

Deconstructing Data Science

David Bamman, UC Berkeley

Info 290

Lecture 8: Naive Bayes

Feb 17, 2016

elements of probability in many of these methods



Linear regression

Decision trees

Ordinal regression

Probabilistic graphical models

Random forests

Logistic regression

Topic models

Survival models

Neural networks

K-means clustering

Perceptron

Random variable

- A variable that can take values within a fixed set (discrete) or within some range (continuous).

$$X \in \{1, 2, 3, 4, 5, 6\}$$

$$X \in \{the, a, dog, cat, runs, to, store\}$$

$$P(X = x)$$

Probability that the random variable X takes the value x (e.g., 1)

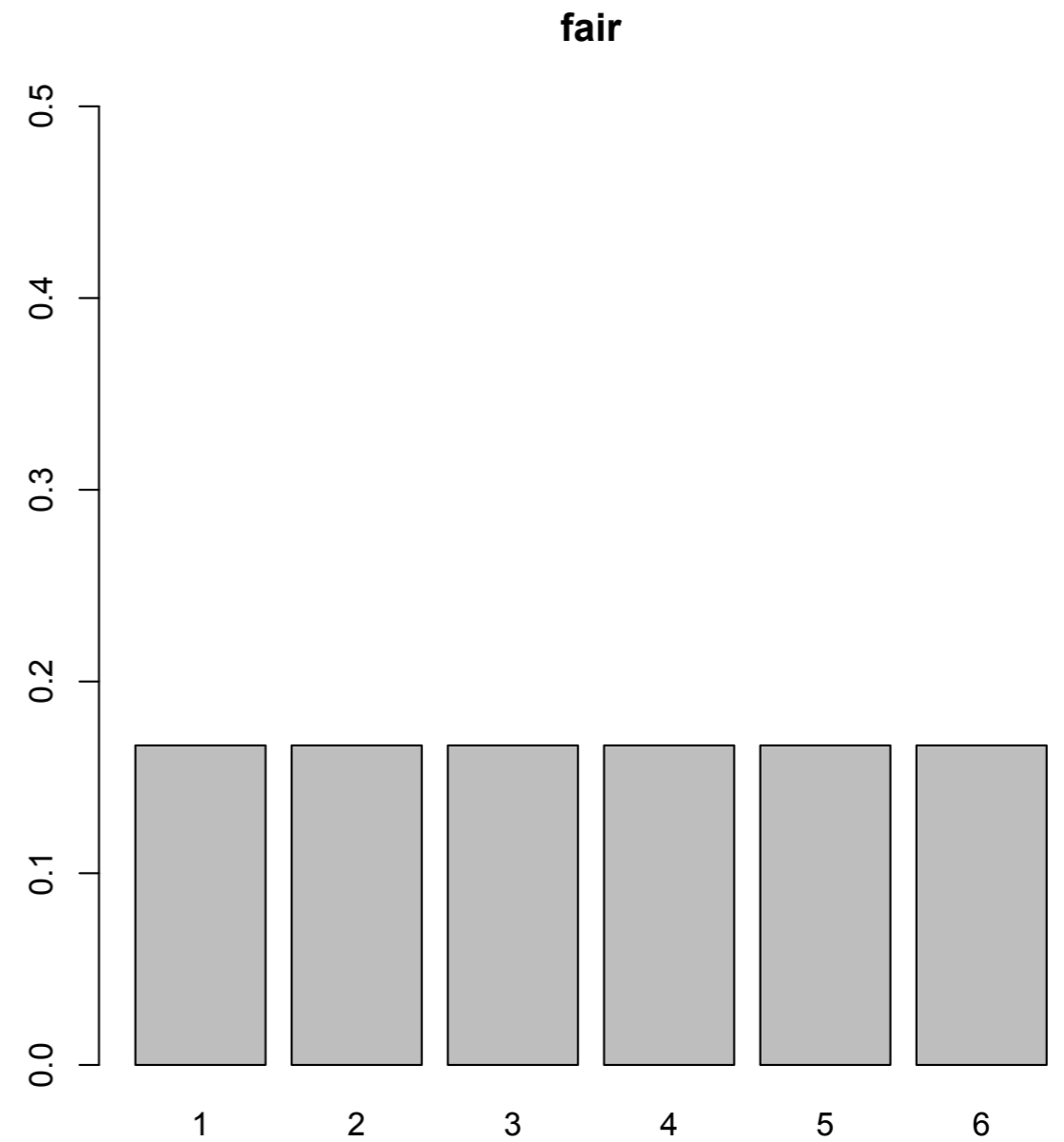
$$X \in \{1, 2, 3, 4, 5, 6\}$$

Two conditions:

1. Between 0 and 1: $0 \leq P(X = x) \leq 1$
2. Sum of all probabilities = 1 $\sum_x P(X = x) = 1$

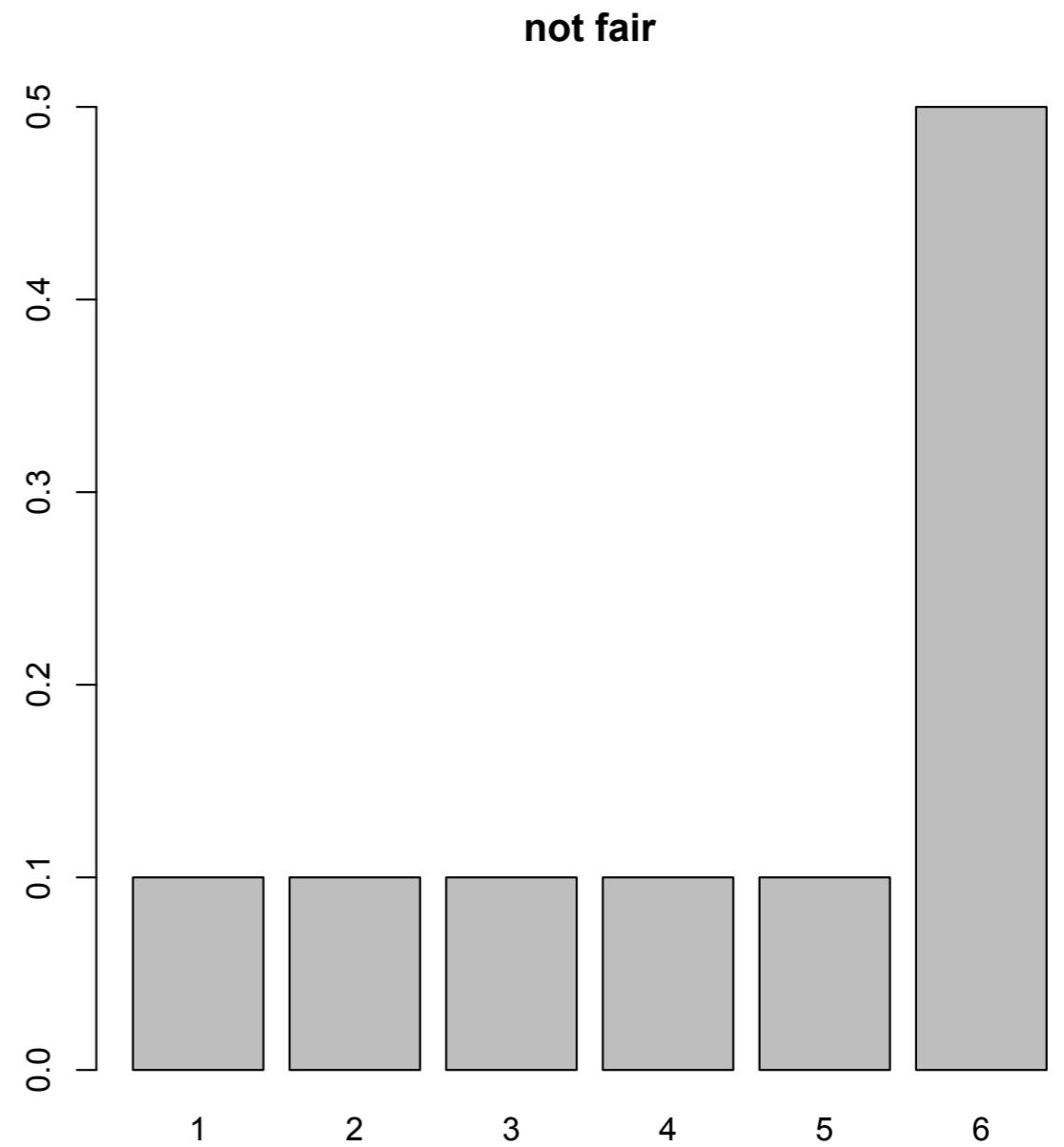
Fair dice

$$X \in \{1, 2, 3, 4, 5, 6\}$$



Weighted dice

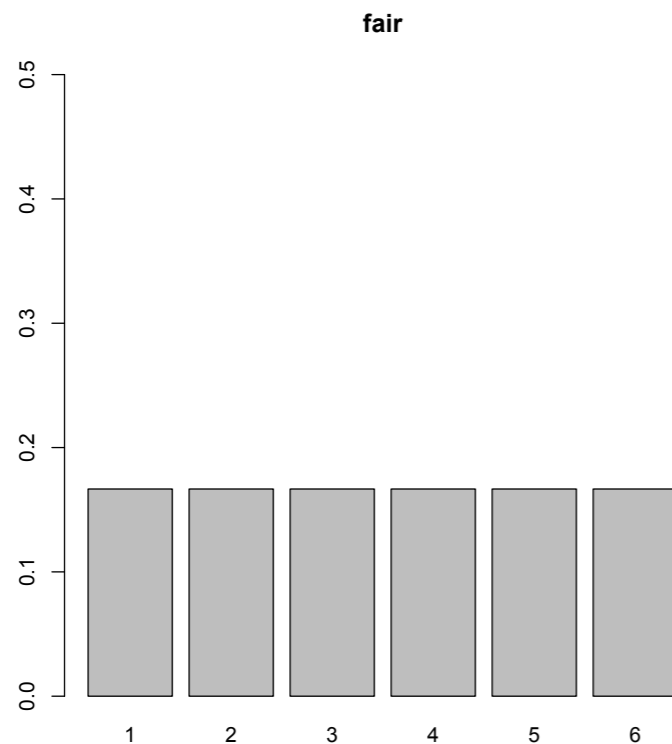
$$X \in \{1, 2, 3, 4, 5, 6\}$$



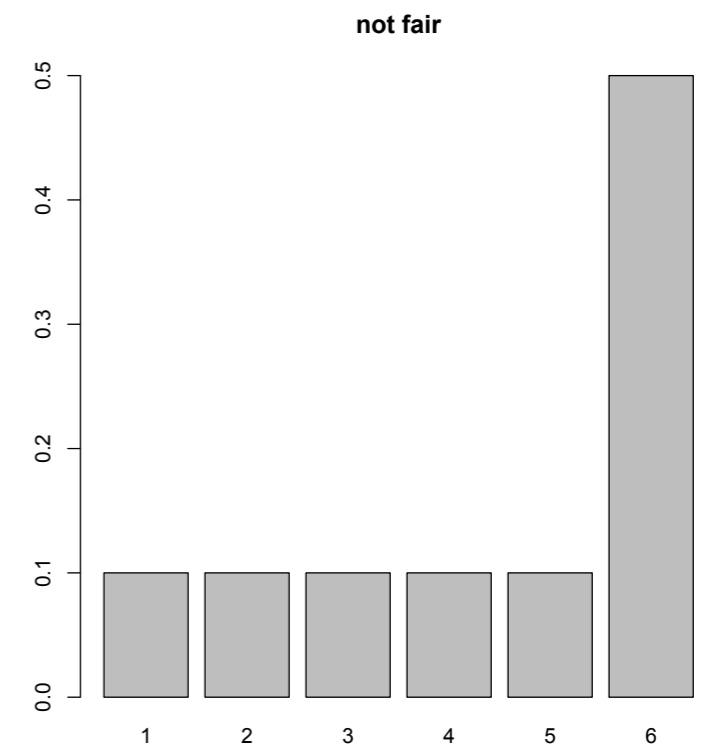
Inference

$$X \in \{1, 2, 3, 4, 5, 6\}$$

We want to *infer* the probability distribution that generated the data we see.

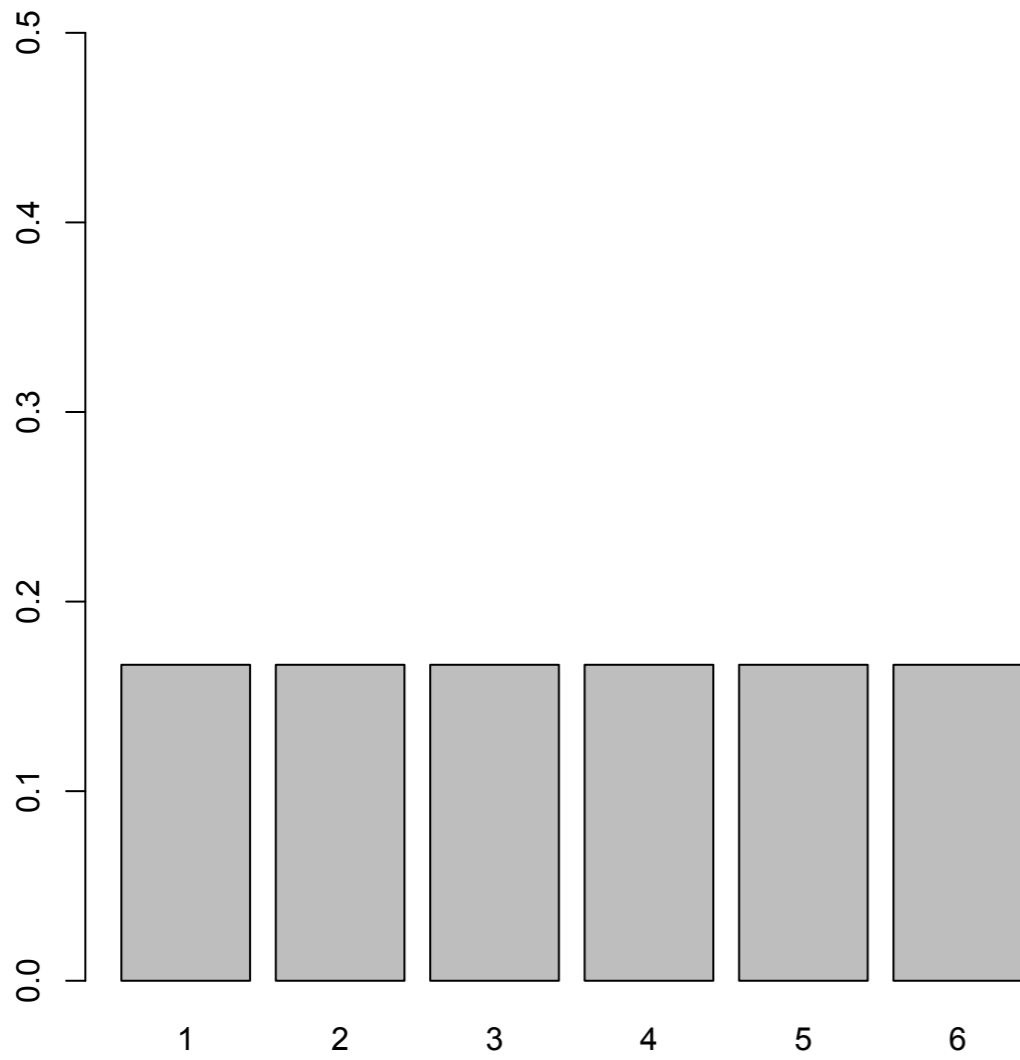


?

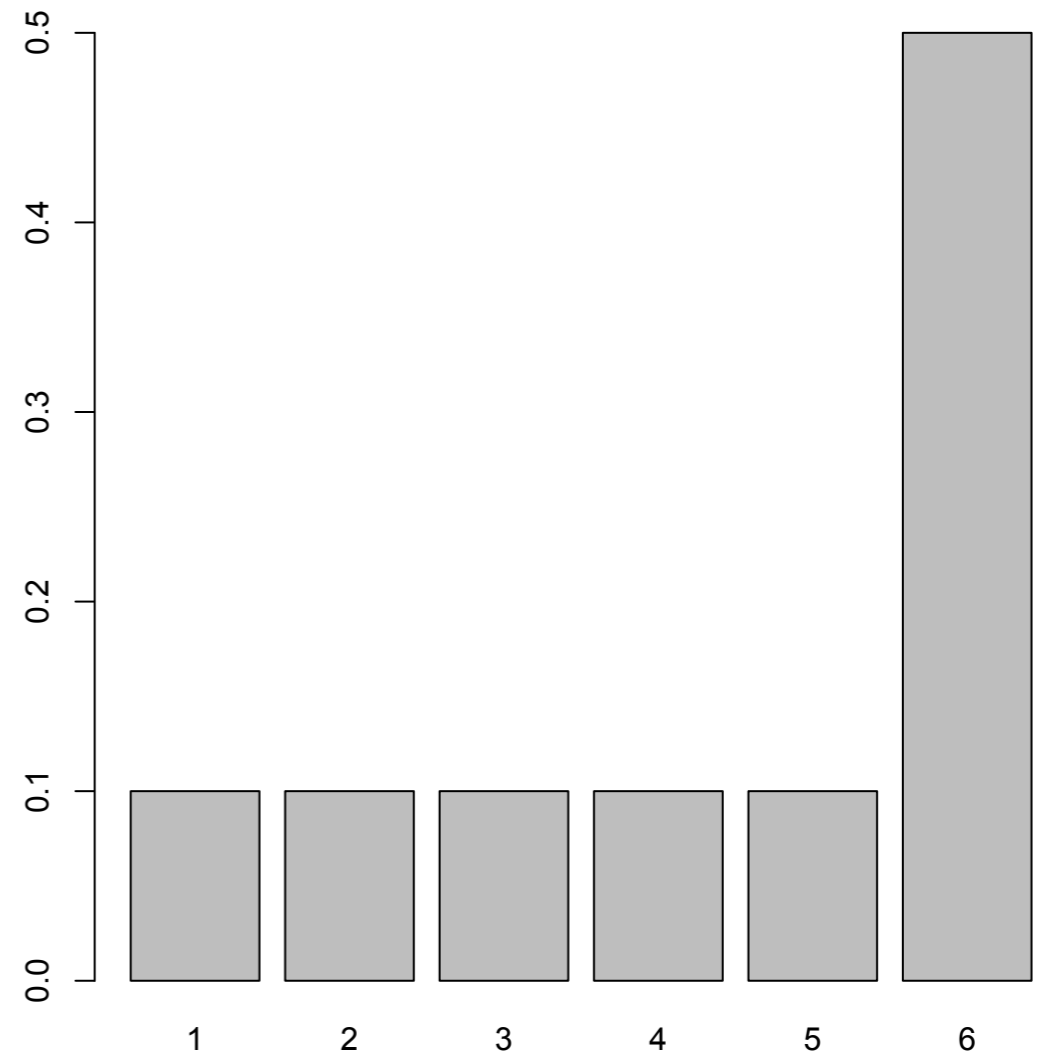


Probability

fair

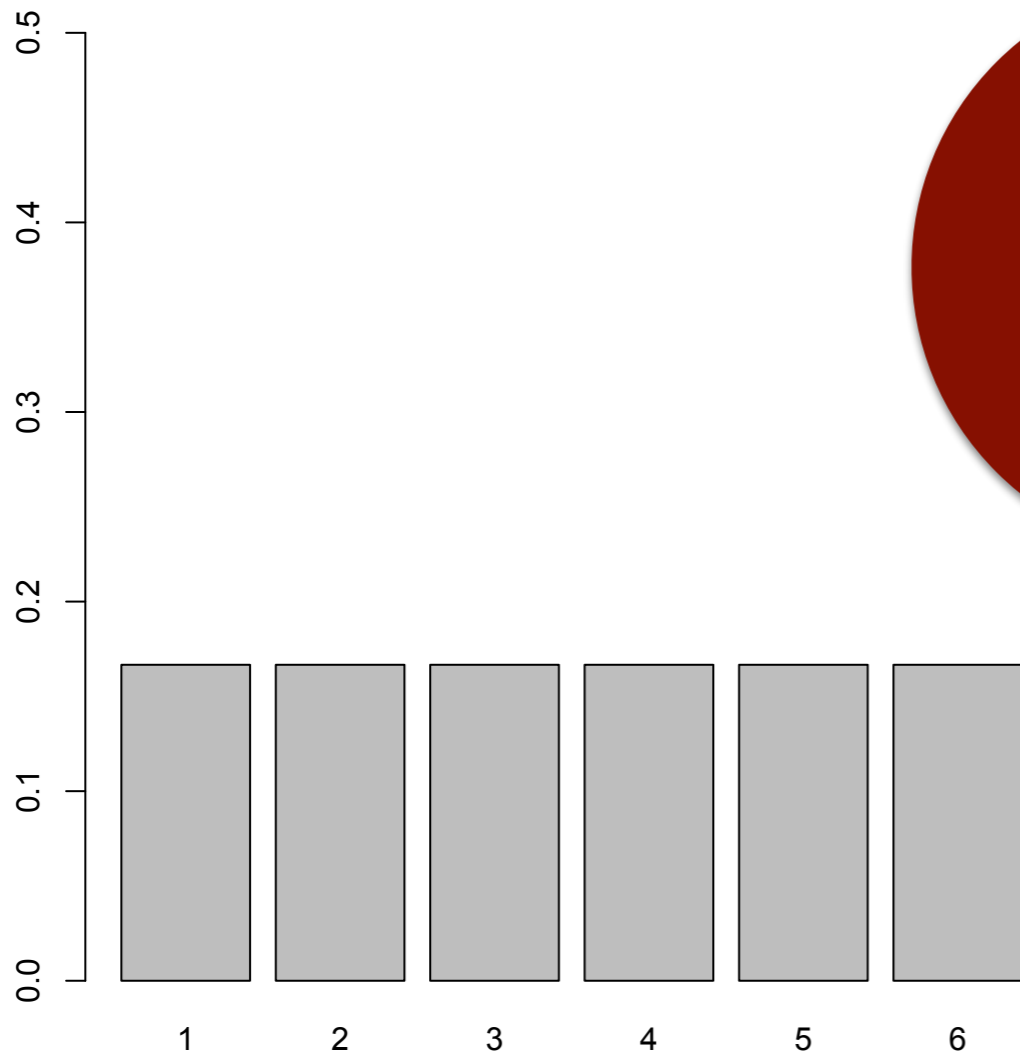


not fair

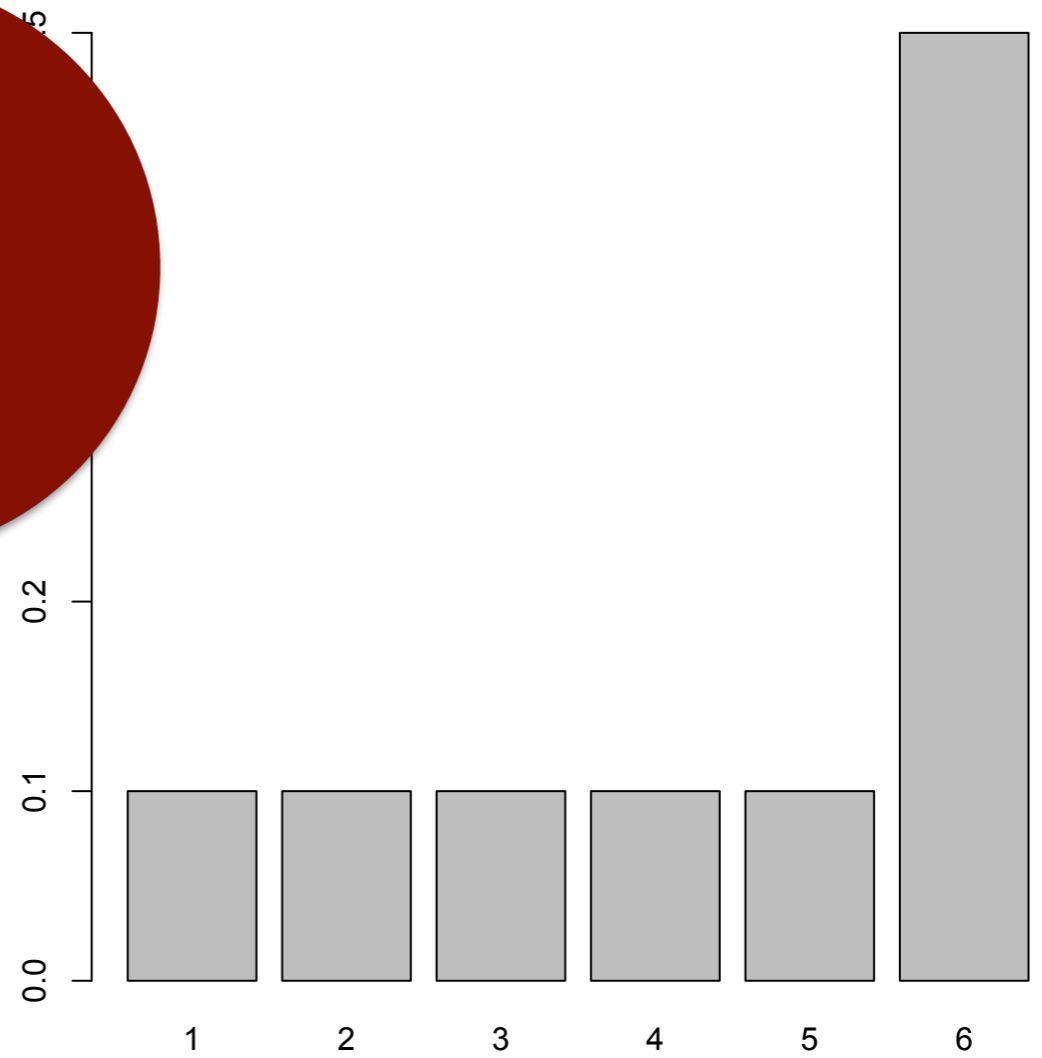


Probability

fair



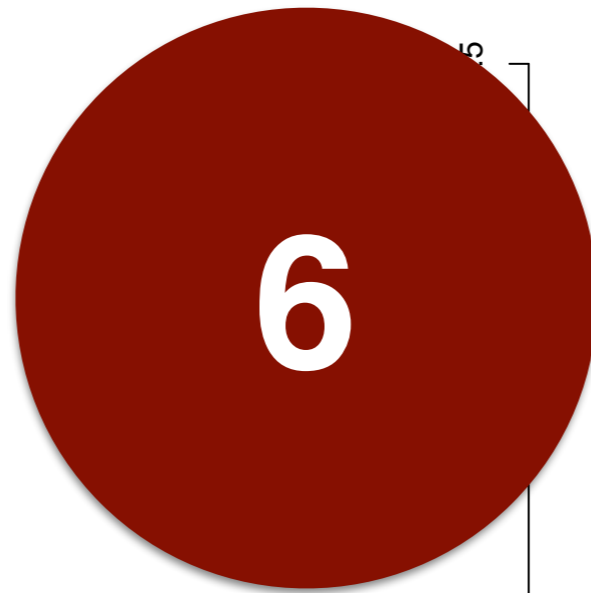
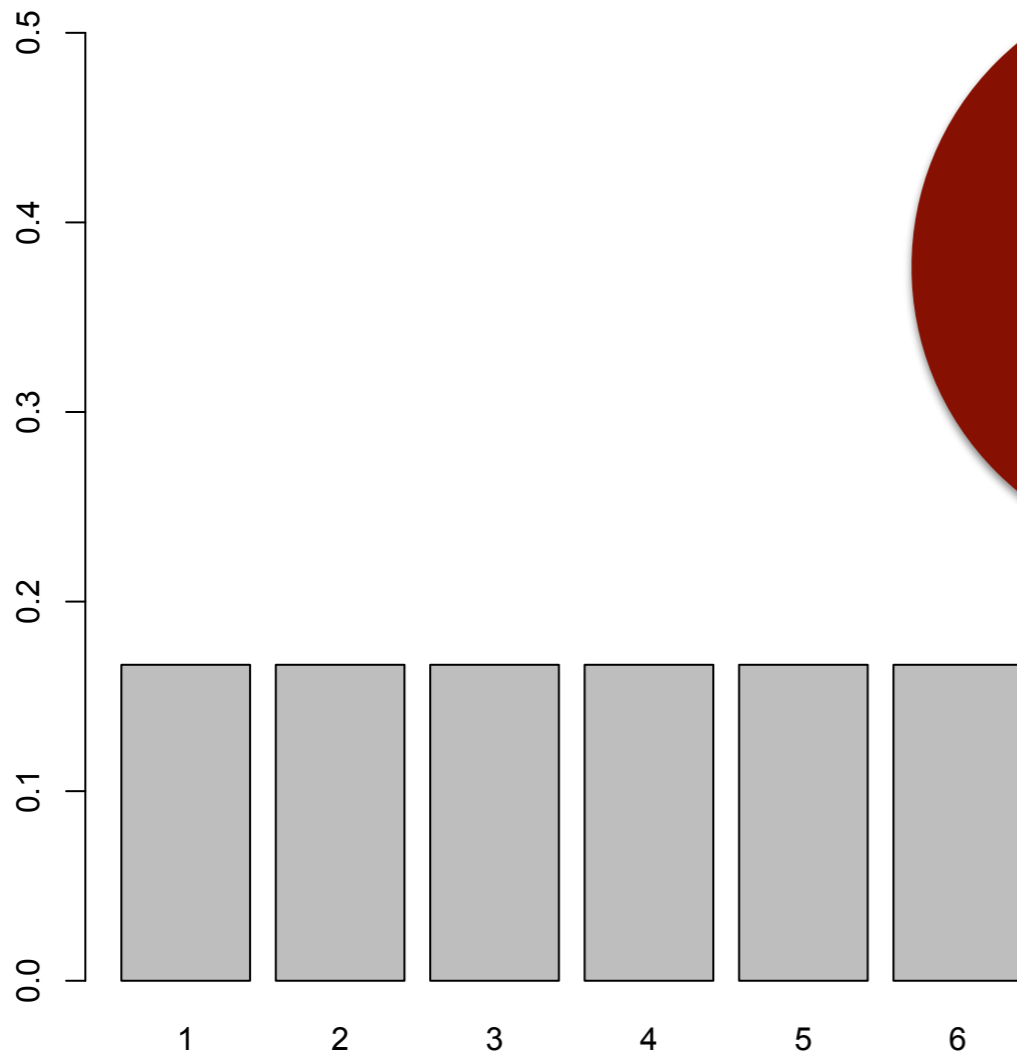
not fair



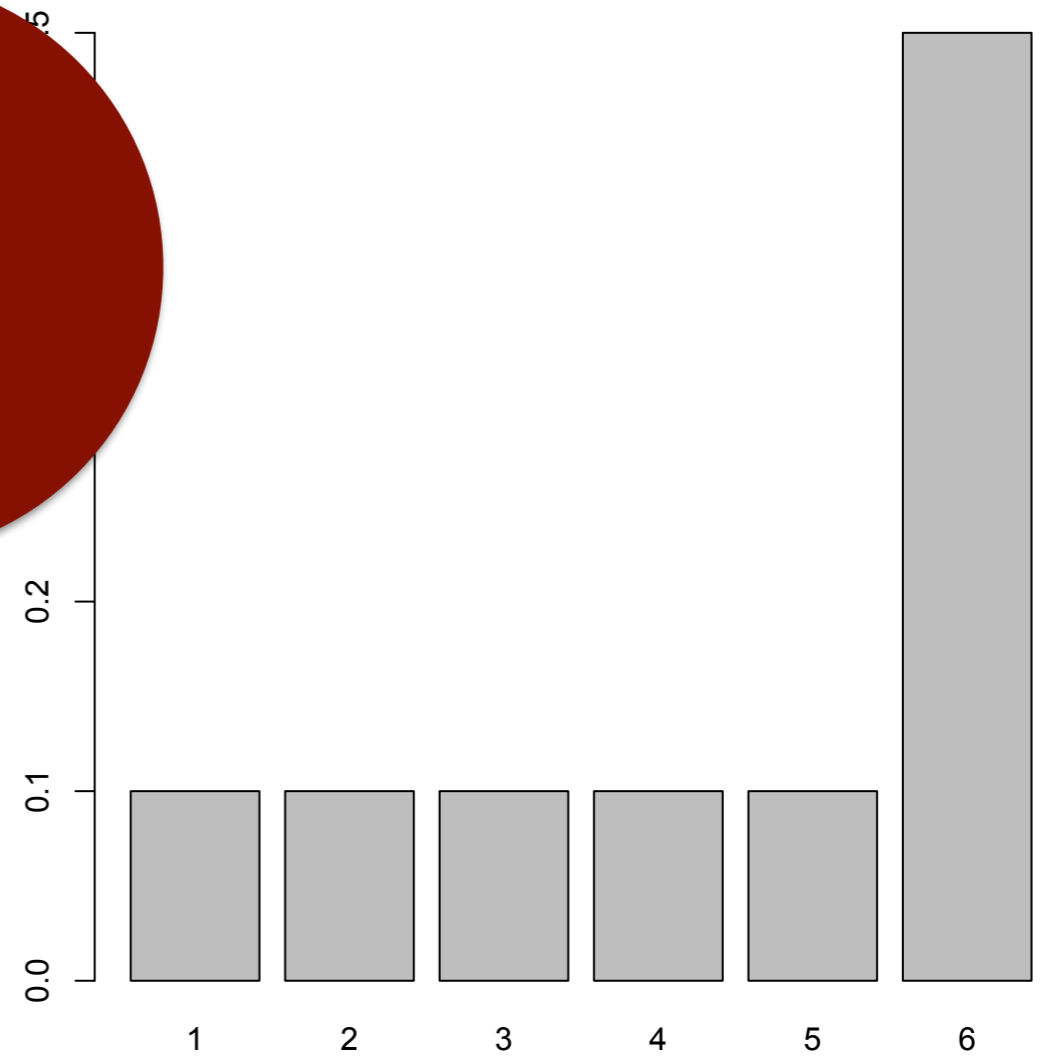
Probability

2

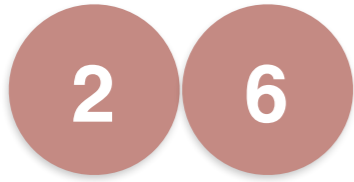
fair



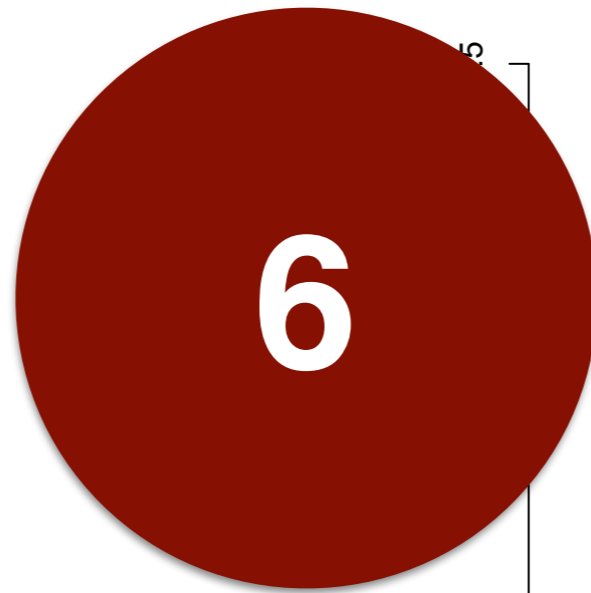
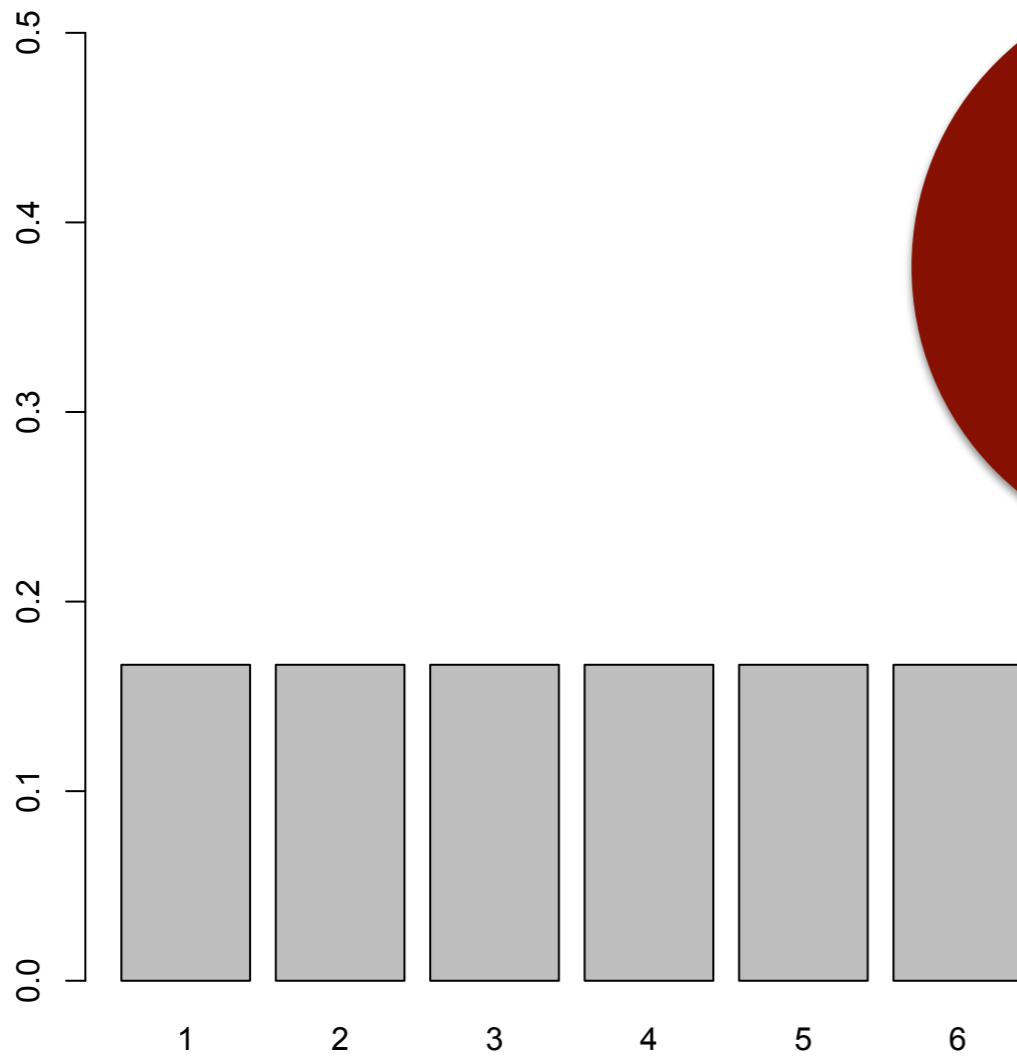
not fair



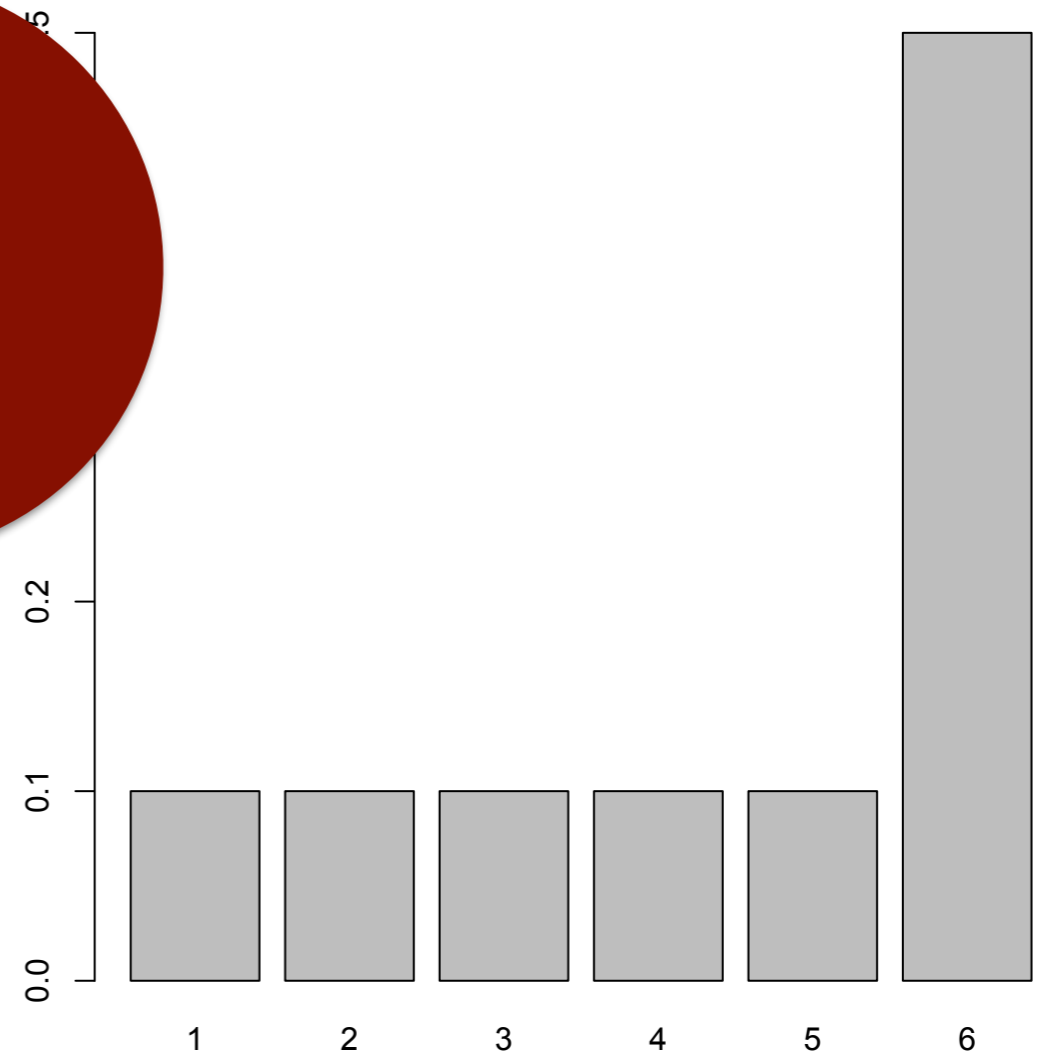
Probability



fair



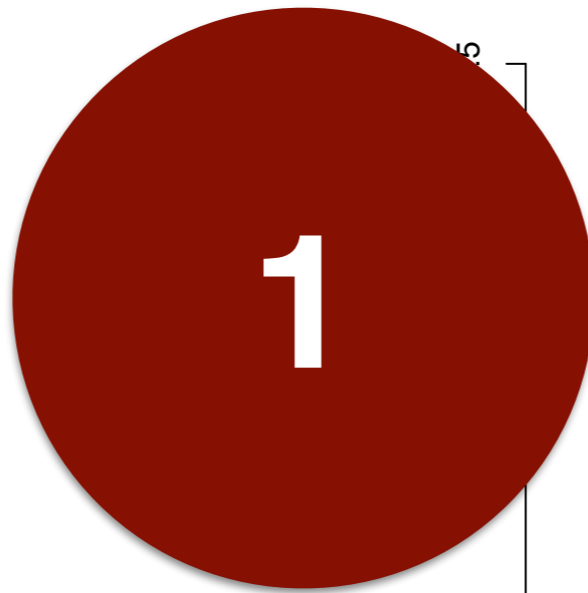
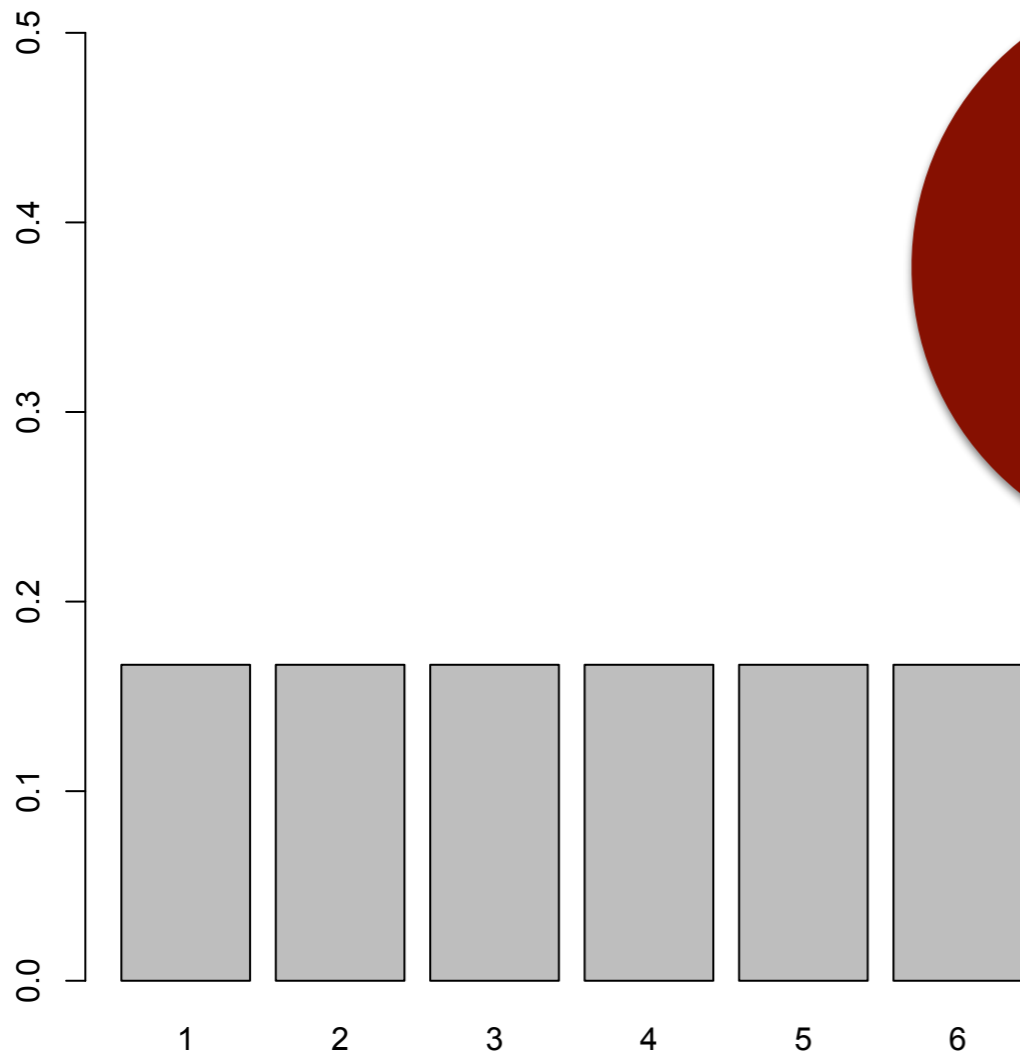
not fair



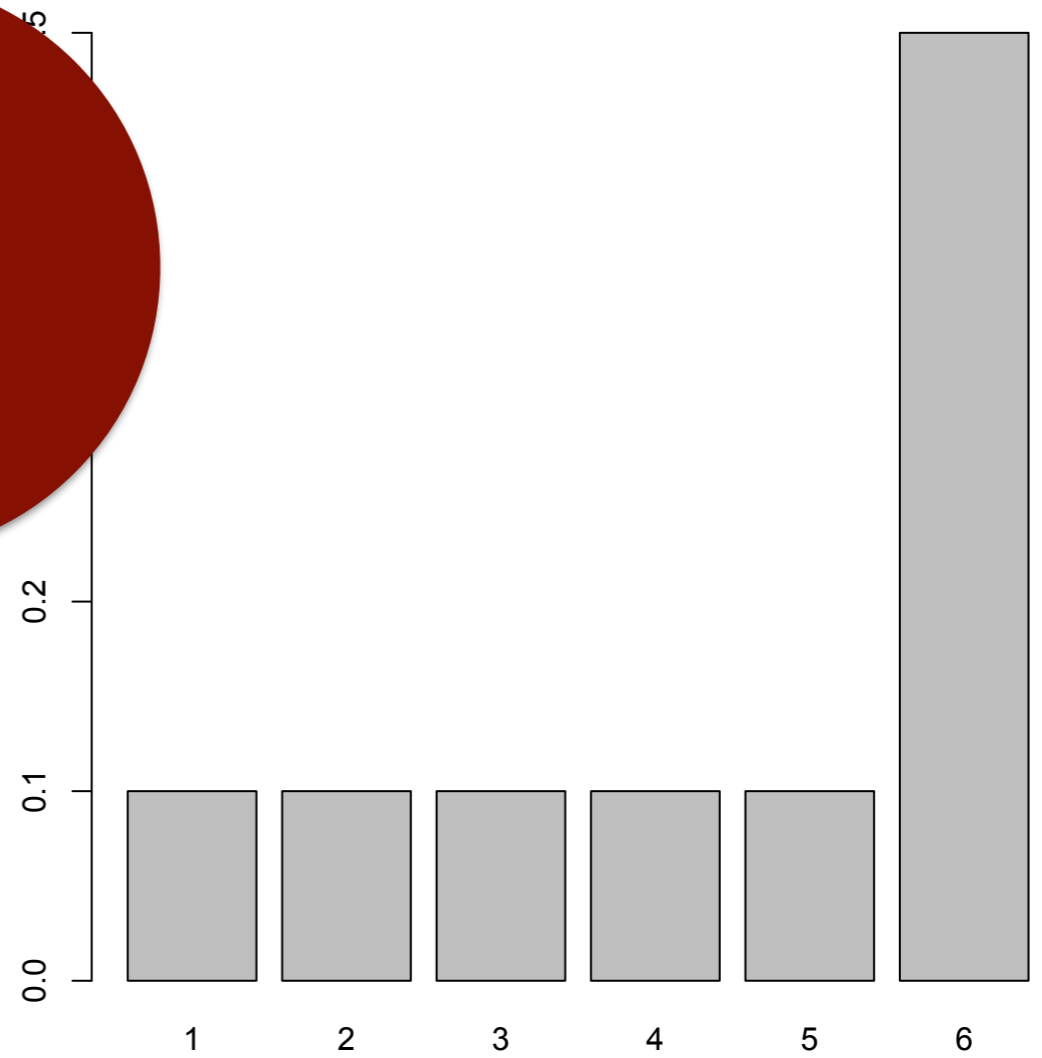
Probability



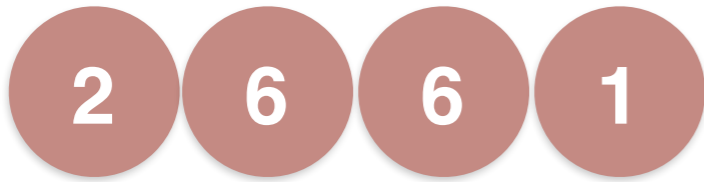
fair



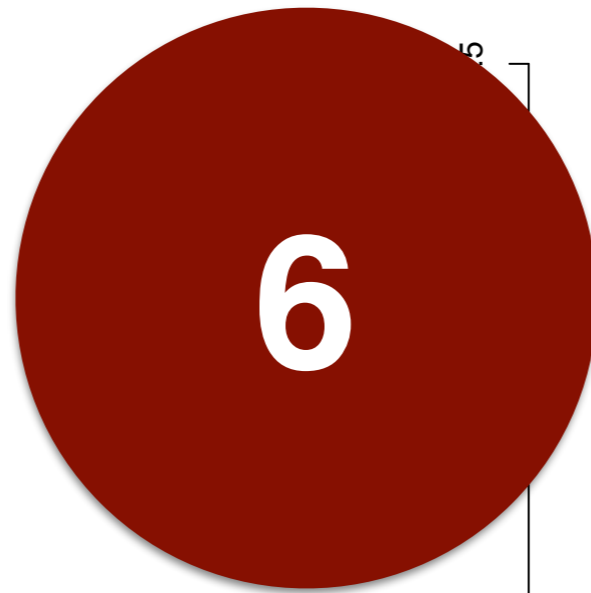
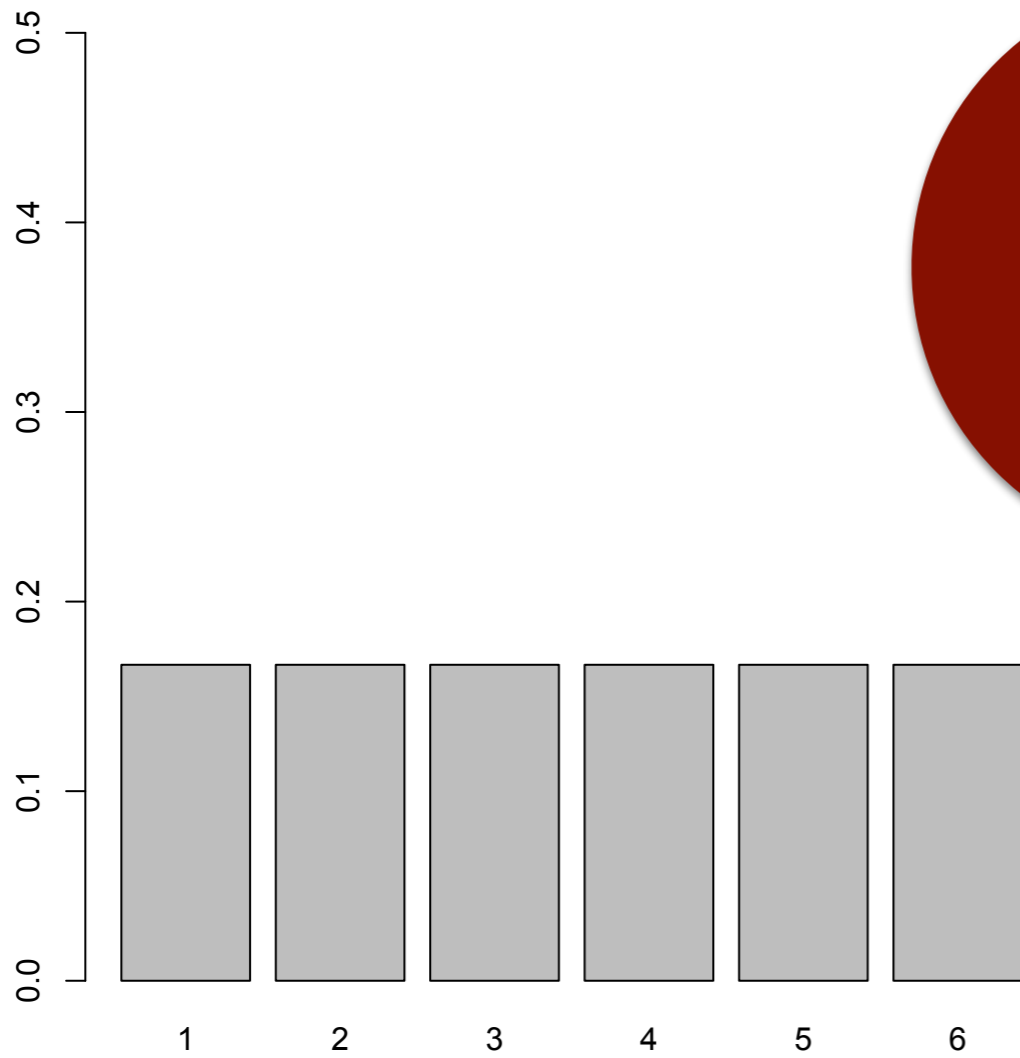
not fair



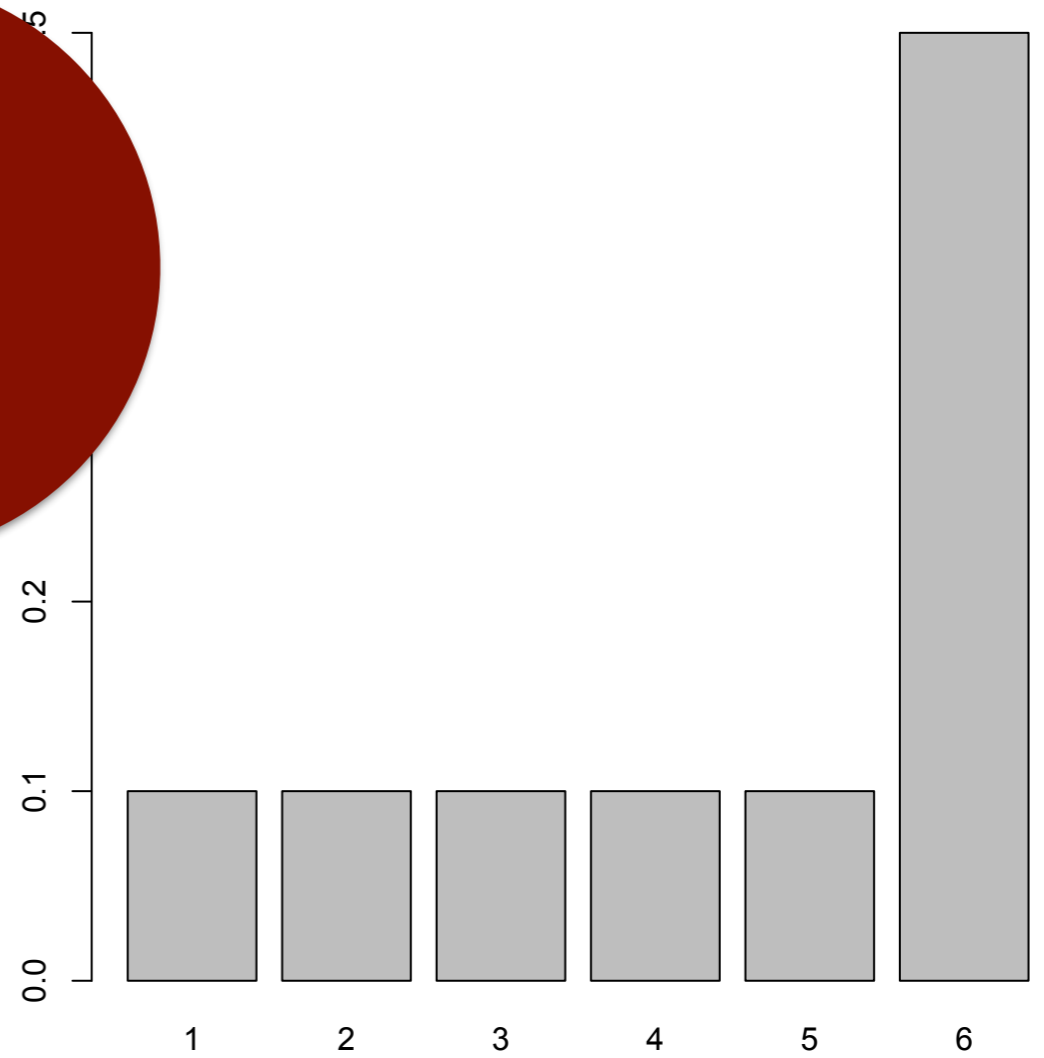
Probability



fair



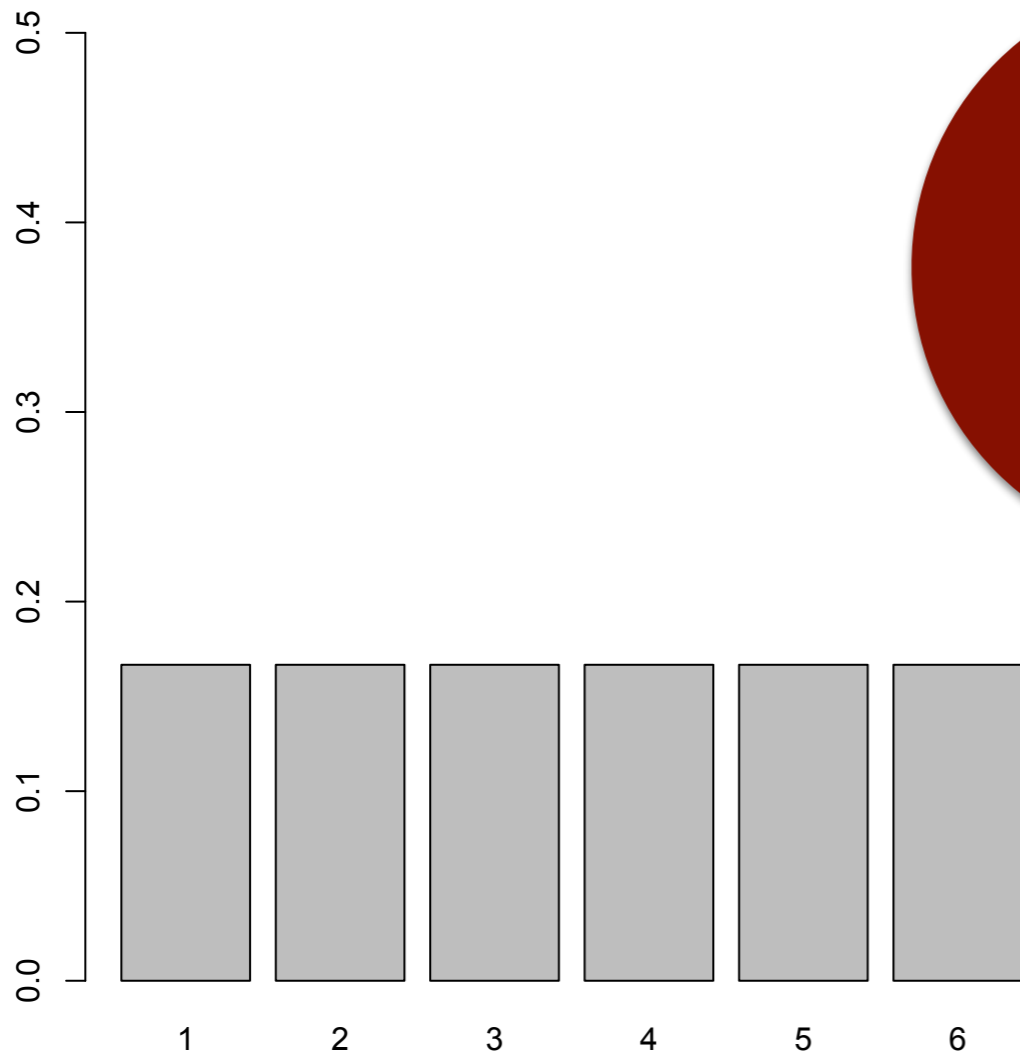
not fair



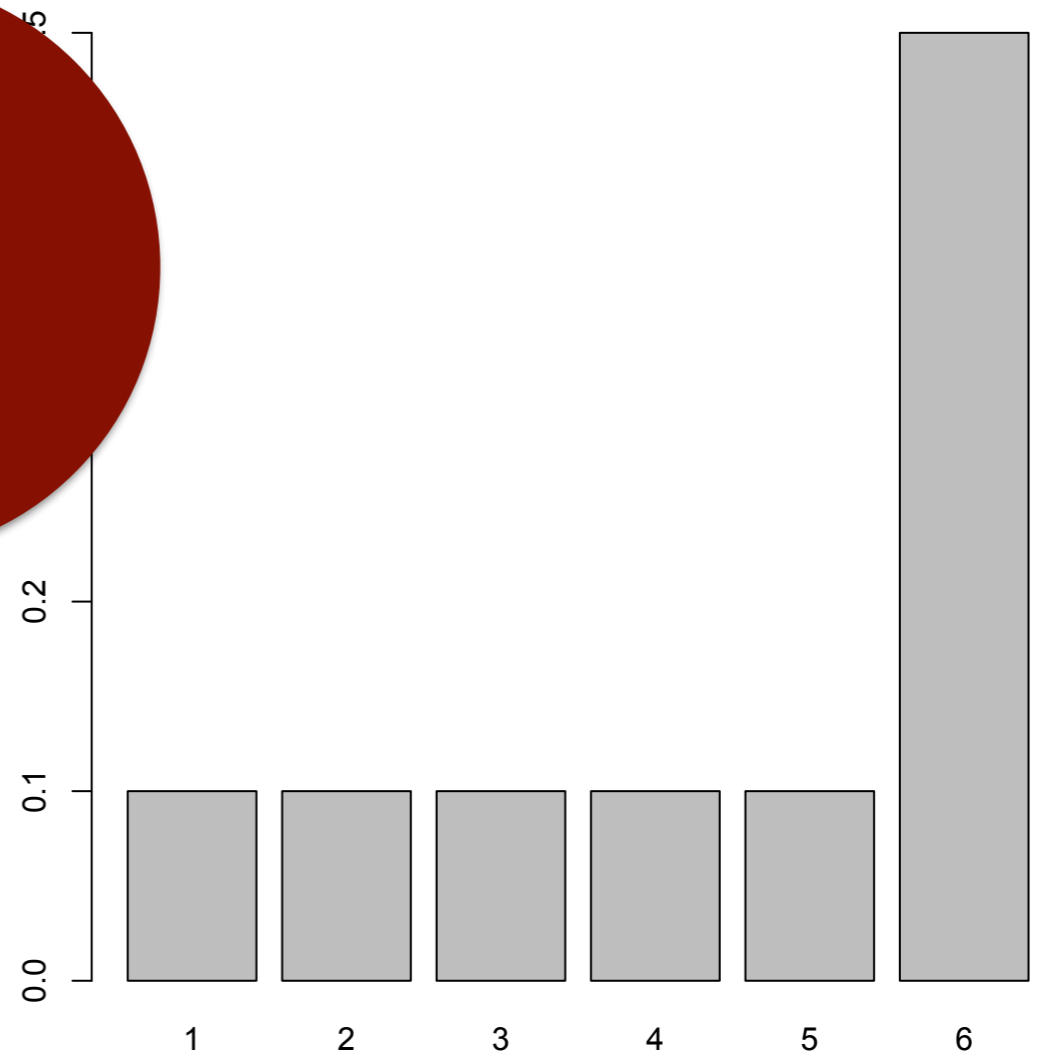
Probability



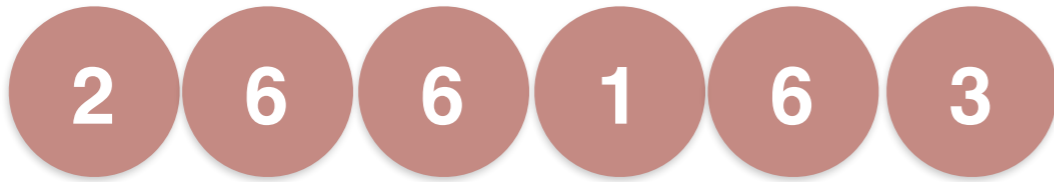
fair



not fair

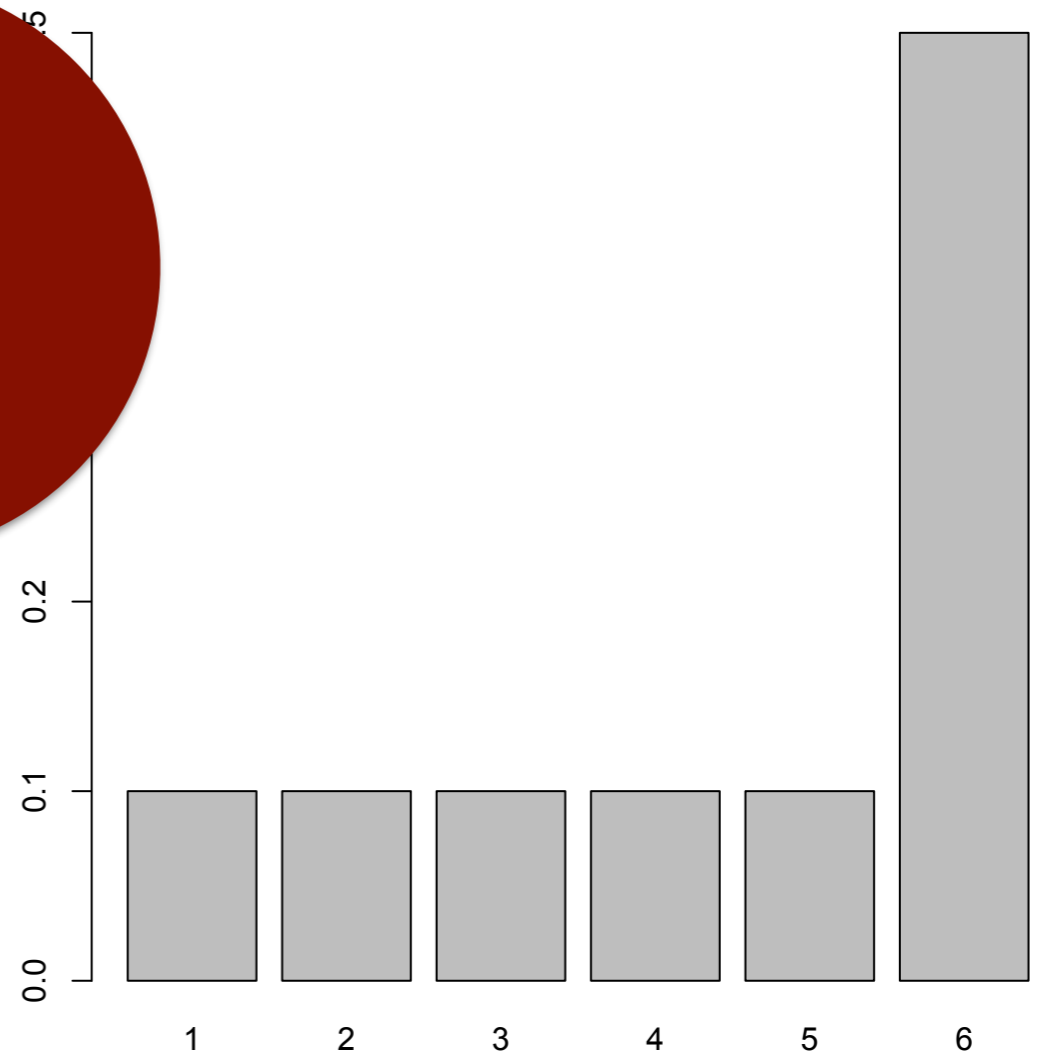
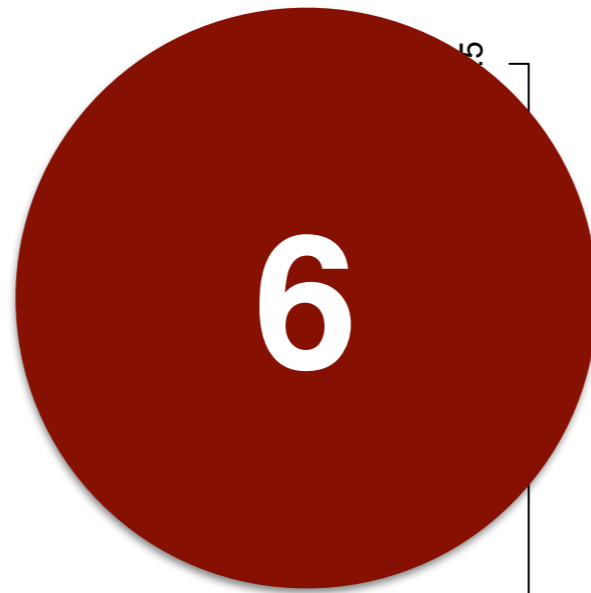
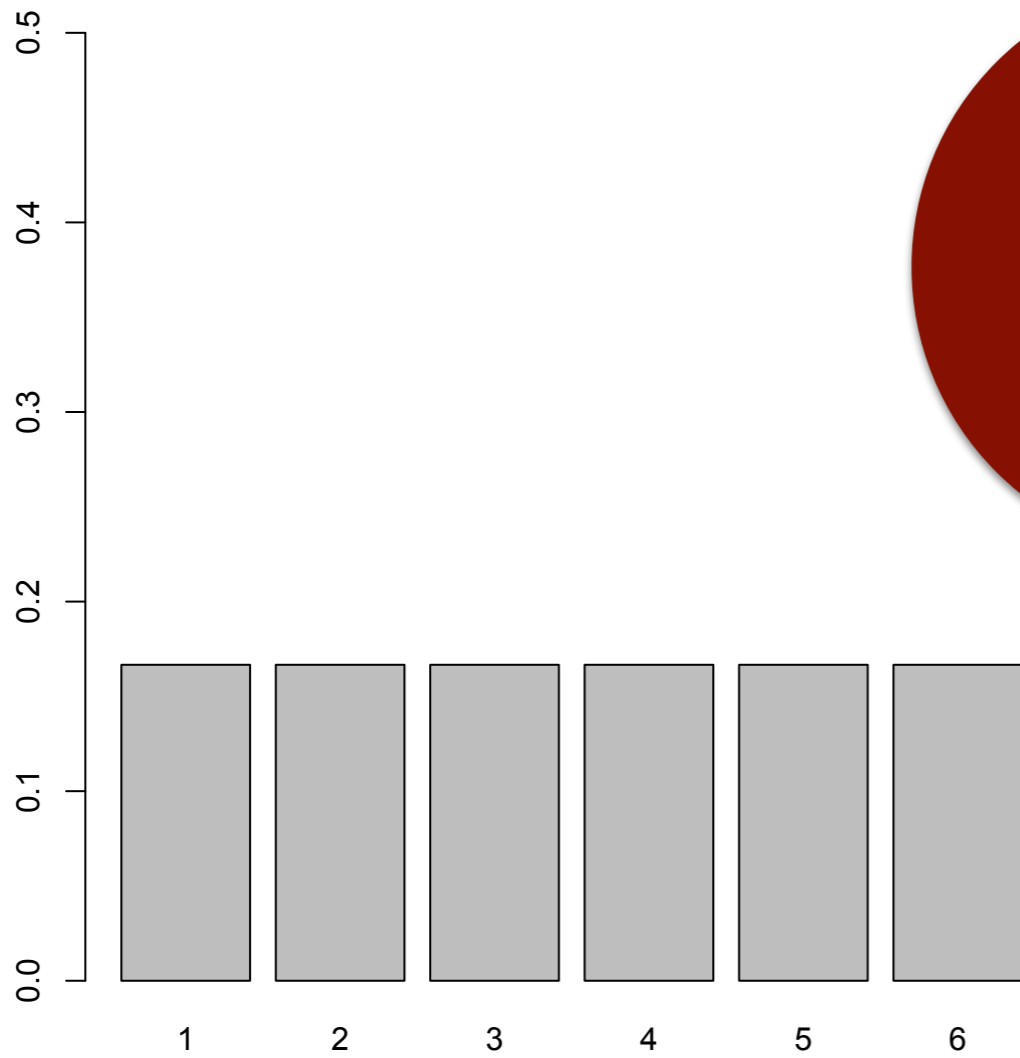


Probability



fair

not fair

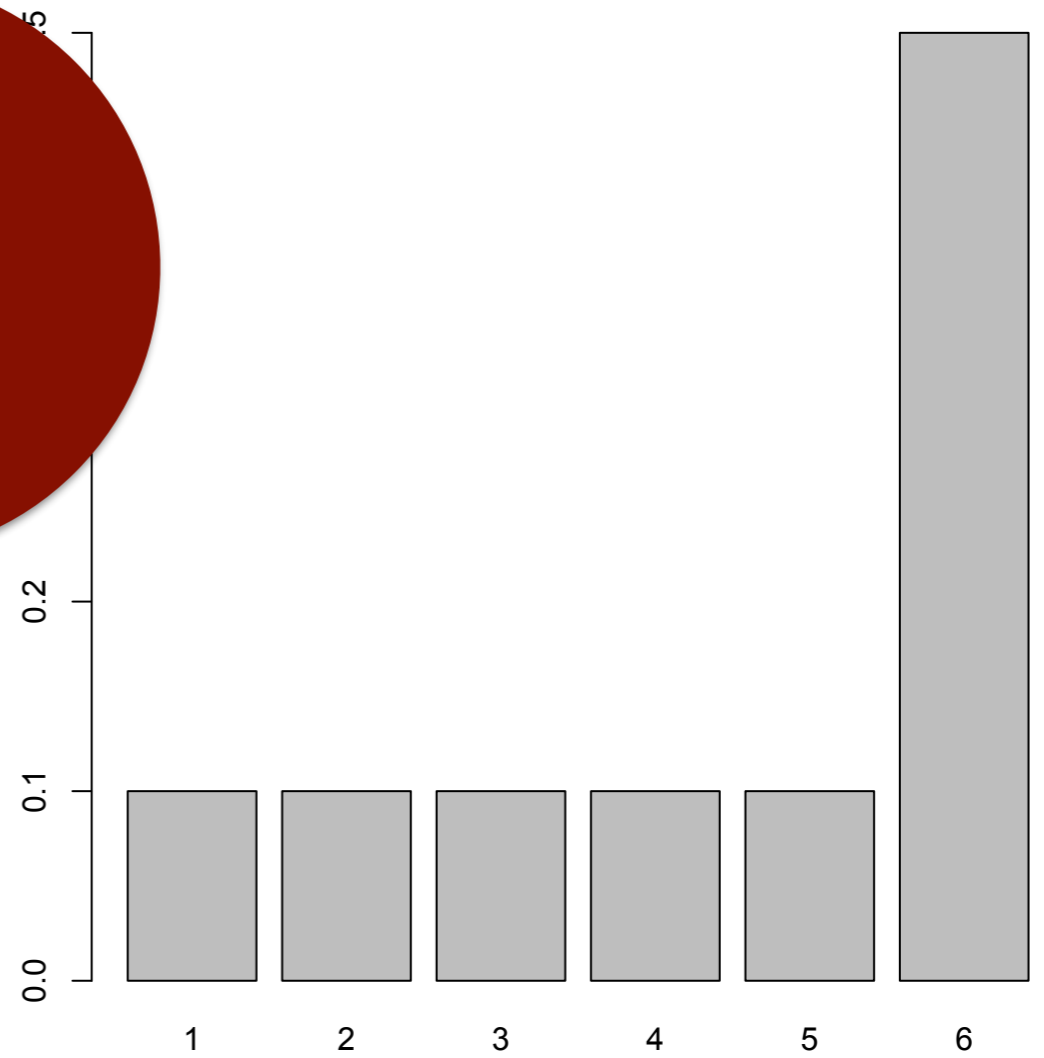
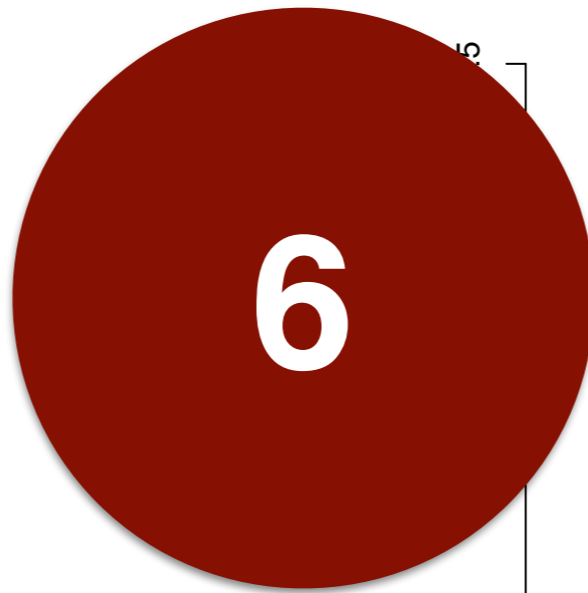
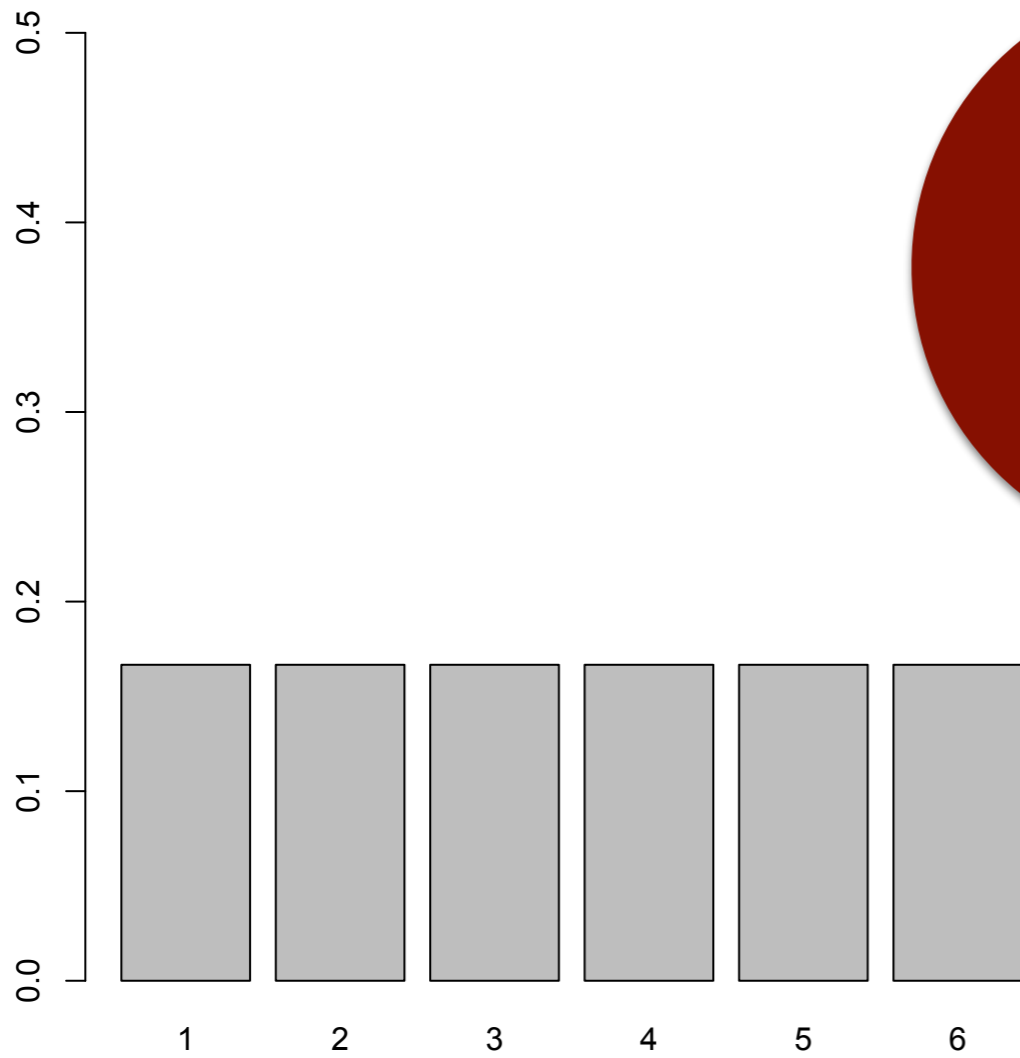


Probability

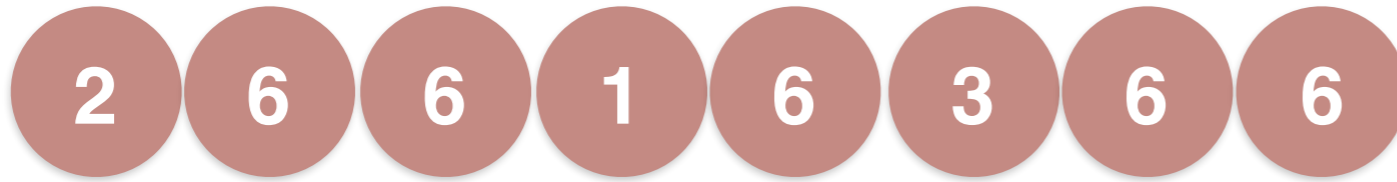


fair

not fair

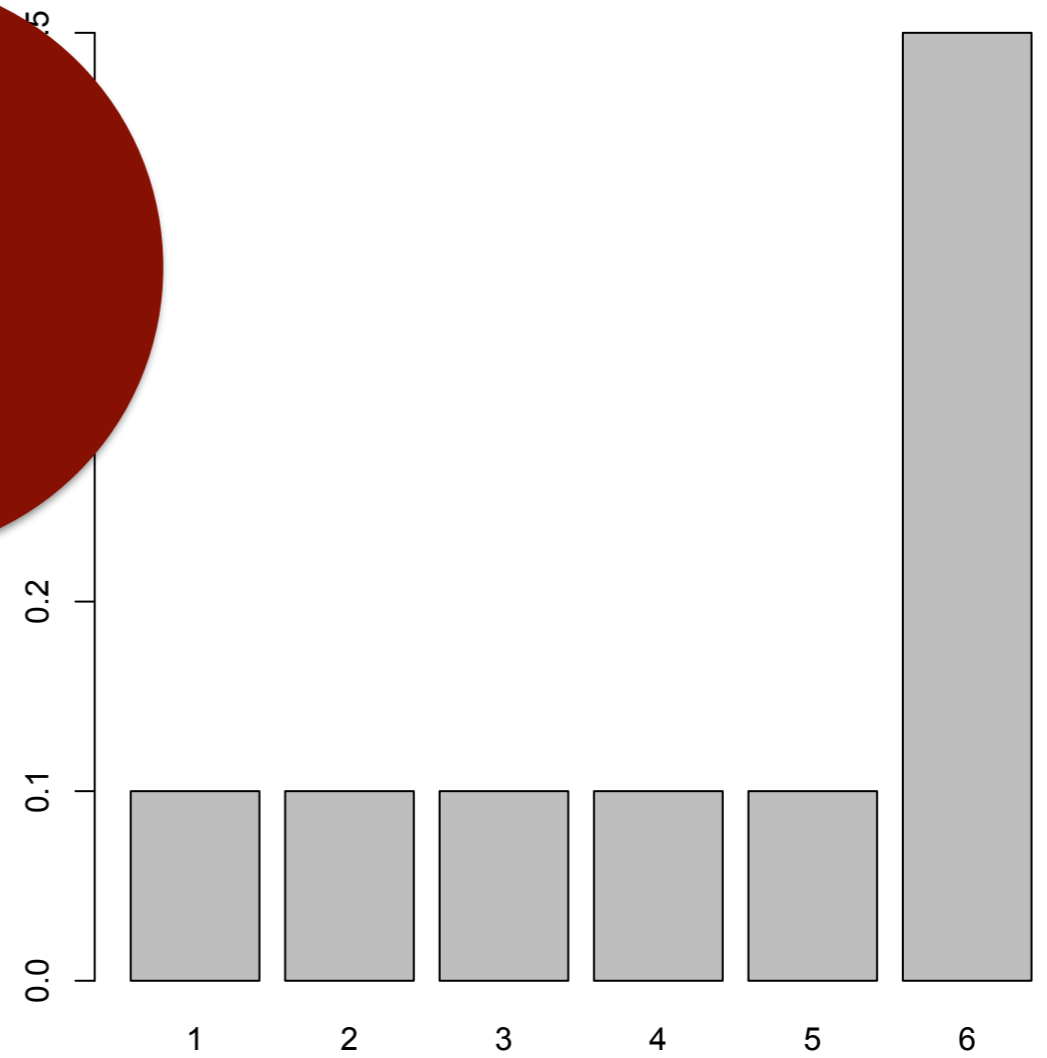
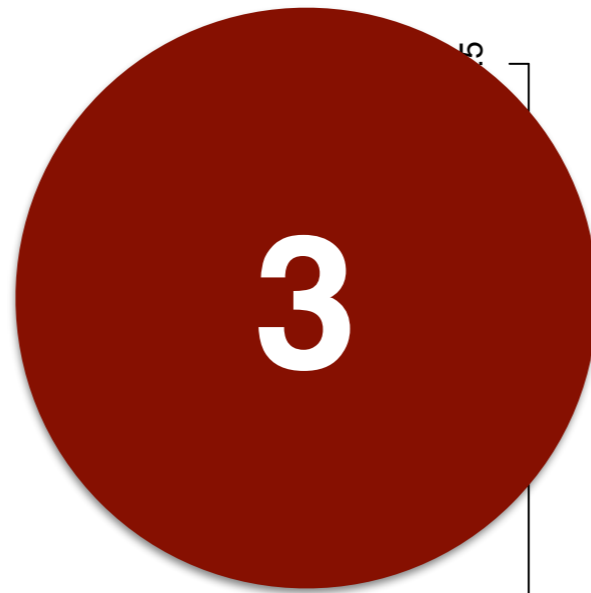
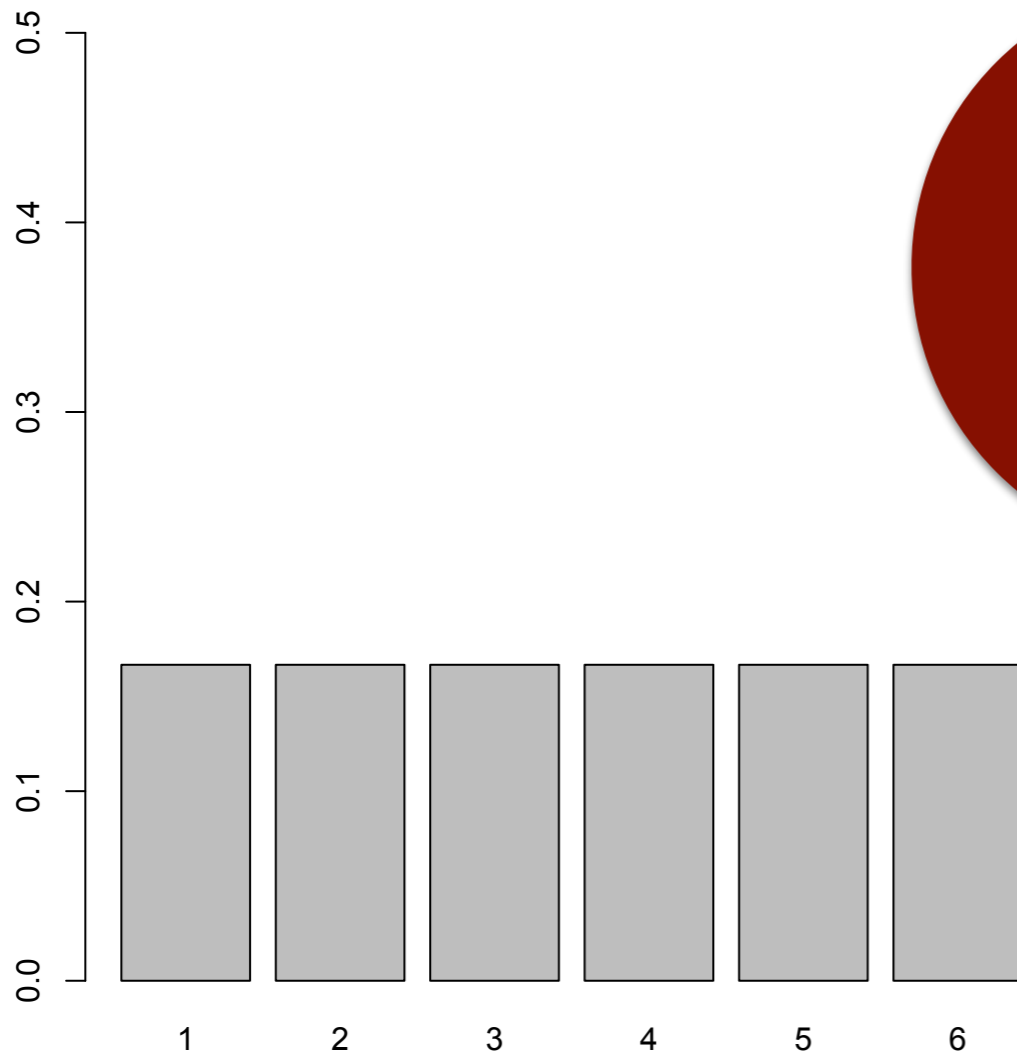


Probability



fair

not fair

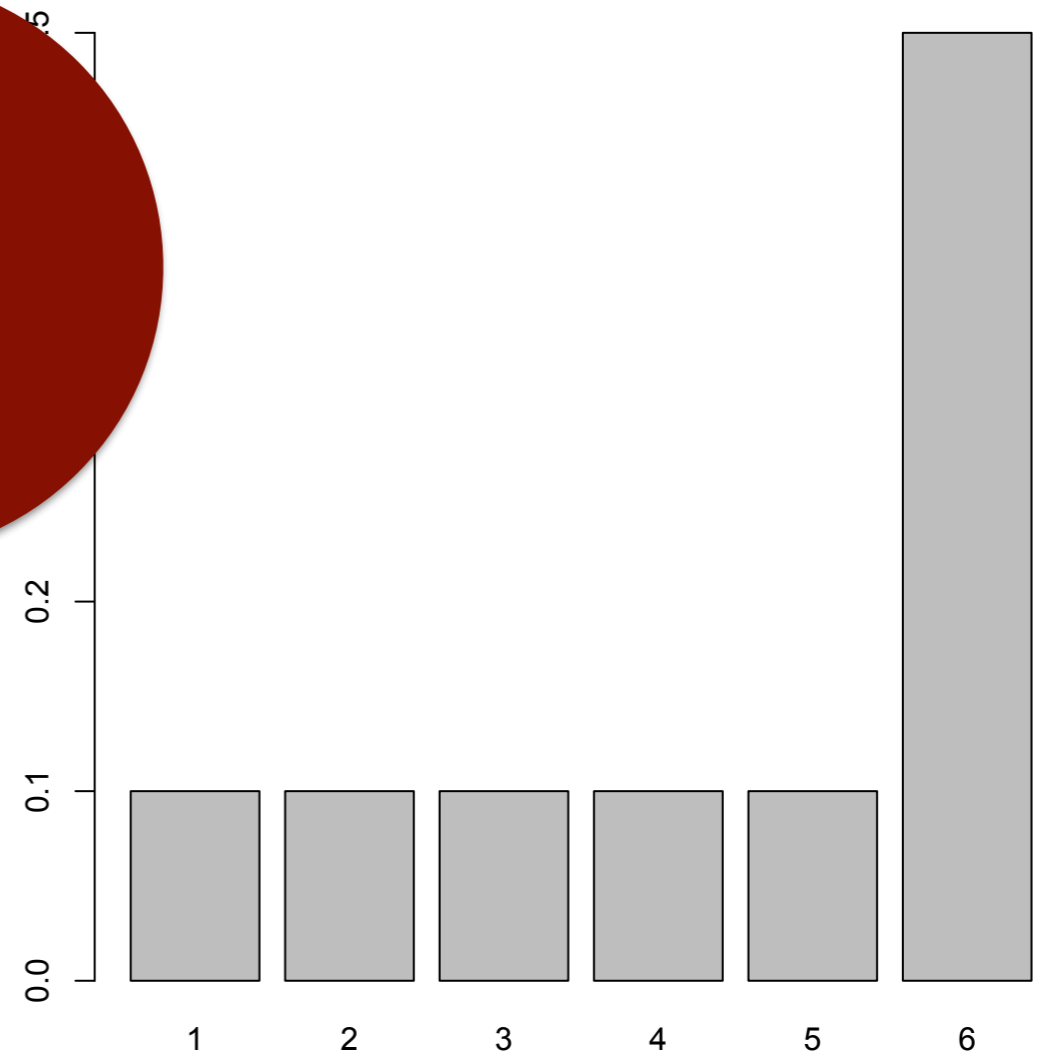
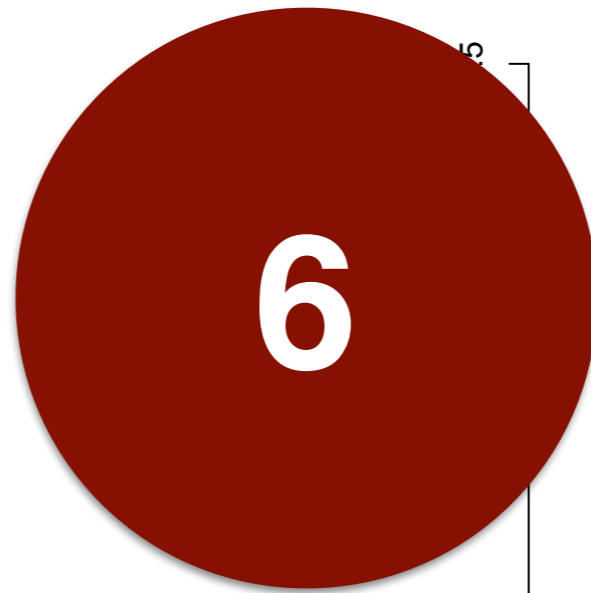
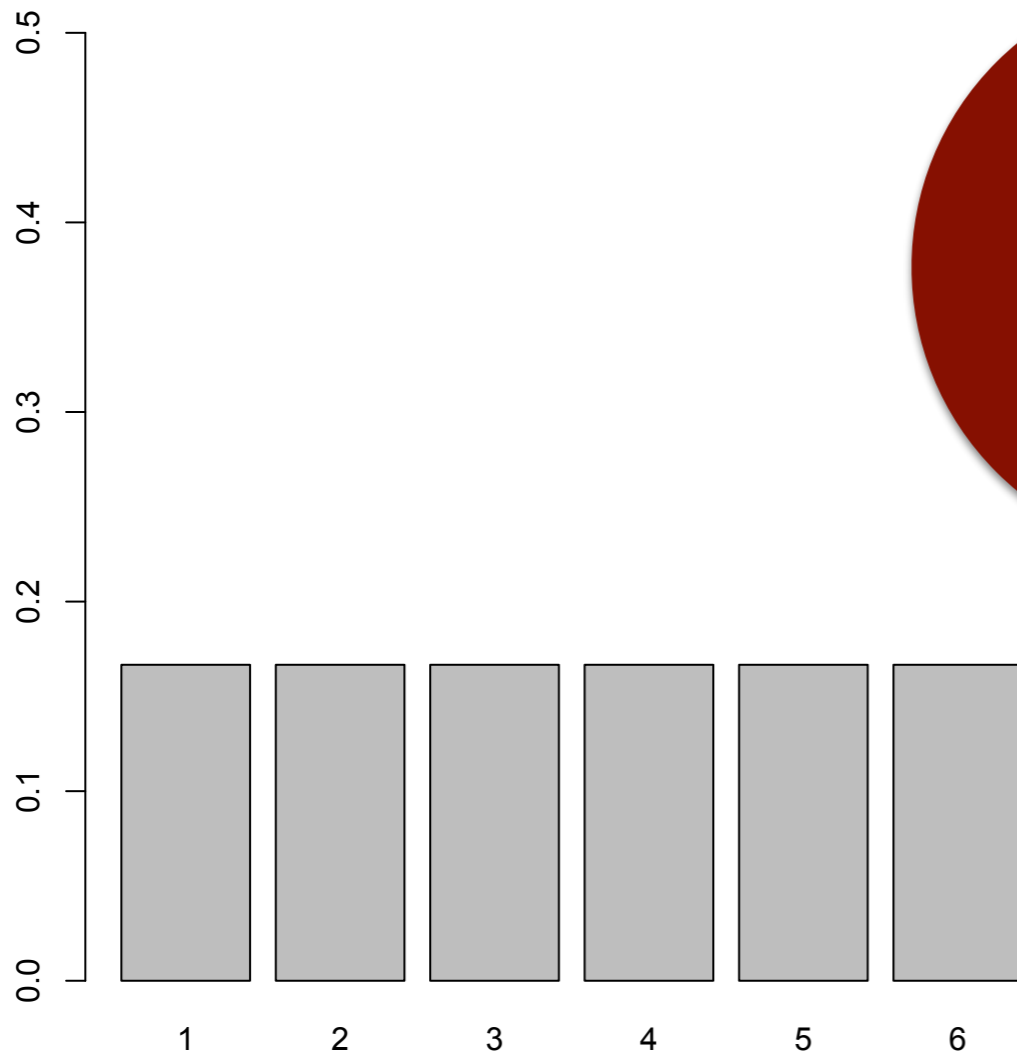


Probability

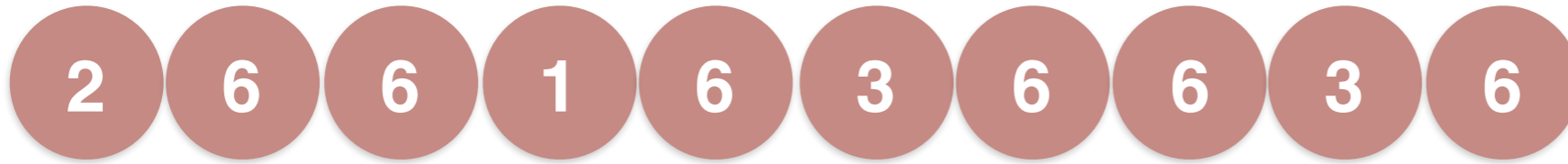


fair

not fair

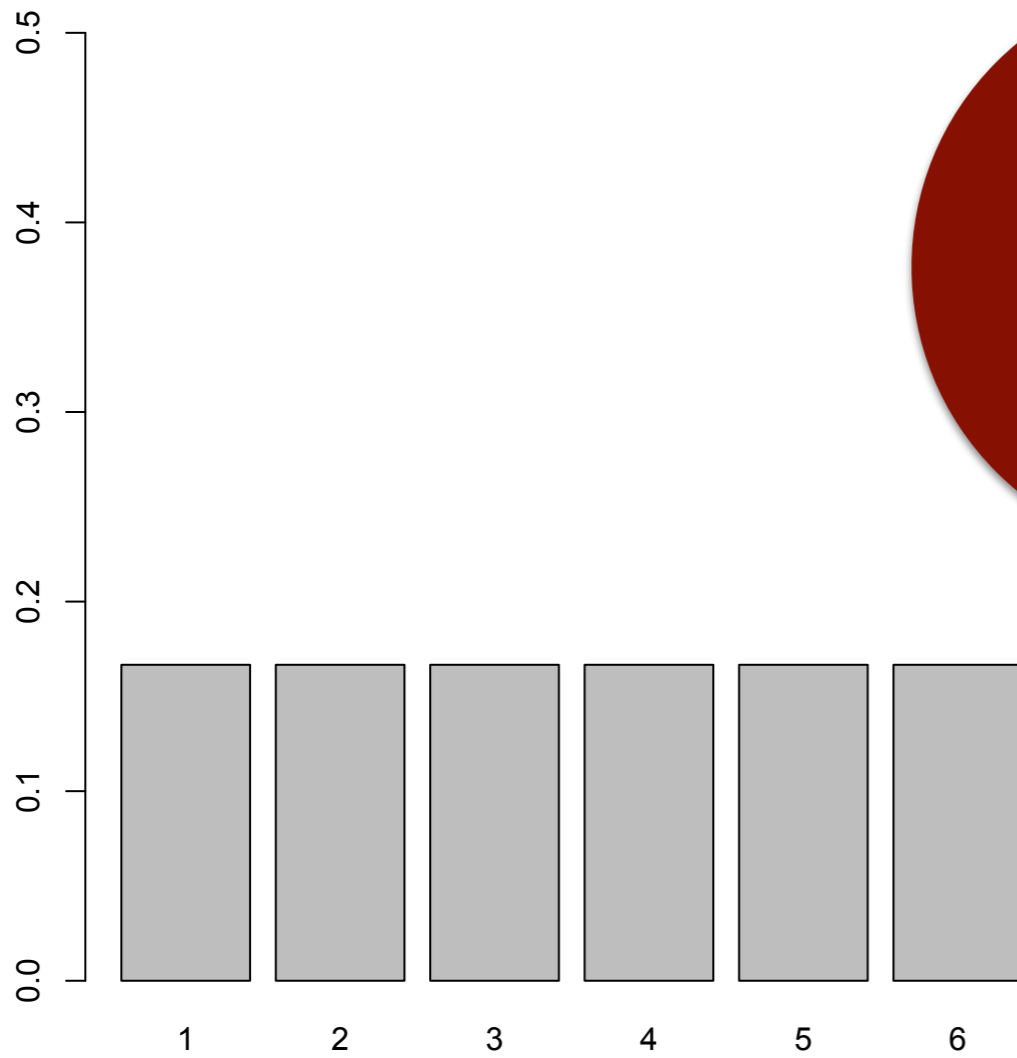
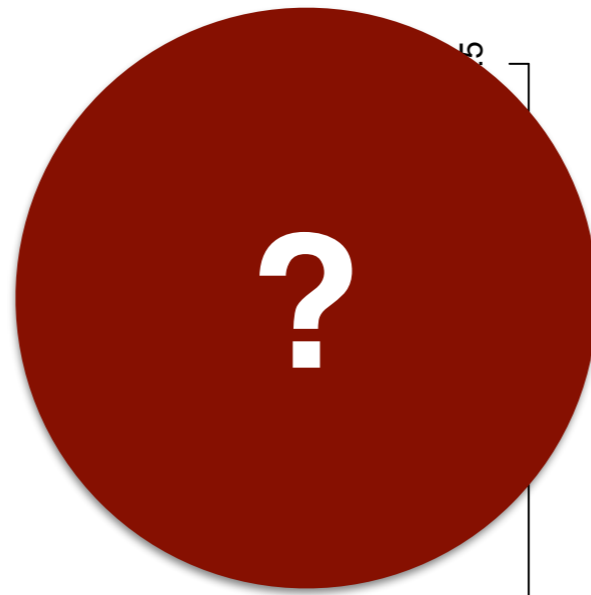


Probability

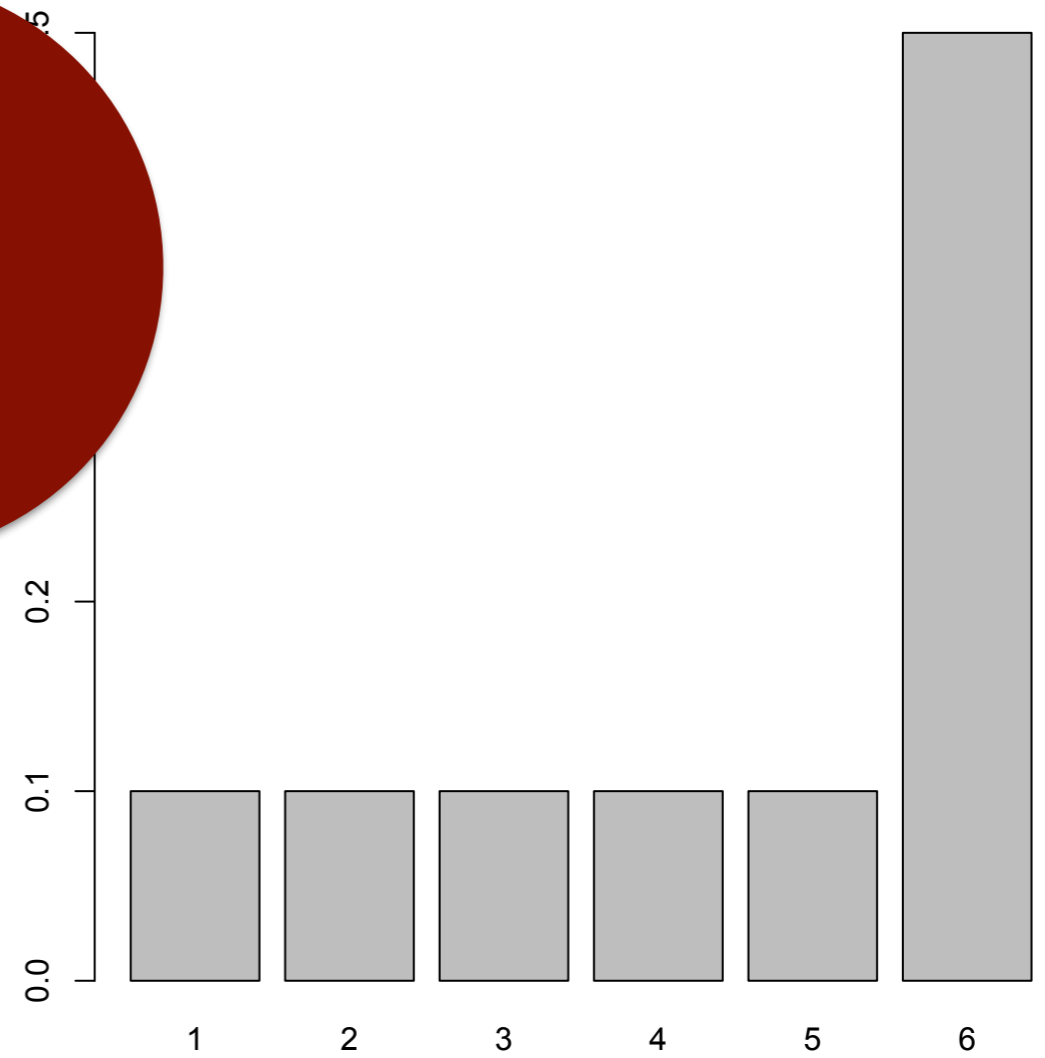


fair

not fair



1



15,625

Independence

- Two random variables are independent if:

$$P(A, B) = P(A) \times P(B)$$

- In general:

$$P(x_1, \dots, x_n) = \prod_{i=1}^n P(x_i)$$

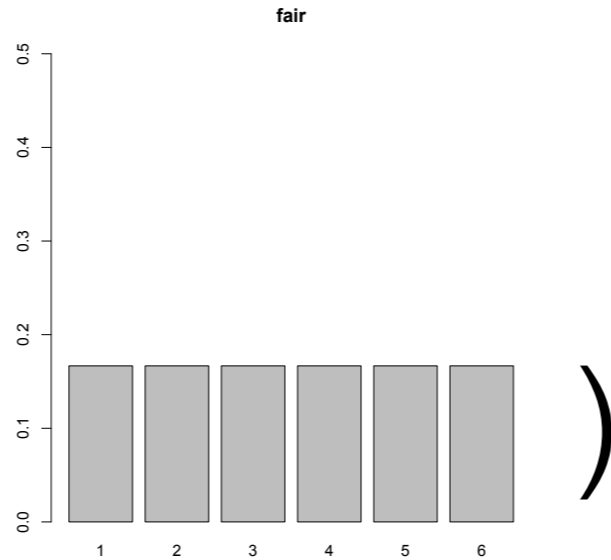
- Information about one random variable (B) gives no information about the value of another (A)

$$P(A) = P(A | B)$$

$$P(B) = P(B | A)$$

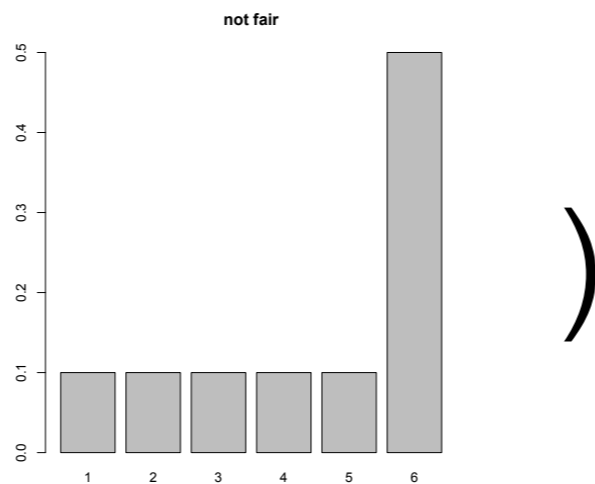
Data Likelihood

P(2 6 6 |



$$= .17 \times .17 \times .17$$
$$= 0.004913$$

P(2 6 6 |



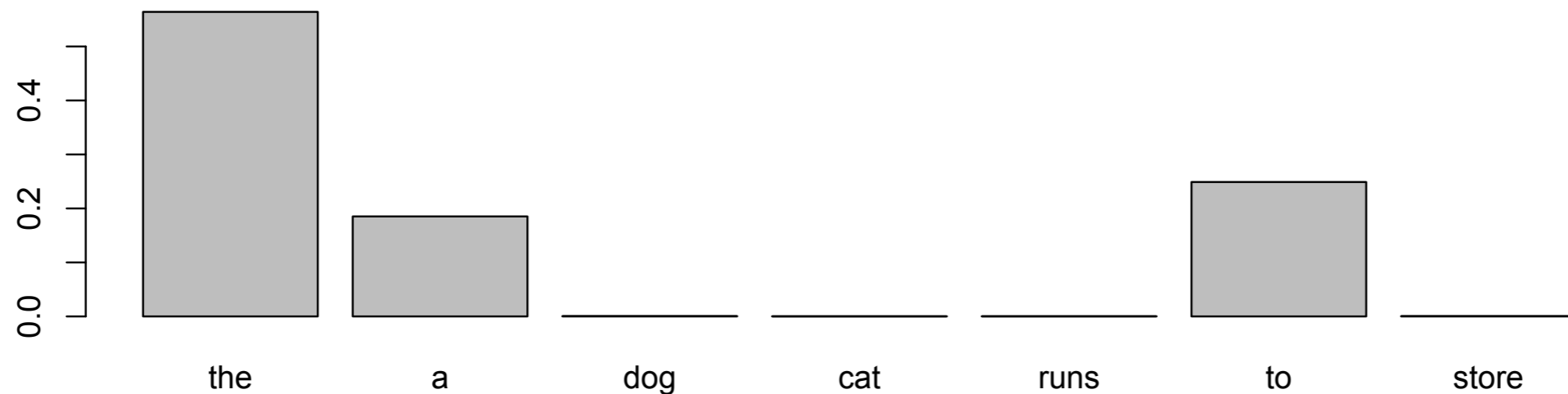
$$= .1 \times .5 \times .5$$
$$= 0.025$$

Data Likelihood

- The likelihood gives us a way of discriminating between possible alternative parameters, but also a strategy for picking a single best* parameter among all possibilities

Unigram probability

$X \in \{the, a, dog, cat, runs, to, store\}$



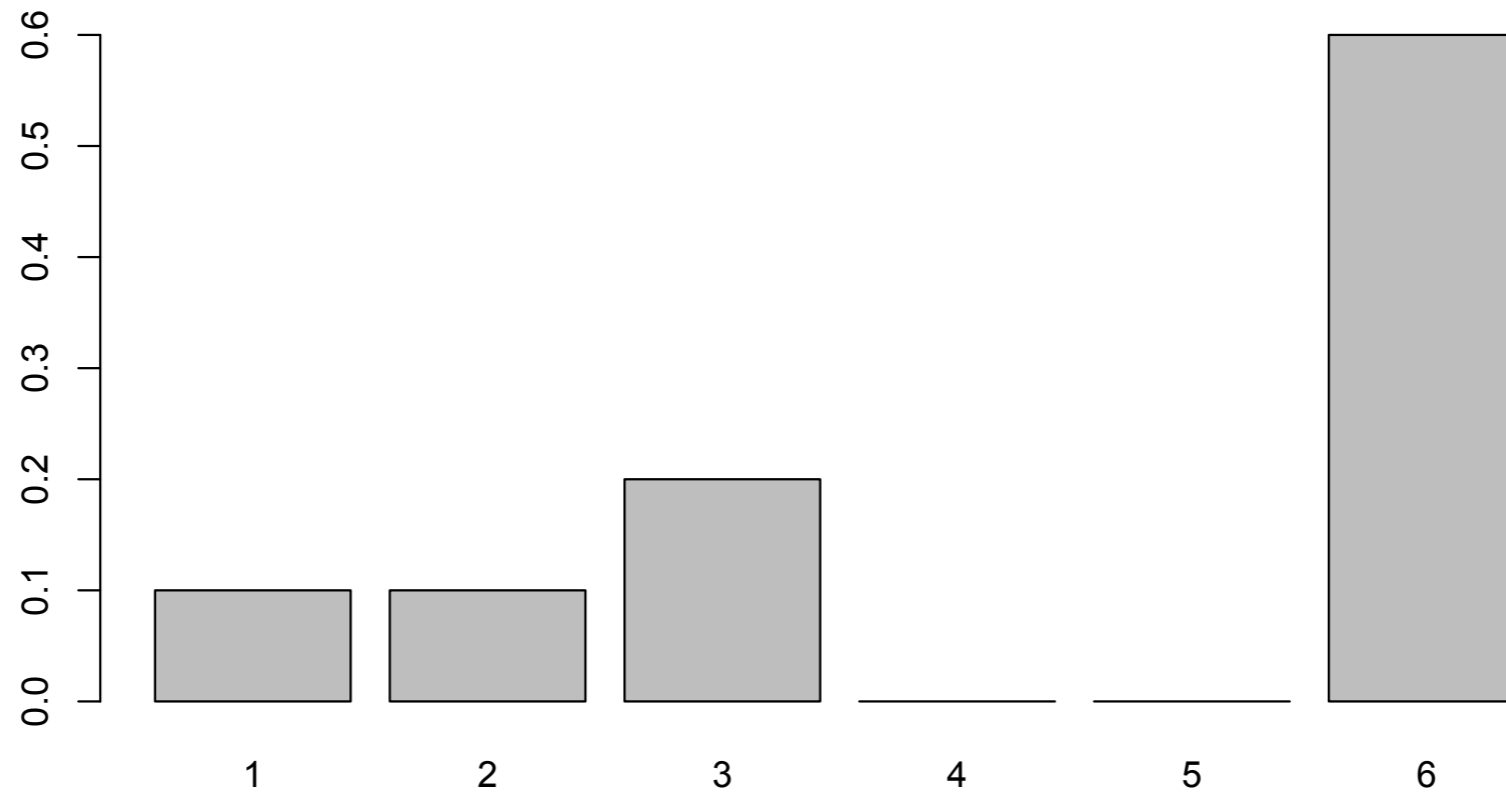
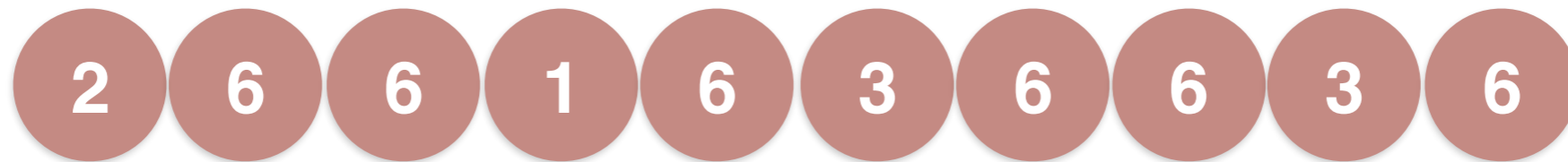
How do we calculate this?

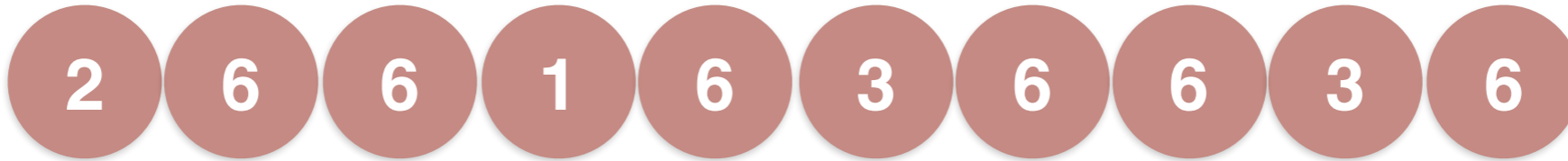
$$P(X=\text{"the"}) = 28/536 = .052$$

Maximum Likelihood Estimate

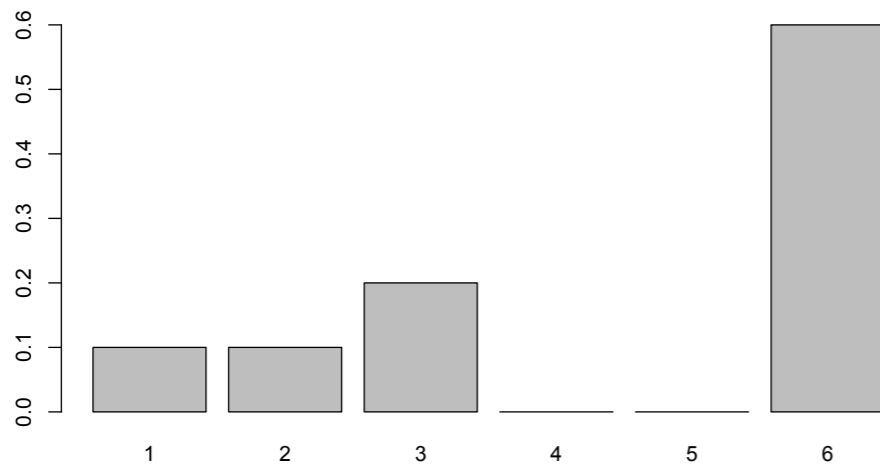
- This is a maximum likelihood estimate for $P(X)$; the parameter values for which the data we observe (X) is *most likely*.

Maximum Likelihood Estimate



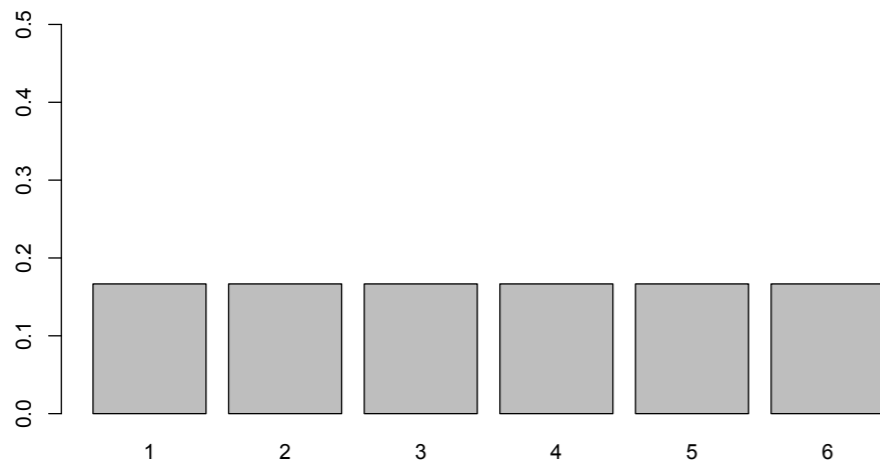


θ_1



$$P(X | \theta_1) = 0.0000311040$$

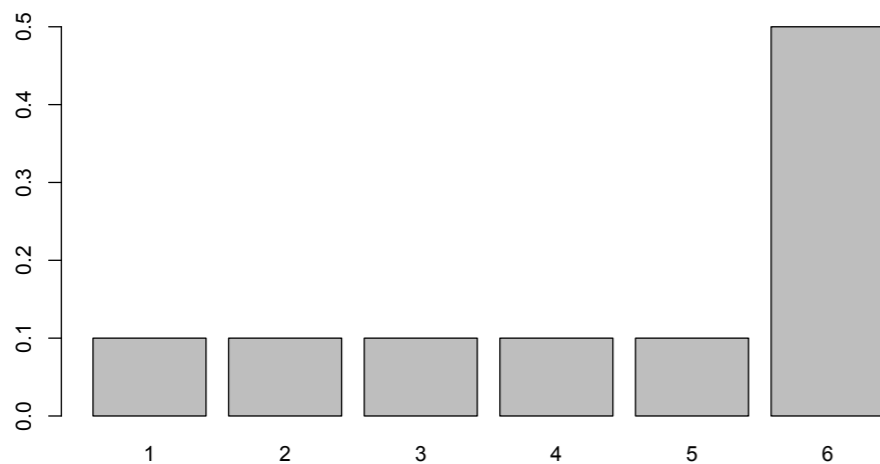
θ_2



$$P(X | \theta_2) = 0.00000000992$$

(313x less likely)

θ_3



$$P(X | \theta_3) = 0.0000031250$$

(10x less likely)

Conditional Probability

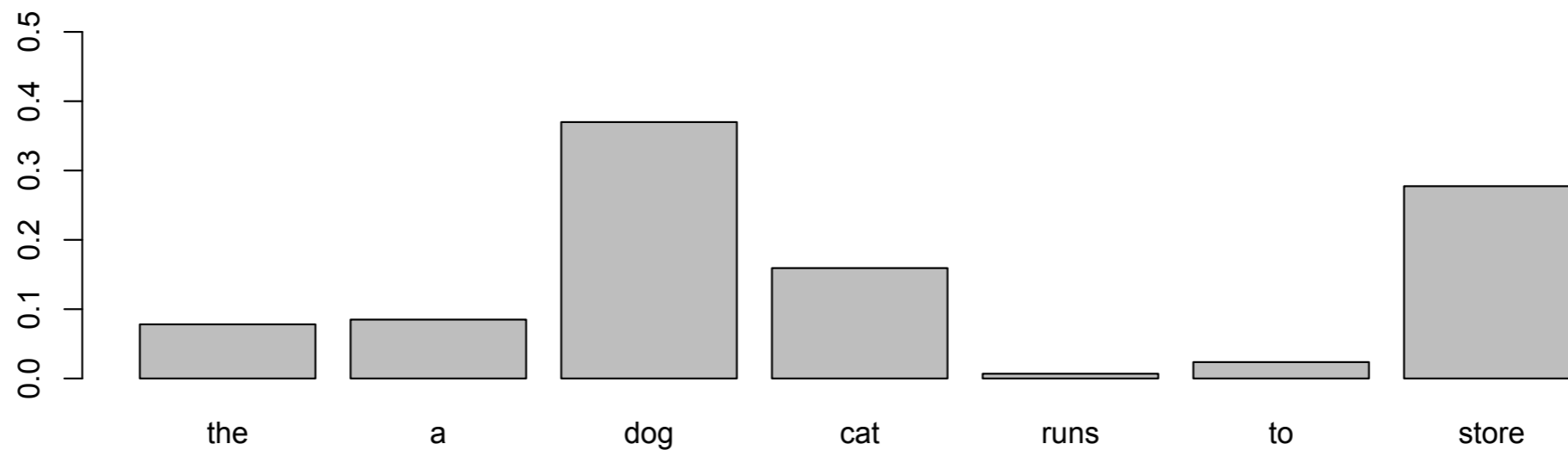
$$P(X = x|Y = y)$$

- Probability that one random variable takes a particular value *given* the fact that a different variable takes another

$$P(X_i = \text{dog}|X_{i-1} = \text{the})$$

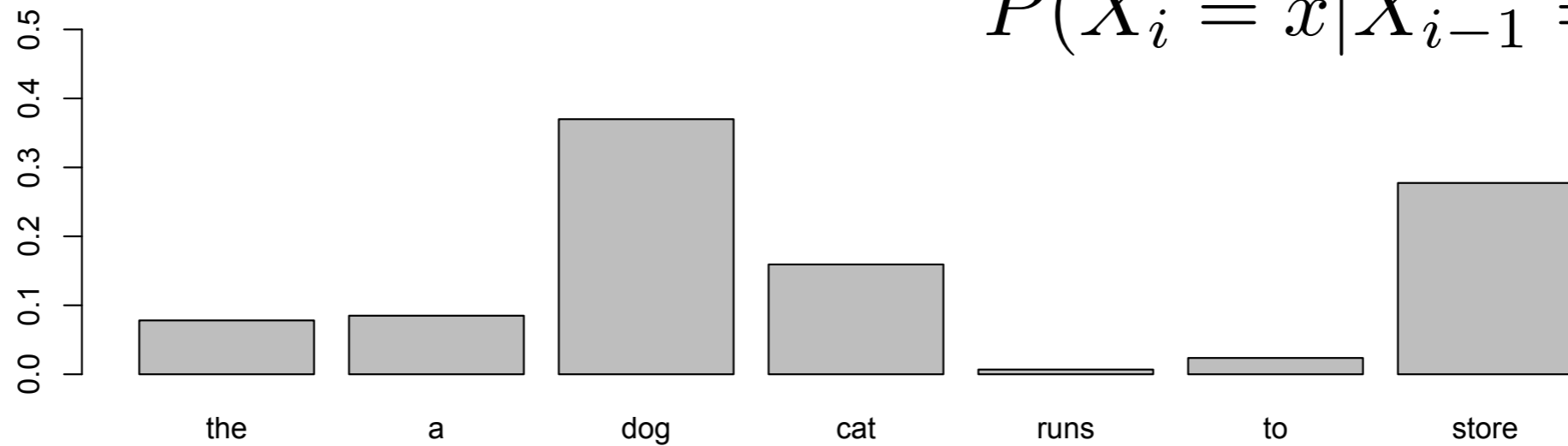
Conditional Probability

$$P(X_i = \text{dog} | X_{i-1} = \text{the})$$

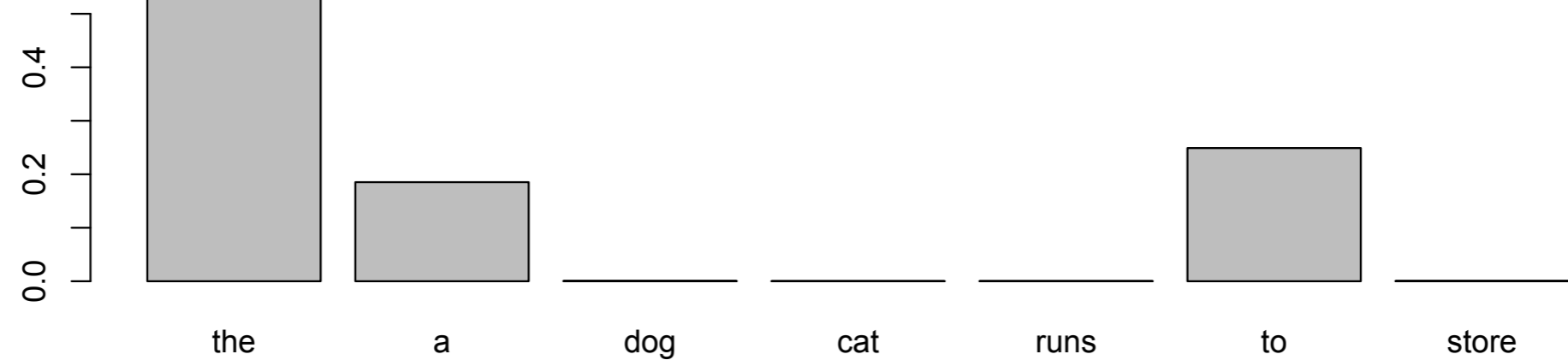


Conditional Probability

$$P(X_i = x | X_{i-1} = \textit{the})$$



$$P(X_i = x)$$



$$P(X_i = \text{"room"} | X_{i-1} = \text{"the"}) = 2/28 = .071$$

Conditional Probability

$P(X = \textit{vampire})$ vs. $P(X = \textit{vampire} | Y = \textit{horror})$

$P(X = \textit{manners} | Y = \textit{austen})$ vs. $P(X = \textit{whale} | Y = \textit{austen})$
0.00036 **0**

$P(X = \textit{manners} | Y = \textit{austen})$ vs. $P(X = \textit{manners} | Y = \textit{dickens})$
0.00036 = **6.7x times more than** **0.000053**

Authorship Attribution

“Mr. Collins was not a sensible man”



Independence Assumption

“Mr. Collins was not a sensible man”



$$P(x_1 = \text{Mr.}, x_2 = \text{Collins}) = P(x_1 = \text{Mr.}) \times P(x_2 = \text{Collins})$$

This is certainly untrue in this case, because the presence of **Mr.** makes **Collins** more likely (they are dependent)

Independence Assumption

“Mr. Collins was not a sensible man”



We will assume the features are independent:

$$P(x_1, x_2, x_3, x_4, x_6, x_7 \mid c) = P(x_1 \mid c)P(x_2 \mid c) \dots P(x_7 \mid c)$$

$$P(x_i \dots x_n \mid c) = \prod_{i=1}^N P(x_i \mid c)$$

A simple classifier

“Mr. Collins was not a sensible man”

Austen		Dickens	
$P(X=\text{Mr.} \mid Y=\text{Austen})$	0.0084	$P(X=\text{Mr.} \mid Y=\text{Dickens})$	0.00421
$P(X=\text{Collins} \mid Y=\text{Austen})$	0.00036	$P(X=\text{Collins} \mid Y=\text{Dickens})$	0.000016
$P(X=\text{was} \mid Y=\text{Austen})$	0.01475	$P(X=\text{was} \mid Y=\text{Dickens})$	0.015043
$P(X=\text{not} \mid Y=\text{Austen})$	0.01145	$P(X=\text{not} \mid Y=\text{Dickens})$	0.00547
$P(X=\text{a} \mid Y=\text{Austen})$	0.01591	$P(X=\text{a} \mid Y=\text{Dickens})$	0.02156
$P(X=\text{sensible} \mid Y=\text{Austen})$	0.00025	$P(X=\text{sensible} \mid Y=\text{Dickens})$	0.00005
$P(X=\text{man} \mid Y=\text{Austen})$	0.00121	$P(X=\text{man} \mid Y=\text{Dickens})$	0.001707

A simple classifier

“Mr. Collins was not a sensible man”

$P(X = \text{“Mr. Collins was not a sensible man”} \mid Y = \text{Austen})$

$$\begin{aligned} &= P(\text{“Mr”} \mid \text{Austen}) \times P(\text{“Collins”} \mid \text{Austen}) \times \\ &P(\text{“was”} \mid \text{Austen}) \times P(\text{“not”} \mid \text{Austen}) \dots \\ &= 0.000000022507322 (\approx \mathbf{2.3 \times 10^{-8}}) \end{aligned}$$

$P(X = \text{“Mr. Collins was not a sensible man”} \mid Y = \text{Dickens})$

$$\begin{aligned} &P(\text{“Mr”} \mid \text{Dickens}) \times P(\text{“Collins”} \mid \text{Dickens}) \times \\ &P(\text{“was”} \mid \text{Dickens}) \times P(\text{“not”} \mid \text{Dickens}) \dots \\ &= 0.000000002078906 (\approx \mathbf{2.1 \times 10^{-9}}) \end{aligned}$$

A simple classifier

- The classifier we just specified is a maximum likelihood classifier, where compare the **likelihood** of the data under each class and choose the class with the highest likelihood

Likelihood: probability of data
(here, under class y)

$$P(X = x_1 \dots x_n \mid Y = y)$$

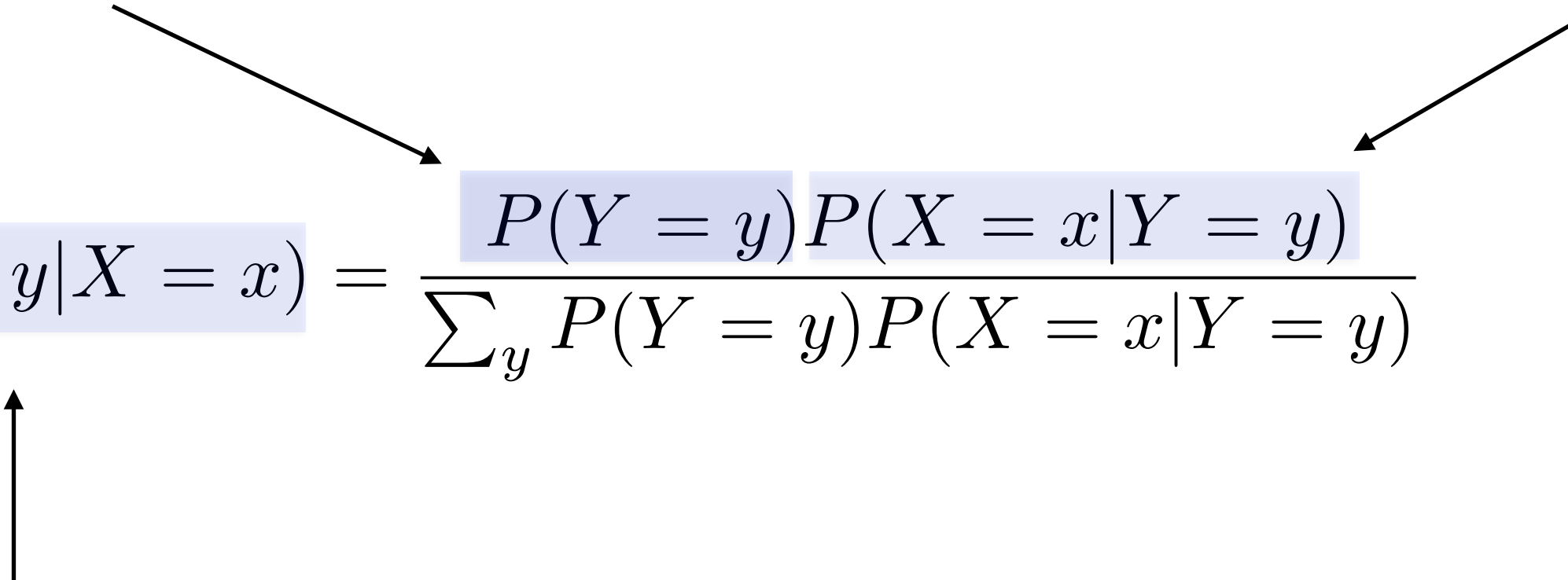
Prior probability of class y

$$P(Y = y)$$

Bayes' Rule

Prior belief that $Y = y$
(before you see any data)

Likelihood of the data
given that $Y=y$


$$P(Y = y|X = x) = \frac{P(Y = y)P(X = x|Y = y)}{\sum_y P(Y = y)P(X = x|Y = y)}$$

Posterior belief that $Y=y$ given that $X=x$

Bayes' Rule

Prior belief that $Y = \text{Austen}$
(before you see any data)

Likelihood of "Mr. Collins
was not a sensible man"
given that $Y = \text{Austen}$

$$P(Y = y | X = x) = \frac{P(Y = y)P(X = x | Y = y)}{\sum_y P(Y = y)P(X = x | Y = y)}$$

Posterior belief that $Y = \text{Austen}$ given that
 $X = \text{"Mr. Collins was not a sensible man"}$

This sum ranges over
 $y = \text{Austen} + y = \text{Dickens}$
(so that it sums to 1)

Likelihood: probability of data
(here, under class y)

$$P(X = x_1 \dots x_n \mid Y = y)$$

Prior probability of class y

$$P(Y = y)$$

Posterior belief in the probability
of class y after seeing data

$$P(Y = y \mid X = x_1 \dots x_n)$$

Naive Bayes Classifier

$$\frac{P(Y = Austen)P(X = \text{“Mr...”}|Y = Austen)}{P(Y = Austen)P(X = \text{“Mr...”}|Y = Austen) + P(Y = Dickens)P(X = \text{“Mr...”}|Y = Dickens)}$$

Let's say $P(Y=Austen) = P(Y=Dickens) = 0.5$
(i.e., both are equally likely a priori)

$$= \frac{0.5 \times (2.3 \times 10^{-8})}{0.5 \times (2.3 \times 10^{-8}) + 0.5 \times (2.1 \times 10^{-9})}$$

$$P(Y = Austen|X = \text{“Mr...”}) = 91.5\%$$

$$P(Y = Dickens|X = \text{“Mr...”}) = 8.5\%$$

Taxicab Problem

“A cab was involved in a hit and run accident at night. Two cab companies, the Green and the Blue, operate in the city. You are given the following data:

- 85% of the cabs in the city are Blue.
 - A witness identified the cab involved in the accident as Blue. The witness was 80% accurate in identifying the color of the cab.
- “Base rate fallacy”
Don't ignore prior information!

What is the probability that the cab involved in the accident was Blue rather than Green knowing that this witness identified it as Blue?”

(Tversky & Kahneman 1981)

Prior Belief

- Now let's assume that Dickens published 1000 times more books than Austen.
- $P(Y = \text{Austen}) = 0.000999$
- $P(Y = \text{Dickens}) = 0.999001$

$$\frac{0.000999 \times (2.3 \times 10^{-8})}{0.000999 \times (2.3 \times 10^{-8}) + 0.999001 \times (2.1 \times 10^{-9})}$$

$$P(Y = \text{Austen} | X) = 0.011$$

$$P(Y = \text{Dickens} | X) = 0.989$$

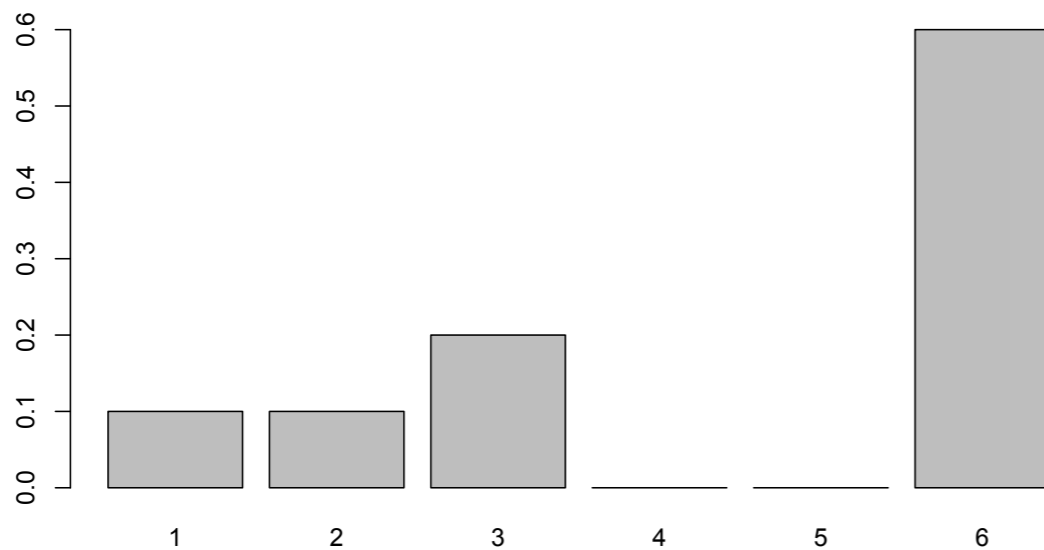
Priors

- Priors can be informed (reflecting expert knowledge) but in practice, but priors in Naive Bayes are often simply estimated from training data

$$P(Y = \text{Austen}) = \frac{\# \text{ of Austen texts}}{\# \text{ of total texts}}$$

Smoothing

- Maximum likelihood estimates can fail miserably when features are never observed with a particular class.



What's the probability of:



Smoothing

- One solution: add a little probability mass to every element.

maximum likelihood
estimate

$$P(x_i | y) = \frac{n_{i,y}}{n_y}$$

$n_{i,y}$ = count of word i in class y
 n_y = number of words in y
 V = size of vocabulary

smoothed estimates

$$P(x_i | y) = \frac{n_{i,y} + a}{n_y + Va}$$

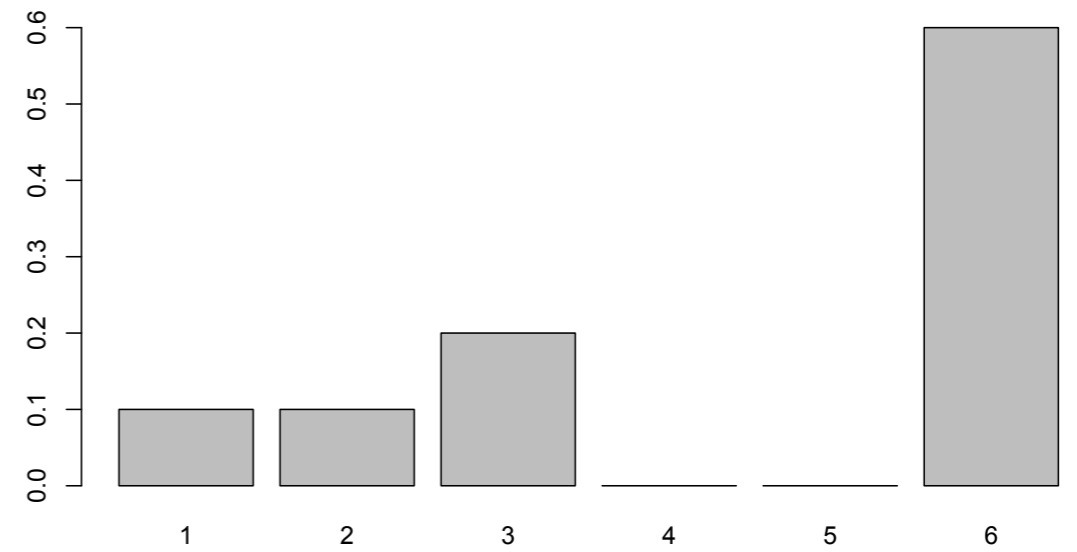
same a for all x_i

$$P(x_i | y) = \frac{n_{i,y} + a_i}{n_y + \sum_{j=1}^V a_j}$$

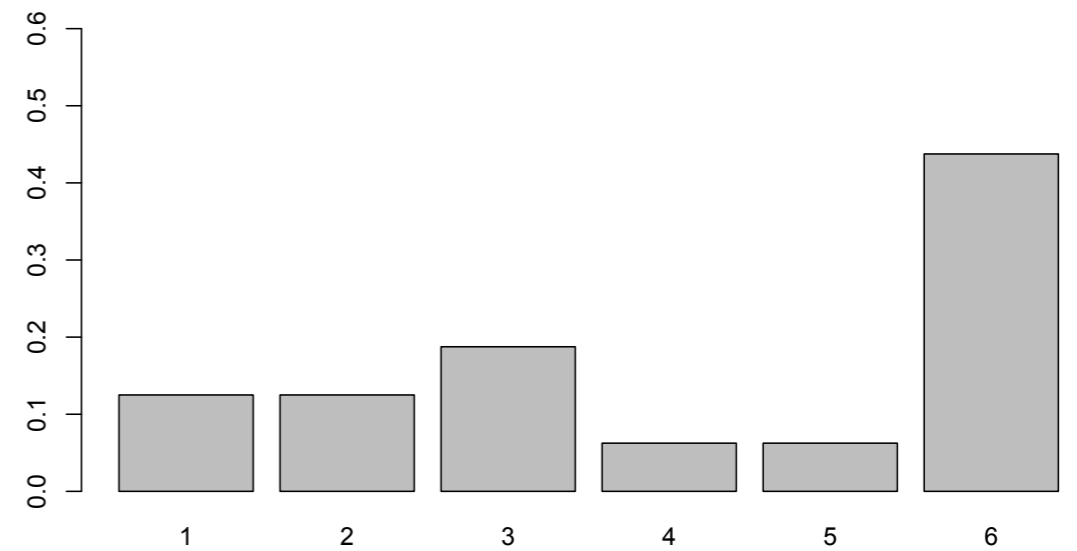
possibly different a for each x_i

Smoothing

MLE



smoothing with $\alpha = 1$



Naive Bayes training

Training a Naive Bayes classifier consists of estimating these two quantities from training data for all classes y

$$P(Y = y|X = x) = \frac{P(Y = y)P(X = x|Y = y)}{\sum_y P(Y = y)P(X = x|Y = y)}$$

At test time, use those estimated probabilities to calculate the posterior probability of each class y and select the class with the highest probability

Naive Bayes

- We've just described Naive Bayes with a multinomial distribution, but any probability distribution can be modeled as well.

Probability distributions

Normal

Gamma

Poisson

Geometric

Exponential

Multinomial

Bernoulli

Beta

Binomial

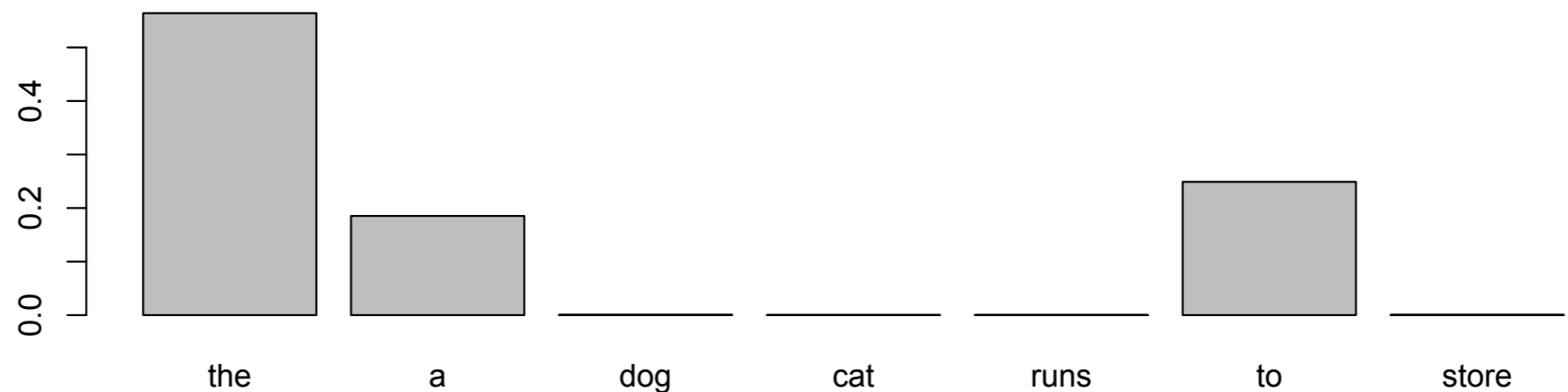
Uniform

Dirichlet

Multinomial

Discrete distribution for modeling count data (e.g., word counts; single parameter θ)

$\theta =$



the	a	dog	cat	runs	to	store
3	1	0	1	0	2	0
531	209	13	8	2	331	1

Multinomial

Maximum likelihood parameter estimate

$$\hat{\theta}_i = \frac{n_i}{N}$$

	the	a	dog	cat	runs	to	store
count n	531	209	13	8	2	331	1
θ	0.48	0.19	0.01	0.01	0.00	0.30	0.00

Bernoulli

- Binary event (true or false; $\{0, 1\}$) $P(x = 1 | p) = p$
- One parameter: p (probability of an event occurring) $P(x = 0 | p) = 1 - p$

Examples:

- Probability of a particular feature being true (e.g., self-reported location = Berkeley)

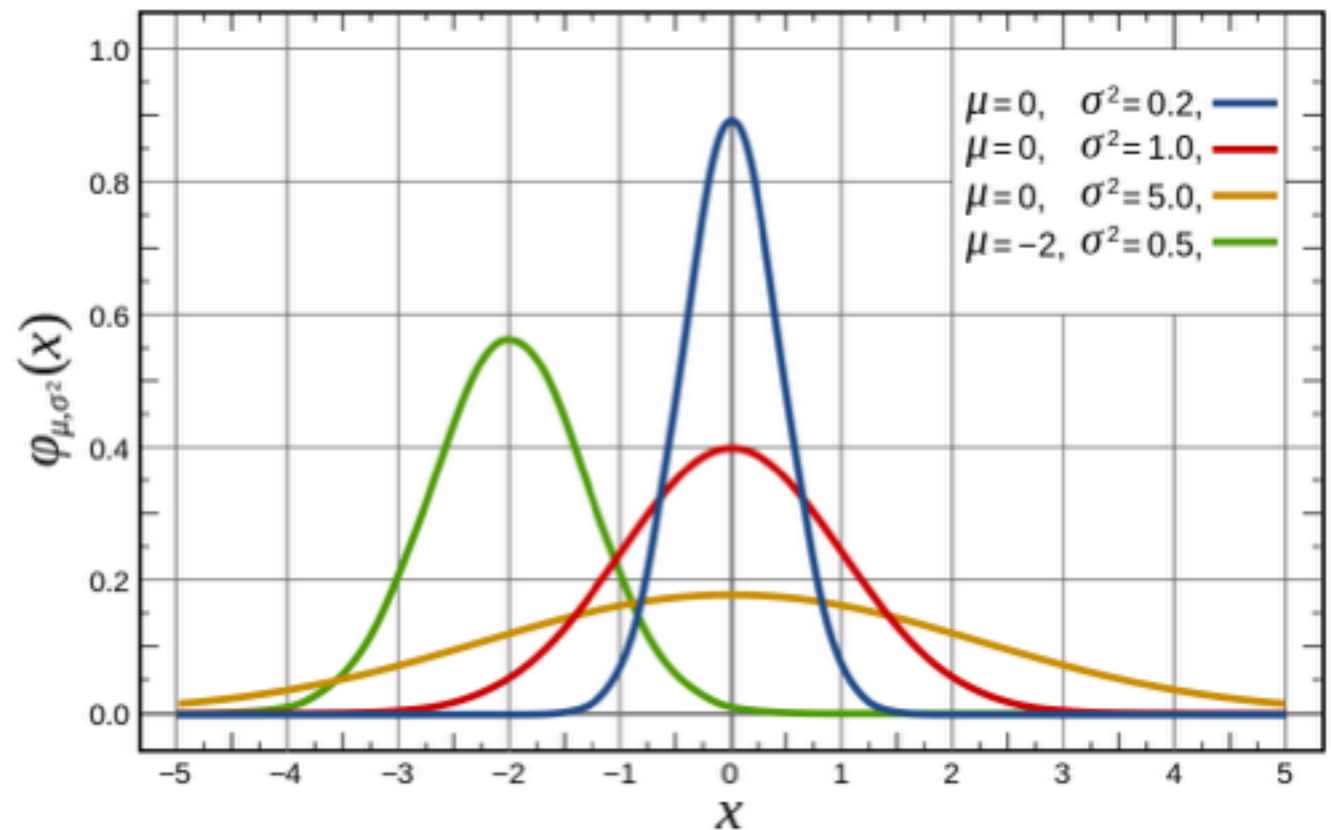
$$\hat{p}_{mle} = \frac{1}{N} \sum_{i=1}^N x_i$$

Normal

- continuous $(-\infty, \infty)$
- μ (mean) $(-\infty, \infty)$
- σ^2 (variance) > 0

Examples:

- Age
- Height



$$P(x = -2 \mid \mu = -2, \sigma^2 = 0.5) = 0.56$$

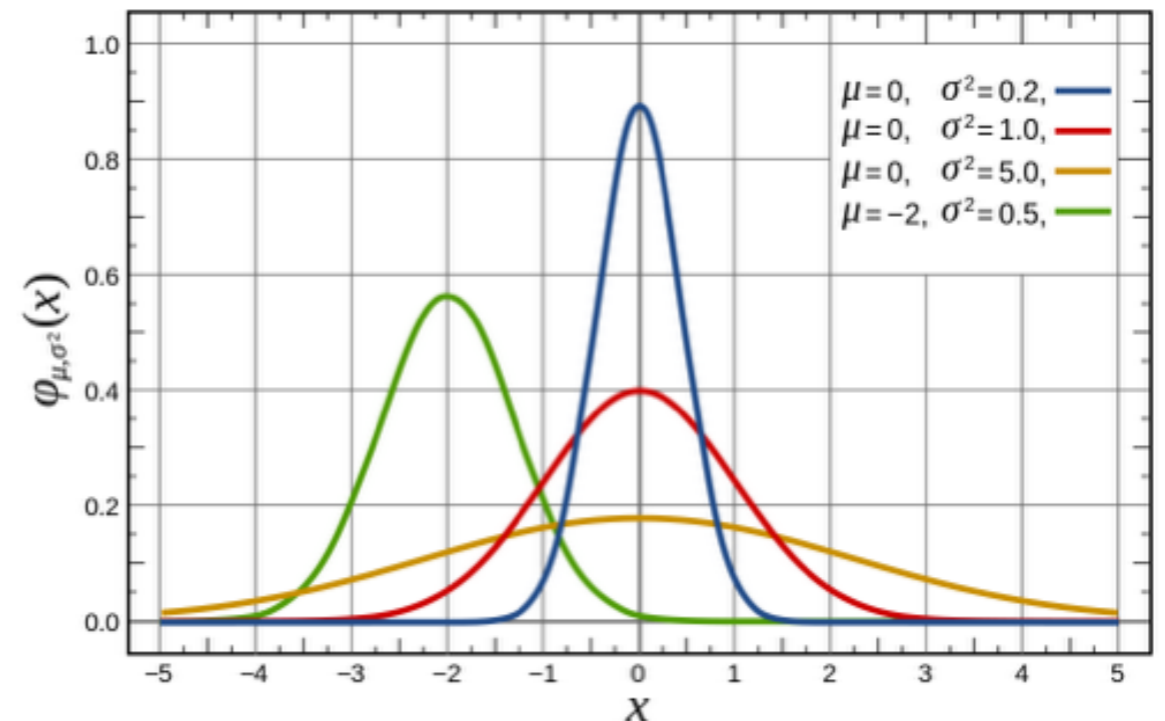
$$P(x = -2 \mid \mu = 0, \sigma^2 = 1) = 0.05$$

Normal

Maximum likelihood parameter estimates

$$\hat{\mu}_{mle} = \frac{1}{N} \sum_{i=1}^N x_i$$

$$\hat{\sigma}_{mle}^2 = \frac{1}{N} \sum_{i=1}^N (x_i - \bar{x})^2$$



Normal

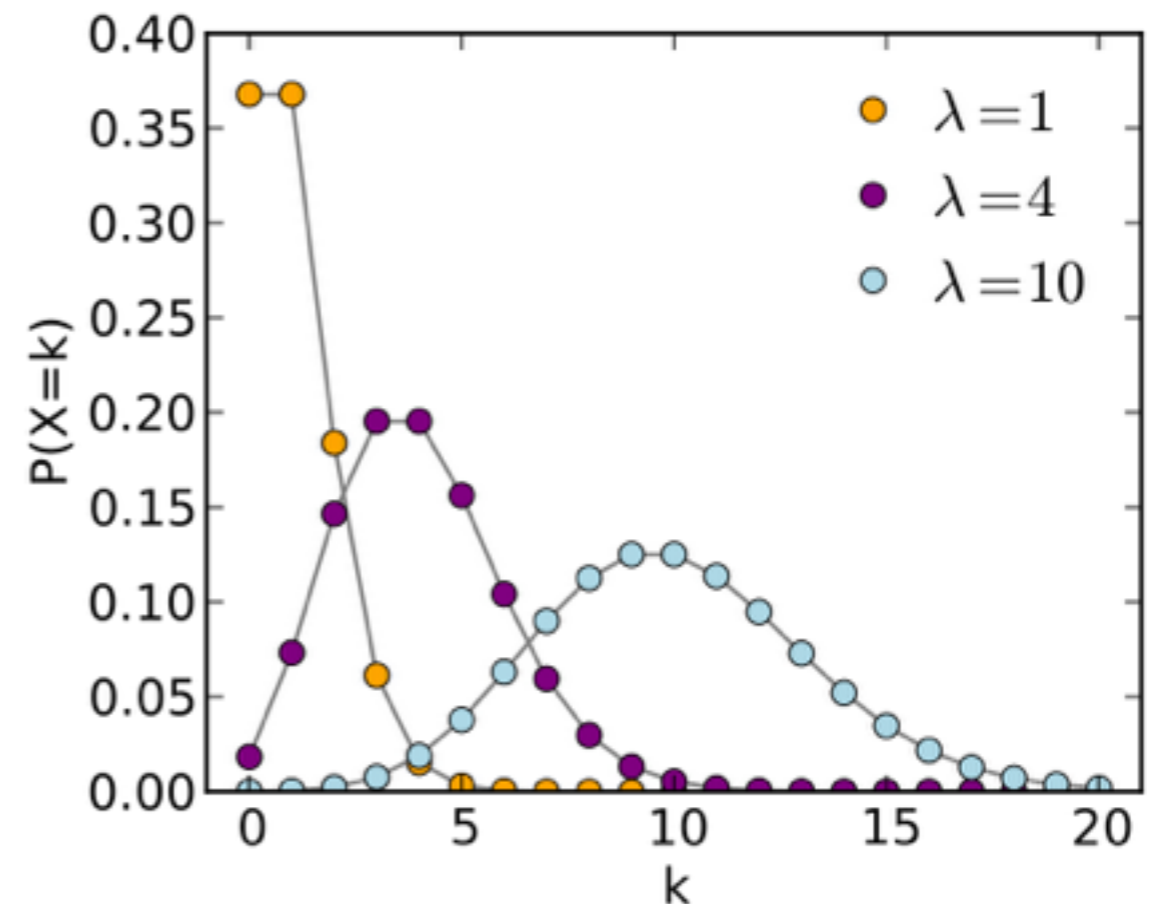
	Republican				Democrat				$\mu_{MLE,R}$	$\mu_{MLE,D}$
	X ₁	X ₂	X ₃	X ₄	X ₅	X ₆	X ₇	X ₈		
f ₁	3.4	-2.1	5.2	7.6	11.6	9.1	9.7	10.8	3.5	10.3
f ₂	-0.3	8.5	5.6	11.5	5.4	6.2	3.1	12.7	6.3	6.8
f ₃	-0.6	3.7	1.2	5.6	3.4	-4.4	8.0	6.2	2.5	3.3
f ₄	2.5	6.7	0.5	2.6	13.2	6.1	13.7	7.7	3.1	10.2
f ₅	7.0	5.0	5.6	16.3	15.4	14.9	2.3	6.3	8.5	9.7

Poisson

- discrete (0, 1, 2, ...)
- $\lambda > 0$
- Models the number of events within a fixed interval of time

Examples:

- Number of emails in one hour
- Number of children in family



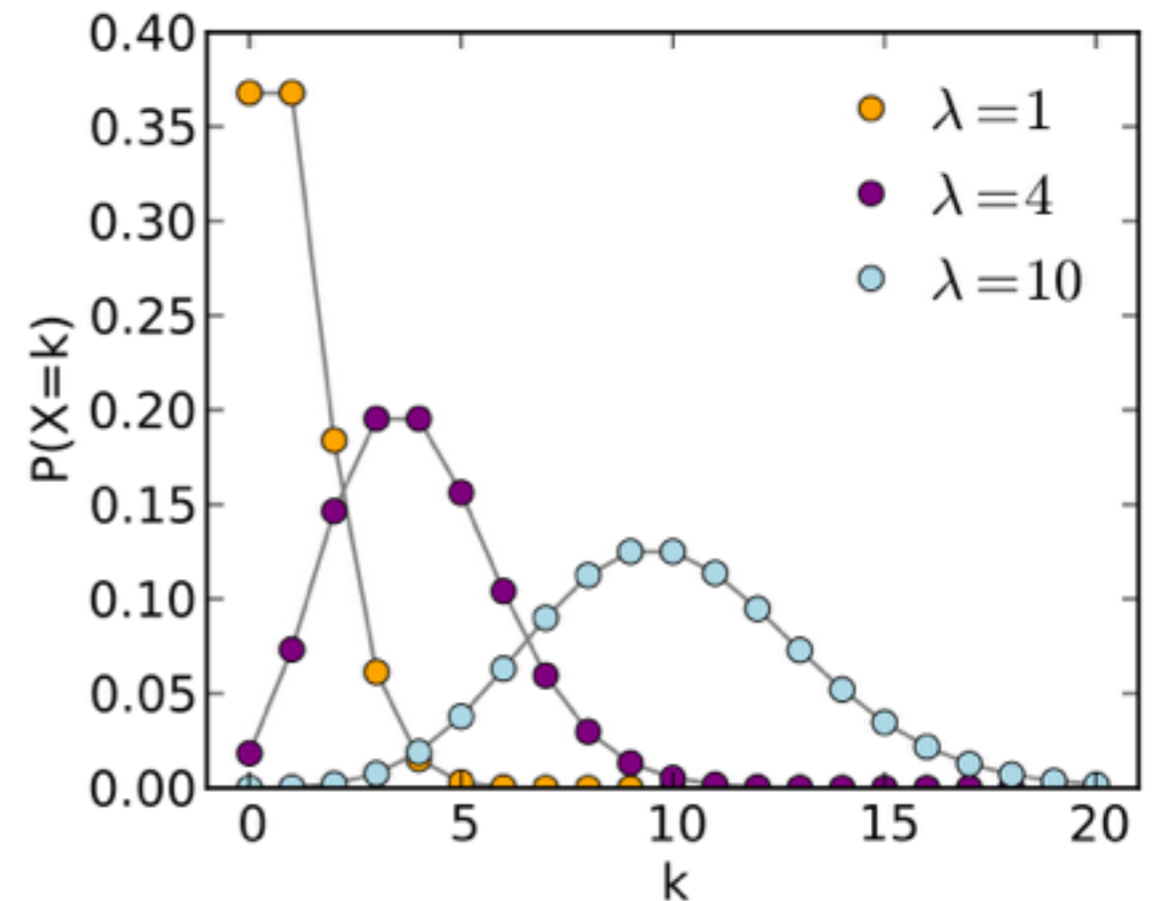
$$P(x = 4 | \lambda = 10) = 0.02$$

$$P(x = 4 | \lambda = 4) = 0.20$$

Poisson

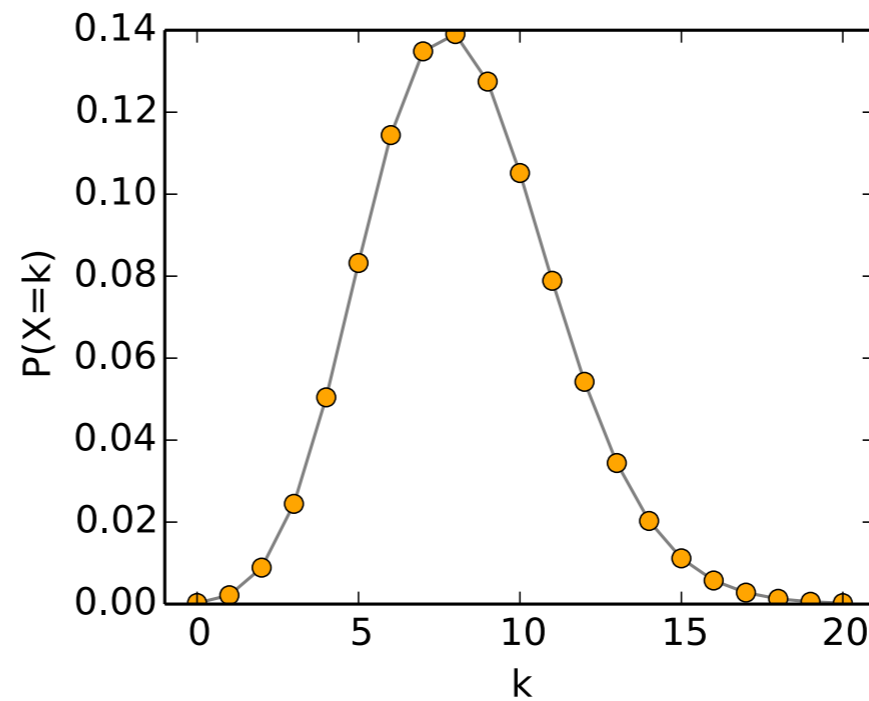
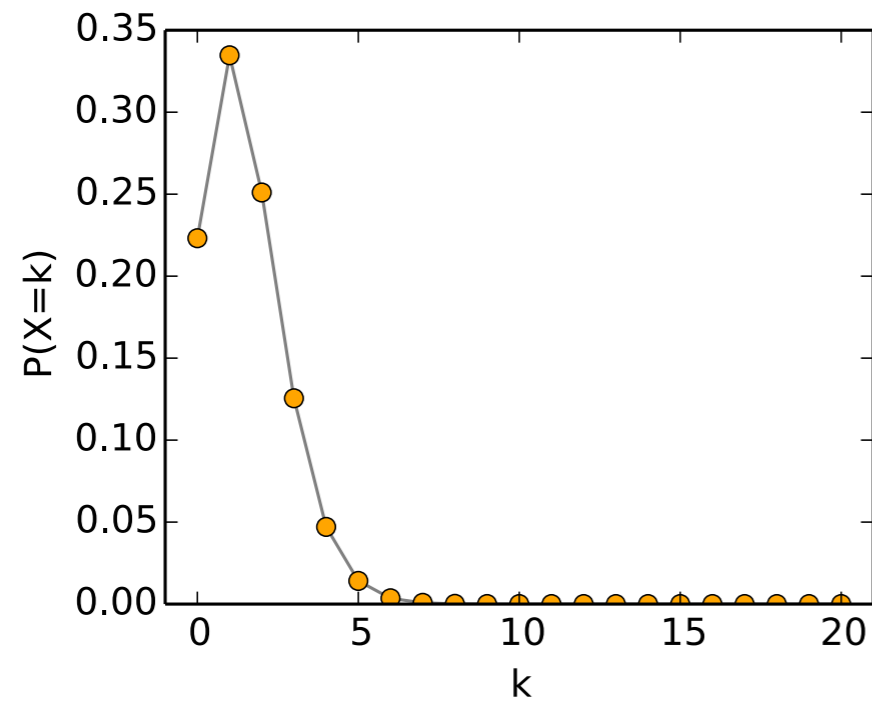
Maximum likelihood parameter estimate

$$\hat{\lambda} = \frac{1}{N} \sum_{i=1}^N x_i$$



Poisson

	Republican				Democrat				$\lambda_{MLE,R}$	$\lambda_{MLE,D}$
f_1	X_1	X_2	X_3	X_4	X_5	X_6	X_7	X_8	1.5	8.25
	1	2	2	1	6	10	8	9		



Feature

Value

Distribution?

follow clinton

0

follow trump

0

age

24

word counts in profile

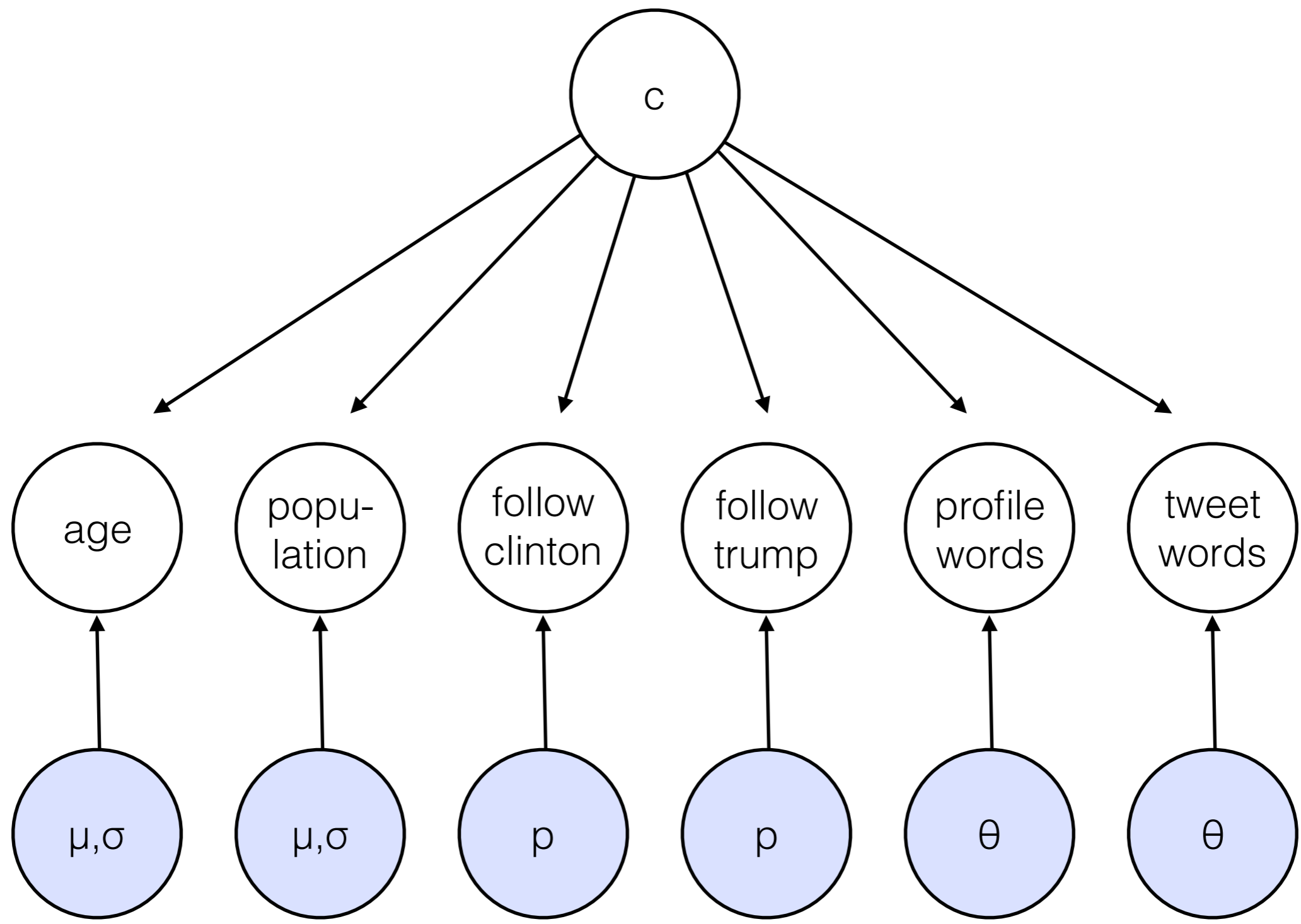
Berkeley, liberal,
runner

word counts in profile

the, election, a,
data, movies

population size of your city

116,000



Normal

Normal

Bernoulli

Bernoulli

Multinomial

Multinomial

$$P(X | c = \text{Dem}) = \prod_{i=1}^N P(X_i | c = \text{Dem})$$

$$= \text{Norm}(\textit{age} | \mu_{\textit{age},\textit{dem}}, \sigma_{\textit{age},\textit{dem}}^2)$$

$$\times \text{Norm}(\textit{population} | \mu_{\textit{population},\textit{dem}}, \sigma_{\textit{population},\textit{dem}}^2)$$

$$\times \text{Bernoulli}(\textit{followClinton} | p_{\textit{followClinton},\textit{dem}})$$

$$\times \text{Bernoulli}(\textit{followTrump} | p_{\textit{followTrump},\textit{dem}})$$

$$\times \text{Multinomial}(w_{\textit{profile}} | \theta_{\textit{profile},\textit{dem}})$$

$$\times \text{Multinomial}(w_{\textit{tweets}} | \theta_{\textit{tweets},\textit{dem}})$$

$$P(c = \text{Dem} \mid X) = \frac{P(c = \text{Dem}) \times P(X \mid c = \text{Dem})}{P(c = \text{Dem}) \times P(X \mid c = \text{Dem}) + P(c = \text{Rep}) \times P(X \mid c = \text{Rep})}$$

Authorship Attribution

Koppel et al. (2009), Computational Methods in Authorship Attribution (JASIST)

Representation

FW

A list of 512 function words, including conjunctions, prepositions, pronouns, modal verbs, determiners, and numbers (purely stylistic)

POS

Thirty-eight part-of-speech unigrams and 1,000 most common bigrams using the Brill (1992) part-of-speech tagger (purely stylistic)

SFL

All 372 nodes in SFL trees for conjunctions, prepositions, pronouns, and modal verbs (purely stylistic)

CW

The 1,000 words with highest information gain (Quinlan, 1986) in the training corpus among the 10,000 most common words in the corpus

CNG

The 1,000 character trigrams with highest information gain in the training corpus among the 10,000 most common trigrams in the corpus (cf. Keselj, 2003)

Models

NB WEKA's implementation (Witten & Frank, 2000) of Naïve Bayes (Lewis, 1998) with Laplace smoothing

J4.8 WEKA's implementation of the J4.8 decision tree method (Quinlan, 1986) with no pruning

RNW Our implementation of a version of Littlestone's (1988) Winnow algorithm, generalized to handle real-valued features and more than two classes (Schler, 2007)

BMR Genkin et al.'s (2006) implementation of Bayesian multiclass regression

SMO Weka's implementation of Platt's (1998) SMO algorithm for SVM with a linear kernel and default settings

Accuracy

TABLE 2. Accuracy on test set attribution for a variety of feature sets and learning algorithms applied to authorship classification for the e-mail corpus.

Features/learner	NB (%)	J4.8 (%)	RMW (%)	BMR (%)	SMO (%)
FW	60.2	58.7	66.1	68.2	63.8
POS	61.0	59.0	66.1	66.3	67.1
FW + POS	65.9	61.6	68.0	67.8	71.7
SFL	57.2	57.2	65.6	67.2	62.7
CW	67.1	66.9	74.9	78.4	74.7
CNG	72.3	65.1	73.1	80.1	74.9
CW + CNG	73.2	68.9	74.2	83.6	78.2

TABLE 4. Accuracy test set attribution for a variety of feature sets and learning algorithms applied to authorship classification for the blog corpus.

Features/learner	NB (%)	J4.8 (%)	RMW (%)	BMR (%)	SMO (%)
FW	38.2	30.3	51.8	63.2	63.2
POS	34.0	30.3	51.0	63.2	60.6
FW + POS	47.0	34.3	62.3	70.3	72.0
SFL	35.4	36.3	61.4	69.2	71.7
CW	56.4	51.0	62.9	72.5	70.5
CNG	65.0	48.9	67.1	80.4	80.9
CW + CNG	69.9	51.6	75.4	86.1	85.7

Homework 2: Validity

HW 2, part I (everyone)

- Pick any of the academic papers assigned throughout this course (i.e., any text except ML and NCM) and discuss the ways in which it establishes (or fails to establish) the nine types of validity outlined in Krippendorff (2004):
 - Face validity
 - Social validity
 - Sampling validity
 - Semantic validity
 - Structural validity
 - Functional validity
 - Convergence validity
 - Discriminant validity
 - Predictive validity

Deliverable: one-page paper

HW 2, part IIa (implementation)

- The permutation test is a robust hypothesis test that doesn't require the parametric or large-sample assumptions of classical tests.
- The GitHub repository contains a dataset mapping movies (featurized through their genres and major actors who performed in them) to a binary decision of whether or not it was among the 25% highest grossing movies in that set.
- For each of the features x , consider the hypothesis "Movies with x are more likely to have a higher box office than those that do not." Code and execute a permutation test evaluating this hypothesis. Can the null hypothesis (that movies featuring x are not likely to have a higher box office than those that do not) be rejected with $p < 0.01$?

HW 2, part 1b (critique)

- The nine forms of validity outlined above represent a detailed taxonomy of the different ways in which an analysis can be judged for the extent to which it is valid. What other possible forms of validity are missing from this taxonomy that should be represented within it? Present an argument for a single form of validity—a.) why it captures an important dimension that should be assessed, b.) why you believe it's missing from Krippendorff's taxonomy, and c.) tangible ways in which an analysis could be assessed according to this dimension.
- Deliverable: one-page paper (single-spaced)