

# 5

## Interpolation

# 插值

分段插值函数，通过已知数据点



人们思考皆，浮皮潦草，泛泛而谈；现实世界却，盘根错节，千头万绪。

*We think in generalities, but we live in details.*

—— 阿尔弗雷德·怀特海 (Alfred Whitehead) | 英国数学家、哲学家 | 1861 ~ 1947



- ◀ `scipy.interpolate.interp1d()` 一维插值
- ◀ `scipy.interpolate.lagrange()` 拉格朗日多项式插值
- ◀ `scipy.interpolate.interp2d()` 二维插值，网格化数据
- ◀ `matplotlib.pyplot.pcolormesh()` 绘制填充颜色网格数据
- ◀ `scipy.interpolate.griddata()` 二维插值，散点化数据
- ◀ `matplotlib.pyplot.imshow()` 绘制数据平面图像

本 PDF 文件为作者草稿，发布目的为方便读者在移动终端学习，终稿内容以清华大学出版社纸质出版物为准。

版权归清华大学出版社所有，请勿商用，引用请注明出处。

代码及 PDF 文件下载：<https://github.com/Visualize-ML>

本书配套微课视频均发布在 B 站——生姜 DrGinger：<https://space.bilibili.com/513194466>

欢迎大家批评指教，本书专属邮箱：[jiang.visualize.ml@gmail.com](mailto:jiang.visualize.ml@gmail.com)

## 5.1 插值

插值是通过已知数据点之间的值来估计未知点的值的方法，它可以用于填补数据缺失或者进行数据平滑处理。插值方法通常基于已知数据点之间的关系，通过数学函数或者曲线拟合等方法来预测未知数据点的值。

如图 1 所示的蓝色点为已知数据点，插值就是根据这几个离散的数据点估算其他点对应的  $y$  值。

插值可分为**内插** (interpolation) 和**外插** (extrapolation)。内插是在已知数据点之间进行插值，估计出未知点的值。而外插则是在已知数据点的范围之外进行插值，从而预测超出已知数据点范围的未知点的值。在进行外插时，需要考虑插值函数是否能够正确地拟合未知数据点，并且需要注意不要过度依赖插值函数来进行预测，以免导致不可靠的预测结果。

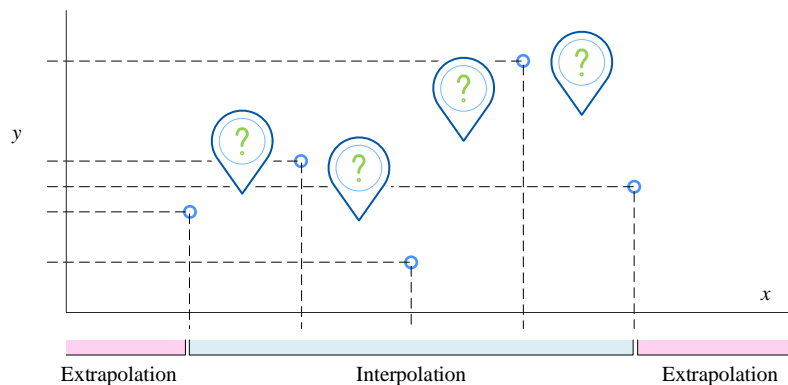


图 1. 插值的意义

### 常见插值方法

图 2 总结常用的插值的算法。图 2 相当于本章的思维导图。

本章主要介绍如下几种方法：

- ◀ **常数插值** (constant interpolation)，比如**向前** (previous 或 forward)、**向后** (next 或 backward)、**最近** (nearest)；
- ◀ **线性插值** (linear interpolation)；
- ◀ **二次插值** (quadratic interpolation)，本章不做介绍；
- ◀ **三次插值** (cubic interpolation)；
- ◀ **拉格朗日插值** (Lagrange polynomial interpolation)。

本章最后还要介绍**二维插值** (bivariate interpolation)，二维插值将一元插值的方法推广到二维。

此外，对于时间序列，处理缺失值或者获得颗粒度更高的数据，都可以使用插值。图 3 所示为利用线性插值插补时间序列数据中的缺失值。

➡ 《可视之美》介绍的贝塞尔曲线本质上也是插值。贝塞尔曲线是一种通过一系列控制点来定义曲线形状的数学函数。在计算机图形学和计算机辅助设计中，常使用贝塞尔曲线来生成平滑的曲线形状。

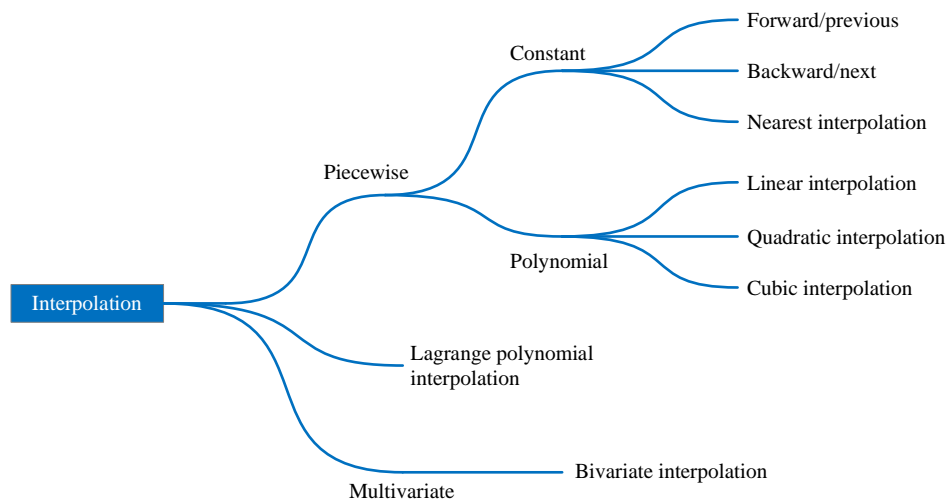


图 2. 插值的分类

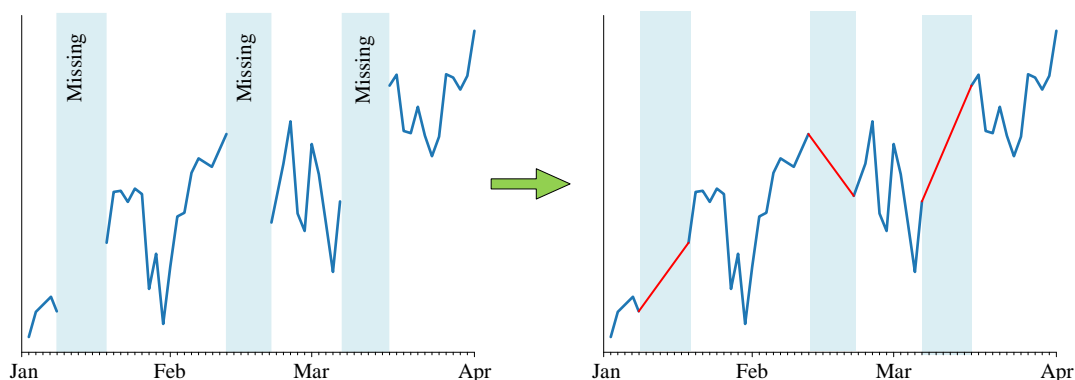


图 3. 时间序列插值

## 分段函数

虽然，一些插值分段函数构造得到的曲线整体看上去平滑。但是绝大多数情况，插值函数是分段函数，因此插值也称**分段插值** (piecewise interpolation)。



《数学要素》第 11 章介绍过分段函数。

对于一元函数  $f(x)$ ，分段函数是指自变量  $x$  在不同取值范围对应不同解析式的函数。

每两个相邻的数据点之间便对应不同解析式：

$$f(x) = \begin{cases} f_1(x) & x^{(1)} \leq x < x^{(2)} \\ f_2(x) & x^{(2)} \leq x < x^{(3)} \\ \dots & \dots \\ f_{n-1}(x) & x^{(n-1)} \leq x < x^{(n)} \end{cases} \quad (1)$$

其中， $n$  为已知点个数。

**▲ 注意**，上式中  $f_i(x)$  代表一个特定解析式。分段函数虽然由一系列解析式构成，但是分段函数还是一个函数。

如图 4 所示，已知数据点一共有五个—— $(x^{(1)}, y^{(1)})$ 、 $(x^{(2)}, y^{(2)})$ 、 $(x^{(3)}, y^{(3)})$ 、 $(x^{(4)}, y^{(4)})$ 、 $(x^{(5)}, y^{(5)})$ 。比如，分段函数  $f(x)$  在  $[x^{(1)}, x^{(2)}]$  区间的解析式为  $f_1(x)$ 。 $f_1(x)$  通过  $(x^{(1)}, y^{(1)})$ 、 $(x^{(2)}, y^{(2)})$  两个已知数据点。图 4 实际上就是线性插值。

(1) 还告诉我们，对于内插， $n$  个已知点可以构成  $n-1$  个区间，即分段函数有  $n-1$  个解析式。

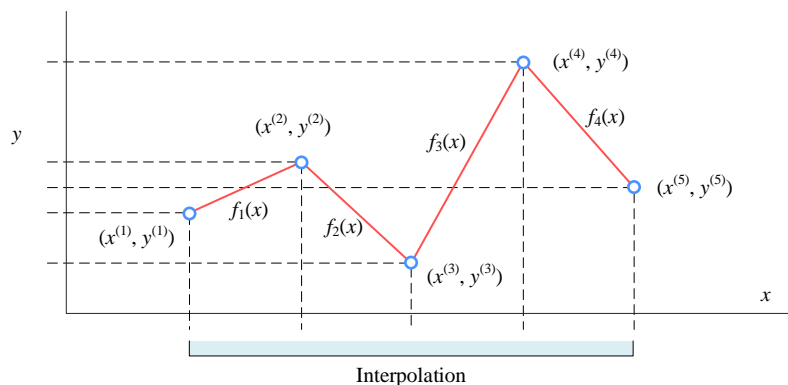


图 4. 分段函数

## 拟合、插值

大家经常混淆拟合和插值这两种方法。插值和拟合有一个相同之处，它们都是根据已知数据点，构造函数，从而推断得到更多数据点。

插值和回归都是用于对数据进行预测的方法，但两者有明显的区别。插值是用于填补已有数据点之间的空缺，预测未知点的值。回归则是预测自变量和因变量之间的关系。插值通常使用插值函数，如多项式插值；而回归则通过拟合数据点的回归方程来预测因变量的值。插值通常用于数据平滑处理、数据填补等。回归则常用于预测和建模。插值要求原始数据点之间要有一定的连续性和平滑性；而回归则对数据点的分布没有明显要求。插值得到的是精确的函数值，但在超出已有数据范围时可能不准确；而回归得到的是变量之间的大致关系，可以预测未来的趋势。

需要根据具体情况选择合适的方法。当数据缺失或需要平滑处理时，可以使用插值方法；当需要建立模型并预测未来趋势时，可以使用回归方法。

插值一般得到分段函数，分段函数通过所有给定的数据点，如图 5 (a)、(b) 所示。回归拟合得到的函数尽可能靠近样本数据点，如图 5 (c)、(d) 所示。图 6 比较二维插值和二维回归。

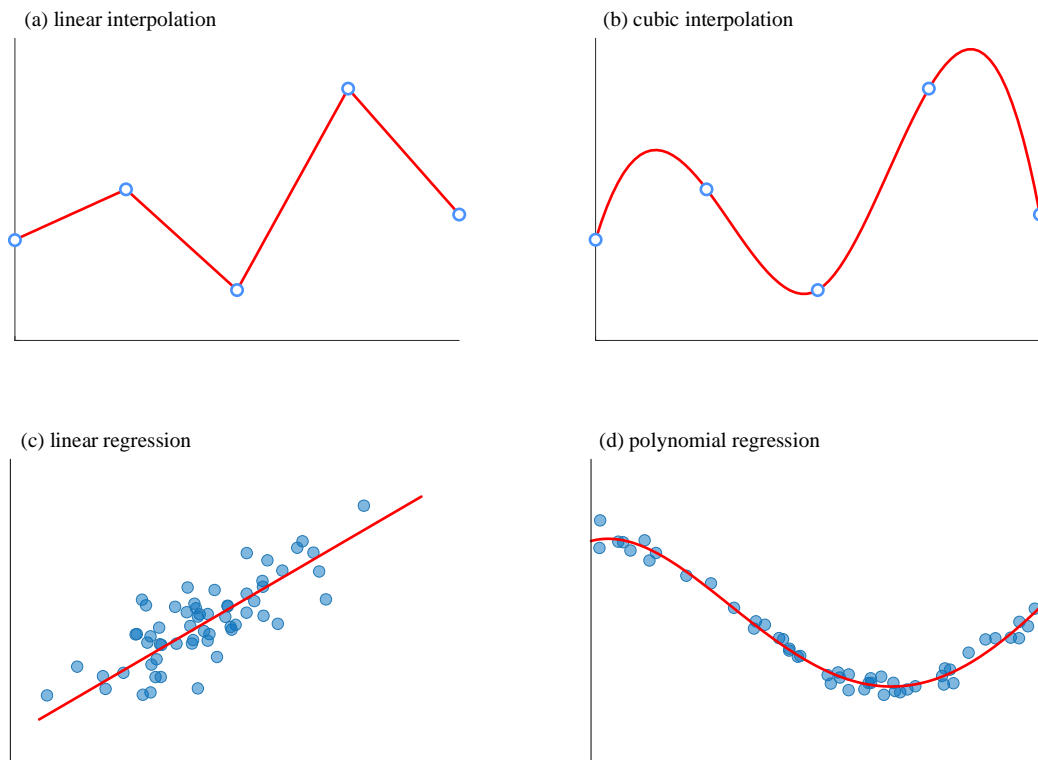


图 5. 比较一维插值和回归

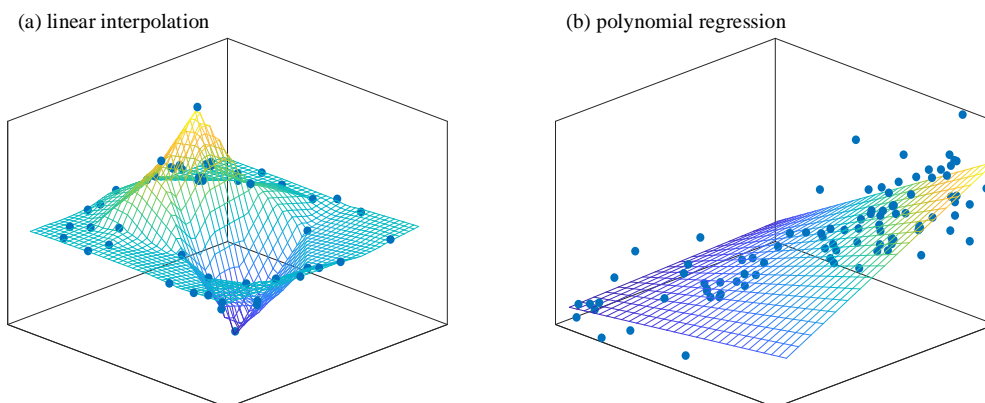


图 6. 比较二维插值和二维回归

## 5.2 常数插值：分段函数为阶梯状

本节介绍常用的三种常数插值方法。

### 向前

向前常数插值对应的分段函数为：

$$f(x) = \begin{cases} f_1(x) = x^{(1)} & x^{(1)} \leq x < x^{(2)} \\ f_2(x) = x^{(2)} & x^{(2)} \leq x < x^{(3)} \\ \dots & \dots \\ f_{n-1}(x) = x^{(n-1)} & x^{(n-1)} \leq x < x^{(n)} \end{cases} \quad (2)$$

如图 7 所示，向前常数插值用区间  $[x^{(i)}, x^{(i+1)}]$  左侧端点，即  $x^{(i)}$ ，对应的  $y^{(i)}$ ，作为常数函数的取值。图 7 中红色划线为真实函数取值。

对于数据帧 `df`，如果存在 NaN 的话，`df.fillna(method = 'ffill')` 便对应向前常数插补。

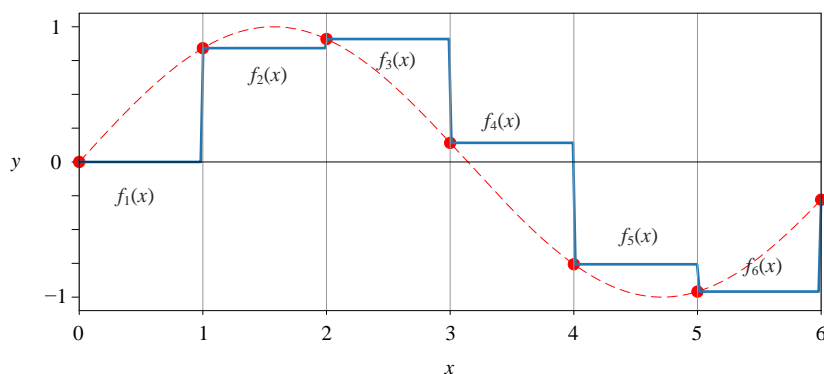


图 7. 向前常数插值

## 向后

向后常数插值对应的分段函数为：

$$f(x) = \begin{cases} f_1(x) = x^{(2)} & x^{(1)} \leq x < x^{(2)} \\ f_2(x) = x^{(3)} & x^{(2)} \leq x < x^{(3)} \\ \dots & \dots \\ f_{n-1}(x) = x^{(n)} & x^{(n-1)} \leq x < x^{(n)} \end{cases} \quad (3)$$

如图 8 所示，向后常数插值和图 7 正好相反。

对于数据帧 df，如果存在 NaN 的话，df.fillna(method = 'bfill') 对应向后常数插补。

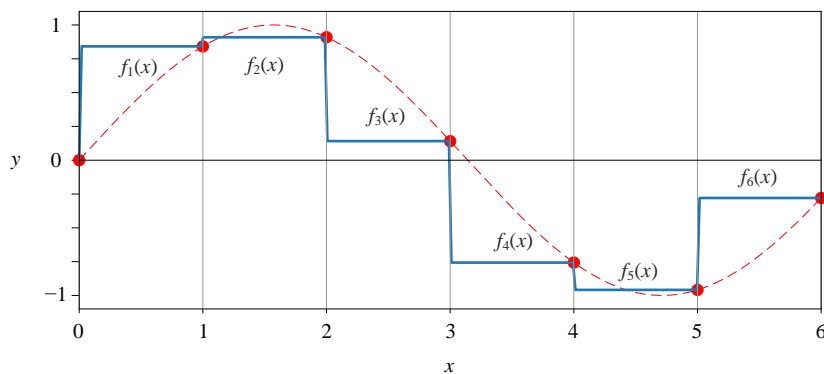


图 8. 向后常数插值

## 最邻近

最邻近插值的分段函数为：

$$f(x) = \begin{cases} f_1(x) = x^{(1)} & x^{(1)} \leq x < \frac{x^{(1)} + x^{(2)}}{2} \\ f_2(x) = x^{(2)} & \frac{x^{(1)} + x^{(2)}}{2} \leq x < \frac{x^{(2)} + x^{(3)}}{2} \\ \dots & \dots \\ f_n(x) = x^{(n)} & \frac{x^{(n-1)} + x^{(n)}}{2} \leq x < x^{(n)} \end{cases} \quad (4)$$

如图9所示，最邻近常数插值相当于“向前”和“向后”常数插值的“折中”。分段插值函数同样是阶梯状，只不过阶梯发生在两个相邻已知点中间处。

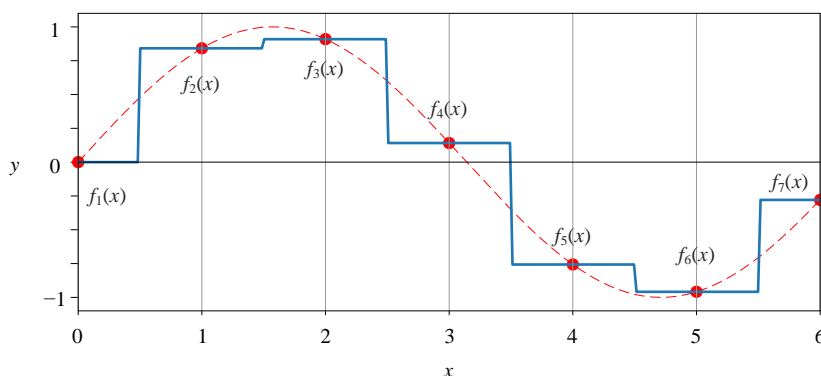


图 9. 最邻近常数插值

## 5.3 线性插值：分段函数为线段

对于线性插值，区间  $[x^{(i)}, x^{(i+1)}]$  对应的  $f_i(x)$  为：

$$f_i(x) = \underbrace{\left( \frac{y^{(i)} - y^{(i+1)}}{x^{(i)} - x^{(i+1)}} \right)}_{\text{slope}} (x - x^{(i+1)}) + y^{(i+1)} \quad (5)$$



容易发现，上式就是《数学要素》第 11 章介绍的一元函数的点斜式。

也就是说，不考虑区间的话，上式代表通过  $(x^{(i)}, y^{(i)})$ 、 $(x^{(i+1)}, y^{(i+1)})$  两点的一条直线。

图 10 所示为线性插值结果。白话说，线性插值就是用任意两个相邻已知点连接成的线段来估算其他未知点的值。



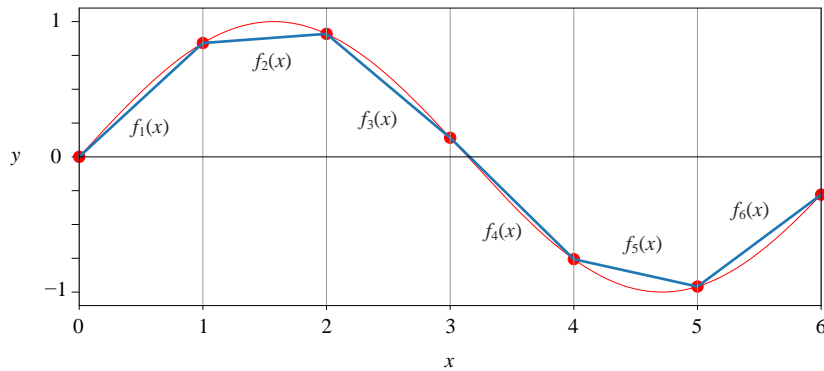


图 10. 线性插值

## 5.4 三次样条插值：光滑曲线拼接

图 11 所示为三次样条插值的结果。虽然，整条曲线看上去连续、光滑，实际上它是由四个函数拼接起来的分段函数。

对于三次样条插值，每一段的分段函数是三次多项式：

$$f_i(x) = a_i x^3 + b_i x^2 + c_i x + d_i \quad (6)$$

其中， $a_i$ 、 $b_i$ 、 $c_i$ 、 $d_i$  为需要求解的系数。

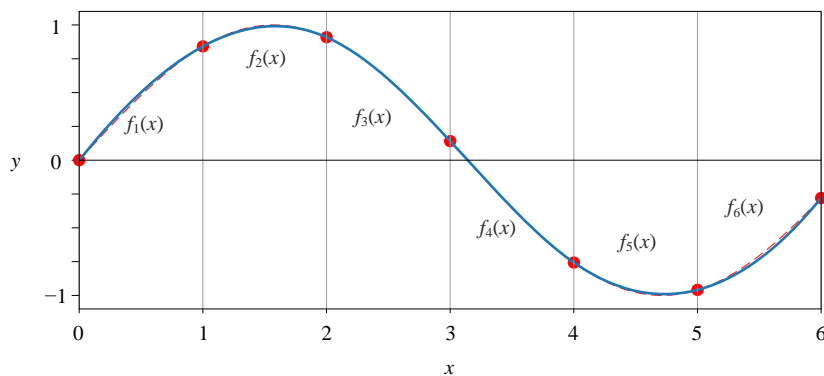


图 11. 三次样条插值

为了求解系数，我们需要构造一系列等式。类似线性插值，每一段三次函数通过区间  $[x^{(i)}, x^{(i+1)}]$  左右两点，即：

$$\begin{cases} f_i(x^{(i)}) = y^{(i)} & i = 1, 2, \dots, n-1 \\ f_i(x^{(i+1)}) = y^{(i+1)} & i = 1, 2, \dots, n-1 \end{cases} \quad (7)$$

曲线之所以看起来很平滑是因为，除两端样本数据点以外，内部数据点处，一阶和二阶导数等值：

$$\begin{cases} f'_i(x^{(i+1)}) = f'_{i+1}(x^{(i+1)}) & i = 1, 2, \dots, n-2 \\ f''_i(x^{(i+1)}) = f''_{i+1}(x^{(i+1)}) & i = 1, 2, \dots, n-2 \end{cases} \quad (8)$$

对于三次样条插值，一般还设定两端样本数据点处二阶导数为 0：

$$\begin{cases} f''_1(x^{(1)}) = 0 \\ f''_{n-1}(x^{(n)}) = 0 \end{cases} \quad (9)$$

插值中系数求解一般都是用矩阵运算完成。举个例子，在三次样条插值中，需要解出一个三对角线方程组，这个方程组可以用矩阵形式表示。具体来说，需要先确定每个小区间内的多项式系数，然后利用这些系数和每个小区间的边界点，构造一个三对角矩阵方程组，利用三对角矩阵求解方法，可以得到每个小区间内的多项式系数，从而得到整个分段函数。本章不展开讲解。



Bk6\_Ch05\_01.py 完成插值并绘制图 7 ~ 图 11。Python 进行一维插值函数为 `scipy.interpolate.interp1d()`，二维插值的函数为 `scipy.interpolate.interp2d()`。

## 5.5 拉格朗日插值

**拉格朗日插值** (Lagrange interpolation) 不同于本章前文介绍的插值方法。前文介绍的插值方法得到的都是分段函数，而拉格朗日插值得到的是一个高次多项式函数  $f(x)$ 。 $f(x)$  相当是由若干多项式函数叠加而成：

$$f(x) = \sum_{i=1}^n f_i(x) \quad (10)$$

其中，

$$f_i(x) = y^{(i)} \cdot \prod_{k=1, k \neq i}^n \frac{x - x^{(k)}}{x^{(i)} - x^{(k)}} \quad (11)$$

$f_i(x)$  展开来写：

$$f_i(x) = y^{(i)} \cdot \frac{(x-x^{(1)})(x-x^{(2)})\dots(x-x^{(i-1)})(x-x^{(i+1)})\dots(x-x^{(n)})}{(x^{(i)}-x^{(1)})(x^{(i)}-x^{(2)})\dots(x^{(i)}-x^{(i-1)})(x^{(i)}-x^{(i+1)})\dots(x^{(i)}-x^{(n)})} \quad (12)$$

比如,  $f_1(x)$  展开来写:

$$f_1(x) = y^{(1)} \cdot \frac{(x-x^{(2)})(x-x^{(3)})\dots(x-x^{(n)})}{(x^{(1)}-x^{(2)})(x^{(1)}-x^{(3)})\dots(x^{(1)}-x^{(n)})} \quad (13)$$

$f_2(x)$  展开来写:

$$f_2(x) = y^{(2)} \cdot \frac{(x-x^{(1)})(x-x^{(3)})\dots(x-x^{(n)})}{(x^{(2)}-x^{(1)})(x^{(2)}-x^{(3)})\dots(x^{(2)}-x^{(n)})} \quad (14)$$

### 举个例子

比如,  $n=3$ , 也就是有三个样本数据点  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)})\}$  的时候,  $f(x)$  为:

$$f(x) = y^{(1)} \cdot \underbrace{\frac{(x-x^{(2)})(x-x^{(3)})}{(x^{(1)}-x^{(2)})(x^{(1)}-x^{(3)})}}_{f_1(x)} + y^{(2)} \cdot \underbrace{\frac{(x-x^{(1)})(x-x^{(3)})}{(x^{(2)}-x^{(1)})(x^{(2)}-x^{(3)})}}_{f_2(x)} + y^{(3)} \cdot \underbrace{\frac{(x-x^{(1)})(x-x^{(2)})}{(x^{(3)}-x^{(1)})(x^{(3)}-x^{(2)})}}_{f_3(x)} \quad (15)$$

观察上式,  $f(x)$  相当于三个二次函数叠加得到。

将三个数据点  $\{(x^{(1)}, y^{(1)}), (x^{(2)}, y^{(2)}), (x^{(3)}, y^{(3)})\}$ , 逐一代入上式, 可以得到:

$$f(x^{(1)}) = y^{(1)}, \quad f(x^{(2)}) = y^{(2)}, \quad f(x^{(3)}) = y^{(3)} \quad (16)$$

也就是说, 多项式函数  $f(x)$  通过给定的已知点。

图 12 所示为拉格朗日插值结果。

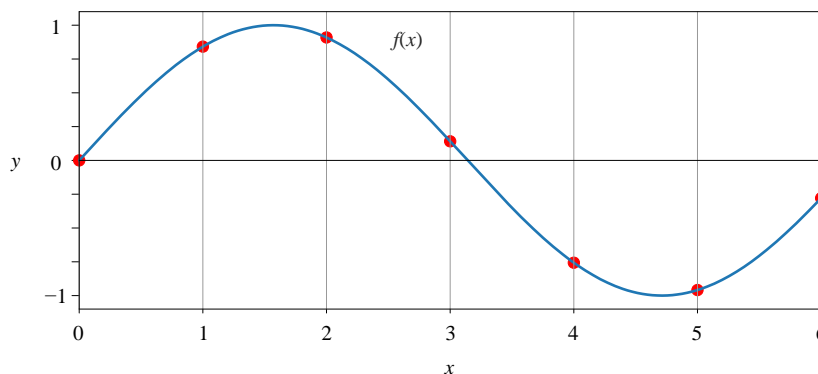


图 12. 拉格朗日插值

## 龙格现象

有一点需要大家注意的是，已知点数量  $n$  不断增大，拉格朗日插值函数多项式函数次数不断提高，插值多项式的插值逼近效果未必好。如图 13 所示，插值多项式（红色曲线）区间边缘处出现振荡问题，这一现象叫做**龙格现象** (Runge's phenomenon)。

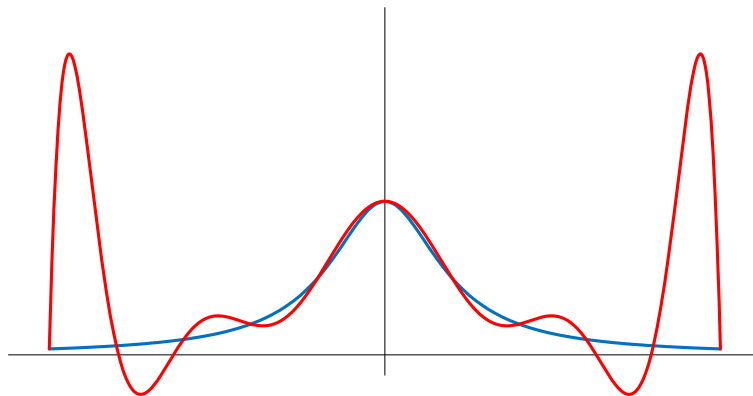


图 13. 龙格现象



Bk6\_Ch05\_02.py 完成拉格朗日插值，并绘制图 12。

## 5.6 二维插值

如图 14 所示，以二维线性插值为例，二维线性插值相当于处理了三个一维线性插值。

对于二维线性插值，先将二维坐标系中的点分别按照横坐标和纵坐标排序。然后，找到待插值点所在的四个相邻的点。分别对这四个点在横坐标和纵坐标上进行一维线性插值，得到在横向和纵向上的两个插值结果。将上述两个插值结果加权平均，作为待插值点的二维线性插值结果。其中，权重的计算基于待插值点相对于四个相邻点在横向和纵向上的距离，距离越远的点权重越小。

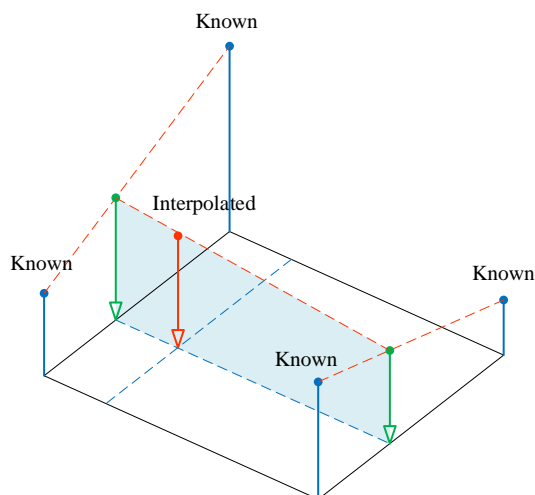


图 14. 二维线性插值原理

### 举个例子

图 15 中 × 为给定的已知数据。图 16 和图 17 所示为分别通过线性插值、三次样条插值完成的二维插值结果。二维插值用到的函数是 `scipy.interpolate.interp2d()`。

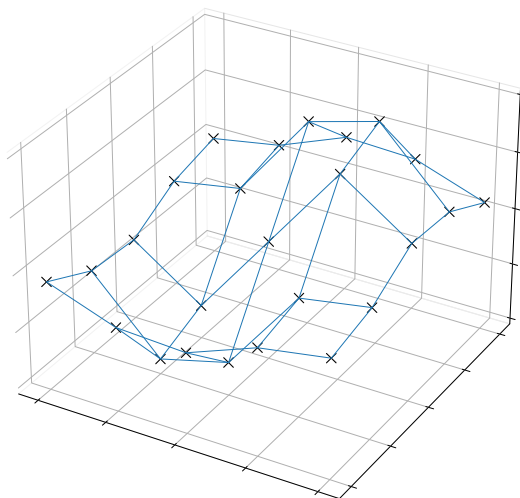


图 15. 已知数据点

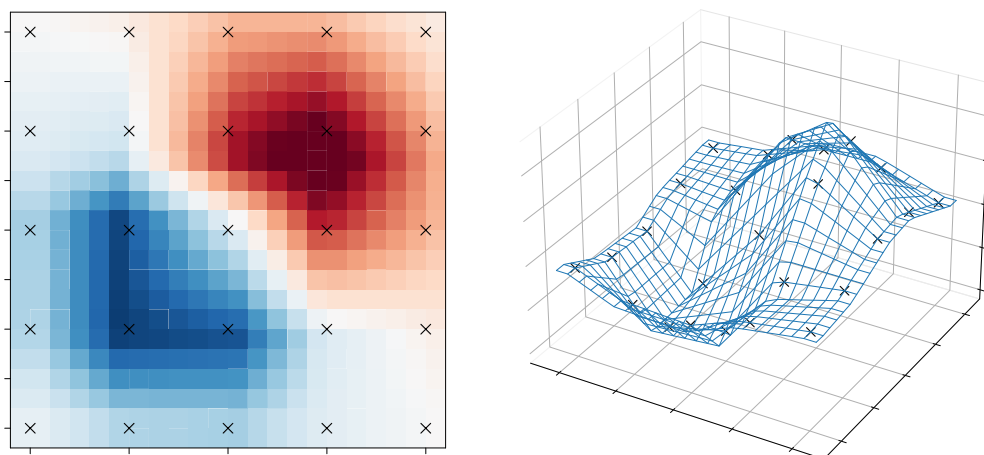


图 16. 二维插值，规则网格，线性插值

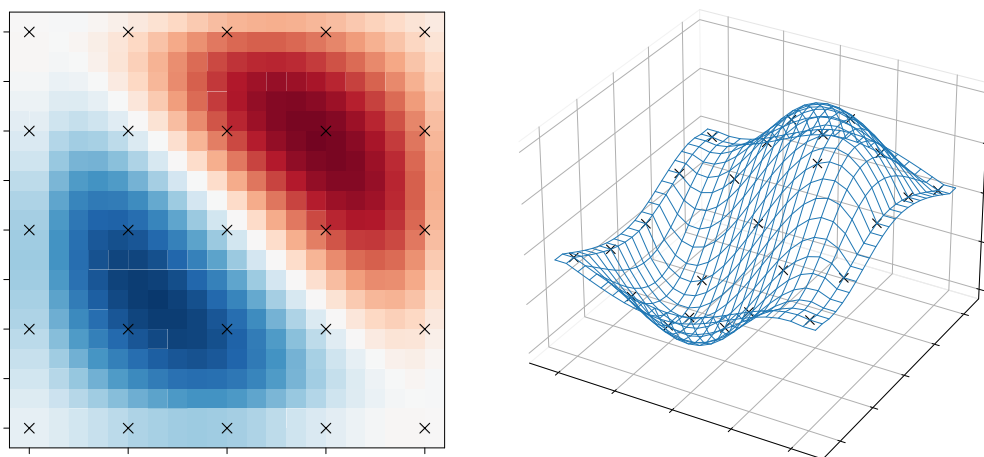


图 17. 二维插值，规则网格，三次样条



Bk6\_Ch05\_03.py 完成二维插值，并绘制图 16 和图 17。

## 不规则散点

大家可能已经注意到，图 15 给定的已知数据是规整的网格数据。当数据并不是规整的网格数据，而是不规则的散点时，我们也可以用 `scipy.interpolate.griddata()` 完成二维插值。图 18、图 19、图 20 分别所示为利用最邻近、线性、三次样条方法完成不规则散点的二维插值。

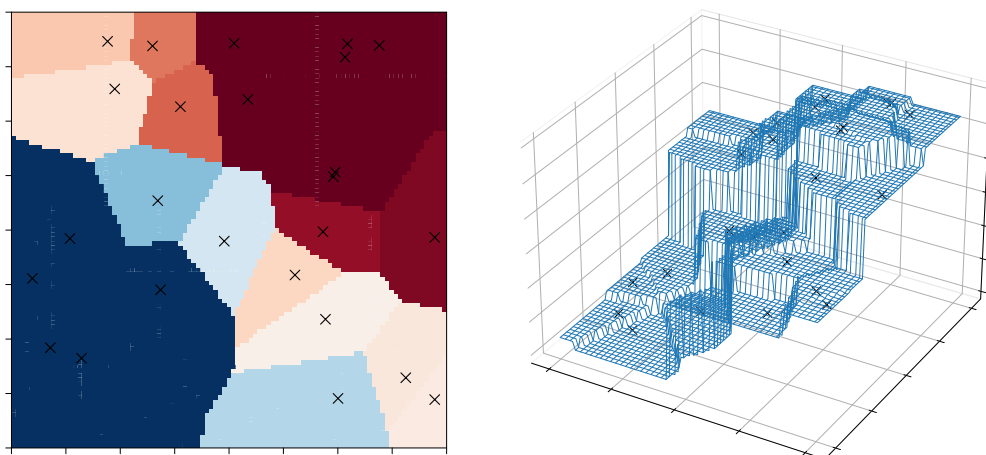


图 18. 二维插值，不规则散点，最近邻

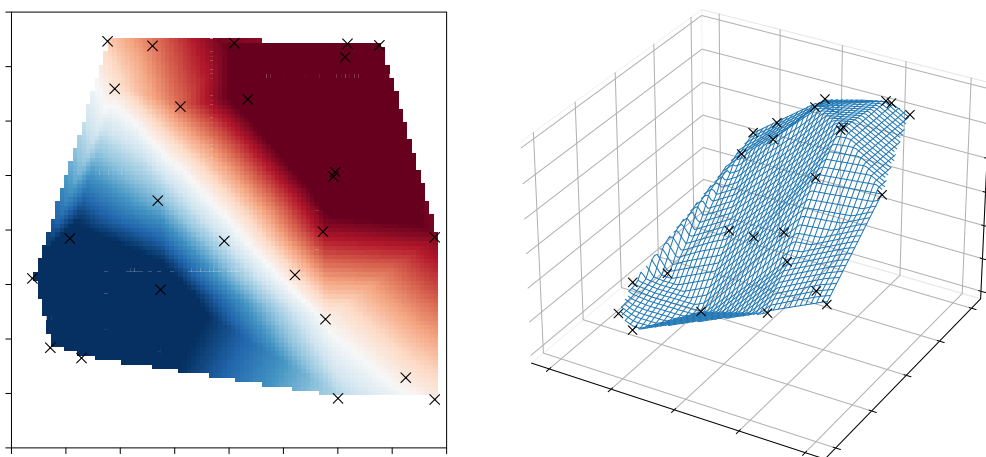


图 19. 二维插值，不规则散点，线性插值

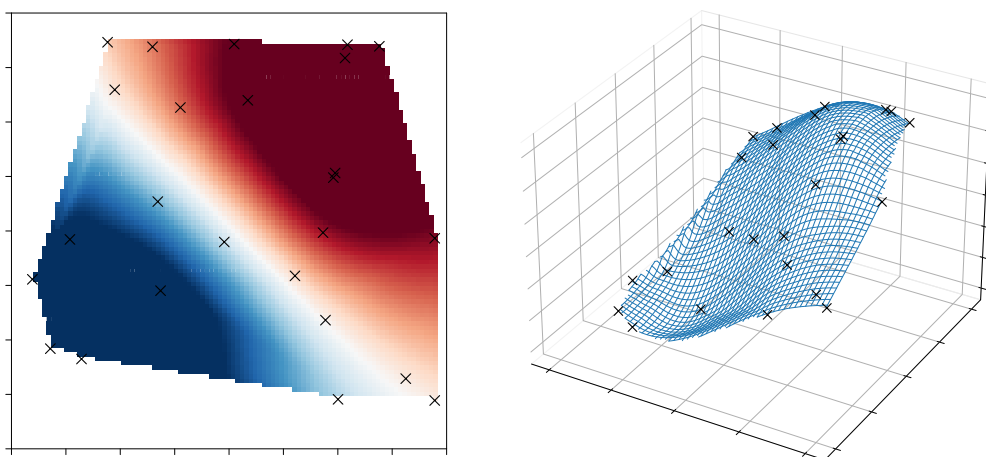
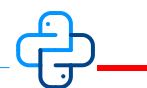


图 20. 二维插值，不规则散点，三次样条插值



Bk6\_Ch05\_04.py 完成不规则散点插值，并绘制图 18、图 19、图 20。

## 更多插值方法

matplotlib.pyplot.imshow() 绘图函数自带大量二维插值方法，请大家参考图 21。

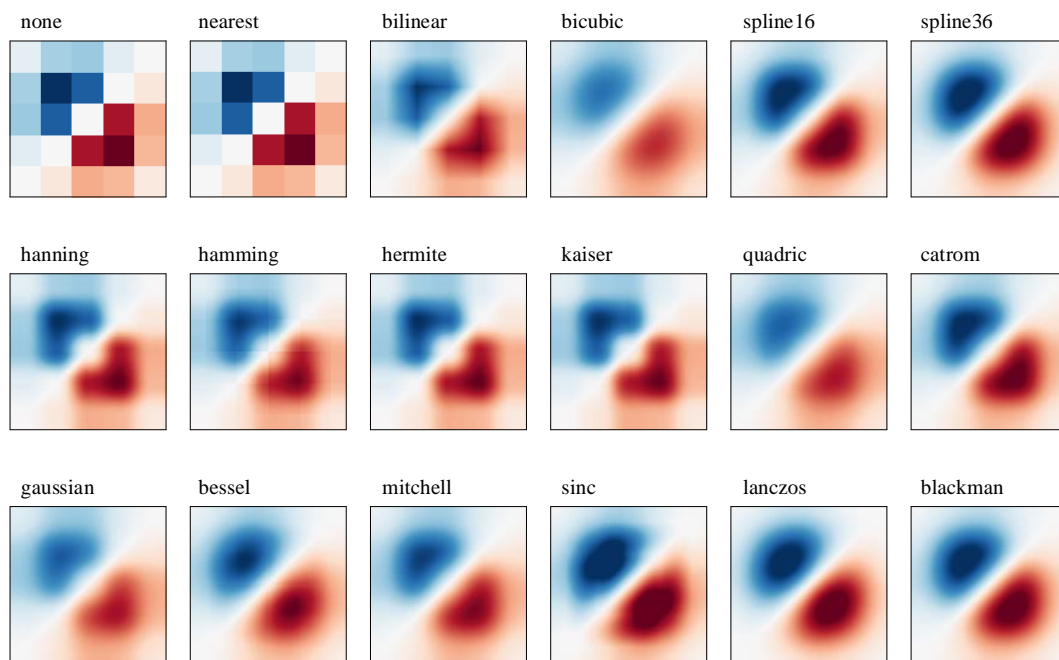


图 21. imshow() 函数插值方法



Bk6\_Ch05\_05.py 绘制图 21。



插值是一种通过已知数据点推断出连续函数在其他位置上取值的方法。在实际问题中，我们常常只知道一些离散的数据点，但需要通过这些数据点来推断出函数在其他位置的取值。插值可以通过拟合一条曲线、平面或者高维曲面来达到这个目的，从而实现对函数的估计。在机器学习中，插值可以用于对数据进行处理和预处理。插值可以通过拟合一条平滑的曲线或者曲面来填充数据中的缺失值，从而获得完整的数据集，这可以提高模型的准确性和可靠性。



请大家格外注意，插值和回归都是处理数据的方法，但插值是通过已知的数据点之间的值来估计未知点的值，而回归是通过已知的数据点来拟合一个函数，预测未知点的值。插值的目的是将数据点之间的缺失值或噪声进行平滑处理，而回归的目的是对数据进行预测和建模。虽然两者都是通过已知数据点来估计未知点的值，但它们的目的是使用场景是不同的。