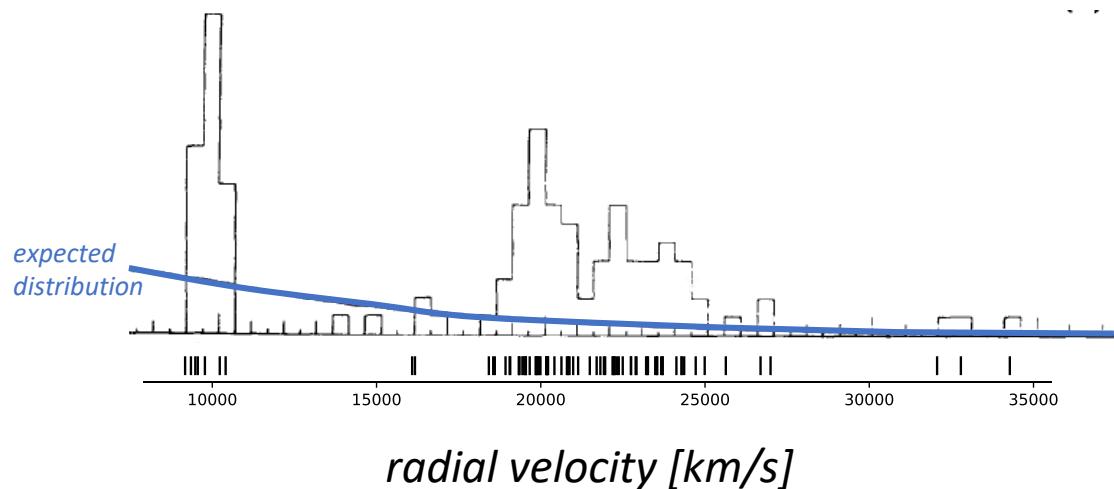


Clustering



This chart shows the distribution of the speeds of 82 galaxies, from a survey of the Corona Borealis region.

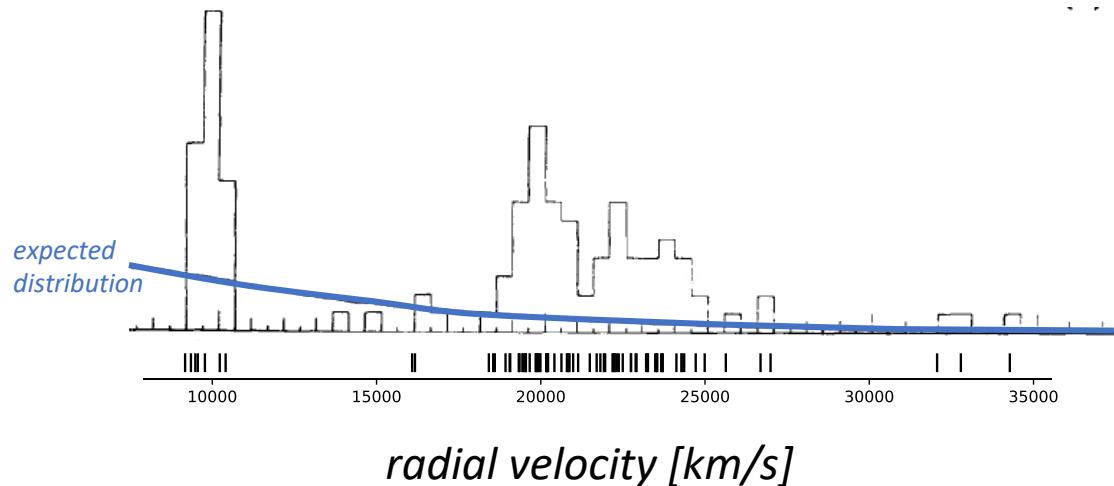
The smooth line shows the expected distribution for a uniform universe. The data however is clumpy, indicating voids and clusters of galaxies.

How can we pick out the clusters?

Since the early 1980s, multi-object spectrographs, CCD detectors, and some dedicated telescopes have allowed the mass production of galaxy redshifts. These large surveys have revealed a very surprising picture of the luminous matter in the Universe. Many astronomers had imagined roughly spherical galaxy clusters floating amongst randomly scattered field galaxies, like meatballs in sauce. Instead, they saw galaxies concentrated into enormous walls and streamers, surrounding huge voids that appear largely empty. The galaxy distribution has been compared to walls of soapy water, surrounding bubbles of air in a basinful of suds; linear filaments appear where the walls of different soap bubbles join, and rich clusters where three or more walls run into each other. A more accurate metaphor is that of a sponge; the voids are interlinked by low-density 'holes' in the walls.

We must remember that in dynamical terms, this structure is very young. Galaxies are not highly concentrated into groups and clusters, in the way that luminous stars are largely confined within galaxies. Averaged over a rich cluster like Coma or Perseus, the density is only 10-100 times greater than the cosmic mean. So the gravitational forces binding a galaxy cluster or group together are relatively weak, and the time for an individual galaxy to move across the system is long. We saw in Section 6.5 that since the Big Bang, galaxies in the outskirts of the Coma cluster have not yet had time to travel to its core; in the Local Group, the large galaxies are falling together for the first time.

Clustering



This chart shows the distribution of the speeds of 82 galaxies, from a survey of the Corona Borealis region.

The smooth line shows the expected distribution for a uniform universe. The data however is clumpy, indicating voids and clusters of galaxies.

How can we pick out the clusters?

```
def rx(p1, p2, p3, μ1, μ2, μ3, σ1, σ2, σ3):  
    u = random.random()  
    if u <= p1:  
        μ, σ = μ1, σ1  
    else if u <= p1+p2:  
        μ, σ = μ2, σ2  
    else:  
        μ, σ = μ3, σ3  
    return numpy.random.normal(loc=μ, scale=σ)
```

Method

1. Invent a probability model i.e. a simulator, parameterized by what we want to learn
2. Write out the likelihood
3. Maximize it

```

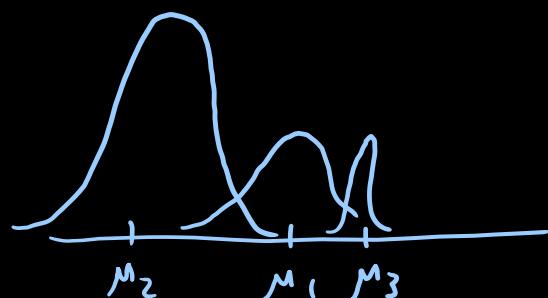
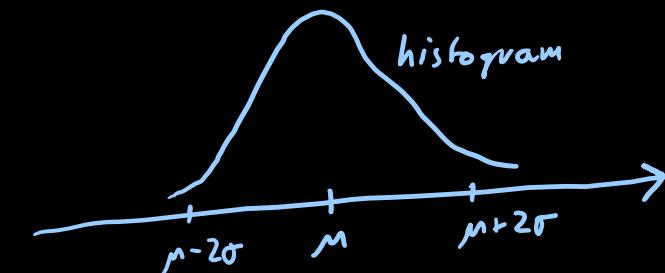
def rx(p1, p2, p3, μ1, μ2, μ3, σ1, σ2, σ3):
    u = random.random()
    if u <= p1:
        μ, σ = μ1, σ1
    else if u <= p1+p2:
        μ, σ = μ2, σ2
    else:
        μ, σ = μ3, σ3
    return numpy.random.normal(loc=μ, scale=σ)

```

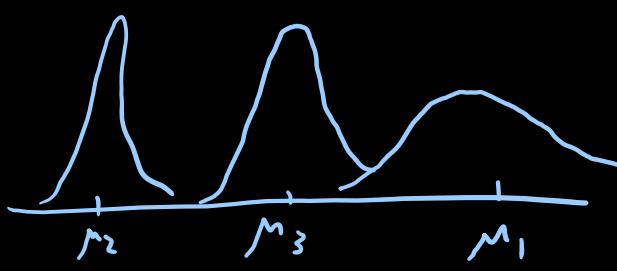
Where did this simulator come from?

`numpy.random.normal(loc=μ, scale=σ)` generates values like it seems fine for describing a single cluster.

If looks like the data might have three clusters, each with different μ and σ.
So, let's pick the μ, σ randomly from a set of choices.



if p₂ is large



if p₁, p₂, p₃ roughly equal

1. Specifying and fitting models

1.3, 1.4. Random variables in code, and in maths

tl;dr. A *random variable* is a function that can give different answers, e.g. a random number generator.

Two random variables are *independent* if knowing the value of one gives no information about the value of the other.

```
def ry(a,b):
    x = random.random()
    y = a * x + b
    return y
```

$$X \sim U[0,1]$$

$$Y = aX + b$$

Generate X from $U[0,1]$
 Sample X from "
 upper case - r.v.
 lower case - constants

```
def ry():
    x1 = random.random()
    x2 = random.random()
    return x1 * math.log(x2)
```

$$X_1, X_2 \sim U[0,1]$$

$$Y = X_1 \log X_2$$

generate x_1 and x_2
 independently from $U[0,1]$

```
def ry():
    (x1,x2) = my_pair_rng()
    return x1 * math.log(x2)
```

$$(X_1, X_2) \sim \dots$$

$$Y = X_1 \log X_2$$

generate x_1 and x_2
 (perhaps) not independently

```
(x1,x2) = [random.random()
            for _ in range(2)]
```

$$X_1, X_2 \sim U[0,1]$$

x_1, x_2 independent.

```
x1 = random.random()
x2 = 1 - random.random()
```

$$X_1, X_2 \sim U[0,1]$$

$$X_1 \sim X_2$$

x_1, x_2 are distributed as $U[0,1]$
 x_1, x_2 are identically distributed

DISCRETE RANDOM VARIABLES

Uniform	$X \sim U\{a, \dots, b\}$	An integer, uniformly distributed in $\{a, \dots, b\}$.	<code>numpy.random.randint(a,b+1)</code>
Geometric	$X \sim \text{Geom}(p)$	Counts the number of failures before a success, where success happens with probability p	<code>numpy.random.geometric(p) - 1</code>
Binomial	$X \sim \text{Bin}(n, p)$	Counts the number of heads in n tosses of a biased coin, where p is the probability of heads	<code>numpy.random.binomial(n,p)</code>
Poisson	$X \sim \text{Pois}(\lambda)$	Used for modelling counts such as the number of buses that arrive in an interval of time.	<code>numpy.random.poisson(\lambda)</code>

CONTINUOUS RANDOM VARIABLES

Uniform	$X \sim U[a, b]$	A floating point value, uniformly distributed in the interval $[a, b]$.	<code>numpy.random.uniform(a,b)</code>
Exponential	$X \sim \text{Exp}(\lambda)$	Used for modelling lifetimes, e.g. the time until the next bus arrives.	<code>numpy.random.exponential(\lambda)</code>
Normal / Gaussian	$X \sim N(\mu, \sigma^2)$	Used to model magnitudes, e.g. the height of a person.	<code>numpy.random.normal(\mu,\sigma)</code>
Beta	$X \sim \text{Beta}(a, b)$	Arises in Bayesian inference.	<code>numpy.random.beta(a,b)</code>

DISCRETE RANDOM VARIABLESUniform
Geometric $X \sim U\{a, \dots, b\}$ $X \sim \text{Geom}(p)$

**the return type is integer
(or, more generally, a
value from some finite /
countable set)**

An integer, uniformly distributed in $\{a, \dots, b\}$.
Counts the number of failures before a success,
where success happens with probability p

Counts the number of heads in n tosses of a
fair coin, where p is the probability of heads
for modelling counts such as the number of events
that arrive in an interval of time.

`numpy.random.randint(a,b+1)``numpy.random.geometric(p) - 1`

Binomial

`numpy.random.binomial(n,p)`

Poisson

`numpy.random.poisson(λ)`**CONTINUOUS RANDOM VARIABLES**

Uniform

 $X \sim U[a, b]$

**the return type is floating
point (and there's also a
smoothness condition...)**

A floating point value, uniformly distributed in
the interval $[a, b]$.

for modelling lifetimes, e.g. the time until
the next bus arrives.

to model magnitudes, e.g. the height of a
person.

in Bayesian inference.

`numpy.random.uniform(a,b)``numpy.random.exponential(λ)``numpy.random.normal(μ,σ)``numpy.random.beta(a,b)`

Exponential

Normal / Gaus

Beta

DISCRETE RANDOM VARIABLES

Uniform

Geometric

Binomial

Poisson

CONTINUOUS

Uniform

Exponential

Normal / Gaussian

Beta

The Normal distribution is used very widely. It's easy to work with, and it's good as an approximation in very many settings.

Useful fact: if

$$X \sim N(\mu, \sigma^2)$$

then

$$aX + b \sim N(a\mu + b, a^2\sigma^2)$$



$\text{Exp}(\lambda)$

$$X \sim N(\mu, \sigma^2)$$

$$X \sim \text{Beta}(a, b)$$

distributed in $\{a, \dots, b\}$.
failures before a success,
with probability p
heads in n tosses of a
the probability of heads
nts such as the number of
interval of time.

uniformly distributed in

`numpy.random.randint(a,b+1)`

`numpy.random.geometric(p) - 1`

`numpy.random.binomial(n,p)`

`numpy.random.poisson(\lambda)`

`numpy.random.uniform(a,b)`

`numpy.random.exponential(\lambda)`

`numpy.random.normal(\mu,\sigma)`

`numpy.random.beta(a,b)`

Specifying numerical random variables

Any random variable X can be specified by its *probability distribution*,

$$\mathbb{P}(X \in A) \text{ for every set } A.$$

Numerical random variables can equivalently be specified by their *cumulative distribution function*,

$$F(x) = \mathbb{P}(X \leq x) \text{ for every } x.$$

If X is a discrete random variable (integer-valued), $F(x)$ is a step function,

$$F(x) = \mathbb{P}(X \leq x) = \sum_{y=-\infty}^{\text{floor}(x)} \mathbb{P}(X = y).$$

X is a continuous random variable if $F(x)$ is differentiable. If so, the *density function* for X is defined to be $f(x) = F'(x)$, and

$$F(x) = \mathbb{P}(X \leq x) = \mathbb{P}(X < x) = \int_{y=-\infty}^x f(y) dy.$$

For a continuous random variable X ,
 $\mathbb{P}(X = x) = 0$ for all x .

$$\begin{aligned} \text{So } \mathbb{P}(X \leq x) &= \mathbb{P}(X < x \text{ or } X = x) \\ &= \mathbb{P}(X < x) + \mathbb{P}(X = x) \\ &= \mathbb{P}(X < x). \end{aligned}$$

Exercise (The uniform distribution).

For X generated by `random.uniform(a, b)`, which generates a floating point number uniformly distributed in $[a, b]$, sketch $\mathbb{P}(X \leq x)$ and derive the density.

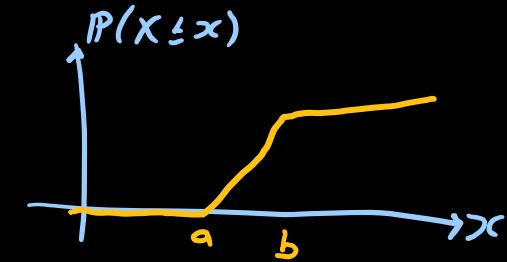
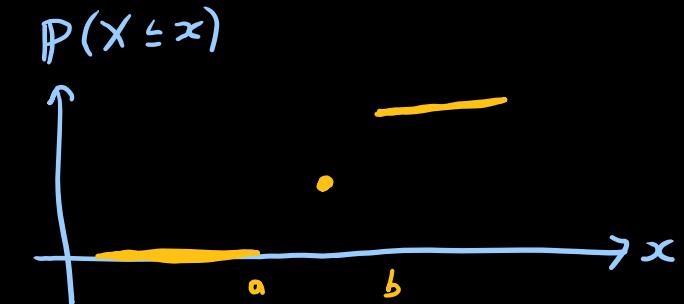
There's no general method for this. You just have to use common sense, and think through all the cases carefully -

Let's try some cases:

- If $x < a$ then $\mathbb{P}(X \leq x) = 0$ since values $< a$ are impossible
- If $x > b$ then $\mathbb{P}(X \leq x) = 1$ since all return values are $\leq b$
- If $x = \frac{a+b}{2}$ then $\mathbb{P}(X \leq x) = \frac{1}{2}$ since all values in $[a, b]$ are equally likely

Generalizing, we can spot that $\mathbb{P}(X \leq x) = \begin{cases} \frac{x-a}{b-a} & \text{for } x \in [a, b] \\ 0 & \text{for } x < a \\ 1 & \text{for } x > b \end{cases}$

The density is $f(x) = \frac{d}{dx} \mathbb{P}(X \leq x) = \begin{cases} \frac{1}{b-a} & \text{for } x \in [a, b] \\ 0 & \text{else} \end{cases}$



Exercise (The uniform distribution).

For X generated by `random.uniform(a, b)`, which generates a floating point number uniformly distributed in $[a, b]$, sketch $\mathbb{P}(X \leq x)$ and derive the density.

We derived the density

$$f(x) = \begin{cases} \frac{1}{b-a} & \text{if } x \in [a, b] \\ 0 & \text{otherwise.} \end{cases}$$

This is commonly written with an indicator function

$$f(x) = \frac{1}{b-a} \mathbf{1}_{x \in [a, b]}$$

The indicator function maps boolean to integer: $\mathbf{1}_{\text{true}} = 1, \quad \mathbf{1}_{\text{false}} = 0$.

We'll use a unified density notation:

$$\Pr_X(x) = \begin{cases} \mathbb{P}(X = x) & \text{if } X \text{ is a discrete r.v.} \\ f(x) & \text{if } X \text{ is continuous, density } f \end{cases}$$

With this notation, it's cleaner to write the various rules...

In lecture 1, we said "maximize the probability of the observed data". But for continuous random variables, $\mathbb{P}(\text{observed data} | \theta) = 0$ for any θ , so we need to use \Pr not \mathbb{P} .

Maximum likelihood estimation.

The maximum likelihood estimator for an unknown parameter θ is that value which maximizes

$$\Pr(\text{observed data} | \theta)$$

Joint density and independence

Two random variables X and Y have a *joint density*, written $\Pr_{X,Y}(x,y)$.

For discrete random variables, it is defined as

$$\Pr_{X,Y}(x,y) = \mathbb{P}(X = x \text{ and } Y = y)$$

For continuous random variables, ... (it's trickier — left to a later lecture)

Independence. Two random variables X and Y are independent if

$$\Pr_{X,Y}(x,y) = \Pr_X(x) \Pr_Y(y) \text{ for all } x \text{ and } y$$

or equivalently if

$$\mathbb{P}(X \in A \text{ and } Y \in B) = \mathbb{P}(X \in A) \mathbb{P}(Y \in B) \text{ for all sets } A \text{ and } B$$

Handling numerical random variables

- To fit a model (i.e. to learn its parameters from data), we need the likelihood, which needs $\Pr_X(x)$
- If the random variable is continuous, then $\Pr_X(x)$ is defined to be the derivative of $\mathbb{P}(X \leq x)$, so we should work out $\mathbb{P}(X \leq x)$ first.

Exercise 3.3 (Transforming random variables).

Let $X \sim U[0,1]$ and let $Y = X^2$. Find the density of Y .

we want a formula for $F(y) = \mathbb{P}(Y \leq y)$.

$$\begin{array}{ll} \text{If } y < 0: \quad \mathbb{P}(Y \leq y) = 0. \\ \text{If } y > 1: \quad \mathbb{P}(Y \leq y) = 1 \end{array} \quad \left. \begin{array}{l} \\ \end{array} \right\} \text{by common sense}$$

$$\begin{aligned} \text{If } y \in [0,1]: \quad \mathbb{P}(Y \leq y) &= \mathbb{P}(X^2 \leq y) = \mathbb{P}(X \leq \sqrt{y}) && \leftarrow \text{good strategy: rewrite the probability you want to work out in terms of simpler random variables, e.g. the standard random variables listed earlier} \\ &= \frac{\sqrt{y} - 0}{1 - 0} && \leftarrow \text{we've already worked out the c.d.f. for a uniform } [a,b] \text{ random variable} \\ &= \sqrt{y} \end{aligned}$$

So the density is $\Pr_Y(y) = \frac{d}{dy} \mathbb{P}(Y \leq y) = \frac{1}{2} y^{-\frac{1}{2}} \mathbf{1}_{y \in [0,1]}$

1. Specifying and fitting models

1.5. Learning generative models

tl;dr. Given a dataset x_1, \dots, x_n can we design a random number generator (i.e. a probability model) that might have generated it?

1. Choose a distribution with tuneable parameters, call it X . We want x_1, \dots, x_n to look like independent samples from X .
2. Write out the likelihood

$$\text{lik}(\theta|x_1, \dots, x_n) = \Pr_X(x_1|\theta) \times \dots \times \Pr_X(x_n|\theta)$$
3. Fit the model—i.e. find the maximum likelihood estimator for θ .

The likelihood is the \Pr of the observed data:

$$\text{lik}(\theta|x_1, \dots, x_n)$$

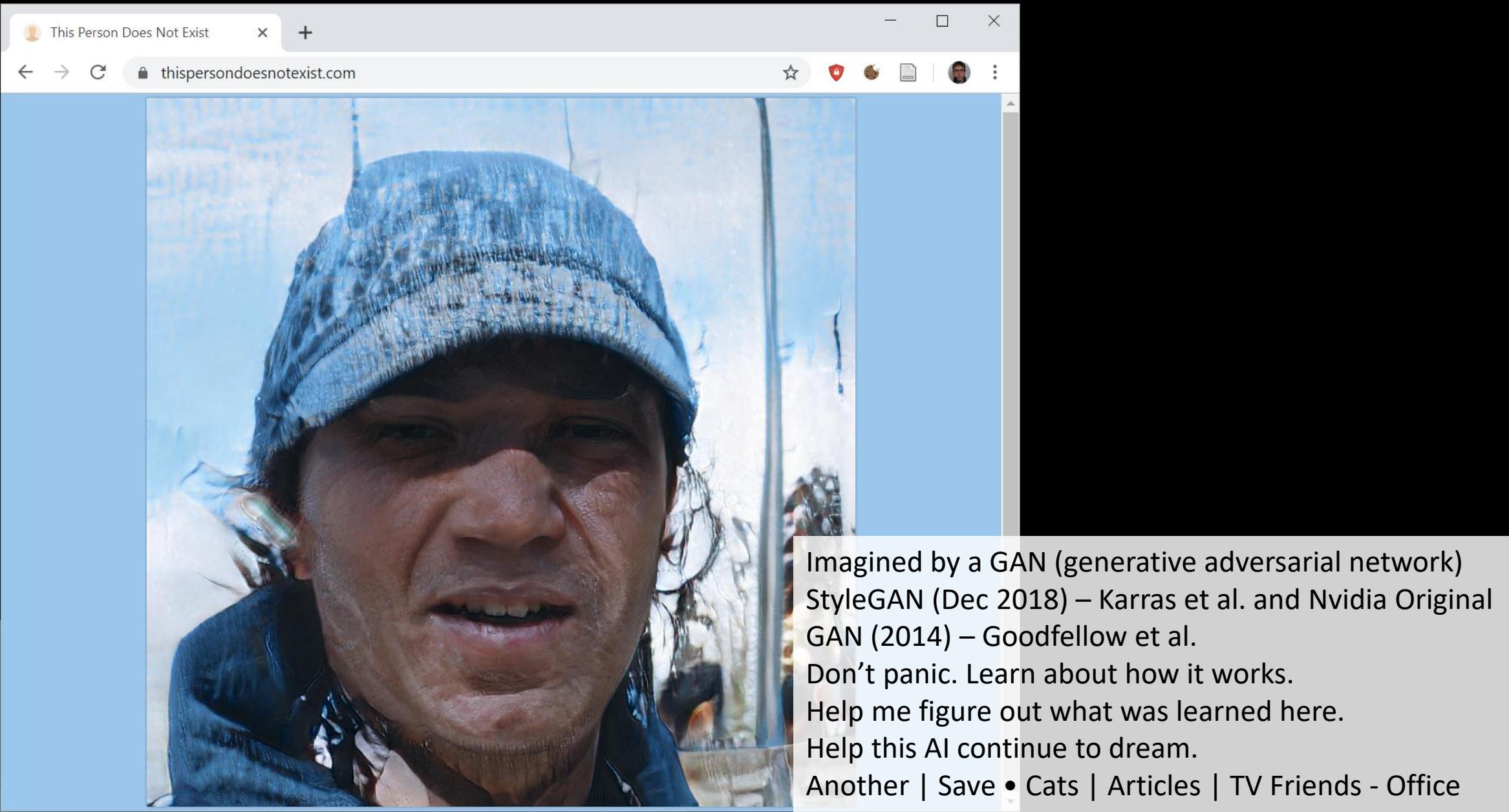
$$= \Pr(x_1, \dots, x_n | \theta)$$

$$= \Pr(x_1|\theta) \times \dots \times \Pr(x_n|\theta)$$

by independence.



<https://thispersondoesnotexist.com/>



Exercise 1.7 (coin tosses).

We take a biased coin and toss it $n = 10$ times, and observe the outcomes

$$(x_1, \dots, x_{10}) = (H, H, t, t, H, H, t, t, H, H).$$

Fit the probability model

$$X = \begin{cases} H & \text{with probability } p \\ t & \text{with probability } 1 - p \end{cases}$$

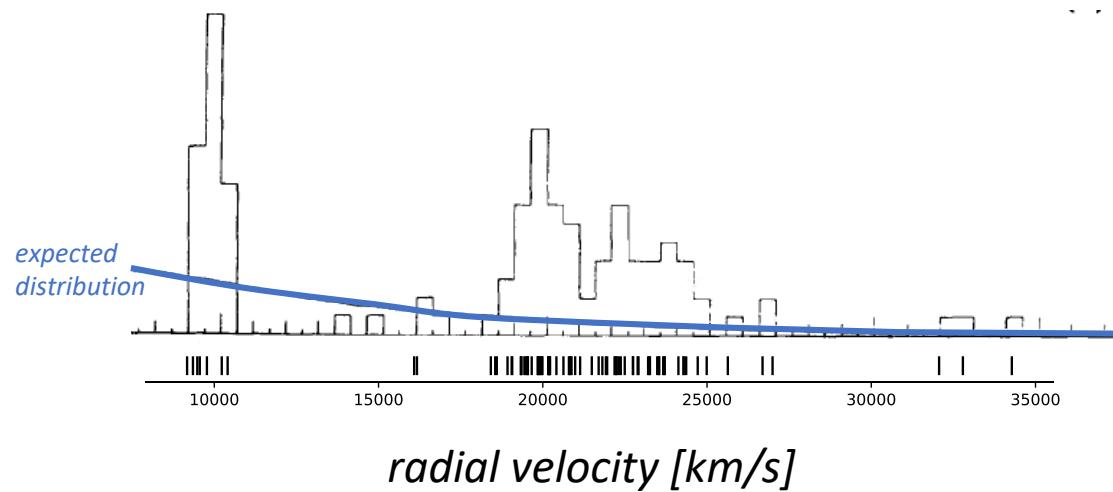
For a single observation: $\Pr_X(x|p) = \begin{cases} p & \text{if } x=H \\ 1-p & \text{if } x=t \end{cases}$

$$\begin{aligned} \text{So } \text{lik}(p|x_1, \dots, x_n) &= \Pr(x_1, \dots, x_n|p) = \Pr(x_1|p) \times \dots \times \Pr(x_n|p) = \prod_{i=1}^n \Pr_X(x_i|p) = \prod_{i=1}^n \begin{cases} p & \text{if } x_i=H \\ 1-p & \text{if } x_i=t \end{cases} \\ &= p^y (1-p)^{n-y} \quad \text{where } y = \# \text{ heads}. \end{aligned}$$

$$\Rightarrow \log \text{lik}(p|x_1, \dots, x_n) = y \log p + (n-y) \log (1-p).$$

To fit the model, i.e. find the maximum likelihood estimator \hat{p} ,

$$\frac{d}{dp} \log \text{lik} = \frac{y}{p} - \frac{n-y}{1-p} = 0 \Rightarrow \hat{p} = \frac{y}{n} = 0.6 \quad \text{for this dataset.}$$



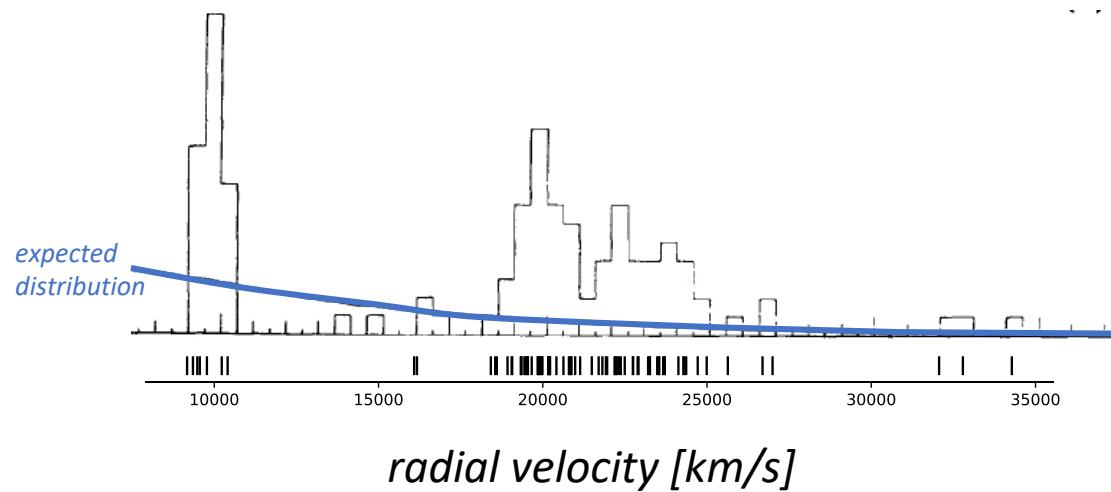
This chart shows the distribution of the speeds of 82 galaxies, from a survey of the Corona Borealis region.

The smooth line shows the expected distribution for a uniform universe. The data however is clumpy, indicating voids and clusters of galaxies.

How can we pick out the clusters?

Method

1. Invent a probability model i.e. a simulator, parameterized by what we want to learn
2. Write out the likelihood
3. Maximize it



This chart shows the distribution of the speeds of 82 galaxies, from a survey of the Corona Borealis region.

The smooth line shows the expected distribution for a uniform universe. The data however is clumpy, indicating voids and clusters of galaxies.

How can we pick out the clusters?

```
def rx(p1,p2,p3, μ1,μ2,μ3, σ1,σ2,σ3):
    u = random.random()
    if u <= p1:
        μ,σ = μ1,σ1
    else if u <= p1+p2:
        μ,σ = μ2,σ2
    else:
        μ,σ = μ3,σ3
    return numpy.random.normal(loc=μ, scale=σ)
```

In maths notation,

$$K = \begin{cases} 1 & \text{with probability } p_1 \\ 2 & \text{with probability } p_2 \\ 3 & \text{with probability } p_3 \end{cases}$$

$$X \sim N(\mu_K, \sigma_K^2)$$

- Invent a probability model, i.e. a simulator, parameterized by what we want to learn

$$K = \begin{cases} 1 & \text{with probability } p_1 \\ 2 & \text{with probability } p_2 \\ 3 & \text{with probability } p_3 \end{cases}$$

$$X \sim N(\mu_K, \sigma_K^2)$$

- Write out the likelihood $= \Pr(x_1, \dots, x_n \mid \text{parameters})$
 $= \prod_{i=1}^n \Pr(x_i \mid \text{params})$ by independence
 Figure out \Pr by calculating $\frac{d}{dx} \Pr(X \leq x)$ (see example sheet 1)

- Maximize it

Use numerical optimization (`scipy.optimize.fmin`) to maximize over the nine unknown parameters

$$\mu_1, \mu_2, \mu_3 \in \mathbb{R}$$

$\sigma_1, \sigma_2, \sigma_3 > 0$ so use the $\sigma = e^\tau$ trick and optimize over $\tau \in \mathbb{R}$

$p_1, p_2, p_3 \in [0, 1]$, $p_1 + p_2 + p_3 = 1$, so use the softmax trick from Lecture 1 / exercise 1.5

```

import scipy.stats
Pr = scipy.stats.norm.pdf

def loglik(theta, x):
    xi1, xi2, mu1, mu2, mu3, tau1, tau2, tau3 = theta
    p = numpy.exp([xi1, xi2, 0])
    [p1, p2, p3] = p / sum(p)
    sigma1, sigma2, sigma3 = numpy.exp([tau1, tau2, tau3])
    loglik = p1 * Pr(x, loc=mu1, scale=sigma1) + p2 * Pr(x, loc=mu2, scale=sigma2) + p3 * Pr(x, loc=mu3, scale=sigma3)
    return numpy.log(loglik)

```

Watch out for sensitivity on the initial guess! Later on in the course, we'll learn strategies for reasoning about confidence in the answer.

```

initial_guess = [-1, -1, 10000, 20000, 24000, math.log(1000), math.log(5000), math.log(8000)]
theta_hat = scipy.optimize.fmin(lambda theta: -numpy.sum(loglik(theta, galaxies)), initial_guess, maxiter=5000)

```

```

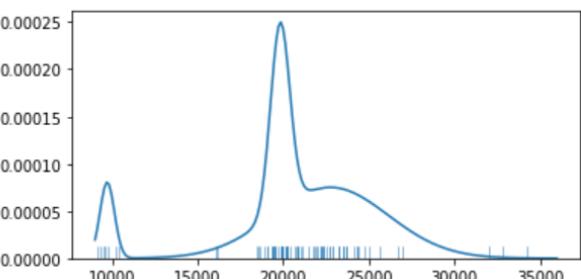
Optimization terminated successfully.
    Current function value: 776.169426
    Iterations: 1360
    Function evaluations: 2008

```

```

fig, ax = plt.subplots(figsize=(6, 3))
x = numpy.linspace(9000, 36000, 200)
f = numpy.exp(loglik(theta_hat, x))
ax.plot(x, f)
for x in galaxies:
    ax.plot([x, x], [0, .05 * max(f)], color='#1f77b4', lw=0.5)
ax.set_xlim(left=0)
plt.show()

```



(Code snippet on Azure Notebooks.)

Implementation note:
 my $\loglik(\theta | \mathbf{x})$ function returns a
 list $[\loglik(\theta | x_1), \dots, \loglik(\theta | x_n)]$
 and the sum is done later.

That's handy when it
 comes to plotting the $\Pr(x | \theta)$
 function — my \loglik/\logPr
 function is already vectorized.