# Histogram Binning with Bayesian Blocks

Brian Pollack, Northwestern University

8/3/17

*Coauthors: Sapta Bhattacharya, Michael Schmitt*

*arXiv:* *1708.00810*
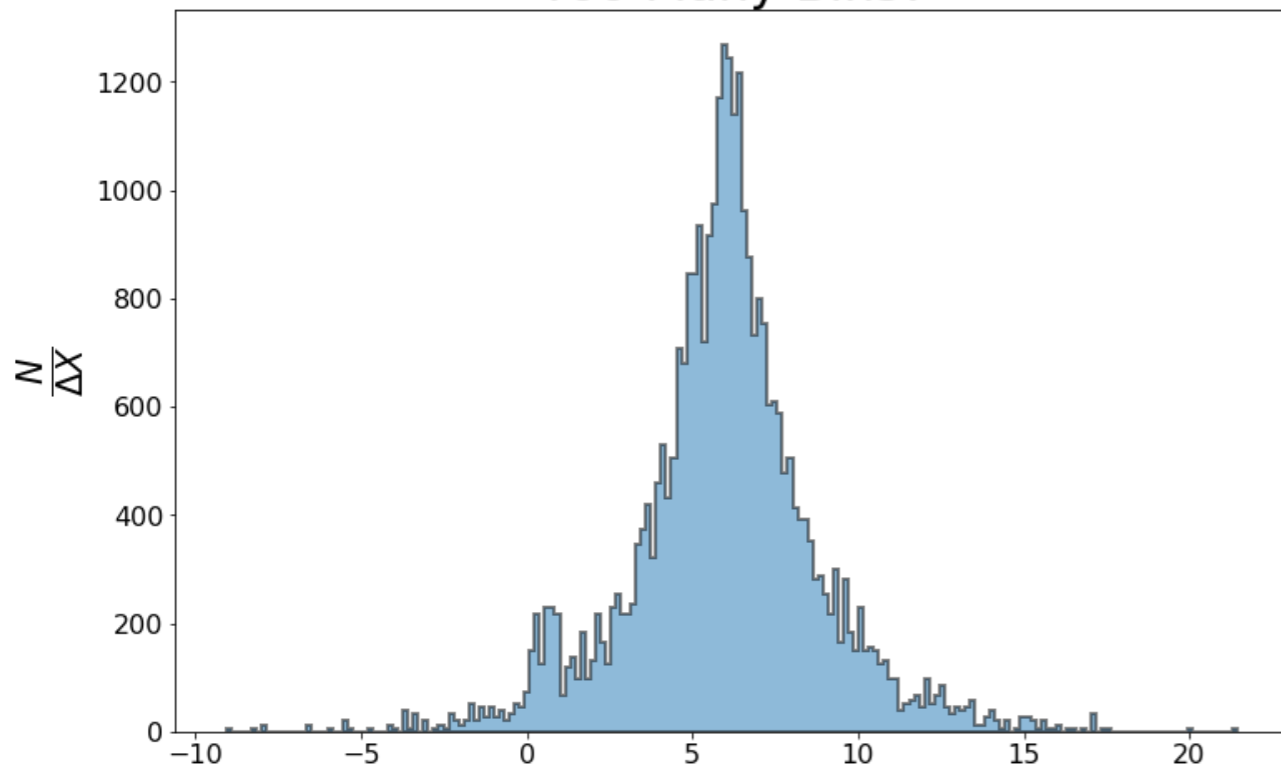
# How Do We Bin?

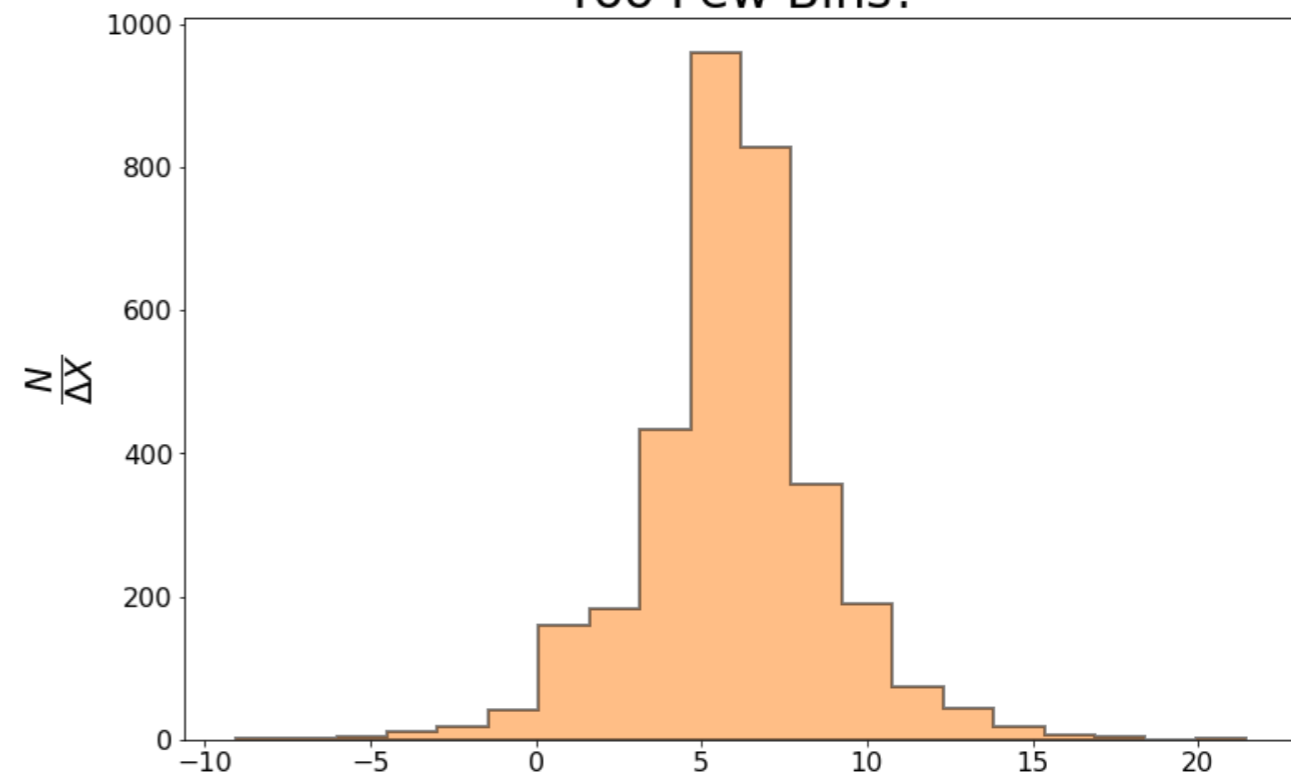★ **Histogram binning is usually *arbitrary*.**

- Number of bins → Whatever seems to look reasonable.

- Too many bins → Statistical fluctuations obscure structure.

- Too few bins → Small structures are swallowed by background.

★ **Bayesian Blocks (BB) chooses 'best' number of blocks (bins), and 'best' choice for bin edges.**
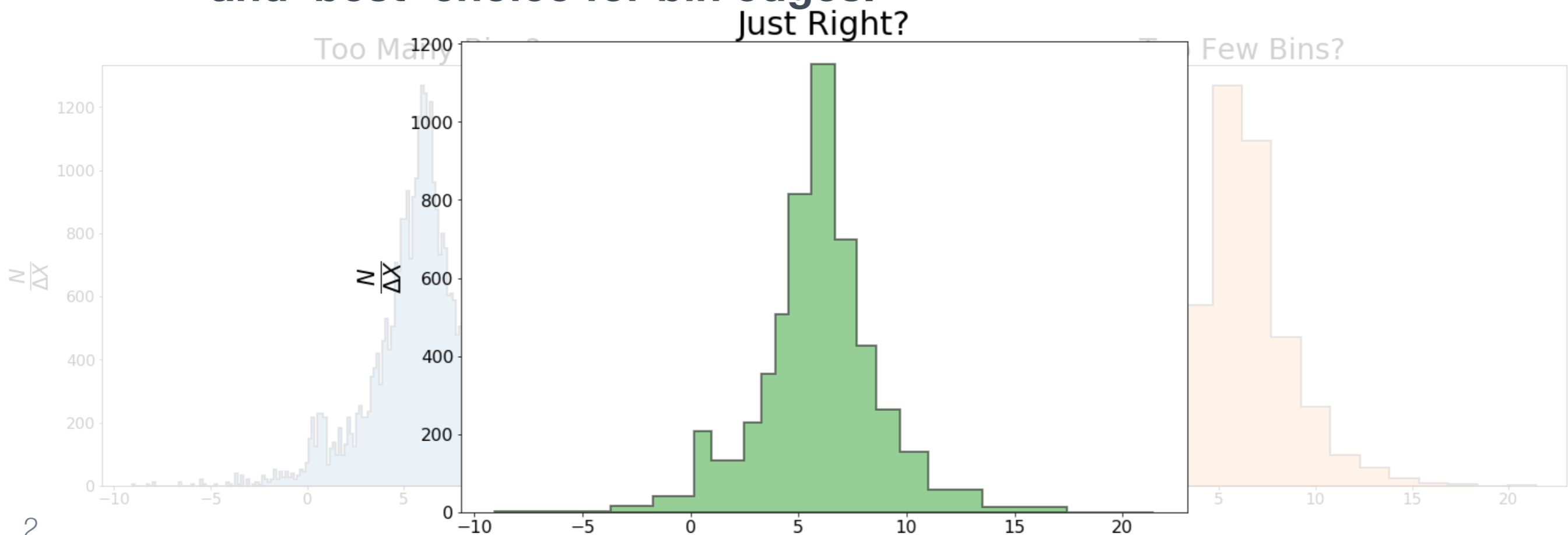
# How Do We Bin?

★ **Histogram binning is usually *arbitrary*.**

- Number of bins → Whatever seems to look reasonable.

- Too many bins → Statistical fluctuations obscure structure.

- Too few bins → Small structures are swallowed by background.

★ **Bayesian Blocks (BB) chooses 'best' number of blocks (bins), and 'best' choice for bin edges.**
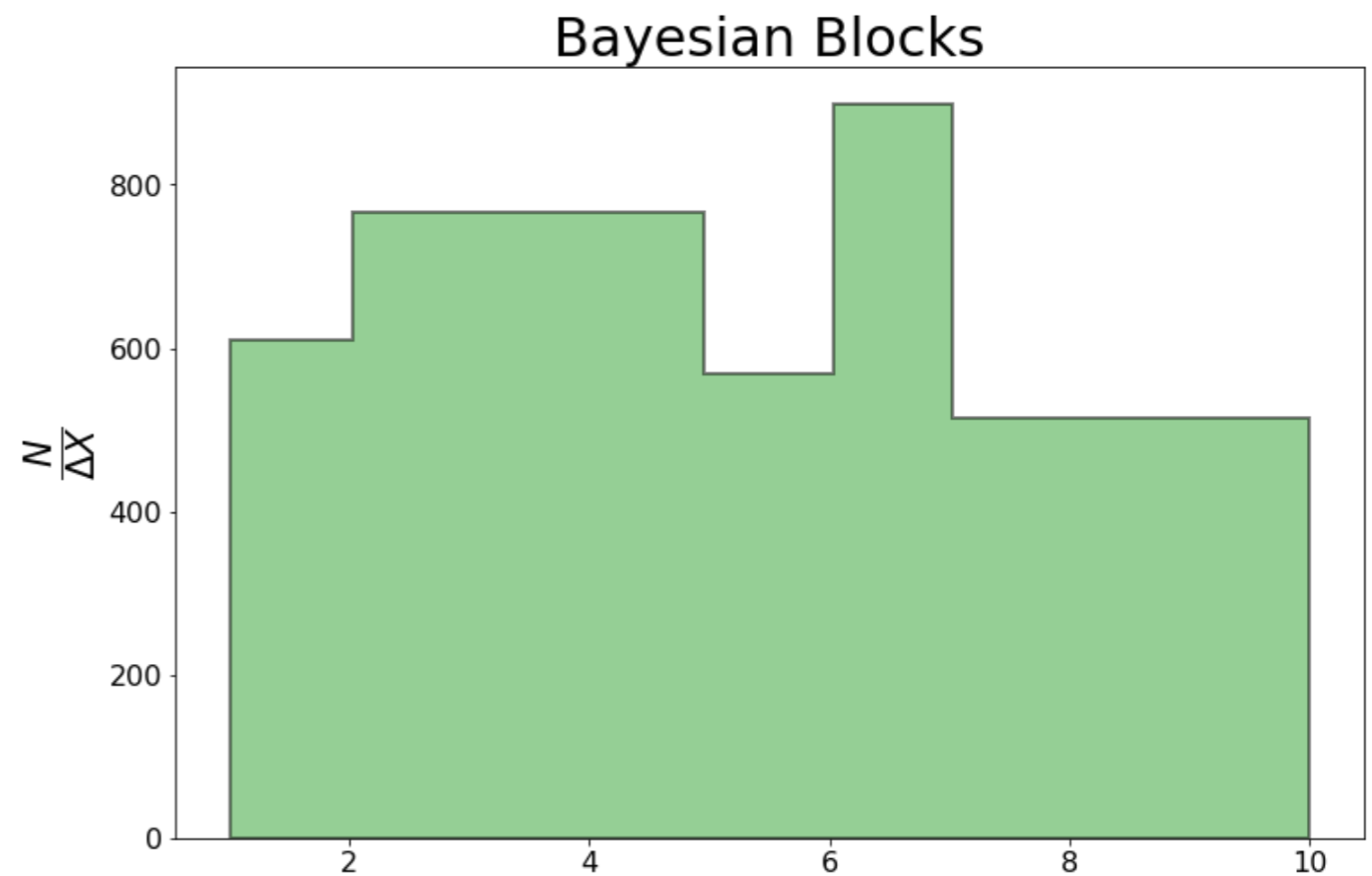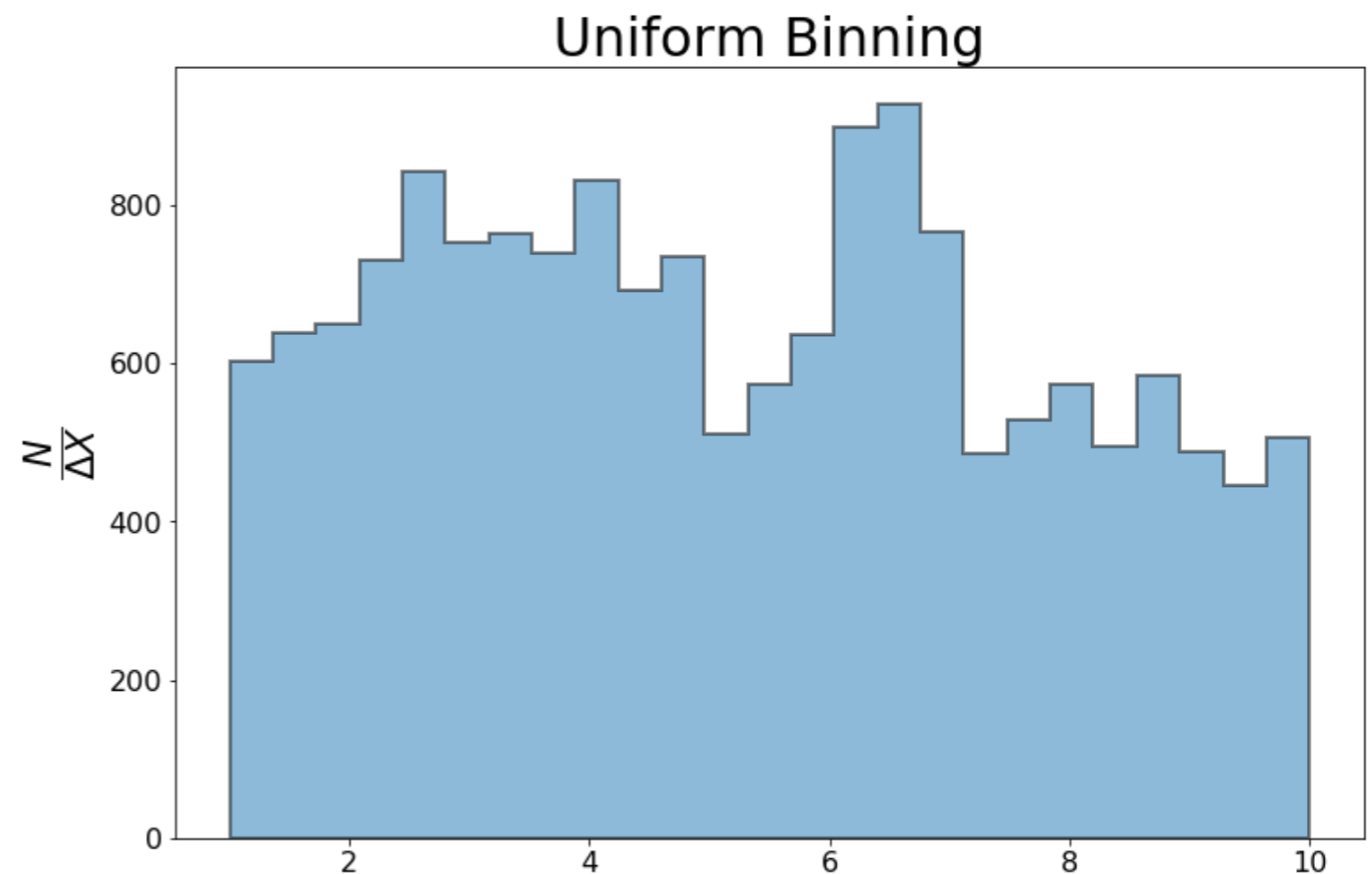
# Bayesian Blocks

★ **Input:**

- Data

- False-positive rate (tuning parameter)

★ **Output:**

- Bin Edges

★ **Each edge is statistically significant**

- New edge → change in underlying pdf



*Underlying pdfs: 3 Uniform distributions*

# Bayesian Blocks
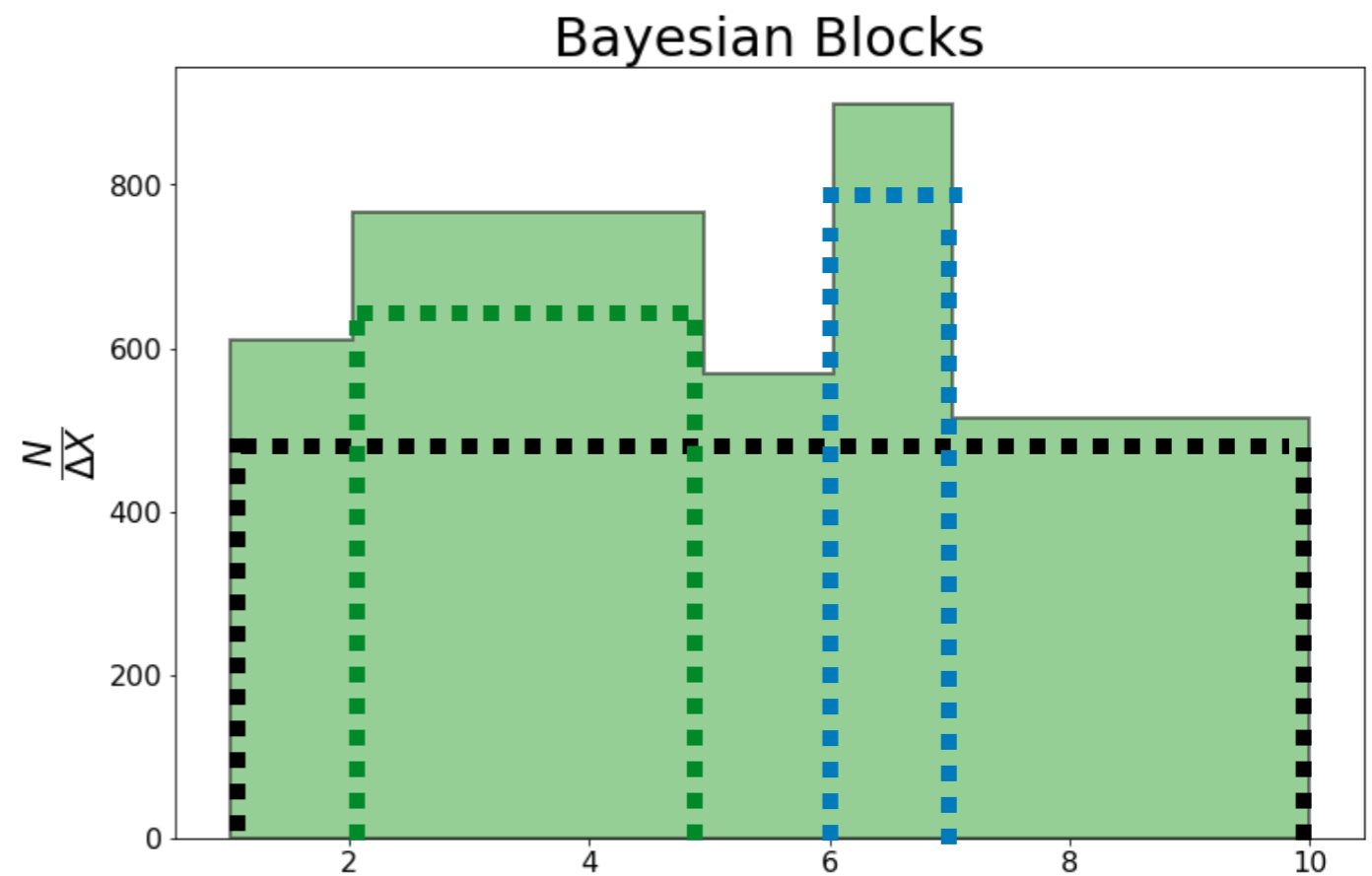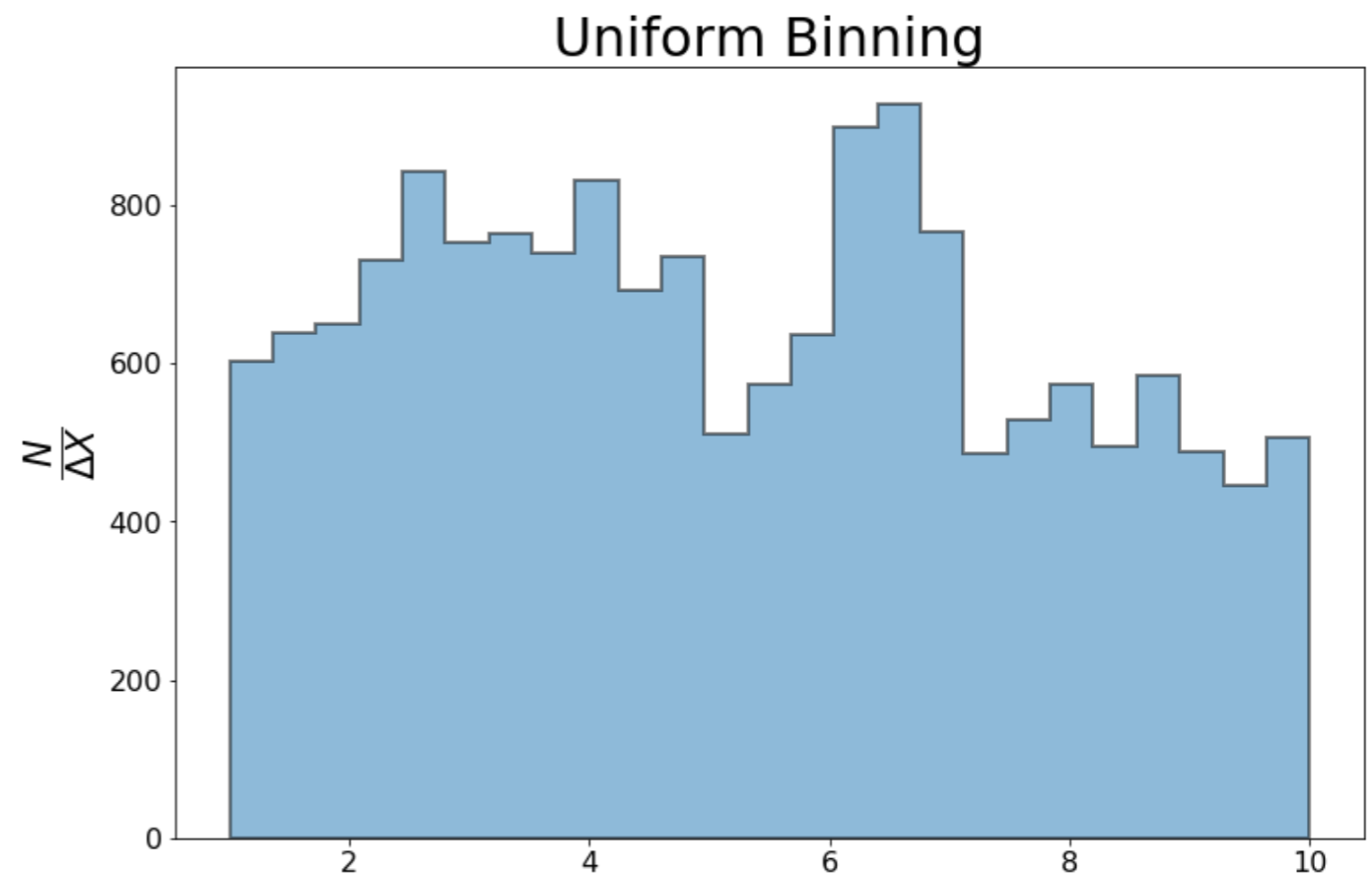

Uniform Binning


Bayesian Blocks

★ **Input:**

- Data

- False-positive rate (tuning parameter)

★ **Output:**

- Bin Edges

★ **Each edge is statistically significant**

- New edge → change in underlying pdf
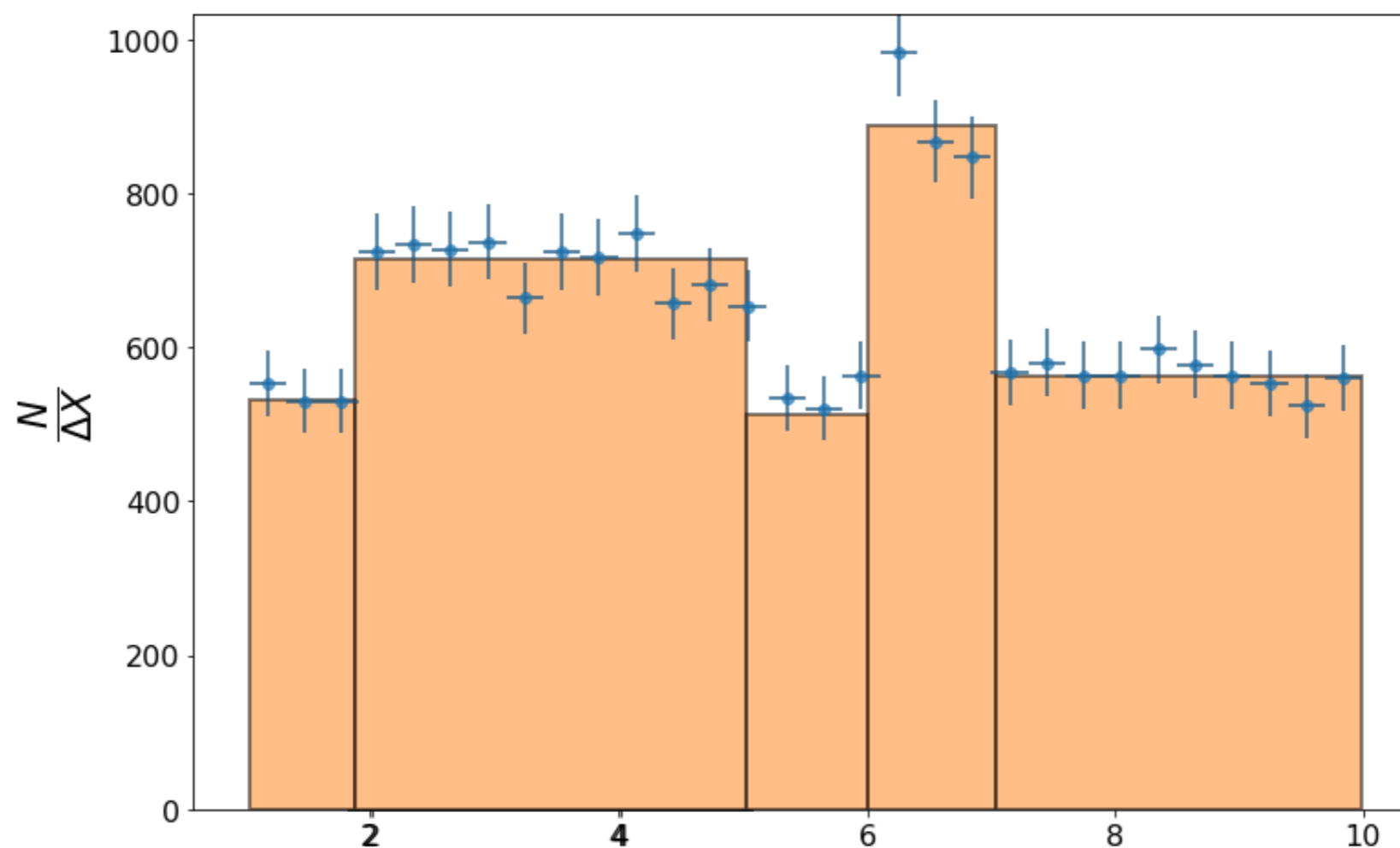
*Underlying pdfs: 3 Uniform distributions*

# Bayesian Blocks

★ **Developed by J. D. Scargle et. al.\*, for use with time-series data in astronomy.**

★ **Goal: characterize statistically significant variations in data.**

- Accomplish via <u>optimal segmentation</u> using non-parametric modeling.

  ✦ Each segment treated as histogram bin (bins have variable widths).

  ✦ Each segment associated with uniform distribution.

  ✦ Combination of data and uniform distributions → calculation of **fitness function**.

★ **Finding maximal fitness function requires clever programming, not feasible to use naive (brute force) methods.**

- For $N$ data points, $2^N$ possible binnings → untenable for large $N$
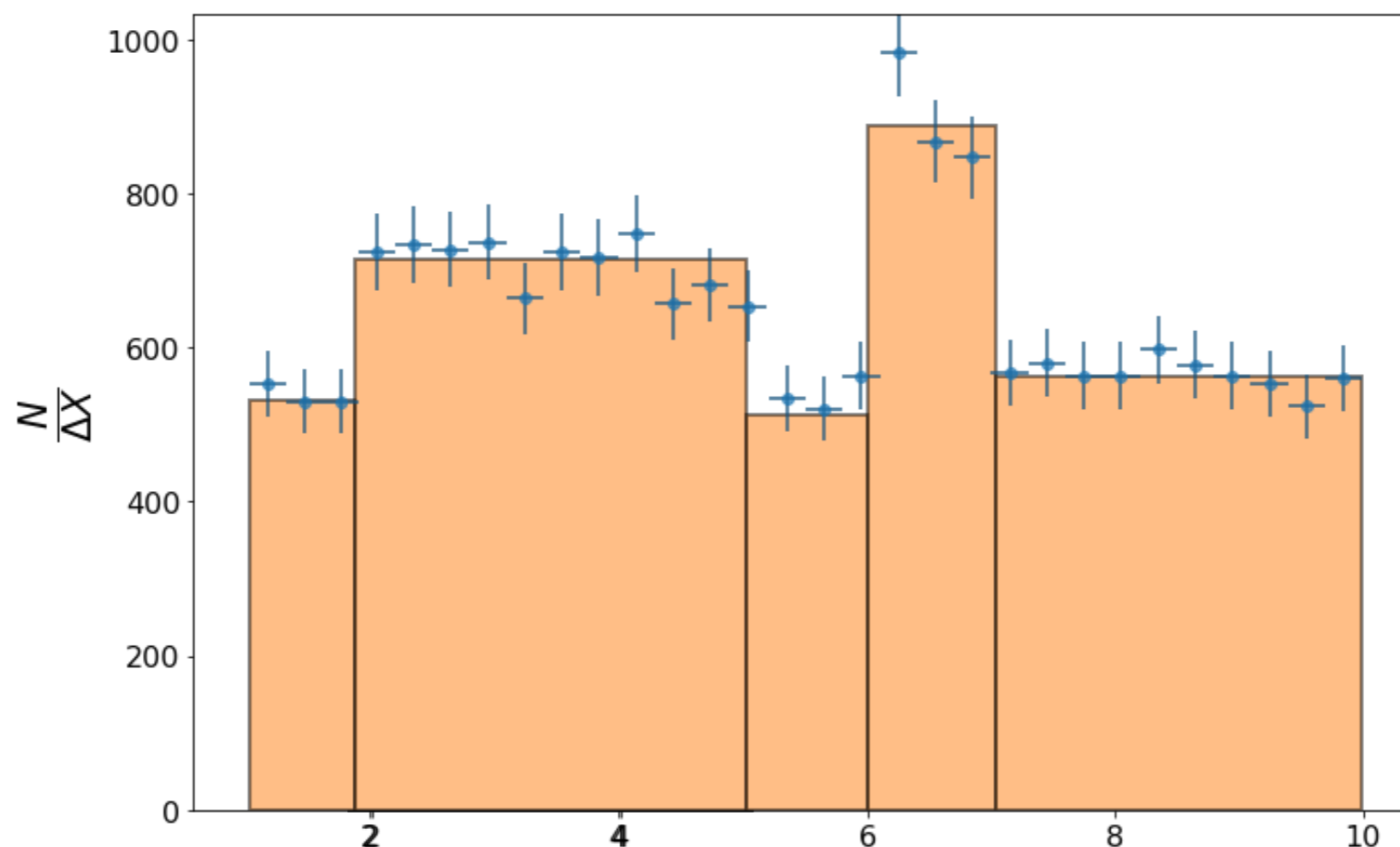
# The Fitness Function

★ **The <u>Fitness Function</u> is a quantity that is maximized when the optimal segmentation of a dataset is achieved.**

# The Fitness Function

★ **The <u>Fitness Function</u> is a quantity that is maximized when the optimal segmentation of a dataset is achieved.**

★ **For *K* bins, the total fitness, $F_{total}$, can be defined as the sum of the fitnesses of each bin, *f(Bᵢ)*:**
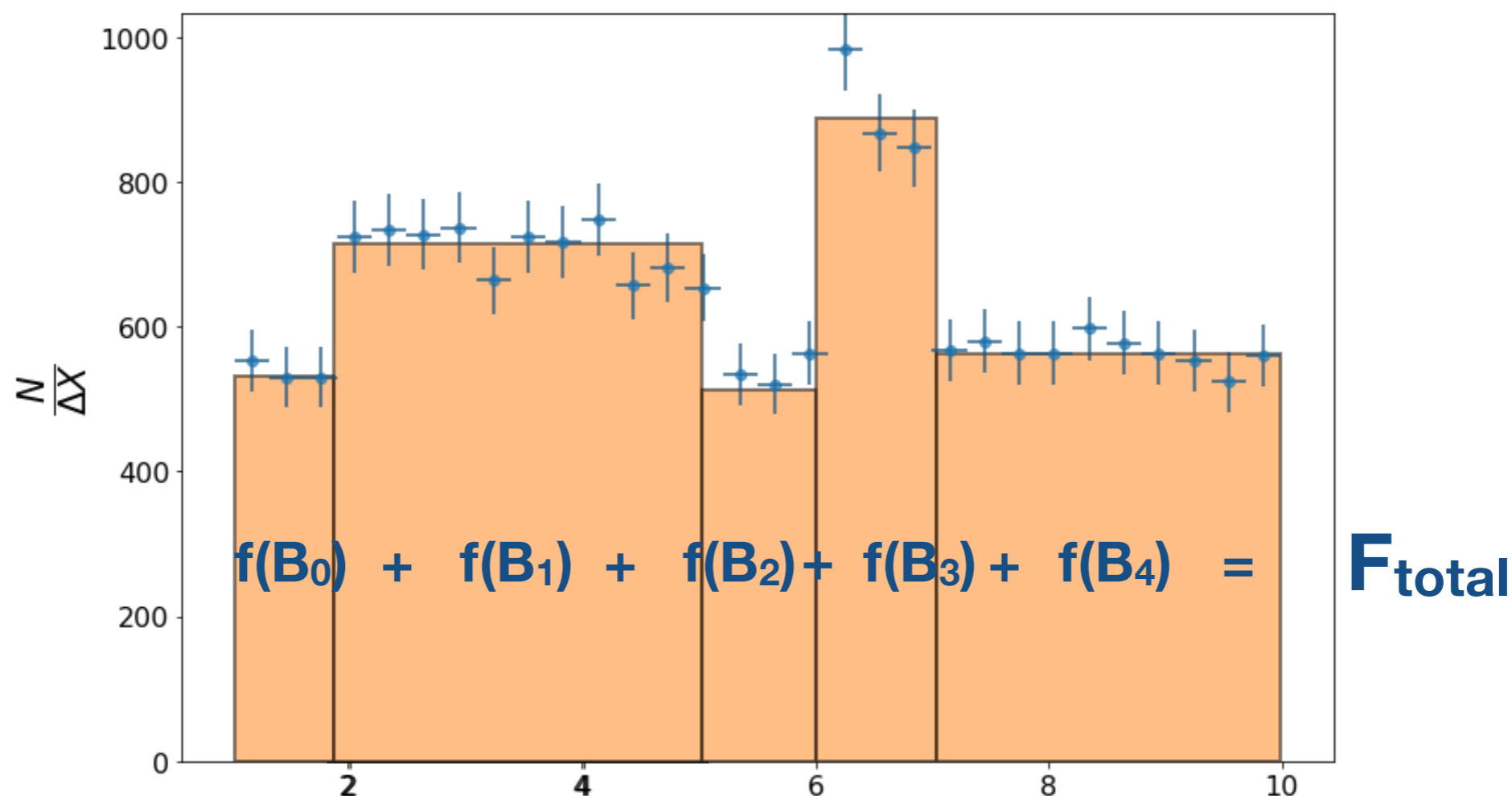
$$F_{total} = \sum_{i=0}^{K} f(B_i)$$
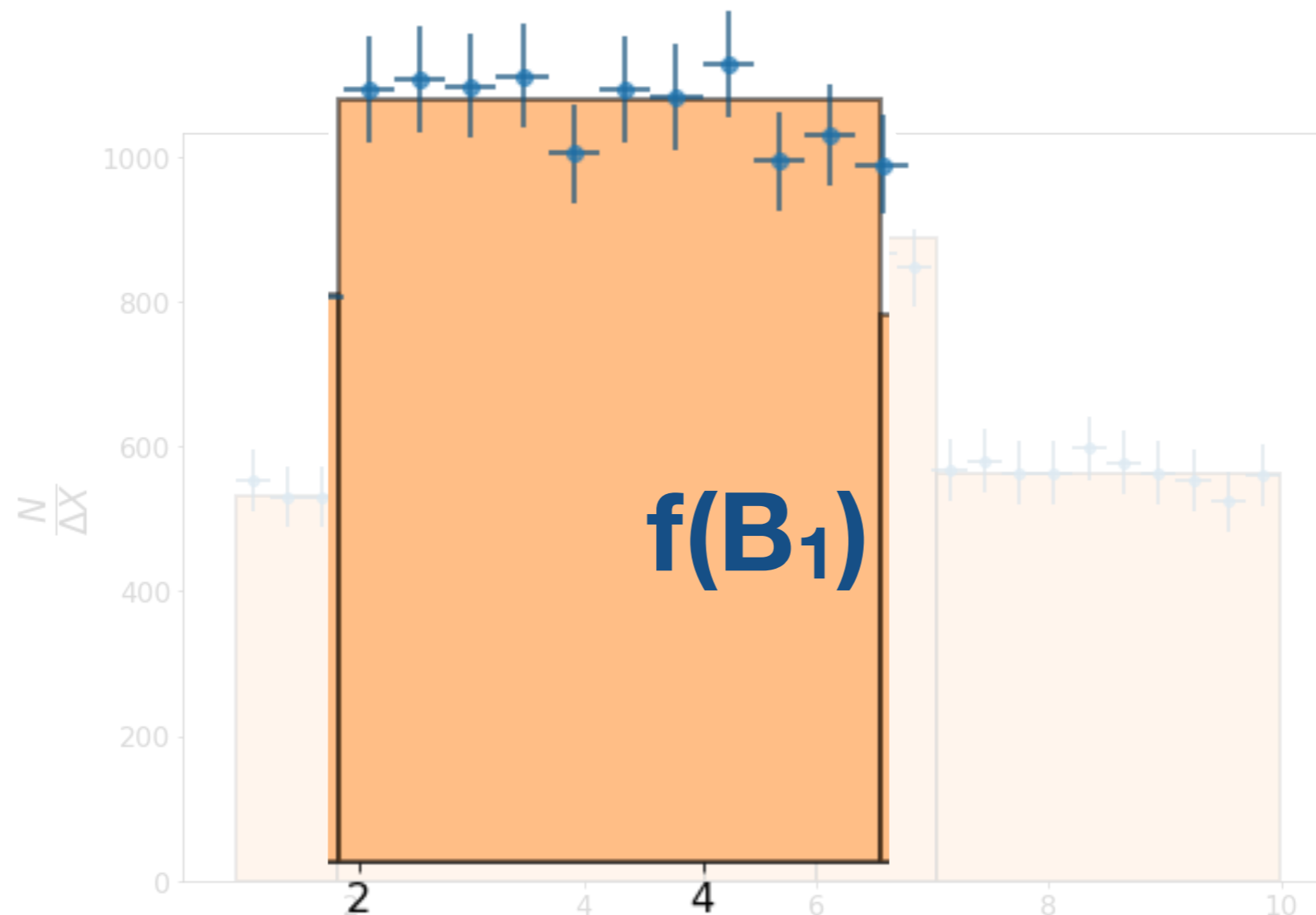
# The Fitness Function

★ **The <u>Fitness Function</u> is a quantity that is maximized when the optimal segmentation of a dataset is achieved.**

★ **For *K* bins, the total fitness, *F*<sub>total,</sub> can be defined as the sum of the fitnesses of each bin, *f(B<sub>i</sub>)*:**

$$F_{total} = \sum_{i=0}^{K} f(B_i)$$



f(B₀)  +  f(B₁)  +  f(B₂)+ f(B₃) +  f(B₄)   =   **F**<sub>total</sub>

# The Fitness Function

**The fitness, $f(B_i)$, of <u>each bin</u> can be treated as a log-likelihood, assuming the events in each bin follow a Poisson distribution.**

# The Fitness Function

**The fitness, $f(B_i)$, of <u>each bin</u> can be treated as a log-likelihood, assuming the events in each bin follow a Poisson distribution.**

$$P_{dx} = \lambda(x)dx \times e^{-\lambda(x)dx} \rightarrow \text{probability for an infinitesimal bin.}$$

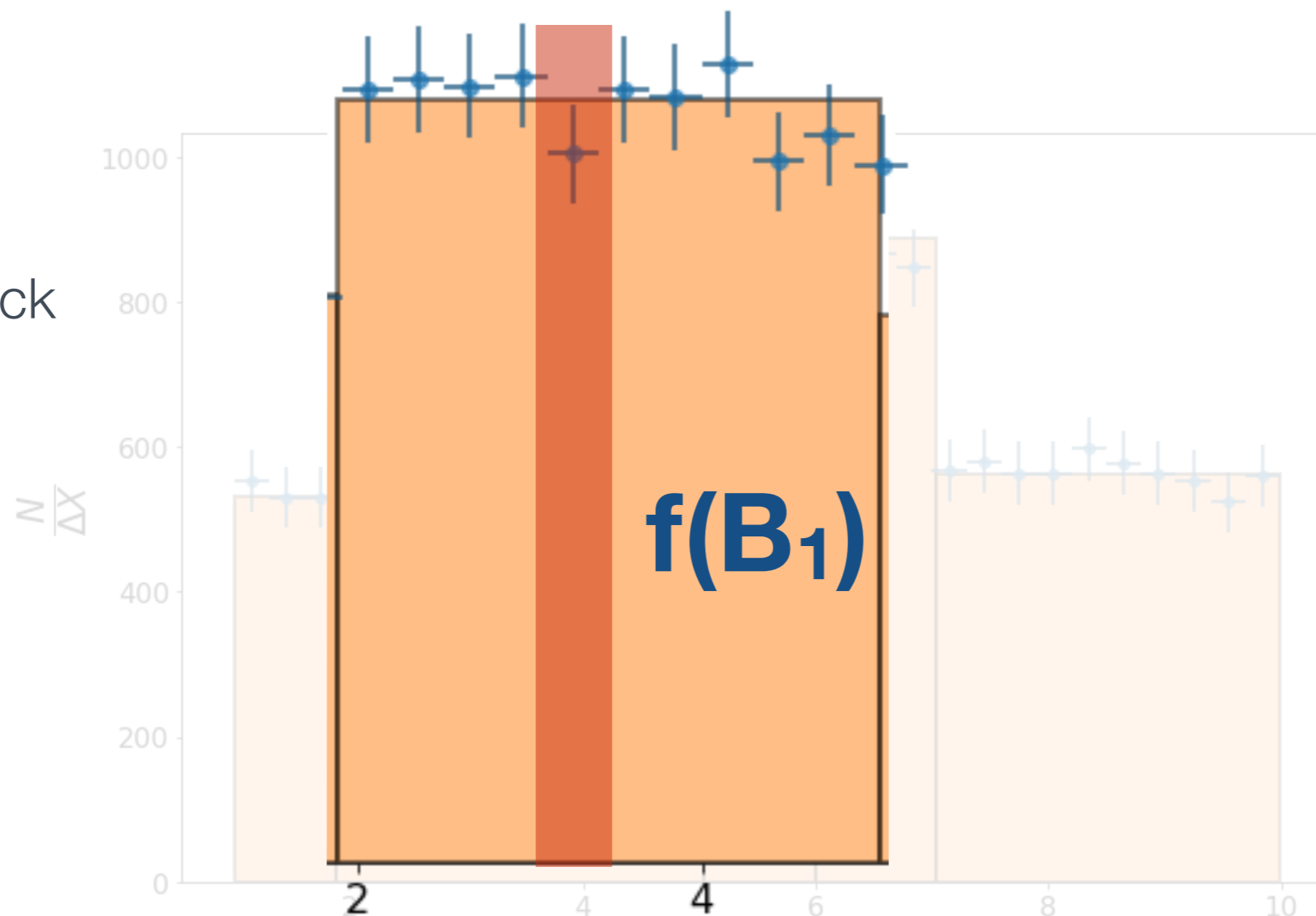$\lambda$: amplitude

$x$: width of block

# The Fitness Function

**The fitness, *f(B$_i$)*, of <u>each bin</u> can be treated as a log-likelihood, assuming the events in each bin follow a Poisson distribution.**

$$P_{dx} = \lambda(x)dx \times e^{-\lambda(x)dx} \rightarrow \text{probability for an infinitesimal bin.}$$

$$\ln L_B = \sum^n \ln \lambda(x) + \sum^n \ln dx - \int \lambda(x)dx \rightarrow \text{log-likelihood for an entire bin.}$$

λ: amplitude

*x*: width of block

*n*: number of events
   in a bin



**f(B₁)**

# The Fitness Function

**The fitness, $f(B_i)$, of <u>each bin</u> can be treated as a log-likelihood, assuming the events in each bin follow a Poisson distribution.**

$$P_{dx} = \lambda(x)dx \times e^{-\lambda(x)dx} \rightarrow \text{probability for an infinitesimal bin.}$$
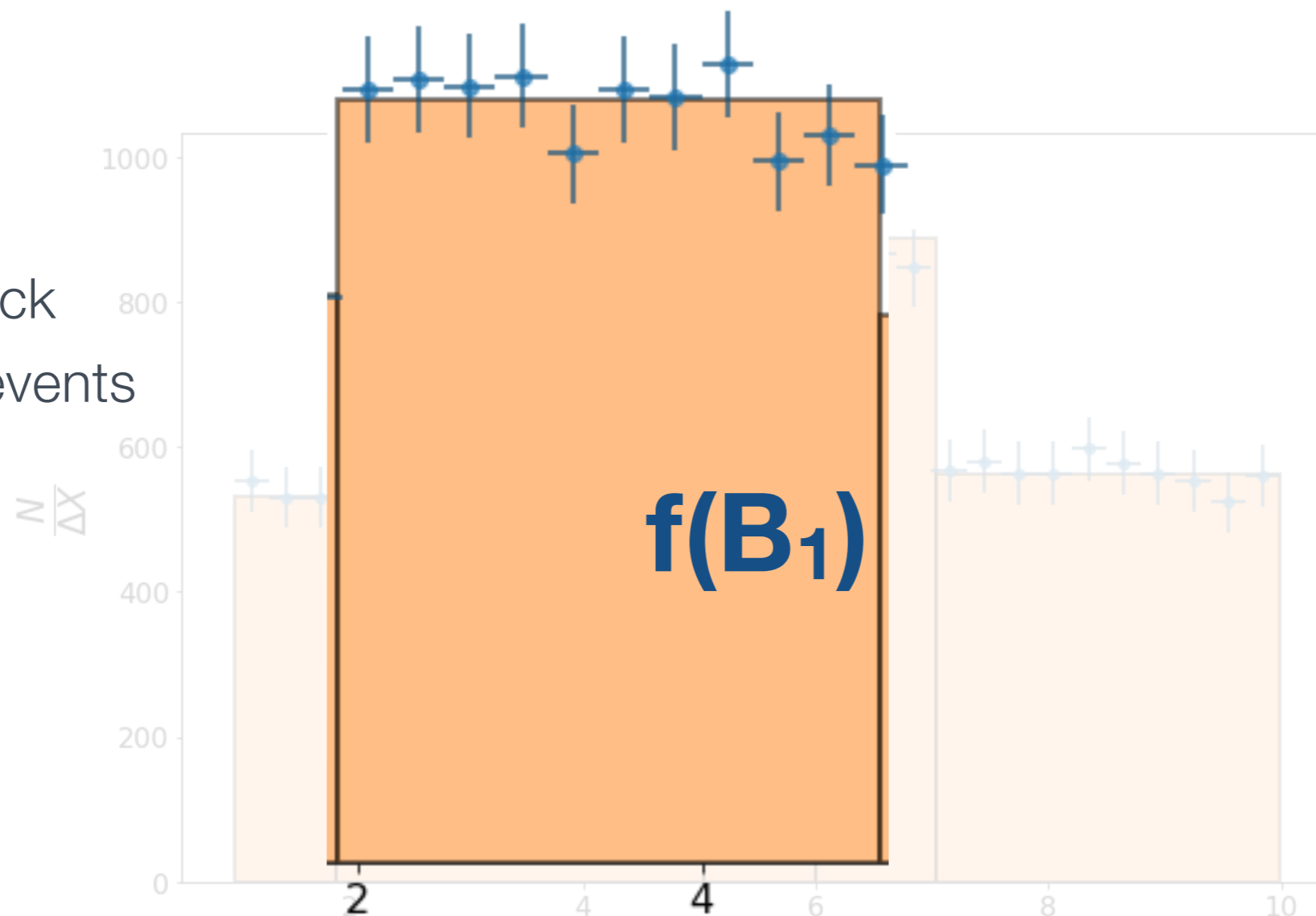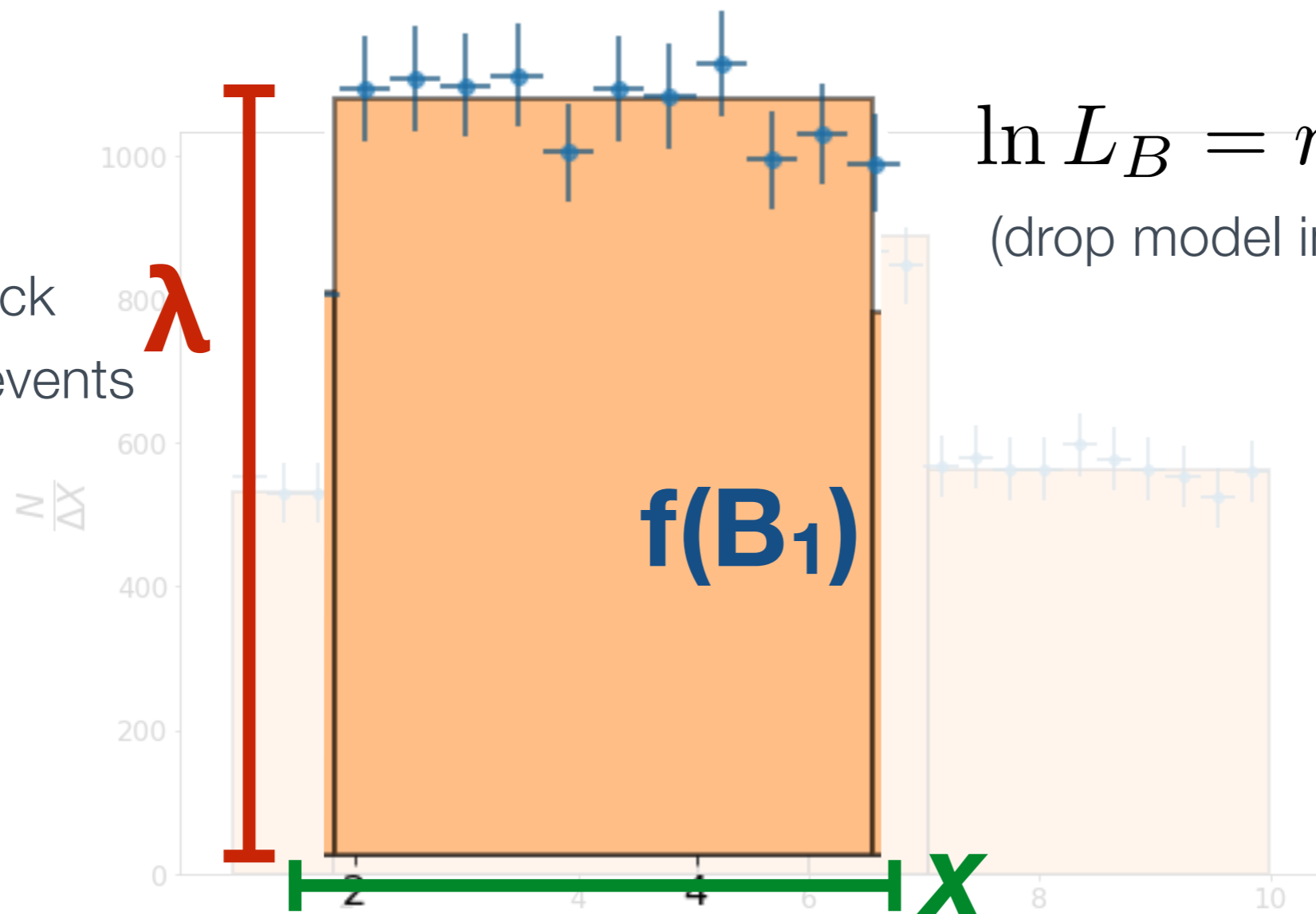
$$\ln L_B = \sum^n \ln \lambda(x) + \sum^n \ln dx - \int \lambda(x)dx \rightarrow \text{log-likelihood for an entire bin.}$$

λ: amplitude

$x$: width of block

$n$: number of events
    in a bin

$$\ln L_B = n \ln \lambda - \lambda x$$

(drop model independent terms)

**f(B₁)**

# The Fitness Function

**The fitness, *f(B$_i$)*, of __each bin__ can be treated as a log-likelihood, assuming the events in each bin follow a Poisson distribution.**

$$P_{dx} = \lambda(x)dx \times e^{-\lambda(x)dx} \rightarrow \text{probability for an infinitesimal bin.}$$

$$\ln L_B = \sum^n \ln \lambda(x) + \sum^n \ln dx - \int \lambda(x)dx \rightarrow \text{log-likelihood for an entire bin.}$$

λ: amplitude

*x*: width of block

*n*: number of events
   in a bin

$$\ln L_B = n \ln \lambda - \lambda x$$
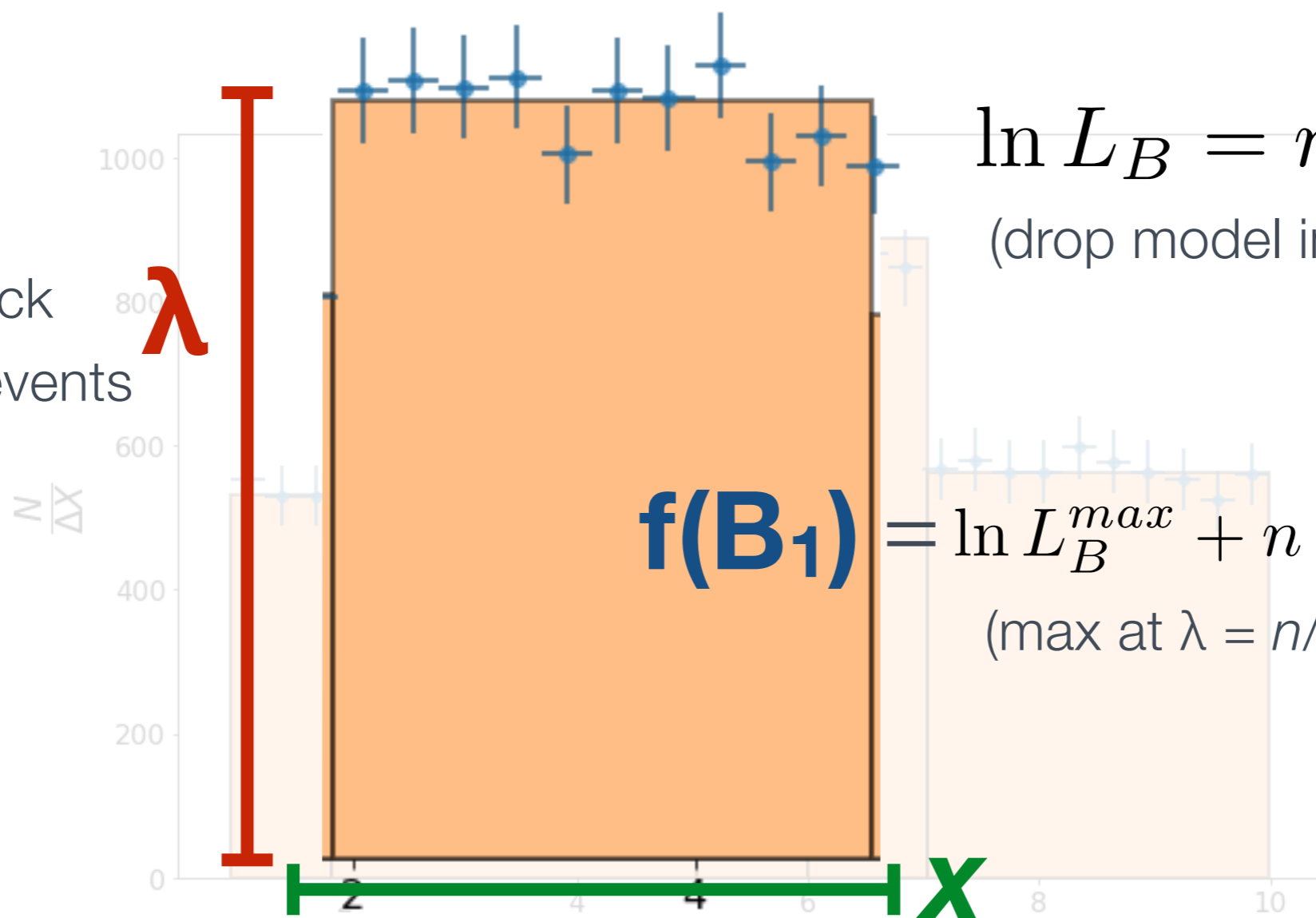
(drop model independent terms)

$$\mathbf{f(B_1)} = \ln L_B^{max} + n = \boxed{n(\ln n - \ln x)}$$

(max at *λ = n/x*)

# Penalty Term

★ **Given the previous definitions, the total fitness, $F_{total}$, will be maximal when the <u>number of bins, K, is equal to the number of data points.</u>**

- This is not desirable!

★ **A penalty term, *g(K)*, is introduced such that:**

$$F_{total} = \sum_{i=0}^{K} f(B_i) \rightarrow \sum_{i=0}^{K} f(B_i) - g(K)$$

★ **<u>Term reduces $F_{total}$ as K increases.</u>**

★ **This term is user defined, and should be tuned on signal-free data.**

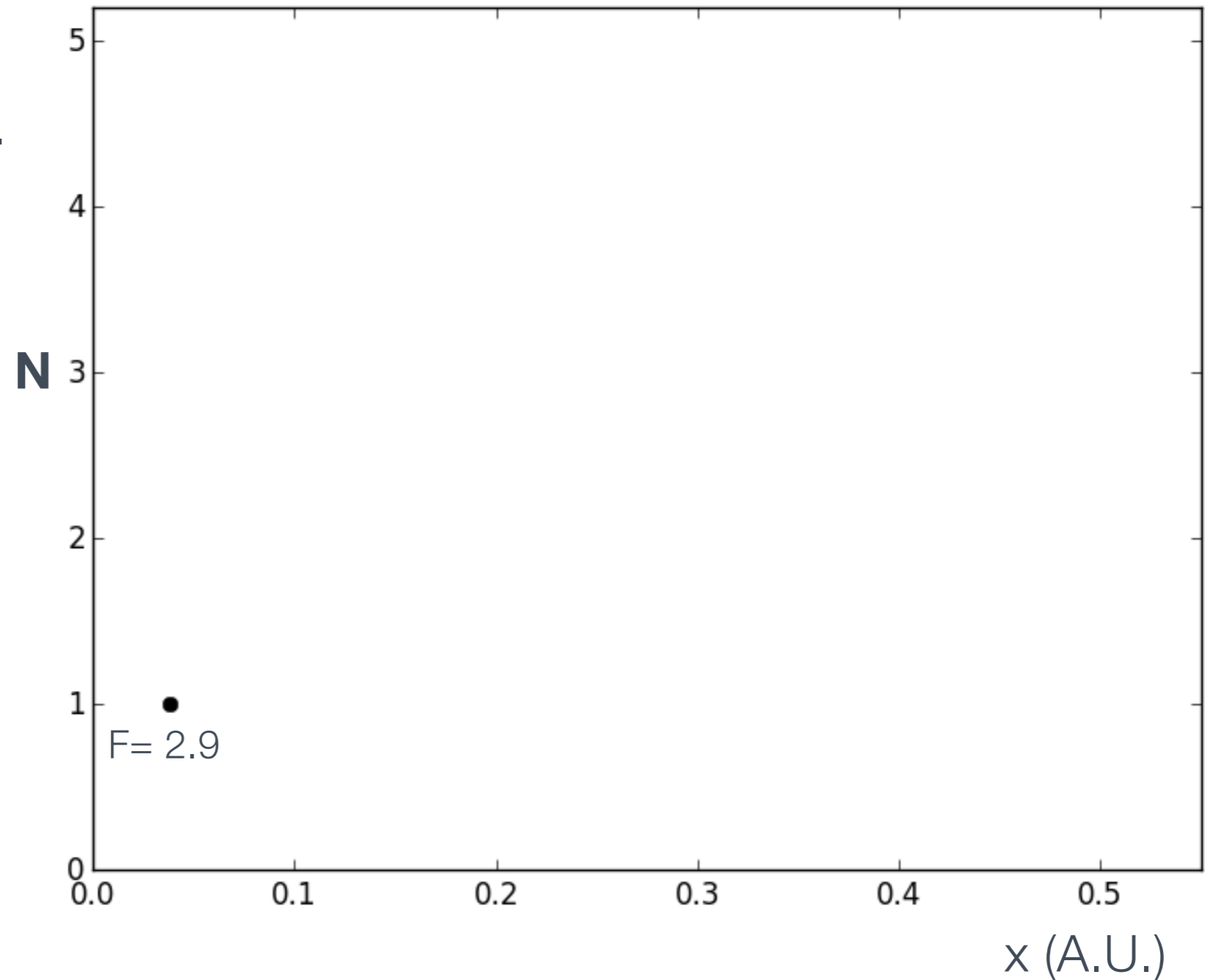# Algorithm Overview

★ **For N data points, there are $2^N$ total bin combinations.**

★ **BB algo finds optimal binning in $O(N^2)$.**

- Start: Ordered, unbinned data.

- Iterate over data:

  - Calculate fitness for all new potential bins *("New bins" = set of all bins that include newest data point).*

  - Determine current maximum total fitness *(Use cached results of previous iterations with new best bin).*

- Finish iteration, return bin edges associated with max fitness.
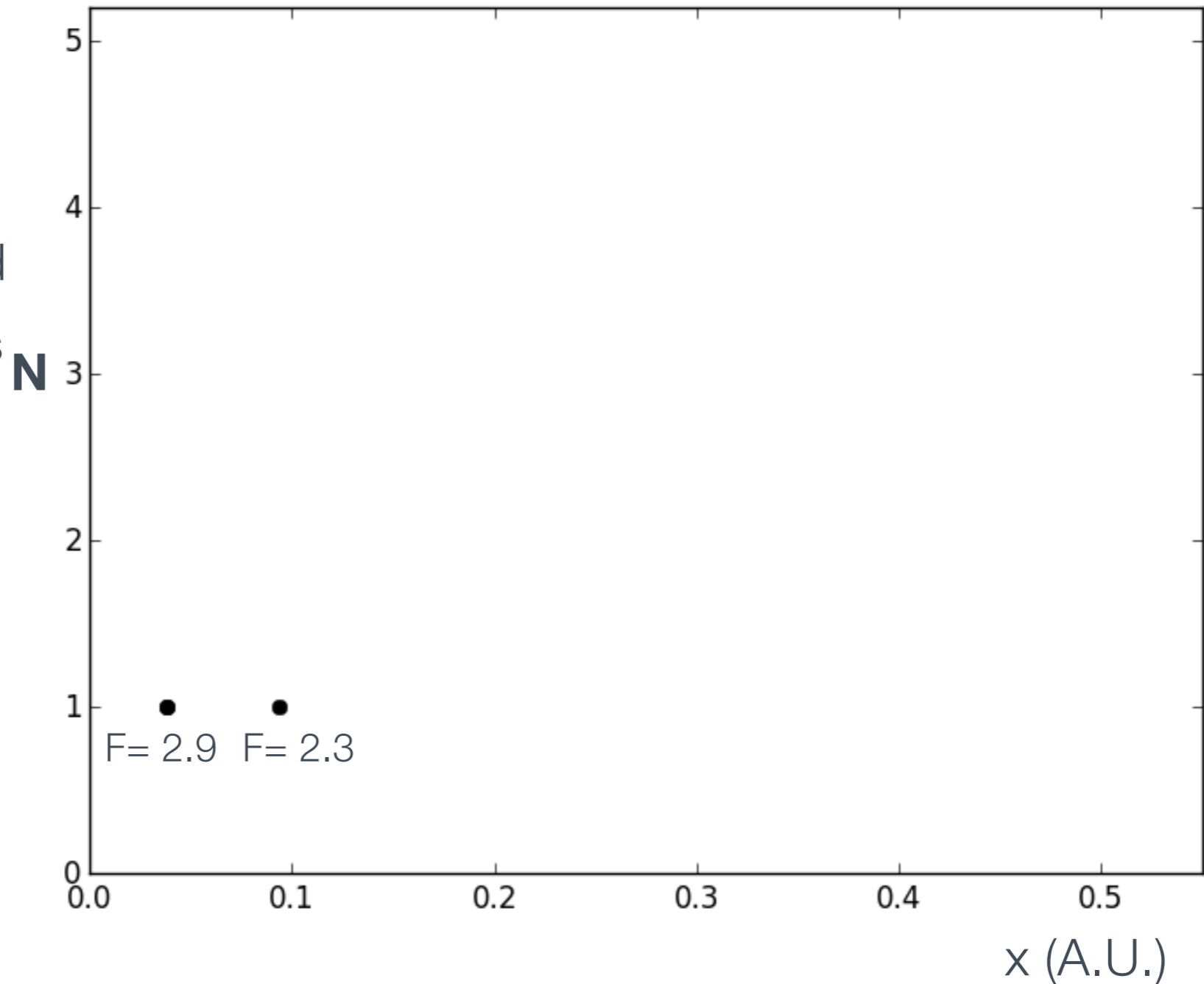
# Algorithm Example

- First data point added.
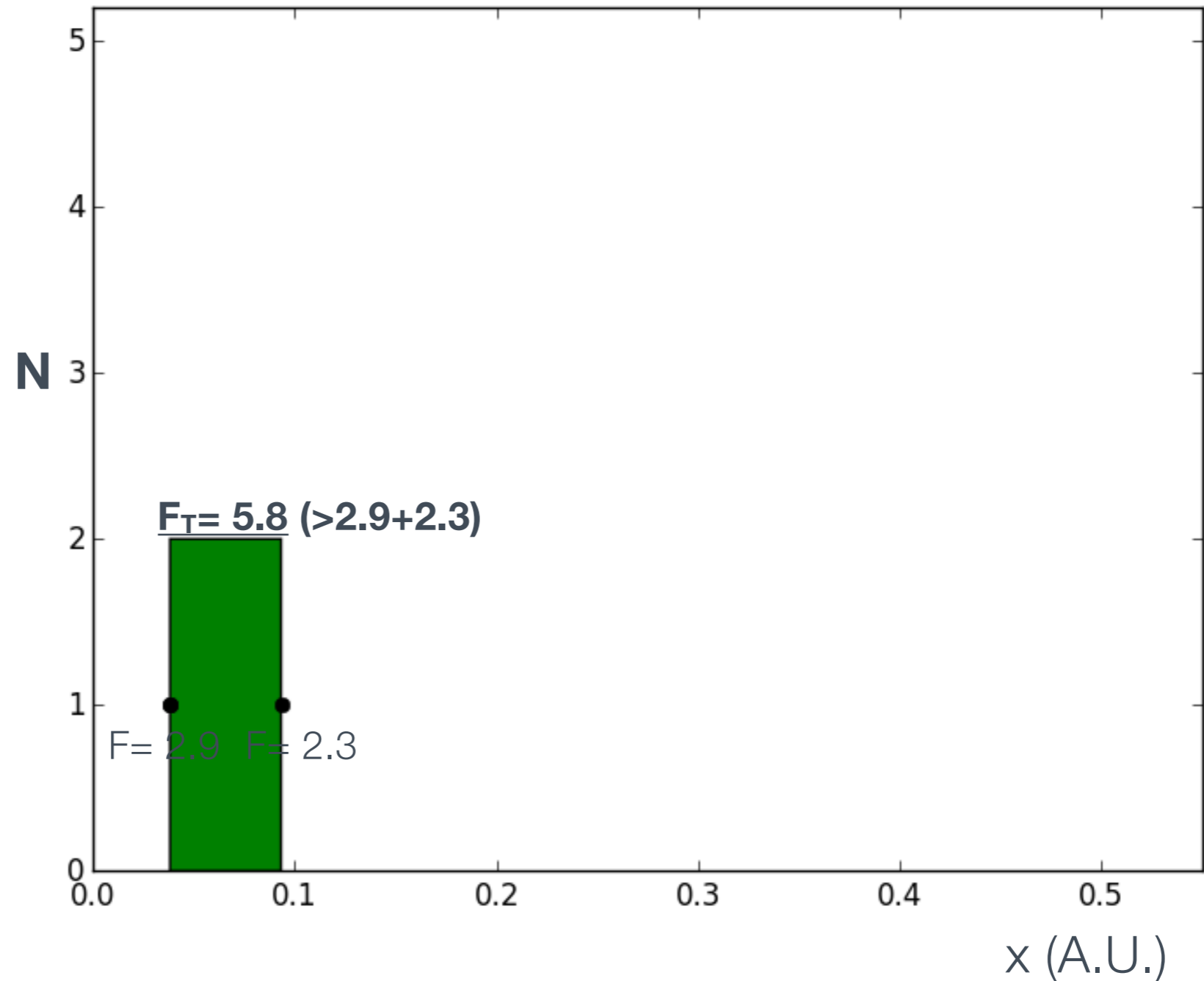- Fitness Function (F) is trivial, only one point considered.

# Algorithm Example

- Second data point added.
- Total fitness calculated ($F_T$ is sum of the fitness of all potential blocks)
- For 2 bins, $F_T = 5.2$



F= 2.9  F= 2.3
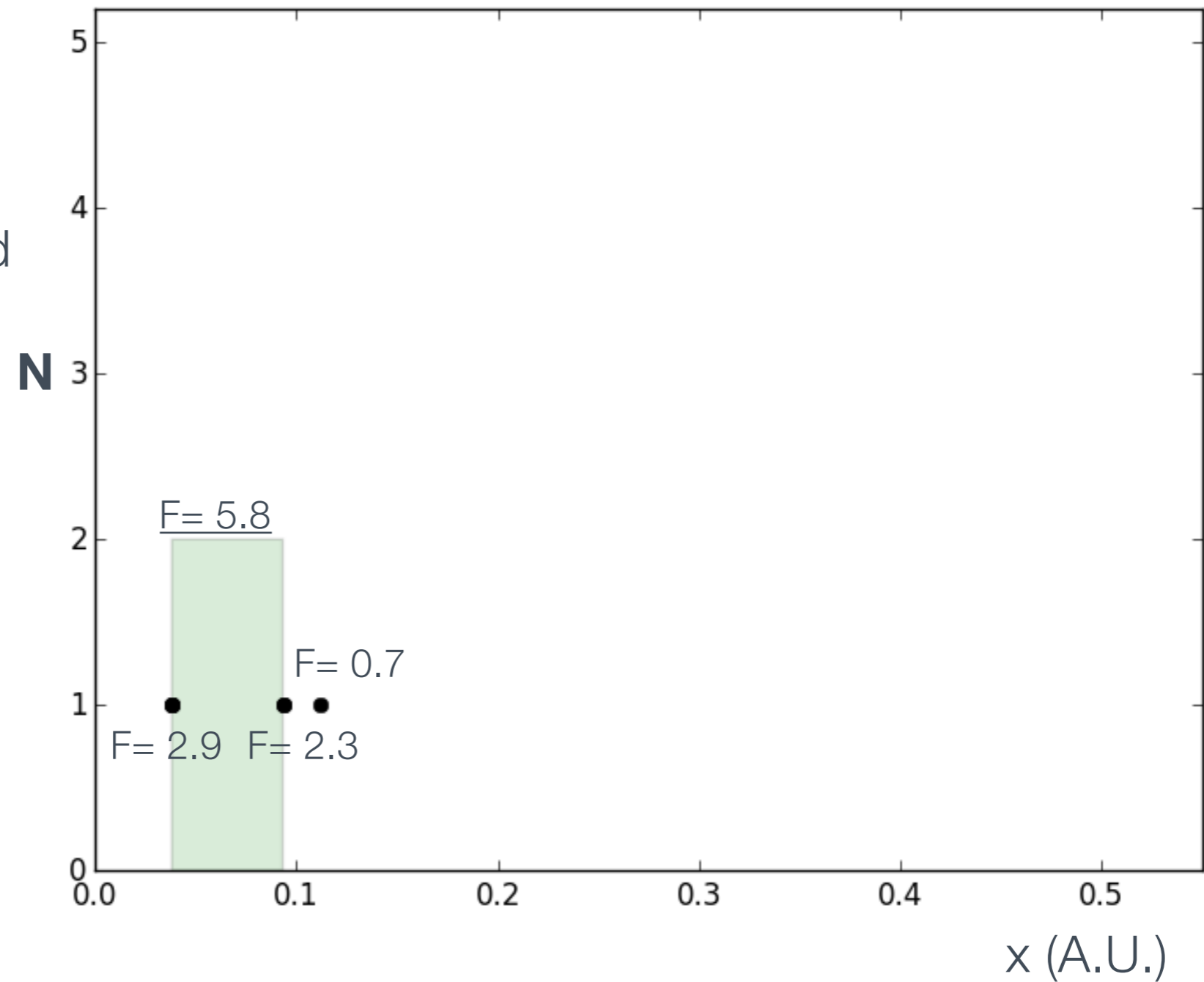
N

x (A.U.)

# Algorithm Example

- $F_T$ of single bin > $F_T$ of two bins.
- Single bin is chosen.



$F_T= 5.8 \ (>2.9+2.3)$
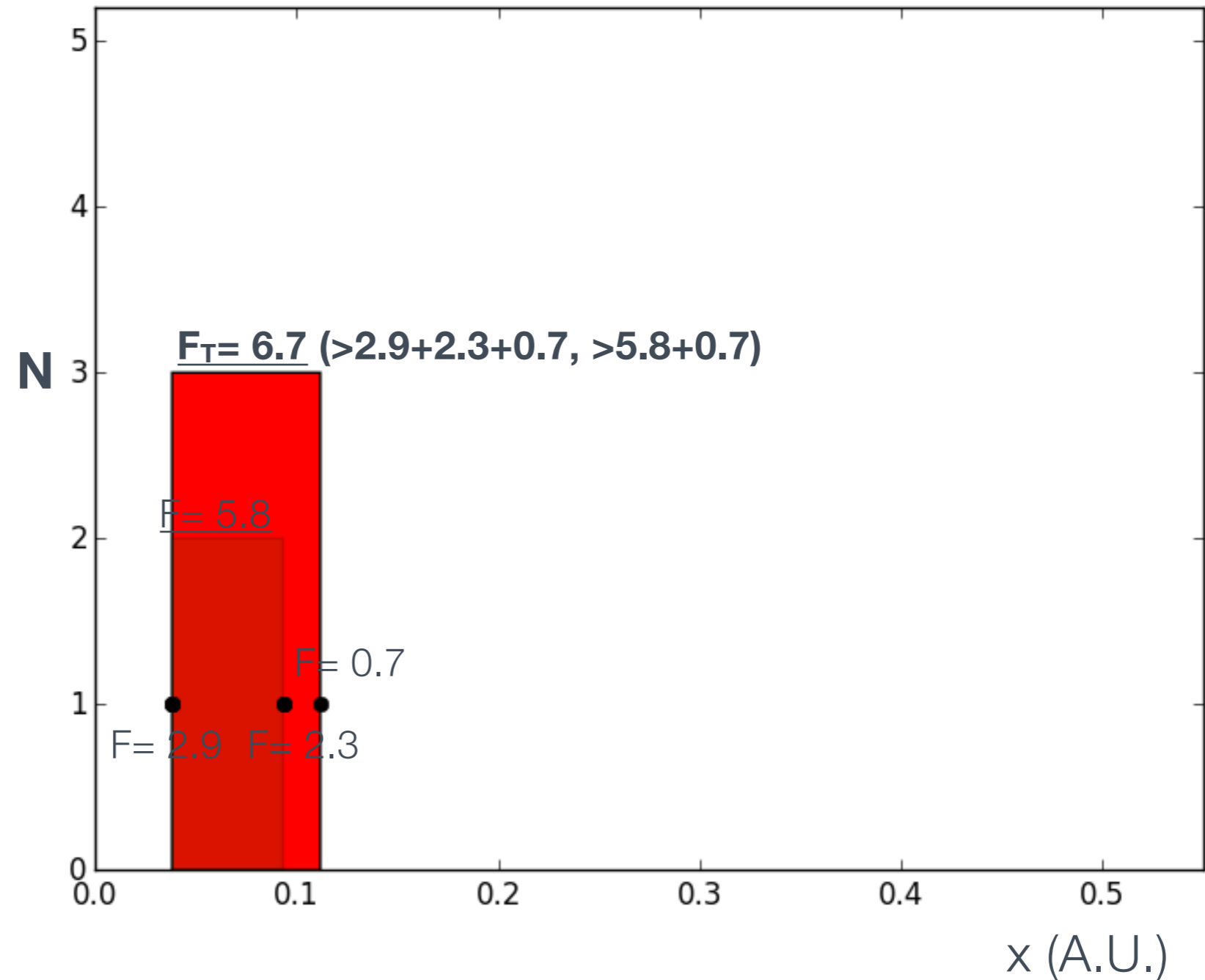
$F= 2.9 \quad F= 2.3$

N

x (A.U.)

# Algorithm Example

- Third data point added

# Algorithm Example

- $F_T$ of single bin > $F_T$ of all other combos (using stored F values from previous iterations)



$F_T$= 6.7 (>2.9+2.3+0.7, >5.8+0.7)

F= 5.8

F= 0.7

F= 2.9   F= 2.3
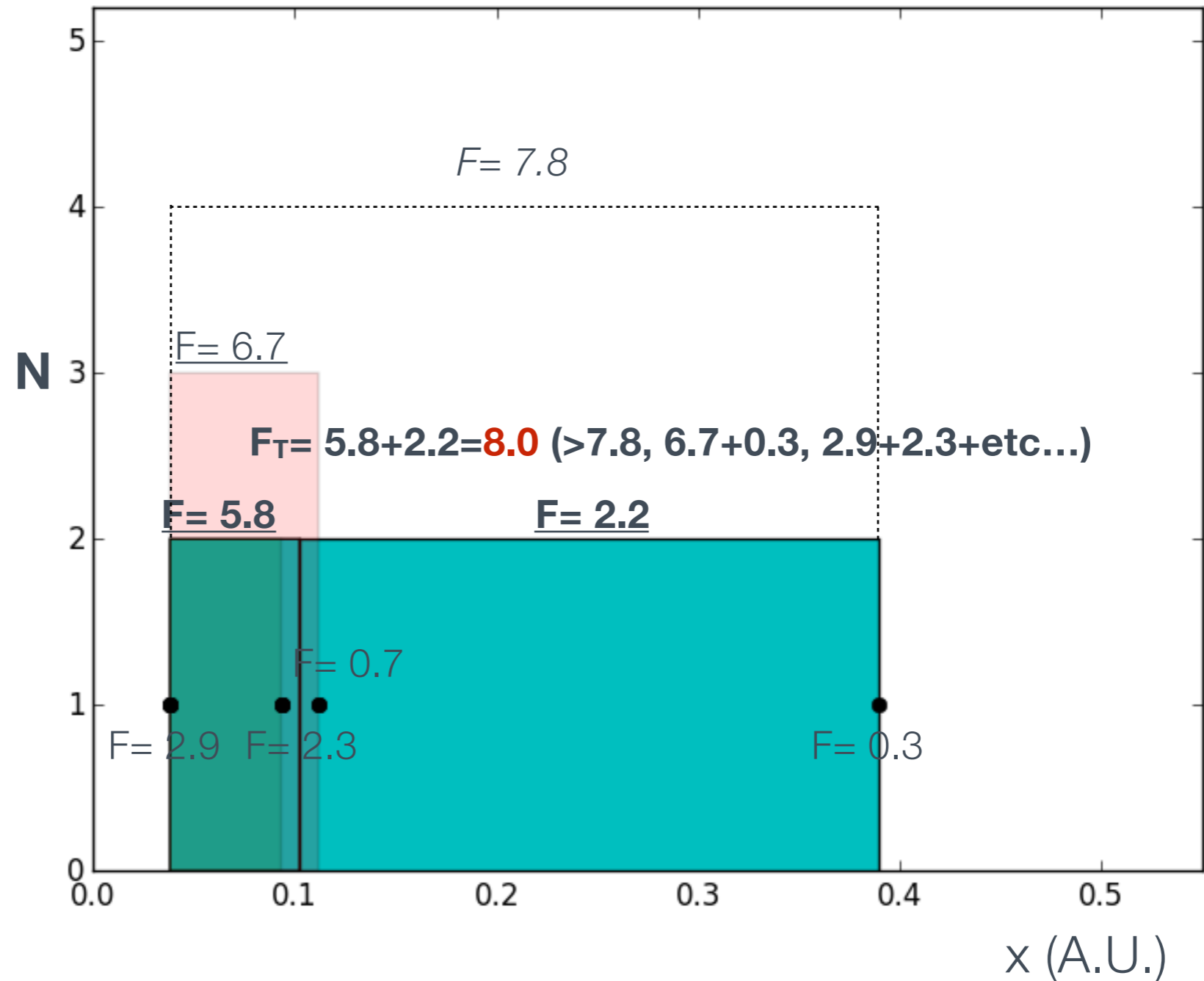
N

x (A.U.)

# Algorithm Example

- Fourth data point added
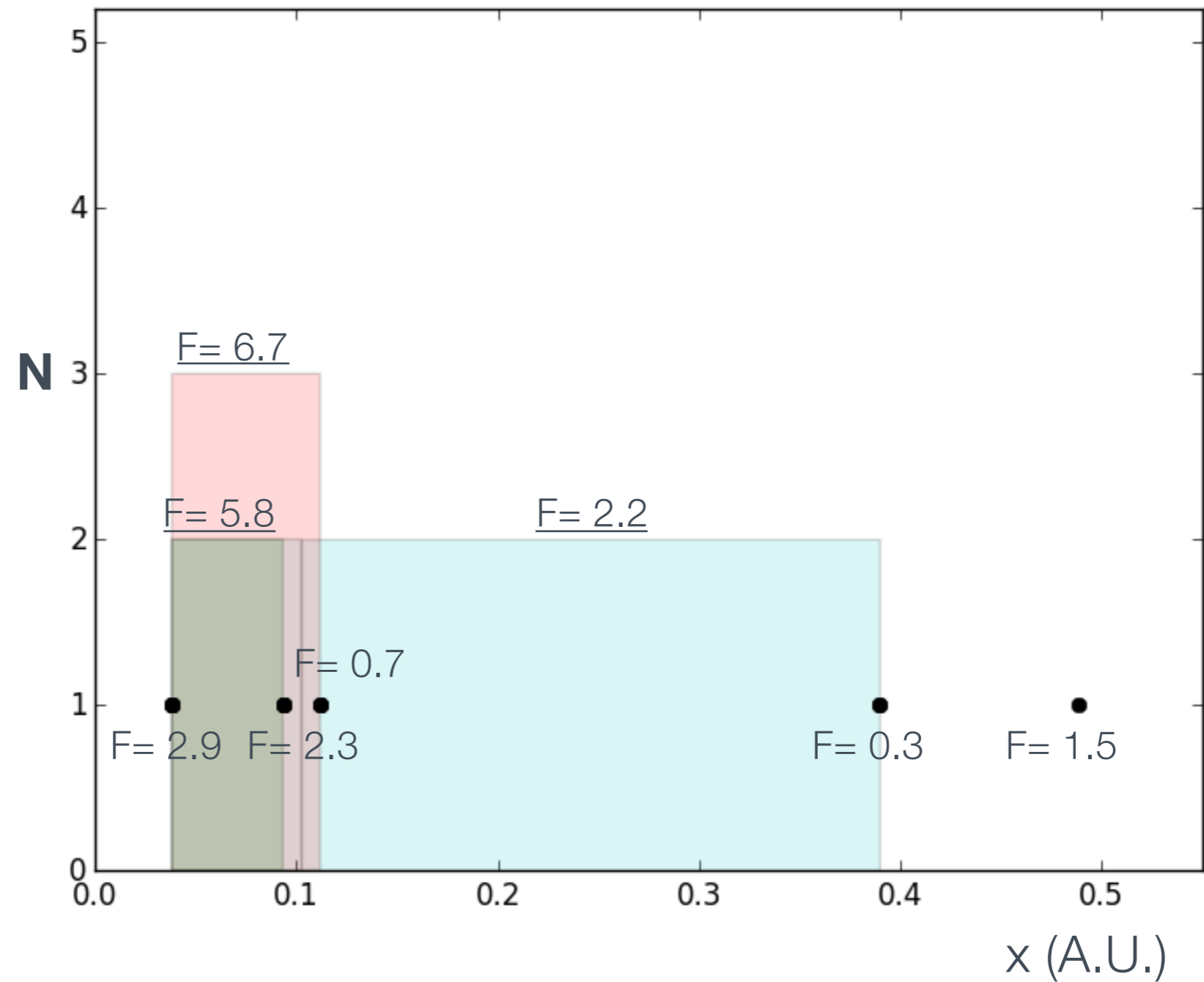
# Algorithm Example

- Maximum $F_T$ is for 2 bins
  - *F value of first bin was stored from previous iteration*
- New change-point is determined between pts 2 and 3
- Change-point is saved along with $F_T$ value



$F= 7.8$

$F= 6.7$

**N** 3

$F_T= 5.8+2.2=8.0$ (>7.8, 6.7+0.3, 2.9+2.3+etc...)

**F= 5.8**          **F= 2.2**

$F= 0.7$

F= 2.9   F= 2.3                    F= 0.3
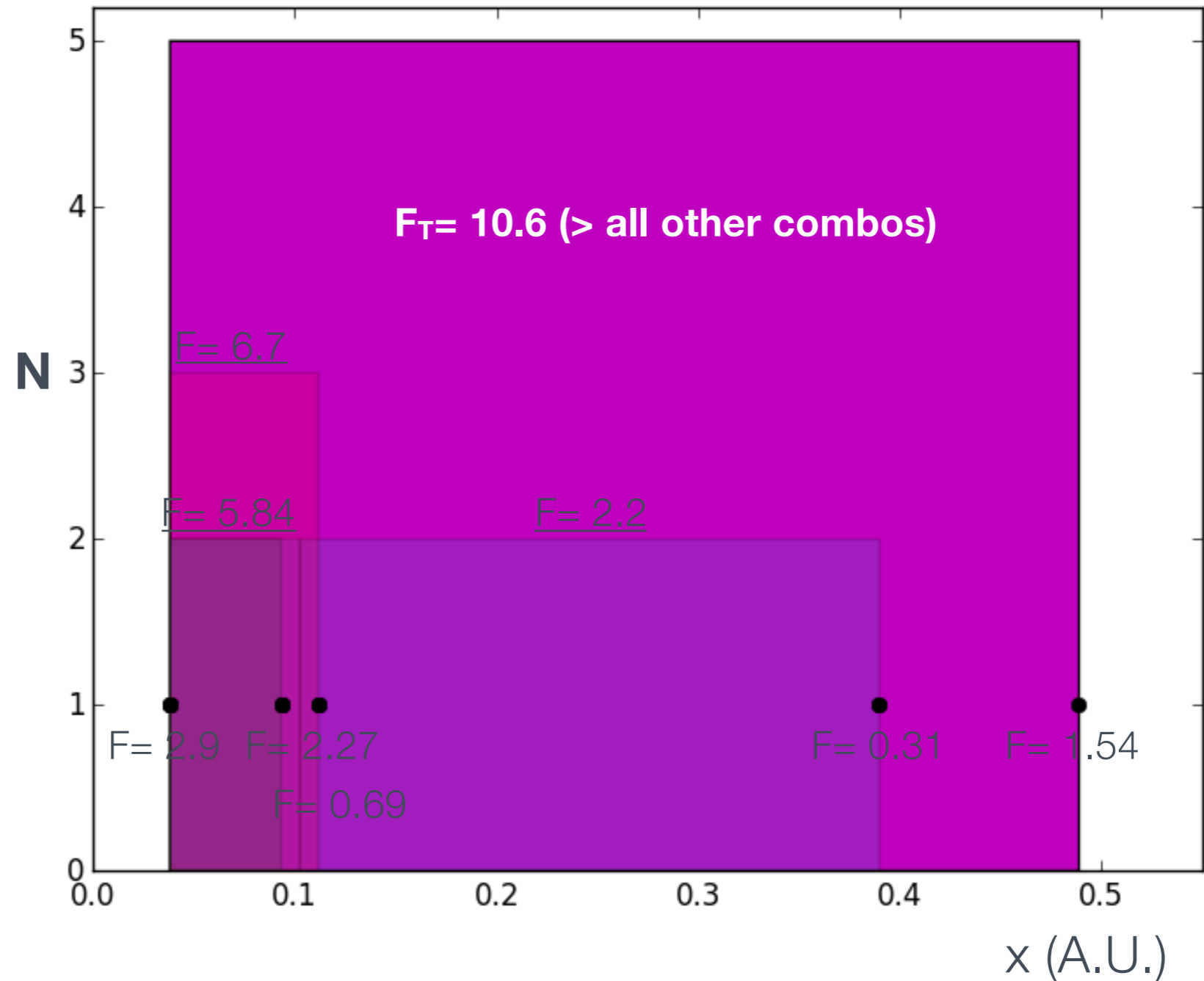
x (A.U.)

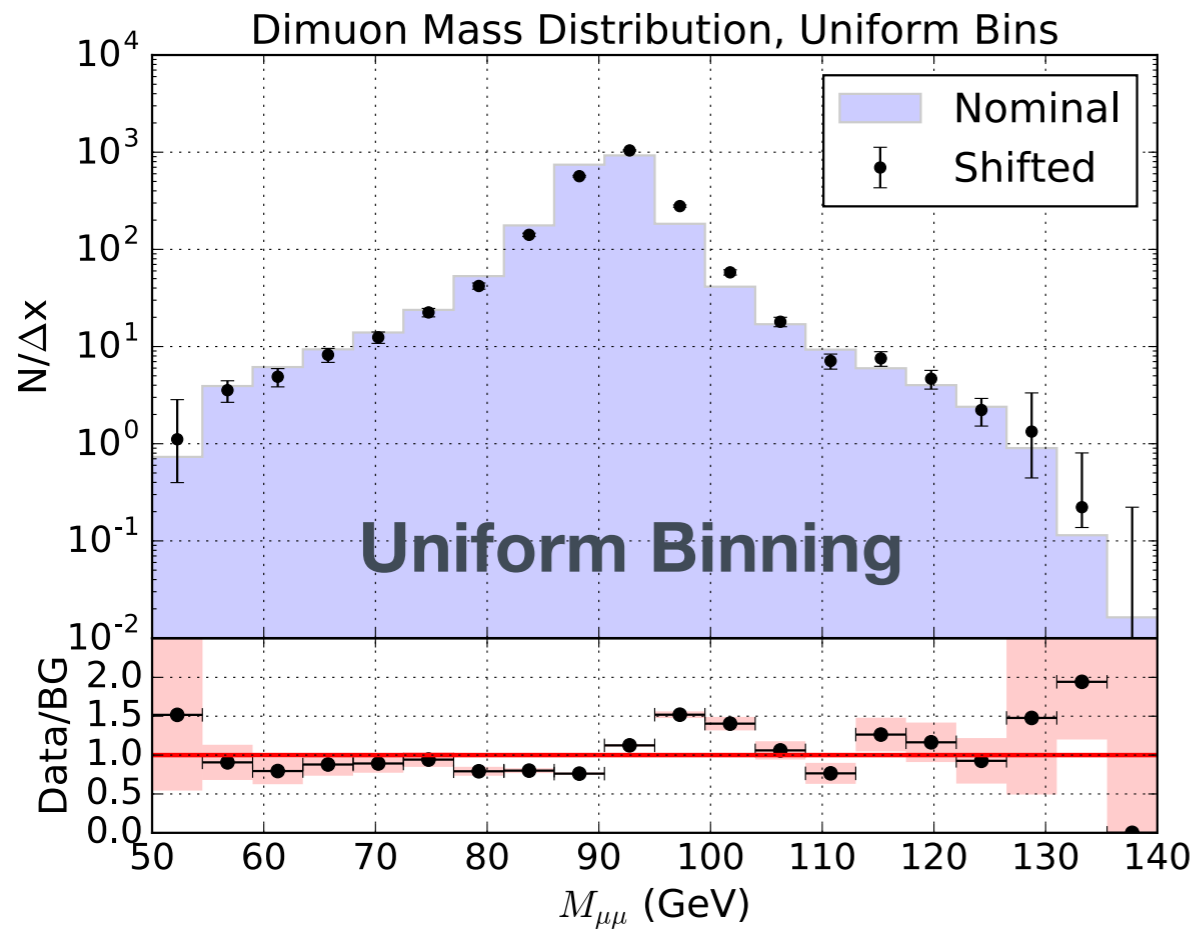# Algorithm Example
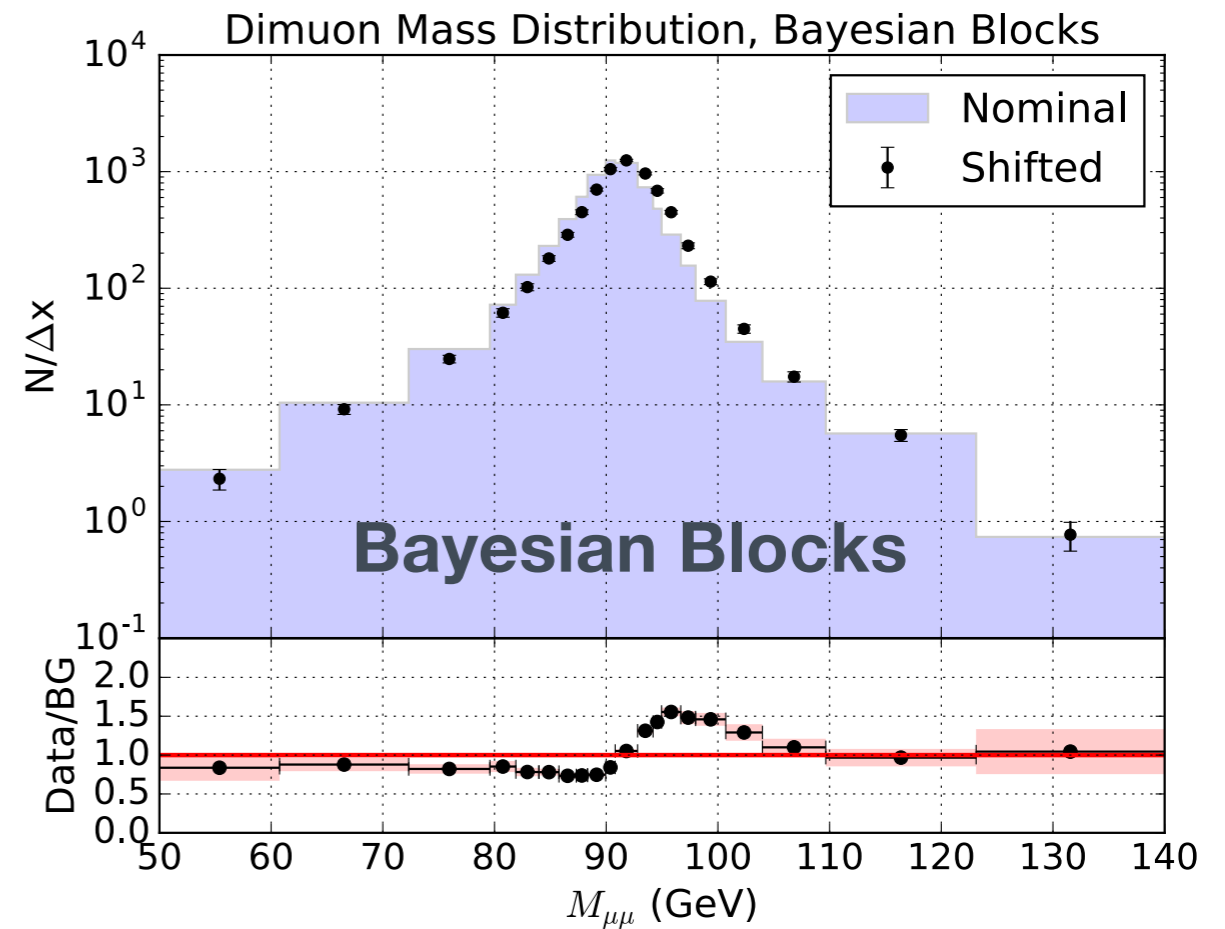
- Final data point added

# Algorithm Example

- Maximum $F_T$ is determined to be single bin
- Previous change-point is ignored because of sub-optimal value
- Final result yields bin edges at [1,5]



$F_T$= 10.6 (> all other combos)

F= 6.7

F= 5.84    F= 2.2

F= 2.9    F= 2.27         F= 0.31    F= 1.54

F= 0.69

N

x (A.U.)

# Visual Impact



(a) Fixed-width binning.

(b) BB binning.

★ **Simulated Z→μμ example.**

- One distribution is slightly shifted w.r.t. other → typical HEP scenario before muon scale corrections are applied.
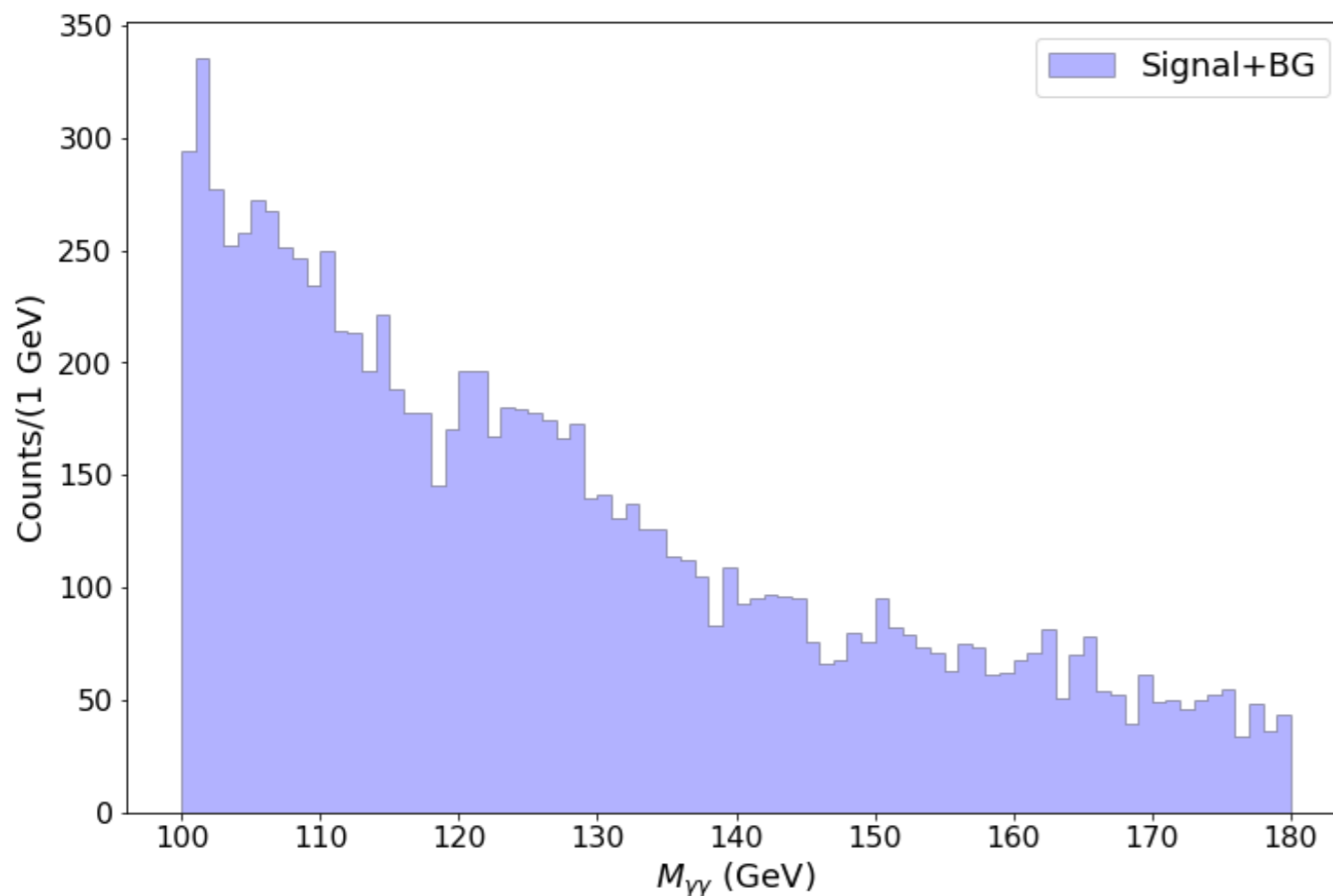
★ **Bayesian Blocks example shows more detail in peak, smooths out statistical fluctuation in tails.**

# Bump Hunting

★ **The bin edges determined by Bayesian Blocks are statistically significant.**

  • Can they assist with analyses, outside of purely visual?
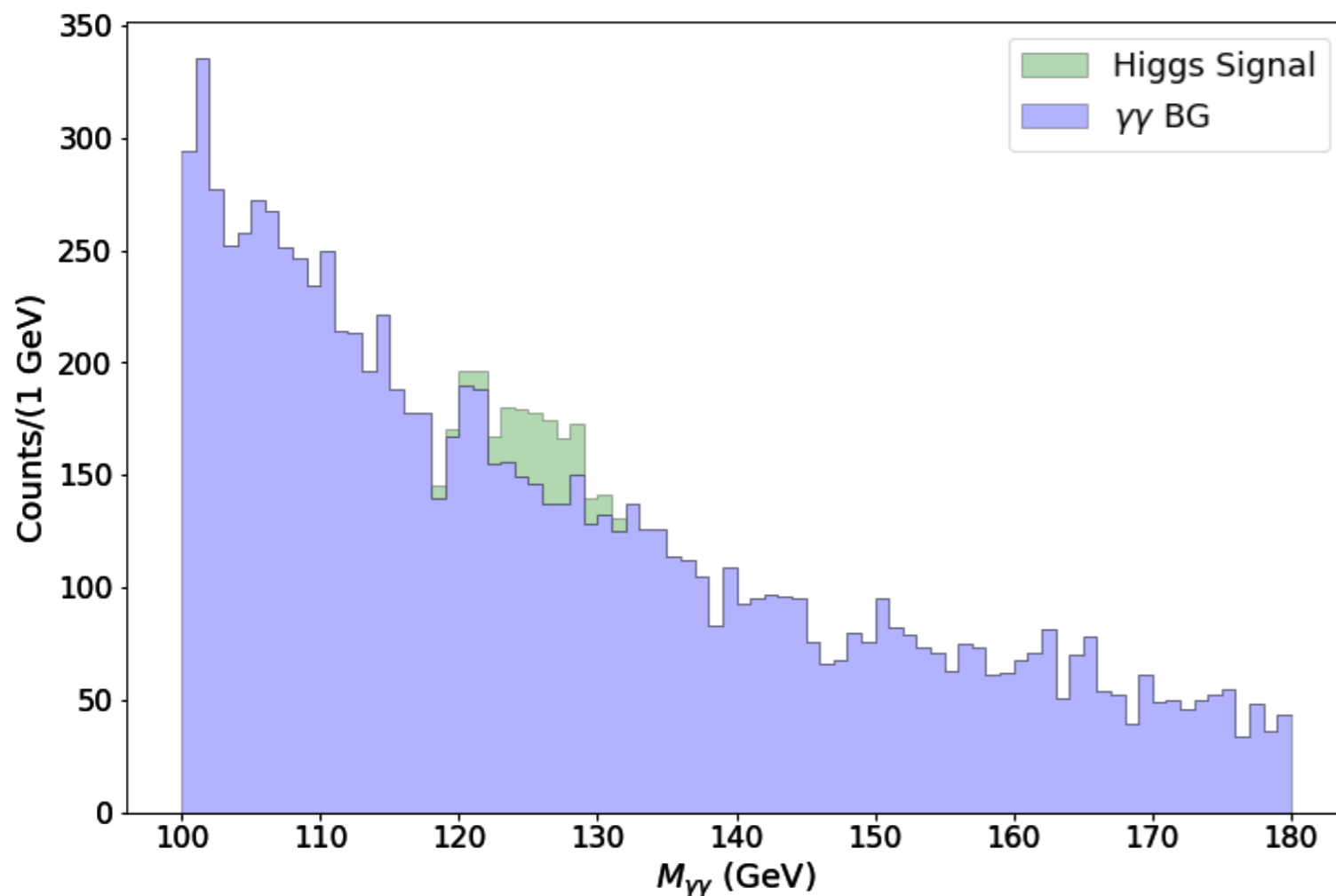
★ **Consider the H→γγ discovery (simulated):**



• Falling diphoton BG, ~10k events.
• ~230 Higgs signal events at $M_{\gamma\gamma}$=125 GeV (~5 σ excess)

# Bump Hunting

★ **The bin edges determined by Bayesian Blocks are statistically significant.**

   • Can they assist with analyses, outside of purely visual?

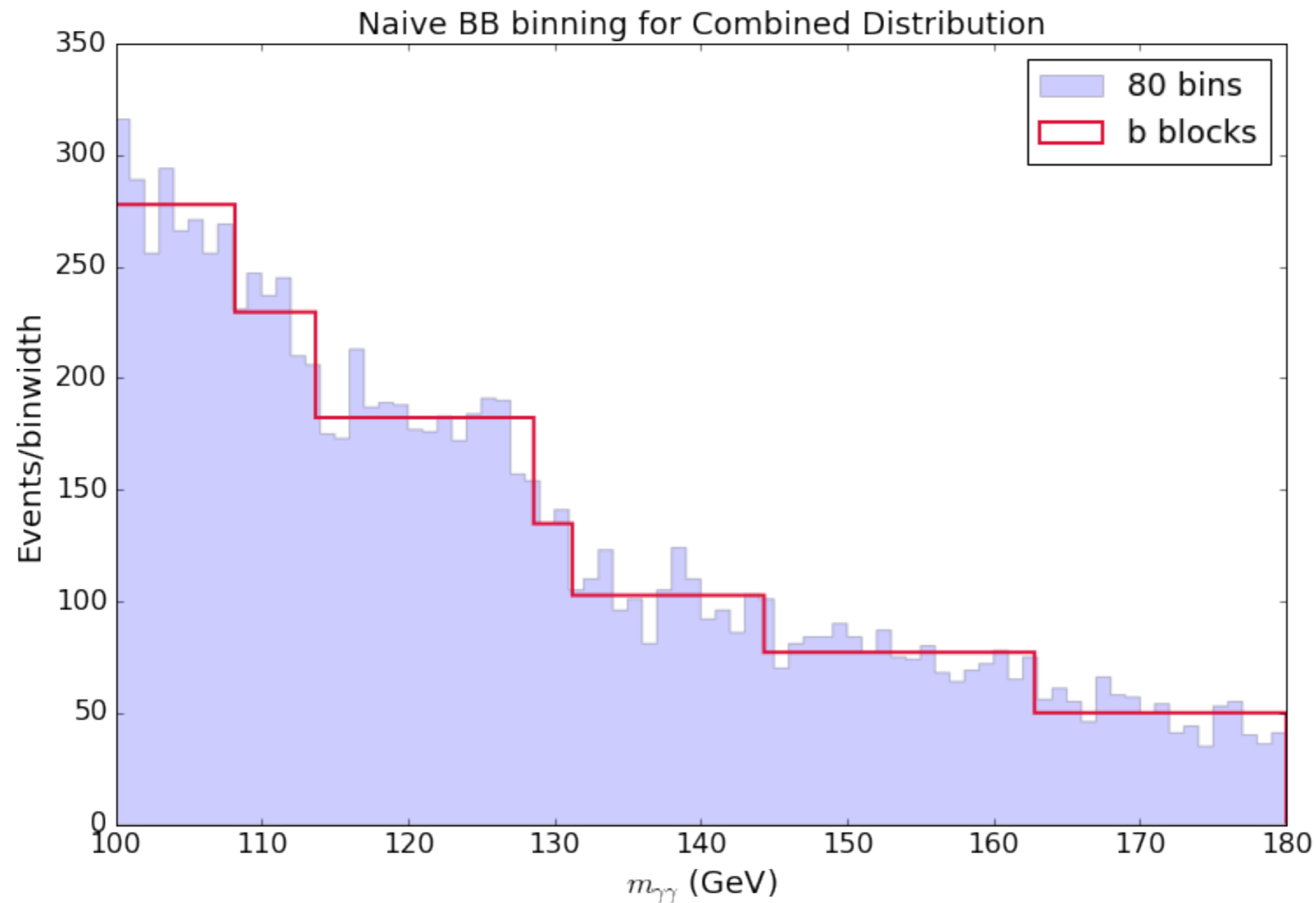★ **Consider the H→γγ discovery (simulated):**



   • Falling diphoton BG, ~10k events.
   • ~230 Higgs signal events at $M_{\gamma\gamma}$=125 GeV (~5 σ excess)

**Significant excess, difficult to discern by eye.**

# Bump Hunting

**First try, naive binning of signal+background:**



Naive BB binning for Combined Distribution

# Bump Hunting

**First try, naive binning of signal+background:**



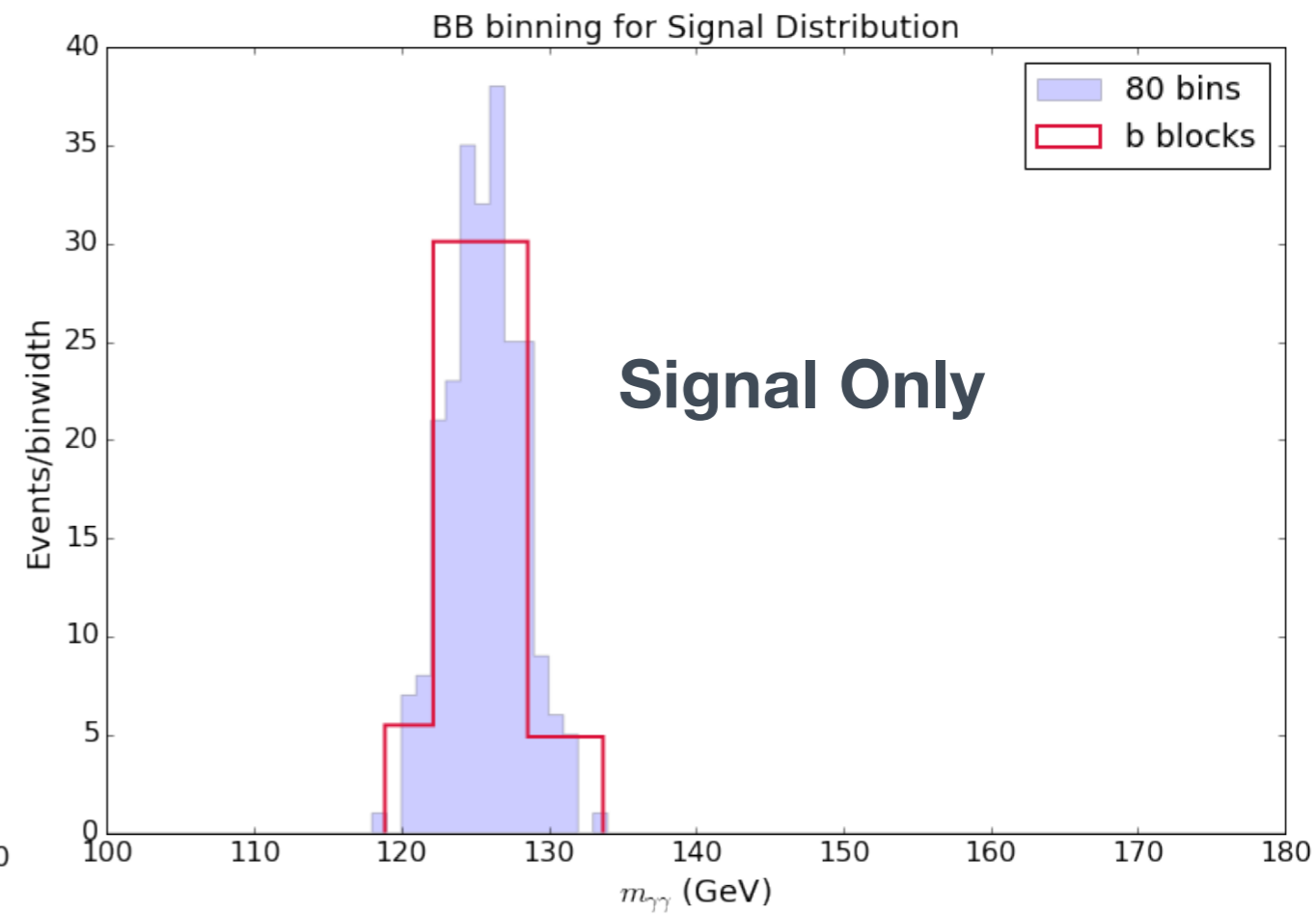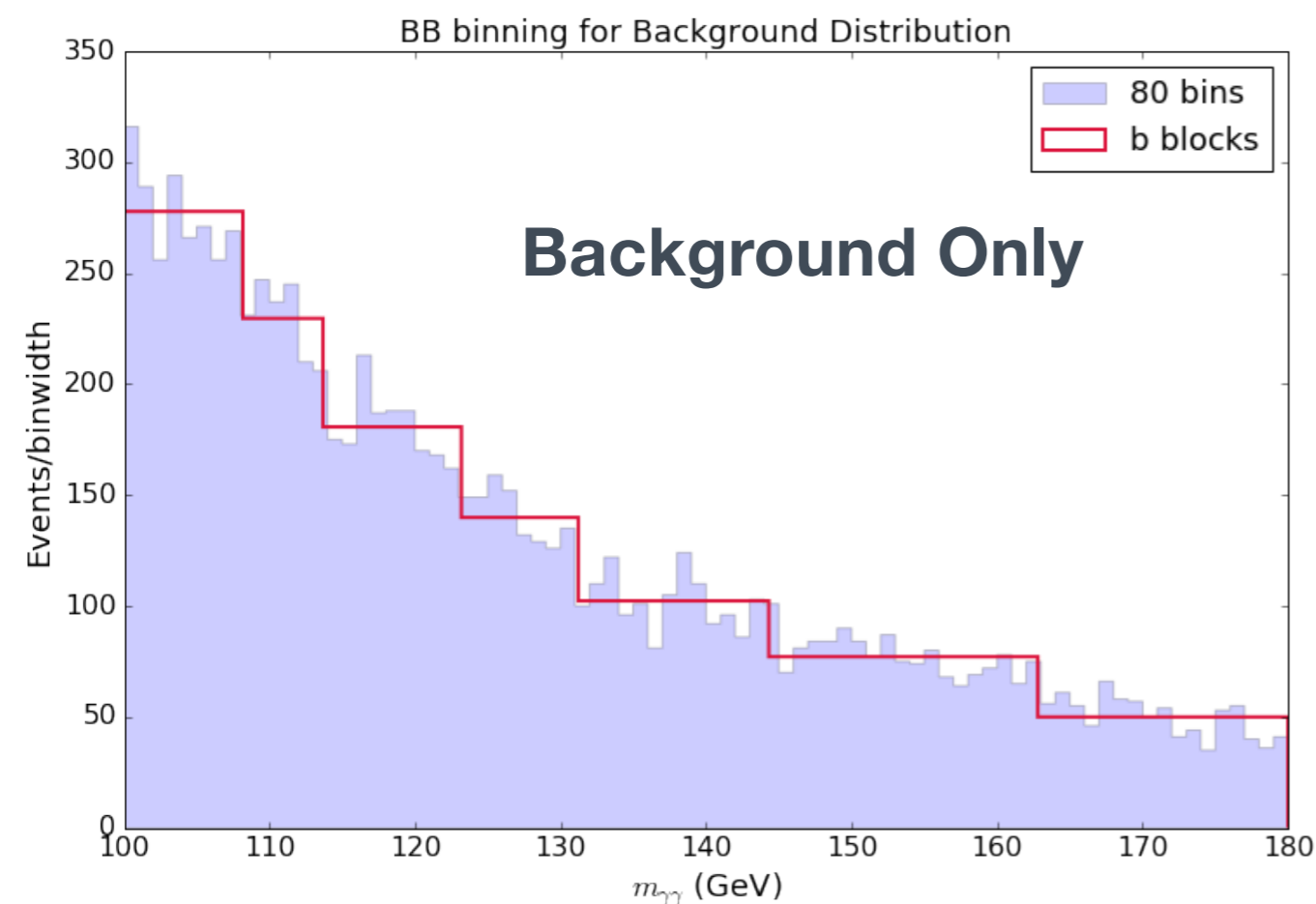Naive BB binning for Combined Distribution

**Results not great.**

**Falling background + rising signal = one large bin.**

# Bump Hunting

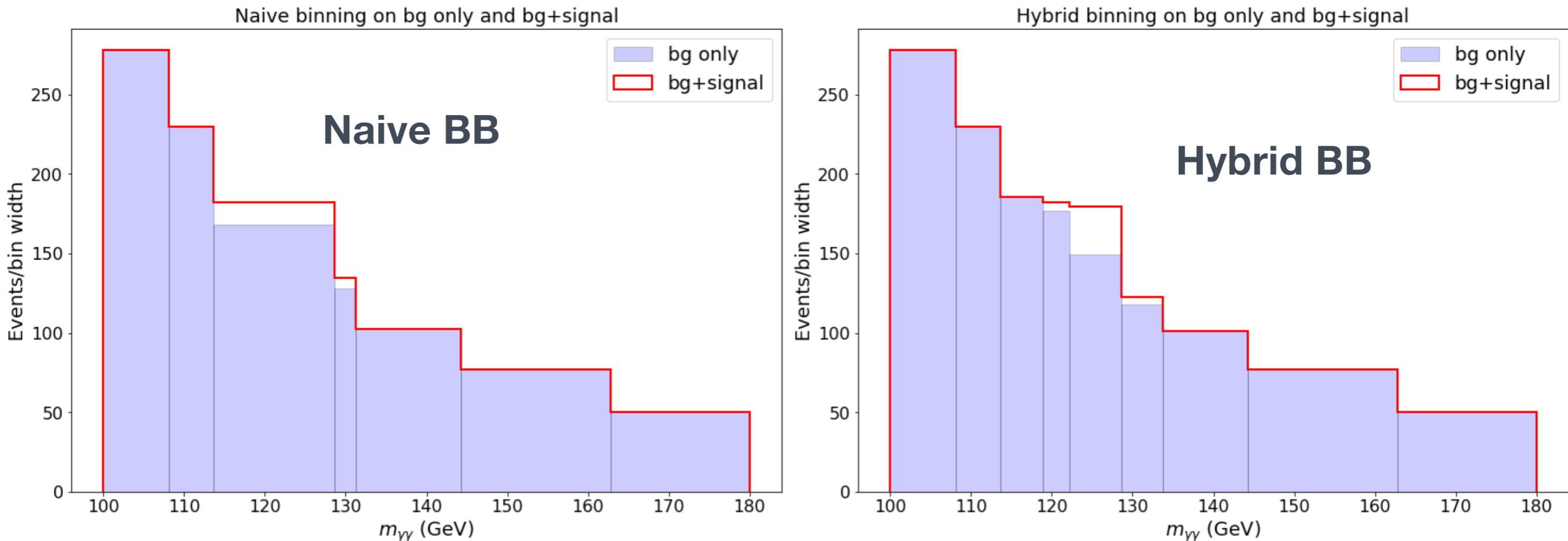★ **Generate a "hybrid" binning, leveraging knowledge of signal shape:**

- Use Bayesian Blocks on simulated signal and background templates.

- Combine the bin edges (background bin edges in signal region replaced by signal bin edges)

# Bump Hunting

★ **Signal excess much more apparent with hybrid binning:**



No parametric models used to generate binning, completely MC dependent.

What is the sensitivity of this excess?

# Bump Hunting

★ **Calculate Gaussian Z-score (# of σ excess) for 1000 simulations, and compare to unbinned likelihood from known underlying pdfs.**

• Z-score from unbinned likelihood are the <u>upper-bound</u>.



**Mean Z-scores:**

**Bayesian Blocks Template: 5.35 σ**

**Unbinned likelihood: 5.57 σ**

*Hybrid binning is only slightly less sensitive than unbinned pdf, and is completely non-parametric!*

# Software



★ **Python histogramming package developed for HEP:**

- Wraps matplotlib, adds automatic error bars, scaling, <u>Bayesian Blocks binning</u>, and more!

★ **Install with pip:**

- `$ pip install histogram_plus`

Documentation (in progress):

https://brovercleveland.github.io/histogram_plus/

```python
In [1]:  import numpy as np
         import matplotlib.pyplot as plt
         from histogram_plus.hist_funcs import hist

         %matplotlib inline
         plt.rcParams['figure.figsize'] = (12,8)
```
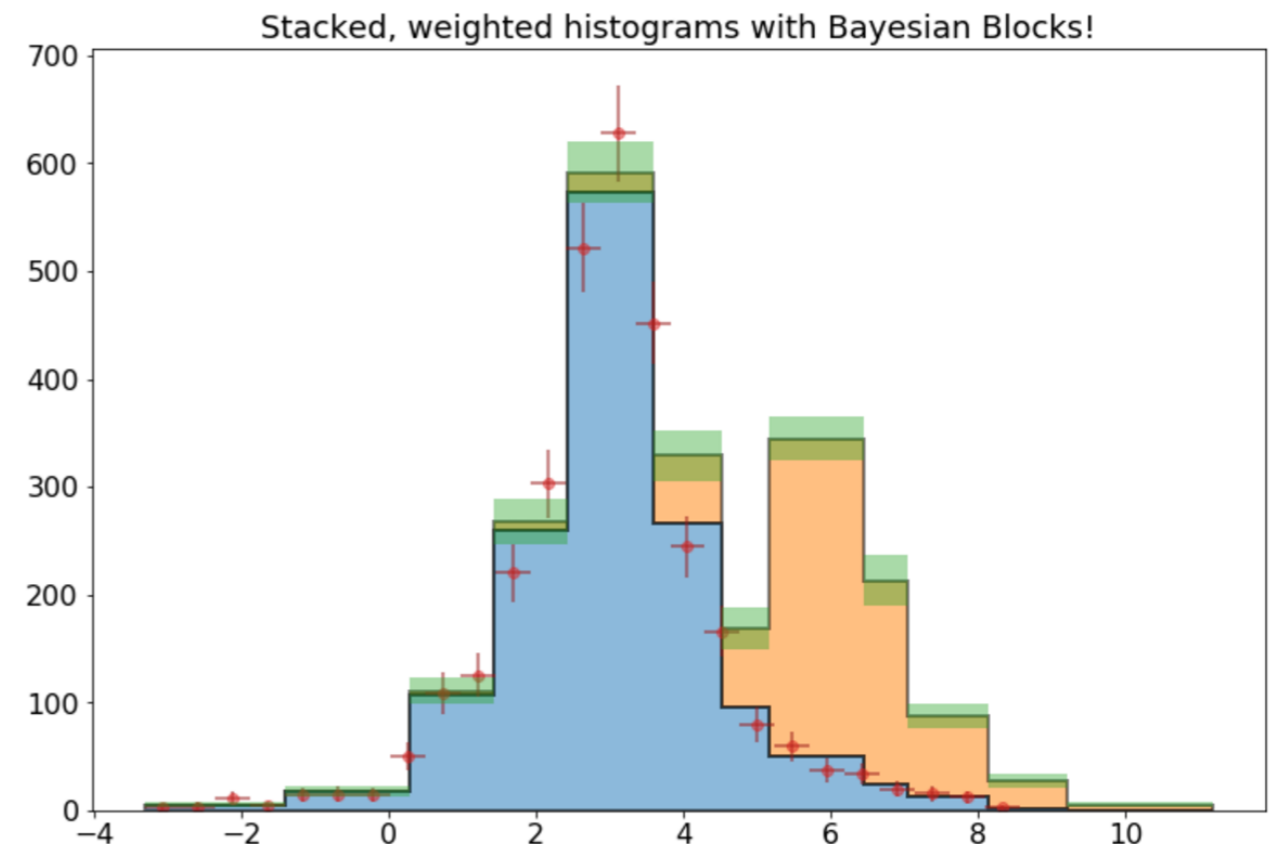
```python
In [2]:  np.random.seed(8008)
         data = np.random.laplace(3, 1, size=1000)
         data2 = np.random.laplace(6, 1, size=500)
         weights = np.random.uniform(1,2, size=1000)
         weights2 = np.random.uniform(1,2, size=500)
```

```python
In [3]:  stacked_hists = hist([data,data2], weights=[weights, weights2], bins='blocks',
                              stacked=True, errorbars=True, scale='binwidth', p0=0.01)

         marker_hists = hist(data, bins=25, weights=weights, histtype='marker',
                             errorbars=True, scale='binwidth')

         plt.title('Stacked, weighted histograms with Bayesian Blocks!')
```

```
Out[3]: <matplotlib.text.Text at 0x113b93b50>
```

# Summary

★ **The Bayesian Blocks algorithm is a data-driven, model-independent method for binning.**

- Bins are variable-width, edges represent statistically significant changes in data.

- Improves visualization of distributions, even with dense peaks and sparse tails.

★ **Bayesian Blocks can also assist in template-based analyses.**

- Provides a non-parametric way of modeling distributions in histograms, with minimal loss in sensitivity when compared to unbinned methods.

★ **New paper on HEP application for Bayesian Blocks:**

- https://arxiv.org/abs/1708.00810