# 1 Week 8: Probability: Part I

## 1.1 Data Fundamentals (H)

## 1.2 Probability theory, Bayes' Rule, discrete and continuous random variables

DF(H) - University of Glasgow - John H. Williamson - 2018

# 2 Summary

By the end of this unit you should know:

- what probability is, and different philosophical interpretations of it
- what inverse and forward probability are
- what a random variable, distribution, probability mass/density function are
- what the empirical distribution is and how it is computed from data
- what expectation/expected value is
- the axioms of probability theory
- conditional, marginal and joint distribution
- Bayes' rule and how to apply it problems
- what bigrams are and how they can be used
- representation of probability as odds, log-odds and log-probabilities and the numerical effects of these representations
- the specific problems of continuous random variables as compared to discrete random variables
- the normal distribution and the central limit theorem
- what multivariate distributions are

```
In [ ]:  import IPython.display
         IPython.display.HTML("""
         <script>
           function code_toggle() {
             if (code_shown){
               $('div.input').hide('500');
               $('#toggleButton').val('Show Code')
             } else {
               $('div.input').show('500');
               $('#toggleButton').val('Hide Code')
             }
             code_shown = !code_shown
           }

           $( document ).ready(function(){
             code_shown=false;
             $('div.input').hide()
           });
         </script>
         <form action="javascript:code_toggle()"><input type="submit" id="toggleE
```

```
In [4]:  import numpy as np
         import matplotlib as mpl
         import matplotlib.pyplot as plt
         from jhwutils.float_inspector import print_shape_html, print_float, prir
         from jhwutils.matrices import show_boxed_tensor_latex, print_matrix
         import jhwutils.image_audio as ia
         import numpy as np
         %matplotlib inline
         plt.rc('figure', figsize=(8.0, 4.0), dpi=140)
```

$x\mathbb{R}$  [♠]

# 3 Example: the lost submarine

The *USS Scorpion* was a nuclear armed submarine that disappeared on 30th June 1968 somewhere in the Atlantic. [This was a year in which *four* submarines were inexplicably lost at sea -- **the 1968 submarine mystery.**]
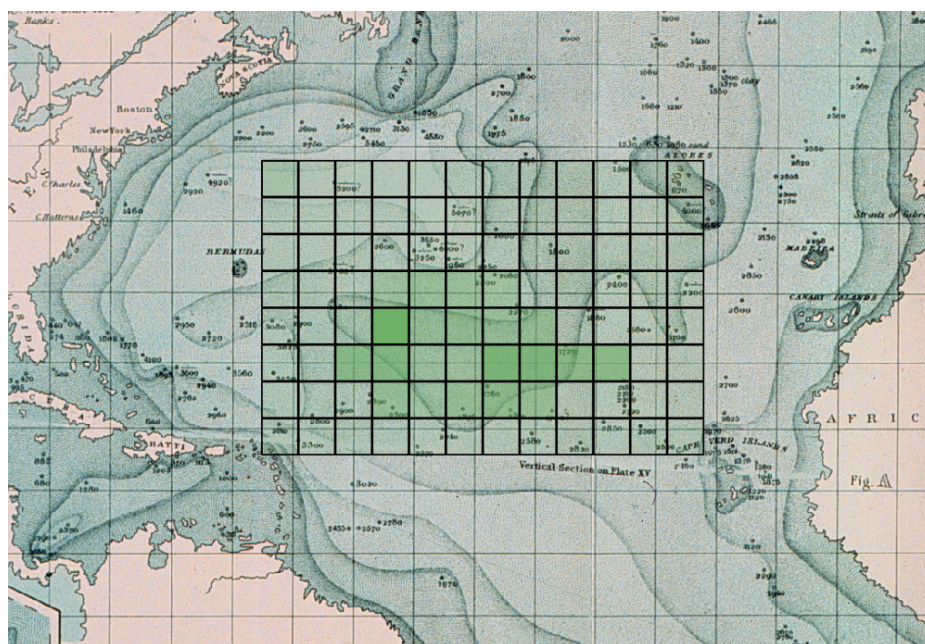
[Image: USS Scorpion at sea. Credit: US Navy, public domain]

The search for the submarine was one of the first times **probabilistic methods** were used for searching. The US Navy needed to find the submarine and recover it as fast as possible. Probabilistic methods allow us to reason *precisely* about things we are uncertain about -- like where the submarine is. Probability gives a concrete, computable representations for an uncertain world.

Imagine dividing the possible search areas into grid squares. Then, for example, it makes sense to ask *what is the probability that the submarine lies within a specific square?*. A map like this might be produced, where squares are colored according to their probability:



## 3.1 But what happened to the Scorpion?

# 4 Probability

This section of the course is concerned with **stochastic elements**; the role of uncertainty, randomness and statistics in computation. The fundamental mathematical principles are drawn from **probability theory**, which gives us simple and powerful ways of manipulating uncertain values, and lets us do useful operations like inferring the most likely hypotheses given some observations. Probability theory is a simple, consistent, and effective way to manipulate uncertainty.

## 4.1 What is probability?

[The following definitions are directly adapted from Peter Norvig's Probability notebook: http://nbviewer.jupyter.org/url/norvig.com/ipython/Probability.ipynb] (http://nbviewer.jupyter.org/url/norvig.com/ipython/Probability.ipynb%5D)

- **Experiment** (or **trial**) An occurrence with an uncertain outcome.

- For example, losing a submarine -- the location of the submarine is now unknown.
- **Outcome** The result of an experiment; one particular state of the world.
  - For example: the submarine is in ocean grid square [2,3].
- **Sample Space** The set of *all possible* outcomes for an experiment.
  - For example, ocean grid squares {[0,0], [0,1], [0,2], [0,3], ..., [8,7], [8,8], [8,9], [9,9]}.
- **Event** A *subset* of possible outcomes with some common property.
  - For example, the grid squares which are south of the Equator.
- **Probability** The probability of an event *with respect to a sample space* is the number of outcomes from the sample space that are in the event, divided by the total number of outcomes in the sample space. Since it is a ratio, probability will always be a real number between 0 (representing an impossible event) and 1 (representing a certain event).
  - For example, the probability of the submarine being below the equator, or the probability of the submarine being in grid square [0,0] (in this case the event is just a single outcome).
- **Frequency** a number describing how often an outcome occurs in a sequence of experiments.
  - which could either be a count of occurrences (the last 10 times, we found submarines below the equator after losing them) or a proportion of the total experiments (the 0.512 of the time, lost submarines are found below the equator).
  - For example, $X$ is a random variable representing the location of the submarine.
- **Distribution** A mapping from outcome to frequency for each outcome in a sample space.
  - For example, grid square [0,0] occurs 0.1 of the time; grid square[0,1] occurs 0.01 of the time, etc.
- **Probability distribution** A distribution that has been *normalized* so that the sum of the frequencies is 1. This is because an outcome must happen from a trial (with probability 1) so the sum of all possible outcomes together will be 1. A random variable has a probability distribution which maps each outcome to a probability.
  - For example $P(X = x)$, the probability that the submarine is in a specific grid square **x**.
- **Random variable** A variable representing an unknown value, whose probability distribution we *do* know. The variable is associated outcomes of a trial
- **Probability density/mass function** A function that *defines* a probability distribution by mapping each outcome to a probability $f_X(x), x \rightarrow \mathbb{R}$. This could be a continuous function over $x$ (density) or discrete function over $x$ (mass).
  - For example $f_X(x)$ would be a probability mass function for the submarine, which maps each grid square to real number representing its probability.
- **Observation** An outcome that we have directly observed; i.e. data.
  - For example, a submarine was found in grid square [0,5]
- **Sample** An outcome that we have simulated according a probability distribution. We say we have **drawn** a sample from a distribution.
  - For example, if we believe that the submarine was distributed according to some pattern, generate possible concrete grid positions that follow this pattern.
- **Expectation/expected value** The "average" value of a random variable.
  - The submarine was on average in grid square [3.46, 2.19]

# 5  Philosophy of probability

There are two schools of thought regarding probability and its uses. We will be (largely) following the **Bayesian interpretation**, but its worth understanding what that entails.

## 5.1  Bayesian/Laplacian view on probability

**Bayesians** treat probability as a **calculus of belief**; in this model of thought, probabilities are measures of *degrees of belief*. $P(A) = 0$ means a belief that event $A$ cannot be true and $P(A) = 1$ is a belief that event $A$ is absolutely certain. In the Bayesian perspective, it makes sense to say "the probability it is raining outside is 0.3" (the probability quantifies our belief about the weather given the information we have). **Note that is not a statement that we believe that the weather is 0.3 rainy (whatever that means)**

*[Image: A portrait of a man who isn't Thomas Bayes, who didn't propose Bayesian probability theory.]*

**Bayesians** allow for belief in states to be combined and manipulated via the rules of probability. The key process in Bayesian logic is *updating of beliefs*. Given some:

- **prior** belief (it's Glasgow, it's not likely to be sunny) and some
- new **evidence** (there seems to be a bright reflection inside) we can
- update our belief to calculate the **posterior** -- our new probability that it is sunny outside.

Bayesian inference requires that we accept priors over events, i.e. that we must explicitly quantify our assumptions with probability distributions. It is an extension of logic to uncertain information.

### 5.1.1 USS Scorpion

For example, in the submarine search, the prior might be that the submarine is probably in the south Atlantic (given the last radio broadcast received). Evidence might be the result of sonar surveys from search ships. After each survey, the posterior probability of the submarine being in the survey area could be updated. This represents our belief about where the vessel might be.

## 5.2 Frequentist view of probability

There is an alternative school of thought that considers probabilities to *only* be the long-term behaviour of repeated events (e.g. the probability of a coin coming up heads in 0.5 because over the long term this will be the average proportion of times this occurs).

A **frequentist** does not accept phrases like "what is the probability it is sunny just now?" as there is no long term behaviour involved (it is only "now" once). It does not make sense in this world view to talk about the probability of events that can only happen once. It *does* make sense in a frequentist view to ask things like "what is the probability it will be sunny on any given day?" since we can measure this event (sunny or not) for many different days. For example, frequentists would not assign a probability to the USS Scorpion being in a specific grid square; this is not an experiment that can be repeated.

### 5.2.1 Objectivity and subjectivity

Frequentist versus Bayesian debates quickly enter philosophical territory. The diversity of viewpoints and depth of arguments cannot be done justice here.

Very briefly, Bayesian probability theory is sometimes said to be **subjective** because it requires the specification of prior belief, whereas frequentist models of probability do not admit the concept of priors and thus is **objective**.

An alternative view is that the Bayesian model explicitly encodes uncertain knowledge and states universal formal rules for manipulating that knowledge, as formal logic does for definite knowledge. Frequentist methods are objective in the sense that they make statements about universal truths (e.g. asymptotic behaviour), but they do not form a calculus of belief, and thus can't answer many questions of importance directly.

#### 5.2.1.1 Bayesian

- Includes **priors**
- Probability is a **degree of belief**
- (Parameters of population considered to be random variables, data to be known)

#### 5.2.1.2 Frequentist

- No **priors**
- Probability is the *long-term frequency of events*
- (Parameters of population assumed to be fixed, data to be random)
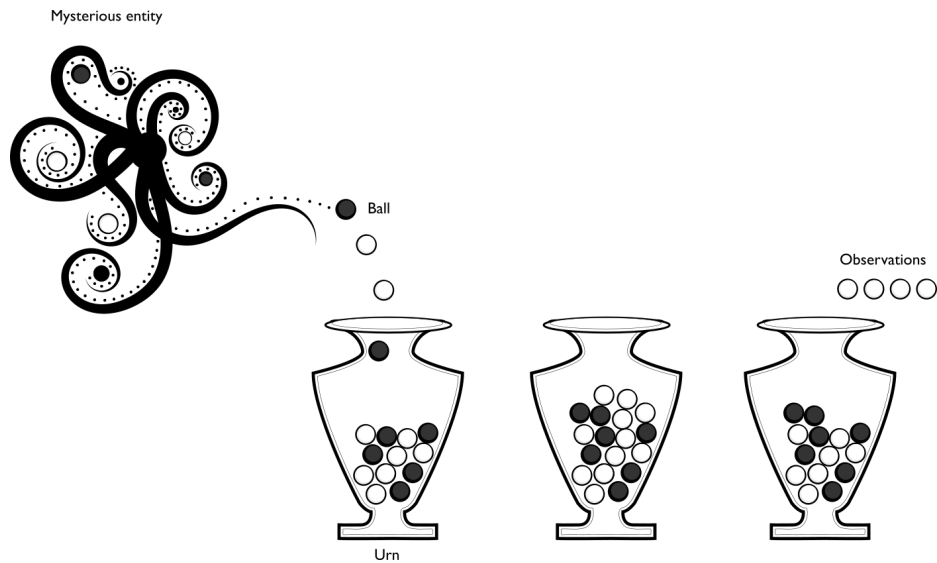
## 5.3 Superiority of probabilistic models

Regardless of the philosophical model you subscribe to, there is one thing you can be sure of: *probability is the best*.

There are other models of uncertainty than probability theory that are sometimes used. However, all other representations of uncertainty are *strictly inferior* to probabilistic methods *in the sense that* a person, agent, computer placing "bets" on future events using probabilistic models has the best possible return out of all decision systems when there is uncertainty.

> *Any theory with as good a gambling outcome as would be achieved using probability theory is equivalent to probability theory.*

## 5.4 Generative models: forward and inverse probability

A key idea in probabilistic models is that of a **generative process**; the idea that there is some unknown process going on, the results of which can be observed. The process itself is governed by unobserved variables that we do not know but which we can **infer**.

The classic example is an **urn problem**. Consider an urn, into which a number of balls have been poured (by some mysterious entity, say). Each ball can be either black or white.

You pull out four random balls from the urn and observe their colour. You get four white balls.

There are lots of questions you can ask now:

- What is the probability that the next ball that is drawn will be white?
  - This is a **forward probability** question. It asks questions related to the distribution of the observations.
- What is the distribution of white and black balls in the urn?
  - This is an **inverse probability** question. It asks questions related to unobserved variables that govern the process that generated the observations.
- Who is the mysterious entity?
  - This is an unknowable question. The observations we make cannot resolve this question.

There are a huge number of processes that can be framed as urn problems (urns where balls are replaced, problems where there are multiple urns and you don't know which urn the balls came from, problems where balls can move between urns, and so on).

---

# 6 A formal basis for probability theory
## 6.1 Axioms of probability

There are only a few basic axioms of probability, from which everything else can be derived. Writing $P(A)$ to mean the probability of event $A$:

- **Boundedness**
$$0 \leq P(A) < 1 \ [\spadesuit]$$
  all possible outcomes $A$ -- probabilities are 0, or positive and less than 1.
- **Unitarity**
$$\sum_A P(A) = 1 \ [\spadesuit]$$
  for the complete set of possible **outcomes** (not events!) $A \in \sigma$ in a sample space $\sigma$ -- something always happens.
- **Sum rule**
$$P(A \vee B) = P(A) + P(B) - P(A \wedge B), \ \ [\spadesuit]$$
  i.e. the probability of either outcome $A$ or $B$ happening is the sum of the independent probabilities minus the probability of both happening. (notation note: $\vee$ means "or" and $\wedge$ means "and")
- **Conditional probability** The conditional probability $P(A|B)$ is defined to be the probability that outcome $A$ will happen *given that we already know $B$ to have happened.*
$$P(A|B) = \frac{P(A \wedge B)}{P(B)} \ [\spadesuit]$$

# 7 Random variables and distributions

A **random variable** is a variable that can take on different values, but we do not know what value it has; i.e. one that is "unassigned". However, we have some knowledge which captures the possible states the variable could take on, and their corresponding probabilities. Probability theory allows us to manipulate random variables without having to assign them a specific value.

A random variable is written with a capital letter, like $X$.

A random variable might represent:

- the outcome of dice throw (discrete);
- whether or not it is raining outside (discrete: binary);
- the latitude of the USS Scorpion (continuous);
- the height of person we haven't met yet (continuous).

In the next unit, we will see how random variables can be implemented as first-class values in programming languages, just like floats or integers.

# 7.1 Distributions

A **probability distribution** defines how likely different states of a random variable are.

We can see $X$ as the the the *experiment* and $x$ as the *outcome*, with a function mapping every possible outcome to a probability. We write $P(X = x)$(note the case!), and use the shorthand notations:

$$P(X = x), \text{ the probability of random variable X taking on value x}$$
$$P(X), \text{ shorthand for probability of X=x}$$
$$P(x), \text{ shorthand for probability of specific value X=x}$$

We can see an outcome as a random variable taking on a specific value i.e. $P(X = x)$. Note that by convention we use $P(A)$ to mean the probability of **event** $A$, not a random variable $A$ (an **event** is a *set* of **outcomes**; **random variables** only assign probabilities to **outcomes**).

## 7.1.1 Discrete and continuous

Random variables can be continuous (e.g. the height of a person) or discrete (the value showing on the face of a dice).

- **Discrete variables** The distribution of a discrete random variable is described with a **probability mass function** (PMF) which gives each outcome a specific value; imagine a Python dictionary mapping outcomes to probabilities. The PMF is usually written $f_X(x)$, where $P(X = x) = f_X(x)$.
- **Continuous variables** A continuous variable has a **probability density function** (PDF) which specifies the spread of the probability over outcomes as a *continuous function* $f_X(x)$. It is **not** the case that $P(X = x) = f_X(x)$ for PDFs.

### 7.1.1.1 Integration to unity

A probability mass function or probability density function *must* sum/integrate to exactly 1, as the random variable under consideration must take on *some* value; this is a consequence of unitarity. Every repetition of an experiment has exactly one outcome.

$$\sum_i f_X(x_i) = 1 \quad \text{for PMFs of discrete RVs}$$

$$\int_x f_X(x)\,dx = 1 \quad \text{for PDFs of continuous RVs}$$

# 7.2 PMF example: sum of dice rolls

A very simple discrete PMF is the expected value of the sum of two six-faced dice. $P(X = x) = f_X(x)$ takes on values for each possible outcome $x \in \{2, 3, 4, 5, 6, 7, 8, 9, 10, 11, 12\}$

```python
# the PMF of the sum of two dice rolls
from IPython.display import display, Latex
def two_dice():
    # form the sum of the cross product of these possibilities
    roll_two = [i+j for i in range(1,7) for j in range(1,7)]


    # now plot the histogram
    fig = plt.figure()
    ax = fig.add_subplot(1,1,1)

    pmf, edges, patches = ax.hist(roll_two, density=True, bins=range(1,1

    for outcome, p in zip(roll_two, pmf[1:]):
        display(Latex("$P(X={d}) = {p:.2f}$".format(d=outcome, p=p)))

    print("Sum of PMF = %.2f" % np.sum(pmf)) # sum of probability shoulc
    ax.set_title("PMF of sum of 2d6 dice")
    ax.set_xlabel("Sum of rolls x")
    ax.set_frame_on(False)
    ax.set_ylabel("P(X=x)")
```

`two_dice()`

$P(X = 2) = 0.03$

$P(X = 3) = 0.06$

$P(X = 4) = 0.08$

$P(X = 5) = 0.11$
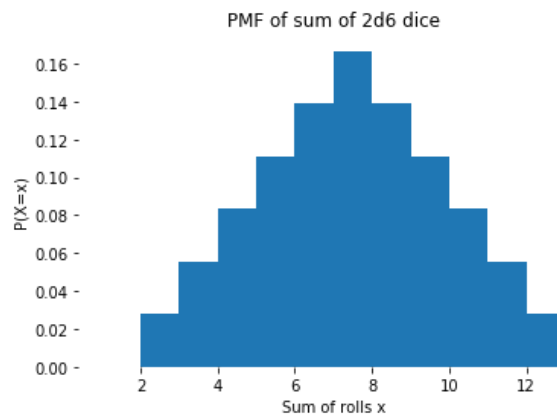
$P(X = 6) = 0.14$

$P(X = 7) = 0.17$

$P(X = 3) = 0.14$

$P(X = 4) = 0.11$

$P(X = 5) = 0.08$

$P(X = 6) = 0.06$

$P(X = 7) = 0.03$

```
Sum of PMF = 1.00
```



PMF of sum of 2d6 dice

## 8 Expectation

If a random variable takes on numerical values, then we can defined the **expectation** or **expected value** of a random variable $\mathbb{E}[X]$ as:

$$\mathbb{E}[X] = \int_x x f_X(x) dx \quad [\spadesuit]$$

For a discrete random variable with probability mass function $P(X = x) = f_X(x)$, we would write this as a summation:

$$\mathbb{E}[X] = \sum_x f_X(x) x \quad [\spadesuit]$$

If there are only a finite number of possibilities, then this is:

$$\mathbb{E}[X] = P(X = x_1)x_1 + P(X = x_2)x_2 + \cdots + P(X = x_n)x_n$$

The expectation is the "average" of a random variable. Informally, it represents what we'd "expect to happen"; the most likely overall "score". It can be thought of as a *weighted sum* of all the possible outcomes of an experiment, where each outcome is weighted by the probability of that outcome occurring.

For example, in the "pair of dice" scenario, we can compute the expected value of the number of "dots" showing in total after a roll. We compute the probability of each number of dots showing and multiply by that number of dots, and summing the result. This is the expected number of dots showing, on average, or the expectation.

```python
from collections import Counter
dice_values = [i+j for i in range(1,7) for j in range(1,7)]

# count the occurence of each possible value
counts = Counter(dice_values)
total_possibilities = len(dice_values)

pmf = np.array([counts[i]/total_possibilities for i in counts])
print("PMF of two dice:", pmf)
values = np.array([i for i in counts])
print("Face value of two dice:", values)

# expectation is each value weighted by probability.
expected_value = np.sum(pmf * values)
print("Expected value of sum of two dice: %.2f" % expected_value)
```

This is an intuitive property.

> Imagine you meet a street hustler, who asks you to play the two dice game. He offers you a chance to buy in for £8; you win as many pounds as the show on the top faces of the dice after throwing them. Is this a fair game?

No. The expected return is only £7, and you have to put in £8 to play, so the expected result is a £1 loss (sometimes stated as "negative expected value" or -ve EV). If it was £7 to buy in, the game would be fair, in the sense that you and the hustler would not transfer money on average.

## 8.1 Expectation and means

Expectation corresponds to the idea of a **mean** or **average** result. The expected value of a random variable is the **true average** of the value of all outcomes that would be observed if we ran the experiment an infinite number of times. This is the **population mean** -- the mean of the whole, possibly infinite, population of a random variable.

Many important properties of random variables can be defined in terms of expectation.

- The mean of a random variable $X$ is just $\mathbb{E}[X]$. It is a measure of **central tendency**.
- The variance of a random variable $X$ is $\text{var}(X) = \mathbb{E}[(X - \mathbb{E}[X])^2]$. It is a measure of **spread**.

### 8.1.1 Expectations of functions of X

We can apply functions to random variables, for example, the square of a random variable. The expectation of any function $g(X)$ of a random variable $X$ is defined as:

$$\mathbb{E}[g(X)] = \int_x f_X(x)g(x)dx \ [\spadesuit]$$

That is, we just take the sum/integral of each outcome, passed through the function $g(x)$, weighted by the probability of the outcome $x$. $g(x)$ could be thought of a "scoring" function, which assigns a real number to every outcome of $X$. Note that $g(x)$ has no effect on the probability density/mass function $f_X(x)$. It doesn't affect the probability of the outcomes, just the *value assigned to those outcomes*. For example, if we played the a dice game where the score of a throw was the *square* of the number of dots showing, (so 2 showing would be worth 4 points, 8 showing would be worth 64 points, etc.) we would compute $\mathbb{E}[X^2]$ as:

```
sqr_expected_value = np.sum(pmf * (values ** 2))
print("Expected value of square of sum of two dice: %.2f" % sqr_expected
```

```
--------------------------------------------------------------------------
--
NameError                                 Traceback (most recent call las
t)
<ipython-input-13-e69a3537d85c> in <module>()
----> 1 sqr_expected_value = np.sum(pmf * (values ** 2))
      2 print("Expected value of square of sum of two dice: %.2f" % sqr_e
xpected_value)

NameError: name 'pmf' is not defined
```

Or we could compute the expectation of a two-dice game where you get twice the face value for any pair of identical rolls, and zero otherwise ("pairs or nothing"):

```python
# 6x6 matrix of outcomes
from jhwutils.matrices import print_matrix
p_faces = np.ones((6,6))/36.0
values = np.diag(np.arange(1,7)*4)
print_matrix("f_X(d_1, d_2)", p_faces)
print_matrix("g(d_1, d_2)", values)
# expected value of playing this game
print("Expected value of 'pairs or nothing' dice game", np.sum(p_faces*v
```

$$f_X(d_1, d_2) = \begin{bmatrix} 0.03 & 0.03 & 0.03 & 0.03 & 0.03 & 0.03 \\ 0.03 & 0.03 & 0.03 & 0.03 & 0.03 & 0.03 \\ 0.03 & 0.03 & 0.03 & 0.03 & 0.03 & 0.03 \\ 0.03 & 0.03 & 0.03 & 0.03 & 0.03 & 0.03 \\ 0.03 & 0.03 & 0.03 & 0.03 & 0.03 & 0.03 \\ 0.03 & 0.03 & 0.03 & 0.03 & 0.03 & 0.03 \end{bmatrix}$$

$$g(d_1, d_2) = \begin{bmatrix} 4 & 0 & 0 & 0 & 0 & 0 \\ 0 & 8 & 0 & 0 & 0 & 0 \\ 0 & 0 & 12 & 0 & 0 & 0 \\ 0 & 0 & 0 & 16 & 0 & 0 \\ 0 & 0 & 0 & 0 & 20 & 0 \\ 0 & 0 & 0 & 0 & 0 & 24 \end{bmatrix}$$

```
Expected value of 'pairs or nothing' dice game 2.3333333333333335
```

Expected values are essential in making **rational decisions**, the central problem of **decision theory**. They combine scores (or **utility**) with uncertainty (**probability**).

The expected value gives us a way of deciding, for example, how much it would be worth paying to play a dice game. If the units were pounds, we'd break even if we paid £2.33 to play the game, and make a profit (on average) if we paid £2.00 to play the game. The expected average profit is just the stake we pay to play each game minus the expected value of one round of the game: $\mathbb{E}[\text{game}] = \text{stake} - \mathbb{E}[X]$

## 8.2 Samples and sampling

**Samples** are observed outcomes of an experiment; we will use the term **observations** synonymously, though samples usually refer to simulations and observations to concrete real data.

We can **sample** from a distribution; this means simulating outcomes according to the probability distribution of those variables. We can also **observe** data which comes from an external source, that might believe is generated by some probability distribution. For example, we can sample from the sum of dice PMF by rolling two dice and summing the result. This is a *sample* or a *draw* from this distribution. For discrete random variables, this is easy: we simply produce samples by drawing each outcome according to its probability. We will discuss sampling strategies in the next unit.

## 8.3 The empirical distribution

For discrete data, we can estimate the probability mass function that might be generating **observations** by counting each outcome seen divided by the total number of trials. This is called the **empirical distribution.**

This can be thought of as the **normalized histogram** of counts of occurrences of outcomes.

```python
## PMF for dice rolls here
## simulate dice rolls
def simulate_two_dice(n):
    d1 = np.random.randint(1,7,n)
    d2 = np.random.randint(1,7,n)
    return d1+d2


### Run some random sampling trials, plotting the empirical PMF
### and the true PMF (that generated those samples)
n_samples = 200 # number of samples in one repetition
trials = 1  # number of times to repeat


dice_values =  [i+j for i in range(1,7) for j in range(1,7)]


# compute the PMF (again)
counts = Counter(dice_values)
total_possibilities = len(dice_values)
pmf = np.array([0]+[counts[i]/total_possibilities for i in counts]+[0])


### Plot the figure
fig = plt.figure()
ax = fig.add_subplot(1,1,1)
ax.step(np.arange(2,15), pmf, label="True PMF")


for j in range(trials):
    # simulate a bunch of pairs of dice throws
    two_dice_samples = simulate_two_dice(n_samples)
    label = "Empirical PDF" if j==0 else None
    # compute the normalized histogram (empirical PMF)
    ax.hist(two_dice_samples, bins = np.arange(0,15), density=True, labe

ax.legend()
ax.set_frame_on(False)
ax.set_title("Two dice empirical PMF, $N={n}$".format(n=n_samples))
ax.set_xlabel("Dice value")
ax.set_ylabel("P(X)")
```

```
----------------------------------------------------------------------
--
NameError                                 Traceback (most recent call las
t)
<ipython-input-16-c0119ab2a58f> in <module>()
     14
     15 # compute the PMF (again)
---> 16 counts = Counter(dice_values)
     17 total_possibilities = len(dice_values)
     18 pmf = np.array([0]+[counts[i]/total_possibilities for i in counts
]+[0])

NameError: name 'Counter' is not defined
```

## 8.3.1 Computing the empirical distribution

For discrete random variables, we can always compute the empirical distribution from a series of observations; for example from the counts of a specific word in a *corpus* of text (e.g. in every newspaper article printed in 1994). We just count the number of times each word is seen and divide by the total

number of words.

$$P(X = x) = \frac{n_x}{N}, \quad [\spadesuit]$$

where $n_x$ is the number of time outcome $x$ was observed, and $N$ is the total number of trials.

Note that the empirical distribution is a distribution which *approximates* an unknown true distribution. For very large samples of discrete variables, the empirical distribution will increasingly closely approximate the **true PMF**, assuming the samples we see are drawn in an unbiased way. However, this approach does not work usefully for continuous random variables, since we will only ever see each observed value once (think about why!).

## 8.4 A more interesting PMF

The following code loads a text file (in this case, *Romeo and Juliet*), and converts into a vector of ASCII codes. It keeps only letters and space, and converts everything to lowercase.

> Rom. Give me a torch. I am not for this ambling.
> Being but heavy, I will bear the light.
>
> Mer. Nay, gentle Romeo, we must have you dance.
>
> Rom. Not I, believe me. You have dancing shoes
> With nimble soles; I have a soul of lead
> So stakes me to the ground I cannot move.

In [46]:
```python
alphabet = 'abcdefghijklmnopqrstuvwxyz'
#alphabet = 'eaoiutnshrdlcmfwypvbgkjqxz'


def numerify(fname):
    with open(fname) as f:
        return np.array([alphabet.index(c) for c in f.read().lower() if
```

We can compute the empirical distribution of letters in Romeo and Juliet. Its probability mass function is just given by the normalised count of each character.

In [39]:
```python
def empirical_pmf_plot(vector, alphabet, title):
    fig = plt.figure()
    ax = fig.add_subplot(1,1,1)
    probs,_,_ = ax.hist(vector, bins=np.arange(len(alphabet)+1), density

    ax.set_xticks(np.arange(len(alphabet)+1)+0.5)
    ax.set_xticklabels(alphabet)
    ax.set_xlabel("Character $x$")
    ax.set_ylabel("Empirical probability $P(X=x)$")
    ax.set_title("Empirical PMF of {0}".format(title))
    return probs
```

In [40]:
```python
rj_chars = numerify("data/romeo_juliet.txt")
rj_pmf = empirical_pmf_plot(rj_chars, alphabet, "Romeo and Juliet")
```



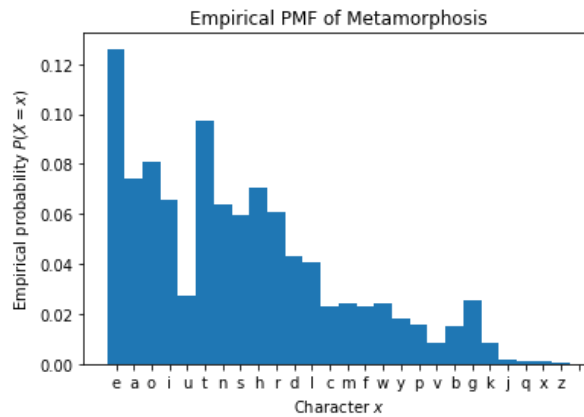## 8.4.1 Kafka

We can repeat the process with Kafka's *Metamorphosis*; an English text, but one with a different writing style.

In [41]:
```
metamorphosis = numerify("data/metamorphosis.txt")
metamorphosis_pmf = empirical_pmf_plot(metamorphosis, alphabet, "Metamor
```



The results are similar (both are English texts), but there are some subtle differences; there slightly more q characters in Kafka compared to Shakespeare; the frequencies of h and i are reversed.

### 8.4.2 Joint, conditional, marginal

The **joint probability** of two random variables is written

$$P(X, Y)$$

and gives the probability that $X$ and $Y$ take the specific values *simultaneously* (i.e. $P(X = x) \wedge P(Y = y)$).

The **marginal probability** is the derivation of $P(X)$ from $P(X, Y)$ by integrating (summing) over all the possible outcomes of $Y$:

$$P(X) = \int_y P(X, Y) dy \text{ for a PDF.} \quad [\spadesuit]$$

$$P(X) = \sum_y P(X, Y) \text{ for a PMF.}$$

This allows us to compute a distribution over one random variable from a joint distribution by summing over all the possible outcomes of the other variable involved.

**Marginalisation** just means integration over one or more variables from a joint distribution: it *removes* those variables from the distribution.

Two random variables are **independent** if the they do not have any dependence on each other. If this is the case then the joint distribution is just the product of the individual distributions: $P(X, Y) = P(X)P(Y)$. This is **not true in the general case where the variables have dependence**.

The **conditional probability** of a random variable $X$ *given* a random variable $Y$ is written as

$$P(X|Y)$$

and can be computed as

$$P(X|Y) = \frac{P(X, Y)}{P(Y)}. \quad [\spadesuit]$$

This tells us how likely the outcomes of $X$ are *if we already know* (or fix) the outcomes of $Y$. Read as "probability of X taking on the value x given that Y has taken on the value y". The conditional probability is $P(X|Y) = P(X)$ and $P(Y|X) = P(Y)$ if $X$ and $Y$ are independent.

## 8.5 Bigrams

We can look at these in the case of the character model of text we saw earlier. From the vector of character codes that make up "Metamorphosis" we can take every *pair of characters*, in order that they appear. That is, we consider two characters $c_{i-1}$ and $c_i$ at some index $i$. This is called a **bigram** model, and

there are unigram, trigram, n-gram generalisations of the idea. A "gram" refers to a unit like a character or word, and we are discussing a character **bigram model**.
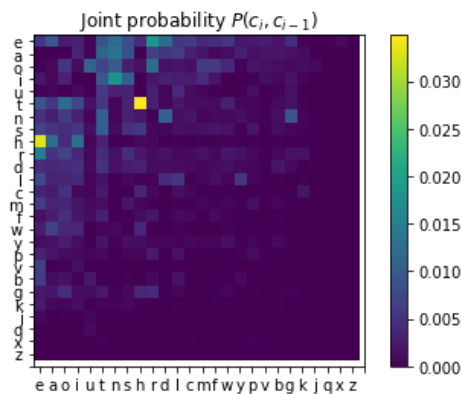
- The joint distribution of bigrams $P(C_i = c_i, C_{i-1} = c_{i-1})$ is given by the normalized count of each character pair.
- The marginal distribution $P(C_i = c_i)$ can be computed from $P(C_i = c_i, C_{i-1} = c_{i-1})$ by summing over every possible character $c_{i-1}$, and likewise to marginalise to find $P(C_{i-1} = c_{i-1})$
- The conditional distribution $P(C_i = c_i | C_{i-1} = c_{i-1})$ is given by the joint distribution, divided by the counts of $P(C_{i-1} = c_{i-1})$. It tells us how likely we are to observe a specific character $c_i$ *given* that we have observed a character $c_{i-1}$ just beforehand.

In [42]:
```python
meta_joint = np.zeros([len(alphabet), len(alphabet)])
for c_1, c_2 in zip(metamorphosis[1:], metamorphosis[:-1]):
    meta_joint[c_2, c_1] += 1

# normalise by total number of pairs
meta_joint = meta_joint / np.sum(meta_joint)

def show_2d_histogram(hist):
    fig = plt.figure()
    ax = fig.add_subplot(1,1,1)
    img = ax.imshow(hist)
    fig.colorbar(img)
    ax.set_xticks(np.arange(len(alphabet)+1))
    ax.set_xticklabels(alphabet)
    ax.set_yticks(np.arange(len(alphabet)+1))
    ax.set_yticklabels(alphabet)
    ax.set_title("Joint probability $P(c_{i},c_{i-1})$")

# rows = c_i-1, columns = c_i
show_2d_histogram(meta_joint)
```



This tells us how likely each possible **pair** of characters is. Some combinations are much more likely than others.

In [43]:
```python
# compute the marginial probabilities, by summing over the relevant axes
# it's this easy!
meta_c1 = np.sum(meta_joint, axis=0)
meta_c2 = np.sum(meta_joint, axis=1)
```
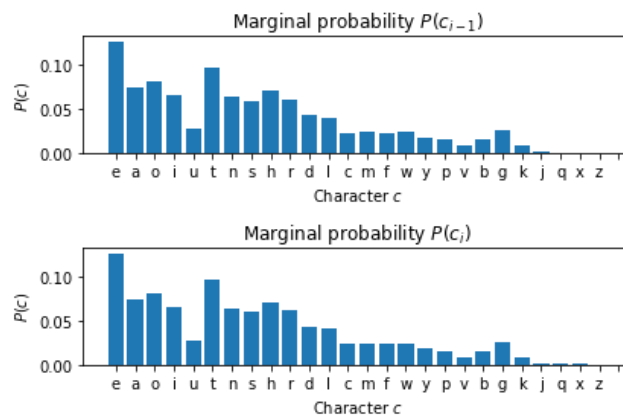
```python
# plot them
fig = plt.figure()

def show_1d_histogram(hist):
    ax.bar(np.arange(len(alphabet)), hist)
    ax.set_xticks(np.arange(len(alphabet)+1))
    ax.set_xticklabels(alphabet)
    ax.set_ylabel("$P(c)$")
    ax.set_xlabel("Character $c$")

ax = fig.add_subplot(2,1,1)
show_1d_histogram(meta_c1)
ax.set_title("Marginal probability $P(c_{i-1})$")

ax = fig.add_subplot(2,1,2)
show_1d_histogram(meta_c2)
ax.set_title("Marginal probability $P(c_{i})$")
plt.tight_layout()
```
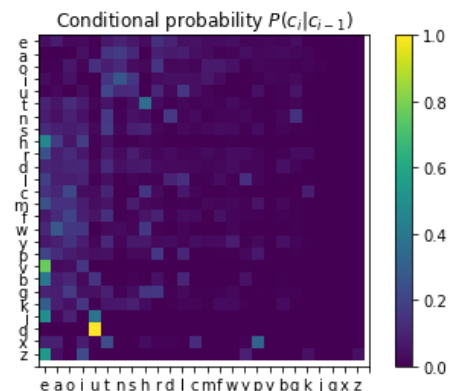


```python
# p(c2|c1) = p(c1,x2) / p(c1)
# have to transpose to normalise across columns, instead of rows
meta_conditional = (meta_joint.T / meta_c1).T

# plot it
show_2d_histogram(meta_conditional)
plt.gca().set_title("Conditional probability $P(c_i|c_{i-1})$")
```

```
Text(0.5,1,'Conditional probability $P(c_i|c_{i-1})$')
```



This tells us what character to expect, *given the character we know we saw before*. It lets us *predict* the next character. We can look along the rows to do this look up: each row sums to 1.0. For example, look at row "q": the probability is concentrated on the entry in column "u" -- "u" always follows "q" in English.

# 9 Quiz

# 10 Bayes' Rule

## 10.1 ☣ Plague ☣

- A new test is developed that can detect *Plague X* with 95% accuracy.
- We'll assume that 95% accuracy means:
    - a 5% **false positive rate**, i.e. 5% of the time people who don't have the disease test positive
    - a 5% **false negative rate**, i.e. 5% of the time people have the disease test negative
- One in a ten thousand people are known to have *Plague X*
    - we might know this because it's universally fatal within one year, and one in ten thousand people died of it last year, for example.
- You go for the test, and it comes back *positive* for *Plague X*.
- How likely are you to have *Plague X*?

### 10.1.1 Options

- A: about 95 in 100
- B: about 5 in 100
- C: about 1 in 100
- D: about 1 in 5000
- E: about 1 in 10000

## 10.2 Prior, likelihood, posterior

### 10.2.1 Inverting conditional distributions

We often want to know the probability of a some outcome $A$ given some other outcome $B$; that is $P(A|B)$.

But we are often in the situation that we can only compute $P(B|A)$.

In general $P(A|B) \neq P(B|A)$; and the two expressions can be completely different.

Typically, this type of problem occurs where we:

- want to know the probability of some event given some *evidence* (*how likely is it that I have a disease given that my blood test came back positive?*)
- but we only know the probability of observing evidence given the event (*if you have this disease, the blood test will come back positive 95% of the time*).

**Bayes' rule** gives the correct way to invert the probability distribution:

$$P(A|B) = \frac{P(B|A)P(A)}{P(B)} \quad [\spadesuit]$$

This follows directly from the axioms of probability. Bayes' Rule is a very important rule, and has some surprising consequences.

### 10.2.2 Nomenclature

- $P(A|B)$ is called the **posterior** -- what we want to know, or will know after the computation

- $P(B|A)$ is called the **likelihood** -- how likely the event $A$ is to produce the evidence we see

- $P(A)$ is the **prior** -- how likely the event $A$ is regardless of evidence

- $P(B)$ is the **evidence** -- how likely the evidence $B$ is regardless of the event.

Bayes' rule gives a consistent rule to take some prior belief and combine it with observed data to estimate a new distribution which combines them.

We often phrase this as some **hypothesis** $H$ we want to know, given some **data** $D$ we observe, and we write Bayes' Rule as:

$$P(H|D) = \frac{P(D|H)P(H)}{P(D)} \quad [\spadesuit]$$

$H$ and $D$ are random variables in this expression.

(the probability of the hypothesis given the data) is equal to (the probability of the data given the hypothesis) times (the probability of the hypothesis) divided by (the probability of the data). In other words, if we want to work out how likely a hypothesis is to be true given observations, but we only know how likely we are to have seen those observations if that hypothesis *was* true, we can use Bayes' rule to solve the problem.

## 10.3 Integration over the evidence

We can say that the posterior probability is *proportional* to the product of the prior and the likelihood.

But to evaluate its value, we need to compute $P(D)$, **the evidence**.

It is difficult to see what this represents at first. But one way to think of it is as the result of marginalising the $P(D)$ from the joint distribution $P(H, D)$; that is integrating $P(H, D)$ over every possible outcome of $H$, for each possible $D$.

Because probabilities must add up to 1, we can write $P(B)$ as:

$$P(D) = \sum_i P(D|H_i)P(H_i) \quad [\spadesuit]$$

for a set of discrete outcomes $A_i$ or

$$P(D) = \int_A P(D|H)P(H)dA \quad [\spadesuit]$$

for a continuous distribution of outcomes.

This trick is essential in understanding Bayes Rule!

In general this can be difficult to compute. For binary simple cases where there are only two possible outcomes ($H$ can only be 0 or 1), Bayes' rule can be written as:

$$P(H = 1|D) = \frac{P(D|H = 1)P(H = 1)}{P(D|H = 1)P(H = 1) + P(D|H = 0)D(H = 0)},$$

Bayes Rule for *Plague X*

```
P(Plague|Test) = [P(Test|Plague) P(Plague)] / P(Test)

# integrate over the two possible
#states of Plague
P(Test) =  [P(Test|not Plague)P(not Plague) +
            P(Test|Plague)P(Plague)]
```

In [ ]:

```
# one in a ten thousand have the plague
P_Plague = 1.0 / 100000
# if you have the plague, test is positive with 5% error
P_Test_Plague = 1-0.05
# If you don't have the plague, the test is positive 5% of the time
P_Test_not_Plague = 0.05
# chance you don't have the plague, before we see the test
P_not_Plague = 1-P_Plague

# integrate over evidence; only two possibilities here
P_Test = (P_Test_Plague * P_Plague +
            P_Test_not_Plague * P_not_Plague)

# what is the probability you have the plague,
# given this 95% accurate test?
# (hint: it is not 95%!)
P_Plague_Test = P_Test_Plague * P_Plague / P_Test

print("After testing positive, you have a 1:%.0f chance of having the pl
```
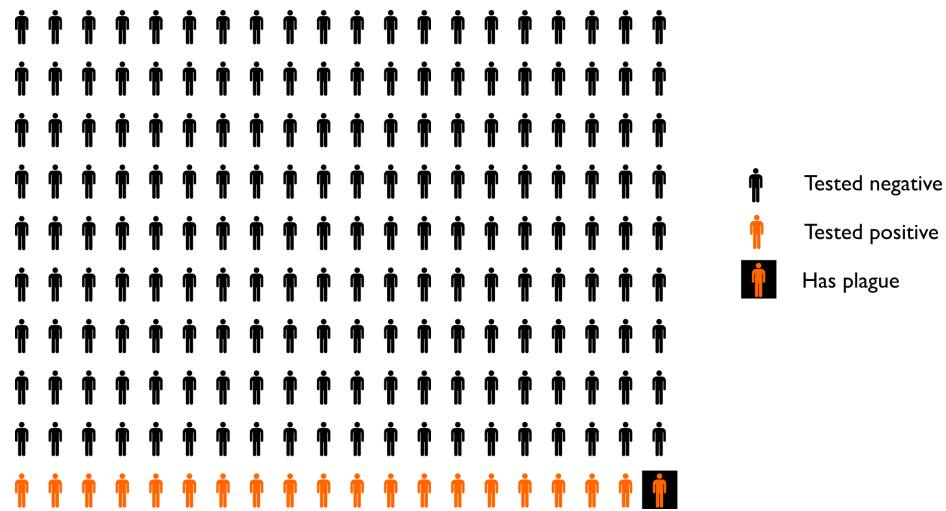
## 10.4 Natural frequency

There is an approach to explaining problems like this which makes it much less likely to make poor judgements. **Natural frequency** explanations involve imagining concrete populations of a fixed size (10000 people in a room, for example), and considering the proportions of the populations as *counts* (how many people in the room have the plague?).

We can use this to visualise the problem above and explain the apparent paradox. The graph below shows the case $P(\text{plague}) = 0.005$ (1 in 200), again with a 5% accurate test.
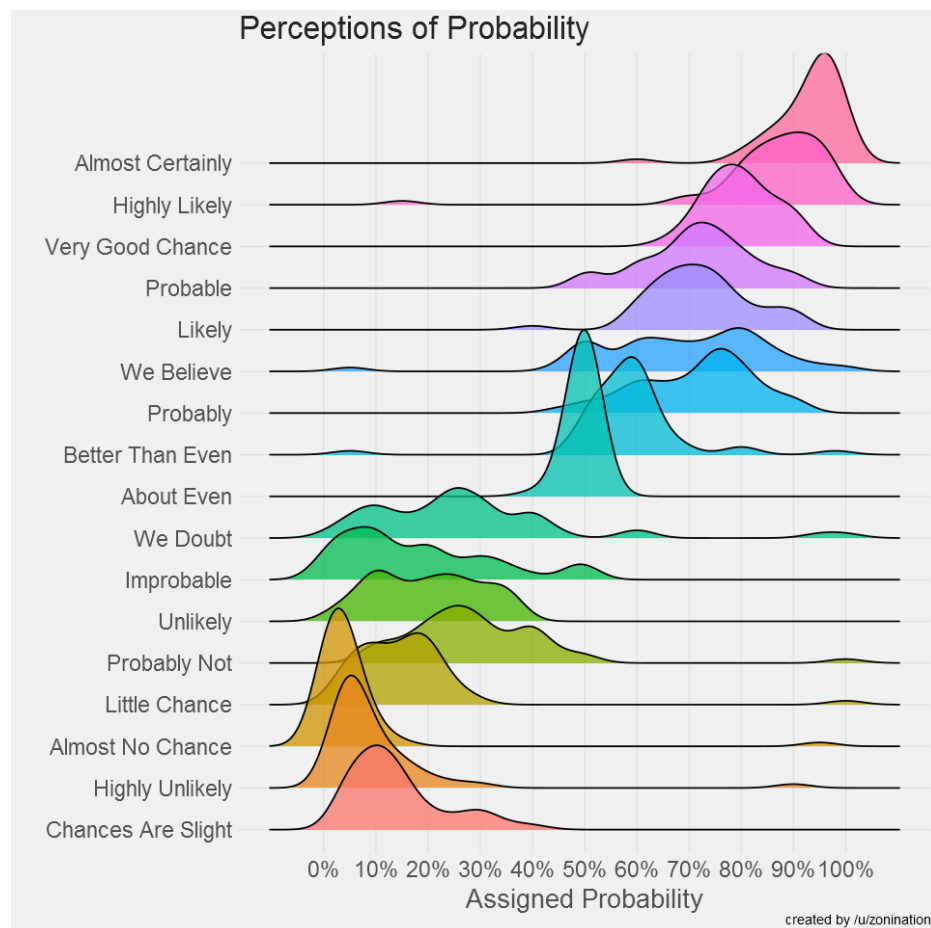
## 10.5 Bayes' rule for combining evidence

Bayes' rule is the correct way to combine prior belief and observation to update beliefs. We always transform from one probability distribution (prior) to a new belief (posterior) using some observed evidence. This can be used to "learn", where "learning" means updating a probability distribution based on observations. It has enormous applications anywhere uncertain information must be fused together, whether from multiple sources (e.g. sensor fusion) or over time (e.g. probabilistic filtering).

# 11  Writing and manipulation of probabilities

Probabilities can be used to represent belief. But the raw numbers (e.g. $P(X = x) = 0.9999$) are not always a useful to make judgments, communicate results, or in some cases, even to do computations.

The graph below shows the probability that a sample of people imagined, given specific verbal cues. The results are very interesting, but notice that the graph has almost no points close 0.0 or 1.0 -- what about 1 in a million events? What about the probability of the sun rising tomorrow?



*[Image credit: Graph by Zonination, from https://github.com/zonination/perceptions, license [MIT] (https://github.com/zonination/perceptions/blob/master/LICENSE)]*

One of the problems is that even if people had responded accurately, the linear visualisation from 0% to 100% makes it very hard to see extreme values.

### 11.0.1 Odds, log odds

The **odds** of an event with probability $p$ is defined by:

$$\text{odds} = \frac{1-p}{p} \quad [\spadesuit]$$

The odds are a more useful unit for discussing unlikely scenarios (odds of 999:1 is easier to understand than $p = 0.001$).

```
In [ ]:  p = 0.001

         def odds(p):
             return (1-p) / (p)

         print("p %.4f => odds %.2f:1" %(p,odds(p)))
```

**Log-odds** or **logit** are particularly useful for very unlikely scenarios:

$$\text{logit}(p) = \log\left(\frac{p}{1-p}\right) \quad [\spadesuit]$$

The logit scales proportionally to the number of zeros in the numerator of the odds.

```
In [ ]:  def logit(p):
             return np.log((p)/(1-p))

         def print_logit(p):
             print("p {p: 12.8f} odds {odds: 12.1f}   log odds {log_odds:12.4f}".

         print_logit(0.1)
         print_logit(0.01)
         print_logit(0.001)
         print_logit(0.00000003)
```

Both of these are typically used to *display* results, rather than to do computations. But **log-probabilities** are widely used for computation as well as display. They help solve *numerical problems* in probability calculations.

### 11.0.2 Log probabilities

The probability of multiple *independent* random variables taking on a set of values can be computed from the product:

$$P(X, Y, Z) = P(X)P(Y)P(Z)$$

and in general

$$P(X_1 = x_i, \ldots, X_n = x_n) = \prod_{i=1}^{n} P(X_i = x_i) \quad [\spadesuit]$$

We often have to have to compute such products, but to multiply lots of values $< 1$ leads to numerical issues: we will get floating point underflow. Instead, it is numerically more reliable to manipulate **log probabilities**, which can be summed instead of multiplied:

$$\log P(x_1, \ldots, x_n) = \sum_{i=1}^{n} \log P(x_i) \quad [\spadesuit]$$

This uses the identity $\log(AB) = \log(A) + \log(B)$

This is simply a numerical convenience which avoids underflow. The **log-likelihood** is just $\log P(B|A)$, and is often more convenient to work with than the raw likelihood.

When talking about **likelihood**, we often write $\mathcal{L}(x_i)$ to mean the likelihood of $x_i$. The likelihood is not a probability. It is a function of data, and $\mathcal{L}(x_i) = f_X(x_i)$

For example, consider the empirical PMF of Romeo and Juliet. This gives the probability of seeing any given character. The likelihood of seeing all of the characters, given our per-character probability model, is:

$$\mathcal{L}(c_1, c_2, \ldots \ c_n) = \prod_{i} P(c_i)$$

This is the likelihood of the text.

```
In [ ]:   def likelihood(sequence, p):
              result = 1.0
              for char in sequence:
                  # this is bad: multiplying lots of small numbers
                  # will underflow very quickly
                  result = result * p[char]
              return result
```

```
In [ ]:   print(likelihood(rj_chars, rj_pmf)) # huh?
```

The log-likelihood does not have this problem with underflow:

$$\log \mathcal{L}(c_1, c_2, \dots \ c_n) = \sum_i \log P(c_i)$$

```
In [59]:   def log_likelihood(sequence, p):
               result = 0.0
               for char in sequence:
                   ## summing logs is stable numerically
                   result = result + np.log(p[char])
               return result
```

```
In [60]:   #
           print(log_likelihood(rj_chars, rj_pmf))
```

```
-462284.1554165428
```

### 11.0.3 Comparing log-likelihoods

We could imagine that writing plays and novels is an activity that mysterious entities do by generating random characters according to a PMF. Under this (very simplified!) assumption, we could now take an "unknown" text (in this case *Macbeth*) and then look at how likely it would have been to have been generated under two models:

- A* It was generated by a mysterious entity using the PMF for *Romeo and Juliet
- B* It was generated by a mysterious entity using the PMF for *Metamorphosis

Neither of these will be exactly true, but we can precisely quantify to what extent *Macbeth* appears to have been generated by a similar process to these two reference texts.

This is a very rough proxy for whether or not they were generated by the same mysterious entity -- i.e. author. Our model is just the distribution of characters, so is a fairly weak model of different styles. However, it is sufficient to do the comparison here.

```
In [71]:   macbeth_chars = numerify("data/macbeth.txt")

           ## A positive value indicates more likely to be "like"
           # Romeo and Juliet; negative is evidence
           ## in favour of Metamorphosis (i.e. written by Kafka)
           print(log_likelihood(macbeth_chars, rj_pmf)-
                 log_likelihood(macbeth_chars, metamorphosis_pmf))
```

```
2744.9718195123132
```

```python
## We could repeat this with another text, in this case
## "The Trial", by Kafka
trial_chars = numerify("data/the_trial.txt")

## A positive value indicates more likely to be "like" Romeo and Juliet
print(log_likelihood(trial_chars , rj_pmf)-
      log_likelihood(trial_chars, metamorphosis_pmf))
```

```
---------------------------------------------------------------------------
NameError                                 Traceback (most recent call last)
<ipython-input-4-3dcfbdb3bb4c> in <module>()
      1 ## We could repeat this with another text, in this case
      2 ## "The Trial", by Kafka
----> 3 trial_chars = numerify("data/the_trial.txt")
      4
      5 ## A positive value indicates more likely to be "like" Romeo and
 Juliet

NameError: name 'numerify' is not defined
```

# 12 Continuous random variables

## 12.1 Problems with continuous variables

Continuous random variables are defined by a PDF, rather than a PMF. A PMF essentially just a vector of values, but a PDF is a function mapping *any* input in its domain to a probability. This seems like a subtle difference (as a PMF has more and more "bins" it gets closer and closer to a PDF, right?), but it has a number of complexities.

- The probability of any specific value is P(X=x)=0$ *zero*, yet any value in the *support* of the distribution function (everywhere the PDF is non-zero) is possible.
- There is no direct way to sample from the PDF in the same way as we did for the PMF. But there are several tricks for sampling from continuous distributions.

- We cannot estimate the true PDF from simple counts of observations like we did for the empirical distribution. This can never "fill in" the PDF, because it will only apply to a single value with zero "width".
- Bayes' Rule is easy to apply to discrete problems. But how do we do computations with continuous PDFs using Bayes' Rule?
- Simple discrete distributions don't have a concept of dimension. But we can have continuous values in $\mathbb{R}$, or in vector spaces $\mathbb{R}^n$, representing the probability of a random variable taking on a vector value, or indeed distributions over other generalisations.
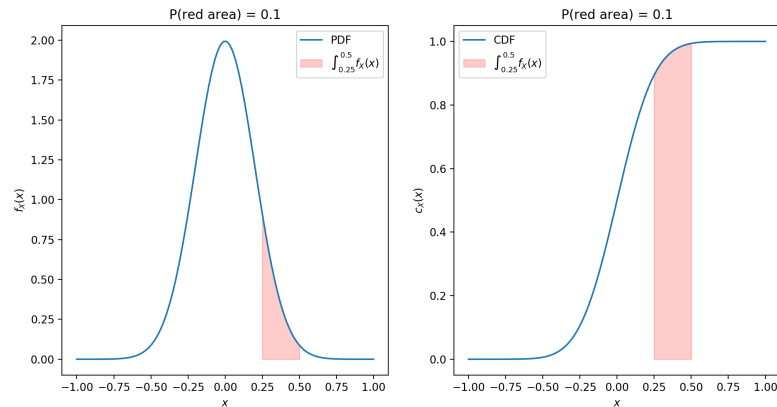
### 12.1.1 Probability distribution functions

The PDF $f_X(x)$ of a random variable $X$ maps a value $x$ (which might be a real number, or a vector, or any other continuous value) to a single number, the density at the point. It is a function $\mathbb{R}^N \to \mathbb{R}^+$, where $\mathbb{R}^+$ is the positive real numbers, and $\int_x f_X(x) = 1$.

- While a PMF can have outcomes with a probability of at most 1, it is *not* the case that the maximum value of a PDF is $f_X(x) \leq 1$ -- *just that the integral of the PDF be 1.*

The value of the PDF at any point is **not** a probability, because the probability of a continuous random variable taking on any specific number must be zero. Instead, we can say that the probability of a continuous random variable $X$ lying in a range $(a, b)$ is:

$$P(X \in (a, b)) = (a < X < b) = \int_a^b f_X(x) \ [\spadesuit]$$

### 12.1.1.1 Support

The **support** of a PDF is the domain it maps from where the density is non-zero.

$$\text{supp}(x) = x \text{ such that } f_X(x) > 0$$

Some PDFs have density over a fixed interval, and have zero density everywhere else. This is true of the uniform distribution, which has a fixed range where the density is constant, and is zero everywhere. This is called **compact support**. We know that sampling from a random variable with this PDF will always give us values in the range of this support. Some PDFs have non-zero density over an infinite domain. This is true of the normal distribution. A sample from a normal distribution could take on *any* real value at all; it's just much more likely to be close to the mean of the distribution than to be far away. This is **infinite support**.

## 12.1.2 Cumulative distribution functions

The **cumulative distribution function** or CDF of a random variable is

$$c_X(x) = \int_{-\infty}^{x} f_X(x) = P(X \le x) \ [\spadesuit]$$

Unlike the PDF, the CDF always maps $x$ to the codomain [0,1]. For any given value $c_X(x)$ tells us how much probability mass there is that is less than or equal to $x$. Given a CDF, we can now answer questions, like: what is the probability that random variable X takes on a value between 3.0 and 4.0?

$$P(3 \le X \le 4) = c_X(4) - c_X(3)$$

This *is* a probability. Sometimes it is more convenient or efficient to do computations with the CDF than with the PDF.

# 12.2 PDF example: the normal disribution

The most ubiquitous of all continuous PDFs is the **normal** or **Gaussian** distribution. It assigns probabilities to real values $x \in \mathbb{R}$ (in other words, a sample space consisting of all of the real numbers). It has a density given by the PDF:

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

You do not need to remember this formula, but it is very useful to know that it is essentially just $e^{-x^2}$ with some scaling factors to normalise it -- this is called the **squared exponential function**.

We use a shorthand notation to refer to the distribution of continuous random variables, which for a normal distribution is:

$$X \sim \mathcal{N}(\mu, \sigma^2) \ [\spadesuit],$$

which is read as

> "Random variable X is distributed as [N]ormal with mean $\mu$ and variance $\sigma^2$"

There are various other symbols used for other continuous distributions, including $\Gamma(\alpha, \beta), \beta(\alpha, \beta), t(\nu), \chi^2(k), \ldots$, which we will not cover in this course. See the references if you want to learn more.

## 12.2.1 Location and scale

The normal distribution places most density close to its center $\mu$ (the "mean"), with a spread defined by $\sigma^2$ (the "variance"). This can be though of as the **location** and **scale** of the density function. Most standard continuous random variable PDFs have a location (where density is concentrated) and scale (how spread out the density is).

```python
import scipy.stats as stats
# Plot the PDF of the normal distribution
def plot_normal():
    # plot the normal (Gaussian distibution) along with a set of points
    x = np.linspace(-4,4,100)
    y = stats.norm.pdf(x,0,1) # mean 0, std. dev. 1
    fig = plt.figure()
    ax = fig.add_subplot(1,1,1)

    ax.plot(x,y, label="PDF")
    ax.axhline(0, color='k', linewidth=0.2) # axis line

    # mark the mean
    ax.text(0, 0.51, '$\mu$')
    ax.axvline(0, color='r')
    # highlight one std. dev. to the right
    ax.axvspan(0,1, facecolor='b', alpha=0.1, label="1 std. dev.")
    ax.text(1.2, 0.3, '$\sigma$')
    # take 1000 random samples and scatter plot them
    samples = stats.norm.rvs(0,1,1000)
    ax.scatter(samples, np.full(samples.shape, .2), s=448, c='b', alpha=
    ax.set_xlabel("$x$")
    ax.set_ylabel("$f_X(x)$")
    ax.legend()

    fig = plt.figure()
    ax = fig.add_subplot(1,1,1)
    ax.plot(x, stats.norm.cdf(x,0,1))
    ax.set_title("Cumulative distribution function")
    ax.set_xlabel("x")
    ax.set_ylabel("$c_X(x)$")
```
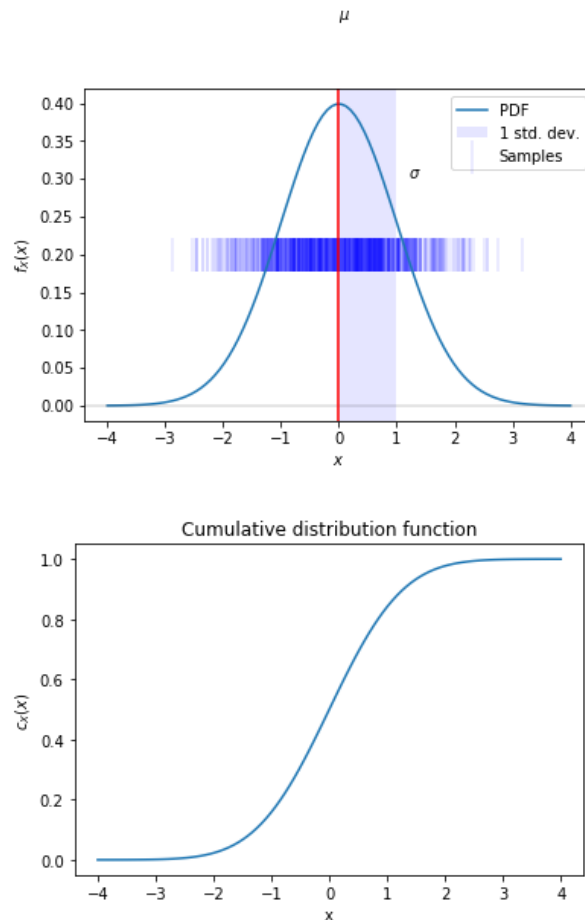
```
plot_normal()
```





#### 12.2.1.1 Normal modelling

It seems that this might be a very limiting choice but there are two good reasons for this:

1. Normal variables have very nice mathematical properties and are easy to work with analyitically (i.e. without relying on numerical computation).
2. The *central limit theorem* tells us that any sum of random variables (however they are distributed) will tend to a *Levy stable distribution* as the number of variables being summed increases. For most random variables encountered, this means the normal distribution (one specific Levy stable distribution).

### 12.2.2 Central limit theorem

If we form a sum of many random variables

$$Y = X_1 + X_2 + X_3 + \ldots,$$

then for almost any PDF that each of $X_1, X_2, \ldots$ might have, the PDF of $Y$ will be approximately normal, $Y \sim \mathcal{N}(\mu, \sigma)$. This means that any process that involves a mixture of many random components will tend to be Gaussian under a wide variety of conditions.

```
def clt():
    # demonstrate the central limit theorem
    for i in [1,2,3,4,8,16,32]:
        x = np.zeros((20000,))
        # add i copies of samples drawn from uniform (flat) distribution
        for j in range(i):
            x += np.random.uniform(-0.5,0.5, x.shape)
        fig = plt.figure()
        ax = fig.add_subplot(1,1,1)
        ax.hist(x/np.sqrt(i), bins=np.linspace(-1.5,1.5,30), density=True
        ax.set_title("(normalised) sum of %d uniform variables" % (i))
        ax.set_xlabel("$x$")
        ax.set_ylabel("$P(x)$")
```

## 12.3 Multivariate distributions: distributions over $\mathbb{R}^n$

Continuous distributions generalise discrete variables (probability mass functions) (e.g. over $\mathbb{Z}$) to continuous spaces over $\mathbb{R}$ via probability density functions.

Probability densities can be further generalised to vector spaces, particularly to $\mathbb{R}^n$. This extends PDFs to assign probability across an entire vector space, under the constraint that the (multidimensional integral)

$$\int_{\mathbf{x}\in\mathbb{R}^n} f_X(\mathbf{x}) = 1, \mathbf{x} \in \mathbb{R}^n.$$

Distributions with PDFs over vector spaces are called **multivariate distributions** (which isn't a very good name; vector distributions might be clearer). In many respects, they work the same as **univariate** continuous distributions. However, they typically require more parameters to specify their form, since they can vary over more dimensions.

### 12.3.1 Multivariate uniform

The multivariate uniform distribution is particularly simple to understand. It assigns equal probability to some (axis-aligned) box in a vector space $\mathbb{R}^n$, such that

$$\int_{\mathbf{x}} f_X(\mathbf{x}) = 1, \mathbf{x} \in \text{a box}.$$

It is trivial to sample from; we just sample *independently* from a one-dimensional uniform distribution in the range [0,1] to get each element of our vector sample. This is a draw from a n-dimensional uniform distribution in the unit box.

```python
def uniform_nd(n, A = None, b=None, d=2):
    if A is None:
        A = np.eye(d)
    if b is None:
        b = np.zeros(d)
    return np.random.uniform(0,1, (n,d)) @ A + b, A, b

def draw_pts_2d(pts, A, b):
    # draw the points
    fig = plt.figure()
    ax = fig.add_subplot(1,1,1)
    ax.plot(pts[:,0], pts[:,1], '.', markersize=4)
    # set sensible scales
    ax.set_xlim(-1,2)
    ax.set_ylim(-1,2)
    ax.axvline(0)
    ax.axhline(0)
    ax.set_aspect("equal")
    # draw a bounding box
    bounding = np.array([[0,0],
                         [0,1],
                         [1,1],
                         [1,0],
                         [0,0]]) @ A + b


    ax.plot(bounding[:,0], bounding[:,1], 'k')
    ax.plot(bounding[0,0], bounding[0,1], 'C1o')

pts, A, b = uniform_nd(100, d=2)
draw_pts_2d(pts, A, b)
plt.gca().set_title("Uniform samples in the range $x\in \\mathbb{R}^2, 6
```
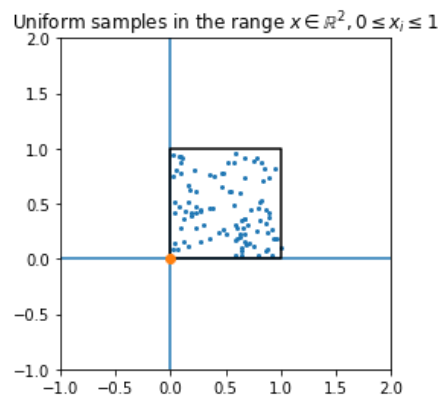
```
Text(0.5,1,'Uniform samples in the range $x\\in \\mathbb{R}^2, 0 \\leq x
_i \\leq 1 $')
```

Uniform samples in the range $x \in \mathbb{R}^2, 0 \leq x_i \leq 1$
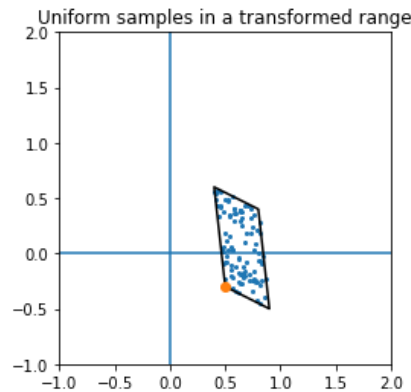


## 12.4 Transformed uniform distribution

If we want to define a distribution over any box, we can simply transform the vectors with a matrix $A$ and shift by adding an offset vector $\mathbf{b}$

```
In [103]:    # A is "spread"
             A = np.array([[0.4, -0.2], [-0.1, 0.9]])
             # b is "offset" or "centre"
             b = np.array([0.5, -0.3])
             pts, A, b = uniform_nd(100, A=A, b=b, d=2)
             draw_pts_2d(pts, A=A, b=b)
             plt.gca().set_title("Uniform samples in a transformed range")
```

```
Text(0.5,1,'Uniform samples in a transformed range')
```



## 12.5 Normal distribution

The normal distribution (above) is very widely used as the distribution of continuous random variables. It can be defined for a random variable of *any dimension*; a **multivariate normal** in statistical terminology. In Unit 5, we saw the idea of a **mean vector** $\mu$ and a **covariance matrix** $\Sigma$ which captured the "shape" of a dataset in terms of an ellipse. *These are in fact the parameterisation of the multivariate normal distribution.*
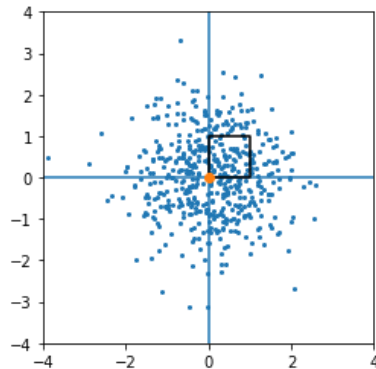
A multivariate normal is fully specified by a mean vector $\mu$ and the covariance matrix $\Sigma$. If you imagine the normal distribution to be a ball shaped mass in space, the mean *translates* the mass, and covariance applies a transformation matrix (scale, rotate and shear) to the ball.

Just like the uniform distribution, we can think of drawing samples from a "unit ball" with an independent normal distribution in each dimension. These samples are transformed linearly by the covariance matrix $\Sigma$ and the mean vector $\mu$, just like $A$ and $b$ above (though $\Sigma$ is actually $A^{-\frac{1}{2}}$ for technical reasons)

In [104]:
```python
def normal_nd(n, A = None, b=None, d=2):
    if A is None:
        A = np.eye(d)
    if b is None:
        b = np.zeros(d)
    return np.random.normal(0,1, (n,d)) @ A + b, A, b

pts, A, b = normal_nd(500, A=None, b=None, d=2)
draw_pts_2d(pts, A=A, b=b)
plt.gca().set_xlim(-4,4)
plt.gca().set_ylim(-4,4)
```
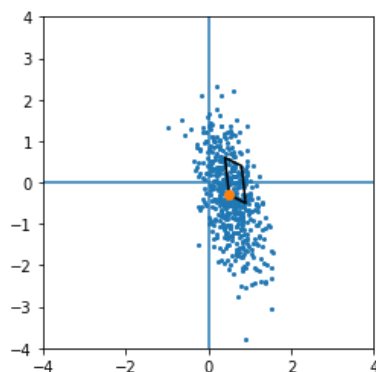
$(-4, \quad 4)$



In [105]:
```python
# A is "spread"
A = np.array([[0.4, -0.2], [-0.1, 0.9]])
# b is "offset" or "centre"
b = np.array([0.5, -0.3])
pts, A, b = normal_nd(500, A=A, b=b, d=2)
draw_pts_2d(pts, A=A, b=b)
plt.gca().set_xlim(-4,4)
plt.gca().set_ylim(-4,4)
```

$(-4, \quad 4)$



### 12.5.1 Joint and marginal PDFs

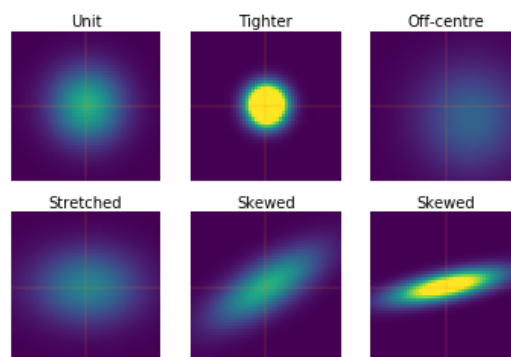We can look at the PDF of a multivariate normals for different covariances and mean vector (centres and spreads).

```python
import scipy.stats
def demo_normal(ax, mean, cov, title):
    x,y = np.meshgrid(np.linspace(-3,3,50), np.linspace(-3,3,50))
    pos = np.empty(x.shape + (2,))
    pos[:,:,0] = x
    pos[:,:,1] = y
    joint_pdf = scipy.stats.multivariate_normal.pdf(pos, mean, cov)
    ax.pcolor(x,y,joint_pdf, cmap='viridis', vmin=0, vmax=0.25)

    ax.axhline(0, color='C1', linewidth=0.2)
    ax.axvline(0, color='C1', linewidth=0.2)
    ax.text(0, 3.2, title, ha='center')
    ax.axis("off")
    ax.axis("image")


fig = plt.figure()
ax = fig.add_subplot(2,3,1)
demo_normal(ax, [0,0], [[1,0],[0,1]], "Unit")
ax = fig.add_subplot(2,3,2)
demo_normal(ax, [0,0], [[0.25,0],[0,0.25]], "Tighter")
ax = fig.add_subplot(2,3,3)
demo_normal(ax, [1,-0.5], [[2,0],[0,2]], "Off-centre")
ax = fig.add_subplot(2,3,4)
demo_normal(ax, [0,0], [[2,0],[0,1]], "Stretched")
ax = fig.add_subplot(2,3,5)
demo_normal(ax, [0,0], [[2,0.1],[1,1]], "Skewed")
ax = fig.add_subplot(2,3,6)
demo_normal(ax, [0,0], [[2,-0.9],[0.4,0.2]], 'Skewed')
```



### 12.5.2 Joint and marginal distributions

We can now talk about the **joint probability density function** (density over all dimensions) and the **marginal probability density function** (density over some sub-selection of dimensions).

For example, consider $X \sim N(\mu, \Sigma), X \in \mathbb{R}^2$, a two dimensional ("bivariate") normal distribution. We can look at some examples of the PDF, showing:

- Joint $P(\mathbf{X})$

- Marginal $P(X_1)$ and $P(X_2)$

- Conditionals $P(X_1|X_2)$ and $P(X_2|X_1)$

```python
import scipy.stats
def joint_marginal(cov):
    # create an independent 2D normal distribution
    x,y = np.meshgrid(np.linspace(-3,3,50), np.linspace(-3,3,50))
    pos = np.empty(x.shape + (2,))
    pos[:,:,0] = x
    pos[:,:,1] = y
    joint_pdf = scipy.stats.multivariate_normal.pdf(pos, [0,0], cov)
    fig = plt.figure()
    # plot the joint
    ax = fig.add_subplot(2,2,1)
    ax.axis('equal')
    ax.set_xlabel("$x_1$")
    ax.set_ylabel("$x_2$")
    plt.title("Joint $P(\\vec{x})=P(x_1,x_2)$")
    ax.pcolor(x,y,joint_pdf, cmap='viridis')
    # plot the marginals
    ax = fig.add_subplot(2,2,3)
    ax.axis('equal')
    plt.title("Marginal $P(x_1) = \int\  P(x_1,x_2) dx_2$")
    ax.plot(x[0,:], np.sum(joint_pdf, axis=0))
    ax = fig.add_subplot(2,2,2)
    ax.axis('equal')
    plt.title("Marginal $P(x_2) = \int\  P(x_1,x_2) dx_1$")
    ax.plot(np.sum(joint_pdf, axis=1), x[0,:])
    plt.tight_layout()

    # plot p(x|y) and p(y|x)
    fig = plt.figure()
    # plot the joint
    ax = fig.add_subplot(2,2,1)
    ax.axis('equal')
    ax.set_xlabel("$x_1$")
    ax.set_ylabel("$x_2$")
    plt.title("Joint $P(\\vec{x})=P(x_1,x_2)$")
    ax.pcolor(x,y,joint_pdf, cmap='viridis')
    ax = fig.add_subplot(2,2,3)
    ax.axis('equal')
    plt.title("Conditional $P(x_1|x_2) = \\frac{P(x_1,x_2)}{p(x_2)}$")
    marginal = np.tile(np.sum(joint_pdf, axis=0), (joint_pdf.shape[0],1))
    ax.pcolor(x,y,joint_pdf/marginal, cmap='viridis')
    plt.tight_layout()
    ax = fig.add_subplot(2,2,2)
    ax.axis('equal')
    plt.title("Conditional $P(x_2|x_1) = \\frac{P(x_1,x_2)}{p(x_1)}$")
    marginal = np.tile(np.sum(joint_pdf, axis=1), (joint_pdf.shape[0],1))
    ax.pcolor(x,y,joint_pdf/marginal.T, cmap='viridis')
    plt.tight_layout()
```
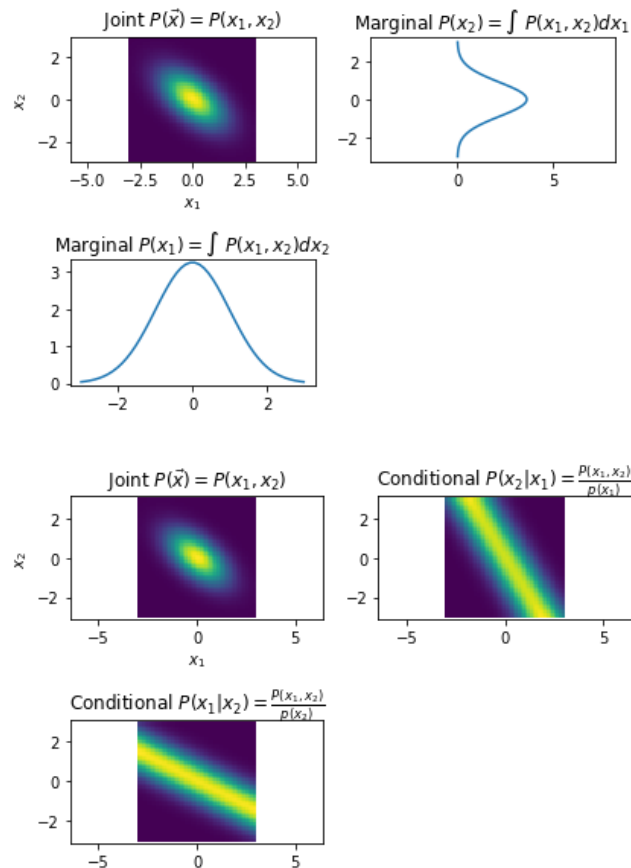
```python
joint_marginal([[1,-0.3],[-0.5,0.8]])
print_matrix("\mu", [0.0,0.0])
print_matrix("\Sigma", [[1,0],[0.5,1]])
plt.tight_layout()
```

$$\mu = \begin{bmatrix} 0.0 & 0.0 \end{bmatrix}$$

$$\Sigma = \begin{bmatrix} 1.0 & 0.0 \\ 0.5 & 1.0 \end{bmatrix}$$



### 12.5.3 Ellipses

When we spoke informally about the covariance matrix "covering" the data with an ellipsoidal shape, we more precisely meant that *if* we represented the data as being generated with a normal distribution, and chose the mean vector and covariance matrix that best approximated that data, then the contours of the density of the PDF corresponding to equally spaced isocontours would be ellipsoidal.

## 12.6 Resources

### 12.6.1 Recommended

- **Probability** by Peter Norvig http://nbviewer.jupyter.org/url/norvig.com/ipython/Probability.ipynb (http://nbviewer.jupyter.org/url/norvig.com/ipython/Probability.ipynb)

- **Chapters 2 and 3 of Information Theory, Inference, and Learning Algorithms** by David Mackay (freely available online)
  - http://www.inference.org.uk/mackay/itprnn/ps/22.40.pdf (http://www.inference.org.uk/mackay/itprnn/ps/22.40.pdf)
  - http://www.inference.org.uk/mackay/itprnn/ps/47.59.pdf (http://www.inference.org.uk/mackay/itprnn/ps/47.59.pdf)

### 12.6.2 Other useful resources

- **Veritasium explains Bayes' Theorem** https://www.youtube.com/watch?v=R13BD8qKeTg (https://www.youtube.com/watch?v=R13BD8qKeTg)
- **Think Bayes** by Allen Downey http://greenteapress.com/wp/think-bayes/ (http://greenteapress.com/wp/think-bayes/) (a freely available book on Bayesian probability theory)
  - **Video by same author** https://www.youtube.com/watch?v=TpgiFIGXcT4 (https://www.youtube.com/watch?v=TpgiFIGXcT4)

- A collection of very readable articles by Count Bayesie
  https://www.countbayesie.com/blog/2016/5/1/a-guide-to-bayesian-statistics
  (https://www.countbayesie.com/blog/2016/5/1/a-guide-to-bayesian-statistics)
- Khan Academy materials on probability https://www.khanacademy.org/math/statistics-probability/probability-library (https://www.khanacademy.org/math/statistics-probability/probability-library) (more in depth than we cover, but high quality stuff)

## 12.7 Beyond this course

- The Probability and Statistics Cookbook (http://pages.cs.wisc.edu/~tdw/files/cookbook-en.pdf)
  which covers most of probability and statistics in 28 incredibly dense pages.
- Probability Theory: The Logic of Science 1st Edition by E. T. Jaynes (this is an excellent but hardcore book; also it is very heavily biased to Bayesian interpretations of probability!)
- A First Course in Probability * *by Sheldon Ross (standard textbook on probability)
- Introduction to statistical learning (http://www-bcf.usc.edu/~gareth/ISL/) (outstanding introduction to statistical learning, including a book, video and course notes)

```
In [ ]:
```