

MATH 3610 – 02

Single population growth models

Introduction to modelling

Julien Arino
University of Manitoba
julien.arino@umanitoba.ca
Fall 2024

The University of Manitoba campuses are located on original lands of Anishinaabeg, Ininew, Anisininew, Dakota and Dene peoples, and on the National Homeland of the Red River Métis.

We respect the Treaties that were made on these territories, we acknowledge the harms and mistakes of the past, and we dedicate ourselves to move forward in partnership with Indigenous communities in a spirit of Reconciliation and collaboration.

Outline

The data – US census

Fitting a curve to the data

Least squares problems

Population growth – Logistic equation

The delayed logistic equation

The logistic map

Objective of this part

In this set of slides, we introduce some of the basic concepts of population growth models

We introduce some of the basic concepts of mathematical modelling and some of the questions that will be considered during the course

The data – US census

Fitting a curve to the data

Least squares problems

Population growth – Logistic equation

The delayed logistic equation

The logistic map

Objective

We are given a table with the population census at different time intervals between a date a and a date b , and want to get an expression for the population. This allows us to:

- ▶ compute a value for the population at any time between the date a and the date b (**interpolation**)
- ▶ predict a value for the population at a date before a or after b (**extrapolation**)

PROCEEDINGS
OF THE
NATIONAL ACADEMY OF SCIENCES

Volume 6

JUNE 15, 1920

Number 6

*ON THE RATE OF GROWTH OF THE POPULATION OF THE
UNITED STATES SINCE 1790 AND ITS MATHEMATICAL
REPRESENTATION¹*

BY RAYMOND PEARL AND LOWELL J. REED

DEPARTMENT OF BIOMETRY AND VITAL STATISTICS, JOHNS HOPKINS UNIVERSITY

Read before the Academy, April 26, 1920

SHOWING THE DATES OF THE TAKING OF THE CENSUS AND THE RECORDED POPULATIONS
FROM 1790 TO 1910

Year	DATE OF CENSUS	RECORDED POPULATION (REVISED FIGURES FROM STATISTICAL ABST., 1918)
	Month and Day	
1790	First Monday in August	3,929,214
1800	First Monday in August	5,308,483
1810	First Monday in August	7,239,881
1820	First Monday in August	9,638,453
1830	June 1	12,866,020
1840	June 1	17,069,453
1850	June 1	23,191,876
1860	June 1	31,443,321
1870	June 1	38,558,371
1880	June 1	50,155,783
1890	June 1	62,947,714
1900	June 1	75,994,575
1910	April 15	91,972,266

USA census from 1790 to 1910

Although we have data up to 2020, we use the data up to 1910 like Pearl & Reed
(note that there were some corrections to the census since the paper of Pearl & Reed)

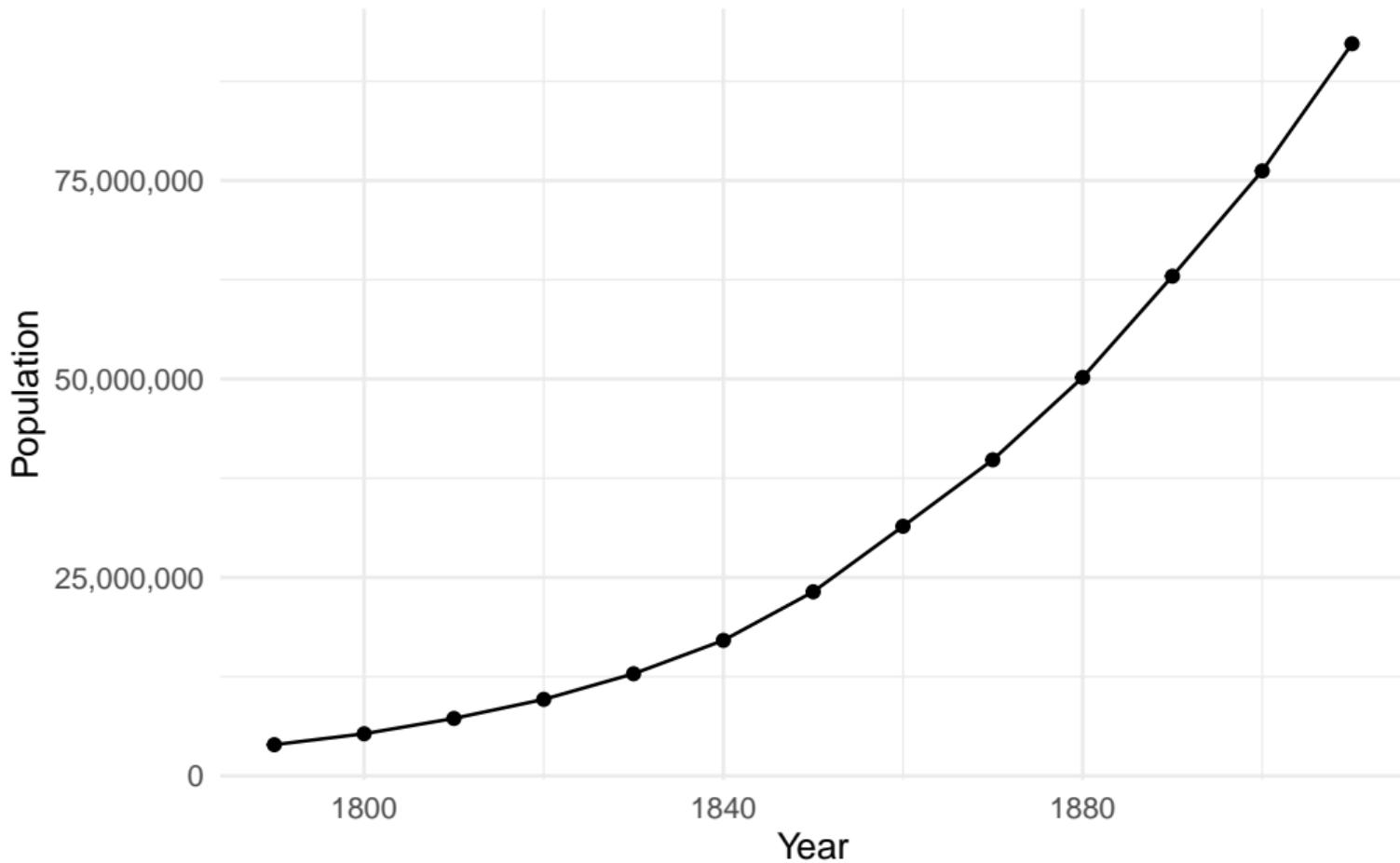
Year	Population	Year	Population	Year	Population
1790	3,929,326	1840	17,069,458	1890	62,947,714
1800	5,308,483	1850	23,191,876	1900	76,212,168
1810	7,239,881	1860	31,443,321	1910	92,228,496
1820	9,638,453	1870	39,818,449		
1830	12,866,020	1880	50,189,209		

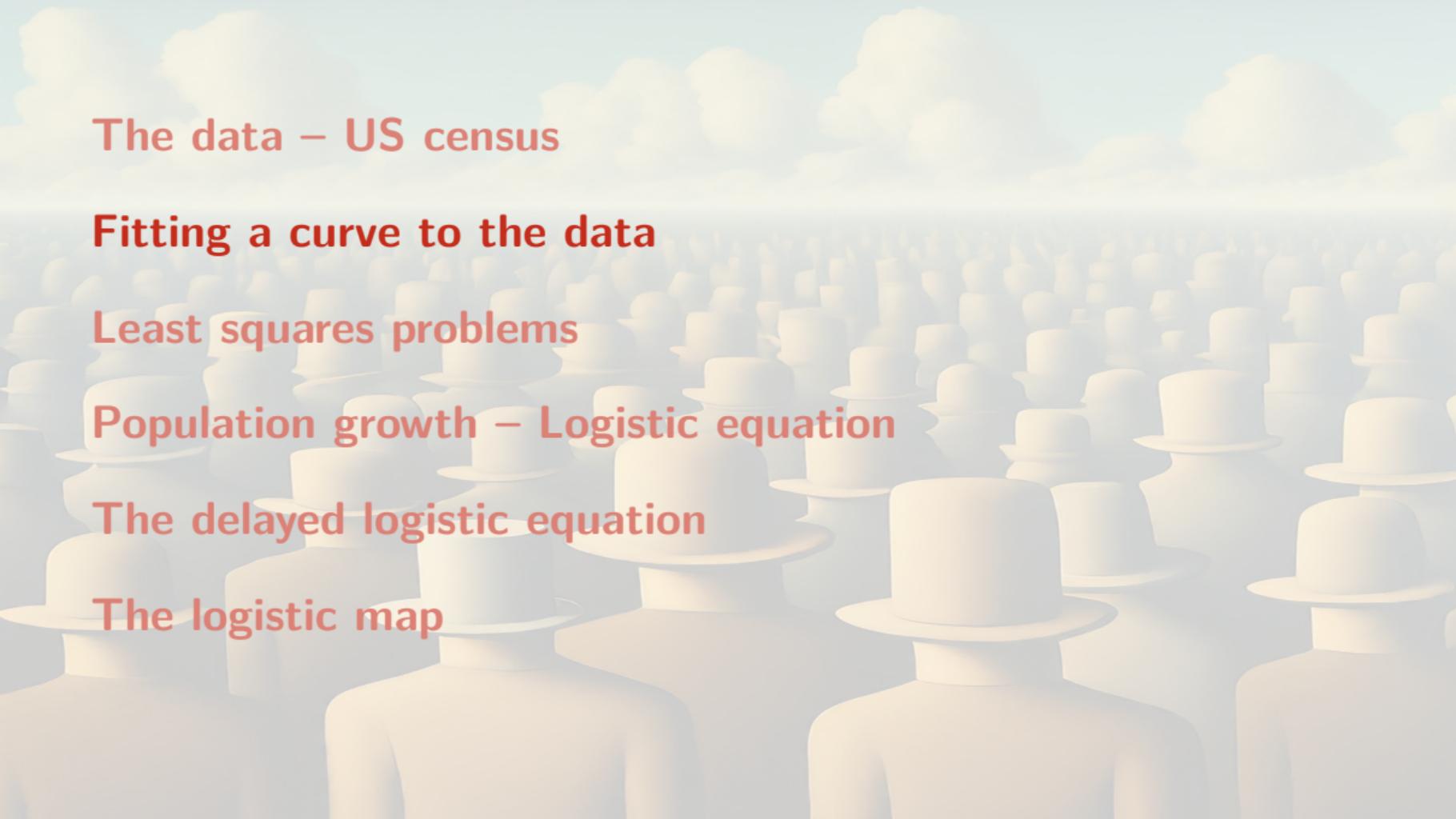
Plot the data !!!

It is always a good idea to plot the data before trying to do anything with it

```
plot_USA_census_to_1910 =  
  ggplot(USA_census_to_1910, aes(x=Year, y=Population)) +  
    geom_line() +  
    geom_point() +  
    labs(title="US population from 1790 to 1910",  
         x="Year",  
         y="Population") +  
    theme_minimal()  
  print(plot_USA_census_to_1910)
```

US population from 1790 to 1910





The data – US census

Fitting a curve to the data

Least squares problems

Population growth – Logistic equation

The delayed logistic equation

The logistic map

The background of the slide features a dense, sprawling landscape filled with numerous white human skulls. The skulls are scattered across a dry, cracked ground, stretching towards a range of mountains in the distance under a pale, overcast sky.

Fitting a curve to the data

Fitting a quadratic curve to the data

Some similar curves

Population curves – Gompertz

First idea – This looks quadratic!

The curve looks like a piece of a parabola. So let us “fit” a curve of the form

$$P(t) = a + bt + ct^2$$

This means we want to find coefficients a, b, c such that the curve $P(t)$ is as close as possible to the data points

The data points

Year	Population	Year	Population	Year	Population
1790	3,929,326	1840	17,069,458	1890	62,947,714
1800	5,308,483	1850	23,191,876	1900	76,212,168
1810	7,239,881	1860	31,443,321	1910	92,228,496
1820	9,638,453	1870	39,818,449		
1830	12,866,020	1880	50,189,209		

We have 13 data points (t_k, P_k) , $k = 1, \dots, 13$, e.g., $(t_1, P_1) = (1790, 3929214)$, $(t_2, P_2) = (1800, 5308483)$, etc.

Some of you are familiar with this problem

If you have taken MATH 2740 (Math of Data Science), you have seen this before!

See the notes on the course website for a refresher on this problem here and the corresponding videos [here](#), [here](#), [here](#), [here](#) and [here](#)

(Sorry about the number of videos, I need to reorganize them!)

To do this, we want to minimize

$$S = \sum_{k=1}^{13} (P(t_k) - P_k)^2$$

where t_k are the known dates, P_k are the known populations, and $P(t_k) = a + bt_k + ct_k^2$

The t_k and P_k are known, a, b, c are to be found, so we write S as a function of a, b, c :
 $S(a, b, c)$

Recall your multivariable calculus:

$$S = S(a, b, c) = \sum_{k=1}^{13} (a + bt_k + ct_k^2 - P_k)^2$$

is maximal if (necessary condition) $\partial S / \partial a = \partial S / \partial b = \partial S / \partial c = 0$

We have

$$\frac{\partial S}{\partial a} = 2 \sum_{k=1}^{13} (a + bt_k + ct_k^2 - P_k)$$

$$\frac{\partial S}{\partial b} = 2 \sum_{k=1}^{13} (a + bt_k + ct_k^2 - P_k) t_k$$

$$\frac{\partial S}{\partial c} = 2 \sum_{k=1}^{13} (a + bt_k + ct_k^2 - P_k) t_k^2$$

Thus, we want

$$\frac{\partial S}{\partial a} = 0 \iff 2 \sum_{k=1}^{13} (a + bt_k + ct_k^2 - P_k) = 0$$

$$\frac{\partial S}{\partial b} = 0 \iff 2 \sum_{k=1}^{13} (a + bt_k + ct_k^2 - P_k)t_k = 0$$

$$\frac{\partial S}{\partial c} = 0 \iff 2 \sum_{k=1}^{13} (a + bt_k + ct_k^2 - P_k)t_k^2 = 0$$

that is

$$\sum_{k=1}^{13} (a + bt_k + ct_k^2 - P_k) = 0$$

$$\sum_{k=1}^{13} (a + bt_k + ct_k^2 - P_k)t_k = 0$$

$$\sum_{k=1}^{13} (a + bt_k + ct_k^2 - P_k)t_k^2 = 0$$

Rearranging the system

$$\sum_{k=1}^{13} (a + bt_k + ct_k^2 - P_k) = 0$$

$$\sum_{k=1}^{13} (a + bt_k + ct_k^2 - P_k) t_k = 0$$

$$\sum_{k=1}^{13} (a + bt_k + ct_k^2 - P_k) t_k^2 = 0$$

we get

$$\sum_{k=1}^{13} (a + bt_k + ct_k^2) = \sum_{k=1}^{13} P_k$$

$$\sum_{k=1}^{13} (at_k + bt_k^2 + ct_k^3) = \sum_{k=1}^{13} P_k t_k$$

$$\sum_{k=1}^{13} (at_k^2 + bt_k^3 + ct_k^4) = \sum_{k=1}^{13} P_k t_k^2$$

$$\sum_{k=1}^{13} (a + bt_k + ct_k^2) = \sum_{k=1}^{13} P_k$$

$$\sum_{k=1}^{13} (at_k + bt_k^2 + ct_k^3) = \sum_{k=1}^{13} P_k t_k$$

$$\sum_{k=1}^{13} (at_k^2 + bt_k^3 + ct_k^4) = \sum_{k=1}^{13} P_k t_k^2$$

after a bit of tidying up, takes the form

$$\left(\sum_{k=1}^{13} 1 \right) a + \left(\sum_{k=1}^{13} t_k \right) b + \left(\sum_{k=1}^{13} t_k^2 \right) c = \sum_{k=1}^{13} P_k$$

$$\left(\sum_{k=1}^{13} t_k \right) a + \left(\sum_{k=1}^{13} t_k^2 \right) b + \left(\sum_{k=1}^{13} t_k^3 \right) c = \sum_{k=1}^{13} P_k t_k$$

$$\left(\sum_{k=1}^{13} t_k^2 \right) a + \left(\sum_{k=1}^{13} t_k^3 \right) b + \left(\sum_{k=1}^{13} t_k^4 \right) c = \sum_{k=1}^{13} P_k t_k^2$$

So the aim is to solve the linear system

$$\begin{pmatrix} 13 & \sum_{k=1}^{13} t_k & \sum_{k=1}^{13} t_k^2 \\ \sum_{k=1}^{13} t_k & \sum_{k=1}^{13} t_k^2 & \sum_{k=1}^{13} t_k^3 \\ \sum_{k=1}^{13} t_k^2 & \sum_{k=1}^{13} t_k^3 & \sum_{k=1}^{13} t_k^4 \end{pmatrix} \begin{pmatrix} a \\ b \\ c \end{pmatrix} = \begin{pmatrix} \sum_{k=1}^{13} P_k \\ \sum_{k=1}^{13} P_k t_k \\ \sum_{k=1}^{13} P_k t_k^2 \end{pmatrix}$$

With R, this is easy to solve..?

```
> t = as.double(USA_census_to_1910$Year)
> pop = as.double(USA_census_to_1910$Population)
> A = matrix(c(13, sum(t), sum(t^2),
+               sum(t), sum(t^2), sum(t^3),
+               sum(t^2), sum(t^3), sum(t^4)),
+               nrow=3, byrow=TRUE)
> b = c(sum(pop),
+       sum(pop * t),
+       sum(pop * t^2))
> sol = try(solve(A,b))
> writeLines(sol)

Error in solve.default(A, b) :
  system is computationally singular: reciprocal condition number = 1.11839e-20
```

So we need to do some “time shifting”: the problem is that some of the entries are too large

```
> t = t - 1790
> A = matrix(c(13, sum(t), sum(t^2),
+              sum(t), sum(t^2), sum(t^3),
+              sum(t^2), sum(t^3), sum(t^4)),
+              nrow=3, byrow=TRUE)
> b = c(sum(pop),
+        sum(pop * t),
+        sum(pop * t^2))
> sol = try(solve(A,b))
> print(sol)
[1] 5544964.000 -109242.513     6849.346
```

Thus

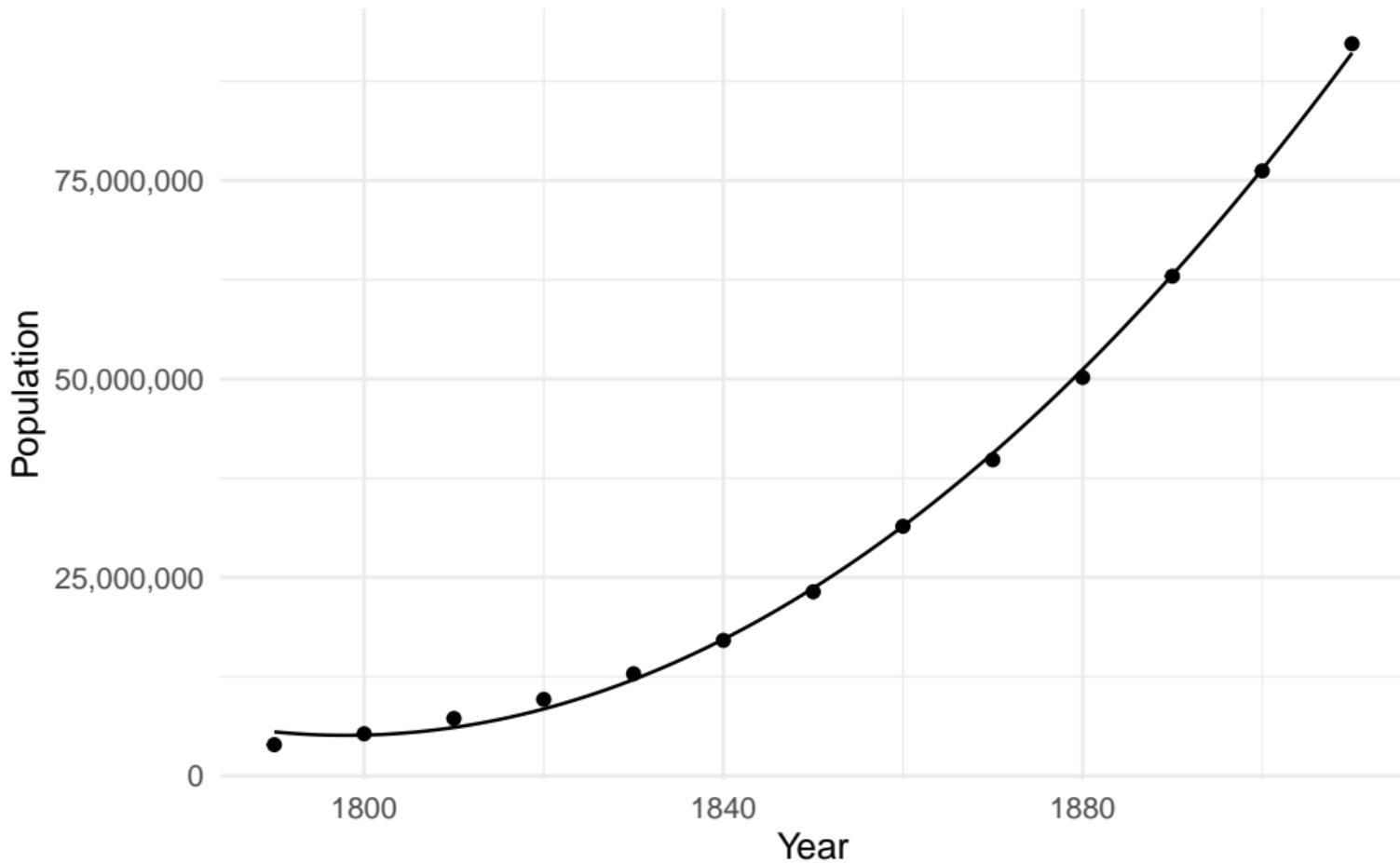
$$P(t) = 5544964 - 109243t + 6849t^2$$

(keeping in mind that time is here shifted and starts at 0)

So we define the function

```
> sol_plot = function(t, sol) {  
+   t = t - 1790  
+   return(sol[1] + sol[2]*t + sol[3]*t^2)  
+ }
```

US population from 1790 to 1910



Form the vector of errors, and compute sum of errors squared:

```
> t = USA_census_to_1910$Year  
> P = USA_census_to_1910$Population  
> E = sum((P - sol_plot(t, sol))^2)
```

Quite a large error (9,256,979,482,173), which is normal since we have used actual numbers, not thousands or millions of individuals, and we are taking the square of the error

Now for the big question...

How does our formula do for present times?

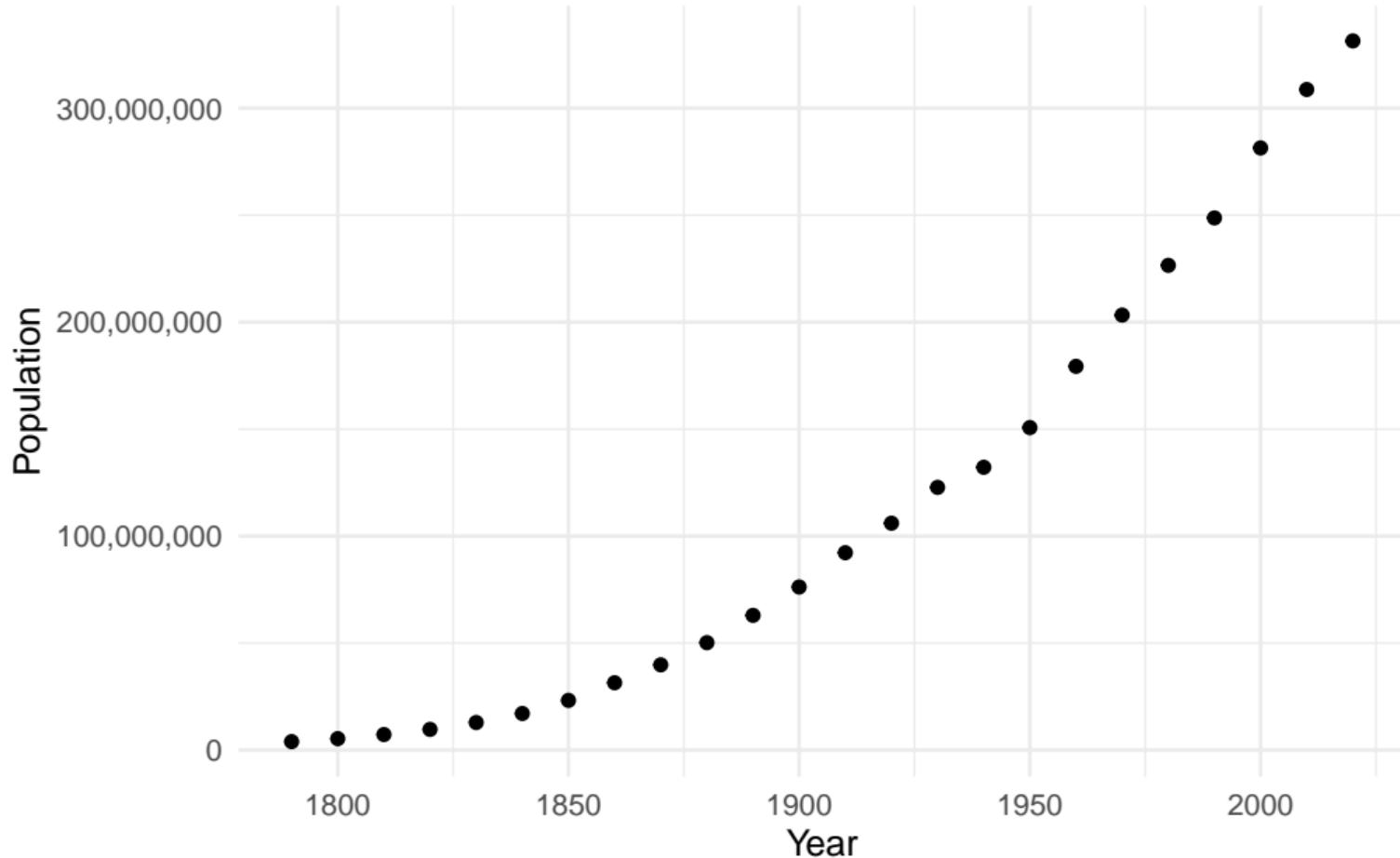
```
> format(sol_plot(2024, sol), big.mark = ",")  
[1] "355,024,999"
```

Actually, quite well: 355,024,999, compared to the 345,786,196 September 2024 estimate, overestimates the population by 9,238,803, a relative error of 2.67%

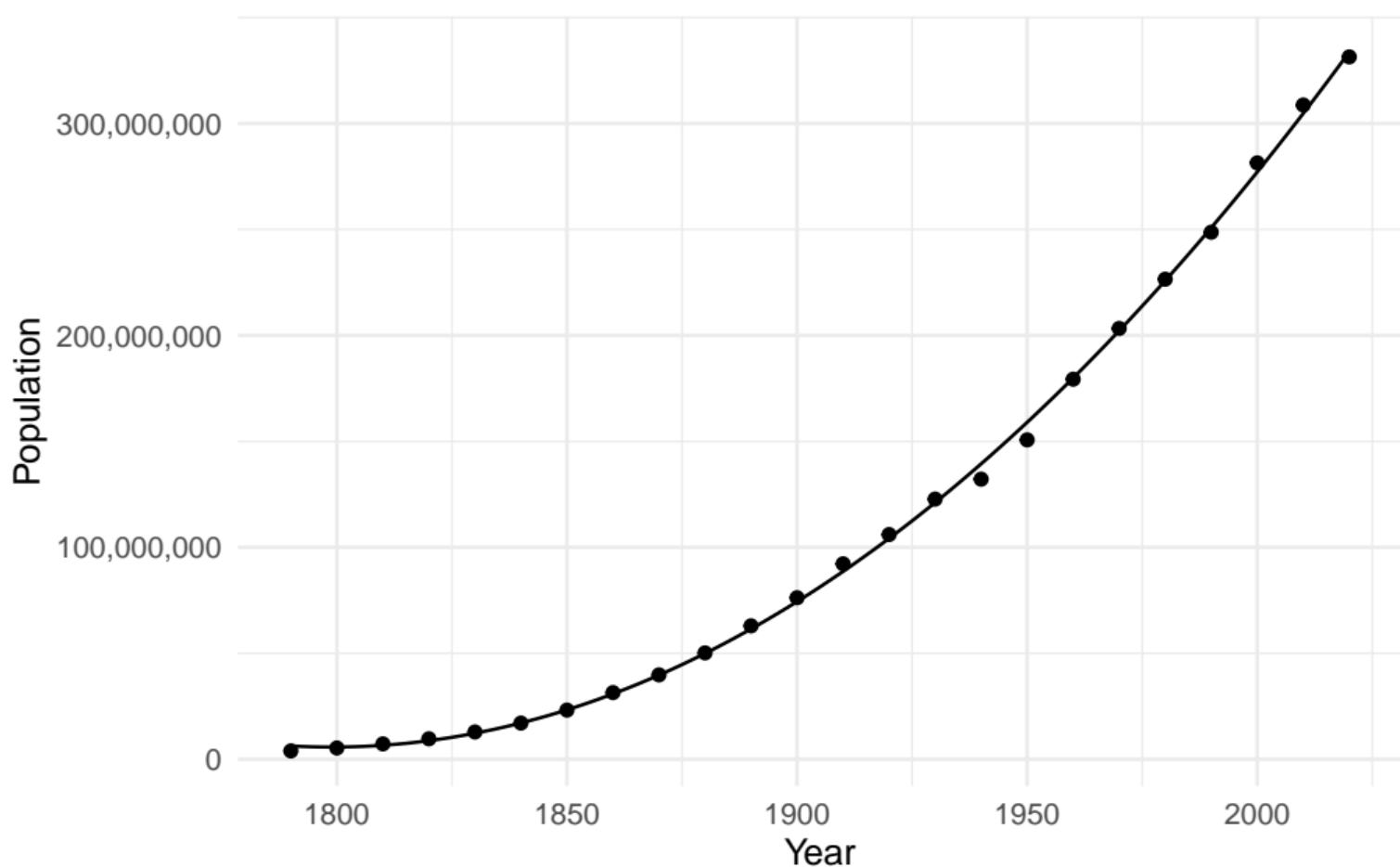
The US population from 1790 to 2020 (revised numbers)

Year	Population	Year	Population	Year	Population
1790	3,929,326	1890	62,947,714	1990	248,709,873
1800	5,308,483	1900	76,212,168	2000	281,421,906
1810	7,239,881	1910	92,228,496	2010	308,745,538
1820	9,638,453	1920	106,021,537	2020	331,449,281
1830	12,866,020	1930	122,775,046		
1840	17,069,458	1940	132,164,569		
1850	23,191,876	1950	150,697,361		
1860	31,443,321	1960	179,323,175		
1870	39,818,449	1970	203,302,031		
1880	50,189,209	1980	226,545,805		

US population from 1790 to 2020



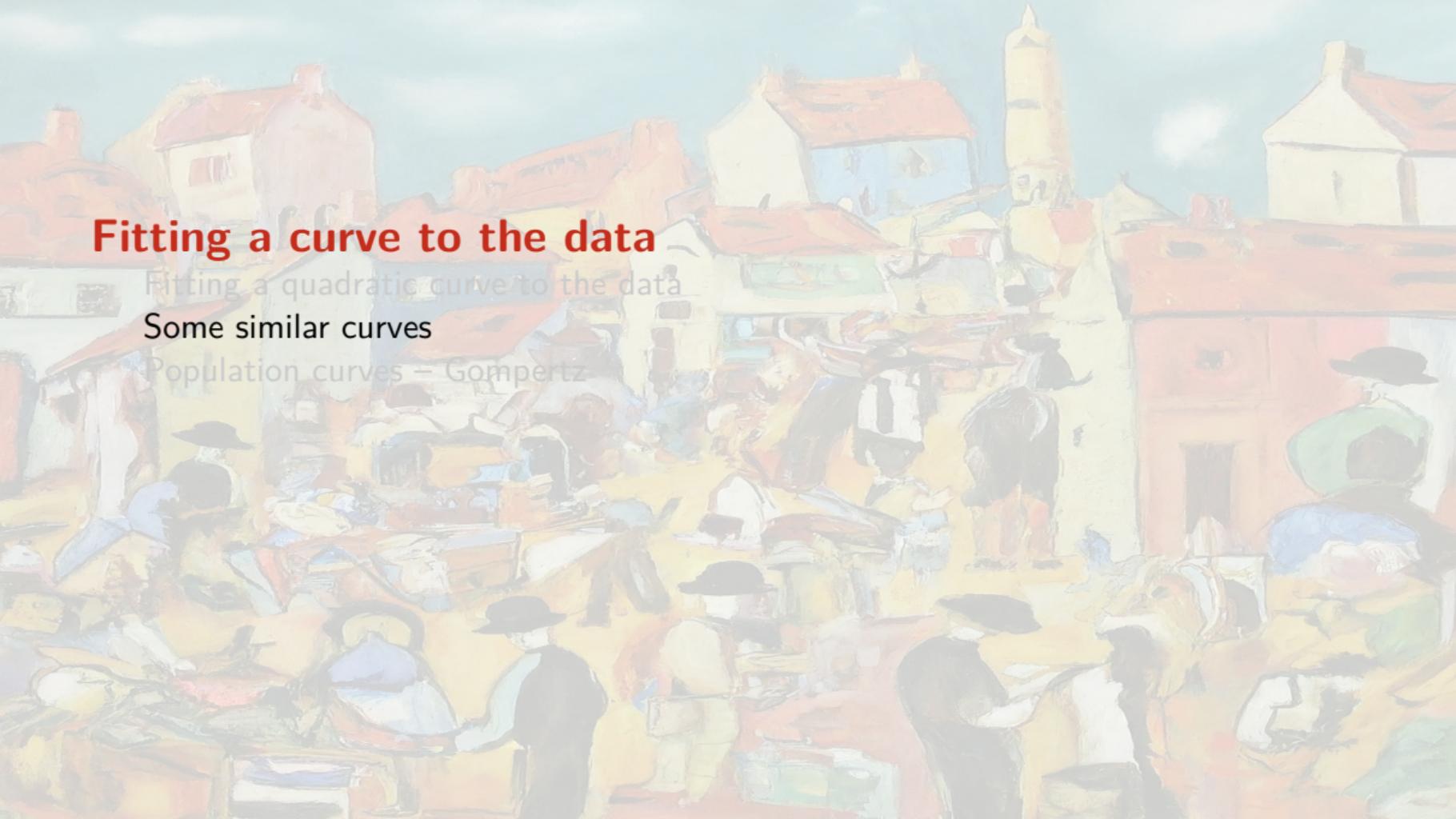
US population from 1790 to 2020 (with fit)



How does our formula do for present times?

```
> format(sol_plot(2024, sol_full), big.mark = ",")  
[1] "345,749,152"
```

Actually, quite well: 345,749,152, compared to the 345,786,196 September 2024 estimate, underestimates the population by -37,043.51, a relative error of -0.01%

A painting of a crowded outdoor market scene. In the foreground, many people wearing hats are gathered around tables with goods. In the background, there are several buildings with red-tiled roofs under a cloudy sky.

Fitting a curve to the data

Fitting a quadratic curve to the data

Some similar curves

Population curves – Gompertz

Other similar approaches

Pritchett, 1891:

$$P = a + bt + ct^2 + dt^3$$

(we have done this one, and found it to be quite good too)

Pearl, 1907:

$$P(t) = a + bt + ct^2 + d \ln t$$

Finds

$$P(t) = 9,064,900 - 6,281,430t + 842,377t^2 + 19,829,500 \ln t.$$

SHOWING (a) THE ACTUAL POPULATION¹ ON CENSUS DATES, (b) ESTIMATED POPULATION
 FROM PRITCHETT'S THIRD-ORDER PARABOLA, (c) ESTIMATED POPULATION FROM
 LOGARITHMIC PARABOLA, AND (d) (e) ROOT-MEAN SQUARE ERRORS
 OF BOTH METHODS

CENSUS YEAR	(a) OBSERVED POPULATION	(b). PRITCHETT ESTIMATE	(c) LOGARITHMIC PARABOLA ES- TIMATE	(d) ERROR OF (b)	(e) ERROR OF (c)
1790	3,929,000	4,012,000	3,693,000	+ 83,000	- 236,000
1800	5,308,000	5,267,000	5,865,000	- 41,000	+ 557,000
1810	7,240,000	7,059,000	7,293,000	- 181,000	+ 53,000
1820	9,638,000	9,571,000	9,404,000	- 67,000	- 234,000
1830	12,866,000	12,985,000	12,577,000	+ 119,000	- 289,000
1840	17,069,000	17,484,000	17,132,000	+ 415,000	+ 63,000
1850	23,192,000	23,250,000	23,129,000	+ 58,000	- 63,000
1860	31,443,000	30,465,000	30,633,000	- 978,000	- 810,000
1870	38,558,000	39,313,000	39,687,000	+ 755,000	+ 1,129,000
1880	50,156,000	49,975,000	50,318,000	- 181,000	+ 162,000
1890	62,948,000	62,634,000	62,547,000	- 314,000	- 401,000
1900	75,995,000	77,472,000	76,389,000	+ 1,477,000	+ 394,000
1910	91,972,000	94,673,000	91,647,000	+ 2,701,000	- 325,000
				935,000 ²	472,000 ²
1920		114,416,000	108,214,000		

¹ To the nearest thousand.

² Root-mean square error.

The logistic curve

Pearl and Reed try

$$P(t) = \frac{be^{at}}{1 + ce^{at}}$$

or

$$P(t) = \frac{b}{e^{-at} + c}$$

What is wrong with the logistic equation here?

- ▶ The carrying capacity is constant
- ▶ The model does not take immigration into account (for the US, this is an important component)

PROCEEDINGS
OF THE
NATIONAL ACADEMY OF SCIENCES

Volume 18

January 15, 1932

Number 1

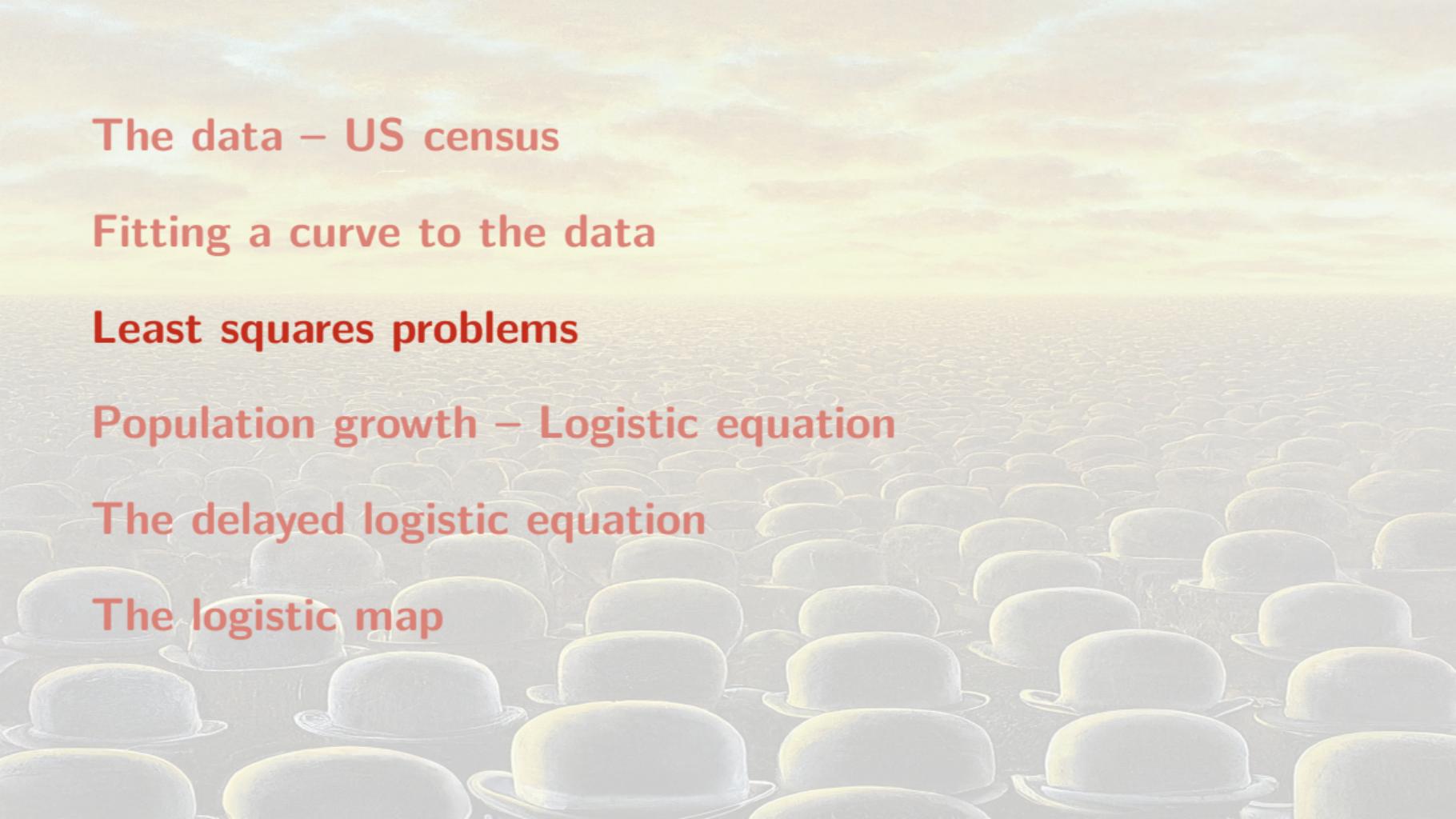
THE GOMPERTZ CURVE AS A GROWTH CURVE

By CHARLES P. WINSOR

DEPARTMENT OF BIOLOGY, SCHOOL OF HYGIENE AND PUBLIC HEALTH, JOHNS HOPKINS
UNIVERSITY

Communicated December 2, 1931

PROPERTY	GOMPERTZ	LOGISTIC
Equation	$y = ke^{-e^{a-bx}}$	$y = \frac{k}{1 + e^{a-bx}}$
Number of constants	3	3
Asymptotes	$\begin{cases} y = 0 \\ y = k \end{cases}$	$\begin{cases} y = 0 \\ y = k \end{cases}$
Inflection	$\begin{cases} x = \frac{a}{b} \\ y = \frac{k}{e} \end{cases}$	$\begin{cases} x = \frac{a}{b} \\ y = \frac{k}{2} \end{cases}$
Straight line form of equation	$\log \log \frac{k}{y} = a - bx$	$\log \frac{k-y}{y} = a - bx$
Symmetry	Assymetrical	Symmetrical about inflection
Growth rate	$\frac{dy}{dx} = bye^{a-bx} = by \log \frac{k}{y}$	$\frac{dy}{dx} = \frac{b}{k} y(k-y)$
Maximum growth rate	$\frac{bk}{e}$	$\frac{bk}{4}$
Relative growth rate as function of time	$\frac{1}{y} \frac{dy}{dx} = be^{a-bx}$	$\frac{1}{y} \frac{dy}{dx} = \frac{b}{1+e^{-a+bx}}$
Relative growth rate as function of size	$\frac{1}{y} \frac{dy}{dx} = b(\log k - \log y)$	$\frac{1}{y} \frac{dy}{dx} = \frac{b}{k}(k-y)$

The background of the slide features a wide, flat landscape filled with numerous small, white bowler hats, suggesting a large crowd of people. The sky above is a pale yellow with scattered, wispy clouds.

The data – US census

Fitting a curve to the data

Least squares problems

Population growth – Logistic equation

The delayed logistic equation

The logistic map

A.k.a. if the Math Dept was less #\\$%&, you'd know this

The following are a brief extract from MATH 2740 slides...

The least squares problem (simplest version)

Definition 1

Given a collection of points $(x_1, y_1), \dots, (x_n, y_n)$, find the coefficients a, b of the line $y = a + bx$ such that

$$\|\mathbf{e}\| = \sqrt{\varepsilon_1^2 + \dots + \varepsilon_n^2} = \sqrt{(y_1 - \tilde{y}_1)^2 + \dots + (y_n - \tilde{y}_n)^2}$$

is minimal, where $\tilde{y}_i = a + bx_i$ for $i = 1, \dots, n$

We just saw how to solve this by brute force using a genetic algorithm to minimise $\|\mathbf{e}\|$, let us now see how to solve this problem “properly”

For a data point $i = 1, \dots, n$

$$\varepsilon_i = y_i - \tilde{y}_i = y_i - (a + bx_i)$$

So if we write this for all data points,

$$\varepsilon_1 = y_1 - (a + bx_1)$$

⋮

$$\varepsilon_n = y_n - (a + bx_n)$$

In matrix form

$$\mathbf{e} = \mathbf{b} - A\mathbf{x}$$

with

$$\mathbf{e} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}, A = \begin{pmatrix} 1 & x_1 \\ \vdots & \vdots \\ 1 & x_n \end{pmatrix}, \mathbf{x} = \begin{pmatrix} a \\ b \end{pmatrix} \text{ and } \mathbf{b} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

The least squares problem (reformulated)

Definition 2 (Least squares solutions)

Consider a collection of points $(x_1, y_1), \dots, (x_n, y_n)$, a matrix $A \in \mathcal{M}_{mn}$, $\mathbf{b} \in \mathbb{R}^m$. A **least squares solution** of $A\mathbf{x} = \mathbf{b}$ is a vector $\tilde{\mathbf{x}} \in \mathbb{R}^n$ s.t.

$$\forall \mathbf{x} \in \mathbb{R}^n, \quad \|\mathbf{b} - A\tilde{\mathbf{x}}\| \leq \|\mathbf{b} - A\mathbf{x}\|$$

Needed to solve the problem

Definition 3 (Best approximation)

Let V be a vector space, $W \subset V$ and $\mathbf{v} \in V$. The **best approximation** to \mathbf{v} in W is $\tilde{\mathbf{v}} \in W$ s.t.

$$\forall \mathbf{w} \in W, \mathbf{w} \neq \tilde{\mathbf{v}}, \quad \|\mathbf{v} - \tilde{\mathbf{v}}\| < \|\mathbf{v} - \mathbf{w}\|$$

Theorem 4 (Best approximation theorem)

Let V be a vector space with an inner product, $W \subset V$ and $\mathbf{v} \in V$. Then $\text{proj}_W(\mathbf{v})$ is the best approximation to \mathbf{v} in W

Let us find the least squares solution

$\forall \mathbf{x} \in \mathbb{R}^n$, $A\mathbf{x}$ is a vector in the **column space** of A (the space spanned by the vectors making up the columns of A)

Since $\mathbf{x} \in \mathbb{R}^n$, $A\mathbf{x} \in \text{col}(A)$

\implies least squares solution of $A\mathbf{x} = \mathbf{b}$ is a vector $\tilde{\mathbf{y}} \in \text{col}(A)$ s.t.

$$\forall \mathbf{y} \in \text{col}(A), \quad \|\mathbf{b} - \tilde{\mathbf{y}}\| \leq \|\mathbf{b} - \mathbf{y}\|$$

This looks very much like Best approximation and Best approximation theorem

Putting things together

We just stated: The least squares solution of $Ax = \mathbf{b}$ is a vector $\tilde{\mathbf{y}} \in \text{col}(A)$ s.t.

$$\forall \mathbf{y} \in \text{col}(A), \quad \|\mathbf{b} - \tilde{\mathbf{y}}\| \leq \|\mathbf{b} - \mathbf{y}\|$$

We know (reformulating a tad):

Theorem 5 (Best approximation theorem)

Let V be a vector space with an inner product, $W \subset V$ and $\mathbf{v} \in V$. Then $\text{proj}_W(\mathbf{v}) \in W$ is the best approximation to \mathbf{v} in W , i.e.,

$$\forall \mathbf{w} \in W, \mathbf{w} \neq \text{proj}_W(\mathbf{v}), \quad \|\mathbf{v} - \text{proj}_W(\mathbf{v})\| < \|\mathbf{v} - \mathbf{w}\|$$

$$\implies W = \text{col}(A), \mathbf{v} = \mathbf{b} \text{ and } \tilde{\mathbf{y}} = \text{proj}_{\text{col}(A)}(\mathbf{b})$$

So if $\tilde{\mathbf{x}}$ is a least squares solution of $A\mathbf{x} = \mathbf{b}$, then

$$\tilde{\mathbf{y}} = A\tilde{\mathbf{x}} = \text{proj}_{\text{col}(A)}(\mathbf{b})$$

We have

$$\mathbf{b} - A\tilde{\mathbf{x}} = \mathbf{b} - \text{proj}_{\text{col}(A)}(\mathbf{b}) = \text{perp}_{\text{col}(A)}(\mathbf{b})$$

and it is easy to show that

$$\text{perp}_{\text{col}(A)}(\mathbf{b}) \perp \text{col}(A)$$

So for all columns \mathbf{a}_i of A

$$\mathbf{a}_i \cdot (\mathbf{b} - A\tilde{\mathbf{x}}) = 0$$

which we can also write as $\mathbf{a}_i^T(\mathbf{b} - A\tilde{\mathbf{x}}) = 0$

For all columns \mathbf{a}_i of A ,

$$\mathbf{a}_i^T(\mathbf{b} - A\tilde{\mathbf{x}}) = 0$$

This is equivalent to saying that

$$A^T(\mathbf{b} - A\tilde{\mathbf{x}}) = \mathbf{0}$$

We have

$$\begin{aligned} A^T(\mathbf{b} - A\tilde{\mathbf{x}}) = \mathbf{0} &\iff A^T\mathbf{b} - A^TA\tilde{\mathbf{x}} = \mathbf{0} \\ &\iff A^T\mathbf{b} = A^TA\tilde{\mathbf{x}} \\ &\iff A^TA\tilde{\mathbf{x}} = A^T\mathbf{b} \end{aligned}$$

The latter system constitutes the **normal equations** for $\tilde{\mathbf{x}}$

Least squares theorem

Theorem 6 (Least squares theorem)

$A \in \mathcal{M}_{mn}$, $\mathbf{b} \in \mathbb{R}^m$. Then

1. $A\mathbf{x} = \mathbf{b}$ always has at least one least squares solution $\tilde{\mathbf{x}}$
2. $\tilde{\mathbf{x}}$ least squares solution to $A\mathbf{x} = \mathbf{b} \iff \tilde{\mathbf{x}}$ is a solution to the normal equations
 $A^T A \tilde{\mathbf{x}} = A^T \mathbf{b}$
3. A has linearly independent columns $\iff A^T A$ invertible.
In this case, the least squares solution is unique and

$$\tilde{\mathbf{x}} = (A^T A)^{-1} A^T \mathbf{b}$$

We have seen 1 and 2, we will not show 3 (it is not hard)

Suppose we want to fit something a bit more complicated..

For instance, instead of the affine function

$$y = a + bx$$

suppose we want to do the quadratic

$$y = a_0 + a_1x + a_2x^2$$

or even

$$y = k_0 e^{k_1 x}$$

How do we proceed?

Fitting the quadratic

We have the data points $(x_1, y_1), (x_2, y_2), \dots, (x_n, y_n)$ and want to fit

$$y = a_0 + a_1 x + a_2 x^2$$

At (x_1, y_1) ,

$$\tilde{y}_1 = a_0 + a_1 x_1 + a_2 x_1^2$$

⋮

At (x_n, y_n) ,

$$\tilde{y}_n = a_0 + a_1 x_n + a_2 x_n^2$$

In terms of the error

$$\varepsilon_1 = y_1 - \tilde{y}_1 = y_1 - (a_0 + a_1 x_1 + a_2 x_1^2)$$

⋮

$$\varepsilon_n = y_n - \tilde{y}_n = y_n - (a_0 + a_1 x_n + a_2 x_n^2)$$

i.e.,

$$\mathbf{e} = \mathbf{b} - A\mathbf{x}$$

where

$$\mathbf{e} = \begin{pmatrix} \varepsilon_1 \\ \vdots \\ \varepsilon_n \end{pmatrix}, A = \begin{pmatrix} 1 & x_1 & x_1^2 \\ \vdots & \vdots & \vdots \\ 1 & x_n & x_n^2 \end{pmatrix}, \mathbf{x} = \begin{pmatrix} a_0 \\ a_1 \\ a_2 \end{pmatrix} \text{ and } \mathbf{b} = \begin{pmatrix} y_1 \\ \vdots \\ y_n \end{pmatrix}$$

Theorem 6 applies, with here $A \in \mathcal{M}_{n3}$ and $\mathbf{b} \in \mathbb{R}^n$

Fitting the exponential

Things are a bit more complicated here

If we proceed as before, we get the system

$$y_1 = k_0 e^{k_1 x_1}$$

⋮

$$y_n = k_0 e^{k_1 x_n}$$

$e^{k_1 x_i}$ is a nonlinear term, it cannot be put in a matrix

However: take the \ln of both sides of the equation

$$\ln(y_i) = \ln(k_0 e^{k_1 x_i}) = \ln(k_0) + \ln(e^{k_1 x_i}) = \ln(k_0) + k_1 x_i$$

If $y_i, k_0 > 0$, then their \ln are defined and we're in business..

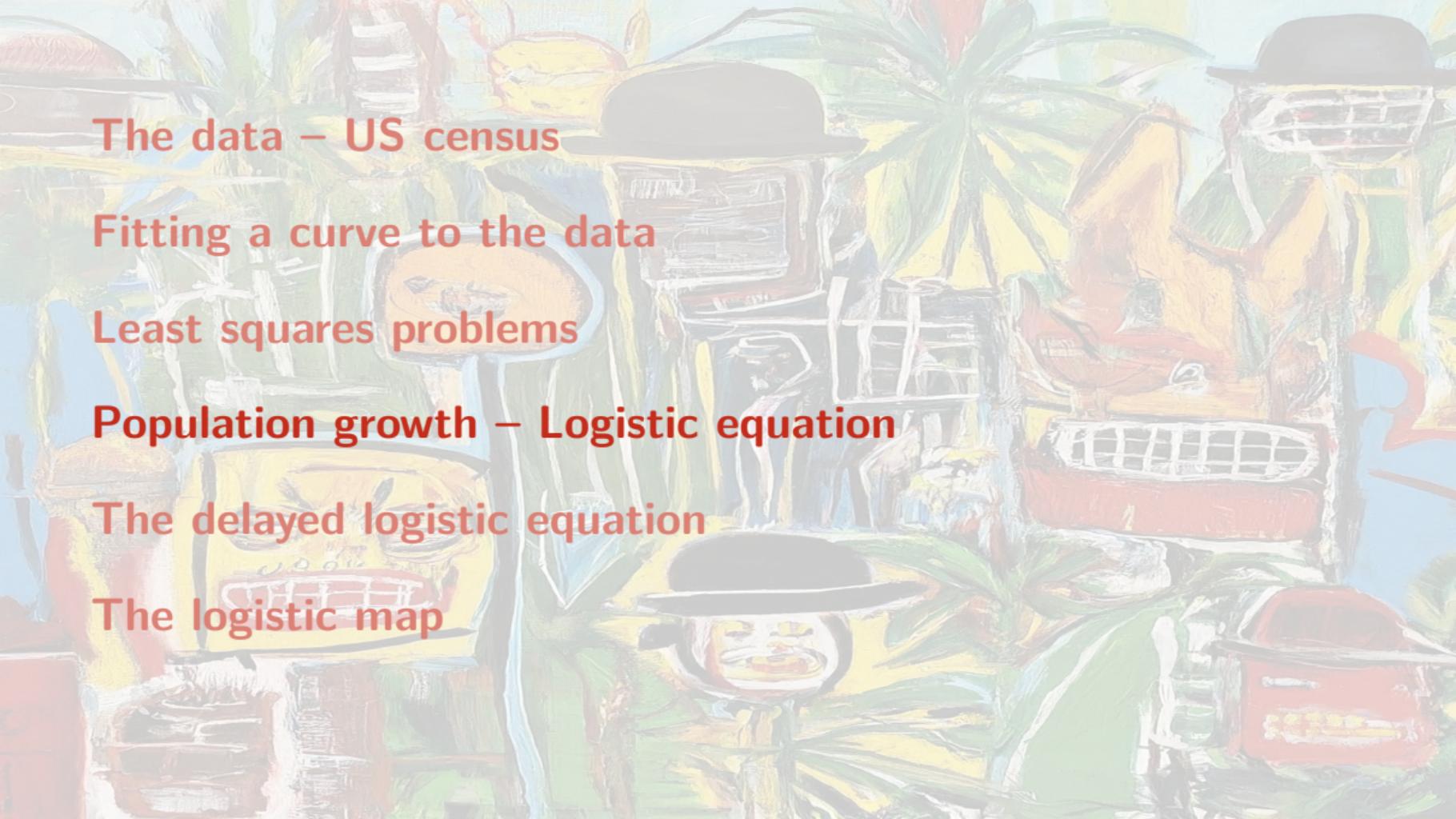
$$\ln(y_i) = \ln(k_0) + k_1 x_i$$

So the system is

$$\mathbf{y} = A\mathbf{x} + \mathbf{b}$$

with

$$A = \begin{pmatrix} x_1 \\ \vdots \\ x_n \end{pmatrix}, \mathbf{x} = (k_1), \mathbf{b} = (\ln(k_0)) \text{ and } \mathbf{y} = \begin{pmatrix} \ln(y_1) \\ \vdots \\ \ln(y_n) \end{pmatrix}$$



The data – US census

Fitting a curve to the data

Least squares problems

Population growth – Logistic equation

The delayed logistic equation

The logistic map

Population growth – Logistic equation

Formulating the logistic equation

Qualitative analysis of the logistic equation

The logistic equation

The logistic curve is the solution to the ordinary differential equation

$$N' = rN \left(1 - \frac{N}{K}\right)$$

which is called the **logistic equation**

r is the **intrinsic growth rate**, K is the **carrying capacity**

This equation was introduced by Pierre-François Verhulst (1804-1849) in 1844

Deriving the logistic equation

The idea is to represent a population with the following components:

- ▶ birth, at the **per capita** rate b
- ▶ death, at the **per capita** rate d
- ▶ competition of individuals with other individuals reduces their ability to survive, resulting in death

This gives

$$N' = bN - dN - \text{competition}$$

Accounting for competition

Competition describes the mortality that occurs when two individuals meet

- ▶ In chemistry, if there is a concentration X of one product and Y of another product, then XY , called **mass action**, describes the number of interactions of molecules of the two products
- ▶ Here, we assume that X and Y are of the same type (individuals). So there are N^2 contacts
- ▶ These N^2 contacts lead to death of one of the individuals at the rate c

Therefore, the **logistic** equation is

$$N' = bN - dN - cN^2$$

Reinterpreting the logistic equation

The equation

$$N' = bN - dN - cN^2$$

is rewritten as

$$N' = (b - d)N - cN^2$$

- ▶ $b - d$ represents the rate at which the population increases (or decreases) in the absence of competition. It is called the **intrinsic growth rate** of the population
- ▶ c is the rate of **intraspecific** competition. The prefix **intra** refers to the fact that the competition is occurring between members of the same species, that is, within the species

[We will see later examples of **interspecific** competition, that is, between different species]

Another (...) interpretation of the logistic equation

We have

$$N' = (b - d)N - cN^2$$

Factor out an N :

$$N' = ((b - d) - cN)N$$

This gives us another interpretation of the logistic equation. Writing

$$\frac{N'}{N} = (b - d) - cN$$

we have N'/N , the **per capita growth rate** of N , given by a constant, $b - d$, minus a **density dependent inhibition** factor, cN

Equivalent equations

$$\begin{aligned}N' &= (b - d)N - cN^2 \\&= ((b - d) - cN)N \\&= \left(r - \frac{r}{r}cN\right)N \quad \text{with } r = b - d \\&= rN\left(1 - \frac{c}{r}N\right) \\&= rN\left(1 - \frac{N}{K}\right)\end{aligned}$$

with

$$\frac{c}{r} = \frac{1}{K}$$

i.e., $K = r/c$

3 ways to tackle this equation

1. The equation is separable [explicit method]
2. The equation is a Bernoulli equation [explicit method]
3. Use qualitative analysis

Population growth – Logistic equation

Formulating the logistic equation

Qualitative analysis of the logistic equation

Studying the logistic equation qualitatively

We study

$$N' = rN \left(1 - \frac{N}{K}\right) \quad (\text{ODE1})$$

For this, write

$$f(N) = rN \left(1 - \frac{N}{K}\right)$$

Consider the initial value problem (IVP)

$$N' = f(N), \quad N(0) = N_0 > 0 \quad (\text{IVP1})$$

- ▶ f is C^1 (differentiable with continuous derivative) so solutions to (IVP1) exist and are unique

Equilibria of (ODE1) are points such that $f(N) = 0$ (so that $N' = f(N) = 0$, meaning N does not vary). So we solve $f(N) = 0$ for N . We find two points:

- ▶ $N = 0$
- ▶ $N = K$

By uniqueness of solutions to (IVP1), solutions cannot cross the lines $N(t) = 0$ and $N(t) = K$

Several cases

- ▶ $N = 0$ for some t , then $N(t) = 0$ for all $t \geq 0$, by uniqueness of solutions
- ▶ $N \in (0, K)$, then $rN > 0$ and $N/K < 1$ so $1 - N/K > 0$, which implies that $f(N) > 0$. As a consequence, $N(t)$ increases if $N \in (0, K)$
- ▶ $N = K$, then $rN > 0$ but $N/K = 1$ so $1 - N/K = 0$, which implies that $f(N) = 0$. As a consequence, $N(t) = K$ for all $t \geq 0$, by uniqueness of solutions
- ▶ $N > K$, the $rN > 0$ and $N/K > 1$, implying that $1 - N/K < 0$ and in turn, $f(N) < 0$. As a consequence, $N(t)$ decreases if $N \in (K, +\infty)$

Therefore,

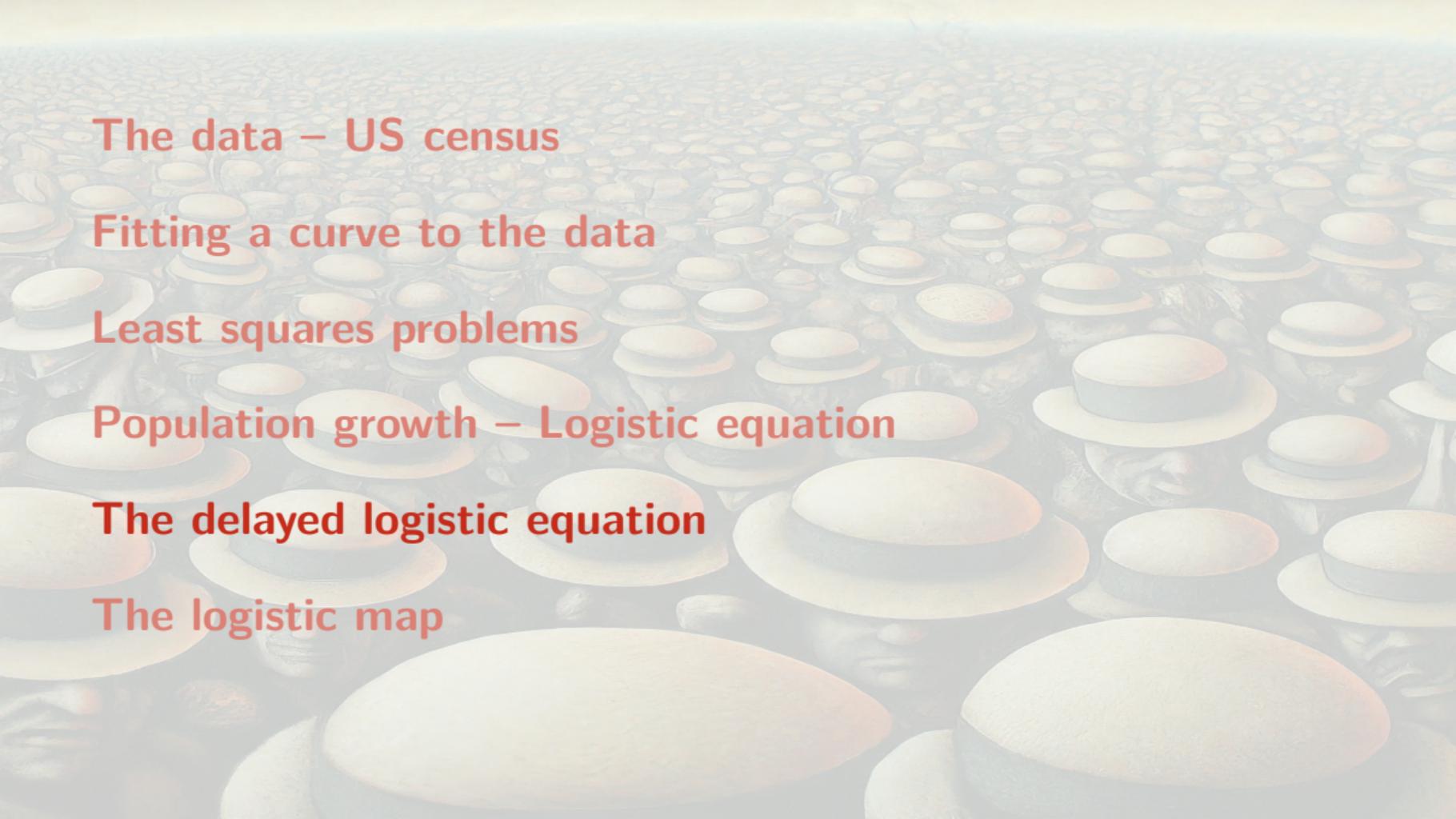
Theorem 7

Suppose that $N_0 > 0$. Then the solution $N(t)$ of (IVP1) is such that

$$\lim_{t \rightarrow \infty} N(t) = K$$

so that K is the number of individuals that the environment can support, the **carrying capacity** of the environment

If $N_0 = 0$, then $N(t) = 0$ for all $t \geq 0$



The data – US census

Fitting a curve to the data

Least squares problems

Population growth – Logistic equation

The delayed logistic equation

The logistic map

The delayed logistic equation

Consider the equation as

$$\frac{N'}{N} = (b - d) - cN$$

that is, the per capita rate of growth of the population depends on the net growth rate $b - d$, and some density dependent inhibition cN (resulting of competition)

Suppose that instead of instantaneous inhibition, there is some delay τ between the time the inhibiting event takes place and the moment when it affects the growth rate

For example, two individuals fight for food, and one later dies of the injuries sustained during this fight

The delayed logistic equation

In the case of a time τ between inhibiting event and inhibition, the equation would be written as

$$\frac{N'}{N} = (b - d) - cN(t - \tau)$$

Using the change of variables introduced earlier, this is written

$$N'(t) = rN(t) \left(1 - \frac{N(t - \tau)}{K}\right) \quad (\text{DDE1})$$

Such an equation is called a **delay** differential equation. It is much more complicated to study than (ODE1). In fact, some things remain unknown about (DDE1)

Delayed initial value problem

The IVP takes the form

$$\begin{aligned}N'(t) &= rN(t) \left(1 - \frac{N(t-\tau)}{K}\right) \\ N(t) &= \phi(t) \text{ for } t \in [-\tau, 0]\end{aligned}\tag{IVP2}$$

where $\phi(t)$ is some continuous function

Hence, initial conditions (called initial data in this case) must be specified on an interval, instead of being specified at a point, to guarantee existence and uniqueness of solutions

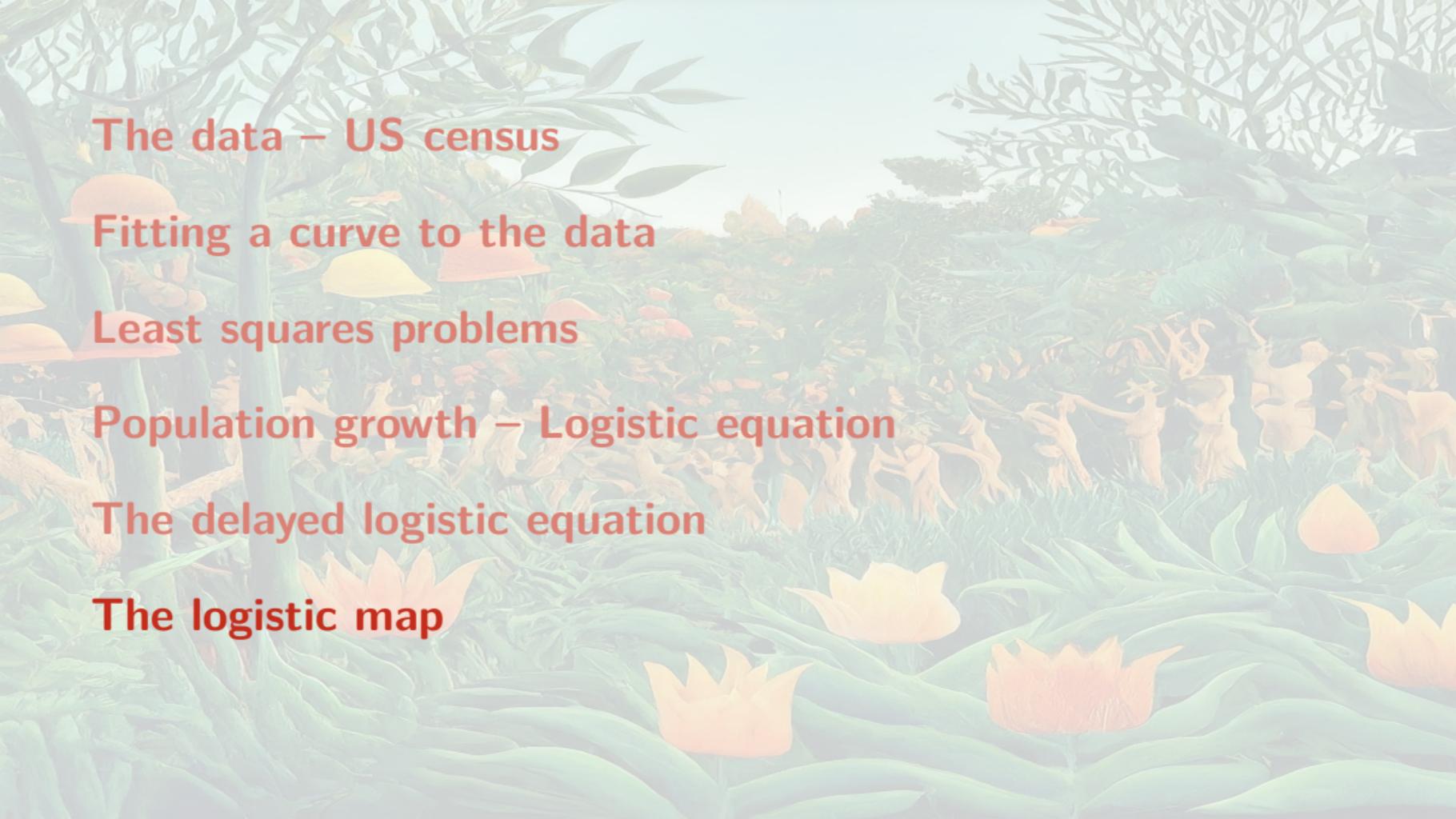
We will not learn how to study this type of equation (this is graduate level mathematics). I will give a few results

To find equilibria, remark that delay should not play a role, since N should be constant. Thus, equilibria are found by considering the equation with no delay, which is (ODE1)

Theorem 8

Suppose that $r\tau < \pi/2$. Then solutions of (IVP2) with positive initial data $\phi(t)$ starting close enough to K tend to K . If $r\tau < 37/24$, then all solutions of (IVP2) with positive initial data $\phi(t)$ tend to K . If $r\tau > \pi/2$, then K is an unstable equilibrium and all solutions of (IVP2) with positive initial data $\phi(t)$ on $[-\tau, 0]$ are oscillatory

There is a gray zone between $37/24$ ($\simeq 1.5417$) and $\pi/2$ ($\simeq 1.5708$). The global aspect was proved for $r\tau < 37/24$ in 1945 by Wright. Although there is very strong numerical evidence that this is in fact true up to $\pi/2$, nobody has yet managed to prove it [Edit: now done!]

A dense field of orange tulips with green leaves, serving as the background for the slide.

The data – US census

Fitting a curve to the data

Least squares problems

Population growth – Logistic equation

The delayed logistic equation

The logistic map

Discrete-time systems

So far, we have seen continuous-time models, where $t \in \mathbb{R}_+$. Another way to model natural phenomena is by using a discrete-time formalism, that is, to consider equations of the form

$$x_{t+1} = f(x_t)$$

where $t \in \mathbb{N}$ or \mathbb{Z} , that is, t takes values in a discrete valued (countable) set

Time could for example be days, years, etc.

The logistic map

The logistic **map** is, for $t \geq 0$,

$$N_{t+1} = rN_t \left(1 - \frac{N_t}{K}\right) \quad (\text{DT1})$$

To transform this into an initial value problem, we need to provide an initial condition $N_0 \geq 0$ for $t = 0$

Some mathematical analysis

Suppose we have a system in the form

$$x_{t+1} = f(x_t)$$

with initial condition given for $t = 0$ by x_0 . Then,

$$x_1 = f(x_0)$$

$$x_2 = f(x_1) = f(f(x_0)) \stackrel{\Delta}{=} f^2(x_0)$$

⋮

$$x_k = f^k(x_0)$$

The $f^k = \underbrace{f \circ f \circ \cdots \circ f}_{k \text{ times}}$ are the **iterates** of f

Fixed points

Definition 9 (Fixed point)

Let f be a function. A point p such that $f(p) = p$ is called a **fixed point** of f

Theorem 10

Consider the closed interval $I = [a, b]$. If $f: I \rightarrow I$ is continuous, then f has a fixed point in I

Theorem 11

Let I be a closed interval and $f: I \rightarrow \mathbb{R}$ be a continuous function. If $f(I) \supset I$, then f has a fixed point in I .

Periodic points

Definition 12 (Periodic point)

Let f be a function. If there exists a point p and an integer n such that

$$f^n(p) = p, \quad \text{but} \quad f^k(p) \neq p \text{ for } k < n,$$

then p is a periodic point of f with (least) period n (or a n -periodic point of f).

Thus, p is a n -periodic point of f iff p is a 1-periodic point of f^n .

Stability of fixed points, of periodic points

Theorem 13

Let f be a continuously differentiable function (that is, differentiable with continuous derivative, or C^1), and p be a fixed point of f .

1. If $|f'(p)| < 1$, then there is an open interval $\mathcal{I} \ni p$ such that $\lim_{k \rightarrow \infty} f^k(x) = p$ for all $x \in \mathcal{I}$.
2. If $|f'(p)| > 1$, then there is an open interval $\mathcal{I} \ni p$ such that if $x \in \mathcal{I}$, $x \neq p$, then there exists k such that $f^k(x) \notin \mathcal{I}$.

Definition 14

Suppose that p is a n -periodic point of f , with $f \in C^1$.

- If $|(f^n)'(p)| < 1$, then p is an **attracting** periodic point of f .
- If $|(f^n)'(p)| > 1$, then p is an **repelling** periodic point of f .

Parametrized families of functions

Consider the equation (DT1), which for convenience we rewrite as

$$N_{t+1} = rN_t(1 - N_t), \quad (\text{DT2})$$

where r is a parameter in \mathbb{R}_+ , and N will typically be taken in $[0, 1]$. Let

$$f_r(x) = rx(1 - x).$$

The function f_r is called a **parametrized family** of functions.

Bifurcations

Definition 15 (Bifurcation)

Let f_μ be a parametrized family of functions. Then there is a **bifurcation** at $\mu = \mu_0$ (or μ_0 is a bifurcation point) if there exists $\varepsilon > 0$ such that, if $\mu_0 - \varepsilon < a < \mu_0$ and $\mu_0 < b < \mu_0 + \varepsilon$, then the dynamics of $f_a(x)$ are “different” from the dynamics of $f_b(x)$.

An example of “different” would be that f_a has a fixed point (that is, a 1-periodic point) and f_b has a 2-periodic point.

Back to the logistic map

Consider the simplified version (DT2),

$$N_{t+1} = rN_t(1 - N_t) \stackrel{\Delta}{=} f_r(N_t).$$

Are solutions well defined?

Suppose $N_0 \in [0, 1]$, do we stay in $[0, 1]$? f_r is continuous on $[0, 1]$, so it has a extrema on $[0, 1]$. We have

$$f'_r(x) = r - 2rx = r(1 - 2x),$$

which implies that f_r increases for $x < 1/2$ and decreases for $x > 1/2$, reaching a maximum at $x = 1/2$.

$f_r(0) = f_r(1) = 0$ are the minimum values, and $f_r(1/2) = r/4$ is the maximum. Thus, if we want $N_{t+1} \in [0, 1]$ for $N_t \in [0, 1]$, we need to consider $r \leq 4$.

- ▶ Note that if $N_0 = 0$, then $N_t = 0$ for all $t \geq 1$.
- ▶ Similarly, if $N_0 = 1$, then $N_1 = 0$, and thus $N_t = 0$ for all $t \geq 1$.
- ▶ This is true for all t : if there exists t_k such that $N_{t_k} = 1$, then $N_t = 0$ for all $t \geq t_k$.
- ▶ This last case might occur if $r = 4$, as we have seen.
- ▶ Also, if $r = 0$ then $N_t = 0$ for all t .

For these reasons, we generally consider

$$N \in (0, 1)$$

and

$$r \in (0, 4).$$

Fixed points: existence

Fixed points of (DT2) satisfy $N = rN(1 - N)$, giving:

- ▶ $N = 0$;
- ▶ $1 = r(1 - N)$, that is, $p \stackrel{\Delta}{=} \frac{r-1}{r}$.

Note that $\lim_{r \rightarrow 0^+} p = 1 - \lim_{r \rightarrow 0^+} 1/r = -\infty$, $\frac{\partial}{\partial r} p = 1/r^2 > 0$ (so p is an increasing function of r), $p = 0 \Leftrightarrow r = 1$ and $\lim_{r \rightarrow \infty} p = 1$. So we come to this first conclusion:

- ▶ 0 always is a fixed point of f_r .
- ▶ If $0 < r < 1$, then p takes negative values so is not relevant.
- ▶ If $1 < r < 4$, then p exists.

Stability of the fixed points

Stability of the fixed points is determined by the (absolute) value f'_r at these fixed points. We have

$$|f'_r(0)| = r,$$

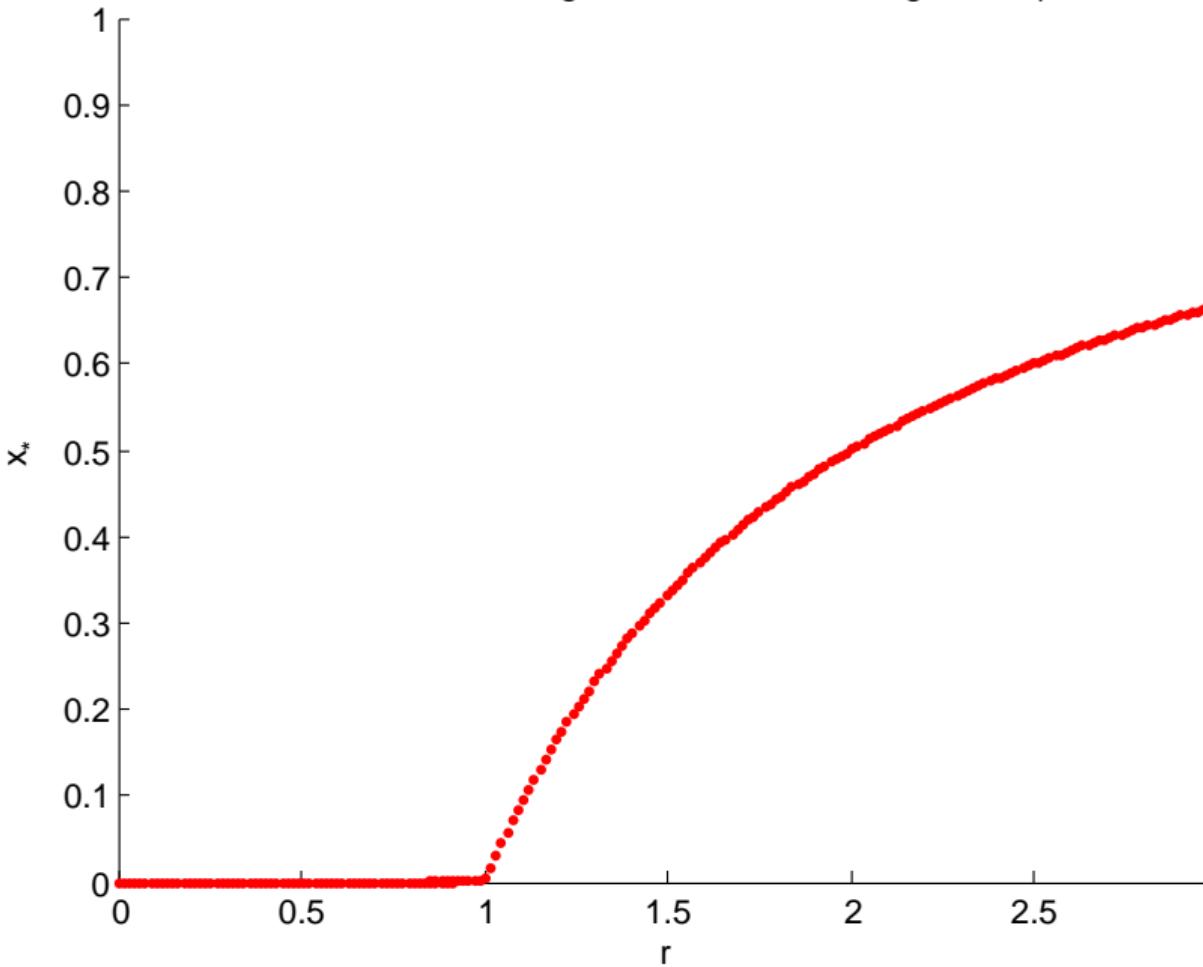
and

$$\begin{aligned}|f'_r(p)| &= \left| r - 2r\frac{r-1}{r} \right| \\&= |r - 2(r-1)| \\&= |2 - r|\end{aligned}$$

Therefore, we have

- ▶ if $0 < r < 1$, then the fixed point $N = p$ does not exist and $N = 0$ is attracting,
- ▶ if $1 < r < 3$, then $N = 0$ is repelling, and $N = p$ is attracting,
- ▶ if $r > 3$, then $N = 0$ and $N = p$ are repelling.

Bifurcation diagram for the discrete logistic map



Another bifurcation

Thus the points $r = 1$ and $r = 3$ are bifurcation points. To see what happens when $r > 3$, we need to look for period 2 points.

$$\begin{aligned}f_r^2(x) &= f_r(f_r(x)) \\&= rf_r(x)(1 - f_r(x)) \\&= r^2x(1 - x)(1 - rx(1 - x)).\end{aligned}\tag{1}$$

0 and p are points of period 2, since a fixed point x^* of f satisfies $f(x^*) = x^*$, and so, $f^2(x^*) = f(f(x^*)) = f(x^*) = x^*$.

This helps localizing the other periodic points. Writing the fixed point equation as

$$Q(x) \stackrel{\Delta}{=} f_r^2(x) - x = 0,$$

we see that, since 0 and p are fixed points of f_μ^2 , they are roots of $Q(x)$. Therefore, Q can be factorized as

$$Q(x) = x(x - p)(-r^3x^2 + Bx + C),$$

Substitute the value $(r - 1)/r$ for p in Q , develop Q and (1) and equate coefficients of like powers gives

$$Q(x) = x \left(x - \frac{r-1}{r} \right) (-r^3x^2 + r^2(r+1)x - r(r+1)). \quad (2)$$

We already know that $x = 0$ and $x = p$ are roots of (2). So we search for roots of

$$R(x) := -r^3x^2 + r^2(r+1)x - r(r+1).$$

Discriminant is

$$\begin{aligned}\Delta &= r^4(r+1)^2 - 4r^4(r+1) \\ &= r^4(r+1)(r+1-4) \\ &= r^4(r+1)(r-3).\end{aligned}$$

Therefore, R has distinct real roots if $r > 3$. Remark that for $r = 3$, the (double) root is $p = 2/3$. For $r > 3$ but very close to 3, it follows from the continuity of R that the roots are close to $2/3$.

Descartes' rule of signs

Theorem 16 (Descartes' rule of signs)

Let $p(x) = \sum_{i=0}^m a_i x^i$ be a polynomial with real coefficients such that $a_m \neq 0$. Define v to be the number of variations in sign of the sequence of coefficients a_m, \dots, a_0 . By 'variations in sign' we mean the number of values of n such that the sign of a_n differs from the sign of a_{n-1} , as n ranges from m down to 1. Then

- ▶ the number of positive real roots of $p(x)$ is $v - 2N$ for some integer N satisfying $0 \leq N \leq \frac{v}{2}$,
- ▶ the number of negative roots of $p(x)$ may be obtained by the same method by applying the rule of signs to $p(-x)$.

Example of use of Descartes' rule

Example 17

Let

$$p(x) = x^3 + 3x^2 - x - 3.$$

Coefficients have signs $++--$, i.e., 1 sign change. Thus $v = 1$. Since $0 \leq N \leq 1/2$, we must have $N = 0$. Thus $v - 2N = 1$ and there is exactly one positive real root of $p(x)$.

To find the negative roots, we examine $p(-x) = -x^3 + 3x^2 + x - 3$. Coefficients have signs $-++-$, i.e., 2 sign changes. Thus $v = 2$ and $0 \leq N \leq 2/2 = 1$. Thus, there are two possible solutions, $N = 0$ and $N = 1$, and two possible values of $v - 2N$. Therefore, there are either two or no negative real roots. Furthermore, note that $p(-1) = (-1)^3 + 3 \cdot (-1)^2 - (-1) - 3 = 0$, hence there is at least one negative root. Therefore there must be exactly two.

Back to the logistic map and the polynomial R ..

We use Descartes' rule of signs.

- ▶ R has signed coefficients $- + -$, so 2 sign changes implying 0 or 2 positive real roots.
- ▶ $R(-x)$ has signed coefficients $---$, so no negative real roots.
- ▶ Since $\Delta > 0$, the roots are real, and thus it follows that both roots are positive.

To show that the roots are also smaller than 1, consider the change of variables $z = x - 1$. The polynomial R is transformed into

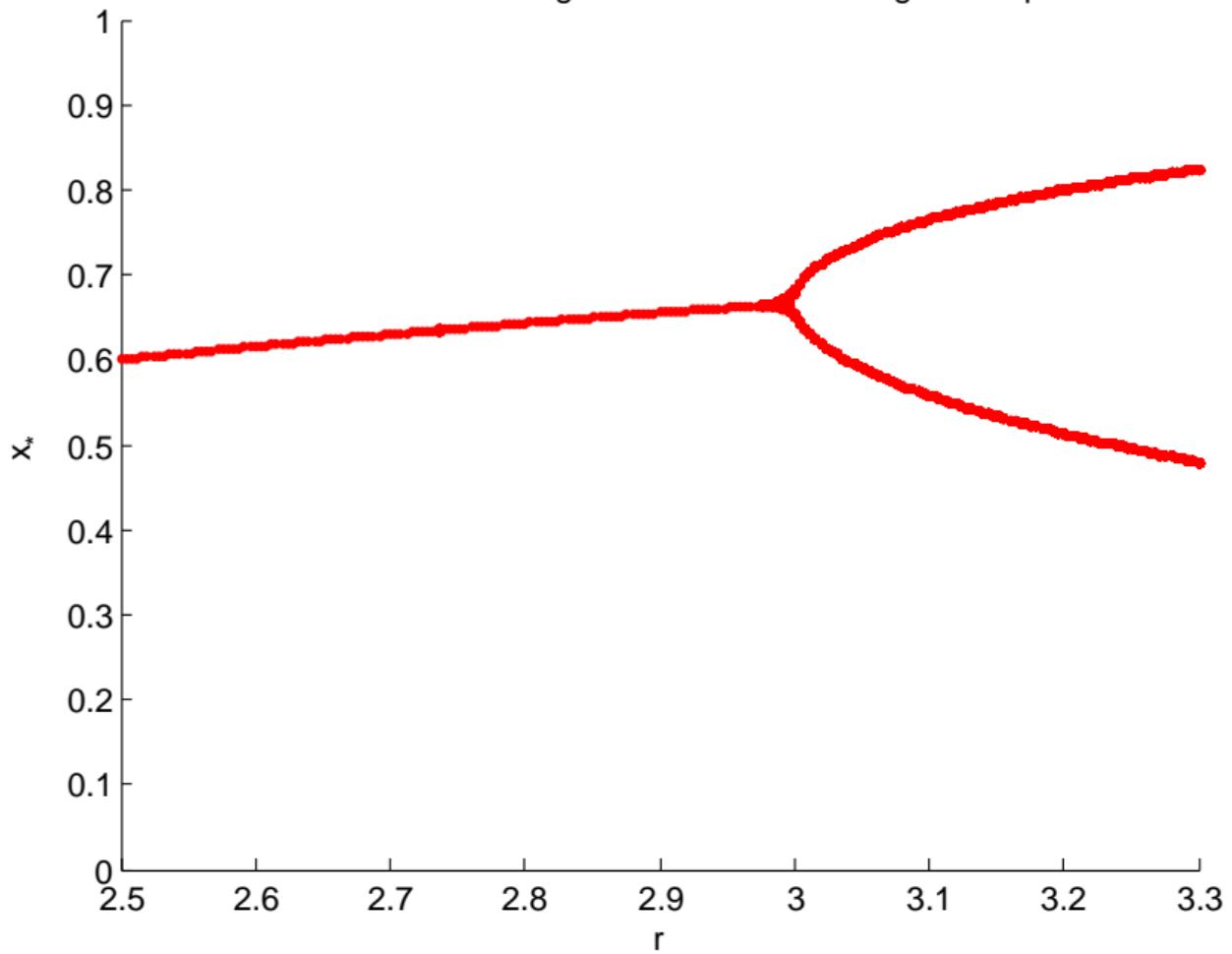
$$\begin{aligned}R_2(z) &= -r^3(z+1)^2 + r^2(r+1)(z+1) - r(r+1) \\&= -r^3z^2 + r^2(1-r)z - r.\end{aligned}$$

For $r > 1$, the signed coefficients are $---$, so R_2 has no root $z > 0$, implying in turn that R has no root $x > 1$.

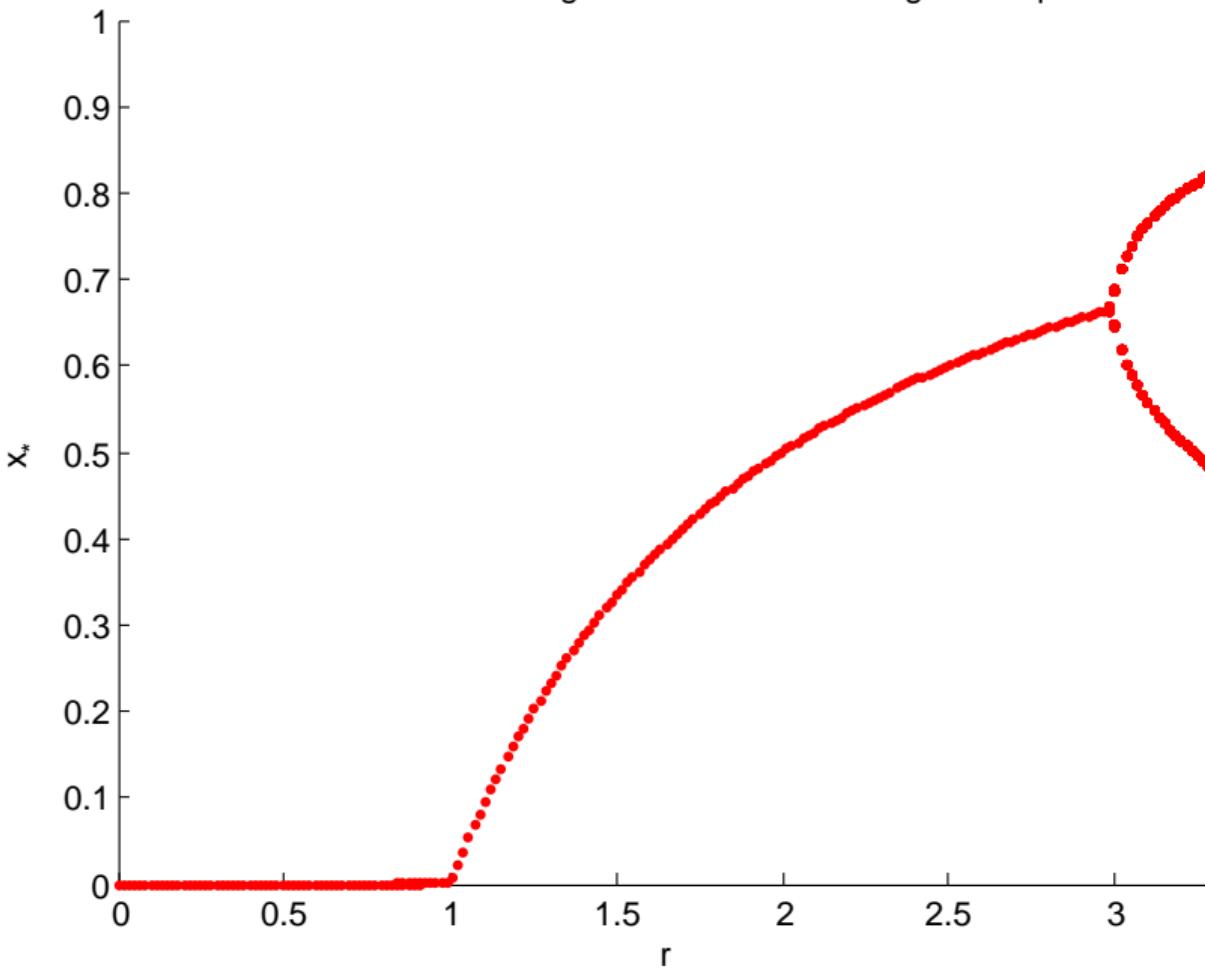
Summing up

- ▶ If $0 < r < 1$, then $N = 0$ is attracting, p does not exist and there are no period 2 points.
- ▶ At $r = 1$, there is a bifurcation (called a **transcritical** bifurcation).
- ▶ If $1 < r < 3$, then $N = 0$ is repelling, $N = p$ is attracting, and there are no period 2 points.
- ▶ At $r = 3$, there is another bifurcation (called a **period-doubling** bifurcation).
- ▶ For $r > 3$, both $N = 0$ and $N = p$ are repelling, and there is a period 2 point.

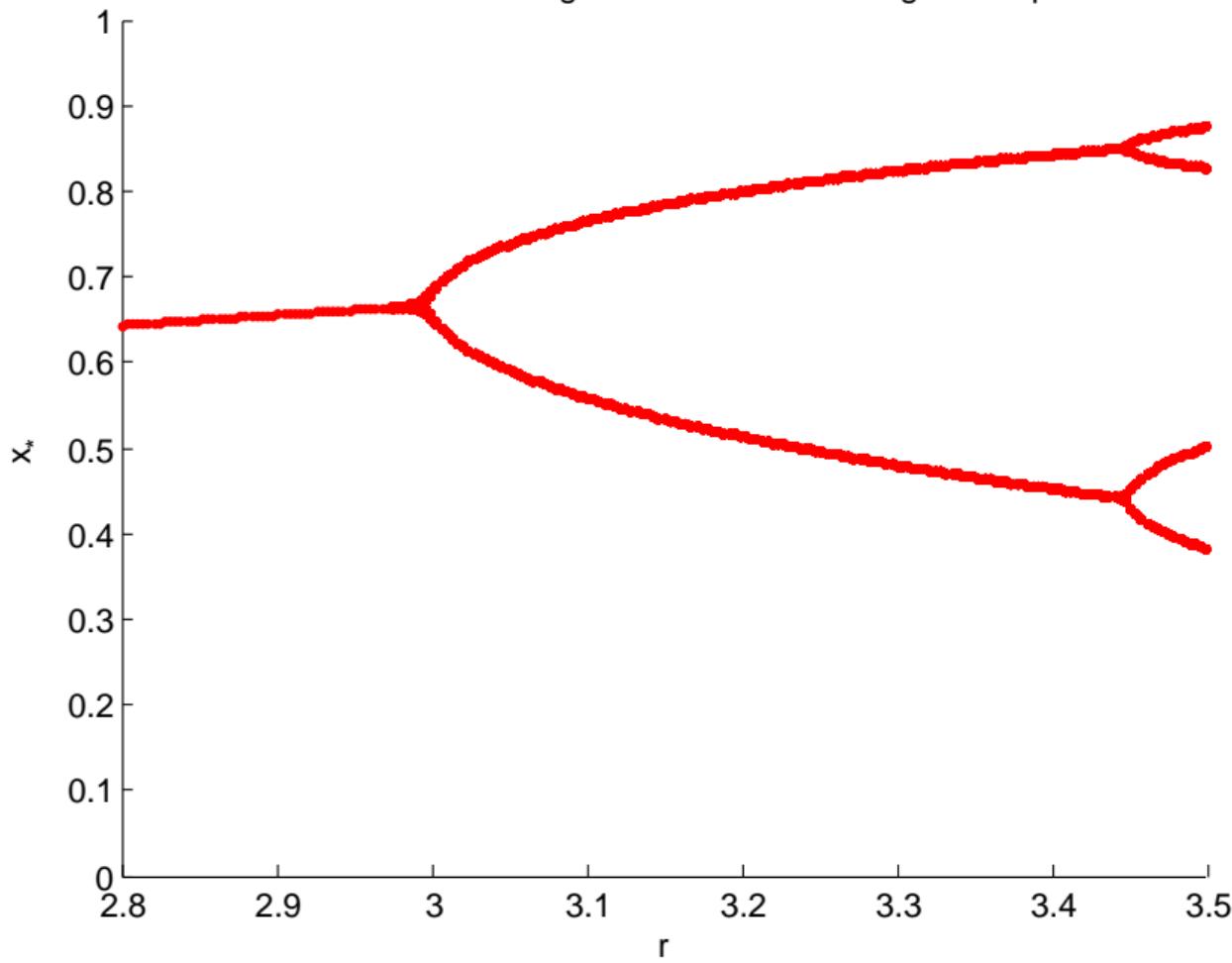
Bifurcation diagram for the discrete logistic map



Bifurcation diagram for the discrete logistic map



Bifurcation diagram for the discrete logistic map



The period-doubling cascade to chaos

The logistic map undergoes a sequence of period doubling bifurcations, called the **period-doubling cascade**, as r increases from 3 to 4.

- ▶ Every successive bifurcation leads to a doubling of the period.
- ▶ The bifurcation points form a sequence, $\{r_n\}$, that has the property that

$$\lim_{n \rightarrow \infty} \frac{r_n - r_{n-1}}{r_{n+1} - r_n}$$

exists and is a constant, called the Feigenbaum constant, equal to 4.669202...

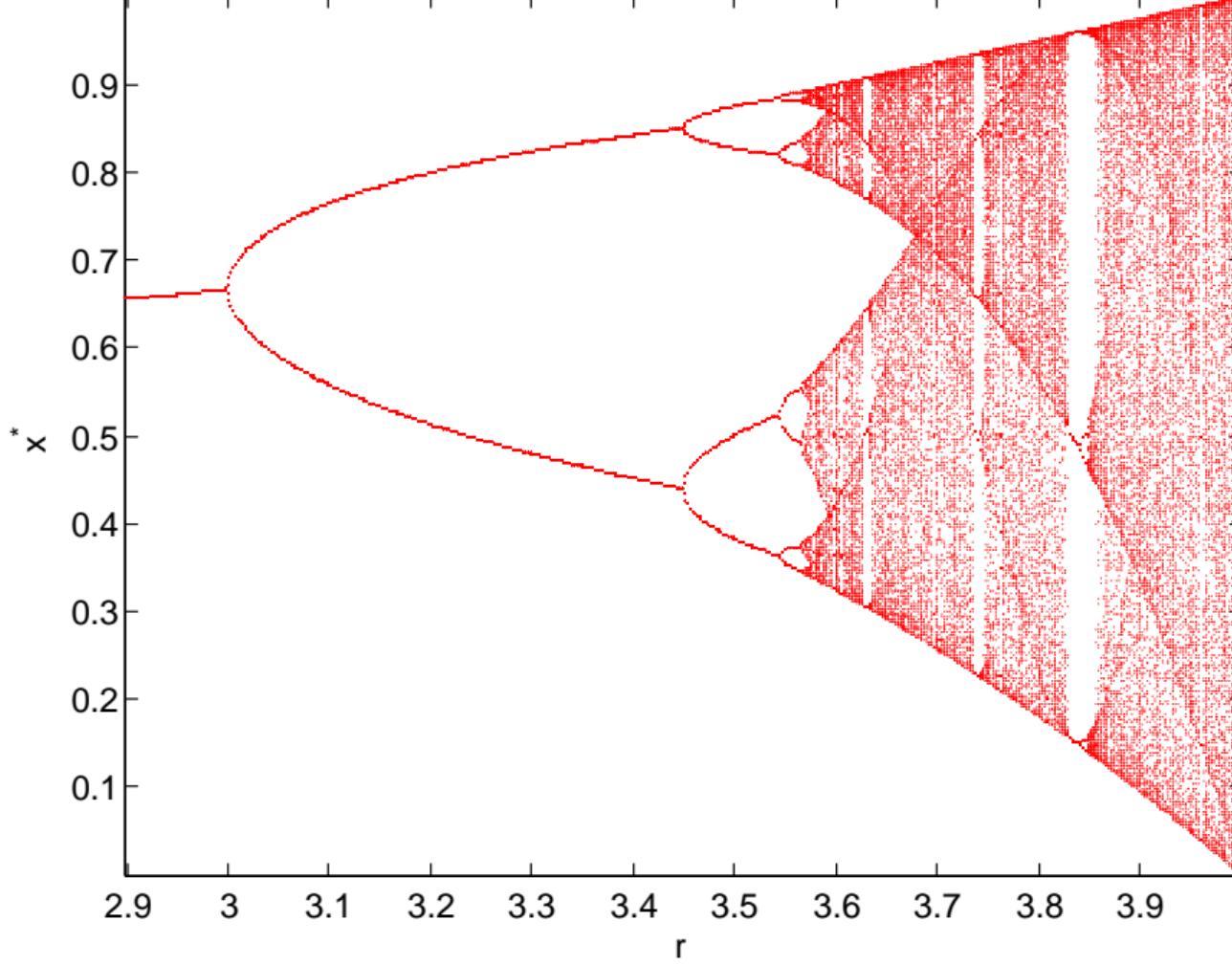
- ▶ This constant has been shown to exist in many of the maps that undergo the same type of cascade of period doubling bifurcations.

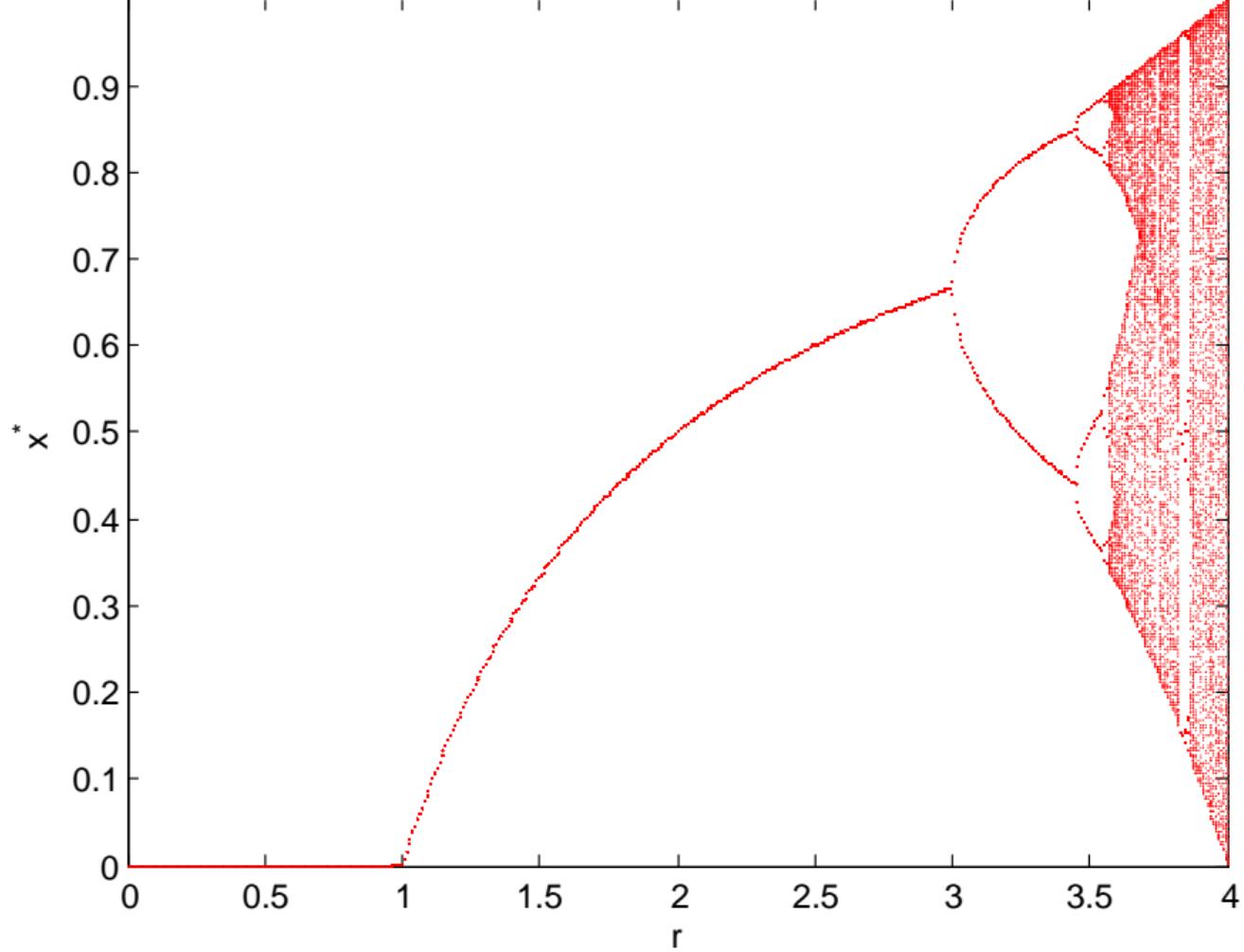
Chaos

After a certain value of r , there are periodic points with all periods. In particular, there are periodic points of period 3.

By a theorem (called **Sarkovskii's theorem**), the presence of period 3 points implies the presence of points of all periods.

At this point, the system is said to be in a **chaotic regime**, or **chaotic**.





Conclusion – A word of caution

We have used three different modelling paradigms to describe the growth of a population in a **logistic** framework:

- ▶ The ODE version has monotone solutions converging to the carrying capacity K
- ▶ The DDE version has oscillatory solutions, either converging to K or, if the delay is too large, periodic about K
- ▶ The discrete time version has all sorts of behaviors, and can be chaotic

The **choice of modelling method** is almost **as important** in the outcome of the model as the precise formulation/hypotheses of the **model**