| | MNIST | | YTF | | 3D-cars | | 3D-chairs | |
|---|---|---|---|---|---|---|---|---|
| | NMI ↑ | ENT ↓ | NMI ↑ | ENT ↓ | NMI ↑ | ENT ↓ | NMI ↑ | ENT ↓ |
| JointVAE ($\beta = 10$) | $0.536 \pm 0.13$ | $1.032 \pm 0.29$ | $0.372 \pm 0.06$ | $1.751 \pm 0.03$ | $\mathbf{0.452 \pm 0.24}$ | $1.026 \pm 0.43$ | $0.392 \pm 0.28$ | $2.053 \pm 0.46$ |
| JointVAE ($\beta = 20$) | $\mathbf{0.704 \pm 0.08}$ | $\mathbf{0.661 \pm 0.13}$ | $0.421 \pm 0.04$ | $1.687 \pm 0.02$ | $0.441 \pm 0.31$ | $\mathbf{1.022 \pm 0.34}$ | $0.431 \pm 0.26$ | $\mathbf{1.817 \pm 0.42}$ |
| JointVAE ($\beta = 30$) | $0.680 \pm 0.07$ | $0.701 \pm 0.13$ | $0.447 \pm 0.03$ | $1.717 \pm 0.02$ | $0.391 \pm 0.23$ | $1.082 \pm 0.53$ | $0.448 \pm 0.21$ | $1.909 \pm 0.44$ |
| JointVAE ($\beta = 40$) | $0.676 \pm 0.09$ | $0.713 \pm 0.14$ | $\mathbf{0.479 \pm 0.02}$ | $\mathbf{1.662 \pm 0.03}$ | $0.324 \pm 0.45$ | $1.151 \pm 0.51$ | $\mathbf{0.480 \pm 0.24}$ | $1.986 \pm 0.41$ |
| JointVAE ($\beta = 50$) | $0.649 \pm 0.09$ | $0.774 \pm 0.15$ | $0.435 \pm 0.03$ | $1.695 \pm 0.02$ | $0.376 \pm 0.40$ | $1.193 \pm 0.25$ | $0.377 \pm 0.31$ | $2.120 \pm 0.30$ |
| InfoGAN + $L_{ntxent}$ | $0.838 \pm 0.05$ | $0.351 \pm 0.09$ | $0.712 \pm 0.01$ | $0.831 \pm 0.02$ | $0.617 \pm 0.28$ | $0.835 \pm 0.52$ | $0.438 \pm 0.15$ | $1.237 \pm 0.32$ |
| Ours | $\mathbf{0.889 \pm 0.04}$ | $\mathbf{0.213 \pm 0.09}$ | $\mathbf{0.792 \pm 0.01}$ | $\mathbf{0.636 \pm 0.01}$ | $\mathbf{0.850 \pm 0.07}$ | $\mathbf{0.303 \pm 0.15}$ | $\mathbf{0.650 \pm 0.08}$ | $\mathbf{0.765 \pm 0.18}$ |

Thank you for the helpful comments. We are encouraged that the reviewers found the problem setting important & unexplored (R2,3,4), and our solution effective & reasonable (R1,2,4) in overcoming issues of existing work (R1,3).

**[R1] Clarifications regarding the VAE-based baseline** We apologize for not including the hyperparameter details of JointVAE. We use the KL term for both continuous as well as discrete variables, to follow the standard normal ($\mathcal{N}(0,1)$) and uniform categorical distribution ($\text{Cat}(p = 1/k)$), respectively. We use uniform categorical because of the unsupervised nature of the problem (L53-5, 83-4), and our full approach itself starts from uniform initialization (L127-30 supp). We use the same weight ($\beta$) for both KL loss terms (similar to JointVAE paper), the value of which was first decided empirically based on image reconstruction quality - we observed that a value in 100s (e.g. 100-300) resulted in poor reconstruction quality. We ultimately went with $\beta = 30$ (the value chosen by JointVAE for MNIST) for all datasets. We present an ablation study on the effect of strength of the KL term ($\beta$) on disentanglement in the table above. While we agree that a lesser weight on the KL term might imply lesser restriction for inference model to follow the uniform prior, it might result in reduced disentanglement as well. So, starting from a particular value (say $\beta = 30$), it is not clear whether increasing (towards uniform) or decreasing (towards less disentanglement) it would help from disentanglement's point of view. We observe this in the ablation study as well, where low and high values of the weight ($\beta = 10$ and $\beta = 50$ respectively) usually result in low disentanglement scores. We can, however, get slightly better results with alternate $\beta$s (different $\beta$ for different datasets) than the ones reported in the main paper, and we'll update them with these in the final version. Note that for our approach we don't perform an exhaustive search for $\lambda_2$ (in $L_{final}$); it's set as 10 for all datasets (L124-5 in supp). Finally, comparison to FactorVAE is not directly applicable, as it can only capture continuous factors, and the paper itself mentions the inability to capture discrete factors as a limitation.

**[R2] [R3] Concerns regarding technical novelty** We agree we leverage existing techniques (L166-8). However, they have previously been used in orthogonal areas: Gumbel-softmax was introduced for differentiable sampling of one-hot like variables, and identity preserving transformations have been used as part of data augmentation, avoiding overfitting, representation learning, etc. In this work, we've integrated these techniques in a coherent framework to address an important problem in a *novel* setting of learning disentangled representations in class-imbalanced data (L169-70).

**[R3] [R4] Concerns regarding datasets** We'd like to point out that seminal works in learning disentangled representations (e.g. InfoGAN, $\beta$-VAE) present results on such small datasets, where it is relatively easy to account for, and capture the factors of variations. R3 states "..model tends to mode collapse on YTF" - we respectfully disagree: the categories in YTF consist of video frames of the same person, resulting in very similar *real images* themselves. Faithfully modeling such image distribution hence results in similar generated images.


Uniform-InfoGAN


Ours

R3 further suggested to try the ShapeNet dataset, which is naturally imbalanced. The original dataset was too big to operate during the rebuttal phase ($\sim$ 600k images, with 10 renderings per model), so we created a subset consisting of 5 categories - cars, airplanes, bowl, can, rifle - in a way which maintains the original imbalance between categories. Due to time constraints, we're only able to compare Uniform-InfoGAN and our final method. We can see that InfoGAN (NMI: 0.545, ENT: 0.687) mixes up different categories more frequently (rows 1/3) than our method (NMI: 0.781, ENT: 0.432), which is more consistent when grouping same category instances together. We'll include more analysis.

**[R2] Clarifications regarding training components** The first component, Gumbel-Softmax, should be thought of as making the latent distribution flexible, so that the model can capture *any* discrete factor having $k$ modes (not necessarily object identity) present in any ratio. $L_{ntxent}$'s role is to push the discovered factor to better correspond with object identity (L162-3), in balanced or imbalanced case (as R2 points out). We include the results for Uniform-InfoGAN + $L_{ntxent}$ in the table above. The results improve compared to Uniform-InfoGAN, but having the rigid uniform prior still results in worse performance compared to our approach. Hence, Gumbel-softmax alone *shouldn't* be thought of as an improvement, as it works best when used along with $L_{ntxent}$ (L229-33). Furthermore, we cannot use class labels since this is an unsupervised task. Finally, we'll discuss the mentioned related works; we thank R2 for pointing them out.

**[R4] Clarification regarding input for $L_{ntxent}$** We agree and will discuss this in detail. *ntxent* is described in L159.