
Generating High Fidelity Synthetic Data via Coreset selection and Entropic Regularization

Anonymous Author(s)

Affiliation

Address

email

Abstract

1 Generative models have the ability to synthesize data points drawn from the data
2 distribution, however, not all generated samples are high quality. In this paper,
3 we propose using a combination of coresets selection methods and “entropic
4 regularization” to select the highest fidelity samples. We leverage an Energy-Based
5 Model which resembles a variational auto-encoder with an inference and generator
6 model for which the latent prior is complexified by an energy-based model. In a
7 semi-supervised learning scenario, we show that augmenting the labeled data-set,
8 by adding our selected subset of samples, leads to better accuracy improvement
9 rather than using all the synthetic samples.

10 1 Introduction

11 In machine learning, augmenting data-sets with synthetic data has become a common practice
12 which potentially provides significant improvements in downstream tasks such as classification. For
13 example, in the case of images, recent methods like MixMatch, FixMatch and Mean Teacher [1] [12]
14 [13] have proposed data augmentation techniques which rely on simple pre-defined transformations
15 such as cropping, resizing, etc.

16 However, generating augmentations is not as straightforward in all modalities. Hence, one suggestion
17 is to use samples from generative models to augment the data-sets. One issue that arises is that simply
18 augmenting a data-set using a generative model can often lead to the degradation of classification
19 accuracy due to some poor samples drawn from the generator. The question arises: can we filter the
20 lower quality generated samples to avoid degradation in accuracy? In our method we select a subset
21 of synthetic samples which have high fidelity to the underlying data-set via CRAIG [6], additionally
22 we introduce “entropic regularization” by filtering samples with low entropy over the latent classifier.

23 In semi-supervised learning, the goal is to learn a classifier model which maintains high classification
24 accuracy while reducing the number of labeled observed examples. Generative modeling and
25 especially likelihood-based learning is a principled formulation for unsupervised and semi-supervised
26 learning. Within this family of models, energy-based models (EBM) are particularly convenient for
27 semi-supervised learning, as they may be interpreted as generative classifiers. That is, we not only
28 have access to the class predictions but may also draw samples from the model.

29 Another direction in supervised learning is to reduce the amount of computation involved in training
30 a model by reducing the data-set to a smaller subset. Such sets are coined *coresets* as a smaller set of
31 representative points attempts to approximate the geometry of a larger point set under some metric.
32 Recent art [6] introduces a novel algorithm CRAIG which constructs a weighted coreset such that the
33 gradient over the full training data-set is closely estimated, which allows for gradient descent on the
34 smaller coreset with considerable improvement in the sample- and computational-efficiency.

35 In this work, we show that semi-supervised learning and coreset subset selection are complementary
 36 and improve generalization as well as generation quality. First, a generative classifier is learned on a
 37 large set of unlabelled data and a small set of labeled data pairs. Then, the generative model is utilized
 38 to draw class conditional samples which augment the labeled data pairs. As such augmentation might
 39 be a considerably large set, in fact, we can draw infinite samples from the generative model, we
 40 recruit CRAIG to reduce the conditional samples to a much smaller coreset while approximately
 41 maintaining the full gradient over the cross-entropy term. As the generative model might synthesize
 42 conditional samples of low quality or even incorrect class identity, we apply an entropic filter to
 43 remove noisy samples. By learning a joint generative classifier we learn a generator that can produce
 44 samples that improve classification accuracy as well as a classifier that can boost generative capacity
 45 and quality.

46 This method may be interpreted as a learned (and filtered) data augmentation as opposed to classical
 47 data augmentation in which the set of augmentation functions (e.g., convolution with Gaussian noise,
 48 horizontal or vertical flipping, etc.) is pre-defined and could be specific to a data-set or modality. We
 49 demonstrate the efficacy of the method by a significant improvement in classification performance.

50 2 Synthetic Data Generation for Semi-Supervised Learning

51 **Notation** Let $x \in \mathcal{R}^D$ be an observed example. Let y be a K -dimensional one-hot vector as the
 52 label for classification with K categories. Suppose $\mathcal{L} = \{(x_i, y_i) \in \mathbb{R}^D \times \{k\}_{k=1}^K, i = 1, \dots, M\}$
 53 denotes a set of labeled examples where K indicates the number of categories and $\mathcal{U} = \{x_i \in$
 54 $\mathbb{R}^D, i = M + 1, \dots, M + N\}$ denotes a set of unlabeled examples.

55 **Semi-Supervised Learning** Let $p_\theta(y | x)$ denote a soft-max classifier with parameters θ . The goal
 56 of semi-supervised learning is to learn θ with “good” generalization while decreasing the number of
 57 labeled examples M .

58 2.1 Latent Energy Based Model

59 Let $z \in \mathbb{R}^d$ be the latent variables, where $D \gg d$. We assume a Markov chain $y \rightarrow z \rightarrow x$. Then the
 60 joint distribution of (y, z, x) is

$$p_\theta(y, z, x) = p_\alpha(y, z) p_\beta(x|z), \quad (1)$$

61 where $p_\alpha(y, z)$ is the prior model with parameters α , $p_\beta(x|z)$ is the top-down generation model with
 62 parameters β , and $\theta = (\alpha, \beta)$. Then, the prior model $p_\alpha(y, z)$ is formulated as an energy-based
 63 model [10],

$$p_\alpha(y, z) = Z(\alpha)^{-1} \exp(F_\alpha(z)[y]) p_0(z). \quad (2)$$

64 where $p_0(z)$ is a reference distribution, assumed to be isotropic Gaussian. $F_\alpha(z) \in \mathbb{R}^K$ is param-
 65 eterized by a multi-layer perceptron. $F_\alpha(z)[y]$ is y th element of $F_\alpha(z)$, indicating the conditional
 66 negative energy. $Z(\alpha)$ is the partition function. In the case where the label y is unknown, the prior
 67 model $p_\alpha(z) = \sum_y p_\alpha(y, z) = Z(\alpha)^{-1} \sum_y \exp(F_\alpha(z)[y]) p_0(z)$. Taking log of both sides:

$$\log p_\alpha(z) = \log \sum_y \exp(F_\alpha(z)[y]) + \log p_0(z) - \log Z(\alpha), \quad (3)$$

68 The prior model can be interpreted as an energy-based correction or exponential tilting of the reference
 69 distribution, p_0 . The correction term is $F_\alpha(z)[y]$ conditional on y , while it is $\log \sum_y \exp(F_\alpha(z)[y])$
 70 when y is unknown. Denote

$$f_\alpha(z) = \log \sum_y \exp(F_\alpha(z)[y]), \quad (4)$$

and then $-f_\alpha(z)$ is the free energy [2]. The soft-max classifier is $p_\alpha(y|z) \propto \exp(\langle y, F_\alpha(z) \rangle) =$
 $\exp(F_\alpha(z)[y])$.

71 The generation model is the same as the top-down network in VAE [4], $x = g_\beta(z) + \epsilon$, where
 72 $\epsilon \sim \mathcal{N}(0, \sigma^2 I_D)$, so that $p_\beta(x|z) \sim \mathcal{N}(g_\beta(z), \sigma^2 I_D)$.

73 We use variational inference to learn our latent space EBM by minimizing the evidence lower bound
 74 (ELBO) over our energy, encoder, and generator models jointly. Refer to appendix B for more details
 75 about learning the model.

76 In summary, we can use the above model to i) classify data points ii) generate class-conditional
 77 samples iii) compute entropy for each generated sample. We will leverage these properties in the
 78 later sections to get better augmentation for our data-set.

79 2.2 Sampling Synthetic data from the EBM

80 Naturally, increasing the cardinality of the set of labeled samples \mathcal{L} may improve the classification
 81 accuracy of soft-max classifier $p_\theta(y|x)$. In the case of image models, traditional methods recruit a set
 82 of transformations or permutations of x such as convolution with Gaussian noise or random flipping.
 83 Instead we leverage the learned top-down generator $p_\beta(x|z)$ to augment \mathcal{L} with class conditional
 84 samples. This is beneficial as (1) the generative path is readily available as an auxiliary model of
 85 learning the variational posterior $q_\phi(z|x)$ by auto-encoding variational Bayes, (2) hand-crafting of
 86 data augmentation is domain and modality-specific, and (3) in principle the number of conditional
 87 ancestral samples is infinite and might capture the underlying data distribution well.

88 We may construct the augmented set of L labelled samples $\mathcal{L}^+ = \{(x_i, y_i)\}$ by drawing conditional
 89 latent samples from the energy-based prior model $p_\alpha(y, z)$ in the form of Markov chains. Then, we
 90 obtain data space samples by sampling from the generator $p_\beta(x|z)$.

91 First, for each label y , we draw an equal number of samples $\mathcal{Z} = \{z_i\}$ in latent space. One convenient
 92 MCMC is the overdamped Langevin dynamics, which we run for T_{LD} steps with target distribution
 93 $p_\alpha(y, z)$,

$$z \sim p_0(z), \quad (5)$$

$$z_{t+1} = z_t + s \nabla_z [f_\alpha(z)[y] - \|z\|^2/2] + \sqrt{2s} \epsilon_t, t = 1, \dots, T_{LD} \quad (6)$$

94 with negative conditional energy $f_\alpha(z)[y]$, discretization step size s , and isotropic $\epsilon_t \sim N(0, I)$.

95 Then, we draw conditional samples $\{x_i\}$ in data space given $\{z_i\}$ from the top-down generator
 96 model $p_\beta(x|z)$,

$$\mathcal{L}^+ = \{(x_i \sim p_\beta(x|z_i), y_i) \mid i = M + N, \dots, M + N + L\} \quad (7)$$

97 which results in an augmented data-set of L class conditional samples.

98 2.3 Entropic Regularization

99 When learning the generative classifier on both labelled samples \mathcal{L} and the above naive construction
 100 of augmentation \mathcal{L}^+ , the classification accuracy tends to be worse than solely learning from \mathcal{L} .
 101 As depicted in Figure 1a, a few conditional samples suffer from either low visual fidelity or even
 102 incorrect label identity. This reveals the implicit assumption of our method is that $p_\beta(x|z)p_\alpha(z|y)$ is
 103 reasonable “close” to the true class conditional distribution $p(x|y)$ under some measure of divergence,
 104 which is not guaranteed.

105 To address the issue of outliers, we propose to exclude conditional samples for which the entropy in
 106 logits $\mathcal{H}(p_\theta(y|z))$ exceeds some threshold \mathcal{T} . We propose the following criteria for outlier detection,

$$\mathcal{H}(z) = - \sum_y p_\theta(z|y) \log p_\theta(z|y). \quad (8)$$

107 Note, (8) is the classical Shannon entropy of over the soft-max normalized logits of the classifier.
 108 Then, we may construct a more faithful data augmentation as follows,

$$\mathcal{Z}_\mathcal{T} = \{z_i \sim p(z|y_i) \mid \mathcal{H}(z_i) < \mathcal{T}, i = M + N, \dots, M + N + L\}, \quad (9)$$

$$\mathcal{L}_\mathcal{H}^+ = \{(x_i \sim p_\beta(x|z_i), y_i) \mid i = M + N, \dots, M + N + L\}. \quad (10)$$

109 Figure 1b depicts conditional samples sorted by $\mathcal{H}(z)$ for which samples with relatively large Shannon
 110 entropy suffer from low visual fidelity.

111 The learning and sampling algorithm is described in Algorithm 1 (appendix) as an extension of [10].

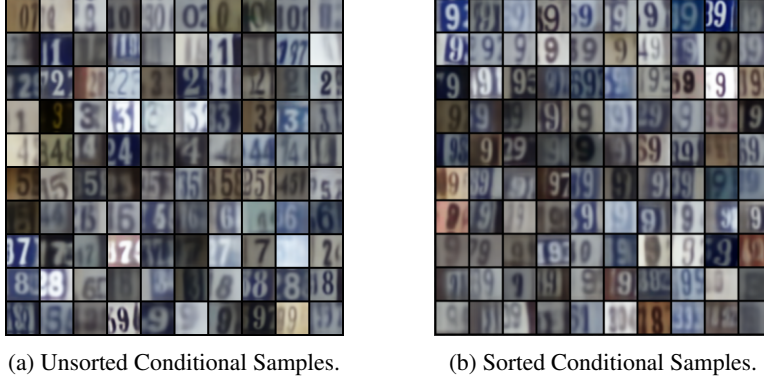


Figure 1: Class conditional samples drawn from $p_\beta(x|z)p_\alpha(z|y)$. (a) Outliers suffer from low visual fidelity (e.g., the last sample in the row of “ones”) or wrong label identity (e.g., the last image of row of “sevenths”). (b) Conditional samples sorted by increasing Shannon entropy $\mathcal{H}(z)$ over the logits.

112 2.4 Coreset Selection

113 Training machine learning models on large data-sets incur considerable computational costs. There
 114 has been substantial effort to develop subset selection methods that can carefully select a subset of the
 115 training samples that generalize on par with the entire training data [6] [11]. Since we can generate
 116 virtually infinite amount of synthetic samples, we must select the best subset of points to augment
 117 our base data-set with. Intuitively CRAIG selects a subset that can best cover the gradient space of the
 118 full data-set. It does this by selecting exemplar medoids from clusters of datapoints in the gradient
 119 space. As a bi-product, CRAIG robustly rejects noisy and even poisoned datapoints. The subset
 120 coreset algorithm ADACORE improves on CRAIG’s results by selecting diverse subsets [11]. Utilizing
 121 coreset methods allows us to select samples from the generator that is representative of the ground
 122 truth data-set while rejecting points that may negatively impact our network performance.

123 Formally, the CRAIG [6] algorithm aims to identify the smallest subset $S \subset V$ and corresponding
 124 per-element stepsizes $\gamma_j > 0$ that approximate the full gradient with an error at most $\epsilon > 0$ for all the
 125 possible values of the optimization parameters $w \in \mathcal{W}$.

$$S^* = \arg \min_{S \subset V, \gamma_j \geq 0 \forall j} |S|, \text{ s.t. } \max_{w \in \mathcal{W}} \left\| \sum_{i \in V} \nabla f_i(w) - \sum_{j \in S} \gamma_j \nabla f_j(w) \right\| \leq \epsilon \quad (11)$$

126 For deep neural networks it is more costly to calculate the above metric than to calculate vanilla SGD,
 127 In deep neural networks, the variation of the gradient norms is mostly captured by the gradient of the
 128 loss w.r.t the inputs of the last layer L . [6] shows that the normed gradient difference between data
 129 points can be efficiently bounded approximately by

$$\|\nabla f_i(w) - \nabla f_j(w)\| \leq c_1 \left\| \Sigma'_L \left(z_i^{(L)} \right) \nabla f_i^{(L)}(w) - \Sigma'_L \left(z_j^{(L)} \right) \nabla f_j^{(L)}(w) \right\| + c_2 \quad (12)$$

130 where $z_i^{(l)} = w^{(l)} x_i^{(l-1)}$. This upper bound is only slightly more expensive than calculating the loss.
 131 In the case of cross entropy loss with soft-max as the last layer, the gradient of the loss w.r.t. the
 132 i -th input of the soft-max is simply $p_i - y_i$, where p_i are logits and y is the one-hot encoded label.
 133 As such, for this case CRAIG does not need a backward pass or extra storage. This makes CRAIG
 134 practical and scalable tool to select higher quality generated synthetic data points.

135 2.5 Implicit learned data augmentation

136 In the following, we will re-interpret the above explicit data augmentation and entropic regularization
 137 into an implicit augmentation which can be merged into a simple term of the learning objective
 138 function.

139 The assumed Markov chain underlying the model is $y \rightarrow z \rightarrow x$. Let $\hat{z} \sim q_\phi(z|x)$ de-
 140 note the conditional sample \hat{z} from the approximate posterior given an observation x . Let
 141 $\hat{y} \sim p_\theta(y|\hat{z})$ denote the predicted label for which the logits of C classes are given as $F_\alpha(z) =$
 142 $(F_\alpha(z)[1], F_\alpha(z)[2], \dots, F_\alpha(z)[C])$.

143 The factorization which recruits the log-sum-exp lifting (3) as exponential tilting of the the reference
 144 distribution $p_0(z)$ so that the conditional $p_\alpha(y|z)$ is defined, and, amortized inference (19) with
 145 variational approximation of the posterior $q_\phi(z|x)$. These conditional distributions allow us to
 146 express learned data augmentation as the chain,

$$y \xrightarrow{q_T(z|y)} z \xrightarrow{p_\theta(x|z)} x \xrightarrow{q_\phi(z|x)} \hat{z} \xrightarrow{F_\alpha(z)[y]} \hat{y}. \quad (13)$$

147 in which the conditional $z|y$ is given as a MCMC dynamics. Specifically, we define $q_T(z|y)$ as K -
 148 steps of an overdamped Langevin dynamics on the learned energy-based prior $\exp(F_\alpha(z)[y])p_0(z)$,
 149 which iterates

$$z_{k+1} = z_k + s\nabla_z \log p(z_k|y) + \sqrt{2Ts}\epsilon_k, \quad k = 0, \dots, K-1, \quad (14)$$

150 with discretization step-size s , temperature T and isotropic noise $\epsilon_k \sim N(0, I)$.

151 For the (labeled) data distribution p_{data} the labels y are known. For the data augmentation, we
 152 assume a discrete uniform distribution over labels $y \sim U\{1, C\}$. Then, we define augmentation of
 153 synthesized examples as the marginal distribution

$$p_{\text{aug}}(x) = E_y E_{z|y} [p(x|z)p(z|y)]. \quad (15)$$

154 Then, we may introduce an augmented data-distribution as the mixture of the underlying labeled
 155 data-distribution p_{data} and the augmentation p_{aug} and mixture coefficient λ ,

$$p_\lambda(x) = \lambda p_{\text{data}}(x) + (1 - \lambda)p_{\text{aug}}(x). \quad (16)$$

156 As we have access to $p_\theta(y|x) = E_{p_\theta(z|x)} p_\theta(y|z)$ and can extend the objective to minimize the KL
 157 divergence under the augmented data distribution such that the labels y of (labeled) p_{data} and p_{aug}
 158 are recovered under the model,

$$E_{p_\lambda(x)} [KL(p(y|x) \| p(\hat{y}|x))]. \quad (17)$$

159 In information theory, the Kraft-McMillian theorem relates the relative entropy $KL(p\|q) =$
 160 $E_p[\log p/q]$ to the Shannon entropy $H(p)$ and cross entropy $H(p, q)$,

$$KL(p\|q) = H(p, q) - H(p). \quad (18)$$

161 In our case, the first term reduces to soft-max cross entropy over the (labeled) data distribution p_{data}
 162 and sampled labels $y \sim U\{1, C\}$. Hence, to minimize the above divergence, we must minimize the
 163 cross entropy which is consistent with classical learning of discriminative models. However, note that
 164 in our case the steps in (13) are fully differentiable, so that the data augmentation itself turns into an
 165 implicit term in the unified objective function rather than an explicitly constructed set of examples.

166 Lastly, we wish to re-introduce the entropic regularization for implicit data augmentation. Note,
 167 the entropic filter can be interpreted as a hard threshold on $H(p(\hat{y}|x)) < \mathcal{T}$. Here the Langevin
 168 dynamics q_T on z maximizes the logit $F_\alpha(z)[y]$, i.e. minimizes $H(p(\hat{y}|x))$, for which the Wiener
 169 process materialized in the noise term $\sqrt{2Ts}\epsilon_k$ with temperature T introduces randomness, or,
 170 smoothens the energy potential such that the dynamics converges towards the correct stationary
 171 distribution. High temperature T leads to Brownian motion, while low T leads to gradient descent.
 172 We realize that T controls $H(p(\hat{y}|x))$ as it may be interpreted as a soft or stochastic relaxation of \mathcal{T} .
 173 That is, we can express the entropic filter in terms of the temperature T of q_T and only need to lower
 174 T to obtain synthesized samples with associated low entropy in the class logits.

175 3 Experiments: Learning data augmentation

176 We evaluate our method on standard semi-supervised learning benchmarks for image data. Specif-
 177 ically, we use the street view house numbers (SVHN) [8] data-set with 1,000 labeled images and
 178 64,932 unlabeled images. The inference network is a standard Wide ResNet [14]. The generator
 179 network is a standard 4-layer de-convolutional network as regularly used in DC-GAN. The energy-
 180 based model is a fully connected network with 3 layers. Adam [3] is adopted for optimization with
 181 batch-sizes $n = m = l = 100$. The models are trained for $T = 1, 200, 000$ steps with augmentation
 182 after $T_a = 600, 000$ steps. The short-run MCMC dynamics in (6) is run for $T_{LD} = 60$ steps.

183 At iteration T_a , we take L class conditional samples from the generator with an equal amount of
 184 samples ($L/10$ for each digit). We filter conditional samples based on \mathcal{H} as described in Section 2.3
 185 for which the threshold $\mathcal{T} = 1e-6$ was determined by grid search. Next, we run CRAIG on the
 186 generated samples to keep a subset of 10% of the samples. For these additional examples, we compute
 187 the soft-max cross-entropy gradient with per-example weights obtained by CRAIG and update the
 188 model with step size $\eta_3 < \eta_2$ or a loss coefficient of 0.1 to weaken the gradient of $\mathcal{L}_{\mathcal{H}}^+$ relative to
 189 the original labeled data \mathcal{L} . Additionally, for every 10,000 iteration, we rerun CRAIG to choose an
 190 updated subset of generated samples.

Method	L				
	0	10,000	40,000	100,000	200,000
Baseline	92.0 ± 0.1	88.1 ± 0.1	-	-	-
\mathcal{H}	-	93.5 ± 0.1	93.8 ± 0.1	-	-
\mathcal{H} & CRAIG	-	93.0 ± 0.1	93.5 ± 0.1	93.9 ± 0.1	93.9 ± 0.1
\mathcal{H} & CRAIG & PL	-	-	94.5 ± 0.1	-	-

Table 1: Test accuracy with varied number of conditional samples L on SVHN [8].

191 Table 1 depicts results for the test accuracy on SVHN for a varied number of conditional samples L .
 192 First, we learned the model without data augmentation as a baseline. Then, we draw L conditional
 193 samples without an entropic filter and observe worse classification performance. As described earlier,
 194 we introduce the entropic filter \mathcal{H} to eliminate conditional samples of low quality which leads to a
 195 significant improvement in classification performance with increasing L . Finally, we combine both
 196 the entropic filter \mathcal{H} and coreset selection by CRAIG to further increase L . For $L = 10,000$ there
 197 is a significant improvement in classification accuracy when introducing CRAIG, which however
 198 decreases with increasing L . Lastly, to further boost accuracy we pseudo-label unlabeled data points
 199 from the SVHN data-set using the latent classifier. We reject data points whose entropy over the
 200 latent classifier is above 10^{-6} .

201 4 Conclusion

202 In the setting of semi-supervised learning, we have investigated the idea of combining generative
 203 models with a coreset selection algorithm, CRAIG. Such a combination is appealing as a generative
 204 model can in theory sample an infinite amount of labeled data, while a coreset algorithm can reduce
 205 such a large set to a much smaller informative set of synthesized examples. Moreover, learned
 206 augmentation is useful as many discrete data modalities such as text, audio, graphs, and molecules do
 207 not allow the definition of hand-crafted semantically invariant augmentations (such as rotations for
 208 images) easily.

209 We illustrated that a naive implementation of this simple result deteriorates the performance of the
 210 classifier in terms of accuracy over a baseline without such data augmentation. The underlying issue
 211 here was isolated to being related to the Shannon entropy in the predicted logits over classes for a
 212 synthesized example. High entropy indicates samples with low visual fidelity or wrong class identity,
 213 which may confuse the discriminative component of the model and lead to a loop in which uncertainty
 214 in the predictions leads to worse synthesis. In the first attempt, we constraint the class entropy in the
 215 set of augmented examples by taking a subset of the generated data-set with a hard threshold on the
 216 Shannon entropy. This resulted in significant empirical improvement of classification accuracy of
 217 two percentage points on SVHN. Moreover, we introduced pseudo labels which further improved
 218 performance.

219 Then, we show that the latent energy-based model with symbol-vector couplings has conditional
 220 distributions for end-to-end training of learned augmentations readily available. We formulate learned
 221 data augmentation as the KL-divergence between two known conditional distributions, show the
 222 relation to cross-entropy, and relax the entropy regularization into the temperature of the associated
 223 Langevin dynamics. This not only allows learning data augmentations as an alteration of the learning
 224 objective function but also paves the way toward a theoretical analysis.

225 **References**

- 226 [1] David Berthelot, Nicholas Carlini, Ian Goodfellow, Nicolas Papernot, Avital Oliver, and Colin A
227 Raffel. Mixmatch: A holistic approach to semi-supervised learning. In *Advances in Neural*
228 *Information Processing Systems*, pages 5049–5059, 2019.
- 229 [2] Will Grathwohl, Kuan-Chieh Wang, Joern-Henrik Jacobsen, David Duvenaud, Mohammad
230 Norouzi, and Kevin Swersky. Your classifier is secretly an energy based model and you should
231 treat it like one. *arXiv preprint arXiv:1912.03263*, 2019.
- 232 [3] Diederik P Kingma and Jimmy Ba. Adam: A method for stochastic optimization. *arXiv preprint*
233 *arXiv:1412.6980*, 2014.
- 234 [4] Diederik P Kingma and Max Welling. Auto-encoding variational bayes. *arXiv preprint*
235 *arXiv:1312.6114*, 2013.
- 236 [5] Dong-Hyun Lee. Pseudo-label: The simple and efficient semi-supervised learning method for
237 deep neural networks. In *Workshop on challenges in representation learning, ICML*, volume 3,
238 2013.
- 239 [6] Baharan Mirzasoleiman, Jeff Bilmes, and Jure Leskovec. Coresets for data-efficient training of
240 machine learning models. *arXiv preprint arXiv:1906.01827*, 2019.
- 241 [7] Takeru Miyato, Shin-ichi Maeda, Masanori Koyama, and Shin Ishii. Virtual adversarial training:
242 a regularization method for supervised and semi-supervised learning. *IEEE transactions on*
243 *pattern analysis and machine intelligence*, 41(8):1979–1993, 2018.
- 244 [8] Yuval Netzer, Tao Wang, Adam Coates, Alessandro Bissacco, Bo Wu, and Andrew Y Ng.
245 Reading digits in natural images with unsupervised feature learning. 2011.
- 246 [9] Erik Nijkamp, Bo Pang, Tian Han, Linqi Zhou, Song-Chun Zhu, and Ying Nian Wu. Learning
247 multi-layer latent variable model via variational optimization of short run mcmc for approximate
248 inference. *stat*, 1050:17, 2020.
- 249 [10] Bo Pang, Tian Han, Erik Nijkamp, Song-Chun Zhu, and Ying Nian Wu. Learning latent space
250 energy-based prior model. *arXiv preprint arXiv:2006.08205*, 2020.
- 251 [11] Omead Pooladzandi, David Davini, and Baharan Mirzasoleiman. Adaptive second order
252 coresets for data-efficient machine learning. In Kamalika Chaudhuri, Stefanie Jegelka, Le Song,
253 Csaba Szepesvari, Gang Niu, and Sivan Sabato, editors, *Proceedings of the 39th International*
254 *Conference on Machine Learning*, volume 162 of *Proceedings of Machine Learning Research*,
255 pages 17848–17869. PMLR, 17–23 Jul 2022.
- 256 [12] Colin Raffel. Fixmatch: Simplifying semi-supervised learning with consistency and confidence.
257 2020.
- 258 [13] Antti Tarvainen and Harri Valpola. Mean teachers are better role models: Weight-averaged
259 consistency targets improve semi-supervised deep learning results. In *Advances in neural*
260 *information processing systems*, pages 1195–1204, 2017.
- 261 [14] Sergey Zagoruyko and Nikos Komodakis. Wide residual networks. *arXiv preprint*
262 *arXiv:1605.07146*, 2016.

Algorithm 1: Semi-supervised learning of generative classifier with coreset selection.

input : Learning iterations T , augmentation iteration T_a , learning rates $(\eta_0, \eta_1, \eta_2, \eta_3)$, initial parameters $(\alpha_0, \beta_0, \phi_0)$, observed unlabelled examples $\{x_i\}_{i=1}^M$, observed labelled examples $\{(x_i, y_i)\}_{i=M+1}^{M+N}$, unlabelled, labelled and augmented batch sizes (n, m, l) , number of augmented samples L , entropy threshold \mathcal{T} , and number of Langevin dynamics steps T_{LD} .

output : $(\alpha_T, \beta_T, \phi_T)$.

for $t = 0 : T - 1$ **do**

1. **Mini-batch:**

Sample $\{x_i\}_{i=1}^m \subset \mathcal{U}$, $\{x_i, y_i\}_{i=m+1}^{m+n} \subset \mathcal{L}$, and $\{x_i, y_i\}_{i=m+n+1}^{m+n+l} \subset \mathcal{L}_{\mathcal{H}}^+$.

2. **Prior sampling:**

For each unlabelled x_i , initialize a Markov chain $z_i^- \sim q_\phi(z|x_i)$ and update by MCMC with target distribution $p_\alpha(z)$ for T_{LD} steps.

264

3. **Posterior sampling:**

For each x_i , sample $z_i^+ \sim q_\phi(z|x_i)$ using the inference network and reparameterization trick.

4. **Unsupervised learning of prior model:**

$\alpha_{t+1} = \alpha_t + \eta_0 \frac{1}{m} \sum_{i=1}^m [\nabla_\alpha F_{\alpha_t}(z_i^+) - \nabla_\alpha F_{\alpha_t}(z_i^-)]$.

5. **Unsupervised learning of inference and generator models:**

$\psi_{t+1} = \psi_t + \eta_1 \frac{1}{m} \sum_{i=1}^m [\nabla_\psi [\log p_{\beta_t}(x|z_i^+)] - \nabla_\psi \text{KL}(q_{\phi_t}(z|x_i) \| p_0(z)) + \nabla_\psi [F_{\alpha_t}(z_i^+)]]$.

6. **Supervised learning of prior and inference model:**

$\theta_{t+1} = \theta_t + \eta_2 \frac{1}{n} \sum_{i=m+1}^{m+n} \sum_{k=1}^K y_{i,k} \log(p_{\theta_t}(y_{i,k}|z_i^+))$.

7. **Augment at iteration T_a :**

$\mathcal{Z}_{\mathcal{T}} = \{z_i \sim p(z|y_i) \mid \mathcal{H}(z_i) < \mathcal{T}, i = M + N, \dots, M + N + L\}$,

$\mathcal{L}_{\mathcal{H}}^+ = \{(x_i \sim p_\beta(x|z_i), y_i) \mid i = M + N, \dots, M + N + L\}$.

8. **Approximate the gradient below with CRAIG after iteration T_a according to (12):**

$\theta_{t+1} = \theta_{t+1} + \eta_3 \frac{1}{n} \sum_{i=n+m+1}^{m+n+l} \sum_{k=1}^K y_{i,k} \log(p_{\theta_t}(y_{i,k}|z_i^+))$.

265 **B Learning the model with variational inference**

266 Given a data point in the unlabeled set, $x \in \mathcal{U}$, the the log-likelihood $\log p_\theta(x)$ is lower bounded by
267 the evidence lower bound (ELBO),

$$\text{ELBO}(\theta) = \mathbb{E}_{q_\phi(z|x)}[\log p_\beta(x|z)] - D_{KL}[q_\phi(z|x) \| p_\alpha(z)] \quad (19)$$

268 where $\theta = \{\alpha, \beta, \phi\}$ is overloaded for simplicity and $q_\phi(z|x)$ is a variational posterior, an approxi-
269 mation to the intractable true posterior $p_\theta(z|x)$.

270 For the prior model, the learning gradient for an example is

$$\nabla_\alpha \text{ELBO}(\theta) = \mathbb{E}_{q_\phi(z|x)}[\nabla_\alpha f_\alpha(z)] - \mathbb{E}_{p_\alpha(z)}[\nabla_\alpha f_\alpha(z)] \quad (20)$$

271 where $f_\alpha(z)$ is the negative free energy defined in equation (4), $\mathbb{E}_{q_\phi(z|x)}$ is approximated by samples
272 from the variational posterior and $\mathbb{E}_{p_\alpha(z)}$ is approximated with short-run MCMC chains [9] initialized
273 from the variational posterior $q_\phi(z|x)$.

274 Let $\psi = \{\beta, \phi\}$ collects parameters of the inference and generation models, and the learning gradients
275 for the two models are,

$$\nabla_\psi \text{ELBO}(\theta) = \nabla_\psi \mathbb{E}_{q_\phi(z|x)}[\log p_\beta(x|z)] - \nabla_\psi D_{KL}[q_\phi(z|x) \| p_0(z)] + \nabla_\psi \mathbb{E}_{q_\phi(z|x)} f_\alpha(z) \quad (21)$$

276 where $D_{KL}[q_\phi(z|x) \| p_0(z)]$ is tractable and the expectation in the other two terms is approximated
277 by samples from the variational posterior distribution $q_\phi(z|x)$.

278 For one example of labeled data, $(x, y) \in \mathcal{L}$, the log-likelihood can be decomposed $\log p_\theta(x, y) =$
279 $\log p_\theta(x) + \log p_\theta(y|x)$. While we optimize $\log p_\theta(x)$ as the unlabeled data, we maximize $\log p_\theta(y|x)$
280 by minimizing the cross-entropy as in standard classifier training. Notice that given the Markov chain
281 assumption $y \rightarrow z \rightarrow x$, we have

$$p_\theta(y|x) = \int p_\theta(y|z)p_\theta(z|x)dz = \mathbb{E}_{p_\theta(z|x)}p_\theta(y|z) \approx \mathbb{E}_{q_\phi(z|x)} \frac{\exp(F_\alpha(x)[y])}{\sum_k \exp(F_\alpha(x)[k])}. \quad (22)$$

282 In the last step, the true posterior $p_\theta(z|x)$ which requires expensive MCMC is approximated by the
283 amortized inference $q_\phi(z|x)$.

284 C Related Work

285 **Data augmentation.** Semi-supervised models with purely discriminative learning mostly rely on
286 data augmentation which exploit the class-invariance properties of images. Pseudo-labels [5] train a
287 discriminative classifier on a small set of labelled data and sample labels for a large set of unlabelled
288 data, which in turns is used to further train the classifier supervised. MixMatch [1] applies stochastic
289 transformations to an unlabeled image and each augmented image is fed to a classifier for which
290 the average logit distribution is sharpened by lowering the soft-max temperature. FixMatch [12]
291 strongly distorts an unlabeled image and trains the model such that the cross-entropy between the
292 one-hot pseudo-labels of the original image and the logits of the distorted image is minimized. Mean
293 teacher [13] employs a teacher model which parameters are the running mean of a student model
294 and trains the student such that a discrepancy between teacher and student predictions of augmented
295 unlabeled examples is minimized. Virtual Adversarial Training (VAT) [7] finds an adversarial
296 augmentation to an unlabeled example within an ϵ -ball with respect to some norm such that the
297 distance between the class distribution conditional on the unlabeled example and the one on the
298 adversarial example is maximized.

299 The methods of MixMatch, FixMatch and Mean teacher rely on pre-defined data augmentations,
300 which are readily available in the modality of images as the semantic meaning is invariant to
301 transforms such as rotation or flipping, but are difficult to construct in modalities such as language or
302 audio modalities. Our method is agnostic to the data modality. Pseudo-labeling is closely related
303 in that labels are sampled given unlabeled examples, whereas our method samples examples given
304 labels. VAT is close to our method as it is modality agnostic and leverages the learned model to
305 sample labeled examples, albeit of “adversarial” nature while our samples are “complementary.”
306 DAPPER is closest to our method as it employs a generative model to augment the data-set, but it
307 misses the coreset reduction.