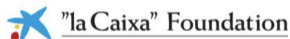


# Anomaly Awareness Estimation in $b \rightarrow s\ell^+\ell^-$ with VAEs

B. Capdevila University of Barcelona, ICCUB



Institut de Ciències del Cosmos  
UNIVERSITAT DE BARCELONA



*In collaboration with:*

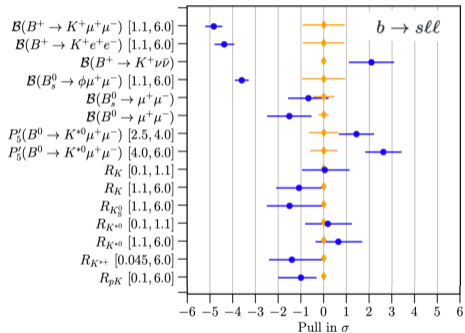
**A. Scaffidi; arXiv:241x.xxxx**

# Anomaly Detection in Low-Energy Observables

- ▶ **Anomalies:** Deviations between **experimental data** and **theoretical predictions** under a **null hypothesis** (SM hypothesis)
- ▶ Searching for anomalies in **low-energy observables** is critical for **NP exploration**
  - ⇒ Observables at **low energies** receive contributions from **higher scales** due to **quantum corrections**
- ▶ **Goal:** *Properly* quantify the **statistical significance** of observed anomalies
- ▶ Examples of recent anomalies
  - ⇒  $B$  anomalies ( $b \rightarrow sl^+\ell^-$ ,  $b \rightarrow cl\nu$ )
  - ⇒  $(g - 2)_\mu$
  - ⇒  $V_{cb}$ ,  $V_{ub}$  puzzle
  - ⇒ ...

# Anomaly Detection in $b \rightarrow sl^+l^-$

## Data



## Stats



$p$  - value

# Experimental Data in Phenomenological Analyses

► Experimental data:

⇒ Released from the experiments as a **vector of means**  $\mu^{\text{exp}}$  and a **covariance matrix**  $\Lambda^{\text{exp}}$

⇒ **Implicitly assumes a Gaussian distribution** for the experimental measurements

$$p(\mathbf{x}^{\text{exp}}) = \mathcal{N}(\mathbf{x}^{\text{exp}}; \boldsymbol{\mu}^{\text{exp}}, \boldsymbol{\Lambda}^{\text{exp}})$$

$4.0 < q^2 < 6.0 \text{ GeV}^2/c^4$			$F_L$	$P_1$	$P_2$	$P_3$	$P'_4$	$P'_5$	$P'_6$	$P'_8$
$P_1$	$0.088 \pm 0.235 \pm 0.029$	$F_L$	1.00	0.04	0.05	-0.10	-0.04	-0.14	-0.17	0.14
$P_2$	$0.105 \pm 0.068 \pm 0.009$	$P_1$		1.00	0.06	0.07	-0.06	-0.10	-0.03	0.02
$P_3$	$-0.090 \pm 0.139 \pm 0.006$	$P_2$			1.00	-0.02	-0.14	-0.09	-0.03	-0.01
$P'_4$	$-0.312 \pm 0.115 \pm 0.013$	$P_3$				1.00	-0.01	0.07	0.19	-0.01
$P'_5$	$-0.439 \pm 0.111 \pm 0.036$	$P'_4$					1.00	0.02	0.04	0.01
$P'_6$	$-0.293 \pm 0.117 \pm 0.004$	$P'_5$						1.00	0.09	0.00
$P'_8$	$0.166 \pm 0.127 \pm 0.004$	$P'_6$							1.00	0.02
		$P'_8$								1.00

LHCb, arXiv:2003.04831

# Theoretical Predictions in Phenomenological Analyses

## ► Theoretical predictions

- ⇒ For the observables in the analysis  $\mathbf{x} = (x_1, \dots, x_n)$ , we have functions representing their **theoretical predictions** in terms of several **input parameters**  $\boldsymbol{\nu} = (\nu_1, \dots, \nu_m)$

$$x_i = x_i(\nu_1, \dots, \nu_m)$$

- ⇒ The **input parameters** are distributed according to some distribution, usually Gaussian

$$\boldsymbol{\nu} \sim \mathcal{N}(\boldsymbol{\nu}; \boldsymbol{\mu}_\nu, \boldsymbol{\Lambda}_\nu)$$

where  $\boldsymbol{\mu}_\nu$  and  $\boldsymbol{\Lambda}_\nu$  means and covariance of the distribution of underlying parameters

## ► Implications

- ⇒ Even if the distribution of parameters is Gaussian, observables with complex structures **do not distribute normally**
- ⇒ Except for observables with a **linear dependence** on the underlying parameters
- ⇒ This means the **likelihood**  $p(\mathbf{x}|H_i)$  will generally be distributed under a **non-Gaussian distribution**

## Gaussian Likelihoods in Phenomenological analyses

- ▶ In **conventional frequentist  $b \rightarrow sl^+\ell^-$  analyses**, Gaussian likelihoods are assumed

$$p(\mathbf{x}|H_i) = \mathcal{N}(\mathbf{x}; \mathbf{x}_{H_i}, \mathbf{\Lambda}_{H_i}) = \frac{1}{\sqrt{(2\pi)^n |\mathbf{\Lambda}_{H_i}|}} \exp\left(-\frac{1}{2}(\mathbf{x} - \mathbf{x}_{H_i})^T \mathbf{\Lambda}_{H_i}^{-1}(\mathbf{x} - \mathbf{x}_{H_i})\right)$$

- ⇒  $\mathbf{\Lambda}_{H_i}$  is the **sum of theoretical and experimental covariances**

$$\mathbf{\Lambda}_{H_i} = \mathbf{\Lambda}_{H_i}^{\text{th}} + \mathbf{\Lambda}^{\text{exp}}$$

- ⇒ Theory side: Covariance estimated from samples of theoretical predictions
- ⇒ Experimental side: Covariance read directly from data released by experiments

- ▶ Measuring goodness-of-fit

- ⇒ Use  $-\log p(\mathbf{x}|H_i)$  as a **statistic** to measure **agreement between data and hypothesis  $H_i$**

$$-\log p(\mathbf{x}|H_i) = \frac{\chi^2}{2} + \text{const}$$

where  $\chi^2$  is the chi-squared function

- ⇒ If theoretical predictions are normally distributed,  $-\log p(\mathbf{x}|H_i)$  follows a  $\chi^2$ -distribution with  $n_{\text{dof}} = n_{\text{obs}}$  in the analysis
- ⇒ If **not normally distributed**,  $-\log p(\mathbf{x}|H_i)$  is **only asymptotically  $\chi^2$**  due to the central limit theorem
- ⇒ Assuming  $\chi^2$ -distribution when it is not **creates biases** in calculating  $p$ -values from  $-\log p(\mathbf{x}_{\text{exp}}|H_i)$

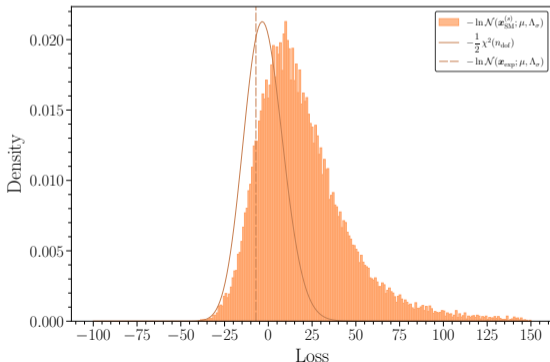
## Distribution of $-\log p(\mathbf{x}|H_i)$ vs Asymptotic $\chi^2$

### ► Distribution of $-\log p(\mathbf{x}|H_i)$

- ⇒ Calculated for the set of observables in the  $b \rightarrow s\ell\ell$  dataset under the SM hypothesis ( $H_0 = \text{SM}$ )
- ⇒ Obtained by calculating  $-\log p(\mathbf{x}|H_0)$  for each  $\mathbf{x}$  in a sample  $\mathbf{x}^s = (\mathbf{x}^1, \dots, \mathbf{x}^{n_{\text{sample}}})$
- ⇒ Sample size:  $n_{\text{sample}} = 10000$
- ⇒ Each  $\mathbf{x}^s$  generated by sampling underlying parameters  $\boldsymbol{\nu}$  from  $\mathcal{N}(\boldsymbol{\nu}; \boldsymbol{\mu}_\nu, \boldsymbol{\Lambda}_\nu)$

### ► Comparison

- ⇒ Distribution of  $-\log p(\mathbf{x}|H_i)$  vs asymptotic  $\chi^2$  distribution with corresponding degrees of freedom
- ⇒ Difference observed between the two distributions



# Estimating Likelihoods

## ► Understanding $p(\mathbf{x}|H_i)$

- ⇒ Goal: Determine the likelihood  $p(\mathbf{x}|H_i) = p(\mathbf{x})$  for enhanced statistical rigor in hypothesis testing
- ⇒ Known: Distribution of underlying inputs  $p(\boldsymbol{\nu})$ , which informs us about the prior probability of the model parameters
- ⇒ Computable:  $p(\mathbf{x}|\boldsymbol{\nu})$ , achievable by simulating observables  $\mathbf{x}$  using sampled parameters  $\boldsymbol{\nu}$  from their known distributions

## ► Obtaining $p(\mathbf{x})$ as a marginal likelihood

$$p(\mathbf{x}) = \int d\boldsymbol{\nu} p(\mathbf{x}|\boldsymbol{\nu})p(\boldsymbol{\nu})$$

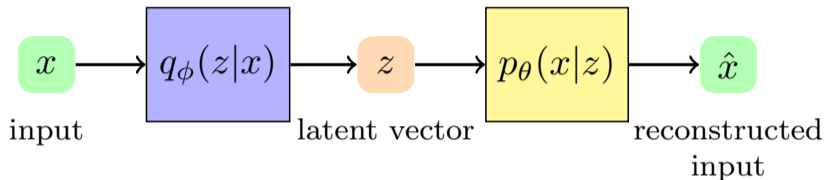
- ⇒ In most real-life applications  $\boldsymbol{\nu}$  is usually high-dimensional
- ⇒ Challenge: a direct computation the  $d\boldsymbol{\nu}$  integral is generally computationally prohibitive
- ⇒ Hence, the likelihood  $p(\mathbf{x})$  is typically **intractable**

## ► Estimation using Variational Autoencoders (VAEs)

- ⇒ VAEs provide a feasible approach to approximate  $p(\mathbf{x})$  with **arbitrary precision**



# Introducing Variational Autoencoders



## ► VAE framework

- ⇒ **Model:** Pairs a probabilistic decoder  $p_{\theta}(\mathbf{x}|\mathbf{z})$  with a probabilistic encoder  $q_{\phi}(\mathbf{z}|\mathbf{x})$
- ⇒ **Latent variable  $\mathbf{z}$**  approximates underlying parameters  $\boldsymbol{\nu}$
- ⇒ VAEs **do not** map inputs to a deterministic latent variable, but to a **probability space**  $p(\mathbf{z})$
- ⇒  $\theta$ : parameters of the decoder
- ⇒  $\phi$ : parameters of the encoder

# The Variational Lower Bound

- ▶ The Variational Lower Bound (ELBO) relates to two joint probability density functions:  $p_\theta$  and  $q_\phi$

$$\mathcal{L}(\theta, \phi; \mathbf{x}) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} \left[ \log \frac{p_\theta(\mathbf{x}, \mathbf{z})}{q_\phi(\mathbf{z}|\mathbf{x})} \right]$$

- ⇒  $p_\theta(\mathbf{x}, \mathbf{z})$ : joint distribution of  $\mathbf{x}$  and  $\mathbf{z}$
- ⇒  $q_\phi(\mathbf{z}|\mathbf{x})$ : approximate encoder posterior
- ⇒ Simplifies to

$$\mathcal{L}(\theta, \phi; \mathbf{x}) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}))$$

- ▶ Includes Kullback-Leibler divergence (KL-div) which is a distance in distribution space

$$D_{KL}(q_\phi(w)||p_\theta(w)) = \mathbb{E}_{q_\phi(w)} \left[ \log \frac{q_\phi(w)}{p_\theta(w)} \right]$$

- ▶ Implications and optimisation objective

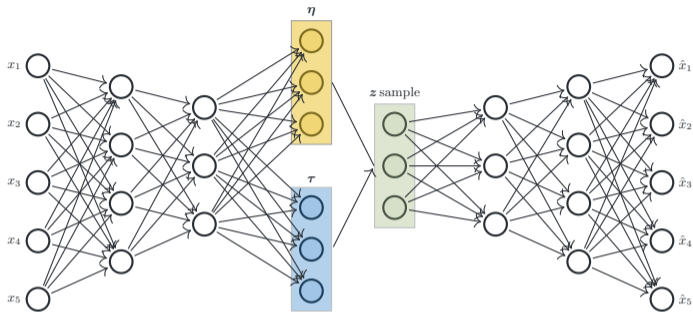
- ⇒ Log-likelihood relation

$$\log p_\theta(\mathbf{x}) \geq \mathcal{L}(\theta, \phi; \mathbf{x})$$

- ⇒ Objective function to train VAEs

- ⇒ Maximize ELBO to approximate the true log-likelihood.
- ⇒ Equivalent to minimising the negative ELBO (-ELBO)

# Deep Learning Implementation of Variational Autoencoders



## ► Implementation and parametrisation

⇒ **Neural Networks as parametrisers**

⇒  $p_{\theta}(\mathbf{x}|\mathbf{z})$  and  $q_{\phi}(\mathbf{z}|\mathbf{x})$  parameterised using deep neural networks

$$\theta = \{W_{l_1}, \dots, W_{l_L}, b_{l_1}, \dots, b_{l_L}\}$$

$$\phi = \{V_{l_1}, \dots, V_{l_L}, c_{l_1}, \dots, c_{l_L}\}$$

where  $W_l$  ( $V_l$ ),  $b_l$  ( $c_l$ ) are the weights and biases of the encoder (decoder) network

⇒ This setup enables the modeling of complex, non-linear relationships between observed data and latent variables

# Preparing Training Data with Theoretical and Experimental Inputs

## ► Generating Theoretical Predictions

- ⇒ Start by sampling the distribution of underlying inputs under hypothesis  $H_0$ ,  $p(\boldsymbol{\nu}|H_0)$
- ⇒ Compute the vector of observables  $\boldsymbol{x}^s$  for these values to obtain a sample of theoretical predictions:  
 $\boldsymbol{x}^s = (\boldsymbol{x}^1, \dots, \boldsymbol{x}^{n_{\text{sample}}})$

## ► Incorporating Experimental Uncertainties

- ⇒ Smear the samples with experimental uncertainties to simulate realistic observational data

$$\boldsymbol{x}'^s = \boldsymbol{x}^s + \boldsymbol{L}_{\Lambda^{\text{exp}}} \boldsymbol{w}$$

- ⇒  $\boldsymbol{L}_{\Lambda^{\text{exp}}}$ : Cholesky decomposition of the experimental covariance matrix  $\boldsymbol{\Lambda}^{\text{exp}}$
- ⇒  $\boldsymbol{w} \sim \mathcal{N}(\boldsymbol{w}; \mathbf{0}, \mathbf{1})$ : Normal noise vector simulating experimental noise

# Training the Variational Autoencoder with Experimental Uncertainties

## ► Training the VAE

- ⇒ Divide the smeared dataset into **training and testing datasets**
- ⇒ Use the training dataset to **optimise the parameters** of the VAE, **minimising** the -ELBO

$$\mathcal{L}(\theta, \phi; \mathbf{x}) = \mathbb{E}_{q_\phi(\mathbf{z}|\mathbf{x})} [\log p_\theta(\mathbf{x}|\mathbf{z})] - \beta D_{KL}(q_\phi(\mathbf{z}|\mathbf{x})||p_\theta(\mathbf{z}))$$

- ⇒ Notice the parameter  $\beta$  in the ELBO. It allows us to test the **generative properties** of the model
- ⇒ Objective: **approximate the full log-likelihood distribution** of the observables under  $H_0$  as closely as possible

## ► Model Assumptions

- ⇒ Assume distributions for **model simplicity**:

$$\begin{aligned} p_\theta(\mathbf{z}) &= \mathcal{N}(\mathbf{z}; \mathbf{0}, \mathbf{1}) \\ p_\theta(\mathbf{x}|\mathbf{z}) &= \mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \mathbf{\Lambda}_{\hat{\mathbf{x}}}) \text{ with } \mathbf{\Lambda}_{\hat{\mathbf{x}}} = \text{diag}(\sigma_{\hat{\mathbf{x}}}^2) \\ q_\phi(\mathbf{z}|\mathbf{x}) &= \mathcal{N}(\mathbf{z}; \boldsymbol{\eta}, \mathbf{\Lambda}_\tau) \text{ with } \mathbf{\Lambda}_\tau = \text{diag}(\tau^2) \end{aligned}$$

- ⇒ These **assumptions do not imply a Gaussian likelihood** but approximate relationships within the VAE structure

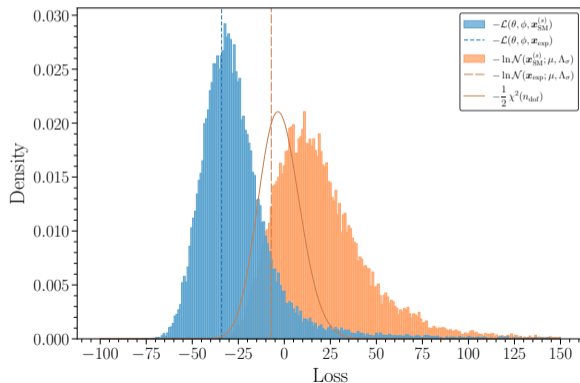
# Statistical Analysis Using Trained VAE for $b \rightarrow sl^+\ell^-$

## ► Analysing the test dataset

- ⇒ Compute the -ELBO distribution using the test dataset, approximating the full -log-likelihood under hypothesis  $H_0$ .
- ⇒ Evaluate -ELBO for the experimental data to compute the  $p$ -value

## ► Preliminary results

- ⇒ Performed for the  $b \rightarrow sl\ell$  dataset with promising preliminary outcomes



# Refining Model Tuning and Validity in VAE Training

## ► Challenges in Model Tuning

- ⇒ How do we determine the optimal dimensionality for the DNNs of the encoder and decoder, or the correct value of  $\beta$ ?
- ⇒ Could these choices bias the  $p$ -value?
- ⇒ The choice of neural networks' architecture and  $\beta$  significantly affects the model's performance and the fidelity of the statistical results

## ► Testing and Validating Model Parameters

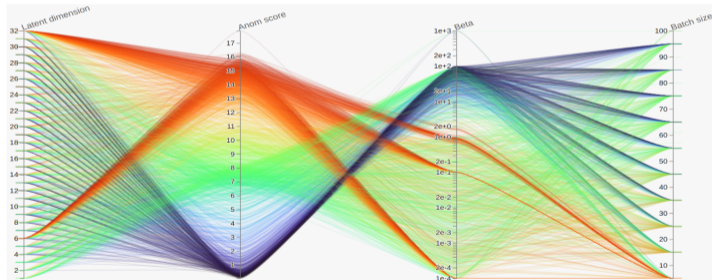
- ⇒ Ongoing research and empirical testing are essential to optimise these parameters while minimising biases

## ► Strategies for Validation and Hyperparameter Optimisation

- ⇒ Employ validation techniques to ensure model outputs are stable and reliable across various parameter configurations
- ⇒ Use synthetic datasets to evaluate the impact of hyperparameter adjustments on model performance

# Optimising Hyperparameters

- ▶ Generating artificially anomalous data
  - ⇒ Used to train the VAE across different configurations to optimise anomaly detection
- ▶ Tuning  $\beta$  in the ELBO
  - ⇒ Exploring the impact of  $\beta$  on anomaly detection and VAE's generative accuracy
- ▶ Pre-experimental blind analysis
  - ⇒ Ensures that the final measurement of experimental data is unbiased





# Outlook and Continued Research

- ▶ Deepening understanding of VAE parameters
  - ⇒ Exploring how different parameters influence the anomaly score and VAE's generative properties
- ▶ Ongoing hyperparameter optimisation
  - ⇒ Continuously refining the model to enhance its predictive accuracy and anomaly detection
- ▶ Addressing sparse covariance matrices
  - ⇒ Using random matrix theory techniques to sample observables and covariance matrices at the same time (LKJ distribution, Wishart distribution)
  - ⇒ Will allow us to quantify the uncertainty attached to many unnatural zeros in the experimental covariance matrix
- ▶ Expanding application scope
  - ⇒ Applying methodologies to SMEFT fits beyond  $b \rightarrow s\ell\ell$

Thank You!

# Backup Slides

# Dual-Branch VAE Architecture

- ⇒ We need the VAE to estimate  $p_{\theta}(\mathbf{x}|\mathbf{z}) = \mathcal{N}(\mathbf{x}; \hat{\mathbf{x}}, \mathbf{\Lambda}_{\hat{\mathbf{x}}})$
- ⇒ We need two output branches: one for  $\hat{\mathbf{x}}$  and one for  $\mathbf{\Lambda}_{\hat{\mathbf{x}}}$

