

Abstract

106 Keywords: Deep Latent Variable Generative Models; Regularized Autoen-
coders; Latent Space; Few-shot Generation.

107 The rise of deep neural networks has significantly advanced unsupervised generative
108 modeling. Numerous Deep Generative Models (DGMs) have been proposed, including Vari-
109 ational Autoencoders (VAE) [KW14], Generative Adversarial Networks (GAN) [GPAM⁺14],
110 Wasserstein Autoencoders (WAE) [TBGS18], Adversarial Autoencoders (AAE) [MSJG16],
111 Autoregressive Models [vdOKE⁺16, OKK16], Normalizing Flow-based Models [KD18], Energy-
112 based Models [LCH⁺06, DM19], and Diffusion Models [HJA20]. These models may be cate-
113 gorized along various dimensions, such as their architectural variations, presence or absence
114 of explicit latent representation, training methodology and stability, density estimation capa-
115 bility, and time taken for sampling. In this thesis, our focus is on *Deep Latent Variable Gen-*
116 *erative Models* (DLVGMs), which refer to generative autoencoders (VAE, WAE, AAE) and
117 GANs. This class of models uses low-dimensional latent representations of high-dimensional
118 data, allowing learning informative low-dimensional representations for various downstream
119 tasks (clustering, classification, and disentangling generative factors) while facilitating novel
120 data generation. Interestingly, GANs are particularly notable for their high-quality gen-
121 eration but suffer from unstable training, mode collapse, and hyperparameter sensitivity.
122 In contrast, Autoencoders with regularized latent spaces (VAE, WAE, AAE) offer stable
123 training, interpretable inference, and efficient sampling, though their generated images are
124 visually less impressive than those from GANs. In this thesis, we intend to address some of
125 the complementary strengths and limitations of these DLVGMs. The thesis is organized in
126 several parts, as outlined below.

127 In Part I (Prologue) of this thesis, we introduce various deep generative frameworks,
128 define DLVGMs, highlight challenges (such as poor generation quality of RAEs, representa-
129 tion learning issues, and the need for large-scale data) associated with DLVGMs, and review
130 existing solutions.

131 In Part II, ‘Optimizing the Latent Space of RAEs for Improved Generation,’ we diagnose
132 the reasons behind the poor generation quality of generative AE frameworks by exploring
133 two questions: 1. What is the ‘optimal’ latent dimension [MCJ⁺20], and 2. What is the
134 ‘optimal’ latent prior [MASP21a] for a good generation? We hypothesize natural data
135 generation as a two-step process involving a true low-dimensional latent space and a non-
136 linear mapping to a high-dimensional data space. We show that under the assumption

137 of a Gaussian prior, the best generation quality is achieved when the dimensionality of
138 the generative AE’s bottleneck layer matches the true latent dimensionality [MCJ+20]. In
139 [MASP21a], we relax the Gaussian prior assumption to learn the prior flexibly, considering
140 the true latent dimensionality.

141 In Part III, ‘Optimizing the Latent Space of RAEs for Task-Specific Representation
142 Learning,’ we focus on representation learning in an RAE framework. First, we study
143 the impact of bias-variance trade-off due to fixed prior distribution versus learnable pri-
144 ors on representation quality and demonstrate that learning the prior flexibly helps the
145 model discover the actual data structure, improving clustering performance [MASP21b].
146 While [MASP21b] uses uni-modal data, we address disentangled representation learning in
147 a multi-modal setting in [MSSA23]. We decompose the joint latent space into continuous
148 and discrete components, each with domain-specific and domain-invariant representations.
149 We demonstrate the effectiveness of these disentangled joint representations in downstream
150 tasks like classification and generation. In our subsequent work [MST+23], we combine the
151 representation learning aspect with the generation ability of RAEs to develop a framework
152 for class-imbalance mitigation to enhance discriminating performance. Precisely, we propose
153 a minority oversampling method that is distance-metric-free and class-preserving by design.

154 In Part IV, ‘Few-shot Generation Using DLVGMs,’ we address the issue that while
155 DLVGMs offer a plethora of applications, they are data-hungry, limiting their applicability
156 in real-world scenarios with data scarcity. Specifically, we develop techniques to trans-
157 fer a source DLVGM built with large-scale data to a ‘close’ target domain with limited
158 data. In [MTSA23], we perform few-shot ‘generative domain adaptation’ via inference-time
159 latent-code learning by prepending a latent adapter network. While [MTSA23] achieves
160 high-quality generation, it requires considerable time to generate due to inference-time opti-
161 mization. In [MTSA24], we address this by learning to sample the parameters of the latent
162 adapter network using a hypernetwork.

163 Finally, in Part V, we summarize our contributions and propose directions for future
164 work based on the techniques and methods introduced in this thesis.