Quick Recap

• Chapters 1–3: The basics

- Introduction to R
- Basic R commands
- Numerical summaries

Quick Recap

• Chapters 1–3: The basics

- Introduction to R
- Basic R commands
- Numerical summaries

• Chapters 4–5: More advanced R stuff

- Graphics
- Functions
- Loops
- Control statements

Quick Recap

• Chapters 1–3: The basics

- Introduction to R
- Basic R commands
- Numerical summaries

• Chapters 4–5: More advanced R stuff

- Graphics
- Functions
- Loops
- Control statements

Chapters 6–7: Randomness

- Random number generation
- Simulating random observations from discrete probability distributions

Basic idea: Repeatedly generate random numbers/samples to approximate a quantity of interest

Basic idea: Repeatedly generate random numbers/samples to approximate a quantity of interest

▲ロト ▲母 ▶ ▲ ヨ ▶ ▲ ヨ ● つんで

• Monte Carlo integration ("numerical integration")

Basic idea: Repeatedly generate random numbers/samples to approximate a quantity of interest

- Monte Carlo integration ("numerical integration")
- Simulation studies:

Basic idea: Repeatedly generate random numbers/samples to approximate a quantity of interest

- Monte Carlo integration ("numerical integration")
- Simulation studies:
 - Monoploy

Basic idea: Repeatedly generate random numbers/samples to approximate a quantity of interest

- Monte Carlo integration ("numerical integration")
- Simulation studies:
 - Monoploy
 - Replicating rolls on a die: the binomial distribution

Basic idea: Repeatedly generate random numbers/samples to approximate a quantity of interest

- Monte Carlo integration ("numerical integration")
- Simulation studies:
 - Monoploy
 - Replicating rolls on a die: the binomial distribution
 - The distribution of the sample mean: CLT

We generate a random data point from a simulation grid. Let

 $A = \{$ The data point lies below the curve $\}.$

We generate a random data point from a simulation grid. Let

 $A = \{$ The data point lies below the curve $\}.$

Then

$$\Pr[A] = \frac{\text{Area under curve}}{\text{Area of simulation grid}}$$

We generate a random data point from a simulation grid. Let

 $A = \{$ The data point lies below the curve $\}.$

Then

$$\Pr[A] = \frac{\text{Area under curve}}{\text{Area of simulation grid}}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Area under curve =

We generate a random data point from a simulation grid. Let

 $A = \{$ The data point lies below the curve $\}.$

◆□▶ ◆御▶ ◆臣▶ ◆臣▶ 三臣 - の�?

Then

$$\Pr[A] = \frac{\text{Area under curve}}{\text{Area of simulation grid}}$$
Area under curve = $\int_{a}^{b} f(x) dx =$

We generate a random data point from a simulation grid. Let

 $A = \{$ The data point lies below the curve $\}.$

Then

$$\Pr[A] = \frac{\text{Area under curve}}{\text{Area of simulation grid}}$$

Area under curve = $\int_{a}^{b} f(x) dx = \Pr[A] \times \text{Area of simulation grid}$

We generate a random data point from a simulation grid. Let

 $A = \{$ The data point lies below the curve $\}.$

Then

$$\Pr[A] = \frac{\text{Area under curve}}{\text{Area of simulation grid}}$$
Area under curve = $\int_{a}^{b} f(x) dx = \Pr[A] \times \text{Area of simulation grid}$

$$\approx \left[\frac{\text{No. of "hits"}}{\text{No. of points simulated}}\right]$$

imes Area of simulation grid

Part IX

Kernel Density Estimation

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

9.1 Introduction



Figure: (a) Histogram of ten values sampled from a N(0, 1) distribution. (b) Three different Kernel density estimators. The data are the X's.

9.2 Definition

A kernel is a non-negative real-valued integrable function K which satisfies the following two requirements:

$$\int_{-\infty}^{\infty} K(t) \, dt = 1 \tag{9.1}$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○□ のへで

and

$$K(-t) = K(t)$$
 for all values of t . (9.2)

Expression 9.1 ensures that the kernel is a pdf, whilst Expression 9.2 makes the distribution symmetric about 0.

9.2 Definition



Figure: (a) The Epanechnikov and Uniform kernels. (b) The triangular and Gaussian kernel.

<ロト <回ト < 注ト < 注ト = 注

Epanechnikov:

$$K(t) = \begin{cases} \frac{3}{4}(1-t^2) & -1 < t < 1\\ 0 & \text{otherwise.} \end{cases}$$

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

Epanechnikov:

$$K(t) = \begin{cases} \frac{3}{4}(1-t^2) & -1 < t < 1\\ 0 & \text{otherwise.} \end{cases}$$

Uniform:

$$K(t) = \begin{cases} \frac{1}{2} & -1 < t < 1\\ 0 & \text{otherwise.} \end{cases}$$

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

Epanechnikov:

$$K(t) = \begin{cases} \frac{3}{4}(1-t^2) & -1 < t < 1\\ 0 & \text{otherwise.} \end{cases}$$

Uniform:

$$\mathcal{K}(t) = \begin{cases} \frac{1}{2} & -1 < t < 1\\ 0 & \text{otherwise.} \end{cases}$$

Triangular:

$$\mathcal{K}(t) = \begin{cases} 1 - |t| & -1 < t < 1 \\ 0 & \text{otherwise.} \end{cases}$$

Epanechnikov:

$$K(t) = \begin{cases} \frac{3}{4}(1-t^2) & -1 < t < 1\\ 0 & \text{otherwise.} \end{cases}$$

Uniform:

$$K(t) = \begin{cases} \frac{1}{2} & -1 < t < 1\\ 0 & \text{otherwise.} \end{cases}$$

Triangular:

$$\mathcal{K}(t) = \begin{cases} 1 - |t| & -1 < t < 1 \\ 0 & \text{otherwise.} \end{cases}$$

Gaussian:

$$K(t)=\frac{1}{\sqrt{2\pi}}e^{-t^2/2}.$$

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで



Figure: (a) The Epanechnikov and Uniform kernels. (b) The triangular and Gaussian kernel.

・ロト ・四ト ・ヨト ・ヨト

æ



Kernel density estimation can be summarised in four steps:

Kernel density estimation can be summarised in four steps:

• We have some sample data. In Figure 9.3a we have three points, so n = 3.

Kernel density estimation can be summarised in four steps:

• We have some sample data. In Figure 9.3a we have three points, so n = 3.

▲ロト ▲母 ▶ ▲ ヨ ▶ ▲ ヨ ● つんで

Around each of the data points, we draw a kernel.

Kernel density estimation can be summarised in four steps:

- We have some sample data. In Figure 9.3a we have three points, so n = 3.
- Around each of the data points, we draw a kernel. In Figure 9.3b we have used a Gaussian kernel. However, we could have used a Uniform, triangular, or Epanechnikov kernel.

Kernel density estimation can be summarised in four steps:

- We have some sample data. In Figure 9.3a we have three points, so n = 3.
- Around each of the data points, we draw a kernel. In Figure 9.3b we have used a Gaussian kernel. However, we could have used a Uniform, triangular, or Epanechnikov kernel.

▲ロト ▲母 ▶ ▲ ヨ ▶ ▲ ヨ ● つんで

Solution Next we combine the kernels - the blue dashed line in Figure 9.3c.

Kernel density estimation can be summarised in four steps:

- We have some sample data. In Figure 9.3a we have three points, so n = 3.
- Around each of the data points, we draw a kernel. In Figure 9.3b we have used a Gaussian kernel. However, we could have used a Uniform, triangular, or Epanechnikov kernel.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

- Solution Next we combine the kernels the blue dashed line in Figure 9.3c.
- The final step is to normalise the distribution.

Kernel density estimation can be summarised in four steps:

- We have some sample data. In Figure 9.3a we have three points, so n = 3.
- Around each of the data points, we draw a kernel. In Figure 9.3b we have used a Gaussian kernel. However, we could have used a Uniform, triangular, or Epanechnikov kernel.
- Next we combine the kernels the blue dashed line in Figure 9.3c.
- The final step is to normalise the distribution. In our example, since we have three points, the total area under the blue dashed curve is 3.

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで

Kernel density estimation can be summarised in four steps:

- We have some sample data. In Figure 9.3a we have three points, so n = 3.
- Around each of the data points, we draw a kernel. In Figure 9.3b we have used a Gaussian kernel. However, we could have used a Uniform, triangular, or Epanechnikov kernel.
- Next we combine the kernels the blue dashed line in Figure 9.3c.
- The final step is to normalise the distribution. In our example, since we have three points, the total area under the blue dashed curve is 3. Hence, to recover a density we divide by 3 to get the black curve in Figure 9.3d.

▲ロト ▲母 ▶ ▲ 国 ▶ ▲ 国 ● の Q @

▲ロト ▲御 ▶ ▲ 臣 ▶ ▲ 臣 ▶ ○臣 ○ の Q @

Let K be a kernel

Let *K* be a kernel and suppose our sample contains *n* values:

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

Let K be a kernel and suppose our sample contains n values: x_1, \ldots, x_n .

◆□▶ ◆□▶ ◆三▶ ◆三▶ 三三 のへで

Let *K* be a kernel and suppose our sample contains *n* values: x_1, \ldots, x_n . Then our estimate of the true pdf f(x) is

Let *K* be a kernel and suppose our sample contains *n* values: x_1, \ldots, x_n . Then our estimate of the true pdf f(x) is

$$\hat{f}(x) =$$

Let *K* be a kernel and suppose our sample contains *n* values: x_1, \ldots, x_n . Then our estimate of the true pdf f(x) is

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^{n} K(x - x_i) .$$
(9.3)

Let *K* be a kernel and suppose our sample contains *n* values: x_1, \ldots, x_n . Then our estimate of the true pdf f(x) is

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^{n} K(x - x_i)$$
 (9.3)

◆□▶ ◆御▶ ◆臣▶ ◆臣▶ 三臣 - の�?

Notice we use $K(x - x_i)$, since we draw a kernel around *each* x_i .

It's fairly straightforward to see that $\hat{f}(x)$ is also a pdf, namely

It's fairly straightforward to see that $\hat{f}(x)$ is also a pdf, namely

$$\int_{-\infty}^{\infty} \hat{f}(x) \, dx =$$

It's fairly straightforward to see that $\hat{f}(x)$ is also a pdf, namely

$$\int_{-\infty}^{\infty} \hat{f}(x) dx = \frac{1}{n} \sum_{i=1}^{n} \int_{-\infty}^{\infty} K(x - x_i) dx$$

It's fairly straightforward to see that $\hat{f}(x)$ is also a pdf, namely

$$\int_{-\infty}^{\infty} \hat{f}(x) dx = \frac{1}{n} \sum_{i=1}^{n} \int_{-\infty}^{\infty} K(x - x_i) dx$$
$$= \frac{1}{n} \sum_{i=1}^{n} \int_{-\infty}^{\infty} K(y) dy$$

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

It's fairly straightforward to see that $\hat{f}(x)$ is also a pdf, namely

$$\int_{-\infty}^{\infty} \hat{f}(x) dx = \frac{1}{n} \sum_{i=1}^{n} \int_{-\infty}^{\infty} K(x - x_i) dx$$
$$= \frac{1}{n} \sum_{i=1}^{n} \int_{-\infty}^{\infty} K(y) dy$$
$$= \frac{1}{n} \sum_{i=1}^{n} 1$$

It's fairly straightforward to see that $\hat{f}(x)$ is also a pdf, namely

$$\int_{-\infty}^{\infty} \hat{f}(x) dx = \frac{1}{n} \sum_{i=1}^{n} \int_{-\infty}^{\infty} K(x - x_i) dx$$
$$= \frac{1}{n} \sum_{i=1}^{n} \int_{-\infty}^{\infty} K(y) dy$$
$$= \frac{1}{n} \sum_{i=1}^{n} 1$$
$$= \frac{1}{n} (1 + 1 + \dots + 1)$$

It's fairly straightforward to see that $\hat{f}(x)$ is also a pdf, namely

$$\int_{-\infty}^{\infty} \hat{f}(x) dx = \frac{1}{n} \sum_{i=1}^{n} \int_{-\infty}^{\infty} K(x - x_i) dx$$
$$= \frac{1}{n} \sum_{i=1}^{n} \int_{-\infty}^{\infty} K(y) dy$$
$$= \frac{1}{n} \sum_{i=1}^{n} 1$$
$$= \frac{1}{n} (1 + 1 + \dots + 1)$$
$$= \frac{1}{n} \times n = 1.$$

(9.4)

The data shown below are the lengths (to the nearest cm) of 10 Giant Groupers caught by expert angler Jeremy Wade for the TV series *River Monsters*.

This sample was taken in 2013 in a lake near the Chernobyl nuclear disaster of 1986. These fish usually grow to around 75cm in length, but genetic mutations caused by the nuclear explosions are thought to have increased the size of this species.

101	97	99	104	103	94	102	94	102	106
-----	----	----	-----	-----	----	-----	----	-----	-----

▲□▶ ▲□▶ ▲□▶ ▲□▶ □ のへで



Produce a density plot for these data using the Gaussian kernel.

Produce a density plot for these data using the Gaussian kernel.

Let *X* represent the length of a fish. We have

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^{n} K(x - x_i),$$

where

$$K(x-x_i) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(x-x_i)^2}{2}\right\}.$$

Produce a density plot for these data using the Gaussian kernel.

Let X represent the length of a fish. We have

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^{n} K(x - x_i),$$

where

$$K(x-x_i) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(x-x_i)^2}{2}\right\}.$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○□ ○○○

The range of our data is $94 \rightarrow 106$, so let's plot over the range $90 \rightarrow 110$.

Produce a density plot for these data using the Gaussian kernel.

Let X represent the length of a fish. We have

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^{n} K(x - x_i),$$

where

$$K(x-x_i) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(x-x_i)^2}{2}\right\}.$$

The range of our data is $94 \rightarrow 106$, so let's plot over the range $90 \rightarrow 110$. For example,

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○□ ○○○

Produce a density plot for these data using the Gaussian kernel.

Let X represent the length of a fish. We have

$$\hat{f}(x) = \frac{1}{n} \sum_{i=1}^{n} K(x - x_i),$$

where

$$K(x-x_i) = \frac{1}{\sqrt{2\pi}} \exp\left\{-\frac{(x-x_i)^2}{2}\right\}.$$

The range of our data is $94 \rightarrow 106$, so let's plot over the range $90 \rightarrow 110$. For example,

$$\hat{f}(90) = \frac{1}{10} \left[\frac{1}{\sqrt{2\pi}} e^{-(90-101)^2/2} + \ldots + \frac{1}{\sqrt{2\pi}} e^{-(90-106)^2/2} \right]$$

◆□▶ ◆□▶ ◆三▶ ◆三▶ ○□ ○○○

= 0.0000267

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

As another example:

As another example:

$$\hat{f}(\mathbf{94}) = \frac{1}{10} \left[\frac{1}{\sqrt{2\pi}} e^{-(\mathbf{94}-\mathbf{101})^2/2} + \ldots + \frac{1}{\sqrt{2\pi}} e^{-(\mathbf{94}-\mathbf{106})^2/2} \right]$$

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

= 0.08023.

As another example:

$$\hat{f}(\mathbf{94}) = \frac{1}{10} \left[\frac{1}{\sqrt{2\pi}} e^{-(\mathbf{94}-\mathbf{101})^2/2} + \ldots + \frac{1}{\sqrt{2\pi}} e^{-(\mathbf{94}-\mathbf{106})^2/2} \right]$$

= 0.08023.

Also,

$$\hat{f}(102) = \frac{1}{10} \left[\frac{1}{\sqrt{2\pi}} e^{-(102 - 101)^2/2} + \ldots + \frac{1}{\sqrt{2\pi}} e^{-(102 - 106)^2/2} \right]$$
$$= 0.13404.$$

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

This gives:

X	$\hat{f}(x)$	X	$\hat{f}(x)$
90	0.0000267	101	0.0995432
91	0.0000886	102	0.1340384
92	0.0107983	103	0.1183411
93	0.0484075	104	0.0807319
94	0.0802317	105	0.0546929
95	0.0538066	106	0.0457634
96	0.0354386	107	0.0246539
97	0.0461933	108	0.0054126
98	0.0488910	109	0.0004433
99	0.0515926	110	0.0000138
100	0.0600920		

・ロト ・御ト ・ヨト ・ヨト

王

▲□▶ ▲□▶ ▲三▶ ▲三▶ 三三 のへで

Plotting $\hat{f}(x)$ against *x* gives:

Plotting $\hat{f}(x)$ against *x* gives:



Plotting $\hat{f}(x)$ against *x* gives:



Plotting $\hat{f}(x)$ against *x* gives:



・ロト ・日下・ ・ 田下・ ・

æ

ъ