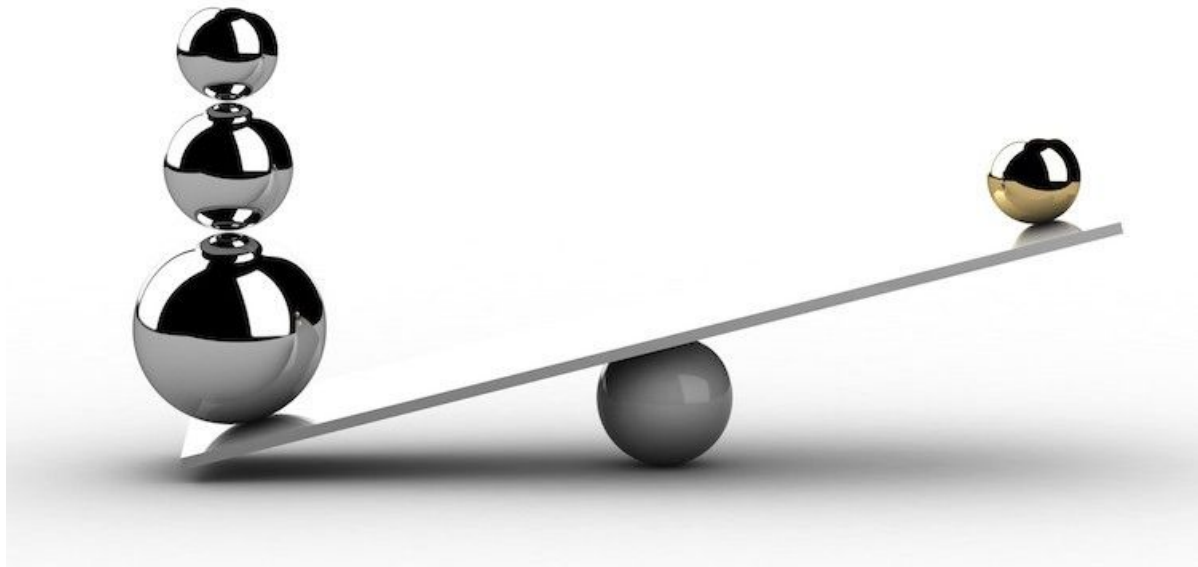
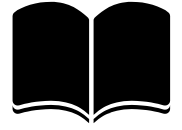


# **Generative Adversarial Networks (GANs) for data imbalance**

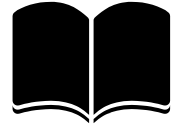


**Armando Vieira & Javier Rodriguez Zaurin**

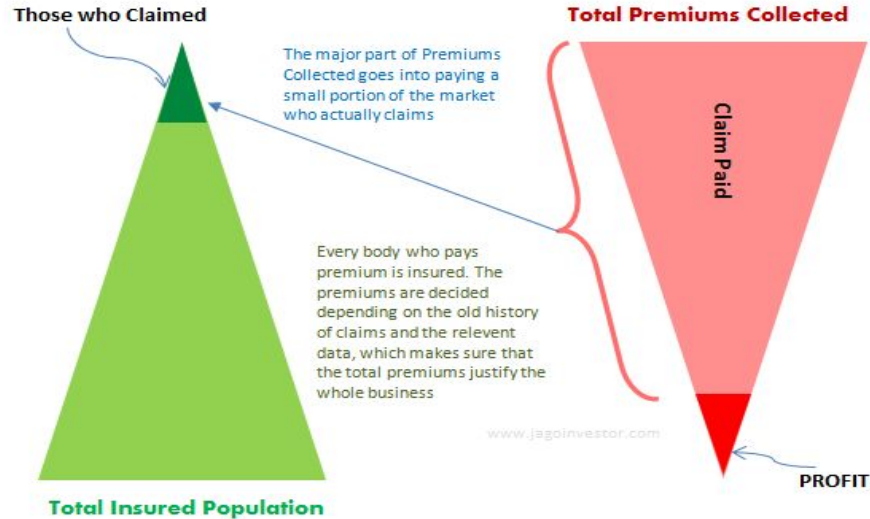
# Data imbalance in insurance - why it matters?



# 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**



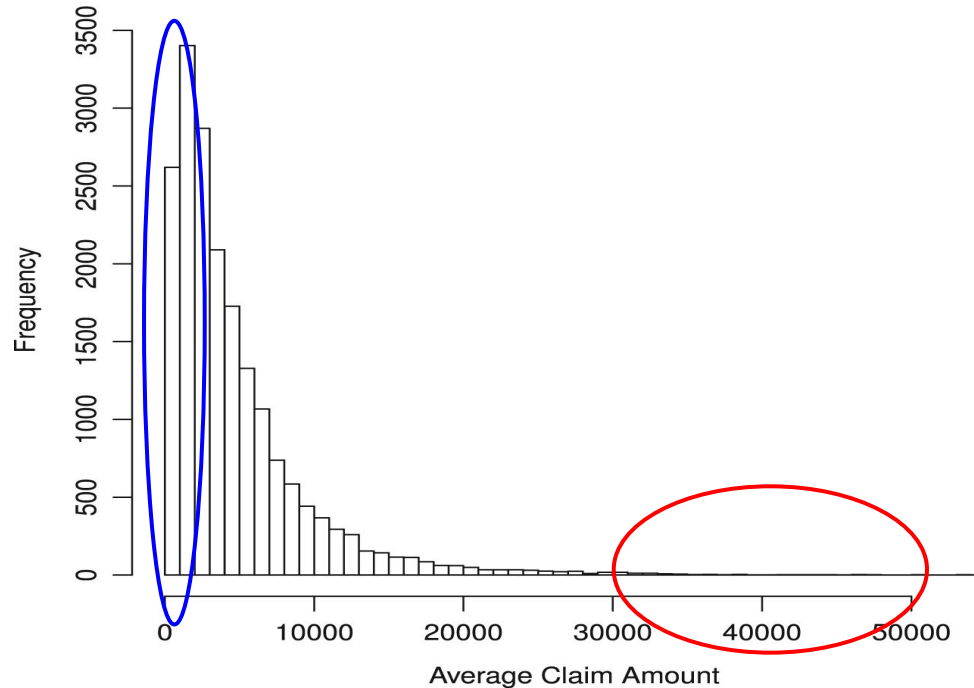
# In Insurance only two types of clients matter

- ❑ 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



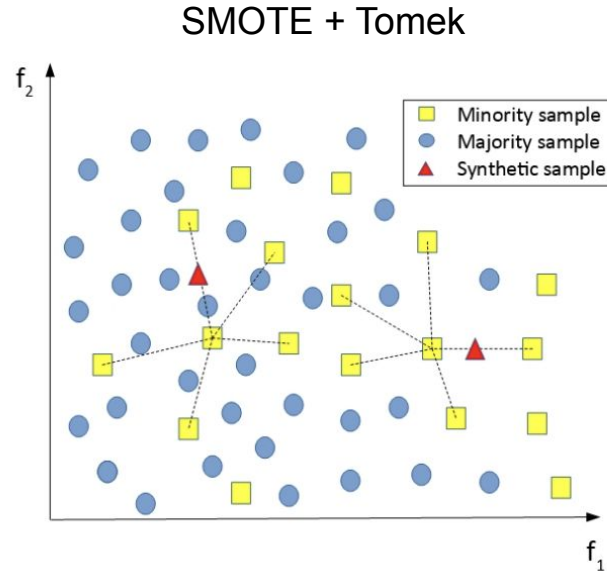
(b) ADASYN

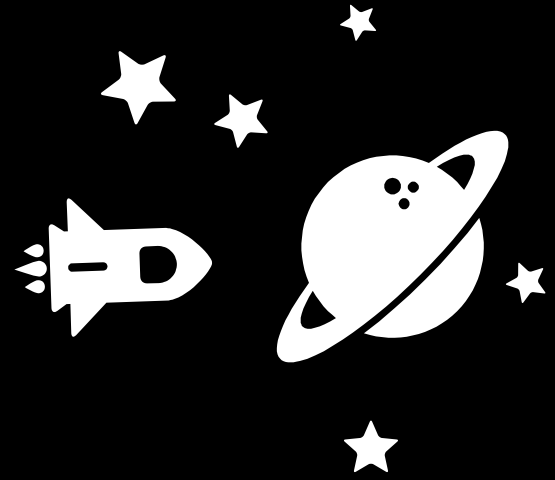


(c) VAE based Method



(d) GAN based Method

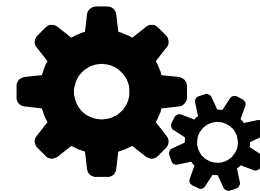




# GANs

Generative Adversarial Networks

# GANs



## 3. Deep Learning

### Deep Learning Achieves Photorealistic Image Generation

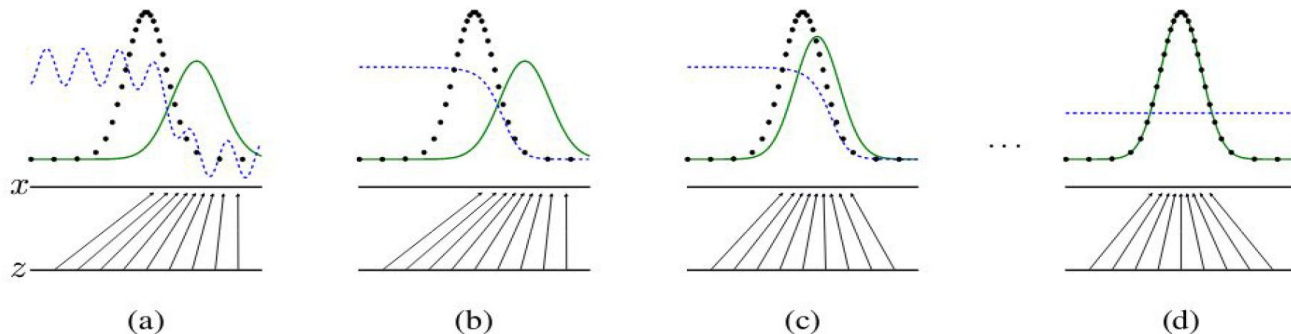


Deep learning can recognize and generate images. Early results were blurry and unconvincing, as seen on the left. The latest results approach photorealism, as seen on the right.

Fake Images Generated Using Deep Learning

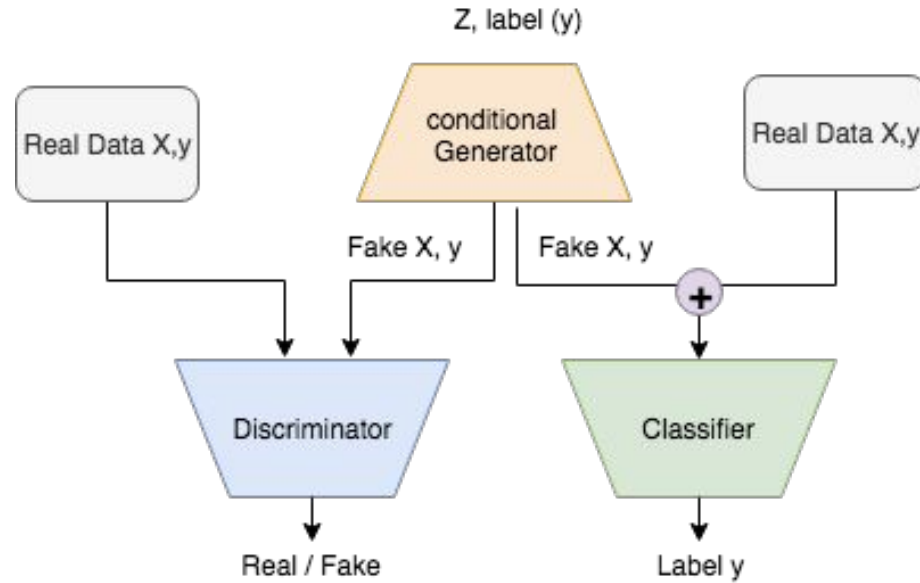
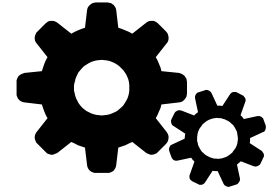


# Generative Adversarial Nets

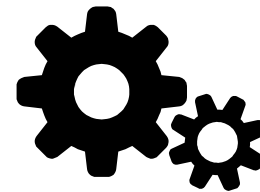


Training stages of a GAN  
Black dotted: True data  
Green solid: Generated data  
Blue dotted: Discriminator loss

# Conditional GAN

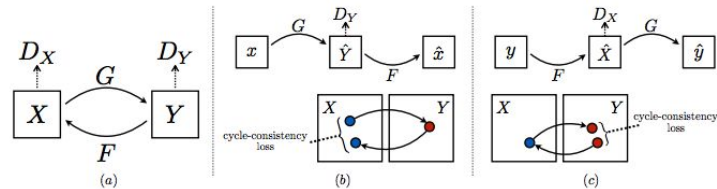
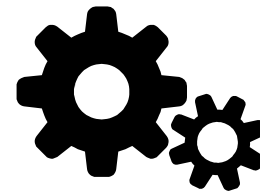


# Conditional GAN



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

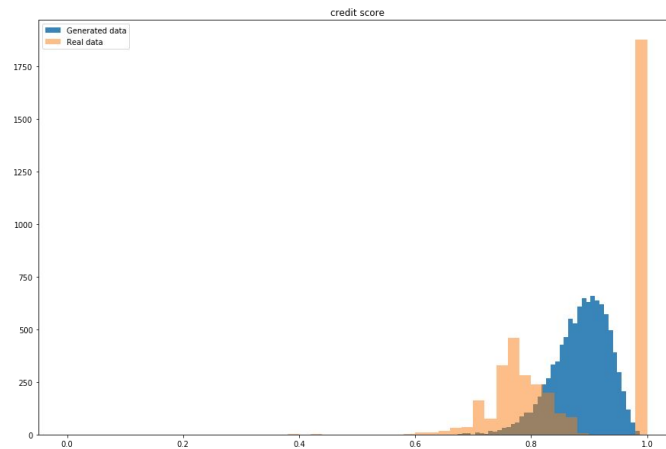
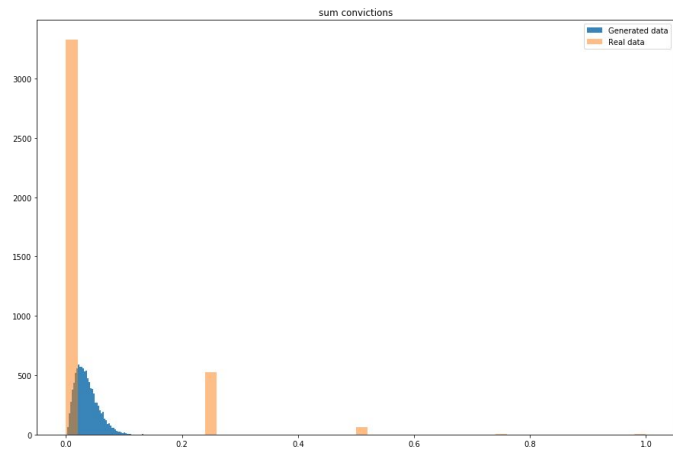
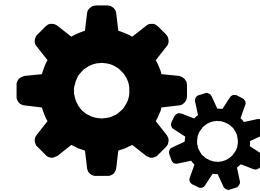


**Horse** → **Zebra**

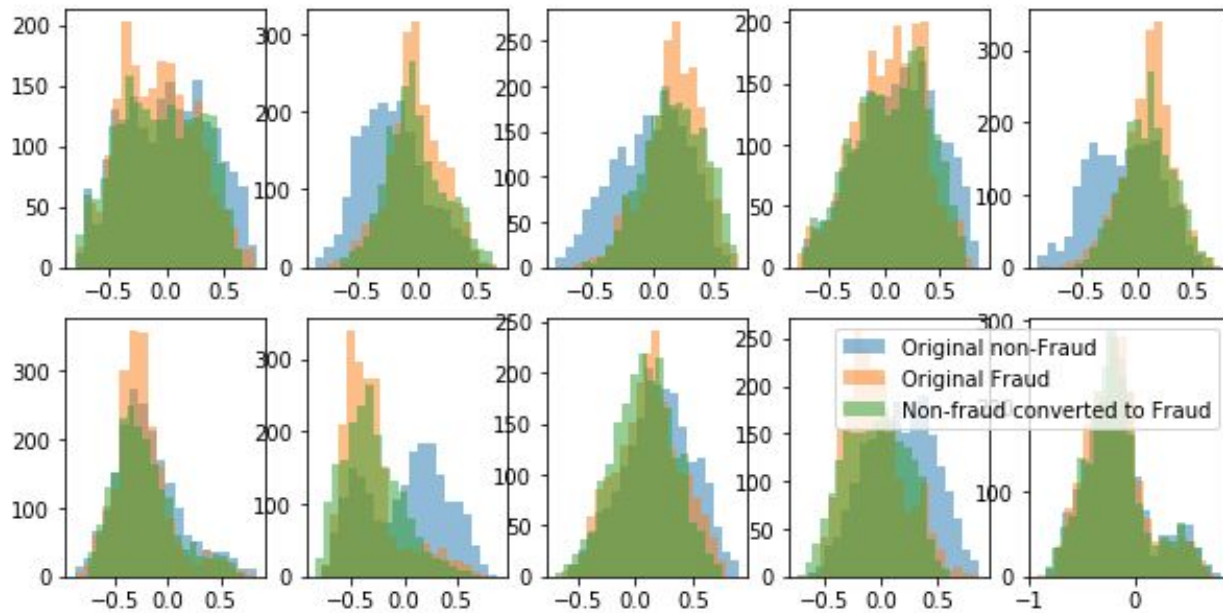
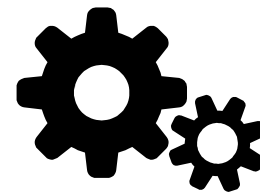
**or ...**

**No Claim** → **Claim**

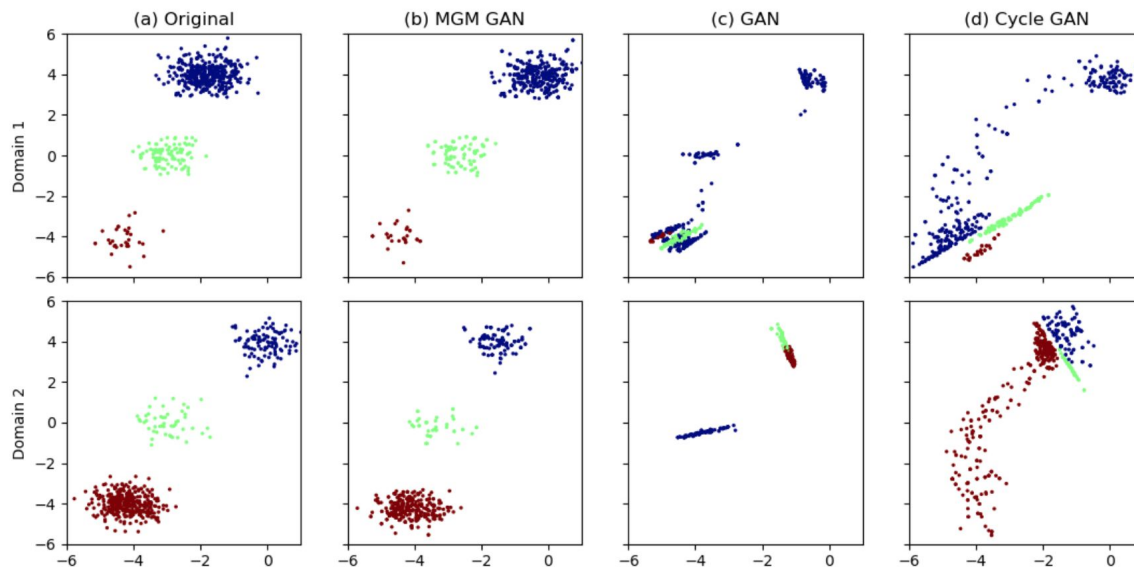
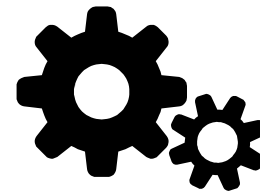
# GANs were very effective for images, but how about tabular data?



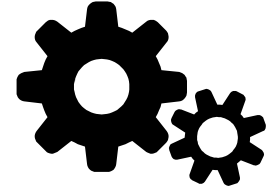
# CycleGAN for fraud detection + sparse Autoencoders



# Topology preserving CycleGAN



# Results

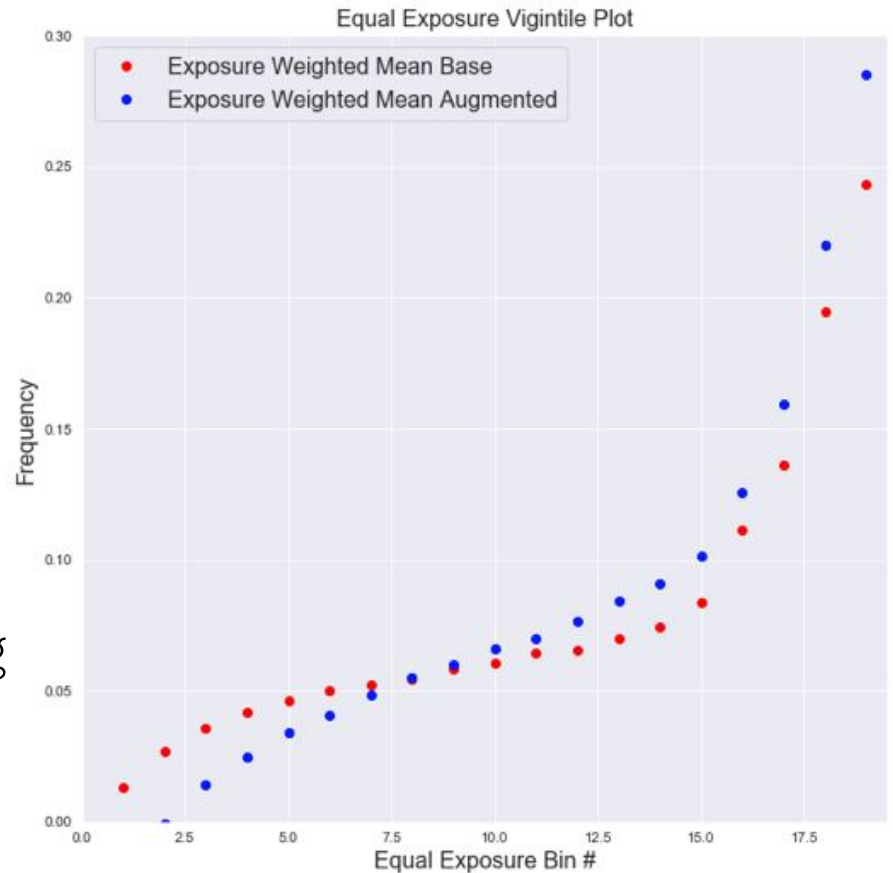


	Precision	Recall	F1-score
No augmentation	0.88	0.77	0.79
SMOTE	0.94	0.79	0.85
ADASYN	0.79	0.76	0.77
cGAN	0.90	0.85	0.87
cycleGAN	0.92	0.83	0.87

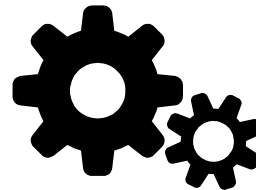
$N_s$	No. of misclassifications for MNIST. Test error rate for others.				
	Supervised	DCGAN	W-DCGAN	DCGAN-10	BayesGAN
MNIST	$N=50k, D = (28, 28)$	14	15	114	32
20	—	$1823 \pm 412$	$1687 \pm 387$	<b><math>1087 \pm 564</math></b>	$1432 \pm 487$
50	—	$453 \pm 110$	$490 \pm 170$	<b><math>189 \pm 103</math></b>	$332 \pm 172$
100	$2134 \pm 525$	$128 \pm 11$	$156 \pm 17$	$97 \pm 8.2$	<b><math>79 \pm 5.8</math></b>
200	$1389 \pm 438$	$95 \pm 3.2$	$91 \pm 5.2$	$78 \pm 2.8$	<b><math>74 \pm 1.4</math></b>
CIFAR-10	$N=50k, D = (32, 32, 3)$	18	19	146	68
1000	$63.4 \pm 2.6$	$58.2 \pm 2.8$	$57.1 \pm 2.4$	<b><math>31.1 \pm 2.5</math></b>	$32.7 \pm 5.2$
2000	$56.1 \pm 2.1$	$47.5 \pm 4.1$	$49.8 \pm 3.1$	$29.2 \pm 1.2$	<b><math>26.2 \pm 4.8</math></b>
4000	$51.4 \pm 2.9$	$40.1 \pm 3.3$	$38.1 \pm 2.9$	$27.4 \pm 3.2$	<b><math>23.4 \pm 3.7</math></b>
8000	$47.2 \pm 2.2$	$29.3 \pm 2.8$	$27.4 \pm 2.5$	$25.5 \pm 2.4$	<b><math>21.1 \pm 2.5</math></b>
SVHN	$N=75k, D = (32, 32, 3)$	29	31	217	81
500	$53.5 \pm 2.5$	$31.2 \pm 1.8$	$29.4 \pm 1.8$	$27.1 \pm 2.2$	<b><math>22.5 \pm 3.2</math></b>
1000	$37.3 \pm 3.1$	$25.5 \pm 3.3$	$25.1 \pm 2.6$	$18.3 \pm 1.7$	<b><math>12.9 \pm 2.5</math></b>
2000	$26.3 \pm 2.1$	$22.4 \pm 1.8$	$23.3 \pm 1.2$	$16.7 \pm 1.8$	<b><math>11.3 \pm 2.4</math></b>
4000	$20.8 \pm 1.8$	$20.4 \pm 1.2$	$19.4 \pm 0.9$	$14.0 \pm 1.4$	<b><math>8.7 \pm 1.8</math></b>
CelebA	$N=100k, D = (50, 50, 3)$	103	98	649	329
1000	$53.8 \pm 4.2$	$52.3 \pm 4.2$	$51.2 \pm 5.4$	$47.3 \pm 3.5$	<b><math>33.4 \pm 4.7</math></b>
2000	$36.7 \pm 3.2$	$37.8 \pm 3.4$	$39.6 \pm 3.5$	$31.2 \pm 1.8$	<b><math>31.8 \pm 4.3</math></b>
4000	$34.3 \pm 3.8$	$31.5 \pm 3.2$	$30.1 \pm 2.8$	<b><math>29.3 \pm 1.5</math></b>	$29.4 \pm 3.4$
8000	$31.1 \pm 4.2$	$29.5 \pm 2.8$	$27.6 \pm 4.2$	$26.4 \pm 1.1$	<b><math>25.3 \pm 2.4</math></b>

# Insurance

Double lift chart with data  
Augmentation + extreme modeling

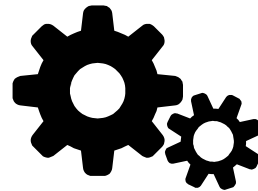


# Implementation Details



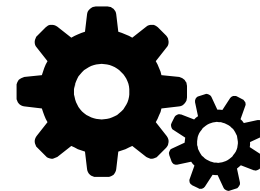
- Data has to be converted into a dense representation
- Domains have to be consistently related - retain same semantic features.
- Loads of tricks (weight clipping, regularizers, training schedule)
- Consistency over time (i.e. stationary)
- Avoid mode collapsing
- Diversity enforcing with conditional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