

# STOR566: Introduction to Deep Learning

## Lecture 12: Generative Models

Yao Li  
UNC Chapel Hill

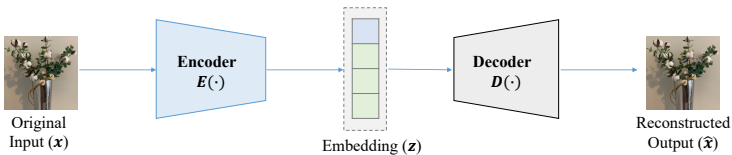
Oct 10, 2024

# Unsupervised Learning

- Working with datasets without a **response** variable
- Some Applications:
  - Clustering
  - Data Compression
  - Exploratory Data Analysis
  - Generating New Examples
  - ...
- Example: PCA, K-means, Autoencoders, GAN, etc

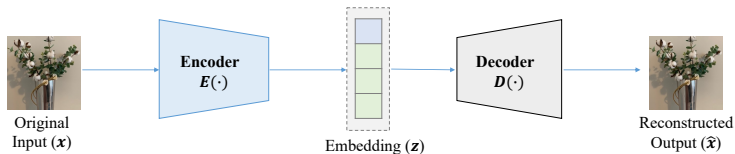
# Autoencoder: Basic Architecture

- Autoencoder: A special type of DNN where the target (response) of each input is the input itself.



# Autoencoder: Basic Architecture

- Autoencoder: A special type of DNN where the target (response) of each input is the input itself.



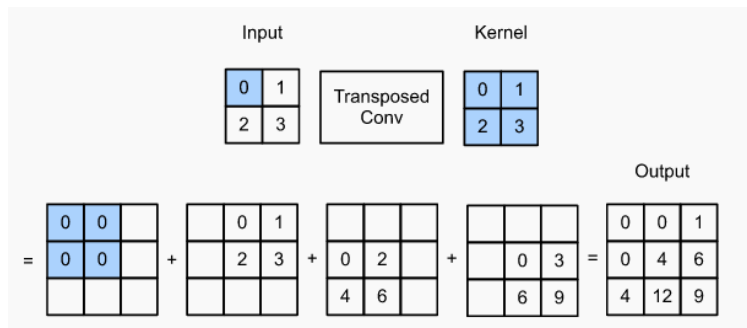
- Objective:

$$\|x - D(E(x))\|^2$$

Encoder:  $E : \mathbb{R}^n \rightarrow \mathbb{R}^d$

Decoder:  $D : \mathbb{R}^d \rightarrow \mathbb{R}^n$

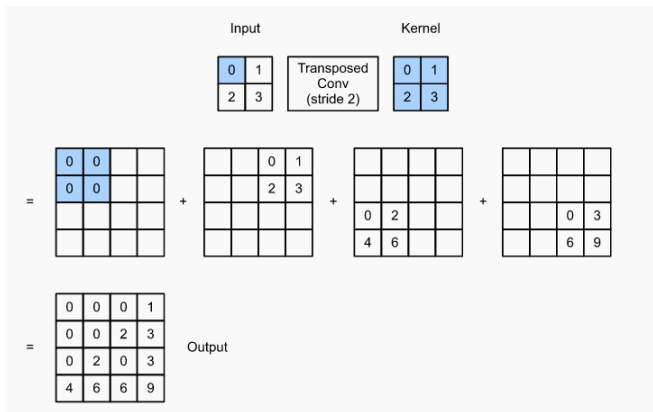
# Transposed Convolution



(Figure from Dive into Deep Learning)

- Multiple input and output channels: works the same as the regular convolution
- Number of weights:  $k_1 \times k_2 \times d_{in} \times d_{out} + d_{out}$

# Transposed Convolution



(Figure from Dive into Deep Learning)

- Strides are specified for the output feature map
- Padding: remove rows and columns from the output

# Overfitting

- Overfitting is a problem
- Solutions:
  - Bottleneck layer: a low-dimensional representation of the data ( $d < n$ )
  - Denoise autoencoder
  - Sparse autoencoder
  - ...

# Regularization

- Objective:

$$L(\mathbf{x}, \hat{\mathbf{x}}) + \text{regularizer},$$



# Regularization

- Objective:

$$L(\mathbf{x}, \hat{\mathbf{x}}) + \text{regularizer},$$

$L(\cdot, \cdot)$ : captures the distance between the input ( $\mathbf{x}$ ) and the output ( $\hat{\mathbf{x}}$ ).

- Example:  $\|\mathbf{x} - \hat{\mathbf{x}}\|^2$

# Regularization

- Objective:

$$L(\mathbf{x}, \hat{\mathbf{x}}) + \text{regularizer},$$

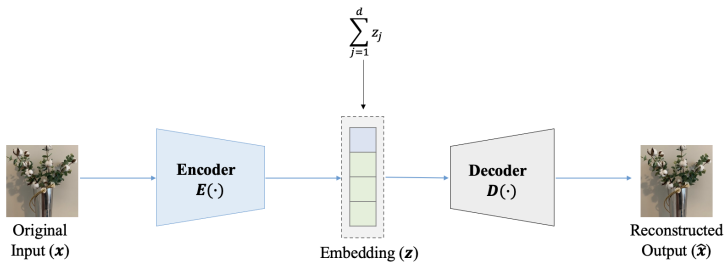
$L(\cdot, \cdot)$ : captures the distance between the input ( $\mathbf{x}$ ) and the output ( $\hat{\mathbf{x}}$ ).

- Example:  $\|\mathbf{x} - \hat{\mathbf{x}}\|^2$

Regularizer example:

- $L_1$  penalty:  $\sum_j |h_j^l|$
- $h_j^l$ : hidden output of  $j$ -th neuron in  $l$ -th layer

# Sparse Autoencoder

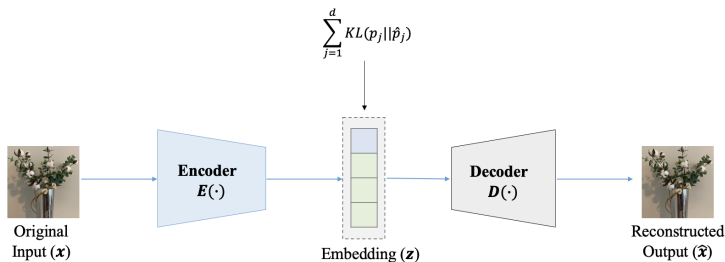


- Objective:

$$\|\mathbf{x} - \mathbf{D}(\mathbf{E}(\mathbf{x}))\|^2 + \lambda \sum_j |z_j|$$

- Iterate over layers.

# Sparse Autoencoder



- Another regularizer:

$$\|\mathbf{x} - \mathbf{D}(\mathbf{E}(\mathbf{x}))\|^2 + \lambda \sum_j KL(p_j || \hat{p}_j)$$

- Convert value of  $z$  to  $[0, 1]$ . (e.g., sigmoid activation)
- $p_j$ : probability of activation for neuron  $j$  in the bottleneck layer
- $\hat{p}_j = \frac{1}{B} \sum_{i=1}^B z_{ij}$

# Denoising Autoencoder

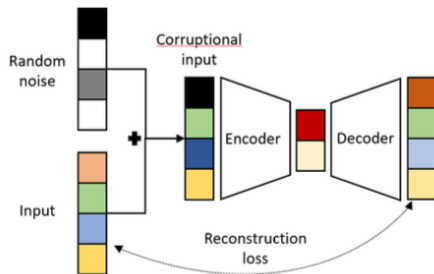


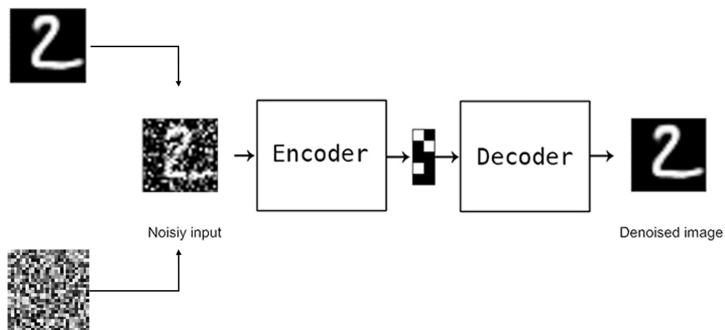
Figure from Bank, Dor, Noam Koenigstein, and Raja Giryes. "Autoencoders." (2020).

- Another regularizer:

$$\|x - D(E(x + \delta))\|^2$$

- $\delta$ : Random noise

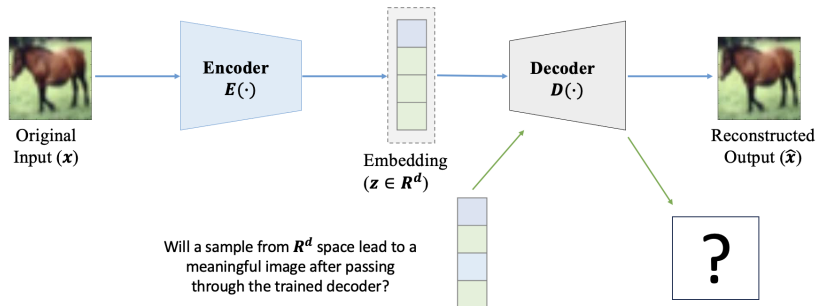
# Denosing Autoencoder



- noisy data → clean data
- Learn to capture valuable features and ignore noise

# Generative Model

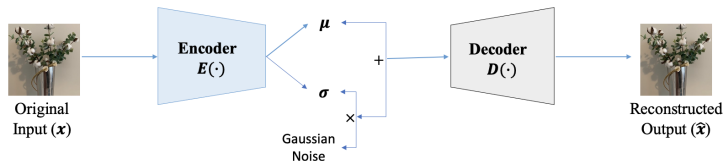
# Generative Problem



- In general, a trained Vanilla auto-encoder cannot be used to generate new data

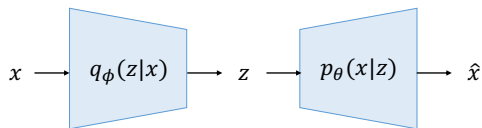
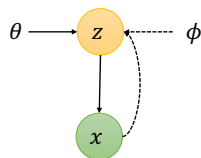


# Variational Autoencoder (VAE)



- Probabilistic model: will let us generate data from the model
- Encoder outputs  $\mu$  and  $\sigma$
- Draw  $\tilde{z} \sim N(\mu, \sigma)$
- Decoder decodes this **latent** variable  $\tilde{z}$  to get the output

# Variational Autoencoder (VAE)



- Maximum likelihood approach:  $\prod_i p(\mathbf{x}_i)$
- Variational lower bound as objective:
  - End-to-End reconstruction loss (e.g., square loss)
  - Regularizer:  $KL(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$
- Objective:

$$L(\mathbf{x}, \hat{\mathbf{x}}) + KL(q_\phi(\mathbf{z}|\mathbf{x})||p(\mathbf{z}))$$

# Variational Lower Bound

- Variational lower bound:

$$\log p(x) \geq E_{q(z|x)} (\log p(x|z)) - KL(q(z|x)||p(z))$$

- How to derive the variational lower bound from the likelihood?
- Suggested reading: Kingma et al. (2013). Auto-encoding variational bayes. ICLR.

# Re-parameterization Trick

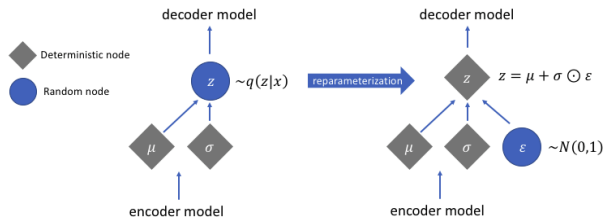
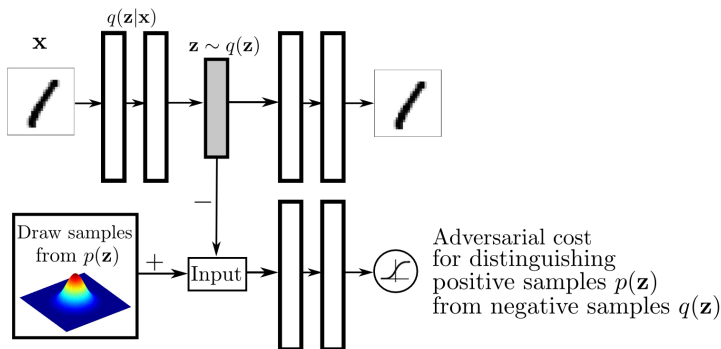


Figure from Jeremy Jordon Blog

- Cannot back-propagate error through random samples
- Reparameterization trick: replace  $\tilde{z} \sim N(\mu, \sigma)$  with  $\epsilon \sim N(0, I)$ ,  
 $z = \epsilon\sigma + \mu$

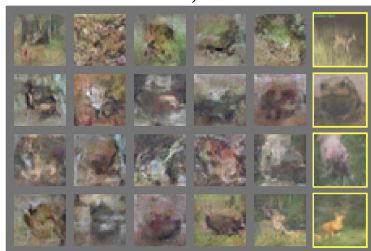
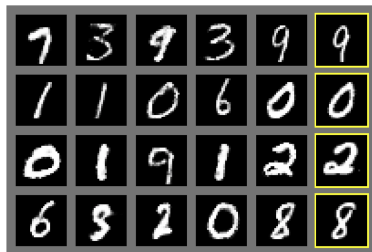
# Adversarial Autoencoder



- The top row is a standard autoencoder
- Force the embedding space distribution towards the prior

# Generated Adversarial Network

- Discriminative models:
  - Given an image  $x$ , predict a label  $y$
  - (by learning  $P(y | x)$ )
- Generative models:
  - Generate new images
  - Learn  $P(x)$  (or  $P(x, y), P(x | y)$ )



(Goodfellow et al., 2014)

# How to represent a distribution

- Define the distribution **implicitly**
- Start from a random vector  $z$ : a simple distribution (e.g., sphere Gaussian)
- Define (the sampling process of) the distribution as a function  $G$ :

$$z \rightarrow G(z) = x$$

- Our goal is to learn this **generator function**  $G$

# How to represent a distribution

- Define the distribution **implicitly**
- Start from a random vector  $z$ : a simple distribution (e.g., sphere Gaussian)
- Define (the sampling process of) the distribution as a function  $G$ :

$$z \rightarrow G(z) = x$$

- Our goal is to learn this **generator function**  $G$

Example:

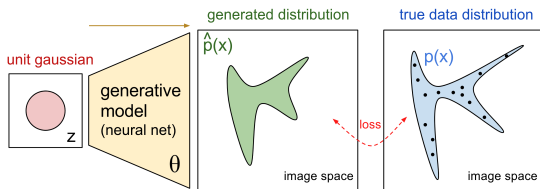
- Gaussian with covariance matrix  $N(0, \Sigma)$

$$z \sim N(0, I) \quad \rightarrow \quad \underbrace{\Sigma^{1/2} z}_{G(z)} \sim N(0, \Sigma)$$



# Neural network as a generator

- Now we assume  $G$  is a neural network parameterized by  $\theta$
- Goal: learn  $\theta$  to make generated distribution similar to the data distribution

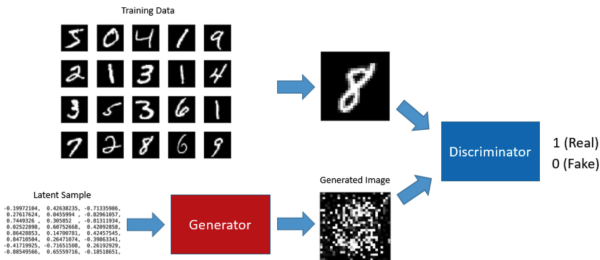


(figure from <https://openai.com/blog/generative-models/>)

- But how to evaluate the quality of generated distribution?

# Generative Adversarial Network (GAN)

- A good measurement: whether there exists a **discriminator** (classifier) to distinguish real/fake images
- Generative Adversarial Network (GAN): Train two networks jointly
  - The **generator network** tries to produce realistic-looking images
  - The **discriminator network** tries to classify real vs fake images



(figure from <https://naokishibuya.medium.com/understanding-generative-adversarial-networks>)

# Training objective

- The **discriminator**'s goal: classify real/fake images

$$L_D = E_{x \sim \text{real data}} [ - \log D(x) ] + E_z [ - \log(1 - D(G(z))) ]$$

- **Generator**'s goal: fool the discriminator
- A simple cost function for generator: the opposite of the **discriminator**'s
- The minmax training objective:

$$\max_G \min_D L_D(G, D)$$

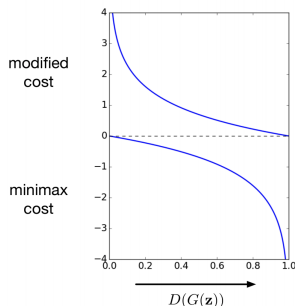
- GAN training: alternatively update  $G$  and  $D$

# Gradient vanishing problem

$$\max_G \min_D E_{x \sim \text{real data}} [-\log D(x)] + E_z [-\log(1 - D(G(z)))]$$

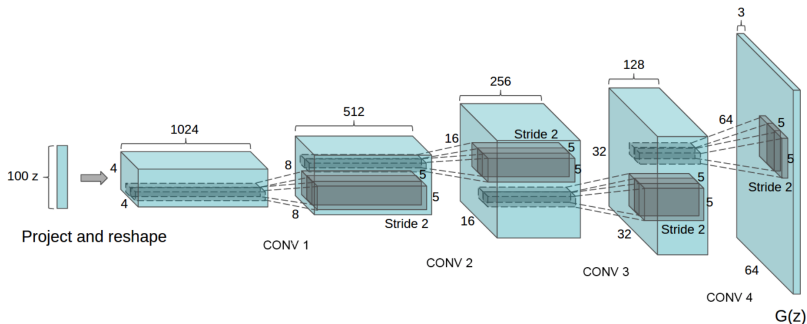
- The discriminator is usually much better than generators ( $D(G(z)) \rightarrow 0$ ), this implies the gradient of generator will vanish
- A modified generator loss:

$$L_G = E_z [\log(1 - D(G(z)))] \Rightarrow L_G = E_z [-\log D(G(z))]$$



# CNN for both generator and discriminator (DC-GAN)

- Discriminator: a regular classification network
- Generator: CNN with **transposed convolution** structure



(Radford et al., 2015)

# DC-GAN results



(Figure from Raford et al., 2015)

## Many improvements have been made

- c-GAN (Mirza and Osindero, 2014): add class label into the generator
- AC-GAN (Odena et al., 2016): discriminator classifies both real/fake and class label
- WGAN (Arjovsky et al., 2017): use Wasserstein distance
- SN-GAN (Miyato et al., 2018): spectral regularization
- Big-GAN (Brock et al., 2018): large batch (2048), bigger model
- Fast-GAN (Liu and Hsieh, 2018), (Zhong et al., 2020): small batch (64) can also work with adversarial training
- Style-GAN<sub>1,2,3</sub> (Karras et al., 2018; Karras et al., 2019; Karras et al., 2021): latent code transformation, progressive growing GAN

# Big-GAN results

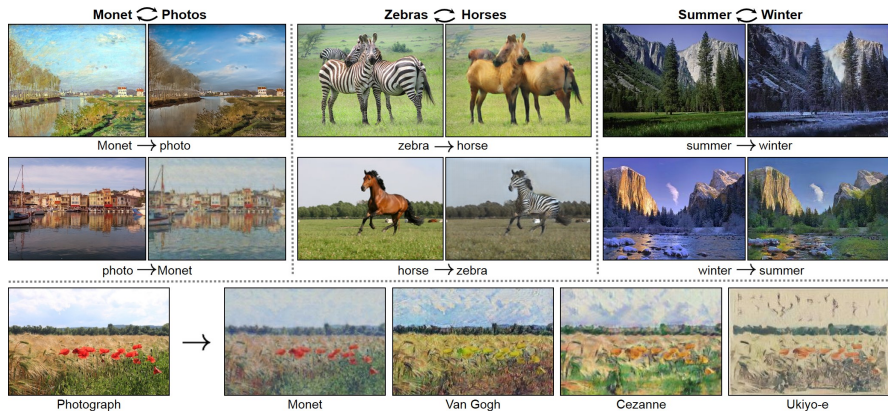


(Figure from Brock et al., 2018)

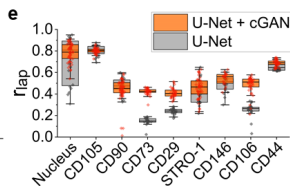
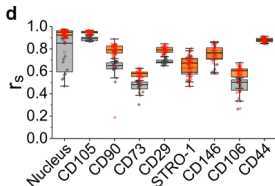
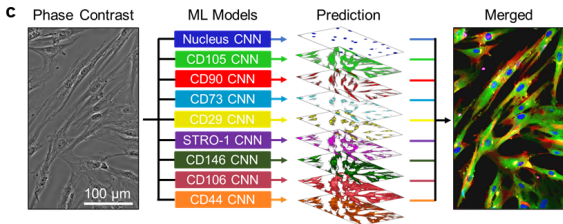
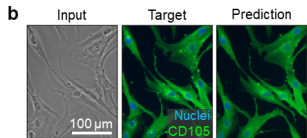
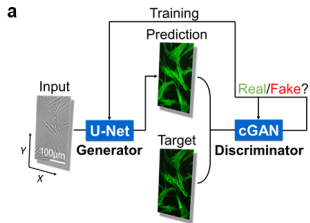


# Image-to-image translation

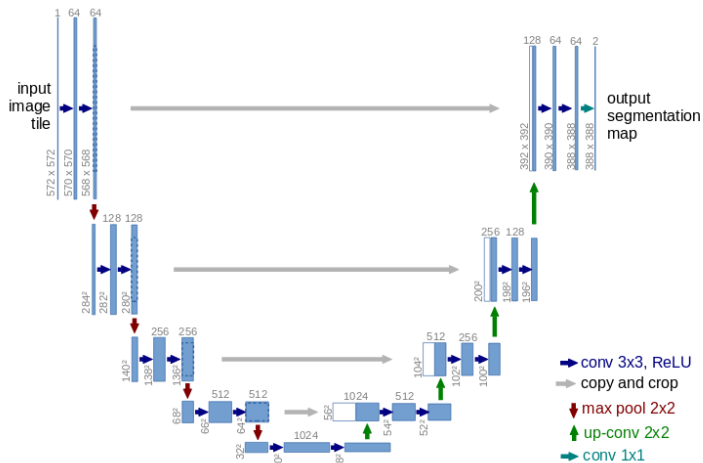
Cycle GAN: Zhu et al., 2017



# Many applications in bioinformatics



# Image-to-image translation



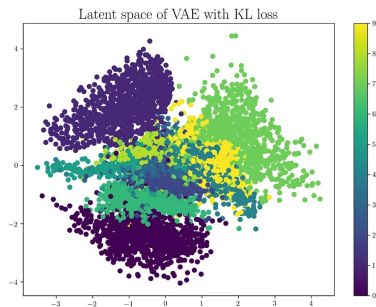
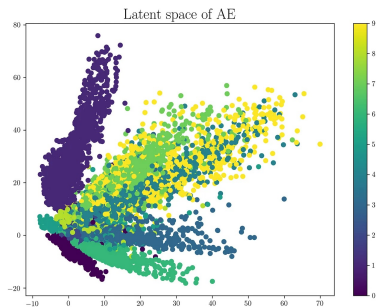
(The unet architecture)

# Embedding Space Visualization

Commonly used visualization tools:

- t-SNE (t-Distributed Stochastic Neighbor Embedding)
  - Van der Maaten et al. (2008). Visualizing data using t-SNE. *Journal of Machine Learning Research*, 9(11).
  - Available: sklearn
- UMAP (Uniform Manifold Approximation and Projection)
  - McInnes et al. (2018). UMAP: Uniform Manifold Approximation and Projection. *Journal of Open Source Software*, 3(29), 861,
  - Available: umap-learn
- PCA (Principal Component Analysis)
  - Available: sklearn

# Examples with tSNE



- Embedding space visualization for a Vanilla autoencoder and a VAE trained on MNIST
- VAE: more compact

# Examples with PCA

- Problem: Game Result Prediction



Figure: Heroes of the Storm and Dota 2 characters

# Assumption

## Assumption

We assume a team's score can be written as

$$s_t^+ = \sum_{i \in I_t^+} w_i + \sum_{i \in I_t^+} \sum_{j \in I_t^+} \mathbf{v}_i^T \mathbf{v}_j$$

- $w_i$ : individual ability of  $i$ -th player
- $\mathbf{v}_i \in R^d$ : teamwork ability of  $i$ -th player
- $I_t^+$ : winning team player index set
- $s_t^+$ : winning team score

# Team Ability Visualization (PCA)

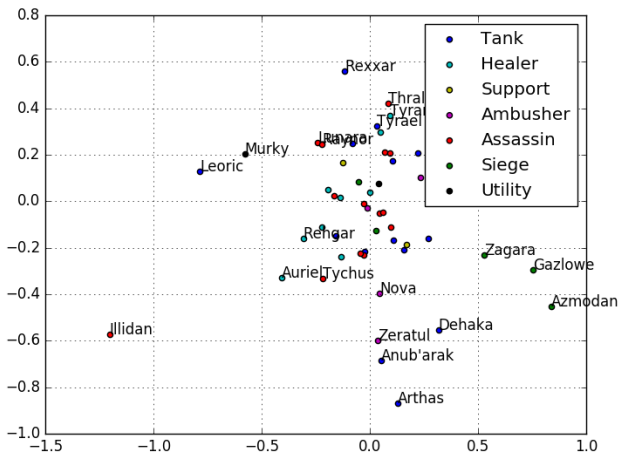


Figure: Projection of team ability vector for each character ( $\mathbf{v}_i$ ) to 2-D space. Colors represents the official categorization for these characters.



# Conclusions

- Autoencoder
- GAN
- Visualization tools

Questions?