

# INNF+ Workshop 2020

Phillip Urquijo, *Justin Tan*

Decorrelation via Disentanglement

October 26, 2020



# Background Sculpting

Classifier output  $f(X; \theta_f) \sim p(\text{signal}|\text{data})$ . Only accept events above a given probability.

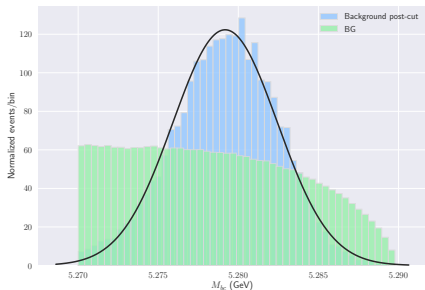


Figure 1: Continuum  $M_{bc}$  before (green) and after (blue) 0.995 suppression

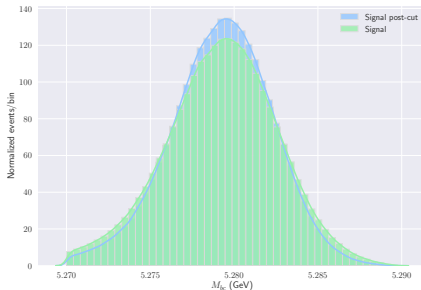


Figure 2: Signal  $M_{bc}$  before (green) and after (blue) 0.995 suppression

e.g. invariant mass. Background looks like signal post-selection - difficult to quantify background contribution via mass sidebands.

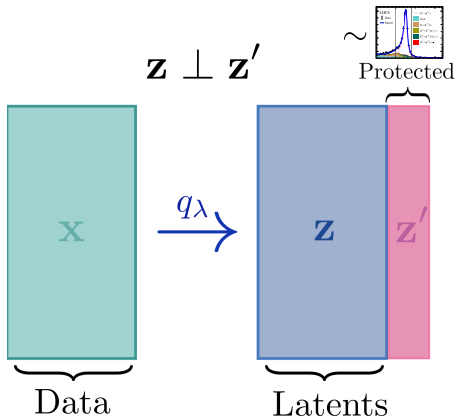
- Most decorrelation approaches:

$$\mathcal{L}_{Decorr.} = \mathcal{L}_{Clf.} + \lambda \cdot \mathcal{L}_? \quad (1)$$

- $\mathcal{L}_?$ : summarizes dependency of classifier output on protected variables.
- We propose to **decouple** the classification task from the decorrelation process.
- Two-stage procedure:
  - ▶ First, learn a representation  $z$  non-informative of protected variables.
  - ▶ Use representation in lieu of original data  $x$  downstream.

# Disentanglement

- Want to find representation of event data that is:
  - ▶ Non-informative of physical observables (e.g.  $M_{inv}$ ).
  - ▶ Maximally informative of the original data.
- Propose to **isolate influence of protected variables** in a subspace of latent representation.





# Latent Variable Models

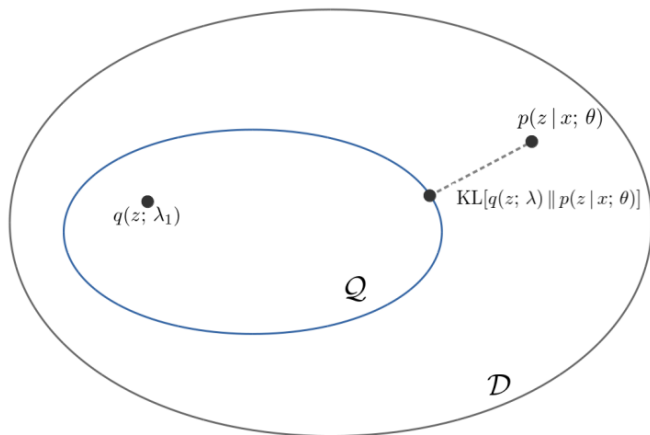
- Want compact description of observed data  $x$ .
- Introduce hidden/latent variables  $z$ , joint density  $p_\theta(x, z) = p_\theta(x|z)p(z)$ .
- $z$  provides low-dimensional representation of  $x$ .
- Want to **infer** latents  $z$  responsible for generating  $x$ :

$$p_\theta(z|x) = \frac{p_\theta(x, z)}{\int_z dz p_\theta(x|z)p(z)} \quad (2)$$

- **Problem:** quadrature takes exponential time in  $\dim(z)$ , likely non-analytic.

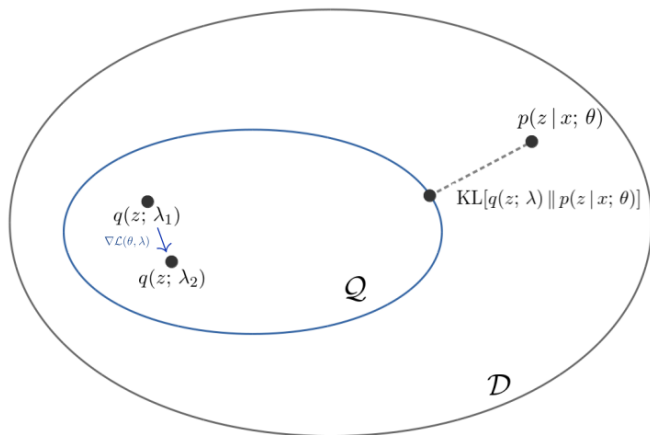
# Inference as Optimization

- Approximate true  $p_\theta(z|x)$  with  $q_\lambda(z; x) \in \mathcal{Q}$ .
  - ▶  $\mathcal{Q}$ : 'Simple' family of parametric distributions.
- Minimize 'distance'  $D_{\text{KL}}(q_\lambda(z; x) \parallel p_\theta(z|x))$ .
- **Inference (hard)**  $\approx$  **optimization (easy)**.



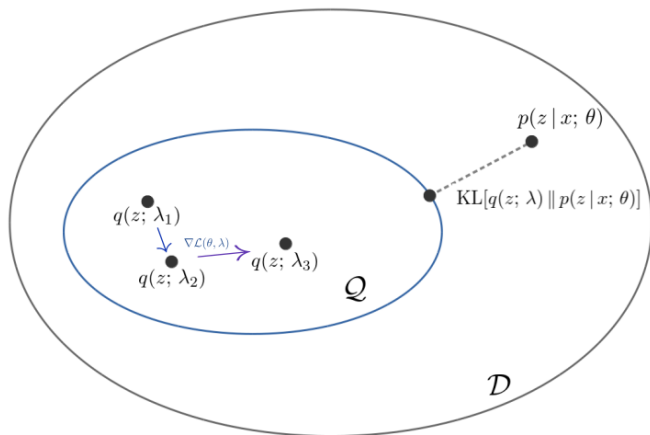
# Inference as Optimization

- Approximate true  $p_\theta(z|x)$  with  $q_\lambda(z; x) \in \mathcal{Q}$ .
  - ▶  $\mathcal{Q}$ : 'Simple' family of parametric distributions.
- Minimize 'distance'  $D_{\text{KL}}(q_\lambda(z; x) \parallel p_\theta(z|x))$ .
- **Inference (hard)**  $\approx$  **optimization (easy)**.



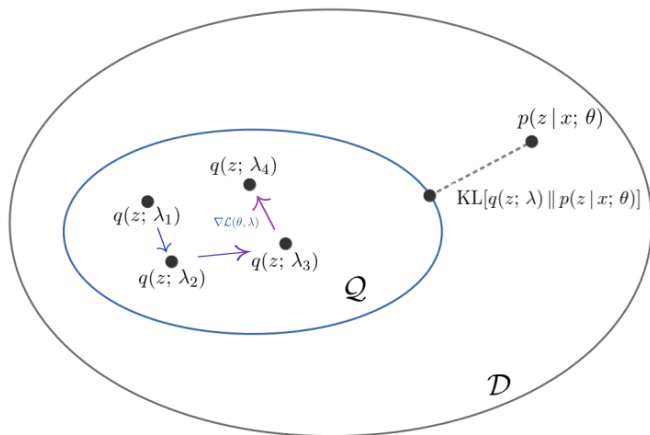
# Inference as Optimization

- Approximate true  $p_\theta(z|x)$  with  $q_\lambda(z; x) \in \mathcal{Q}$ .
  - ▶  $\mathcal{Q}$ : 'Simple' family of parametric distributions.
- Minimize 'distance'  $D_{\text{KL}}(q_\lambda(z; x) \parallel p_\theta(z|x))$ .
- **Inference (hard)**  $\approx$  **optimization (easy)**.



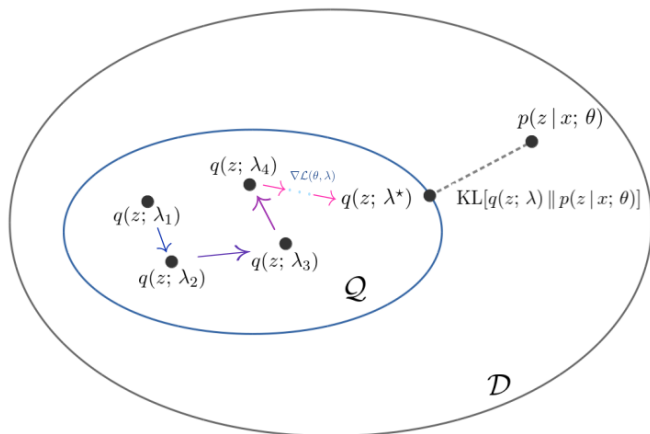
# Inference as Optimization

- Approximate true  $p_\theta(z|x)$  with  $q_\lambda(z; x) \in \mathcal{Q}$ .
  - ▶  $\mathcal{Q}$ : 'Simple' family of parametric distributions.
- Minimize 'distance'  $D_{\text{KL}}(q_\lambda(z; x) \parallel p_\theta(z|x))$ .
- **Inference (hard)**  $\approx$  **optimization (easy)**.



# Inference as Optimization

- Approximate true  $p_\theta(z|x)$  with  $q_\lambda(z;x) \in \mathcal{Q}$ .
  - ▶  $\mathcal{Q}$ : 'Simple' family of parametric distributions.
- Minimize 'distance'  $D_{\text{KL}}(q_\lambda(z;x) \parallel p_\theta(z|x))$ .
- **Inference (hard)**  $\approx$  **optimization (easy)**.



# Inference as Optimization

- Equivalent to maximization of Evidence Lower BOund (ELBO):

$$\text{ELBO}(\mathbf{x}) \triangleq \mathbf{E}_{q_{\lambda}(\mathbf{z}; \mathbf{x})} \left[ \log \frac{p_{\theta}(\mathbf{x}|\mathbf{z})p(\mathbf{z})}{q_{\lambda}(\mathbf{z}; \mathbf{x})} \right] \leq \log p_{\theta}(\mathbf{x}) \quad (3)$$

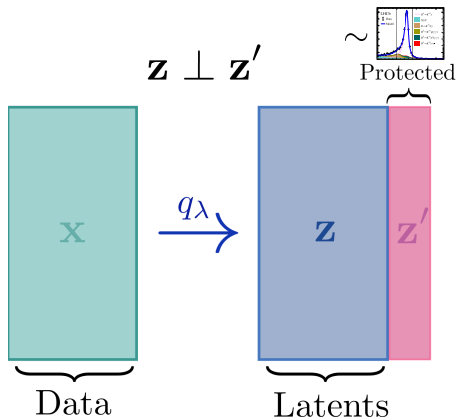
$$\min_{q_{\lambda} \in \mathcal{Q}} D_{\text{KL}}(q_{\lambda}(\mathbf{z}; \mathbf{x}) \| p_{\theta}(\mathbf{z}|\mathbf{x})) \iff \max_{\theta, \lambda} \text{ELBO}(\mathbf{x}) \quad (4)$$

- Variational Autoencoding (VAEs) framework:  $q_{\lambda}$  and conditional model  $p_{\theta}$  parameterized by neural networks.
  - ▶ Instead of generation, interested in inference.

# Disentanglement

- Different components of  $z$  should capture distinct semantic information about  $x$ .
- How to encourage independence between latent dimensions?
- Try to factorize **average** joint variational posterior over  $z, z'$ :

$$\mathbf{E}_{p(x)} [q_{\lambda}(z, z'; x)] \approx \mathbf{E}_{p(x)} [q_{\lambda}(z; x)] \mathbf{E}_{p(x)} [q_{\lambda}(z'; x)]. \quad (5)$$





$$\mathbf{E}_{p(x)} [\text{ELBO}(x)] = \mathbf{E}_{p(x)} [\mathbf{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - D_{\text{KL}}(q_{\lambda}(z;x) \parallel p(z))]$$

$$\begin{aligned} \mathbf{E}_{p(x)} [D_{\text{KL}}(q_{\lambda}(z;x) \parallel p(z))] &= \mathbb{I}_q(x; z) + D_{\text{KL}}(q(z) \parallel p(z)) \\ &= \underbrace{\mathbb{I}_q(x; z)}_{\textcircled{1}} + \underbrace{D_{\text{KL}}\left(q(z) \parallel \prod_j q(z_j)\right)}_{\textcircled{2}} \\ &\quad + \underbrace{\sum_i D_{\text{KL}}(q(z_i) \parallel p(z_i))}_{\textcircled{3}} \end{aligned}$$

- ① Observed-latent mutual information.

$$\mathbf{E}_{p(x)} [\text{ELBO}(x)] = \mathbf{E}_{p(x)} [\mathbf{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - D_{\text{KL}}(q_{\lambda}(z;x) \parallel p(z))]$$

$$\begin{aligned} \mathbf{E}_{p(x)} [D_{\text{KL}}(q_{\lambda}(z;x) \parallel p(z))] &= \mathbb{I}_q(x; z) + D_{\text{KL}}(q(z) \parallel p(z)) \\ &= \underbrace{\mathbb{I}_q(x; z)}_{\textcircled{1}} + \underbrace{D_{\text{KL}}\left(q(z) \parallel \prod_j q(z_j)\right)}_{\textcircled{2}} \\ &\quad + \underbrace{\sum_i D_{\text{KL}}(q(z_i) \parallel p(z_i))}_{\textcircled{3}} \end{aligned}$$

- ① Observed-latent mutual information.
- ② Latent Total Correlation.

$$\mathbf{E}_{p(x)} [\text{ELBO}(x)] = \mathbf{E}_{p(x)} [\mathbf{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - D_{\text{KL}}(q_{\lambda}(z;x) \parallel p(z))]$$

$$\begin{aligned} \mathbf{E}_{p(x)} [D_{\text{KL}}(q_{\lambda}(z;x) \parallel p(z))] &= \mathbb{I}_q(x; z) + D_{\text{KL}}(q(z) \parallel p(z)) \\ &= \underbrace{\mathbb{I}_q(x; z)}_{\textcircled{1}} + \underbrace{D_{\text{KL}}\left(q(z) \parallel \prod_j q(z_j)\right)}_{\textcircled{2}} \\ &\quad + \underbrace{\sum_i D_{\text{KL}}(q(z_i) \parallel p(z_i))}_{\textcircled{3}} \end{aligned}$$

- ① Observed-latent mutual information.
- ② Latent Total Correlation.
- ③ Dimension-wise KL.

$$\min_{q_\lambda \in \mathcal{Q}} \left( \underbrace{\mathbb{I}_q(x; z)}_{\textcircled{1}} + \underbrace{D_{\text{KL}} \left( q(z) \parallel \prod_j q(z_j) \right)}_{\textcircled{2}} + \underbrace{\sum_i D_{\text{KL}} (q(z_i) \parallel p(z_i))}_{\textcircled{3}} \right)$$

**Aggregate Posterior:**  $q(z) = \frac{1}{N} \sum_n q_\lambda(z; x_n)$  empirical average of  $q_\lambda(z; x_n)$ .

① Observed-latent mutual information.

① Predictiveness of observed data from latents.

$$\min_{q_\lambda \in \mathcal{Q}} \left( \underbrace{\mathbb{I}_q(x; z)}_{\textcircled{1}} + \underbrace{D_{\text{KL}} \left( q(z) \parallel \prod_j q(z_j) \right)}_{\textcircled{2}} + \underbrace{\sum_i D_{\text{KL}}(q(z_i) \parallel p(z_i))}_{\textcircled{3}} \right)$$

**Aggregate Posterior:**  $q(z) = \frac{1}{N} \sum_n q_\lambda(z; x_n)$  empirical average of  $q_\lambda(z; x_n)$ .

- ① Observed-latent mutual information.
- ② Latent Total Correlation.

② Statistical dependency between dimensions of  $z$ .

$$\min_{q_\lambda \in \mathcal{Q}} \left( \underbrace{\mathbb{I}_q(x; z)}_{\textcircled{1}} + \underbrace{D_{\text{KL}} \left( q(z) \parallel \prod_j q(z_j) \right)}_{\textcircled{2}} + \underbrace{\sum_i D_{\text{KL}}(q(z_i) \parallel p(z_i))}_{\textcircled{3}} \right)$$

**Aggregate Posterior:**  $q(z) = \frac{1}{N} \sum_n q_\lambda(z; x_n)$  empirical average of  $q_\lambda(z; x_n)$ .

- ① Observed-latent mutual information.
- ② Latent Total Correlation.
- ③ Dimension-wise KL.

③ Complexity penalty on dimensions of  $q(z)$ .

# Disentangling Objective

$$\mathbf{E}_{p(x)} \left[ \underbrace{\mathbf{E}_{q_{\lambda}(z|x)} [\log p_{\theta}(x|z)] - D_{\text{KL}}(q_{\lambda}(z; x) \| p(z))}_{\text{ELBO}(x)} \right] - (\beta - 1) D_{\text{KL}} \left( q(z) \| \prod_{i=1}^d q(z_i) \right)$$

Maximize w.r.t. variational parameters  $(\lambda, \theta)$  for  $\beta > 1$ .

- Total correlation hard to estimate, but multiple (biased) estimators developed:
  - ▶ Chen, et. al. Isolating Sources of Disentanglement in VAEs, 2018. ( $\beta$ -TCVAE).
  - ▶ Kim, et. al. Disentangling by Factorizing, 2018. (Factor-VAE).

# Identifying latent dimensions

Identify selected dimensions  $z'$  of the latent code with the protected variables  $v$  with cross-entropy-like objective:

$$R_s(x, v) = \sum_{i=1}^{|v|} (v_i \log \sigma(r(x)_i) + (1 - v_i) \log (1 - \sigma(r(x)_i))). \quad (6)$$

- $r(x)$  is mean of the variational posterior  $q_\lambda : r(x) = \langle q_\lambda(z; x) \rangle$ .
- We protect the invariant mass:  $v = M_{inv}$ .



$$\mathbf{E}_{p(x)} [\text{ELBO}(x) + \lambda_s \mathbf{E}_{p(v|x)} [R_s(x, v)]] - (\beta - 1) D_{\text{KL}} \left( q(z) \parallel \prod_{i=1}^d q(z_i) \right) \quad (7)$$

- Maximize w.r.t. variational parameters  $\theta, \lambda$ .
- Adopt mean of variational posterior as representation:

$$r(x) = \int_z dz z q_\lambda(z; x)$$

- Excise dimensions  $z'$  to yield final representation for use in arbitrary downstream tasks without exploiting protected vs:

$$r(x)' = \langle q_\lambda(z \setminus z'; x) \rangle \quad (8)$$

# Only disentangle protected factors

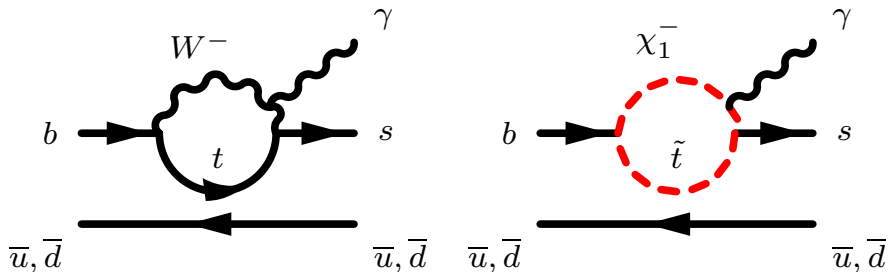
Only disentangle protected dimensions  $z'$  (with indices  $S$ ) instead of full  $z$ :

$$\begin{aligned} \mathcal{L}(\theta, \phi) = & \mathbf{E}_{p(x)} \left[ \mathbf{E}_{q_{\lambda}(z;x)} [\log p_{\theta}(x|z)] + \lambda_s \mathbf{E}_{p(v|x)} [R_s(q_{\lambda}(z;x), v)] \right] \\ & - \left( \alpha \mathbb{I}_q(x; z) + \beta \mathbf{E}_{q(z)} \left[ \log \frac{q(z)}{q(z \setminus z') \prod_{j \in S} q(z'_j)} \right] \right. \\ & \left. + \gamma D_{\text{KL}}(q(z \setminus z') \parallel p(z \setminus z')) + \delta \sum_{j \in S} D_{\text{KL}}(q(z'_j) \parallel p(z'_j)) \right) \end{aligned}$$

- Each term estimated through stratified importance sampling.
- We call this approach  $\beta$ -TCVAE-Sensitive.

# Experiments

- Apply to  $B$ -physics analysis of 'penguin'  $b \rightarrow s\gamma$  transition ( $\text{BF} \approx 10^{-5}$ ).
- Background estimation relies on beam-constrained mass spectrum:  
$$M_{bc}^2 = E_{beam}^2 - p_B^2.$$
  - ▶ Set protected  $v = M_{bc}$ .
- Use output representation  $\langle q_\lambda(z \setminus z'; x) \rangle$  as input to 'downstream' classification task
  - ▶  $10^6$  events, 1 : 10 signal ratio,  $M_{bc}$ -correlated features removed.



- Quantify results:
  - ▶ **Disentanglement:** Mutual Information Gap.

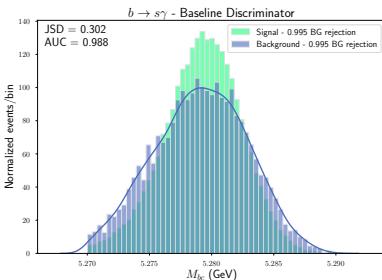
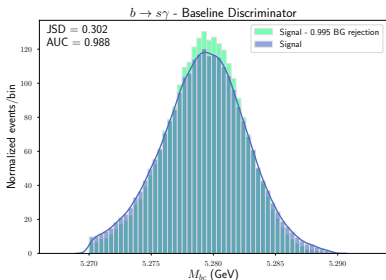
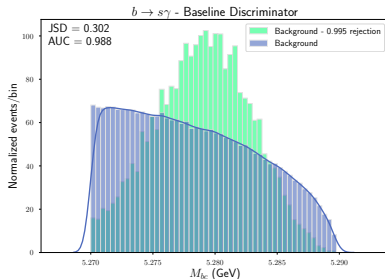
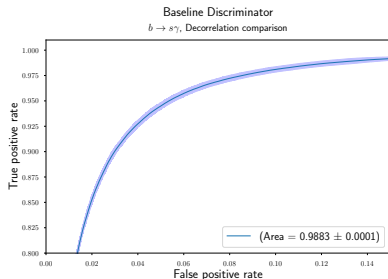
$$\text{MIG}(z, v) = \frac{1}{K} \sum_{k=1}^K \frac{1}{H(v_k)} \left( I(z_{j^{(k)}}; v_k) - \max_{j \neq j^{(k)}} I(z_j; v_k) \right)$$
$$j^{(k)} = \underset{j}{\operatorname{argmax}} I(z_j; v_k)$$

- ▶ **Sculpting extent:** (Discretized) Jensen-Shannon divergence of background mass spectra before/after 0.995 background rejection.

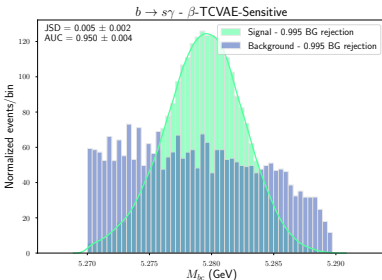
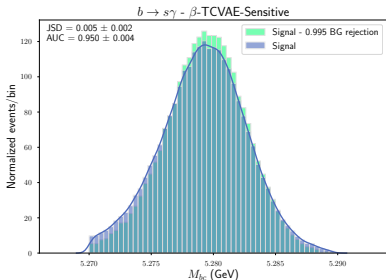
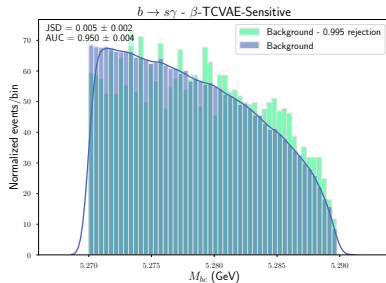
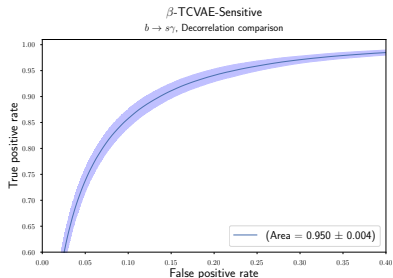
$$\text{JSD} = D_{\text{JS}} \left( N_{\text{bkg}}^{\text{pass}}(m) / \sum_i N_{\text{bkg},i}^{\text{pass},i}(m) \left\| \left\| N_{\text{bkg}}^{\text{fail}}(m) / \sum_i N_{\text{bkg},i}^{\text{fail},i}(m) \right. \right. \right).$$

- ▶ **Classification:** AUC-ROC.

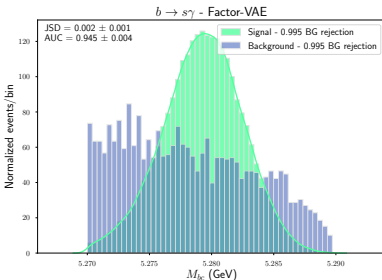
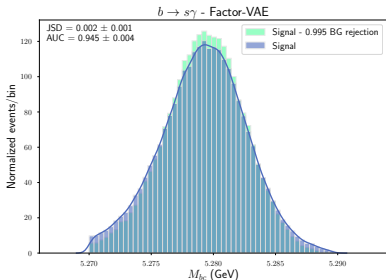
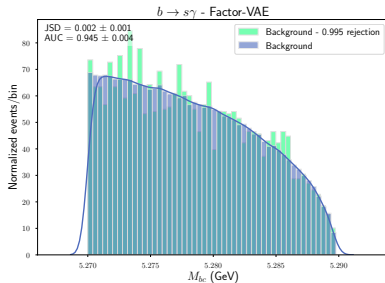
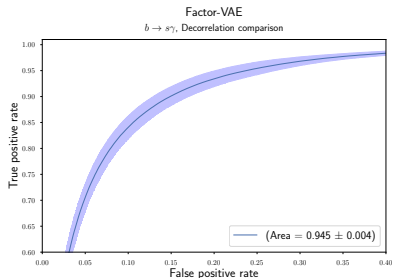
# Baseline Discriminator



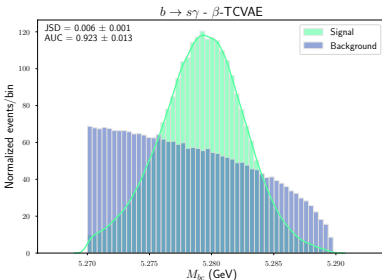
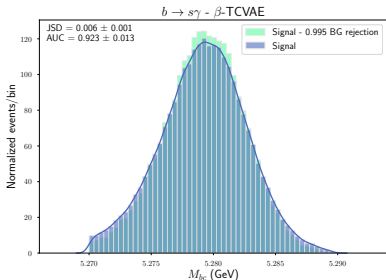
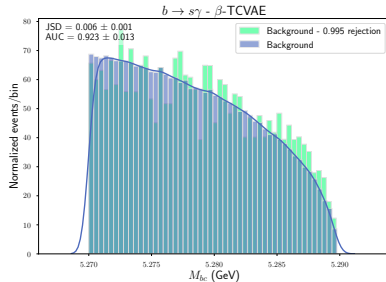
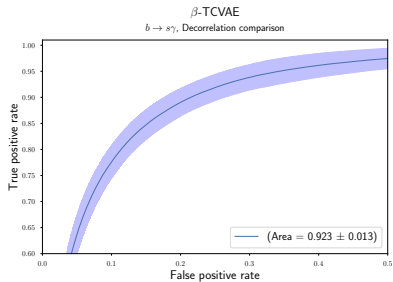
# Decorrelated Representation



# Decorrelated Representation (Factor-VAE)

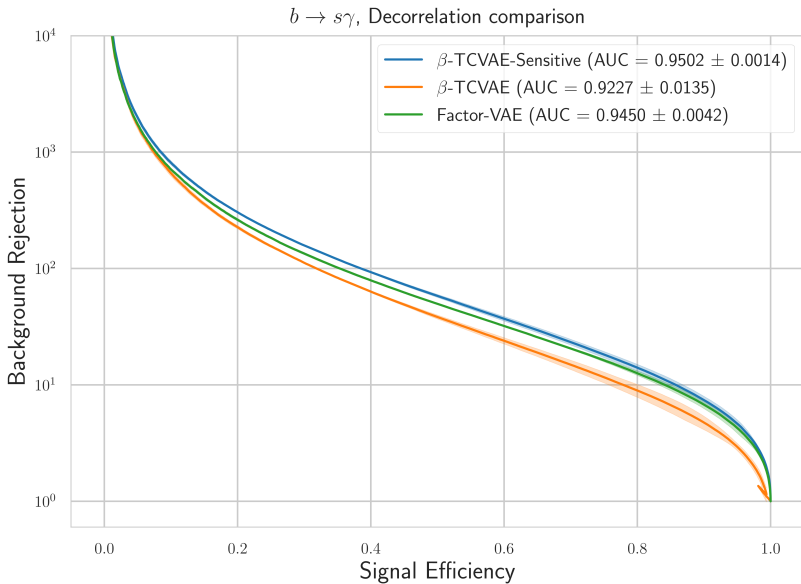


# Decorrelated Representation ( $\beta$ -TCVAE)





# Efficiency Plots

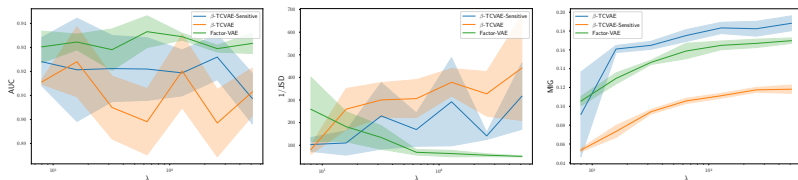


Model	AUC	JSD	MIG
$\beta$ -TCVAE-Sensitive	<b>0.950 <math>\pm</math> 0.002</b>	0.0046 $\pm$ 0.0016	0.071 $\pm$ 0.009
$\beta$ -TCVAE	0.923 $\pm$ 0.014	0.0059 $\pm$ 0.0011	0.16 $\pm$ 0.03
Factor-VAE	0.945 $\pm$ 0.004	<b>0.0021 <math>\pm</math> 0.0007</b>	<b>0.152 <math>\pm</math> 0.005</b>
Annealed-VAE	0.915 $\pm$ 0.006	0.04 $\pm$ 0.02	0.062 $\pm$ 0.003
$\beta$ -VAE	0.903 $\pm$ 0.021	0.03 $\pm$ 0.02	0.064 $\pm$ 0.017
Baseline	0.988 $\pm$ 0.000	0.324 $\pm$ 0.0074	-

**Table 1:** Metrics obtained when the invariant representation  $z \setminus t$  is used as input for the downstream discriminator. Higher is better for AUC and MIG, while lower is better for JSD.

# Hyperparameter Robustness

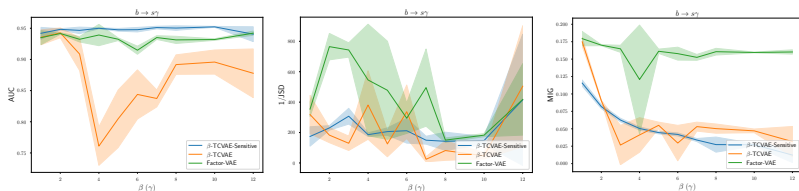
Scan over the weighting  $\lambda_s$  of the supervised regularization term  $R_s$  (Eq. 6).



**Figure 8:** AUC (left), 1/JSD (middle), and MIG (right) versus  $\beta$  ( $\gamma$ ) for  $\beta$ -TCVAE (orange),  $\beta$ -TCVAE-Sensitive (blue) and Factor-VAE (green). Higher is better for each metric.

# Hyperparameter Robustness

Scan over the weighting  $\beta$  of the latent total correlation term  $\mathbb{E}_{q(z)} [\log (q(z) / \prod_i q(z_i))]$  (Eq. 7).



**Figure 9:** AUC (left), 1/JSD (middle), and MIG (right) versus  $\beta$  ( $\gamma$ ) for  $\beta$ -TCVAE (orange),  $\beta$ -TCVAE-Sensitive (blue) and Factor-VAE (green). Higher is better for each metric.

- Generative modelling via *intervention* on disentangled factors.
  - ▶ Manipulate generated event characteristics in latent space.
- Anomaly detection via isolation of invariant mass dimension.
  - ▶ Associate latent dimension with invariant mass and perform approximate posterior inference.
- Discrete-variable decorrelation via discrete latent relaxations.

- Presented a method for decorrelation based on latent variable representations in VAEs.
- Robust, non-adversarial training procedure.
- Decorrelation process decoupled from and agnostic to the downstream task.

Thanks for listening!

[justin.tan@unimelb.edu.au](mailto:justin.tan@unimelb.edu.au)

[github.com/Justin-Tan](https://github.com/Justin-Tan)

Backup



Proposition (*Factorization of  $q(z) \implies$  factorization of  $r(x) = \int dz z q_\lambda(z; x)$ )*)

Proof.

$q(z) \triangleq \mathbf{E}_{p(x)} [q_\lambda(z; x)]$ . Take expectation of  $r(x)$  over  $p(x)$ :

$$\int dx r(x) p(x) = \int dx dz z q_\lambda(z; x) p(x) \quad (9)$$

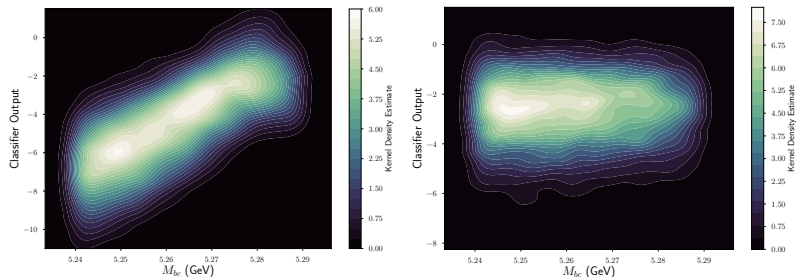
$$= \int dz z \int dx q_\lambda(z; x) p(x) \quad (10)$$

$$= \int dz z q(z) \quad (11)$$

$$= \int dz z \prod_{i=1}^{\dim(\mathcal{Z})} q(z_i) \quad (12)$$

$$= \prod_{i=1}^{\dim(\mathcal{Z})} \int dz_i z_i q(z_i). \quad (13)$$





**Figure 10:** Visualization of invariance of the downstream classification task to the sensitive variable  $s = M_{bc}$ . Bivariate kernel density estimate between  $M_{bc}$  and the classifier output,  $\log\left(\frac{p}{1-p}\right)$ , where  $p$  is the probability that a given event belongs to the positive signal class given by the downstream classifier. The output is a monotonic function of  $p$  - larger values of the classifier output indicate higher confidence by the classifier that an example belongs to the positive class. The baseline classifier output (left) is strongly correlated with  $M_{bc}$ , and that it outputs higher confidence that an event belongs to the signal class around the signal region  $M_{bc} \in (5.27, 5.29)$  GeV, where the signal  $M_{bc}$  distribution is sharply peaked. The output of the invariant classifier (right) does not exhibit significant dependency on  $M_{bc}$ , with the classifier response flat throughout the full  $M_{bc}$  range.