

STATISTICS

A SHORT AND PAINLESS INTRODUCTION

Nils M Holm

Contents

Preface	9
Probability	10
Basics	10
Conditional Probability	13
Permutations and Combinations	17
Probability Functions	20
Discrete Probability Distributions	22
Location	24
Variation	27
Discrete Distributions Revisited	31
Continuous Distributions	33
Quantiles	37
Central Limit Theorem	39
Statistics	43
Point Estimators	43
Sampling Distributions	44
Confidence Intervals	46
Regression and Correlation	50
<i>Covariance</i>	<i>53</i>
<i>Standard Error of the Estimate</i>	<i>55</i>
<i>Final Notes</i>	<i>56</i>
Probability Distributions	58
Uniform Distribution	58
Geometric Distribution	62
Binomial Distribution	66
Hypergeometric Distribution	70
Poisson Distribution	74
Normal Distribution	78
Standard Normal Distribution	82
Chi-Square Distribution	86
Student's t-Distribution	90

Further Applications	95
Hypothesis Testing	95
Contingency Tables	98
Special Functions	104
Gauss Error Function and Normal CDF	104
Gamma Function	104
Beta Function	108
t-Distribution CDF	109
Quantile Functions	110
Probability Tables	112
Mathematical Notation	116
Bibliography	119
Index	120

Preface

Statistics is a fascinating subject but, unfortunately, there is little literature that teaches it in a way that is accessible to mathematical laymen while still being precise and straight to the point.

On the one side, there are math textbooks that provide the interested reader with rigorous proof of every little detail, which may not be necessary for all students of humanities or natural sciences. On the other side, there are lengthy works that use entertaining language and cute comic figures to try to get their point across.

This book has neither rigorous proofs nor cute characters. It strives to build a solid foundation for people who use statistics as a tool. After finishing this book, the reader should be able to digest more complete works in the field of statistics without too much difficulty.

Topics covered in this brief volume include basic probability, probability functions and distributions, confidence intervals, linear regression, correlation, and hypothesis testing.

The final chapters describe numerical methods for computing statistical functions. They are optional for most students, but may be of interest to the mathematically inclined reader.

There are no exercises, but at some points questions are asked and at those points the reader is invited to put aside the book and try to find their own solution before reading on.

The matter of the book progresses rather quickly, so the reader is advised not to skip sentences or paragraphs. Doing so would probably (!) complicate the comprehension of later parts of the text.

Enjoy the tour through the foundations of statistics!

Nils M Holm, July 2016

Probability

Basics

Probability is an *estimate* predicting how often an *event* will occur given a fixed number of *trials*. For example, when a coin is tossed 10 times, “heads” will probably show up *about* 5 times. So the probability of “heads” is 50% or

$$p = 0.5$$

In statistics, probabilities are expressed as a real number $p \in \mathbf{R}$ where $0 \leq p \leq 1$, i.e. p is in the interval $[0, 1]$. *Impossibility* (an event will never occur) is denoted by $p = 0$ and *certainty* (an event will always occur) is represented by $p = 1$.

Note that getting 7 “heads” and 3 “tails” when tossing a coin 10 times does not violate the prediction of getting “heads” half of the time! The probability of $p = 0.5$ for getting “heads” is only the *most probable outcome* of tossing a coin repeatedly.

As the number of trials (coin tosses) increases, the actual distribution of heads and tails will converge towards the expectation of $p = 0.5$ for “heads”.

When tossing two coins at the same time, there are four possible *outcomes* (H indicates “heads” and T indicates “tails”):

HH HT TH TT

Each outcome has the same probability of $p = 0.25$. This can be expressed using the *probability function* P as follows:

$$P(HH) = 0.25$$

$$P(HT) = 0.25$$

$$P(TH) = 0.25$$

$$P(TT) = 0.25$$

Meaning: the probability of two times “heads” is $p = 0.25$, etc.

The set of all possible outcomes is called the *sample space* S . In the above example, this would be

$$S = \{HH, HT, TH, TT\}$$

The probability of S is $P(S) = 1$, because one of the above events *must* occur in every trial (we are assuming an “ideal” coin that cannot get stuck on its edge or disappear in a storm drain).

So if $S = A_1 \cup \dots \cup A_n$, then $P(A_1) + \dots + P(A_n) = P(S) = 1$, where $A \cup B$ denotes the *union*, or logical “or”, of two events, i.e. either A or B or both A and B happens.

The notation A' (sometimes also \bar{A}) is called the *complement* of A . It indicates that an event A does *not* occur. The probability $P(A')$ is $1 - P(A)$ for any A . For example, the probability of not getting two times heads when tossing two coins would be

$$P(\overline{HH}) = 1 - P(HH) = 1 - 0.25 = 0.75$$

and the probability of none of the events in the sample space happening would be $P(S') = 1 - 1 = 0$.

When drawing cards from a standard deck of 32 cards, the probability of drawing an “ace” would be $P(A) = \frac{4}{32} = \frac{1}{8}$, because there are 4 aces in the standard deck. The probability of drawing a red card from the deck would be $P(B) = \frac{16}{32} = \frac{1}{2}$, because there are 16 red and 16 black cards in the deck. The sample space in this case would be the entire deck.

The probability of drawing a “red ace” would be

$$P(A \cap B) = P(A) \cdot P(B) = \frac{1}{8} \cdot \frac{1}{2} = \frac{1}{16} = 0.0625$$

Meaning: the probability of A and B happening *at the same time* (i.e. in the same trial) is 0.0625.

The probability of the *intersection*, or logical “and”, $P(A \cap B)$ of two *independent* events A and B is calculated by multiplying their probabilities. More on this later (pg 14).

The probability of drawing an “ace” *or* a red card is:

$$P(A \cup B) = P(A) + P(B) - P(A \cap B) = \frac{1}{8} + \frac{1}{2} - \frac{1}{16} = 0.5625$$

Meaning: the probability of *either A or B or $A \cap B$* happening in the same trial is 0.5625.

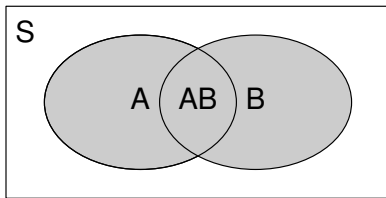
When calculating the probability of the *union $P(A \cup B)$* of two events *A* and *B*, the probability of the intersection $P(A \cap B)$ has to be subtracted, because otherwise it would be duplicated (if it is non-zero).

Of course, if the intersection of *A* and *B* is empty, its probability does not have to be subtracted, so $P(A \cup B) = P(A) + P(B)$, *iff* *A* and *B* are mutually exclusive, i.e. *iff* the events *A* and *B* *cannot* occur in the same trial.

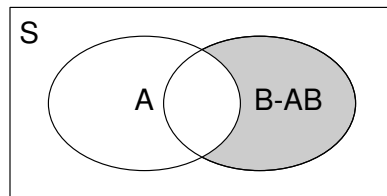
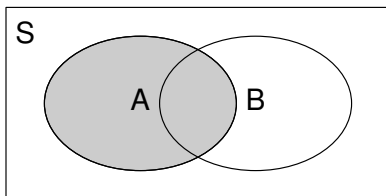
(Note: “iff” is a common abbreviation for the bidirectional “if”, i.e. “if and only if”.)

In the above example, $P(A)$ (aces) includes two red cards and $P(B)$ (red cards) includes two aces, giving an “overlap” of 4 cards. However, there are only 2 red aces in the deck, so half of the overlap has to be eliminated by subtracting $P(A \cap B)$.

This is probably best demonstrated using a Venn diagram (each ellipse denotes an event and the overlap of ellipses denotes the intersection of events):



Both the event *A* and *B* would contribute to the intersection *AB*, thereby duplicating it, so one of the sets (ellipses) has to lose its intersection part before adding it:



Summary

$0 \leq p \leq 1$ for any probability p .

$0 \leq P(A) \leq 1$ for any event A .

$P(A') = 1 - P(A)$ for any event A .

$P(S) = 1$ for any sample space S .

If $S = \bigcup_{i=1}^n A_i$, then $P(S) = \sum_{i=1}^n P(A_i) = 1$.

$P(A \cap B) = P(A) \cdot P(B)$, iff A and B are independent.

$P(A \cup B) = P(A) + P(B) - P(A \cap B)$.

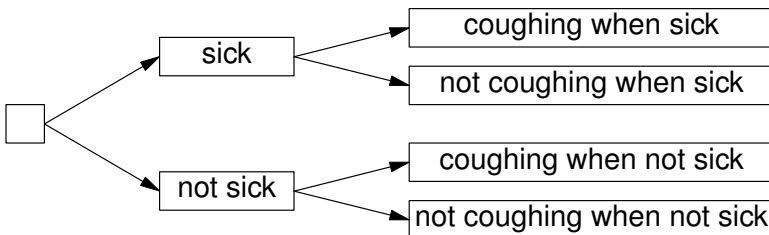
$P(A \cup B) = P(A) + P(B)$, iff A and B are mutually exclusive.

Conditional Probability

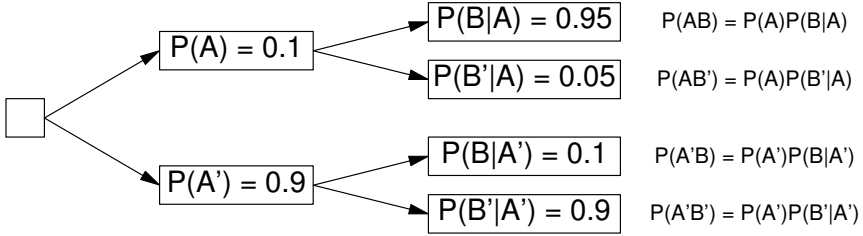
Imagine it is flu season and

- the probability of a random person having a cold is $p = 0.1$
- the probability of a person coughing while having a cold is $p = 0.95$ (5% may not cough and still be sick)
- the probability of a person not coughing while not having a cold is $p = 0.9$ (10% might cough for different reasons)

The following tree diagram can be constructed from this scenario:



The notation $P(B|A)$ denotes the *conditional probability* of “ B given A ”, i.e. the probability of B in the case where it is already known that A has happened. Using this notation, the diagram can be populated with probabilities:



For example, $P(A)$ denotes the probability of a person being sick, $P(B|A)$ denotes the probability of a person coughing *given* they are sick, and $P(AB) = P(A \cap B) = P(A) \cdot P(B|A)$ is the probability of a person coughing *and* being sick ($0.1 \cdot 0.95 = 0.095$).

Note that the probability $P(A \cap B)$ is given as $P(A) \cdot P(B|A)$ here, while it was given as $P(A) \cdot P(B)$ earlier in this text (pg 11). In the example given here, A and B are *dependent*, because $P(B) \neq P(B|A)$.

In fact, two events A and B are *independent* if, and only if, $P(B) = P(B|A)$. That is, $P(B)$ is the same, no matter whether A has is given or not.

The probability $P(B)$ of a person coughing, although not explicitly stated in the data, can be inferred from the diagram. It is the combined probability of a person coughing, no matter if they have a cold or not:

$$\begin{aligned}
 P(B) &= P(B \cap A) + P(B \cap A') \\
 &= P(A) \cdot P(B|A) + P(A') \cdot P(B|A') \\
 &= 0.1 \cdot 0.95 + 0.9 \cdot 0.1 \\
 &= 0.185
 \end{aligned}$$

Note that $P(B \cap A)$ and $P(B \cap A')$ are mutually exclusive, because A and A' cannot happen in the same trial. (E.g. a person cannot be sick and not sick at the same time.)

Reverse Conditional Probability

A more interesting question in the flu season might be: “Given that someone is coughing, what is the probability that they have a cold?” I.e.: what is $P(A|B)$?

The probability of a person coughing *because* they have a cold equals the *proportion* of people who cough while having a cold and those who cough at all (for whatever reason):

$$P(A|B) = \frac{\text{coughing and sick}}{\text{coughing}} = \frac{P(A \cap B)}{P(B)}$$

Inserting the formulae for $P(A \cap B)$ and $P(B)$ from above gives:

$$\frac{P(A) \cdot P(B|A)}{P(A) \cdot P(B|A) + P(A') \cdot P(B|A')}$$

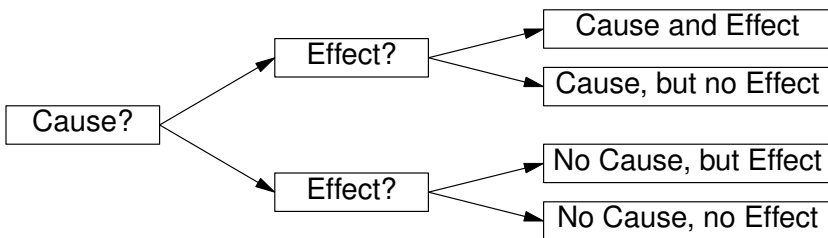
All probabilities that appear in this formula can be extracted from the tree diagram. The above formula is widely known as *Bayes' Theorem* or *Bayes' Rule*.

Substituting values for probability functions finally gives:

$$\frac{P(A) \cdot P(B|A)}{P(A) \cdot P(B|A) + P(A') \cdot P(B|A')} = \frac{0.1 \cdot 0.95}{0.1 \cdot 0.95 + 0.9 \cdot 0.1} = \frac{19}{37}$$

The probability of someone coughing because they have a cold in said flu season is just $p = 0.514$ — so results from the above test would only be slightly more significant than tossing a coin.

The basic idea behind *reverse conditional probability* is as follows:



We may observe an *effect*, and that effect may or may not have a specific *cause*. Given the probability of “effect given cause” ($P(B|A)$) and the probability of “no effect given no cause” ($P(B'|A')$) as well as the probability of the cause in general ($P(A)$), what is the probability of “cause given effect”? I.e. what is the probability that the observation of the effect was triggered by the specific cause? This question is answered by Bayes' Theorem.

When using reverse conditional probability (RCP) to evaluate test results, $P(B|A)$ is called the *sensitivity* of the test and $P(B'|A')$ is called its *specificity*. A test is *sensitive*, if it catches a lot of positives (i.e. has few false negatives). It is *specific*, if it catches few negatives (i.e. has few false positives).

A *false positive* occurs when a test wrongly delivers a positive result. Analogously, a *false negative* is a wrong negative result.

The above “cold test” has a sensitivity of $p = 0.95$ and a specificity of $p = 0.9$. This sounds good, but the reliability of the test is low, because a positive only indicates a $p = 0.514$ chance for the specific cause.

RCP depends a lot on the *prior probability* (or just “prior”), i.e. the probability of the cause in general, $P(A)$. When $P(A) = 1$, then

$$\frac{P(A) \cdot P(B|A)}{P(A) \cdot P(B|A) + P(A') \cdot P(B|A')} = \frac{1 \cdot P(B|A)}{1 \cdot P(B|A) + 0 \cdot P(B|A')} = \frac{P(B|A)}{P(B|A)}$$

So the test result is always $p = 1$. Similarly, when the prior is $P(A) = 0$, the test will always be negative.

Above cold test would not be as bad, if the *prevalence* (the medical term for the prior) was higher, i.e. if more people had a cold in the first place. Given $P(A) = 0.5$:

$$\frac{P(A) \cdot P(B|A)}{P(A) \cdot P(B|A) + P(A') \cdot P(B|A')} = \frac{0.5 \cdot 0.95}{0.5 \cdot 0.95 + 0.5 \cdot 0.1} = \frac{19}{21}$$

So given a prevalence of 50%, a coughing person would indicate a cold in about 90% of the observed cases.

The smaller the prior is, the greater the sensitivity and the specificity of a test have to be in order for the test result to be significant.

Summary

$P(A \cap B) = P(A) \cdot P(B|A)$ if A and B are dependent.

$P(B) = P(B|A)$ iff A and B are independent.

$P(B) = P(A) \cdot P(B|A) + P(A') \cdot P(B|A')$

Probability Distributions

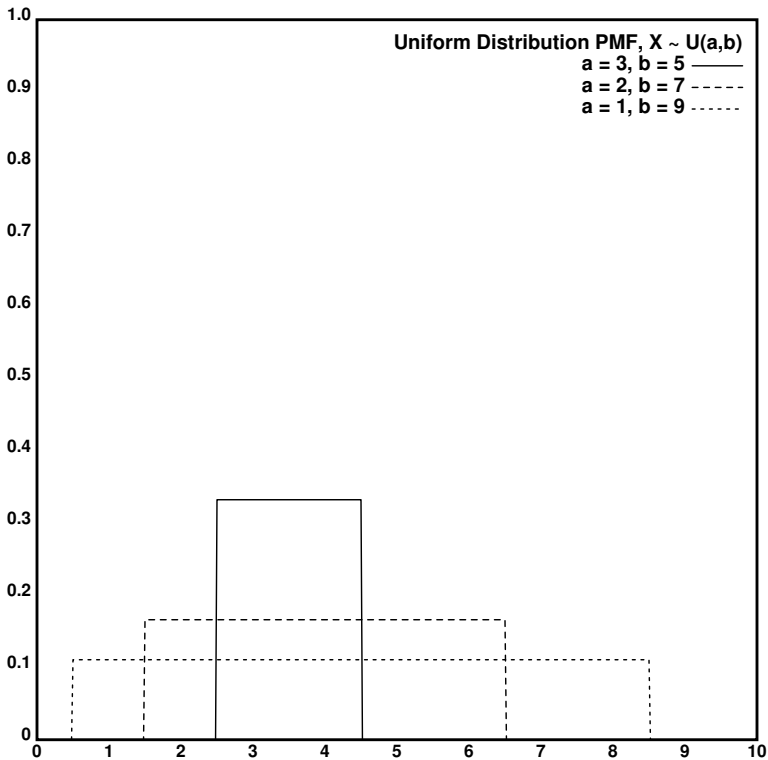
Uniform Distribution

$X \sim U(a, b)$	
PMF	$\begin{cases} 0 & \text{if } x < a \\ \frac{1}{b-a+1} & \text{if } a \leq x \leq b \\ 0 & \text{if } x > b \end{cases}$
CDF	$\begin{cases} 0 & \text{if } x < a \\ \frac{x-a+1}{b-a+1} & \text{if } a \leq x \leq b \\ 1 & \text{if } x > b \end{cases}$
Parameters	$a, b \in \mathbf{Z}, a \leq b$: range
μ	$\frac{a+b}{2}$
σ^2	$\frac{(b-a+1)^2 - 1}{12}$
Skewness (γ_1)	0

Question answered

PMF: what is the probability of an event x happening, given a constant probability?

CDF: what is the probability of any event in the range from a to x happening?



Examples

The probability of getting a specific face when rolling a six-sided die follows the uniform distribution $X \sim U(1, 6)$, so the probability of getting a “3” is

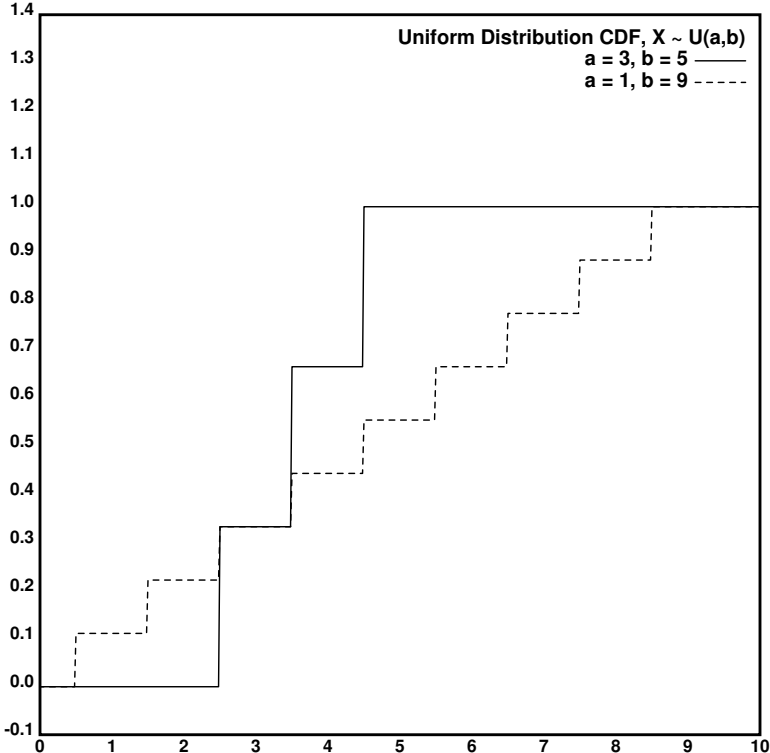
$$P(X = 3) = \frac{1}{b - a + 1} = \frac{1}{6 - 1 + 1} = \frac{1}{6}$$

The probability of getting a “1”, “2”, or “3” is:

$$P(X \leq 3) = \frac{x - a + 1}{b - a + 1} = \frac{3 - 1 + 1}{6 - 1 + 1} = \frac{1}{2}$$

The probability of getting at least a “3” is:

$$P(X \geq 3) = 1 - P(X \leq 2) = 1 - \frac{x - a + 1}{b - a + 1} = 1 - \frac{2 - 1 + 1}{6 - 1 + 1} = \frac{2}{3}$$



Geometric Distribution

$X \sim \text{Geo}(p)$	
PMF	$q^{x-1} \cdot p$
CDF	$1 - q^x$
Parameters	$p \in [0, 1]$: probability of success
	$q: 1 - p$
	$x \in \mathbf{N}_0$: number of trials
μ	$\frac{1}{p}$
σ^2	$\frac{q}{p^2}$
Skewness (γ_1)	$\frac{2-p}{\sqrt{q}}$

Questions answered

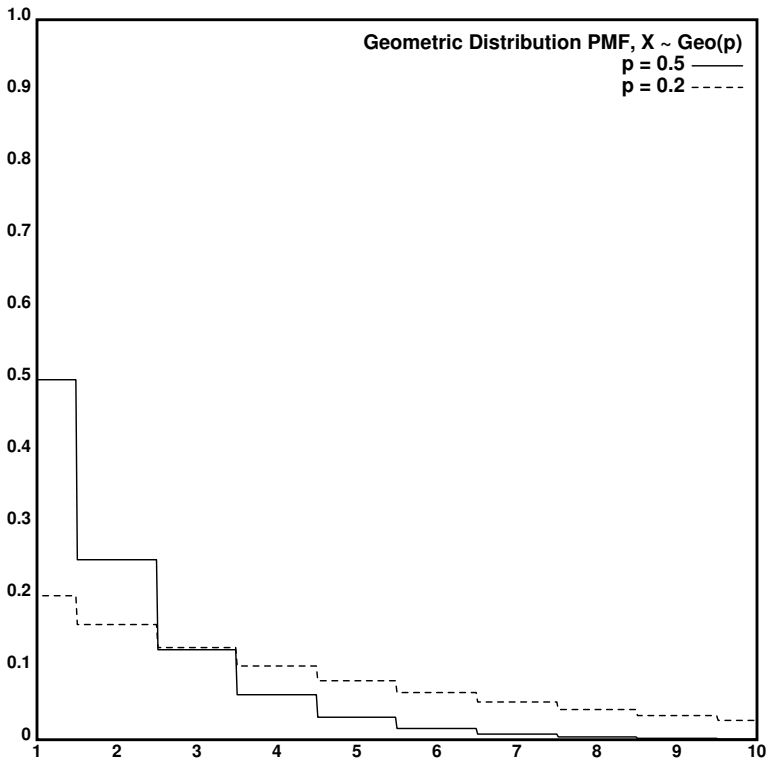
PMF: given x independent trials, all with equal probability of success p , what is the probability of *exactly* one success after $x - 1$ failures?

CDF: given x independent trials, what is the probability of *at least* one success?

Examples

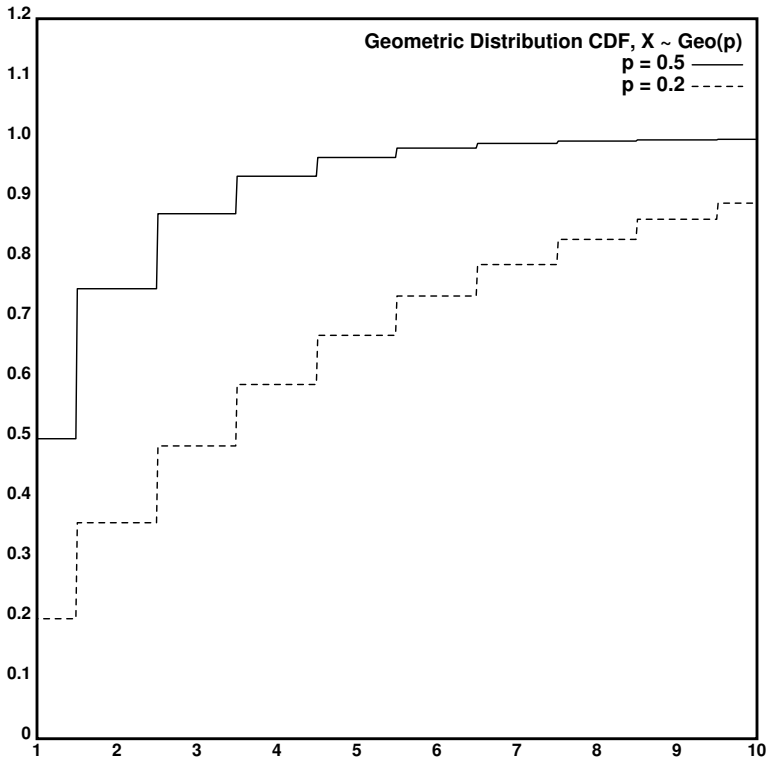
The probability of getting the first “six” in the n^{th} subsequent roll of a six-sided die follows the geometric distribution $X \sim \text{Geo}(\frac{1}{6})$. The probability of getting a six ($p = \frac{1}{6}$) in the $x = 3^{\text{rd}}$ toss of a die is:

$$\begin{aligned}
 P(X = 3) &= q^{x-1} \cdot p = \left(1 - \frac{1}{6}\right)^2 \cdot \frac{1}{6} = \left(\frac{5}{6}\right)^2 \cdot \frac{1}{6} \\
 &= \frac{25}{36} \cdot \frac{1}{6} = \frac{25}{216} \approx 0.116
 \end{aligned}$$



The probability of getting at least one six in three tosses is:

$$\begin{aligned} P(X \leq 3) &= 1 - q^x = 1 - \left(1 - \frac{1}{6}\right)^3 = 1 - \left(\frac{5}{6}\right)^3 \\ &= 1 - \frac{125}{216} = \frac{91}{216} \approx 0.421 \end{aligned}$$



Binomial Distribution

$X \sim B(n, p)$	
PMF	$\binom{n}{x} \cdot p^x \cdot q^{n-x}$
CDF	$\sum_{i=0}^x \binom{n}{i} \cdot p^i \cdot q^{n-i}$
	$I_q(n-x, 1+x)$
Parameters	$n \in \mathbf{N}_0$: number of trials
	$p \in [0, 1]$: probability of success
	$q: 1 - p$ (probability of failure)
	$x \in \mathbf{N}_0, x \leq n$: number of successes
μ	np
σ^2	npq
Skewness (γ_1)	$\frac{q-p}{\sqrt{npq}}$
Approximations	$N(np, npq)$ for $np > 5, nq > 5$
	$Poi(np)$ for $n \geq 50, p < 0.1$

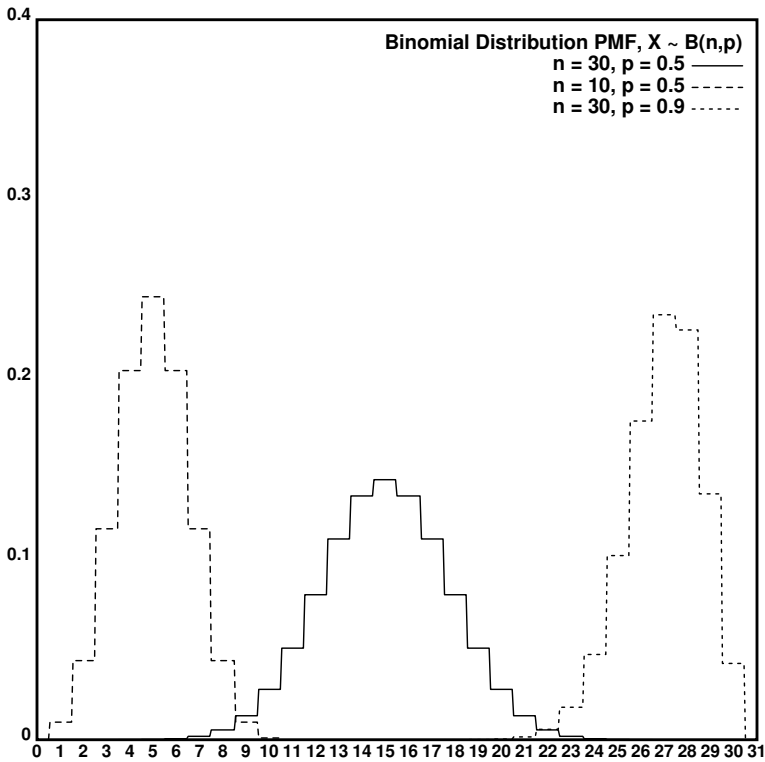
Questions answered

PMF: what is the probability of *exactly* x successes in n independent trials with a success probability of p ?

CDF: what is the probability of *up to* x successes in n independent trials?

Examples

The probability for x out of 5 children being girls (or boys) follows the binomial distribution $X \sim B(5, 0.5)$ given equal chances for a child being a girl or boy ($p = 0.5$).



Given this distribution, the probability for a couple having exactly $x = 3$ girls is:

$$\begin{aligned}
 P(X = 3) &= \binom{n}{x} \cdot p^x \cdot q^{n-x} = \binom{5}{3} \cdot 0.5^3 \cdot (1 - 0.5)^{5-3} \\
 &= \binom{5}{3} \cdot 0.5^3 \cdot 0.5^2 = 10 \cdot 0.125 \cdot 0.25 = 0.3125
 \end{aligned}$$

All other factors being equal, the probability of *up to* three of the children being girls is:

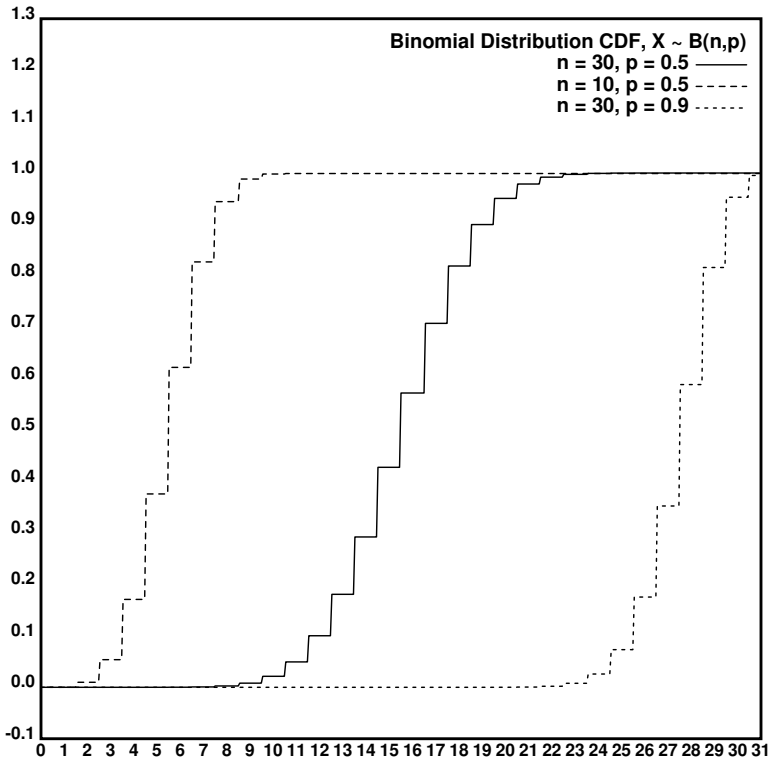
$$\begin{aligned}
 P(X \leq 3) &= \sum_{i=0}^x \binom{n}{i} \cdot p^i \cdot q^{n-i} = \sum_{i=0}^3 \binom{5}{i} \cdot 0.5^i \cdot 0.5^{5-i} \\
 &= \binom{5}{0} \cdot 0.5^0 \cdot 0.5^5 + \binom{5}{1} \cdot 0.5^1 \cdot 0.5^4 + \binom{5}{2} \cdot 0.5^2 \cdot 0.5^3 + \binom{5}{3} \cdot 0.5^3 \cdot 0.5^2
 \end{aligned}$$

$$= 0.03125 + 0.15625 + 0.3125 + 0.3125 = 0.8125$$

A more effective way to compute the above would be:

$$P(X \leq 3) = I_q(n - x, 1 + x) = I_{0.5}(2, 4) = 0.8125$$

where $I_q(a, b)$ denotes the regularized incomplete B (beta) function (see page 108).



Index

- 3-sigma, rule of 29
- 5-sigma certainty 49
- alpha level 97
- alternative hypothesis 95
- anticorrelation 50
- Bayes' Theorem 15
- bell curve 34
- Bernoulli trial 32
- best guess 43
- beta function 118
 - incomplete 108
 - regularized incomplete 68
- binomial coefficient 32
- binomial distribution 32, 66
- cause 15
- CDF 21, 34
- central limit theorem 39, 45
- central tendency 25
- certainty 10
 - 5-sigma 49
- chi-square distribution 86, 106
- chi-square statistic 86, 100
- chi-square test 86, 100
- choice 19
- CLT 39, 45
- combination 18, 31
- complement 11
- confidence interval 46, 91, 110
- confidence level 46, 91, 110
- contingency table 99
- continuity correction 80
- correlation coefficient 53
- correlation 50, 53
- covariance 53
- critical region 46, 88, 96
- cumulative distribution function
 - 21, 34
- degrees of freedom 86, 102
- distribution
 - binomial 32, 66
 - chi-square 86, 106
 - geometric 62
 - hypergeometric 70
 - normal 78
 - Poisson 74, 106
 - standard normal 82
 - uniform 58
- distribution function 21, 34
- effect 15
- error function 79, 83
 - computation 104
- Euler's constant 34
- event 10
- expectation 24, 86
 - contingency table 100
- explained variable 50
- explanatory variable 50
- factorial 17
- failure 32
- false negative 15
- false position method 111
- false positive 15
- frequency table 98
- function
 - beta 118
 - distribution 21, 34
 - error 79, 83
 - gamma 104
 - phi 35, 104
 - root of 110

- function
 - quantile 37, 38, 110
 - stochastic 20
- gamma function 104
 - incomplete 106
 - regularized incomplete 75, 88
- Gauss error function 79, 83
 - computation 104
- geometric distribution 62
- goodness of fit 86
- head 44
- hypergeometric distribution 70
- hypothesis testing 95
- iid 39
- impossibility 10
- income inequality 29
- incomplete beta function 108
- incomplete gamma function 106
- independence 14, 100
- inference 43
- intercept 51
- interquartile range 28
- intersection 11
 - probability 14
- interval, two-tailed 47
- inverse CDF 37, 38, 110
- k-combination 18
- k-permutation 18
- level of confidence 46, 91, 110
- level of significance 95
- linear regression 51
- location 25
- logical AND 11
- logical OR 11
- mean of sample 43
- mean 24
- median 26, 26, 37
- mode 26
- multivariate data 99
- mutual exclusion 12
- negative correlation 50
- Newton's method 110
- non-linear regression 56
- normal distribution CDF 104
- normal distribution 34, 78
- null hypothesis 95
- observation 86, 102
- outcome 10
- outlier 26
- P-value 0
- PDF 33
- Pearson's r 53
- percentile 37
- permutation 17
- phi function 35, 104
- pi 34
- PMF 21
- point estimate 43
- point estimator 43
- Poisson distribution 74, 106
- polynomial regression 56
- population 43
- positive correlation 50
- prevalence 16
- prior 16
- probability 10
 - complementary 31
 - conditional 13
 - prior 16
- probability density function 33
- probability distribution
 - continuous 33
 - discrete 22
- probability mass function 20, 21
- proportion 14
- q-quantile 37
- QF 37, 38, 110
- quantile function 37, 38, 110
- quantile 37
- quartile 28, 37

- random distribution 50
- random variable 20, 39
- random variable, discrete 21
- range 27
- RCP 15
- regressand 50
- regression line 51
- regression 50
- regression, linear 51
- regressor 50
- regula falsi method 111
- regularized incomplete
 - beta function 68
- regularized incomplete
 - gamma function 75, 88
- replacement 70
- residual sum of squares 55
- reverse conditional probability 15
- root of a function 110
- RSS 55
- rule of the three sigma 29
- sample mean 43
- sample point 45
- sample space 20
- sample 39, 43
 - variance 10
- sampling distribution 44
 - of means 45, 90
- scatter diagram 50
- score 34
- sensitivity 15
- significance 88
- significance level 95
- skewness 26, 44
- slope 51
- specificity 15
- SSE 55
- standard deviation 28
- standard error 48
- standard error of estimate 55
- standard normal CDF 35
- standard normal distribution
 - 34, 40, 46, 82
- standard score 34, 39
- statistic 34, 43
- statistic inference 43
- stochastic function 20
- Student's t-distribution 90, 108
- success 32
- sum of squared errors 55
- t-distribution 90, 108, 109
- t-score 91
- tail 44
- three sigma 29
- trial 10
 - Bernoulli 32
- type-I error 97
- uniform distribution 23, 58
- union 11
 - probability 12
- variables
 - explained 50
 - explanatory 50
 - iid 39
- variance 28, 53
 - sample 43
- variation 27
- Z-distribution 34, 82
- z-score 34, 39, 82
 - adjusted 47
- z-statistic 34