

How to change to normal distribution

Gaussian و Gaussian-Like

قد تكون هناك مناسبات عندما تعمل بتوزيع غير غاوسي ، لكنك ترغب في استخدام الأساليب الإحصائية البارامترية بدلاً من الأساليب اللامعلمية.

على سبيل المثال ، قد يكون لديك عينة بيانات لها شكل الجرس المألوف ، مما يعني أنها تبدو هذا يشير إلى .غاوسية ، لكنها فشلت في واحد أو أكثر من اختبارات الحالة الطبيعية الإحصائية قد تفضل استخدام الإحصائيات البارامترية في هذه الحالة . Gaussian. أن البيانات قد تكون مثل أفضل ولأن البيانات من الواضح أنها غاوسية ، أو يمكن أن تكون ، بعد [لقوة إحصائية](#) نظراً التحويل الصحيح للبيانات

في هذا المنشور ، هناك العديد من الأسباب التي قد تجعل مجموعة البيانات غير غاوسية تقنياً سنلقي نظرة على بعض الأساليب البسيطة التي قد تكون قادراً على استخدامها لتحويل عينة Gaussian. إلى توزيع Gaussian بيانات بتوزيع يشبه

قد تكون هناك حاجة إلى بعض التجريب والحكم .لا يوجد حل سحري لهذه العملية

حجم العينة

هو أن حجم عينة البيانات صغير جداً Gaussian أحد الأسباب الشائعة لكون عينة البيانات غير

قد .ومن ثم ، الحد الأدنى .تم تطوير العديد من الأساليب الإحصائية حيث كانت البيانات شحيحة يصل عدد العينات للعديد من الطرق إلى 20 أو 30 ملاحظة

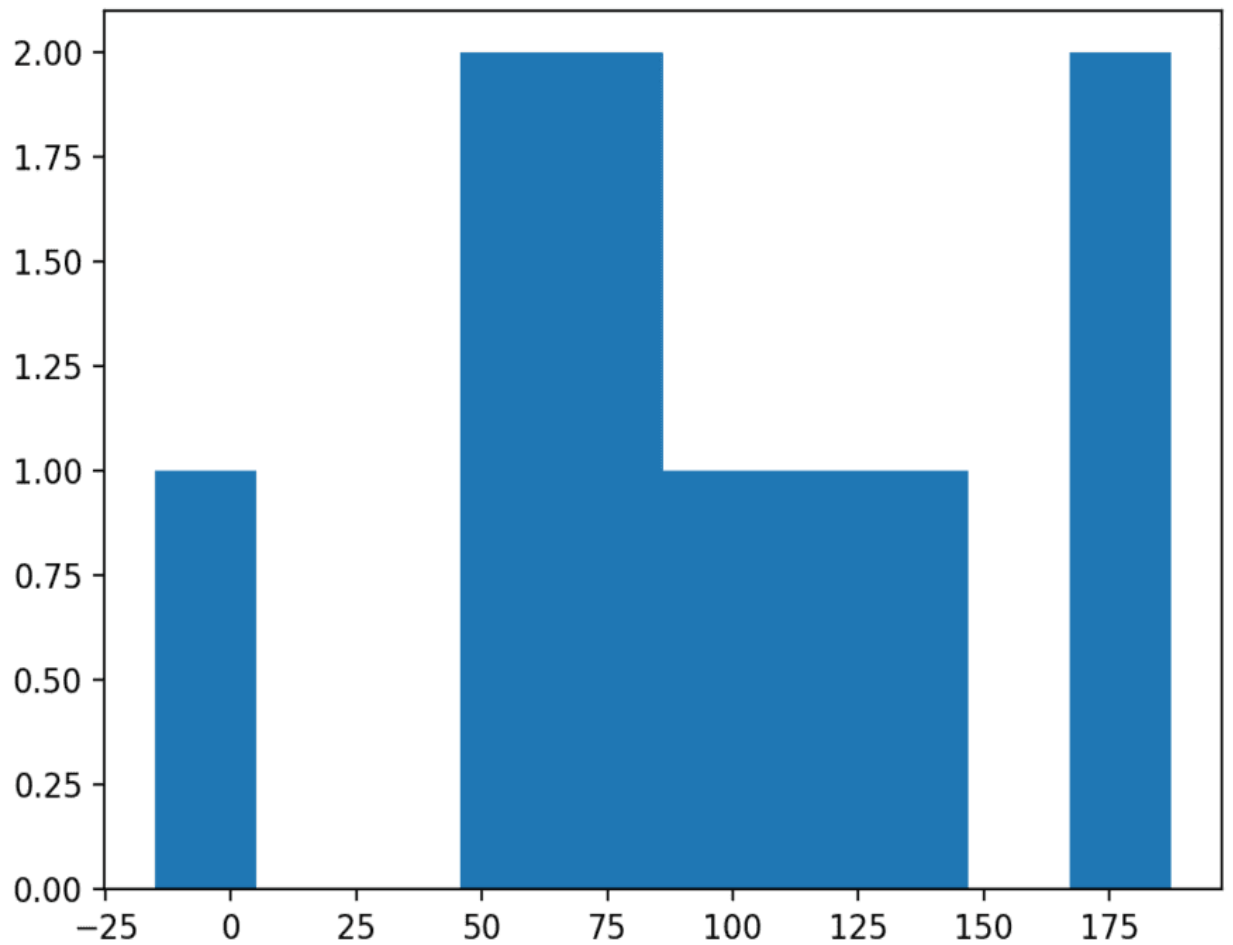
ومع ذلك ، نظراً للوضاء في بياناتك ، فقد لا ترى شكل الجرس المألوف أو تفشل في اختبارات الحالة الطبيعية مع عدد متواضع من العينات ، مثل 50 أو 100. إذا كان هذا هو الحال ، فربما بفضل قانون الأعداد الكبيرة ، كلما زادت البيانات التي تجمعها ، يمكنك جمع المزيد من البيانات زادت احتمالية استخدام بياناتك لوصف التوزيع الأساسي للسكان

لجعل هذا ملموساً ، فيما يلي مثال على مخطط لعينة صغيرة من 50 ملاحظة مأخوذة من توزيع غاوسي بمتوسط 100 وانحراف معياري قدره 50

- 1 # histogram plot of a small sample
- 2 from numpy.random import seed

```
3         from numpy.random import randn
4         from matplotlib import pyplot
5         # seed the random number generator
6         seed(1)
7         # generate a univariate data sample
8         data = 50 * randn(50) + 100
9         # histogram
10        pyplot.hist(data)
11        pyplot.show()
```

يؤدي تشغيل المثال إلى إنشاء رسم بياني للبيانات يظهر عدم وجود توزيع غاوسي واضح ، ولا حتى توزيع غاوسي.



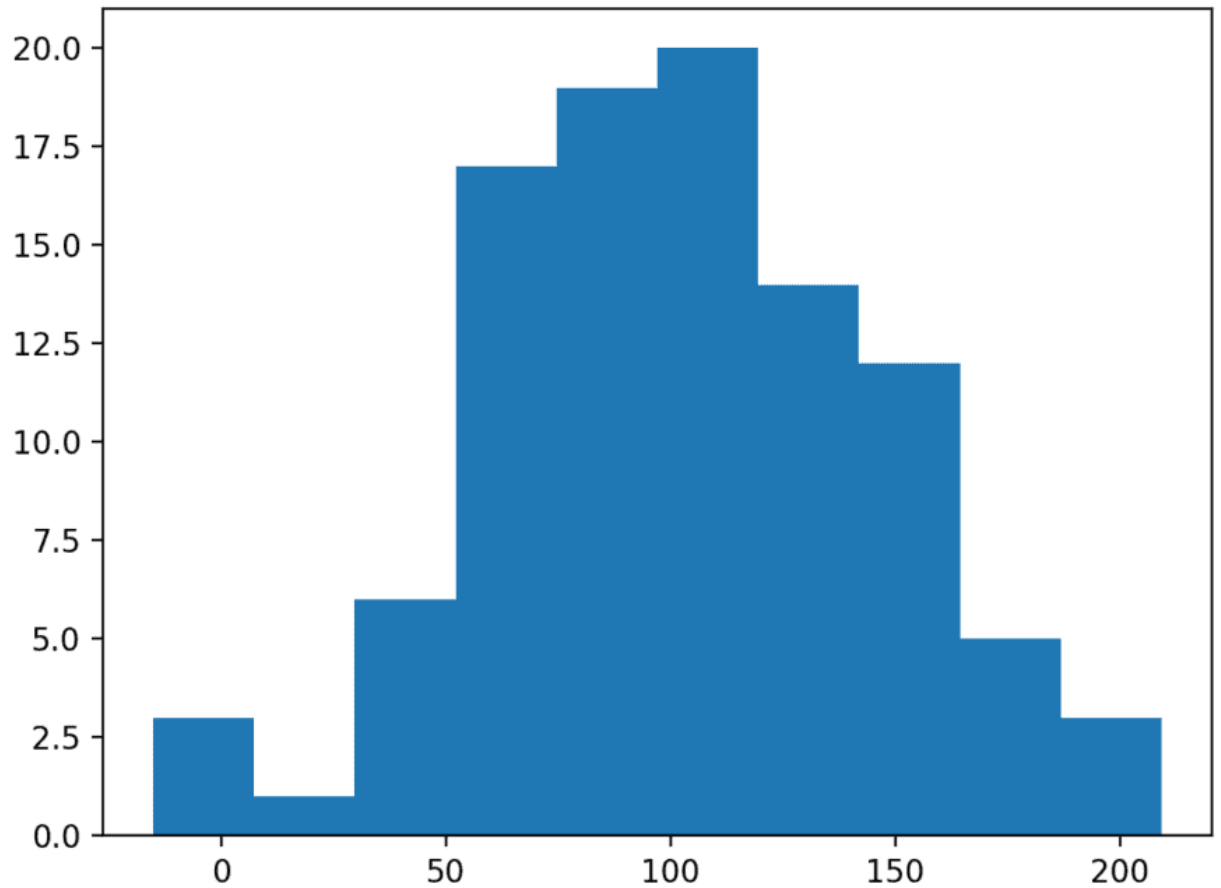
رسم بياني لعينة بيانات صغيرة جدًا

يمكن أن تساعد زيادة حجم العينة من 50 إلى 100 في عرض الشكل الغاوسي لتوزيع البيانات بشكل أفضل.

- 1 # histogram plot of a small sample
- 2 from numpy.random import seed
- 3 from numpy.random import randn
- 4 from matplotlib import pyplot

```
5 # seed the random number generator
6 seed(1)
7 # generate a univariate data sample
8 data = 50 * randn(100) + 100
9 # histogram
10 pyplot.hist(data)
11 pyplot.show()
```

باستخدام هذا المثال ، يمكننا أن نرى بشكل أفضل التوزيع الغاوسي للبيانات التي ستجتاز كلاً من الاختبارات الإحصائية وفحوصات كرة العين.



رسم بياني لعينة بيانات أكبر

دقة البيانات

ربما تتوقع توزيعًا غاوسيًا من البيانات ، ولكن بغض النظر عن حجم العينة التي تجمعها ، فإنها لا تتحقق.

قد يتم حجب توزيع البيانات من السبب الشائع لذلك هو الدقة التي تستخدمها لتجميع الملاحظات قد يكون هناك العديد من الأسباب لتعديل دقة . خلال الدقة المختارة للبيانات أو دقة الملاحظات: البيانات قبل النمذجة ، مثل

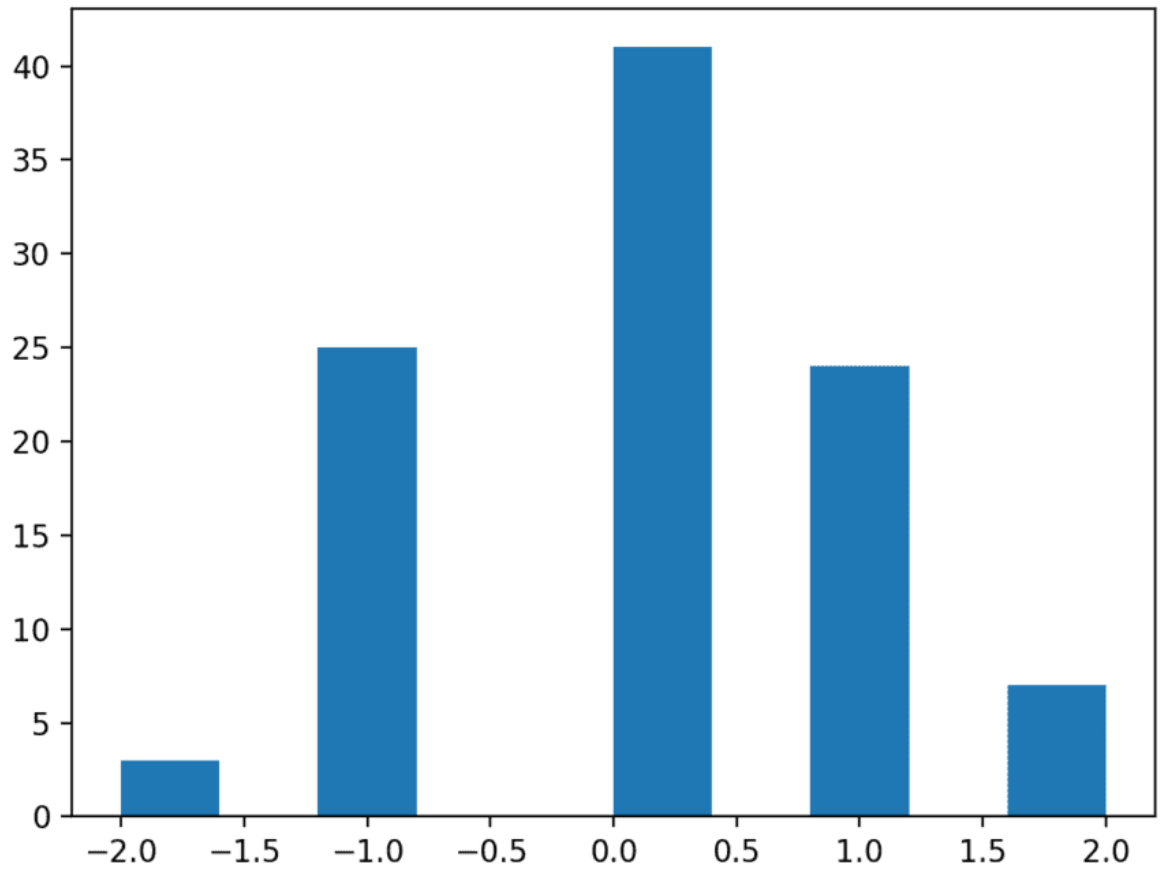
- تكوين آلية عمل الملاحظة .
- البيانات تمر عبر عملية مراقبة الجودة .

- دقة قاعدة البيانات المستخدمة لتخزين البيانات

بمتوسط 0 وانحراف رقم غاوسي عشوائي لجعل هذا الأمر ملموسًا ، يمكننا عمل عينة من 100 معياري قدره 1 وإزالة جميع المنازل العشرية

```
1 # histogram plot of a low res sample
2
3 from numpy.random import seed
4
5 from numpy.random import randn
6
7 from matplotlib import pyplot
8
9 # seed the random number generator
10 seed(1)
11
12 # generate a univariate data sample
13 data = randn(100)
14
15 # remove decimal component
16 data = data.round(0)
17
18 # histogram
19 pyplot.hist(data)
20
21 pyplot.show()
```

ستؤدي إضافة Gaussian. تشغيل المثال ينتج عنه توزيع يبدو منفصلاً على الرغم من أنه يشبه الدقة إلى الملاحظات إلى توزيع أكثر اكتمالاً للبيانات



رسم بياني لعينة بيانات منخفضة الدقة

القيم المتطرفة

قد يكون لعينة البيانات توزيع غاوسي ، ولكن قد يتم تشويهها لعدد من الأسباب

يمكن أن توجد القيم المتطرفة لعدد من الأسباب الشائع هو وجود قيم متطرفة عند حافة التوزيع :الأسباب ، مثل

- خطأ في القياس .
- بيانات مفقودة .
- تلف البيانات .
- أحداث نادرة .

في مثل هذه الحالات ، يمكن تحديد القيم القصوى وإزالتها من أجل جعل التوزيع أكثر
غالبًا ما تسمى هذه القيم المتطرفة القيم المتطرفة .غاوسيًا

قد يتطلب ذلك خبرة في المجال أو استشارة أحد خبراء المجال من أجل تصميم معايير لتحديد القيم
المتطرفة ثم إزالتها من عينة البيانات وجميع عينات البيانات التي تتوقع أنت أو نموذجك العمل
معها في المستقبل

يمكننا توضيح مدى سهولة وجود قيم متطرفة تعطل توزيع البيانات

يُنشئ المثال أدناه عينة بيانات تحتوي على 100 رقم غاوسي عشوائي تم قياسه بمتوسط 10
وانحراف معياري قدره 5. ثم تتم إضافة 10 ملاحظات إضافية ذات قيمة صفرية إلى
هذا سلوك شائع .يمكن أن يحدث هذا إذا تم تعيين القيم المفقودة أو التالفة بقيمة صفر. التوزيع
فمثلا في مجموعات بيانات التعلم الآلي المتاحة للجمهور ؛

```
1 # histogram plot of data with outliers
2
3 from numpy.random import seed
4
5 from numpy.random import randn
6
7 from numpy import zeros
8
9 from numpy import append
10
11 from matplotlib import pyplot
12
13 # seed the random number generator
14
15 seed(1)
16
17 # generate a univariate data sample
18
19 data = 5 * randn(100) + 10
20
21 # add extreme values
22
23 data = append(data, zeros(10))
24
25 # histogram
```

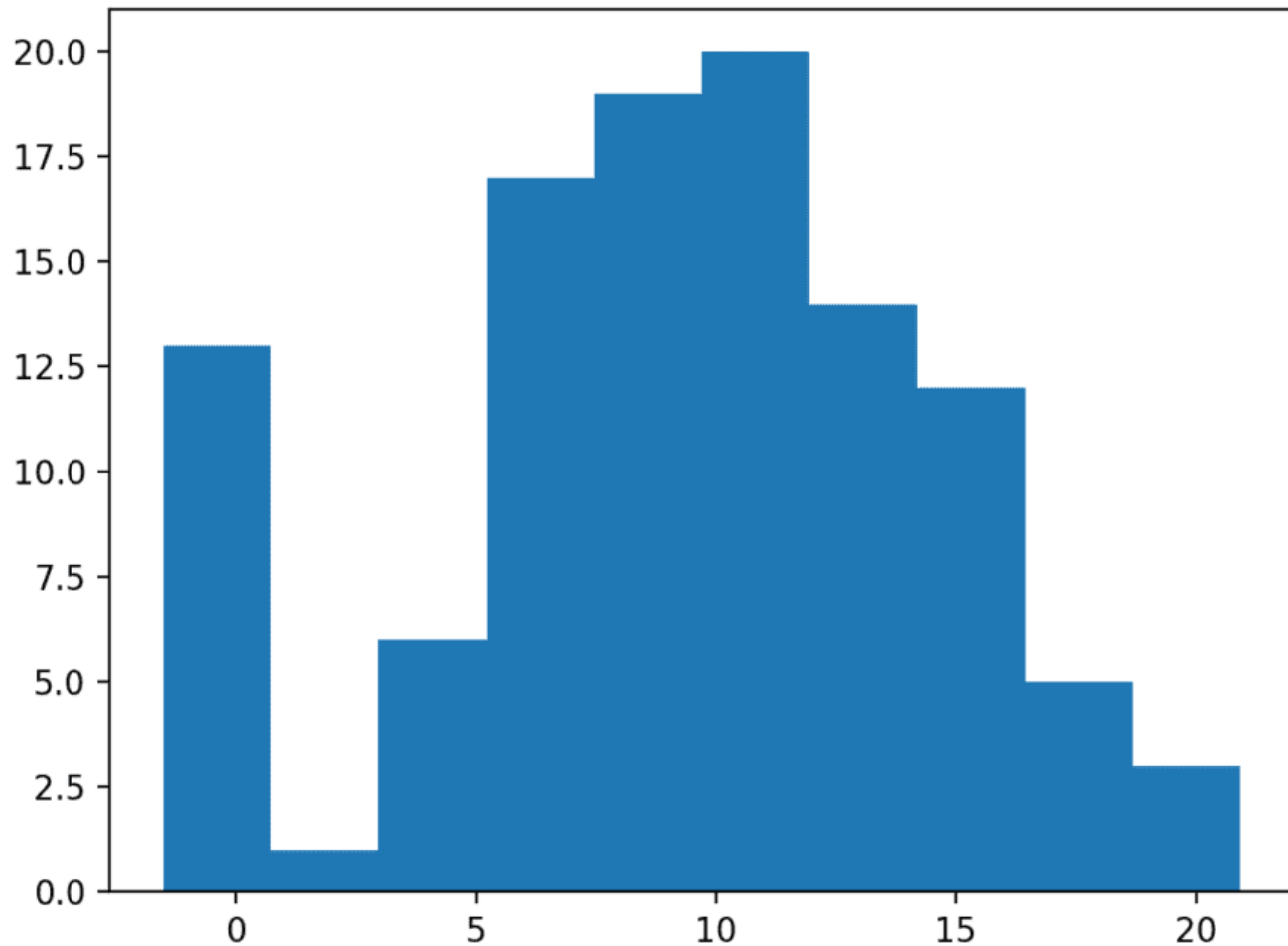

14

`pyplot.hist(data)`

15

`pyplot.show()`

يمكنك أن ترى بوضوح كيف أن التردد العالي يؤدي تشغيل المثال إلى إنشاء ورسم عينة البيانات غير المتوقع للملاحظات ذات القيمة الصفرية يعطل التوزيع.



رسم بياني لعينة البيانات ذات القيم القصوى

ذبول طويلة

بالإضافة إلى وفرة الأحداث النادرة على حافة التوزيع ، يمكن أن تظهر القيم المتطرفة بعدة طرق. قد ترى ذيلًا طويلًا للتوزيع في أحد الاتجاهين أو كلاهما

في المؤامرات ، يمكن أن يجعل هذا التوزيع يبدو وكأنه أسي ، بينما في الواقع قد يكون غاوسيًا مع وفرة من الأحداث النادرة في اتجاه واحد

يمكنك استخدام قيم حد بسيطة ، ربما بناءً على عدد الانحرافات المعيارية عن المتوسط ، لتحديد قيم الذيل الطويل وإزالتها

يمكننا إثبات ذلك بمثال مفتعل. تحتوي عينة البيانات على 100 رقم عشوائي غاوسي بمتوسط 10 وانحراف معياري 5. تمت إضافة 50 قيمة عشوائية منتظمة في النطاق من 10 إلى 110. هذا يخلق ذيل طويل على التوزيع

```
1 # histogram plot of data with a long tail
2
3 from numpy.random import seed
4
5 from numpy.random import randn
6
7 from numpy.random import rand
8
9 from numpy import append
10
11 from matplotlib import pyplot
12
13 # seed the random number generator
14
15 seed(1)
16
17 # generate a univariate data sample
18
19 data = 5 * randn(100) + 10
20
21 tail = 10 + (rand(50) * 100)
22
23 # add long tail
24
25 data = append(data, tail)
26
27 # histogram
```

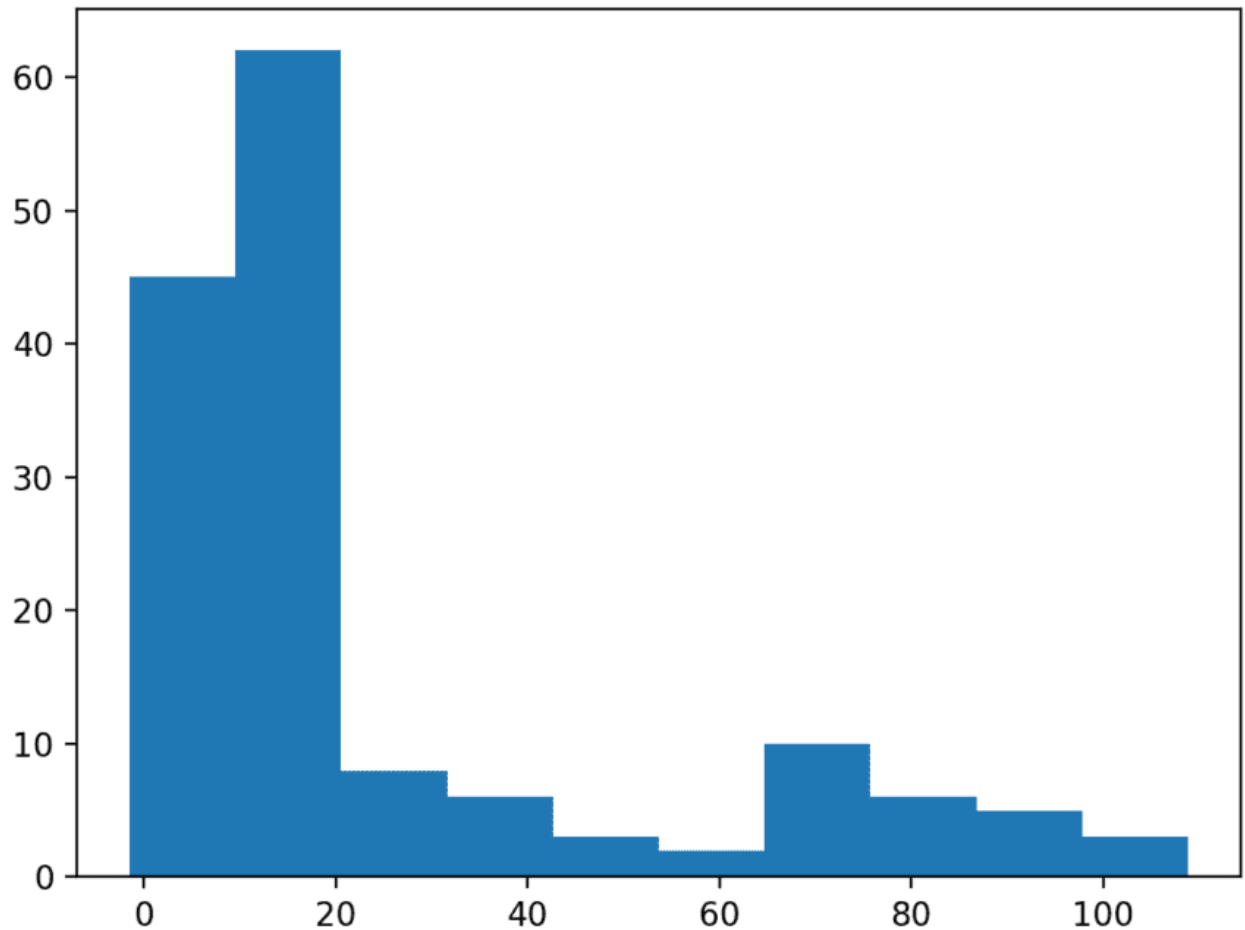
15

`pyplot.hist(data)`

16

`pyplot.show()`

عند تشغيل المثال ، يمكنك أن ترى كيف يشوه الذيل الطويل التوزيع الغاوسي ويجعله يبدو أسياً (تقريباً أو ربما حتى ثنائي النسق (نتوءان).



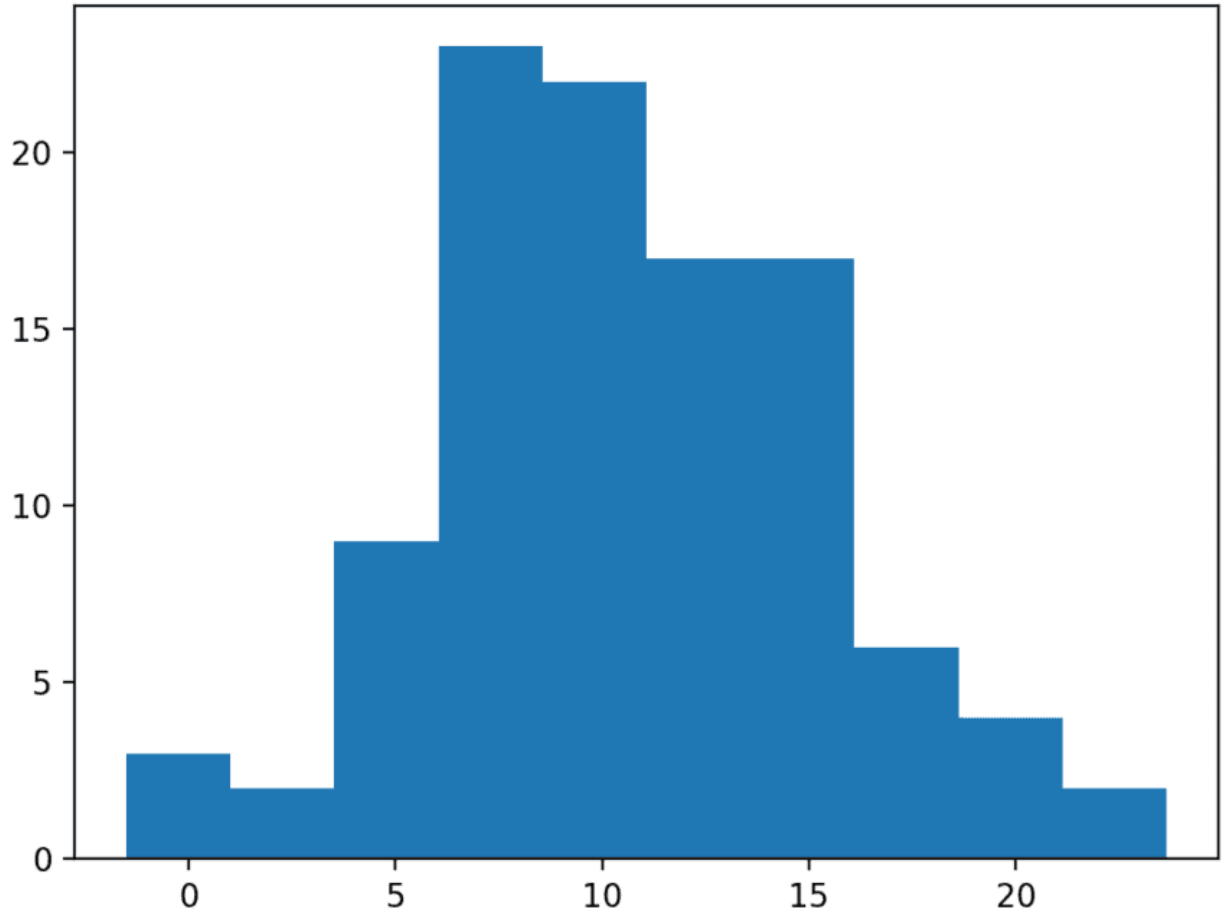
رسم بياني لعينة البيانات ذات الذيل الطويل

يمكننا استخدام حد بسيط ، مثل القيمة 25 ، في مجموعة البيانات هذه كقطع وإزالة جميع
لقد اخترنا هذا الحد بمعرفة مسبقة عن كيفية إنشاء عينة .الملاحظات الأعلى من هذا الحد

البيانات ، ولكن يمكنك تخيل اختبار عتبات مختلفة على مجموعة البيانات الخاصة بك وتقييم تأثيرها.

```
1 # histogram plot of data with a long tail
2
3 from numpy.random import seed
4
5 from numpy.random import randn
6
7 from numpy.random import rand
8
9 from numpy import append
10
11 from matplotlib import pyplot
12
13 # seed the random number generator
14
15 seed(1)
16
17 # generate a univariate data sample
18
19 data = 5 * randn(100) + 10
20
21 tail = 10 + (rand(10) * 100)
22
23 # add long tail
24
25 data = append(data, tail)
26
27 # trim values
28
29 data = [x for x in data if x < 25]
30
31 # histogram
32
33 pyplot.hist(data)
34
35 pyplot.show()
```

يُظهر تشغيل الكود كيف يؤدي هذا القص البسيط للذيل الطويل إلى إرجاع البيانات إلى توزيع غاوسي.



رسم بياني لعينة البيانات مع ذيل طويل مبتور

تحويلات الطاقة

قد يكون توزيع البيانات طبيعيًا ، ولكن قد تتطلب البيانات تحويلًا للمساعدة في كشفها

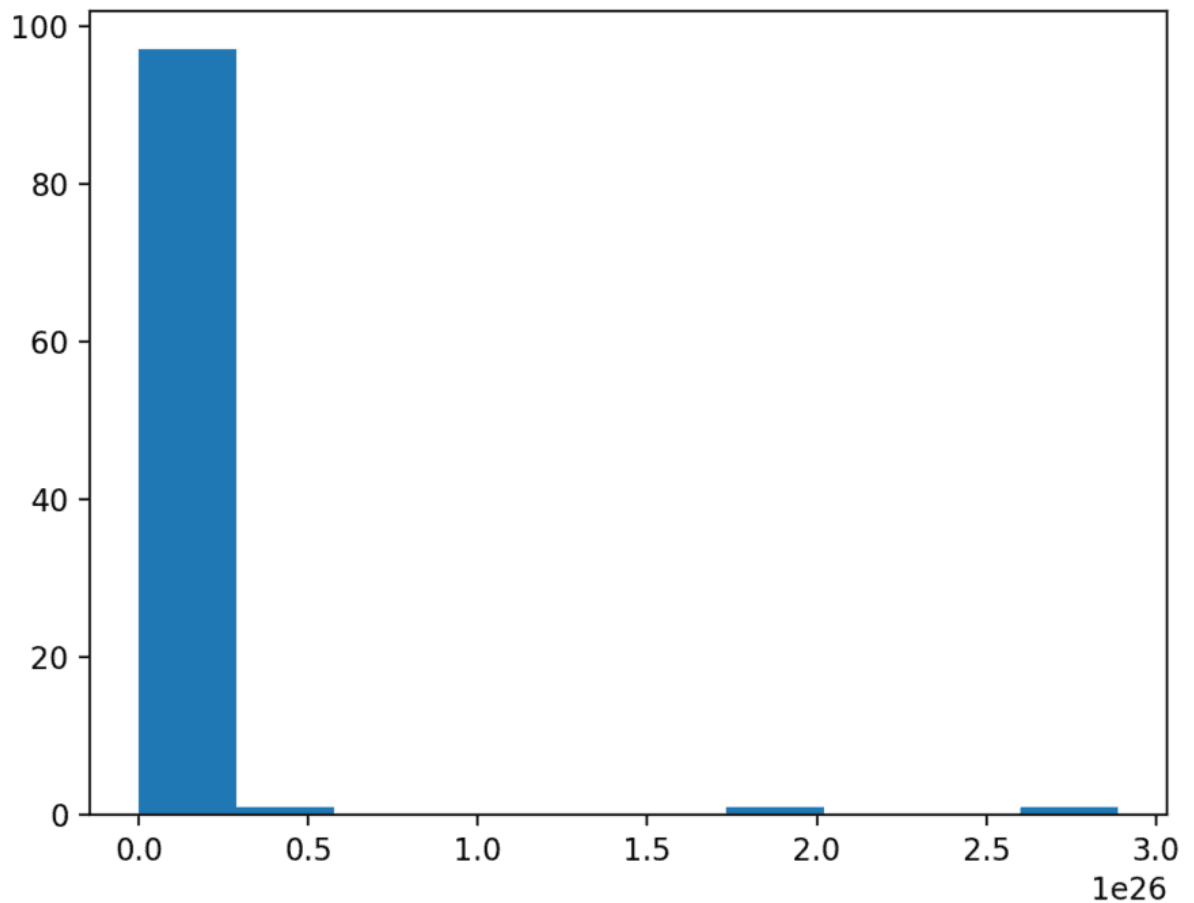
على سبيل المثال ، قد تحتوي البيانات على انحراف ، مما يعني أن الجرس في شكل الجرس قد في بعض الحالات ، يمكن تصحيح ذلك عن طريق تحويل البيانات عن . يتم دفعه بطريقة أو بأخرى طريق حساب الجذر التربيعي للملاحظات

بالتناوب ، قد يكون التوزيع أسياً ، ولكن قد يبدو طبيعيًا إذا تم تحويل الملاحظات عن طريق أخذ تسمى البيانات مع هذا التوزيع السجل العادي . اللوغاريتم الطبيعي للقيم

لجعل هذا ملموسًا ، يوجد أدناه مثال لعينة من الأرقام الغوسية التي تم تحويلها ليكون لها توزيع أسّي.

```
1 # log-normal distribution
2 from numpy.random import seed
3 from numpy.random import randn
4 from numpy import exp
5 from matplotlib import pyplot
6 # seed the random number generator
7 seed(1)
8 # generate two sets of univariate observations
9 data = 5 * randn(100) + 50
10 # transform to be exponential
11 data = exp(data)
12 # histogram
13 pyplot.hist(data)
14 pyplot.show()
```

ليس من الواضح أن البيانات .يؤدي تشغيل المثال إلى إنشاء رسم بياني يوضح التوزيع الأسّي هي في الواقع لوغاريتمي عادي



الرسم البياني للتوزيع العادي للسجل

أخذ الجذر التربيعي ولو غار يتم الملاحظة من أجل جعل التوزيع طبيعيًا ينتمي إلى فئة من هي طريقة لتحويل البيانات قادرة على أداء Box-Cox طريقة. التحويلات تسمى تحويلات الطاقة تم تسمية هذه الطريقة باسم مجموعة من تحويلات الطاقة ، بما في ذلك السجل والجذر التربيعي جورج بوكس وديفيد كوكس.

يمكن أكثر من ذلك ، يمكن تهيئته لتقييم مجموعة من التحويلات تلقائيًا واختيار أفضل ملائمة قد تكون عينة . اعتبره أداة قوية لتسوية التغيير المعتمد على الطاقة في عينة البيانات الخاصة بك البيانات الناتجة أكثر خطية وستمثل بشكل أفضل التوزيع الأساسي غير للطاقة ، بما في ذلك Gaussian.

، تتحكم lambda يتطلب الأمر حجة ، تسمى Box-Cox طريقة [SciPy \(\) boxcox](#) وظيفة تنفذ في نوع التحويل المطلوب إجراؤه.

lambda: فيما يلي بعض القيم الشائعة لـ

- هو تحويل متبادل . لامدا = -1
- تحويل جذري تربيعي متبادل لامدا = -0.5
- تحويل سجل لامدا = 0.0
- تحويل جذر تربيعي لامدا = 0.5
- لا يوجد تحويل لامدا = 1.0

لإجراء Box-Cox على سبيل المثال ، نظرًا لأننا نعلم أن البيانات غير طبيعية ، يمكننا استخدام صراحةً على $\lambda = 0$ تحويل السجل عن طريق تعيين

```
1 # power transform
2 data = boxcox(data, 0)
```

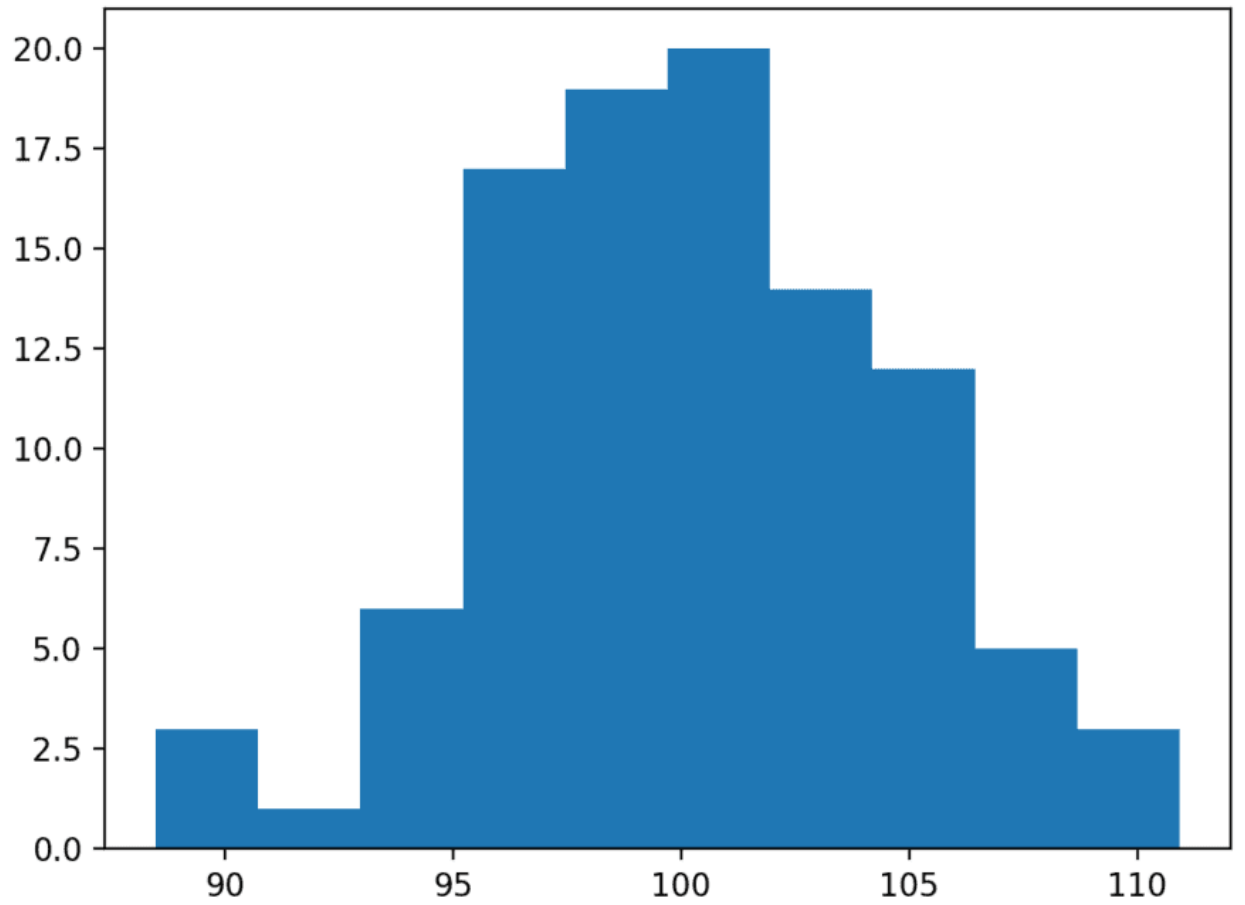
على عينة البيانات الأسية مُدرج أدناه Box-Cox المثال الكامل لتطبيق تحويل

```
1 # box-cox transform
2 from numpy.random import seed
3 from numpy.random import randn
4 from numpy import exp
5 from scipy.stats import boxcox
6 from matplotlib import pyplot
7 # seed the random number generator
8 seed(1)
9 # generate two sets of univariate observations
10 data = 5 * randn(100) + 100
```



```
11 # transform to be exponential
12 data = exp(data)
13 # power transform
14 data = boxcox(data, 0)
15 # histogram
16 pyplot.hist(data)
17 pyplot.show()
```

على عينة البيانات ورسم النتيجة ، مع إظهار Box-Cox يؤدي تشغيل المثال إلى إجراء تحويل
التوزيع الغاوسي بوضوح



مخطط الرسم البياني لمربع كوكس عينة البيانات الأسية المحولة

في أنه يفترض أن جميع القيم في عينة البيانات موجبة Box-Cox يتمثل أحد قيود تحويل

Yeo-Johnson. طريقة بديلة لا تجعل هذا الافتراض هو تحول