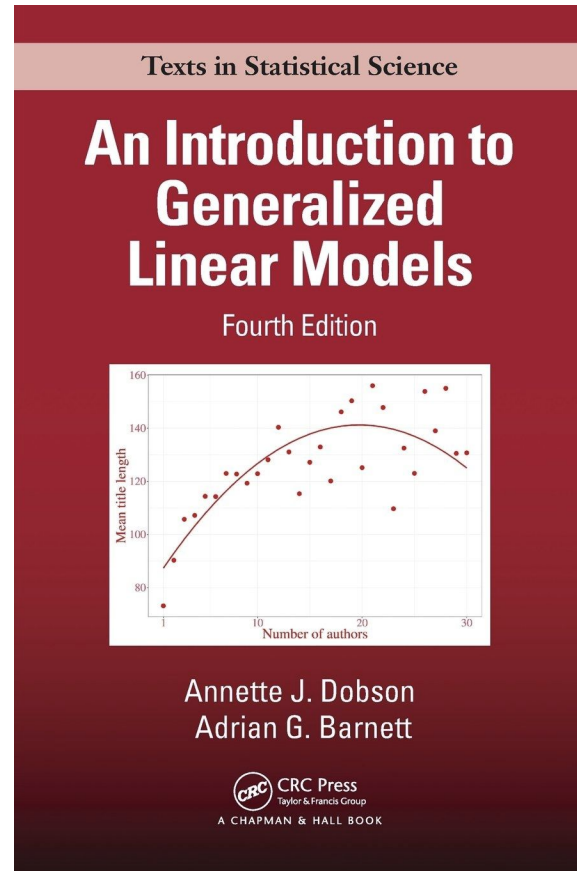


Introduction to Data Science

Distributions

Prof. Dr. Ralf Lämmel & M.Sc. **Johannes Härtel**
(johanneshaertel@uni-koblenz.de)

A book recommendation for this lecture: Especially the first part of this book provides a **formal and concise** introduction to distributions.



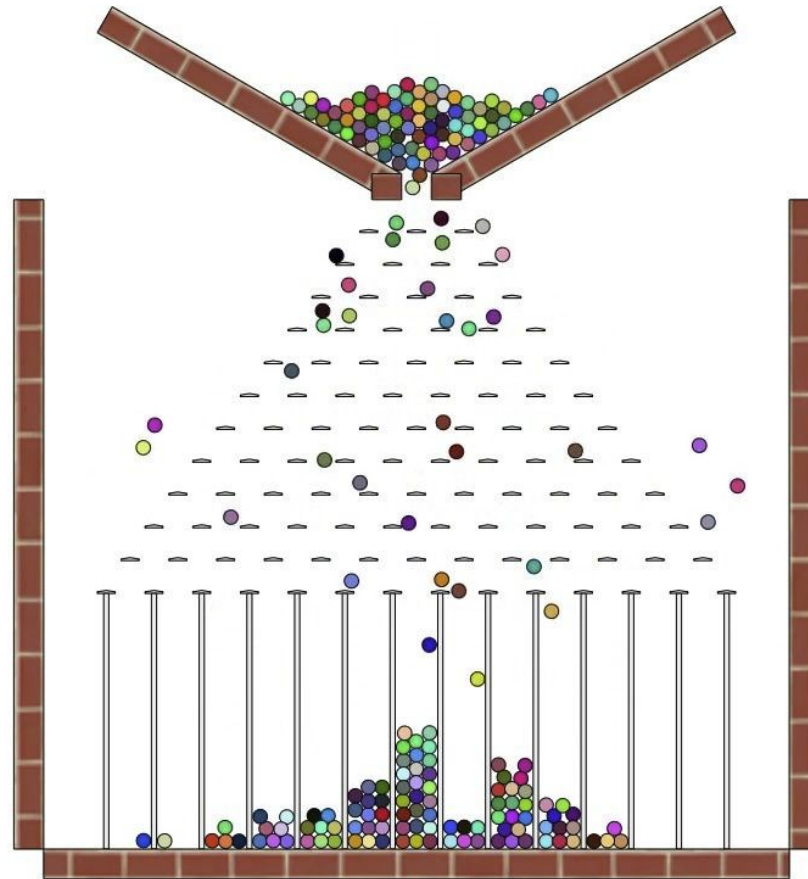
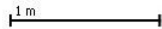
Some distributions in real life



A coin flip: head or number

Binomial Distribution (special case of Bernoulli distribution if just 0 or 1)

[Source](#)

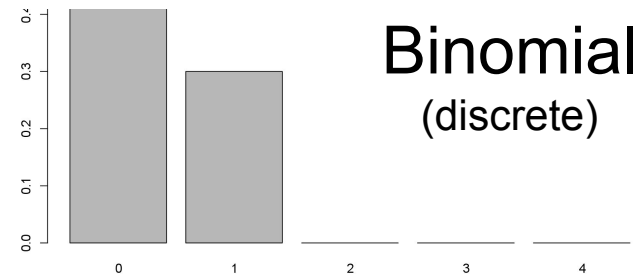
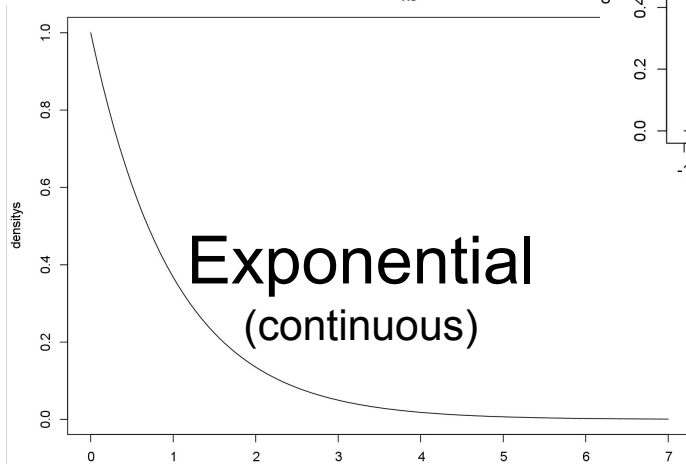
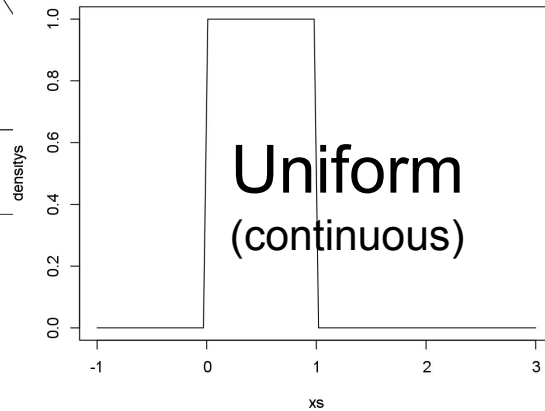
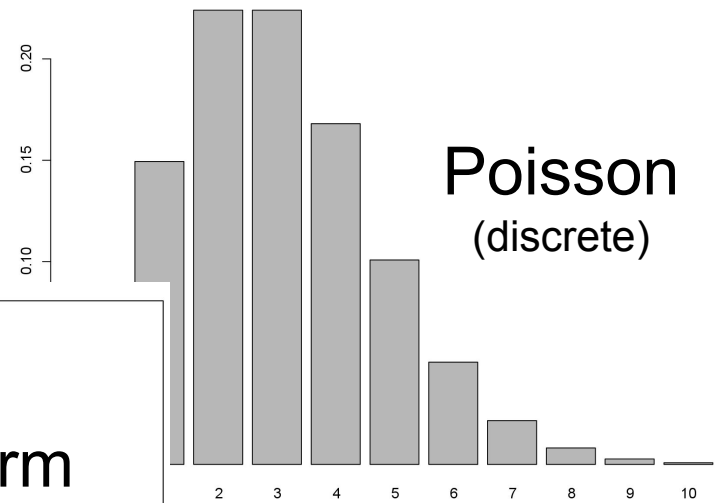
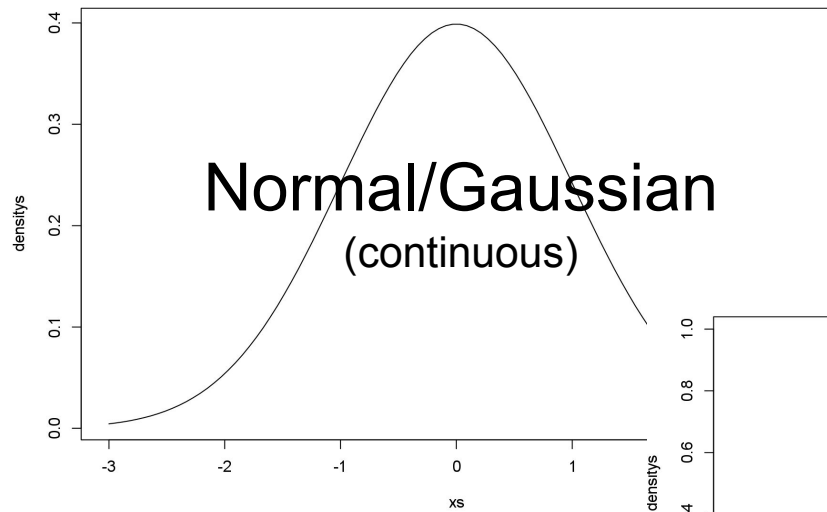


A (discretized) normal distribution

| Filter messages by: Sender Recipient | | |
|--|---------------------------------------|---------------------|
| Subject | Correspondents | Date |
| [Mitarbeiter] Sparkasse, Volksbank, etc.: Phishing-Versuche herausfiltern | • Christoph Litauer via Mitarbeiter | 10/25/2021, 2:44 PM |
| [Mitarbeiter] neue Spam-Phishingwelle, Hinweise | • Konstantin Root via Mitarbeiter | 8/17/2021, 1:31 PM |
| Re: [Mitarbeiter] Spam, Phishing und Konfiguration der Mailinglisten | • Christoph Litauer via Mitarbeiter | 4/11/2021, 11:10 AM |
| [Mitarbeiter] Spam, Phishing und Konfiguration der Mailinglisten | • Christoph Litauer via Mitarbeiter | 4/9/2021, 1:38 PM |
| [Mitarbeiter] Phishing-Mail (?) | • Christoph Litauer via Mitarbeiter | 3/3/2021, 10:53 AM |
| Achtung vor Phishing: So prüfen Sie die Echtheit einer STRATO E-Mail | • STRATO | 2/22/2021, 3:35 PM |
| [Mitarbeiter] Achtung: Sogo-Phishing Mail im Umlauf | • Christoph Litauer via Mitarbeiter | 5/28/2020, 11:46 AM |
| [Mitarbeiter] Hochwertige Spear-Phishing-Angriffe gegen deutsche Universitäten und Forschungseinrichtungen | • Christian Schneider via Mitarbeiter | 12/4/2019, 11:38 AM |
| [Mitarbeiter] Phishing und die Folgen | • Christoph Litauer via Mitarbeiter | 6/4/2019, 4:40 PM |
| [Mitarbeiter] Vorsicht Phishing Mail / Das gleiche wie das letzte mal!!!! | | 6/1/2019, 8:32 PM |
| [Mitarbeiter] Vorsicht Phishing Mail | | 5/29/2019, 5:10 PM |
| [Mitarbeiter] Wieder mal ein PHISHING Mail | | 1/3/2019, 10:52 AM |
| [Mitarbeiter] Phishing/Virus-Welle | | 10/6/2017, 9:18 AM |
| [Studierende] Phishing erkennen und vorbeugen | • Christoph Litauer | 2/10/2016, 1:11 PM |
| [Studierende] Achtung Phishing! | • Christoph Litauer | 7/18/2014, 10:41 AM |
| [Studierende] Neue Phishing Mails | • Uwe Arndt | 7/9/2013, 4:05 PM |

Number of “Phishing mail” mails in each year.

A Poisson distribution



Ingredients of Distributions

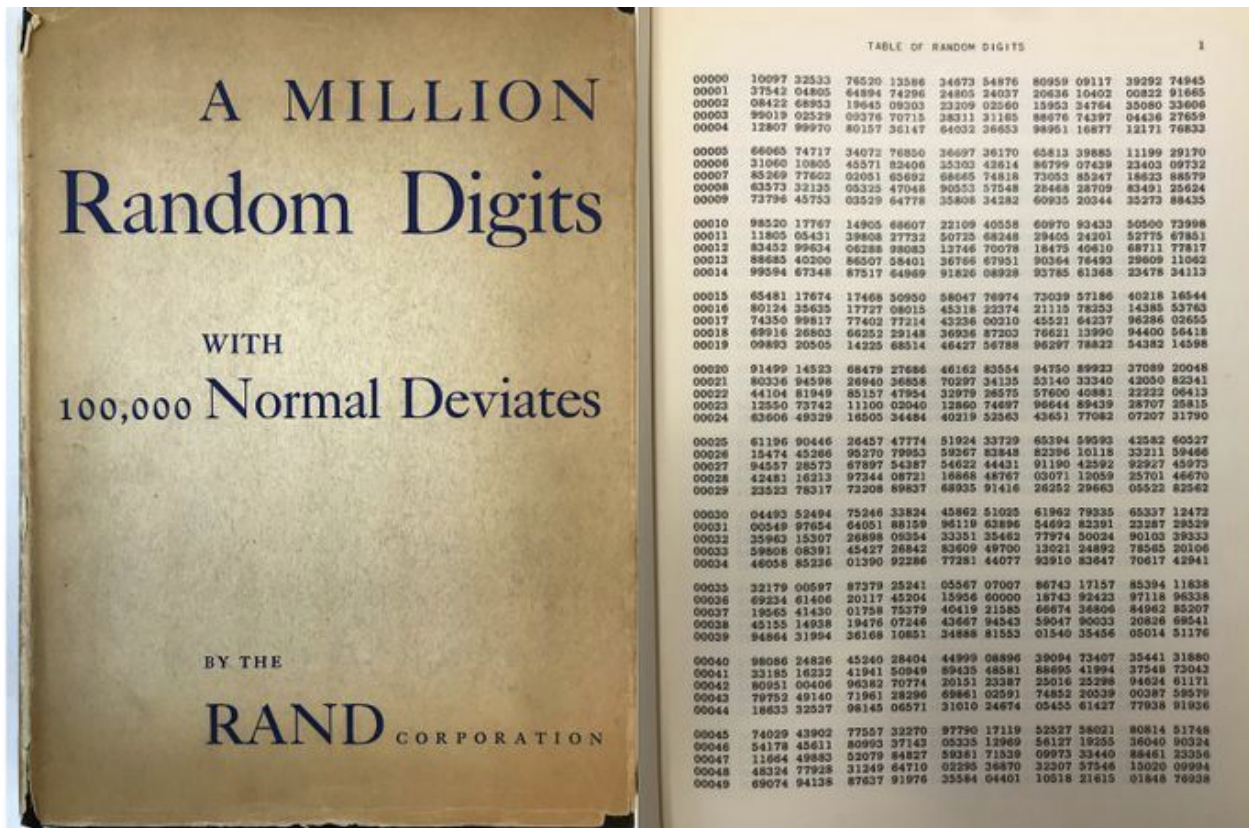
- **Type:** normal, exponential, uniform, Poisson, Binomial (discrete/continuous) ...
- **Parameters:** mean, standard deviation, lambda, rate, prob, min max ...
- **Functions:** (probability) density/mass function (PDF), cumulative distribution function (CDF), quantile function (inverse CDF) and random generation function.





The **random generation** function





Some early work on random values.

We want to generate random values (efficiently).

Management /
Data Base Systems

R. Benjamin
Editor

Computer Methods for Sampling from the Exponential and Normal Distributions

J. H. Ahrens
Nova Scotia Technical College
and
U. Dieter
Universität Karlsruhe

We see that it might also be
called '*sampling from a
distribution*'.

You will find publications on **efficient method**.
For us, it is sufficient to know that such method
exist.

A simple random generation function

This code shows a Linear congruential generator ([Link](#))

```
n <- 2000

# I borrowed these constants from a paper.
m <- 2^32
a <- 1103515245
c <- 12345

# d starts with some sort of seed.
d <- 1000
results <- vector(length = n)
for (i in 1:n) {
  d <- (a * d + c) %% m
  results[i] <- d / m
}
```

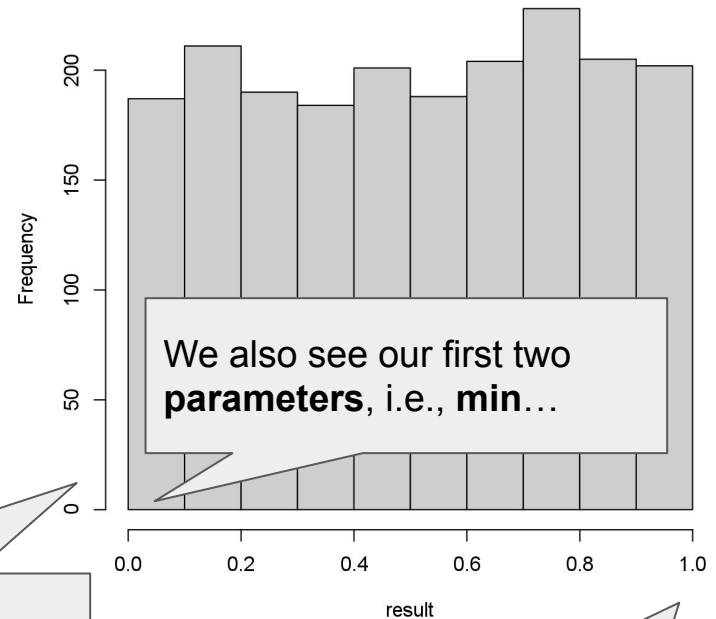
Plotting the results as histogram

```
[1] 0.932167650 0.121204376 0.920264259  
[8] 0.775016561 0.357110500 0.270018101  
[15] 0.815303341 0.601030588 0.734871149  
[22] 0.916099012 0.558927298 0.739606619 ...  
[29] 0.368894354 0.891547620 0.640638113  
[36] 0.144315481 0.578168631 0.047849417  
[43] 0.577669621 0.810473919 0.197372913
```

The generated
random values

⋮

We have been sampling from a
uniform distribution.



... and
max.

Modifying the parameters

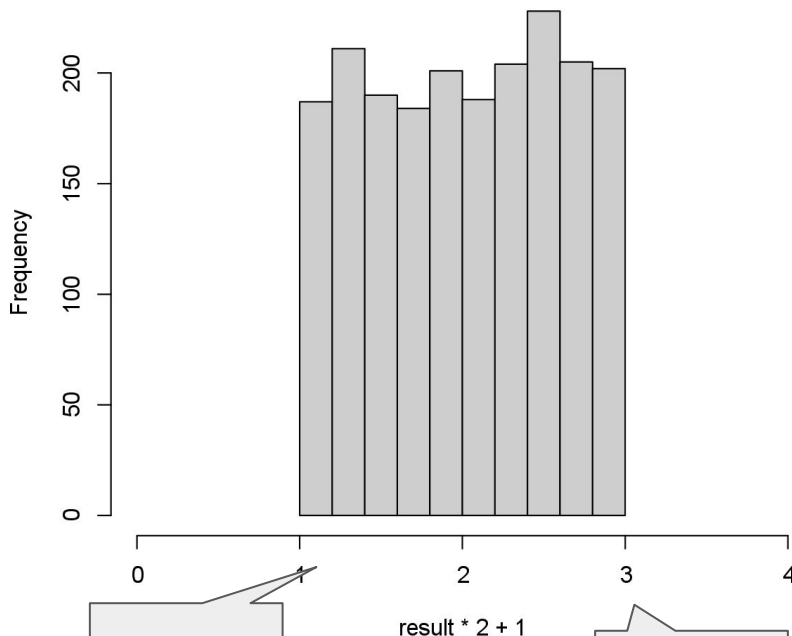
In most case, we want to have different parameters.

```
hist(result * 2 + 1, xlim = c(0, 4))
```

Plotting
histogram
in R

Modification

Limits for
plot x-axis.



new min

new max

R's API

We have comparable packages in python, but R is convenient here ([link](#)).

Bars are added to keep the tension!



R Documentation

Uniform {stats}

The Uniform Distribution

Description

These functions provide information about the uniform distribution on the interval from min to max.

and runif generates random deviates.

Usage

Here it is called
'*generates
random deviates*'
($\cup \square \cup$) \cap \perp

```
runif(n, min = 0, max = 1)
```

Changeling the type of a distribution

We can convert a uniform to a normal distribution.

```
results <- NULL
for (x in 1:1000) {
  parts <- runif(n = 1000, min = 0, max = 1)
  results <- c(results, sum(parts))
}
```

We can employ the [central limit theorem](#), and **sum up random (uniform) variables**.

There are much more efficient way, but we don't care here.

Some normalization to center on 0 and make standard deviation (sd) equals 1.

```
results <- results - mean(results)
```

```
results <- results / sd(results)
```

Plot the data.

```
hist(results)
```

Plotting the results as histogram

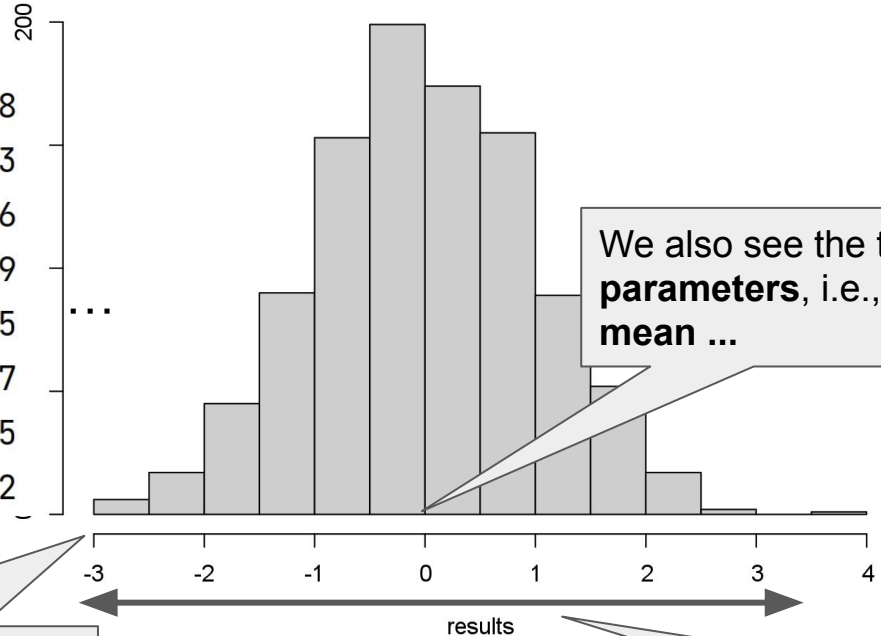
```
[1] 0.450729024 -0.809649934 0.764605888
[7] 1.736888830 0.059764177 -0.576108443
[13] 0.419742595 -0.439459404 -1.177983766
[19] -2.586490148 -0.127243277 0.475076599
[25] 1.709972580 0.202269372 -0.471295545
[31] 0.612078333 0.998179839 -1.227958907
[37] -0.096025296 1.528426825 -0.065776675
[43] 0.243227469 -0.489929905 0.012433712
```

The random values

⋮

We have sampled from a **normal distribution**.

Histogram of results



We also see the two **parameters**, i.e., **mean** ...

... and **standard deviation**.


R's API

Normal {stats}

R Documentation

The Normal Distribution

Description

 random generation for the normal distribution with mean equal to mean and standard deviation equal to sd.

Usage



```
rnorm(n, mean = 0, sd = 1)
```

Here it is called
random generation
for...
(°□°) ^ _

Random number generation for other types of distributions

- Binomial distribution: `rbinom(n = ..., size = ..., prob = ...)`
- Poisson distribution: `rpois(n = ..., lambda = ...)`
- Exponential distribution: `rexp(n = ..., rate = ...)`
- ...

Short Demo

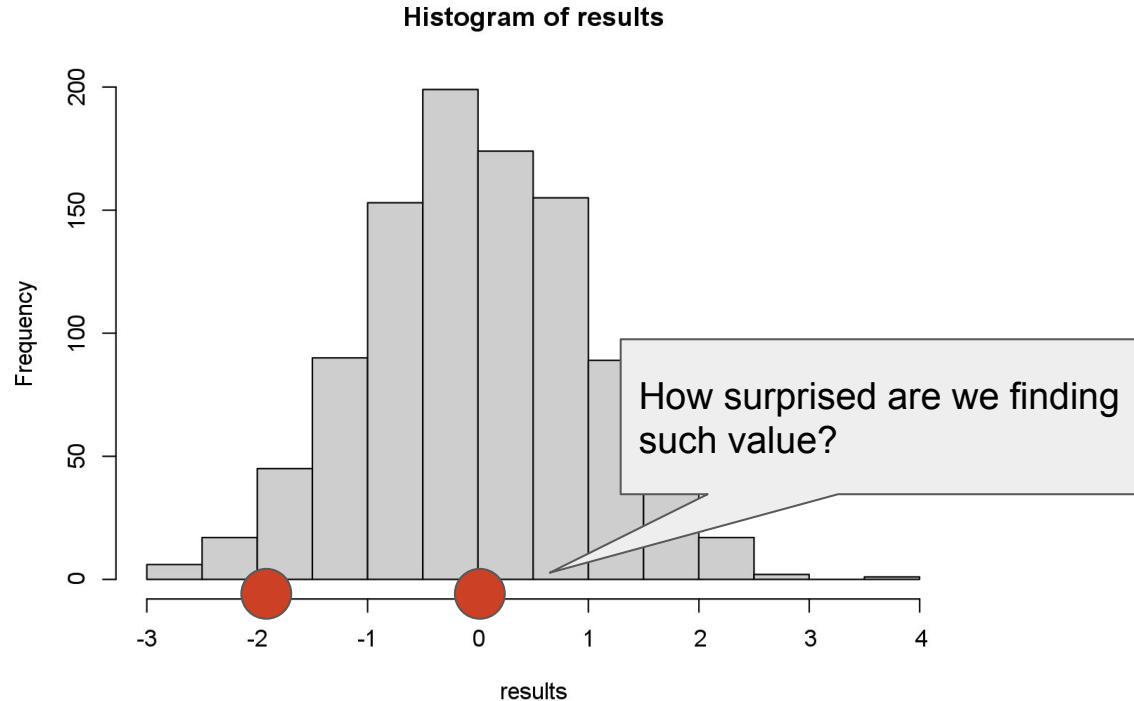
(The best way for understanding distributions is using them)



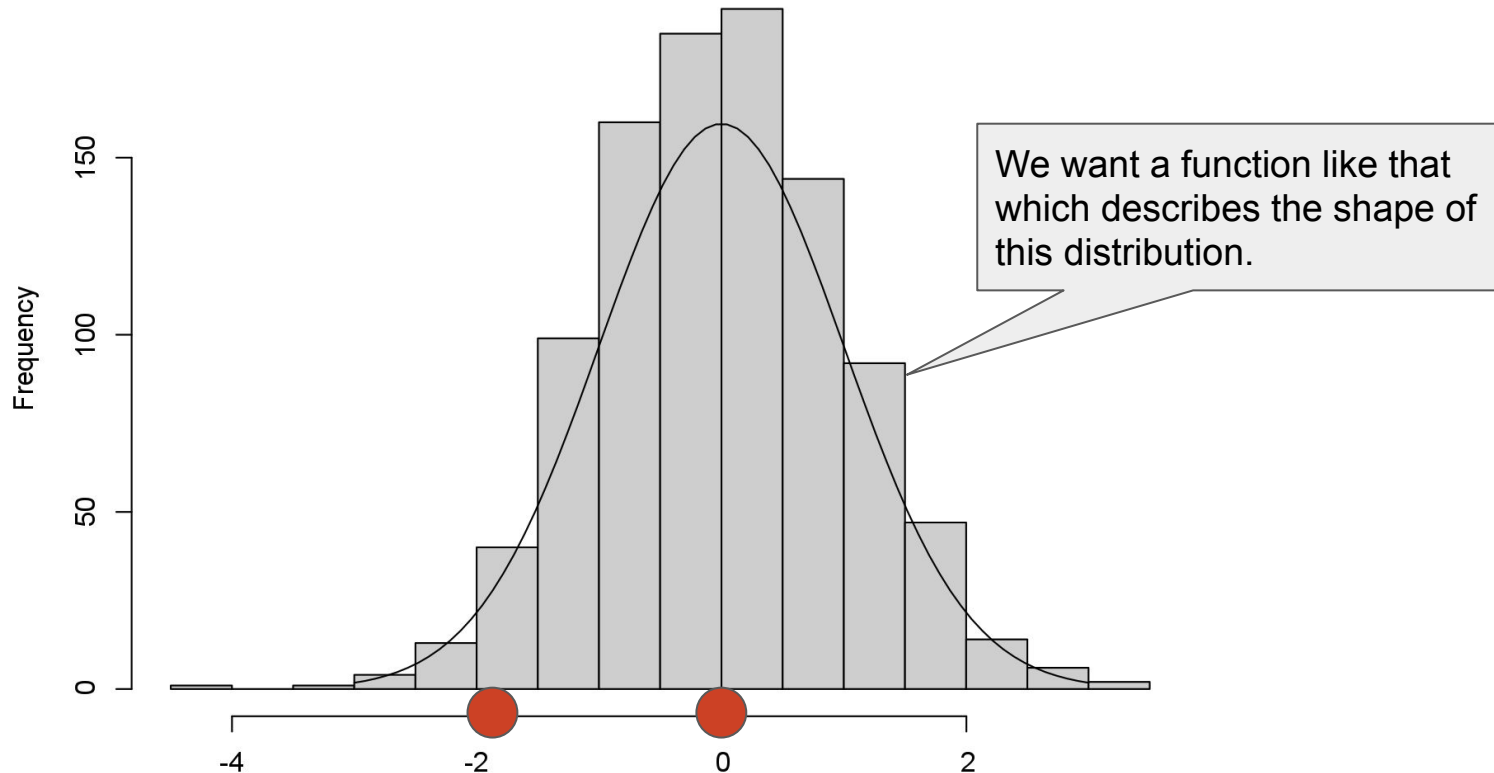
The **(probability) density** function (PDF)



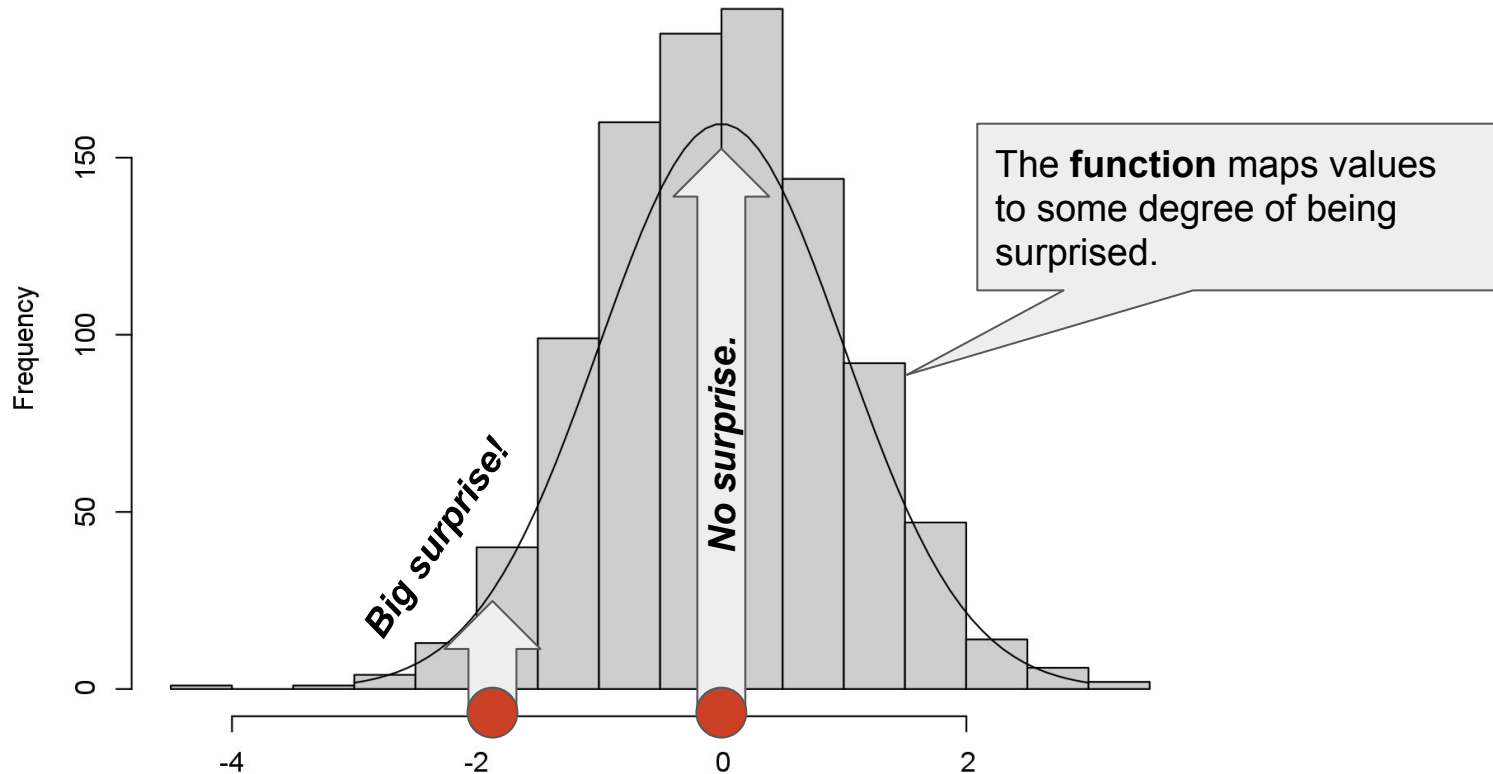
We need the degree of being “surprised” by a value



We need the degree of being “surprised” by a value (cont)



We need the degree of being “surprised” by a value (cont)



The (probability) density function (short PDF)

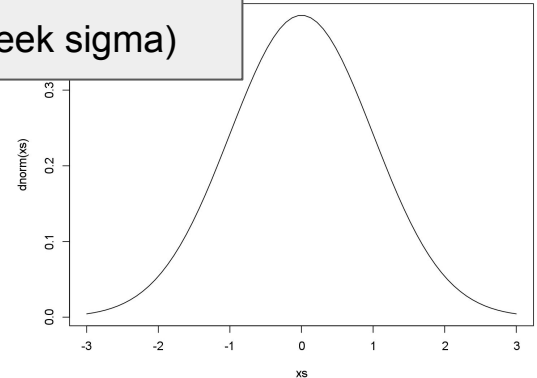
In this case, for the normal distribution (see, we have parameters again)

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

μ (greek mu) is the mean

σ (greek sigma)

Parameters are fixed again:
 σ (greek sigma) refers to the
standard deviation.



R's API


We will never implement this by hand.

Normal {stats}

The Normal Distribution
Description

Just called density, but it's the probability density function (PDF).

R Documentation

Density,  and random generation for the normal distribution with mean equal to mean and standard deviation equal to sd.

Usage

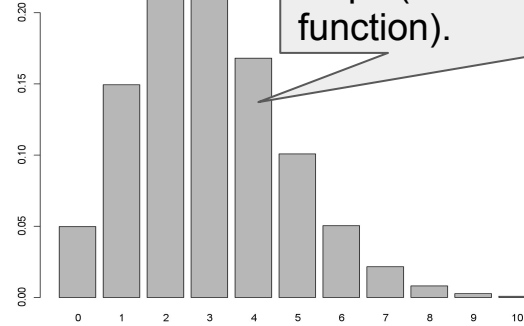
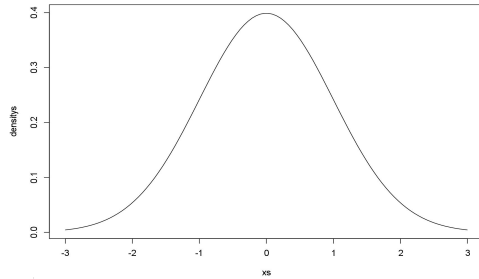
```
dnorm(x, mean = 0, sd = 1, log = FALSE)
```



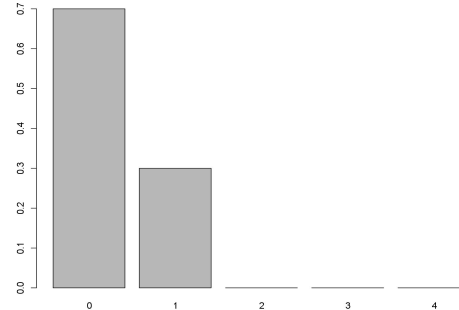
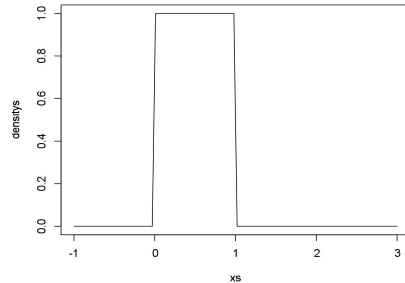
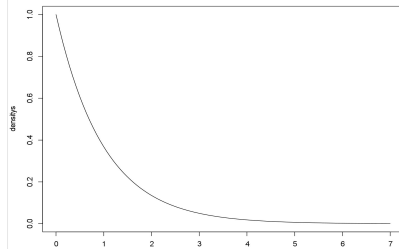
```
rnorm(n, mean = 0, sd = 1)
```

Probability density functions (PDF) for other types

You might start to notice the symmetry (and where to find this functions).



For discrete distributions, we have steps (and we say probability mass function).



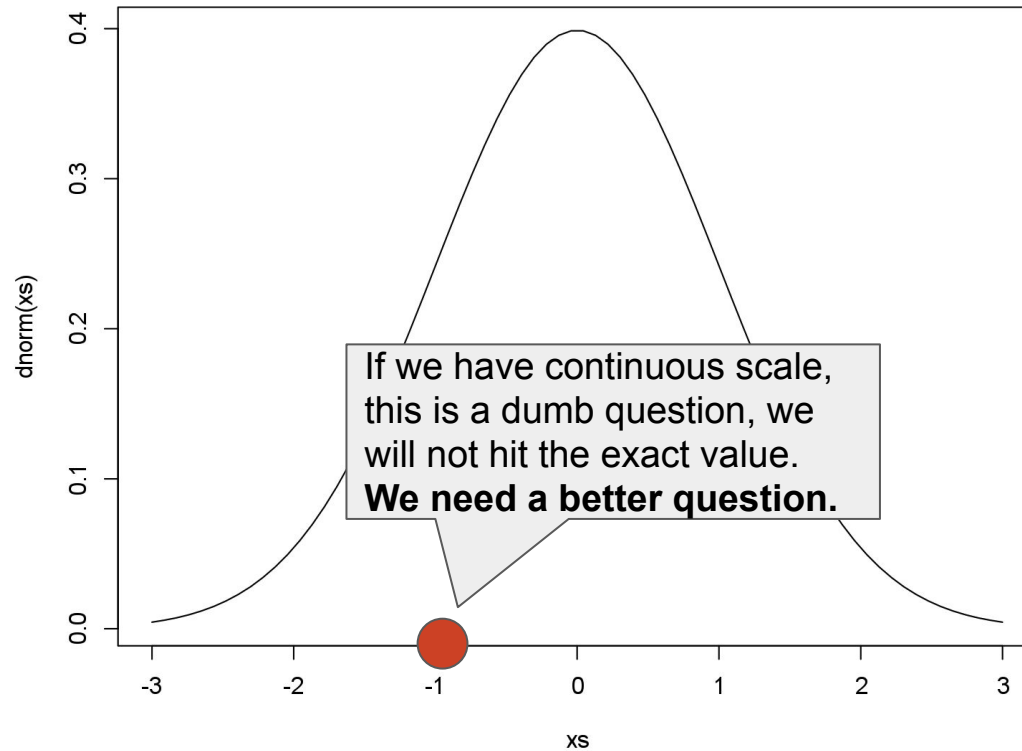
Short Demo



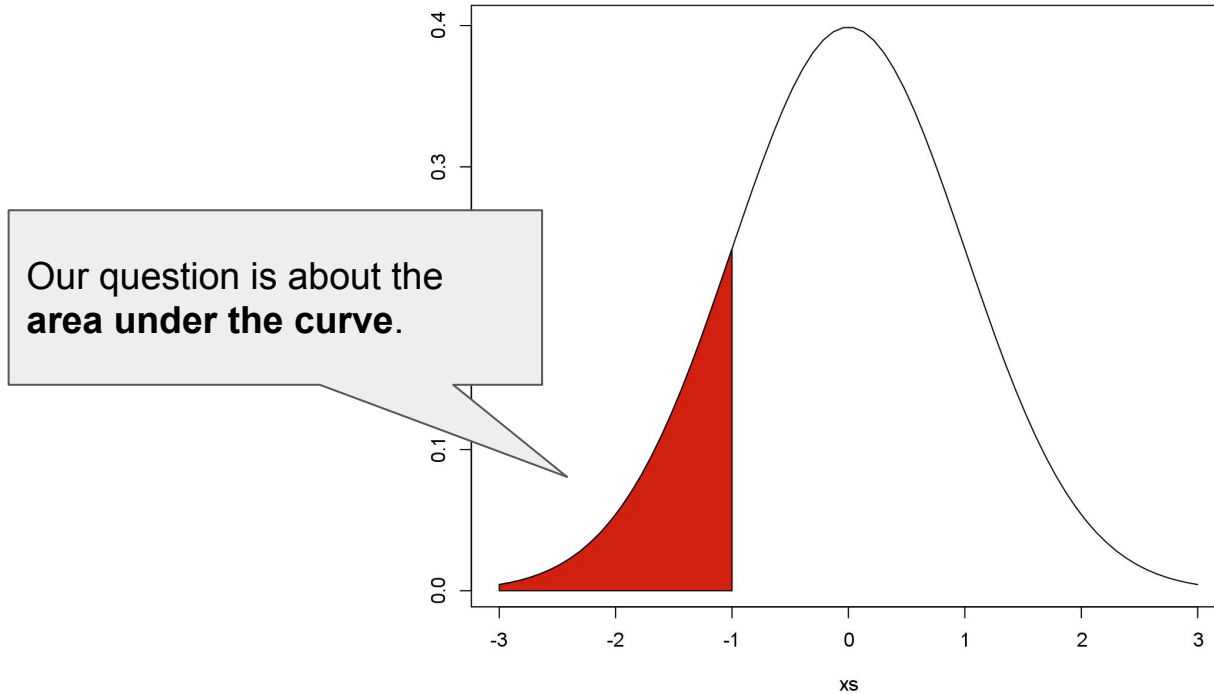
The **cumulative distribution** function (CDF)



What is the probability of facing an exact value?



What is the probability of facing a value within a certain range (e.g., $X \leq -1$)?



The cumulative distribution function (CDF)

for a normal distribution

new

PDF (in this example of normal distribution)

No **closed form solution** for this case; we need to approximate it numerically.

μ (greek mu) is the mean

$$F_X(x) = \int_{-\infty}^x \frac{1}{\sigma\sqrt{2\pi}} e^{-\frac{1}{2}\left(\frac{x-\mu}{\sigma}\right)^2}$$

σ (greek sigma)

If you have a discrete distribution, the integral can be replaced by a sum.

Parameters are fixed again: (greek sigma) refers to the standard deviation.

Since we may have no closed form solution, we approximate it.
Seen an example for the CDF of a normal distribution

Checking the cumulative distribution function (CDF) for x smaller or equal to 0.

```
x <- 0
```

Just generate 2000 random values and then do simple counting.

```
values <- rnorm(n = 2000, sd = 1, mean = 0)
```

```
print(sum(values < x) / 2000)
```

≈ 0.5

R can sum up boolean values
(True = 1, False = 0)

Exploring the CDF with respect to X

Plotting the cumulative distribution function (CDF) for x smaller or equal.

```
values <- rnorm(n = 2000, sd = 1, mean = 0)
```

```
ys <- NULL
```

```
xs <- seq(-3, 3, length.out = 100)
```

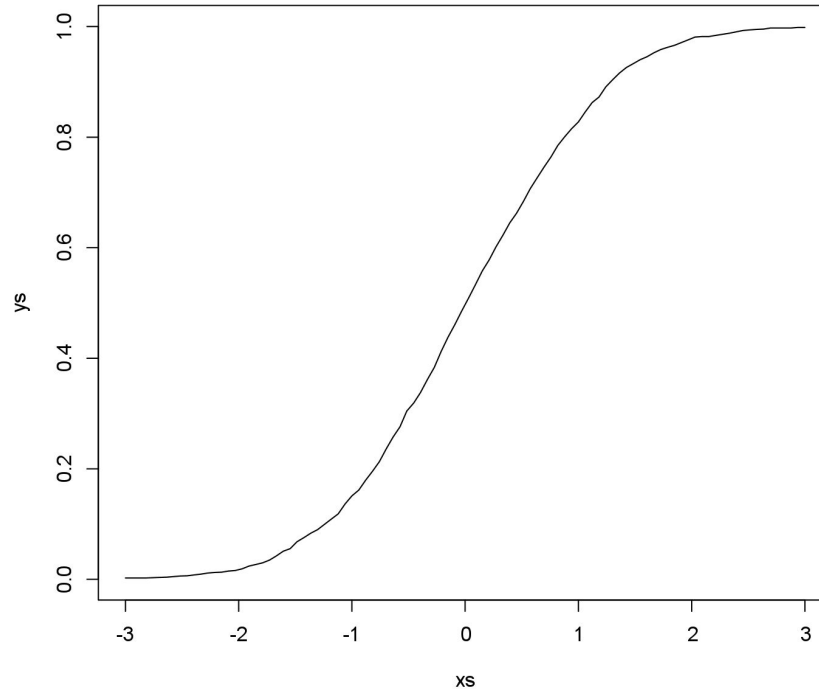
```
for (x in xs) {
```

```
  ys <- c(ys, sum(values < x) / 2000)
```

```
}
```

```
plot(xs, ys, type = "l")
```

The classical plot of the cumulative distribution function (CDF) for the normal distribution




R's API

Normal {stats}

R Documentation

The Normal Distribution

Description

Density, distribution function  and random generation for the normal distribution with mean equal to mean and standard deviation equal to sd.

Usage

```
dnorm(x, mean = 0, sd = 1, log = FALSE)
```

```
pnorm(q, mean = 0, sd = 1, lower.tail = TRUE, log.p = FALSE)
```

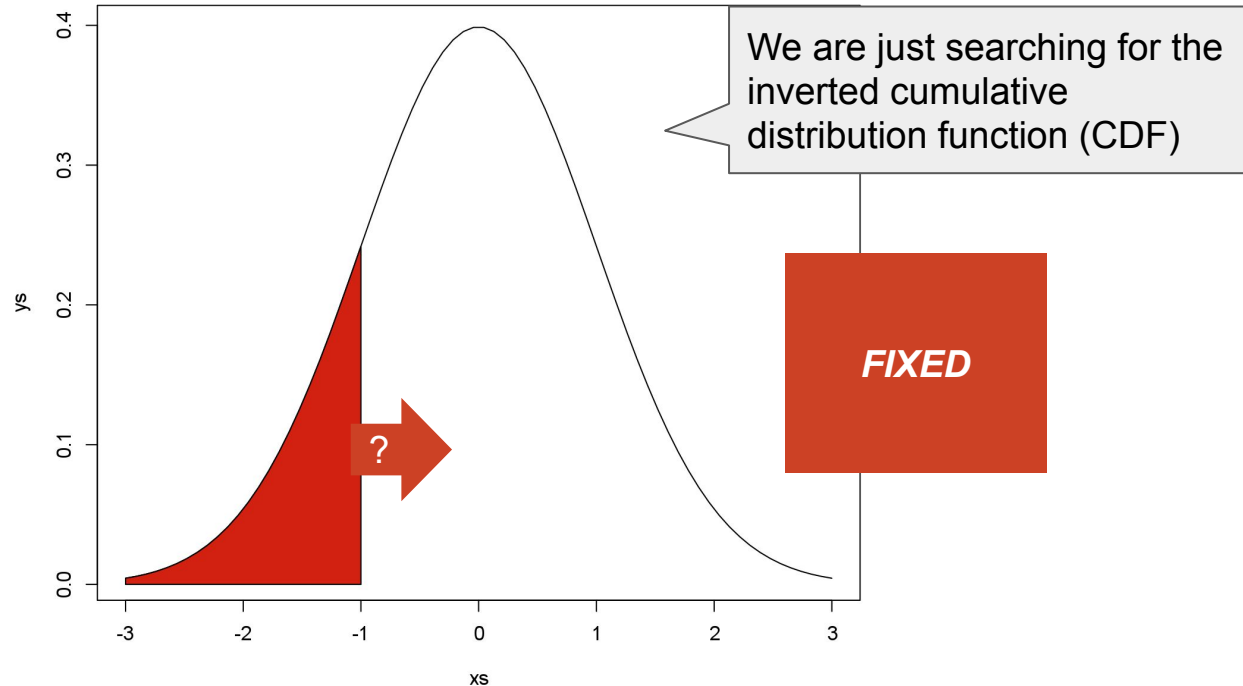
```
  
rnorm(n, mean = 0, sd = 1)
```



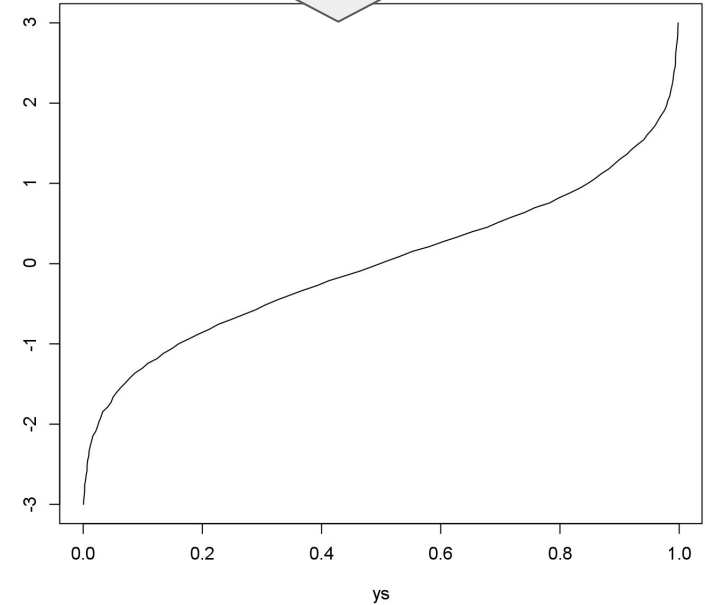
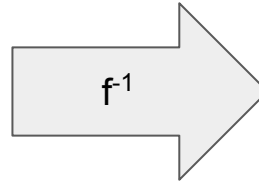
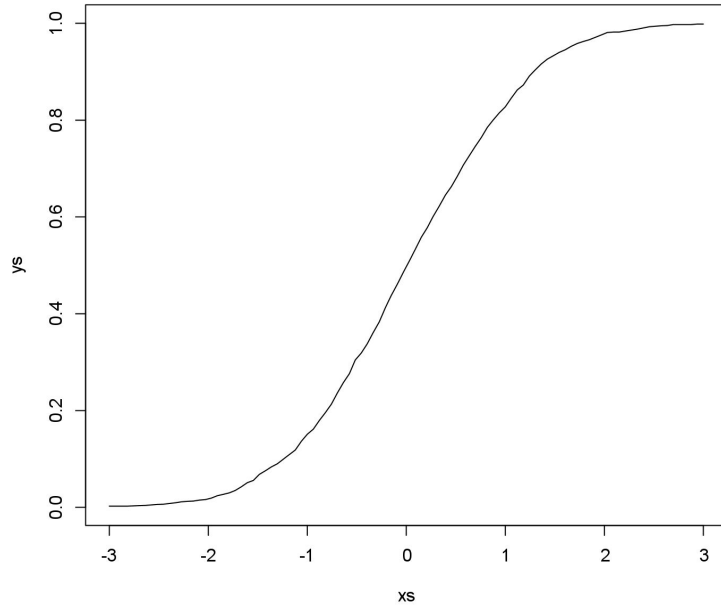
The **quantile** function



Requirement: Given the area, find the boundary:



Just inverting the CDF



BTW: x and y are
just switched in the
plot. `plot(ys,xs,...)`

R's API

Normal {stats}

R Documentation

The Normal Distribution

Description

Density, distribution function, quantile function and random generation for the normal distribution with mean equal to mean and standard deviation equal to sd.

Usage

```
dnorm(x, mean = 0, sd = 1, log = FALSE)
```

```
pnorm(q, mean = 0, sd = 1, lower.tail = TRUE, log.p = FALSE)
```

```
qnorm(p, mean = 0, sd = 1, lower.tail = TRUE, log.p = FALSE)
```

```
rnorm(n, mean = 0, sd = 1)
```

Example discrete distribution (binomial)

Binomial {stats}

R Documentation

The Binomial Distribution

Description

Density, distribution function, quantile function and random generation for the binomial distribution with parameters `size` and `prob`.

This is conventionally interpreted as the number of 'successes' in `size` trials.

Usage

```
dbinom(x, size, prob, log = FALSE)
```

```
pbinom(q, size, prob, lower.tail = TRUE, log.p = FALSE)
```

```
qbinom(p, size, prob, lower.tail = TRUE, log.p = FALSE)
```

```
rbinom(n, size, prob)
```

Summary

- We have covered distributions in terms of **types, parameters and some interesting function.**
- What we are missing is a way to **come from fixed data to unknown parameters.** We will cover this next week.

We will meet all these distributions again in the remainder of this course.

