

Generating Hands with Reinforcement Learning

Written by Chase Mutzig

February 3, 2024

Abstract

From its inception, AI has been able to do incredible things, from modern art, to web searching and creation, to faking an entire person into existence. However, something that has always been difficult for machines are hands, specifically of the human variety. Always too many fingers, too many segments, or no hand at all, it is incredibly rare that a hand is properly formed via a machine. In this paper, we will continue an experiment that used a variational auto encoder to generate hands by adding reinforcement learning. The goal is to generate a human hand that is anatomically accurate 70-80% of the time. To do this, we will use a dataset of around 11k hands in specific orientations and positions and feed them to the modified neural network. This paper will summarize the end results gained from using a VAE, explain how reinforcement learning is projected to help accuracy, and the expected implementation of reinforcement learning to the VAE.

Introduction

AI mimicry in the form of vocals, photos, and videos is growing in popularity. With this growth in popularity, it is also becoming increasingly more difficult to differentiate between what is fake and real. Take a person's face as an example. In the early generations, an AI would have all the parts of the face; eyes, nose, ears, eyebrows, mouth, etc. However, these parts may have had mismatched sizes, locations, or other impossibilities (Kalpokas and Kalpokiene, 2022). In current times, these mistakes have been ironed out and it can often be very difficult to spot the malformed details. If you take a human or any humanoid figure, one of the most reliably messed up details would be the hands (Meg, 2023).

Hands are typically attached via a wrist, have a palm, and five digits. Each of these digits are made up of sections that have fingernails, fingerprints, and joints. In AI generated material, there are often incorrect numbers of such details. Currently, the main issue relates to quantity of features, rather than quality of features, as many current generators are able to make good quality and high fidelity images. In an effort to correctly generate hands 70-80% of the time, we are going to be modifying a neural network specialized in focusing on these easy-to-miss details by adding reinforcement learning. The reinforcement learning agent will work to improve the baseline understanding of the neural network in an effort to improve the results enough as to reach the 70-80% accuracy goal.

Background

To summarize the preceding experiments, a variational auto encoder is a neural network structure that forms a latent understanding of the input, then tries to recreate the input from the created understanding. In this case, the input is roughly 11,000 hands all in similar positions, and the output is a *strictly new* hand that was not previously in the input data. This newly generated hand should be visibly recognizable as a hand, and should be anatomically correct. This means that the hand should have a palm, four fingers, and a thumb. Each digit also must have the correct amount of segments, if visible.

A variational auto encoder was chosen specifically due to the nature of this task. VAEs generate probabilistic latent spaces, allowing them to generate many different samples and are thus well suited for tasks where we want to interpolate between different data; in this case, the skin tones, positions, and structure of the input hands (Pu, et. al. 2016).

A standard auto encoder has a deterministic latent space, meaning they do not capture variations in the input data well (Bank, et. al., 2023).

As for GANs, or generative adversarial networks, they implicitly learn a latent space through a generator network which may not be as understandable or controllable as a VAE latent space (El-Kaddoury, et. al., 2019).

Previous Results

The results from the VAE section will be considered the preceding and baseline results for this experiment.



Figure. 1: Best baseline results from pure VAE training.

The minimum goal will be to generate results better than

these preceding results, and the overall goal will be to generate anatomically correct hands 70-80% of the time. To clarify, if a batch of 20 hands are created, then roughly 15 of those hands must be anatomically correct.

Related Work

Specifically with references to hands, there has been work mainly on reading hand motions and gestures. In one such case, Stanford University has made “[a] new AI learning scheme combined with a spray-on smart skin [that] can decipher the movements of human hands to recognize typing, sign language, and even the shape of simple familiar objects” (Patel, 2023).

In another similar project, people are using AI to learn and read gestures in an effort to make a hands-free, no external devices blackboard, similar in style to Iron-Man’s holo-desk. By reading hand gestures, a person can write numbers and letters or create more blackboards, lead meetings, and many other applications (Soroni et al., 2021).

Lastly, although there are many more examples of AI using hand reading, there is a physics simulation that tries to map hand movements to a virtual environment, similar to the aforementioned project. However, this project also attempts *accurate* haptic feedback by using reinforcement and imitation learning. This means they don’t want any clipping of fingers into objects, and for the person doing the manipulation to feel the opposing forces physically, from virtual objects (Garcia-Hernando, et. al., 2021).

Methodology

For this experiment, we want to change the latent space of the VAE by adjusting latent variables to improve the understanding of what a hand is, and thus improve the generated hands.

Due to its suitability, the Soft Actor-Critic (SAC) algorithm will be used. SAC is an algorithm that is well suited for continuous action spaces, which the RL agent will be dealing with since the base of this model is a VAE that generates a continuous latent space (Haarnoja, et. al., 2019). Continuing with the SAC algorithm, the agent must also have clearly defined states, actions, and rewards.

The states represent information about the environment that the agent will use to make decisions regarding any available actions it may take. Each state will be the collection of latent variables that encode information on the latent space, or understanding of the model. There will be one state per episode.

The actions will be adjustments and transformations applied to the latent variables in the VAE latent space. The agent will use SAC to learn a policy that maps the states to the actions, allowing it to modify the latent variables and hopefully improve the latent understanding of what a hand may be.

The reward for the RL agent will be based on the change in loss after each training epoch. If the loss decreases between epochs, then the agent is rewarded a positive reward in proportion to the change. If the loss instead increases, the

proportional reward will be negative. The algorithm used is as follows:

$R_t = -\lambda * \Delta L_t$ if $t > 1$, where R_t is the reward at epoch t , ΔL_t is the change in loss from the previous epoch, defined as $L_t = L_t - L_{t-1}$, and λ is a positive constant scaling factor. The use of λ allows for the control of reward scaling, making rewards scale more or less depending on the loss change. λ is not strictly necessary for this algorithm to work, yet is included due to the flexibility in experimentation it provides.

SAC Algorithm

The Soft Actor-Critic algorithm has four main parts, each of which are responsible for a major portion of what makes the SAC algorithm good (Derman, et. al., 2018).

Critic Update The critic update evaluates the quality of the actions taken by the agent. This update ensures an accurate estimation of the cumulative reward, incorporating the reward, discounted future rewards, and the entropy-adjusted policy. The critic update is given as the following:

$$\mathcal{L}_Q(\phi) = E_{(s,a,r,s') \sim \mathcal{D}} \left[\frac{1}{2} (Q_\phi(s,a) - (r + \gamma \min_{a'} Q_{\phi'}(s',a') - \alpha \log \pi_\psi(a'|s'))^2 \right] \quad (1)$$

- s : Current state
- s' : Next state
- a : Taken action
- a' : Next taken action
- r : Reward
- \mathcal{D} : Data distribution representing experiences collected from the environment
- γ : Discount factor for future rewards
- α : Temperature parameter
- $\pi_\psi(a' | s')$: Policy representing the probability of taking action a' in state s'

Value Function Update The value function update is responsible for estimating the cumulative reward and aligning the policy’s entropy-adjusted distribution with a target entropy. This encourages a balance between trying new things (exploration) or doing something already done (exploitation), leading to diverse yet high-quality policies. The VFU is given as:

$$\mathcal{L}_V(\theta) = E_{s \sim \mathcal{D}} \left[\text{KL} \left[\pi_\psi(\cdot | s) \parallel \exp \left(\frac{Q_\phi(s, \cdot)}{\alpha} \right) \right] \right] \quad (2)$$

- $\mathcal{L}_V(\theta)$: Loss function for the value function update.
- $E_{s \sim \mathcal{D}}$: Expectation over the data distribution \mathcal{D}
- **KL**: Kullback-Leibler divergence
- $\pi_\psi(\cdot | s)$: Policy representing the probability distribution of actions given state s
- $Q_\phi(s, \cdot)$: Q-function representing the expected cumulative reward given state s and various actions
- α : Entropy temperature parameter

Policy Update The policy update updates the policy to maximize the expected reward along with the entropy term. Doing this encourages the agent to explore its environment and to take actions that yield high rewards. The policy update is given as:

$$\mathcal{L}_\pi(\psi) = E_{s \sim \mathcal{D}} \left[E_{a \sim \pi_\psi(\cdot|s)} \left[\alpha \log(\pi_\psi(a|s)) - Q_\phi(s, a) \right] \right] \quad (3)$$

- $\mathcal{L}_\pi(\psi)$: Loss function for the policy update
- $E_{s \sim \mathcal{D}}$: Expectation over the data distribution \mathcal{D}
- $E_{a \sim \pi_\psi(\cdot|s)}$: Expectation over the policy's action distribution given state s .
- α : Entropy temperature parameter
- $\log(\pi_\psi(a|s))$: Log probability of the action a given state s under the policy
- $Q_\phi(s, a)$: Q-function representing the expected cumulative reward given state s and action a

Temperature Entropy Update The entropy temperature update regulates exploration vs. exploitation. This update makes sure the agent maintains the wanted amount of exploration or exploitation by changing the temperature parameter based on the policy's entropy. JuMathematically, the update is given as:

$$\alpha \leftarrow \text{clip}(\alpha + \eta * E_{s \sim \mathcal{D}} [-\log(\pi_\psi(a|s))] - \text{Target Entropy}, \alpha_{\min}, \alpha_{\max}) \quad (4)$$

- α : Entropy temperature parameter (to be updated)
- $\text{clip}(\cdot)$: Clip function ensuring the value stays within a specified range
- η : Learning rate for the update
- $E_{s \sim \mathcal{D}}$: Expectation over the data distribution \mathcal{D}
- $-\log(\pi_\psi(a|s))$: Negative log probability of the action a given state s under the policy
- Target Entropy: Variable to be specified in practical application
- α_{\min} and α_{\max} : Minimum and maximum values for the entropy temperature

Implementation

In practical application, the SAC algorithm is not a sequence of mathematical formulas (Haarnoja, et. al., 2019) and instead uses miniature neural networks to function as actor and critic. The preceding math is used to define the theoretical application of this algorithm while in reality neural networks are used to approximate the purpose of SAC algorithmically. In the case of this paper, both the actor and critic are dense neural networks with ReLu activation applied to each dense layer allowing for non-linear learning. The input data has all RGB values normalized, meaning that the RGB scale is from zero to one, instead of 0-255. Due

to this normalization, an action modifying any latent space with a Delta greater than 1 would lead to extreme changes, and thus the actor has an added dense output layer using tanh to regulate any actions to a range of [-1, 1]. The critic has an added input layer that takes the current state and actor's action and concatenates them, leading to a modified state. The critic then passes this through two dense layers, compressing it into a final single node dense layer. This process estimates the Q value of the action taken by providing a linear combination of both the input state and the action applied to it, simplified down to a single value. This single value can then be compared to other values to determine if the current state-action pair will progress or regress policy optimization. The single node output layer of the critic neural network does not have an activation function applied due to the goal being to predict a continuous Q-value rather than trying to generate a classification, in which case the non-linearity provided from activation functions would be useful. The actor's action is the prediction generated from the actor neural network given the current state.

For this project, there will be a single actor determining actions based off of training not yet implemented. There will also be two critic networks due to the use of multiple discriminators allowing for better results and decreasing the likelihood of an exploding gradient occurring. (Fujimoto, et. al., 2018). With some training, the hope is that the actor will only predict changes that are useful in policy optimization and the critics will prevent future poor decisions.

As for training the actor and critic, there will be a sample_latent function which generates a sample latent space. This sample space then gets used in making a duplicate VAE model that is nearly identical to the originally described base model. The only difference being the latent space and subsequent changes to handle a differently sized space. The environment in this case is a simple usage of the temporary model, a reset back to zero weights, tracking of actions taken and modified states, sample image generation, and giving rewards based on the loss values. Loss as stated previously is calculated using KL_reconstruction. The actor and critic networks as a whole will function as a reinforcement learning agent, and thus will be referred to as 'the agent'.

The Q-learning method for improving agents and their corresponding policies will be used, however, there is no current implementation of a training loop. There is only the described agent, temporary model, and environment. This does not mean that the agent as a whole is meaningless until it is trained, merely that it has much room for improvement. For example, the agent is currently not basing its actions off of anything and thus any actions taken are completely random and act as adding noise to the representation. Adding noise is a common practice when preprocessing data, however, it is now being used in the model training loop, leading to a VAE that is better able to understand what is noise and what is part of the original images. There was no enhancement image preprocessing done, so even this relatively simple application of the SAC algorithm has made a large difference in results.

Issues and Solutions

When training the agent or updating a latent space, it is computationally expensive to remake a model from scratch each time if the only change is the latent space. However, this issue arises from the size of the latent space and how keras compiles models. For the first issue, when creating a model the latent space, with a dimensionality of 300, will have a shape of (0, 300). The x represents the number of active samples being understood and the y represents dimensionality. After a single epoch of training, when taking the latent space out of a model, the shape will be (num_images, 300) and thus incompatible with the original model, which was made to be compiled with a shape of (0,300). This means a new model needs to be created that is compatible with the new latent space. Technically speaking, after the first recreation of a model, all future changes to the latent space will not make the space incompatible with the model. However, due to how keras model compiling works, of which this project uses, a models latent space cannot be 'updated'. While a model in its totality may be updated with new weights or pre-training, since the latent space is a *part* of a model and not considered weights or training, it cannot be specifically updated. To get around this issue, the 'updating' of a model is in actuality the recreation of the model using the updated latent space, leading to many models and a large use of resources.

When making changes to a latent space, taking a latent space out of a model, or putting it back in, the data type will change. When leaving or creating a model, the type will be a symbolic tensor which stores data and is not conducive to data editing. When making changes to a space or running it through the agent networks, it is necessary to have a data type that is modifiable and keeps any relevant data. To solve this, the latent space is taken via a prediction from the model *encoder*, which returns an easy to visualize and edit numpy array. After making changes to this array it is passed through a Lambda function which converts it back into a symbolic tensor and thus able to be used in a new model.

Preliminary Results

All preliminary results shown are done using a single epoch to increase the frequency at which tests can be done. Due to the following results having been generated from single epochs, the results are not expected to be any indicator of quality and instead should be seen as small-scale tests to ensure that there is improvement from the addition of an agent.



Figure. 2: VAE model with no agent, 8 training images and 2 testing images



Figure. 3: VAE model with agent, 8 training images and 2 testing images

After many quick tests using very little of the dataset, it quickly became clear that while the agent was improving the results, there was a limit to how much improvement could be gained from 10 images total. In the results gained without the agent, there are no definable features, while with the agent a vague shadow of a hand is apparent in 3/8 of the shown examples.

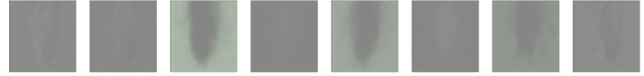


Figure. 4: VAE model with no agent, 80 training images and 20 testing images

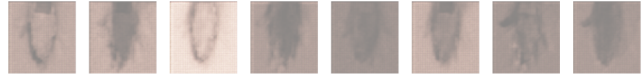


Figure. 5: VAE model with agent, 80 training images and 20 testing images

By expanding the available images by a factor of 10, the improvement from the simple and untrained agent becomes readily evident. In the example without the agent, the images still have extreme faint shadows of a hand or a green hue, while the images generated with the agent are showing progress towards a skin tone being used and zero images being mostly grey squares.

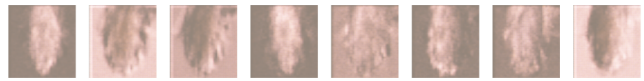


Figure. 6: VAE model with no agent, 800 training images and 200 testing images



Figure. 7: VAE model with agent, 800 training images and 200 testing images

Unfortunately, because the agent is not being trained and is simply adding in random noise, there is a turning point at which it becomes a hindrance instead of an assistance. As showcased in figures 6. and 7., there is more definition being attempted and shown in the results generated without the agent, while results with the agent almost seem to revert back to the solid squares of color. The only notable difference between the results in favor of the agent is that the main color across all results is closer to a normal skin tone than the mostly pink hue found in the agent-less results.

Evaluation

Aside from being anatomically correct, the hand must be the correct color (within any reasonable *human* skin tone range), have the correct orientation, be a clear image with little blur, etc. Essentially, the evaluation will be on a visual basis of whether the output is closer to what a hand is than the previous output. Out of the roughly 20 images per output,

about 15 of them need to be highly accurate for this project to be deemed a success and reach the 70-80% benchmark.

After preliminary tests in which the agent was added and experimented with, this first test of using the SAC algorithm is deemed a partial success. It is not a partial failure as the results are mostly poor, with a small amount of success. While the agent did improve results with a clear visual indication, the main evaluation criteria is to have the current output be closer to what a hand is than the previous output, and since most of the current output doesn't remotely resemble a hand, this first experiment is mostly a failure. The number of epochs will be increased dramatically in future phases which will lead to a major increase in quality.

Timeline

March 20: Midterm Paper and Presentation

April 3: Implement agent training loop

April 3 - 17: Train the agent using progressively larger sections of the dataset for more episodes. Improve training loop as necessary.

April 18 - 30: Have trained the model using the agent on the full dataset, for a minimum of 100 epochs. Increase epoch number.

May 1: Consolidate findings and complete final Paper, Presentation, and Demo

Conclusion

In this paper we talked about the continuation of a neural network experimentation that tried to generate images of hands and how we will work to improve it using reinforcement learning. This previous experimentation used a variational auto encoder and got poor results in color. The agent spoken of replicated the theoretical application of an algorithm known as soft actor-critic or SAC in a practical environment. The agent was able to improve the latent understanding and output generation of the base VAE model without training, and will require further refinement to achieve the goal of 70-80% accurate and visually acceptable results.

References

Matthias, Meg. "Why does AI art screw up hands and fingers?". Encyclopedia Britannica, 25 Aug. 2023

Kalpokas, Ignas, and Julija Kalpokiene. "From Gans to Deepfakes: Getting the Characteristics Right." Springer-Link, Springer International Publishing, 2022

Prachi Patel. "Spray-on Smart Skin Reads Typing and Hand Gestures." IEEE Spectrum, IEEE Spectrum, 3 Mar. 2023

F. Soroni, S. a. Sajid, M. N. H. Bhuiyan, J. Iqbal and M. M. Khan, Hand Gesture Based Virtual Blackboard Using Webcam, 06 Dec. 2021 IEEE 12th Annual Information Technology, Electronics and Mobile Communication Conference (IEMCON), Vancouver, BC, Canada, 2021

Guillermo Garcia-Hernando, Edward Johns, & Tae-Kyun Kim, Physics-Based Dexterous Manipulations with Estimated Hand Poses and Residual Reinforcement Learning. IEEE Spectrum, IEEE Spectrum, 2021

Pu, Y., Gan, Z., Hénao, R., Yuan, X., Li, C., Stevens, A., & Carin, L. (n.d.). Variational Autoencoder for Deep Learning of Images, Labels, and Captions, NeurIPS Proceedings, 2016

Bank, Dor, Noam Koenigstein, and Raja Giryes. Autoencoders. Machine learning for data science handbook: data mining and knowledge discovery handbook Feb. 2023

El-Kaddoury, M., Mahmoudi, A., Himmi, M.M. "Deep Generative Models... and Generative Adversarial Networks", Springer, 2019.

Tuomas Haarnoja, Aurick Zhou, Kristian Hartikainen, George Tucker, Sehoon Ha, Jie Tan, Vikash Kumar, Henry Zhu, Abhishek Gupta, Pieter Abbeel, and Sergey Levine. "Soft Actor-Critic Algorithms and Applications". Cornell University arXiv, 2019

Esther Derman, Daniel J. Mankowitz, Timothy A. Mann, and Shie Mannor. "Soft-Robust Actor-Critic Policy Gradient". Cornell University arXiv, 2018

Scott Fujimoto, Herke van Hoof, David Meger, "Addressing Function Approximation Error in Actor-Critic Methods", Proceedings of Machine Learning Research, 22 Oct. 2018