

Cauchy Distribution

The Cauchy distribution has PDF given by:

$$f(x) = \frac{1}{\pi} \frac{1}{1+x^2}$$

for $x \in (-\infty, \infty)$.

Note that $\int \frac{1}{1+x^2} dx = \arctan(x)$.

The Cauchy distribution has infinite mean and variance.

$$E(|x|) = \frac{2}{\pi} \int_0^\infty \frac{x}{1+x^2} dx = \frac{1}{\pi} \int_0^\infty \frac{1}{1+y} dy = \frac{1}{\pi} \log(1+y) \Big|_0^\infty$$

Normal Random Variables

The normal distribution behaves well under addition.

Consider $z_1 \sim N(0; v_1)$ and $z_2 \sim N(0; v_2)$ are two independent random variables, and scalar a . Then:

$$\begin{aligned} z_1 + z_2 &\sim N(0; v_1 + v_2) \\ az_1 &\sim N(0; a^2 v_1) \end{aligned}$$

Or more generally, for fixed but arbitrary real numbers r_1, \dots, r_k and random variables z_1, \dots, z_k all iid $N(0; 1)$ we have:

$$\begin{aligned} r_1 z_1 + \dots + r_k z_k &=_D N(0, r_1^2 + \dots + r_k^2) \\ &=_D \sqrt{r_1^2 + \dots + r_k^2} N(0; 1) \end{aligned}$$

So

$$\begin{aligned} \left(\sum r_i z_i\right)^2 &=_D \left(N(0; \sum r_i^2)\right)^2 \\ &=_D \left(\sum r_i^2\right) (N(0; 1))^2 \end{aligned}$$

P-Stability

A distribution f is p-stable if for any $K > 0$, any real numbers r_1, \dots, r_K , and given K iid random variables z_1, \dots, z_K following distribution f , then

$$\sum_{i=1}^K r_i z_i =_D \left(\sum_{i=1}^K |r_i|^p\right)^{1/p} Y$$

where Y has distribution f .

Notes:

- For any $p \in (0, 2]$ there exists some p -stable distribution.
- The Cauchy distribution is 1-stable.
- The Normal distribution is 2-stable.
- The CLT suggests that no other distribution is 2-stable

F2 Estimation

$$F_2(t) = \sum_{a \in U} (f_t(a))^2$$

This looks similar to computing a variance.

Define the consistent normal random variable $h_i(a) \sim N(0; 1)$ such that $h_i(a)$ and $h_j(b)$ are independent if $i \neq j$ or $a \neq b$.

In practice:

```
h(i,a) = { rseed(<i,a>)
           return rand_normal()
           }
```

Then for a multiset S (a set where the same element may occur multiple times) we define the following sketch:

$$\sigma = \langle \sum_{a \in S} h_1(a), \dots, \sum_{a \in S} h_J(a) \rangle$$

With τ being coordinate wise addition:

$$\tau(\sigma(S_1), \sigma(S_2)) = \langle \dots, \sum_{a \in S_1} h_i(a) + \sum_{a \in S_2} h_i(a), \dots \rangle$$

Now given $\sigma = \langle \sigma_1, \dots, \sigma_J \rangle$ we wish to estimate F_2 . Notice that:

$$\begin{aligned} \sigma_i &= \sum_{a \in S} h_i(a) = \sum_{a \in U} h_i(a) f_s(a) \\ \sigma_i &=_D \left(\sum_{a \in U} (f_s(a))^2 \right)^{1/2} N(0; 1) \end{aligned}$$

Then:

$$\sigma_i^2 =_D \sum_{a \in U} (f_s(a))^2 (N(0; 1))^2$$

So

$$\begin{aligned} \text{mean}(\sigma^2) &= F_2 \\ \text{median}(\sigma^2) &= 0.4705 F_2 \end{aligned}$$

Where the number 0.4705 arises as the median of a chi-squared distribution with 1 degree of freedom.

F2 Estimation in Stream

Suppose at time t you have a sketch $\langle \sigma_1, \dots, \sigma_J \rangle$. You then wish to update once a_t arrives according to:

$$\langle \sigma_1 + h_1(a_t), \dots, \sigma_J + h_J(a_t) \rangle$$

Suppose now your data is arranged in key-value pairs $\langle a_t, v_t \rangle$ such that $f_t(a) = \sum_{t:a=a_t} v_t$. Then the update becomes:

$$\langle \sigma_1 + v_t h_1(a_t), \dots, \sigma_J + v_t h_J(a_t) \rangle$$

F1 Estimation

In the simple case $F_1(t) = t$.

However, when the data is arranged in key-value pairs $\langle a_t, v_t \rangle$ where the v_t may take on negative values, computing $F_1(t)$ is no longer trivial.

Consider the following example:

$$\langle 1, -3 \rangle, \langle 1, 7 \rangle, \langle 2, 1 \rangle, \langle 3, -1 \rangle, \langle 2, 2 \rangle$$

Produces

$$f_t(1) = 4 \quad f_t(2) = 3 \quad f_t(3) = -1$$

This has $F_1(t) = \sum_i |f_t(i)| = 8$.

To estimate $F_1(t)$ we follow the same technique as for $F_2(t)$ but replacing the consistent normal random variable by a consistent Cauchy random variable. So now let $h_i(a) \sim \text{Cauchy}$. Then:

$$\sigma_1 = \sum_{a \in U} h_1(a) f_t(a) =_D F_1(t) C$$

So

$$\text{median}(|\sigma_1|, \dots, |\sigma_J|) =_D F_1(t) \text{median}(|C_1|, \dots, |C_J|)$$

Where C and C_i are Cauchy distributed random variables.

$|C_i|$ has PDF $\frac{2}{\pi} \frac{1}{1+x^2}$ for $x \in [0, \infty)$. The median can be calculated by solving:

$$\frac{1}{2} = \frac{2}{\pi} \int_0^x \frac{1}{1+t^2} dt$$

This gives solution $x = \tan(\frac{\pi}{4}) = 1$.

Therefore the median of σ is a good estimator for $F_1(t)$.

Machine Learning Application

Suppose you have N points in \mathbb{R}^D where N and D are large, and you want to produce a summary or reduction of these points while preserving distance: $d(x, y) = \|X - Y\|_2$.

1. Pick D random variables z_1, \dots, z_D iid for $N(0; 1)$. Fix these values.

2. Map $X = \langle x_1, \dots, x_D \rangle$ to $\sum_{i=1}^D x_i z_i$. The resulting value will follow a normal distribution multiplied by a constant.

Then for two points X and Y we have:

$$\sum_{i=1}^D x_i z_i - \sum_{i=1}^D y_i z_i = \sum_{i=1}^D (x_i - y_i) z_i \sim N(0; 1) \|X - Y\|_2$$

So the L2-norm is (approximately) preserved.