# Network Traffic Measurements and Analysis

## Lecture III: Probability models and measurements

Markku Liinaharja (slides originally made by Esa Hyytiä),

Department of Communications and Networking
Aalto University, School of Electrical Engineering

Version 0.2, October 4, 2017

# Contents

- Introduction
- Probability
- Distributions for network measurements
- Parameter estimation
- Model validation

Aalto University
School of Electrical
Engineering

# Measurements and modeling

- Exploratory approach
  Problem ⇒ Data ⇒ Analysis ⇒ (Model) ⇒ Conclusions
- Measurement analysis is often intertwined with traffic modeling
- If the observations can be described using an idealized mathematical model their implications are often easier to understand
- E.g., input to a queue
- "All models are wrong, but some models are useful"

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
3/51

# Why models?

- **Descriptive models** for measurements
  - Efficient summary of observed data
  - E.g., Gaussian with mean 100ms and standard deviation 10ms
- **Constructive models** for what-if scenarios
  - Model that could have produced the observed data
  - E.g., The trace could have been produced by a certain stochastic process
  - We are interested in the underlying phenomena instead of details of data
- In this lecture, we focus on descriptive probability models

**Aalto University**
School of Electrical
Engineering

**ELEC-E7130 - Internet Traffic Measurements and Analysis**
October 4, 2017
4/51

# Probability models

- Probability captures features that are unknown or difficult to characterize
  - Exact user behavior
  - Immense amount of functionalities in the Internet
- Probability allows us to model, reason, and proceed with inference in an uncertain environment

**Aalto University**
School of Electrical
Engineering

**ELEC-E7130 - Internet Traffic Measurements and Analysis**
October 4, 2017
5/51

# Modeling process

1. Model selection
   - ▶ Prefer models with a few parameters over those that have more parameters (Occam's razor)
   - ▶ Model should be parsimonious to avoid over-fitting
2. Parameter estimation
   - ▶ Choose the parameters that best describe the observed data
3. Validation
   - ▶ Descriptive: Compare the distributions
   - ▶ Constructive: Confirm that the observed data is relatively likely outcome of the model

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
6/51

# Contents

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
7/51

# Probability

- Consider an experiment with random outcomes
- Event is a set of outcomes
- **Probability**
  - How likely an event is?
  - Long-run proportion of an event in a series of experiments
- Mathematically defined by three axioms, for an event A
  - i) $0 \leq P\{A\} \leq 1$
  - ii) $P\{S\} = 1$, where $S$ is the sample space
  - iii) $A$ and $B$ disjoint $\Rightarrow P\{A \text{ or } B\} = P\{A\} + P\{B\}$
- Everything else follows from these!

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
8/51

# Useful notions

- **Conditional probability**:

$$P\{B \mid A\} = \frac{P\{A \text{ and } B\}}{P\{A\}} = \frac{P\{A \cap B\}}{P\{A\}}.$$

- **Independence**: If the two events are independent, then

$$P\{B \mid A\} = P\{B\}.$$

- **Bayes' rule**: assume $S = B_1 \cup B_2 \cup \ldots \cup B_n$, where $B_i \cap B_j = \emptyset$ for $i \neq j$. Then,

$$P\{B_i \mid A\} = \frac{P\{A \mid B_i\} \cdot P\{B_i\}}{\sum_j P\{A \mid B_j\} \cdot P\{B_j\}}.$$

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
9/51

# Random variables

- Random variable $X$ is a variable whose value is a numerical outcome of a random phenomenon
  - Coin toss: "heads" $X = 1$, "tails" $X = 0$

- **Discrete random variable** $X$ takes discrete values $\{x_1, x_2, x_3, \ldots\}$ with probabilities $\{p_1, p_2, p_3, \ldots\}$

- **Continuous random variable** $X$ takes continuous values $\{x\}$ according to a probability density function $f(x)$ (pdf)
  - Probability of event $x \in (a, b)$ is the area under pdf,

$$P\{a < X < b\} = \int_a^b f(x)\,dx$$

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
10/51

# Expectation and variance of a random variable

▶ **Expectation** (mean of a random variable) $\mu = \mathsf{E}(X)$

$$\mu = \sum_i p_i x_i, \qquad \mu = \int x \cdot f(x)\, dx.$$

▶ **Variance** $\mathsf{V}(X) = \mathsf{E}((X - \mu)^2)$,

$$\mathsf{V}(X) = \sum_i p_i (x_i - \mu)^2, \qquad \mathsf{V}(X) = \int (x - \mu)^2 \cdot f(x)\, dx.$$

▶ **Covariance** of two random variables

$$\mathrm{Cov}[X, Y] = \mathsf{E}((X - \mu_x)(Y - \mu_y)).$$

**Aalto University**
School of Electrical
Engineering

**ELEC-E7130 - Internet Traffic Measurements and Analysis**
October 4, 2017
11/51

# Properties of expectation and variance

- For constants $a, b$, and random variables $X$ and $Y$

$$\mathsf{E}(aX + bY) = a\,\mathsf{E}(X) + b\,\mathsf{E}(Y).$$

- For constants $a, b$, and a random variable $X$

$$\mathsf{V}(aX + b) = a^2\,\mathsf{V}(X).$$

- For *independent* random variables $X$ and $Y$

$$\mathsf{V}(X + Y) = \mathsf{V}(X) + \mathsf{V}(Y).$$

$$\mathsf{E}(XY) = \mathsf{E}(X) \cdot \mathsf{E}(Y).$$

**Aalto University**
School of Electrical
Engineering

**ELEC-E7130 - Internet Traffic Measurements and Analysis**
October 4, 2017
12/51

# Contents

- Introduction
- Probability
- Distributions for network measurements
- Parameter estimation
- Model validation

# Distributional models

Distributional (analytic) models for measurement data have certain benefits

- ▶ Models can be manipulated mathematically, leading to improved understanding
- ▶ Models are concise and easily communicated (only a few parameters)
- ▶ Values of model's parameters can give insight into the nature of the underlying data (distribution varies predictably, when its parameters are varied)
- ▶ Models can take into account features that have not been observed or external knowledge

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
14/51

# Discrete distributions

► Discrete distributions are characterized by the probability mass function (pmf)

$$p_i = p(x_i) = \mathsf{P}\{X = x_i\}.$$

► Cumulative distribution $\mathsf{P}\{X \leq x\}$ is given by
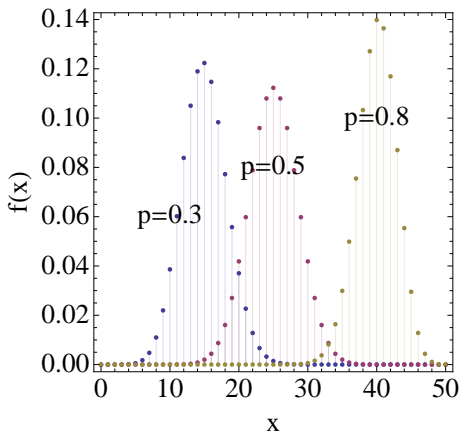
$$F(x) = \sum_{i:x_i \leq x} p(x_i).$$

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
15/51

# Binomial distribution

Number of successes (each with probability *p*) in *n* attempts takes values $\{0, 1, 2, \ldots, n\}$.

Bin(*n*, *p*):

$$P\{X = i\} = \binom{n}{i} p^i (1 - p)^{n-i}.$$

| Mean $\mu$: | $np$ |
| Variance $\sigma^2$: | $np(1 - p)$ |

Aalto University
School of Electrical
Engineering

# Poisson distribution

When $n \to \infty$ and $p \to 0$ in binomial distribution so that $np = \lambda$ takes values $\{0, 1, 2, \ldots\} \Rightarrow$ Poisson distribution.

Poisson($\lambda$):

$$P\{X = i\} = \frac{\lambda^i}{i!}\, e^{-\lambda}.$$

| Mean $\mu$: | $\lambda$ |
|---|---|
| Variance $\sigma^2$: | $\lambda$ |

Aalto University
School of Electrical
Engineering

# Zipf's Law

▶ Consider a set of categorical variables, e.g., URLs of web pages sorted in decreasing number of references made to each page
  ▶ $R$ number of references to a page
  ▶ $n$ rank of the page
▶ Then, for some constants $c$ and $\beta$, Zipf's law states that

$$R = c\, n^{-\beta}.$$

▶ "Discrete power-law distribution"
▶ Linear in log-log plot

$$\log R = \log c - b \log n.$$

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
18/51

# Example: Zipf's Law Applied To WWW Documents



Source: C.Cunha, A. Bestavros, M. Crovella, Characteristics of WWW Client-based Traces, Tech. Report BU-CS-95-010, Boston University, 1995.

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
19/51

# Continuous distributions

- A continuous random variable has a probability density function (pdf), denoted by $f(x)$
- Cumulative distribution function (cdf) defines the probability $P\{X \leq x\}$ and it is denoted by $F(x)$,

$$F(x) = \int_{-\infty}^{x} f(x)\, dx.$$

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
20/51

# Normal distribution $N(\mu, \sigma^2)$

Normal (Gaussian) distribution is denoted by $N(\mu, \sigma^2)$ with mean $\mu$ and variance $\sigma^2$.

Probability density function:

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma} e^{-\frac{(x-\mu)^2}{2\sigma^2}}.$$

| Mean: | $\mu$ |
|---|---|
| Variance: | $\sigma^2$ |

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
21/51

# Exponential distribution Exp($\lambda$)

Exponential distribution with intensity $\lambda$ is denoted by Exp($\lambda$).

Probability density function:

$$f(x) = \lambda e^{-\lambda x}, \quad x \geq 0.$$

| Mean $\mu$: | $1/\lambda$ |
|---|---|
| Variance $\sigma^2$: | $1/\lambda^2$ |

Memorylessness property:
$$P\{X > x+t \mid X > t\} = P\{X > x\}.$$
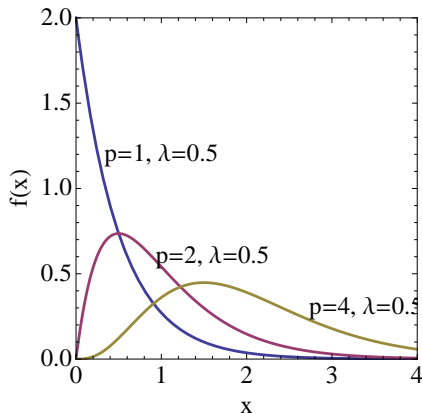
# Gamma distribution Gamma($p, \lambda$)

Probability density function:

$$f(x) = \frac{(\lambda x)^{p-1}}{\Gamma(p)} \, \lambda \, e^{-\lambda x}, \quad x \geq 0,$$

where

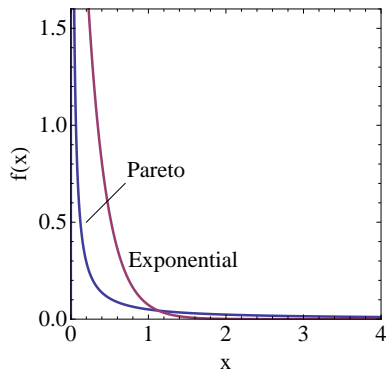$$\Gamma(p) = \int_0^\infty t^{p-1} e^{-t} \, dt.$$

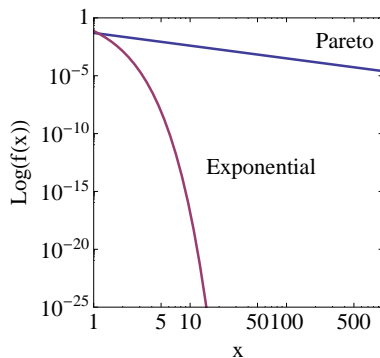| Mean: | $p/\lambda$ |
|---|---|
| Variance: | $p/\lambda^2$ |

Aalto University
School of Electrical
Engineering

# Heavy-tailed distributions

- Heavy-tailed distributions are distributions with a right tail that decays slower than exponentially
- Evidence found, e.g., in sizes of
  - Files stored on Web servers
  - Data transferred through the Internet
  - Files stored in general-purpose Unix file systems
  - I/O traces of file system, disk, and tape activity

**Aalto University**
School of Electrical
Engineering

**ELEC-E7130 - Internet Traffic Measurements and Analysis**
October 4, 2017
24/51

# Visual comparison



(a) Linear scale

(b) Log scale

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
25/51

# Definitions

- Distribution has a **heavy tail** if for all $\gamma > 0$

$$\lim_{x \to \infty} e^{\gamma x} G(x) \to \infty,$$

where $G(x) = 1 - F(x)$, i.e., the ccdf.

- Distribution has a long tail if

$$G(x + t) \sim G(x), \quad \text{as } x \to \infty.$$

- Distribution has a **power tail** if for some $\alpha$ and $\beta > 0$

$$G(x) \sim \alpha x^{-\beta}, \quad \text{as } x \to \infty.$$

- In a nutshell:
  - Large values likely
  - High variability

Aalto University
School of Electrical
Engineering

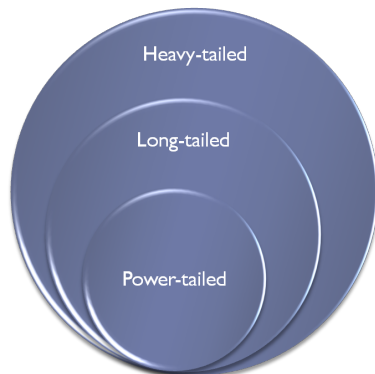ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
26/51

# Relations

Power-tailed $\subset$ Long-tailed $\subset$ Heavy-tailed.

Distribution with a short tail has

$$\lim_{x \to \infty} e^{\gamma x} G(x) \to 0,$$

for some $\gamma > 0$.

Aalto University
School of Electrical
Engineering

# Effects of "heavy tails"

- Expectation paradox:
  - The longer we have waited for an event the longer we have to wait
- Aggregate size of small variables is negligible compared to the largest one

$$\lim_{x \to \infty} \frac{\mathsf{P}\{X_1 + X_2 + \ldots + X_n > x\}}{\mathsf{P}\{\max\{X_1, X_2, \ldots, X_n\} > x\}} = 1, \quad \forall \, n \geq 2.$$

  - Typical flow is small, but typical transferred byte belongs to a large flow

**Aalto University**
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
28/51

# Examples of utilizing heavy tails

- Load balancing in distributed systems
  - Only a few flows are redirected with a significant effect on load distribution
- Scheduling in web servers
  - Shortest-remaining-processing-time scheduling lets the small tasks interrupt larger ones and with heavy tails the benefit becomes large
- Routing and switching in the Internet
  - Shortcuts established only for large flows (cf. data centers)

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
29/51

# Log-normal distribution, LogNormal($\mu, \sigma^2$)

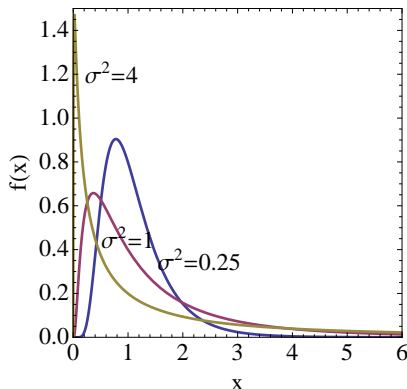Def.: *Random variable X follows the LogNormal($\mu, \sigma^2$) distribution if* $\log(X)$ *is distributed as N($\mu, \sigma^2$).*

Probability density function:

$$f(x) = \frac{1}{\sqrt{2\pi}\sigma \cdot x} e^{-\frac{(\log x - \mu)^2}{2\sigma^2}}.$$

Mean:   $e^{\mu + \sigma^2/2}$

Variance:   $\left(e^{\sigma^2} - 1\right) e^{2\mu + \sigma^2}.$

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
30/51

# Pareto distribution Pareto($\alpha, k$)
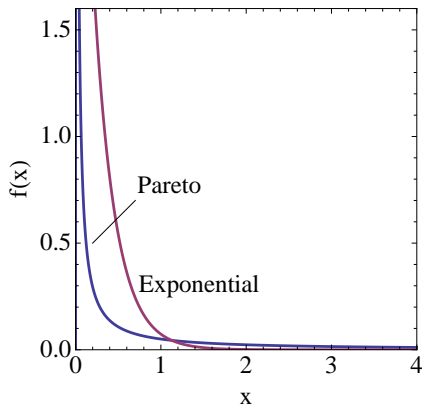
Probability density function:

$f(x) = \alpha k^{\alpha} x^{-\alpha-1}$.

- $\alpha$ is the **shape parameter**
- k is the **scale parameter**

Pareto has a **power tail**,

$$G(x) = \left(\frac{k}{x}\right)^{\alpha}, \quad x \geq k.$$

Mean: $\dfrac{\alpha k}{\alpha - 1}$

Variance: $\left(\dfrac{k}{\alpha - 1}\right)^2 \dfrac{\alpha}{\alpha - 2}$



Note: mean is infinite for $\alpha \leq 1$, and variance for $\alpha \leq 2$.

Aalto University
School of Electrical
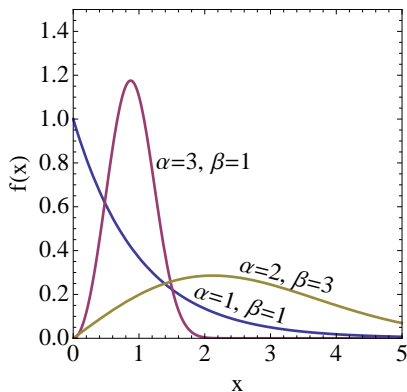Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
31/51

# Weibull distribution, Weibull($\alpha, \beta$)

In Weibull($\alpha, \beta$) distribution $\alpha$ is the shape parameter and $\beta$ the scale parameter.

- ► $\alpha < 1$ "failure rate decreases in time"
- ► $\alpha > 1$ "failure rate increases in time"

The pdf has the form

$$f(x) = \alpha \beta^{-\alpha} x^{\alpha-1} e^{-(x/\beta)^\alpha}.$$

| Mean: | $\beta\, \Gamma\left(1 + 1/\alpha\right)$ |
|---|---|
| Variance: | $\beta^2\, \Gamma\left(1 + 2/\alpha\right)$ |
| | $-\beta^2\, \Gamma\left(1 + 1/\alpha\right)^2$ |



$\alpha=3$, $\beta=1$

$\alpha=2$, $\beta=3$

$\alpha=1$, $\beta=1$

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
32/51

# Choosing the distribution

- Distribution can be chosen according to "best fit"
- Distribution function can be adopted from a similar situation
  - Often assumed that the distribution function remains valid in other "similar conditions", only parameters vary
- Distribution functions can emerge from generative processes
  - E.g., CLT and Gaussian distribution
  - Allows taking into account external knowledge on the variable
- As a result probabilistic modeling of network measurements is seldom purely objective process ...

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
33/51

# Contents

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
34/51

# Parameter estimation

- Given a distribution (or distribution family) and data, the next step is to fit the parameters of the distribution to match the data
- **Estimator** is a function of (sample) data that attempts to estimate an unknown (population) parameter
- We try to **optimize** the parameters of a density with respect to some **measure of fit**

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
35/51

# Estimation

- ▶ We briefly outline two practical methods of parameter estimation
    1. Method of moments
    2. Maximum likelihood (ML)
- ▶ Note that a lot of literature exists on parameter estimation
    - ▶ Quality of the estimators ignored here . . .
    - ▶ E.g., sample mean and sample variance of data are "best" estimates for Normal distribution

**Aalto University**
School of Electrical
Engineering

**ELEC-E7130 - Internet Traffic Measurements and Analysis**
October 4, 2017
36/51

# Method of moments

- $k$:th moment of a random variable: $E(X^k)$
- k:th sample moment of data: $m_k = \frac{1}{n} \sum_i x_i^k$
- **Method of moments**:
  1. Derive as many moments of the distribution as there are parameters
  2. Compute the corresponding sample moments from data
  3. Solve the parameters so that the moments and sample moments are equal
- Simple, but not always available
- The estimates are not necessarily "optimal"

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
37/51

# Maximum likelihood

▶ Idea is to maximize the probability/likelihood that the selected distribution has produced the observations

▶ Likelihood function for independent observations

$$L(x_1, x_2, \ldots, x_n; \theta) = \prod_{i=1}^{n} f(x_i; \theta).$$

▶ Likelihood is a function of the parameter(s) of the distribution, for which the estimate is

$$\arg\max_{\theta} \, L(x_1, x_2, \ldots, x_n; \theta)$$

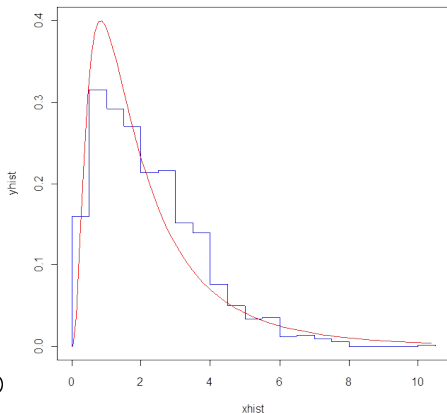▶ Maximum likelihood estimates have many favorable properties, but require often complex non-linear optimization

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
38/51

# Other estimates

- Non-linear least-squares optimization of the density
  - Non-linear optimization methods more generally available in mathematical software
- Statistical software offer many direct ways of fitting parameters for given densities
  - Often based on maximum likelihood
  - General optimization methods can be sensitive to selected starting values
  - Results are affected by outliers

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
39/51

# Example

```
> require(MASS)
> fitdistr(logftp,"lognormal")
    meanlog         sdlog
 0.52062324      0.84232334
(0.02663660)    (0.01883492)
```



```
> h<-hist(logftp,n=20)
> xhist<-c(min(h$breaks),h$breaks)
> yhist<-c(0,h$density,0)
> xfit<-seq(min(logftp),max(logftp),length=100)
> yfit<-dlnorm(xfit,meanlog=0.52062,sdlog=0.84232)
> plot(xhist,yhist,type="s",ylim=c(0,max(yhist,yfit)),col=4)
> lines(xfit,yfit,col="red")
```
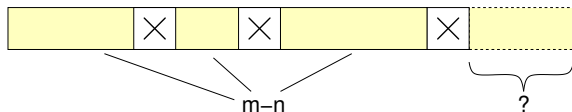
Aalto University
School of Electrical
Engineering

# Example: German tank problem

Suppose that a measurement has given us

- A set of serial numbers of some device
  - How many devices has been sold?
- A list of user-id's of people using some service
  - How many users the given service has (say globally)?

These questions are examples of the *German tank problem*.

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
41/51

# Example: German tank problem (2)

Assumptions:

- ▶ The existing serial numbers are $1, \ldots, N$ (all equally likely)
- ▶ Sampling without replacement, $\{X_i\}$, $i = 1, \ldots, n$
- ▶ Task is to estimate $N$ based on the given $n$ samples

Clearly, $n \leq m \leq N$, where $m = \max\{X_i\}$, but ...?

1. Bayesian estimate,

$$\hat{N} = \frac{(m-1)(n-1)}{n-2}.$$

2. Minimum-variance unbiased estimator,

$$\hat{N} = m + \frac{m}{n} - 1.$$

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
42/51

# Off-line vs. On-line estimation

**Off-line estimation**

- ▶ Collect samples and store them
- ▶ Estimate the quantities of interest
- ▶ No (strict) memory or time constraints

**On-line estimation**

- ▶ Collect samples in real-time (streaming data)
  - ▶ Update estimates at the same time
- ▶ Both memory and time constraints (typically)

The latter is important, e.g., for real-time monitoring systems.

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
43/51

# Example: On-line estimation

**Mean:**

Init:
1. $S \leftarrow 0$
2. $n \leftarrow 0$

Per sample:
1. $S \leftarrow S + x_i$
2. $n \leftarrow n + 1$
3. $\hat{m} \leftarrow S/n$

- Two state variables
- Fast constant computation time
- Unfortunately, e.g., median or mode are more difficult! (Why?)

Sometimes(!) sufficient:

$$\begin{cases} \text{mean} & \leftarrow & \text{mean} & + & \eta \times (x_i - \text{mean}) & \text{(cf. EWMA)} \\ \text{median} & \leftarrow & \text{median} & + & \eta \times \text{sgn}(x_i - \text{median}) \end{cases}$$

**Aalto University**
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
44/51

# Example: On-line estimation (2)

- Moving average (MA):

$$y \leftarrow \frac{x_{i-k} + \ldots + x_i}{k+1}$$

- Exponentially Weighted Moving Average (EWMA):

$$y \leftarrow \alpha x_i + (1 - \alpha) y$$

In general, on-line estimation of the streaming data is an interesting topic itself, and there are advanced algorithms for different scenarios. For example,

R. Jain and I. Chlamtac, "*The P-Square Algorithm for Dynamic Calculation of Percentiles and Histograms without Storing Observations*", Communications of the ACM, October 1985.
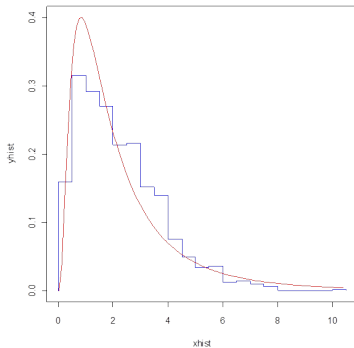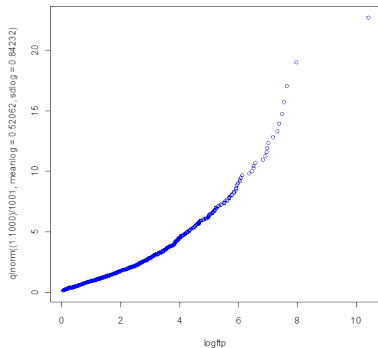
**Aalto University**
School of Electrical
Engineering

**ELEC-E7130 - Internet Traffic Measurements and Analysis**
October 4, 2017
45/51

# Contents

- Introduction
- Probability
- Distributions for network measurements
- Parameter estimation
- Model validation

Aalto University
School of Electrical
Engineering

# Checking for model fit

- Visual tools
  - Plot density and histogram in the same figure
  - Plot cdf and ecdf in the same figure
  - Compare data with the distribution in a QQ-plot
  - For highly variable data, log-log complementary distribution could be considered
  - Plot $\log(1 - F(x))$ against $\log x$ for CDF and ECDF
- Statistical tests
  - $\chi^2$-test
  - Kolmogorov-Smirnov -test

**Aalto University**
School of Electrical
Engineering

**ELEC-E7130 - Internet Traffic Measurements and Analysis**
October 4, 2017
47/51

# Example

- Our log-normal fit
  - QQ-plot does not support the fit!



```
>qqplot(logftp,qlnorm((1:1000)/1001,meanlog=0.52062,sdlog=0.84232),col=4)x))
```

Aalto University
School of Electrical
Engineering

# Example Continued

Let's try with Weibull distribution . . .

```
> fitdistr(logftp,"weibull")
   shape        scale
1.49738316    2.47066018
(0.03691000) (0.05496849)
```

Aalto University
School of Electrical
Engineering

# Statistical testing for goodness of fit

- As the model selection is usually subjective, visual tools are generally sufficient for validating the model
- However, there are also statistical tests available for the **goodness of fit**
  - A $p$-value is computed for null hypothesis
    
    "*Sample comes from a population with a given distribution*"
  - $p$-value is roughly the probability that given the null hypothesis, we actually observe the data
  - If $p$-value is small, there is only a small probability that the data is from the distribution and null hypothesis is rejected
  - A large $p$-value does not automatically mean that the distribution is correct

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
50/51

# Literature

- Mark Crovella, Balachander Krishnamurthy, **Internet Measurement: Infrastructure, Traffic & Applications**, John Wiley and Sons Ltd, 2006.
- David J. Marchette, **Computer Intrusion Detection and Network Monitoring: A Statistical Viewpoint**, Springer, 2001.
- Pierre Baldi, Paolo Frasconi, Padhraic Smyth, **Modeling the Internet and the Web, John Wiley and Sons Ltd**, 2003.
- NIST/SEMATECH, Engineering Statistics Handbook, `http://www.itl.nist.gov/div898/handbook/index.htm`

Aalto University
School of Electrical
Engineering

ELEC-E7130 - Internet Traffic Measurements and Analysis
October 4, 2017
51/51