

















# From Pixe's to Profits



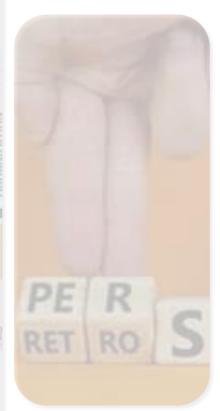
Business Understanding

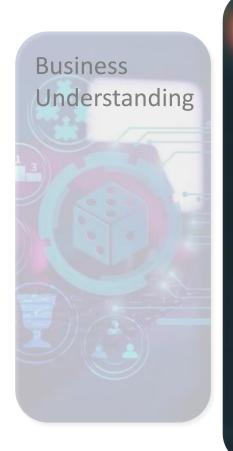
- Background Intro
- Project Objectives











Data Source Preprocessing

- Data Source
- Repository
- Preprocessing









Modelling Text Analytics

- PredictiveModelling
- Text Analytics



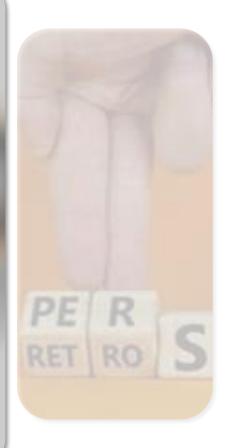




Modelling Text Analytics

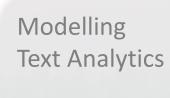
Model Evaluation

- Assessment
- Interpretability







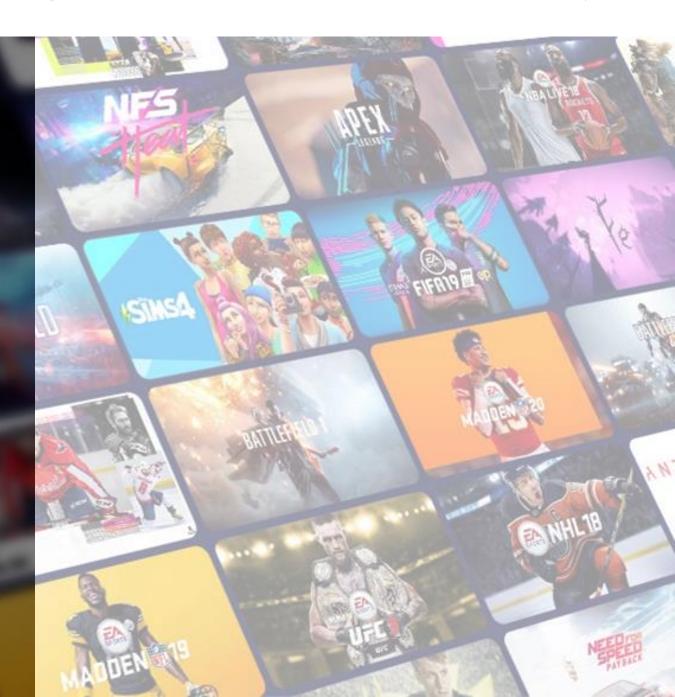


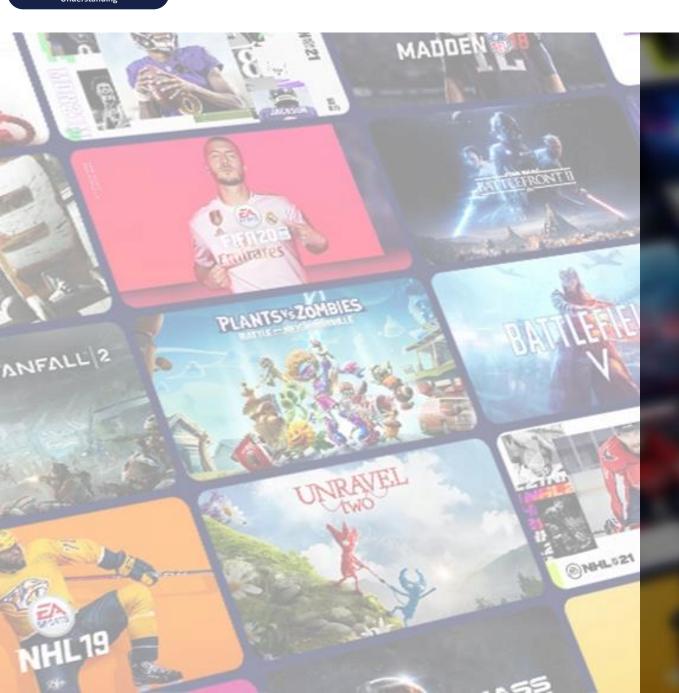




Retrospective

- Limitations
- Learning Points





## Historical data



Guiding patterns and insights to better business decisions

# Project Objectives

**Key Revenue Indicators** 

### Sales Revenue Prediction Model

# Project Objectives



Key Revenue Indicators

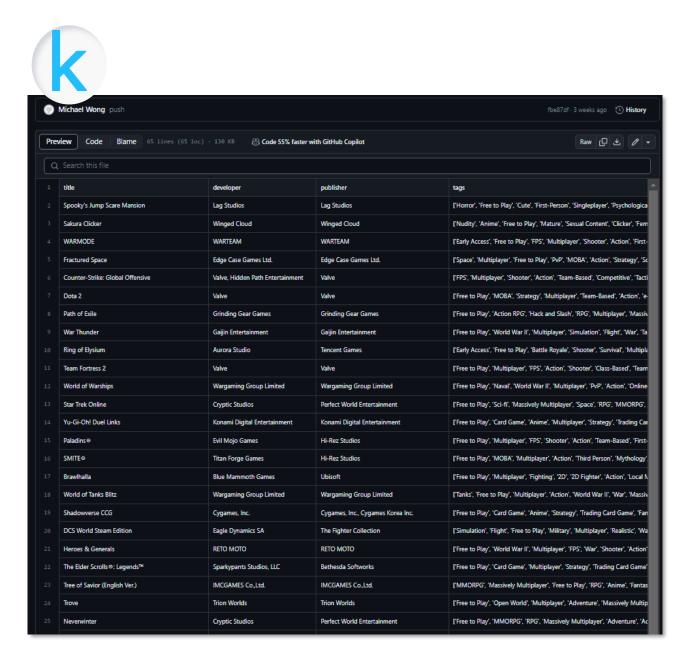
# **Project Objectives**

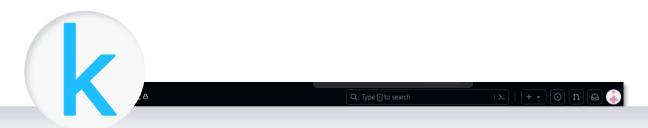


### Sales Revenue Prediction Model

More informed business decisions

Text Analysis of Players' Reviews





A	В	0	D	E	FG	н	N		0   P	Q Q	j K	5		
	platform -	genre *	publisher -	vg_score 💌	critic_score 💌 user_scor	re 💌 total_shipp	ed 💌 Production Co	st 💌 releas	e_year 🕝 Developer_metac 🕶	Genre_metacritic_game_info	<sup>▼</sup> Metascore_metacritic_game_ir ▼	Avg_Userscore_metacritic_game_ =	No. of players_metacritic_game_info 💌	releasedate_qua
J1 A.D.	PC	Simulation	Sony Computer Entertainment	8	7.7	5	1	6.96	2006 RelatedDesigns	Strategy;Real-Time;Historic;General;Historic	79	8.4	1 to 4	26/10/2006
2002 FIFA World Cup	PC	Sports	Sony Computer Entertainment	2	6	7	277	6.06	2002 EASports	Sports;Traditional;Team;Soccer;Sim;Sim	80	8	Mutliplayer	26/4/2002
2010 FIFA World Cup South Afri	ic PS3	Sports	Ubisoft	7	8.2	9	53	4.24	2010 EACanada	Sports;Traditional;Team;Soccer;Sim;Sim	82	7.8	1	27/4/2010
ABZU	PC	Adventure	Activision	9	8	9	0.58	3.75	2016 GiantSquid	Action Adventure;Adventure;General;3D;Third-Person	83	6.9	1	2/8/2016
ABZU	PS4	Adventure	Electronic Arts	8	7	7	117	3.24	2017 GiantSquid	Fantasy; Action Adventure; Adventure; General; 3D; Third-Person	78	7.6	1	2/8/2016
Aggressive Inline	PS2	Sports	Activision	6	8.7	10	117	3.21	2002 Z-Axis,Ltd.	Sports; Alternative; Skating	85	8.7	1 to 2	28/5/2002
Alice in Wonderland	DS	Adventure	Electronic Arts	8	7	3	59	0.03	2010 EtrangesLibellules	Action; Action Adventure; General; General; Fantasy	78	tbd	1	2/3/2010
Alien Isolation	PS4	Adventure	Frontier Developments	10	4	8	189	2.61	2014 CreativeAssembly	Action Adventure;Sci-Fi;General;Survival	79	8.2	1	6/10/2014
Alien Hominid	PS2	Shooter	Frontier Developments	3	4	8	296	2.58	2004 TheBehemoth	Action;Shooter;Scrolling	78	8.3	1	21/11/2004
Aliens versus Predator 2	PC	Shooter	Sony Computer Entertainment	5	8.2	2	275	2.07	2001 MonolithProduction	Action;Shooter;Shooter;First-Person;Sci-Fi;Sci-Fi;Arcade	85	8.7	Mutliplayer	22/10/2001
Allegiance	PC	Simulation	Activision	1	10	10	231	2.81	2000 MicrosoftGameStu-	Simulation;Sci-Fi;Space;Small Spaceship;Small Spaceship;Combat	86	8.3	1	16/3/2000
Alone in the Dark The New Nightr	m: PS	Adventure	Electronic Arts	3	7.2	4	162	5.65	2001 Darkworks	Action Adventure;Horror	77	8.2	1	8/5/2001
Animal Crossing Wild World	DS	Simulation	Amazon Game Studios	6	8.5	4	11.75	5.09	2005 Nintendo	Simulation;Miscellaneous;Virtual Life;Virtual;Virtual Life	86	8.6	1 to 4	23/11/2005
Anno 1701 Dawn of Discovery	DS	Simulation	Ubisoft	9	8	5	221	4.05	2008 SunflowersInteraction	Simulation;General;General	78	7.4	1 to 4	8/6/2007
Apollo Justice Ace Attorney	DS	Adventure	Ubisoft	7.9	8	10	23	5.95	2008 Capcom	Adventure;General;General;Visual Novel	78	8.2	1	12/4/2007
Aqua Aqua	PS2	Puzzle	Frontier Developments	10	4	2	299	4.23	2001 ZedTwoLimited	Miscellaneous;Puzzle;Puzzle;General	79	tbd	1 to 2	2/11/2000
Aquaria	PC	Adventure	Activision	4	4	10	0.16	2.24	2007 BitBlot	Action Adventure;Fantasy;General;Fantasy	82	8.2	1	7/12/2007
ArmA II	PC	Shooter	Electronic Arts	3	8	1	214	4.22	2009 BohemiaInteractive	Action;Shooter;Shooter;First-Person;Modern;Modern;Arcade	77	7.5	1	17/6/2009
Armed and Dangerous	PC	Shooter	Electronic Arts	9	8	9	112	2.78	2003 PlanetMoonStudios	Action;Shooter;Shooter;Third-Person;Sci-Fi;Sci-Fi;Arcade	78	7.8	1	2/12/2003
Armored Core 2	PS2	Simulation	Sony Computer Entertainment	1	10	9	249	6.26	2000 FromSoftware	Simulation;Sci-Fi;Mech	78	8.6	1 to 2	3/8/2000
ATV Offroad Fury	PS2	Racing	Activision	7	8.4	2	208	3.34	2001 RainbowStudios	Driving;Racing;Rally / Offroad	82	8.3	1 to 2	6/2/2001
ATV PHONE			<u> </u>		0.7		4.0		0000 0 11 0 1				4. 4	4014110000

Title Platform Userscore Comment This game could have been a lot better. The campaign was way too short and it didn't really explain much of the story. The multiplayer is just downright bad. Every time I play it, i remember why i stopped. The game is fun if u have a bunch of friends who play and you guys can screw around of forge. This is to Oscar G.: yes a boring history..even though this game takes place in 500 This 5.0 game could have been a lot better. The campaign was way too short and it didn't really explain much of the story. The multiplayer is just downright bad. Every time I play it, i remember why i Halo 3 Xbox360 stopped. The game is fun if u have a bunch of friends who play and you guys can screw around of forge. This is to Oscar G.: yes a boring history..even though this game takes place in 500 years. And of course the gameplay didn't change. What are you going to be killing this game? Zombies? Overall, the campaign was fun for the 2 hours and the multiplayer is way too tournament based. The Legend of This game is a masterpiece, it would definitely feature in my Top 10 favourite games of all times. Simple controls, beautiful environments, great sense of freedom and your discoveries really feel Zelda: Breath WiiU 10.0 like the result of your own work. of the Wild Every year the same. If u r fan its ok;) the line up is good but the glitches, pace and pass system is nt what we want again and again; (i would prefer every 2 years a new game and a live team that 6.0 FIFA 16 PlayStation4

keeps updating the game over the years until the next one release.

Loll'm playing on nemisis difficulty, I never paid a cent in, games brilliant and I'm super powerful now. Microtransactions didn't ruin nothing. I earned 6k mithril last night 5hr session. If you are putting money into this game you probably aren't any good at it.

One of the best of the CS series. Sure the first CS is the leader of them all, but CS:GO just puts the series onto a whole new level. It proves that Source can give amazing graphics as well as incredible gameplay, and it shows that CS will never die out. I think this is an amazing release made by VALVe and it's worth every penny.

Offensive

Middle-earth:

Strike: Global

Shadow of

War Counter-

XboxOne

PC

10.0

10.0



Predictive Modelling

Modelling

- 1. Platform
- 2. Genre
- 3. Publisher
- 4. Developer
- 5. No. of players
- 6. VG Score
- 7. Critic Score
- 8. User Score
- 9. Meta Score
- 10. Production Cost
- 11. Release Date



### **Causal Prediction Model**

Linear Regression Model

Neural Network Potential Revenue

## **Data Pre-processing**

#### **Apply Exclusion Rules**



#### **Remove Missing Values**



### **Addressed Data Inconsistencies**



#### **Group by Top 10 Developer**

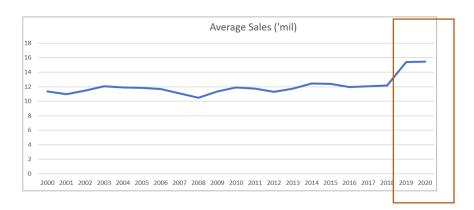


#### **Derived Time Dimension**

- There is a spike in the average sales from 2019 onwards due to Covid-19
- Remove those data points to eliminate Covid-19 impact
- 8 missing values identified
- Remove those data points due to not missing completely at random based on MCAR test
- Standardized the description for Game Developer
- E.g. EA LA; ElectronicArts .. Standardized to Electronic Arts

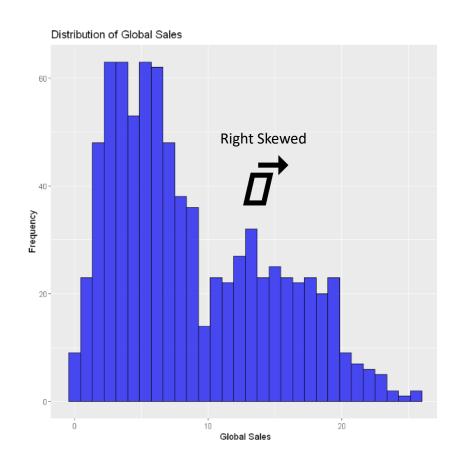
• 283 Developers had been grouped into Top 10 Developers & Others

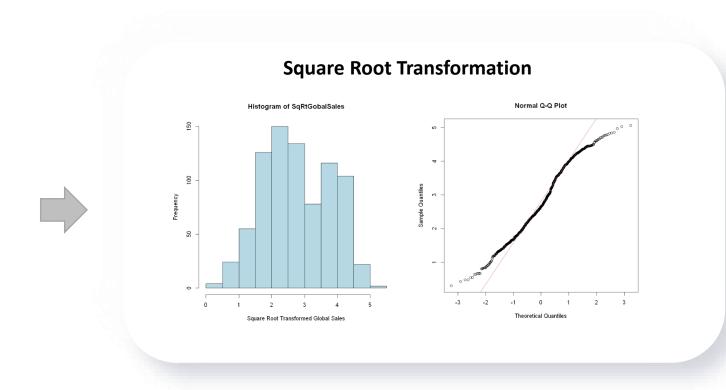
Release Year; Release Quarter; Release Month; Release Weekday; Release Weekend; Release End of Month



# **Data Pre-processing**

### **Addressed Non-Normality**





Business
Understanding

Data

Modelling

Evaluation

Retrospective

Note:

# **Linear Regression**

 Observed consistently poor R-square results across all models, ranging from 30% to 42%

Note: Please note that this simplified version of the stargazer output is provided for presentation purposes. For comprehensive and detailed information, kindly refer to the project documentation

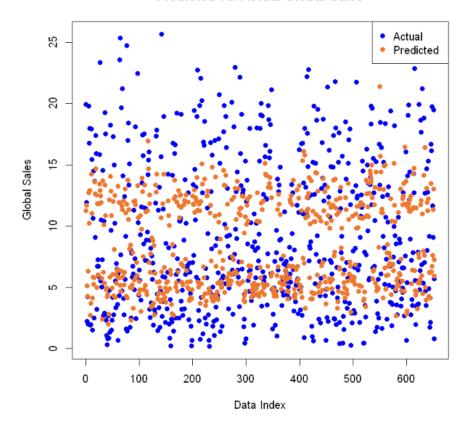
		Dependent	Variable	
		Response	Variable	
	(1)		(3)	(4)
PlatformPC	1.710*** (0.173)	1.365*** (0.165)	1.614* (0.825)	1.506* (0.854)
PlatformPS	0.900*** (0.284)	0.545* (0.281)	1.546* (0.846)	1.823** (0.863)
GenreAdventure	0.608*** (0.180)	0.550*** (0.184)	0.801 (0.631)	0.930 (0.651)
GenreSports	0.078 (0.165)	-0.001 (0.172)	1.639 (1.010)	1.467 (1.044)
PublisherActivision	0.228** (0.108)	0.182* (0.110)		0.192* (0.109)
PublisherIBM	0.310** (0.155)	0.272* (0.156)		0.363** (0.152)
CriticScore	-0.015 (0.012)	-0.011 (0.013)		
UserScore	-0.022* (0.012)	-0.017 (0.012)		
ProductionCost	-0.001 (0.019)	-0.0005 (0.020)		
ReleaseMonth	0.325* (0.195)	0.287 (0.198)		
ReleaseQuarter	-0.077 (0.130)	-0.082 (0.131)		
ReleaseWeekendWeekend	0.370 (0.237)	0.408* (0.239)	0.410** (0.205)	0.456** (0.211)
ReleaseWeek	-0.067 (0.043)	-0.057 (0.044)		
MetaScore	0.021** (0.008)	0.025*** (0.008)		
NoOfPlayers1	-0.347 (0.856)	-0.360 (0.865)		
TopDeveloper_SumNamco	1.046*** (0.354)		1.190*** (0.360)	
TopDeveloper_SumNintendo	2.566*** (0.446)		2.389*** (0.451)	
TopDeveloper_CountNamco		1.236*** (0.388)		
TopDeveloper_CountSportsInteracti	ve	1.019** (0.411)		
PlatformPS3:GenrePuzzle			2.558** (1.073)	2.306** (1.106)
PlatformPS:GenreSports			-2.187* (1.219)	-2.237* (1.255)
Constant	-11.673 (19.707)	-15.346 (19.885)	0.666 (0.629)	1.261** (0.603)
Observations	652	652	652	652
R2 Adjusted R2	0.403 0.357	0.391 0.344	0.419 0.368	0.382 0.331

\*p<0.1; \*\*p<0.05; \*\*\*p<0.01

### **Actual vs Predicted Global Sales**

Model 3

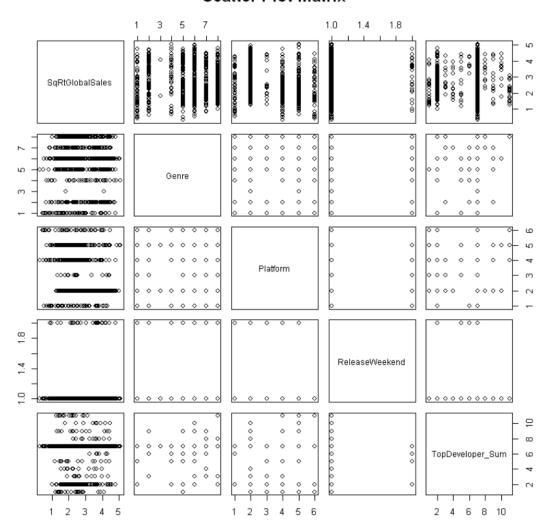
#### Predicted vs. Actual Global Sales

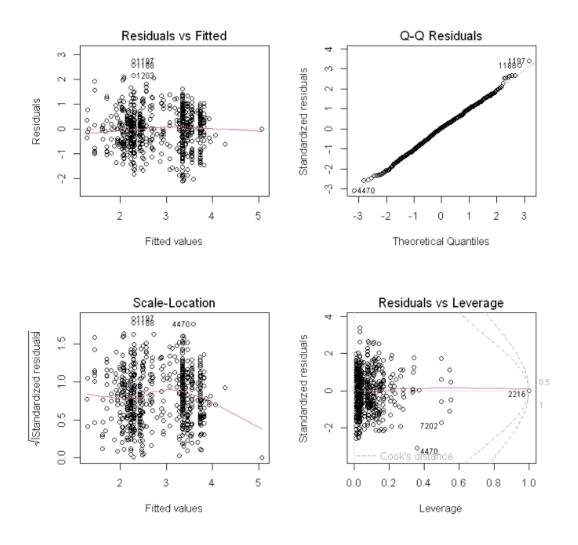


	Training Dataset (80%)	Testing Dataset (20%)
Mean Absolute Deviation	3.55	3.60
Mean Squared Error	21.56	22.30
Standard Error	4.64	4.72

Business Data Modelling Evaluation Retrospects

#### **Scatter Plot Matrix**

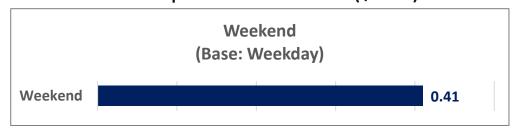


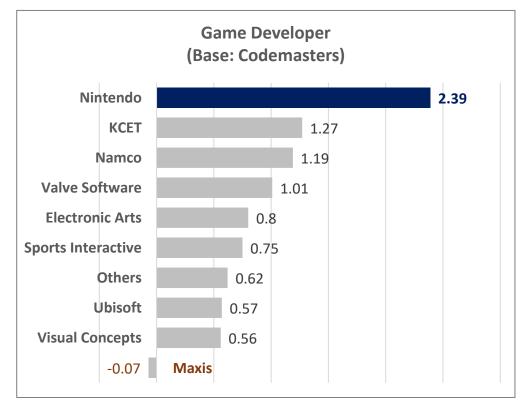


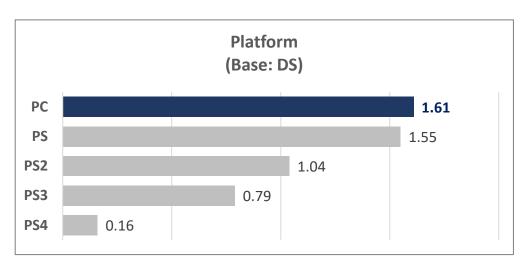
Business Data Modelling Evaluation Retrospective

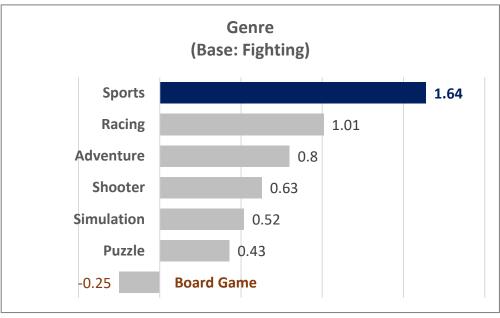
# **Key Revenue Contributors**

### Net impact on Revenue (\$'mil)









### **Neural Network**

Normalized all the **Neural Network Evaluate the** Revenue input variables Model **Models Dataset** Model 2 Model 1 **Architecture:**  Activation Function: Relu 1 Output 1 Output • Loss Function: Mean squared Variable Variable 2 Nodes 3 Nodes error 5 Nodes 4 Nodes **Training Configuration:** • Epochs: 100 16 Input 16 Input Variables **Variables** 

# **Models Comparison**

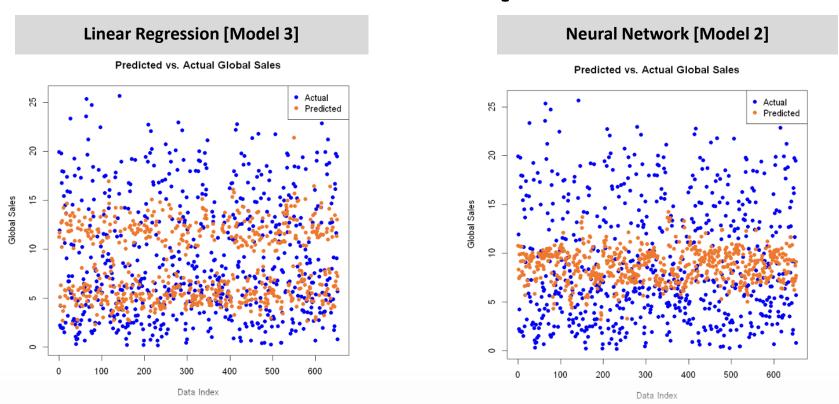
			el 1]	Neural Network [Model 2] Training Testing	
Training	Testing	Training	Testing	Training	Testing
3.67	3.66	5.11	4.91	5.10	4.98
22.87	22.52	37.03	32.84	36.97	33.98
4.78	4.75	6.09	5.73	6.08	5.83
	3.67 22.87	3.67 3.66 22.87 22.52	3.67       3.66       5.11         22.87       22.52       37.03	3.67       3.66       5.11       4.91         22.87       22.52       37.03       32.84	3.67     3.66     5.11     4.91     5.10       22.87     22.52     37.03     32.84     36.97

### **Linear Regression Outperforms Neural Network:**

- Highly likely attributed to the limited dataset size, which consists of around 815 data points.
- Limited dataset restricted the complexity of the neural network model, allowing only a simple architecture with few nodes and layers to be set up.

Business Data Modelling Evaluation Retrospective

## **Models Comparison**



- Both Models Underfitting: Indicating limitations in capturing the complexity of the relationship within the data
- Consider Categorical Prediction Models: These models aim to predict the hit or miss category rather than attempting to predict a precise point estimate of revenue

# Classification model

"Hit" or "Miss"

### **Model variables**

### Input variables

- Numerical:
  - Critic score, User score, Production Cost
- Categorical:
  - Genre, platform, publisher, number of players, Console, Developer
- Output from Text Mining TDM:
  - Features and improvements in beta players' user comments

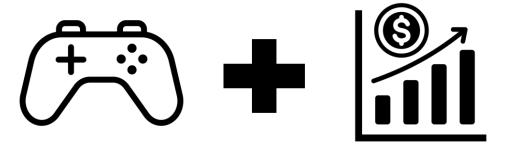
### **Output variables**

• Game's Hit or Miss (Based on Net profit)





# **Data Preparation**





Inner Join game sales dataset with user review dataset

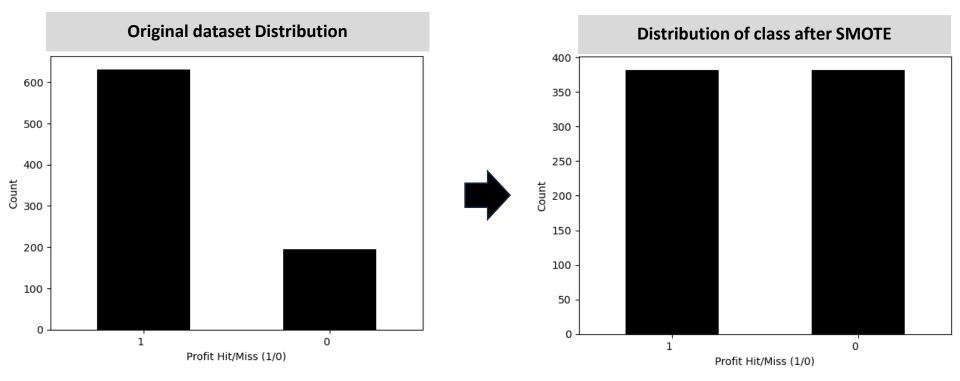
Derived "Game hit/miss" column

platform	genre	publisher	Developer_metacritic_game_info
PC	Simulation	Sony Computer Entertainment	Related Designs
PS3	Sports	Ubisoft	EACanada
PC	Adventure	Activision	GiantSquid
PS2	Sports	Activision	Z-Axis,Ltd.
DS	Adventure	Electronic Arts	Etranges Libellules
PS	Adventure	Electronic Arts	Darkworks
DS	Simulation	Amazon Game Studios	Nintendo
DS	Simulation	Ubisoft	SunflowersInteractive
DS	Adventure	Ubisoft	Capcom
PS2	Puzzle	Frontier Developments	ZedTwoLimited
PC	Adventure	Activision	BitBlot



	genre_Board Game	genre_Fighting	genre_Puzzle	genre_Racing	genre_Shooter	genre_Simulation	genre_Sports	publisher_Amazon Game Studios
	0	0	0	0	0	1	0	0
	0	0	0	0	0	0	1	0
	0	0	0	0	0	0	1	0
•	0	0	0	0	0	0	0	0
	0	0	0	0	0	0	0	0
	0	0	0	0	0	0	1	0
	0	0	0	0	0	0	0	0

# **Data Imbalance Handling**



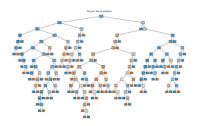
- Dataset exhibits imbalance with the profit hit/miss
  - Class 1: Profit 76%
  - Class 2: Loss 24%
- Applied **SMOTE** to training set to generate synthetic samples for minority class
- Applied train test split at 80:20 rule

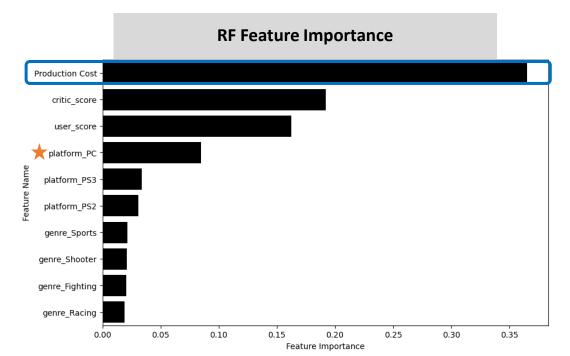
Business Data Modelling Evaluation Retrospects

## **Parameters & Tuning Factor**

### Random Forest

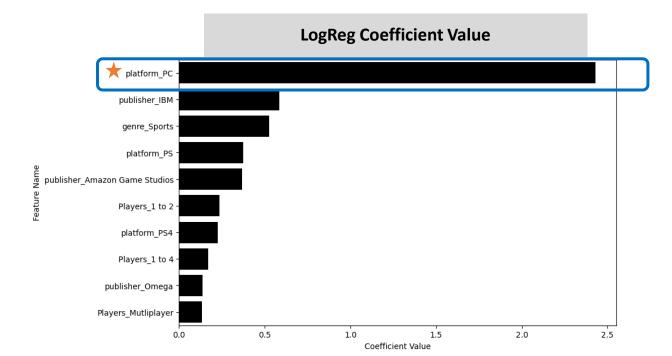
- Hyperparameter Tuning: Grid search
- Features selected: Genre, Platform, critic score, user score, Production cost





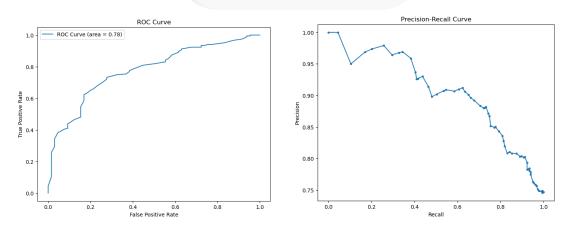
### Logistic regression

- Maximum number of iterations: 1000
- Features selected: Platform, No. of players, critic score, user score, Production cost

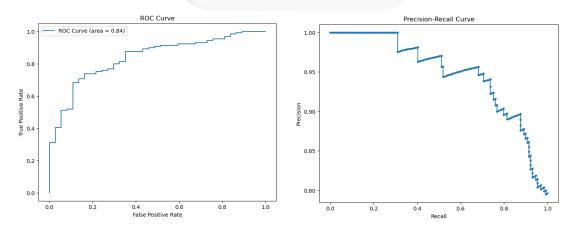


### **Model Performance**





### Logistic regression



	Random Forest	Logistic regression	Neural Network	SVM
Accuracy	0.76	0.81	0.72	0.73
ROC AUC	0.78	0.84	0.56	0.71
Cross Entropy	0.48	0.43	9.57	0.46

- LogReg outperforms RF in accuracy and ROC AUC.
- Lower cross-entropy suggests a higher level of confidence in its predictions.

## **Next steps**



Further refine Logistic Regression model.



Evaluate the interpretability of the model and ensure it aligns with business needs.



Leveraging domain knowledge in the gaming industry dynamics



Adopt an iterative approach to model improvement



Text Analytics

# Text Analytics

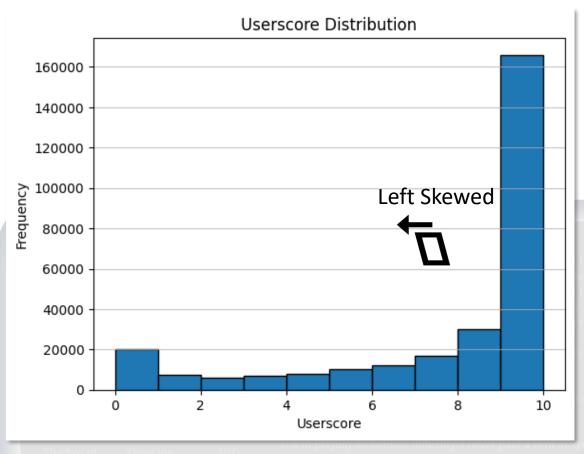
- "How would you rate the stated game, played via this stated platform, on a Likert scale of 0 to 10?"
- 2 "Please justify why you gave this rating score."

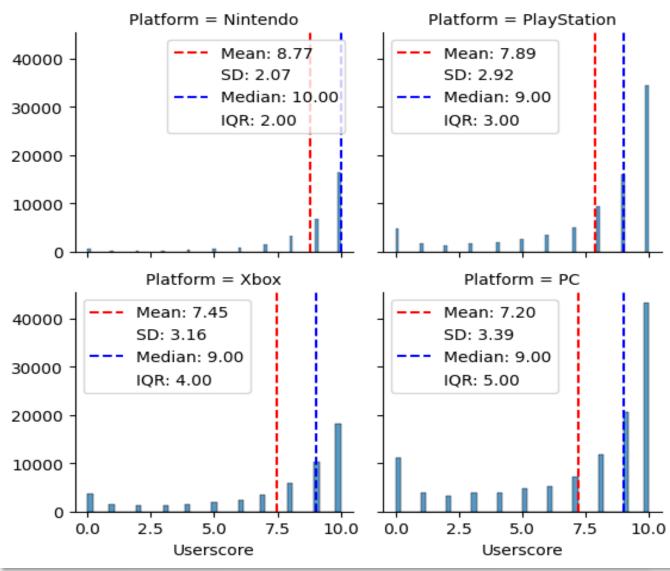
### Objective:

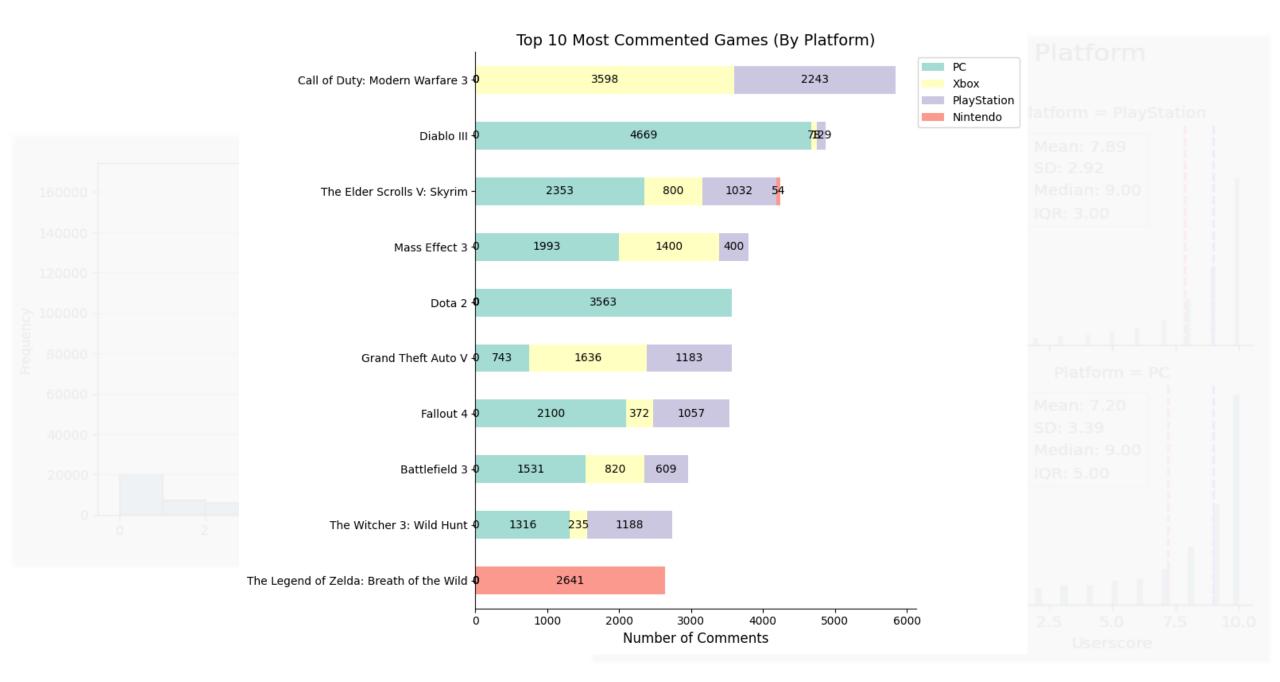
To use these two questions to find out 3 to 5 key 'features' and main 'improvements' that players look for in Games per Gaming Platform

orm Userscore Com	Userscore	Platform	Title
This game could have been a lot better. The campaign was way too short and it didn't really explain much of the story. The multiplayer is just downright bad. Every time I play it, i remember stopped. The game is fun if u have a bunch of friends who play and you guys can screw around of forge. This is to Oscar G.: yes a boring historyeven though this game takes place in 500 game could have been a lot better. The campaign was way too short and it didn't really explain much of the story. The multiplayer is just downright bad. Every time I play it, i remember stopped. The game is fun if u have a bunch of friends who play and you guys can screw around of forge. This is to Oscar G.: yes a boring historyeven though this game takes place in 500 yand of course the gameplay didn't change. What are you going to be killing this game? Zombies? Overall, the campaign was fun for the 2 hours and the multiplayer is way too tournament be	5.0	Xbox360	Halo 3
WiiU 10.0 This game is a masterpiece, it would definitely feature in my Top 10 favourite games of all times. Simple controls, beautiful environments, great sense of freedom and your discoveries really like the result of your own versions.	10.0	WiiU	The Legend of Zelda: Breath of the Wild
Every year the same.If u r fan its ok ;) the line up is good but the glitches ,pace and pass system is nt what we want again and again :( i would prefer every 2 years a new game and a live team keeps updating the game over the years until the next one rel	6.0	PlayStation4	FIFA 16
One 10.0 Loll'm playing on nemisis difficulty, I never paid a cent in, games brilliant and I'm super powerful now.Microtransactions didn't ruin nothing.I earned 6k mithril last night 5hr session.If yo putting money into this game you probably aren't any good	10.0	XboxOne	Middle-earth: Shadow of War
PC 10.0 One of the best of the CS series. Sure the first CS is the leader of them all, but CS:GO just puts the series onto a whole new level. It proves that Source can give amazing graphics as w incredible gameplay, and it shows that CS will never die out. I think this is an amazing release made by VALVe and it's worth every personal contents.	10.0	PC	Counter- Strike: Global Offensive
One of the best of the CS series. Sure the first CS is the leader of them all, but CS:GO just puts the series onto a whole new level. It proves that Source can give amazing graphics as w incredible gameplay, and it shows that CS will never die out. I think this is an amazing release made by VALVe and it's worth every possible.	10.0	ЬС	Counter- Strike: Global Offensive

### Userscore Distribution by Platform







Q&A method

Topic Modelling method POS
Tagging
and TF-IDF
method
(chosen)

- 1. 'What are the best features of the game?'
- 2. 'What are the main improvements needed for the game?'

#### Plan:

Test if the sentiments drawn from both questions are different. If yes, it means that Q&A method is able to discern key features vs. main improvements well and may be used. Else, cannot be used.

- 1. 'What are the best features of the game?'
- 2. 'What are the main improvements needed for the game?'

### **Key Features**



#### Plan:

Test if the sentiments drawn from both questions are different. If yes, it means that Q&A method is able to discern key features vs. main improvements well and may be used. Else, cannot be used.

### Improvement Needed









#### Plan:

If the words in each group points to a consistent topic, we may use that topic as the 'features' (for high score) or 'improvements' (for low score). Else, cannot be used.

```
High Score:
[(0, '0.002*"gamecube" + 0.002*"metroid" + 0.002*"owns" + 0.002*"decade" + 0.002*"fifa" + 0.002*"1.6" + 0.001*"soccer" + 0.001*"sexy" + 0.001*"pe" + 0.001*"B"'),
(1, '0.005*"de" + 0.004*"e" + 0.004*"sweet" + 0.004*"que" + 0.003*"un" + 0.002*"la" + 0.002*"juego" + 0.002*"comment" + 0.002*"en" + 0.002*"jogo"'),
(2 '0.006*"best") + 0.005*"ever" + 0.004*"played" + 0.004*"good" + 0.004*"story" + 0.004*"like" + 0.003*"graphic" + 0.003*"amazing + 0.003*"really" + 0.003*"get"')]
Low Score:
[(0, '0.004*"volvo" + 0.002*"B" + 0.002*"hat" + 0.002*"hat" + 0.001*"wer" + 0.001*"a" + 0.001*"a" + 0.001*"halloween" + 0.001*"muκροτραμβακιμи"'),
(1, '0.003*"like + 0.003*"good") + 0.002*"money" + 0.002*"get" + 0.002*"give" + 0.002*"even" + 0.002*"player" + 0.002*"story" + 0.002*"boring"),
(2, '0.003*"pokemon" + 0.002*"uplay" + 0.002*"blade" + 0.001*"batmobile + 0.001*"blacklist" + 0.001*"pet" + 0.000*"washing" + 0.000*"batman" + 0.000*"smashing" + 0.000*"spirit"')]
```

Topic Modelling

```
filter_list = ["would", "could", "left", "right", "a.m.", "p.m.", "'s", "! ! ! !", "...", ":", ";", "n't",
                game", "games", "play", "fun", "much", "one", "great", "perfect", "time", "year", "lot", "thing", "etc","
               "hour", "hours", "ways", "ways", "everything", "anything", "thing", "things", "review", "reviews", "years", "years",
               "feel", "feels", "thing", "nothing", "problem", "end", "begin", "kind", "piece", "work", "call", "anyone",
               "minute", "minutes", "waste", "crap", "garbage", "masterpiece"]
def preprocess text(tokens, needtokenizeBoolean = True, grams = False, ngramsNumber = 2, furtherPreProcessNgrams = False):
   if needtokenizeBoolean:
        tokens = nltk.word tokenize(tokens)
            tokens = list(ngrams(tokens, ngramsNumber))
   if grams:
        tokens = [' '.join(gram) for gram in tokens]
       if furtherPreProcessNgrams == False:
           return tokens
   tokens = [t.lower() for t in tokens]
   tokens = [t for t in tokens if t not in stopwords.words('english') + filter list]
   tokens = [t for t in tokens if t not in string.punctuation]
   tokens = [t for t in tokens if not t.isnumeric()]
   tokens = [nltk.WordNetLemmatizer().lemmatize(w) for w in tokens]
   return tokens
 odf in userCommentsTESTExtreme_list:
  for platform in platformCondensed list:
         if df.scoreBin.max() == 'High':
           features improvements = 'Key Features'
        elif df.scoreBin.max() == 'Low':
           features_improvements = 'Improvements Needed'
            wc(df = df[df.platformCondensed == platform], columnName = 'Comment',
            preProcessingFunctionBoolean = True, vectorizerMinDf = 2, vectorizerMaxDf = 0.7,
            countVectorizerBinary = True, tfidfVectorizerBoolean = True,
            ngrams = False, ngramsNumber = 3, furtherPreProcessNgrams = False,
            nounTaggingBoolean = True, universalNounTagsetBoolean = False,
            top = 20, features_improvements = features_improvements, platform = platform, titleFontSize = 20)
            print(f"\n\nNo-Word-Cloud for '{features improvements}'({platform})' due to insufficient sample size (No. of comments = 0 or < vectorizerMinDf).\n\n")
```

```
nounTaggingBoolean = False, universalNounTagsetBoolean = False, ngrams = False, ngramsNumber = 2, furtherPreProcessNgrams = False top = 10, features improvements = 'Nord Cloud', platform = 'All', titleFontSize = 20):
 f preProcessingFunctionBoolean -- True & ngrams -- True
         If furtherPreProcessNgrams -- True
                                  min df - vectorizerMinDf max df - vectorizerMaxDf
           tdm = TfidfVectorizer(tokenizer = lambda text: preprocess_text(text, grams = True, ngramsNumber-ngramsNumber
min df = vectorizerMinOf, max df = vectorizerMaxOf)
        If furtherPreProcessignams == Irue:

tdm - CountVectorizer(binary = countVectorizerBinary, tokenizer = lambda text: preprocess text(text, grams = True, ngramsNumber-ngramsNumber, furtherPreProcessignams = True
            tdm = CountVectorizer(binary = countVectorizerBinary, tokenizer = lambda text; preprocess_text(text, grams = True, ngramsNumber=ngramsNumber)
    min_df = vectorizerMinOf, max_df = vectorizerMaxOf)
tdmMatrix = tdm.fit_transform(df[columnName])
  f nounTaggingBoolean == True and ngrams == False
            noun = [word for word, tag in tagged_value if tag == 'NOUN']
           tagged_value = pos_tag(word_tokenize(value))
noun = [word for word, tag in tagged_value if tag == 'NN' or tag == 'NNS']
       df['Text_NounOnly'] = df['Text_NounOnly'].apply(lambda x: preprocess_text(tokens = x, needtokenizeBoolean = False))
df['Text_NounOnly'] = df['Text_NounOnly'].apply(lambda x: ', '.join(x))
       tdm = CountVectorizer(binary = countVectorizerBinary, min_df = vectorizerMinDf, max_df = vectorizerMaxDf)
    tdmMatrix = tdm.fit transform(dff Text NounOnly
if nounTaggingBoolean != True and ngrams == False
     if tfidfVectorizerBoolean ==
               tdm = TfidfVectorizer(tokenizer = preprocess_text, min_df = vectorizerMinDf, max_df = vectorizerMaxDf)
               tdm = TfidfVectorizer(min df = vectorizerMinDf, max df = vectorizerMaxDf)
         if preProcessingFunctionBoolean == True:
              tdm = CountVectorizer(binary = countVectorizerBinary, tokenizer = preprocess_text, min_df = vectorizerMinDf, max_df = vectorizerMaxDf)
              tdm = CountVectorizer(binary = countVectorizerBinary, min df = vectorizerMinDf, max df = vectorizerMaxDf)
     tdmMatrix = tdm.fit transform(df[columnName])
array = tdmMatrix.toarray()
feature_names = tdm.get_feature_names_out()
word = dict(zip(feature_names, array.sum(axis=0)))
wc = WordCloud(background_color="white").generate_from_frequencies(fd)
plt.figure()
plt.suptitle(f"{features_improvements} ((platform))", fontsize = titleFontSize, x = 0.5, y = 0.85, fontweight = 'bold', fontname = 'Calibri')
plt.imshow(wc. interpolation='bilinear')
plt.axis("off"
displayList = []
for x,y in fd.most_common(top):
```

# Chosen plan:

### **Preprocessing**

- capitalization
- stop words
- custom filter list
- punctuation
- numeric removals

Unigram
Lemmatization
POS Tagging for Nouns
TF-IDF

Business Data

#### **Key Features (Nintendo)**



#### **Key Features (PlayStation)**



#### **Key Features (PC)**



#### **Key Features (Xbox)**



#### **Key Features (Others)**



#### Evaluation Retrosp

#### **Key Features (All platforms)**



### Similarities:

Graphic, Story (storyline), Character (character development)

### Differences:

Music in *Nintendo,*Multiplayer in *XBox* 



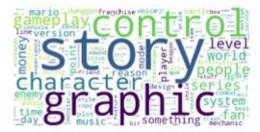






Business Data

#### Improvements Needed (Nintendo)



#### Improvements Needed (PlayStation)



#### Improvements Needed (PC)



#### Improvements Needed (Xbox)



#### Improvements Needed (Others)



#### Evaluation Retro

#### Improvements Needed (All platforms)



### Similarities:

Money, Graphic, Story (storyline)

### Differences:

Level (levelling difficulty) in *Nintendo,*Map (limited world map) in *Xbox,*Camera (POV camera angle) in *Others* 









# Applying back to business objective

- Information from Word Cloud:
  - Better advisory to clients
    - Recommend must-have key features and pitfalls (improvements) to avoid based on gaming platform type
    - Appeal to gamers

# Conclusion



Business Data Modelling Evaluation Retrospectiv

# Model



## 1. Logistic Regression \*\*\*

- ROC and AUC
- Business Use Cases (Sales Revenue vs Hit/Miss)
- Availability and Relevancy of Data



# 2. POS Tagging

- Term Frequency Distribution
- Critical Game attributes
  - Character
  - Gameplay
  - Story



Business Data Modelling Evaluation Retrospective anderstanding

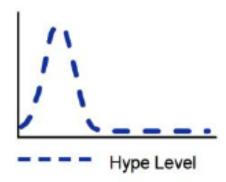
# Retrospective

### 1. Imbalance Dataset

• SMOTE, under sampling etc

### 2. Sales Revenue Data

- Assumption games sales tapering (older vs newer games)
- Gartner Hype Cycle





## 1. Gaming Industry

- Dynamic and fast product sale lifecycles
- Changing consumer preferences

# 2. Respondent Language

Languages other than English in review comments

### 3. Limited usable data size

Not adequate for a NN model



### 1. Enhancing Sentiment Analysis

- Positive vs Negative Categorization
- Eliciting useful response thru targeted questions

## 2. Language Handling

- Respondent Demographics (российский & español)
- Google Translate API

## 3. Potential Model Improvements

Gradient Boosting / Ensemble model

