

UNIMORE

UNIVERSITÀ DEGLI STUDI DI  
MODENA E REGGIO EMILIA



# Simulation

# Simulation overview

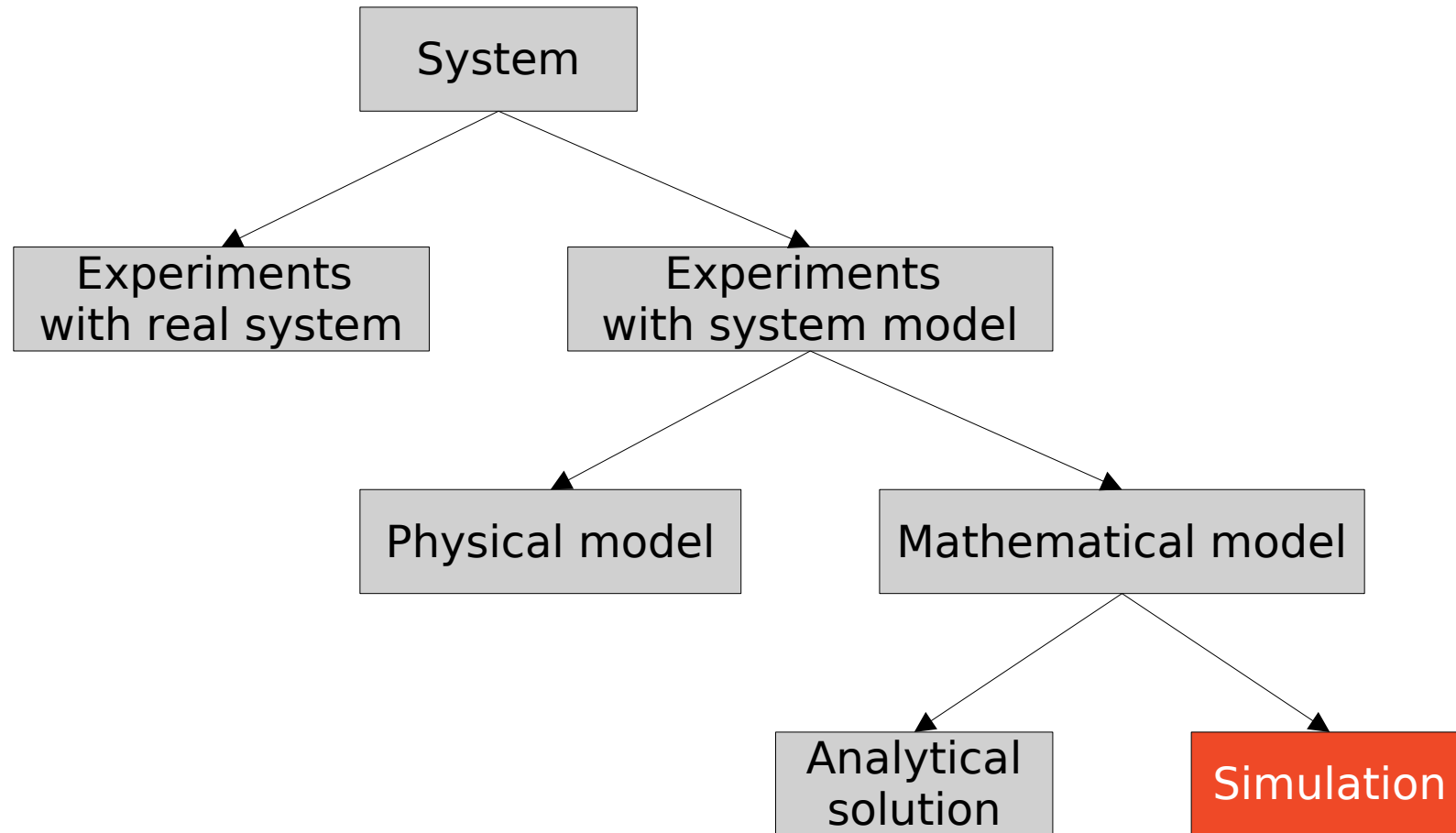
# Goal of simulation

---

- **Build model of a system that captures its behavior**
  - Can test under different conditions
  - Can evaluate modification to the system
- **System**: set of parts cooperating for a common goal
- **Model**: abstract representation of system
  - Typically simplified version
  - Level of simplification depends on goal of the analysis
- Model typically implemented as computer program

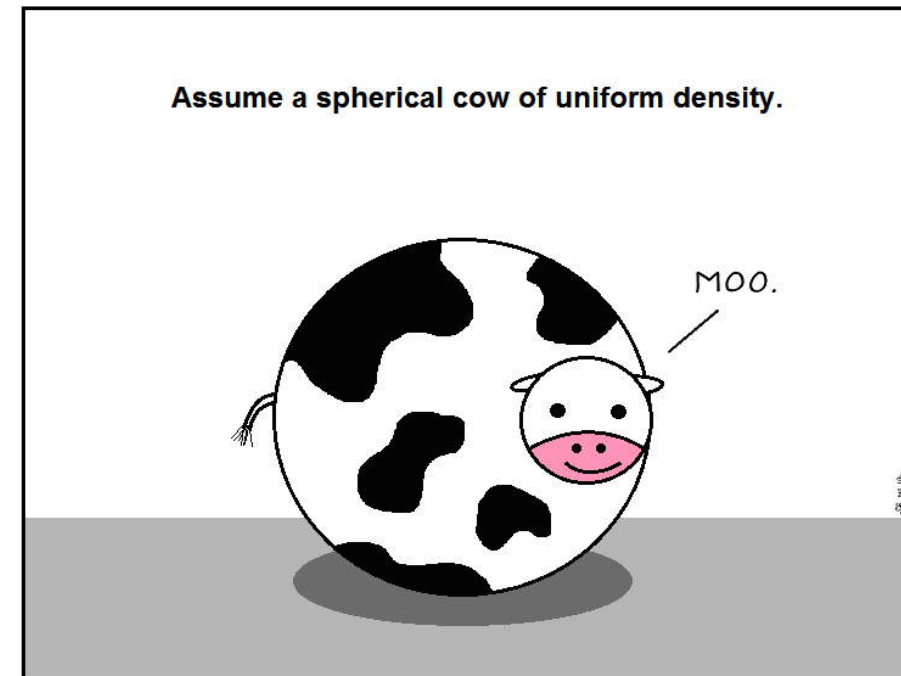
# Possible approaches

---



# Motivations for simulation

- Actual system vs. **System model**
  - Easier to do/Less dangerous (if system exists)
  - Less expensive than building a real system
- Physical vs. **Mathematical model**
  - Easier and less expensive to implement
  - Easier to modify
  - Mathematical model must be **accurate enough**
  - Remember the Spherical Cow
- Analytical solution vs. **Simulation**
  - Analytical is more precise
  - Simulation is more **flexible**



# Types of simulations

---

- Static vs. **Dynamic**
  - **Time** is considered or not?
- Deterministic vs. **Stochastic model**
  - Typically **random variables** are used in the simulation
- Discrete vs. Continuous model
  - **Discrete** → **event-based** simulation
  - **Continuous** → time cannot be split into events
  - **Depends on the system and what we are looking for**
  - Example: car traffic
    - Model each car → model is discrete (movement of car)
    - Model traffic → flow-based model (differential equations)
  - See also fluid models in large systems

# Discrete simulator

---

- Several components:
- **System state**
  - Includes system clock (“now” variable)
- Statistical counters (**data collectors**)
- **Routines:**
  - Initialization (setup for  $t=0$ )
  - Timing (select next event)
  - Event management
- Report generation
- Main()
  - Startup and shutdown

# Steps for good simulation

---

- Problem and study formulation
  - Define the **problem** of interest
  - Define:
    - **Goal** of the study
    - **Questions** to be answered
    - **KPI** to measure
    - **Scope** of model
    - Configurations to consider (**scenarios**)
- Model definition
  - Collect **information on system**
    - Several people to involve
    - Check for inaccurate information
    - Check for not-formalized aspects of problem



# Steps for good simulation

---

- Model definition (continued)
  - Collect **data** for **model parameters** (if possible)
  - **Document** the model assumptions
  - Collect **data** on performance of existing system (for **validation**)
  - Define **level of detail** (remember the Spherical Cows)
- **Model validation**
  - Discuss with other people (experts, management, ...)
- Software development
- **Simulator validation**
  - Compare with **data** (if available) or **theoretical model**
  - Run **sensitivity analyses** on parameters

# Steps for good simulation

---

- Experiments design
  - Run lengths
  - Warm-up period (avoid transients)
  - Number of runs
- Experiments
- Data analysis
- (Good) reporting
  - Simulation presents data, scenarios, ...
  - Management takes decisions
  - Point out assumptions and limits of validity
  - Get paid

# Advantages of simulation

---

- Can deal with **complex systems**
  - Hard to manage **random variables** with analytical models
  - Hard to find **analytical models** to describe them
  - Analytical models are hard to solve in **closed form**
- Can analyze **several scenarios**
  - Existing system with different operating parameters (**what-if**)
  - Variations of existing system (**capacity planning**)
- **Easier to control** than prototype systems
  - External interferences like network status
- Explore **large time windows**

# Drawbacks of simulation

---

- Use of stochastic processes
  - Simulation provides **just an estimation**
- (Good) simulators are **complex to develop**
- **Simulators are as good as the model they implement**
  - GIGO systems → Garbage In Garbage Out
  - A word of wisdom from the past:



*“ I have been asked: ‘Mr. Babbage, if you put into the machine wrong figures, will the right answers come out?’ I am not able rightly to apprehend the kind of confusion of ideas that could provoke such a question.”*

# Some common pitfalls

---

- Not having well-defined **objectives** for the simulation
  - “everything” is not a well-defined objective
- Wrong **level of detail**
  - Too much detail → simulation not scalable (common newbie mistake, often related to previous point)
  - Too little detail → inaccurate/wrong results
- **Wrong model**
  - Input variables models?
  - KPI?
- Considering simulation an exercise in computer programming
  - Not considering simulation as a **complex set of tasks**
  - Error in **statistical analysis** (e.g., just one run)

# Recall of statistics

- **Experiment**: a **sample** of a random variable  $X$
- **Sample space**  $S$ : the set of **possible outcomes** (Co-dominion of random variable)
  - Each sample belongs to  $S$
  - Example:
    - Coin toss  $\rightarrow S=\{H, T\}$
    - 1d6 roll  $\rightarrow S=\{1, 2, \dots, 6\}$
  - $S$  can be a **set** of finite values, an **interval** or  **$\mathbb{R}$**
- **Probability density**:
  - $P(x) = P(X=x)$
  - The **sum or probabilities is 1**
    - $\sum_{x \in S} p(x)=1$
    - $\int p(x) dx = 1$

# Relevant metrics

---

- **Average** value  $\mu = E(X)$ 
  - $E(X) = \sum x p(x)$
  - Relevant properties:
    - $E(cX) = cE(X)$
    - $E(\sum_i c_i X_i) = \sum_i c_i E(X_i)$
- **Median** value
  - Value  $x$ :  $P(X \leq x) = 0.5$
- **Quantiles** are defined in a similar way
- **Variance**  $\sigma^2 = E[(X - \mu)^2] = E(X^2) - \mu^2$ 
  - $E[(X - E(X))^2] = E[X^2 - 2X E(X) + E(X)^2] =$   
 $= E(X^2) - 2E(X)E(X) + E(X)^2 = E(X^2) - 2E(X)^2 + E(X)^2$
- **Standard Deviation**  $\sigma = \sqrt{\sigma^2}$



# Variable independence

---

- Variables can be:
  - Independent
  - Correlated
- Independent variables can be studied separately
- For correlated variables the probability distribution are intertwined
  - $P(X=x)=f(P(Y=y))$
- In simulation variables are often correlated
- To obtain statistically valid results we need to achieve variable independence

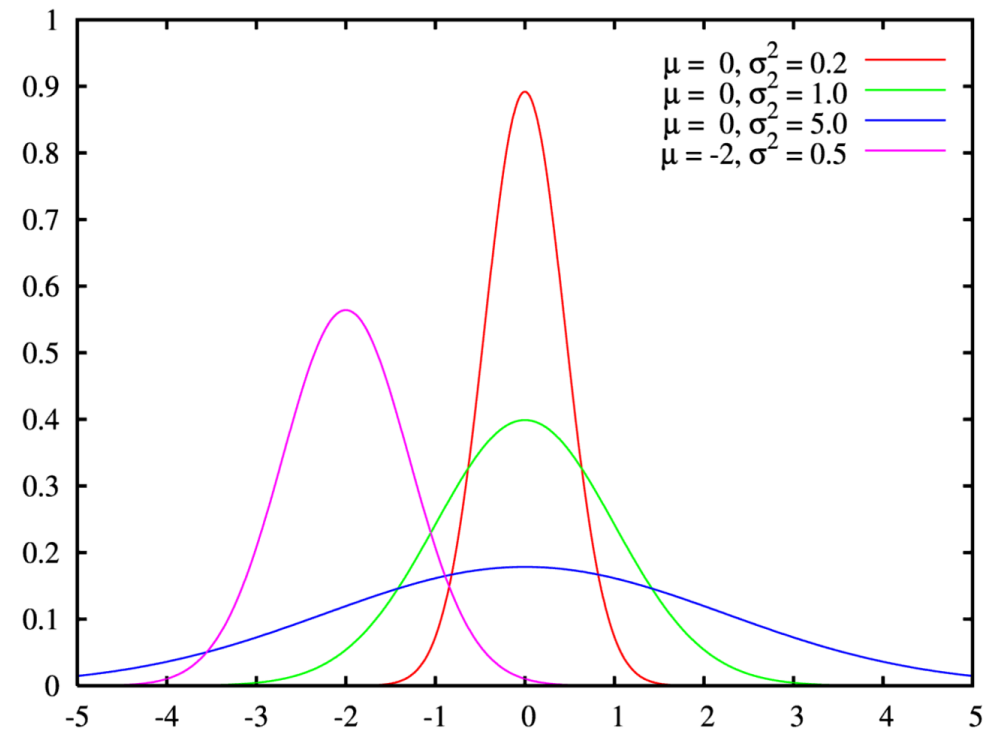
# Repeated experiments

---

- In simulation we **repeat experiments**:
  - Each experiments starts with **a different seed for RNG**
  - Experiments are **statistically independent**
- From **central limit theorem**
  - $F(X) = P(X \leq x)$
  - $F(X) \rightarrow G(X)$  for  $N \rightarrow \infty$
  - $G(X)$  = Gaussian function
- **Law of large numbers**
  - $N$  = number of observations
  - $N \rightarrow \infty$  means that sample average converges to  $E(X)$

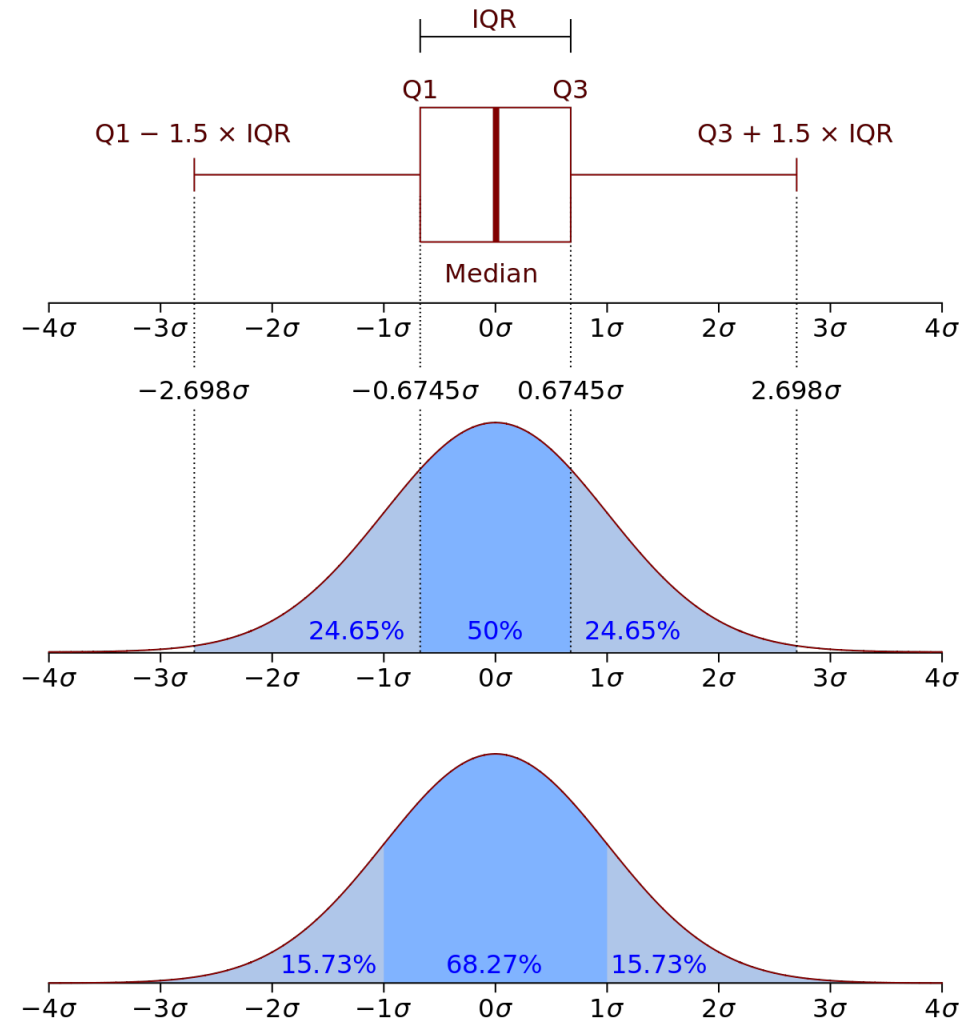
# Gaussian functions

- Meaning of **variance**:
- **Low** variance:  
→ High and narrow peak
- **High** variance:  
→ Low and wide peak



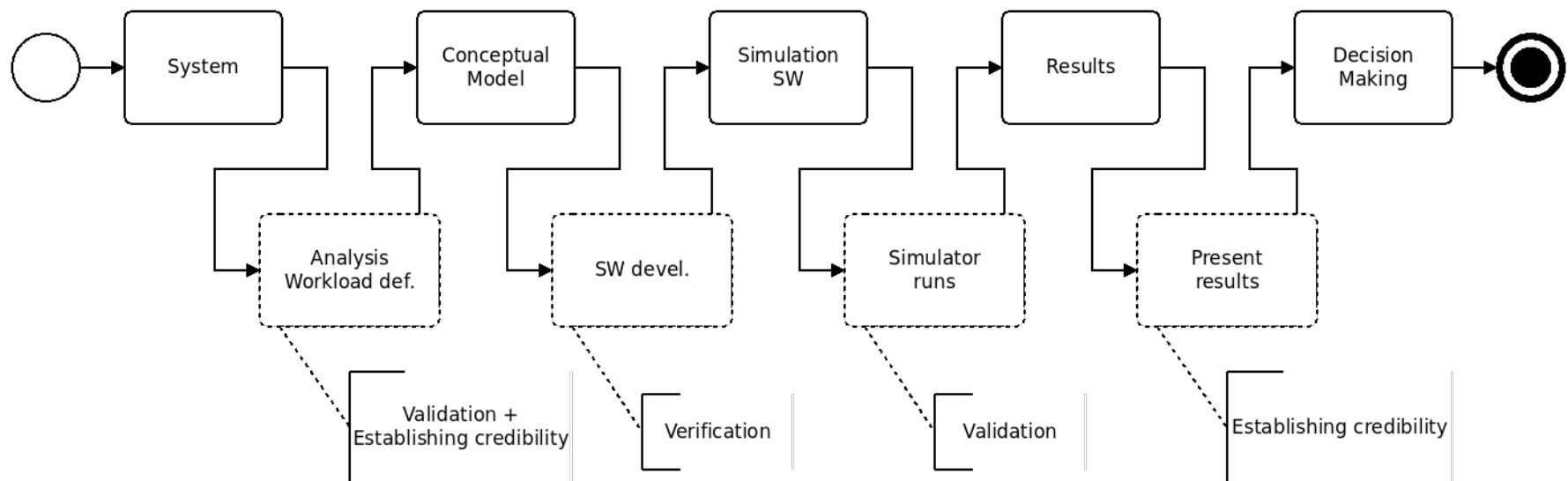
# Confidence intervals

- Analysis of a **Gaussian distribution**
- Median=Average
- **Confidence:**
  - $\mu \pm \sigma \rightarrow 68\%$  confidence
  - $\mu \pm 3\sigma \rightarrow 99\%$  confidence
- Need to estimate  $\mu$  and  $\sigma$  from repeated experiments



# Building models

- **Validation**
  - Model is an accurate representation of the system,
- **Verification**
  - Concepts maps into SW implementation
- Establishing **credibility**
  - Information accepted as credible from experts/management



- **Example:**
  - Accreditation with DoD in USA
- Key elements to be considered
  - Verification and validation correctly carried out
  - Quality of SW development (including historical data)
  - Quality of available data
  - Quality of documentation
  - Known problems and limitation in simulator or model
- Sources of errors → potentially inaccurate results
  - Approximations in **system/model**
  - Approximations in **simulator runs** (e.g., # of runs, length)

# Some remarks

---

- Common newbie **mistake**: **add every detail**
  - **Un-scalable** simulation
  - **Huge** parameter **space to explore**
- Detail depends on the data available
  - No need to model input variables parameters if no data are available
  - Rely on simple, standard assumptions



# Verification techniques

---

- Use **modular** approach
  - Unit testing of components
- Use **external code review**
  - Walk-through sessions with other developers
- Explore large **parameter space**
  - See if the results make sense
  - Compare with theoretical models when possible
  - Compare with historical results if possible
- Use **traces**
  - Easier to reproduce experiments and to analyze problems

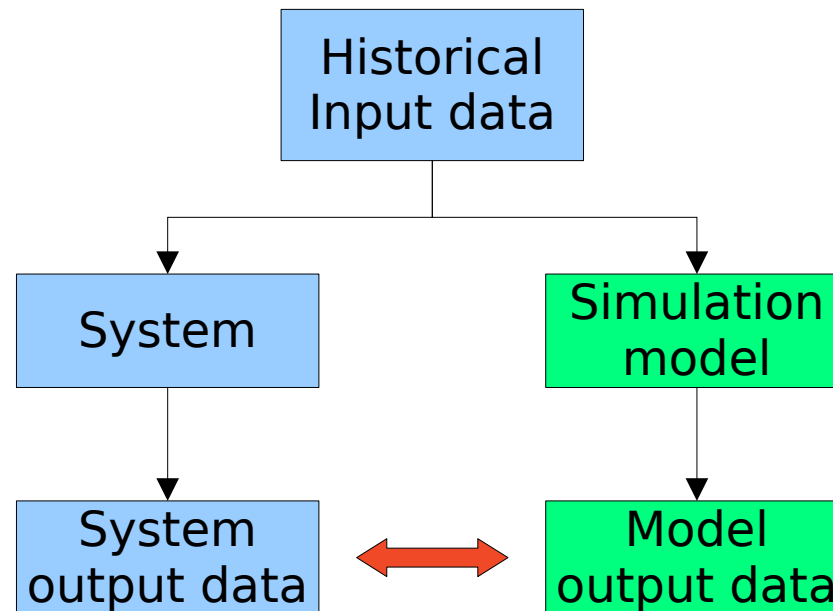
# Validation / Credibility-improving techniques

---

- Seek wisdom from **experts**
  - No need to make a fool of yourself
  - Especially if you are not an expert on the subject
- **Observe the system**
  - Collect data
  - Evaluate data Quality (Missing points, variance issues, ...)
- Rely on **existing** theory/studies
  - Two weeks in lab can save you a 30-min visit to the library...
  - Consider previous simulation projects

# Validation / Credibility-improving techniques

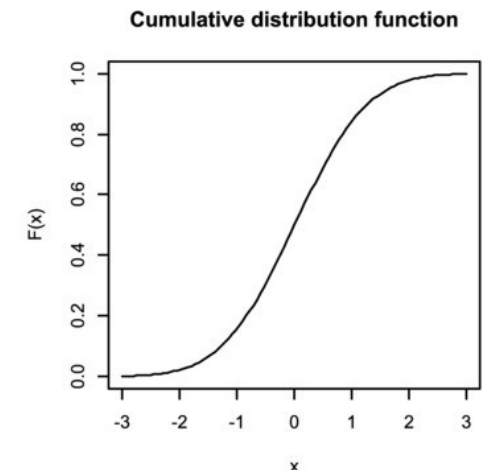
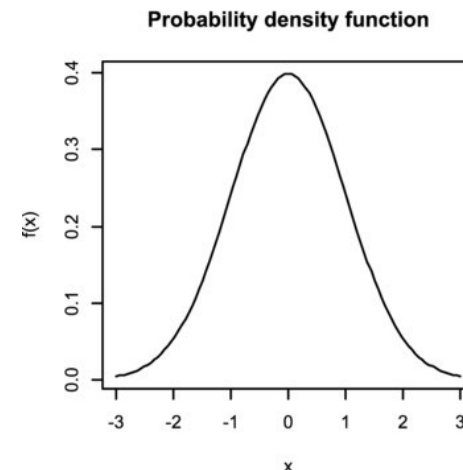
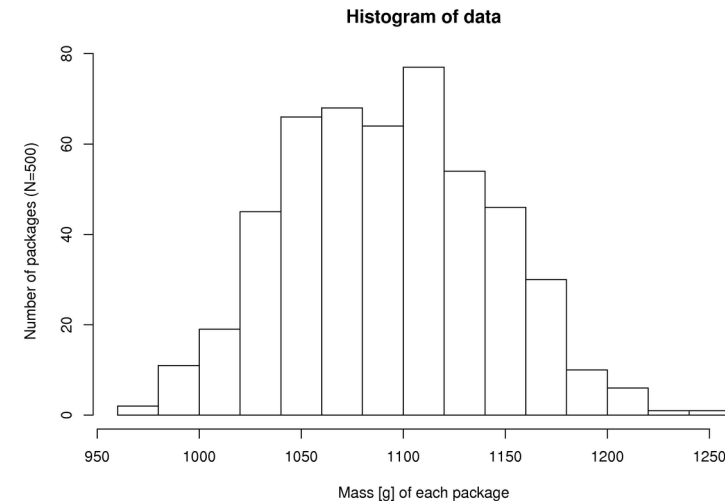
- Correlation-inspection approach
- Strong validation for existing/known scenario



# Describing input

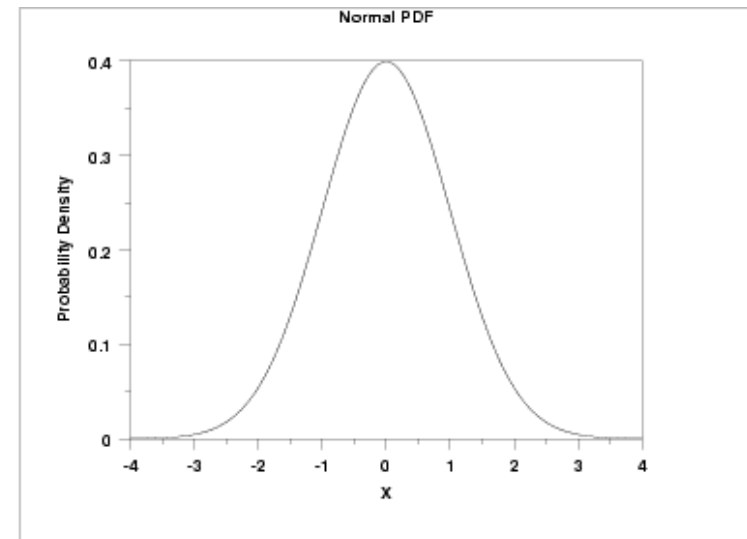
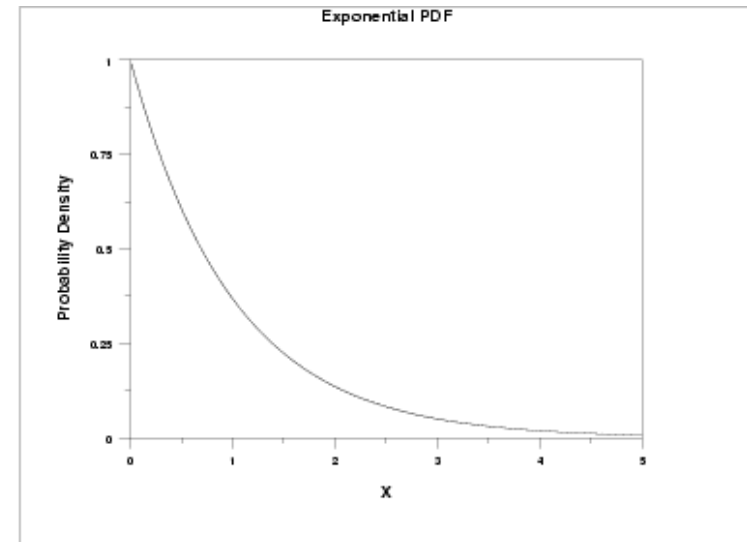
# Input variables

- Input variables are **random** variables
- Typically described **statistically**
- Use of a probability distributions
  - Derived from real observation (**histogram**)
  - **Mathematical** model



# Some useful distributions

- **Exponential**
  - Typical inter-arrival time
  - Poisson process
- **Gaussian** (Normal)
  - Independent samples of same process



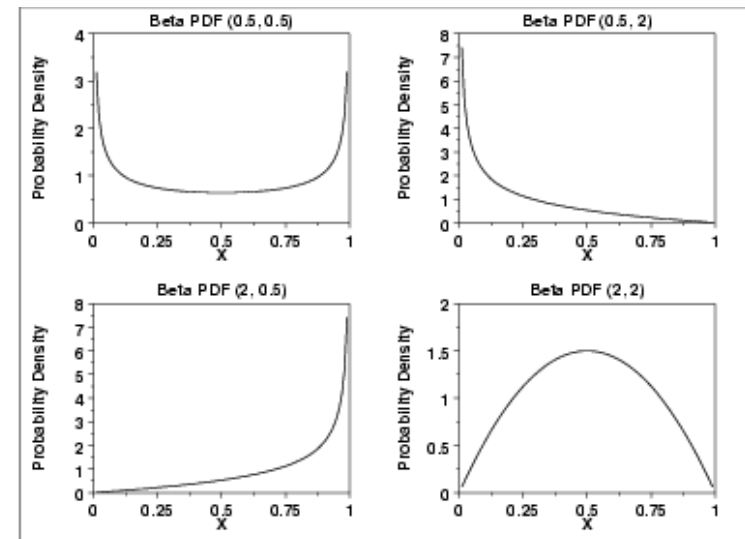
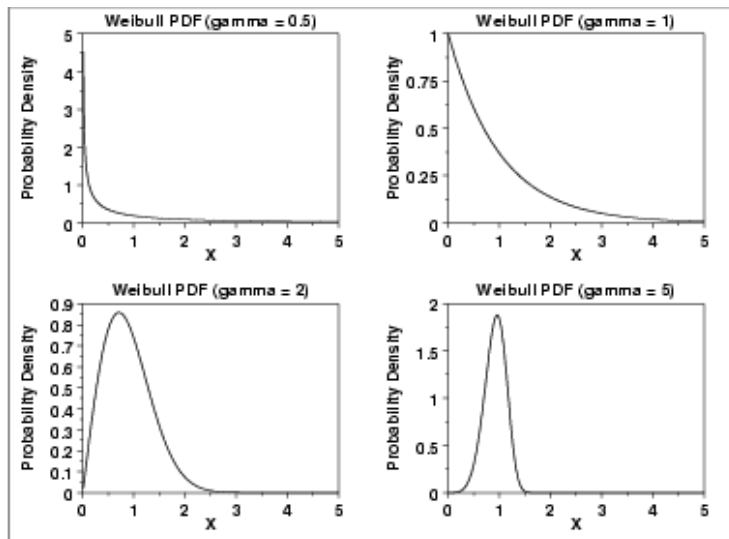
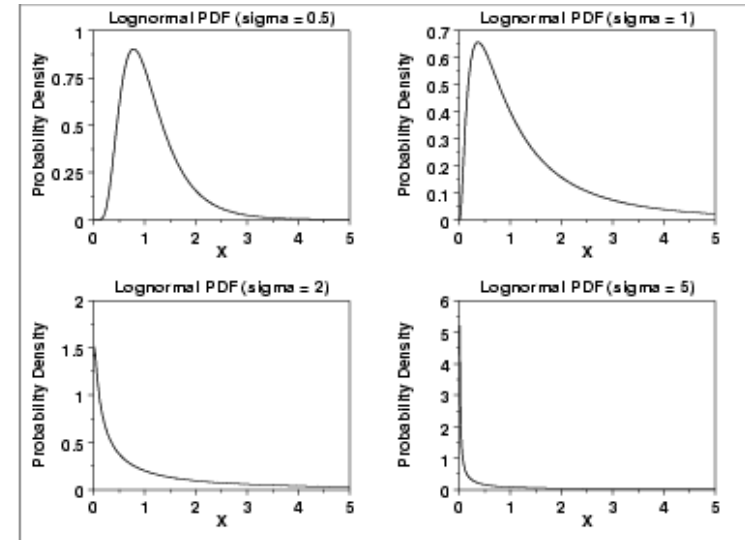
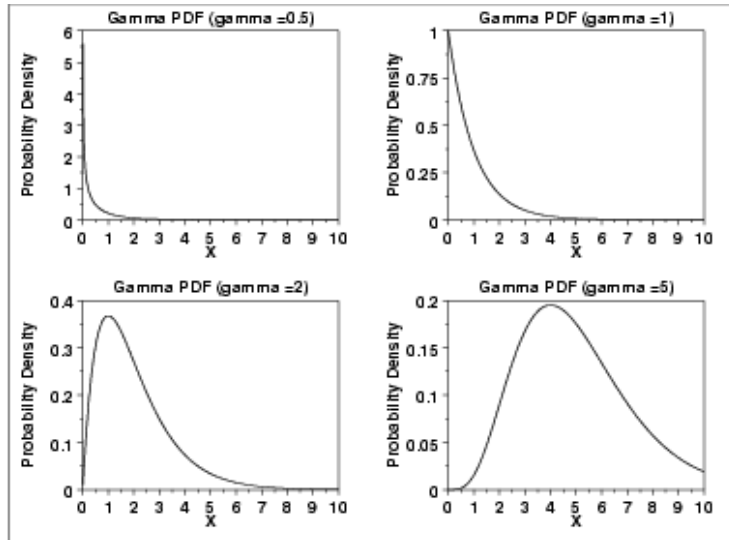
# Some useful distributions

UNIMORE

UNIVERSITÀ DEGLI STUDI DI  
MODENA E REGGIO EMILIA

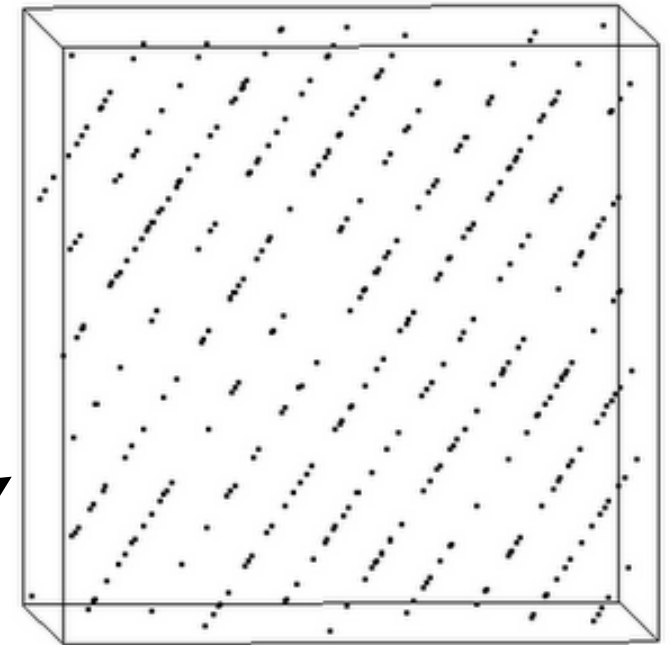


1175



# A remark on RNG

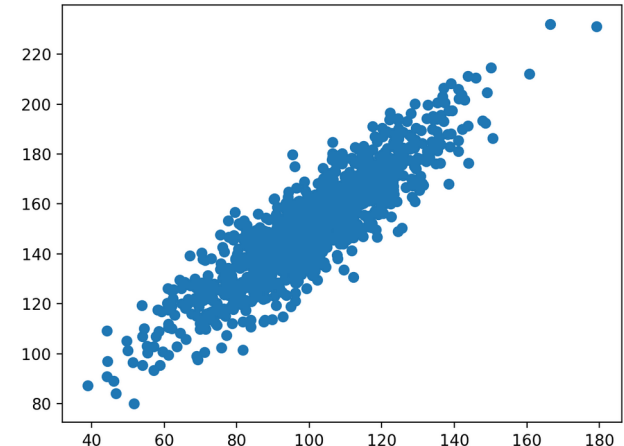
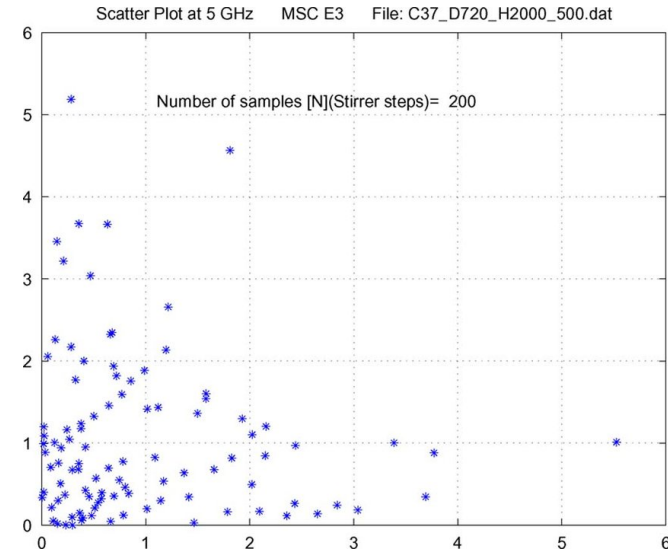
- **RNG**: Random Number Generators
  - Used to produce input variables
  - Starts from a **seed** value like `srand()`
- **Not really random!**
- The **good**:
  - **Reproducible** experiments
  - Multiple seeds → independent experiments
- The **bad**:
  - Can be really **un-random-like**
  - Need to select the RNG
  - Need to select the seed





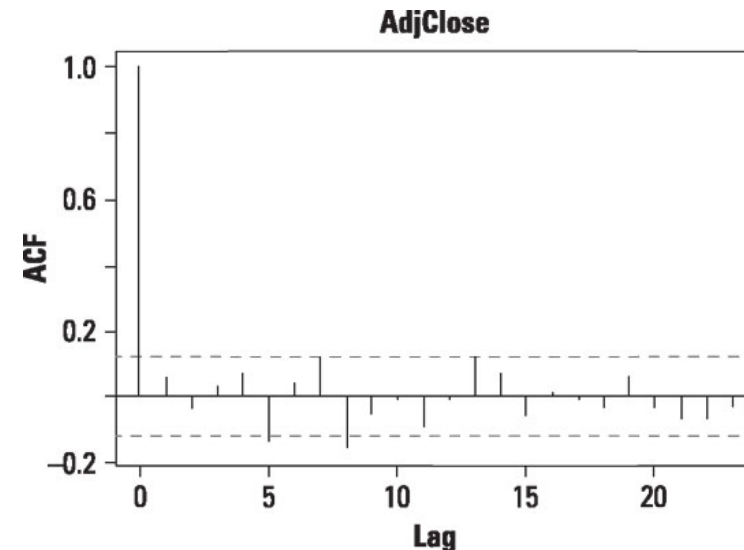
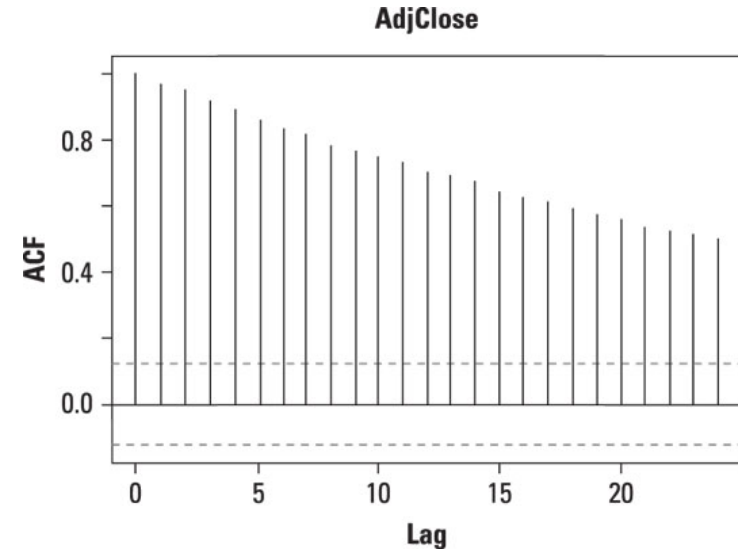
# Assessing sample independence

- Variable **independence**
  - **Among** variables
  - Correlation analysis
- Techniques:
  - **Correlation** value
  - **Visual** inspection
- Correlation (**scatter**) plot
  - Two **examples**
  - Uncorrelated data
  - Correlated data



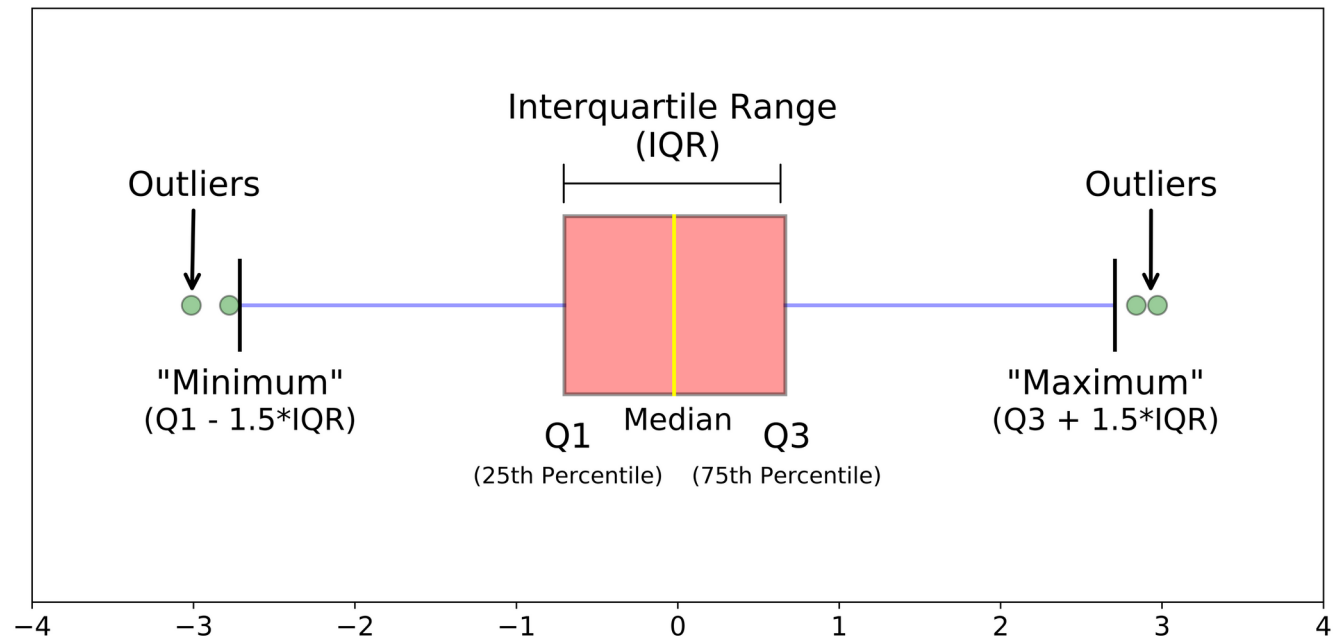
# Assessing sample independence

- **Auto-correlation**
- Correlation between
  - Current series
  - Time-shifted series
- Examples of:
  - Auto-correlated values
  - Non auto-correlated values



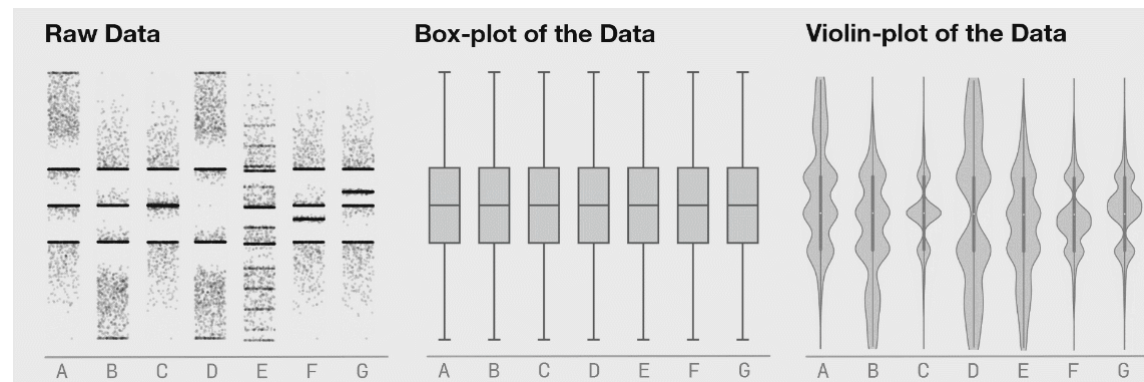
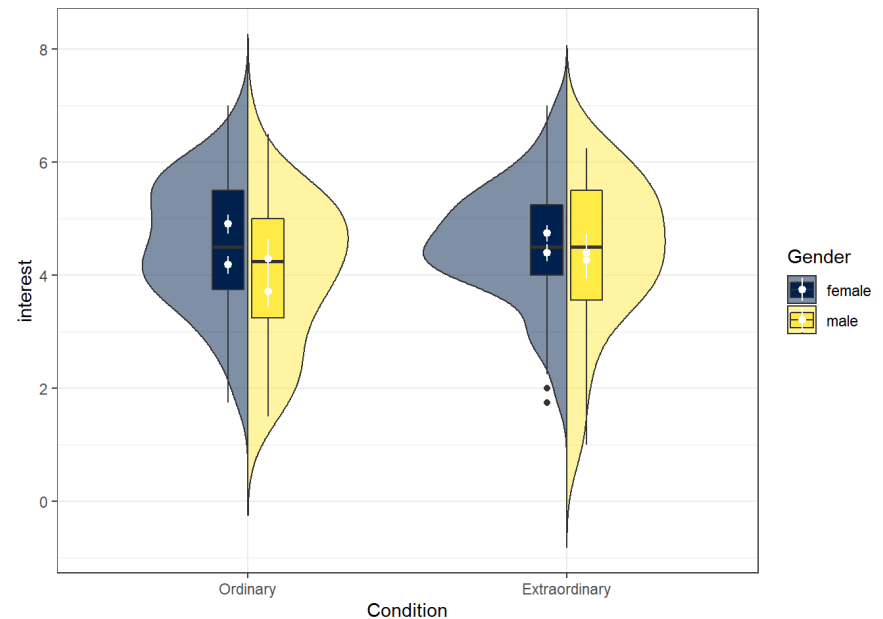
# Short representation of random variable

- **Boxplot**
- Compact representation of a random variable
- Main key points of
  - PDF function
  - Samples
- **Median**
- **Quartiles** (25-perc, 75-perc)
- Estimation of **min/max**
  - Useful for Gaussian var.
- **Outliers**



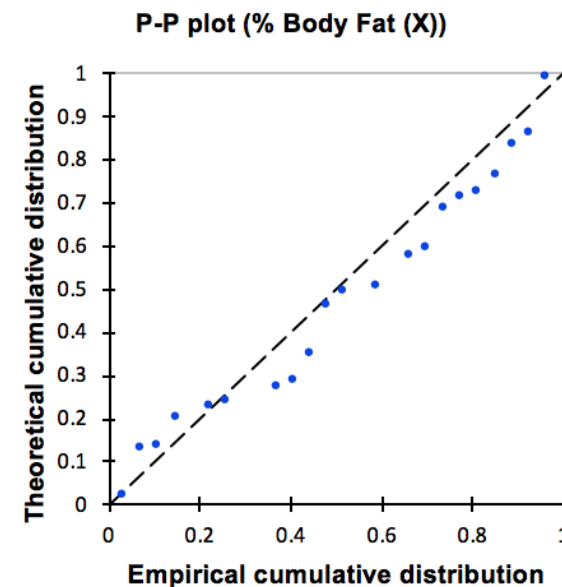
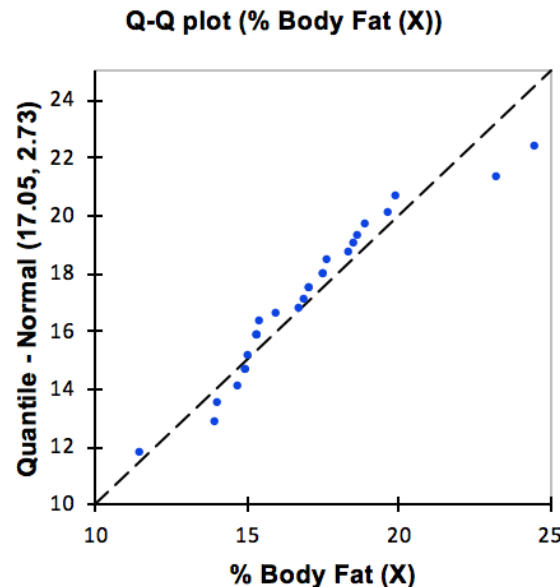
# Short representation of random variable

- **Violin plot**
- Composite plot
  - Boxplot +
  - Probability density
- Useful to:
  - Provide **detailed** statistical description
  - Compare different variables
- Can be valuable when distribution is irregular



# Variable comparison

- Given two variables with PDF
  - $F(x) = P(F \leq x)$
  - $G(x) = P(G \leq x)$
- **P-P plot**
  - Plot  $(F(x), G(x))$  for all  $x$
- Percentile view
  - $Q(i, F(x)) \rightarrow i$ -th percentile
- **Q-Q plot**
  - Plot  $(Q(i, F(x)), Q(i, G(x)))$  for  $i$  in  $[1..100]$
- More emphasis on **head** and **tail**



# Data analysis

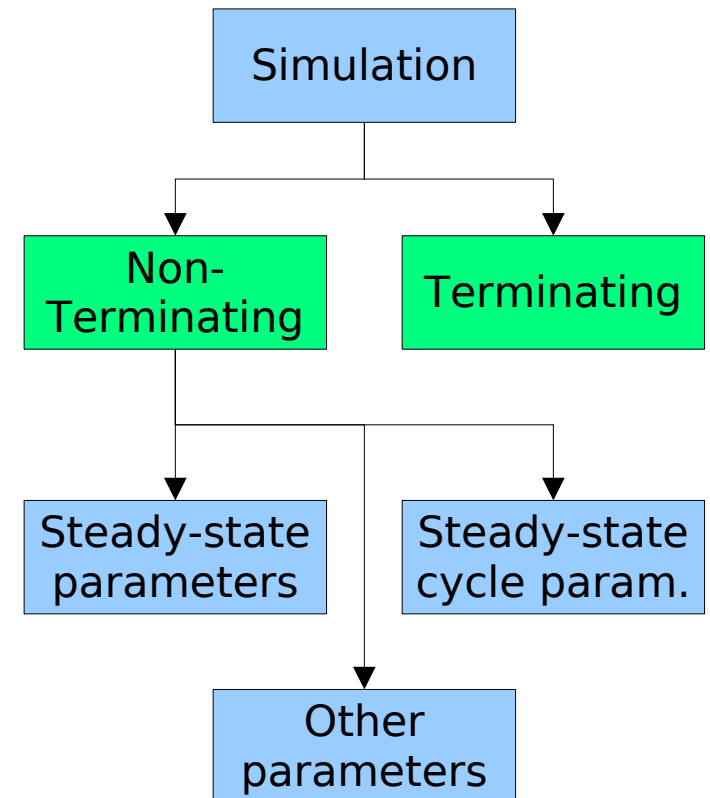
# Transient and steady-state

---

- Simulation goals:
  - Analysis of **steady** state
  - Analysis of **transient** behavior
- Examples of steady state analyses:
  - Throughput of a computing system under a known workload
  - Steady state is not necessarily **constant**
    - Workload changes over the time of day
    - Workload evolution is much slower than client requests
  - Transient effect are present but not wanted
- Examples of transient analysis:
  - Throughput evolution during traffic surge and cloud scale-up

# Types of simulation

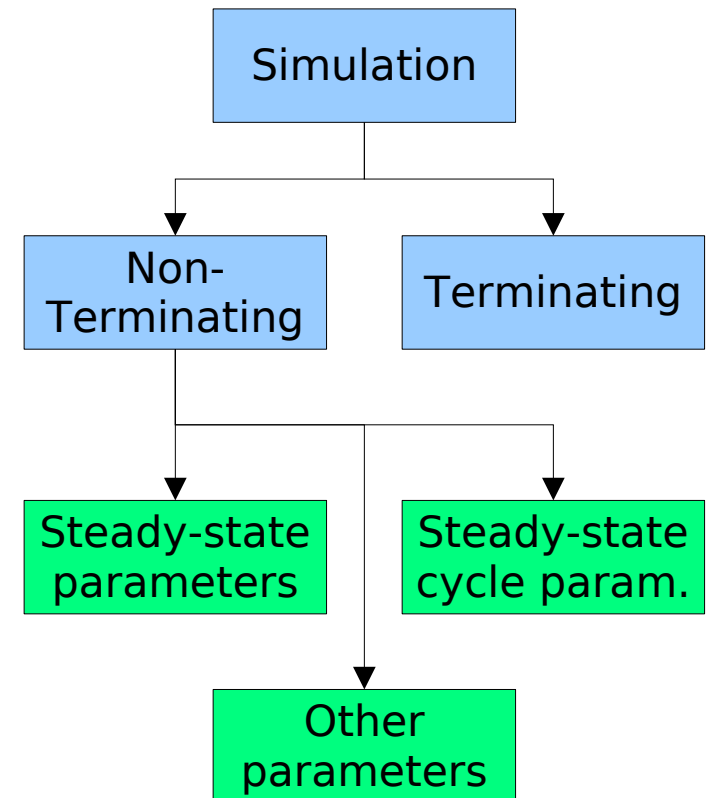
- Types of simulation based on output analyses
- **Terminating** analyses
  - Natural termination event
  - Focus on transient behavior
  - Examples:
    - A day in a retail establishment
    - Manufacturing a batch of items
- **Non-terminating** simulations
  - No termination event
  - Normal (long run) operation setup





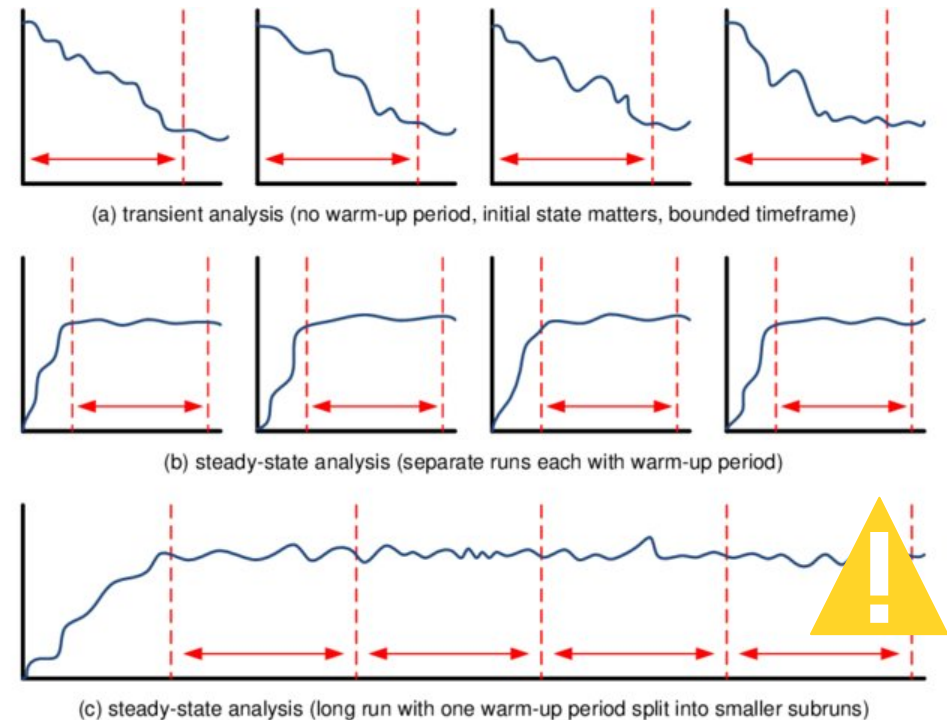
# Types of simulation

- **Steady state** parameters analyses
  - Avoids transient effects
  - We assume system status constant over time
- Steady state **cycle** parameters analyses
  - Focus on **returning patterns**
  - E.g., scaling behavior of cloud over daily patterns
    - Not interested in single requests, **fluid models**
- Other parameters
  - Everything else (es, transient)



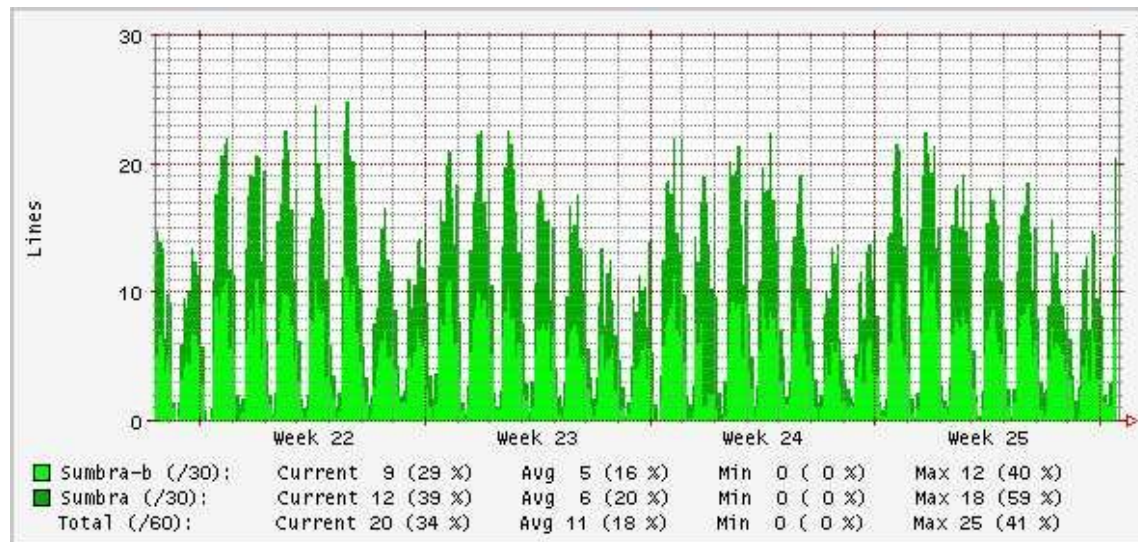
# Steady state analysis

- Some critical elements to remember
- Presence of **initial transient** periods (**cold system**)
  - Reducing impact of transient using **long** simulations
  - **Avoid collecting data** during the initial period (warm-up)
- **Simulation length:**
  - Must collect enough system data
  - Remember **law of large numbers**



# Steady state cycle analysis

- Dual focus:
  - **Variance analysis** (focus on oscillations)
  - **Trend analysis** (focus on long term patterns)
- Use of **filters/ICA** techniques to separate patterns
- **Smoothing** (EWMA) for finding trends



# Variance reduction

---

- Increasing accuracy of simulations
- Some approaches:
- CRN: **Common Random Numbers**
  - When comparing alternative setup use the same seeds
  - Environment as similar as possible between configurations
  - Even RNG settings!
  - Better if we use separate sequences for different parameter

- More complex approaches
- AV: **Antithetic Variates**
  - Use RNG such that time series are negatively correlated
  - Average is more stable
- **Control variates**
  - Usable if we know two variables to be correlated
  - We can use this correction to reduce variance
  - Use correlation to estimate the shift from the real average
  - Needs knowledge of statistical properties of output