

COSC343: Artificial Intelligence

Lecture 4: Probability Theory: introduction

Lech Szymanski

Dept. of Computer Science, University of Otago

In today's lecture

- Mathematical framework for dealing with uncertainty
- Probability distributions
- Conditional probability
- Independence
- Expectation

Probability Theory

- Fundamental concept underlying all machine learning is **uncertainty**
- Probability theory = mathematical framework for quantification and manipulation of uncertainty
- What's the *best* action to take, when the outcome is uncertain?

Defining a sample space

A **sample space** is a model of 'all possible ways the world can be'.

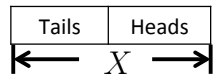
- Formally, it's the space of all possible values of the input and outputs to the function
- Each of these defines one dimension of the samples space
- Each possible combination is called a **sample point**

Formally, a **probability model** assigns a probability to each sample point in a sample space.

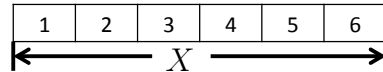
- Each probability is between 0 and 1 inclusive
- Probabilities for all points in the space sum to 1

Examples of sample spaces

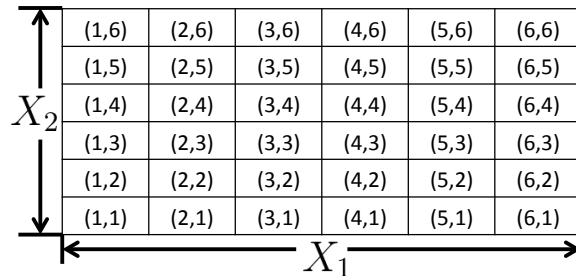
Coin toss



Dice roll



Double Dice roll



Probability distribution

Imagine we roll a single die. Our sample space has a single **random variable** (call it X), which has 6 possible values.

n	1	2	3	4	5	6
	$P(X = 1)$	$P(X = 2)$	$P(X = 3)$	$P(X = 4)$	$P(X = 5)$	$P(X = 6)$

We can estimate the probability at each point by generating a training set of N die rolls and using **relative frequencies** of events in this set

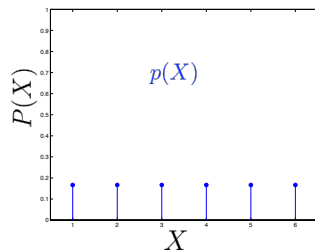
$$P(X = n) = \frac{\text{count}(X = n)}{N}$$

A simple probability model

A probability model induces a **probability distribution** for each possible value of the random variable.

- This distribution is a function, whose domain is all possible value for the random variable, which returns probability for each possible value
- The distribution must sum to 1

n	1	2	3	4	5	6
$P(X)$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$	$\frac{1}{6}$



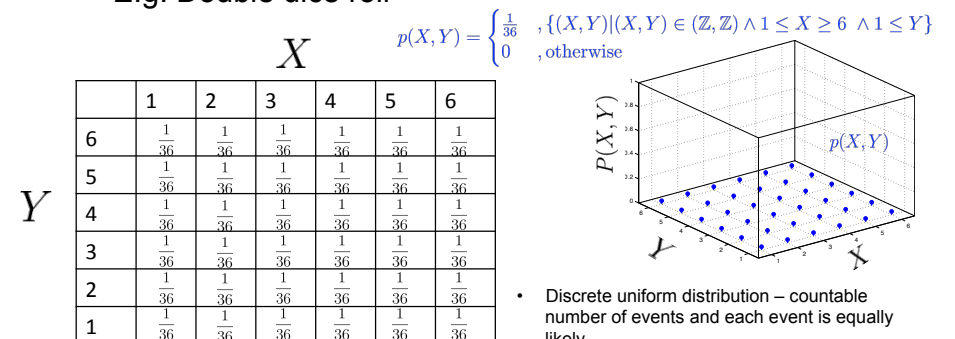
$$p(X) = \begin{cases} \frac{1}{6} & , \{X | X \in \mathbb{Z} \wedge 1 \leq X \leq 6\} \\ 0 & , \text{otherwise} \end{cases}$$

- Discrete uniform distribution – countable number of events and each event is equally likely

Joint distribution

A distribution function over two, or more, random variables is called a **joint distribution**

- E.g. Double dice roll



- Discrete uniform distribution – countable number of events and each event is equally likely

Some terminology

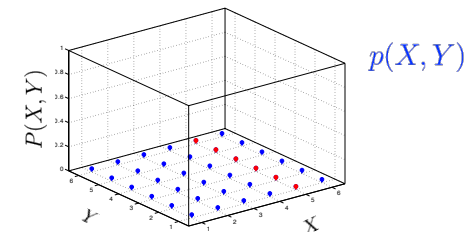
- An **event** is any subset of points in a sample space.
- The probability of an event E is the sum of probabilities of each sample point it contains.

$$P(E) = \sum_{\{n \in E\}} P(X = n)$$

Events

	X					
	1	2	3	4	5	6
6	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$
5	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$
4	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$
3	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$
2	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$
1	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$

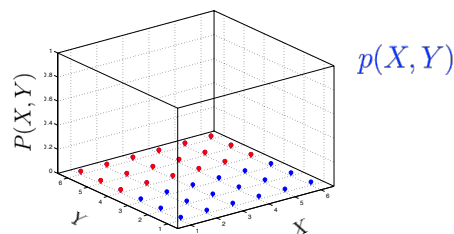
- Double dice roll
- What's $P(X = 5)$?



Events

	X					
	1	2	3	4	5	6
6	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$
5	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$
4	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$
3	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$
2	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$
1	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$	$\frac{1}{36}$

- Double dice roll
- What's $P(Y \geq 4)$?



A simple medical example

Consider a medical scenario, with 3 Boolean variables

- cavity* (does the patient have a cavity or not?)
- toothache* (does the patient have a toothache or not?)
- catch* (does the dentist's probe catch on the patient's tooth?)

Here's an example probability model: the joint probability distribution $p(\text{Toothache}, \text{Cavity}, \text{Catch})$

	<i>toothache</i>		\neg <i>toothache</i>	
	<i>catch</i>	\neg <i>catch</i>	<i>catch</i>	\neg <i>catch</i>
<i>cavity</i>	.108	.012	.072	.008
\neg <i>cavity</i>	.016	.064	.144	.576

Inference from a joint distribution

Given a full joint distribution, we can compute the probability of any event simply by summing the probabilities of the relevant sample points.

E.g. how to calculate $P(\text{toothache})$?

$$P(\text{toothache}) = 0.108 + 0.012 + 0.016 + 0.06 = 0.2$$

	toothache		$\neg \text{toothache}$	
	catch	$\neg \text{catch}$	catch	$\neg \text{catch}$
cavity	.108	.012	.072	.008
$\neg \text{cavity}$.016	.064	.144	.576

Inference from a joint distribution

Given a full joint distribution, we can compute the probability of any event simply by summing the probabilities of the relevant sample points.

E.g. how to calculate $P(\text{cavity} \vee \text{toothache})$?

	toothache		$\neg \text{toothache}$	
	catch	$\neg \text{catch}$	catch	$\neg \text{catch}$
cavity	.108	.012	.072	.008
$\neg \text{cavity}$.016	.064	.144	.576

Inference from a joint distribution

Given a full joint distribution, we can compute the probability of any event simply by summing the probabilities of the relevant sample points.

E.g. how to calculate $P(\text{cavity} \vee \text{toothache})$?

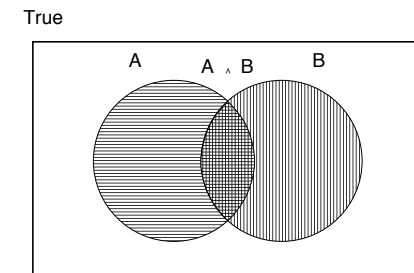
$$P(\text{cavity} \vee \text{toothache}) = 0.108 + 0.012 + 0.072 + 0.008 + 0.016 + 0.064 = 0.28$$

	toothache		$\neg \text{toothache}$	
	catch	$\neg \text{catch}$	catch	$\neg \text{catch}$
cavity	.108	.012	.072	.008
$\neg \text{cavity}$.016	.064	.144	.576

Set-theoretic relationships in probability

Given a full joint distribution, we can compute the probability of any event simply by summing the probabilities of the relevant sample points.

For instance: $P(a \vee b) = P(a) + P(b) - P(a \wedge b)$

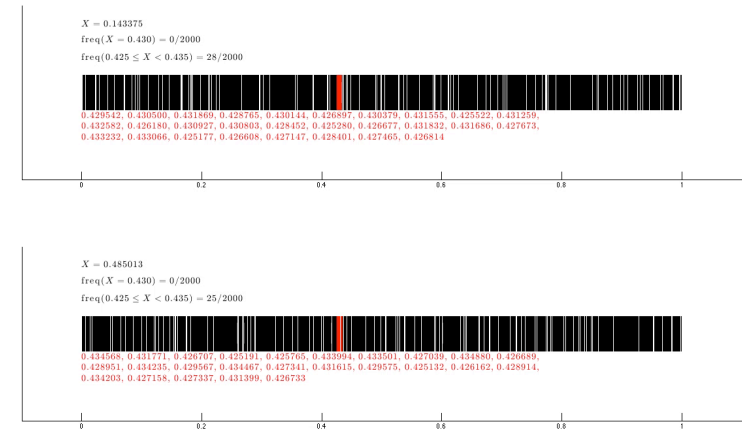


Continuous random variables

The sample spaces we've seen so far have been built from discrete random variables. But you can build probability models using **continuous random variables** too.

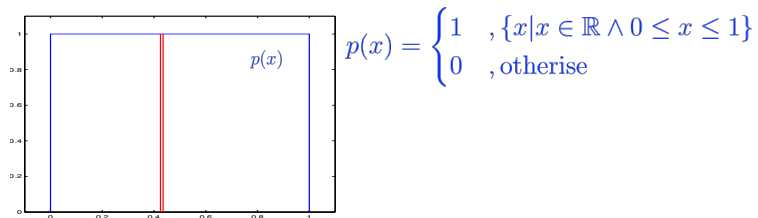
- E.g. we can define a random variable *Temperature*, whose domain is the real numbers.
- In the real domain (even if it's bounded) domain there is an infinite number of samples. Probability of continuous random variable hitting a specific value is 0.
- However, we can talk about probability of value being in certain range.

Continuous random variables



Probability density function

- For continuous variables, probability distributions are continuous, and are referred to as **probability density functions**
- E.g. here's a function which gives uniform probability for values between 0 and 1

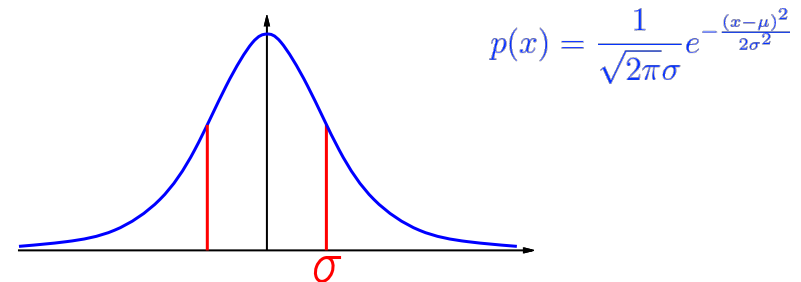


- This function is a density; it integrates to 1. So:

$$P(0.425 \leq x < 0.435) = \int_{0.425}^{0.435} p(x) dx = 0.01$$

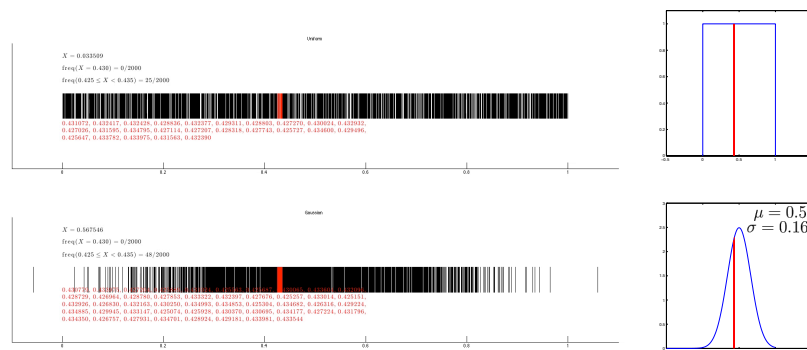
Gaussian distribution

- A particularly useful probability function for continuous variables is the **Gaussian** function (often referred to as the **normal** distribution)

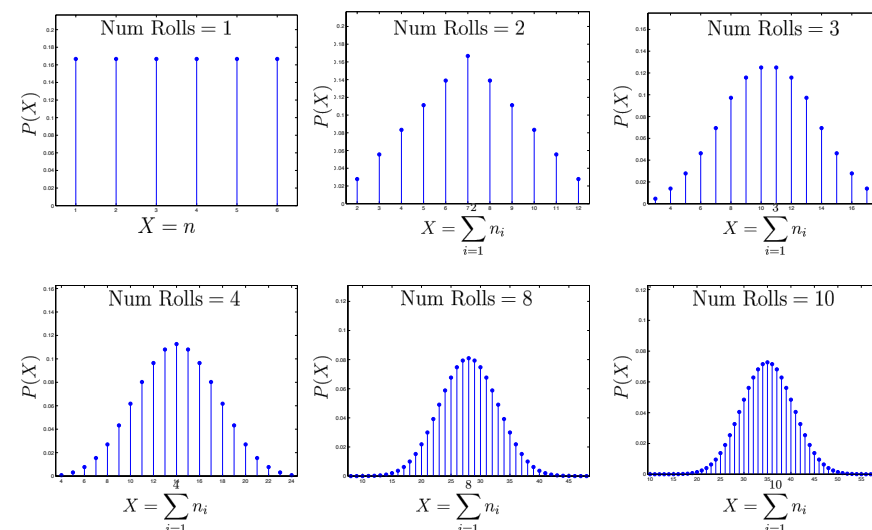


Lots of real-world variables have this distribution

Gaussian distribution



Central Limit Theorem



Expectation

- Probability weighted value of all possible values of a function dependent on a random variable
- “Average” result expected

Discrete distribution

$$E[g(x)] = \sum_i p(x_i)g(x_i)$$

Continuous distribution

$$E[g(x)] = \int p(x)g(x)dx$$

Mean and variance

- The expected value of the random variable itself

$$\mu = E[x]$$

- The expected value of the squared deviation of random variable from its mean (measures the spread of a probability distribution).

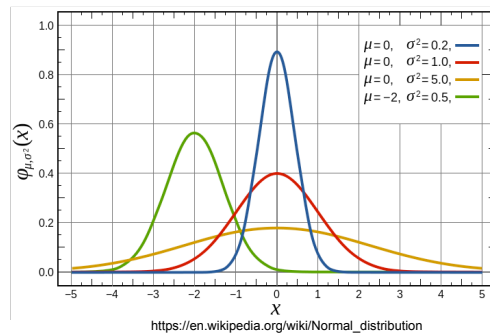
$$\sigma^2 = E[(x - \mu)^2]$$

An example: mean and variance of the normal distribution

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

$$E[x] = \mu$$

$$E[(x - \mu)^2] = \sigma^2$$



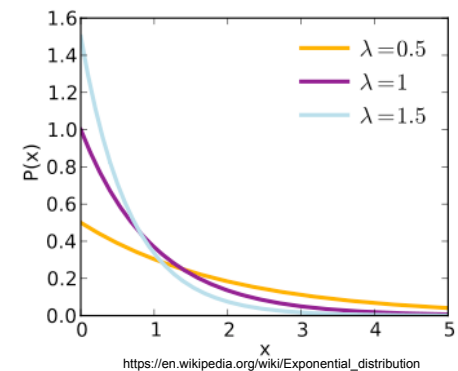
- Gaussian distribution is completely parametrised by its mean and variance
- σ - standard deviation

An example: mean and variance of the exponential distribution

$$p(x) = \begin{cases} \lambda e^{-\lambda x} & x \geq 0, \\ 0 & x < 0. \end{cases}$$

$$\mu = E[x] = \lambda^{-1}$$

$$E[(x - \mu)^2] = \lambda^{-2}$$



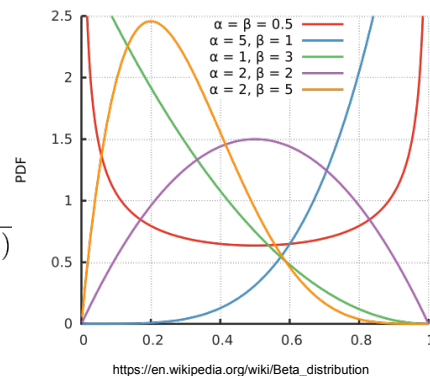
- Mean is the point splitting the probability density, such that area under curve is exactly 0.5 on either side

An example: mean and variance of the Beta distribution

$$p(x) = \frac{x^{\alpha-1}(1-x)^{\beta-1}}{\int_0^1 u^{\alpha-1}(1-u)^{\beta-1} du}$$

$$\mu = E[x] = \frac{\alpha}{\alpha + \beta}$$

$$E[(x - \mu)^2] = \frac{\alpha\beta}{(\alpha + \beta)^2(\alpha + \beta + 1)}$$



Summary and reading

Probability theory is the foundation for many learning algorithms.

- Key concepts: sample space, random variable, probability distribution, probability density, expectation

Reading for the lecture: AIMA Chapter 13 Sections 1-2

Reading for next lecture: AIMA Chapter 13 Section 3-6