is no strictly confidential data being used. Almost all of the data used is publically available on our web store. The only potentially sensitive data in there is revenue generated by each of the top selling products. If accessed in its entirety, one could imagine a potential risk of a competitor being able to peek behind the curtain and try to replicate - say - our Top 1000 selling products and undercutting our pricing. To avoid that, I have sorted revenue in ascending order and never show the column content - only few rows at a time.

## Loading libraries

```
In [1]: # for Loading and manipulating dataframes
2 import pandas as pd
3
4 # for additional mathematical functionality
5 import numpy as np
6
7 # Seaborn and matplotlib for visualisation
8 import matplotlib.pyplot as plt
9 import seaborn as sns
10
11 # for building and evaluating our classification models
12 from sklearn import tree, metrics # This is for us to simply evaluate our mo
13 from sklearn.model_selection import train_test_split
14 from sklearn.ensemble import RandomForestClassifier # This provides the impl
15 from sklearn.metrics import confusion_matrix #This is for a more thorough ev
16 from sklearn.metrics import accuracy_score, precision_score, recall_score, f
17
18 # for plotting our tree:
19 from sklearn.externals.six import StringIO
20 from IPython.display import Image
21 from sklearn.tree import export_graphviz
22 import pydotplus
```

C:\Users\oszer\anaconda3\lib\site-packages\sklearn\externals\six.py:31: FutureWarning: The module is deprecated in version 0.21 and will be removed in version 0.23 since we've dropped support for Python 2.7. Please rely on the official ve  
rsion of six (<https://pypi.org/project/six/>).  
"(<https://pypi.org/project/six/>).", FutureWarning)

## Sourcing data

```
In [2]: 1 best_sell = pd.read_excel("../Top Sellers Digital Data_KP.xlsx")
```

```
In [3]: 1 best_sell.tail()
```

Out[3]:

		Article	Description	FSTD Global	Category	Cell	Technology	Visits	Product Views	Stock View
11453	7609005	ALUMINIUM PAINT ENCLOSURE 160X160X90MM	9999- 12-31 00:00:00	IA&C	Mach Guard & Enc	Enclosures	390	383	0.7%	
11454	4992470	PIPE CLEANING BRUSH FOR 15MM DIA PIPE	9999- 12-31 00:00:00	TCFM	Tools Storage & FM	Hand Tools	1194	1125	0.9%	
11455	1257961	INDUSTRIAL CABLE 10MM2 10M	9999- 12-31 00:00:00	II&T	Cable, Anti-Vib & Ac	Elec & Data Cable	291	335	0.2%	
11456	3971315	MINERAL INSULATED PT 100 SENSOR,3X150MM	9999- 12-31 00:00:00	IA&C	Sensors & Switches	Proc Auto Sen	392	408	0.9%	
11457	2414782	ADHESIVE BASE MICRO WIRE SADDLE,9X7.3MM	9999- 12-31 00:00:00	II&T	Cable, Anti-Vib & Ac	Cable Accs	591	561	0.9%	

## Exploring, Cleaning & Transforming data

### Exploring

Let's first have a quick look at what type data and its shape we are dealing with..

```
In [4]: 1 best_sell.shape
```

Out[4]: (11458, 13)

As we can see, our dataframe has almost 11,500 rows and 13 columns

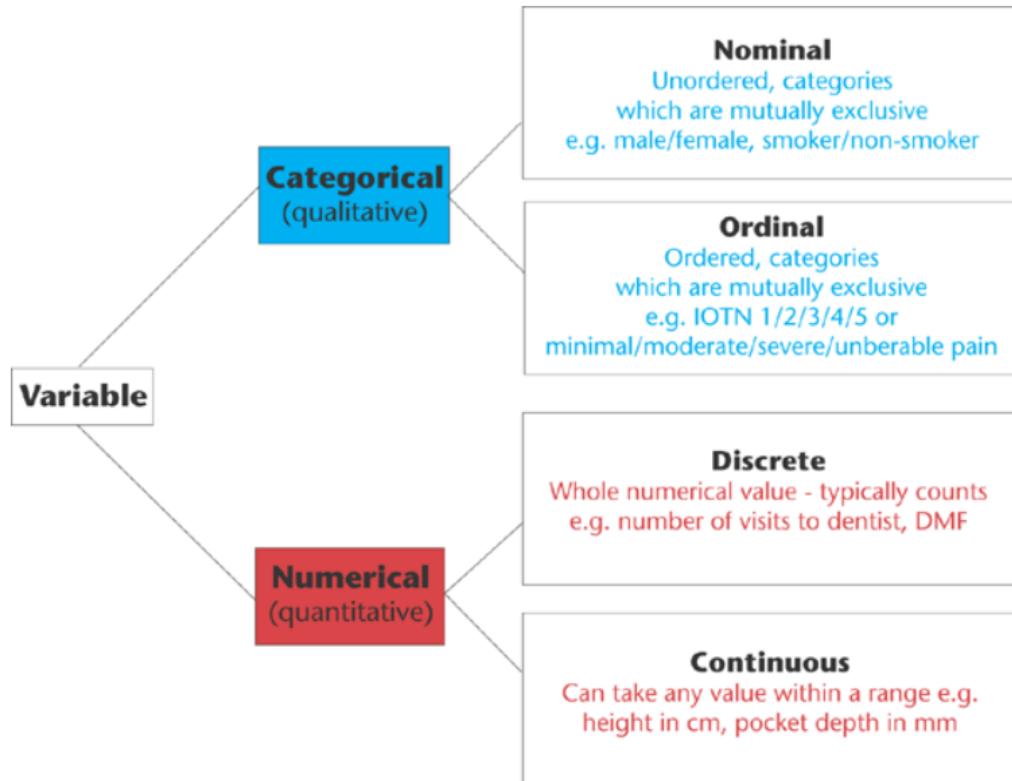
In [5]: 1 best\_sell.info()

```
<class 'pandas.core.frame.DataFrame'>
RangeIndex: 11458 entries, 0 to 11457
Data columns (total 13 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Article          11458 non-null   int64  
 1   Description      11458 non-null   object  
 2   FSTD Global      11458 non-null   object  
 3   Category          11458 non-null   object  
 4   Cell              11458 non-null   object  
 5   Technology        11458 non-null   object  
 6   Visits            11458 non-null   int64  
 7   Product Views    11458 non-null   int64  
 8   In Stock Views   11458 non-null   float64 
 9   Cart Additions   11458 non-null   int64  
 10  Orders            11458 non-null   int64  
 11  Conversion Rate  11458 non-null   float64 
 12  Revenue           11458 non-null   float64 
dtypes: float64(3), int64(5), object(5)
memory usage: 1.1+ MB
```

Good news - no null values in any of the rows - that doesn't happen often. I guess my colleague made sure the dataset has been of proper quality.

We can see that majority of the columns contain categorical values, while remaining 5 columns (1st column, 7th, 8th, 10th and 11th col) have numerical records.

Here is a quick refresher on how we classify types of variables:



As per above dichotomy, we can see that columns 6 through 11 contain discrete numerical values, the same for Article column and while Revenue columns appear to contain continuous numerical variables. Remaining columns possess nominal variables.

In [6]: 1 best\_sell.dtypes

Out[6]:

Article	int64
Description	object
FSTD Global	object
Category	object
Cell	object
Technology	object
Visits	int64
Product Views	int64
In Stock Views	float64
Cart Additions	int64
Orders	int64
Conversion Rate	float64
Revenue	float64
dtype:	object

We can further see that columns Article, Visits, Product Views, Cart Additions and Orders contain whole numbers (aka integral = integers) while columns In Stock Views, Conversion rate and Revenue contain fractionals (in this case 64-bit floating point data types).

### Interesting trivia - Nonintegral Numeric Types:

Nonintegral data types are those that represent numbers with both integer and fractional parts. The nonintegral numeric data types are Decimal (128-bit fixed point), Single Data Type (32-bit floating point), and Double Data Type (64-bit floating point). They are all signed types. If a variable can contain a fraction, declare it as one of these types. Decimal is not a floating-point data type. Decimal numbers have a binary integer value and an integer scaling factor that specifies what portion of the value is a decimal fraction. You can use Decimal variables for money values. The advantage is the precision of the values. The Double data type is faster and requires less memory, but it is subject to rounding errors. The Decimal data type retains complete accuracy to 28 decimal places. Floating-point (Single and Double) numbers have larger ranges than Decimal numbers but can be subject to rounding errors. Floating-point types support fewer significant digits than Decimal but can represent values of greater magnitude.

## Cleaning

Let's check for unique values now in columns with categorical record - to make sure there are redundancies in naming (which could negatively impact our model):

```
In [7]: 1 best_sell["Category"].unique()
```

```
Out[7]: array(['III&T', 'TCFM', 'IA&C', 'BLE', 'OKdo'], dtype=object)
```

```
In [8]: 1 best_sell["Cell"].unique()
```

```
Out[8]: array(['Cable', 'Anti-Vib & Ac', 'Consumables & Safety',
   'Tools Storage & FM', 'I-Conn', 'Sensors & Switches',
   'Process Control', 'Mach Guard & Enc', 'E-Conn',
   'Test & Measurement', 'Passives', 'EMECH', 'Mech & Flu Pwr',
   'BLE E-Mech', 'Opto', 'Analog Disc Power', 'SBC & IOT'],
  dtype=object)
```

```
In [9]: 1 best_sell["Technology"].unique()
```

```
Out[9]: array(['Cable Accs', 'Battery', 'Hand Tools', 'Adhesives, Clean & M',
   'Fasteners', 'Measurement & Insp', 'Crimp Term & Power', 'HVAC',
   'Lighting', 'Proc Auto Sen', 'Power Management', 'Enclosures',
   'Pow Tools & Abras', 'Site Safety & Janit', 'IT', 'Fact Auto Sen',
   'Rect Connectors', 'Electronic T&M', 'Resistors',
   'Push But & Swit', 'Electrical T&M', 'Soldering & ESD',
   'Indicators', 'Acc, Stor & Mat Hand', 'Pow Tran & Struc',
   'Printing & Office', '3D Printing', 'Pwr & Ind Cable',
   'RF Datacomm Conns', 'Motors & Control', 'Ind & AV Conns',
   'Capacitors', 'PPE', 'Plumbing', 'Level & Flow Con',
   'Engineering Mat', 'Cable Mgt & Elec Acc', 'Hydraulics',
   'Thermal Management', 'Hose & Pumps', 'Pneumatics',
   'Elec & Data Cable', 'Security', 'Ctrl & Ind Relay',
   "DIN Rail PSU's", 'Timers, Counters', 'BLE E-Mech Tech',
   'Temp Control', "PLC's & HMI's", 'Inductors', 'Switches-EMECH',
   'Fuses', 'DIN Rail & Term', 'PCB Connectors', 'Circuit Protection',
   'MIL Ind & Test Conns', 'Sounders/Beacons', 'Anti-Vibration',
   'Relays-EMECH', 'Contactors', 'Heavy Duty Conn', 'Machine Safety',
   'LED', 'Power PCB', 'Displays', 'SBC IOT Dev. Tools', 'Sensors'],
  dtype=object)
```

I am purposely skipping Article column as I expect there to be 1000's of unique product names.

We are lucky - there doesn't seem to be a need for further data cleaning at this point. If that need arises, I'll do it on the fly.

## Investigating

Let's see now what are the most commonly ordered products:

```
In [10]: 1 best_sell.sort_values(by=['Orders'], ascending=False).head(1)
```

Out[10]:

	Article	Description	FSTD Global	Category	Cell	Technology	Visits	Product Views	In Stock Views	A
14	8660672	LITHIUM COIN CELL CR2032 5PK	9999-12-31 00:00:00	TCFM	Consumables & Safety	Battery	45475	28009	0.977	

Oh, the revenue... I'm going to hide it, as not to reveal too much at this stage:

```
In [11]: 1 best_sellhid=best_sell.drop(['Revenue'], axis=1)
```

```
In [12]: 1 best_sellhid.head(7)
```

Out[12]:

	Article	Description	FSTD Global	Category	Cell	Technology	Visits	Product Views	I Stoc View
0	233487	BLACK CABLE TIE, 300X4.8MM, PACK 100	9999-12-31 00:00:00	II&T	Cable, Anti-Vib & Ac	Cable Accs	43371	28655	0.99
1	5375488	RS SEALED LEAD-ACID BATTERY,12V 7AH	9999-12-31 00:00:00	TCFM	Consumables & Safety	Battery	40171	41566	0.99
2	7348885	88PC 1/2IN. SOCKET & VDE TOOL KIT	9999-12-31 00:00:00	TCFM	Tools Storage & FM	Hand Tools	79286	89157	0.97
3	1681644	DUST REMOVER CLEANER 400ML	9999-12-31 00:00:00	TCFM	Consumables & Safety	Adhesives, Clean & M	8704	8340	0.93
4	6191506	SEALING STRIP,PVC,WIRE INSERT,1-2,9X6.	9999-12-31 00:00:00	TCFM	Consumables & Safety	Fasteners	9605	9731	0.99
5	8412518	ELECTRONIC CALIPER 150MM/6"	9999-12-31 00:00:00	TCFM	Tools Storage & FM	Measurement & Insp	35240	34067	0.99
6	100282	BRASS SLOT'D SCREW,NUT & WASHER KIT	9999-12-31 00:00:00	TCFM	Consumables & Safety	Fasteners	1609	1760	0.99

ok, now we can safely go back to seeing most frequently ordered products:

In [13]: 1 best\_sellhid.sort\_values(by=[ 'Orders' ], ascending=False).head(10)

Out[13]:

		Article	Description	FSTD Global	Category	Cell	Technology	Visits	Product Views
14	8660672	RS LITHIUM COIN CELL CR2032 5PK		9999-12-31 00:00:00	TCFM	Consumables & Safety	Battery	45475	28009
89	7757233	LR44 ALKALINE COIN CELL 1.5V 158MAH		9999-12-31 00:00:00	TCFM	Consumables & Safety	Battery	51456	37928
0	233487	BLACK CABLE TIE, 300X4.8MM, PACK 100		9999-12-31 00:00:00	II&T	Cable, Anti-Vib & Ac	Cable Accs	43371	28655
29	512238	PTFE THREAD SEAL TAPE,12M L X 12MM W		9999-12-31 00:00:00	TCFM	Consumables & Safety	Adhesives, Clean & M	36293	26000
157	233455	CABLE TIE,100 X 2.5,BLACK,PACK 100		9999-12-31 00:00:00	II&T	Cable, Anti-Vib & Ac	Cable Accs	27824	17143
72	458702	RED INSUL BOOTLACE FERRULE,8MMPIN 1MMSQ.		9999-12-31 00:00:00	II&T	I-Conn	Crimp Term & Power	29405	19737
33	7442199	NON-RECHARGEABLE AA ALKALINE BATTERY		9999-12-31 00:00:00	TCFM	Consumables & Safety	Battery	22215	13742
36	458718	BLK INSUL BOOTLACE FERRULE,1.5MMSQ. WIRE		9999-12-31 00:00:00	II&T	I-Conn	Crimp Term & Power	29259	19750
26	593423	RS SR44 COIN CELL BATTERY, 1.55V		9999-12-31 00:00:00	TCFM	Consumables & Safety	Battery	23043	17434
138	467406	STANDARD TOP HAT PUNCHED DIN RAIL,0.5M		9999-12-31 00:00:00	IA&C	Mach Guard & Enc	Enclosures	32740	25221

It looks like batteries and cables are most popular among our customers.

Now let's find out which products got most eyeballs (but not necessarily most orders):

In [14]: 1 best\_sellhid.sort\_values(by=[ 'Product Views' ], ascending=False).head()

Out[14]:

		Article	Description	FSTD Global	Category	Cell	Technology	Visits	Product Views	In Stock Views
2	7348885	88PC 1/2IN. SOCKET & VDE TOOL KIT	9999-12-31 00:00:00	TCFM	Tools Storage & FM	Hand Tools	79286	89157	0.970	
1293	8179236	55PC BIT AND SOCKET SET 1/4" DRIVE	9999-12-31 00:00:00	TCFM	Tools Storage & FM	Hand Tools	55287	58233	0.730	
1	5375488	RS SEALED LEAD-ACID BATTERY,12V 7AH	9999-12-31 00:00:00	TCFM	Consumables & Safety	Battery	40171	41566	0.992	
89	7757233	LR44 ALKALINE COIN CELL 1.5V 158MAH	9999-12-31 00:00:00	TCFM	Consumables & Safety	Battery	51456	37928	0.997	
397	7613319	PROFILE 8 40X40 LIGHT 1M	9999-12-31 00:00:00	IA&C	Mech & Flu Pwr	Pow Tran & Struc	28691	37330	0.995	

Pandas is telling us that it wasn't batteries or cables but, actually, sockets. Interesting. It would appear that conversion rate of sockets is noticeably lower than that of the 2 former groups. Perhaps it's driven by price. Either way - interesting fact...

Finally, let's learn which product was **least** often in stock, when viewed:

In [15]: 1 best\_sellhid.sort\_values(by=['In Stock Views']).head()

Out[15]:

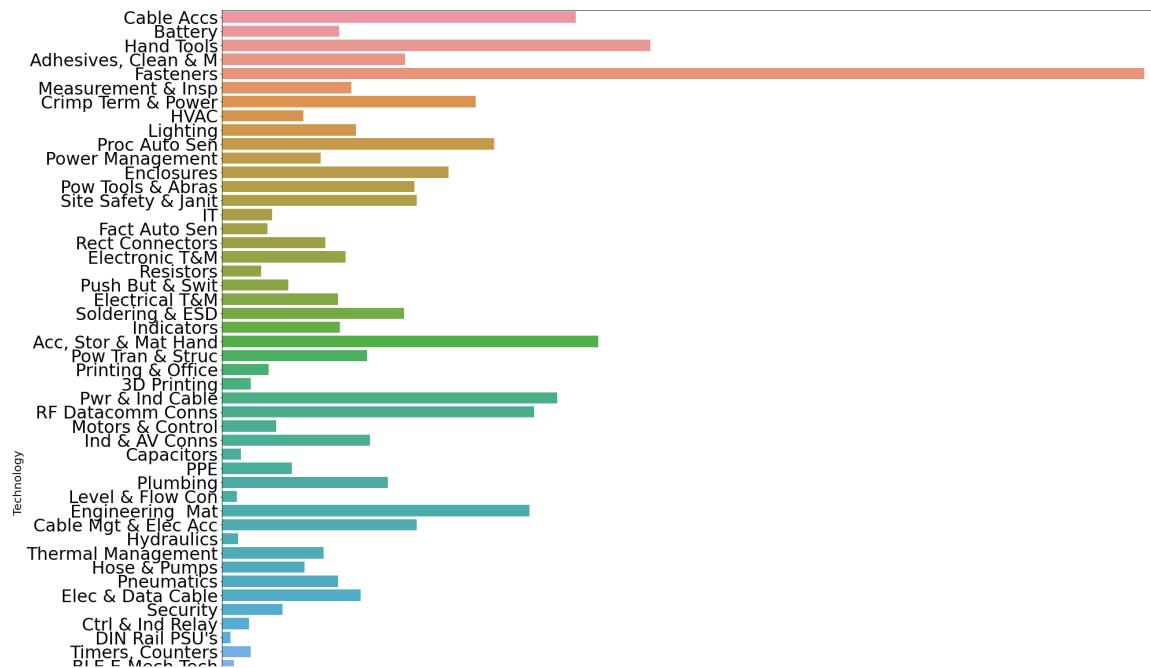
		Article	Description	FSTD Global	Category	Cell	Technology	Visits	Product Views	In Stock	Views	Add
1622	6693677		THREE STEP MOBILE PLATFORM	9999-12-31 00:00:00	TCFM	Tools Storage & FM	Acc, Stor & Mat Hand	0	0	0	0.000	
8995	1777087		OIL FILLED RADIATOR 2KW 9FIN UK PLUG	9999-12-31 00:00:00	TCFM	Tools Storage & FM	HVAC	8	0	0	0.000	
6164	9076422		5/2 PILOT OPERATED SOLENOID VALVE, G1/2"	9999-12-31 00:00:00	IA&C	Mech & Flu Pwr	Pneumatics	190	222	0.009		
10254	8938534		FUSE KIT 6.3X32 FAST AND SLOW GLASS 250V	9999-12-31 00:00:00	II&T	EMECH	Fuses	1921	1922	0.051		
10273	7761989		ORANGE 3 CORE ARCTIC CABLE 2.5MM 100M	2024-12-19 00:00:00	II&T	Cable, Anti-Vib & Ac	Pwr & Ind Cable	520	542	0.174		

I won't pretend to be able to explain what's so special about solenoid valve but it looks like it was viewed 222 times and was only in stock less than %1 of those times. And fuse kit was clicked on to view 1921 times and yet it was available for purchase only 5.1% of the time...

## Visualizing the data

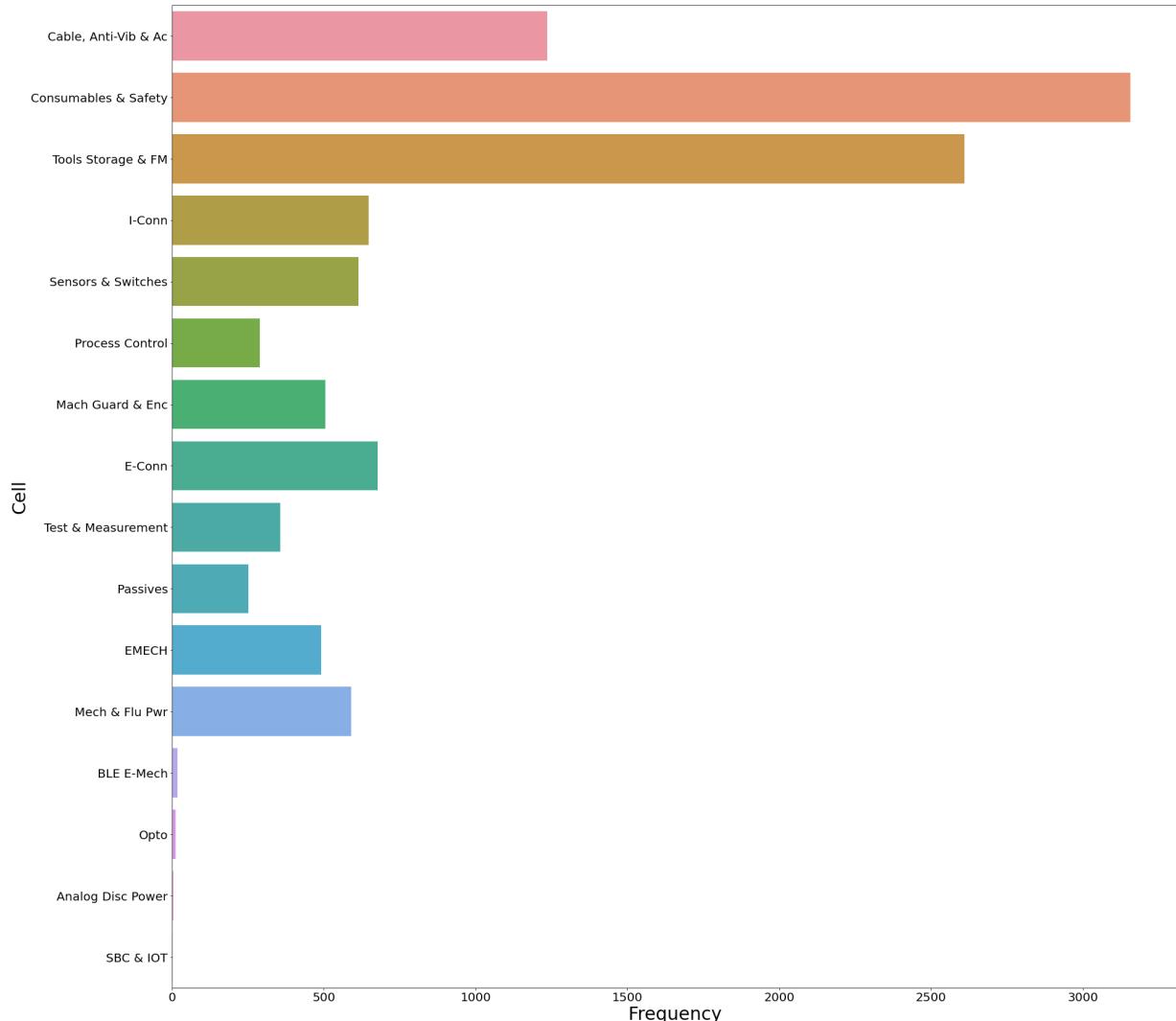
In [16]:

```
1 # bar-plotting distribution of Technology
2 fig1 = sns.countplot(y = 'Technology', data = best_sellhid)
3
4 fig1.figure.set_size_inches(30,30)
5
6
7 fig1.set_ylabel("Technology", fontsize=20)
8 fig1.set_xlabel("Frequency", fontsize=30)
9 fig1.tick_params(labelsize = 30)
10
11 plt.show()
```



In [17]:

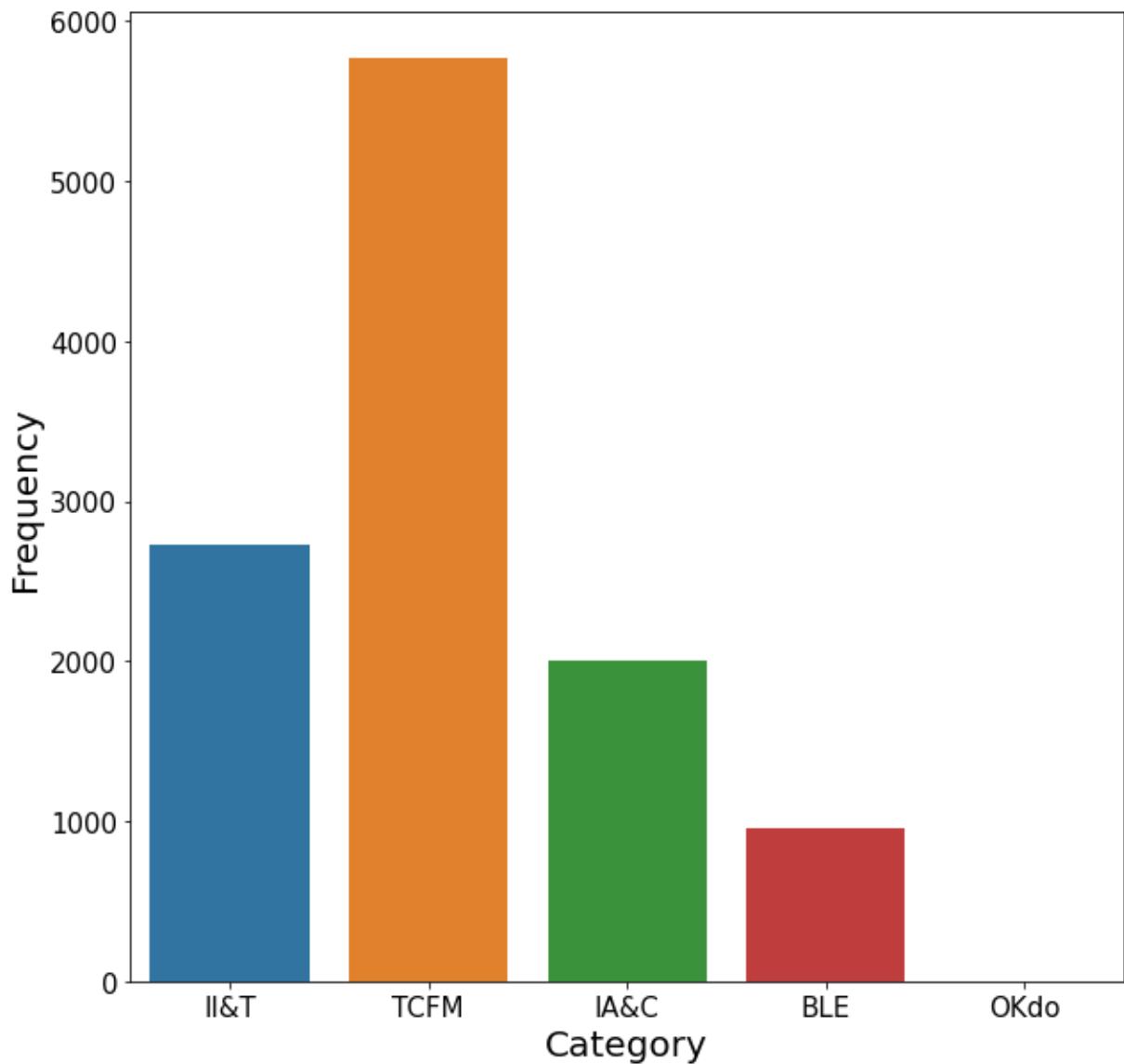
```
1 # bar-plotting distribution of Technology
2 fig1 = sns.countplot(y = 'Cell', data = best_sellhid)
3
4 fig1.figure.set_size_inches(30,30)
5
6
7 fig1.set_ylabel("Cell", fontsize=30)
8 fig1.set_xlabel("Frequency", fontsize=30)
9 fig1.tick_params(labelsize = 20)
10
11 plt.show()
```





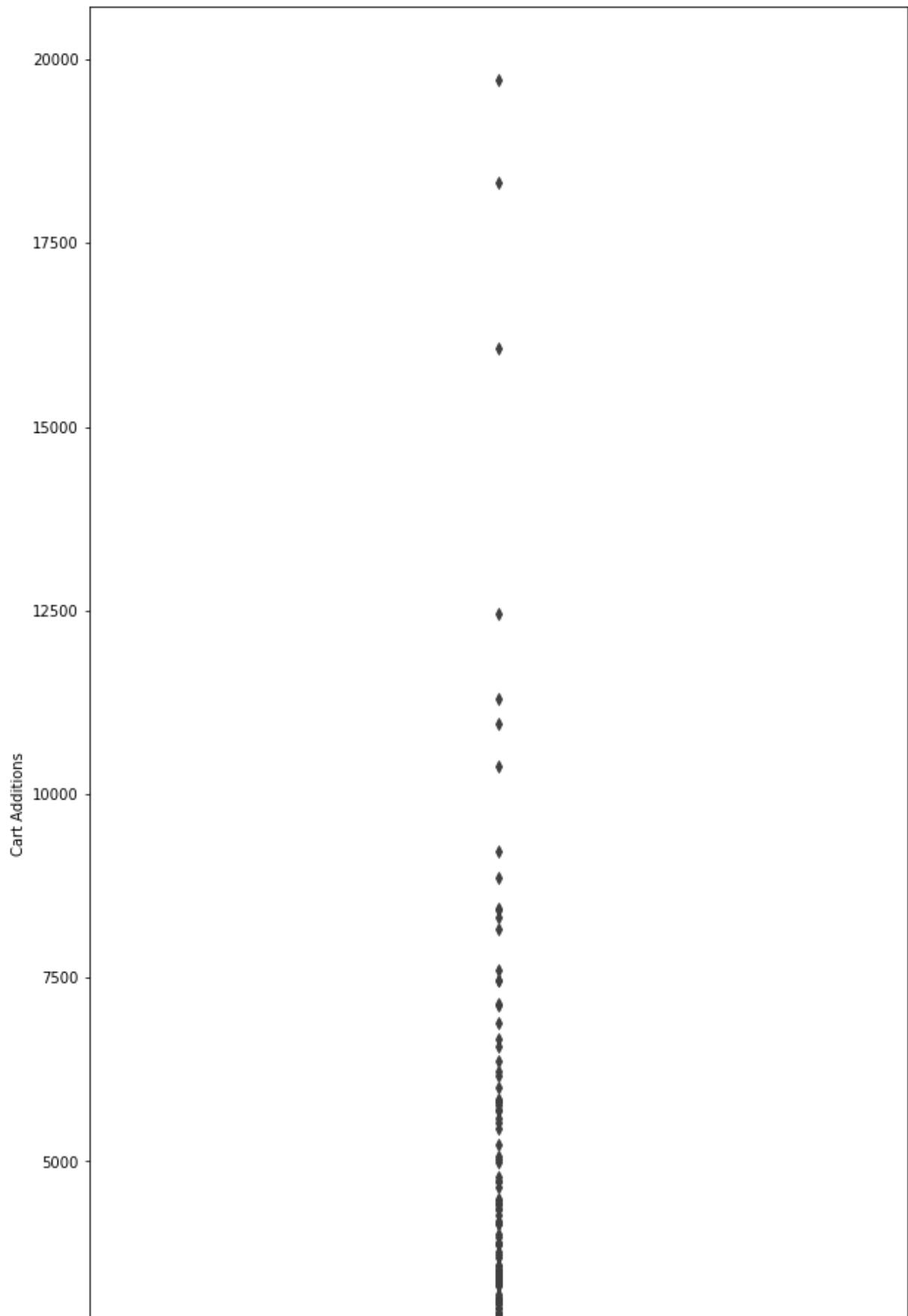
In [18]:

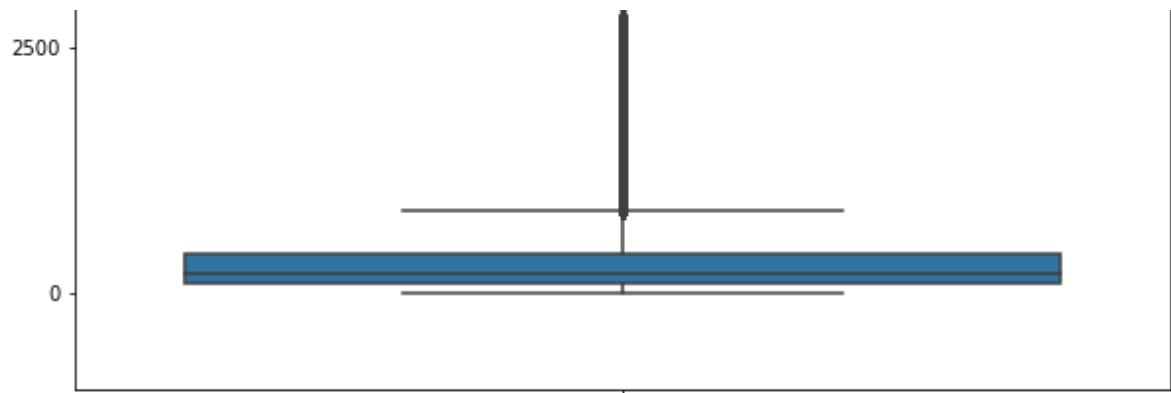
```
1 # bar-plotting distribution of Technology
2 fig1 = sns.countplot(x = 'Category', data = best_sellhid)
3
4 fig1.figure.set_size_inches(10,10)
5
6
7 fig1.set_xlabel("Category", fontsize=20)
8 fig1.set_ylabel("Frequency", fontsize=20)
9 fig1.tick_params(labelsize = 15)
10
11 plt.show()
```



In [19]:

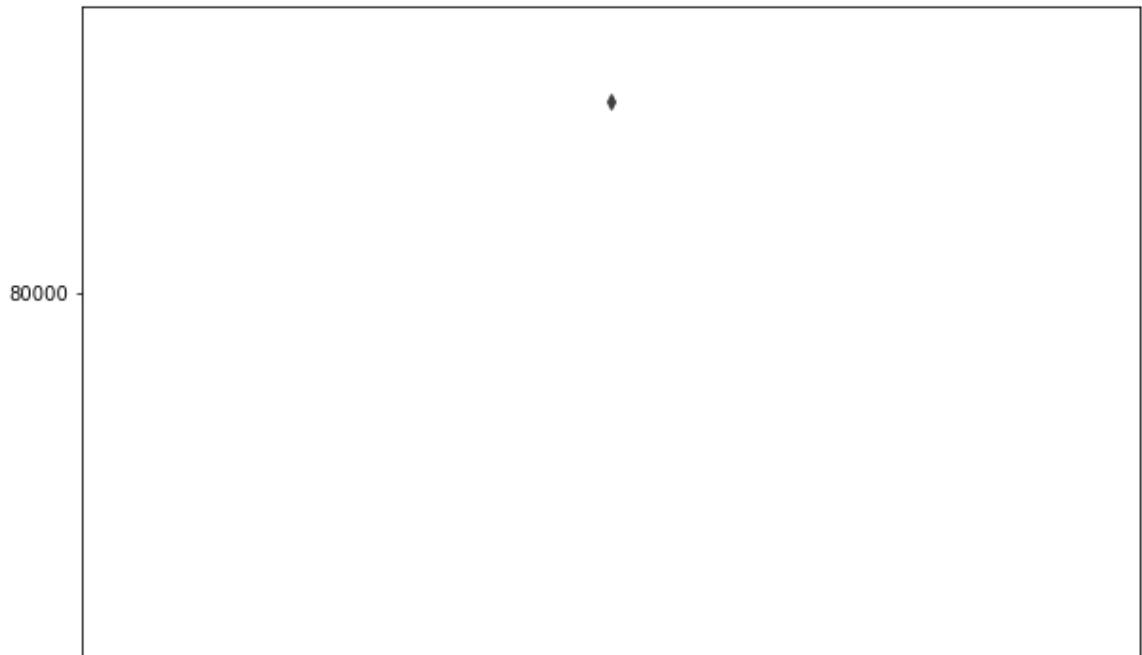
```
1 # boxplotting to check the variation of Cart Additions
2 fig2 = sns.boxplot(y = 'Cart Additions', data = best_sellhid)
3 fig2.figure.set_size_inches(10,20)
4 plt.show()
```





In [20]:

```
1 # boxplotting to check the variation of Product Views
2 fig3 = sns.boxplot(y = 'Product Views', data = best_sellhid)
3 fig3.figure.set_size_inches(10,20)
4 plt.show()
```



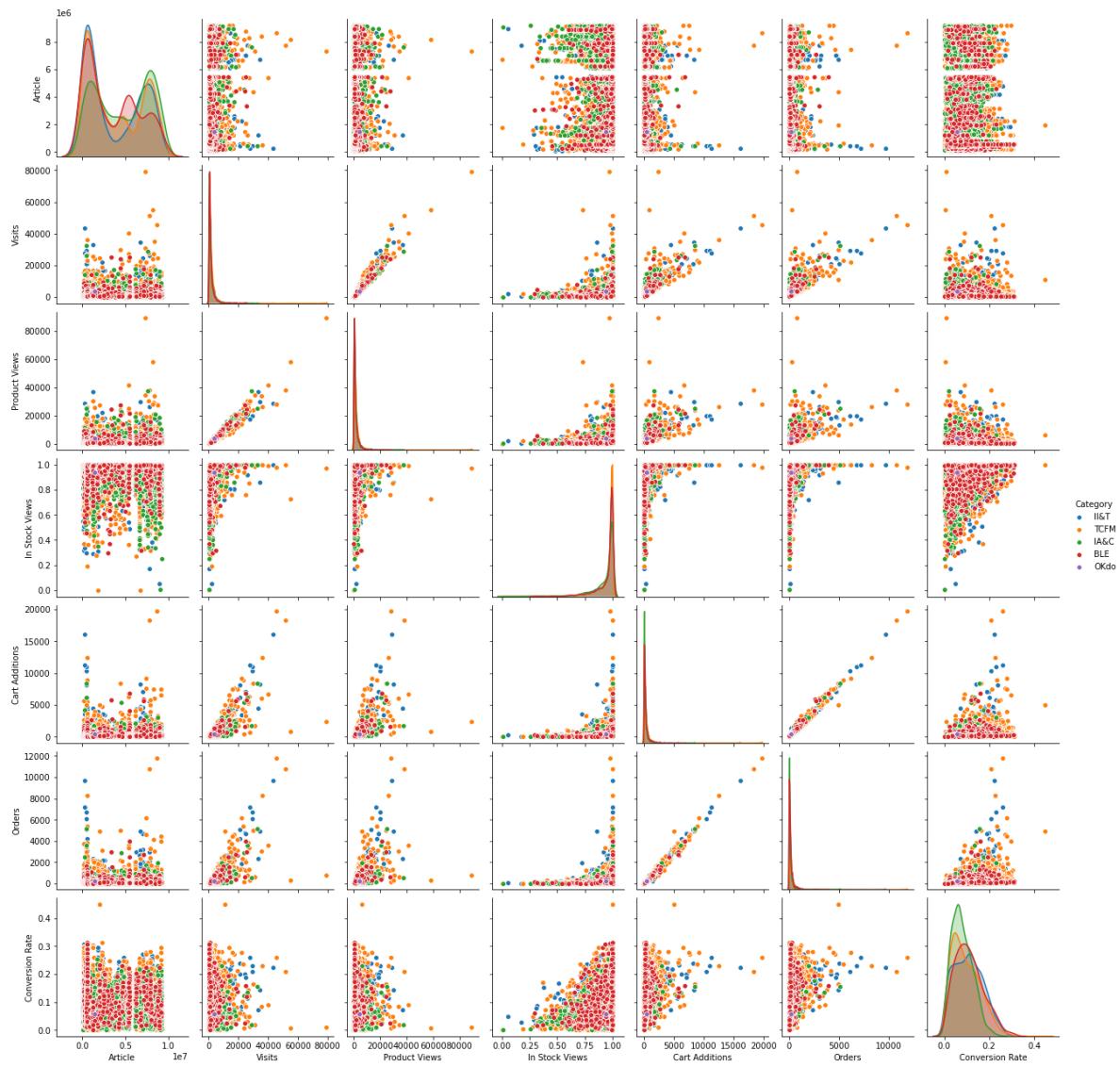
Let's use seaborn's pairplot to see if anything pops up, looking at Category feature:

In [21]:

```
1 pair = sns.pairplot(best_sellhid, hue="Category")
2 pair
3 # Run this cell twice if you have a warning message to see the plot
```

```
C:\Users\oszer\anaconda3\lib\site-packages\seaborn\distributions.py:288: UserWa
rning: Data must have variance to compute a kernel density estimate.
    warnings.warn(msg, UserWarning)
C:\Users\oszer\anaconda3\lib\site-packages\seaborn\distributions.py:288: UserWa
rning: Data must have variance to compute a kernel density estimate.
    warnings.warn(msg, UserWarning)
C:\Users\oszer\anaconda3\lib\site-packages\seaborn\distributions.py:288: UserWa
rning: Data must have variance to compute a kernel density estimate.
    warnings.warn(msg, UserWarning)
C:\Users\oszer\anaconda3\lib\site-packages\seaborn\distributions.py:288: UserWa
rning: Data must have variance to compute a kernel density estimate.
    warnings.warn(msg, UserWarning)
C:\Users\oszer\anaconda3\lib\site-packages\seaborn\distributions.py:288: UserWa
rning: Data must have variance to compute a kernel density estimate.
    warnings.warn(msg, UserWarning)
C:\Users\oszer\anaconda3\lib\site-packages\seaborn\distributions.py:288: UserWa
rning: Data must have variance to compute a kernel density estimate.
    warnings.warn(msg, UserWarning)
C:\Users\oszer\anaconda3\lib\site-packages\seaborn\distributions.py:288: UserWa
rning: Data must have variance to compute a kernel density estimate.
    warnings.warn(msg, UserWarning)
```

Out[21]: &lt;seaborn.axisgrid.PairGrid at 0x1d6d924a088&gt;



That's too noisy to discern much. But looking at OkDo category there seemed to have been hardly any. let's confirm it:

```
In [22]: 1 okdo = best_sell["Category"].value_counts()
2 okdo
```

```
Out[22]: TCFM      5768
II&T      2730
IA&C      1998
BLE       961
Okdo      1
Name: Category, dtype: int64
```

Indeed there is only 1 product in OkDo category. It's safe to drop it and see if that helps...

```
In [23]: 1 #dropping OKdo from Category column:
2 best_sellhidwithoutOk = best_sellhid[best_sellhid['Category'] != 'OKdo']
```

```
In [24]: 1 best_sellhidwithoutOk["Category"].value_counts()
```

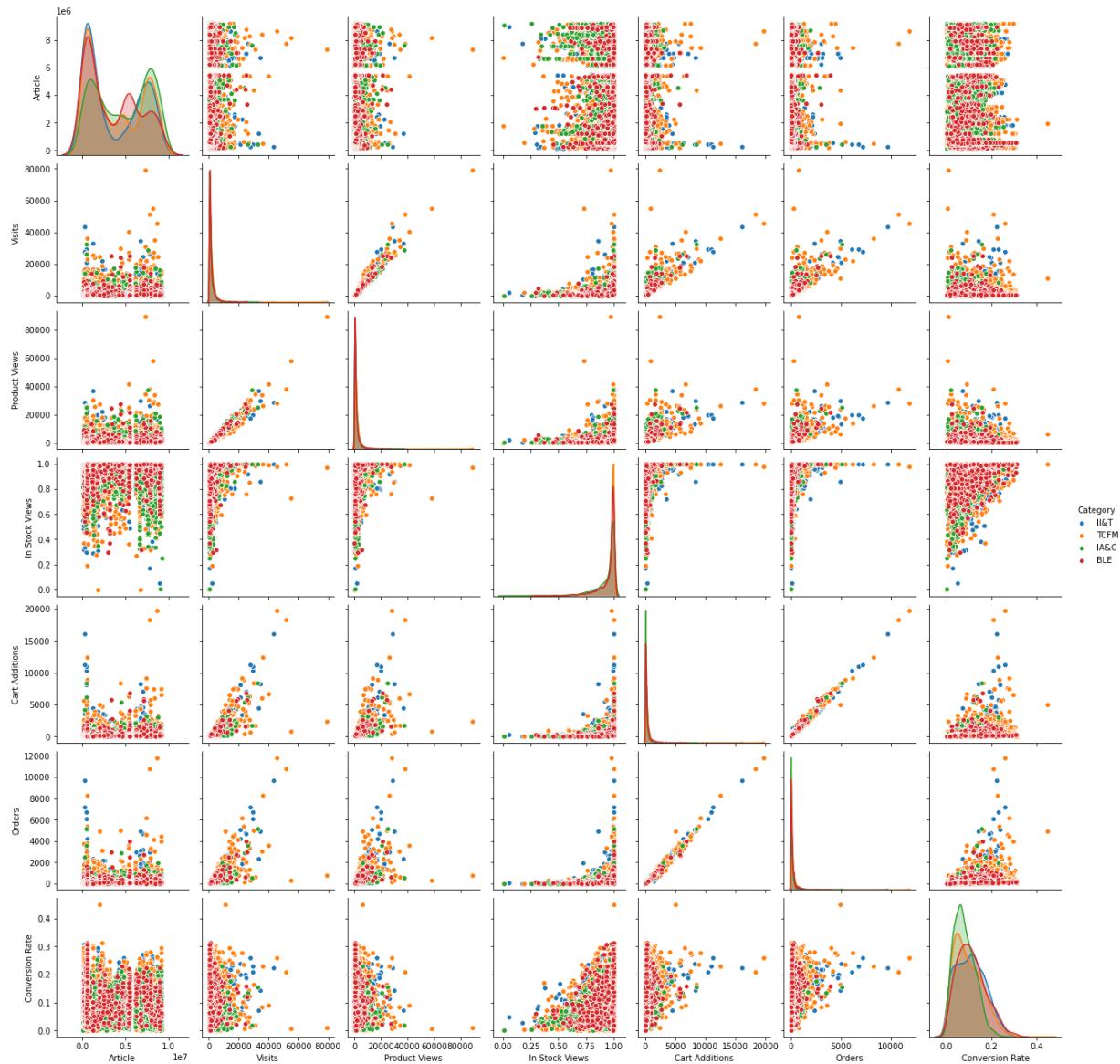
```
Out[24]: TCFM      5768
II&T     2730
IA&C     1998
BLE       961
Name: Category, dtype: int64
```

```
In [25]: 1 best_sellhidwithoutOk.shape
```

```
Out[25]: (11457, 12)
```

```
In [26]: 1 #Trying this again - this time with less categories (without Okdo):
2 pair2 = sns.pairplot(best_sellhidwithoutOk, hue="Category")
3 pair2
4 # Run this cell twice if you have a warning message to see the plot
```

```
Out[26]: <seaborn.axisgrid.PairGrid at 0x1d6db6a9ec8>
```



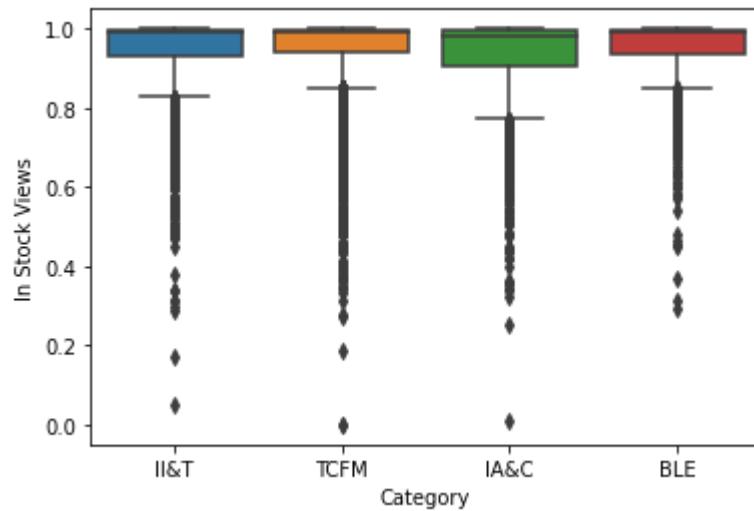
It's still too much overlap but we can't drop any more categories - they have too many records.

In [27]:

```

1 # box plot of In Stock Views based on Category:
2 sns.boxplot(y = 'In Stock Views', x = 'Category', data = best_sellhidwithout
3 plt.show()

```



It appears we stock BLE best as observations don't reach the bottom of the graph.

In [28]:

```

1 # Assigning dataframe new name - shorter:
2 data = best_sellhidwithoutOk

```

In [29]:

```
1 data.sample(3)
```

Out[29]:

		Article	Description	FSTD Global	Category	Cell	Technology	Visits	Product Views	In Stock Views
5420	8117689		PE SPIRAL WRAP 15MM BLACK 10M	9999-12-31 00:00:00	II&T	Cable, Anti-Vib & Ac	Cable Accs	3286	3022	0.813
8407	7887144		TYPE K IMMERSION PROBE	9999-12-31 00:00:00	II&T	Test & Measurement	Electrical T&M	1502	1733	0.989
4526	2830746		5 WAY DOUBLE SCREW EARTH TERMINAL BLOCK	9999-12-31 00:00:00	TCFM	Tools Storage & FM	Cable Mgt & Elec Acc	1062	1042	0.968



## Preparing dataframe for modelling:

1) Since we don't have any binary type data, we will need to create new column(s) with dychotomy split. It could sense to look at key statistical measures and use one of them as a cut-off value to cut dataset in 2...

```
In [30]: 1 data['Product Views'].describe()
```

```
Out[30]: count    11457.000000
          mean     1946.953915
          std      2752.089032
          min      0.000000
          25%     632.000000
          50%     1153.000000
          75%     2245.000000
          max     89157.000000
          Name: Product Views, dtype: float64
```

We can see that mean value of the distribution is 1946 (that's how many times a product among Top Sellers was viewed) so it could be an option to consider...

```
In [31]: 1 data['Cart Additions'].describe()
```

```
Out[31]: count    11457.000000
          mean     361.568561
          std      657.753388
          min      0.000000
          25%     95.000000
          50%     195.000000
          75%     390.000000
          max     19733.000000
          Name: Cart Additions, dtype: float64
```

For Cart Additions the average value is 361.

**Important:** We can see (from "min" value) in both columns examined above that value of 0 (zero) makes an appearance. It's only few instances but it will cause problems down the line (in fact it did - producing infinity error in my first attempts) so they have to go.

```
In [32]: 1 #Filtering out zero value sfrom columns Product Views and Cart Additions:
2 nonzero = data[~((data['Product Views'] == 0) | (data['Cart Additions'] == 0))]
```

In [33]: 1 nonzero.head()

Out[33]:

	Article	Description	FSTD Global	Category	Cell	Technology	Visits	Product Views	In Stock Views
0	233487	BLACK CABLE TIE, 300X4.8MM, PACK 100	9999-12-31 00:00:00	II&T	Cable, Anti-Vib & Ac	Cable Accs	43371	28655	0.996
1	5375488	RS SEALED LEAD-ACID BATTERY,12V 7AH	9999-12-31 00:00:00	TCFM	Consumables & Safety	Battery	40171	41566	0.992
2	7348885	88PC 1/2IN. SOCKET & VDE TOOL KIT	9999-12-31 00:00:00	TCFM	Tools Storage & FM	Hand Tools	79286	89157	0.970
3	1681644	DUST REMOVER CLEANER 400ML	9999-12-31 00:00:00	TCFM	Consumables & Safety	Adhesives, Clean & M	8704	8340	0.933
4	6191506	SEALING STRIP,PVC,WIRE INSERT,1-2,9X6.	9999-12-31 00:00:00	TCFM	Consumables & Safety	Fasteners	9605	9731	0.997



In [34]: 1 #checking the shape of thus filtered out dataframe:  
2 nonzero.shape

Out[34]: (11454, 12)

We can see - comparing to result of .shape function used earlier (11457,12) that we've only dropped 3 additional rows - minuscule amount of data.

In [35]:

```
1 #Double-checking if zeros are truly gone, by sorting Product Views in ascending order
2 nonzero.sort_values(by=['Product Views']).head()
```

Out[35]:

		Article	Description	FSTD Global	Category	Cell	Technology	Visits	Product Views	I	Stock View
7440	521358	INDUSTRIAL LOCKER 1 DOOR RED 450MM	U/BRIGHT RED LED BLK CHR RECESSED,12VDC	9999-12-31 00:00:00	TCFM	Tools Storage & FM	Acc, Stor & Mat Hand	40	40	0.97	
4353	212547	CLAD EXTENSION BAY,500MM 24KG	U/BRIGHT GRN LED BLK CHR RECESSED,24VDC	9999-12-31 00:00:00	II&T	EMECH	Indicators	45	47	1.00	
7386	519469	110V 40W/M CONSTANT WATT HEATER TAPE-50M	110V 40W/M CONSTANT WATT HEATER TAPE-50M	9999-12-31 00:00:00	TCFM	Tools Storage & FM	Acc, Stor & Mat Hand	48	48	0.78	
4278	212171	110V 40W/M CONSTANT WATT HEATER TAPE-50M	110V 40W/M CONSTANT WATT HEATER TAPE-50M	9999-12-31 00:00:00	II&T	EMECH	Indicators	77	58	0.93	
10605	7033123	110V 40W/M CONSTANT WATT HEATER TAPE-50M	110V 40W/M CONSTANT WATT HEATER TAPE-50M	9999-12-31 00:00:00	II&T	Cable, Anti-Vib & Ac	Pwr & Ind Cable	62	61	0.95	

2)...I came up with a different approach, though. Why not compare 2 variables as a ratio and then split data in 2 according to the mean value of such ratio:

In [36]:

```
1 nonzero['CartAdd_to_ProductView'] = nonzero['Cart Additions']/nonzero['Produ
```

C:\Users\oszer\anaconda3\lib\site-packages\ipykernel\_launcher.py:1: SettingWithCopyWarning:

A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy) ([https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy))

"""Entry point for launching an IPython kernel.

This is SettingWithCopyWarning - basically pandas advising me I'm working on a copy of a dataframe and not original one. In this case it's fine - that's what I want to do.

In [37]: 1 nonzero.head()

Out[37]:

	Article	Description	FSTD Global	Category	Cell	Technology	Visits	Product Views	In Stock Views
0	233487	BLACK CABLE TIE, 300X4.8MM, PACK 100	9999-12-31 00:00:00	II&T	Cable, Anti-Vib & Ac	Cable Accs	43371	28655	0.996
1	5375488	RS SEALED LEAD-ACID BATTERY,12V 7AH	9999-12-31 00:00:00	TCFM	Consumables & Safety	Battery	40171	41566	0.992
2	7348885	88PC 1/2IN. SOCKET & VDE TOOL KIT	9999-12-31 00:00:00	TCFM	Tools Storage & FM	Hand Tools	79286	89157	0.970
3	1681644	DUST REMOVER CLEANER 400ML	9999-12-31 00:00:00	TCFM	Consumables & Safety	Adhesives, Clean & M	8704	8340	0.933
4	6191506	SEALING STRIP,PVC,WIRE INSERT,1,2,9X6.	9999-12-31 00:00:00	TCFM	Consumables & Safety	Fasteners	9605	9731	0.997



Now we can calculate the average (mean) of the newly created column:

In [38]: 1 # mean calculation of CartAdd\_to\_ProductView:  
2 nonzero['CartAdd\_to\_ProductView'].mean()

Out[38]: 0.20035606227139058

In [39]: 1 # Checking stats of newly created ratio:  
2 nonzero['CartAdd\_to\_ProductView'].describe()

Out[39]: count 11454.000000  
mean 0.200356  
std 0.130551  
min 0.011609  
25% 0.102109  
50% 0.171446  
75% 0.267959  
max 0.986318  
Name: CartAdd\_to\_ProductView, dtype: float64

We can see that the average ratio of Cart Addition(s) to Product View(s) is around **20% - or 1:5**, which intuitively feels right.

I'm curious, though to find out what that one product with super high ratio of over 98% is... Let's find out:

```
In [40]: 1 max = nonzero[nonzero['CartAdd_to_ProductView'] > 0.9863]
```

```
In [41]: 1 max.head()
```

Out[41]:

	Article	Description	FSTD Global	Category	Cell	Technology	Visits	Product Views	In Stock Views
272	9176588	RS ALKALINE AA 20 PK	9999-12-31 00:00:00	TCFM	Consumables & Safety	Battery	13527	6651	0.995

it's the batteries. Looks like clients really like them and once they view them, they virtually always add them to their carts.

I can now create a new column, where there will be indication if our CartAdd\_to\_ProductView ratio is above or below the average.

```
In [42]: 1 #Creating and assigning new variable for our average:  
2 a_mean = nonzero['CartAdd_to_ProductView'].mean()  
3 a_mean
```

Out[42]: 0.20035606227139058

```
In [43]: 1 #Using Lambda function for comparison to mean and creating new column with b  
2 nonzero['aboveMean'] = nonzero['CartAdd_to_ProductView'].apply(lambda x: x >
```

C:\Users\oszer\anaconda3\lib\site-packages\ipykernel\_launcher.py:2: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame.  
Try using .loc[row\_indexer,col\_indexer] = value instead

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy) ([https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy))

In [44]: 1 nonzero.head()

Out[44]:

	Article	Description	FSTD Global	Category	Cell	Technology	Visits	Product Views	In Stock Views
0	233487	BLACK CABLE TIE, 300X4.8MM, PACK 100	9999-12-31 00:00:00	II&T	Cable, Anti-Vib & Ac	Cable Accs	43371	28655	0.996
1	5375488	RS SEALED LEAD-ACID BATTERY,12V 7AH	9999-12-31 00:00:00	TCFM	Consumables & Safety	Battery	40171	41566	0.992
2	7348885	88PC 1/2IN. SOCKET & VDE TOOL KIT	9999-12-31 00:00:00	TCFM	Tools Storage & FM	Hand Tools	79286	89157	0.970
3	1681644	DUST REMOVER CLEANER 400ML	9999-12-31 00:00:00	TCFM	Consumables & Safety	Adhesives, Clean & M	8704	8340	0.933
4	6191506	SEALING STRIP,PVC,WIRE INSERT,1-2,9X6.	9999-12-31 00:00:00	TCFM	Consumables & Safety	Fasteners	9605	9731	0.997

We now have a new column with binary values and we can treat it as a class for predictions.

## Correlations:

In [45]: 1 #Assigning new name to the dataframe:  
2 dataset = nonzero

In [46]: 1 dataset.sample(3)

Out[46]:

		Article	Description	FSTD Global	Category	Cell	Technology	Visits	Product Views	In Stock Views	Ad
6042	7840085	G 1/2 FRL PORTABLE AIR CRADLE		9999-12-31 00:00:00	IA&C	Mech & Flu Pwr	Pneumatics	716	807	0.791	
8666	7004523	RED HEATSHRINK TUBE 24/8MM I/D		9999-12-31 00:00:00	II&T	Cable, Anti-Vib & Ac	Cable Accs	934	744	0.999	
9757	2512224	YELLOW SEALED LOW PROFILE FOOT SWITCH,3A		9999-12-31 00:00:00	IA&C	Sensors & Switches	Push But & Swit	302	321	0.994	



We can now check correlations:

In [47]: 1 dataset.corr()

Out[47]:

	Article	Visits	Product Views	In Stock Views	Cart Additions	Orders	Conversion Rat
<b>Article</b>	1.000000	0.005626	0.021034	-0.110331	-0.033665	-0.050572	-0.16737
<b>Visits</b>	0.005626	1.000000	0.977554	0.106543	0.808881	0.772378	0.04915
<b>Product Views</b>	0.021034	0.977554	1.000000	0.092866	0.677227	0.633791	-0.04917
<b>In Stock Views</b>	-0.110331	0.106543	0.092866	1.000000	0.126527	0.135813	0.24137
<b>Cart Additions</b>	-0.033665	0.808881	0.677227	0.126527	1.000000	0.991103	0.33049
<b>Orders</b>	-0.050572	0.772378	0.633791	0.135813	0.991103	1.000000	0.37159
<b>Conversion Rate</b>	-0.167378	0.049156	-0.049179	0.241371	0.330492	0.371596	1.00000
<b>CartAdd_to_ProductView</b>	-0.132046	0.033293	-0.079390	0.169494	0.332880	0.354950	0.91464
<b>aboveMean</b>	-0.101098	0.013253	-0.071804	0.176931	0.248793	0.266624	0.79104



In [48]:

```

1 #Plotting correlations above as a seaborn heatmap:
2 plt.figure(figsize = (20, 10))
3 fig1 = sns.heatmap(dataset.corr(), annot = True)
4 fig1.set_xlabel("Correlations", fontsize=20)
5 plt.show()

```



In [49]:

```

1 #let's remind ourselves of the structure of dataframe:
2 dataset.sample()

```

Out[49]:

	Article	Description	FSTD Global	Category	Cell	Technology	Visits	Product Views	In Stock Views	A
2916	7755362	MILD STEEL IP66 WALL BOX, 800X600X300MM	9999-12-31 00:00:00	IA&C	Mach Guard & Enc	Enclosures	992	1049	0.827	

Since we will need to convert categorical columns' values into numerical, we need to get rid of mostly useless values in 2 columns: Description and FSTD Global. They don't help us much but make next step very time-consuming.

In [50]:

```

1 dataset.drop(['Description', 'FSTD Global'], axis=1, inplace=True)

```

C:\Users\oszer\anaconda3\lib\site-packages\pandas\core\frame.py:3997: SettingWithCopyWarning:  
A value is trying to be set on a copy of a slice from a DataFrame

See the caveats in the documentation: [https://pandas.pydata.org/pandas-docs/stable/user\\_guide/indexing.html#returning-a-view-versus-a-copy](https://pandas.pydata.org/pandas-docs/stable/user_guide/indexing.html#returning-a-view-versus-a-copy)  
errors=errors,

In [51]: 1 dataset.info()

```
<class 'pandas.core.frame.DataFrame'>
Int64Index: 11454 entries, 0 to 11457
Data columns (total 12 columns):
 #   Column           Non-Null Count  Dtype  
--- 
 0   Article          11454 non-null   int64  
 1   Category         11454 non-null   object  
 2   Cell              11454 non-null   object  
 3   Technology       11454 non-null   object  
 4   Visits            11454 non-null   int64  
 5   Product Views    11454 non-null   int64  
 6   In Stock Views   11454 non-null   float64 
 7   Cart Additions   11454 non-null   int64  
 8   Orders            11454 non-null   int64  
 9   Conversion Rate  11454 non-null   float64 
 10  CartAdd_to_ProductView 11454 non-null   float64 
 11  aboveMean        11454 non-null   bool    
dtypes: bool(1), float64(3), int64(5), object(3)
memory usage: 1.1+ MB
```

## Dummy data (variables):

To evaluate the correlation between categorical variables it is necessary to first perform one-hot encoding to create numerical dummy variables - one binary variable for each level of the original categorical variable.

In [52]:

```

1 # Generating new data set with dummy variables:
2 data_dummy = pd.get_dummies(dataset)
3
4 data_dummy.info()

```

```

<class 'pandas.core.frame.DataFrame'>
Int64Index: 11454 entries, 0 to 11457
Data columns (total 94 columns):
 #   Column           Non-Null Count Dtype
 ---  -- 
 0   Article          11454 non-null  int64
 1   Visits           11454 non-null  int64
 2   Product Views    11454 non-null  int64
 3   In Stock Views   11454 non-null  float64
 4   Cart Additions   11454 non-null  int64
 5   Orders           11454 non-null  int64
 6   Conversion Rate  11454 non-null  float64
 7   CartAdd_to_ProductView 11454 non-null  float64
 8   aboveMean        11454 non-null  bool
 9   Category_BLE     11454 non-null  uint8
 10  Category_IA&C   11454 non-null  uint8
 11  Category_II&T   11454 non-null  uint8
 12  Category_TCFM   11454 non-null  uint8
 13  Cell_Analog Disc Power 11454 non-null  uint8
 14  Cell_BLE E-Mech  11454 non-null  uint8
 15  Cell_Cable, Anti-Vib & Ac 11454 non-null  uint8
 16  Cell_Consumables & Safety 11454 non-null  uint8
 17  Cell_E-Conn      11454 non-null  uint8
 18  Cell_EMECH       11454 non-null  uint8
 19  Cell_I-Conn      11454 non-null  uint8
 20  Cell_Mach Guard & Enc 11454 non-null  uint8
 21  Cell_Mech & Flu Pwr  11454 non-null  uint8
 22  Cell_Opto        11454 non-null  uint8
 23  Cell_Passives   11454 non-null  uint8
 24  Cell_Process Control 11454 non-null  uint8
 25  Cell_Sensors & Switches 11454 non-null  uint8
 26  Cell_Test & Measurement 11454 non-null  uint8
 27  Cell_Tools Storage & FM 11454 non-null  uint8
 28  Technology_3D Printing 11454 non-null  uint8
 29  Technology_Acc, Stor & Mat Hand 11454 non-null  uint8
 30  Technology_Adhesives, Clean & M 11454 non-null  uint8
 31  Technology_Anti-Vibration 11454 non-null  uint8
 32  Technology_BLE E-Mech Tech 11454 non-null  uint8
 33  Technology_Battery 11454 non-null  uint8
 34  Technology_Cable Accs 11454 non-null  uint8
 35  Technology_Cable Mgt & Elec Acc 11454 non-null  uint8
 36  Technology_Capacitors 11454 non-null  uint8
 37  Technology_Circuit Protection 11454 non-null  uint8
 38  Technology_Contactors 11454 non-null  uint8
 39  Technology_Crimp Term & Power 11454 non-null  uint8
 40  Technology_Ctrl & Ind Relay 11454 non-null  uint8
 41  Technology_DIN Rail & Term 11454 non-null  uint8
 42  Technology_DIN Rail PSU's 11454 non-null  uint8
 43  Technology_Displays 11454 non-null  uint8
 44  Technology_Elec & Data Cable 11454 non-null  uint8
 45  Technology_Electrical T&M 11454 non-null  uint8
 46  Technology_Electronic T&M 11454 non-null  uint8

```

```

47 Technology_Enclosures           11454 non-null  uint8
48 Technology_Engineering_Mat     11454 non-null  uint8
49 Technology_Fact_Auto_Sen       11454 non-null  uint8
50 Technology_Fasteners          11454 non-null  uint8
51 Technology_Fuses               11454 non-null  uint8
52 Technology_HVAC                11454 non-null  uint8
53 Technology_Hand_Tools          11454 non-null  uint8
54 Technology_Heavy_Duty_Conn    11454 non-null  uint8
55 Technology_Hose_&_Pumps        11454 non-null  uint8
56 Technology_Hydraulics          11454 non-null  uint8
57 Technology_IT                  11454 non-null  uint8
58 Technology_Ind_&_AV_Conns     11454 non-null  uint8
59 Technology_Indicators          11454 non-null  uint8
60 Technology_Inductors          11454 non-null  uint8
61 Technology_LED                 11454 non-null  uint8
62 Technology_Level_&_Flow_Con   11454 non-null  uint8
63 Technology_Lighting            11454 non-null  uint8
64 Technology_MIL_Ind_&_Test_Conns 11454 non-null  uint8
65 Technology_Machine_Safety      11454 non-null  uint8
66 Technology_Measurement_&_Insp  11454 non-null  uint8
67 Technology_Motors_&_Control    11454 non-null  uint8
68 Technology_PCB_Connectors     11454 non-null  uint8
69 Technology_PLC's_&_HMI's       11454 non-null  uint8
70 Technology_PPE                 11454 non-null  uint8
71 Technology_Plumbing            11454 non-null  uint8
72 Technology_Pneumatics          11454 non-null  uint8
73 Technology_Pow_Tools_&_Abras   11454 non-null  uint8
74 Technology_Pow_Tran_&_Struc    11454 non-null  uint8
75 Technology_Power_Management   11454 non-null  uint8
76 Technology_Power_PCB           11454 non-null  uint8
77 Technology_Printing_&_Office   11454 non-null  uint8
78 Technology_Proc_Auto_Sen       11454 non-null  uint8
79 Technology_Push_But_&_Swit     11454 non-null  uint8
80 Technology_Pwr_&_Ind_Cable    11454 non-null  uint8
81 Technology_RF_Datacomm_Conns   11454 non-null  uint8
82 Technology_Rect_Connectors    11454 non-null  uint8
83 Technology_Relays-EMECH       11454 non-null  uint8
84 Technology_Resistors          11454 non-null  uint8
85 Technology_Security            11454 non-null  uint8
86 Technology_Sensors             11454 non-null  uint8
87 Technology_Site_Safety_&_Janit 11454 non-null  uint8
88 Technology_Soldering_&_ESD     11454 non-null  uint8
89 Technology_Sounders/Beacons   11454 non-null  uint8
90 Technology_Switches-EMECH     11454 non-null  uint8
91 Technology_Temp_Control        11454 non-null  uint8
92 Technology_Thermal_Management 11454 non-null  uint8
93 Technology_Timers,_Counters   11454 non-null  uint8
dtypes: bool(1), float64(3), int64(5), uint8(85)
memory usage: 1.7 MB

```

We can see now that one-hot encoding resulted in creation of additional columns.

In [53]: 1 data\_dummy.head()

Out[53]:

	Article	Visits	Product Views	In Stock Views	Cart Additions	Orders	Conversion Rate	CartAdd_to_ProductView	above
0	233487	43371	28655	0.996	16065	9658	0.2227		0.560635
1	5375488	40171	41566	0.992	6666	3611	0.0899		0.160371
2	7348885	79286	89157	0.970	2373	789	0.0100		0.026616
3	1681644	8704	8340	0.933	2245	1517	0.1743		0.269185
4	6191506	9605	9731	0.997	2152	1556	0.1620		0.221149

5 rows × 94 columns



## Correlation matrix

We can now calculate the correlation coefficients again using the function `.corr` as before:

In [54]: 1 cor\_matrix = data\_dummy.corr()

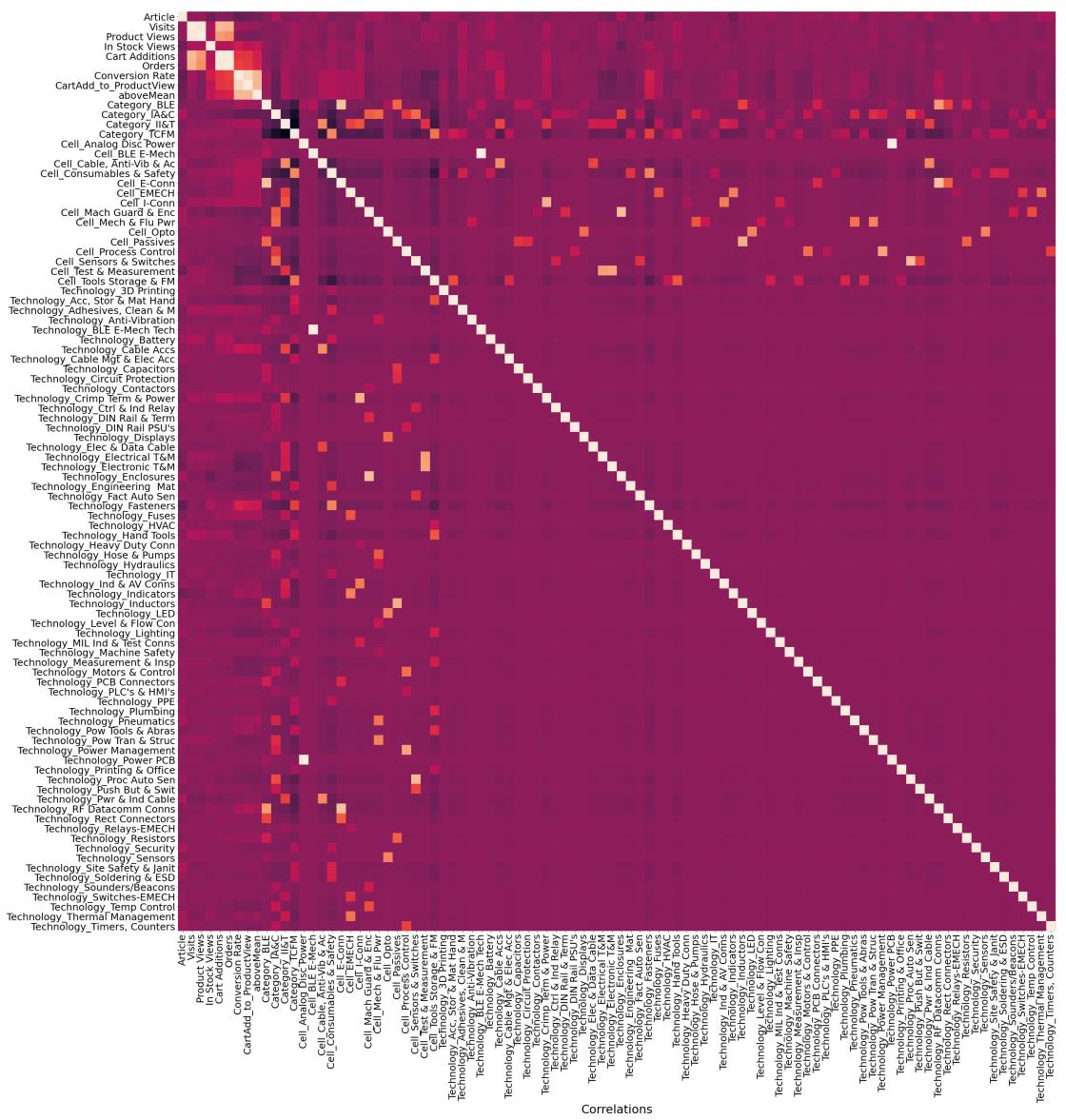
Since it's huge, printing this matrix wouldn't make much sense. Let's try to see the whole picture by displaying it as a heatmap:

In [55]:

```

1 plt.figure(figsize=(50,43))
2 fig1=sns.heatmap(cor_matrix)
3 fig1.set_xlabel("Correlations", fontsize=30)
4 fig1.tick_params(labelsize = 25)

```



Since the scale on the right hand side is difficult to read, it's worth pointing out that the scale extremes are from -1 (light color) to +1 (the darkest coloring).

Along the diagonal there are a couple of small squares of fully white and black cells - this is just due to the dummy variables we produced, where the generated features representing different values (e.g. "above average") will always have perfect anti-correlation. In the rest of the heatmap, there are some non-zero correlations too however. Let's explore these more thoroughly.

## Calculating $r$ -coefficients and $p$ -values

*The important thing to determine for us is, not just what the correlation coefficient is, but whether a correlation is statistically significant or not i.e. not just due to chance.*

We can assess this using the Pearson correlation test. This calculates the p-value, which recall is the probability of observing the current result (in this case the r coefficient) given that there is in fact no significant correlation i.e. the result is due to chance. If that probability is tiny, then the correlation is indeed significant and it is not due to chance. Note that the p-value is not related to the strength of the correlation as measured by the r-coefficient - it only tells us whether a correlation is due to chance or not. Sometimes you can observe r-coefficients very close to zero but with p-values that suggest that the correlation is not due to chance, whereas in other cases r-coefficients can be very large in magnitude (close to -1 or 1), but with non-significant p-values.

Time to import the library `researchpy` which will help a lot in our correlation analysis - in particular, the functions `corr_case` and `corr_pair`:

In [56]:

```
1 import researchpy as rp
2
3 corr_type, corr_matrix, corr_ps = rp.corr_case(data_dummy)
```

The function `.corr_case` generates 3 matrices - one about the type of correlation, the 2nd matrix showing the r-coefficient, and the last showing the p-value of the correlation.

In [57]:

```
1 # Running the last subfunction (i.e. matrix) of .corr_case - corr_ps - showing
2 corr_ps.head(13)
```

Out[57]:

	Article	Visits	Product Views	In Stock Views	Cart Additions	Orders	Conversion Rate	CartAdd_to_ProductView
Article	0.0000	0.5471	0.0244	0.0000	0.0003	0.0000	0.0000	0.0000
Visits	0.5471	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
Product Views	0.0244	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
In Stock Views	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
Cart Additions	0.0003	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
Orders	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
Conversion Rate	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
CartAdd_to_ProductView	0.0000	0.0004	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
aboveMean	0.0000	0.1561	0.0000	0.0000	0.0000	0.0000	0.0000	0.0000
Category_BLE	0.0000	0.0009	0.0014	0.3254	0.1292	0.0820	0.0000	0.0000

In [58]:

```
1 # 2nd matrix of .corr_case - type of correlation:
2 corr_type.head()
```

Out[58]:

Pearson correlation test using list-wise deletion

0	Total observations used = 11454
---	---------------------------------

In [59]:

```
1 # Checking the 2nd matrix that shows the r-coefficient:  
2 corr_matrix.head(15)
```

Out[59]:

	Article	Visits	Product Views	In Stock Views	Cart Additions	Orders	Conversion Rate	CartAd
<b>Article</b>	1	0.0056	0.021	-0.1103	-0.0337	-0.0506	-0.1674	
<b>Visits</b>	0.0056	1	0.9776	0.1065	0.8089	0.7724	0.0492	
<b>Product Views</b>	0.021	0.9776	1	0.0929	0.6772	0.6338	-0.0492	
<b>In Stock Views</b>	-0.1103	0.1065	0.0929	1	0.1265	0.1358	0.2414	
<b>Cart Additions</b>	-0.0337	0.8089	0.6772	0.1265	1	0.9911	0.3305	
<b>Orders</b>	-0.0506	0.7724	0.6338	0.1358	0.9911	1	0.3716	
<b>Conversion Rate</b>	-0.1674	0.0492	-0.0492	0.2414	0.3305	0.3716	1	
<b>CartAdd_to_ProductView</b>	-0.132	0.0333	-0.0794	0.1695	0.3329	0.3549	0.9146	
<b>aboveMean</b>	-0.1011	0.0133	-0.0718	0.1769	0.2488	0.2666	0.791	
<b>Category_BLE</b>	-0.0403	-0.031	-0.0298	0.0092	-0.0142	-0.0162	0.0659	
<b>Category_IA&amp;C</b>	0.1137	-0.0527	-0.0332	-0.0571	-0.0826	-0.0819	-0.132	
<b>Category_II&amp;T</b>	-0.0455	-0.0006	-0.0154	-0.007	0.0409	0.0451	0.09	
<b>Category_TCFM</b>	-0.0252	0.0576	0.0548	0.0442	0.0357	0.0328	-0.013	
<b>Cell_Analog Disc Power</b>	0.0145	0.0162	0.0216	-0.0314	0.0014	-0	-0.0157	
<b>Cell_BLE E-Mech</b>	-0.0052	0.0181	0.0179	-0.0049	0.0228	0.0236	0.0223	

15 rows × 94 columns

This is what researchpy documentation says about .corr\_pair function:

"Conducts Pearson (default method), Spearman rank, or Kendall's Tau-b correlation analysis using pair wise deletion. Returns the relevant information and results in 1 DataFrame for easy exporting".

In [60]:

```
1 corr_table = rp.corr_pair(data_dummy)
```

In [61]:

```
1 # Checking what that correlations' table looks like:
2 corr_table.head(20)
```

Out[61]:

		r value	p-value	N
	<b>Article &amp; Visits</b>	0.0056	0.5471	11454
	<b>Article &amp; Product Views</b>	0.0210	0.0244	11454
	<b>Article &amp; In Stock Views</b>	-0.1103	0.0000	11454
	<b>Article &amp; Cart Additions</b>	-0.0337	0.0003	11454
	<b>Article &amp; Orders</b>	-0.0506	0.0000	11454
	<b>Article &amp; Conversion Rate</b>	-0.1674	0.0000	11454
	<b>Article &amp; CartAdd_to_ProductView</b>	-0.1320	0.0000	11454
	<b>Article &amp; aboveMean</b>	-0.1011	0.0000	11454
	<b>Article &amp; Category_BLE</b>	-0.0403	0.0000	11454
	<b>Article &amp; Category_IA&amp;C</b>	0.1137	0.0000	11454
	<b>Article &amp; Category_II&amp;T</b>	-0.0455	0.0000	11454
	<b>Article &amp; Category_TCFM</b>	-0.0252	0.0069	11454
	<b>Article &amp; Cell_Analog Disc Power</b>	0.0145	0.1212	11454
	<b>Article &amp; Cell_BLE E-Mech</b>	-0.0052	0.5745	11454
	<b>Article &amp; Cell_Cable, Anti-Vib &amp; Ac</b>	0.0435	0.0000	11454
	<b>Article &amp; Cell_Consumables &amp; Safety</b>	-0.0678	0.0000	11454
	<b>Article &amp; Cell_E-Conn</b>	-0.0292	0.0018	11454
	<b>Article &amp; Cell_EMECH</b>	-0.0343	0.0002	11454
	<b>Article &amp; Cell_I-Conn</b>	-0.0292	0.0018	11454
	<b>Article &amp; Cell_Mach Guard &amp; Enc</b>	0.0711	0.0000	11454

In [62]:

```
1 corr_table.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 4371 entries, Article & Visits to Technology_Thermal Management & Technology_Timers, Counters
Data columns (total 3 columns):
 #   Column    Non-Null Count  Dtype  
--- 
 0   r value   4371 non-null   object 
 1   p-value   4371 non-null   object 
 2   N         4371 non-null   int64  
dtypes: int64(1), object(2)
memory usage: 296.6+ KB
```

In [63]:

```
1 # Let's sort now by r value (correlation coefficient):
2 corr_table.sort_values(by=['r value'], ascending=False).head(21)
```

Out[63]:

		r value	p-value	N
	<b>Cell_Analog Disc Power &amp; Technology_Power PCB</b>	1.0000	0.0000	11454
	<b>Cell_BLE E-Mech &amp; Technology_BLE E-Mech Tech</b>	1.0000	0.0000	11454
	<b>Cart Additions &amp; Orders</b>	0.9911	0.0000	11454
	<b>Visits &amp; Product Views</b>	0.9776	0.0000	11454
	<b>Conversion Rate &amp; CartAdd_to_ProductView</b>	0.9146	0.0000	11454
	<b>Category_BLE &amp; Cell_E-Conn</b>	0.8282	0.0000	11454
	<b>Cell_E-Conn &amp; Technology_RF Datacomm Conns</b>	0.8217	0.0000	11454
	<b>Cell_Mach Guard &amp; Enc &amp; Technology_Enclosures</b>	0.8119	0.0000	11454
	<b>Visits &amp; Cart Additions</b>	0.8089	0.0000	11454
	<b>Cell_Sensors &amp; Switches &amp; Technology_Proc Auto Sen</b>	0.8062	0.0000	11454
	<b>CartAdd_to_ProductView &amp; aboveMean</b>	0.7916	0.0000	11454
	<b>Conversion Rate &amp; aboveMean</b>	0.7910	0.0000	11454
	<b>Visits &amp; Orders</b>	0.7724	0.0000	11454
	<b>Cell_I-Conn &amp; Technology_Crimp Term &amp; Power</b>	0.7560	0.0000	11454
	<b>Cell_Passives &amp; Technology_Inductors</b>	0.7497	0.0000	11454
	<b>Cell_Test &amp; Measurement &amp; Technology_Electronic T&amp;M</b>	0.7124	0.0000	11454
	<b>Cell_Process Control &amp; Technology_Power Management</b>	0.7075	0.0000	11454
	<b>Cell_Test &amp; Measurement &amp; Technology_Electrical T&amp;M</b>	0.6904	0.0000	11454
	<b>Category_BLE &amp; Technology_RF Datacomm Conns</b>	0.6805	0.0000	11454
	<b>Product Views &amp; Cart Additions</b>	0.6772	0.0000	11454
	<b>Product Views &amp; Orders</b>	0.6338	0.0000	11454

We can see that Cart Additions & Orders have **r value (correlation coefficient)** close to 1 (0.9911), denoting very strong positive linear relationship, which means in practical terms that act of adding a product to the shopping cart almost always means ordering (buying) that product.

Product Views & Cart Additions correlation coefficient value of 0.6772 also shows us relatively strong positive linear relationship - it tells us that on typically with a product viewed there is 67.72% chance that it will be added to the cart.

It bears observing that the r-coefficient and p-value are not numeric values by default. We can use pd.to\_numeric to rectify this however:

```
In [64]: 1 corr_table['p-value'] = pd.to_numeric(corr_table['p-value'])
          2 corr_table['r value'] = pd.to_numeric(corr_table['r value'])
```

```
In [65]: 1 # Checking data types again:
          2 corr_table.info()
```

```
<class 'pandas.core.frame.DataFrame'>
Index: 4371 entries, Article & Visits to Technology_Thermal Management & Technology_Timers, Counters
Data columns (total 3 columns):
 #   Column    Non-Null Count  Dtype  
--- 
 0   r value   4371 non-null   float64
 1   p-value   4371 non-null   float64
 2   N         4371 non-null   int64  
dtypes: float64(2), int64(1)
memory usage: 296.6+ KB
```

We can see now that r and p - values have been converted to numbers.

```
In [66]: 1 corr_table.tail(4)
```

		r value	p-value	N
<b>Technology_Switches-EMECH &amp; Technology_Timers, Counters</b>		-0.0044	0.6393	11454
<b>Technology_Temp Control &amp; Technology_Thermal Management</b>		-0.0104	0.2667	11454
<b>Technology_Temp Control &amp; Technology_Timers, Counters</b>		-0.0055	0.5566	11454
<b>Technology_Thermal Management &amp; Technology_Timers, Counters</b>		-0.0071	0.4462	11454

Now let's look at only those correlations that are high (positive or negative) and significant:

In [67]:

```

1 #Searching for correlations that are high positive and significant:
2 corr_table[(corr_table['r value'] > 0.5) & (corr_table['p-value'] < 0.05)]

```

Out[67]:

		r value	p-value	N
	<b>Visits &amp; Product Views</b>	0.9776	0.0	11454
	<b>Visits &amp; Cart Additions</b>	0.8089	0.0	11454
	<b>Visits &amp; Orders</b>	0.7724	0.0	11454
	<b>Product Views &amp; Cart Additions</b>	0.6772	0.0	11454
	<b>Product Views &amp; Orders</b>	0.6338	0.0	11454
	<b>Cart Additions &amp; Orders</b>	0.9911	0.0	11454
	<b>Conversion Rate &amp; CartAdd_to_ProductView</b>	0.9146	0.0	11454
	<b>Conversion Rate &amp; aboveMean</b>	0.7910	0.0	11454
	<b>CartAdd_to_ProductView &amp; aboveMean</b>	0.7916	0.0	11454
	<b>Category_BLE &amp; Cell_E-Conn</b>	0.8282	0.0	11454
	<b>Category_BLE &amp; Technology_RF Datacomm Conn</b>	0.6805	0.0	11454
	<b>Category_IA&amp;C &amp; Cell_Mech &amp; Flu Pwr</b>	0.5070	0.0	11454
	<b>Category_IA&amp;C &amp; Cell_Sensors &amp; Switches</b>	0.5173	0.0	11454
	<b>Category_II&amp;T &amp; Cell_Cable, Anti-Vib &amp; Ac</b>	0.6214	0.0	11454
	<b>Category_TCFM &amp; Cell_Consumables &amp; Safety</b>	0.6129	0.0	11454
	<b>Category_TCFM &amp; Cell_Tools Storage &amp; FM</b>	0.5393	0.0	11454
	<b>Cell_Analog Disc Power &amp; Technology_Power PCB</b>	1.0000	0.0	11454
	<b>Cell_BLE E-Mech &amp; Technology_BLE E-Mech Tech</b>	1.0000	0.0	11454
	<b>Cell_Cable, Anti-Vib &amp; Ac &amp; Technology_Cable Accs</b>	0.6323	0.0	11454
	<b>Cell_Cable, Anti-Vib &amp; Ac &amp; Technology_Pwr &amp; Ind Cable</b>	0.6146	0.0	11454
	<b>Cell_Consumables &amp; Safety &amp; Technology_Fasteners</b>	0.5989	0.0	11454
	<b>Cell_E-Conn &amp; Technology_RF Datacomm Conn</b>	0.8217	0.0	11454
	<b>Cell_EMECH &amp; Technology_Indicators</b>	0.5903	0.0	11454
	<b>Cell_EMECH &amp; Technology_Thermal Management</b>	0.5480	0.0	11454
	<b>Cell_I-Conn &amp; Technology_Crimp Term &amp; Power</b>	0.7560	0.0	11454
	<b>Cell_I-Conn &amp; Technology_Ind &amp; AV Conn</b>	0.5733	0.0	11454
	<b>Cell_Mach Guard &amp; Enc &amp; Technology_Enclosures</b>	0.8119	0.0	11454
	<b>Cell_Mech &amp; Flu Pwr &amp; Technology_Pneumatics</b>	0.5314	0.0	11454
	<b>Cell_Mech &amp; Flu Pwr &amp; Technology_Pow Tran &amp; Struc</b>	0.5949	0.0	11454
	<b>Cell_Opto &amp; Technology_Displays</b>	0.5221	0.0	11454
	<b>Cell_Opto &amp; Technology_LED</b>	0.6028	0.0	11454
	<b>Cell_Opto &amp; Technology_Sensors</b>	0.6028	0.0	11454
	<b>Cell_Passives &amp; Technology_Inductors</b>	0.7497	0.0	11454
	<b>Cell_Process Control &amp; Technology_Motors &amp; Control</b>	0.5236	0.0	11454

	r value	p-value	N
Cell_Process Control & Technology_Power Management	0.7075	0.0	11454
Cell_Sensors & Switches & Technology_Proc Auto Sen	0.8062	0.0	11454
Cell_Test & Measurement & Technology_Electrical T&M	0.6904	0.0	11454
Cell_Test & Measurement & Technology_Electronic T&M	0.7124	0.0	11454

In [68]:

```

1 #Searching for correlations that are high negative and significant:
2 corr_table[(corr_table['r value'] < -0.45) & (corr_table['p-value'] < 0.05)

```

Out[68]:

	r value	p-value	N
Category_IA&C & Category_TCFM	-0.4627	0.0	11454
Category_II&T & Category_TCFM	-0.5631	0.0	11454

Highly correlated features can cause lower performance in certain classification algorithms (though non-parametric models such as decision trees/random forests are not so badly affected) because the feature redundancy, relative to how much data there is, might increase the chance of overfitting. As such, it can sometimes be useful to remove some of the correlated features, or use dimensionality reduction techniques like **Principal Component Analysis (PCA)** - which I have used in my Clustering projects. We will leave the dataset as it is for now, and see how our model performs.

## Model & Evaluate

*During the Classification workshop we've evaluated a classification model (i.e. a Decision Tree) - below we are going to follow similar steps, but in this case we will use an "ensemble method" known as Random Forests. As the name suggests, these are closely related to Decision Trees.*

Random Forests are an example of an "ensemble" method. In essence, a Random Forest is simply a "forest" (collection or ensemble) of multiple decision trees that each use a different random subsample of features - within this, each decision tree "votes" for the best solution, and the average is taken as the final prediction.

First - as is an accepted custom - let's separate the features from the class. Since typing them up one by one would've been way too time-consuming, I'm gonna "cheat" a bit:

```
In [69]: 1 list(data_dummy.columns)
```

```
Out[69]: ['Article',
'Visits',
'Product Views',
'In Stock Views',
'Cart Additions',
'Orders',
'Conversion Rate',
'CartAdd_to_ProductView',
'aboveMean',
'Category_BLE',
'Category_IA&C',
'Category_II&T',
'Category_TCFM',
'Cell_Analog Disc Power',
'Cell_BLE E-Mech',
'Cell_Cable, Anti-Vib & Ac',
'Cell_Consumables & Safety',
'Cell_E-Conn',
'Cell_EMECH',
'Cell_I-Conn',
'Cell_Mach Guard & Enc',
'Cell_Mech & Flu Pwr',
'Cell_Opto',
'Cell_Passives',
'Cell_Process Control',
'Cell_Sensors & Switches',
'Cell_Test & Measurement',
'Cell_Tools Storage & FM',
'Technology_3D Printing',
'Technology_Acc, Stor & Mat Hand',
'Technology_Adhesives, Clean & M',
'Technology_Anti-Vibration',
'Technology_BLE E-Mech Tech',
'Technology_Battery',
'Technology_Cable Accs',
'Technology_Cable Mgt & Elec Acc',
'Technology_Capacitors',
'Technology_Circuit Protection',
'Technology_Contactors',
'Technology_Crimp Term & Power',
'Technology_Ctrl & Ind Relay',
'Technology_DIN Rail & Term',
"Technology_DIN Rail PSU's",
'Technology_Displays',
'Technology_Elec & Data Cable',
'Technology_Electrical T&M',
'Technology_Electronic T&M',
'Technology_Enclosures',
'Technology_Engineering Mat',
'Technology_Fact Auto Sen',
'Technology_Fasteners',
'Technology_Fuses',
'Technology_HVAC',
'Technology_Hand Tools',
'Technology_Heavy Duty Conn',
```

```
'Technology_Hose & Pumps',
'Technology_Hydraulics',
'Technology_IT',
'Technology_Ind & AV Conns',
'Technology_Indicators',
'Technology_Inductors',
'Technology_LED',
'Technology_Level & Flow Con',
'Technology_Lighting',
'Technology_MIL Ind & Test Conns',
'Technology_Machine Safety',
'Technology_Measurement & Insp',
'Technology_Motors & Control',
'Technology_PCB Connectors',
"Technology_PLC's & HMI's",
'Technology_PPE',
'Technology_Plumbing',
'Technology_Pneumatics',
'Technology_Pow Tools & Abras',
'Technology_Pow Tran & Struc',
'Technology_Power Management',
'Technology_Power PCB',
'Technology_Printing & Office',
'Technology_Proc Auto Sen',
'Technology_Push But & Swit',
'Technology_Pwr & Ind Cable',
'Technology_RF Datacomm Conns',
'Technology_Rect Connectors',
'Technology_Relays-EMECH',
'Technology_Resistors',
'Technology_Security',
'Technology_Sensors',
'Technology_Site Safety & Janit',
'Technology_Soldering & ESD',
'Technology_Sounders/Beacons',
'Technology_Switches-EMECH',
'Technology_Temp Control',
'Technology_Thermal Management',
'Technology_Timers, Counters']
```

In [70]:

```
1 feature_cols = ['Article',
2   'Visits',
3   'Product Views',
4   'In Stock Views',
5   'Cart Additions',
6   'Orders',
7   'Category_BLE',
8   'Category_IA&C',
9   'Category_IIT&T',
10  'Category_TCFM',
11  'Cell_Analog Disc Power',
12  'Cell_BLE E-Mech',
13  'Cell_Cable, Anti-Vib & Ac',
14  'Cell_Consumables & Safety',
15  'Cell_E-Conn',
16  'Cell_EMECH',
17  'Cell_I-Conn',
18  'Cell_Mach Guard & Enc',
19  'Cell_Mech & Flu Pwr',
20  'Cell_Opto',
21  'Cell_Passives',
22  'Cell_Process Control',
23  'Cell_Sensors & Switches',
24  'Cell_Test & Measurement',
25  'Cell_Tools Storage & FM',
26  'Technology_3D Printing',
27  'Technology_Acc, Stor & Mat Hand',
28  'Technology_Adhesives, Clean & M',
29  'Technology_Anti-Vibration',
30  'Technology_BLE E-Mech Tech',
31  'Technology_Battery',
32  'Technology_Cable Accs',
33  'Technology_Cable Mgt & Elec Acc',
34  'Technology_Capacitors',
35  'Technology_Circuit Protection',
36  'Technology_Contactors',
37  'Technology_Crimp Term & Power',
38  'Technology_Ctrl & Ind Relay',
39  'Technology_DIN Rail & Term',
40  "Technology_DIN Rail PSU's",
41  'Technology_Displays',
42  'Technology_Elec & Data Cable',
43  'Technology_Electrical T&M',
44  'Technology_Electronic T&M',
45  'Technology_Enclosures',
46  'Technology_Engineering Mat',
47  'Technology_Fact Auto Sen',
48  'Technology_Fasteners',
49  'Technology_Fuses',
50  'Technology_HVAC',
51  'Technology_Hand Tools',
52  'Technology_Heavy Duty Conn',
53  'Technology_Hose & Pumps',
54  'Technology_Hydraulics',
55  'Technology_IT',
56  'Technology_Ind & AV Conns',
```

```

57  'Technology_Indicators',
58  'Technology_Inductors',
59  'Technology_LED',
60  'Technology_Level & Flow Con',
61  'Technology_Lighting',
62  'Technology_MIL Ind & Test Conns',
63  'Technology_Machine Safety',
64  'Technology_Measurement & Insp',
65  'Technology_Motors & Control',
66  'Technology_PCB Connectors',
67  "Technology_PLC's & HMI's",
68  'Technology_PPE',
69  'Technology_Plumbing',
70  'Technology_Pneumatics',
71  'Technology_Pow Tools & Abras',
72  'Technology_Pow Tran & Struc',
73  'Technology_Power Management',
74  'Technology_Power PCB',
75  'Technology_Printing & Office',
76  'Technology_Proc Auto Sen',
77  'Technology_Push But & Swit',
78  'Technology_Pwr & Ind Cable',
79  'Technology_RF Datacomm Conns',
80  'Technology_Rect Connectors',
81  'Technology_Relays-EMECH',
82  'Technology_Resistors',
83  'Technology_Security',
84  'Technology_Sensors',
85  'Technology_Site Safety & Janit',
86  'Technology_Soldering & ESD',
87  'Technology_Sounders/Beacons',
88  'Technology_Switches-EMECH',
89  'Technology_Temp Control',
90  'Technology_Thermal Management',
91  'Technology_Timers, Counters']
92
93 X = data_dummy[feature_cols]
94 y = dataset.aboveMean

```

Now let's perform a train/test split (recall that it is important that we did the one-hot encoding to dummy variables *before* the train/test split):

```
In [71]: 1 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25,
```

**We can easily train and fit a Random Forest classifier with certain hyperparameters as follows:**

```
In [72]: 1 clf2 = RandomForestClassifier(n_estimators = 500, criterion = 'entropy', max_
2
3 clf2.fit(X_train, y_train)
```

```
Out[72]: RandomForestClassifier(bootstrap=True, ccp_alpha=0.0, class_weight=None,
                                 criterion='entropy', max_depth=5, max_features='auto',
                                 max_leaf_nodes=None, max_samples=None,
                                 min_impurity_decrease=0.0, min_impurity_split=None,
                                 min_samples_leaf=1, min_samples_split=2,
                                 min_weight_fraction_leaf=0.0, n_estimators=500,
                                 n_jobs=None, oob_score=False, random_state=123,
                                 verbose=0, warm_start=False)
```

But how do I pick the best values of the hyperparameters?

The best way to do this is with a **grid search** and **cross-validation**:

```
In [73]: 1 from sklearn.model_selection import GridSearchCV
2
3 rfc = RandomForestClassifier(random_state = 123)    # we set the random_state
4
5 param_grid = [{max_depth': [1,5,10], 'n_estimators': [1, 10, 100], 'criteri
6
7 grid_search = GridSearchCV(rfc, param_grid, cv = 5, scoring = 'accuracy', re
8
9 grid_search.fit(X_train, y_train)
10
11 grid_search.best_params_
```

```
Out[73]: {'criterion': 'gini', 'max_depth': 10, 'n_estimators': 100}
```

Let's pick the best estimator from the grid search as our final model:

```
In [74]: 1 clf2 = grid_search.best_estimator_
```

## Evaluating model:

Now that we have built a model (and validated it using cross-validation), let's see how well it performs on the (so far unused) test set:

```
In [75]: 1 # Making predictions on the test set
2 y_pred = clf2.predict(X_test)
```

In [76]:

```

1 # Evaluating performance on the test by computing several metrics
2 print("Accuracy (test):", metrics.accuracy_score(y_test, y_pred))
3 print("Precision (test)", metrics.precision_score(y_test, y_pred))
4 print("Recall (test)", metrics.recall_score(y_test, y_pred))
5 print("F1-score (test)", metrics.f1_score(y_test, y_pred))

```

```

Accuracy (test): 0.8997905027932961
Precision (test) 0.8955996548748921
Recall (test) 0.8621262458471761
F1-score (test) 0.8785442234447737

```

In [77]:

```

1 #I'm getting elementwise comparison warning and - as eper advice from StackO
2 import warnings
3 with warnings.catch_warnings():
4     warnings.filterwarnings('ignore', r'elementwise comparison failed; return

```

In [78]:

```

1 # for comparison, see performance on the train set too
2 y_pred_train = clf2.predict(X_train)
3
4 print("Accuracy (train):", metrics.accuracy_score(y_train, y_pred_train))
5 print("Precision (train)", metrics.precision_score(y_train, y_pred_train))
6 print("Recall (train)", metrics.recall_score(y_train, y_pred_train))
7 print("F1-score (train)", metrics.f1_score(y_train, y_pred_train))

```

```

Accuracy (train): 0.9514551804423749
Precision (train) 0.9558350394852296
Recall (train) 0.9247311827956989
F1-score (train) 0.9400258881058536

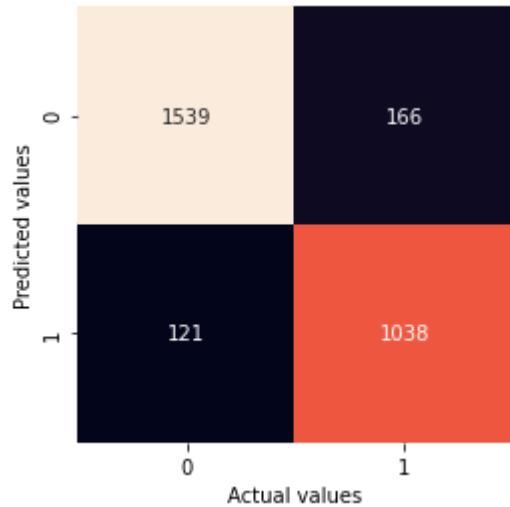
```

We can see better scores for training set but that is to be expected.

The model is performing pretty well - an accuracy of upper 80-90% is quite good, especially considering our dataset was fairly balanced (since it has been split along the mean) to begin with so our benchmark guess would only be approximately 50% accurate. The precision and recall are also similarly high, so the model is performing well for both "below the mean" and "above the mean" predictions. The model also has similar performance on the train and test set, which means it is not drastically overfitting.

In [79]:

```
1 # Look at the confusion matrix to analyse exactly what the errors are
2 conf_mat = metrics.confusion_matrix(y_test, y_pred, labels = clf2.classes_)
3
4 sns.heatmap(conf_mat.T, square = True, annot = True, fmt = 'd', cbar = False)
5 plt.xlabel('Actual values')
6 plt.ylabel('Predicted values')
7 plt.show()
```

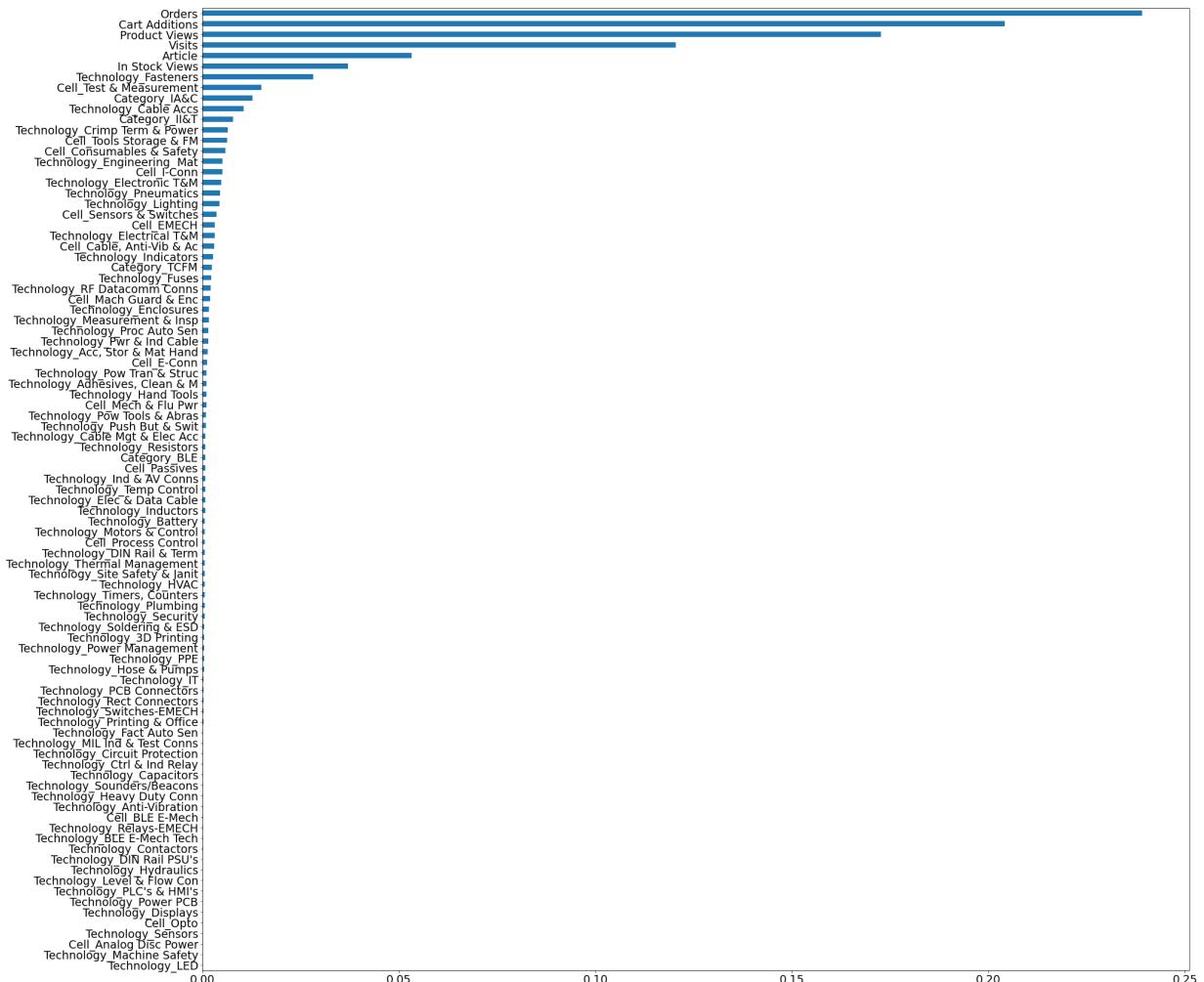


Reading the above confusion matrix, we can see that there were 1539 True Negatives and 1038 True Positives and 166 False Negatives, with 121 False Positives - in the training set (let's remember test set was 25% of all the dataset).

Let's check which features are most **important** in classifying the data - this is determined by which features appear higher up in the trees (on average in the ensemble):

```
In [80]: 1 feature_importances = pd.Series(clf2.feature_importances_, index = X.columns
2 feature_importances = feature_importances.sort_values()
3 feature_importances.plot(kind='barh', figsize = (30,30), fontsize=19)
```

Out[80]: <matplotlib.axes.\_subplots.AxesSubplot at 0x1d6e1455f08>



Not surprisingly, Orders have been indicated as Nr 1 and Cart Additions as Nr 2 when it comes to features' importance in predicting the class.

Then we have Product Views and Visits. They matter - just not as much as Orders and Cart Additions.

**Alternative approach:**

All top 5 features appear to have much higher importance than other. The challenge to the above approach is that given we are trying to predict a binary class of Above-Below mean of CartAdd\_to\_ProductView ratio, one could argue that ratio can be calculated - directly or indirectly - from the values of those top features (columns).

Let's try different approach and ignore them - dropping 'Visits', 'Product Views', 'Cart Additions' and 'Orders' from our training set features. We can certainly expect noticeably lower model scores but perhaps this approach will highlight patterns easier to miss at first glance.

In [81]:

```
1 feature_cols = ['Article',
2   'In Stock Views',
3   'Category_BLE',
4   'Category_IA&C',
5   'Category_II&T',
6   'Category_TCFM',
7   'Cell_Analog Disc Power',
8   'Cell_BLE E-Mech',
9   'Cell_Cable, Anti-Vib & Ac',
10  'Cell_Consumables & Safety',
11  'Cell_E-Conn',
12  'Cell_EMECH',
13  'Cell_I-Conn',
14  'Cell_Mach Guard & Enc',
15  'Cell_Mech & Flu Pwr',
16  'Cell_Opto',
17  'Cell_Passives',
18  'Cell_Process Control',
19  'Cell_Sensors & Switches',
20  'Cell_Test & Measurement',
21  'Cell_Tools Storage & FM',
22  'Technology_3D Printing',
23  'Technology_Acc, Stor & Mat Hand',
24  'Technology_Adhesives, Clean & M',
25  'Technology_Anti-Vibration',
26  'Technology_BLE E-Mech Tech',
27  'Technology_Battery',
28  'Technology_Cable Accs',
29  'Technology_Cable Mgt & Elec Acc',
30  'Technology_Capacitors',
31  'Technology_Circuit Protection',
32  'Technology_Contactors',
33  'Technology_Crimp Term & Power',
34  'Technology_Ctrl & Ind Relay',
35  'Technology_DIN Rail & Term',
36  "Technology_DIN Rail PSU's",
37  'Technology_Displays',
38  'Technology_Elec & Data Cable',
39  'Technology_Electrical T&M',
40  'Technology_Electronic T&M',
41  'Technology_Enclosures',
42  'Technology_Engineering Mat',
43  'Technology_Fact Auto Sen',
44  'Technology_Fasteners',
45  'Technology_Fuses',
46  'Technology_HVAC',
47  'Technology_Hand Tools',
48  'Technology_Heavy Duty Conn',
49  'Technology_Hose & Pumps',
50  'Technology_Hydraulics',
51  'Technology_IT',
52  'Technology_Ind & AV Conns',
53  'Technology_Indicators',
54  'Technology_Inductors',
55  'Technology_LED',
56  'Technology_Level & Flow Con',
```

```
57 |     'Technology_Lighting',
58 |     'Technology_MIL Ind & Test Conns',
59 |     'Technology_Machine Safety',
60 |     'Technology_Measurement & Insp',
61 |     'Technology_Motors & Control',
62 |     'Technology_PCB Connectors',
63 |     "Technology_PLC's & HMI's",
64 |     'Technology_PPE',
65 |     'Technology_Plumbing',
66 |     'Technology_Pneumatics',
67 |     'Technology_Pow Tools & Abras',
68 |     'Technology_Pow Tran & Struc',
69 |     'Technology_Power Management',
70 |     'Technology_Power PCB',
71 |     'Technology_Printing & Office',
72 |     'Technology_Proc Auto Sen',
73 |     'Technology_Push But & Swit',
74 |     'Technology_Pwr & Ind Cable',
75 |     'Technology_RF Datacomm Conns',
76 |     'Technology_Rect Connectors',
77 |     'Technology_Relays-EMECH',
78 |     'Technology_Resistors',
79 |     'Technology_Security',
80 |     'Technology_Sensors',
81 |     'Technology_Site Safety & Janit',
82 |     'Technology_Soldering & ESD',
83 |     'Technology_Sounders/Beacons',
84 |     'Technology_Switches-EMECH',
85 |     'Technology_Temp Control',
86 |     'Technology_Thermal Management',
87 |     'Technology_Timers, Counters']  
88  
89 X = data_dummy[feature_cols]
90 y = dataset.aboveMean
91
92
93 X_train, X_test, y_train, y_test = train_test_split(X, y, test_size = 0.25,
94
95 rfc = RandomForestClassifier(random_state = 123)
96
97 param_grid = [{'max_depth': [1,5,10], 'n_estimators': [1, 10, 100], 'criteri
98
99 grid_search = GridSearchCV(rfc, param_grid, cv = 5, scoring = 'accuracy', re
100 grid_search.fit(X_train, y_train)
101
102 clf2 = grid_search.best_estimator_
103
104 y_pred = clf2.predict(X_test)
```

In [82]:

```
1 # Evaluating performance on the test by computing several metrics
2 print("Accuracy (test):", metrics.accuracy_score(y_test, y_pred))
3 print("Precision (test)", metrics.precision_score(y_test, y_pred))
4 print("Recall (test)", metrics.recall_score(y_test, y_pred))
5 print("F1-score (test)", metrics.f1_score(y_test, y_pred))
```

```
Accuracy (test): 0.7081005586592178
Precision (test) 0.7233009708737864
Recall (test) 0.4950166112956811
F1-score (test) 0.5877712031558185
```

In [83]:

```
1 # for comparison, let's see performance on the train set too
2 y_pred_train = clf2.predict(X_train)
3
4 print("Accuracy (train):", metrics.accuracy_score(y_train, y_pred_train))
5 print("Precision (train)", metrics.precision_score(y_train, y_pred_train))
6 print("Recall (train)", metrics.recall_score(y_train, y_pred_train))
7 print("F1-score (train)", metrics.f1_score(y_train, y_pred_train))
```

```
Accuracy (train): 0.7628637951105938
Precision (train) 0.8094253823894171
Recall (train) 0.5540464063384267
F1-score (train) 0.6578195867629766
```

We can see that testing set scores are a bit worse than those of training set, e.g. testing Accuracy seems to be circa 6% points lower than in the training set while Precision - 7.6% points worse. It's noticeable but not drastically different.

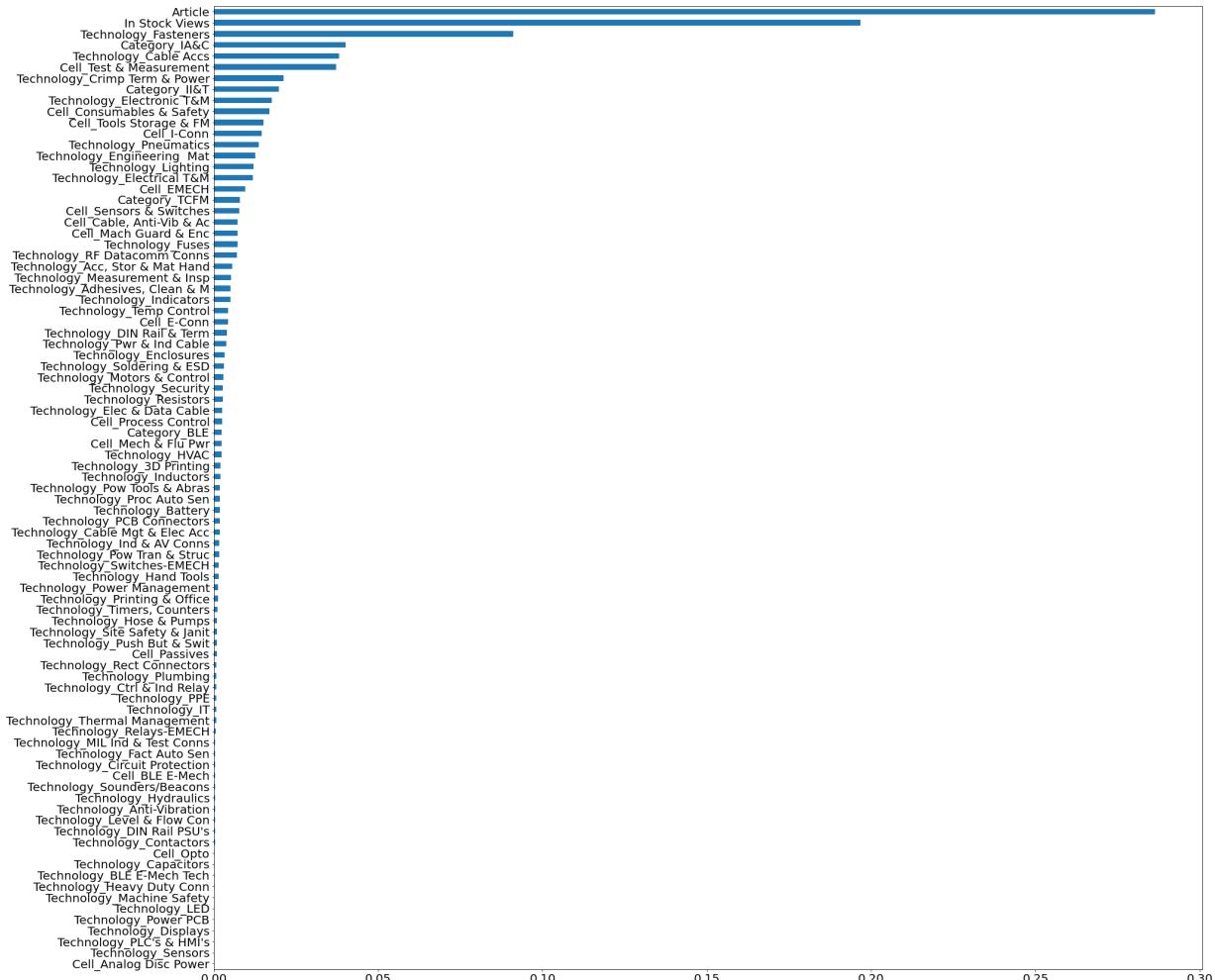
In [84]:

```

1 #Checking features importance in the alternative approach:
2 feature_importances = pd.Series(clf2.feature_importances_, index = X.columns
3 feature_importances = feature_importances.sort_values()
4 feature_importances.plot(kind='barh', figsize = (30,30), fontsize=20)

```

Out[84]: &lt;matplotlib.axes.\_subplots.AxesSubplot at 0x1d6d9093c88&gt;



Since Articles and In-Stock Views were Nr 5 and 6th accordingly in our first model, not surprisingly they rose here to top 2 spots when it comes to features' importance. But we can now see Fasteners in Technology and IA&C in category rising to prominent positions when it comes to how important each feature is.

## Visualizing decision trees in our classifier models:

We can also visualise the individual trees (as we did in the workshop).

**Drawing a decision tree in a single decision tree model (simpler and easier to compute):**

In [85]:

```
1 # Telling system to use a single decision tree from scikit
2 model=tree.DecisionTreeClassifier()
3 #Training model:
4 model.fit(X_train, y_train)
```

Out[85]: DecisionTreeClassifier(ccp\_alpha=0.0, class\_weight=None, criterion='gini',
max\_depth=None, max\_features=None, max\_leaf\_nodes=None,
min\_impurity\_decrease=0.0, min\_impurity\_split=None,
min\_samples\_leaf=1, min\_samples\_split=2,
min\_weight\_fraction\_leaf=0.0, presort='deprecated',
random\_state=None, splitter='best')

In [86]:

```
1 #Drawing decision tree from a single tree classifier using graphviz:
2 dot_data=StringIO()
3 tree.export_graphviz(model, out_file=dot_data, precision = 3, filled=True,
4 rounded=True, special_characters=False, feature_names=X_
5 class_names=True)
6 graph=pydotplus.graph_from_dot_data(dot_data.getvalue())
7 Image(graph.create_png())
```

dot: graph is too large for cairo-renderer bitmaps. Scaling by 0.384305 to fit

Out[86]:



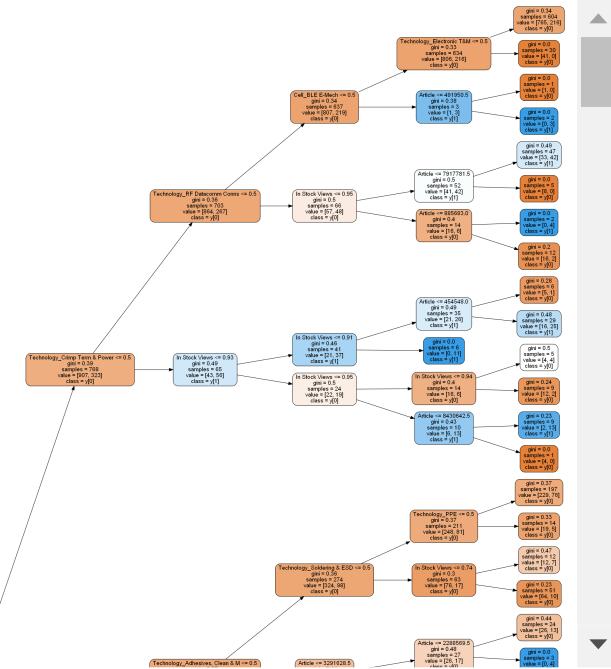
**Drawing a selected decision tree from Random Forest classifier model (more complex and compute-intensive but more robust model):**

Each tree represents a subset of the random forest ensemble model. For example, we might want to visualise tree number 39, but we could also explore any of the trees from tree 0 to tree 50 in our 100 tree ensemble.

**Let's draw a tree number 39 from our random forest ensemble model:**

```
In [87]: 1 # Extracting single tree from our rfc:
2 estimator = clf2.estimators_[39]
3
4 # Exporting tree as a dot file:
5 dot_data = StringIO()
6 export_graphviz(estimator, out_file = dot_data,
7                  special_characters = False, feature_names = X_train.columns,
8                  rounded = True, proportion = False,
9                  precision = 2, filled = True, rotate = True)
10
11 graph = pydotplus.graph_from_dot_data(dot_data.getvalue())
12 Image(graph.create_png())
```

Out[87]:



**Now I'm going to visualize a decision tree number 50 from our random forest model:**

This is alternative approach that I have discovered while dealing with an unforeseen error (having to do with Python running into problems when trying to concatenate np.bool\_() with strings (I have since managed to resolve it but in the process discovered that multi-step approach - detailed below):

```
In [88]: 1 # Extracting single tree from our rfc:
2 estimator = clf2.estimators_[50]
3
4 # Exporting tree as a dot file:
5 export_graphviz(estimator,
6                  out_file='tree.dot',
7                  feature_names = X_train.columns,
8                  class_names = True,
9                  rounded = True, proportion = False,
10                 precision = 2, filled = True, impurity=True)
```

```
In [89]: 1 # Convert to png
2 from subprocess import call
3 call(['dot', '-Tpng', 'tree.dot', '-o', 'tree.png', '-Gdpi=600'])
```

Out[89]: 0

```
In [90]: 1 # Display in jupyter notebook
2 from IPython.display import Image
3 Image(filename = 'tree.png')
```

Out[90]:



## BONUS MATERIAL: Catboost

**CatBoost is a machine learning algorithm that uses gradient boosting on decision trees. It is available as an open source library. It has been developed by Yandex and more information is available here:**

<https://catboost.ai/docs/concepts/about.html> (<https://catboost.ai/docs/concepts/about.html>)

```
In [91]: 1 # Importing necessary components of catboost:
2 from catboost import CatBoostClassifier, Pool, cv
3 # for visualizing metrics
4 from catboost import MetricVisualizer
```

**Now let's create the model itself:**

I'm going to go here with default parameters (as they provide a really good baseline almost all the time). The only thing I'd like to specify here is custom\_loss parameter, as this would give us an ability to see what's going on in terms of key metric - accuracy, as well as to be able to watch for logloss, as it would be more smooth on dataset of this size.

In [92]:

```

1 #Creating catboost model:
2 model3 = CatBoostClassifier(
3     custom_loss=['Accuracy'],
4     random_seed=49,
5     logging_level='Silent'
6 )

```

Below I'll train our model in catboost gradient boosting, using default parameters (which are quite useful As Is), plotting the logloss of Accuracy.

After Wikipedia: "**Logarithmic loss (related to cross-entropy) measures the performance of a classification model where the prediction input is a probability value between 0 and 1. The goal of our machine learning models is to minimize this value. A perfect model would have a log loss of 0. Log loss increases as the predicted probability diverges from the actual label. So predicting a probability of .015 when the actual observation label is 1 would be bad and result in a high log loss.**"

It's worth pointing out that Catboost allows for customizing loss function so instead of Accuracy, one could e.g. plot Recall. But let's stick to Accuracy.

In the output we can interact with the graph, by (un)checking the boxes - e.g by clicking on Smooth and adjust the value of smoothing between 0 to 1 or applying logarithmic scale:

In [93]:

```

1 #Training our model:
2 model3.fit(
3     X_train, y_train,
4     cat_features=None,
5     eval_set=(X_test, y_test),
6     #     Logging_Level='Verbose', # you can uncomment this for text output but
7     plot=True
8 )

```

MetricVisualizer(layout=Layout(align\_self='stretch', height='500px'))

Out[93]: &lt;catboost.core.CatBoostClassifier at 0x1d6855e5708&gt;

Learn line is basically plotting of logarithmic loss of Accuracy of our training set while Eval line - for Evaluation - denotes log-loss of Accuracy in our training set. The lower the log-loss the better.

**\*We can see that the best result has been achieved by catboost in 986th boosting step in evaluation (testing) dataset - 13 steps sooner than in training set.**

In [94]:

```

1 print('Simple model validation accuracy: {:.4}'.format(
2     accuracy_score(y_test, model.predict(X_test))
3 ))

```

Simple model validation accuracy: 0.7336

Catboost enables one to choose from many supported metrics for overfitting detection and best model selection. Let's do the same, but this time highlighting Precision as key metric to watch for as a criterium to avoid overfitting:

In [95]:

```
1 #Creating catboost model:
2 model4 = CatBoostClassifier(
3     custom_loss=['Precision'],
4     random_seed=36,
5     logging_level='Silent'
6 )
```

In [96]:

```
1 #Training our model:
2 model4.fit(
3     X_train, y_train,
4     cat_features=None,
5     eval_set=(X_test, y_test),
6     #logging_level='Verbose',
7     plot=True,
8 )
```

```
MetricVisualizer(layout=Layout(align_self='stretch', height='500px'))
```

Out[96]: <catboost.core.CatBoostClassifier at 0x1d6858a0ec8>

This time best Precision log-loss was achieved by 998th boosting step.

In [97]:

```
1 print('Simple model validation precision: {:.4}'.format(precision_score(y_te
```

```
Simple model validation precision: 0.6864
```

### **Model cross-validation in Catboost:**

It's valuable to validate our model, but to cross-validate it - even better:

In [98]:

```
1 cv_params = model3.get_params()
2 cv_params.update({
3     'loss_function': 'Logloss'
4 })
5 cv_data = cv(
6     Pool(X, y),
7     cv_params,
8     plot=True
9 )
```

```
MetricVisualizer(layout=Layout(align_self='stretch', height='500px'))
```

Now we have values of our loss functions at each boosting step averaged by 3 folds, which should provide us with a more accurate estimation of our model performance.

```
In [99]: 1 print('Best validation accuracy score: {:.2f}±{:.2f} on step {}'.format(
2     np.max(cv_data['test-Accuracy-mean']),
3     cv_data['test-Accuracy-std'][np.argmax(cv_data['test-Accuracy-mean'])]),
4     np.argmax(cv_data['test-Accuracy-mean']))
5 ))
```

Best validation accuracy score: 0.75±0.00 on step 999

```
In [100]: 1 # Printing validation Accuracy score:
2 print('Precise validation accuracy score: {}'.format(np.max(cv_data['test-Ac')))
```

Precise validation accuracy score: 0.7489091469740589

```
In [101]: 1 # Cross-validating 2nd catboost model (with Precision as a key metric to mon
2 cv_params = model4.get_params()
3 cv_params.update({
4     'loss_function': 'Logloss'
5 })
6 cv_data = cv(
7     Pool(X, y),
8     cv_params,
9     plot=True
10 )
```

MetricVisualizer(layout=Layout(align\_self='stretch', height='500px'))

```
In [102]: 1 # Printing validation Precision score:
2 print('Precise validation Precision score: {}'.format(np.max(cv_data['test-P')))
```

Precise validation Precision score: 0.7089567224477711

## Applying the model (predictions):

In [103]:

```

1 predictions = model3.predict(X_test)
2 predictions_probs = model3.predict_proba(X_test)
3 # Printing first 10 predictions:
4 print(predictions[:10])
5 print(predictions_probs[:10])

```

['False' 'True' 'False' 'False' 'False' 'True' 'True' 'False' 'False' 'False'  
'True']  
[[0.99402407 0.00597593]  
[0.08929608 0.91070392]  
[0.77264062 0.22735938]  
[0.61959351 0.38040649]  
[0.99608816 0.00391184]  
[0.08776633 0.91223367]  
[0.13130309 0.86869691]  
[0.92669961 0.07330039]  
[0.76693001 0.23306999]  
[0.31261148 0.68738852]]

We can get better ones by using some catboost's built in features:

Let's define some params and create Pool for more convenience. It stores all information about dataset (features, labeles, categorical features indices, weights and and much more):

In [104]:

```

1 params = {
2     'iterations': 1500,
3     'learning_rate': 0.1,
4     'eval_metric': 'Accuracy',
5     'random_seed': 55,
6     'logging_level': 'Silent',
7     'use_best_model': False
8 }
9 train_pool = Pool(X_train, y_train)
10 validate_pool = Pool(X_test, y_test)

```

### Using best model:

If we essentially have a validation set, it's always better to use the use\_best\_model parameter during training. By default, this parameter is enabled. If it is enabled, the resulting trees ensemble is shrinking to the best iteration.

```
In [105]: 1 model = CatBoostClassifier(**params)
2 model.fit(train_pool, eval_set=validate_pool)
3
4 best_model_params = params.copy()
5 best_model_params.update({
6     'use_best_model': True
7 })
8 best_model = CatBoostClassifier(**best_model_params)
9 best_model.fit(train_pool, eval_set=validate_pool)
```

Out[105]: <catboost.core.CatBoostClassifier at 0x1d684ebe908>

Another brilliant feature of catboost is MetricVisualizer which functionality allow one to plot multiple evaluation metrics on 1 chart for quick comparison, saving space and lines of code. One can simply click on different metrics on the same graph (on top of the graph scale):

```
In [106]: 1 eval_metrics = best_model.eval_metrics(validate_pool, ['Accuracy', 'Precision'])
```

MetricVisualizer(layout=Layout(align\_self='stretch', height='500px'))

### Feature importance:

Finally using catboost get\_feature\_importance function we can easily get a sorted ranking of importance of features from our training dataset in the model:

In [107]:

```
1 feature_importances = model.get_feature_importance(train_pool)
2 feature_names = X_train.columns
3 for score, name in sorted(zip(feature_importances, feature_names), reverse=True):
4     print('{}: {}'.format(name, score))
```

Article: 38.58037746569475  
In Stock Views: 19.945221434607994  
Cell\_Test & Measurement: 4.453229995694893  
Technology\_Fasteners: 3.075684543833669  
Category\_IA&C: 1.6238240967871782  
Technology\_Cable Accs: 1.4707558703879104  
Technology\_Engineering Mat: 1.1884922718997997  
Technology\_Lighting: 1.147151798085583  
Technology\_Pwr & Ind Cable: 0.9592125696495065  
Technology\_Measurement & Insp: 0.9299113084553916  
Technology\_Acc, Stor & Mat Hand: 0.9189732693412925  
Category\_II&T: 0.8977574611686182  
Category\_TCFM: 0.8925073446495475  
Technology\_Hand Tools: 0.8246671364422077  
Technology\_Elec & Data Cable: 0.8176029254335158  
Cell\_Tools Storage & FM: 0.7664729863487418  
Technology\_Pneumatics: 0.7436176399754864  
Technology\_Temp Control: 0.7099025402698625  
Technology\_Pow Tools & Abras: 0.7016134327066847  
Technology\_RF Datacomm Conns: 0.6966923380765181  
Technology\_Electronic T&M: 0.6857253391616259  
Cell\_Consumables & Safety: 0.6822667058639351  
Technology\_Crimp Term & Power: 0.6780522582571002  
Technology\_Pow Tran & Struc: 0.6603148434657569  
Cell\_Cable, Anti-Vib & Ac: 0.631736602587004  
Technology\_Adhesives, Clean & M: 0.6256506845186643  
Technology\_Security: 0.5827222591984218  
Technology\_Inductors: 0.5677006330463804  
Technology\_Motors & Control: 0.5228057572782586  
Technology\_Cable Mgt & Elec Acc: 0.5176601404728891  
Cell\_I-Conn: 0.5076422052768433  
Cell\_Sensors & Switches: 0.4910705998732789  
Technology\_Fuses: 0.4792971519701038  
Technology\_Enclosures: 0.4773347663579847  
Technology\_Power Management: 0.47541735577385286  
Technology\_Soldering & ESD: 0.47486220140071844  
Technology\_Rect Connectors: 0.46903053700430714  
Cell\_E-Conn: 0.4569292387099711  
Technology\_Site Safety & Janit: 0.45292488442642465  
Technology\_HVAC: 0.43287257318519134  
Technology\_Electrical T&M: 0.4272260616941882  
Technology\_Proc Auto Sen: 0.4001105898200931  
Technology\_Battery: 0.3910137579678977  
Technology\_PPE: 0.3809122384251792  
Technology\_Indicators: 0.357662521528442  
Technology\_Printing & Office: 0.3431593597501616  
Technology\_Ind & AV Conns: 0.34178870467270256  
Cell\_Process Control: 0.33726930743886985  
Technology\_Plumbing: 0.3178627058294164  
Technology\_DIN Rail & Term: 0.30867118446173974  
Technology\_Resistors: 0.2539185766786817  
Cell\_Passives: 0.25371366606773055

Technology\_3D Printing: 0.2535336477392696  
Technology\_Push But & Swit: 0.24886900189267627  
Technology\_Thermal Management: 0.2400101878899823  
Category\_BLE: 0.23513256991463735  
Technology\_MIL Ind & Test Conns: 0.22962205686294399  
Technology\_IT: 0.22689345679450904  
Technology\_Ctrl & Ind Relay: 0.22290589300734792  
Technology\_Timers, Counters: 0.20274748617712016  
Technology\_Hose & Pumps: 0.19149955195612337  
Technology\_Switches-EMECH: 0.1681194915506174  
Technology\_PCB Connectors: 0.16125821398442616  
Cell\_Mech & Flu Pwr: 0.14268532546053342  
Technology\_Fact Auto Sen: 0.11478949659434402  
Technology\_Level & Flow Con: 0.11358785023799228  
Cell\_EMECH: 0.11358387193298426  
Technology\_Hydraulics: 0.10441735285386344  
Cell\_Mach Guard & Enc: 0.10177369985271964  
Technology\_Anti-Vibration: 0.08657259395555832  
Technology\_Sounders/Beacons: 0.08163294097334729  
Technology\_Capacitors: 0.07695762536498478  
Cell\_BLE E-Mech: 0.07654191146334507  
Cell\_Opto: 0.06016903645760899  
Technology\_DIN Rail PSU's: 0.050758276842717304  
Technology\_Circuit Protection: 0.047242602696098174  
Technology\_Relays-EMECH: 0.032711571615089426  
Technology\_Machine Safety: 0.02280257891368375  
Technology\_Contactors: 0.02203406502316675  
Technology\_BLE E-Mech Tech: 0.020719409859415387  
Technology\_Heavy Duty Conn: 0.011327756138073224  
Cell\_Analog Disc Power: 0.005022242293213834  
Technology\_Power PCB: 0.002460753797098796  
Technology\_PLC's & HMI's: 0.0018232699393543673  
Technology\_Displays: 0.0006637554983654682  
Technology\_LED: 0.00013661279379845789  
Technology\_Sensors: 0.0

Finally, catboost allows us to draw statistics of selected features on a graph:

In [108]:

```
1 stats = model.calc_feature_statistics(X_train,
2                               y_train,
3                               feature = [1,4,6],
4                               plot=True,
5                               max_cat_features_on_plot=10,
6                               thread_count=-1,
7                               plot_file=None)
```

## Summary:

So far there isn't anything obvious popping up and begging for a deep dive. That could be due to multi-dimensional analysis performed above. Let's refresh our memory by having another look at correlation table of our features, this time filtering all features correlated in some way with our prediction class, namely aboveMean. As a reminder, that class covers a ratio of Cart Additions to Product Views that is above mean value for all values of that column.

In [109]:

```
1 #Filtering our previous correlation table by rows that have a correlation pa
2 corr_aboveavg = corr_table.filter(like='aboveMean', axis=0)
```

In [110]:

```
1 corr_aboveavg.sample()
```

Out[110]:

		r value	p-value	N
aboveMean & Technology_Electronic T&M		-0.1017	0.0	11454

In [111]:

```
1 #Sorting our filtered correlation table in descending order (to see highest
2 corr_aboveavg.sort_values(by=['r value'], ascending=False).head(15)
```

Out[111]:

		r value	p-value	N
CartAdd_to_ProductView & aboveMean		0.7916	0.0	11454
Conversion Rate & aboveMean		0.7910	0.0	11454
Orders & aboveMean		0.2666	0.0	11454
Cart Additions & aboveMean		0.2488	0.0	11454
aboveMean & Technology_Fasteners		0.2104	0.0	11454
In Stock Views & aboveMean		0.1769	0.0	11454
aboveMean & Technology_Cable Accs		0.1366	0.0	11454
aboveMean & Cell_Consumables & Safety		0.1232	0.0	11454
aboveMean & Technology_Crimp Term & Power		0.1073	0.0	11454
aboveMean & Cell_I-Conn		0.1032	0.0	11454
aboveMean & Category_II&T		0.0960	0.0	11454
aboveMean & Cell_EMECH		0.0831	0.0	11454

In [112]:

```
1 #Sorting correlation table by r value in ascending order:
2 corr_aboveavg.sort_values(by=['r value'], ascending=True).head(10)
```

Out[112]:

		r value	p-value	N
aboveMean & Category_IA&C		-0.1595	0.0	11454
aboveMean & Cell_Test & Measurement		-0.1435	0.0	11454
aboveMean & Cell_Tools Storage & FM		-0.1118	0.0	11454
aboveMean & Technology_Electronic T&M		-0.1017	0.0	11454
Article & aboveMean		-0.1011	0.0	11454
aboveMean & Cell_Mach Guard & Enc		-0.1001	0.0	11454
aboveMean & Technology_Electrical T&M		-0.0997	0.0	11454
aboveMean & Cell_Sensors & Switches		-0.0958	0.0	11454
aboveMean & Technology_Lighting		-0.0957	0.0	11454
aboveMean & Technology_Enclosures		-0.0857	0.0	11454

## CONCLUSIONS:

As we can see above, that wasn't a straight line journey. I have explored multiple approaches to classification problem, from a single decision tree through a random forest ensemble all the way to relatively little known Catboost gradient boosting model approach. None of those managed to reveal anything ground-breaking, not because of any major inherent flaws but more so because our data set hasn't been as rich as one would have liked. Out of necessity, I have elected to try different features to predict aboveMean class to find out which of those model features would be the best predictors of a ratio of cart additions to product views being above average.

Not surprisingly, given it's a ratio, the top 2 were Cart Additions and Product Views and just behind them - Conversion Rate. The latter because our ratio can be treated as a proxy for Conversion Rate (which is calculated by dividing Orders to Product Views). Since typically a customer will proceed from adding product to a shopping cart to ordering and purchasing items, it was to be expected.

There are some interesting observations to be made, though. Among our main categories, Industrial Automation & Control (IA & C) has the highest negative r value in a correlation pair with aboveMean which tells us there exists an inverse relationship between products in that category being viewed and added to the cart and above average proportion of those 2 activities. It might sound a bit confusing but what it basically means is if a product being viewed belongs to Industrial Automation & Control category, chances are client is less likely to add it to the cart and thus our conversion rate will be below average (when it comes to all 10,000+ products we sell).

The same can be said about Cell\_Test & Measurement (which our random forest ensemble model has picked as 5th most important features when it comes to predicting if our ratio will be above or below mean) or Cell\_Tools Storage & FM product ranges. While r values of -0.1435 or -0.1118 are not very high, they are among the top negative relationships with aboveMean so for some reason our clients are less likely to add them to the shopping cart than virtually all other best selling product ranges. It could be due to pricing, insufficient product information being available, poor picture quality of myriad other reasons, but the pattern exists. It would be interesting to learn what profit margin those products provide, because if it's significant and sold in large enough volumes, it might make sense to try to improve their conversion ratio to bring it closer to our aggregate average.

On the positive side, some of our technologies, like Fasteners and Cable Accessories have slightly positive correlation with above average ratio (and thus conversion rate). For whatever reason our clients are slightly more likely to add those items to the cart and - presumably - pull the trigger when it comes to shopping for them. The same can be said about Cell\_Consumables & Safety. Those products seem to be more tempting to our clients so perhaps there are some learning around the way we do promoting and pricing those products so that they get higher than average conversion rates.

While I can't see a full path yet to implementing these findings into Inventory Management function's operations or to directly boost our Marketing department efforts, it further reinforces my view that we, as a global organization, would be well served to undertake to collect more data on our product features and customer journey across our online sales channel, so that my team has

more features to work with. Those we seem to be collecting are too often either relatively highly correlated (thus risking overfitting) or don't provide much insight. Overall, we would definitely benefit from an ability to notice otherwise unseen patterns and to come up with a model to generalize those patterns onto observations and it's precisely why tools like major classification algorithms are essential in any data scientist's toolbox and of value to us as a corporation stocking hundreds of thousands of products and hoping to predict which ones are going to sell better or foresee which ones will sell out faster . That said, even the best algorithm on its own won't help us much unless the quality and breadth of data we collect gets expanded so that we draw insights from it.

In [ ]:

1