Efficient statistics for future neutrino oscillation experiments

2022-06-19, Erice School of Subnuclear Physics Lukas Berns, Tohoku University

first few slides are with T2K and SK Collaborations





- Study oscillation of neutrino beam from J-PARC accelerator
- ~500 collaborators from institutions in

ν -oscillation

(interaction) (propagation) For neutrinos flavor basis \neq Hamiltonian basis.

→ Flavor ($\nu_e | \nu_\mu | \nu_\tau$) oscillates over $L \times \Delta m^2 / E$, amplitude controlled by (PMNS) mixing matrix *U*:



Neutrino beam

- 30 GeV protons produce π, K in 90 cm graphite target
- Three magnetic horns selectively focus π^+, K^+ or π^-, K^- to produce ν_{μ} or $\bar{\nu}_{\mu}$ beam (decay in-flight). $\int_{X}^{\frac{1}{2}} 0.5$ $\lim_{X_{2}^{n^{2}} 2.4 \times 10^{3} \text{ cv}^{2}}$ • Narrowband beam than 1 ks 3 to off-axis technique.

2.5°

2

E_v (GeV)

ND280

INGRID

to SuperK



The near detectors







INGRID on-axis detector

 Iron-scintillator sandwich detectors monitor neutrino beam direction and intensity ND280 off-axis detector

- Active scintillator + passive water targets
- Tracking with time projection chambers
- Magnetized for charge and momentum measurement

WAGASCI + BabyMIND

- Latest addition at intermediate **1.5**° off-axis flux
- Water target with cuboid lattice scintillators for high angle acceptance
- Compact **magnetized** iron muon range detector
- First xsec meas. published:
 <u>PTEP, ptab014 (2021)</u> <u>New!</u>



Analysis strategy

- Beam monitors + hadron production experiments
 → neutrino flux
- ND280 measurements

 interaction model
 external constraints
 unoscillated flux × xsec
 - 5 samples at SK $\rightarrow \nu_{\mu}$ disappearance + ν_{e} appearance



Thesis: a joint fit between SK atmospheric and T2K accelerator ν



SK + T2K Joint fit

CP and mass ordering sensitivity

-SK Atmospheric -

-T2K Accelerator



- Resonance in Earth mantle & core sensitive to mass ordering
- Weakly sensitive to $\delta_{\rm CP}$ via normalization of sub-GeV e-like



SK + T2K Joint fit



Systematic correlations

- Overlapped true energy region
 - → coherent interaction model to capture correlations
 - → Bonus: ND constraint for atmospherics!

Same Super-K detector
 used by both experiments
 → estimate contribution from
 detector syst. correlations

Sensitivities



 $\sin \delta > 0$

 χ^2

 δ

11

 $\sin \delta < 0$

 Both experiments complement each othegin many parameters, - T2K — SK (+ND) effect of degeneracy-resolution is quite nice 2.5 • At current statistics, contribution from systematic correlation on sensitivity seems limited 0.5 • This kind of analysis very 0 important for future HyperK True values: $\delta_{CP} = x \text{ axis}$ $\sin^2 \theta_{12} = 0.307 |\Delta m_{31}^2| = 2.509 \times 10^{-10}$ Inverted ordering $\sin^2 \theta_{23} = 0.528 \sin^2 \theta_{13} = 0.0218 \Delta m_{21}^2 = 7.53 \times 10^{-10}$ experiment χ^2

 χ^2

from Neutrino 2022

Efficient statistics for future neutrino oscillation experiments

personal work

hopefully useful for a wider group than just our experiment

Inference

• Toy example: threw coin n = 100 times, got heads x = 40 times,

what is the probability θ of this coin to give heads?

- Probability distribution to generate the data $p(x \mid \theta) = \binom{n}{x} \theta^{x} (1 - \theta)^{n-x}$ seen as function of θ this is called the likelihood $L(\theta \mid x) := p(x \mid \theta)$
- Best agreement = maximum likelihood at $\hat{\theta} = x/n$

better: interval over $\boldsymbol{\theta}$ values to account for statistical fluctuations



Confidence intervals

- Confidence interval: Whatever the true value is, the interval would cover it with at least 68%, 90%, ...
- Define $\chi^2(\theta) := -2 \log L(\theta \mid x)$. Small $\chi^2(\theta)$ in general means better agreement with data.
- Study how much worse than at the minimum we are: $\Delta \chi^2(\theta) := \chi^2(\theta) \chi^2_{\min}$
- Strategy: create a confidence interval using $\Delta \chi^2$

Choose range of θ such that $\Delta \chi^2(\theta) < \Delta \chi_c^2$ with some critical value $\Delta \chi_c^2$



• In the presence of nuisance parameters η , choose e.g. the best agreement

$$\chi_p^2(\theta) = \min_{\eta} \chi^2(\theta, \eta)$$

- Under certain conditions
 - regardless of $p(x \mid \theta)$
 - regardless of $\theta_{\rm true}$ (in the case of profiling also regardless of true nuisance $\eta_{\rm true}$)
 - regardless of *n*
 - in n → ∞ limit, Δχ² evaluated at the true value becomes a chi-squared distribution,
 i.e. p(Δχ²(θ_{true})) = p(χ_k²) where k is dimension of θ
- So can just construct interval by choosing e.g. $\Delta \chi^2(\theta) < 1$ for a 68% interval if k = 1.





- Under certain conditions
 - regardless of $p(x \mid \theta)$
 - regardless of $\theta_{\rm true}$ (in the case of profiling also regardless of true nuisance $\eta_{\rm true}$)
 - regardless of *n*
 - in n → ∞ limit, Δχ² evaluated at the true value becomes a chi-squared distribution, i.e. p(Δχ²(θ_{true})) = p(χ_k²) where k is dimension of θ
- So can just construct interval by choosing e.g. $\Delta \chi^2(\theta) < 1$ for a 68% interval if k = 1.





- Under certain conditions
 - regardless of $p(x \mid \theta)$
 - regardless of $\theta_{\rm true}$ (in the case of profiling also regardless of true nuisance $\eta_{\rm true}$)
 - regardless of *n*
 - in $n \to \infty$ limit, $\Delta \chi^2$ evaluated at the true value becomes a chi-squared distribution, i.e. $p(\Delta \chi^2(\theta_{\text{true}})) = p(\chi_k^2)$ where k is dimension of θ
- So can just construct interval by choosing e.g. $\Delta \chi^2(\theta) < 1$ for a 68% interval if k = 1.





- Under certain conditions
 - regardless of $p(x \mid \theta)$
 - regardless of $\theta_{\rm true}$ (in the case of profiling also regardless of true nuisance $\eta_{\rm true}$)
 - regardless of *n*
 - in n → ∞ limit, Δχ² evaluated at the true value becomes a chi-squared distribution,
 i.e. p(Δχ²(θ_{true})) = p(χ_k²) where k is dimension of θ
- So can just construct interval by choosing e.g. $\Delta \chi^2(\theta) < 1$ for a 68% interval if k = 1.





- Under certain conditions
 - regardless of $p(x \mid \theta)$
 - regardless of $\theta_{\rm true}$ (in the case of profiling also regardless of true nuisance $\eta_{\rm true}$)
 - regardless of *n*
 - in $n \to \infty$ limit, $\Delta \chi^2$ evaluated at the true value becomes a chi-squared distribution, i.e. $p(\Delta \chi^2(\theta_{\text{true}})) = p(\chi_k^2)$ where k is dimension of θ
- So can just construct interval by choosing e.g. $\Delta \chi^2(\theta) < 1$ for a 68% interval if k = 1.







• Small amount of data \rightarrow no $\overleftrightarrow{}$ ν -mode: 119, $\overline{\nu}$ -mode: 16 $\overleftrightarrow{}$ 24



- Small amount of data \rightarrow no ν -mode: 119, $\overline{\nu}$ -mode: 16
- Parameters with boundaries → no
 e.g. -1 ≤ sin δ ≤ 1



- Small amount of data \rightarrow no ν -mode: 119, $\overline{\nu}$ -mode: 16
- Parameters with boundaries → no
 e.g. -1 ≤ sin δ ≤ 1
- (approximate) degeneracies → no
 e.g. mostly sensitive to sin δ
 = degenerate in ±sgn cos δ

In fact an 8-fold degeneracy exists.



- Small amount of data \rightarrow no ν -mode: 119, $\overline{\nu}$ -mode: 16
- Parameters with boundaries → no
 e.g. -1 ≤ sin δ ≤ 1
- (approximate) degeneracies → no
 e.g. mostly sensitive to sin δ
 = degenerate in ±sgn cos δ

In fact an 8-fold degeneracy exists.

 θ contains (effectively) discrete parameters / hypotheses here: mass ordering → no ⁶



Feldman-Cousins

- generate many toy experiments at each $\theta_{\rm true}$
- compute $\Delta\chi^2(\theta_{\rm true})$ for all to obtain its distribution
- get "critical values" s.t. $P(\Delta \chi^2_{\rm true} < \Delta \chi^2_c \mid \theta_{\rm true}) = 68 \%, 90 \%, \cdots$ basically the percentiles as function of $\theta_{\rm true}$



Feldman-Cousins

- generate many toy experiments at each $\theta_{\rm true}$
- compute $\Delta\chi^2(\theta_{\rm true})$ for all to obtain its distribution
- get "critical values" s.t. $P(\Delta \chi^2_{\text{true}} < \Delta \chi^2_c \mid \theta_{\text{true}}) = 68 \%, 90 \%, \cdots$ basically the percentiles as function of θ_{true}



Feldman-Cousins

- generate many toy experiments at each θ_{true}
- compute $\Delta \chi^2(\theta_{\rm true})$ for all to obtain its distribution
- get "critical values" s.t. $P(\Delta \chi_{\text{true}}^2 < \Delta \chi_c^2 \mid \theta_{\text{true}}) = 68\%, 90\%, \cdots$ basically the percentiles as function of $\theta_{\rm true}$



Feldman-Cousins

Neyman-Construction using $\Delta \chi^2$ as ordering principle

- generate many toy experiments at each $\theta_{\rm true}$
- compute $\Delta \chi^2(\theta_{\rm true})$ for all to obtain its distribution
- get "critical values" s.t. $P(\Delta \chi_{\text{true}}^2 < \Delta \chi_c^2 \mid \theta_{\text{true}}) = 68\%, 90\%, \cdots$ basically the percentiles as function of $\theta_{\rm true}$



T2K Preliminary



true

Feldman-Cousins

- generate many toy experiments at each θ_{true}
- compute $\Delta \chi^2(\theta_{\rm true})$ for all to obtain its distribution
- get "critical values" s.t. $P(\Delta \chi_{\text{true}}^2 < \Delta \chi_c^2 \mid \theta_{\text{true}}) = 68\%, 90\%, \cdots$ basically the percentiles as function of $\theta_{\rm true}$
- construct confidence interval by $\Delta \chi^2(\theta) < \Delta \chi_c^2(\theta)$





Feldman-Cousins

Neyman-Construction using $\Delta \chi^2$ as ordering principle

- generate many toy experiments at each θ_{true}
- compute $\Delta \chi^2(\theta_{\rm true})$ for all to obtain its distribution
- get "critical values" s.t. $P(\Delta \chi_{\text{true}}^2 < \Delta \chi_c^2 \mid \theta_{\text{true}}) = 68\%, 90\%, \cdots$ basically the percentiles as function of $\theta_{\rm true}$
- construct confidence interval by $\Delta \chi^2(\theta) < \Delta \chi^2_c(\theta)$

30





Critical values for δ_{CP}



from L. Berns, Moriond 2021

Problem

• Current experiments (T2K) started excluding $\delta_{\rm CP}$ ranges at 3σ CL. Future experiments (HK, DUNE) aiming for 5σ discovery of CP violation in leptons.



J. Wilson, Neutrino 2022

M. Muether, Neutrino 2022

Problem

- Current experiments (T2K) started excluding $\delta_{\rm CP}$ ranges at 3σ CL. Future experiments (HK, DUNE) aiming for 5σ discovery of CP violation in leptons.
- By definition, to obtain high confidence-level critical values in Feldman-Cousins method, need many toy experiments:
 - $3\sigma : \gg 370$ $4\sigma : \gg 16k$ $5\sigma : \gg 1.7M$

each fit is complicated with $\mathcal{O}(100)$ parameters often having non-linear responses, so this is not very practical.

Problem

- Current experiments (T2K) started excluding $\delta_{\rm CP}$ ranges at 3σ CL. Future experiments (HK, DUNE) aiming for 5σ discovery of CP violation in leptons.
- By definition, to obtain high confidence-level critical values in Feldman-Cousins method, need many toy experiments:
 - $3\sigma : \gg 370$ $4\sigma : \gg 16k$ $5\sigma : \gg 1.7M$

each fit is complicated with $\mathcal{O}(100)$ parameters often having non-linear responses, so this is not very practical.

- Now even if we produce say 100M toy experiments to get 5σ critical values, most of the toy experiments will be around $1\sigma \sim 2\sigma$ and "wasted".
- Can we do better? Can we specifically generate toys that matter at high CL?

• Let's preferentially sample all low-probability toys "increase the temperature T "

e.g.
$$p_{\text{smpl}}(x) := [p(x \mid \theta_{\text{true}})]^{1/T}$$

and weight each toy as $w(x) = \frac{p(x \mid \theta_{\text{true}})}{p_{\text{smpl}}(x)}$

• So e.g. for each bin *i*, sample toys with $T\sqrt{\lambda_i}$ error instead of usual Poisson error $\sqrt{\lambda_i}$

 (\times)

• Let's preferentially sample all low-probability toys "increase the temperature T "

e.g.
$$p_{\text{smpl}}(x) := [p(x \mid \theta_{\text{true}})]^{1/T}$$

and weight each toy as $w(x) = \frac{p(x \mid \theta_{\text{true}})}{p_{\text{smpl}}(x)}$

• So e.g. for each bin *i*, sample toys with $T\sqrt{\lambda_i}$ error instead of usual Poisson error $\sqrt{\lambda_i}$



e.g.
$$p_{\text{smpl}}(x) := [p(x \mid \theta_{\text{true}})]^{1/T}$$

and weight each toy as $w(x) = \frac{p(x \mid \theta_{\text{true}})}{p_{\text{smpl}}(x)}$

- So e.g. for each bin *i*, sample toys with $T\sqrt{\lambda_i}$ error instead of usual Poisson error $\sqrt{\lambda_i}$
- This gives you toys with larger values of χ^2



e.g.
$$p_{\text{smpl}}(x) := [p(x \mid \theta_{\text{true}})]^{1/T}$$

and weight each toy as $w(x) = \frac{p(x \mid \theta_{\text{true}})}{p_{\text{smpl}}(x)}$

- So e.g. for each bin *i*, sample toys with $T\sqrt{\lambda_i}$ error instead of usual Poisson error $\sqrt{\lambda_i}$
- This gives you toys with larger values of χ^2
- But! almost all of them have small $\Delta \chi^2(\theta_{\rm true})$, because there are $N_{\rm bins}$ directions in which we can increase the temperature, but only *k* (non-linear) directions of them contribute to a change of $\Delta \chi^2(\theta)$, the rest will only increase $\chi^2_{\rm min}$



e.g.
$$p_{\text{smpl}}(x) := [p(x \mid \theta_{\text{true}})]^{1/T}$$

and weight each toy as $w(x) = \frac{p(x \mid \theta_{\text{true}})}{p_{\text{smpl}}(x)}$

- So e.g. for each bin *i*, sample toys with $T\sqrt{\lambda_i}$ error instead of usual Poisson error $\sqrt{\lambda_i}$
- This gives you toys with larger values of χ^2
- But! almost all of them have small $\Delta \chi^2(\theta_{\rm true})$, because there are $N_{\rm bins}$ directions in which we can increase the temperature, but only *k* (non-linear) directions of them contribute to a change of $\Delta \chi^2(\theta)$, the rest will only increase $\chi^2_{\rm min}$
- In fact, because now we have to weight each toy x with w(x), any variation in w(x) for the same value of $\Delta \chi^2$ will actually "cost" us toy statistics, so this is even worse than before, not only at low $\Delta \chi^2$, but also at high $\Delta \chi^2$



e.g.
$$p_{\text{smpl}}(x) := [p(x \mid \theta_{\text{true}})]^{1/T}$$

and weight each toy as $w(x) = \frac{p(x \mid \theta_{\text{true}})}{p_{\text{smpl}}(x)}$

- So e.g. for each bin *i*, sample toys with $T\sqrt{\lambda_i}$ error instead of usual Poisson error $\sqrt{\lambda_i}$
- This gives you toys with larger values of χ^2
- But! almost all of them have small $\Delta \chi^2(\theta_{\rm true})$, because there are $N_{\rm bins}$ directions in which we can increase the temperature, but only *k* (non-linear) directions of them contribute to a change of $\Delta \chi^2(\theta)$, the rest will only increase $\chi^2_{\rm min}$
- In fact, because now we have to weight each toy x with w(x), any variation in w(x) for the same value of $\Delta \chi^2$ will actually "cost" us toy statistics, so this is even worse than before, not only at low $\Delta \chi^2$, but also at high $\Delta \chi^2$



• What does high $\Delta \chi^2(\theta_{\rm true})$ mean?

→ there exists a different value $\hat{\theta}$ from which one is more likely to obtain the data *x* than from θ_{true}

• What does high $\Delta \chi^2(\theta_{\rm true})$ mean?

→ there exists a different value $\hat{\theta}$ from which one is more likely to obtain the data *x* than from θ_{true}



- What does high Δχ²(θ_{true}) mean?
 → there exists a different value θ̂ from which one is more likely to obtain the data x than from θ_{true}
- Just use all toys sampled from many $\theta_{\rm smpl}$ at once

$$p_{\text{smpl}}(x) = p(x \mid \{\theta\}) := \frac{1}{S} \sum_{s=1}^{S} p(x \mid \theta_s)$$

(mixture model)



44

- What does high Δχ²(θ_{true}) mean?
 → there exists a different value θ̂ from which one is more likely to obtain the data x than from θ_{true}
- Just use all toys sampled from many θ_{smpl} at once

$$p_{\text{smpl}}(x) = p(x \mid \{\theta\}) := \frac{1}{S} \sum_{s=1}^{S} p(x \mid \theta_s)$$

(mixture model)

As long as θ_s are close enough, there will always be some θ_s close enough to $\hat{\theta}$

$$\begin{aligned} \forall x, \Delta \chi^2(\theta_{\text{target}}) \geq \Delta \chi_0^2 \\ \frac{dN_{\text{toys}}(x \mid \{\theta\})}{dN_{\text{toys}}(x \mid \theta_{\text{target}})} &= \frac{\sum_s p(x \mid \theta_s)}{p(x \mid \theta_{\text{target}})} \\ \geq \frac{p(x \mid \hat{\theta})}{p(x \mid \theta_{\text{target}})} \\ &= \exp\left[\frac{1}{2}\Delta \chi_0^2\right] \end{aligned}$$



all toys will be sampled more often, high- $\Delta \chi^2$ exponentially more often

- What does high Δχ²(θ_{true}) mean?
 → there exists a different value θ̂ from which one is more likely to obtain the data *x* than from θ_{true}
- Just use all toys sampled from many $\theta_{\rm smpl}$ at once

$$p_{\text{smpl}}(x) = p(x \mid \{\theta\}) := \frac{1}{S} \sum_{s=1}^{S} p(x \mid \theta_s)$$

(mixture model)

• As long as θ_s are close enough, there will always be some θ_s close enough to $\hat{\theta}$, such that the weights w will be bounded from above resulting in numerical stability (at least for the $\Delta \chi^2$ values we care about).

$$w(x) \lesssim S \exp\left[-\frac{1}{2}\Delta\chi^2(\theta_{\text{target}})\right]$$

with

$$w(x) = \frac{p(x \mid \theta_{\text{target}})}{p(x \mid \{\theta\})}$$



- What does high Δχ²(θ_{true}) mean?
 → there exists a different value θ̂ from which one is more likely to obtain the data x than from θ_{true}
- Just use all toys sampled from many $\theta_{\rm smpl}$ at once

$$p_{\text{smpl}}(x) = p(x \mid \{\theta\}) := \frac{1}{S} \sum_{s=1}^{S} p(x \mid \theta_s)$$

(mixture model)

• As long as θ_s are close enough, there will always be some θ_s close enough to $\hat{\theta}$, such that the weights w will be bounded from above resulting in numerical stability (at least for the $\Delta \chi^2$ values we care about).

$$w(x) \lesssim S \exp\left[-\frac{1}{2}\Delta\chi^2(\theta_{\text{target}})\right]$$

with

$$w(x) = \frac{p(x \mid \theta_{\text{target}})}{p(x \mid \{\theta\})}$$



Critical value distributions

- The same toys are used to produce both results!
- Errors for critical values significantly reduced (e.g. 3σ)
- 5σ were impossible to estimate before, now easy





Critical value distributions

- The same toys are used to produce both results!
- Errors for critical values significantly reduced (e.g. 3σ)
- 5σ were impossible to estimate before, now easy
- Bonus: can now interpolate critical values for free!



Nuisance parameters

- Feldman-Cousins paper itself does not deal with presence of nuisance parameters, various approaches exist
- As prior distribution (prior Highland-Cousins)

Just define $p(x \mid \theta) := \int d\eta \, p(\eta) \, p(x \mid \theta, \eta)$

and use this both for fitting and throwing toys, i.e. fit by marginalizing.

In most cases using profile likelihood is similar to marginal likelihood so should be ok for using in fit too (with constraint term corresponding to prior).

Performance improvement same.

If systematics are not well constrained by data,

can also compute toy weight at fixed η instead of marginalizing, which is much cheaper.

i.e. $w(x, \eta_{\text{true}}) = \frac{p(x \mid \theta_{\text{target}}, \eta_{\text{true}})}{p(x \mid \{\theta\}, \eta_{\text{true}})}$ instead of $w(x) = \frac{p(x \mid \theta_{\text{target}})}{p(x \mid \{\theta\})}$

 For prior unconstrained parameters and parameters that get very constrained by the data this is not ideal. Use posterior Highland-Cousins conditioned on the true value of the parameter of interest θ:

i.e. define toy distribution by $p(x \mid \theta) := \int d\eta \, p(\eta \mid x_{data}, \theta) \, p(x \mid \theta, \eta)$

and fit by **profiling**. Here the constraint term should be thrown for each toy experiment (included in *x*) as though it is external data.

In this case we get the same performance improvement *if* for fixed θ the η -posterior is **sufficiently gaussian**.

I do not know yet, if one can prove any properties if η-posterior is not sufficiently gaussian, and method may not be directly applicable. However, in this case the Feldman-Cousins method will be dependent on the nuisance parameter treatment choice (similar to priors in Bayesian analysis) so more careful analysis would be necessary anyway. If the number of such non-gaussian nuisances is small, they could be included as part of the parameters of interest.

Summary

- Neutrino oscillation experiments aiming to discover CP violation in leptonic sector
- Frequentist analysis of neutrino oscillation requires expensive Feldman-Cousins procedure because Wilk's theorem cannot be used
- Proposed a simple technique that efficiently generates high- $\Delta \chi^2$ toys, resulting in exponential reduction of uncertainty for high- $\Delta \chi^2$
 - + provides correct interpolation of critical values for free
- Hopefully this can help many experiments, not only ν -osc.
- Exploring various other techniques as well, hopefully Feldman-Cousins will not be considered an expensive calculation soon

"Simply use a mixture of toys generated at many θ values and reweight"

Similar to "multiple histogram reweighting" in statistical mechanics / lattice QCD calculations see e.g. Kari Rummukainen, "Monte Carlo simulation methods" <u>lecture notes on reweighting</u>

I think some Higgs p-value studies (by K. Cranmer's group) also used similar methods.

to be posted on arXiv soon!