Simultaneous extraction of spin-dependent parton distributions and fragmentation functions in the MC framework

Nobuo Sato University of Connecticut/JLab Seminar at BNL Dec 15, 2017 In collaboration with: J. Ethier, W. Melnitchouk





Outline

• The fitting methodology

 \blacksquare QCD analysis of $\Delta {\rm PDFs}$ and FFs

The fitting methodology

The parent distribution

"If we could make an infinite number of measurements, then we could describe exactly the distribution of the data points. This is not possible in practice, but we can hypothesize the existence of such a distribution that determines the probability of getting any particular observation in a single measurement. This distribution is called **parent distribution**. Similarly we can hypothesize that the measurements we have make are samples from the parent distribution and they form the sample distribution. In the limit of an infinite number of measurements, the sample distribution becomes the parent distribution"

Data reduction and error analysis for the physical sciences Bevington and Robison

 Consider a quantity <u>f</u> for which we want to determine its parent distribution

$\mathcal{P}(f)$

We are interested in the case where f cannot be measured directly, but instead it is inferred from experimental data. In this case the parent distribution is conditioned to the evidence, and mathematically this is written as



$\mathcal{P}(f|data)$

• How do we compute $\mathcal{P}(f|data)$? \rightarrow Bayes theorem:

 $\mathcal{P}(f|data) = \frac{1}{Z}\mathcal{L}(data|f)\pi(f)$

 $\mathcal{L}(data|f)$: Likelihood $\pi(f)$: prior Z: evidence

The likelihood function is chosen to describe the probability of the data to be drawn from a model with a given *f*. e.g Gaussian likelihood

$$\mathcal{L}(data|f) = \exp\left[-\frac{1}{2}\sum_{i}\left(\frac{d_i - \text{model}_i(f)}{\delta d_i}\right)^2\right]$$



The prior function allows us to restrict unphysical regions of f. We make the priors to be as flat as possible to avoid biases (uninformative priors) i.e.

$$\pi(f) = \begin{cases} 1 & \operatorname{condition}(f) == \operatorname{True} \\ 0 & \operatorname{condition}(f) == \operatorname{False} \end{cases}$$

$$\mathcal{P}(f|d) = \frac{1}{Z}\mathcal{L}(d|f)\pi(f)$$

 In practice f needs to be represented mathematically e.g

$$f(x) = Nx^{a}(1-x)^{b}(1+c\sqrt{x}+dx+...)$$

$$f(x) = Nx^{a}(1-x)^{b}NN(x; \{w_{i}\})$$

$$f(x) = NN(x; \{w_{i}\}) - NN(1; \{w_{i}\})$$



• The parent distribution for *f* becomes

$$\boldsymbol{a} = (N, a, b, c, d, ...)$$
$$\mathcal{P}(\boldsymbol{a}|d) = \frac{1}{Z} \mathcal{L}(d|\boldsymbol{a}) \pi(\boldsymbol{a})$$
$$\mathcal{L}(d|\boldsymbol{a}) = \exp\left[-\frac{1}{2} \sum_{i} \left(\frac{d_{i} - \text{model}_{i}(\boldsymbol{a})}{\delta d_{i}}\right)^{2}\right]$$
$$\pi(\boldsymbol{a}) = \prod_{i} \theta(a_{i} - a_{i}^{min}) \theta(a_{i}^{max} - a_{i})$$

$$\mathcal{P}(f|d) = \frac{1}{Z}\mathcal{L}(d|f)\pi(f)$$

$$\downarrow$$

$$\mathcal{P}(\boldsymbol{a}|d) = \frac{1}{Z}\mathcal{L}(d|\boldsymbol{a})\pi(\boldsymbol{a})$$

Having the parent distribution we can compute

$$E[\mathcal{O}] = \int d^{n}a \ \mathcal{P}(\boldsymbol{a}|data) \ \mathcal{O}(\boldsymbol{a})$$
$$V[\mathcal{O}] = \int d^{n}a \ \mathcal{P}(\boldsymbol{a}|data) \ (\mathcal{O}(\boldsymbol{a}) - E[\mathcal{O}])^{2}$$

O is any function of a. e.g

$$\mathcal{O}(\boldsymbol{a}) = f(x; \boldsymbol{a})$$
$$\mathcal{O}(\boldsymbol{a}) = \int_{x}^{1} \frac{d\xi}{\xi} C(\xi) f\left(\frac{x}{\xi}; \boldsymbol{a}\right)$$

• How do we compute $E[\mathcal{O}], V[\mathcal{O}]$?

- Maximum likelihood
- Monte Carlo approach



Attention:

- typically $n \gg 1$
- $\mathcal{P}(\boldsymbol{a}|data)$ is computationally expensive
- for \$\mathcal{O} == f(x)\$, an n-dim integration is needed for each x. Not practical!

Maximum Likelihood

Estimation of expectation value

$$\mathbf{E}[\mathcal{O}] = \int d^n a \ \mathcal{P}(\boldsymbol{a}|data) \ \mathcal{O}(\boldsymbol{a}) \simeq \mathcal{O}(\boldsymbol{a}_0)$$

• a_0 is estimated from optimization algorithm

$$\max \left[\mathcal{P}(\boldsymbol{a}|data) \right] = \mathcal{P}(\boldsymbol{a}_0|data)$$
$$\max \left[\mathcal{L}(data|\boldsymbol{a})\pi(\boldsymbol{a}) \right] = \mathcal{L}(data|\boldsymbol{a}_0)\pi(\boldsymbol{a}_0)$$

equivalently

$$\min \left[-2 \log \left(\mathcal{L}(data|\boldsymbol{a})\pi(\boldsymbol{a})\right)\right] = -2 \log \left(\mathcal{L}(data|\boldsymbol{a}_0)\pi(\boldsymbol{a}_0)\right)$$
$$= \sum_{i} \left(\frac{d_i - \text{model}_i(\boldsymbol{a}_0)}{\delta d_i}\right)^2 - 2 \log \left(\pi(\boldsymbol{a}_0)\right)$$
$$= \chi^2(\boldsymbol{a}_0) - 2 \log \left(\pi(\boldsymbol{a}_0)\right)$$
$$\text{this is Chi-squared}$$
minimization

Maximum Likelihood + Hessian method

Estimation of variance

$$\mathbf{V}[\mathcal{O}] = \int d^{n}a \ \mathcal{P}(\boldsymbol{a}|data) \ (\mathcal{O}(\boldsymbol{a}) - \mathbf{E}[\mathcal{O}])^{2}$$

 \blacksquare Eigen direction decomposition of $\mathcal{P}(\pmb{a}|data)$

$$\begin{split} \mathcal{P}(\boldsymbol{a}|data) &\propto \exp\left(-\frac{1}{2}\chi^{2}(\boldsymbol{a})\right) \propto \exp\left(-\frac{1}{2}\chi^{2}(\boldsymbol{a}_{0}) - \frac{1}{2}\Delta\chi^{2}(\boldsymbol{a})\right) \right) \\ &\propto \exp\left(-\frac{1}{2}\Delta\chi^{2}(\boldsymbol{a})\right) \right) \\ &\propto \exp\left(-\frac{1}{2}\Delta\boldsymbol{a}^{T} H \Delta\boldsymbol{a}\right) + O(\Delta a^{3}) \\ &\propto \exp\left(-\frac{1}{2}\sum_{k}\left(t_{k}\frac{\hat{\boldsymbol{e}}_{k}^{T}}{\sqrt{w_{k}}}\right) H \sum_{l}\left(t_{l}\frac{\hat{\boldsymbol{e}}_{l}}{\sqrt{w_{l}}}\right)\right) + O(\Delta a^{3}) \\ &\propto \exp\left(-\frac{1}{2}\sum_{k}t_{k}^{2}\right) + O(\Delta a^{3}) \\ &\propto \prod_{k}\exp\left(-\frac{1}{2}t_{k}^{2}\right) + O(\Delta a^{3}) \\ &\propto \prod_{k}\exp\left(-\frac{1}{2}t_{k}^{2}\right) + O(\Delta a^{3}) \end{split}$$
 The probability distribution "factorizes" along each eigen direction

Maximum Likelihood + Hessian method

Estimation of variance

$$\mathbf{V}[\mathcal{O}] = \int d^{n}a \ \mathcal{P}(\boldsymbol{a}|data) \ (\mathcal{O}(\boldsymbol{a}) - \mathbf{E}[\mathcal{O}])^{2}$$

 \blacksquare Linear approximation of $\mathcal{O}(\boldsymbol{a})$

$$\left[\mathcal{O}(\boldsymbol{a}) - \mathbf{E}[\mathcal{O}]\right]^2 = \left[\sum_i \frac{\partial \mathcal{O}}{\partial a_i}(a_i - a_0) + O(a^2)\right]^2 = \left[\sum_k \frac{\partial \mathcal{O}}{\partial t_k} t_k\right]^2 + O(a^3)$$

 \blacksquare Combining with factorized $\mathcal{P}(\pmb{a}|data)$ we get

$$\begin{split} \mathbf{V}[\mathcal{O}] &\simeq \prod_{k} \int dt_{k} \frac{e^{-\frac{1}{2}t_{k}^{2}}}{\sqrt{2\pi}} \sum_{lm} \frac{\partial \mathcal{O}}{\partial t_{l}} \frac{\partial \mathcal{O}}{\partial t_{m}} t_{l} t_{m} \\ &= \sum_{k} \left(\frac{\partial \mathcal{O}}{\partial t_{k}} \right)^{2} \simeq \sum_{k} \left[\frac{\mathcal{O}(t_{k}=1) - \mathcal{O}(t_{k}=-1)}{2} \right]^{2} \end{split}$$

Maximum Likelihood + Hessian method

pros

- $\rightarrow\,$ Very practical. Most of the PDF groups use this method
- $\rightarrow\,$ It is computationally inexpensive
- $\rightarrow~f$ and its eigen directions can be precalculated/tabulated

cons

- ightarrow Assumes local gaussian approximation of the likelihood
- ightarrow Assumes linear approximation of the observables ${\cal O}$ around $oldsymbol{a}_0$
- ightarrow The assumptions are strictly valid for linear models.
- $\rightarrow\,$ Computation of the hessian matrix is numerically unstable if flat directions are present

examples

$$\rightarrow$$
 if $f(x) = a + bx + cx^2$ then $\mathbf{E}[f(x)] = \mathbf{E}[a] + \mathbf{E}[b]x + \mathbf{E}[c]x^2$

 \rightarrow but $f(x)=Nx^a(1-x)^b$ then $\mathrm{E}[f(x)]\neq\mathrm{E}[N]x^{\mathrm{E}[a]}(1-x)^{\mathrm{E}[b]}$

Monte Carlo Methods

Recall that we are interested in computing

$$E[\mathcal{O}] = \int d^{n}a \ \mathcal{P}(\boldsymbol{a}|data) \ \mathcal{O}(\boldsymbol{a})$$
$$V[\mathcal{O}] = \int d^{n}a \ \mathcal{P}(\boldsymbol{a}|data) \ (\mathcal{O}(\boldsymbol{a}) - E[\mathcal{O}])^{2}$$

 Any MC method attempts to do this using MC sampling

$$\begin{split} \mathbf{E}[\mathcal{O}] &\simeq \sum_{k} w_k \mathcal{O}(\boldsymbol{a}_k) \\ \mathbf{V}[\mathcal{O}] &\simeq \sum_{k} w_k (\mathcal{O}(\boldsymbol{a}_k) - \mathbf{E}[\mathcal{O}])^2 \end{split}$$

 Here {w_k, a_k} is the sample distribution of the parent distribution P(a|data)

■ Given the $\mathcal{P}(\boldsymbol{a}|data)$ the sample distribution is unique, regardless of the MC method

$$\rightarrow \sum_k w_k = 1$$

 \rightarrow unweighted sampling

$$w_1 = w_2 = \dots$$

 \rightarrow weighted sampling $w_1 \neq w_2 \neq \dots$

MC Method 1: data resampling

 Construct pseudo data sets where each data point is sampled using Gaussian distribution with mean and variance given by the original data

$$d_{k,i}^{\text{(pseudo)}} = d_i^{\text{(exp)}} + \sigma_i^{\text{(exp)}} R_{k,i}$$

- i: i-th data point
- k: k-th pseudo data set index

 $R_{k,i}$: random number from normal distribution

Fit each pseudo data sample k = 1, ..., N to obtain parameter vectors a_k The sample distribution of P(a|data) is approximately

$$\{w_k = 1/N, \boldsymbol{a}_k\}$$

here "fit" means Chi-square minimization

MC Method 1+: data resampling+cross validation

Issues with number of parameters

- $\rightarrow\,$ Ideally one should not be worried about the number of parameters to be used.
- $\rightarrow\,$ This is an issue for Hessian method due to the flat directions.
- → However flat directions are typically only a local feature of the parent distribution.

Over-fitting

- $\label{eq:parameters} \begin{array}{l} \rightarrow & \mbox{If there are too many parameters there} \\ & \mbox{would be regions in the parameter space} \\ & \mbox{where } \mathcal{P}(\boldsymbol{a}|data) \mbox{ develops "spikes"} \\ & \rightarrow \mbox{ signal of over-fitting} \end{array}$
- $\rightarrow~$ Use cross-validation to tame the "spikes"

Procedure

- → For each pseudo data sample k split randomly the data set in 50/50 and label them as "training" and "validation" respectively
- → Fit the "training" set and stop the fitting whenever the description of the "validation" set deteriorates to avoids over-fitting

MC Method 1+: data resampling+cross validation

Procedure

- \rightarrow For each pseudo data sample k split randomly the data set in 50/50 and label them as "training" and "validation" respectively
- $\rightarrow\,$ Fit the "training" set and stop the fitting whenever the description of the "validation" set deteriorates $\rightarrow\,$ it avoids over-fitting

Caveat

 \rightarrow the resulting sample distribution is sensitive to the partition. Possible solutions include to rescale the uncertainties of the training and validation set to compensate for the splitting

MC Method 1+++: data resampling+cross validation

One vs. multiple minima

- ightarrow It is possible that $\mathcal{P}(oldsymbol{a}|data)$ is multi modal.
- $\rightarrow\,$ Hence it is important to start the scan from many different starting points

Caveat

- \rightarrow Optimization algorithms are based on gradient descent search. It is possible that in a given run with N independent scans the sample distribution does not represent accurately the "true" parent distribution
- → To solve this, we start a new run by sampling guessing parameters from the prior iteration

 $+a^{(ext{guess})}$ randomization

+iterative runs





MC Method 2: Hybrid Markov Chain Monte Carlo

The basic idea

- \rightarrow This is an MCMC based algorithm (random walks + rejection sampling)
- \rightarrow The random walks are optimized by solving Hamilton's equations.
- ightarrow The parameters $m{a}$ are the "coordinates" and a conjugate vector $m{p}$ e.g. "momentum" is defined
- ightarrow An initial "state" is defined by a random coordinate vector $m{a}_0$ and a random momentum vector $m{p}_0.$
- \rightarrow A new state is proposed by solving a Hamiltonian using the leap frog method

$$H(\boldsymbol{p}, \boldsymbol{a}) = \frac{\boldsymbol{p}^2}{2m} - \log(\mathcal{L}(\boldsymbol{a}))$$

pros

→ It provides a faithful sampling distribution

cons

- \rightarrow the number of steps and step size of the leap frog must be tuned.
- \rightarrow Cannot be parallelized

MC Method 3: nested sampling

The basic idea: compute

$$Z = \int \mathcal{L}(\text{data}|\boldsymbol{a}) \pi(\boldsymbol{a}) d^{n} \boldsymbol{a} = \int_{0}^{1} \mathcal{L}(X) dX$$

- \rightarrow The algorithm traverses ordered isolikelihood contours in the variable X such that X follows the progression $X_i = t_i X_{i-1}$
- \rightarrow The variable t_i is estimated statistically
- → The algorithm can be optimized iteration to iteration. One can sample only in the regions where the likelihood is larger → "importance sampling"
- $\rightarrow\,$ The nested sampling is parallelizable





Toy example

- \rightarrow We generate events from f(x) to mimic realistic counting experiment
- → The fits and the error bands are performed with four different algorithms
- → This is expected as all the methods uses same likelihood



QCD analysis \triangle **PDF**s and **FF**s

Polarized PDFs: inclusive polarized DIS NS, Melnitchouk, Kuhn, Ethier, Accardi (PRD 93,074005)



Fragmentation Functions: SIA

NS, Ethier, Melnitchouk, Hirai, Kumano, Accardi (PRD 94, 114004)





 $\rightarrow~\pi$ and K global data from Belle and Babar up to LEP data at $Q=M_z$

JAM and DSS $D_{s^+}^K$ consistent

The ΔS^+ puzzle



- \blacksquare Constraints on Δs^+
- \rightarrow JAM: Δ DIS + SU3
- \rightarrow DSSV: Δ DIS + SU3, Δ SIDIS

Note

- $\rightarrow\,$ DSSV analysis shows no violation of SU3 due to penalties
- $\rightarrow\,$ In DSSV, FF is extracted independently from SIA, SIDIS and pp data
- $\rightarrow~{\rm In}~{\rm JAM}$ negative Δs^+ comes only from SU3

Questions

- ightarrow What controls the sign of Δs^+ ?
- $\rightarrow\,$ What are the actual uncertainties on Δs^+ ?

Combined $\triangle PDF$ and FF: $\triangle DIS + \triangle SIDIS + SIA$

Ethier, NS, Melnitchouk (PRL 119, 132001)

Setup

- $\rightarrow\,$ Simultaneous extraction of polarized ΔPDFs and FFs
- ightarrow Data: Δ DIS, Δ SIDIS, SIA
- \rightarrow No SU(3) constraints

Results

- $\rightarrow\,$ Sea polarization consistent with zero
- $\rightarrow\,$ The current precision of ΔSIDIS data is not sufficient to determined the sea polarization
- $\rightarrow \ D_{s^+}^K$ consistent with SIA only analysis



What determines the sign of Δs^+ ?

case 1

- $ightarrow \sim 5 \text{ COMPASS } d \text{ data points at} x < 0.002 \text{ favor small } \Delta s^+(x)$
- $\rightarrow~$ To generate $\Delta s^{+(1)}(Q_0^2)\sim -0.1$ a peak at $x\sim 0.1$ is generated

case 2

- $\rightarrow\,$ In the absence of x<0.002 data, the negative $\Delta s^{+(1)}(Q_0^2)\sim -0.1$ is mostly generated at small x.
- $\rightarrow~$ No need for negative $\Delta s^+(x)$ at $x\sim 0.1$
 - case 3
- $\rightarrow \ \Delta s^+(x\sim 0.1) < 0$ disfavored by HERMES $A_{1d}^{K^-}$
- $\rightarrow \mbox{ Smaller } \Delta s^{+(1)}(Q_0^2)$ but larger uncertainties

case	data	sign change	$\Delta s^{+(1)}(Q_0^2)$
1	$\Delta DIS+SU(3)$	No	-0.1
2	$\Delta \text{DIS}+\text{SU}(3) \ (x > 0.02)$	Possible	-0.1
3	$\Delta DIS + \Delta SIDIS + FF$	Possible	-0.03(10)



Updates on the moments

- \rightarrow We construct flat priors that gives flat a_8 in order to have an unbias extraction of a_8
- $\rightarrow\,$ Data prefers smaller values for $a_8 \rightarrow 25\%$ larger total spin carried by quarks.
- $\label{eq:a3} \begin{array}{l} \rightarrow \ a_3 \ {\rm which} \ {\rm is} \ {\rm in} \ {\rm a} \ {\rm good} \\ {\rm agreement} \ {\rm with} \ {\rm values} \ {\rm from} \ \beta \\ {\rm decays} \ {\rm within} \ 2\%. \end{array}$
- $\label{eq:alpha} \begin{array}{l} \rightarrow & \mbox{Data indicates possible} \\ & \Delta \bar{u} > \Delta \bar{d} \mbox{ consistent with} \\ & \mbox{measurements of } W^{\pm}(Z) \\ & \mbox{asymmetries from PHENIX and} \\ & \mbox{STAR} \end{array}$



obs.	JAM15	JAM17
g_A	1.269(3)	1.24(4)
g_8	0.586(31)	0.46(21)
$\Delta\Sigma$	0.28(4)	0.36(9)
$\Delta \bar{u} - \Delta \bar{d}$	0	0.05(8)

SIDIS+Lattice analysis of nucleon tensor charge Lin, Melnitchouk, Prokudin, NS, Shows (arXiv:1710.09858)



- \rightarrow Extraction of transversity and Collins FFs from SIDIS A_{UT} +Lattice g_T
- $\rightarrow\,$ In the absence of Lattice, SIDIS at present has no significant constraints on $g_T \rightarrow$ this will change with the upcoming JLab12 measurements

Summary and outlook

JAM status

- $\rightarrow\,$ Global analysis methodology based on Bayesian perspectives has been studied and implemented.
- $\rightarrow\,$ Provides faithful representation of uncertainties
- $\rightarrow\,$ First simultaneous extraction of ΔPDF and FF from $\Delta \text{DIS},$ ΔSIDIS and SIA

Ongoing work

- $\rightarrow\,$ Almost done in implementing inclusive DIS analysis to extract PDFs
- \rightarrow Extraction of helicity distributions

Future work

- $\rightarrow\,$ Combined ($\Delta)$ DIS, ($\Delta)$ SIDIS and SIA to extract simultaneously ($\Delta)$ PDFs and FFs.
- $\rightarrow~(\Delta) {\rm Jets},~(\Delta) W^{\pm}$ from Tevatron and RHIC