

Lecture 2.5

Data challenges – missing data, low or high values, and historical/paleo information

Flood Frequency Analysis

Beth Faber, PhD, PE

Hydrologic Engineering Center

1

Goals

- To be aware of the data challenges we're likely to encounter in frequency analysis
- To understand how "challenging data" can **impact** the resulting frequency curve
- To discuss former (17B) and current (17C) methods for handling these issues **to produce robust estimates**

↓
*do well, even when
assumptions are wrong*

2

Bulletin 17B/C methods

- Estimate probability distribution when:
 - Recorded annual data has gaps, missing values or annual peak flow is zero
 - there are low outliers
 - there are high outliers and/or historical information
- Scenarios:
 - *Broken record*
 - *Censored flow records*
 - *Historical Information*

3

These are the difficult data situations we might have to deal with when doing frequency analysis. Bulletin 17B had adjustments for each, and Bulletin 17C has a different way of handling each, but the impacts can be the same, or at least similar.

Outline

- Missing values
- Censored values, zero flows, outliers
 - low flows (PILFs)
 - high flows
- Historical/Paleo information

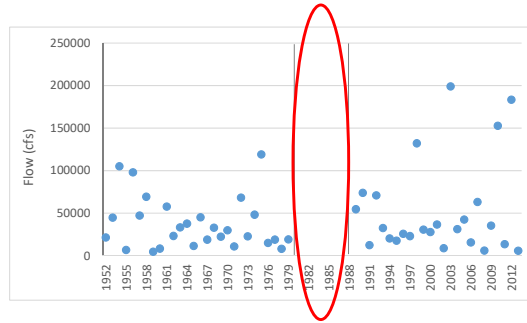
Broken Record

- Two or more periods of systematic record separated by unobserved periods.



- Periods combined and analyzed as a ***single record***.
 - if we have absolutely no information about the unobserved years, they are treated the same as years before or after the systematic record.
- If we have some information, such as knowing flows are below a threshold, EMA lets us use that non-exceedance information...

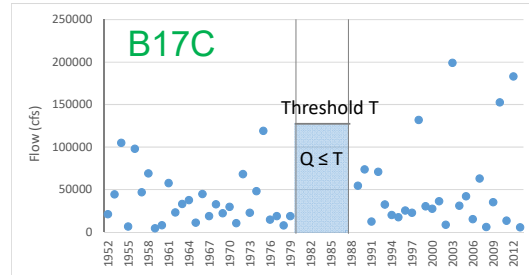
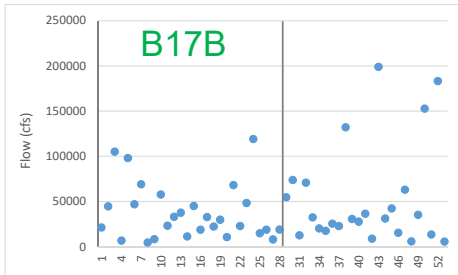
The gage record is missing the years 1980 – 1988



Combine data into a single record for analysis

OR

Define a non-exceedance threshold for the missing years



Or, define perception range as (∞, ∞)

6

Sometimes the gage is out of service causing missing years. If there were known to be relevant flow values (floods) during that time, they should be investigated. But if nothing is known, B17B just closed the record around those years to make a single record of before and after.

The same result can be achieved with B17C by assuming a perception threshold range of infinity to infinity. But if you can say that large events would have been noted, and define a perception threshold, can use a flow range for the missing years.

Outline

- Missing values
- Censored values, zero flows, outliers
 - low flows (PILFs)
 - high flows
- Historical/Paleo information

Censored Records: Flows above or below recording level

- Some gages have a lower bound (*and/or upper bound*) of channel stage that can be measured
 - Flows can result in stages that are below the smallest (*or above the largest*) recording level
 - Flows below (*or above*) the threshold are referred to as **censored** values
 - **unable to measure**
- Unobserved historical flows are also “censored”
 - There can be **evidence in the watershed** that some threshold was exceeded, or knowledge it wasn't
- *Though flows are not measured, knowing they were above or below a threshold is valuable info in estimating the flood distribution*

↙
adjective

8

This slide describes CENSORED as an adjective – they are flows whose values are not known precisely, though some information might be available.

B17B Outliers: values notably different from the rest of the data

Low Outliers are censored → *verb*

- analysis excludes the values, but not the fact they occurred
- B17B used Conditional Probability Adjustment
- B17C uses intervals for excluded values

B17B used Grubbs-Beck test for high and low outlier thresholds

High Outliers are **NOT** censored

- high values are left in the data set
- seek historical information to either add a past large event, or determine that the largest gaged event has a longer return period (*is the largest in a longer period of time.*)
- B17B used Weighted Moments Algorithm
- B17C uses intervals for unobserved years

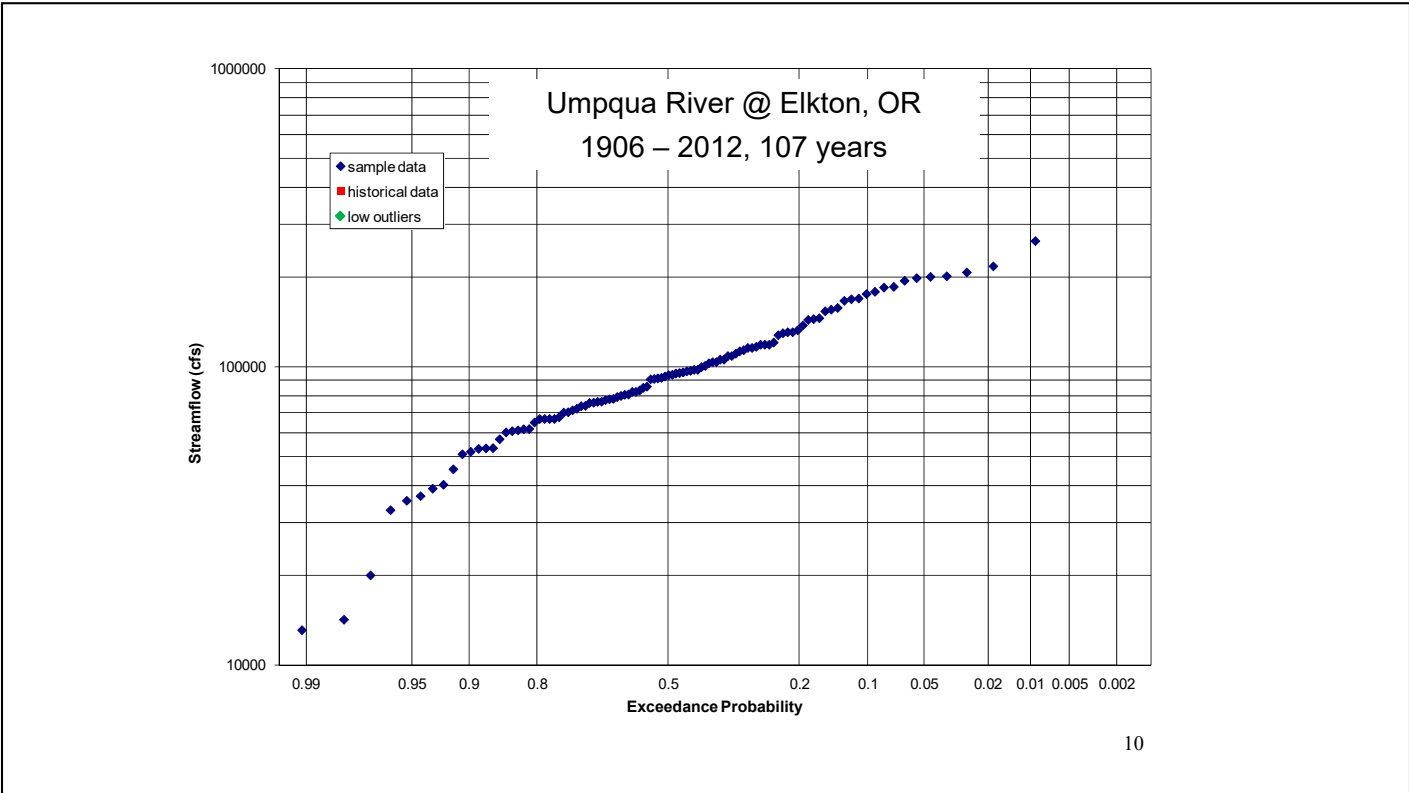
B17C uses Multiple GB

9

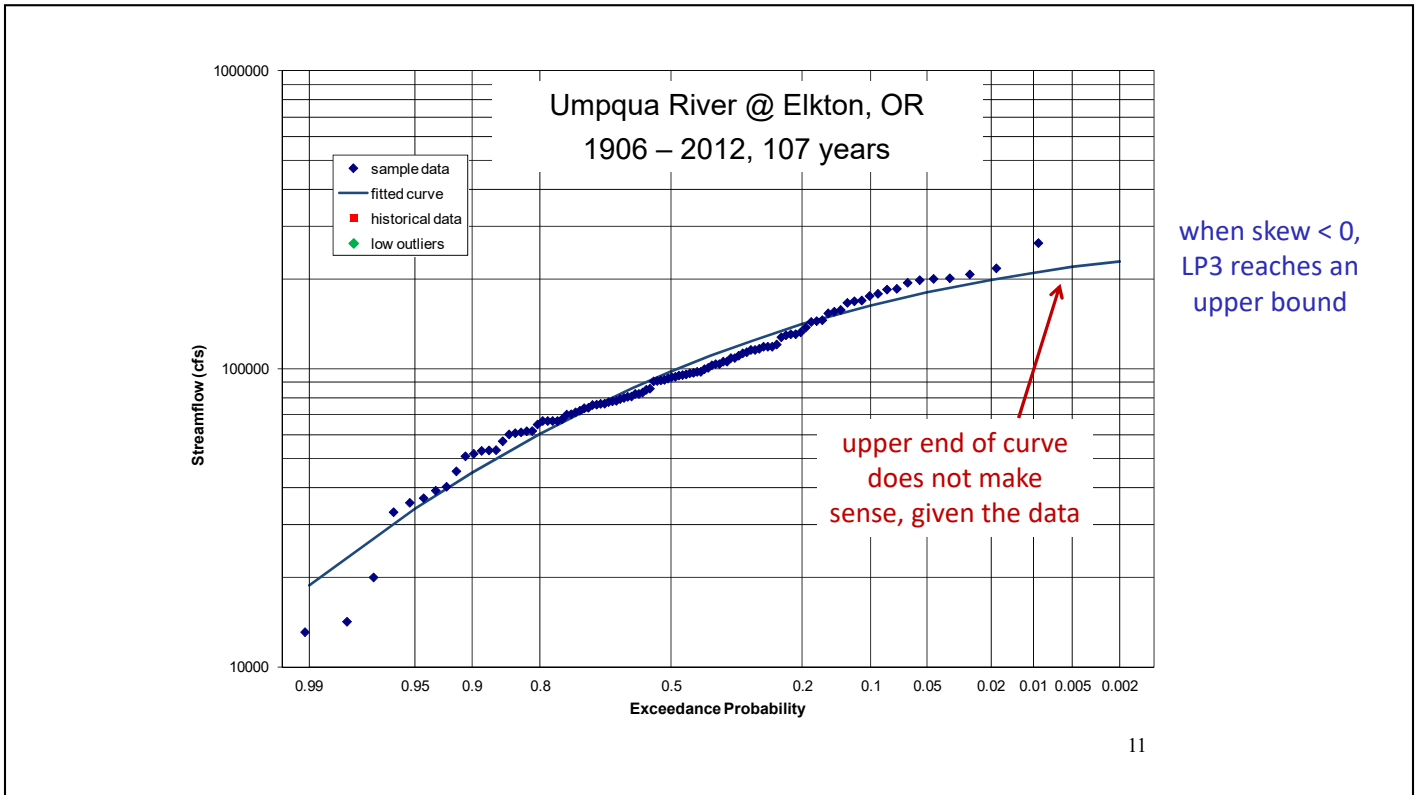
Outliers are values that are notably different from the rest of the data. B17B used a Grubbs-Beck test for both high and low outliers.

Here, CENSORED is a verb, ie to specify some years as censored means to acknowledge they occurred but not use their precise value in computation.

Low outliers are censored in this way. High outlier are not – they simply prompt the analyst to seek historical or paleo data. In Bulletin 17C, the concept of high outliers is not present, and historical data should always be sought.

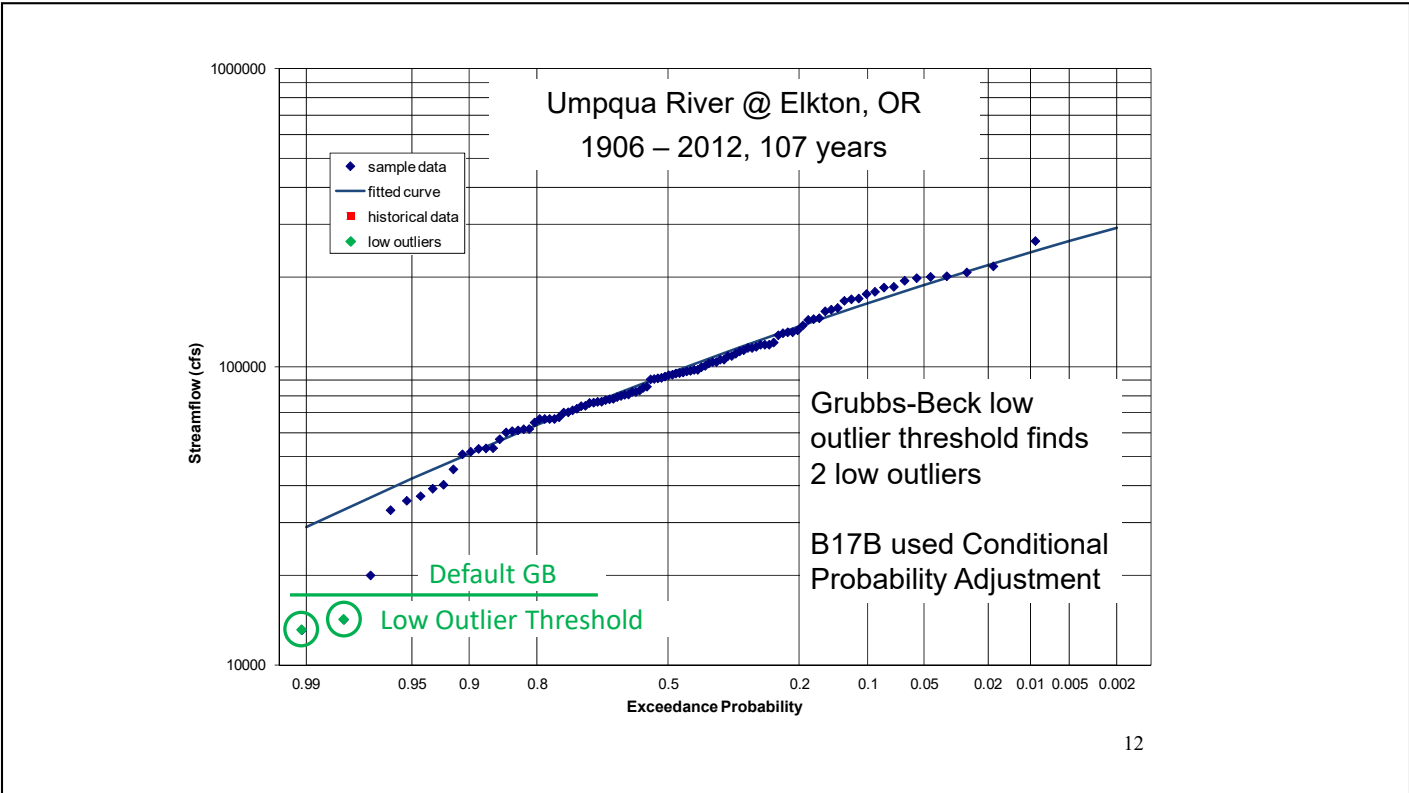


Here is the empirical distribution from plotted annual maximum points for a 107 year record in Oregon.

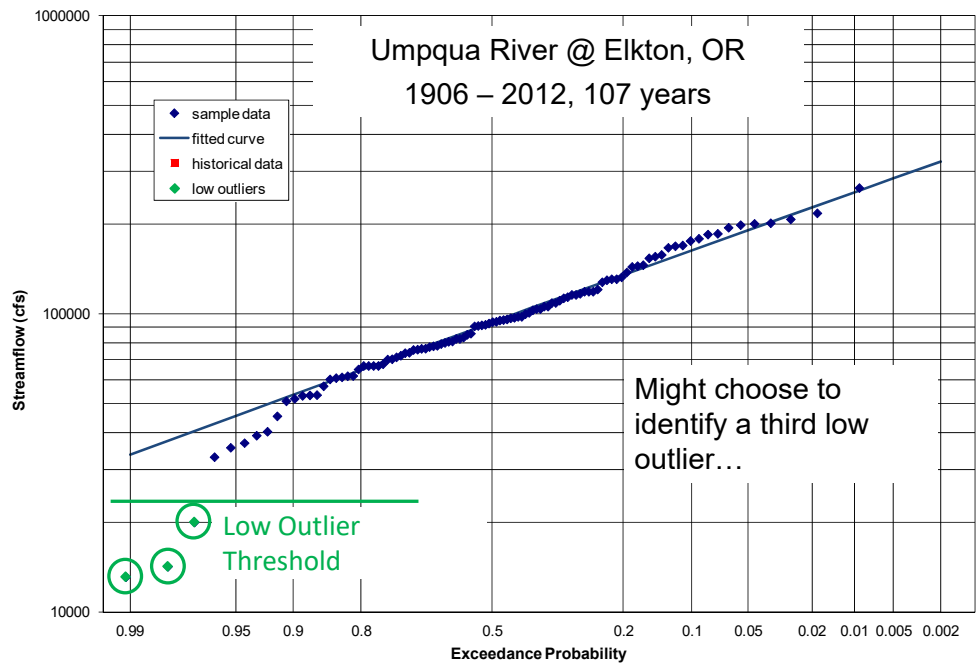


The LP3 estimate is not only a bad fit, but it will reach an upper bound that is below the largest value. (A negatively skewed LP3 distribution has an upper bound. A positively skewed LP3 has a lower bound.)

The problem is that the low values cause a negative skew, which pulls down the high end of the frequency curve as well. A probabilistic model that says the largest event could not have happened is not a good model.

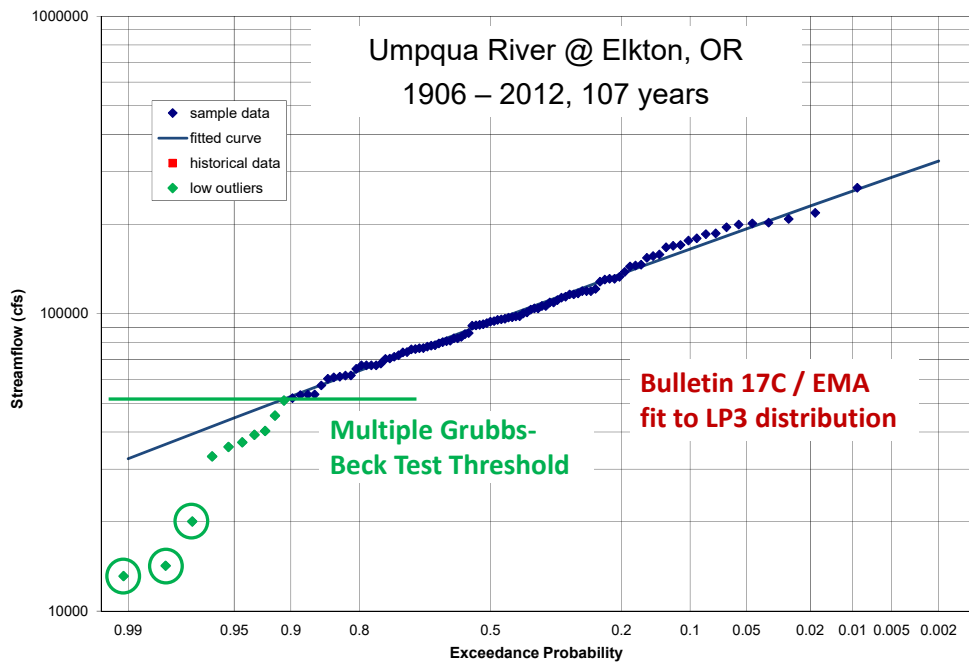


The GB test finds 2 values below the computed low outlier threshold. Censoring them (not using their precise values) increases the skew and raises both upper and lower tails of the frequency curve.



13

But why are these values low outliers? Are they a different flood type, or maybe not a flood at all? Hydrologically, might choose to also include the similar flow and censor 3...



14

Bulletin 17C has a new test that is more aggressive, and it censors at a higher level. Though engineering judgement is still needed.

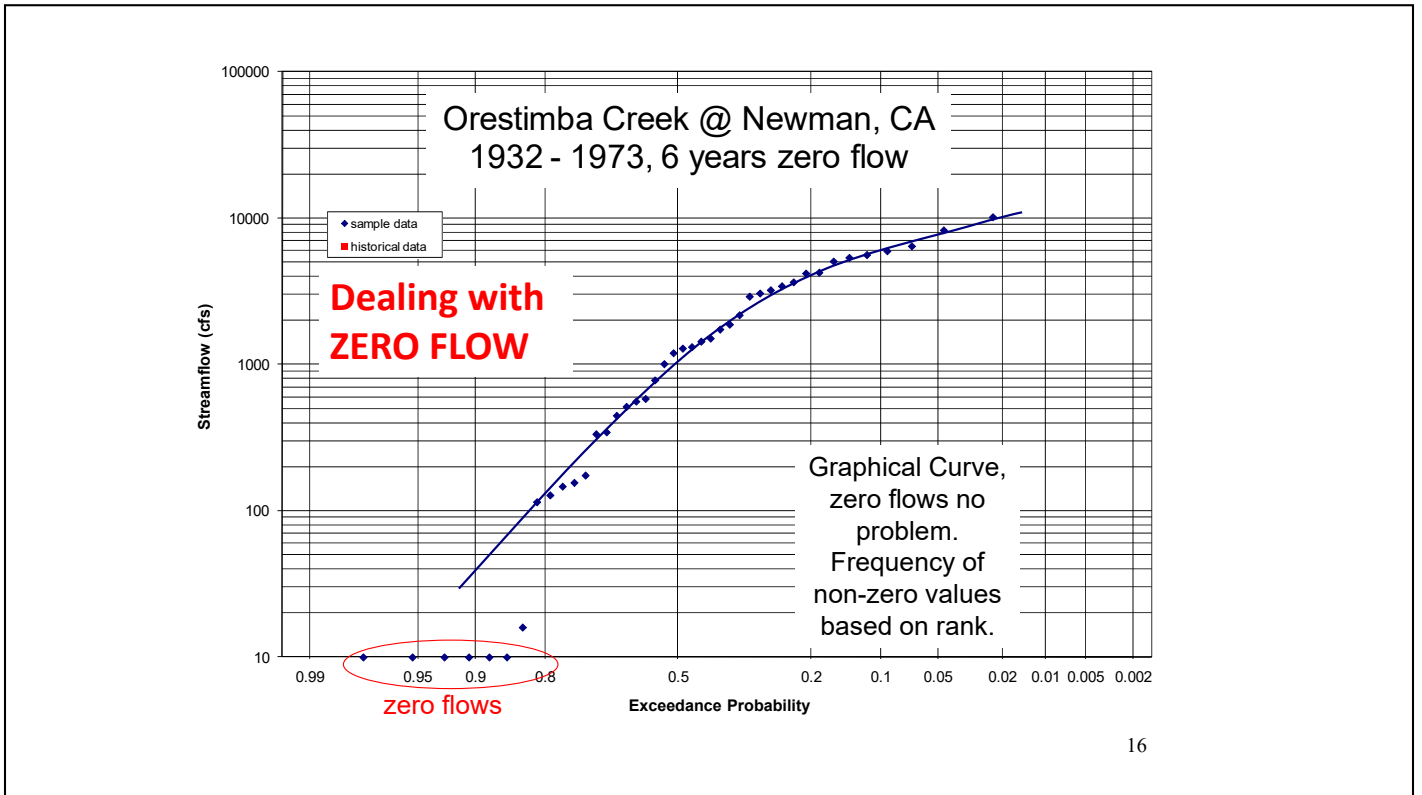
Any Values Below a Threshold

1. Values **below a recording level** are censored, by definition
 - but must be incorporated in the estimate of the flood frequency distribution – ie, they happened, flow was low
2. Flood years with **zero flow**
 - $\log(0)$ is undefined, handled as censored low flows
3. **Low-outliers** intentionally censored and accounted for the same way

In B17B, used the *Conditional Probability Adjustment* to deal with all of these cases

In B17C, these cases are handles with **flow intervals**

15



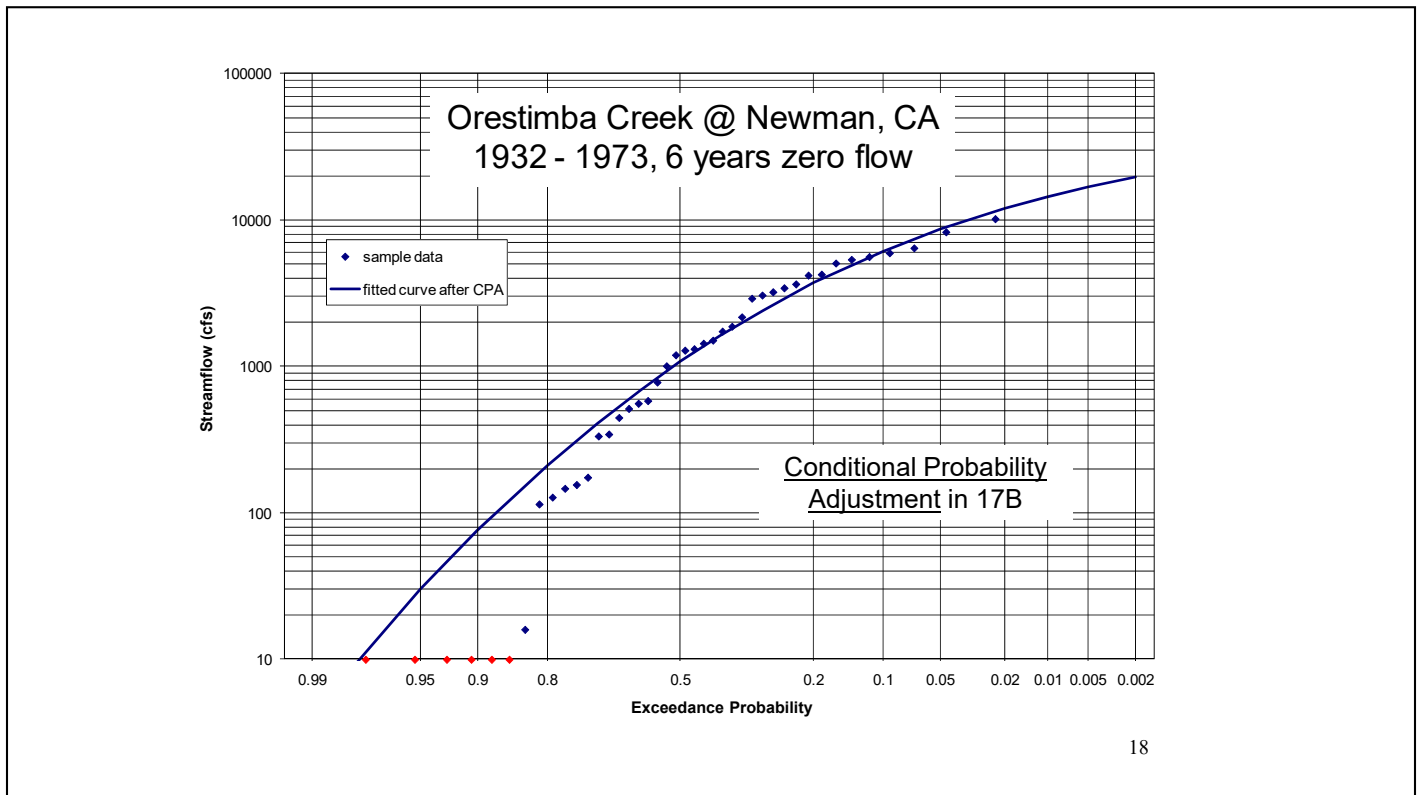
Zero flows are a problem with a log distribution because the parameters are based on log of flow, and the log of zero is undefined.

Defining the curve graphically would not have a problem with zero flows, because they would not affect the fit at the top, which is dependent only on plotting positions of the larger data points.

But an analytical distribution fit must account for all the data points. So the zero flows are treated as if they were censored flows.

Dealing with Zero flows

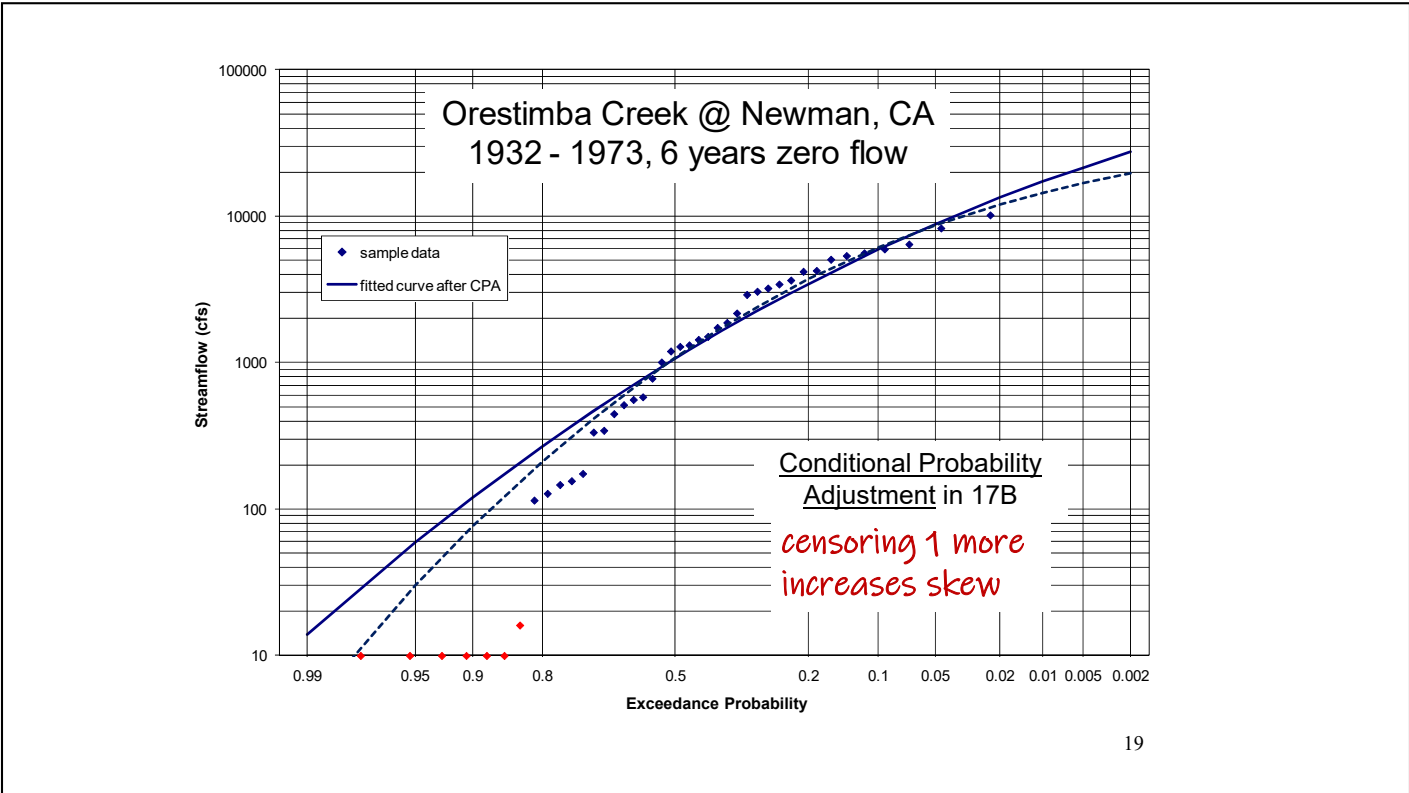
- Graphical Fit
- Adding small increment to flows (does not work well)
 - Standard deviation and skew are sensitive to small values (deviations from mean are squared or cubed)
- Use a different distribution (**no log transform**).
- Conditional Probability Adjustment (B17B)
- Expected Moments Algorithm (B17C)
- Maximum likelihood



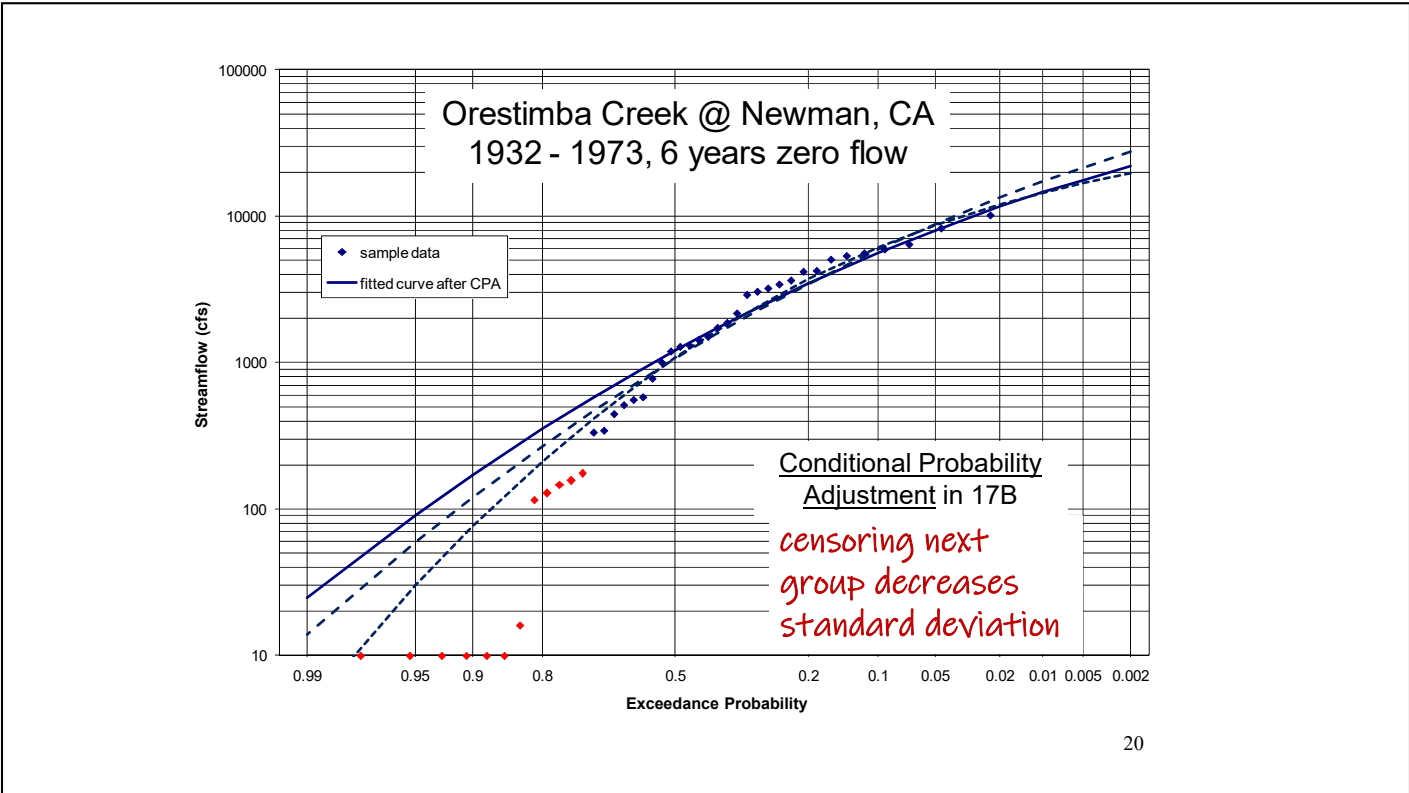
Bulletin 17B used a method called the conditional probability adjustment for accounting for censored flows. This curve results from just censoring the zeros.

The conditional probability adjustment have several steps.

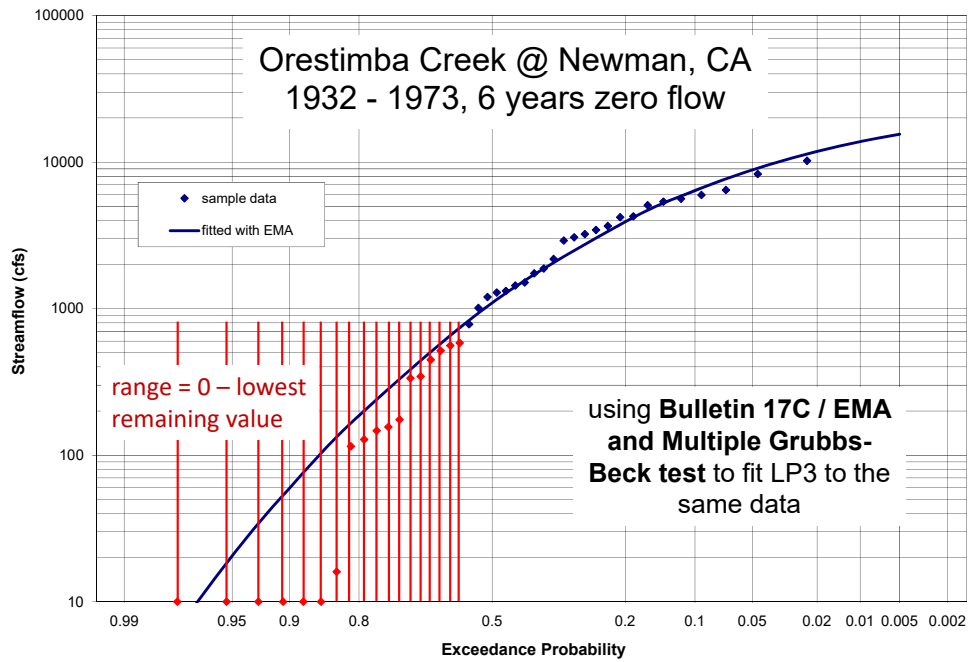
- (1) Fit an LP3 frequency curve with the non-censored data
- (2) Adjust the probability of each flow value based on the ratio of data points included to total data points. For example, if 35 of 43 points were included, a probability of 3/35 should more correctly be 3/43, and so all probability values were multiplied by 35/43 to correct.
- (3) The new frequency curve with adjusted probabilities is not LP, so the parameters of the LP3 are estimated with equation for “synthetic statistics”
 - (1) The synthetic statistics were based on Q(50%), Q(10%), and Q(1%), so the top half of the frequency curve, which is more important than the bottom.



The frequency curve changes to the solid from the dashed line when also censor the extremely low value. Removing the low value raises the skew, and so raises both top and bottom of the curve.



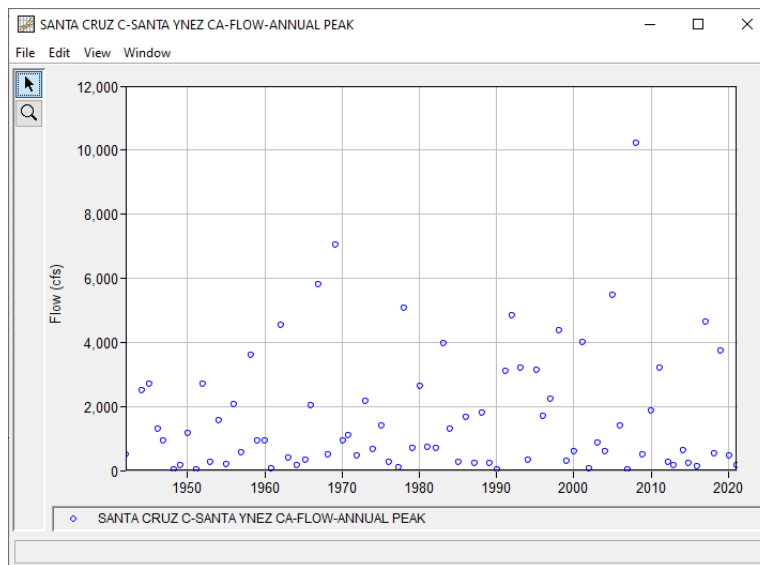
Censoring a few more low values to produce the solid blue line decreases the standard deviation, fitting the upper end more closely to the plotted points.



21

The multiple Grubbs Beck test in EMA censors even more values, replacing them with flow ranges.

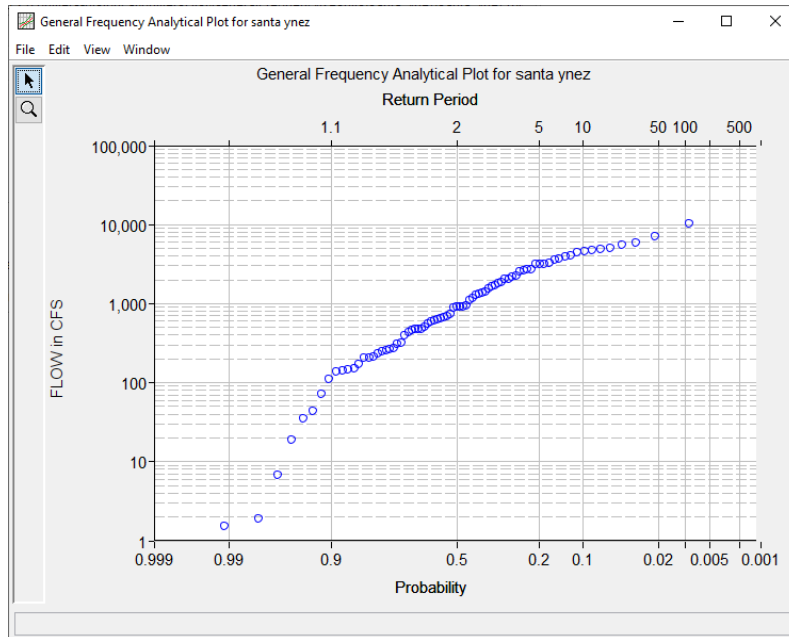
Another Example: Santa Cruz CREEK, CA



23

The Santa Cruz Creek at Santa Inez is another interesting example.

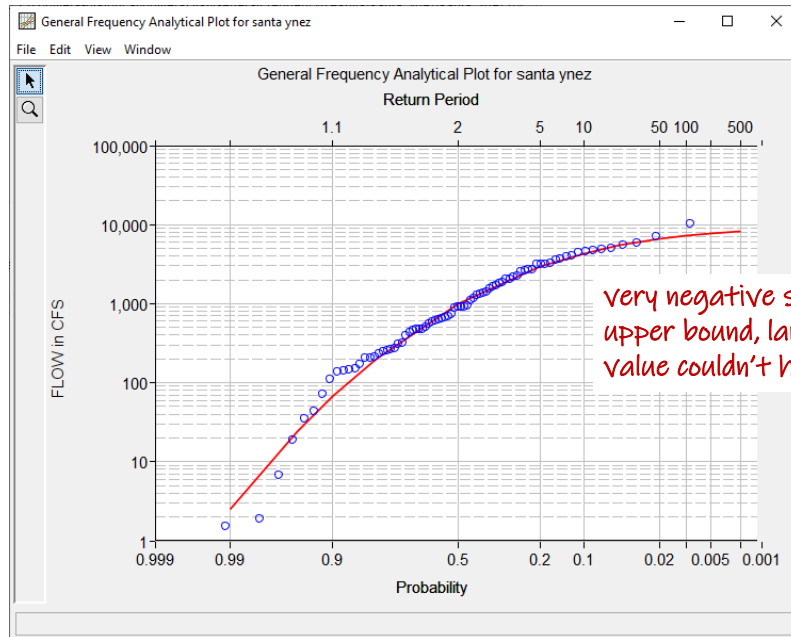
Another Example: Santa Cruz CREEK, CA



24

This figure is the empirical distribution of just the annual maximum flows versus plotting positions.

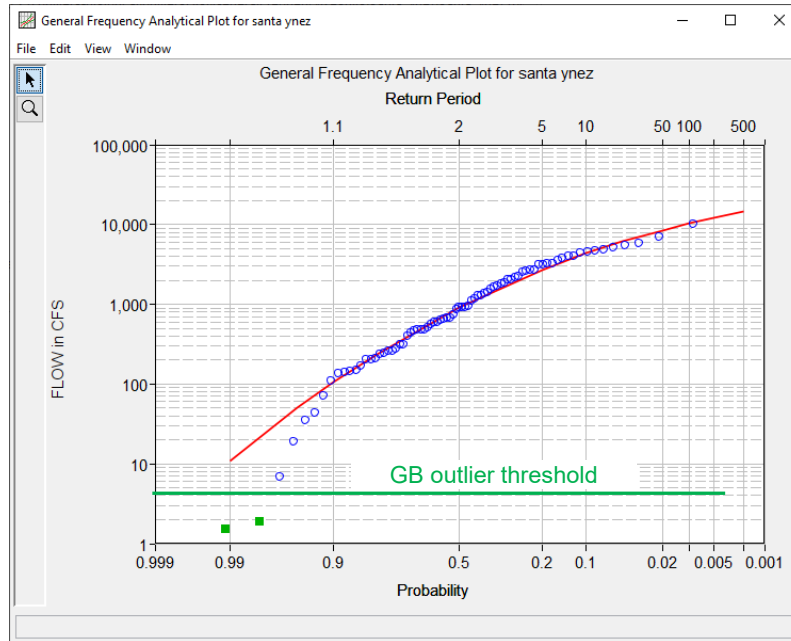
Another Example: Santa Cruz CREEK, CA



25

This example is another that the upper end of the frequency curve doesn't make sense. The upper bound of the negatively skewed distribution is below the two highest values. Clearly more low values must be censored to produce a distribution (model) that does make sense for the data.

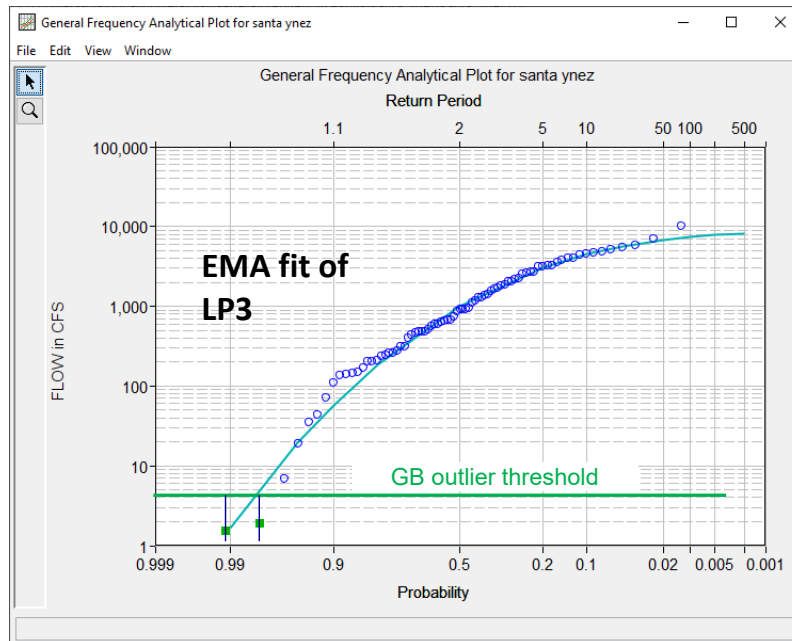
Another Example: Santa Cruz CREEK, CA



26

With the Grubbs-Beck test outlier threshold, two outliers are censored, and a sensible frequency curve results.

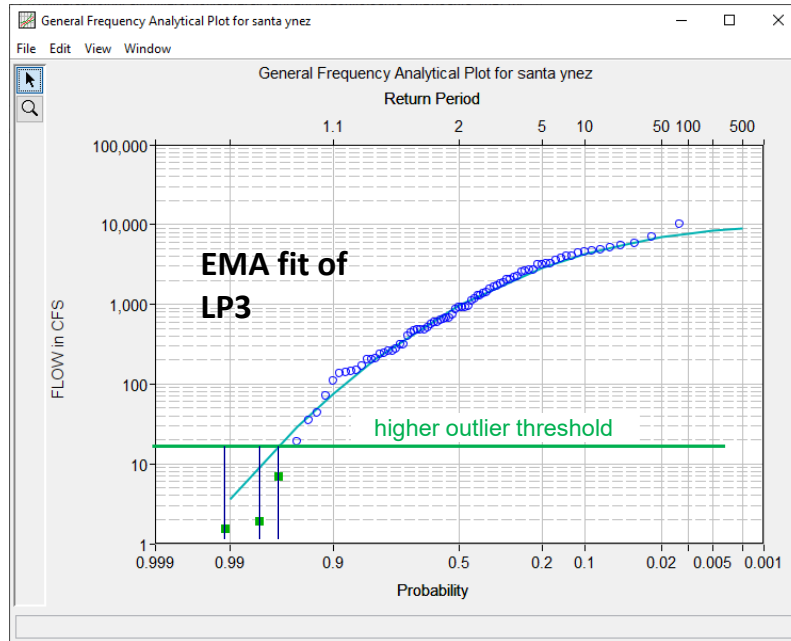
Another Example: Santa Cruz River



27

Using the same threshold choice to censor the lowest 2 values does not work as well with an EMA fit of LP3. EMA is much more sensitive to the values of the remaining low values.

Another Example: Santa Cruz River



28

Even using a higher threshold to censor another value doesn't make it much better.

...the Issue

- **EMA is more sensitive** to the low values
 - Some low values have an unwarranted effect on the upper end of the fitted frequency curve, but are not identified by Grubbs-Beck test
- Developed a **Multiple Grubbs-Beck** test – more aggressive
- Identifies values that depart significantly from the rest of the data, and so are potentially influential
 - **PILFs** *Potentially Influential Low Floods*
- Recodes them as censored observations, interval = (0, PILF-thresh)
- ***We don't want the lowest values to define the upper-tail***

29

Since EMA is more sensitive to the low values, more aggressive censoring is needed. A Multiple Grubbs Beck test was developed that identifies more low values.

Rather than call them outliers, because of their potential for influencing the upper end of the distribution adversely, they are called Potentially Influential Low Floods, or PILFs.

Potentially influential low floods (PILFs)

“PILF” describes small observations that may have an inappropriately large impact on the higher flood quantiles

A PILF might be:

→ high leverage

1. Unusually small flood given sample size and selected distribution
2. A flow that reflects a different physical process than the largest flow in other years
 - perhaps a zero or almost-zero flood year
3. Regular small observation (not “outlier”) whose potential impact upon estimated upper-tail quantiles is just too large

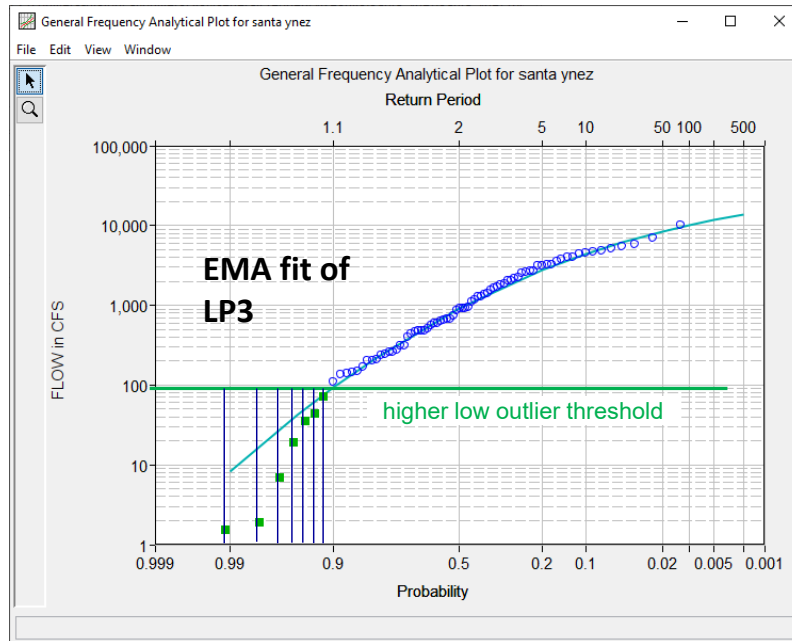
Robustness in Flood Frequency Analysis

- Why do we have this problem?
- Analytical distributions can't really duplicate the true (population) distribution of flood flows
 - Some flood records include ZERO FLOWS in drought years
 - LP3 is not the truth... Neither is GEV or LogNormal
- Need **robust** procedures that provide good estimates when the true physical process isn't captured by the simple model → MGBT

32

In general, we need a method that does a good job even when some assumptions aren't true or the model is too simple to represent the process. This is called a ROBUST method.

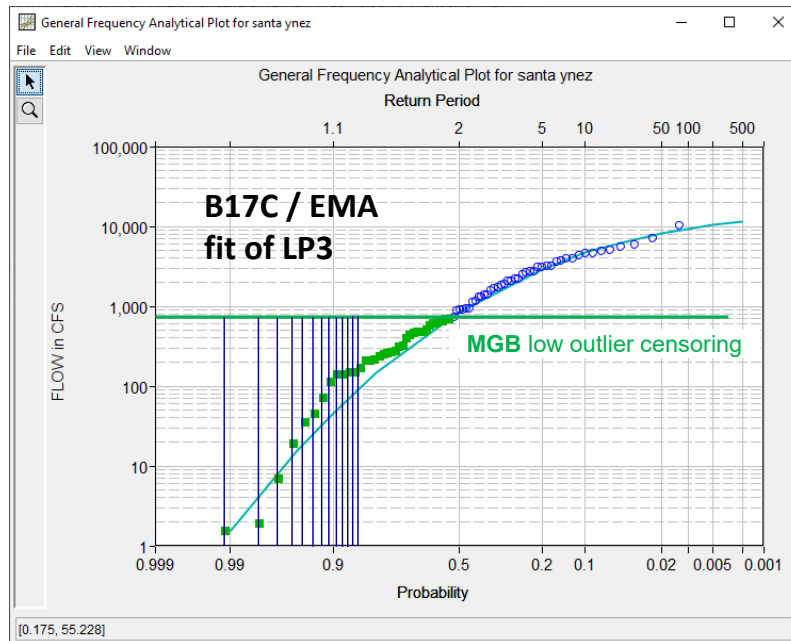
Another Example: Santa Cruz CREEK, CA



36

EMA gets a more successful fit when censoring several more values. It needs a more aggressive low outlier censoring threshold

Another Example: Santa Cruz CREEK, CA



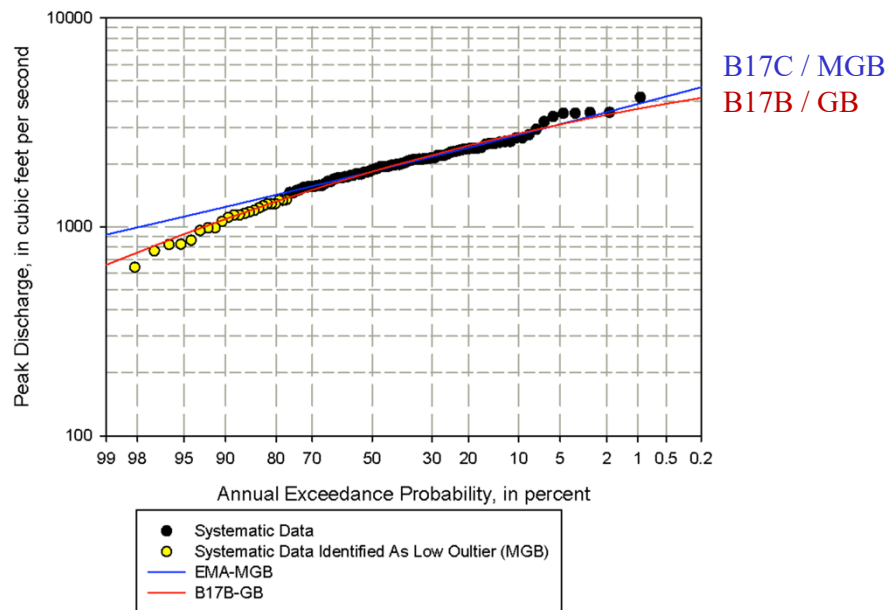
37

The Multiple Grubbs-Beck MGB test censors even more values. The fit is good, but it is up to the analyst to check the sensitivity to the censoring threshold to see if a more defensible value can be chosen while still maintaining a good fit.

The important aspect there is that that MGB test will result in a good fit in the first compute of the data set.

Comparison of B17B/GB (0) and EMA/MGB (24)

Weber River near Oakley, UT
(Station 10128500)



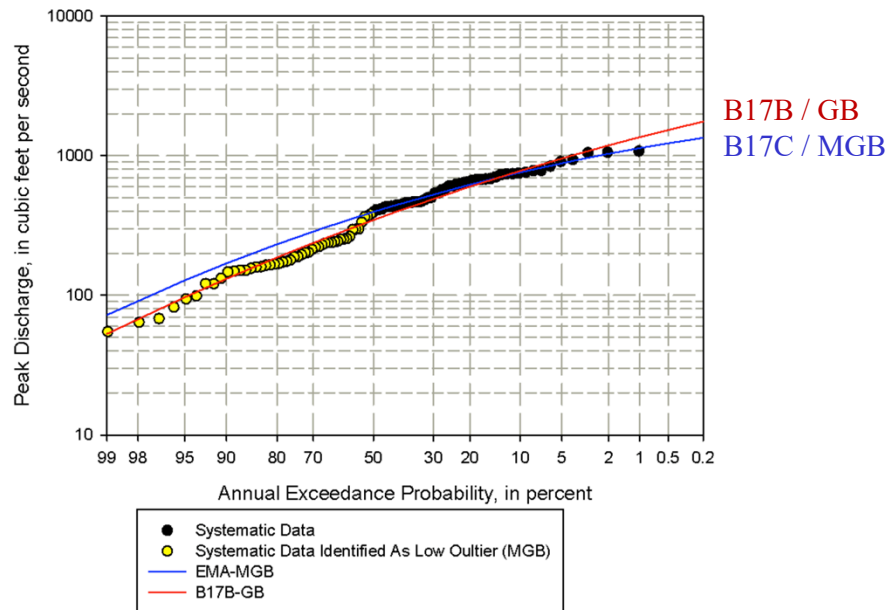
38

Here, largest event would be 1000 year event (or so) on red B17B curve. EMA curve says 200 year event

In this case, the MGB test and the censoring of many PILFs results in a higher upper end.

Comparison of B17B/GB (0) and EMA/MGB (48)

Beaver River near Beaver, UT
(Station 10234500)



39

What's probability that 3 largest events in 100 years do not exceed the 20 year flow?
(B17B fit says this.) Compute using binomial...

In this case, the MGB test and the censoring of many PILFs results in a lower upper end.

Outline

- Missing values
- Censored values, zero flows, outliers
 - low flows (PILFs)
 - high flows
- Historical/Paleo information

Values Above a Threshold

- Since we are interested in extreme flow, we should estimate any flows that were **above** the highest recording level, or know highs **before** the gage.
- B17B noted High Outliers:
 - *If information is available that implies the high value is the **largest in a longer period of time**, the largest events during a systematic period can be given **historical weighting** in the Bulletin 17B computation procedure.*
 - *Therefore, when there are high outliers, should **seek historical information**, even non-exceedance*

41

B17B Outliers: values notably different from the rest of the data

Low Outliers are censored

- analysis excludes the values, but not the fact they occurred
- B17B used Conditional Probability Adjustment
- B17C uses intervals for excluded values

B17B used Grubbs-Beck test for high and low outlier thresholds

High Outliers are **NOT** censored

- high values are left in the data set
- seek historical information to either add a past large event, or determine that the largest gaged event has a longer return period (*is the largest in a longer period of time.*)
- B17B used Weighted Moment Algorithm
- B17C uses intervals for unobserved years

B17C uses Multiple GB

42

Outliers are values that are notably different from the rest of the data. B17B used a Grubbs-Beck test for both high and low outliers.

Bulletin 17C doesn't designate high outliers at all. But there is still a suggestion to seek historical and paleo information if its available.

Historical/Paleo Information

Historical Information = evidence from human records

Paleo Information = physical evidence in the watershed

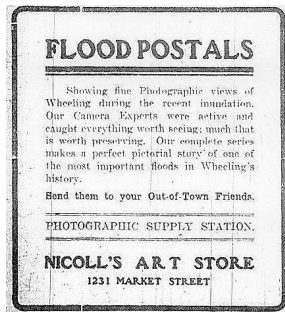
- Systematic record = years in which flow is gaged
- Historical period = includes years that are not gaged, but some information is known, for example:
 - Large flow happened, and can be estimated
 - There were flows that exceeded some threshold (at least 1), or
 - There were **NO** flows that exceeded some threshold
 - For example: might determine that systematic event in 1997 was not exceeded since 1862 = 161 years

43

43

Historic Flood Information

Wheeling (Ohio) Flood (1907)



“...The **worst flood since the memorable 1884** flood now holds sway in the Ohio valley. A new high water record has been established in Pittsburg, and though the mark of '84 was not passed at Wheeling the second flood stage to that destructive water will be attained here this morning. ...”

--*The Intelligencer*, March 15, 1907, p. 1

(Source: Tim Cohn, USGS)

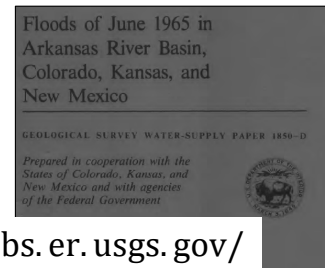
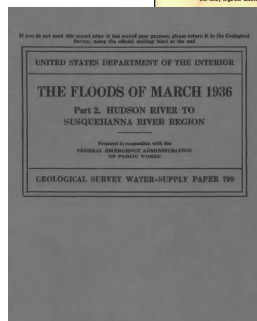
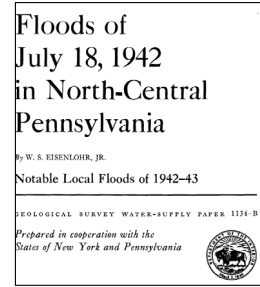
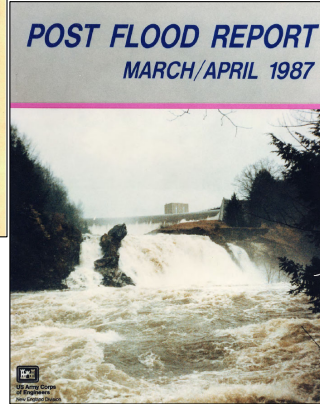
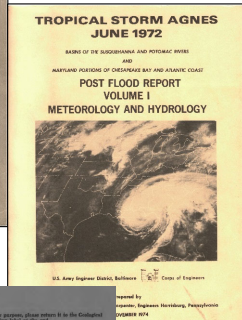
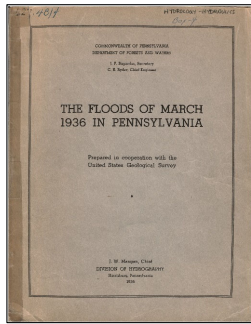
44

44

Historical information might be a newspaper article, such as this one about a flood on the Ohio River in 1907. The article also mentions it being the worst flood since the 1884 flood, so that provides an additional piece of information about the minimum return period of the 1907 flood.

Sources of Data

Post Flood Reports



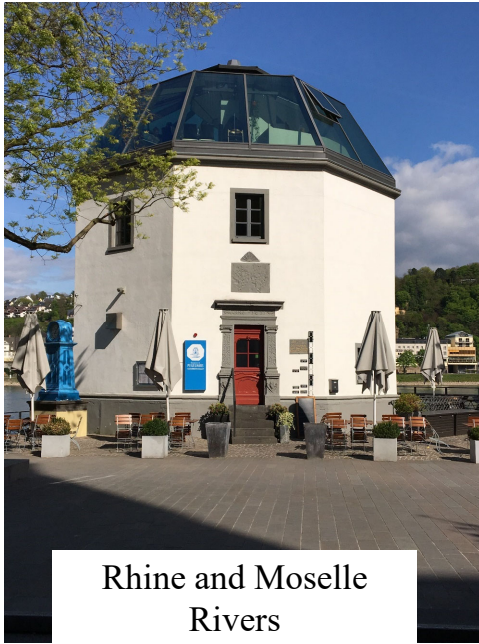
Source:
Mike Bartles

<https://pubs.er.usgs.gov/>

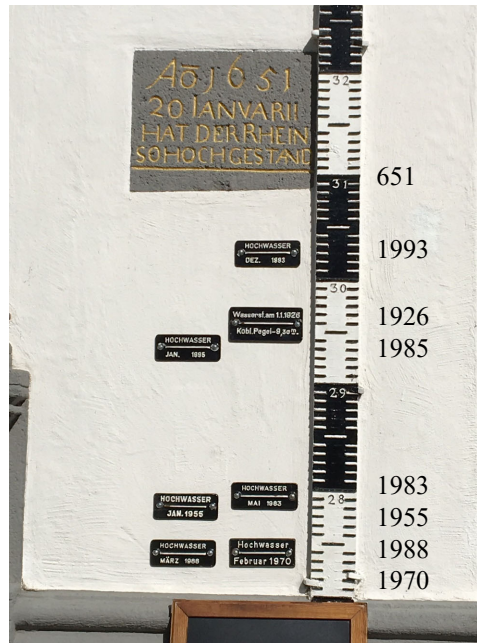
45

- Post Flood Reports are oftentimes prepared following major flood events. Information within these reports can be used for a multitude of purposes.
- These reports are available through a multitude of sources ranging from the U.S. Geological Survey, USACE districts, Bureau of Reclamation, state agencies, etc.
- Google is your friend when searching for these types of documents.
- The following slides present just a small sample of the information that is usually contained within these publications. Getting copies of them and doing research is a MUST!

Historical flood data



Rhine and Moselle
Rivers
Koblenz, Germany



46

Here's an example of historical information shown as high water marks – going back to the year 651.



Potomac River at Great Falls Park, VA

(Source: John England, USBR)

Crooked River near Prineville, OR, 1861 flood

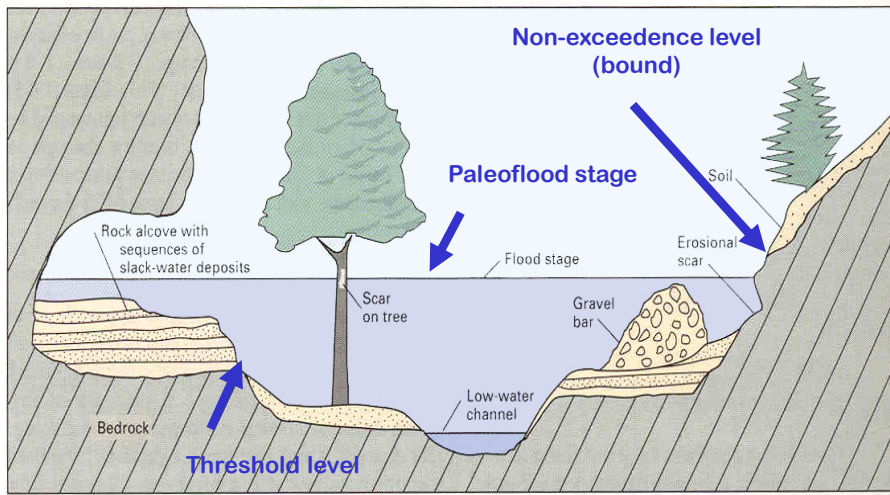


47

These are other forms of historical information. The flood level sign along the Potomac River shows the levels of the large floods.

Paleo: The “schooner” trees in Oregon are physical evidence of a large flood. A flood large enough to knock over the trees happened, but they lived and sent up new trunks vertically. By coring and dating the new trunk, we can determine the data of the flood. This flood is a threshold exceedance, because we only know it was at least to the level of the tree, but not how much higher it was.

Paleoflood Data Sources

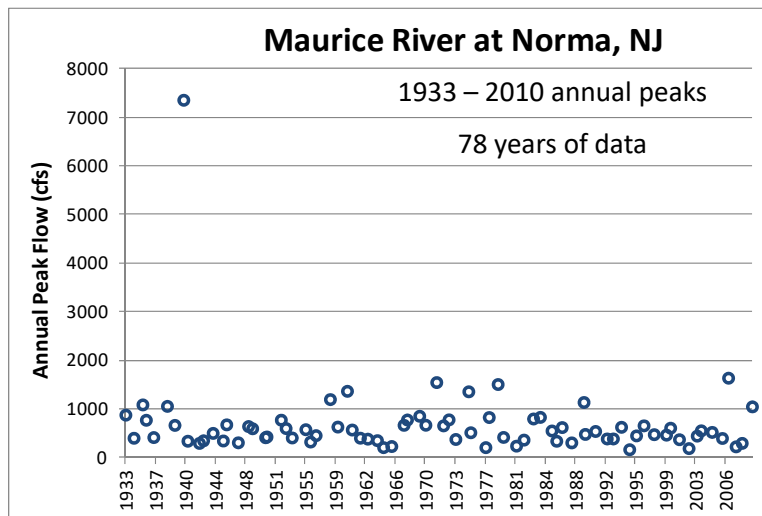


(Source: Jarrett 1991, modified from Baker 1987)

48

These are other types of physical evidence of large floods, such as tree scars, slack water deposits and a non-exceedance bound.

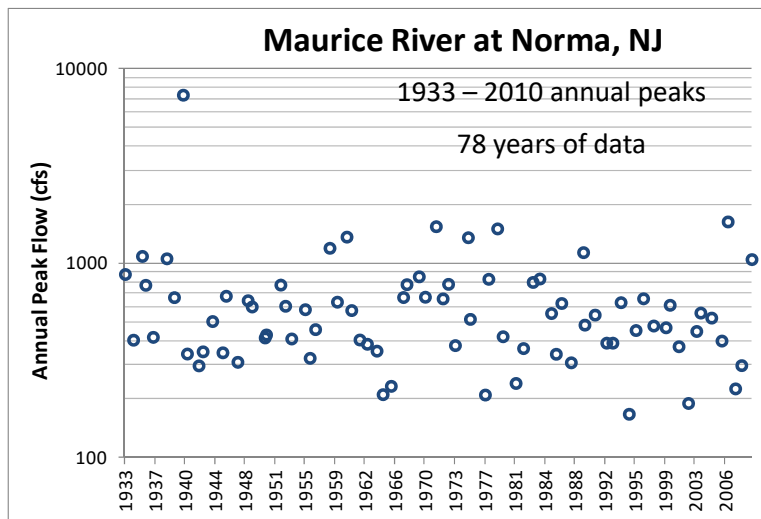
Example: High Outlier, Historical Info



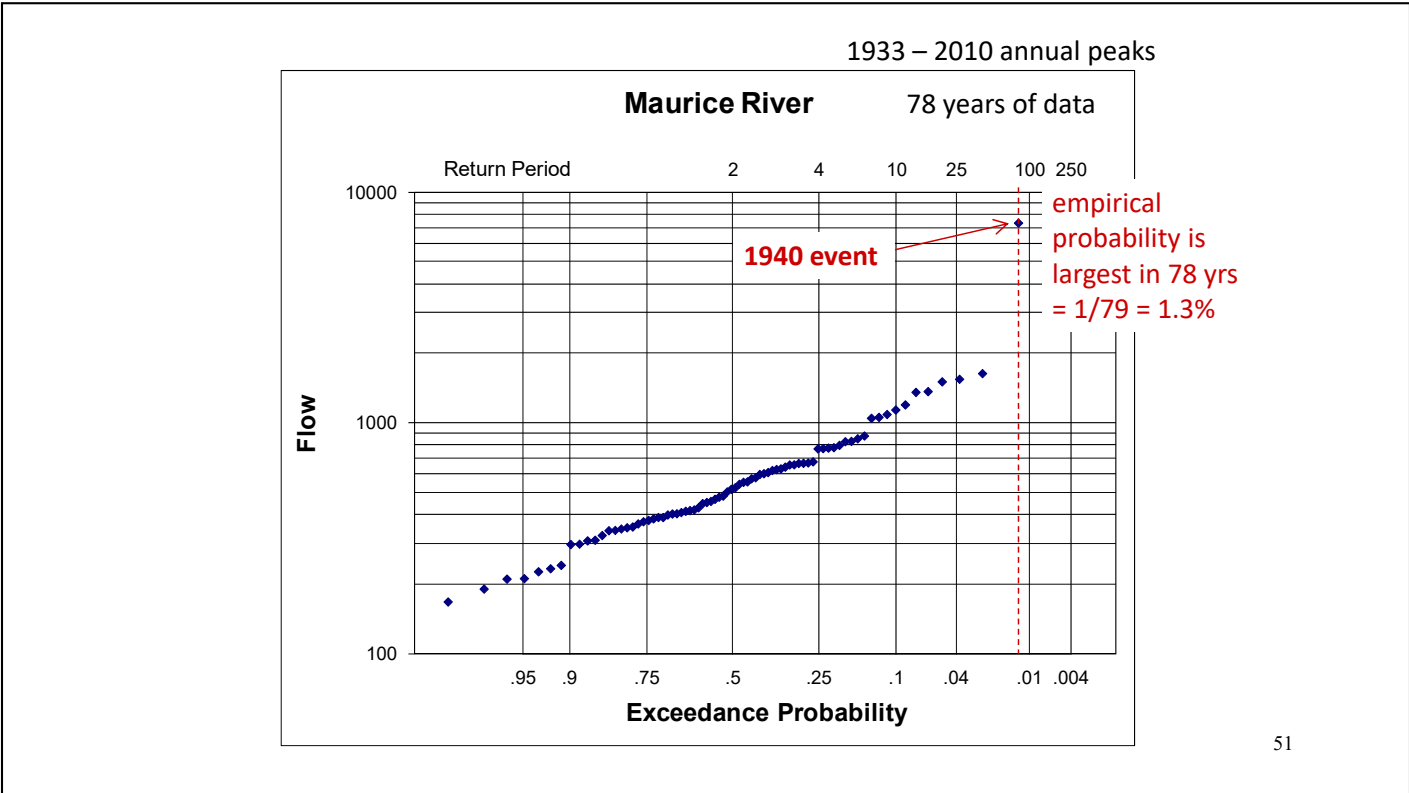
49

Here's an example that includes what Bulletin 17B would have called a high outlier.

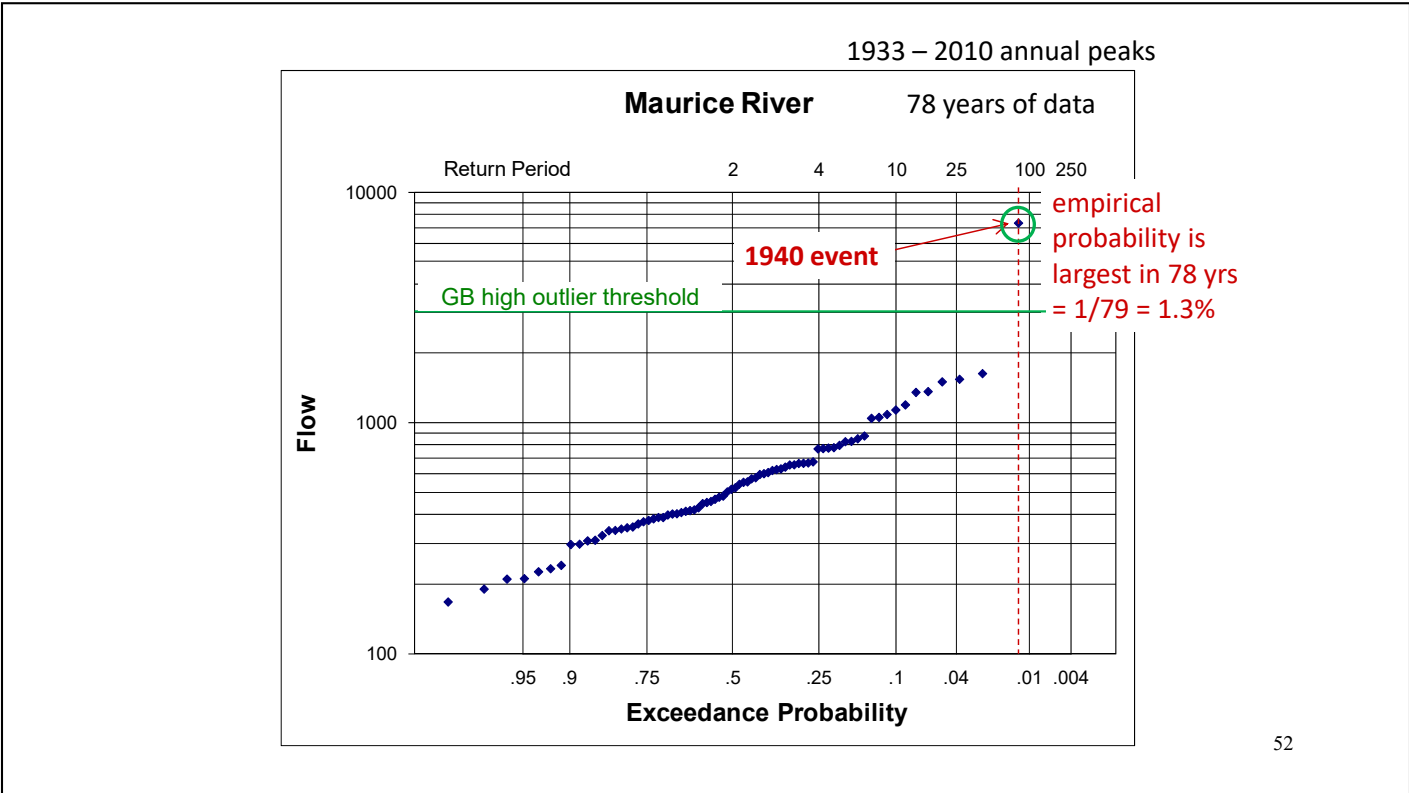
Example: High Outlier, Historical Info



50

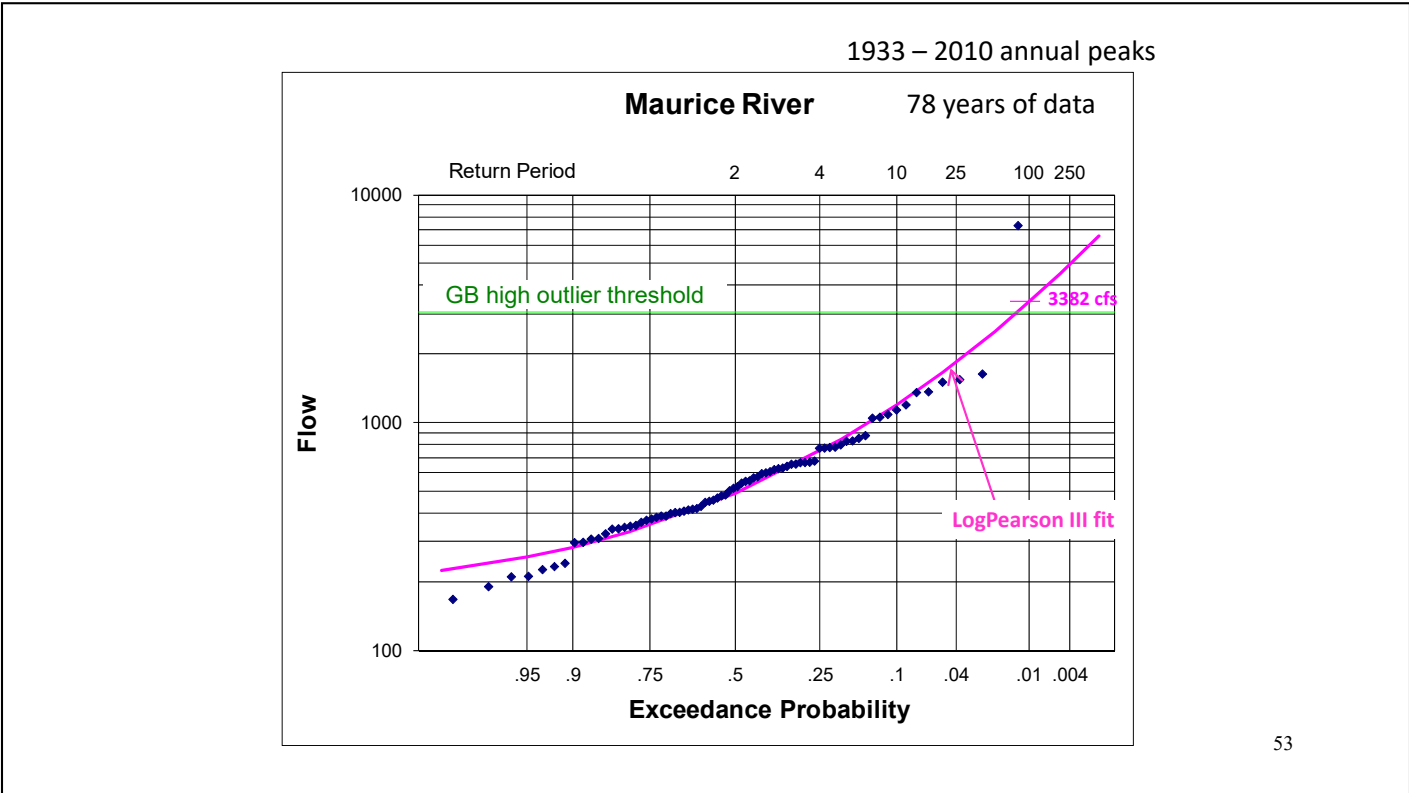


In the original data set, the most we can say about the large flow in 1940 is that it's the largest in 78 years. The fitted LP3 distribution will make that assumption. But, the return period of that event is probably greater than 78 years.



52

The Grubs-Beck test would call this a high outlier, suggesting we seek historical information.



This figure has the simple MoM (Method of Moments) LP3 fit to the gaged data.

The following are excerpts from a newspaper account of the 1940 flood:

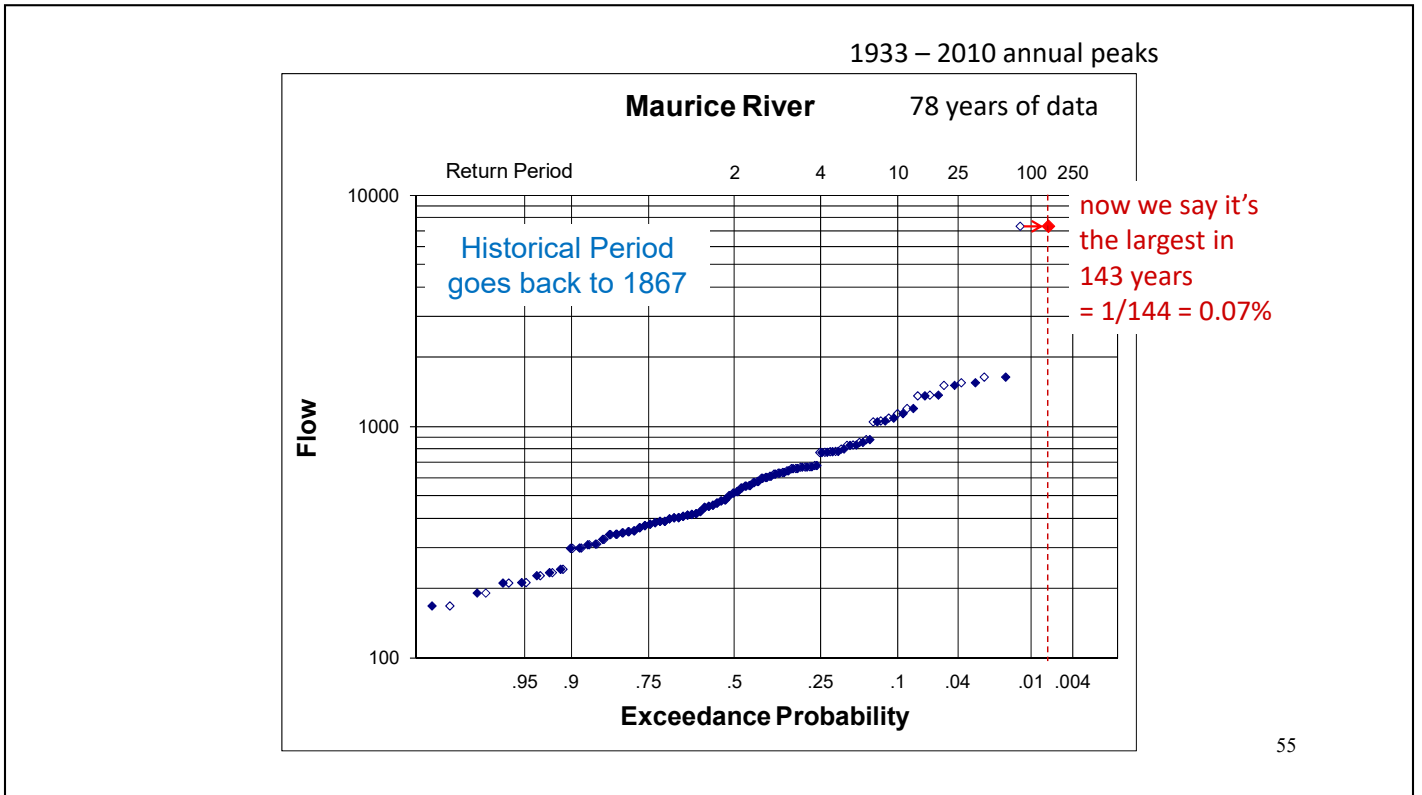
Millville Daily Republican
Tuesday, September 3, [1940](#)

“Yesterday’s flood conditions were the [worst ever](#) to strike Millville. In all the years other communities have suffered from floods, Millville has been unscathed....

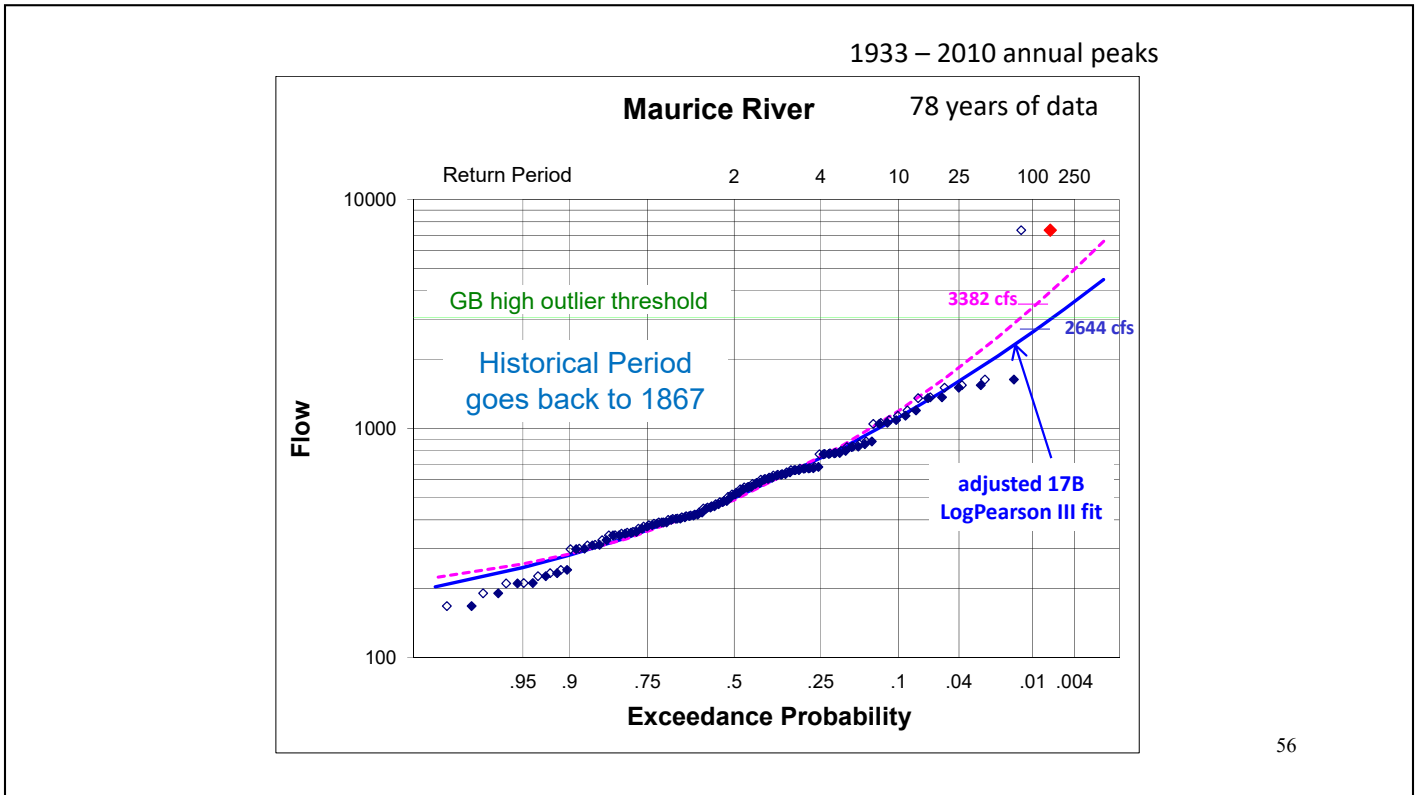
....Dozens of workmen continuously piled sand and sandbags high along the sluiceway inside the Millville Manufacturing Company yard. [Millville’s dam at Union Lake, built in 1867](#), held fast all day yesterday even though there [was more pressure brought to bear against it than at any time in its long and useful career.](#)”

54

This newspaper article from 1940 says that the flood event is the worst known. It mentions a dam build in 1867 which has never experienced a worse flood, meaning the 1940 is the largest back to at least 1867.

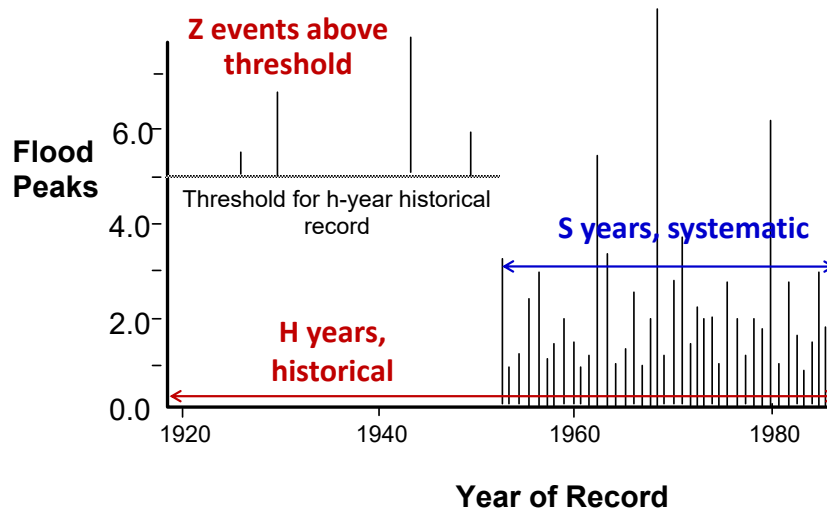


The plotting position of the 1940 event can be adjusted to largest in 143 years. This adjustment does not change the LP3 fit, but does show visually what assumptions we can make about the 1940 event.



The LP3 curve adjusted to include 1940 as largest since 1867 has a lower skew and so lower upper end.

Historical and Gauged Record

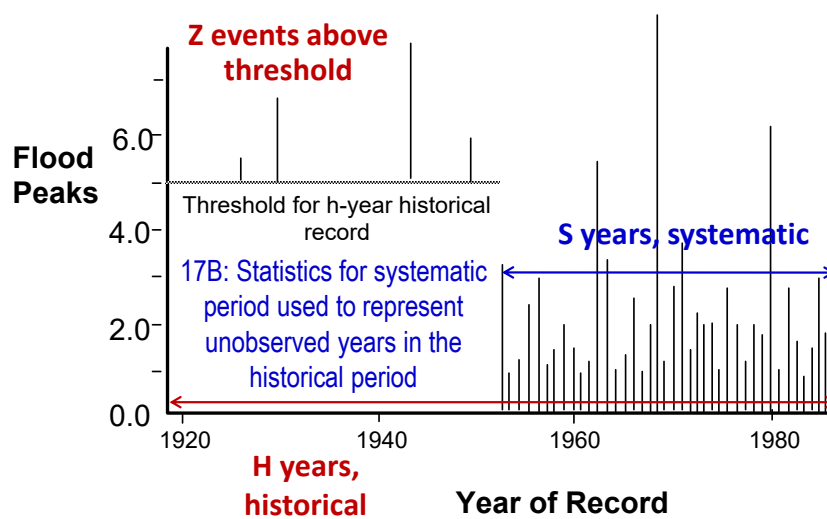


In Bulletin 17B, we adjust the statistics of the LP3 distribution by assuming the characteristics of the systematic record apply across the entire record $S + H - Z$, and then include Z

57

Systematic years plus historical years including some historical events that exceed a perception threshold.

Historical and Gauged Record

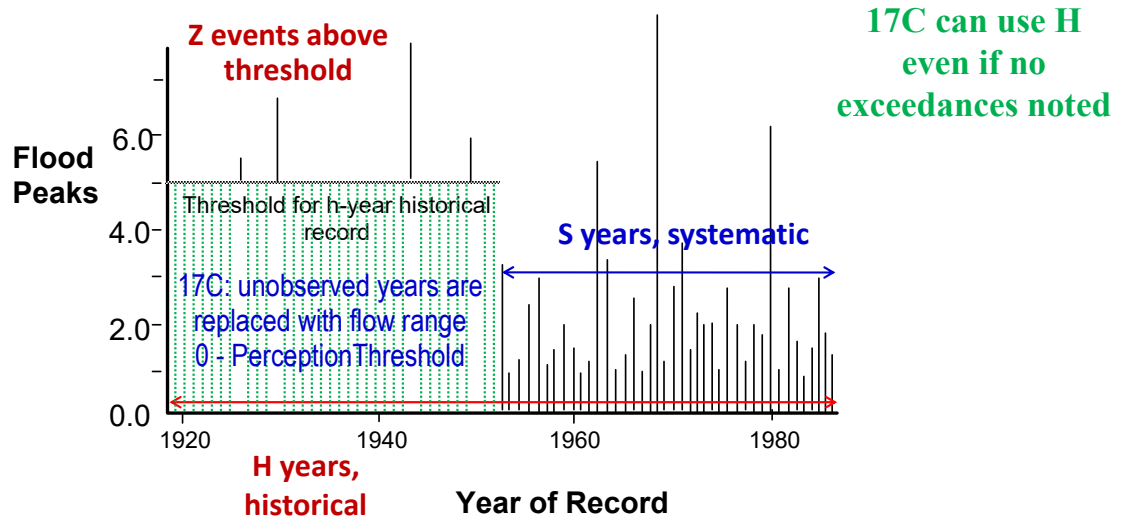


17B: Systematic years were given a weight greater than 1 to represent unobserved historical years when computing sample statistics

58

The 17B weighted moments algorithm used the statistics of the systematic record to represent the unobserved years in the historical period. This worked very well for a limited historical period.

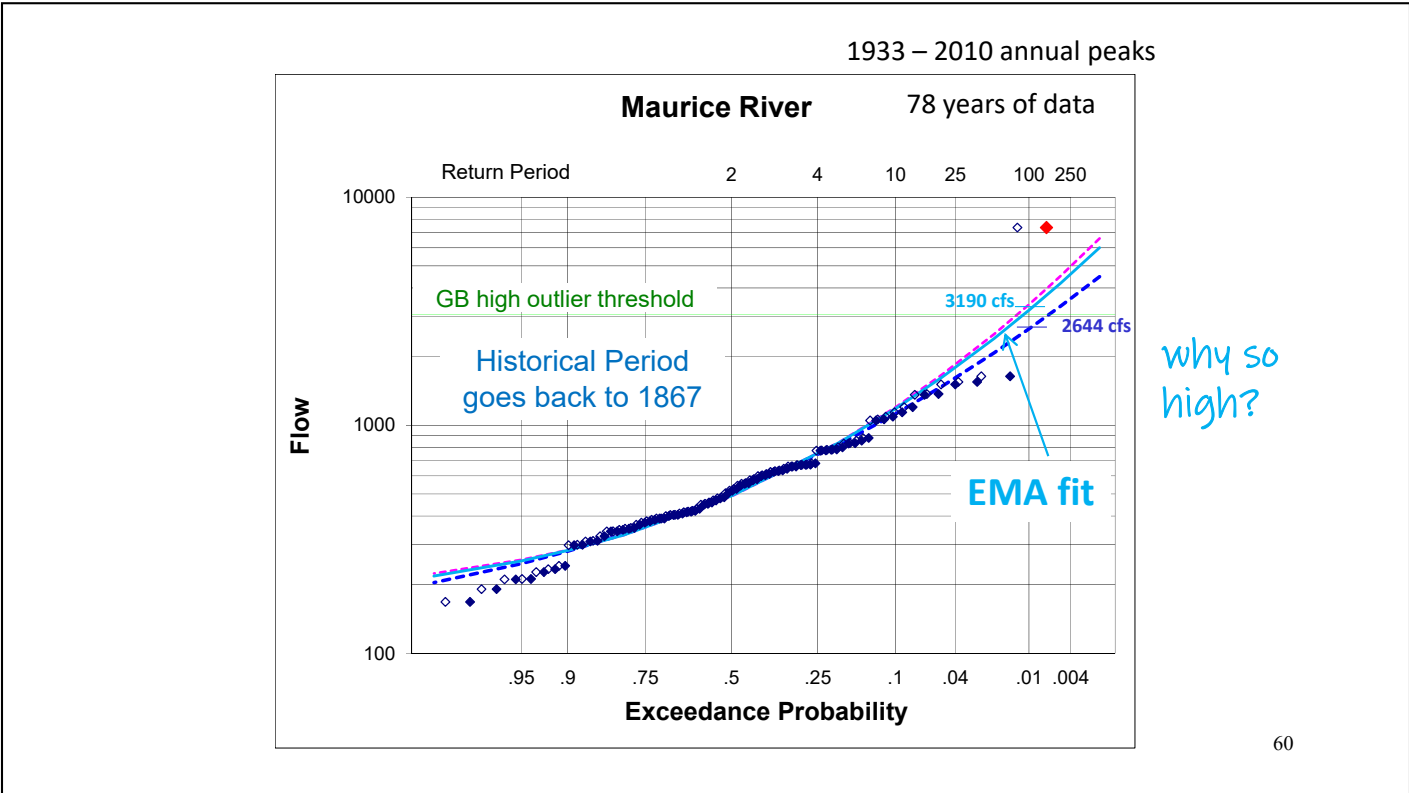
Historical and Gauged Record



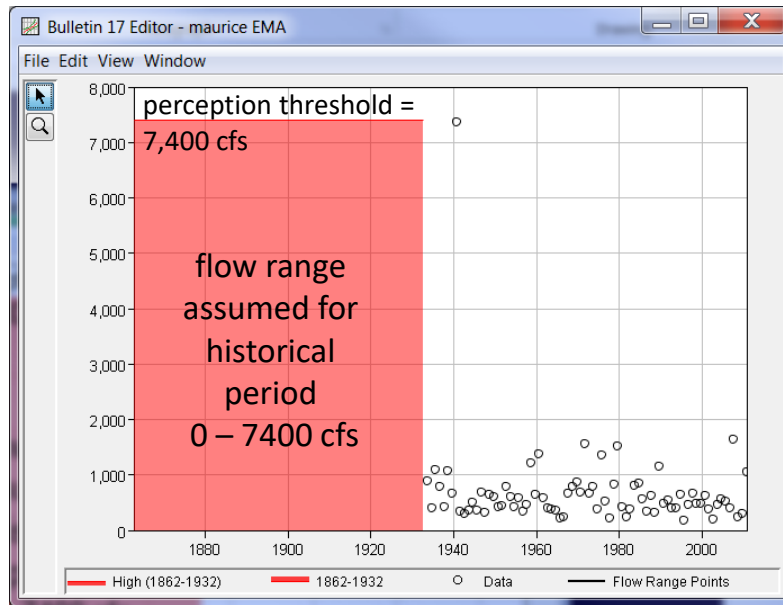
In 17C, include flow ranges for the unobserved years equal to the complement of the perception range

59

Bulletin 17C replaces the unobserved years with a range of flow from 0 to the perception threshold. This method works better for much longer historical periods from paleo information.

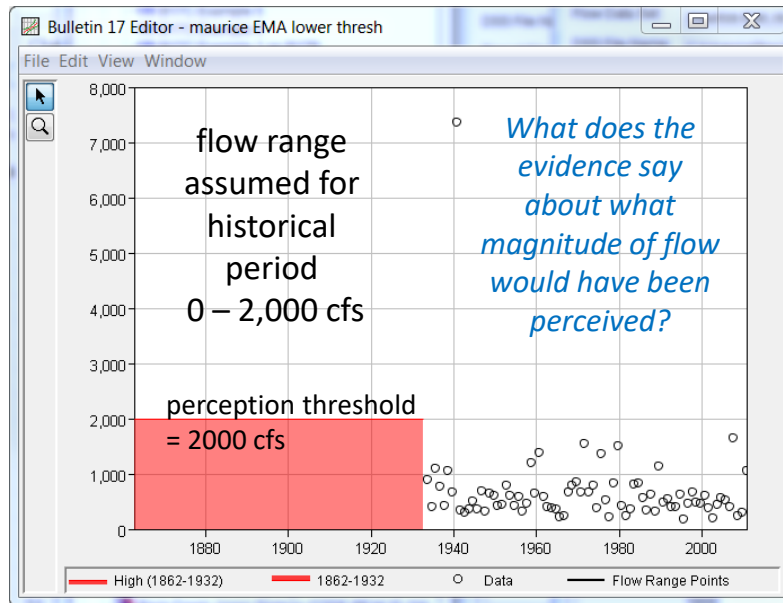


This curve shows the initial EMA fit of the same data, with 1940 largest since 1867. However, the curve seems as high at the original....



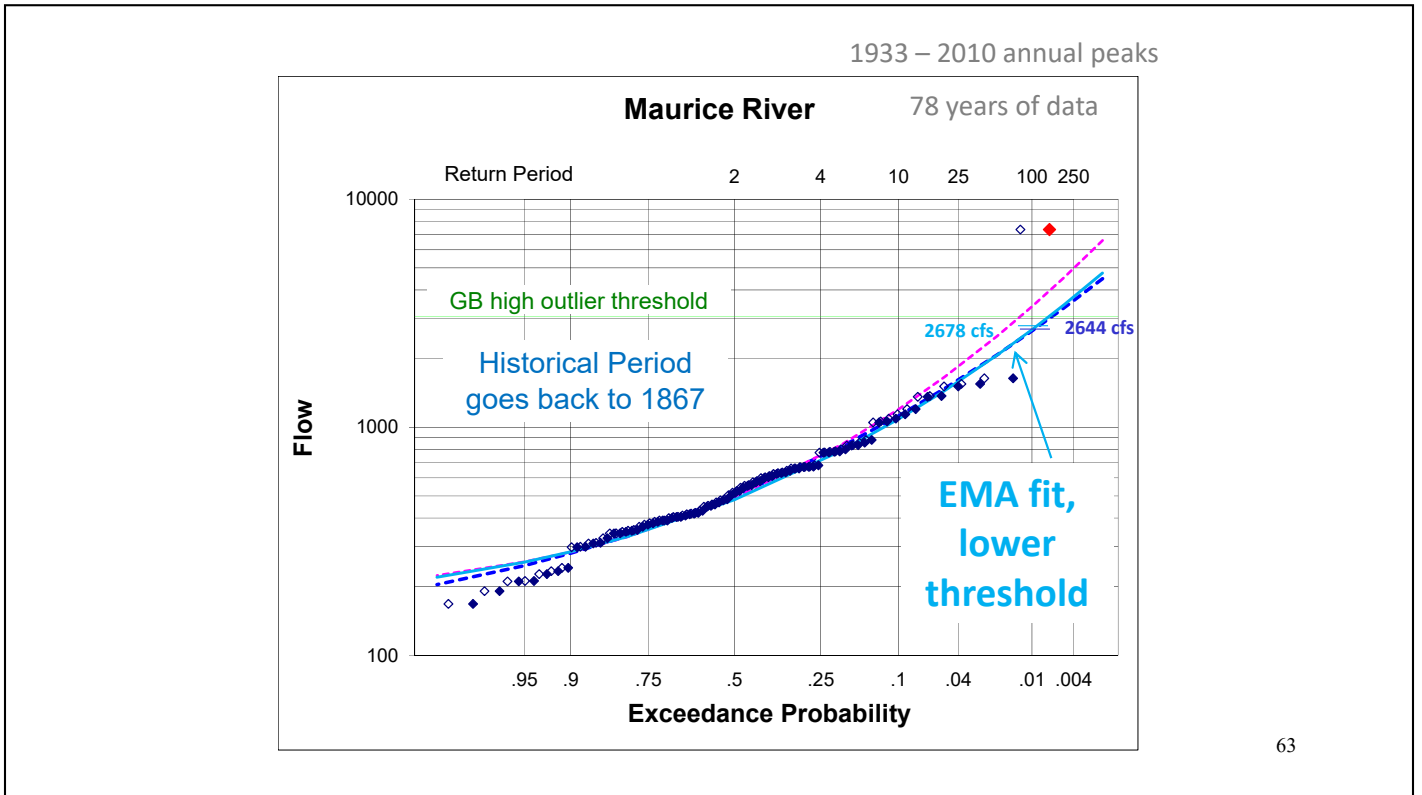
61

The LP3 curve fit with EMA was high because the perception threshold for the historical period was as high as the event itself, implying that anything lower than 7,400 cfs would not have been perceived.



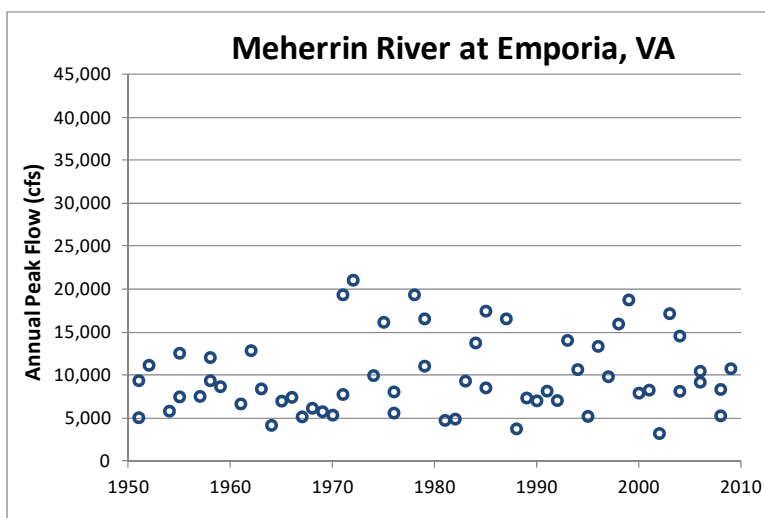
62

It is more likely that events much smaller than 7,400 cfs would have been perceived and recorded if they occurred in that period of time. Perhaps as low as 2,000 cfs. So, a perception threshold is set at 2000 cfs, meaning that flow ranges of 0 to 2000 cfs are used for the unobserved years between 1867 and 1935.



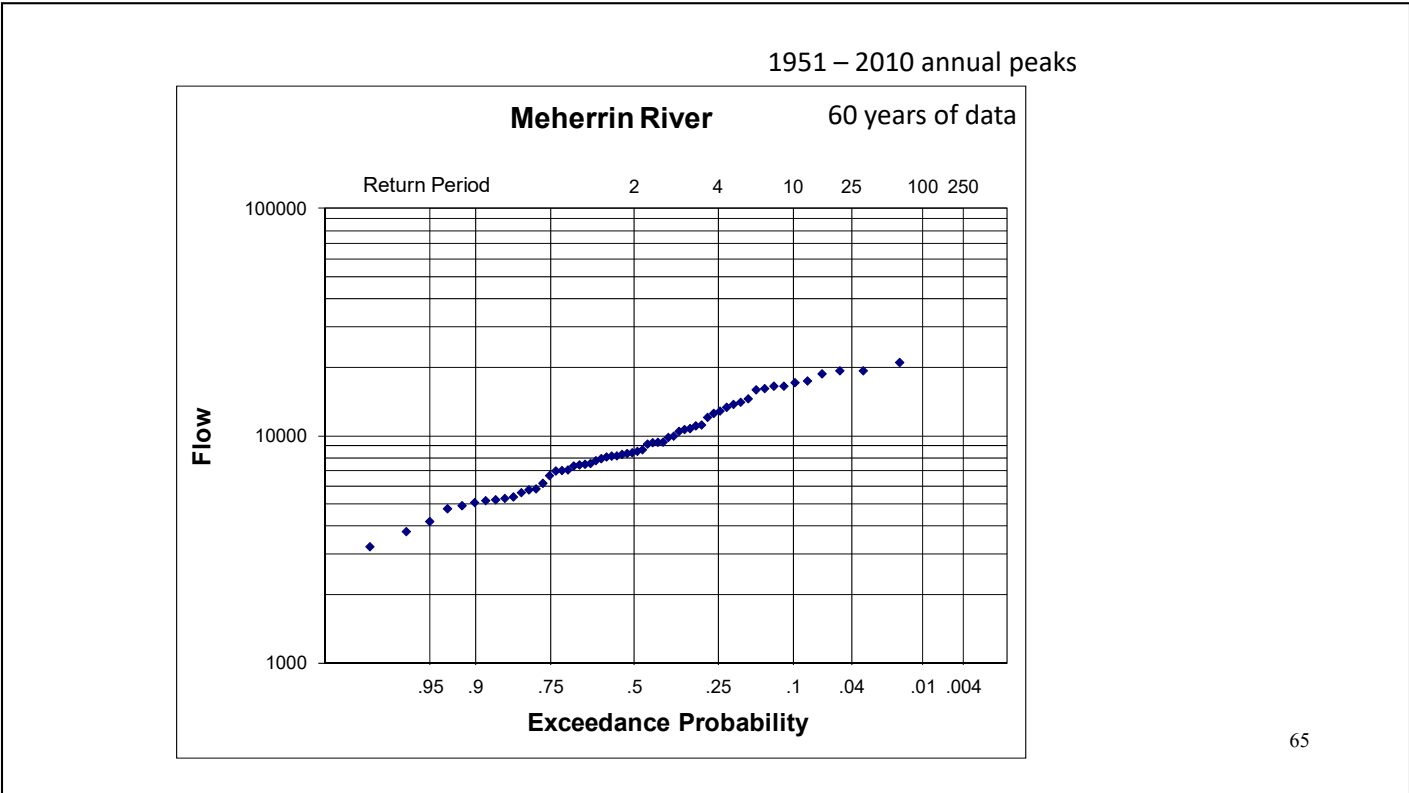
The resulting LP3 fit by EMA is therefore lower at the upper end.

Example: historical information

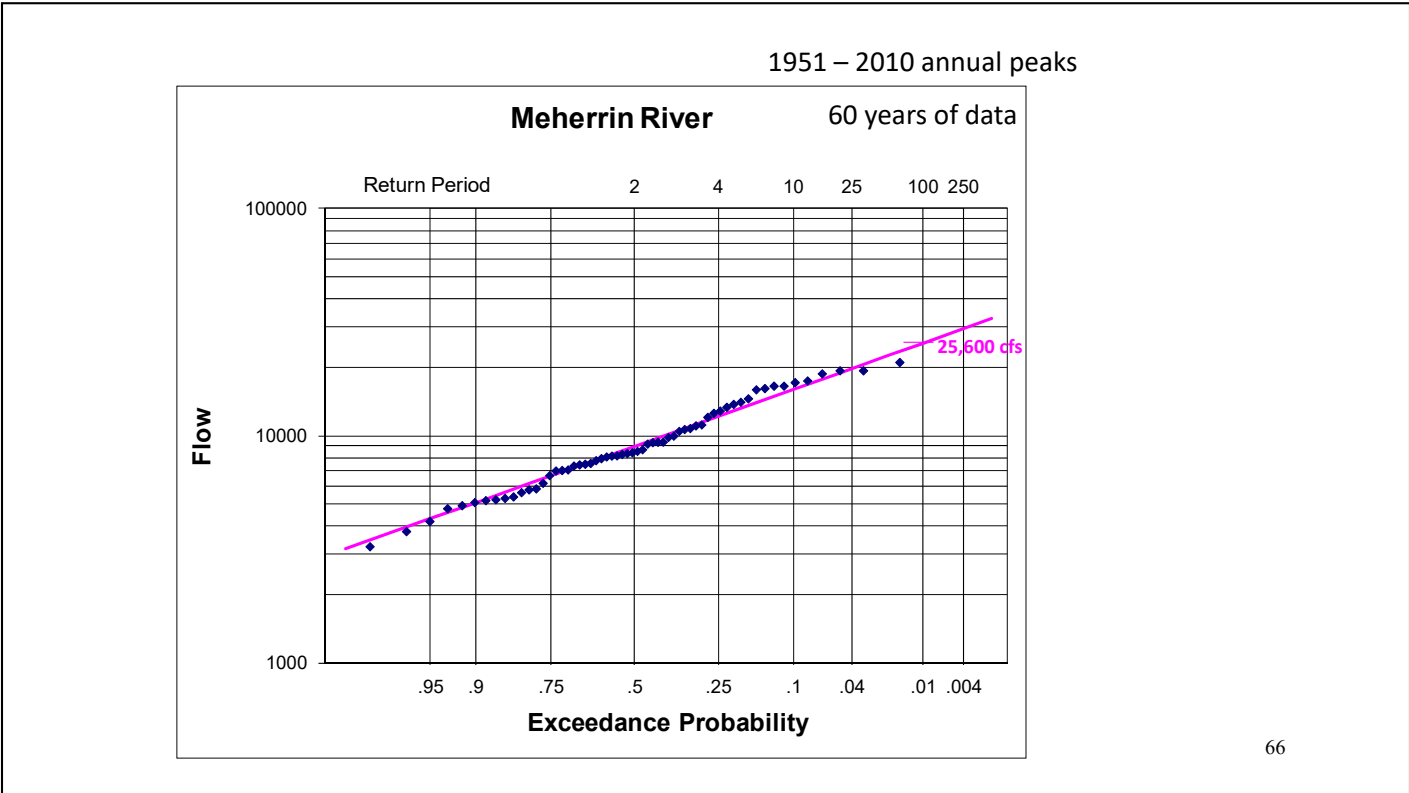


64

Here is another example – the Meherrin River in Virginia



The 60 years of systematic data plot like this.



This curve is the LP3 fit of the systematic data.

520. Meherrin River at Emporia, Va.

Location--Lat 36°41'20", long 77°32'20", on left bank at downstream side of bridge on U. S. Highway 301, in Emporia, Greensville County.

Drainage area--749 sq mi.

Gage--Recording. Altitude of gage is 68 ft (by barometer).

Stage-discharge relation--Defined by current-meter measurements below 11,000 cfs and extended above by logarithmic plotting on basis of record for station near Lawrenceville.

Bankfull stage--13 ft.

Historical data--Flood of Aug. 17, 1940, was greatest since at least 1873.

Remarks--Subsequent to July 1, 1957, records furnished by Virginia Department of Conservation and Economic Development, Division of Water Resources. Information for floods prior to 1929 derived from data reported in Congressional documents: 71st Cong., 2d sess., H. Doc. 446, Meherrin River (1930). Base for partial-duration series, 6,000 cfs.

Peak stages and discharges

Water year	Date	Gage height (feet)	Discharge (cfs)	Water year	Date	Gage height (feet)	Discharge (cfs)
1873	Feb. 10, 1873	(a)	-	1953	Nov. 23, 1952	21.90	11,200
1888	Sept. 13, 1888	-	-		Jan. 26, 1953	19.18	7,640
1889	June 2, 1889	(b)	-	1954	May 21, 1954	17.63	5,860
1893	May 6 or 7, 1893	-	-	1955	Aug. 21, 1955	22.80	12,600
1908	Aug. 28, 1908	28	-	1956	Oct. 3, 1955	19.07	7,520
1912	March 1912	25	-		Oct. 16, 1955	18.82	7,180
1919	July 25, 1919	-	-		Feb. 8, 1956	17.86	6,190
1928	Apr. 27, 1928	26	-		Mar. 18, 1956	18.07	6,410
1940	Aug. 17, 1940	30.0	40,000	1957	July 22, 1956	17.87	6,190
1951	Mar. 21, 1951	18.50	5,100		Feb. 3, 1957	19.78	7,580
1952	Dec. 23, 1951	20.60	9,410		Feb. 28, 1957	18.77	6,500
	Jan. 11, 1952	18.68	7,070	1958	Dec. 11, 1957	19.02	6,700
	Jan. 30, 1952	20.32	8,990		Dec. 22, 1957	18.52	6,300
	Mar. 5, 1952	18.31	6,630		Jan. 27, 1958	18.37	6,100
	Mar. 26, 1952	18.60	6,960		Mar. 1, 1958	19.02	6,700
	Apr. 27, 1952	20.30	8,990		Apr. 1, 1958	18.90	6,630
					May 8, 1958	22.76	12,100
				1959	Dec. 31, 1958	21.18	9,400

a At least 4 ft lower than flood of 1889.
b Slightly lower than flood of 1908 at station "near Lawrenceville."

Meherrin River at Emporia, VA

many additional forms of information available from USGS

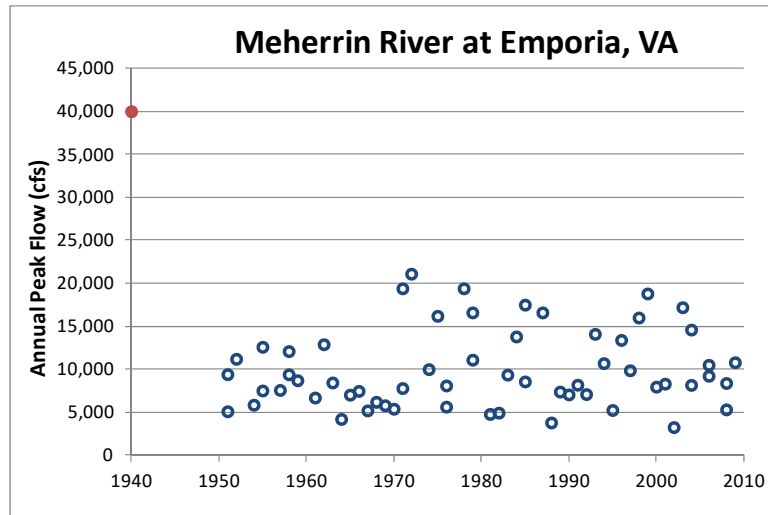
67

But, there was a lot more information available from the USGS than just the systematic starting in 1951. There was a flood of 40,000 cfs in 1840 that was estimated to be the largest since 1873. And there were also stage estimates for several of the years between 1873 and 1940.

This is good information that we should bring into the frequency analysis, if possible.

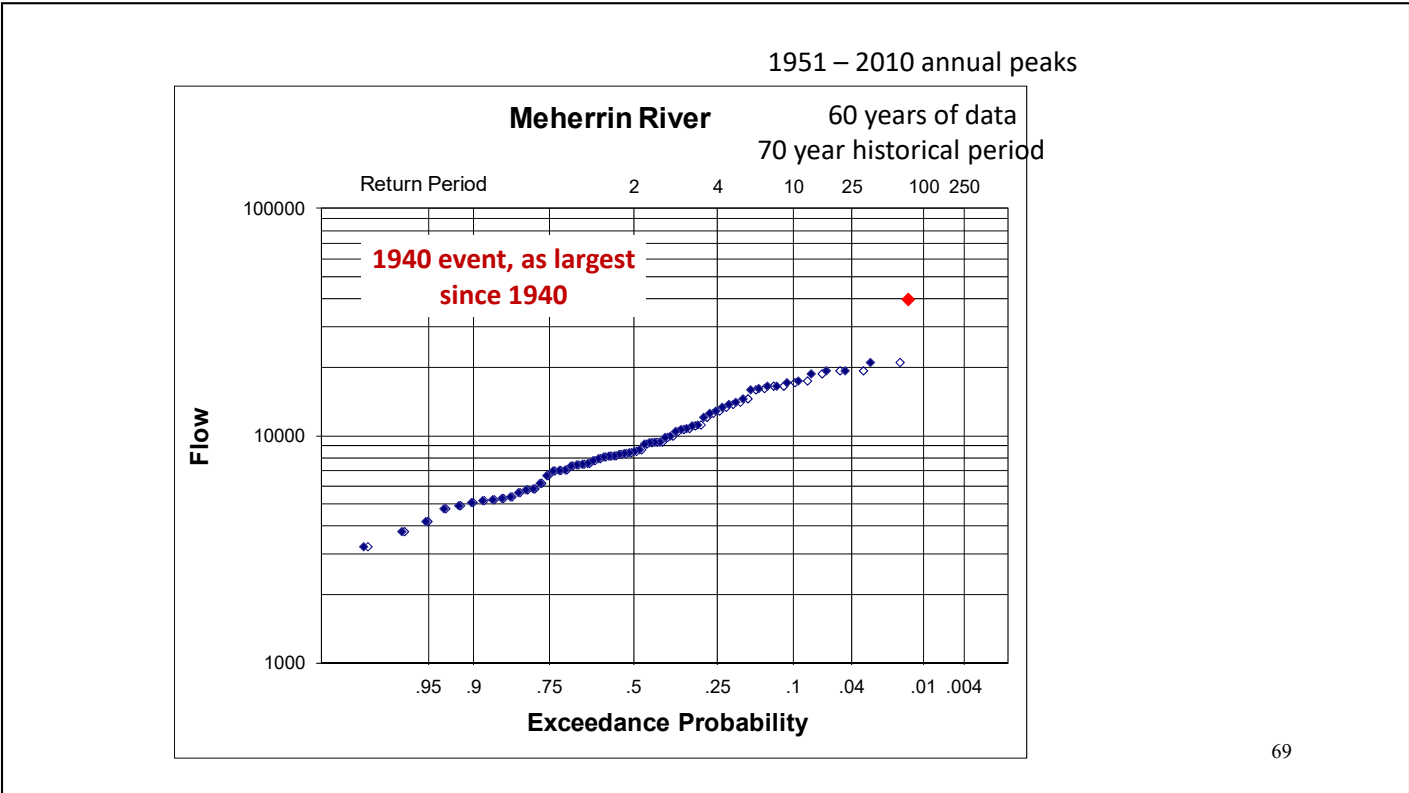
Here, we'll work with this information one piece at a time to see how it affects the frequency curve.

Example: historical information



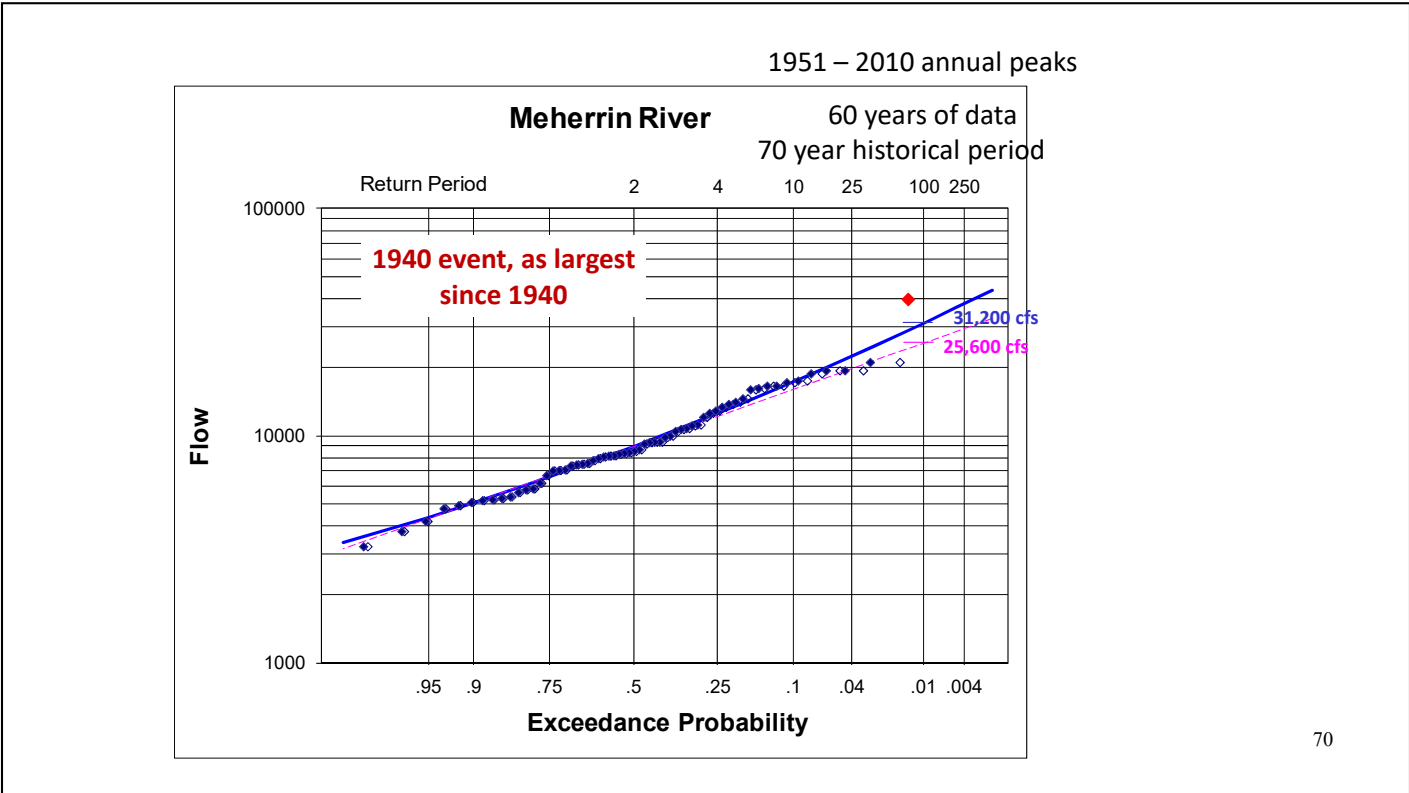
68

Here, we've included the 1940 event estimated at 40,000 cfs.



69

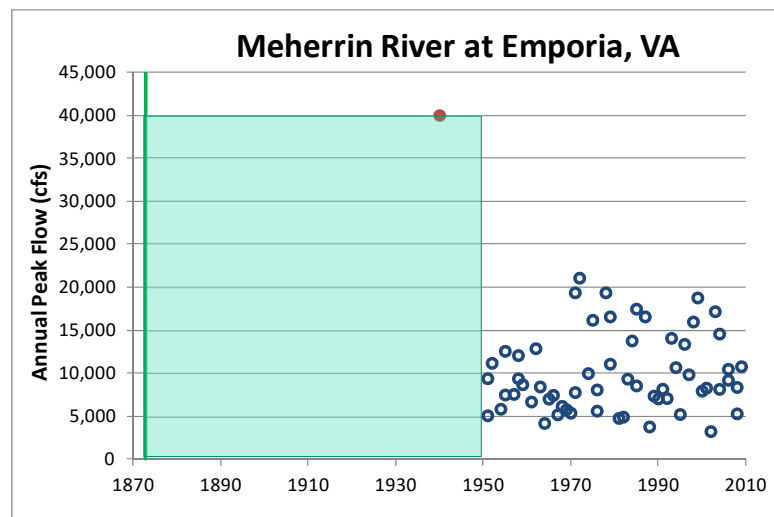
When add the 1940 event, can see from the plotting positions that, at this point, we're considering it the largest since only 1940 – the return period is about 1 in 70. We expect a higher frequency curve will result from this data.



70

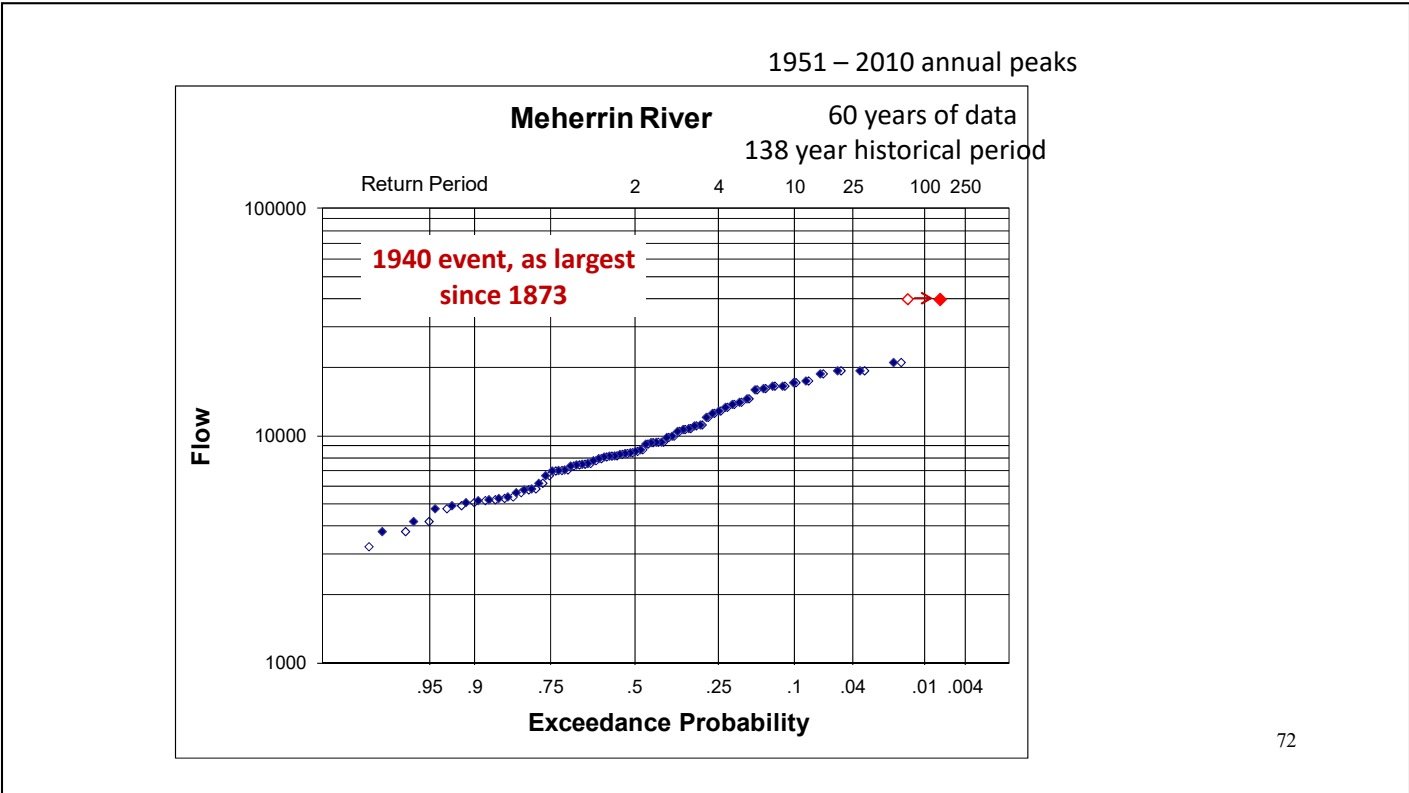
This is the higher LP3 fit including 1940. Note, the plotting position didn't affect the LP3 fit, but it can be seen to be consistent.

Example: historical information

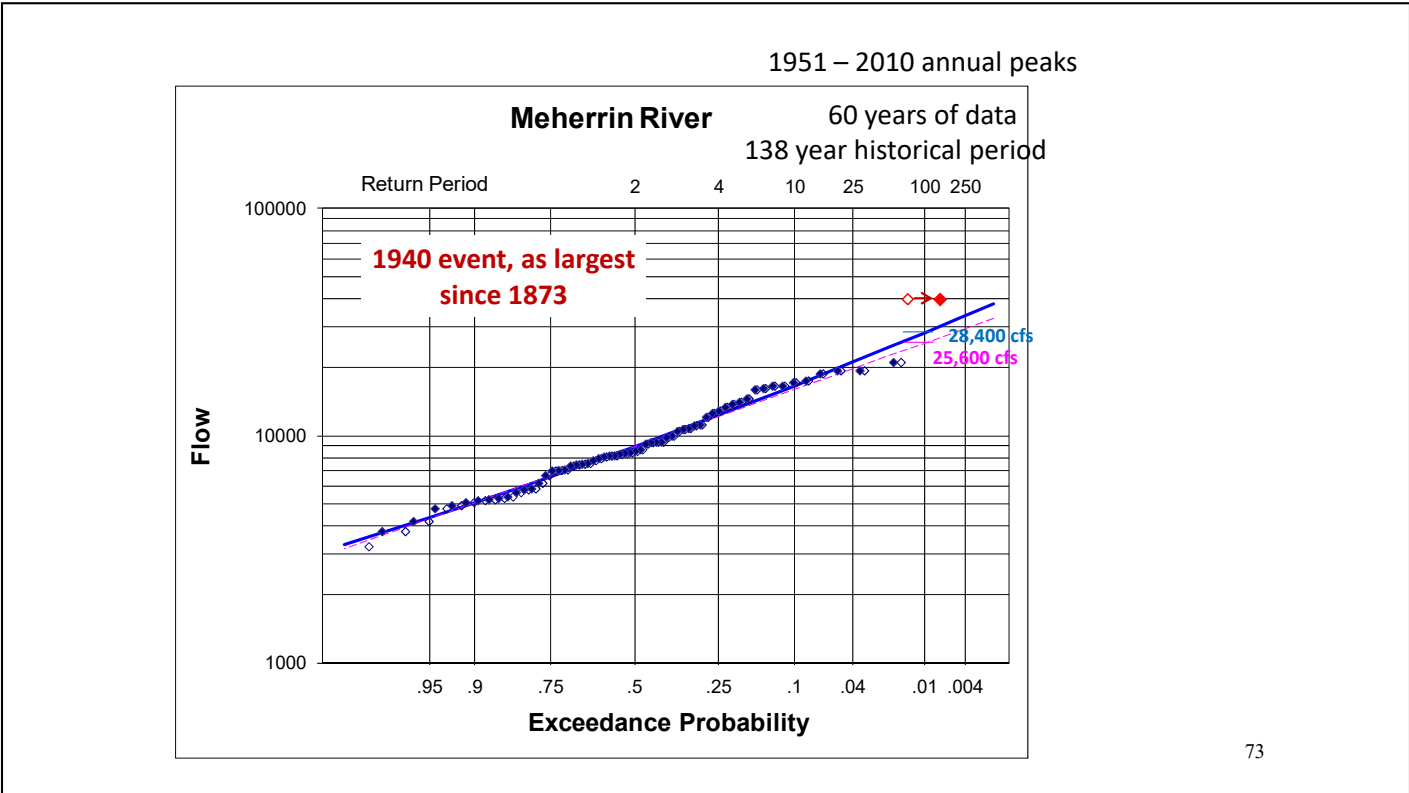


71

Here we include the fact that the 1940 event was the largest since 1873.



We can see the plotting position of the largest event now shows it as the largest since 1873, with a return period of 1 in 138. We expect this information will lower the frequency curve from the previous.



73

The dark blue is the LP3 fit with 1940 as largest since 1873. It's higher than the pink it with no historical data, but lower than the estimate that did not add the "largest since 1873" information.

520. Meherrin River at Emporia, Va.

Location--Lat 36°41'20", long 77°32'20", on left bank at downstream side of bridge on U. S. Highway 301, in Emporia, Greensville County.

Drainage area--749 sq mi.

Gage--Recording. Altitude of gage is 68 ft (by barometer).

Stage-discharge relation--Defined by current-meter measurements below 11,000 cfs and extended above by logarithmic plotting on basis of record for station near Lawrenceville.

Bankfull stage--13 ft.

Historical data--Flood of Aug. 17, 1940, was greatest since at least 1873.

Remarks--Subsequent to July 1, 1957, records furnished by Virginia Department of Conservation and Economic Development, Division of Water Resources. Information for floods prior to 1929 derived from data reported in Congressional documents: 71st Cong., 2d sess., H. Doc. 446, Meherrin River (1930). Base for partial-duration series, 6,000 cfs.

Peak stages and discharges

Water year	Date	Gage height (feet)	Discharge (cfs)	Water year	Date	Gage height (feet)	Discharge (cfs)
1873	Feb. 10, 1873	(a)	-	1953	Nov. 23, 1952	21.90	11,200
1888	Sept. 13, 1888	-	-		Jan. 26, 1953	19.18	7,640
1889	June 2, 1889	(b)	-	1954	May 21, 1954	17.63	5,860
1893	May 6 or 7, 1893	-	-	1955	Aug. 21, 1955	22.80	12,600
1908	Aug. 28, 1908	28	-	1956	Oct. 3, 1955	19.07	7,520
1912	March 1912	25	-		Oct. 16, 1955	18.82	7,180
1919	July 25, 1919	-	-		Feb. 8, 1956	17.86	6,190
1928	Apr. 27, 1928	26	-		Mar. 18, 1956	18.07	6,410
1940	Aug. 17, 1940	30.0	40,000		July 22, 1956	17.87	6,190
1951	Mar. 21, 1951	16.90	5,100	1957	Feb. 3, 1957	19.78	7,580
1952	Dec. 23, 1951	20.60	9,410		Feb. 28, 1957	18.77	6,500
	Jan. 11, 1952	18.68	7,070	1958	Dec. 11, 1957	19.02	6,700
	Jan. 30, 1952	20.32	8,990		Dec. 22, 1957	18.52	6,300
	Mar. 5, 1952	18.31	6,630		Jan. 27, 1958	18.37	6,100
	Mar. 26, 1952	18.60	6,960		Mar. 1, 1958	19.02	6,700
	Apr. 27, 1952	20.30	8,990		Apr. 1, 1958	18.90	6,630
					May 8, 1958	22.76	12,100
				1959	Dec. 31, 1958	21.18	9,400

a At least 4 ft lower than flood of 1889.
b Slightly lower than flood of 1908 at station "near Lawrenceville."

determine
flow/stage
rating from
1951 – 1958

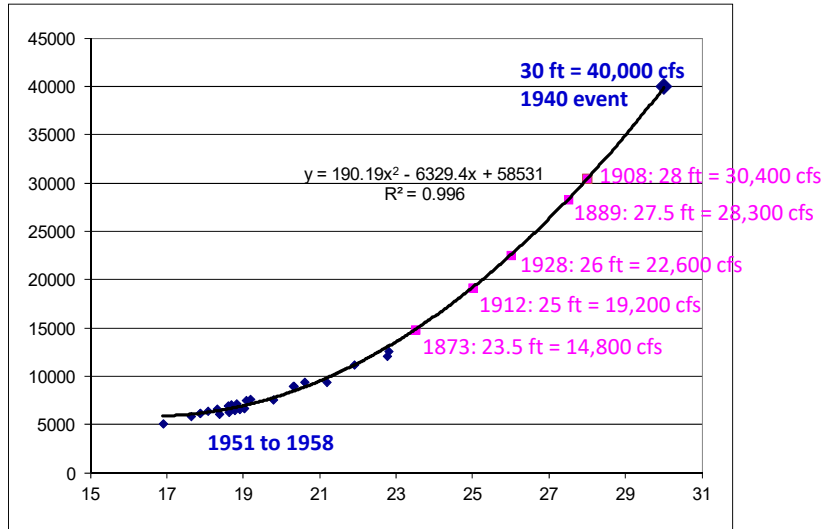
estimate
historical flows
from stage

NOTE: in B17B,
lowest
historical flow
becomes new
high outlier
threshold...

74

Now we can try to use this additional information about peak stages.

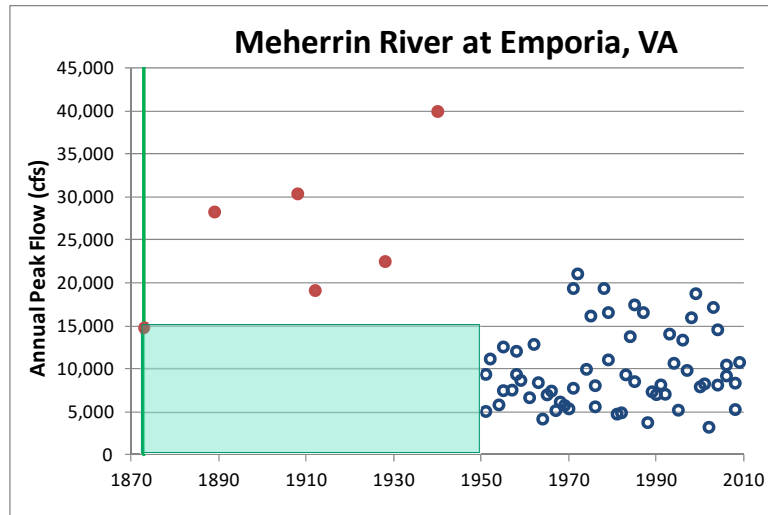
Estimate flow from stage for historical events



75

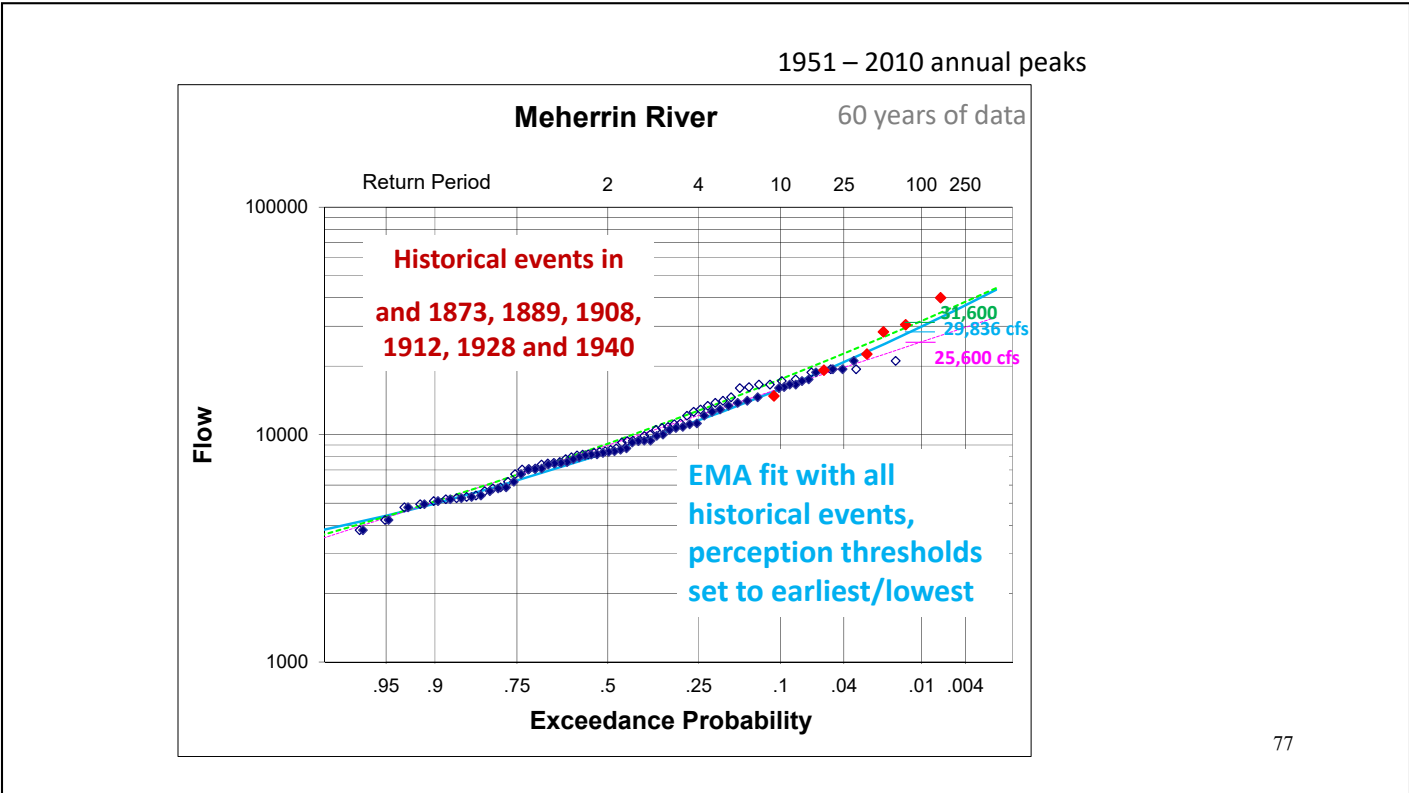
Here's a simple estimate of the historical flows made by plotting 1951 to 1959 flow versus stage along with the 1940 estimate of flow of 40,000 cfs from stage of 30 feet. The polynomial fit is used to estimate flow for the recorded stages.

Example: historical information



76

This image shows the additional historical events in the historical period. We can safely assume that since the flow of 15,000 cfs was observed in 1873, that flow would have been observed at any time since then, and so it can be used as a perception threshold in a B17C EMA fit.



EMA is able to do a better job of using all of the historical information without needing to specify more of the historical events or systematic flows as high outliers.

Assumptions?

- What assumptions are we making by including historical or paleoflood information?
 - **stationarity** – those events are part of the same population as the recent events, i.e., the same physical processes are in place: hydrology, hydraulics, etc
 - **consistency** – that the estimates of the historic flow are equivalent to systematic measurements
- Consider whether these are safe assumptions

78

By using the historical data in the frequency analysis, we are saying we believe it to be from the same flood population as the gaged flows, and thus identically-distributed. Other ways to word this assumption are that the data is homogeneous, and since we collect the data over time, that it is stationary. Using the historical data means we are saying the watershed was similar enough to the current watershed to produce the same distribution of flood events (with the same likelihoods).

We are also using the historical data here as point values, implying the estimates are equally as reliable. However, we could choose to represent them as intervals instead.

Are these assumptions valid?

Summary of how we use historical info

Bulletin 17B method

- can use historical events that have estimated values
- Can use non-exc info only if observed large event
- extend record to length of historical, *assuming syst record also represents unobs years*
- works well for historical length, too weak for paleo length

Bulletin 17C - EMA

- can use a point or an interval for an historical event
- can use intervals for all of the unobserved events, *even if no historical event* (non-exc. info)
- can also use other info:
 - the fact that a flow is below any threshold
 - the fact that a flow is above any threshold

79

Summary of all Data Adjustments

Bulletin 17B

- **Broken record** – missing years excluded
- **Zero flow years** – conditional prob adj
- High and Low outliers:
 - **Low outliers** – conditional prob adj
 - **High outliers** – seek historical info
- **Historical Information** – use Weighed Moment algorithm
 - Historical events
 - Periods of known non-exceedance, if hi-out
 - Note whether systematic events become high outliers

Bulletin 17C (EMA)

- **Broken record** – missing years represented by a range, (∞, ∞) or (T, ∞) if some info
- Low values
 - **Zero flow years**: censored, replaced by range
 - **Low outliers and PILFs**, replaced by range
- **Historical Information**
 - Historical events
 - Periods of known non-exceedance, or exceedance
 - Set thresholds for historical periods