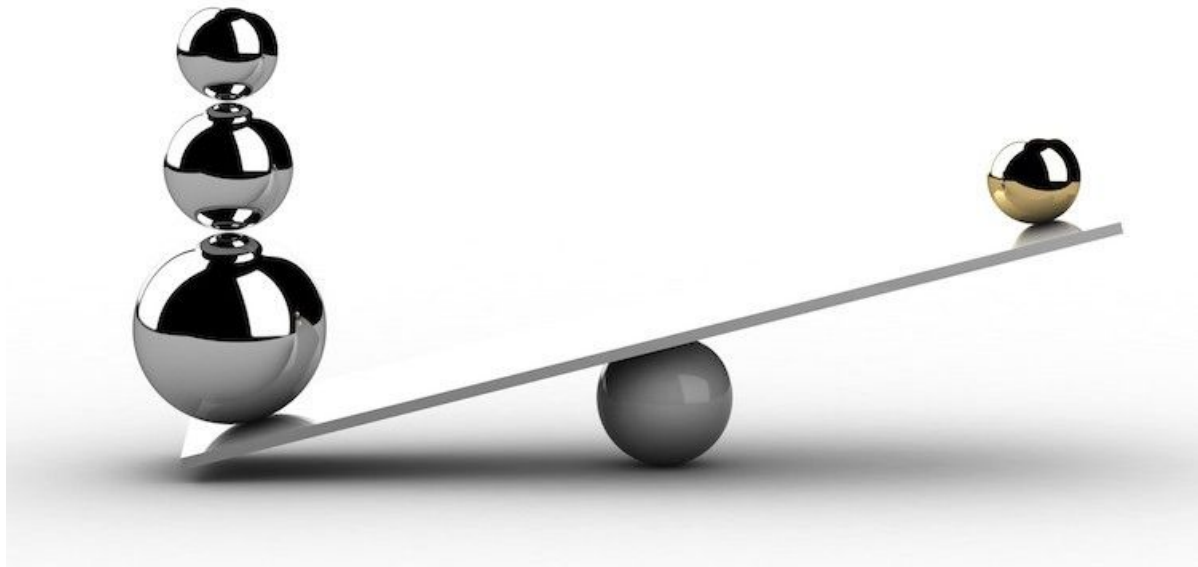
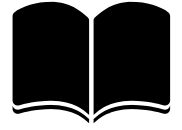


# **Generative Adversarial Networks (GANs) for data imbalance**

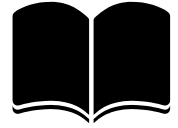


**Armando Vieira & Javier Rodriguez Zaurin**

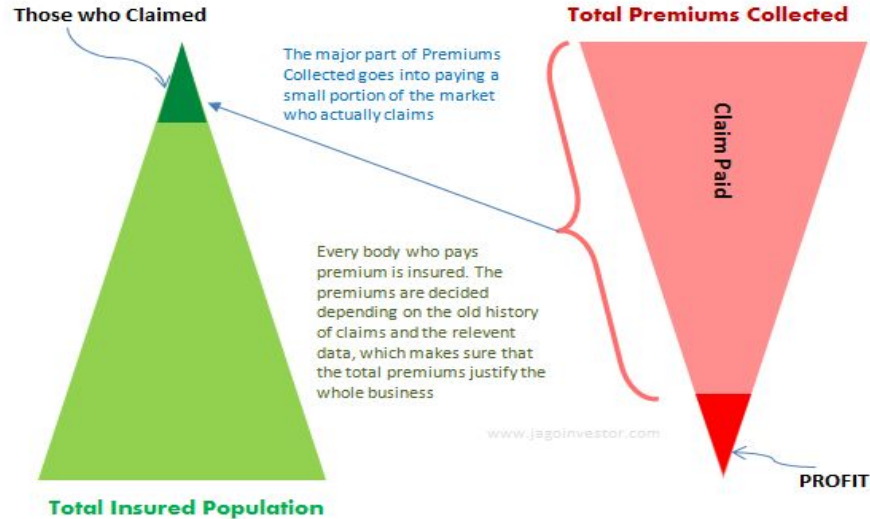
# Data imbalance in insurance - why it matters?



# Data imbalance in insurance - why it matters?



How any Insurance Business Works !



Disclaimer - The Above chart explains an average year. While the above situation might not be true for a particular year. Over a very long term, this is how the insurance business depends on for success. The overall success would depend on various factors like premium pricing, claim handling and other many factors which cant be discussed here.

Created by : [www.jagoinvestor.com](http://www.jagoinvestor.com)

1.

# **Myths Around Insurance modeling**

Statistics are there for a reason...



- Data follows **Poisson** / overdispersed Poisson / Gamma distributions
- There are **well defined quantities** we can always compute
  - mean
  - standard deviation
  - skewness
- Data is **stationary** (i.e. The rate at which events occur is constant)
- Events are **independent** (i.e. The occurrence of one event does not affect the probability that a second event will occur)
- **Pooling** solves the issues

**2.**

**The mean is  
meaningless**



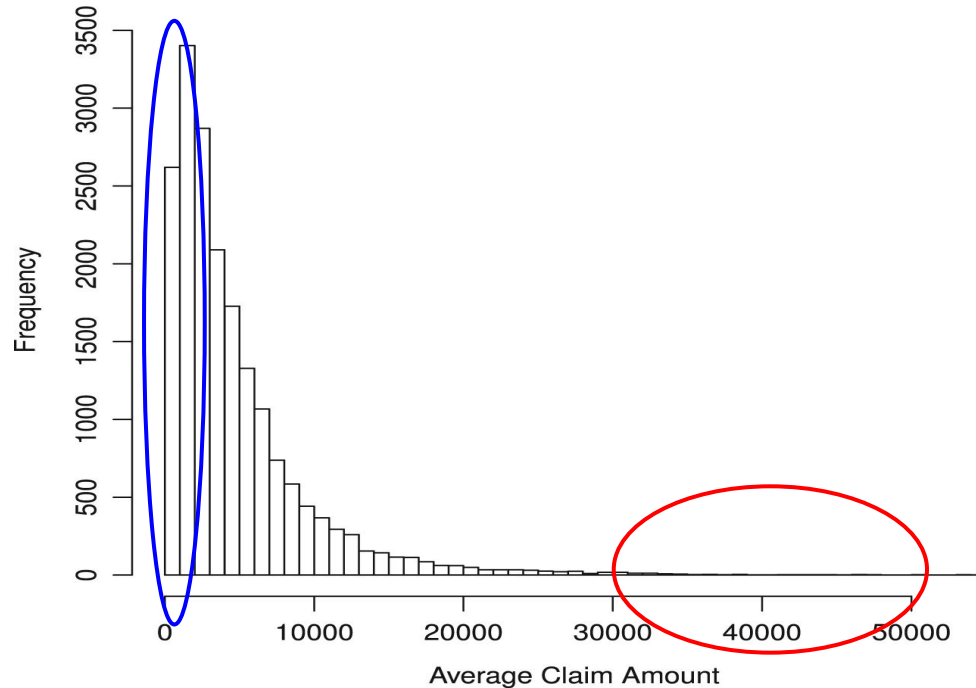
# In Insurance only two types of clients matter

- ❑ Clients that **never** claims
- ❑ Clients with **large** claims

# By focusing on the extremes we can better model the important factors



Severity Histogram





**3.**

**Imbalance  
Datasets**

# Techniques to handle imbalance data



(a) SMOTE



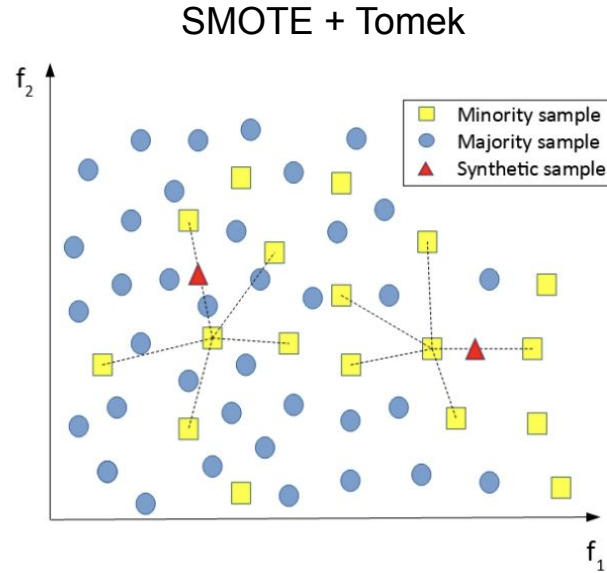
(b) ADASYN

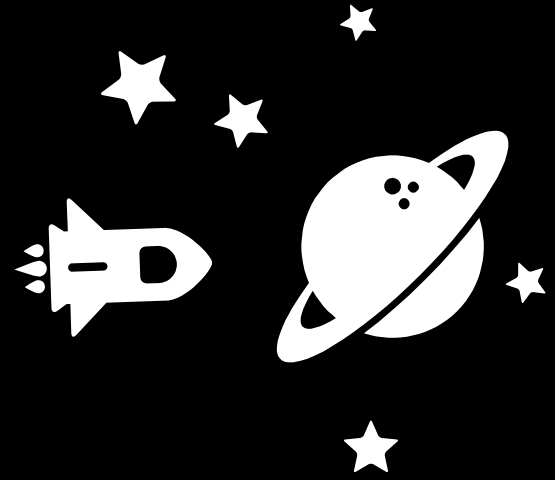


(c) VAE based Method



(d) GAN based Method

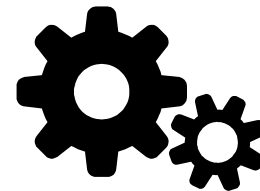




# GANs

Generative Adversarial Networks

# GANs



## 3. Deep Learning

### Deep Learning Achieves Photorealistic Image Generation

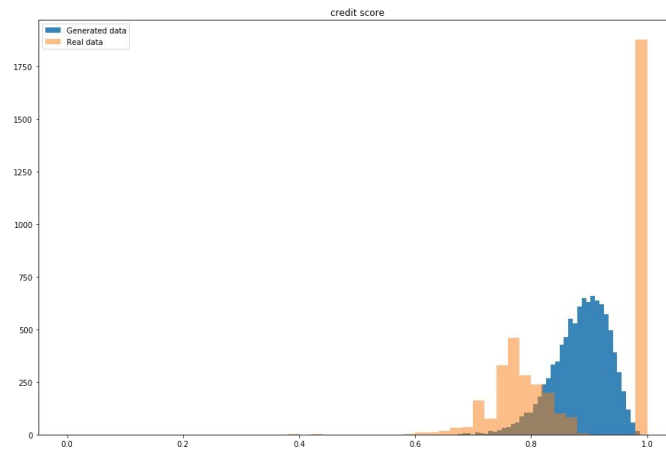
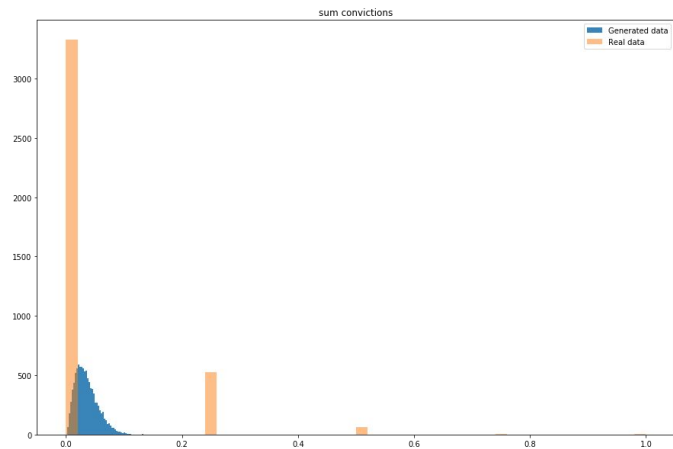
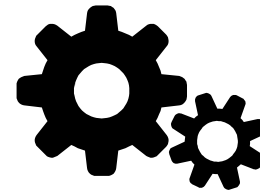


Deep learning can recognize and generate images. Early results were blurry and unconvincing, as seen on the left. The latest results approach photorealism, as seen on the right.

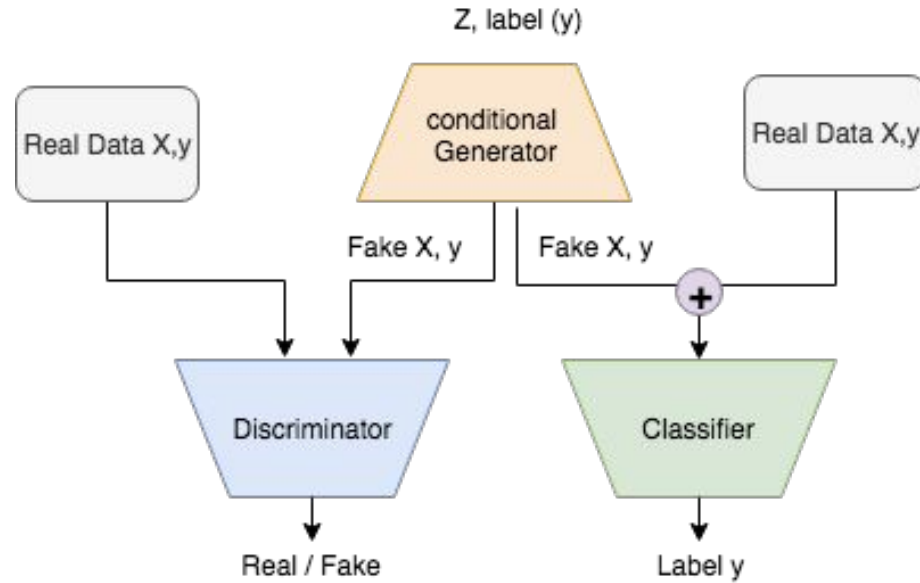
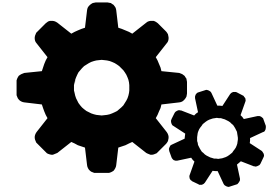
Fake Images Generated Using Deep Learning



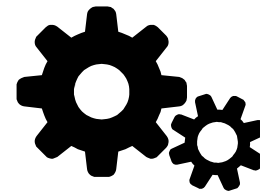
# GANs were very effective for images, but how about tabular data?



# Conditional GAN

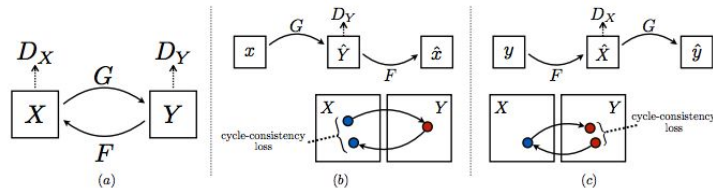
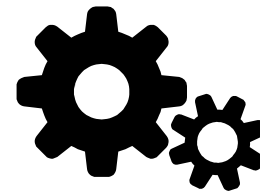


# Conditional GAN



0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9
0	1	2	3	4	5	6	7	8	9

# CycleGAN



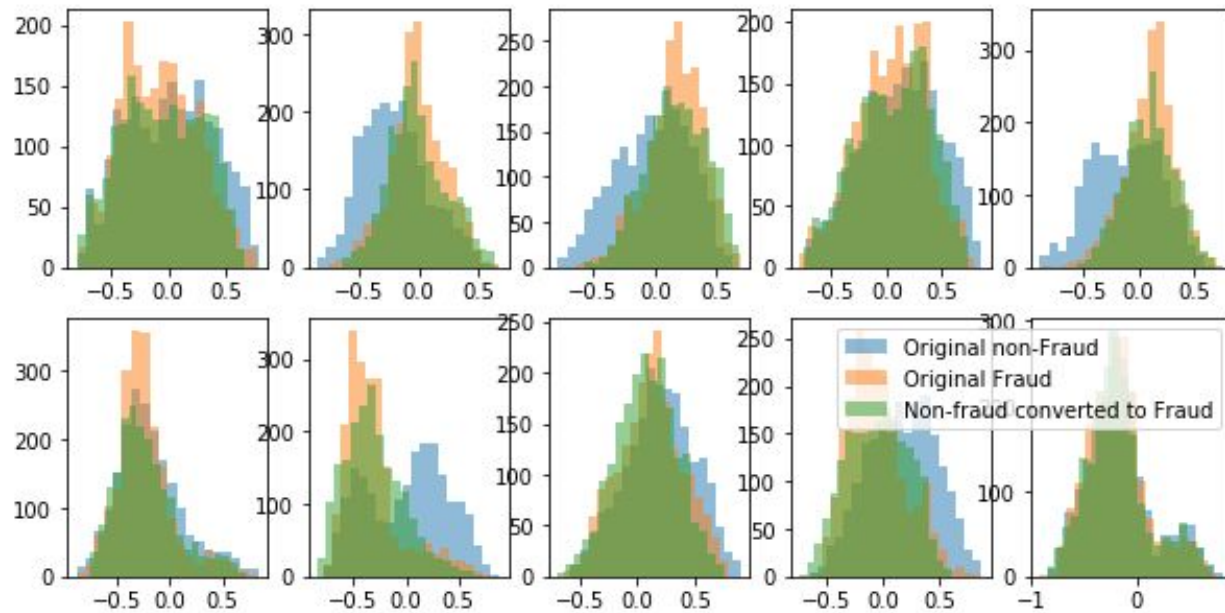
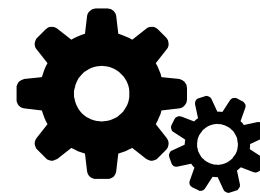
**Horse**  $\longrightarrow$  **Zebra**

**or ...**

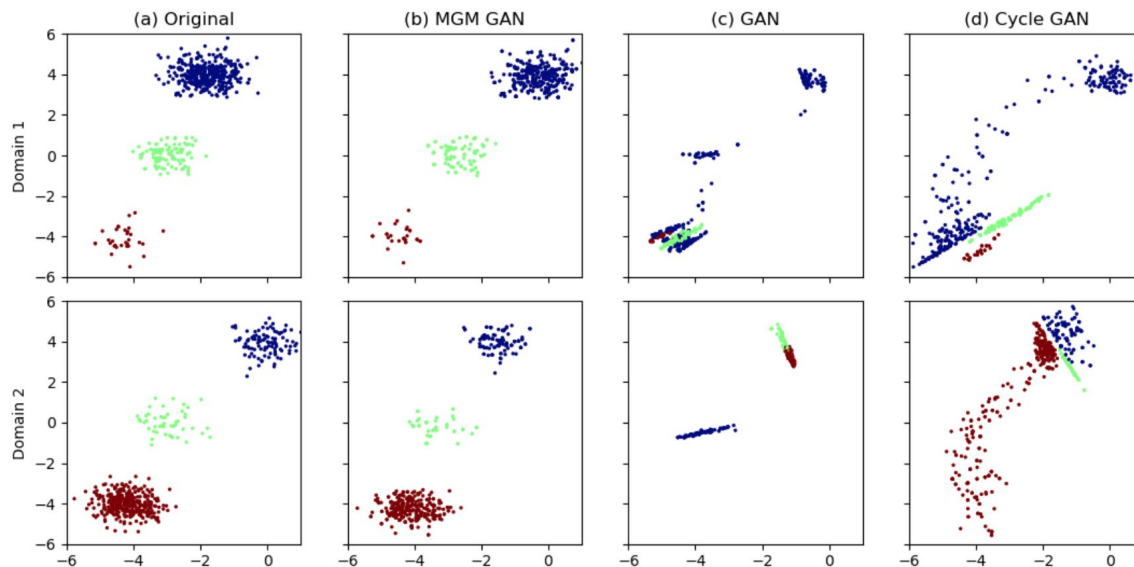
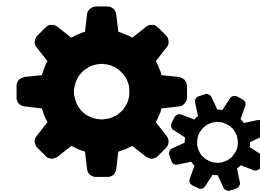
**No Claim**  $\longrightarrow$  **Claim**



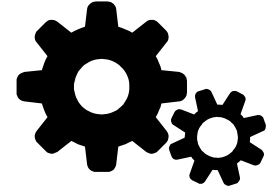
# CycleGAN for fraud detection + sparse Autoencoders



# Topology preserving CycleGAN

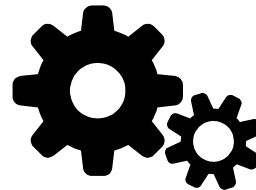


# Results



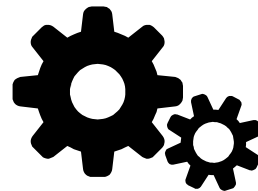
	Precision	Recall	F1-score
No augmentation	0.88	0.77	0.79
SMOTE	0.94	0.79	0.85
ADASYN	0.79	0.76	0.77
cGAN	0.90	0.85	0.87
cycleGAN	0.92	0.83	0.87

# Implementation Details



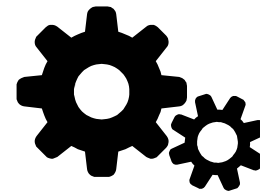
- Data has to be converted into a dense representation
- Domains have to be consistently related - retain same semantic features.
- Loads of tricks (weight clipping, regularizers, training schedule)
- Consistency over time (i.e. stationary)
- Avoid mode collapsing
- Diversity enforcing with conditional batch normalization and noise

# How to train GANs



- Small (decreasing) learning rate
- Different optimizers for encoder and generator
- Train generator multiple times for each discriminator
- Batch normalization and Instance Batch Normalization

# Conclusions:



- Focus on the **extremes** (less data, better results)
- Use **augmentation** to alleviate imbalance data
- **GANS can be effective** under some conditions