

# Why projects go over time?

---

DAVID GUSZEJNOV

# Underestimating project length

---

We have all ran out of time before deadlines or had projects „blow up”

It is common in both business and academia

Somehow we are really bad at estimating the length of projects



# Are we bad at statistics?

---

Hypothesis from the data science blog of Erik Bernhardsson:

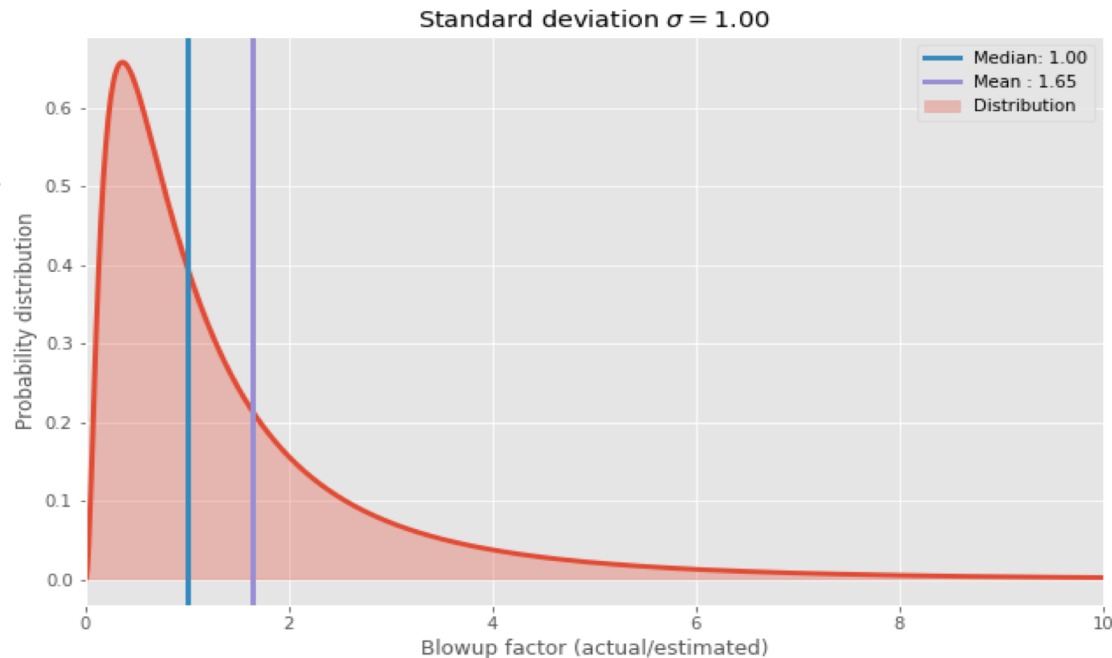
**Humans can only estimate the median, not the mean**

Let's define the **blowup factor** as

$$f = \frac{\Delta t_{actual}}{\Delta t_{estimate}}$$

Let's assume  $f$  is **lognormal**:

$$\log f \sim N(0, \sigma)$$



# How multiple tasks add-up

---

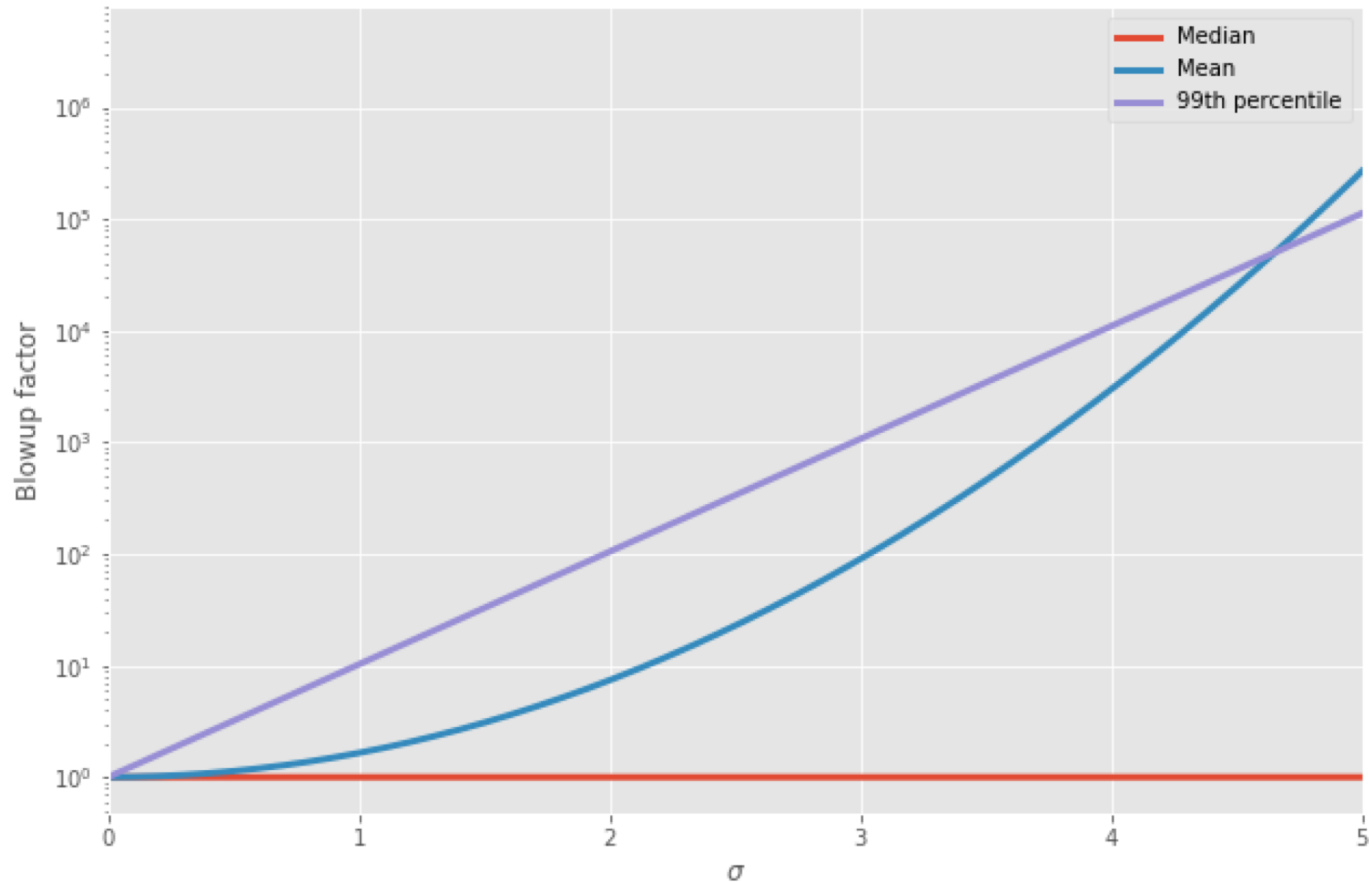
	Median	Mean	99%
Task A	1.00	1.65	10.24
Task B	1.00	1.65	10.24
Task C	1.00	1.65	10.24
SUM	3.98	4.95	18.85

**Even worse if projects are allowed to have different variances for the log f**

	Median	Mean	99%
Task A ( $\sigma = 0.5$ )	1.00	1.13	3.20
Task B ( $\sigma = 1$ )	1.00	1.65	10.24
Task C ( $\sigma = 2$ )	1.00	<b>7.39</b>	<b>104.87</b>
SUM	<b>4.00</b>	<b>10.18</b>	<b>107.99</b>

# Mean vs Median for high $\sigma$

---



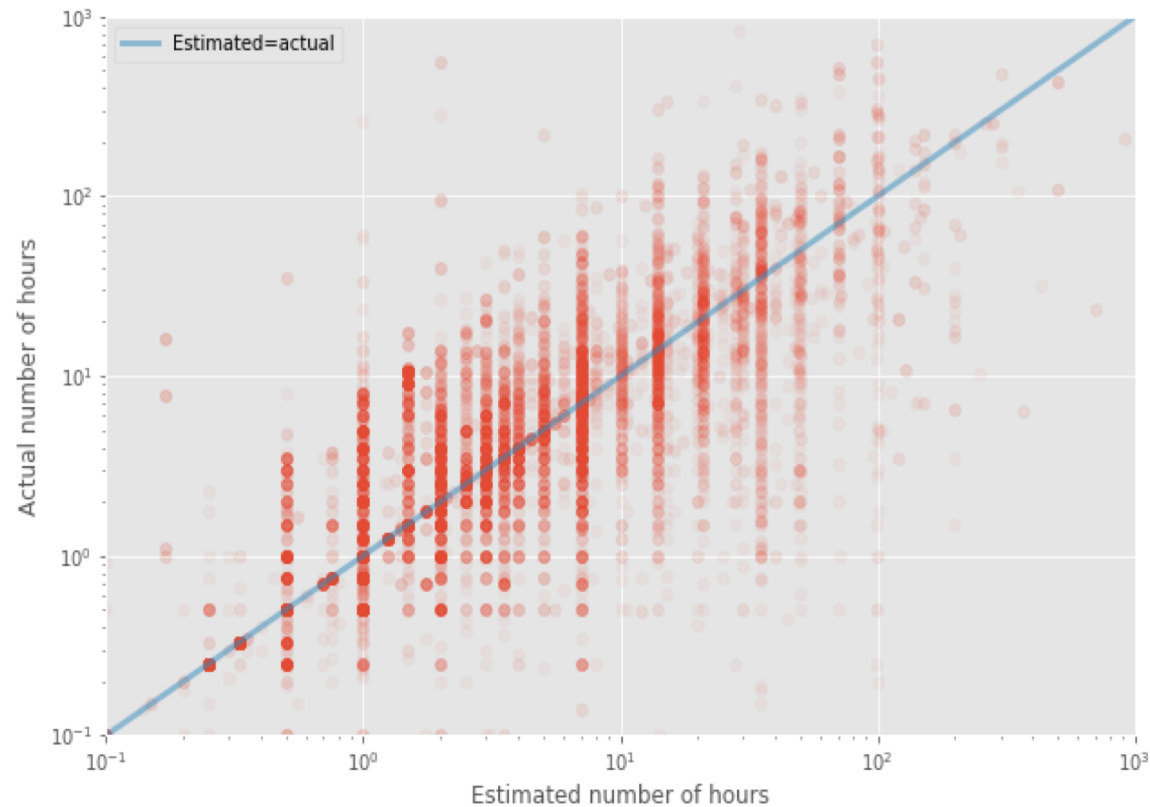
# Real data on blowup factor

So far this was only a hypothesis

Actual data for blowup factor in software development projects:

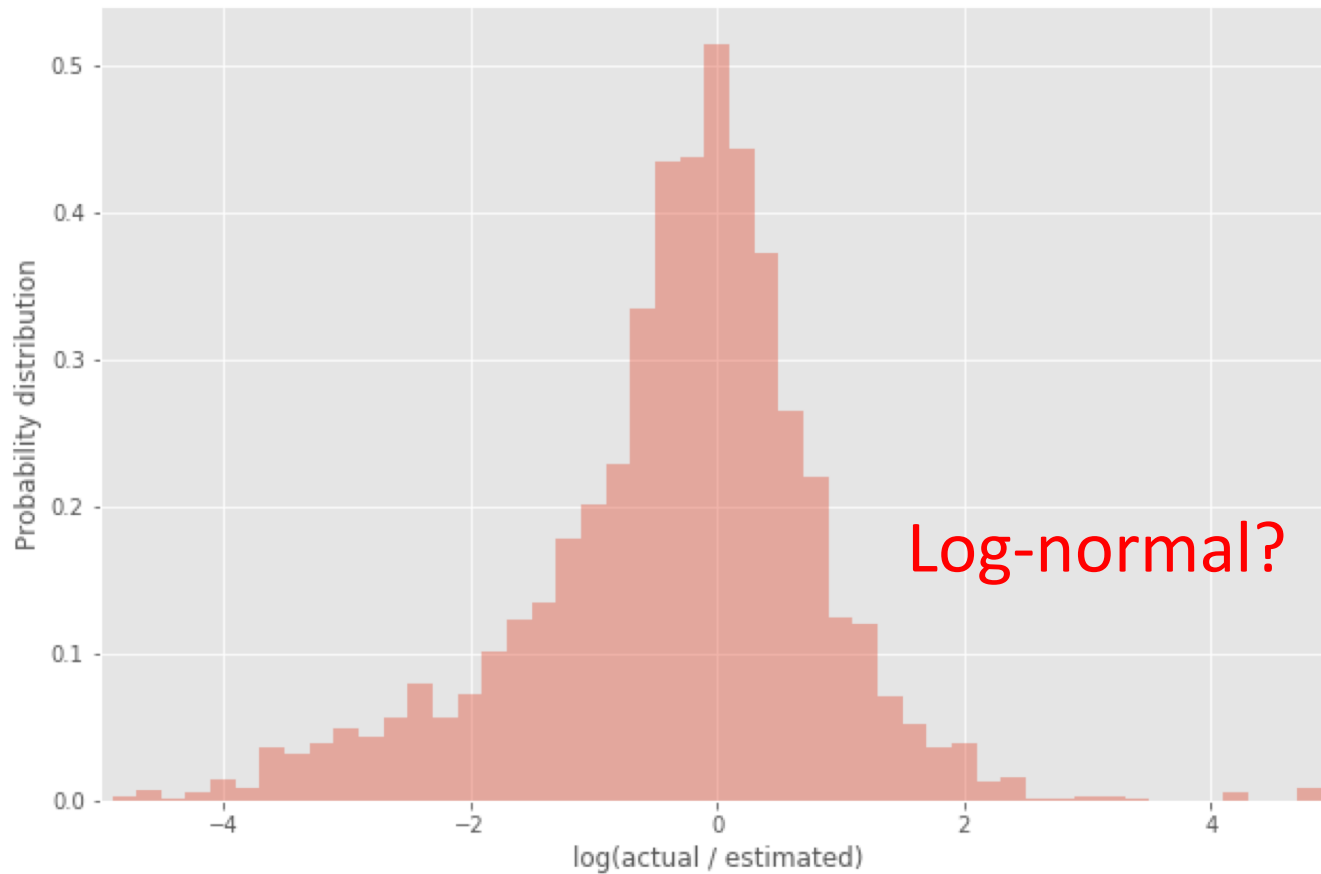
- median = 1
- mean = 1.8

Maybe our hypothesis is not that bad?

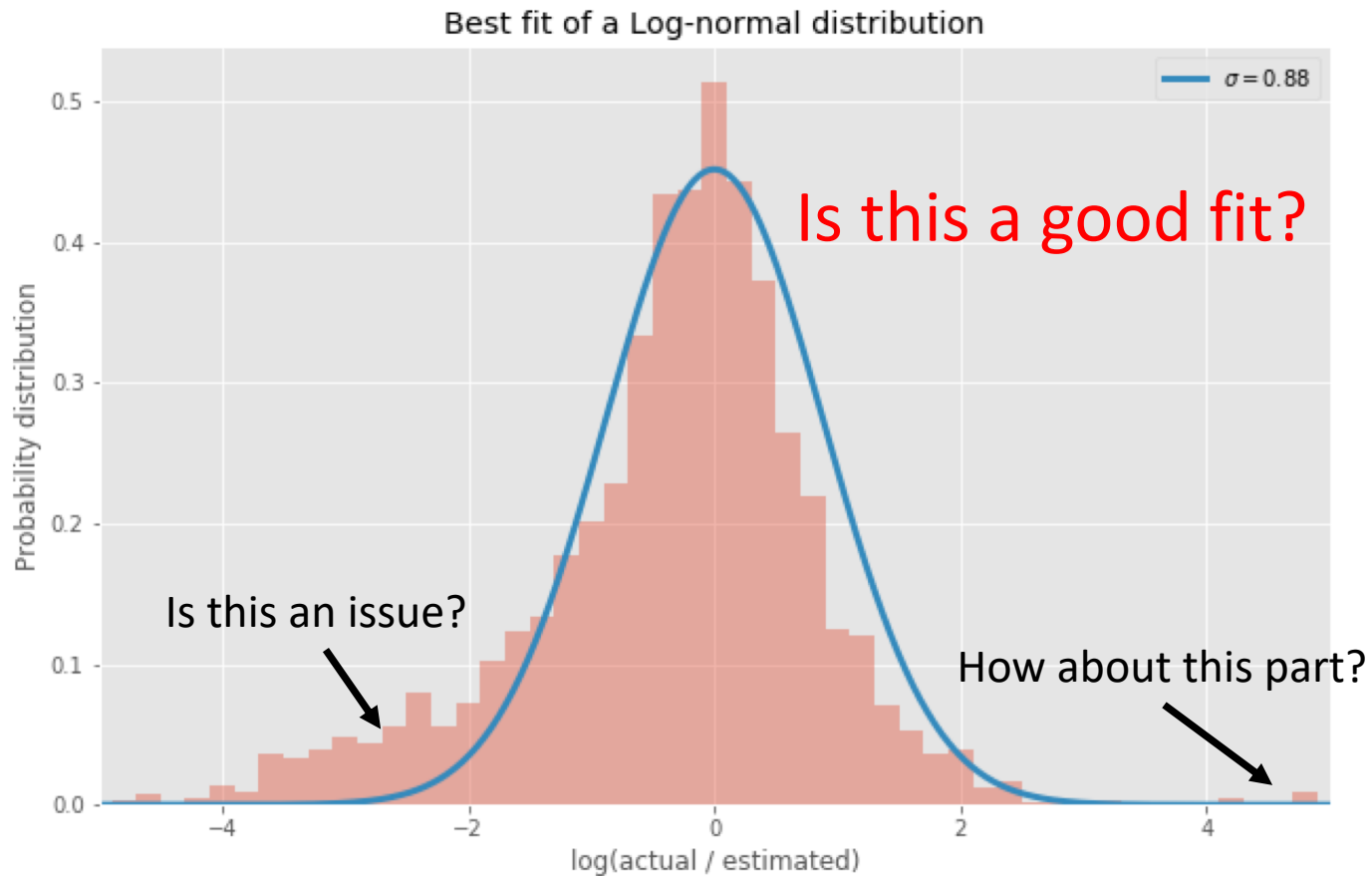


# Distribution of log blowup factor

---

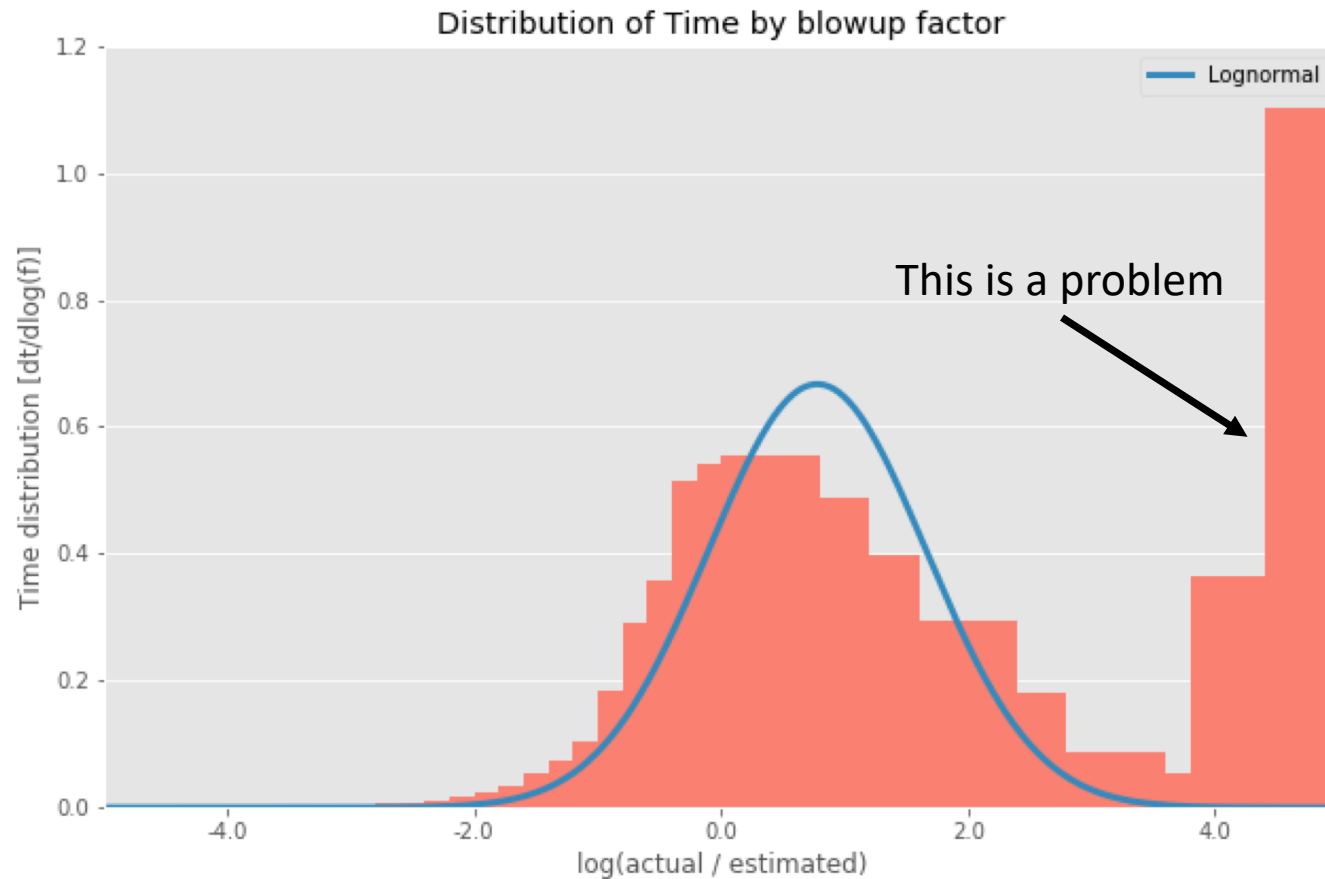


# Is log-normal a good fit?





# Is log-normal a good fit?



# Fitting an analytic PDF

---

Let's fit an analytic PDF to the data!

**But which one?**

**Hypothesis:** Each project has an  $f$  blowup value that is lognormal with unit median and  $\sigma$  variance

**Data:** Mixture of different projects, with some distribution over  $\sigma$

**Data distribution** = compound distribution of  $f_\sigma$  and  $\sigma$

$$\text{PDF of data } p(f) = \int \text{PDF of } \sigma \text{ among projects } p_\sigma(\sigma) p(f|\sigma) d\sigma$$

Blowup factor for each project (lognormal)

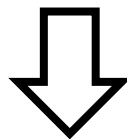
# Fitting an analytic PDF

We have the freedom to choose  $p_\sigma(\sigma)$

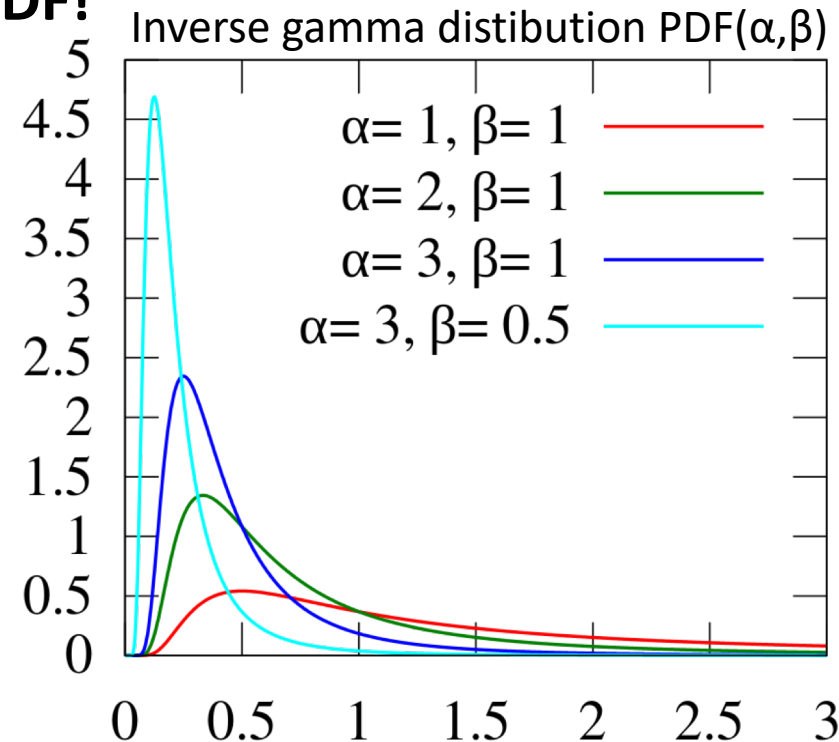
Let's choose the **inverse gamma PDF!**

Advantages:

- Positive values
- Flexible shape (2 parameters)
- **Gives an analytical PDF if compounded with a Gaussian**



**Generalized Student t-distribution**

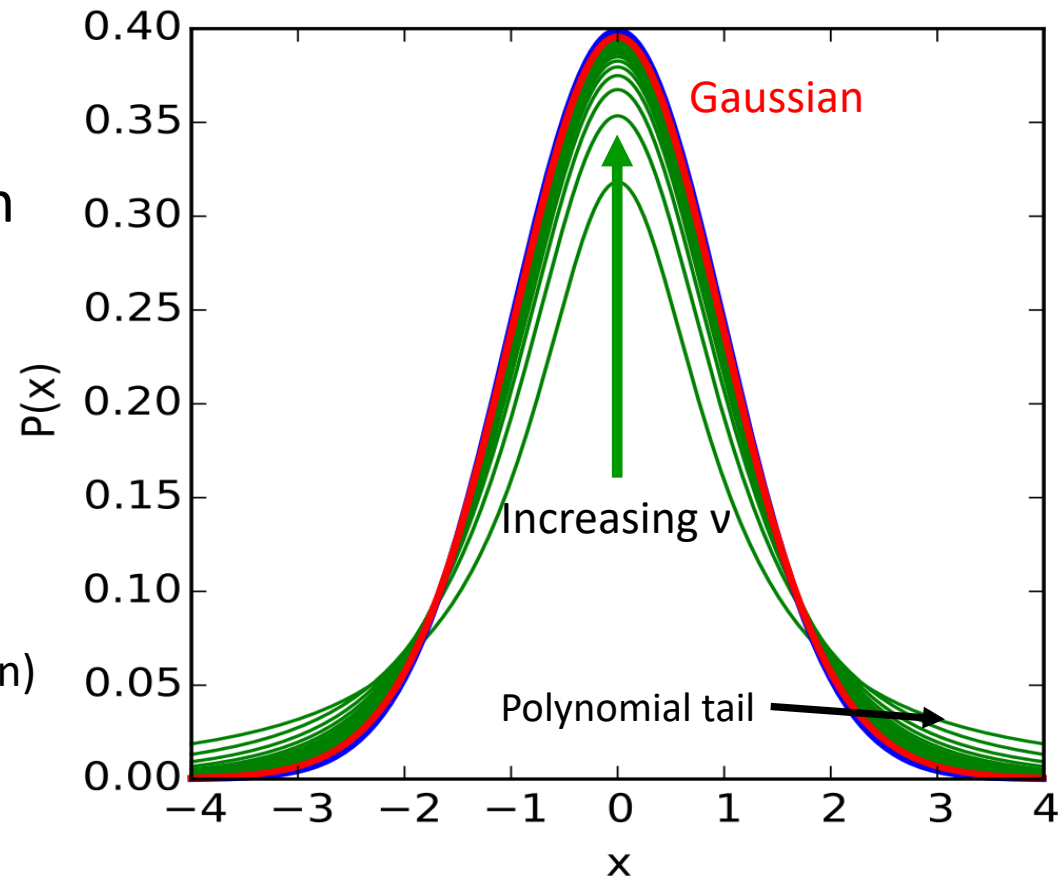


# Generalized Student-t distribution

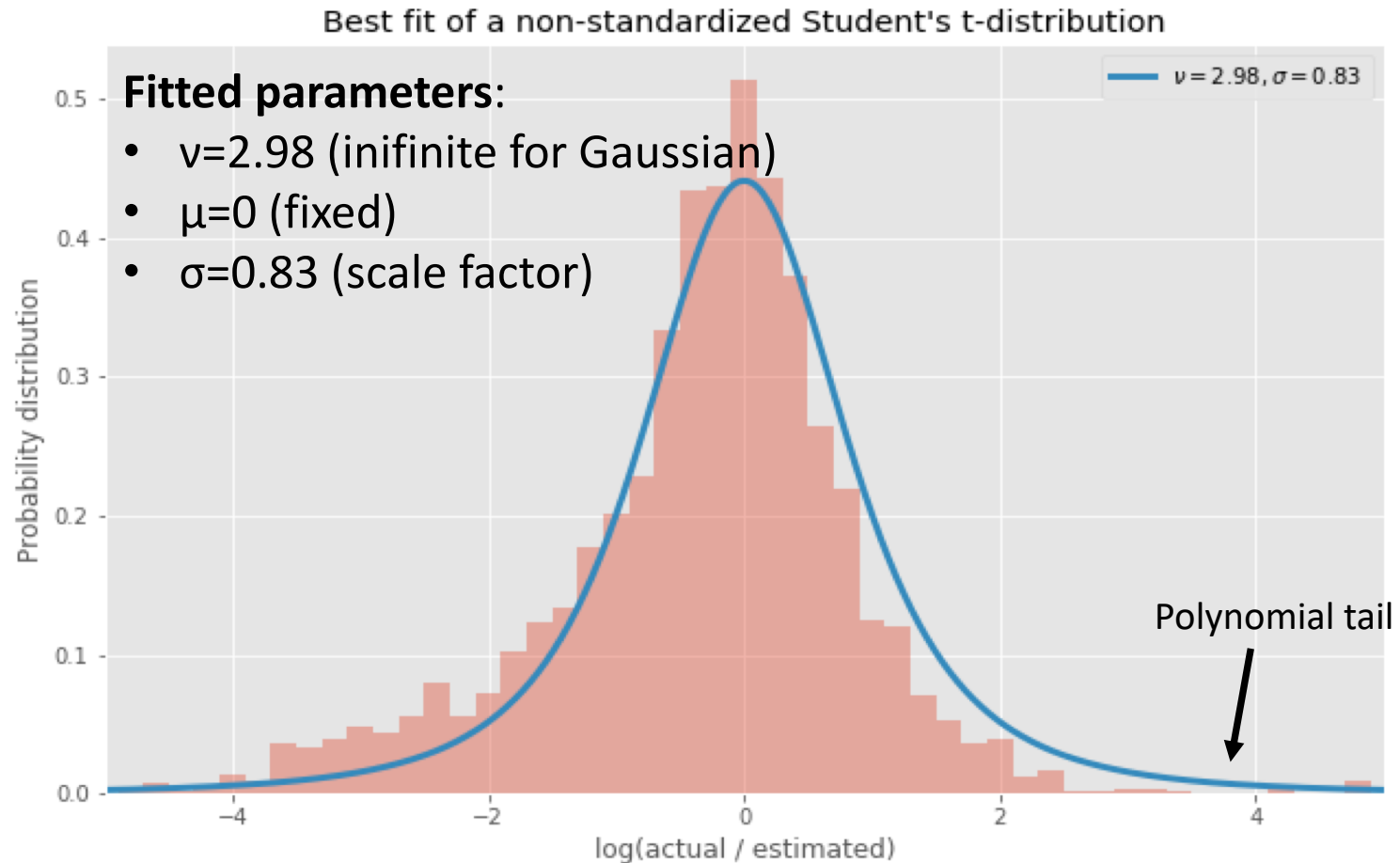
Arises when calculating the mean of a finite samples from a Gaussian distribution

Parameters:

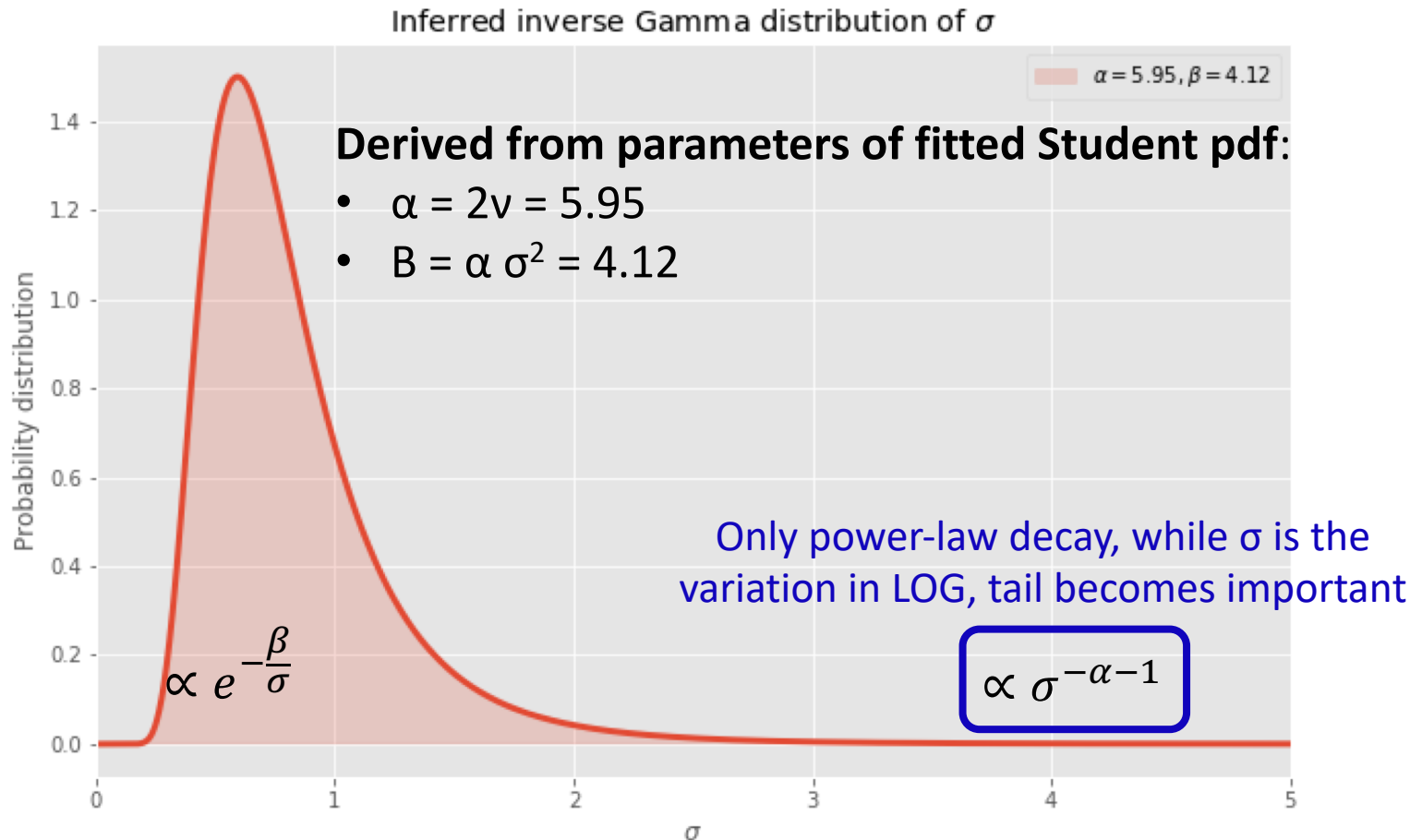
- $\nu$  = degrees of freedom
- $\mu$  = location parameter
- $\sigma$  = scale factor (not variation)



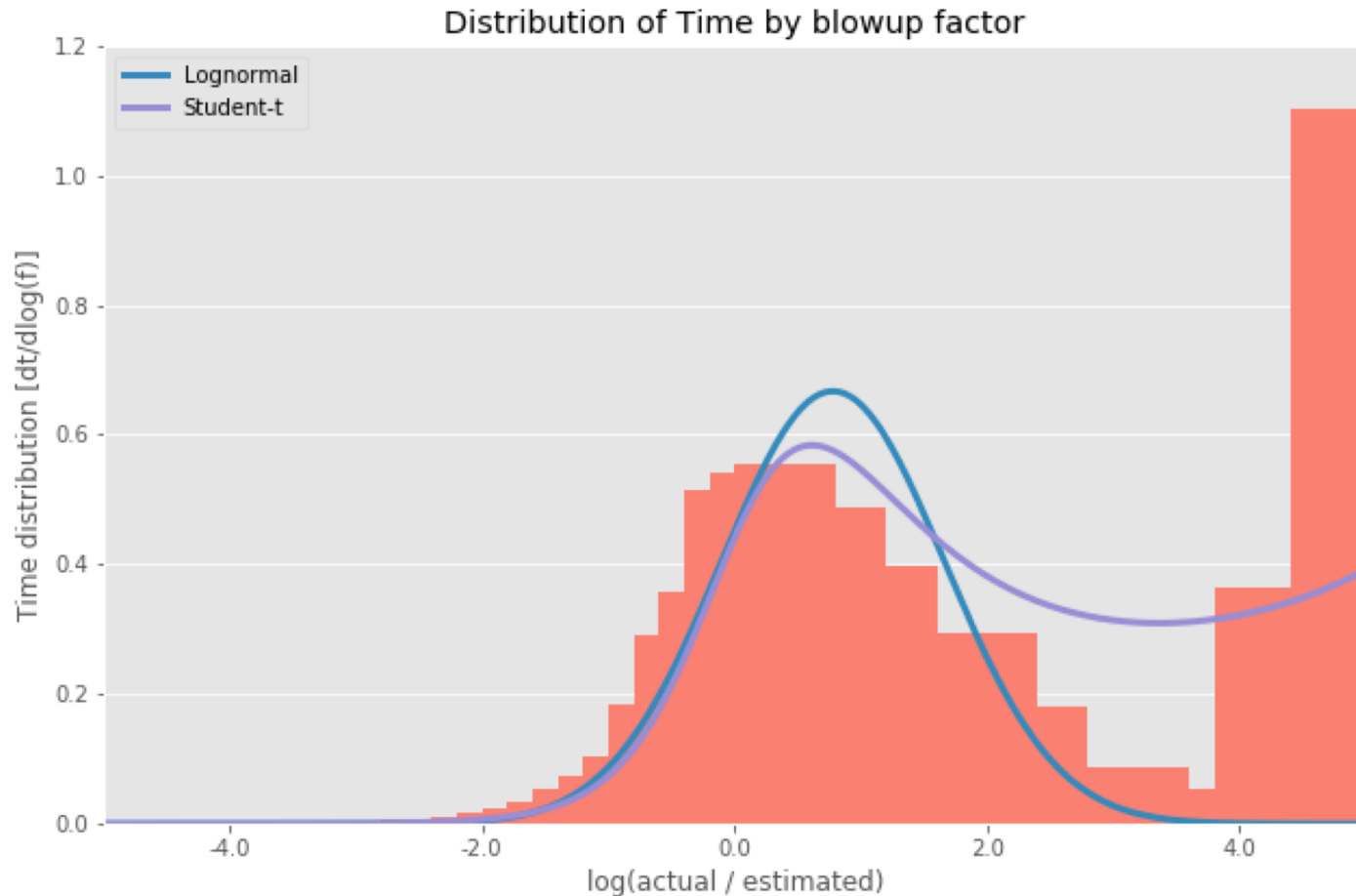
# Fitting an analytic PDF



# Distribution of $\sigma$ for data



# Distribution of project lengths

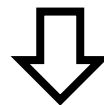


# Derived quantities from fit to data

---

<b>Median value</b>	<b>1</b>
90 <sup>th</sup> percentile	3.92
99 <sup>th</sup> percentile	44
<b>Mean value</b>	<b>Infinite!</b>

**Polynomially rare events require exponentially longer time, which leads to an infinite expected value**



The mean time to solve a completely unknown project is infinite.  
No wonder we can't keep the deadlines...



# Summary

---

## Robust results from data

- People actually estimate the **median** time a task requires, not the mean
- For skewed distributions the **medians don't add up**, unlike the means, leading to people underestimating the total time needed
- High blowup-factor projects dominate the project time length distribution

## Modeling results

- Hypothesis: Blowup factor of projects is well approximated by a lognormal distribution
- Blowup-factor well characterized by a Student-t distribution
- Probability of high  $f$  blowup-factor events decay as polynomial of  $\log f$ , so these dominate the mean

**Best practice: to estimate a project's length always look for the part that can potentially take the longest**