# *Optimising HEP parameter fits: event-by-event sensitivities, weight derivative regression*

Andrea Valassi (CERN IT-DI)

CHEP2019, 7 November 2019 – Adelaide, Australia
https://indico.cern.ch/event/773049/contributions/3476059

# This is a follow-up of my CHEP2018 talk about
# *binned fits of a parameter θ*

**Evaluation and training metrics:**
**Fisher Information Part**

## Previous CHEP2018 talk

Event selection
Binary classification

Bin-by-bin sensitivity to θ

Cross-section fits (FIP1, FIP2)

Medical Diagnostics (AUC),
Information Retrieval (F1)

## This CHEP2019 talk

Event partitioning
Non-binary **regression**

*WEIGHT DERIVATIVE REGRESSION*

**Event-by-event sensitivity to θ**

*MINIMUM ERROR WITH AN IDEAL DETECTOR*

Mass fits, Coupling fits (FIP3)

**Meteorology (MSE, Brier),**
Medical Prognostics

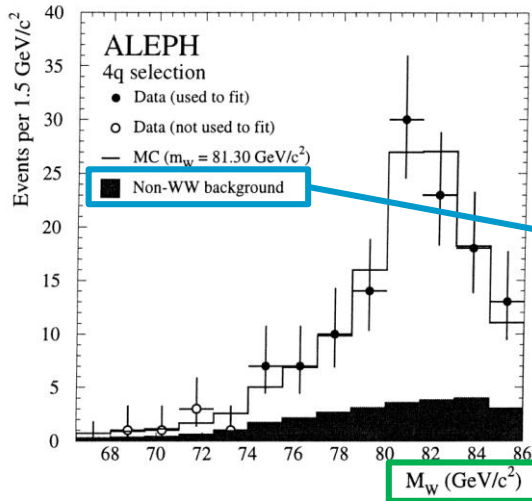**Compare to and learn**
**from other domains**

# Outline

- 1 - HEP parameter fits and Weight Derivative Regression

- 2 - Learning from others

- Conclusions

_This talk only provides some maths and some literature review_

_No toy model or concrete applications are presented_

# 1 – Binned fit of a parameter θ

ALEPH Collaboration, *Measurement of the W mass by direct reconstruction in* $e^+e^-$ *collisions at 172 GeV*, Phys. Lett. B 422 (1998) 384. doi:10.1016/S0370-2693(98)00062-8



ALEPH
4q selection
- • Data (used to fit)
- ○ Data (not used to fit)
- — MC ($m_W = 81.30$ GeV/$c^2$)
- ■ Non-WW background

There are two handles **to minimize the statistical error** $\Delta\theta$ :

## 1. Event selection
Signal-background discrimination

## 2. Event partitioning
Variable(s) for the distribution fit

My CHEP2018 talk: event selection

This CHEP2019 talk: event partitioning
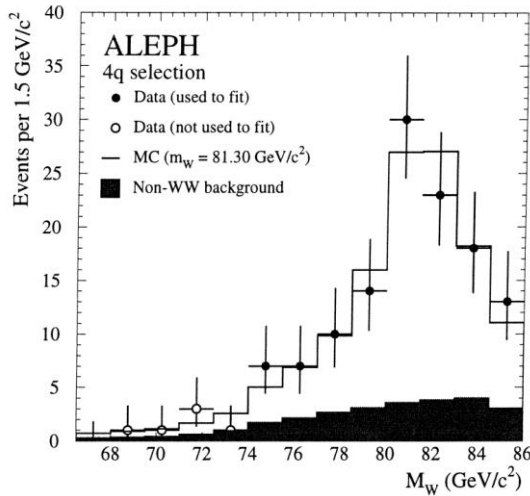(selection is a special case of partitioning)

$$m_{\mathrm{W}} = 81.30 \pm 0.47 (\text{stat.}) \pm 0.11 (\text{syst.}) \, \text{GeV}/c^2$$

*I only discuss the **statistical error** $\Delta\theta$ in this talk*
(I ignore systematic errors, even if at LHC they are the limitation)

# Fisher Information $\frac{1}{(\Delta\theta)^2}$ from bin-by-bin sensitivities

ALEPH Collaboration, *Measurement of the W mass by direct reconstruction in $e^+e^-$ collisions at 172 GeV*, Phys. Lett. B 422 (1998) 384. doi:10.1016/S0370-2693(98)00062-8



ALEPH
4q selection
- ● Data (used to fit)
- ○ Data (not used to fit)
- — MC ($m_w$ = 81.30 GeV/$c^2$)
- ■ Non-WW background

For a given partitioning scheme with K bins
($n_k$ is the number of selected events in bin k):

Statistical errors:
information adds up
(independent bins)

**Bin-by-bin sensitivity to θ**

Recap CHEP2018 talk

$$\mathcal{I}_\theta = \frac{1}{(\Delta\theta)^2} = \sum_{k=1}^{K} \frac{1}{(\Delta\theta)_k^2} = \sum_{k=1}^{K} n_k \left( \frac{1}{n_k} \frac{\partial n_k}{\partial \theta} \right)^2$$
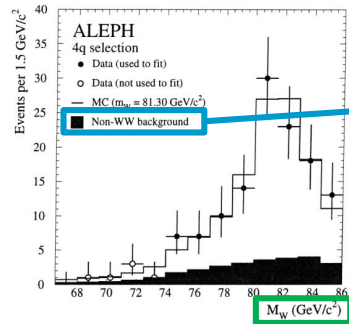
**Minimizing Δθ is equivalent to maximizing I$_θ$**

$$m_{\mathrm{w}} = 81.30 \pm 0.47(\text{stat.}) \pm 0.11(\text{syst.}) \text{ GeV}/c^2$$

# Fisher Information Part (FIP)

ALEPH Collaboration, *Measurement of the W mass by direct reconstruction in $e^+e^-$ collisions at 172 GeV*, Phys. Lett. B 422 (1998) 384. doi:10.1016/S0370-2693(98)00062-8



There are two handles to minimize the statistical error $\Delta\theta$ :

### 1. Event selection
Signal-background discrimination

### 2. Event partitioning
Variable(s) for the distribution fit

## My CHEP2018 talk:
## FIP evaluation of event selection

For a given data set and given partitioning, FIP compares $I_\theta$ to $I_\theta^{(ideal)}$ for the **ideal selection** (select all signal, reject all bkg)

*Recap CHEP2018 talk*

## Fisher Information Part (FIP): the fraction of the information available *"in an ideal case"* retained by a given analysis

$$\text{FIP} = \frac{\mathcal{I}_\theta}{\mathcal{I}_\theta^{(ideal)}} = \frac{(\Delta\theta^{(ideal)})^2}{(\Delta\theta)^2} \leq 100\%$$

**FIP is a metric between 0 and 1 – higher is better**

## *This CHEP2019 talk:*
## *FIP evaluation of event partitioning*

*For a given data set,*
*FIP compares $I_\theta$ to $I_\theta^{(ideal)}$ for the ideal partitioning (and the ideal selection)*

*But what is the smallest statistical error achievable on a given data set with ideal partitioning and selection?*
*Enter event-by-event sensitivities*

# Event-by-event Monte Carlo reweighting



ALEPH Collaboration, *Measurement of the W mass by direct reconstruction in $e^+e^-$ collisions at 172 GeV*, Phys. Lett. B 422 (1998) 384. doi:10.1016/S0370-2693(98)00062-8

$$w_i(m_W, \Gamma_W) = \frac{|\mathcal{M}(m_W, \Gamma_W, p_i^1, p_i^2, p_i^3, p_i^4)|^2}{|\mathcal{M}(m_W^{MC}, \Gamma_W^{MC}, p_i^1, p_i^2, p_i^3, p_i^4)|^2}$$

Fit for θ → Compare data in bin k to
*model prediction $n_k$ as a function of θ*

$$n_k(\theta) = \sum_{i \in k} w_i(\theta) = \overset{\text{Sig}}{\sum_{i \in k}} w_i(\theta) + \overset{\text{Bkg}}{\sum_{i \in k}} w_i = s_k(\theta) + b_k$$

*1. Generate signal sample at $\theta_{ref}$, with $w_i(\theta_{ref})=1$*
(By definition, background does not depend on θ)

*2. Full detector simulation*
(MC truth event properties $\mathbf{x}_i^{(true)}$ → observed event properties $\mathbf{x}_i$)

*3. Reweight each event by matrix element ratio*

$$w_i(\theta) = \frac{\text{Prob}_{(\theta)}(\mathbf{x}_i^{(true)})}{\text{Prob}_{(\theta_{ref})}(\mathbf{x}_i^{(true)})} = \frac{|\mathcal{M}(\theta, \mathbf{x}_i^{(true)})|^2}{|\mathcal{M}(\theta_{ref}, \mathbf{x}_i^{(true)})|^2}$$

*Monte Carlo reweighting: used extensively at LEP*
*Simpler than Matrix Element Method (no integration)*
*[see Gainer2014, Mattelaer2016 for hadron colliders]*

J. S. Gainer, J. Lykken, K. T. Matchev, S. Mrenna, M. Park, *Exploring theory space with Monte Carlo reweighting*, JHEP 2014 (2014) 78. doi:10.1007/JHEP10(2014)078

O. Mattelaer, *On the maximal use of Monte Carlo samples: re-weighting events at NLO accuracy*, Eur. Phys. J. C 76 (2016) 674. doi:10.1140/epjc/s10052-016-4533-7

# Event-by-event sensitivities $\gamma_i$: MC weight derivatives

Bin-by-bin model prediction $n_k(\theta)$

$$n_k(\theta) = \sum_{i \in k} w_i(\theta) = \sum_{i \in k}^{\text{Sig}} w_i(\theta) + \sum_{i \in k}^{\text{Bkg}} w_i = s_k(\theta) + b_k$$

Define the **event-by-event sensitivity $\gamma_i$ to θ** as the *derivative with respect to θ of the MC weight $w_i$*

$$\gamma_i|_\theta = \left( \frac{1}{w_i} \frac{\partial w_i}{\partial \theta} \right)_\theta \longrightarrow \gamma_i = \gamma_i|_{\theta=\theta_{\text{ref}}} = \left( \frac{\partial w_i}{\partial \theta} \right)_{\theta=\theta_{\text{ref}}}$$

(normalized by $1/w_i$, but $w_i(\theta_{\text{ref}})=1$ at the reference $\theta=\theta_{\text{ref}}$)

The **bin-by-bin sensitivity** to θ in bin k is the *average in bin k of the event-by-event sensitivity $\gamma_i$ to θ*

$$\left( \frac{1}{n_k} \frac{\partial n_k}{\partial \theta} \right)_{\theta=\theta_{\text{ref}}} = \frac{1}{n_k} \sum_{i \in k} \gamma_i = \langle \gamma \rangle_k = \frac{1}{n_k} \frac{\partial n_k}{\partial \theta}$$

# Beyond the signal-background dichotomy

***Background events have $\gamma_i=0$***
*because by definition they are insensitive to $\theta$*

$$\gamma_i = \left( \frac{1}{w_i} \frac{\partial w_i}{\partial \theta} \right) = 0 \,, \qquad \text{if } i \in \{\text{Background}\}$$

$$\gamma_i = \left( \frac{1}{w_i} \frac{\partial w_i}{\partial \theta} \right) \in \{-\infty, +\infty\} \,, \qquad \text{if } i \in \{\text{Signal}\}$$

Signal events may have sensitivity $\gamma_i>0$, $\gamma_i=0$ or $\gamma_i<0$
(special case: cross-section fit $\gamma_i=1/\sigma_s$)

*For what concerns
statistical errors in a parameter fit,*
***there is no distinction between
background events and
signal events with low sensitivity ($|\gamma_i|\sim0$)***

Bin-by-bin sensitivity $\phi_k$
of signal events alone:

$$\phi_k = \langle \gamma \rangle_{k,\text{Sig}} = \frac{1}{s_k} \sum_{i \in k}^{(\text{Sig})} \gamma_i = \frac{1}{s_k} \frac{\partial s_k}{\partial \theta}$$

***Bin-by-bin purity $\rho_k \leq 1$:***

$$\delta_i = \begin{cases} 1 & \text{if } i \in \{\text{Signal}\} \\ 0 & \text{if } i \in \{\text{Background}\} \end{cases} \qquad \rho_k = \frac{s_k}{s_k + b_k} = \frac{s_k}{n_k} = \frac{\sum_{i \in k} \delta_i}{n_k} = \langle \delta \rangle_k$$

Bin-by-bin sensitivity $\langle\gamma\rangle_k$
of signal + background:

$$\langle \gamma \rangle_k = \frac{1}{n_k} \frac{\partial n_k}{\partial \theta} = \frac{\rho_k}{s_k} \frac{\partial s_k}{\partial \theta} = \boxed{\rho_k} \phi_k$$

Information from all bins
for signal + background:

$$\mathcal{I}_\theta = \sum_{k=1}^{K} n_k \langle \gamma \rangle_k^2 = \sum_{k=1}^{K} n_k (\rho_k \phi_k)^2 = \sum_{k=1}^{K} s_k \boxed{\rho_k} \phi_k^2$$

**Effect of background:
it dilutes by a factor $\rho_k \leq 1$
the bin-by-bin
sensitivity and information
for signal events alone**

Information $I_\theta$ in terms of average bin-by-bin sensitivities:

$$\mathcal{I}_\theta = \sum_{k=1}^{K} n_k \left( \frac{1}{n_k} \frac{\partial n_k}{\partial \theta} \right)^2 = \sum_{k=1}^{K} n_k \langle \gamma \rangle_k^2$$

There is an **information gain** in partitioning two events $i_1$ and $i_2$ in two 1-event bins rather than one 2-event bin if their sensitivities $\gamma_{i_1}$ and $\gamma_{i_2}$ are different

$$\Delta\mathcal{I}_\theta = \gamma_{i_1}^2 + \gamma_{i_2}^2 - 2\left( \frac{\gamma_{i_1} + \gamma_{i_2}}{2} \right)^2 = \frac{1}{2}(\gamma_{i_1} - \gamma_{i_2})^2$$

**Goal of a distribution fit: partition events by their different MC-truth event-by-event sensitivities $\gamma_i$ to θ**

*How to achieve this in practice: next two slides (WDR)*

*Use $I_\theta^{(ideal)}$ to compute FIP: following two slides*

**Knowing one's limits: maximum achievable information with an ideal detector**
- Ideal acceptance, select all signal events $S_{sel}=S_{tot}$
- Ideal resolution, measured $\gamma_i$ is that from MC truth *(implies ideal rejection of background events, $\gamma_i=0$)*

$$\mathcal{I}_\theta^{(\text{ideal})} = \sum_{i=1}^{N_{tot}} \gamma_i^2 = \sum_{i=1}^{S_{tot}} \gamma_i^2$$

# Weight Derivative Regression (WDR): train $q_i$ for $\gamma_i$

Goal of a distribution fit: separate events with different MC-truth event-by-event sensitivities $\gamma_i$ to θ

***But $\gamma_i$ is not observable on real data events!***

**Weight Derivative Regression:**
**train a regressor $q_i = q(x_i)$**
**on detector-level MC observables $x_i$** $\Longrightarrow$
**against the MC-truth $\gamma_i = \partial w_i / \partial \theta$**
for signal and background MC events

Then determine θ
by the 1-D fit of $q(x_i)$
for real data events $x_i$

*Some of many caveats:*
*- Dependency of weight derivative on reference $\theta_{ref}$:*
  *WDR easier for coupling fits than for mass fits?*
*- How feasible is it to compute and store MC-truth weight derivatives?*
*- How useful is this for measurements limited by systematics?*
*- Train q on signal + background and 1-D fit of q, or*
  *train q on signal alone and 2-D fit on q and scoring classifier?*
*- How to deal with simultaneous fits of many parameters?*

*Training metric: maximize FIP*
*Evaluation metric: maximize FIP*

(or equivalently minimize MSE? see final slides)
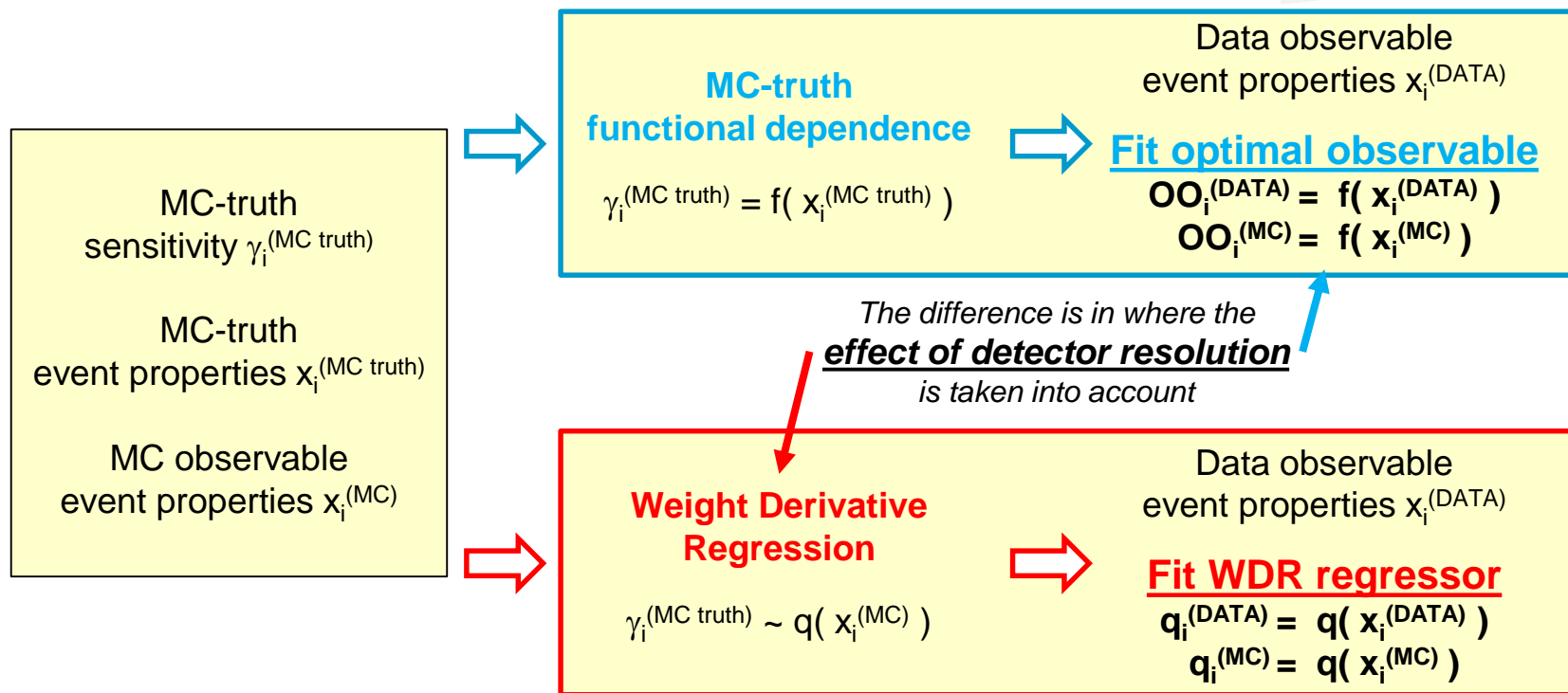
# WDR and Optimal Observables

The WDR idea was inspired by the
**Optimal Observables (OO) method**

*Both OO and WDR partition data by an approximation of a MC-truth sensitivity $\gamma_i$ to θ*
(OO does not use MC weight derivatives but it is similar)

D. Atwood, A. Soni, *Analysis for magnetic moment and electric dipole moment form factors of the top quark via $e^+e^- \to t\bar{t}$*, Phys. Rev. D 45 (1992) 2405. doi:10.1103/PhysRevD.45.2405,

M. Davier, L. Duflot, F. LeDiberder, A. Rougé, *The optimal method for the measurement of tau polarization*, Phys. Lett. B 306 (1993) 411. doi:10.1016/0370-2693(93)90101-M

M. Diehl, O. Nachtmann, *Optimal observables for the measurement of three-gauge-boson couplings in $e^+e^- \to W^+W^-$*, Z. Phys. C 62 (1994) 397. doi:10.1007/BF01555899

O. Nachtmann, F. Nagel, *Optimal observables and phase-space ambiguities*, Eur. Phys. J. C40 (2005) 497. doi:10.1140/epjc/s2005-02153-9

*Like OO, WDR can be useful in coupling/EFT fits*
*(more than in mass fits)*

*Some similarities also with the MadMiner approach See CHEP 2019 contribution #506 "Constraining effective field theories with ML"*

**MC-truth sensitivity $\gamma_i^{(\text{MC truth})}$**

**MC-truth event properties $x_i^{(\text{MC truth})}$**

**MC observable event properties $x_i^{(\text{MC})}$**

**MC-truth functional dependence**

$\gamma_i^{(\text{MC truth})} = f( x_i^{(\text{MC truth})} )$

Data observable event properties $x_i^{(\text{DATA})}$

**Fit optimal observable**
$OO_i^{(\text{DATA})} = f( x_i^{(\text{DATA})} )$
$OO_i^{(\text{MC})} = f( x_i^{(\text{MC})} )$

*The difference is in where the*
**effect of detector resolution**
*is taken into account*

**Weight Derivative Regression**

$\gamma_i^{(\text{MC truth})} \sim q( x_i^{(\text{MC})} )$

Data observable event properties $x_i^{(\text{DATA})}$

**Fit WDR regressor**
$q_i^{(\text{DATA})} = q( x_i^{(\text{DATA})} )$
$q_i^{(\text{MC})} = q( x_i^{(\text{MC})} )$

# FIP decomposition: efficiency, sharpness, purity

Numerator: Information retained by a given analysis using $N_{sel} = \Sigma n_k$ events with the given detector

Denominator: maximum theoretically available information from the given sample of $N_{tot}$ events ($S_{tot}$ signal events) if the true $\gamma_i$ were known for each event (ideal detector)

$$\mathrm{FIP}_3 = \frac{\mathcal{I}_\theta}{\mathcal{I}_\theta^{(\mathrm{ideal})}} = \frac{\sum_{k=1}^{K} n_k \langle\gamma\rangle_k^2}{\sum_{i=1}^{S_{tot}} \gamma_i^2} = \frac{\sum_{k=1}^{K} s_k \rho_k \phi_k^2}{\sum_{i=1}^{S_{tot}} \gamma_i^2}$$

$$\mathrm{FIP}_3 = \frac{\sum_{k=1}^{K} s_k \rho_k \phi_k^2}{\sum_{i=1}^{S_{tot}} \gamma_i^2} = \mathrm{FIP}_{\mathrm{eff}} \times \mathrm{FIP}_{\mathrm{sha}} \times \mathrm{FIP}_{\mathrm{pur}}$$

$$= \frac{\sum_{i=1}^{S_{sel}} \gamma_i^2}{\sum_{i=1}^{S_{tot}} \gamma_i^2} \times \frac{\sum_{k=1}^{K} s_k \phi_k^2}{\sum_{i=1}^{S_{sel}} \gamma_i^2} \times \frac{\sum_{k=1}^{K} s_k \rho_k \phi_k^2}{\sum_{k=1}^{K} s_k \phi_k^2}$$

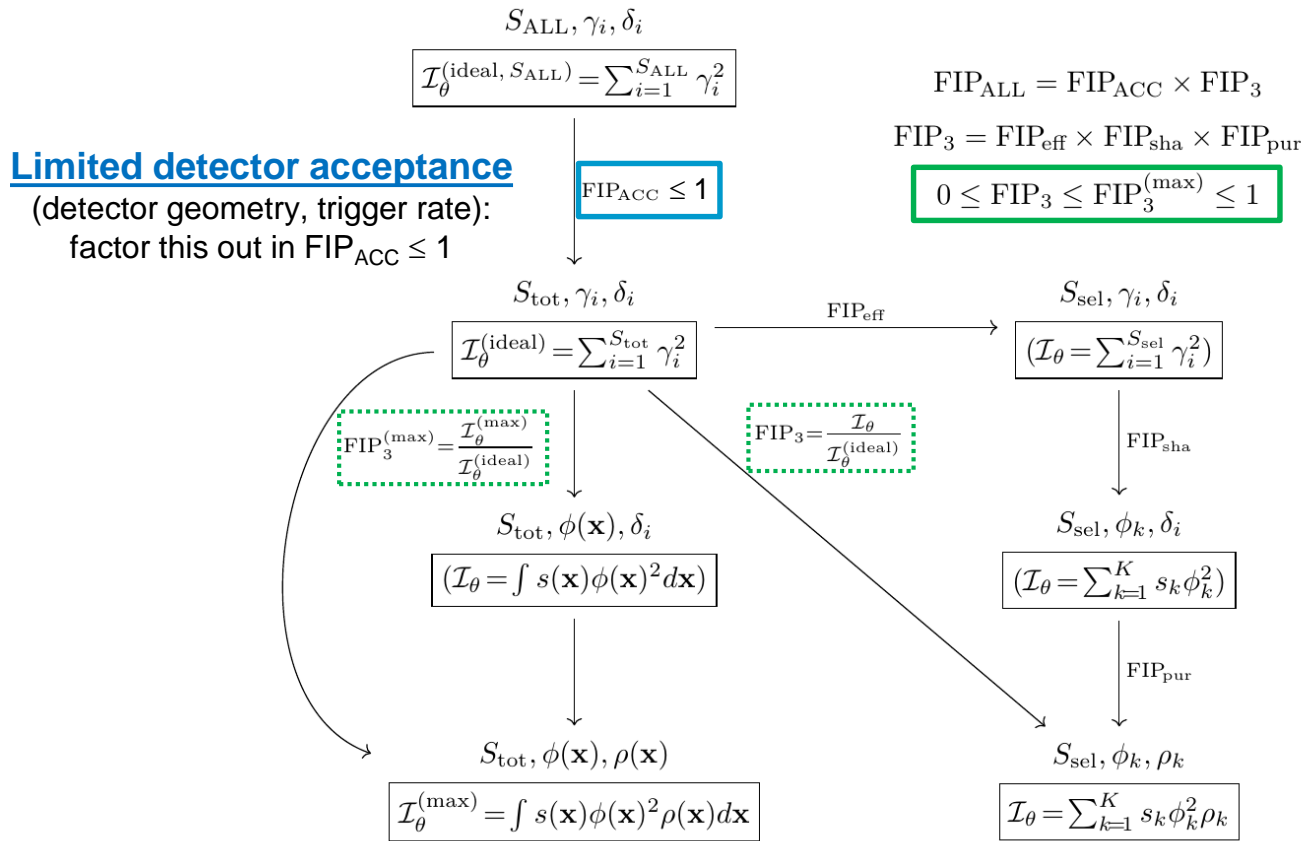Sensitivity-weighted signal **efficiency**: keep $S_{sel}$ of $S_{tot}$ events

**Sharpness** in separating signal events with different sensitivities: partition $S_{sel}$ signal events into K bins

Sensitivity-weighted signal **purity** or equivalently **sharpness** in separating signal events from background events: dilution of signal sensitivity caused by bin-by-bin purity $\rho_k$

*"sharpness" as in meteorology: see later why*

# Limits to knowledge: FIP for a realistic detector

$$S_{\mathrm{ALL}}, \gamma_i, \delta_i$$

$$\boxed{\mathcal{I}_\theta^{(\mathrm{ideal}, S_{\mathrm{ALL}})} = \sum_{i=1}^{S_{\mathrm{ALL}}} \gamma_i^2}$$

$$\mathrm{FIP}_{\mathrm{ALL}} = \mathrm{FIP}_{\mathrm{ACC}} \times \mathrm{FIP}_3$$

$$\mathrm{FIP}_3 = \mathrm{FIP}_{\mathrm{eff}} \times \mathrm{FIP}_{\mathrm{sha}} \times \mathrm{FIP}_{\mathrm{pur}}$$

$$\boxed{0 \le \mathrm{FIP}_3 \le \mathrm{FIP}_3^{(\mathrm{max})} \le 1}$$

**Limited detector acceptance**
(detector geometry, trigger rate):
factor this out in $\mathrm{FIP}_{\mathrm{ACC}} \le 1$

$$\boxed{\mathrm{FIP}_{\mathrm{ACC}} \le 1}$$

**Limited detector resolution**
In the multi-dimensional space
of event observables **x**,
**it is impossible to resolve**:

- signal events
with high sensitivity $\gamma_i$
from signal events
with low sensitivity $\gamma_i$:
average sensitivity is $\phi(\mathbf{x})$

- signal events $\delta_i = 1$
from background events $\delta_i = 0$:
average purity is $\rho(\mathbf{x})$

$$S_{\mathrm{tot}}, \gamma_i, \delta_i$$

$$\boxed{\mathcal{I}_\theta^{(\mathrm{ideal})} = \sum_{i=1}^{S_{\mathrm{tot}}} \gamma_i^2}$$

$$\xrightarrow{\mathrm{FIP}_{\mathrm{eff}}}$$

$$S_{\mathrm{sel}}, \gamma_i, \delta_i$$

$$\left( \mathcal{I}_\theta = \sum_{i=1}^{S_{\mathrm{sel}}} \gamma_i^2 \right)$$

$$\mathrm{FIP}_3^{(\mathrm{max})} = \frac{\mathcal{I}_\theta^{(\mathrm{max})}}{\mathcal{I}_\theta^{(\mathrm{ideal})}}$$

$$\mathrm{FIP}_3 = \frac{\mathcal{I}_\theta}{\mathcal{I}_\theta^{(\mathrm{ideal})}}$$

$$\Big\downarrow \mathrm{FIP}_{\mathrm{sha}}$$

$$S_{\mathrm{tot}}, \phi(\mathbf{x}), \delta_i$$

$$\left( \mathcal{I}_\theta = \int s(\mathbf{x}) \phi(\mathbf{x})^2 d\mathbf{x} \right)$$

$$S_{\mathrm{sel}}, \phi_k, \delta_i$$

$$\left( \mathcal{I}_\theta = \sum_{k=1}^{K} s_k \phi_k^2 \right)$$

$$\Big\downarrow \mathrm{FIP}_{\mathrm{pur}}$$

$$S_{\mathrm{tot}}, \phi(\mathbf{x}), \rho(\mathbf{x})$$

$$\boxed{\mathcal{I}_\theta^{(\mathrm{max})} = \int s(\mathbf{x}) \phi(\mathbf{x})^2 \rho(\mathbf{x}) d\mathbf{x}}$$

$$S_{\mathrm{sel}}, \phi_k, \rho_k$$

$$\boxed{\mathcal{I}_\theta = \sum_{k=1}^{K} s_k \phi_k^2 \rho_k}$$

> **FIP is a metric in [0,1], but**
> **the detector acceptance and resolution**
> **limit it to $0 \le \mathrm{FIP} \le \mathrm{FIP}^{(\mathrm{max})} < 1$**

⇨ *FIP>FIP(max) while training $q_i$*
*implies **overtraining**…*

Reading Room, British Museum
Diliff (own work, unmodified) CC BY 2.5

**Reading is a revolutionary act**
**(Inge Feltrinelli, 1930-2018)**

*Different problems in different domains require different metrics and tools…*

# Evaluating the evaluation metrics

Evaluation metrics of (binary and non-binary) classifiers
have been analysed and compared in many ways

There are two approaches which I find particularly useful:

1. Studying the <u>symmetries</u> and <u>invariances</u> of evaluation metrics

M. Sokolova, G. Lapalme, *A Systematic Analysis of Performance Measures for Classification Tasks*, Information Processing and Management 45 (2009) 427. doi:10.1016/j.ipm.2009.03.002

A. Luque, A Carrasco, A. Martin, J. R. Lama, *Exploring Symmetry of Binary Classification Performance Metrics*, Symmetry 11 (2019) 47. doi:10.3390/sym11010047.

*Example: (ir)relevance of True Negatives:
in my CHEP2018 talk*

2. Separating <u>threshold</u>, <u>ranking</u> and <u>probabilistic</u> metrics

R. Caruana, A. Niculescu-Mizil, *Data mining in metric space: an empirical analysis of supervised learning performance criteria*, Proc. 10th Int. Conf. on Knowledge Discovery and Data Mining (KDD-04), Seattle (2004). doi:10.1145/1014052.1014063

C. Ferri, J. Hernández-Orallo, R. Modroiu, *An Experimental Comparison of Classification Performance Metrics*, Proc. Learning 2004, Elche (2004). http://dmip.webs.upv.es/papers/Learning2004.pdf

C. Ferri, J. Hernández-Orallo, R. Modroiu, *An Experimental Comparison of Performance Measures for Classification*, Pattern Recognition Letters 30 (2009) 27. doi:10.1016/j.patrec.2008.08.010

*Example: AUC (ranking) vs. MSE (probabilistic):
in this CHEP2019 talk (next 3 slides)*

MSE (mean squared error) of regressor prediction $q_i$ versus the true $\gamma_i$ for event $i$:

$$\text{MSE} = \frac{1}{N_{\text{tot}}} \sum_{i=1}^{N_{\text{tot}}} (q_i - \gamma_i)^2$$

MSE is a probabilistic metric for both evaluation and training

MSE decomposition
(if the $N_{\text{tot}}$ events are split into K *partitions*, with $q_i = q_{(k)}\ \forall i \in k$):

***Paraphrases the "Brier score" decomposition in Meteorology***

G. W. Brier, *Verification of forecasts expressed in terms of probability*, Weather Rev. 78 (1950) 1. doi:10.1175/1520-0493(1950)078%3C0001:VOFEIT%3E2.0.CO;2
F. Sanders, *On Subjective Probability Forecasting*, J. Applied Meteorology 2 (1963) 191. https://www.jstor.org/stable/26169573

Validity, Reliability, Calibration

Sharpness, Resolution, Refinement

$$\text{MSE} = \frac{1}{N_{\text{tot}}} \left[ \sum_{k=1}^{K} n_k \left( q_{(k)} - \langle\gamma\rangle_k \right)^2 \right] + \frac{1}{N_{\text{tot}}} \left[ \left( \sum_{i=1}^{N_{\text{tot}}} \gamma_i^2 \right) - \left( \sum_{k=1}^{K} n_k \langle\gamma\rangle_k^2 \right) \right]$$

Validity: in a partition with given true average sensitivity $\langle\gamma_k\rangle$, is the predicted sensitivity $q_{(k)}$ well calibrated?

~0 in training by construction
~0 in evaluation if there are no systematics

**Sharpness: *how well do we separate events with different true sensitivities $\gamma_i$?***

***This is what determines the statistical error on the measurement of $\theta$: related to FIP!***

# FIP is related to Sharpness (MSE)

(Validity, Reliability, Calibration)  **$MSE_{sha}$** (Sharpness, Resolution, Refinement)

$$\text{MSE} = \frac{1}{N_{\text{tot}}} \left[ \sum_{k=1}^{K} n_k \left( q_{(k)} - \langle \gamma \rangle_k \right)^2 \right] + \frac{1}{N_{\text{tot}}} \left[ \left( \sum_{i=1}^{N_{\text{tot}}} \gamma_i^2 \right) - \left( \sum_{k=1}^{K} n_k \langle \gamma \rangle_k^2 \right) \right]$$

$$\frac{1}{N_{\text{tot}}} [\mathcal{I}_\theta^{(\text{ideal})} - \mathcal{I}_\theta] \qquad \mathcal{I}_\theta^{(\text{ideal})} = \sum_{i=1}^{N_{\text{tot}}} \gamma_i^2 = \sum_{i=1}^{S_{\text{tot}}} \gamma_i^2 \qquad \mathcal{I}_\theta = \sum_{k=1}^{K} n_k \left( \frac{1}{n_k} \frac{\partial n_k}{\partial \theta} \right)^2 = \sum_{k=1}^{K} n_k \langle \gamma \rangle_k^2$$

**FIP is related to Sharpness:**

In the ideal case: $MSE_{sha}=0$ and $FIP=1$
(events with different $\gamma_i$ can be resolved)

$$\text{FIP} = \frac{\mathcal{I}_\theta}{\mathcal{I}_\theta^{(\text{ideal})}} = \left( 1 - \frac{N_{\text{tot}} \times \text{MSE}_{\text{sha}}}{\mathcal{I}_\theta^{(\text{ideal})}} \right)$$

*Practical implication for Weight Derivative Regression:*
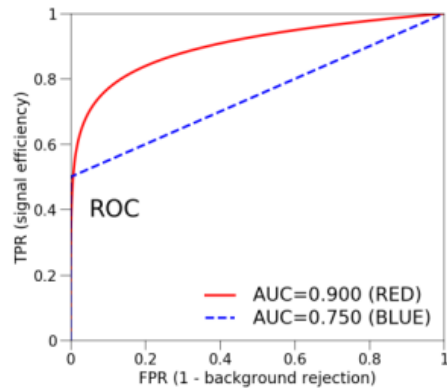*MSE is the most appropriate loss function for training the WDR regressor*

# HEP needs partitioning, and probabilistic metrics

## Ranking, and ranking metrics
Pick two events at random and rank them

**Medical Diagnostics** → *ranking evaluation of diagnostic prediction*
Patient A is diagnosed as more likely sick than B: how often am I right?

TPR (signal efficiency) vs FPR (1 - background rejection)
ROC
AUC=0.900 (RED)
AUC=0.750 (BLUE)

D. M. Green, *General Prediction Relating Yes-No and Forced-Choice Results*, J. Acoustical Soc. Am. 36 (1964) 1042. doi:10.1121/1.2143339
D. J. Goodenough, K. Rossmann, L. B. Lusted, *Radiographic applications of signal detection theory*, Radiology 105 (1972) 199. doi:10.1148/105.1.199
J. A. Hanley, B. J. McNeil, *The meaning and use of the area under a receiver operating characteristic (ROC) curve*, Radiology 143 (1982) 29. doi:10.1148/radiology.143.1.7063747
A. P. Bradley, *The use of the area under the ROC curve in the evaluation of Machine Learning algorithms*, Pattern Recognition 30 (1997) 1145. doi:10.1016/S0031-3203(96)00142-2

*AUC (Area Under the ROC Curve): probability that a randomly chosen diseased subject is correctly rated or ranked with greater suspicion than a randomly chosen non-diseased subject*

**IRRELEVANT FOR HEP PARAMETER FITS?**

## Partitioning, and probabilistic metrics
Group events and make a forecast on each subset

**Meteorology** → *probabilistic evaluation of weather prediction*
Rain forecast was 30% for these 10 days: actual rainy days?

**Medical Prognostics** → *probabilistic evaluation of survival prediction*
5yr survival forecast was 90% for these 10 patients: actual survivors?

**HEP parameter fits** → *probabilistic evaluation of measurement of $\theta$*
MC forecast for #events in this bin is 10 (20) for $\theta$=1 (2): actual data?

Validity, Reliability, Calibration    Sharpness, Resolution, Refinement

$$\text{MSE} = \frac{1}{N_{\text{tot}}} \left[ \sum_{k=1}^{K} n_k \left( q_{(k)} - \langle \gamma \rangle_k \right)^2 \right] + \frac{1}{N_{\text{tot}}} \left[ \left( \sum_{i=1}^{N_{\text{tot}}} \gamma_i^2 \right) - \left( \sum_{k=1}^{K} n_k \langle \gamma \rangle_k^2 \right) \right]$$

*Sharpness (from MSE): how well can I resolve days with 10% and 90% chance of rain? Patients with 10% and 90% 5yr survival rate? Signal events with high sensitivity to $\theta$ from (signal or background) events with low sensitivity?*

**ESSENTIAL FOR HEP PARAMETER FITS!**

# Conclusions – HEP measurement of a parameter $\theta$

- **MC weight derivatives** (event-by-event sensitivities $\gamma_i$ to $\theta$) may be used :
  - To determine the **ideal partitioning strategy**: partition by $\gamma_i$
  - To derive the **minimum error on the measurement of $\theta$** (ideal detector)

  $$\mathcal{I}_\theta^{(\text{ideal})} = \sum_{i=1}^{N_{\text{tot}}} \gamma_i^2 = \sum_{i=1}^{S_{\text{tot}}} \gamma_i^2$$

  - To derive **training and validation metrics** to optimize the measurement

  $$\text{FIP} = \frac{\mathcal{I}_\theta}{\mathcal{I}_\theta^{(\text{ideal})}} = \frac{\sum_{k=1}^{K} n_k \langle\gamma\rangle_k^2}{\sum_{i=1}^{S_{\text{tot}}} \gamma_i^2} = \frac{\sum_{k=1}^{K} s_k \rho_k \phi_k^2}{\sum_{i=1}^{S_{\text{tot}}} \gamma_i^2}$$

  **Evaluation and training metrics: FIP**

  - To train a **regressor $q_i$ of $\gamma_i$ (optimal observable)** for a 1-D fit of $\theta$

- HEP parameter fits are closer to **Meteorology** than to Medical Diagnostics
  - They use **partitioning** and need **probabilistic metrics** (sharpness, MSE)

  $$\text{FIP} = \frac{\mathcal{I}_\theta}{\mathcal{I}_\theta^{(\text{ideal})}} = \left( 1 - \frac{N_{\text{tot}} \times \text{MSE}_{\text{sha}}}{\mathcal{I}_\theta^{(\text{ideal})}} \right)$$

  **Compare to and learn from other domains**

  - They do not use ranking and do not need ranking metrics (AUC)

# Backup slides

# Non-dichotomous truth: examples

- **Medical Diagnostics** → *continuous scale gold standard*
  - The Obuchowski measure, e.g. five stages of liver fibrosis,

  N. A. Obuchowski, *An ROC-Type Measure of Diagnostic Accuracy When the Gold Standard is Continuous-Scale*, Statistics in Medicine 25 (2006) 481. doi:10.1002/sim.2228
  M. J. Pencina, R. B. D'Agostino, *Overall C as a measure of discrimination in survival analysis: model specific population value and confidence interval estimation*, Statistics in Medicine 23 (2004) 2109. doi:10.1002/sim.1802
  J. Lambert et al., *How to Measure the Diagnostic Accuracy of Noninvasive Liver Fibrosis Indices: The Area Under the ROC Curve Revisited*, Clinical Chemistry 54 (2008) 1372. doi:10.1373/clinchem.2007.097923

- **Information Retrieval** → *graded relevance assessment and DCG*
  - Discounted Cumulated Gain $\mathrm{DCG}[k] = \sum_{i=1}^{k} \frac{\mathrm{G}[i]}{\min(1, \log_2 i)}$
    *Response: partitioning + ranking*

  K. Järvelin, J. Kekäläinen, *IR evaluation methods for retrieving highly relevant documents*, Proc. 23rd ACM SIGIR Conf. (SIGIR 2000), Athens (2000). doi:10.1145/345508.345545
  J. Kekäläinen, K. Järvelin, *Using graded relevance assessments in IR evaluation*, J. Am. Soc. Inf. Sci. 53 (2002) 1120. doi:10.1002/asi.10137
  K. Järvelin, J. Kekäläinen, *Cumulated gain-based evaluation of IR techniques*, J. ACM Trans. on Inf. Sys. (TOIS) 20 (2002) 422. doi:10.1145/582415.582418

- **ML (for finance)** → *example-dependent cost-sensitive classification*
  - Payoff matrix for transaction x\$:
    *Response: yes/no decision*

    |         | fraudulent | legitimate |
    |---------|------------|------------|
    | refuse  | \$20       | −\$20      |
    | approve | −$x$       | $0.02x$    |

  B. Zadrozny, C. Elkan, *Learning and making decisions when costs and probabilities are both unknown*, Proc. 7th Int. Conf. on Knowledge Discovery and Data Mining (KDD-01), San Francisco (2001). doi:10.1145/502512.502540
  C. Elkan, *The Foundations of Cost-Sensitive Learning*, Proc. 17th Int. Joint Conf. on Artificial Intelligence (IJCAI-01), Seattle (2001).

- **Meteorology** → *probabilistic evaluation of weather forecasts*
  - Rain forecast was 30% for these 10 days: actual rainy days?

  G. W. Brier, *Verification of forecasts expressed in terms of probability*, Weather Rev. 78 (1950) 1. doi:10.1175/1520-0493(1950)078%3C0001:VOFEIT%3E2.0.CO;2
  F. Sanders, *On Subjective Probability Forecasting*, J. Applied Meteorology 2 (1963) 191. https://www.jstor.org/stable/26169573

- **Medical Prognostics** → *probabilistic evaluation of survival forecasts*
  - 5yr survival forecast was 90% for these 10 patients: actual survivors?

  **HEP-like: probabilistic!**

  D. J. Spiegelhalter, *Probabilistic prediction in patient management and clinical trials*, Statist. Med. 5 (1986) 421. doi:10.1002/sim.4780050506
  F. E. Harrell, K. L. Lee, D. B. Mark, *Multivariable prognostic models: issues in developing models, evaluating assumptions and adequacy, and measuring and reducing errors*, Statist. Med. 15 (1996) 361. 10.1002/(SICI)1097-0258(19960229)15:4<361::AID-SIM168>3.0.CO;2-4

- **HEP measurement of** θ → *evt-by-evt sensitivity to θ*

# Weight Derivative Regression – in practice

- *Compute event-by-event sensitivities $\gamma_i$ from signal MC weight derivatives*
  - Possibly at various reference values of θ

- *Pre-select events to remove most backgrounds*
  - Possibly maximizing a sensitivity-weighted signal efficiency?

- ***Train a regressor $q_i$ for the MC-truth $\gamma_i$ from measured event properties***
  - Possibly using MSE as the loss function in the training (see next slides)

- *Determine θ from a 1-D fit on the optimal observable $q_i$*
  - Or possibly a 2-D fit on $(q_i, D_i)$ including the pre-selection classifier $D_i$

Some of the many **limitations of this approach**
  - MC weight derivative depend on θ: coupling fits easier than mass fits
  - I ignored systematic errors
  - I only discussed fits of a single physics parameter at a time
  - *But I still find this approach better than maximizing an AUC…*

(Note: I did not try a real measurement – I did a few tests with a toy model, but I am not presenting them today)

# Weight derivative regressors and their training
## *(a frequentist dinosaur's view of Machine Learning)*

https://openclipart.org

Classic ML problem: create a model $q(\mathbf{x})=R_\gamma(\mathbf{x})$ to predict the value of $\gamma(\mathbf{x})$ in a multi-dimensional space of variables $\mathbf{x}$

Choosing a ML methodology
mainly implies two choices:

**1. The shape of the function $R_\gamma(x)$:**
i.e. how we choose to model $\gamma(\mathbf{x})$
*Examples: decision tree (sparsely uniform),
neural network (sigmoids), linear discriminant*

*I focus on **Decision Trees**
because of the similarities
to binned distribution fits*

**2. The training metric:** a "distance"
of $R_\gamma(\mathbf{x}_i)$ to $\gamma(\mathbf{x}_i)$ or $\gamma_i$ to minimize, or
a property of $R_\gamma(\mathbf{x}_i)$ to maximize
*Examples: Gini, Shannon entropy/information, MSE*

*I suggest to use $I_\theta$ **or FIP** both
for training and for evaluation*

# Event selection and partitioning: a blurred boundary

(1) The scoring classifier D for signal/background discrimination is related to the average purity $\rho(\mathbf{x})$: it would be a pity to use it only for a yes/no decision

*It can be used both for measuring cross-sections (1-D fit of D) or for measuring a mass or coupling (2-D fit against another variable)*

**Use the scoring classifier D to partition events, not only to accept or reject events**

(2) Signal events with zero or low sensitivity to θ and background events are equally irrelevant

*(As far as statistical errors are concerned)*

**Separating signal events with high sensitivity to θ from background events**

***is as important as***

**Separating signal events with high sensitivity to θ from signal events with low sensitivity to θ**

# Fisher information (about a parameter θ)

- **Fisher information $I_\theta$** is a useful concept because
  - 1. It refers to the parameter θ that is being measured
  - 2. It is additive: the information from independent measurements adds up
  - 3. The higher the information $I_\theta$, the lower the error $\Delta\theta$ achievable on θ

F. James, *Statistical Methods in Experimental Physics*, 2nd edition, World Scientific (2006).

**Cramer-Rao lower bound CRLB**
(lowest achievable variance $\Delta\theta^2$)

$$(\Delta\hat{\theta})^2 = \mathrm{var}(\hat{\theta}) \geq \frac{1}{\mathcal{I}_\theta}$$

- Some estimators achieve the CRLB and are called efficient
  - Example: a maximum likelihood fit (given the event counts in a given partitioning scheme)

- In the following *I will express statistical error $\Delta\theta$ in terms of information $I_\theta$*

*i.e. I will treat errors $\Delta\theta$ and information $I_\theta$ as equivalent concepts*

$$\mathcal{I}_\theta = \frac{1}{(\Delta\theta)^2}$$

# HEP cross-section in a counting experiment

- Measurement of a total cross-section $\sigma_s$ in a counting experiment

- A distribution fit with a single bin

- Well-known since decades if final goal is to minimize statistical error $\Delta\sigma_s$
  - *Maximise $\varepsilon_s*\rho$* ("common knowledge" in the LEP2 experiments) → *"FIP1"*
  - NB: This metric only makes sense for this specific HEP optimization problem!

$$\mathcal{I}_{\sigma_s} = \frac{1}{(\Delta\sigma_s)^2} = \frac{1}{\sigma_s^2}\,\epsilon_s\varrho S_{\text{tot}} = \frac{1}{\sigma_s^2}\left(\frac{S_{\text{sel}}^2}{S_{\text{sel}}+B_{\text{sel}}}\right)$$

$$\mathcal{I}_{\sigma_s}^{(\text{ideal})} = \frac{S_{\text{tot}}}{\sigma_s^2}\,,\text{if } \varrho=1 \text{ and } \epsilon_s=1$$

$$\Longrightarrow \quad \boxed{\text{FIP}_1 = \frac{\mathcal{I}_{\sigma_s}}{\mathcal{I}_{\sigma_s}^{(\text{ideal})}} = \epsilon_s\varrho}$$

By the way: $\rho/\varepsilon_s=1$ where $\partial\text{FIP1}/\partial\rho=\partial\text{FIP1}/\partial\varepsilon_s$ (just like for F1)

# A brief comparison of MD, IR and HEP

- **Medical Diagnostics**
  - *All patients are important, both truly ill (TP) and truly healthy (TN)*
  - e.g. ACC metric depends on all four categories: average over TP+TN+FP+FN

$$ACC = \frac{TP + TN}{TP + TN + FP + FN}$$

- **Information Retrieval**
  - Based on *qualitative distinction between "relevant" and "non relevant" documents*
  - e.g. F1 metric *does not depend on True Negatives*
    - Rejected "irrelevant" documents are utterly irrelevant

$$F_1 = \frac{2\,TP}{2\,TP + FP + FN}$$

- **HEP (cross section measurement by counting)**
  - Based on *qualitative distinction between signal and background*
  - e.g. FIP1 metric *does not depend on True Negatives*
    - Measured cross section cannot depend on how many background events are rejected

$$FIP_1 = \frac{TP^2}{(TP+FN)(TP+FP)}$$

**_HEP is more similar to Information Retrieval than to Medical Diagnostics_**
*(qualitative asymmetry between positives and negatives)*

*Invariance under TN change is only one of many useful symmetries to analyse*
*[Sokolova-Lapalme, Luque et al.]*

M. Sokolova, G. Lapalme, *A Systematic Analysis of Performance Measures for Classification Tasks*, Information Processing and Management 45 (2009) 427. doi:10.1016/j.ipm.2009.03.002

A. Luque, A Carrasco, A. Martin, J. R. Lama, *Exploring Symmetry of Binary Classification Performance Metrics*, Symmetry 11 (2019) 47. doi:10.3390/sym11010047.

# HEP: cross section in a counting experiment
## (maximize FIP1 – the AUC is misleading!)

To minimize the statistical error Δσ:
**Maximize** $\boxed{\mathrm{FIP}_1 = \epsilon_s \varrho}$

**Choice of operating point** is simple:
- Plot $\epsilon_s \times \rho$ as a function of $\epsilon_s$
- Choose the point where $\epsilon_s \times \rho$ is maximum

**Choice between two classifiers** is simple:
- Determine max $(\epsilon_s \times \rho)$ for each
- Choose the classifier with the higher max

*NB1: The choice depends on prevalence*
[which is fixed by physics and approximately known in advance]

*NB2: AUC is misleading and irrelevant in this case*

*But there are better ways
than a counting experiment
to measure a total cross section
in this case…*



|  | FIP1 | AUC |
|---|---|---|
| **Range in [0,1]** | YES | YES |
| **Higher is better** | YES | NO |
| **Numerically meanigful** | YES | NO |

# HEP: cross section by a fit to the score distribution

**Use the scoring classifier D to partition events, not to accept or reject events**

This is the most common method to measure a total cross section (example: a BDT or NN output fit)

Keep all Stot events and partition them in K bins

$$\text{FIP}_2 = \frac{\mathcal{I}_{\sigma_s}}{\mathcal{I}_{\sigma_s}^{(\text{ideal})}} = \frac{\sum_k s_k \rho_k}{\sum_k s_k} = \frac{\sum_k s_k^2 / n_k}{\sum_k s_k} = \frac{\sum_k n_k \rho_k^2}{\sum_k s_k}$$

There is a benefit in partitioning events into subsets with different purities because

$$\Delta \mathcal{I}_{\sigma_s} = \frac{n_1 n_2}{n_1 + n_2} (\rho_1 - \rho_2)^2$$

Better than a counting experiment for two reasons
- All events are used, none are rejected
- Those which were previously in a single bin are now subpartitioned

# FIP2 from the ROC (+prevalence) or from the PRC

- From the previous slide: $\mathrm{FIP2} = \dfrac{\sum_{i=1}^{m} \rho_i s_i}{\sum_{i=1}^{m} s_i}$

FIP2: integrals on ROC and PRC,
more relevant to HEP than AUC or AUCPR!
(well-defined meaning for distribution fits)

- FIP2 from the ROC (+prevalence $\pi_s = \dfrac{S_{\mathrm{tot}}}{S_{\mathrm{tot}} + B_{\mathrm{tot}}}$ ):

$$S_{\mathrm{sel}} = S_{\mathrm{tot}}\,\epsilon_s \qquad s_i = dS_{\mathrm{sel}} = S_{\mathrm{tot}}\,d\epsilon_s \qquad \rho_i = \dfrac{1}{1 + \dfrac{B_{\mathrm{tot}}}{S_{\mathrm{tot}}}\dfrac{d\epsilon_b}{d\epsilon_s}} \qquad \mathrm{FIP2} = \int_0^1 \dfrac{d\epsilon_s}{1 + \dfrac{1-\pi_s}{\pi_s}\dfrac{d\epsilon_b}{d\epsilon_s}}$$

$$B_{\mathrm{sel}} = B_{\mathrm{tot}}\,\epsilon_b \qquad b_i = dB_{\mathrm{sel}} = B_{\mathrm{tot}}\,d\epsilon_b$$

*Compare FIP2(ROC) to AUC*

$$\mathrm{AUC} = \int_0^1 \epsilon_s\, d\epsilon_b = 1 - \int_0^1 \epsilon_b\, d\epsilon_s$$
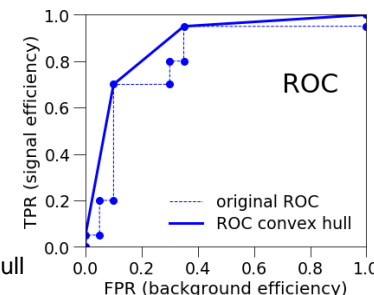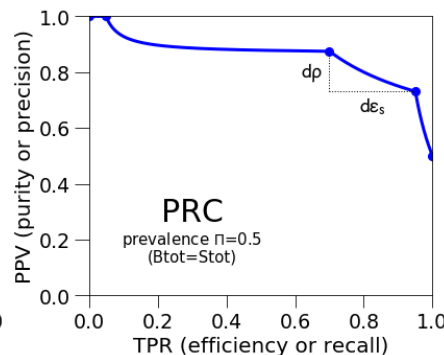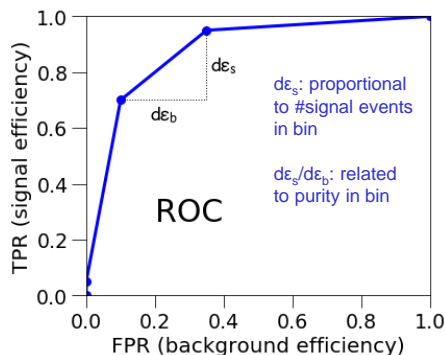
- FIP2 from the PRC:

$$S_{\mathrm{sel}} = S_{\mathrm{tot}}\,\epsilon_s \qquad s_i = dS_{\mathrm{sel}} = S_{\mathrm{tot}}\,d\epsilon_s$$

$$B_{\mathrm{sel}} = S_{\mathrm{sel}}\left(\dfrac{1}{\rho} - 1\right) \qquad b_i = dB_{\mathrm{sel}} = S_{\mathrm{tot}}\left[d\epsilon_s\left(\dfrac{1}{\rho} - 1\right) - \epsilon_s\dfrac{d\rho}{\rho^2}\right] \qquad \rho_i = \dfrac{\rho}{1 - \dfrac{\epsilon_s}{\rho}\dfrac{d\rho}{d\epsilon_s}} \qquad \mathrm{FIP2} = \int_0^1 \dfrac{\rho\, d\epsilon_s}{1 - \dfrac{\epsilon_s}{\rho}\dfrac{d\rho}{d\epsilon_s}}$$

*Compare FIP2(PRC) to AUCPR*

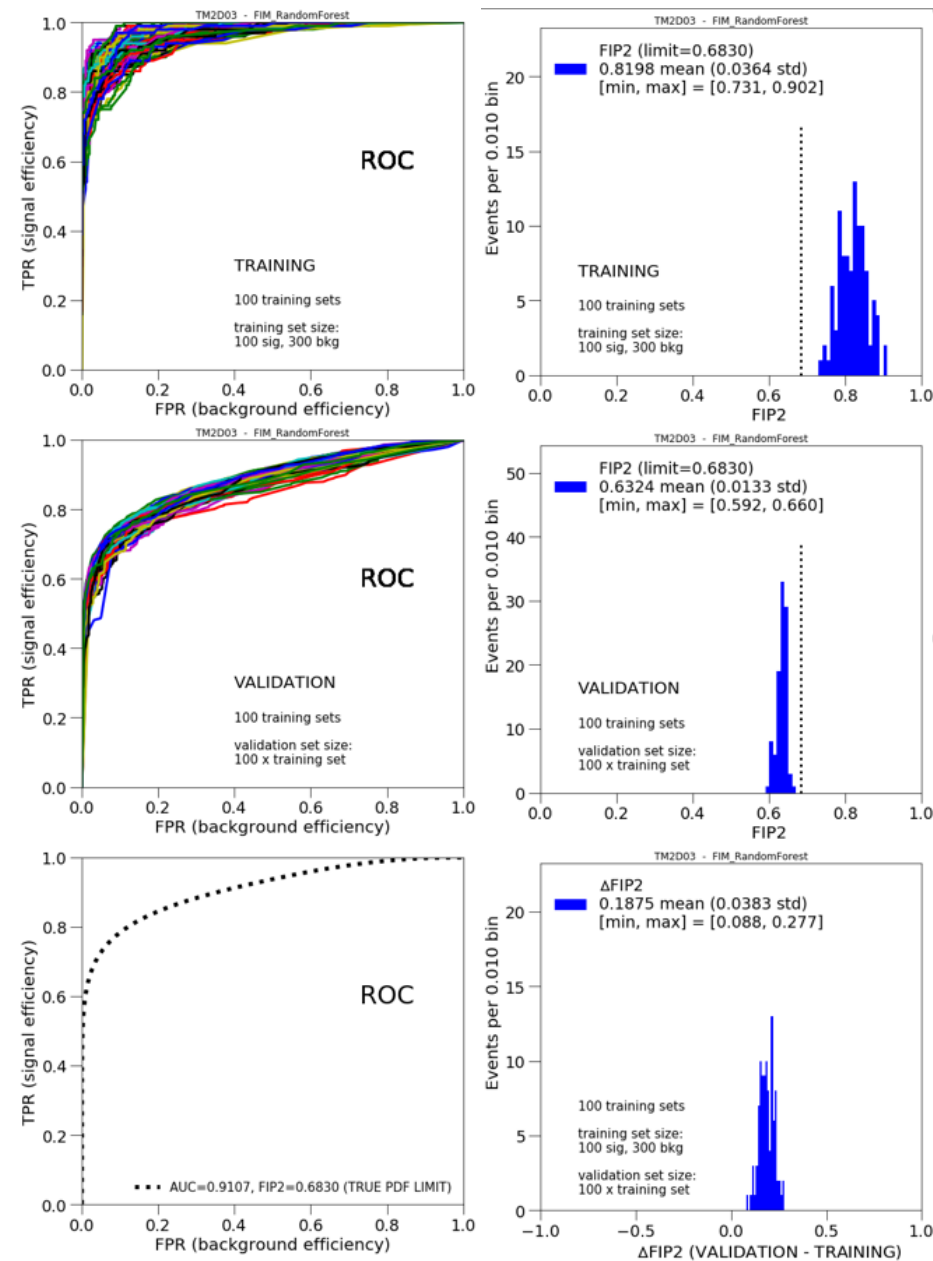$$\mathrm{AUCPR} = \int_0^1 \rho\, d\epsilon_s$$

- Easier calculation and interpretation from ROC (+prevalence) than from PRC
  - *region of constant ROC slope = region of constant signal purity*
  - decreasing ROC slope = decreasing purity
    - technicality (my Python code): convert ROC to convex hull* first

$d\epsilon_s$: proportional to #signal events in bin

$d\epsilon_s/d\epsilon_b$: related to purity in bin

ROC

PRC
prevalence Π=0.5
(Btot=Stot)

ROC

original ROC
ROC convex hull

*Convert ROC to convex hull
- ensure decreasing slope
- avoid staircase effect that would artificially inflate FIP2
  (bins of 100% purity: only signal or only background)

# FIP2$^{(max)}$ example
## (and overtraining)



**FIP2 is a metric in [0,1]
but the detector resolution
effectively determines a FIP2$^{(max)}$ < 1**

# Fisher information $I_\theta$ about θ (statistical errors)

For a given partitioning scheme with K bins
($n_k$ is the number of selected events in bin k)

**Bin-by-bin sensitivity to θ**

$$\mathcal{I}_\theta = \frac{1}{(\Delta\theta)^2} = \sum_{k=1}^{K} \frac{1}{(\Delta\theta)_k^2} = \sum_{k=1}^{K} n_k \left( \frac{1}{n_k} \frac{\partial n_k}{\partial \theta} \right)^2$$
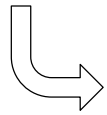
Statistical errors: information adds up
Each bin is an independent measurement with error $(\Delta\theta)_k = \left(\frac{\partial n_k}{\partial\theta}\right)^{-1} \Delta n_k = \left(\frac{\partial n_k}{\partial\theta}\right)^{-1} \sqrt{n_k}$

(Combination more complex with systematic errors, or for searches)

# Optimal partitioning

$$\mathcal{I}_\theta = \frac{1}{(\Delta\theta)^2} = \sum_{k=1}^{K} \frac{1}{(\Delta\theta)_k^2} = \sum_{k=1}^{K} n_k \left( \frac{1}{n_k} \frac{\partial n_k}{\partial \theta} \right)^2$$

Is there a benefit (information inflow) in *splitting bin 0 into two bins 1, 2 with $n_0 = n_1 + n_2$?*

$$\Delta\mathcal{I}_\theta = \frac{1}{n_1} \left( \frac{\partial n_1}{\partial \theta} \right)^2 + \frac{1}{n_2} \left( \frac{\partial n_2}{\partial \theta} \right)^2 - \frac{1}{n_1 + n_2} \left( \frac{\partial (n_1 + n_2)}{\partial \theta} \right)^2$$

$$= \frac{n_1 n_2}{n_1 + n_2} \left[ \left( \frac{1}{n_1} \frac{\partial n_1}{\partial \theta} \right) - \left( \frac{1}{n_2} \frac{\partial n_2}{\partial \theta} \right) \right]^2$$

Information increases if the two new bins have different sensitivities to θ

$$\Delta\mathcal{I}_\theta > 0 \iff \left( \frac{1}{n_1} \frac{\partial n_1}{\partial \theta} \right) \neq \left( \frac{1}{n_2} \frac{\partial n_2}{\partial \theta} \right)$$

**Goal of a distribution fit: partition events into subsets with different bin-by-bin sensitivities to θ**

# Signal and background are not dichotomous classes
## (with one exception: cross section measurements)

Background events by definition are insensitive to θ
Signal events may have positive, zero or negative sensitivity

θ: mass, coupling
NON-DICHOTOMOUS

$$\gamma_i = \left(\frac{1}{w_i}\frac{\partial w_i}{\partial \theta}\right) = 0, \qquad \text{if } i \in \{\text{Background}\}$$

$$\gamma_i = \left(\frac{1}{w_i}\frac{\partial w_i}{\partial \theta}\right) \in \{-\infty, +\infty\}, \qquad \text{if } i \in \{\text{Signal}\}$$

*The distinction between
signal events with low ($|\gamma_i|{\sim}0$) sensitivity
and background events is blurred*
(example: events far from an invariant mass peak)

$$\delta_i = \begin{cases} 1 & \text{if } i \in \{\text{Signal}\} \\ 0 & \text{if } i \in \{\text{Background}\} \end{cases}$$

Changing the signal cross section ~is a
global rescaling of all differential distributions

$$s_k(\sigma_s) = \frac{\sigma_s}{\sigma_{s,\text{ref}}} \times s_k(\sigma_{s,\text{ref}})$$

In a cross section measurement
All background events are equivalent to one another
All signal events are equivalent to one another

$$\gamma_i = \frac{1}{\sigma_s}\delta_i = \begin{cases} \frac{1}{\sigma_s} & \text{if } i \in \{\text{Signal}\}, \\ 0 & \text{if } i \in \{\text{Background}\}, \end{cases} \qquad \text{if } \theta \equiv \sigma_s$$

θ: cross section σ_s
DICHOTOMOUS

# FIP1 and FIP2 revisited

FIP$_{sha}$=1 for both
(dichotomous, all signal events are equivalent)

$$\text{FIP}_3 = \frac{\sum_{k=1}^{K} s_k \rho_k \phi_k^2}{\sum_{i=1}^{S_{\text{tot}}} \gamma_i^2} = \text{FIP}_{\text{eff}} \times \text{FIP}_{\text{sha}} \times \text{FIP}_{\text{pur}}$$

$$= \frac{\sum_{i=1}^{S_{\text{sel}}} \gamma_i^2}{\sum_{i=1}^{S_{\text{tot}}} \gamma_i^2} \times \frac{\sum_{k=1}^{K} s_k \phi_k^2}{\sum_{i=1}^{S_{\text{sel}}} \gamma_i^2} \times \frac{\sum_{k=1}^{K} s_k \rho_k \phi_k^2}{\sum_{k=1}^{K} s_k \phi_k^2}$$

$$\text{FIP}_1 = \epsilon_s \varrho$$

FIP1:
FIP$_{\text{eff}}$ =ε
FIP$_{\text{pur}}$=ρ

$$\text{FIP}_2 = \frac{\mathcal{I}_{\sigma_s}}{\mathcal{I}_{\sigma_s}^{(\text{ideal})}} = \frac{\sum_k s_k \rho_k}{\sum_k s_k} = \frac{\sum_k s_k^2/n_k}{\sum_k s_k} = \frac{\sum_k n_k \rho_k^2}{\sum_k s_k}$$

FIP2:
FIP$_{\text{eff}}$ =1
FIP$_{\text{pur}}$=FIP2

$$\text{FIP}_3 = \frac{\sum_{k=1}^{K} s_k \rho_k \phi_k^2}{\sum_{i=1}^{S_{\text{tot}}} \gamma_i^2} = \text{FIP}_{\text{eff}} \times \text{FIP}_{\text{sha}} \times \text{FIP}_{\text{pur}}$$

$$= \frac{\sum_{i=1}^{S_{\text{sel}}} \gamma_i^2}{\sum_{i=1}^{S_{\text{tot}}} \gamma_i^2} \times \frac{\sum_{k=1}^{K} s_k \phi_k^2}{\sum_{i=1}^{S_{\text{sel}}} \gamma_i^2} \times \frac{\sum_{k=1}^{K} s_k \rho_k \phi_k^2}{\sum_{k=1}^{K} s_k \phi_k^2}$$

# From CRLB to Fisher Information Part (FIP)

| Particles produced in beam collisions |
| Raw data events |
| Analysis object data |
| Event counts in individual bins of a distribution |
| Measured value of the parameter M ± ΔM |

**Detector Trigger**

**Data processing**

*PHYSICS ANALYSIS*
**Event selection (Sig vs Bkg)**
**Event partitioning**

**Max likelihood fit**

$$\mathcal{I}_\theta^{(\text{ideal})}$$

$$\text{FIP} = \frac{\mathcal{I}_\theta}{\mathcal{I}_\theta^{(\text{ideal})}} = \frac{(\Delta\theta^{(\text{ideal})})^2}{(\Delta\theta)^2} \leq 100\%$$

**FIP**

$$\mathcal{I}_\theta$$

$(1/\Delta\theta^2)/\mathcal{I}_\theta$ is 100%

A max likelihood fit is 100% efficient: it achieves the CRLB, *for the given event selection and event partitioning*

$$1/\Delta\theta^2$$

**CRLB**

# Two optimization handles: event selection and partitioning