



WHERE DO WE COME FROM? WHAT ARE WE? WHERE ARE WE GOING?

STEFANO FORTE

MILAN UNIVERSITY & INFN

FOR THE COLLABORATION: R. D. BALL, L. DEL DEBBIO,

S.F., A. GUFFANTI, J. I. LATORRE, J. ROJO, M. UBIALI



PDF4LHC & TH INSITITUTE



CERN, AUGUST 6, 2009

WHAT ARE WE?

WHAT DO WE DELIVER?

A MONTE CARLO SET OF REPLICAS OF PDFS

Example: the gluon distribution in the NNPDF1.0 set



- ENSEMBLE OF REPLICAS \leftrightarrow PROBABILITY DISTRIBUTION OF PDFs
- EXPECTED CENTRAL VALUE \leftrightarrow MEAN; UNCERTAINTY \leftrightarrow STANDARD DEVIATION
- ANY FEATURES OF DISTRIBUTION CAN BE DETERMINED (C.L. INTERVALS, CORRELATIONS...)

WHERE DO WE COME FROM?

WHAT DO WE START FROM?

A MONTE CARLO SET OF REPLICAS OF DATA

1

10

• Space of functions huge 5 bins for 10 pts \times 5 fctns \rightarrow $5^{50} \sim 10^{34}$ bins

• THE DATA TELL US WHICH ONES ARE POPULATED (IMPORTANCE SAMPLING)

replica averages

replica standard dev. vs. uncertainties







UNCERTAINTIES, 1000 FOR CORRELNS



WHERE DO WE COME FROM?

HOW DO WE GET PDF REPLICAS FROM DATA REPLICAS? NEURAL NETWORK PARM+ CROSS-VALIDATION METHOD

- Each PDF \leftrightarrow neural network parametrized by 37 parameters
- NNPDF1.2: $37 \otimes 7 = 259$ parms (COMP. MSTW08 $\rightarrow 28$ free parameters) "INFINITE" NUMBER OF PARAMETERS \Rightarrow CAN REPRESENT ANY FUNCTION
- COMPLEX SHAPES (LARGE NO.OF PARAMETERS) REQUIRE LONGER FITTING
- FIT STOPS WHEN QUALITY OF FIT TO RANDOMLY SELECTED "VALIDATION" DATA (NOT FITTED) STOPS IMPROVING
- CAN OBTAIN A FIT WITH χ^2 LOWER THAN BEST FIT ("OVERLEARNING")

WHERE DO WE COME FROM? HOW DO WE GET PDF REPLICAS FROM DATA REPLICAS? CROSS-VALIDATION

- OPTIMAL FIT OBTAINED WHEN QUALITY OF FIT TO VALIDATION (CONTROL) DATA STOPS IMPROVING
- POSSIBILITY OF OVERFITTING GUARANTESS THAT MINIMUM NOT DRIVEN BY PARAMETRIZATION



WHERE DO WE COME FROM? HOW DO WE GET PDF REPLICAS FROM DATA REPLICAS? CROSS-VALIDATION

- OPTIMAL FIT OBTAINED WHEN QUALITY OF FIT TO VALIDATION (CONTROL) DATA STOPS IMPROVING
- POSSIBILITY OF OVERFITTING GUARANTESS THAT MINIMUM NOT DRIVEN BY PARAMETRIZATION



OVERFITTING

WHERE DO WE COME FROM? WHY ARE WE DOING ALL THIS? ISSUES IN THE STANDARD APPROACH

BENCHMARK FITS (HERALHC 2005-2008)

- PERFORM A FIT TO A CONSISTENT SUBSET OF DATA, USE $\Delta\chi^2=1$
- RESULTS NOT CONSISTENT, BUT IMPROVE W/ MORE GENERAL PARM. (MSTW08 VS MRST01)



- UNCERTAINTY DOES NOT GO DOWN AS DATASET INCREASES
- MUST TUNE PARAMETRIZATION AND STATISTICAL TREATMENT (TOLERANCE) TO DATASET

WHO ARE WE? IS THIS ISSUE SOLVED?

MRST/MSTW: BENCH VS REF

NNPDF: BENCH VS REF

NNPDF BENCH VS MRST/MSTW BENCH



- SINGLE PARAMETRIZATION AND STAT. TREATMENT CAN ACCOMMODATE DIFFERENT DATASETS
- IMPACT OF DATA CAN BE STUDIED INDEPENDENT OF THEORETICAL FRAMEWORK

WHO ARE WE? WHAT IS THIS GOOD FOR? DETERMINATION OF WEAKLY CONSTRAINED QUANTITIES

THE STRANGE PDF

- NNPDF1.0: $s(x, Q_0^2) = \bar{s}(x, Q_0^2)$, $s + \bar{s} = \frac{1}{2}(\bar{u} + \bar{d})$, no dimuon data
- NNPDF1.1: s, \bar{s} (actually s^{\pm}) indep. parametrized, no dimuon data
- NNPDF1.2: s, \bar{s} (actually s^{\pm}) indep. parametrized, dimuon data



STRANGE PDFS

WHO ARE WE? WHAT IS THIS GOOD FOR? COMPATIBILITY CAN BE CHECKED QUANTITATIVELY

THE STRANGE PDF



NNPDF1.2 vs. NNPDF1.1		
	DATA	EXTRAPOLATION
$\Sigma(x,Q_0^2)$	$5 \ 10^{-4} \le x \le 0.1$	$10^{-5} \le x \le 10^{-4}$
$ \begin{array}{c} \langle d[q] \rangle \\ \langle d[\sigma] \rangle \end{array} $	2.7 3.1	1.2 1.8
$g(x,Q_0^2)$	$5 \ 10^{-4} \le x \le 0.1$	$10^{-5} \le x \le 10^{-4}$
$egin{array}{c} \langle d[q] angle \ \langle d[\sigma] angle \end{array} \ (d[\sigma]) \ \langle d[\sigma] angle$	2.4 1.3	2.0 1.4
$T_3(x, Q_0^2)$	$0.05 \le x \le 0.75$	$10^{-3} \le x \le 10^{-2}$
$egin{array}{c} \langle d[q] angle \ \langle d[\sigma] angle \end{array} egin{array}{c} \langle d[\sigma] angle \end{array}$	1.5 1.1	0.9 1.2
$V(x,Q_0^2)$	$0.1 \le x \le 0.6$	$3 \ 10^{-3} \le x \le 3 \ 10^{-2}$
$egin{array}{c} \langle d[q] angle \ \langle d[\sigma] angle \end{array}$	1.1 1.3	1.0 1.4
$\Delta_S(x,Q_0^2)$	$0.1 \le x \le 0.6$	$3 \ 10^{-3} \le x \le 3 \ 10^{-2}$
$ \begin{array}{ c c } & \langle d[q] \rangle \\ & \langle d[\sigma] \rangle \end{array} $	0.8 1.3	0.8 1.1
$s^+(x, Q_0^2)$	$5 \ 10^{-4} \le x \le 0.1$	$10^{-5} \le x \le 10^{-4}$
$ \begin{array}{ c c } & \langle d[q] \rangle \\ & \langle d[\sigma] \rangle \end{array} $	2.0 4.5	1.6 1.8
$s^{-}(x, Q_0^2)$	$0.1 \le x \le 0.6$	$3 \ 10^{-3} \le x \le 3 \ 10^{-2}$
$ \begin{array}{c} \langle d[q] \rangle \\ \langle d[\sigma] \rangle \end{array} $	1.1 6.1	1.3 4.6

DISTANCE OF EXP. VALUES RESCALED BY $\sigma/\sqrt{N_{rep}}$

WHO ARE WE? WHAT IS THIS GOOD FOR? CAN DETERMINE RELIABLY LARGE & ASYMMETRIC UNCERTAINTIES

THE STRANGE MOMENTUM FRACTIONS



AN IMPLICATION: THE "NUTEV ANOMALY" IS GONE

WHO ARE WE? WHAT IS THIS GOOD FOR? CAN DISENTANGLE PHYSICAL PARAMETERS FROM PDFS DETERMINATION OF CKM PARMS FROM DIS

$$F_{2}^{\nu,c} = x \left[C_{2,q} \otimes \left(|V_{cd}|^{2} (u+d) + 2|V_{cs}|^{2} s \right) + C_{2,g} \otimes g \right]$$
$$F_{2}^{\bar{\nu},c} = x \left[C_{2,q} \otimes \left(|V_{cd}|^{2} (\bar{u}+\bar{d}) + 2|V_{cs}|^{2} \bar{s} \right) + C_{2,g} \otimes g \right]$$

1.02			
	ANALYSIS	DETERMINATION	$ V_{CS} $
	NNPDF1.2	DIRECT FROM GLOBAL PDF ANALYSIS	$0.96 \pm 0.07^{\mathrm{tot}}$
-	CDHS	LO FROM $\nu N \rightarrow \mu^+ \mu^- X$	≥ 0.59 (90% C.L.)
CKM unit. fit	CCFR	NLO FROM $\nu N \rightarrow \mu^+ \mu^- X$	≥ 0.74 (90% C.L.)
	PDG08	Average from D decays	1.04 ± 0.06
	HOCKER	AVERAGE FROM $\nu N \rightarrow \mu^+ \mu^- X$	1.04 ± 0.16
	DELPHI	Direct from $W^+ \to c \bar{s}$ decays	$0.94 \frac{+0.32}{-0.26} \pm 0.13$
	PDG08	CKM LINITARITY FIT	0.97334 ± 0.00023
	IDddd		0.91334 ± 0.00023
0.94	ANALYSIS	DETERMINATION	V _{cd}
0.94	ANALYSIS NNPDF1.2	DETERMINATION DIRECT FROM GLOBAL PDF ANALYSIS	$\begin{array}{c c} \hline & & & \\ \hline \\ \hline$
0.94	ANALYSIS NNPDF1.2 CDHS	$\frac{\text{Determination}}{\text{Direct from global PDF analysis}}$ $\text{LO from } \nu N \rightarrow \mu^+ \mu^- X$	$ \begin{array}{c c} & V_{cd} \\ \hline & 0.244 \pm 0.019^{\text{tot}} \\ & 0.24 \pm 0.03 \end{array} $
	ANALYSIS NNPDF1.2 CDHS CCFR	$\frac{\text{Determination}}{\text{Direct from global PDF analysis}}$ $\text{LO from } \nu N \rightarrow \mu^+ \mu^- X$ $\text{NLO from } \nu N \rightarrow \mu^+ \mu^- X$	$\begin{array}{c c} V_{cd} \\ \hline 0.244 \pm 0.019^{\text{tot}} \\ 0.24 \pm 0.03 \\ 0.232^{+0.017}_{-0.019} \end{array}$
0.94	ANALYSIS NNPDF1.2 CDHS CCFR PDG08	DETERMINATION DIRECT FROM GLOBAL PDF ANALYSIS LO FROM $\nu N \rightarrow \mu^+ \mu^- X$ NLO FROM $\nu N \rightarrow \mu^+ \mu^- X$ AVERAGE FROM $\nu N \rightarrow \mu^+ \mu^- X$	$\begin{array}{c c} V_{cd} \\\hline 0.244 \pm 0.019^{\text{tot}} \\0.232 \pm 0.013 \\0.232 \pm 0.017 \\-0.019 \\0.230 \pm 0.011 \end{array}$
0.94 0.92 0.90	ANALYSIS NNPDF1.2 CDHS CCFR PDG08 PDG08	DETERMINATION DIRECT FROM GLOBAL PDF ANALYSIS LO FROM $\nu N \rightarrow \mu^+ \mu^- X$ NLO FROM $\nu N \rightarrow \mu^+ \mu^- X$ AVERAGE FROM $\nu N \rightarrow \mu^+ \mu^- X$ AVERAGE FROM $D \rightarrow K/\pi l\nu$ DECAYS	$\begin{array}{c c} V_{cd} \\\hline 0.244 \pm 0.019^{\rm tot} \\0.232 \pm 0.017 \\- 0.019 \\0.230 \pm 0.011 \\0.218 \pm 0.023 \end{array}$
0.94 0.92 0.90 0.88 Vcd 	ANALYSIS NNPDF1.2 CDHS CCFR PDG08 PDG08 PDG08 PDG08	DETERMINATION DIRECT FROM GLOBAL PDF ANALYSIS LO FROM $\nu N \rightarrow \mu^+ \mu^- X$ NLO FROM $\nu N \rightarrow \mu^+ \mu^- X$ AVERAGE FROM $\nu N \rightarrow \mu^+ \mu^- X$ AVERAGE FROM $D \rightarrow K/\pi l \nu$ DECAYS CKM UNITARITY FIT	$\begin{array}{c c} V_{cd} \\\hline 0.244 \pm 0.019^{\rm tot} \\0.244 \pm 0.03 \\0.232 \pm 0.017 \\-0.019 \\0.230 \pm 0.011 \\0.218 \pm 0.023 \\0.2256 \pm 0.0010 \end{array}$

 \Rightarrow most precise available direct determination of V_{cs}



RELATIVE UNCERTAINTY ON FLUX



- UNCERTAINTIES IN FUTURE (GLOBAL) NNPDF FIT CAN ONLY DECREASE
- SENSITIVITY TO α_s of NNPDF partons negligible

WHERE ARE WE GOING? WHAT IS THIS GOOD FOR? CAN STUDY UNCERTAINTIES AT FUTURE ACCELERATORS

GLUON DETERMINATION AT THE LHEC

IMPACT OF LHEC F_2 data



CAN DISENTANGLE DIFFERENT SCENARIOS FOR SMALL \boldsymbol{x} GLUON BEHAVIOUR

- THE MONTE CARLO SAMPLE CAN BE USED DIRECTLY FOR SUCH STUDIES WITHOUT REFITTING: JUST REWEIGHT SAMPLE ACCORDING TO PROJECTED DATA
- LHC STUDIES ALONG THESE LINES POSSIBLE/IN PREPARATION (SEE MCNULTY'S TALK)

WHERE ARE WE GOING? WHAT NEXT? ROADMAP

- $2007 \rightarrow \text{First}$ (nonsinglet) fit 500 DIS data; one PDF
- $2008 \rightarrow$ First parton set 4000 DIS data; five, then seven PDFs
- $2009 \rightarrow$ First global fit 4000 data: DIS, DY, jets; seven PDFs



NNPDF ROADMAP

WHERE ARE WE GOING? WHAT NEXT? NNPDF2.0



• $800_{\text{DY}} + 200_{\text{JET}} = \mathcal{O}(1000)$ NeW DATA

- FASTDY ALGORITHM \Rightarrow FIRST TRULY NLO FIT
- **DEFECT:** STILL ZM-VFN SCHEME FOR HEAVY QUARKS

SUMMARY

- THE NNPDF APPROACH TRIES TO AVOID AS MUCH AS POSSIBLE
 - PARAMETRIZATION BIAS
 - SUBJECTIVE CHOICES IN UNCERTAINTY ESTIMATE AND COMBINATION
 - ASSUMPTIONS ON THE SHAPE OF THE UNDERLYING PROBABILITY DISTRIBUTION (GAUSSIANITY, SYMMETRY, LINEAR ERROR PROPAGATION...)
- A SET OF PDFS BASED ON THIS ASSUMPTIONS IS AVAILABLE, IT HAS BEEN USED FOR PHYSICAL APPLICATIONS AND IT IS ALREADY USEFUL FOR COMPETITIVE UNCERTAINTY ESTIMATE
- A GLOBAL SET OF PDFS BASED ON THIS METHODOLOGY IS BEHIND THE CORNER
- IT WILL ALLOW STUDIES OF DATA COMPATIBILITY, UNCERTAINTY ASSESSMENT, AND IT WILL VALIDATE OTHER PDF STUDIES



POSITIVITY

- PDFs can go negative, provided physical cross sections remain positive definite
- (IN PRINCIPLE) THIS SHOULD BE TRUE FOR ANY OBSERVABLE
- IN PRACTICE, MOSTLY RELEVANT CLOSE TO THE DATA EDGE
- NNPDF1.X: POSITIVITY OF F_L ENFORCED



STRANGE PDF AND DIMUON XSECT

NNPDF1.2, Q² = 20 GeV², y = 0.4

- STRANGENESS GOES NEGATIVE PRETTY SOON...
- BUT THE CROSS-SECTION REMAINS POSITIVE



COMPARISON PLOTS

REPLICAS: STRANGENESS FLEXIBLE PARAMETRZATION WITH LITTLE INFORMATION



IMPACT ON SOME LHC STANDARD CANDLES: IT COULD BE WORSE

- TOTAL W/Z XSCT DOMINATED BY CENTRAL RAPIDITY REGION \Rightarrow UNCERTAINTY ON STRANGENESS UNDER CONTROL
- AWAY FROM CENTRAL REGION UNCERTAINTY MUCH LARGER







RAPIDITY DISTRIBUTION

THE NEURAL MONTE CARLO



FEATURES OF THE FIT THE DATASET



$$Q^2 > 2 \text{ GeV}^2$$
; $W^2 > 12.5 \text{ GeV}^2$

NAME	DATA POINTS	TARGET
NMC_PD	153	F_2^d/F_2^p
NMC	245	$F_2^{\tilde{p}}$
SLAC	47 (47)	$F_2^{\overline{p}(d)}$
BCDMS	333 (248)	$F_2^{p(d)}$
ZEUS97	240 (29)	$\tilde{\sigma}^{\tilde{N}C(CC),+}$
ZEUS02	92 (26)	$\tilde{\sigma}^{NC(CC),-}$
ZEUS03	90 (30)	$\tilde{\sigma}^{NC(CC),+}$
H1Lx97	135	$\int \tilde{\sigma}^{NC}, +$
H197	130 (25)	$\tilde{\sigma}^{NC(CC),+}$
H199	139 (28)	$\tilde{\sigma}^{NC(CC),-}$
H100	147 (28)	$\tilde{\sigma}^{NC(CC),+}$
H108	8	F_{L}
CHORUS	471 (471)	$\tilde{\sigma}^{\nu(\bar{\nu})}$
TOTAL	3161	

FEATURES OF THE FIT THEORY

- **NLO** EVOLUTION (N SPACE, EXPANDED)
- ZM-VFN SCHEME FOR THRESHOLDS
- $\alpha_s(M_z) = 0.119$, PDFS given at $Q_0^2 = 2 \text{ GeV}^2$
- TARGET-MASS CORRECTIONS INCLUDED UP TO TWIST FOUR

BASIS FUNCTIONS AND PARAMETRIZATION

- FIVE INDEPENDENT PDFs: SINGLET, GLUON, TOTAL VALENCE, TRIPLET, $\bar{d} \bar{u}$.
- Symmetric strange sea $s(x) = \bar{s}(x)$, proportional to non-strange sea, $\bar{s}(x) = \frac{C}{2}(\bar{u}(x) + \bar{d}(x))$, (C = 0.5)
- All PDFS parametrized by a 2-5-3-1 neural network: $37 \times 5 = 185$ parameters
- MOMENTUM AND VALENCE SUM RULES ENFORCED STRICTLY
- POSITIVITY OF F_L ENFORCED for $x \ge 10^{-7}, Q^2 \ge 2 \text{ GeV}^2$



RESULTS PHYSICAL OBSERVABLES



TOTAL CROSS-SECTIONS AT LHC, NLO FROM MCFM

13	9.5
$ \begin{array}{ c c c c c c c c c } \sigma_{W} + \mathcal{B}_{l+\nu} & \Delta\sigma/\sigma & \sigma_{W} - \mathcal{B}_{l-\nu} & \Delta\sigma/\sigma_{\overline{x}} \end{array} \end{array} $	
- $[nb]$ W^+ $[nb]$ $W^{-\frac{c}{d}}$	
NNPDF08 11.96 ± 0.30 2.5% 8.49 ± 0.19 2.3% \vdots	8.5
$\begin{array}{ $	1E CTE065 8 NNPDE08 (mail) CTE061 MPST2001E CTE065
$\begin{bmatrix} \text{CTEQ6.1} & 11.85 \pm 0.28 & 2.4\% & 8.73 \pm 0.23 & 2.6\% \end{bmatrix}$	
$\begin{array}{ c c c c c c c c c c c c c c c c c c c$	7.5
$ \begin{array}{ c c c c c c c c } \hline \sigma_{Z}\mathcal{B}_{l+l-} & \Delta\sigma/\sigma & \sigma_{t\bar{t}} & \Delta\sigma/\sigma & \sigma_{H} & \Delta \sigma/\sigma \\ \hline \end{array} $	$\Delta \sigma / \sigma$ Z ⁰ Cross Section at the LHC [MCFM]
$\begin{bmatrix} nb \end{bmatrix}^{} Z \begin{bmatrix} pb \end{bmatrix} t\overline{t} \begin{bmatrix} pb \end{bmatrix}$	H 23
NNPDF08 2.22 ± 0.04 2.0% 1014 ± 24 2.3% 35.79 ± 1.04 35.79 ± 1.04	3.0%
$ \text{CTEQ6.5} 2.27 \pm 0.05 2.2\% 942 \pm 19 2.0\% 37.51 \pm 0.80 10\%$	2.2% <u>z</u> ²² <u>i</u>
$ $ CTEQ6.1 $ $ 2.12 \pm 0.05 $ $ 2.3% $ $ 970 \pm 18 $ $ 1.9% $ $ 38.50 \pm 0.85 $ $ 5	
MRST01 1.98 ± 0.02 1.0% 1013 ± 13 1.3% 37.52 ± 0.40	1.1% [№] 2

1.9

1.8

NNPDF08 (prel) CTEQ61 MRST2001E CTEQ65

RESULTS GENERAL STATISTICAL FEATURES





- POISSONIAN DISTRIBUTION OF TRAINING LENGTHS
- BEST FIT $\chi^2 = 1.34$: MINOR DATA INCOMPATIBILITIES (?)

PARAMETRIZATION INDEPENDENCE: METHODOLOGY

- EFFECTIVELY INFINITE NUMBER OF PARAMETERS \Rightarrow CAN REPRESENT ANY FUNCTION
- COMPLEX SHAPES (LARGE NO.OF PARAMETERS) REQUIRE LONGER FITTING
- FIT STOPS WHEN QUALITY OF FIT TO RANDOMLY SELECTED "VALIDATION" DATA (NOT FITTED) STOPS IMPROVING
- CAN OBTAIN A FIT WITH χ^2 LOWER THAN BEST FIT ("OVERLEARNING")

PARAMETRIZATION INDEPENDENCE: REDUNDANCY AND OVERLEARNING

- OPTIMAL FIT OBTAINED WHEN QUALITY OF FIT TO VALIDATION (CONTROL) DATA STOPS IMPROVING
- POSSIBILITY OF OVERFITTING GUARANTESS THAT MINIMUM NOT DRIVEN BY PARAMETRIZATION



OPTIMAL FITTING

PARAMETRIZATION INDEPENDENCE: REDUNDANCY AND OVERLEARNING

- OPTIMAL FIT OBTAINED WHEN QUALITY OF FIT TO VALIDATION (CONTROL) DATA STOPS IMPROVING
- POSSIBILITY OF OVERFITTING GUARANTESS THAT MINIMUM NOT DRIVEN BY PARAMETRIZATION



OVERFITTING

- IRREGULAR OR KNOTTY SHAPES ALLOWED IF DATA FLUCTUATE
- STATISTICS SHOW WHETHER THE EFFECT IS REAL



10 REPLICAS

- IRREGULAR OR KNOTTY SHAPES ALLOWED IF DATA FLUCTUATE
- STATISTICS SHOW WHETHER THE EFFECT IS REAL



100 REPLICAS

- IRREGULAR OR KNOTTY SHAPES ALLOWED IF DATA FLUCTUATE
- STATISTICS SHOW WHETHER THE EFFECT IS REAL



200 REPLICAS

- IRREGULAR OR KNOTTY SHAPES ALLOWED IF DATA FLUCTUATE
- STATISTICS SHOW WHETHER THE EFFECT IS REAL



400 REPLICAS

- IRREGULAR OR KNOTTY SHAPES ALLOWED IF DATA FLUCTUATE
- STATISTICS SHOW WHETHER THE EFFECT IS REAL



1000 REPLICAS

PARAMETRIZATION INDEPENDENCE: STATISTICAL STABILITY

COMPARE DISTANCE IN UNITS OF SIGMA OF RESULTS OBTAINED WITH DIFFERENT ASSUMPTIONS

• DISTANCE IN UNITS OF SIGMA

$$\langle d[q] \rangle = \sqrt{\left\langle \frac{\left(\langle q_i \rangle_{(1)} - \langle q_i \rangle_{(2)} \right)^2}{\sigma^2 [q_i^{(1)}] + \sigma^2 [q_i^{(2)}]} \right\rangle_{\text{dat}}}$$

- NOTE σ ⇒ ERROR ON AVERAGE
 = (ERROR ON q_i)/√N with 100 replicas,d = 1 → fits differ by 1/10 of nominal error
- TEST PREDICTIONS FOR CENTRAL VALUES & ERRORS

DISTANCE BETWEEN STANDARD & FIT WITH SMALLER NEURAL NETS 2-4-3-1 VS 2-5-3-1 ARCHITECTURE (31 vs. 37 parms per net)

`	_	- /
	DATA	EXTRAPOLATION
SINGLET	$0.005 \le x \le 0.1$	$10^{-4} \le x \le 10^{-3}$
$\langle d[q] \rangle$	0.96	1.32
$\langle d[\sigma] \rangle$	1.23	1.32
GLUON	$0.005 \le x \le 0.1$	$10^{-4} \le x \le 10^{-3}$
$\langle d[q] \rangle$	1.40	1.13
$\langle d[\sigma] \rangle$	1.17	1.06
VALENCE	$0.1 \le x \le 0.6$	$0.03 \le x \le 0.3$
$\langle d[q] \rangle$	1.40	0.93
$\langle d[\sigma] \rangle$	1.09	0.96
TRIPLET	$0.05 \le x \le 0.75$	$0.01 \le x \le 0.1$
$\left[\langle d[q] \rangle \right]$	1.05	1.09
$\langle d[\sigma] \rangle$	1.68	2.5

PARAMETRIZATION INDEPENDENCE: THE "HERALHC BENCHMARK"

 $Q^2 > 9 \; \mathrm{GeV}^2; \, W^2 > 15 \; \mathrm{GeV}^2$

REDUCED DATASET \Rightarrow WIDER ERROR BAND from 3161 to 773 datapoints reduced info on small x sea (no low Q^2 data) & large x valence (no neutrino data)

UP QUARK

NAME	DATA POINTS	TARGET
NMC_PD	73	F_2^d/F_2^p
NMC	95	$F_2^{\overline{p}}$
BCDMS	322	F_2^{p}
ZEUS97	206	F_2^{p}
H1LX97	77	F_2^{p}
TOTAL	773	





RESULTS COMPATIBLE TO WITHIN LESS THAN TWO SIGMA

PARAMETRIZATION INDEPENDENCE: THE "HERALHC BENCHMARK":INCOMPATIBLE DATA



NO ERROR REDUCTION WHEN DATA IN WIDER DATA SET ARE INCOMPATIBLE

DELIVERY: RESTRICTED SAMPLE OF REPLICAS

- WIDE SAMPLE OF PSEUDODATA ENDURES NO BIAS
- IMPRACTICAL TO AVERAGE OVER THOUSAND(S) OF REPLICAS
- SELECT SUBSET OF REPLICAS WITH APPROXIMATELY SAME STATISTICAL DISTRIBUTION AS FULL SET
- construct histogram for #
 of replicas n sigma away
 from mean
 compare result for subset & singlet at x=0.1
 Singlet at x=0.1
 Singlet at x=0.1
 Yalence at x=0.01
 Yalence at x=0.01
- compare result for subset & s full
- minimize relative entropy of two histograms S ==

$$\sum_{i \text{ bins}} \left(p_i^{(1)} - p_i^{(2)} \right) \ln \frac{p_i^{(1)}}{p_i^{(2)}}$$

• select with genetic algorithm subset which minimizes S





OPTIMAL FITTING VS. OVERLEARNING: AN EXAMPLE THE TRUE FUNCTION



OPTIMAL FITTING VS. OVERLEARNING: AN EXAMPLE UNDERLEARNING



OPTIMAL FITTING VS. OVERLEARNING: AN EXAMPLE OPTIMAL FIT



OPTIMAL FITTING VS. OVERLEARNING: AN EXAMPLE OVERLEARNING



WHAT ARE NEURAL NETWORKS?



MULTILAYER FEED-FORWARD NETWORKS

- Each neuron receives input from neurons in preceding layer and feeds output to neurons in subsequent layer
- Activation determined by weights and thresholds

$$\xi_i = g\left(\sum_j \omega_{ij}\xi_j - \theta_i\right)$$

• Sigmoid activation function $g(x) = \frac{1}{1 + e^{-\beta x}}$

JUST ANOTHER SET OF BASIS FUNCTIONS!

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

ANY FUNCTION CAN BE REPRESENTED BY A SUFFICIENTLY BIG NEURAL NETWORK LESS PARAMETERS \rightarrow SMOOTHER FUNCTIONS

IN A STANDARD FIT, ONE LOOKS FOR MINIMUM χ^2 WITH GIVEN FINITE PARM.

- IF THE BASIS IS TOO LARGE, THE FIT NEVER CONVERGES
- IF THE BASIS IS TOO SMALL, THE FIT IS BIASED

Q: HOW CAN ONE BE SURE THAT THE COMPROMISE IS UNBIASED?

IN A STANDARD FIT, ONE LOOKS FOR MINIMUM χ^2 WITH GIVEN FINITE PARM.

- IF THE BASIS IS TOO LARGE, THE FIT NEVER CONVERGES
- IF THE BASIS IS TOO SMALL, THE FIT IS BIASED

IN A STANDARD FIT, ONE LOOKS FOR MINIMUM χ^2 WITH GIVEN FINITE PARM.

- IF THE BASIS IS TOO LARGE, THE FIT NEVER CONVERGES
- IF THE BASIS IS TOO SMALL, THE FIT IS BIASED



IN A STANDARD FIT, ONE LOOKS FOR MINIMUM χ^2 WITH GIVEN FINITE PARM.

- IF THE BASIS IS TOO LARGE, THE FIT NEVER CONVERGES
- IF THE BASIS IS TOO SMALL, THE FIT IS BIASED



IN A STANDARD FIT, ONE LOOKS FOR MINIMUM χ^2 WITH GIVEN FINITE PARM.

- IF THE BASIS IS TOO LARGE, THE FIT NEVER CONVERGES
- IF THE BASIS IS TOO SMALL, THE FIT IS BIASED



IN A STANDARD FIT, ONE LOOKS FOR MINIMUM χ^2 WITH GIVEN FINITE PARM.

- IF THE BASIS IS TOO LARGE, THE FIT NEVER CONVERGES
- IF THE BASIS IS TOO SMALL, THE FIT IS BIASED

Q: HOW CAN ONE BE SURE THAT THE COMPROMISE IS UNBIASED? IN A NEURAL FIT, SMOOTHNESS DECREASES AS FIT QUALITY IMPROVES:



A: STOP THE FIT BEFORE OVERLEARNING SETS IN!

IN A STANDARD FIT, ONE LOOKS FOR MINIMUM χ^2 WITH GIVEN FINITE PARM.

- IF THE BASIS IS TOO LARGE, THE FIT NEVER CONVERGES
- IF THE BASIS IS TOO SMALL, THE FIT IS BIASED

Q: HOW CAN ONE BE SURE THAT THE COMPROMISE IS UNBIASED? IN A NEURAL FIT, SMOOTHNESS DECREASES AS FIT QUALITY IMPROVES:



A: STOP THE FIT BEFORE OVERLEARNING SETS IN! COULD BE DONE WITH STANDARD PARAMETRIZATIONS, BUT VERY INEFFICIENTLY

MINIMIZE BY GENETIC ALGORITHM: AT EACH GENERATION, THE 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



MINIMIZE BY GENETIC ALGORITHM: AT EACH GENERATION, THE 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)

GO!

• WHEN THE VALIDATION χ^2 STOPS DECREASING, STOP THE FIT



MINIMIZE BY GENETIC ALGORITHM: AT EACH GENERATION, THE 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



STOP!

MINIMIZE BY GENETIC ALGORITHM: AT EACH GENERATION, THE 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)

TOO LATE!

• WHEN THE VALIDATION χ^2 STOPS DECREASING, STOP THE FIT



MONTE CARLO DATA GENERATION

- BCDMS+ NMC PROTON & DEUTERON F_2 DATA (FULL CORRELATED SYSTEMATICS AVAILABLE), TAKEN AT 4 BEAM ENERGIES
- ON TOP OF STAT. ERRORS, 4 SYSTEMATICS + 1 NORMALIZATION (NMC) OR 6 SYSTEMATICS + 1 ABSOLUTE & 2 RELATIVE NORMALIZATIONS (BCDMS), WITH VARIOUS FORMS OF CORRELATION (FULL, OR FOR EACH TARGET, OR FOR EACH BEAM ENERGY)

GENERATE DATA ACCORDING TO A MULTIGAUSSIAN DISTRIBUTION

$$F_{i}^{(art)(k)} = (1 + r_{5}^{(k)} \sigma_{N}) \sqrt{1 + r_{i,6}^{(k)} \sigma_{N_{t}}} \sqrt{1 + r_{i,7}^{(k)} \sigma_{N_{b}}} \left[F_{i}^{(exp)} + \frac{r_{i,1}^{(k)} f_{b} + r_{i,2}^{(k)} f_{i,s} + r_{i,3}^{(k)} f_{i,r}}{100} F_{i}^{(exp)} + r_{i,s}^{(k)} \sigma_{s}^{i} \right]$$

r univariate gaussian random nos., one $r_{i,s}$ for each data, but single $r_{i,j}$ for all correlated data



SCATTER PLOT ART. VS. EXP. FOR 10 (RED) 100 (GREEN) AND 1000 (BLUE) REPLICAS

NEED 1000 REPLICAS TO REPRODUCE CORRELATIONS TO PERCENT ACCURACY

QUESTION: (F. OLNESS) HOW "STRANGE" IS THE STRANGE PDF?

QUESTION: (F. OLNESS) HOW "STRANGE" IS THE STRANGE PDF?

ANSWER: THE NNPDF1.2 SET

THE DATA SET

Х

NMC-pd Х Ж NMC SLAC **10**⁴ • SEVEN PDFS BCDMS ZEUS INDEP. PARAMETRIZED H1 CHORUS 74 EXTRA FREE PARMS FLH108 10³ (GeV²) 2010² (s^{\pm}) IN COMPARISON NTVDMN TO ZEUS-H2 NNPDF1.0 • NUTEV DIMUON DATA USED TO CONSTRAIN **STRANGENESS** 10 1 **10**⁻⁴ **10⁻³** 10⁻² **10**⁻¹ 1

STABILITY COMPARISON TO PREVIOUS NNPDF SETS

- NNPDF1.0: $s(x, Q_0^2) = \bar{s}(x, Q_0^2), s + \bar{s} = \frac{1}{2}(\bar{u} + \bar{d})$
- NNPDF1.1: s, \bar{s} (actually s^{\pm}) indep. parametrized, no dimuon data
- NNPDF1.2: s, \bar{s} indep. parametrized, dimuon data



NONSTRANGE PDFS

DETERMINING STRANGENESS COMPARISON TO PREVIOUS NNPDF SETS

- NNPDF1.0: STRANGENESS UNCERTAINTY UNDERESTIMATED
- NNPDF1.1: STRANGENESS UNCERTAINTY HUGE
- NNPDF1.2: STRANGENESS UNCERTAINTY UNDER CONTROL



STRANGE PDFS

DETERMINING STRANGENESS COMPARISON TO OTHER NNPDF SETS

- CTEQ6.6: $s = \bar{s}$, s^+ parm w. two free parameters
- MSTW08: s^+ & s^- parm w. two free parameters each
- EVERYBODY ENFORCES STRANGENESS SUM RULE



STRANGE PDFS

THE NUTEV ANOMALY

THE NUTEV ANOMALY

IS GONE

$$R_{\rm PW} \equiv \frac{\sigma(\nu \mathcal{N} \to \nu X) - \sigma(\bar{\nu} \mathcal{N} \to \bar{\nu} X)}{\sigma(\nu \mathcal{N} \to \ell X) - \sigma(\bar{\nu} \mathcal{N} \to \bar{\ell} X)}$$
$$= \frac{1}{2} - \sin^2 \theta_{\rm W} + \left(\frac{(U^- - D^-) + (C^- - S^-)}{Q^-} \frac{1}{6} \left(3 - 7\sin^2 \theta_{\rm W}\right)\right),$$

RATIO

Determinations of the weak mixing angle $\sin^2 \theta_W$

