
Generative Adversarial Networks (GAN)

A Gentle Introduction

Su Wang

Department of Statistics and Data Science
University of Texas at Austin

Abstract

In this tutorial, I present an intuitive introduction to the *Generative Adversarial Network* (GAN) [1], invented by Ian Goodfellow of Google Brain, overview the general idea of the model, and describe the algorithm for training it as per the original work. I further briefly introduce the application of GAN in Natural Language Processing to show its flexibility and strong potential as a neural network architecture.

1 Overview

Generative Adversarial Networks (GAN) [1] is a deep learning framework in which two models, a generative model G and a discriminative model D , are trained simultaneously. The objective of G is to capture the distribution of some target data (e.g. distributions of pixel intensity in images). D aids the training of G by examining the data generated by G in reference to “real” data, and thereby helping G learn the distribution that underpins the real data. GAN is fleshed out in Goodfellow et al. (2014) [1] as a pair of simple neural networks. However in practice, the models can in principle be any generative-discriminative pairs¹

In the original work [1,2] and elsewhere [3], GAN has been given the analogy as a process of making counterfeit money: G plays the role of a counterfeiter in-training, while D the bank strives to identify fake bills and in the process (hopefully unintentionally) helps G honing its bill-making skills. More concretely, let $\mathbf{x} \sim p_{data}$ be characterizing features for real bills, and $G(\mathbf{z})$ be features G creates from some noise distribution $\mathbf{z} \sim p_z$. Further let J be some quantitative metric which measures the extent to which a bill is real. Then, D 's job is to lower $J(G(\mathbf{z}))$ (the score of fake bill) while increase the score $J(\mathbf{x})$ (the score of real bill) for more successful identification. G , on the other hand, aims at increasing $J(G(\mathbf{z}))$ (i.e. improving the quality of fake bill) by learning from “observing” how D make differentiations. As the game of “busting fake bill” and “making better fake bill” proceeds, the model distribution p_G draws closer to p_{data} , and eventually reaches an *equilibrium* [2] where D can no longer classify better than chance (i.e. $D(\mathbf{x}) = D(G(\mathbf{z})) = \frac{1}{2}$). Now we say G has arrived at an optimal point² in counterfeiting.

GAN has gained massive attraction in computer vision [4,5,6], feature representation [7], and more recently in *Natural Language Processing* (NLP) tasks: Document Modeling [8], Dialogue Generation [9], Sentiment Analysis [10], and Domain Adaptation [11]. In Section 3 I briefly exemplify the successful application of GAN in NLP [8,9].

¹These can be any generative and discriminative models, in a much wider sense than the term “G-D pair” is used in the literature [12].

²I will show that the optimal point is reached iff $p_G = p_{data}$ (Section 2)

2 Model

Formal description. To begin, I describe a barebone variant of GAN in its original formulation, in which both D and G are simple *Multi-Layer Perceptrons* (MLP). Let p_{data} be our target distribution which the generator G will learn by approximating it with the model distribution p_G . G is associated with a noise prior p_z , from which G draws sample z , and create fake datum³ $G(z; \theta_G)$, where θ_G are model parameters. The discriminator $D(x; \theta_D)$ takes x or $G(z)$ as input, and returns a binary judgment as to whether the input is from p_{data} or p_G .

D evaluates the quality of input x with the following metric:

$$J(D, G) = \mathbb{E}_{x \sim p_{data}} [\log D(x)] + \mathbb{E}_{z \sim p_z} [\log(1 - D(G(z)))] \quad (1)$$

The first term of (1) evaluates the expectation of the log probability that x is drawn from real data; the second term, that x is drawn from fake data (i.e. noise). D 's objective is to maximize J , i.e. pushing the both terms towards 0 (i.e. $D(x) \rightarrow 1$, $D(G(z)) \rightarrow 0$). G , on the other hand, wants to minimize J by lower the second term⁴ (i.e. $D(G(z)) \rightarrow 1$). The ‘‘conflicting objectives’’ of the two models results in a *Minmax Game* [1] which is described as follows:

$$\min_G \max_D J(D, G) \quad (2)$$

Algorithmically, the training takes the form of an alternation process between minimization and maximization, which is described in Algorithm 1. In practice, Eq. 1 often does not bring the model to equilibrium, this is because $\log(1 - D(G(z)))$ rapidly saturates in the early stage of training, where D easily rejects $G(z)$ because G generates fake data of poor quality, such that they conspicuously differ from the real data. Therefore, rather than evaluating how bad fake data are (G 's objective in the second term of Eq. 1), we instead evaluate how good they are by setting G 's goal to maximizing $J_G(G) = \log D(G(z))$. In so doing we end up with two objective functions:

$$\begin{aligned} \max_D J_D(D, G) \\ \max_G J_G(G) \end{aligned} \quad (3)$$

To illustrate the training process graphically, we observe (a) how p_G changes over time, and consequently (b) how the discrimination boundary of D changes accordingly (Figure 1, Figure 1 in [1]).

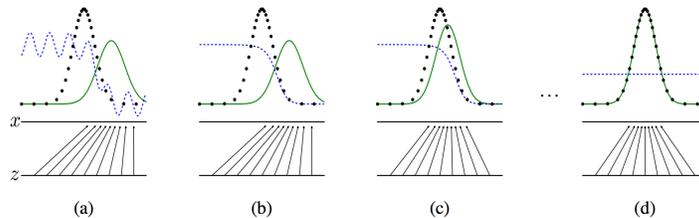


Figure 1: Training Process of GAN (green solid: p_G ; black dotted: p_{data} ; blue dash: D 's discrimination boundary; arrows: generation of fake data)

In (a) through (d), z are sampled uniformly from the noise prior p_z , and p_G draws closer to p_{data} . In the process, the discrimination boundary changes accordingly, and finally morphs into a flat line (i.e. $D(x) = D(G(z)) = \frac{1}{2}$) which indicates D is now unable to tell fake and real data apart. Specifically, Figure 1 shows a scenario where the model is near convergence: (a) p_G is similar to p_{data} , and D is now partially accurate⁵; (b) D is updated: based on the relative distribution of p_G and p_{data} , D converges in the inner loop of Algorithm 1 to $D^*(x) = \frac{p_{data}(x)}{p_{data}(x) + p_G(x)}$; (c) G is updated: p_G

³The generator is essentially a function that maps a noise datum into the space of a real datum, i.e. $G : z \mapsto x$. Note that z and x therefore do not have to be equal in dimensionality.

⁴The first term is independent of G .

⁵In the late stage of training, G starts to generate high quality fakes, to the effect that D 's classification performance suffers (but not entirely down to the level of randomness, i.e. $D(x) = D(G(z)) = \frac{1}{2}$).

Algorithm 1 Minmax Game

- 1: **for** specified # of training iterations **do**
- ▷ TRAINING DISCRIMINATOR
- 2: **for** specified k steps* **do**
- 3: Draw minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\} \sim p_z$.
- 4: Draw minibatch of m data samples $\{\mathbf{x}^{(1)}, \dots, \mathbf{x}^{(m)}\} \sim p_{data}$.
- 5: Update D 's parameters by gradient ascent:

$$\nabla_{\theta_d} \frac{1}{m} \sum_{i=1}^m \left[\log D(\mathbf{x}^{(i)}) + \log \left(1 - D(G(\mathbf{z}^{(i)})) \right) \right]$$

- 6: **end for**
- ▷ TRAINING GENERATOR
- 7: Draw minibatch of m noise samples $\{\mathbf{z}^{(1)}, \dots, \mathbf{z}^{(m)}\} \sim p_z$.
- 8: Update G 's parameters by gradient descent:

$$\nabla_{\theta_g} \frac{1}{m} \sum_{i=1}^m \log \left(1 - D(G(\mathbf{z}^{(i)})) \right)$$

- 9: **end for**

* k is a tunable hyperparameter which is usually set to 1 to lower the training cost of each iteration.

is drawn closer to p_{data} under the guidance of the gradient of D ; (d) final convergence: assuming sufficient model capacity, the adversarial pair reach the equilibrium $p_G = p_{data}$, where $D(\mathbf{x}) = \frac{1}{2}$.

Analysis. We now look at (b)-(d) in Figure 1 analytically to understand why the training scheme works. We begin by addressing (b): how does D converge to $D^*(\mathbf{x})$ in the inner loop?

Statement 1. For fixed G , the optimal discriminator D is

$$D_G^*(\mathbf{x}) = \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_G(\mathbf{x})} \quad (4)$$

Proof. Given a fixed G , the training objective for D is to maximize $J(D, G)$, where

$$\begin{aligned} J(G, D) &= \mathbb{E}_{\mathbf{x} \sim p_{data}} [\log D(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z} [\log(1 - D(G(\mathbf{z})))] \\ &= \int_{(\mathbf{x})} p_{data}(\mathbf{x}) \log(D(\mathbf{x})) dx + \int_{\mathbf{z}} p_z(\mathbf{z}) \log(1 - D(G(\mathbf{z}))) dz \\ &= \int_{\mathbf{x}} p_{data}(\mathbf{x}) \log(D(\mathbf{x})) + p_G(\mathbf{x}) \log(1 - D(\mathbf{x})) dx \end{aligned} \quad (5)$$

The last equality employs change of variable for the second term of the penultimate formula: as $G : \mathbf{z} \mapsto \mathbf{x}$, (i) integrating over $p_z(\mathbf{z})$ is equivalent to integrating over $p_G(\mathbf{x})$, and (ii) integrating over $G(\mathbf{z})$ is equivalent to integrating over \mathbf{x} . Further, we have⁶

$$\begin{aligned} D^* &= \operatorname{argmax}_D J(G, D) \\ &= \operatorname{argmax}_D \int_{\mathbf{x}} p_{data}(\mathbf{x}) \log(D(\mathbf{x})) + p_G(\mathbf{x}) \log(1 - D(\mathbf{x})) dx \end{aligned} \quad (6)$$

where the integrand can be abstracted in the form $f(D) = a \log(D) + b \log(1 - D)$. By setting $\frac{\partial f}{\partial D} \triangleq 0$ and solving for D , we have $D^* = \frac{a}{a+b}$, satisfying Eq. 4, concluding the proof. \square

Next we look at (c,b) and ask: how does updating G get us to $p_G = p_{data}$?

Statement 2. Let $C(G) = \max_D J(G, D)$, i.e. $C(G)$ is G 's minimization objective (cf. Eq. 2). The global minimum of $C(G)$ is reached iff $p_G = p_{data}$, at which point $C(G) = -\log 4$.

⁶The last equality in Eq. 6 is copied from the results of Eq. 5.

Proof. We know that

$$\begin{aligned}
 C(G) &= \max_D J(G, D) \\
 &= \mathbb{E}_{\mathbf{x} \sim p_{data}} [\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{z} \sim p_z} [\log(1 - D_G^*(G(\mathbf{z})))] \\
 &= \mathbb{E}_{\mathbf{x} \sim p_{data}} [\log D_G^*(\mathbf{x})] + \mathbb{E}_{\mathbf{x} \sim p_G} [\log(D_G^*(\mathbf{x}))] \quad (\text{by change of variable (cf. Eq. 5)}) \\
 &= \mathbb{E}_{\mathbf{x} \sim p_{data}} \left[\log \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_G(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_G} \left[\log \frac{p_G(\mathbf{x})}{p_{data}(\mathbf{x}) + p_G(\mathbf{x})} \right] \tag{7}
 \end{aligned}$$

We also know that if $p_{data} = p_G$, then $D_G^*(\mathbf{x}) = \frac{1}{2}$ (by Eq. 4), which we plug in Eq. 7 to obtain $C(G) = \log \frac{1}{2} + \log \frac{1}{2} = -\log 4$. We now show that $C(G) = -\log 4$ is the global minimum, reached when $p_G = p_{data}$:

$$\begin{aligned}
 C(G) &= \mathbb{E}_{\mathbf{x} \sim p_{data}} \left[\log \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_G(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_G} \left[\log \frac{p_G(\mathbf{x})}{p_{data}(\mathbf{x}) + p_G(\mathbf{x})} \right] \\
 &= -\log 4 + \log 4 + \mathbb{E}_{\mathbf{x} \sim p_{data}} \left[\log \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_G(\mathbf{x})} \right] + \mathbb{E}_{\mathbf{x} \sim p_G} \left[\log \frac{p_G(\mathbf{x})}{p_{data}(\mathbf{x}) + p_G(\mathbf{x})} \right] \\
 &= -\log 4 + \mathbb{E}_{\mathbf{x} \sim p_{data}} \left[\log \left(2 \cdot \frac{p_{data}(\mathbf{x})}{p_{data}(\mathbf{x}) + p_G(\mathbf{x})} \right) \right] + \mathbb{E}_{\mathbf{x} \sim p_G} \left[\log \left(2 \cdot \frac{p_G(\mathbf{x})}{p_{data}(\mathbf{x}) + p_G(\mathbf{x})} \right) \right] \\
 &= -\log 4 + \mathbb{E}_{\mathbf{x} \sim p_{data}} \left[\log \frac{p_{data}(\mathbf{x})}{\frac{p_{data}(\mathbf{x}) + p_G(\mathbf{x})}{2}} \right] + \mathbb{E}_{\mathbf{x} \sim p_G} \left[\log \frac{p_G(\mathbf{x})}{\frac{p_{data}(\mathbf{x}) + p_G(\mathbf{x})}{2}} \right] \\
 &= -\log 4 + KL \left(p_{data} \parallel \frac{p_{data}(\mathbf{x}) + p_G(\mathbf{x})}{2} \right) + KL \left(p_G \parallel \frac{p_{data}(\mathbf{x}) + p_G(\mathbf{x})}{2} \right) \\
 &= -\log 4 + KL(p_{data} \parallel p_G) + KL(p_G \parallel p_{data}) \quad (\text{by } p_G = p_{data}) \\
 &= -\log 4 + 2 \cdot \left(\frac{1}{2} \cdot KL(p_{data} \parallel p_G) + \frac{1}{2} \cdot KL(p_G \parallel p_{data}) \right) \\
 &= -\log 4 + 2 \cdot JS(p_{data} \parallel p_G) \tag{8}
 \end{aligned}$$

As the Jensen-Shannon divergence $JS(p_{data} \parallel p_G) \geq 0$, we have shown that $C^*(G) = -\log 4$ is indeed the global minimum, reached on the condition that $p_G = p_{data}$, concluding the proof. \square

Statement 1 and 2 show the optima in the alternating updates in Algorithm 1 lead us to the equilibrium of the Minmax Game. The original work [1] also proves the convergence of Algorithm 1. However the derivation involves advanced knowledge of optimization, I thus refer the interested reader to the proof there.

Implementation note. So far, for the convenience of exposition, we have been describing D as a single MLP. In practice, to implement the objective function (Eq. 1) properly, we actually need two MLPs D_1 and D_2 which share parameters. Concretely, the Python code⁷ for the objective functions for D is as follows [3]:

```

# Objective function for discriminator D
batch=tf.Variable(0)
obj_d=tf.reduce_mean(tf.log(D1)+tf.log(1-D2))
opt_d=tf.train.GradientDescentOptimizer(0.01)
      .minimize(1-obj_d,global_step=batch,var_list=theta_d)

```

3 Applications in NLP

As briefly mentioned in Section 1, recent years have seen applications of GAN in the space of NLP. In this section, I take Document Modeling [8] and Dialogue Generation [9] as examples to demonstrate the potential of GAN as a flexible and versatile learning framework.

⁷The complete code is at https://github.com/ericjang/genadv_tutorial/blob/master/genadv1.ipynb.

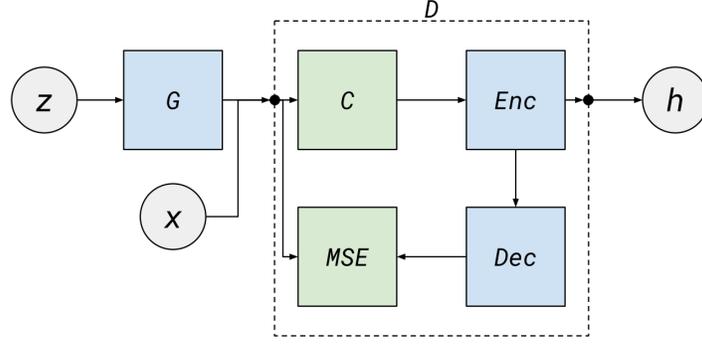


Figure 2: Adversarial Document Model: a variant of Energy-based GAN

Document modeling. Glover (2016) [8] proposes an *Energy-based GAN* [13] (ADM)⁸ where the generator is a regular MLP, while the discriminator is a *Denoising Autoencoder* (DAE) [14] (as the energy function [15]) instead. The model is pitched against a highly well-tuned *Restricted Boltzmann Machine* (RBM) [17] based system DocNADE [17], and a baseline (a stand-alone discriminative DAE classifier). The architecture is illustrated graphically in Figure 2.

In ADM, a document is modeled as a binary bag-of-words vector $\mathbf{x} \in \{0, 1\}^V$, where V is the size of vocabulary. D (the DAE discriminator) takes two inputs: (i) a “real” document vector \mathbf{x} , and (ii) a fake document vector $G(\mathbf{z})$ generated by G (an MLP). In input vector is first processed through a corruption process⁹ C to obtain vector \mathbf{x}^c and feed \mathbf{x}^c to a regular encode-decode component [14] (Enc, Dec).

$$\mathbf{h} = f(\mathbf{W}^e \mathbf{x}^c + \mathbf{b}_e) \quad (9)$$

$$\mathbf{y} = \mathbf{W}^d \mathbf{h} + \mathbf{b}_d \quad (10)$$

where $(\mathbf{W}^e, \mathbf{b}_e), (\mathbf{W}^d, \mathbf{b}_d)$ are the parameters of the encoder/decoder, respectively. The quality of the input is evaluated by *Mean Square Error* (MSE) as a metric for reconstruction error.

$$MSE = \frac{1}{V} \sum_{i=1}^V (\mathbf{x}_i - \mathbf{y}_i)^2 \quad (11)$$

In training, G and D strive to minimize their respective energy functions¹⁰:

$$f_D(\mathbf{x}, \mathbf{z}) = D(\mathbf{x}) + \max(0, m - D(G(\mathbf{z}))) \quad (12)$$

$$f_G(\mathbf{z}) = D(G(\mathbf{z})) \quad (13)$$

The process amounts to the same effect as Minmax Game (Eq. 3), in that D increases the score it assigns the real data (i.e. \mathbf{x}) while lowering the score for the fake data (i.e. $G(\mathbf{z})$) generated by G , and G increases the score of the fake data.

In his experiments, [8] formulates a document classification task with the 20 Newsgroup Dataset [18], where the model takes a query document \mathbf{d} as input and produces a set of output documents that are closest¹¹ to \mathbf{d} . The results are shown in Figure 3 as a precision-recall curve. While producing overall weaker performance than the state-of-the-art DocNADE, ADM demonstrates its power by defeating the strong baseline DAE models by a respectable margin, showing its promise¹² in more sophisticated formulation.

⁸The model is listed in [8] with the name *Adversarial Document Model*.

⁹Vincent et al. (2010) [14] show that formulating an autoencoder reconstruction task as a denoising task (with corrupted input) helps the autoencoder to generalize better.

¹⁰One may consider an energy function as a family of loss function [15].

¹¹By cosine similarity.

¹²In image generation tasks, while sometimes fall short behind VAE [19] in classification tasks [2], GAN often times generates more realistic images by human judgment. This thus leads us to reasonably believe in its potential in document-based tasks.

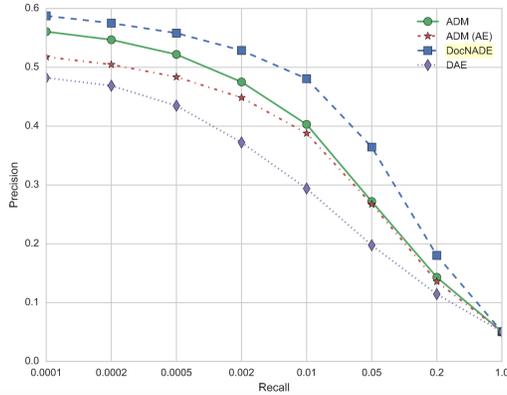


Figure 3: ADM results

| | |
|-------------------------|---|
| Input | tell me ... how long have you had this falling sickness ? |
| Vanilla-MLE | i 'm not a doctor . |
| Vanilla-Sample | well everything you did was totally untrue . |
| REINFORCE | i don 't know how long it 's been . |
| REGS Monte Carlo | A few months, I guess . |
| Input | so i had the doctors test sammy 's response to conditioning . |
| Vanilla-MLE | sammy wrote the test sammy wrote the test . |
| Vanilla-Sample | objects pick objects up objects objects objects objects |
| REINFORCE | i 'm not sure that 's relevant . |
| REGS Monte Carlo | so he took the pills . |
| Input | what are you going to do ? |
| Vanilla-MLE | i 'm going to the movies . |
| Vanilla-Sample | get him outta here first ! |
| REINFORCE | i 'm going to get you . |
| REGS Monte Carlo | i 'm going to get you a new car . |
| Input | they fear your power your intellect . |
| Vanilla-MLE | you 're the only one who knows what 's going on . |
| Vanilla-Sample | when they are conquered and you surrender they will control all of us . |
| REINFORCE | i 'm afraid i 'm not ready yet . |
| REGS Monte Carlo | i 'm not afraid of your power . |

Figure 4: Dialogue Generation (REINFORCE and REGS are GANs)

Dialogue generation. Li et al., (2017) [9] presents an interesting GAN-based dialogue generation system. They formulate a reinforcement learning [20] based *Turning Test* where the goal is for the generator G , under the guidance of a discriminator D , to learn to generate realistic responses to input sentences that are indistinguishable from responses given by humans. D guides G by giving it feedbacks in the form of reinforcement rewards: Positive reward for realistic responses, and negative reward for non-realistic ones.

Concretely, G takes the form of a sequence-to-sequence *Recurrent Neural Network* (RNN) [21] which generates response y based on a dialogue history x . D is an autoencoder-based binary classifier [22] that takes a (dialogue history, response) pair $\{x, y\}$, and produces a label indicating whether the input is generated by a human (i.e. real data) or G (i.e. fake data). It further returns a reward score to guide G : $Q_+(\{x, y\})$, $Q_-(\{x, y\})$ for positive and negative rewards, respectively.

Deferring the details of the model to the original work [9], I only show some sample outputs from its two variants¹³ evaluated therein to demonstrate the quality of the generated responses (Figure 4). Here, the *input* is a dialogue probe, and the following lines are the respective responses of the models give to the probe. The two GAN's produce apparently superior responses by human judgment. In addition to generating good responses, [9] also shows the reliability of their model (see Table 3 in the original work).

¹³REGS improves on REINFORCE by ameliorating its tendency for *mode collapse* — for G to generate the same fake data over and over again [2].

4 Conclusion

In this tutorial, I started by giving an in-detail description of GAN (Section 1) and step-by-step derivation in related proofs (Section 2). I then demonstrated the flexibility of GAN as a novel neural net architecture with examples of its application in NLP (Section 3).

References

- [1] Goodfellow, I.J., Pouget-Abadie, J., Mirza, M., Xu, B., Warde-Farley, D., Ozair, S., Courville, A. & Bengio, Y. (2014) *Generative Adversarial Nets*. NIPS 2014.
- [2] Goodfellow, I.J. (2016) *NIPS 2016 Tutorial: Generative Adversarial Networks*. NIPS 2016. CoRR, arXiv:1701.00160.
- [3] Jang, E. (2015) *Generative Adversarial Nets in Tensorflow*. Eric Jang's Blog: blog.evjang.com/2016/06/generative-adversarial-nets-in.html.
- [4] Lotter, W., Kreiman, G. & Cox, D. (2015) *Unsupervised Learning of Visual Structure Using Predictive Generative Networks*. CoRR, arXiv:1151.06380.
- [5] Ledig, C., Theis, L., Huszar, F., Caballero, J., Aitken, A. P., Tejani, A., Totz, J., Wang, Z., & Shi, W. (2016) *Photo-realistic single Image Super-resolution Using a Generative Adversarial Network*. CoRR, arXiv:1609.04802.
- [6] Isola, P., Zhu, J.-Y., Zhou, T., & Efros, A. A. (2016). *Image-to-image Translation with Conditional Adversarial Networks*. CoRR, arXiv:1611.07004.
- [7] Donahue, J., Krähenbühl, P. & Darrell, T. (2017) *Adversarial Feature Learning*. CoRR, arXiv:1605.09782.
- [8] Glover, J. (2016) *Modeling Documents with Generative Adversarial Networks*. CoRR, arXiv:1612.09122.
- [9] Li, J., Monroe, W., Shi, T., Jean, S., Ritter, A. & Jurafsky, D. (2017) *Adversarial Learning for Neural Dialogue Generation*. CoRR, arXiv:1701.06547.
- [10] Chen X., Athiwaratkun, B., Sun, Y., Weinberger, K. & Cardie, C. (2016) *Adversarial Deep Averaging Networks for Cross-lingual Sentiment Classification*. CoRR, arXiv:1606.01614.
- [11] Zhang, Y., Barzilay, R. & Jaakkola, T. (2017) *Aspect-augmented Adversarial Networks for Domain Adaptation*. CoRR, arXiv:1701.00188.
- [12] Ng, A. & Jordan, M.I. (2002) *On Discriminative vs. Generative Classifiers: A Comparison of Logistic Regression and Naïve Bayes*. NIPS 2002.
- [13] Zhao, J., Mathieu, M. & LeCun Y. (2016) *Energy-based Generative Adversarial Network*. CoRR, arXiv:1609.03126.
- [14] Vincent, P., Larochelle, H., Lajoie, I., Bengio, Y. & Manzagol, P-A. (2010) *Stacked Denoising Autoencoders: Learning Useful Representations in a Deep Network with a Local Denoising Criterion*. JMLR 2010.
- [15] LeCun, Y., Chopra, S., Hadsell, R., Ranzato, M-A. & Huang, F.J. (2006) *A Tutorial on Energy-based Learning*. In Bakir, G., Hofman, T., Schölkopf, B., Smola, A. & Taskar, B. (eds.) *Predicting Structured Data*. MIT Press.
- [16] Hinton, G.E. (2002) *Training Products of Experts by Minimizing Contrastive Divergence*. *Neural Computation*, vol. 14, no. 8, pp. 1607–1614.
- [17] Larochelle, H. & Lauly, S (2012) *A Neural Autoregressive Distribution Estimator*. NIPS 2012.
- [18] Lang, K. (1995) *Newsweeder: Learning to Filter News*. ICML 1995.
- [19] Kingma, D.P. & Welling, M. (2014) *Auto-encoding Variational Bayes*. ICLR 2014.
- [20] Williams, R.J. (1992) *Simple Statistical Gradient-following Algorithms for Connectionist Reinforcement Learning*. *Machine Learning*, vol. 8 (3-4): 229–256.
- [21] Sutskever, I., Vinyals, O. & Quoc, V.L. (2014) *Sequence to Sequence Learning with Neural Networks*. NIPS 2014.
- [22] Li, J., Luong, M-T. & Jurafsky, D. (2015) *A Hierarchical Neural Autoencoder for Paragraphs and Documents*. CoRR, arXiv:1506.01057.