

TREE BOOSTING FOR LEARNING EFT PARAMETERS

S. Chatterjee, R. Schöfbeck, **D. Schwarz** (HEPHY Vienna)

Glühwein Workshop 2.0 - GlühWien

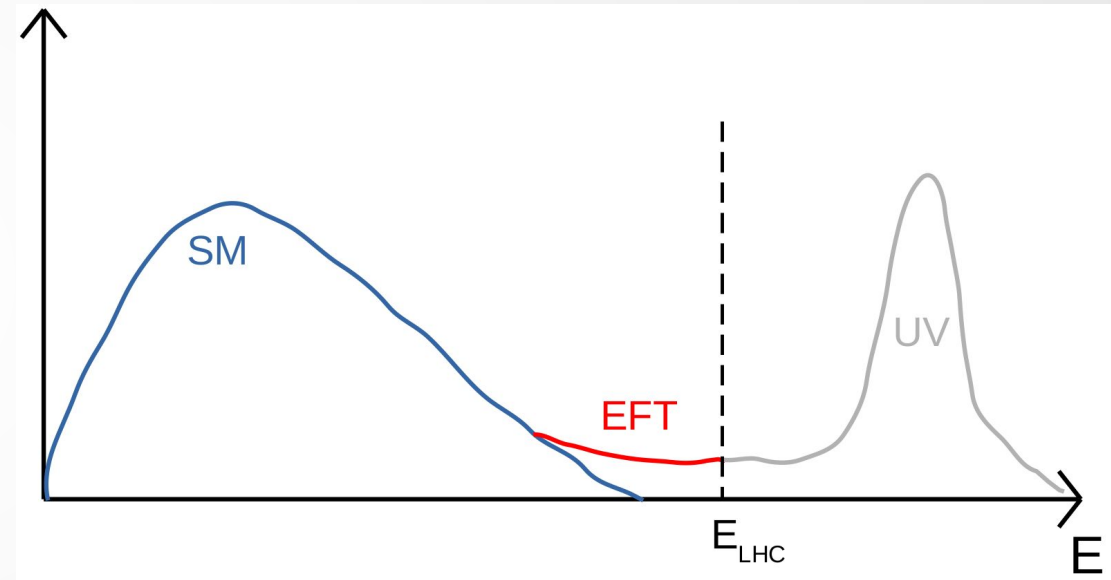
December 13-15, 2023

EFT in a nutshell

- No direct evidence for new physics at the LHC
- Might be out of reach of LHC
- Study small deviations from the SM with effective field theories
- Extend SM Lagrangian by dimension 6 operators:

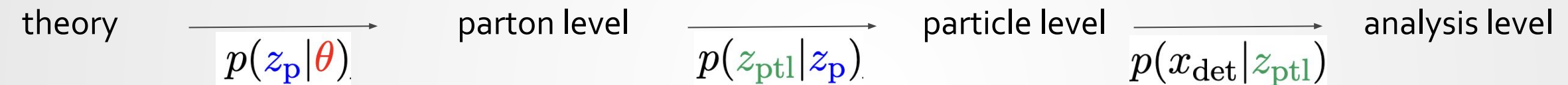
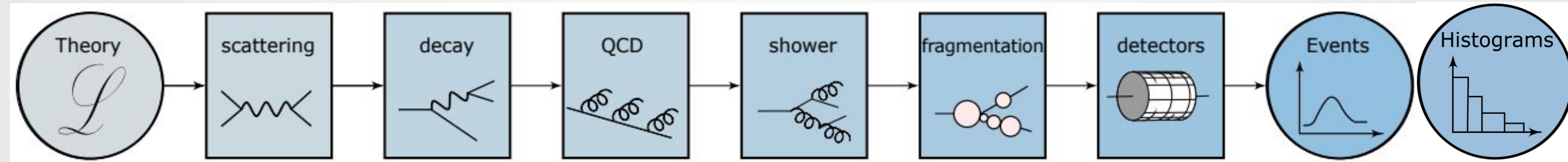
$$\mathcal{L}_{\text{eff}} = \mathcal{L}_{\text{SM}} + \sum_i \frac{c_i}{\Lambda^2} \mathcal{O}_i$$

- Wilson coefficients: coupling strength



Simulation sequence

image from [arXiv:2211.01421](https://arxiv.org/abs/2211.01421)



- In this picture, θ is an EFT parameter (Wilson coefficient)
- pdf $p(z_p|\theta)$ can be interpreted as diff. cross section $1/\sigma d\sigma(\theta)/dz_p$
- Now particle level with non-perturbative effects
- Reconstruction level depends on θ : $p(x_{\text{det}}|\theta) = \int dz_{\text{ptl}} \int dz_p [\dots]$
- Integrated over large latent space!
(observables we don't observe)

Likelihood ratio is the optimal statistic (Neyman-Pearson Lemma) for testing two hypotheses:

$$\text{LR}(x_{\text{det}}|\mathbf{H}_1, \mathbf{H}_2) \equiv \frac{p(x_{\text{det}}|\theta = \mathbf{H}_1)}{p(x_{\text{det}}|\theta = \mathbf{H}_2)}$$

Would like to evaluate for varying θ
Note: latent space integrated in num&den!

Situation in classification

- ML helps in parameterizing the likelihood ratio using “simulation based inference”
- Possible because we have a history of the generation sequence!
- In classification (signal vs background) $\theta \rightarrow \text{isSig} \in \{0,1\}$
- The loss function

$$L = \langle f(x)^2 \rangle_{\text{isSig}=1} + \langle (1 - f(x))^2 \rangle_{\text{isSig}=0} = \sum_{\text{isSig}} \int dx \dots$$

can be evaluated using generator level information (i.e. “isSignal”)

- With this loss, we converge to:

$$f^*(x) = \frac{1}{1 + \frac{p(x|\text{isSig}=1)}{p(x|\text{isSig}=0)}}$$

(here z is already integrated out)

- Latent space integrated in numerator and denominator
- Neyman-Pearson: optimal test statistic for a cross section measurement

-> Exactly what we need!

How we do the same for EFTs?

Can we learn EFT effects “on average”?

$$L = \sum_{\theta \in \mathcal{B}} \int d\mathbf{x} \left(p(\mathbf{x}|\theta) \hat{f}(\mathbf{x})^2 + p(\mathbf{x}|\text{SM})(1 - \hat{f}(\mathbf{x}))^2 \right)$$

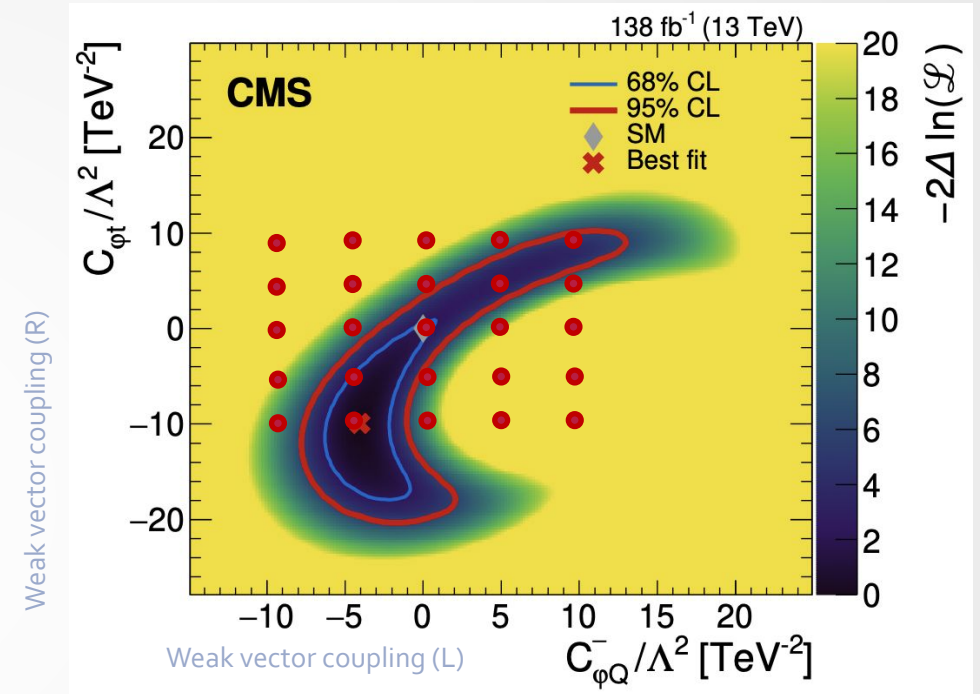
θ - ignorant



mixing signals &
case dependent mixes



$$f^*(\mathbf{x}) = \frac{1}{1 + r_{\mathcal{B}}(\mathbf{x})} \quad r_{\mathcal{B}}(\mathbf{x}) = \frac{\frac{1}{|\mathcal{B}|} \sum_{\theta \in \mathcal{B}} p(\mathbf{x}|\theta)}{p(\mathbf{x}|\text{SM})}$$



[[JHEP 12 \(2021\) 083](#)]

- We can try to learn EFT effects “on average”
 - Sending ‘mixed signals’ to the loss function
 - Averages the training data set -> suboptimal when linear effects dominate -> positive and negative can cancel
 - Classifier does not reflect knowledge on the θ -dependence, which we do know!
1. Go back and inject θ polynomial SMEFT dependence in estimator
 2. Exploit the fact that SMEFT predictions can be weighted: Only need one training data set

SMEFT reweighting

- Matrix element depends on Wilson coefficients:

$$|\mathcal{M}_{\text{SM}} + \theta_i \mathcal{M}_{\text{BSM}}^i|^2 = |\mathcal{M}_{\text{SM}}|^2 + \underbrace{\theta_i \text{Re} \mathcal{M}_{\text{SM}}^* \mathcal{M}_{\text{BSM}}^i}_{\text{interference, linear}} + \underbrace{\theta_i \theta_j \mathcal{M}_{\text{BSM}}^{i*} \mathcal{M}_{\text{BSM}}^j}_{\text{BSM, quadratic}}$$

- In simulated samples we store the ME for each event at various EFT points as weights
- Strictly polynomial structure allows to get function $w(\theta)$ for each event
- Sufficient EFT points needed: $(N + 2)(N + 1)/2$ for N different θ

-> Only one sample needed!

Exploiting the SMEFT reweighting

Start with classifier but now it is aware of θ :
(details later)

$$L = \sum_{\theta \in \mathcal{B}} \left(\langle \hat{f}(x; \theta)^2 \rangle_{\theta} + \langle (1 - \hat{f}(x; \theta))^2 \rangle_{\text{SM}} \right)$$

Exploiting the SMEFT reweighting

Start with classifier but now it is aware of θ :
(details later)

$$L = \sum_{\theta \in \mathcal{B}} \left(\langle \hat{f}(x; \theta)^2 \rangle_{\theta} + \langle (1 - \hat{f}(x; \theta))^2 \rangle_{\text{SM}} \right)$$

Write as integral over joint space:

$$= \sum_{\theta \in \mathcal{B}} \int dx dz \left(p(x, z | \theta) \hat{f}(x; \theta)^2 + p(x, z | \text{SM}) (1 - \hat{f}(x; \theta))^2 \right)$$

Exploiting the SMEFT reweighting

Start with classifier but now it is aware of θ :
(details later)

$$L = \sum_{\theta \in \mathcal{B}} \left(\langle \hat{f}(x; \theta)^2 \rangle_{\theta} + \langle (1 - \hat{f}(x; \theta))^2 \rangle_{\text{SM}} \right)$$

Write as integral over joint space:

$$= \sum_{\theta \in \mathcal{B}} \int dx dz \left(p(x, z | \theta) \hat{f}(x; \theta)^2 + p(x, z | \text{SM}) (1 - \hat{f}(x; \theta))^2 \right)$$

Now pull SM pdf in front:

$$= \sum_{\theta \in \mathcal{B}} \int dx dz p(x, z | \text{SM}) \left(r(x, z | \theta) \hat{f}(x; \theta)^2 + (1 - \hat{f}(x; \theta))^2 \right)$$

“joint” likelihood ratio $r(x, z | \theta)$ appears in the estimator

-> we need to know $r(x, z | \theta)$ for each event

Exploiting the SMEFT reweighting

Start with classifier but now it is aware of θ :
(details later)

$$L = \sum_{\theta \in \mathcal{B}} \left(\langle \hat{f}(x; \theta)^2 \rangle_{\theta} + \langle (1 - \hat{f}(x; \theta))^2 \rangle_{\text{SM}} \right)$$

Write as integral over joint space:

$$= \sum_{\theta \in \mathcal{B}} \int dx dz \left(p(x, z | \theta) \hat{f}(x; \theta)^2 + p(x, z | \text{SM}) (1 - \hat{f}(x; \theta))^2 \right)$$

Now pull SM pdf in front:

$$= \sum_{\theta \in \mathcal{B}} \int dx dz p(x, z | \text{SM}) \left(r(x, z | \theta) \hat{f}(x; \theta)^2 + (1 - \hat{f}(x; \theta))^2 \right)$$

“joint” likelihood ratio $r(x, z | \theta)$ appears in the estimator

-> we need to know $r(x, z | \theta)$ for each event

Let's write it down, assume θ only affects parton level:

Factorizes and PDFs cancel [well established, references in backup]

$$r = \frac{p(x_{\text{det}}, \dots, z_{\text{ptl}}, \dots, z_{\text{p}} | \theta)}{p(x_{\text{det}}, \dots, z_{\text{ptl}}, \dots, z_{\text{p}} | \text{SM})} = \frac{p(x_{\text{det}} | z_{\text{ptl}}) \cdots p(z_{\text{ptl}} | z_{\text{p}}) \cdots p(z_{\text{p}} | \theta)}{p(x_{\text{det}} | z_{\text{ptl}}) \cdots p(z_{\text{ptl}} | z_{\text{p}}) \cdots p(z_{\text{p}} | \text{SM})} = \frac{p(z_{\text{p}} | \theta)}{p(z_{\text{p}} | \text{SM})} \sim \frac{|\mathcal{M}(z_{\text{p}}, \theta)|^2}{|\mathcal{M}(z_{\text{p}}, \text{SM})|^2}$$

This is exactly our weight function $w(\theta)$!

Parametrized Classifiers

Loss:
$$L = \sum_{\theta \in \mathcal{B}} \int dx dz p(x, z | \text{SM}) \left(r(x, z | \theta) \hat{f}(x; \theta)^2 + (1 - \hat{f}(x; \theta))^2 \right)$$
 here: MSE, but cross entropy also works

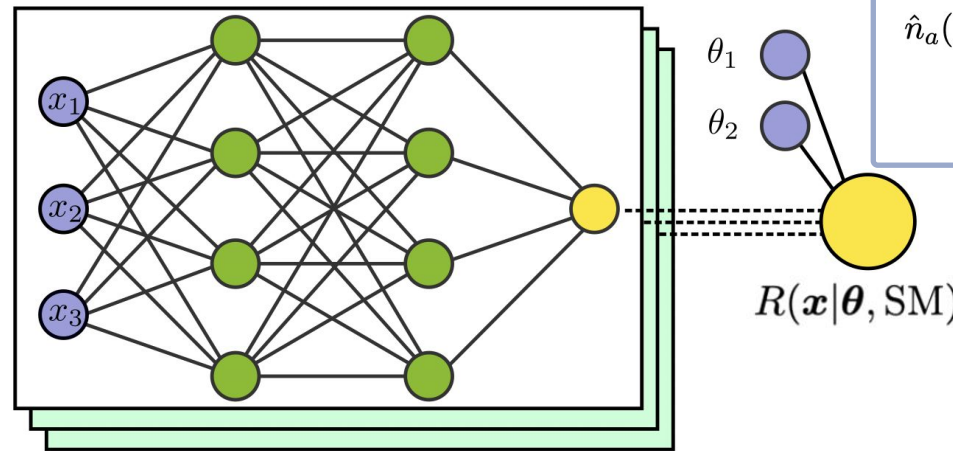
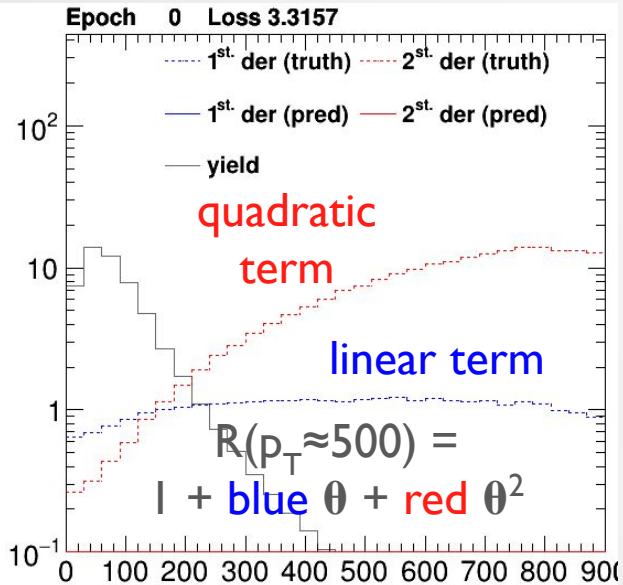
Similar to
 S. Chen, A. Glioti,
 G. Panico, A. Wulzer
[JHEP 05 \(2021\) 247](#)
[arXiv:2308.05704](#)

$$\hat{f}(x; \theta) = \frac{1}{1 + \hat{R}(x; \theta)}$$

likelihood trick, we need to know the likelihood ratio R

Insert knowledge of θ -dependence, assume general polynomial:

$$\hat{R}(x; \theta) = \left(1 + \sum_a \theta_a \hat{n}_a(x) \right)^2 + \sum_a \left(\sum_{b \geq a} \theta_b \hat{n}_{ab}(x) \right)^2$$



$$\hat{n}_a(x) \rightarrow \left. \frac{\partial_a p(x | \theta)}{p(x | \text{SM})} \right|_{\theta = \text{SM}} = \left. \frac{\partial_a \int dz p(x, z | \theta)}{\int dz p(x, z | \text{SM})} \right|_{\theta = \text{SM}}$$

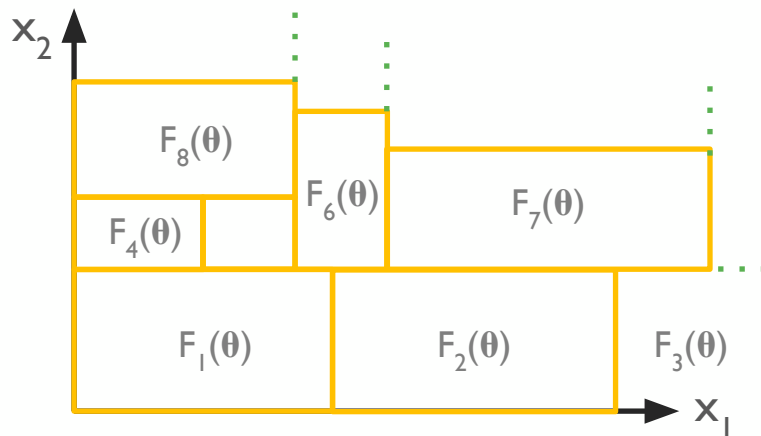
Integrates latent space!

Fit coefficient functions with NN. This is the current status - Why would you want to use trees instead?

The tree algorithm

[[arXiv:2107.10859](https://arxiv.org/abs/2107.10859), [arXiv:2205.12976](https://arxiv.org/abs/2205.12976)]

phase-space partitioning



A simple tree

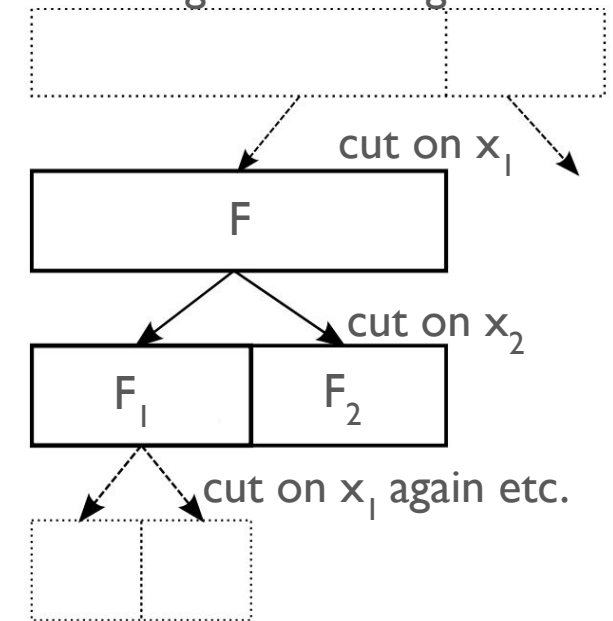
index-function (non-linearity)

$$\hat{F}(\mathbf{x}, \boldsymbol{\theta}) = \sum_{j \in \mathcal{J}} \mathbb{1}_j(\mathbf{x}) F_j(\boldsymbol{\theta})$$

phase space partitioning \mathcal{J} prediction F_j

need to solve for partitioning \mathcal{J} and $\{F_j\}$

training phase:
e.g. "CART" algo



- A tree is a hierarchical phase-space partitioning (\mathcal{J})
- Novelty in the Boosted Information Tree is that we associate each region j with a polynomial $F_j(\boldsymbol{\theta})$
- Fitting tree means optimizing "node split positions" on some loss.
- Training on **ENSEMBLE!**

Trees for SMEFT

[[arXiv:2107.10859](#), [arXiv:2205.12976](#)]

Want to regress in r , making use of θ -dependence:

$$r(x, z|\theta) = \frac{d\sigma(\mathbf{x}, \theta)/d\mathbf{x}}{d\sigma(\mathbf{x}, \text{SM})/d\mathbf{x}}$$

→ will allow to compute the optimal LLR test statistic $q(\mathcal{D})$

Trees for SMEFT

[[arXiv:2107.10859](https://arxiv.org/abs/2107.10859), [arXiv:2205.12976](https://arxiv.org/abs/2205.12976)]

Want to regress in r , making use of θ -dependence:

$$r(x, z|\theta) = \frac{d\sigma(\mathbf{x}, \theta)/d\mathbf{x}}{d\sigma(\mathbf{x}, \text{SM})/d\mathbf{x}}$$

→ will allow to compute the optimal LLR test statistic $q(\mathcal{D})$

$$L = \sum_{\theta \in \mathcal{B}} \int d\mathbf{x} dz p(\mathbf{x}, z|\text{SM}) \left(r(x, z|\theta) - \hat{F}(\mathbf{x}, \theta) \right)^2$$

Tree ansatz

$F_j(\theta)$ polynomial with const. coeff.
(per node)

$$\hat{F}(\mathbf{x}, \theta) = \sum_{j \in \mathcal{J}} \mathbb{1}_j(\mathbf{x}) F_j(\theta)$$

↑
find optimal
partitioning

find optimal
predictor

Trees for SMEFT

[[arXiv:2107.10859](https://arxiv.org/abs/2107.10859), [arXiv:2205.12976](https://arxiv.org/abs/2205.12976)]

Want to regress in r , making use of θ -dependence:

$$r(x, z|\theta) = \frac{d\sigma(\mathbf{x}, \theta)/d\mathbf{x}}{d\sigma(\mathbf{x}, \text{SM})/d\mathbf{x}}$$

→ will allow to compute the optimal LLR test statistic $q(\mathcal{D})$

$$L = \sum_{\theta \in \mathcal{B}} \int d\mathbf{x} dz p(\mathbf{x}, z|\text{SM}) \left(r(x, z|\theta) - \hat{F}(\mathbf{x}, \theta) \right)^2$$

Tree ansatz

$F_j(\theta)$ polynomial with const. coeff.
(per node)

$$\hat{F}(\mathbf{x}, \theta) = \sum_{j \in \mathcal{J}} \underbrace{\mathbb{1}_j(\mathbf{x})}_{\text{find optimal partitioning}} \underbrace{F_j(\theta)}_{\text{find optimal predictor}}$$

Injecting the ansatz, this minimizes the loss:

$$F_j(\theta) = \frac{\sum_{i \in j} w_i(\theta)}{\sum_{i \in j} w_i(\theta_0)} \quad \text{sum up event weights within node } j \text{ \& divide} \\ \rightarrow \text{very fast!}$$

No trainable parameters in the predictor
(Regression by means of classification)

Trees for SMEFT

[arXiv:2107.10859, arXiv:2205.12976]

Want to regress in r , making use of θ -dependence:

$$r(x, z|\theta) = \frac{d\sigma(\mathbf{x}, \theta)/d\mathbf{x}}{d\sigma(\mathbf{x}, \text{SM})/d\mathbf{x}}$$

→ will allow to compute the optimal LLR test statistic $q(\mathcal{D})$

$$L = \sum_{\theta \in \mathcal{B}} \int d\mathbf{x} dz p(\mathbf{x}, z|\text{SM}) \left(r(x, z|\theta) - \hat{F}(\mathbf{x}, \theta) \right)^2$$

Tree ansatz

$F_j(\theta)$ polynomial with const. coeff. (per node)

$$\hat{F}(\mathbf{x}, \theta) = \sum_{j \in \mathcal{J}} \underbrace{\mathbb{1}_j(\mathbf{x})}_{\text{find optimal partitioning}} \underbrace{F_j(\theta)}_{\text{find optimal predictor}}$$

Injecting the ansatz, this minimizes the loss:

$$F_j(\theta) = \frac{\sum_{i \in j} w_i(\theta)}{\sum_{i \in j} w_i(\theta_0)} \quad \text{sum up event weights within node } j \text{ \& divide} \rightarrow \text{very fast!}$$

No trainable parameters in the predictor (Regression by means of classification)

Leading to this loss (minimizing partitioning!):

$$L = - \sum_{\theta \in \mathcal{B}} \sum_{j \in \mathcal{J}} \frac{w_j^2(\theta)}{w_j(\theta_0)} = - \sum_{j \in \mathcal{J}} \sum_{\theta \in \mathcal{B}} \theta^a \theta^b I_{ab}^{(j)} + \mathcal{O}(\theta - \theta_0)^3$$

Taylor expansion shows: We're optimizing the Fisher information!

This is only one tree → boost!

Tree boosting

[[arXiv:2107.10859](https://arxiv.org/abs/2107.10859), [arXiv:2205.12976](https://arxiv.org/abs/2205.12976)]

- Boosting: Fit model iteratively to pseudo-residuals of the preceding iteration with learning rate η

- Ansatz :

$$\hat{F}^{(b)}(\mathbf{x}, \boldsymbol{\theta}) = \underbrace{\hat{f}(\mathbf{x}, \boldsymbol{\theta})}_{\text{current iteration}} + \eta \underbrace{\hat{F}^{(b-1)}(\mathbf{x}, \boldsymbol{\theta})}_{\text{previous iteration}}$$

current
iteration

previous
iteration

1. Fit first tree
2. Multiply output with η
3. Fit next tree to residual

Insert into the loss function:

$$L[\hat{f}^{(b)}] = \sum_{\boldsymbol{\theta} \in \mathcal{B}} \int d\mathbf{x} dz p(\mathbf{x}, z | \text{SM}) \left(r(\mathbf{x}, z | \boldsymbol{\theta}) - \underbrace{\eta \hat{F}^{(b-1)}(\mathbf{x}, \boldsymbol{\theta})}_{\text{previous iteration}} - \underbrace{\hat{f}^{(b)}(\mathbf{x}, \boldsymbol{\theta})}_{\text{current iteration}} \right)^2$$

current
iteration

previous iteration

current iteration

pseudo-residual, expressed
with weights:

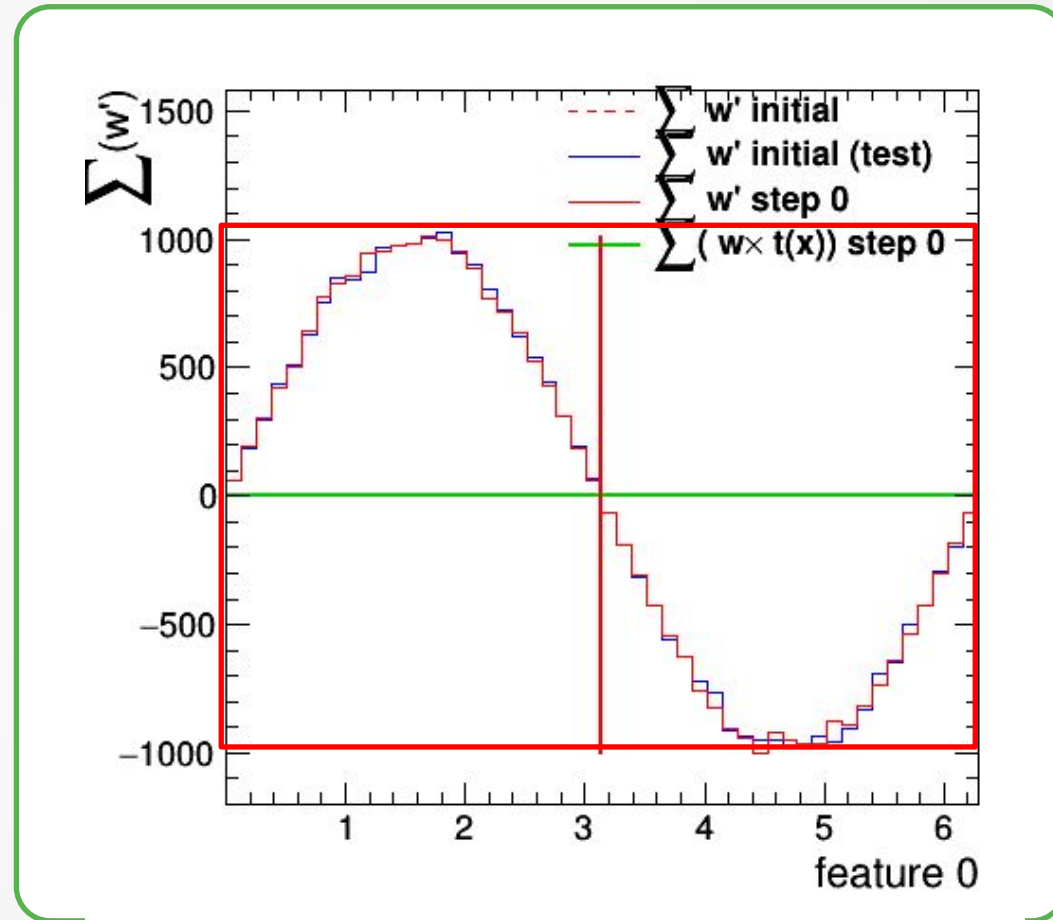
$$w_i^{(b)}(\boldsymbol{\theta}) \rightarrow w_i^{(b-1)}(\boldsymbol{\theta}) - \eta w_i^{(b-1)}(\boldsymbol{\theta}_0) \hat{F}^{(b-1)}(\mathbf{x}_i, \boldsymbol{\theta})$$

-> Iterative fitting and re-weighting (similar to AdaBoost)

.... perform this iteratively
"Boosted Information Tree"

1D TOY EXAMPLE

- $\text{pdf}(x|\theta) = 1/(2\pi) - \theta \sin(x)$
- We predict the first derivative of the pdf wrt. to the parameter



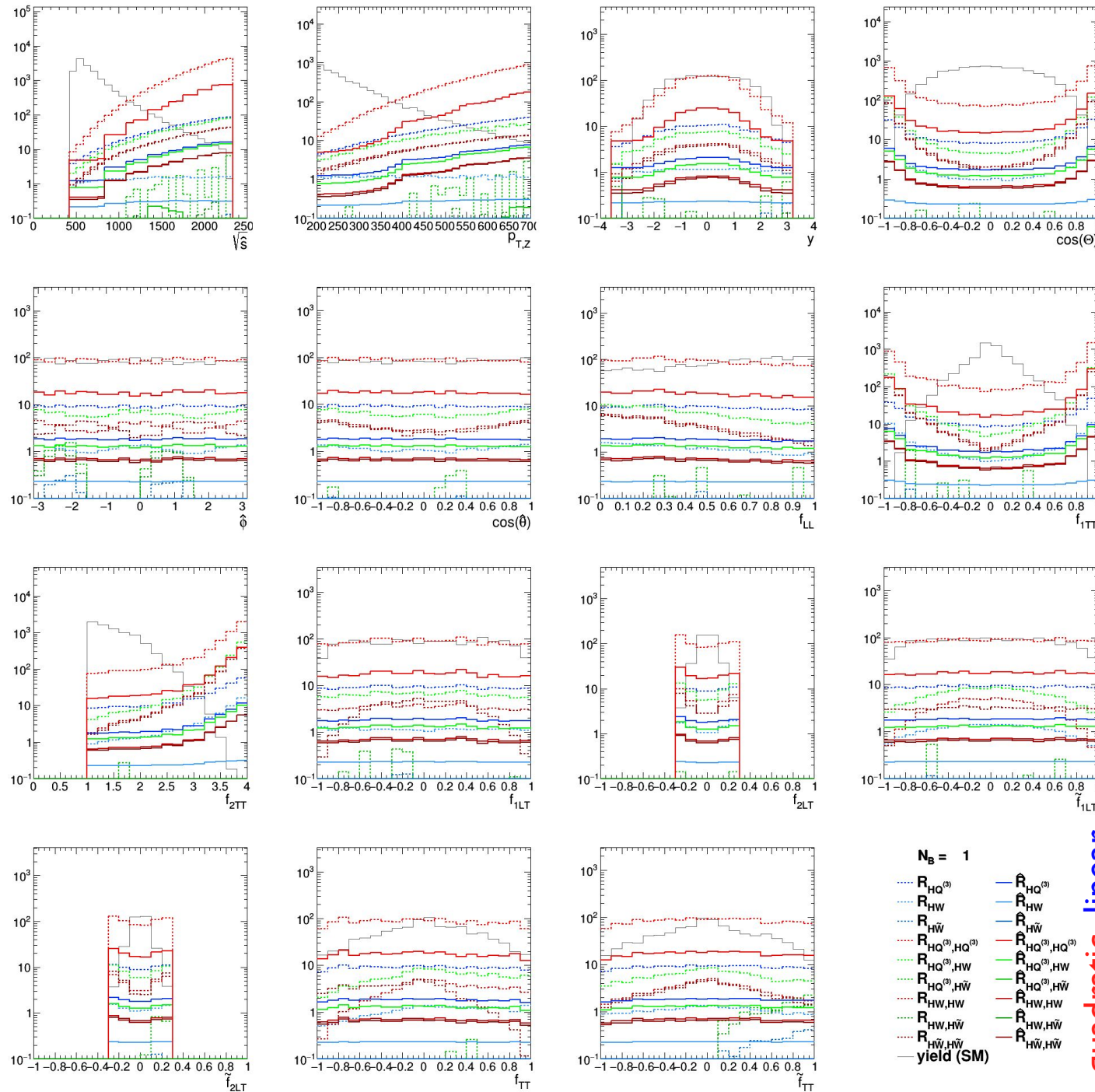
GIF animation not showing in pdf

- Realistic case: model of the ZH process

- “Boosted Information Tree (BIT)”

- 3 WC, 9 DOF, 500k events, ZH
- 200 trees, $D=5$, 9 minutes of training
- also more realistic study, including backgrounds [2107.10859], [2205.12976]

- Learning coefficient functions to compute parametrized optimal observables

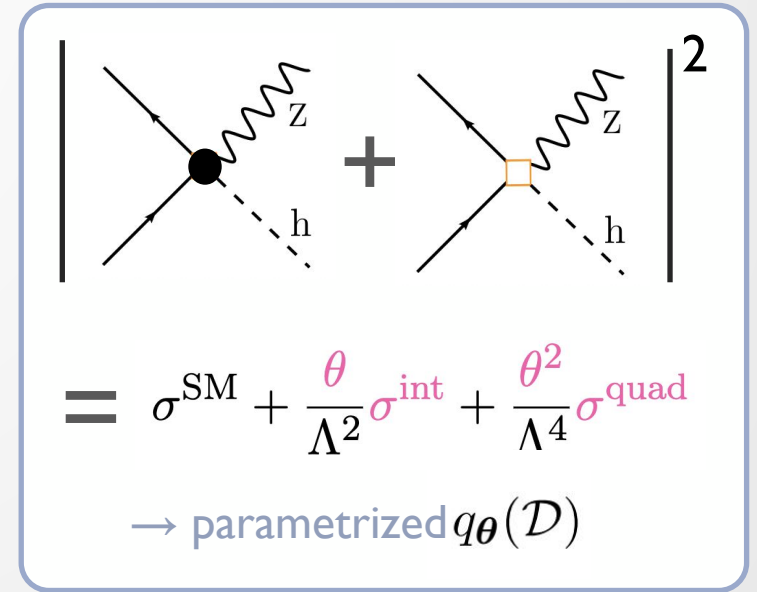


$N_B = 1$

$R_{HQ^{(3)}}$	$\hat{R}_{HQ^{(3)}}$
R_{HW}	\hat{R}_{HW}
$R_{H\tilde{W}}$	$\hat{R}_{H\tilde{W}}$
$R_{HQ^{(3)},HQ^{(3)}}$	$\hat{R}_{HQ^{(3)},HQ^{(3)}}$
$R_{HQ^{(3)},HW}$	$\hat{R}_{HQ^{(3)},HW}$
$R_{HQ^{(3)},H\tilde{W}}$	$\hat{R}_{HQ^{(3)},H\tilde{W}}$
$R_{HW,H\tilde{W}}$	$\hat{R}_{HW,H\tilde{W}}$
$R_{H\tilde{W},H\tilde{W}}$	$\hat{R}_{H\tilde{W},H\tilde{W}}$
$R_{H\tilde{W},HW}$	$\hat{R}_{H\tilde{W},HW}$
$R_{H\tilde{W},H\tilde{W}}$	$\hat{R}_{H\tilde{W},H\tilde{W}}$
yield (SM)	

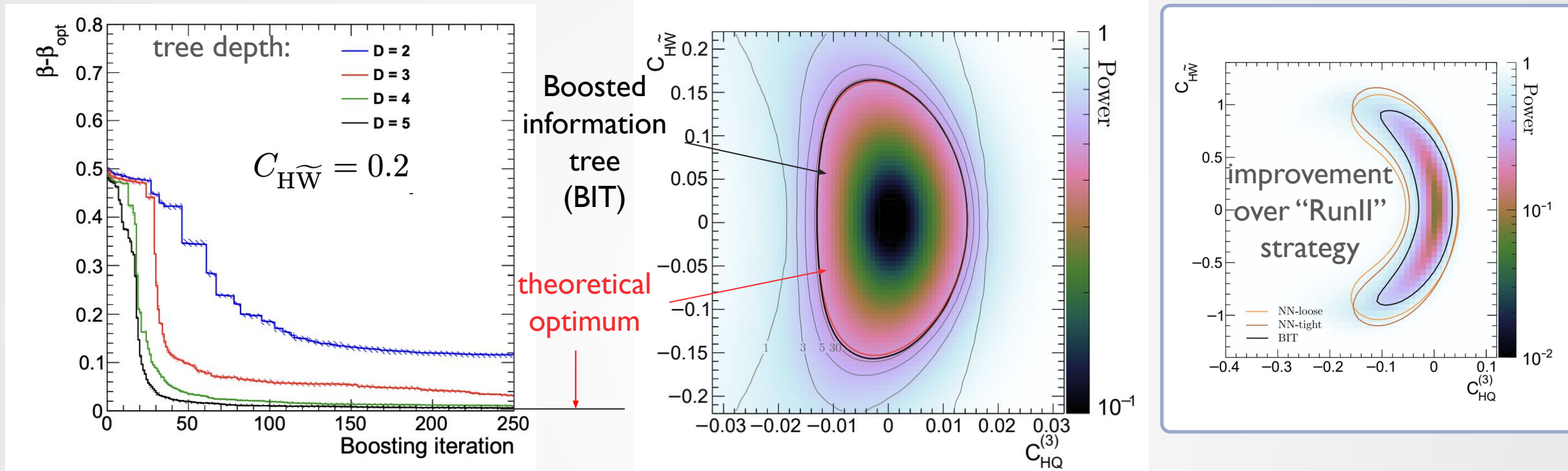
quadratic

linear



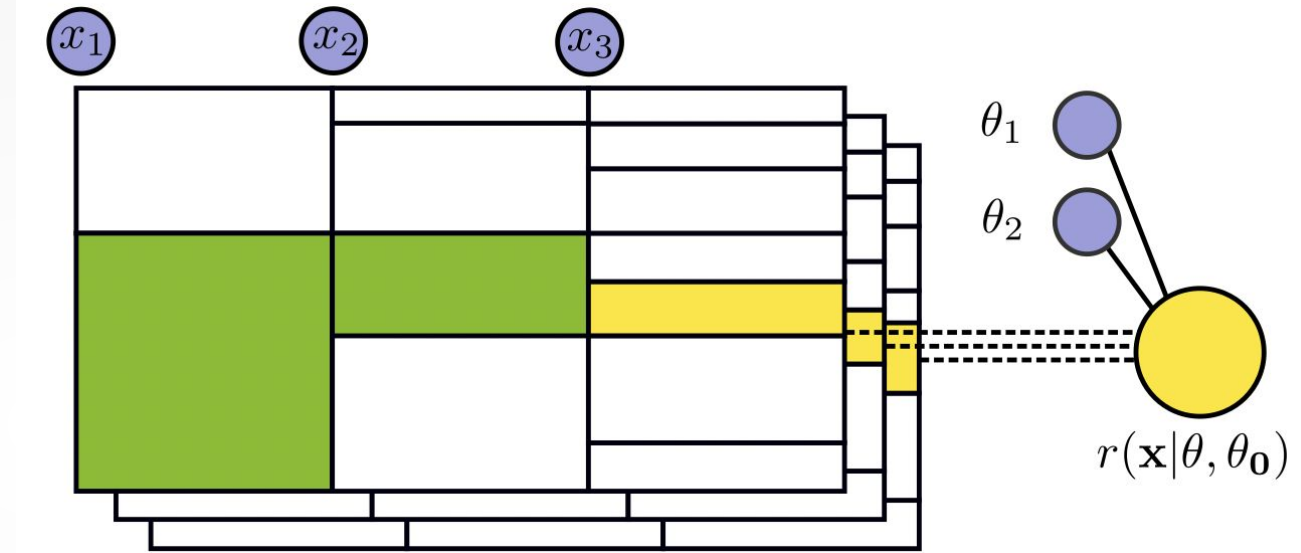
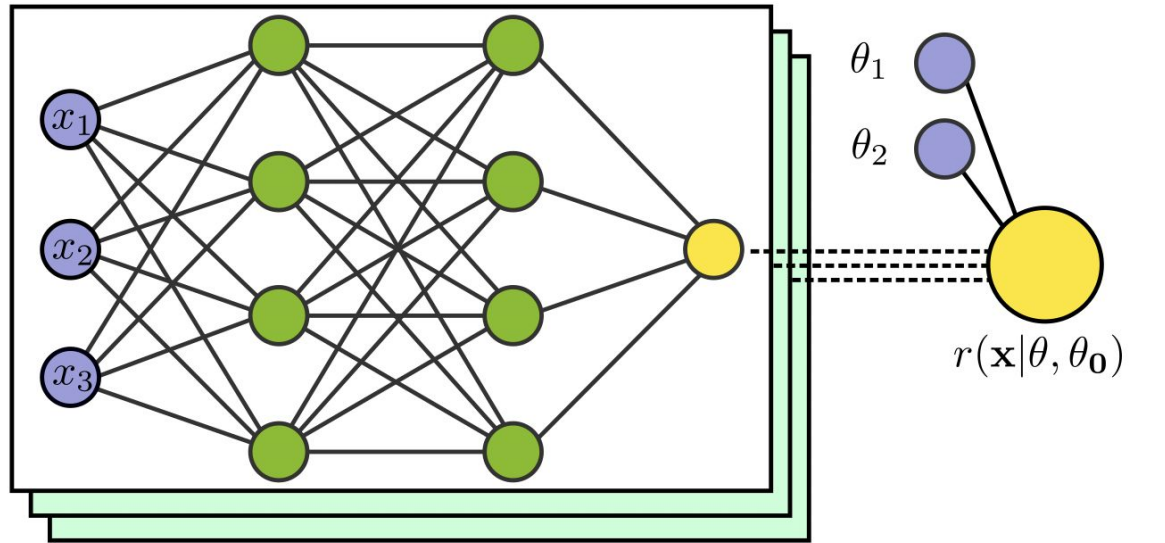
OPTIMALITY IN TEST CASES

[[arXiv:2107.10859](https://arxiv.org/abs/2107.10859), [arXiv:2205.12976](https://arxiv.org/abs/2205.12976)]



- Obtain parametrized classifiers with 20-40% improvements in 2D toy cases (NOT marginalized!)
- Easy to train with backgrounds
- Most gain when performing unbinned analysis (otherwise choices in binning must be made)
- Strong performance for marginalized limits in high dimensions

Network vs Trees



- Given a phase space region with EFT dependence: NN must select & predict
- In the Boosted Information Tree, the weak learner only selects
 - The prediction (F_j) is computed from the boxed events → integrates latent space
 - The regression problem is solved with computational complexity of classification
 - Speed advantage at high operator dimensions!

Summary

- Various algorithms predict SMEFT dependence, mostly capitalizing on weighted simulation
 - Provide NP-optimal observables for hypothesis tests
- Trees **eliminate** the **latent space integration** and are suitable to weight-based SMEFT predictions
- Parametrized classifiers "ask" for **unbinned analysis** -> Needs development in experiments

*Get in touch
with
us for an
in-person
walkthrough*

```
from MultiBoostedInformationTree import MultiBoostedInformationTree
bit = MultiBoostedInformationTree(
    training_features = training_features,
    training_weights = training_weights,
    base_points = base_points,
    feature_names = model.feature_names,
    **model.multi_bit_cfg
)
bit.boost()
test_predictions = bit.vectorized_predict(test_features)
```

REFERENCES

- **Madminer**: Neural networks based likelihood-free inference & related techniques
 - K. Cranmer, J. Pavez, and G. Louppe [[1506.02169](#)]
 - J. Brehmer, K. Cranmer, G. Louppe, J. Pavez [[1805.00013](#)] [[1805.00020](#)] [[1805.12244](#)]
 - J. Brehmer, F. Kling, I. Espejo, K. Cranmer [[1907.10621](#)]
 - J. Brehmer, S. Dawson, S. Homiller, F. Kling, T. Plehn [[1908.06980](#)]
 - A. Butter, T. Plehn, N. Soybelman, J. Brehmer [[2109.10414](#)]
 - established many of the *main ideas* & *statistical interpretation* in various *NN applications*
- **Weight derivative regression** (A.Valassi) [[2003.12853](#)]
- **Parametrized classifiers** for SM-EFT: NN with quadratic structure
- S. Chen, A. Glioti, G. Panico, A. Wulzer [[JHEP 05 \(2021\) 247](#)] [[arXiv:2308.05704](#)]
- **Boosted Information Trees**: Tree algorithms & boosting
 - S. Chatterjee, S. Rohshap, N. Frohner, R. Schöfbeck, D. Schwarz [[2107.10859](#)], [[2205.12976](#)]
- **ML₄EFT** R. Ambrosio, J. Hoeve, M. Madigan, J. Rojo, V. Sanz [[2211.02058](#)]
- All approaches are “SMEFT-specific ML” with differences mostly on the practical side

