



MACHINE LEARNING PDFs

STEFANO FORTE UNIVERSITÀ DI MILANO & INFN



UNIVERSITÀ DEGLI STUDI DI MILANO DIPARTIMENTO DI FISICA



XXVII EPIPHANY CONFERENCE



This project has received funding from the European Union's Horizon 2020 research and innovation programme under grant agreement No 740006

MOTIVATIONS I: ACCURACY PDF4LHC PDFs (2014) NNPDF3.0 NNLO



- GLUON BETTER KNOWN AT SMALL x, VALENCE QUARKS AT LARGE x, SEA QUARKS IN BETWEEN
- TYPICAL UNCERTAINTIES IN DATA REGION $\sim 3-5\%$
- SWEET SPOT: VALENCE Q G; DOWN TO 1%
- UP BETTER KNOWN THAN DOWN; FLAVOR SINGLET BETTER THAN INDIVIDUAL FLAVORS

MOTIVATIONS I: ACCURACY CURRENT PDFs (2017) NNPDF3.1 NNLO



- GLUON BETTER KNOWN AT SMALL x, VALENCE QUARKS AT LARGE x, SEA QUARKS IN BETWEEN
- TYPICAL UNCERTAINTIES IN DATA REGION $\sim 1-3\%$
- SWEET SPOT: VALENCE Q G; 1% OR BELOW
- UP BETTER KNOWN THAN DOWN; FLAVOR SINGLET BETTER THAN INDIVIDUAL FLAVORS



- CT18: PDF SETS RELEASED WITH/WITHOUT ATLAS W/Z DATA INCLUDED
- NNPDF3.1: CONSISTENCY OF ALL DATASETS INCLUDED

EQUALLY PRECISE BUT MORE ACCURATE RESULT!

- CENTRAL VALUE MOVES TOWARDS KNOWN NNLO
- RELATIVE ERROR ϕ ON PREDICTION MILDLY INCREASED
- FIT QUALITY χ^2 IMPROVES



THE WAY AHEAD



ARTIFICIAL INTELLIGENCE: PARADIGMS

"KNOWLEDGE BASED" AI

- LEARN AND IMPLEMENT A SET OF RULES
- GOOD FOR CHESS, **BAD** FOR REAL LIFE



 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0
 0

MACHINE LEARNING

- "INTUITIVE" REPRESENTATION
- THE AI AGENT BUILID UP ITS OWN KNOWLEDGE



ARTIFICIAL INTELLIGENCE: ALGORITHMS



EXTRACT AND OPTIMIZE DATA FEATURES

OPTIMIZE A PROPERTY LEARNING FROM DATA LEARN FROM DATA THE LEARNING STRATEGY

PROTON STRUCTURE AS AN AI PROBLEM: NNPDF



MONTECARLO + NEURAL NETWORKS

THE NNPDF APPROACH COMBINING DATA BY MONTE CARLO

TWO MEASUREMENTS: $\mu_1 \pm \sigma_1$; $\mu_2 \pm \sigma_2$ **MC COMBINATION**: $\bar{\mu} \pm \bar{\sigma}$; $\bar{\mu} = \frac{\frac{\mu_1}{\sigma_1^2} + \frac{\mu_2}{\sigma_2^2}}{\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}}$; $\bar{\sigma}^2 = \frac{1}{\frac{1}{\sigma_1^2} + \frac{1}{\sigma_2^2}}$

MONTE CARLO REPRESENTATION



 $\mu^{(i)} \Leftrightarrow \operatorname{REPLICA} \operatorname{SAMPLE} \Leftrightarrow \operatorname{REPRESENTATION} \operatorname{OF} \operatorname{PROBABILITY} \operatorname{DISTRIBUTION} \operatorname{NEED} \operatorname{ONLY} \operatorname{TO} \operatorname{KNOW} \operatorname{HOW} \operatorname{TO} \operatorname{COMBINE} \operatorname{CENTRAL} \operatorname{VALUES}$

THE NNPDF APPROACH THE FUNCTIONAL MONTE CARLO

REPLICA SAMPLE OF FUNCTIONS ⇔ PROBABILITY DENSITY IN FUNCTION SPACE KNOWLEDGE OF FUNCTIONAL FORM NOT NECESSARY



FINAL PDF SET: $f_i^{(a)}(x,\mu)$; i =up, antiup, down, antidown, strange, antistrange, charm, gluon; $j = 1, 2, ... N_{\text{rep}}$

ARTIFICIAL INTELLIGENCE NEURAL NETWORKS

output layer

ARCHITECTURE



- WEIGHTS ω_{ij}
- THRESHOLDS θ_i



$$F_{\rm out}^{(i)}(\vec{x}_{\rm in}) = F\left(\sum_{j} \omega_{ij} x_{\rm in}^{j} - \theta_{i}\right)$$

SIMPLEST EXAMPLE 1-2-1

 $f(x) = \frac{1}{\substack{\theta_1^{(3)} - \frac{\omega_{11}^{(2)}}{1 + e^{\theta_1^{(2)} - x\omega_{11}^{(1)}} - \frac{\omega_{12}^{(2)}}{1 + e^{\theta_2^{(2)} - x\omega_{21}^{(1)}}}}$

NNPDF: 2-5-3-1 NN for each PDF: $37 \times 8 = 296$ parameters

NEURAL LEARNING

- COMPLEXITY INCREASES AS THE FITTING PROCEEDS
- UNTIL LEARNING NOISE
- WHEN SHOULD ONE STOP?



NEURAL LEARNING

- COMPLEXITY INCREASES AS THE FITTING PROCEEDS
- UNTIL LEARNING NOISE
- WHEN SHOULD ONE STOP?



NEURAL LEARNING

- COMPLEXITY INCREASES AS THE FITTING PROCEEDS
- UNTIL LEARNING NOISE
- WHEN SHOULD ONE STOP?



GENETIC MINIMIZATION: AT EACH GENERATION, χ^2 EITHER UNCHANGED OR DECREASING

- DIVIDE THE DATA IN TWO SETS: TRAINING AND VALIDATION
- MINIMIZE THE χ^2 OF THE DATA IN THE TRAINING SET
- AT EACH ITERATION, COMPUTE THE χ^2 FOR THE DATA IN THE VALIDATION SET (NOT USED FOR FITTING)
- WHEN THE VALIDATION χ^2 STOPS DECREASING, STOP THE FIT



GENETIC MINIMIZATION: AT EACH GENERATION, χ^2 EITHER UNCHANGED OR DECREASING

- DIVIDE THE DATA IN TWO SETS: TRAINING AND VALIDATION
- MINIMIZE THE χ^2 OF THE DATA IN THE TRAINING SET
- AT EACH ITERATION, COMPUTE THE χ^2 FOR THE DATA IN THE VALIDATION SET (NOT USED FOR FITTING)
- WHEN THE VALIDATION χ^2 STOPS DECREASING, STOP THE FIT



GENETIC MINIMIZATION: AT EACH GENERATION, χ^2 EITHER UNCHANGED OR DECREASING

- DIVIDE THE DATA IN TWO SETS: TRAINING AND VALIDATION
- MINIMIZE THE χ^2 OF THE DATA IN THE TRAINING SET
- AT EACH ITERATION, COMPUTE THE χ^2 FOR THE DATA IN THE VALIDATION SET (NOT USED FOR FITTING)
- WHEN THE VALIDATION χ^2 STOPS DECREASING, STOP THE FIT





GENETIC MINIMIZATION: AT EACH GENERATION, χ^2 EITHER UNCHANGED OR DECREASING

- DIVIDE THE DATA IN TWO SETS: TRAINING AND VALIDATION
- MINIMIZE THE χ^2 OF THE DATA IN THE TRAINING SET
- AT EACH ITERATION, COMPUTE THE χ^2 FOR THE DATA IN THE VALIDATION SET (NOT USED FOR FITTING)
- WHEN THE VALIDATION χ^2 STOPS DECREASING, STOP THE FIT

TOO LATE!









- THE METHODOLOGY IS FAITHFUL
- BUT IS IT OPTIMAL?

ML: UNSUPERVISED LEARNING OPTIMIZATION I

- HOW TO MAXIMIZE ACCURACY?
- LARGE (PRIOR) REPLICA SET
- GENETIC SELECTION \Rightarrow OPTIMIZATION OF STATISTICAL INDICATORS (KULLBACK-LEIBLER DIVERGENCE)
- 50 optimizes replicas \Leftrightarrow 1000 starting replicas





- SIGNIFICANT DEPENDENCE ON NUMBER OF REPLICAS
- Asymptotic "tolerance" $T=1.3\pm0.3; \ \Delta\chi^2=1.7\pm0.7$
- For $N_{\rm rep} = 100$, T = 2.3, even for $N_{\rm rep} = 1000$, T = 1.6

DO WE HAVE TO FIT 10000 REPLICAS? DO WE HAVE TO USE 10000 REPLICAS?

ML: SUPERVISED LEARNING **OPTIMIZATION II**

- CAN WE REDUCE THE NUMBER OF COMPRESSED REPLICAS WITHOUT LOSS OF INFORMATION? SOLUTION FOR USER
- CAN WE INCREASE THE NUMBER OF REPLICAS WITHOUT REFITTING? • SOLUTION FOR PDF FITTER



GENERATIVE ADVERSARIAL NETWORKS

- TRAIN A NETWORK TO SIMULATE THE TRUE DISTRIBUTION (GENERATOR)
- TRAIN A NETWORK TO DISCRIMINATE TRUTH FROM SIMULATION (DISCRIMINATOR)
- TRAIN THE GENERATOR TO TRICK THE DISCRIMINATOR



 ● 1D GAN: REPRODUCE THE INFORMATION IN THE UNDERLYING REPLICA SET, BUT NO GAIN (WIGGLY REPLICAS)
 ⇒ REDUCE THE NUMBER OF COMPRESSED REPLICA WITH FIXED NUMBER OF FITTED REPLICAS W/O INFORMATION LOSS



- 2D GAN: COMBINE CORRELATED INFORMATION FROM UNDERLYING REPLICA SET Machine Learning * PDFs * QCD INFERRING THE TRUE UNDERLYING DISTTRIBUTION
 - \Rightarrow REDUCE THE NUMBER OF INPUT REPLICAS W/O INFORMATION LOSS





CLOSURE TEST: A CLOSER LOOK (NNPDF3.1)

ONE σ : ACTUAL/PREDICTED

FOR DATA, BY EXPERIMENT

	NNPDF3.1 ratio
experiment	
NMC	0.882828
SLAC	0.767063
BCDMS	0.730569
CHORUS	0.698907
NTVDMN	0.991090
HERACOMB	0.847359
HERAF2CHARM	1.867597
F2BOTTOM	1.124157
DYE886	0.655955
DYE605	0.585725
CDF	0.961652
D0	0.881199
ATLAS	0.904127
CMS	1.090241
LHCb	1.092194
Total	0.842168

ONE σ VALUE FOR PDFS, VS x



- UNCERTAINTIES OVERESTIMATED
- 1 σ >68% at very small and very large x; 1 σ <68% at intermediate x

FITTING THE METHODOLOGY

THE N3FIT PROJECT



HOW DO WE KNOW THAT THE METHODOLOGY IS THE BEST? "ACCUMULATED WISDOM" INEFFICIENT AND SLOW

CHANGE OF PHILOSOPHY \Rightarrow DETERMINISTIC MINIMIZATION (GRADIENT DESCENT) GO FOR THE ABSOLUTE MINIMUM, AND (HYPER)OPTIMIZE



- PYTHON-BASED KERAS + TENSORFLOW FRAMEWORK
- EACH BLOCK INDEPENDENT LAYER
- CAN VARY ALL ASPECT OF METHODOLOGY



- SCAN PARAMETER SPACE
- OPTIMIZE FIGURE OF MERIT: VALIDATION χ^2
- BAYESIAN UPDATING



• OVERFITTING $\Rightarrow \chi^2_{\text{train}} << \chi^2_{\text{valid}}$!! & WIGGLY PDFS

• CORRELATIONS BETWEEN DATA IN A SET

WHAT HAPPENED?



CROSS-VALIDATION SELECTS THE OPTIMAL MINIMUM

WHAT HAPPENED?

HYPEROPTIMIZATION



WE ARE MISSING A SELECTION CRITERION

MACHINE LEARNING THE SOLUTION

TUNED HYPEROPTIMIZATION



COMPARE TO A A TEST SET (NEW SET OF DATA PREVIOUSLY NOT USED AT AL) TESTS GENERALIZATION POWER

THE TEST SET METHOD

- COMPLETELY UNCORRELATED TEST SET
- OPTIMIZE ON WEIGHTED AVERAGE OF VALIDATION AND TEST \Rightarrow NO OVERLEARNING





- NO OVERFITTING
- COMPARED TO NNPDF3.1
 - MUCH GREATER STABILITY \Rightarrow FEWER REPLICAS FOR EQUAL ACCURACY
 - UNCERTAINTIES SOMEWHAT REDUCED

CLOSURE TESTS AGAIN

ONE σ : ACTUAL/PREDICTED

FOR DATA, BY EXPERIMENT

	NNPDF3.1 ratio	n3fit ratio			
experiment					
NMC	0.882828	0.843427			
SLAC	0.767063	0.690118			
BCDMS	0.730569	0.770704			
CHORUS	0.698907	0.734656			
NTVDMN	0.991090	0.797017			
HERACOMB	0.847359	1.326333			
HERAF2CHARM	1.867597	3.566076			
F2BOTTOM	1.124157	1.532634			
DYE886	0.655955	0.857915			
DYE605	0.585725	0.870151			
CDF	0.961652	0.779424			
D0	0.881199	1.015202			
ATLAS	0.904127	1.132229			
CMS	1.090241	1.017136			
LHCb	1.092194	0.993525			
Total	0.842168	0.940737			





- UNCERTAINTIES WELL ESTIMATED; BUT OVERESTIMATED FOR DIS
- ONE σ PERFECT IN DATA REGION; BUT UNDERESTIMATED IN EXTRAPOLATION

WHAT ARE UNCERTAINTIES WHEN THERE ARE NO DATA?

WHAT IS "PROPER LEARNING"? FORECASTING AN UNKNOWN TRUTH \Rightarrow WHAT IS "OPTIMAL"?

MENU V nature

NEWS • 08 JANUARY 2019

Machine learning leads mathematicians to unsolvable problem

Simple artificial-intelligence problem puts researchers up against a logical parc by famed mathematician Kurt Gödel.

Davide Castelvecchi

Ƴ **f** ■



SOME POSSIBLE ANSWERS/CRITERIA

- PASS A CLOSURE TEST
- PASS A "FUTURE TEST": GENERALIZE TO CURRENT DATA BASED ON PAST DATA
- REPRODUCE THE EXPECTED STATISTICAL PROPERTIES: ONE $\sigma \Leftrightarrow \Delta \chi^2 = 1$
- SATISFY THEORETICAL PREJUDICE?

REINFORCEMENT LEARNING?

THE WORK OF MANY PEOPLE



NNPDF collaboration and N³PDF team meeting, Varenna, Italy, September 2019 "Io stimo più il trovare un vero, benché di cosa leggiera, che il disputar lungamente delle massime questioni senza verità nissuna"

"I am more interested in uncovering a fact, however trifling, than to dispute at length about profound questions devoid of any truth"

Galileo Galilei, letter to Tommaso Campanella



CONTEMPORARY PDF TIMELINE (ONLY PUBLISHED GLOBAL)

	20	08	2009 20		10	2011 2012		2013		2014		2015 2017		17	2019		
SET MONTH	CTEQ6.6 (02)	NNPDF1.00	MSTW (01)	ABKM09 (08)	NNPDF2.00	CT10 (NLO) OT	NNPDF2.16 (NNLO)	ABM11 (02)	NNPDF2.30	CT10 (NNLO) 02	ABM 12 (10)	0 NNPDF3.0	MMHT (12)	CT14 06	ABMP16 O	NNPDF3.10	CT18 (12)
F. T. DIS ZEUS+H1-HI COMB. HI ZEUS+H1-HII HERA JETS	× × ×	× × ×	x x x	× × ×	2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2 2	v v x x	v v some	× × ×	2 2 2 2 2	× × some	× ×		× × ×	v v x x	****		
F. T. DY Tev W+Z LHC W+Z	× × × ×	× × ×	V V X	× × ×	~ ~ ~ ~	 × 	× × ×	× × ×		、 、 、 、 、	× × some		2 2 2	<	✓ × some	~ ~ ~ ~	~ ~ ~ ~
Tev jets LHC jets	× X	× ×	✓ ×	x x	× ×	✓ ×	x x	✓ ×	v v	* ×	× ×	v v	22	~ ~	× ×	~ ~	2 2
TOP TOTAL SINGLE TOP TOTAL TOP DIFFERENTIAL	× × ×	× × ×	× × ×	X X X	× × ×	× × ×	X X X	× × ×	× × ×	× × ×	✓ × ×	× ×	× × ×	× × ×	> > X	× ×	× ×
$ \begin{array}{c} \mathbb{W} \ p_T \\ \mathbb{W} + \mathbb{C} \\ \mathbb{Z} \ p_T \end{array} $	× × ×	× × ×	× × ×	× × ×	× × ×	× × ×	× × ×	× × ×	× × ×	× × ×	× × ×	~ ~ X	× × ×	× × ×	× × ×	× × •	× × •

THEORY PROGRESS:

- MSTW, ABKM: all NNLO; NNPDF NNLO since 07/11 (2.1), CT since 02/13 (CT10); NNPDF THRESHOLD RESUMMATION (3.0RESUM, 07/15), SMALL *x* RESUMMATION (3.1SX, 10/17)
- MSTW, CT, NNPDF all GM-VFN; NNPDF since 01/11 (2.1); ABM FFN+ZM-VFN since 01/17 (ABMP16)
- NNPDF FITTED CHARM since 05/16 (NNPDF3IC)
- PHOTON PDF: (mrst2004qed), NNPDF2.3QED (08/13), NNPDF3.0QED (06/16), NNPDF3.1LUXQED (12/17)



- GLUON BETTER KNOWN AT SMALL x, VALENCE QUARKS AT LARGE x, SEA QUARKS IN BETWEEN
- TYPICAL UNCERTAINTIES IN DATA REGION $\sim 3-5\%$
- SWEET SPOT: VALENCE Q G; DOWN TO 1%
- UP BETTER KNOWN THAN DOWN; FLAVOR SINGLET BETTER THAN INDIVIDUAL FLAVORS
- NO QUALITATIVE DIFFERENCE BETWEEN NLO AND NNLO

DATASET WIDENING NNPDF3.0 vs NNPDF3.1

Kinematic coverage



NEW DATA: (BLACK EDGE)

- HERA COMBINED F_2^b
- D0 W LEPTON ASYMMETRY
- ATLAS *W*, *Z* 2011, HIGH & LOW MASS DY 2011; CMS *W*[±] RAPIDITY 8TEV LHCB *W*, *Z* 7TEV & 8TEV
- ATLAS 7TEV JETS 2011, CMS 2.76TEV JETS
- ATLAS & CMS TOP DIFFERENTIAL RAPIDITY
- ATLAS $Z p_T$ DIFFERENTIAL RAPIDITY & INVARIANT MASS 8TEV, CMS $Z p_T$ DIFFERENTIAL

RAPIDITY 8TEV



- SIGNIFICANT UNCERTAINTY REDUCTION
- MANY PDFs CHANGE BY MORE THAN ONE SIGMA
- BOTH FLAVOR SEPARATION & GLUON SIGNIFICANTLY AFFECTED

DATA VS. THEORY/METHODOLOGY THE STRANGE PDF: DIS VS. W PRODUCTION

- STRANGE PDF CONTROLLED BY NEUTRINO DIS CHARM PRODUCTION + W PRODUCTION
- DIS DATA FAVOR "SUPPRESSED STRANGE" \Rightarrow SMALL $R_s \equiv \frac{s+\bar{s}}{\bar{u}+\bar{d}}$
- ATLAS FAVORS ENHANCED STRANGENESS
- ATLAS IMPACT EXAGGERATED IN XFITTER ANALYSIS
- EVERYTHING CONSISTENT WITHIN UNCERTAINTIES IN GLOBAL FIT



DATA VS. THEORY/METHODOLOGY THE STRANGE PDF: DIS VS. W PRODUCTION

- MASSIVE CORRECTIONS TO CHARGED CURRENT DIS HITERTO INCLUDED TO NLO MASSLESS TO NNLO
- Gao, $2018 \Rightarrow$ NNLO COMPUTED
- STRANGENESS ENHANCED BY NNLO CORRECTIONS





(Gao, 2108)

LESSONS:

- BEWARE OF XFITTER HERA+X FITS
- IN A GLOBAL FIT DIFFERENT DATA ALWAYS PULL IN DIFFERENT DIRECTIONS!
- TENSIONS CAN BE RESOLVED BY BETTER THEORY

DATA VS. THEORY/METHODOLOGY THE CHARM MASS AND TREATMENT CT18 \rightarrow CT18Z

- ATLAS W and Z 7TeV rapidity included
- CHARM MASS INCREASED
- x-dependent factorization scale



DATA VS. THEORY/METHODOLOGY THE CHARM MASS AND TREATMENT CHARM FROM DATA

• CHARM SHOULD NOT DEPEND STRONGLY ON CHARM MASS



• ITS SHAPE SHOULD NOT BE DETERMINED BY FIRST-ORDER MATCHING (NO HIGHER NONTRIVIAL ORDERS KNOWN)

• MIGHT EVEN HAVE A NONPERTURBATIVE COMPONENT

FITTED VS. PERTURBATIVE: SUPPRESSED AT MEDIUM-SMALL x, ENHANCED AT VERY SMALL, VERY LARGE x



- QUARK LUMI AFFECTED BECAUSE OF CHARM SUPPRESSION AT MEDIUM-x
- FLAVOR DECOMPOSITION ALTERED
- UNCERTAINTIES ON LIGHT QUARKS NOT SIGNIFICANTLY INCREASED
- AGREEMENT OF 13TeV W,Z PREDICTED CROSS-SECTIONS IMPROVES!



• W, Z CROSS-SECTIONS AT 13 TEV IN PERFECT AGREEMENT WITH DATA THANKS TO FITTED CHARM!

LESSONS:

- TENSIONS CAN REVEAL METHODOLOGICAL ISSUES
- MORE LIKELY AS DATASET INCREASES, EXPERIMENTAL UNCERTAINTIES DECREASE
- RESOLVED BY MORE COMPLEX METHODOLOGY

DATA vs. METHODOLOGY

- NEW DATA \Rightarrow MAJOR METHODOLOGICAL CHOICES \Rightarrow SIGNIFICANT IMPACT
- NNPDF3.1 vs NNPDF3.0: DATA AND METHODOLOGY HAVE SIMILAR IMPACT





- TO CONVERT HESSIAN INTO MONTECARLO GENERATE MULTIGAUSSIAN REPLICAS IN PARAMETER SPACE
- ACCURATE WHEN NUMBER OF REPLICAS SIMILAR TO THAT WHICH REPRODUCES DATA





(Carrazza, SF, Kassabov, Rojo, 2015)

- TO CONVERT MONTE CARLO INTO HESSIAN, SAMPLE THE REPLICAS $f_i(x)$ AT A DISCRETE SET OF POINTS & CONSTRUCT THE ENSUING COVARIANCE MATRIX
- EIGENVECTORS OF THE COVARIANCE MATRIX AS A BASIS IN THE VECTOR SPACE SPANNED BY THE REPLI-CAS BY SINGULAR-VALUE DECOMPOSITION
- NUMBER OF DOMINANT EIGENVECTORS SIMILAR TO NUMBER OF REPLICAS \Rightarrow ACCURATE REPRESENTATION

TOOLS II NONGAUSSIAN BEHAVIOUR

$\begin{array}{c} \text{MONTE CARLO COMPARED TO HESSIAN} \\ \text{CMS } W + c \text{ production} \end{array}$



- DEFINE KULLBACK-LEIBLER DIVERGENCE $D_{\text{KL}} = \int_{-\infty}^{\infty} P(x) \frac{\ln P(x)}{\ln Q(x)} dx$ BETWEEN A PRIOR P AND ITS REPRESENTATION Q
- $D_{\rm KL}$ between prior and hessian depends on degree of gaussianity
- $D_{\rm KL}$ between prior and compressed MC does not

- DEVIATION FROM GAUSSIANITY E.G. AT LARGE x DUE TO LARGE UNCERTAINTY + POSITIVITY BOUNDS \Rightarrow RELEVANT FOR SEARCHES
- CANNOT BE REPRODUCED IN HESSIAN FRAMEWORK
- Well reproduced by compressed MC



CAN (A) GAUGE WHEN MC IS MORE ADVANTAGEOUS THAN HESSIAN; (B) ASSESS THE ACCURACY OF COMPRESSION

TOOLS III OPTIMIZED PDFS: SMPDF

- OLD ASPIRATION: PDFs OPTIMIZED TO PROCESSES (Pumplin 2009)
- SELECT SUBSET OF THE COVARIANCE MATRIX CORRELATED TO A GIVEN SET OF PROCESSES
- PERFORM SVD ON THE REDUCED COVARIANCE MATRIX, SELECT DOMINANT EIGENVECTOR, PROJECT OUT ORTHOGONAL SUBSPACE
- ITERATE UNTIL DESIRED ACCURACY REACHED
- CAN ADD PROCESSES TO GIVEN SET; CAN COMBINE DIFFERENT OPTIMIZED SETS
- WEB INTERFACE AVAILABLE



w_etmiss_13tev(LO)

(Carrazza, SF, Kassabov, Rojo, 2016)

- EG ggH, $Hb\bar{b}$, $W E_T^{\text{miss}} \Rightarrow 11$ EIGENVECTORS
- STUDY CORRELATIONS OF PDFS TO DATA AND AMONG THEMSELVES!