

Míry polohy a variability

Dominik Heger

Masaryk University

hegerd@chemi.muni.cz

STDT 03 Míry polohy a variability

K čemu sumarizovat data? Není nejlepší mít všechna data?

K čemu sumarizovat data? Není nejlepší mít všechna data?

```
8 2 0 3 1 4 5 8 2 1 7 2 7 3 8 5 5 2 9 0 6 3 1 6 4
0 8 7 3 3 1 9 7 5 2 5 7 6 9 8 0 0 3 6 2 5 1 2 7 5 2
2 3 3 8 6 6 1 4 2 4 0 2 6 1 8 9 5 2 6 9 8 3 4 0 1 0
4 7 5 5 6 3 0 7 7 1 9 1 6 1 7 4 1 7 1 3 7 9 3 3 7
1 9 3 9 5 3 4 9 5 5 2 7 5 8 0 3 4 8 8 1 2 7 5 3 4
2 8 7 8 1 4 1 4 9 4 2 4 1 5 2 2 9 4 6 2 1 5 2 8 1 9
8 4 8 5 1 3 9 6 6 0 7 2 1 9 0 2 2 0 6 7 0 6 0 1 3 0
0 3 8 8 4 7 5 1 5 1 7 3 4 5 2 0 7 4 7 9 6 6 7 7 4
3 5 3 1 9 3 7 4 9 5 0 2 0 1 4 6 2 5 4 5 8 5 0 9 2
3 4 5 9 5 2 7 9 8 9 0 5 5 8 5 1 7 7 3 5 5 4 7 7 2
4 1 5 3 0 9 1 3 7 2 5 8 7 7 1 3 6 3 9 7 8 7 9 1 7
7 2 9 5 6 7 8 5 4 5 3 4 5 4 1 9 8 6 7 5 7 9 3 1 8
5 9 2 8 9 8 6 4 4 1 5 3 7 7 0 8 0 2 5 6 0 6 1 2 0
1 3 3 3 9 0 5 2 8 7 4 0 9 0 3 7 3 1 7 9 4 5 5 2 8
4 6 0 1 0 8 6 2 1 0 0 5 0 3 1 5 4 9 0 3 7 4 7 0 1
7 7 0 6 6 3 2 8 8 5 8 9 5 6 4 0 5 9 1 8 0 5 4 9 4
3 3 8 5 7 5 7 4 3 4 5 7 9 6 9 5 0 7 7 6 6 8 8 5 9
9 1 7 1 3 6 9 2 2 9 1 9 4 2 3 3 0 8 1 8 7 7 6 4 7 2
6 2 2 8 0 9 4 5 3 7 2 5 4 6 6 5 6 6 5 0 4 6 5 6 8
1 7 5 9 0 0 2 0 5 6 5 8 5 1 9 5 3 3 7 4 0 5 8 2 4
0 3 9 6 9 4 7 3 5 7 0 6 5 4 7 1 1 8 5 3 2 8 0 9 8
```

K čemu sumarizovat data? Není nejlepší mít všechna data?

```
8 2 0 3 1 4 5 8 2 1 7 2 7 3 8 5 5 2 9 0 6 3 1 6 4
0 8 7 3 3 1 9 7 5 2 5 7 6 9 8 0 3 6 2 5 1 2 7 5 2
2 3 3 8 6 6 1 4 2 4 0 2 6 1 8 9 5 2 6 9 8 3 4 0 1 0
4 7 5 5 6 3 0 7 7 1 9 1 6 1 7 4 1 7 1 3 7 9 3 3 7
1 9 3 9 5 3 4 9 5 5 2 7 5 8 0 3 4 8 8 1 2 7 5 3 4
2 8 7 8 1 4 1 4 9 4 2 4 1 5 2 9 4 6 2 1 5 2 8 1 9
8 4 8 5 1 3 9 6 6 0 7 2 1 9 0 2 0 6 7 0 6 0 1 3 0
0 3 8 8 4 7 5 1 5 1 7 3 4 5 2 0 7 4 7 9 6 6 7 7 4
3 5 3 1 9 3 7 4 9 5 0 2 0 1 4 6 2 5 4 5 8 5 0 9 2
3 4 5 9 5 2 7 9 8 9 0 5 5 8 5 1 7 7 3 5 5 4 7 7 2
4 1 5 3 0 9 1 3 7 2 5 8 7 7 1 3 6 3 9 7 8 7 9 1 7
7 2 9 5 6 7 8 5 4 5 3 4 5 4 1 9 8 6 7 5 7 9 3 1 8
5 9 2 8 9 8 6 4 4 1 5 3 7 7 0 8 0 2 5 6 0 6 1 2 0
1 3 3 3 9 0 5 2 8 7 4 0 9 0 3 7 3 1 7 9 4 5 5 2 8
4 6 0 1 0 8 6 2 1 0 0 5 0 3 1 5 4 9 0 3 7 4 7 0 1
7 7 0 6 6 3 2 8 8 5 8 9 5 6 4 0 5 9 1 8 0 5 4 9 4
3 3 8 5 7 5 7 4 3 4 5 7 9 6 9 5 0 7 7 6 6 8 8 5 9
9 1 7 1 3 6 9 2 2 9 1 9 4 2 3 3 0 8 1 8 7 7 6 4 7 2
6 2 2 8 0 9 4 5 3 7 2 5 4 6 6 5 6 6 5 0 4 6 5 6 8
1 7 5 9 0 0 2 0 5 6 5 8 5 1 9 5 3 3 7 4 0 5 8 2 4
0 3 9 6 9 4 7 3 5 7 0 6 5 4 7 1 1 8 5 3 2 8 0 9 8
```

Grafický popis

Numerický popis

K čemu sumarizovat data? Není nejlepší mít všechna data?

```
8 2 0 3 1 4 5 8 2 1 7 2 7 3 8 5 5 2 9 0 6 3 1 6 4
0 8 7 3 3 1 9 7 5 2 5 7 6 9 8 0 3 6 2 5 1 2 7 5 2
2 3 3 8 6 6 1 4 2 4 0 2 6 1 8 9 5 2 6 9 8 3 4 0 1 0
4 7 5 5 6 3 0 7 7 1 9 1 6 1 7 4 1 7 1 3 7 9 3 3 7
1 9 3 9 5 3 4 9 5 5 2 7 5 8 0 3 4 8 8 1 2 7 5 3 4
2 8 7 8 1 4 1 4 9 4 2 4 1 5 2 9 4 6 2 1 5 2 8 1 9
8 4 8 5 1 3 9 6 6 0 7 2 1 9 0 2 0 6 7 0 6 0 1 3 0
0 3 8 8 4 7 5 1 5 1 7 3 4 5 2 0 7 4 7 9 6 6 7 7 4
3 5 3 1 9 3 7 4 9 5 0 2 0 1 4 6 2 5 4 5 8 5 0 9 2
3 4 5 9 5 2 7 9 8 9 0 5 5 8 5 1 7 7 3 5 5 4 7 7 2
4 1 5 3 0 9 1 3 7 2 5 8 7 7 1 3 6 3 9 7 8 7 9 1 7
7 2 9 5 6 7 8 5 4 5 3 4 5 4 1 9 8 6 7 5 7 9 3 1 8
5 9 2 8 9 8 6 4 4 1 5 3 7 7 0 8 0 2 5 6 0 6 1 2 0
1 3 3 3 9 0 5 2 8 7 4 0 9 0 3 7 3 1 7 9 4 5 5 2 8
4 6 0 1 0 8 6 2 1 0 0 5 0 3 1 5 4 9 0 3 7 4 7 0 1
7 7 0 6 6 3 2 8 8 5 8 9 5 6 4 0 5 9 1 8 0 5 4 9 4
3 3 8 5 7 5 7 4 3 4 5 7 9 6 9 5 0 7 7 6 6 8 8 5 9
9 1 7 1 3 6 9 2 9 1 9 4 2 3 3 0 8 1 8 7 7 6 4 7 2
6 2 2 8 0 9 4 5 3 7 2 5 4 6 6 5 6 6 5 0 4 6 5 6 8
1 7 5 9 0 0 2 0 5 6 5 8 5 1 9 5 3 3 7 4 0 5 8 2 4
0 3 9 6 9 4 7 3 5 7 0 6 5 4 7 1 1 8 5 3 2 8 0 9 8
```

Grafický popis

Numerický popis

1. - **jedné proměnné** (stupeň vzdělání, příjem, oblíbená barva)

K čemu sumarizovat data? Není nejlepší mít všechna data?

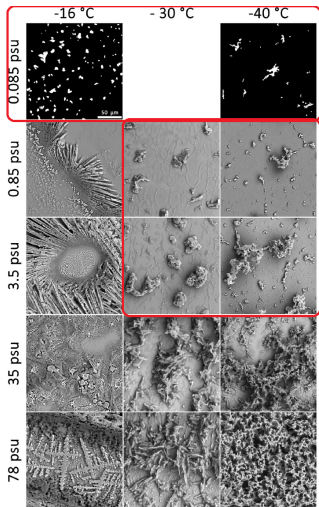
```
8 2 0 3 1 4 5 8 2 1 7 2 7 3 8 5 5 2 9 0 6 3 1 6 4
0 8 7 3 3 1 9 7 5 2 5 7 6 9 8 0 3 6 2 5 1 2 7 5 2
2 3 3 8 6 6 3 1 4 2 4 0 2 6 1 8 9 5 2 6 9 8 3 4 0 1 0
4 7 5 5 6 3 0 7 7 1 9 1 6 1 7 4 1 7 1 3 7 9 3 3 7
1 9 3 9 5 3 4 9 5 5 2 7 5 8 0 3 4 8 8 1 2 7 5 3 4
2 8 7 8 1 4 1 4 9 4 2 4 1 5 2 9 4 6 2 1 5 2 8 1 9
8 4 8 5 1 3 9 6 6 0 7 2 1 9 0 2 0 6 7 0 6 0 1 3 0
0 3 8 8 4 7 5 1 5 1 7 3 4 5 2 0 7 4 7 9 6 6 7 7 4
3 5 3 1 9 3 7 4 9 5 0 2 0 1 4 6 2 5 4 5 8 5 0 9 2
3 4 5 9 5 2 7 9 8 9 0 5 5 8 5 1 7 7 3 5 5 4 7 7 2
4 1 5 3 0 9 1 3 7 2 5 8 7 7 1 3 6 3 9 7 8 7 9 1 7
7 2 9 5 6 7 8 5 4 5 3 4 5 4 1 9 8 6 7 5 7 9 3 1 8
5 9 2 8 9 8 6 4 4 1 5 3 7 7 0 8 0 2 5 6 0 6 1 2 0 0
1 3 3 3 9 0 5 2 8 7 4 0 9 0 3 7 3 1 7 9 4 5 5 2 8
4 6 0 1 0 8 6 2 1 0 0 5 0 3 1 5 4 9 0 3 7 4 7 0 1
7 7 0 6 6 3 2 8 8 5 8 9 5 6 4 0 5 9 1 8 0 5 4 9 4
3 3 8 5 7 5 7 4 3 4 5 7 9 6 9 5 0 7 7 6 6 8 8 5 9
9 1 7 1 3 6 9 2 2 9 1 9 4 2 3 3 0 8 1 8 7 7 6 4 7 2
6 2 2 8 0 9 4 5 3 7 2 5 4 6 6 5 6 6 5 0 4 6 5 6 8
1 7 5 9 0 0 2 0 5 6 5 8 5 1 9 5 3 3 7 4 0 5 8 2 4
0 3 9 6 9 4 7 3 5 7 0 6 5 4 7 1 1 8 5 3 2 8 0 9 8
```

Grafický popis

Numerický popis

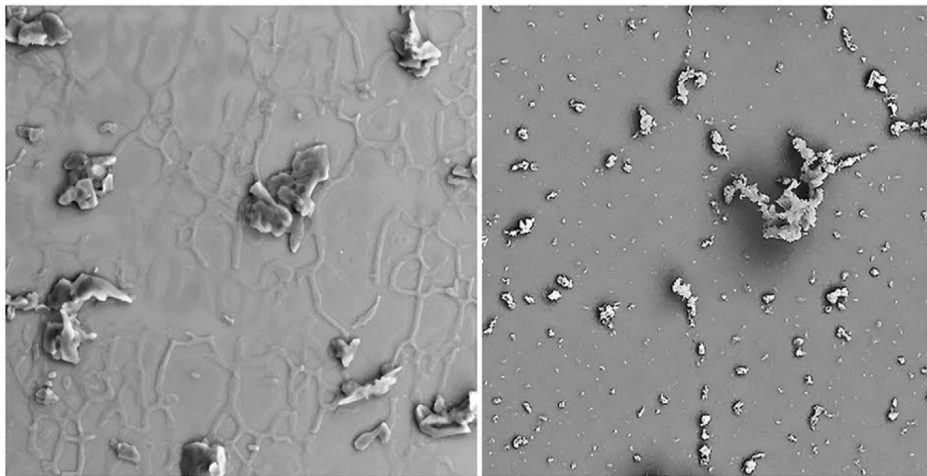
1. - **jedné proměnné** (stupeň vzdělání, příjem, oblíbená barva)
2. - **vzath mezi dvěmi proměnnými** (Jak stupeň vzdělání ovlivní příjem?)

Relationship between the Particle Sizes, Sea Salt Concentration, and Sublimation Temperature



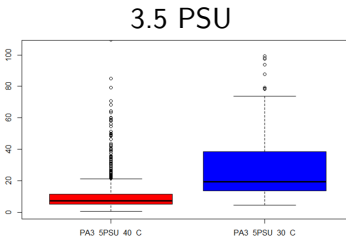
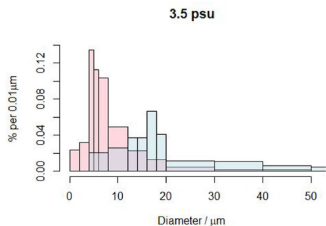
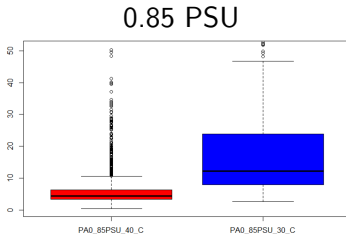
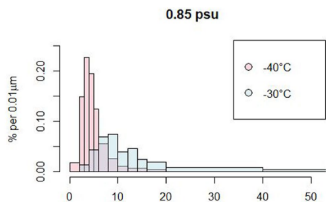
ZÁVACKÁ, K., V. NEDĚLA, M. OLBERT, E. TIHLAŘÍKOVÁ, L. VETRÁKOVÁ, X. YANG AND D. HEGER Temperature and Concentration Affect Particle Size Upon Sublimation of Saline Ice: Implications for Sea Salt Aerosol Production in Polar Regions. *Geophysical Research Letters*, 2022, 49(8) 10.1029/2021GL097098

At $-30\text{ }^{\circ}\text{C}$, Brine is Still Fluid to Become Mostly Solid at $-40\text{ }^{\circ}\text{C}$ (0.85 PSU) ESEM pictures



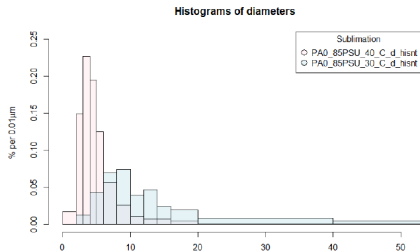
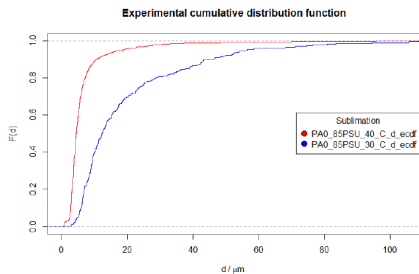
ZÁVACKÁ, K., V. NEDĚLA, M. OLBERT, E. TIHLAŘÍKOVÁ, L. VETRÁKOVÁ, X. YANG AND D. HEGER Temperature and Concentration Affect Particle Size Upon Sublimation of Saline Ice: Implications for Sea Salt Aerosol Production in Polar Regions. *Geophysical Research Letters*, 2022, 49(8),10.1029/2021gl097098.

SSAs: the Size Decreases with the Temperature and Concentration: Histogram and Box-Plot



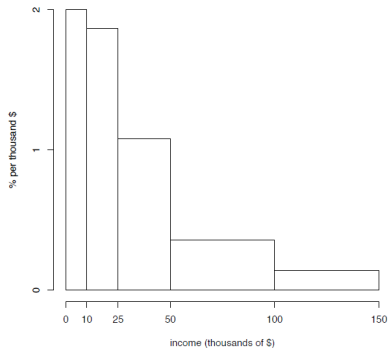
ZÁVACKÁ, K., V. NEDĚLA, M. OLBERT, E. TIHLAŘÍKOVÁ, L. VETRAKOVÁ, X. YANG AND D. HEGER Temperature and Concentration Affect Particle Size Upon Sublimation of Saline Ice: Implications for Sea Salt Aerosol Production in Polar Regions. *Geophysical Research Letters*, 2009, 36(10): 10.1029/2009-102709

Empirická distribuční funkce



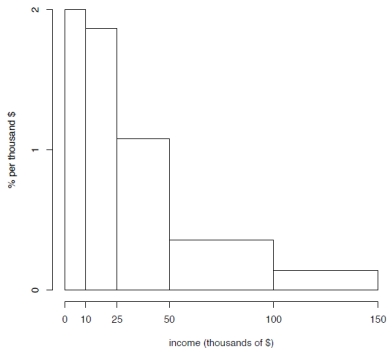
A right-skewed distribution = do prava zšikmělé rozložení

Income in USA 2010



A right-skewed distribution = do prava zšikmělé rozložení

Income in USA 2010



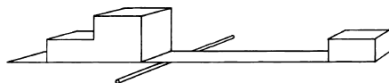
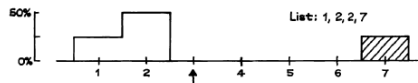
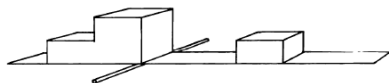
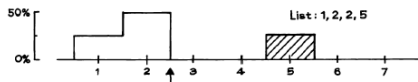
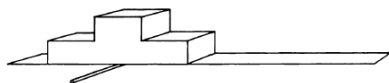
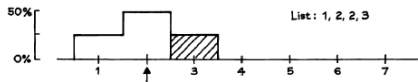
Income in USA 1973

Figure 4. Distribution of families by income in the U.S. in 1973.



- 1 plocha pravoúhelníků představuje pravděpodobnosti
- 2 výška pravoúhelníků představuje hustotu - procenta celkové populace na jednotku horizontální osy.

Average



Comparing averages

age (years)	20-30	30-40	40-50	50-60	60-75	75+
average height(")	69.3	69.5	69.4	69.2	68.3	67.2

Intervals include the left endpoint but not the right.

[National Health and Nutrition Examination Survey, 1999-2002]

Comparing averages

age (years)	20-30	30-40	40-50	50-60	60-75	75+
average height(")	69.3	69.5	69.4	69.2	68.3	67.2

Intervals include the left endpoint but not the right.

[National Health and Nutrition Examination Survey, 1999-2002]

Is this table telling us that as men get older, on average they get a bit taller and then get shorter?

Comparing averages

age (years)	20-30	30-40	40-50	50-60	60-75	75+
average height(")	69.3	69.5	69.4	69.2	68.3	67.2

Intervals include the left endpoint but not the right.

[National Health and Nutrition Examination Survey, 1999-2002]

Is this table telling us that as men get older, on average they get a bit taller and then get shorter?

Data - **longitudinal** × **cross-sectional**

Comparing averages

age (years)	20-30	30-40	40-50	50-60	60-75	75+
average height(")	69.3	69.5	69.4	69.2	68.3	67.2

Intervals include the left endpoint but not the right.

[National Health and Nutrition Examination Survey, 1999-2002]

Is this table telling us that as men get older, on average they get a bit taller and then get shorter?

Data - **longitudinal** × **cross-sectional**

The men in various categories are not the same.

Comparing averages

age (years)	20-30	30-40	40-50	50-60	60-75	75+
average height(")	69.3	69.5	69.4	69.2	68.3	67.2

Intervals include the left endpoint but not the right.

[National Health and Nutrition Examination Survey, 1999-2002]

Is this table telling us that as men get older, on average they get a bit taller and then get shorter?

Data - **longitudinal** × **cross-sectional**

The men in various categories are not the same.

The older men were shorter, on average.

Comparing averages

age (years)	20-30	30-40	40-50	50-60	60-75	75+
average height(")	69.3	69.5	69.4	69.2	68.3	67.2

Intervals include the left endpoint but not the right.

[National Health and Nutrition Examination Survey, 1999-2002]

Is this table telling us that as men get older, on average they get a bit taller and then get shorter?

Data - **longitudinal** × **cross-sectional**

The men in various categories are not the same.

The older men were shorter, on average.

When comparing averages first think:

- 1 How are the groups related to each other?
- 2 Take a look on the numerical averages.

Measures of Spread

How can we quantify your distance from the median and/or mean?

Measures of Spread

How can we quantify your distance from the median and/or mean?

How far from average am I?

Measures of Spread

How can we quantify your distance from the median and/or mean?

How far from average am I?

How much am I deviating?

Measures of Spread

How can we quantify your distance from the median and/or mean?

How far from average am I?

How much am I deviating?

The amount your score is off (from average) is the **deviation**.

Range and interquartile range

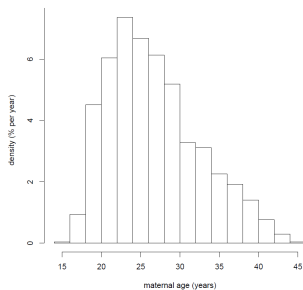
Range and interquartile range

How far from median?

Range and interquartile range

How far from median?

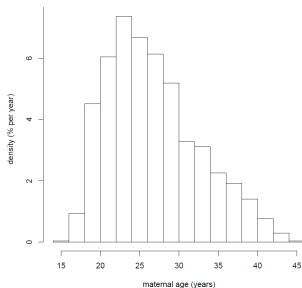
Maternal ages



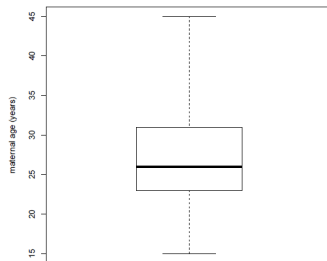
Range and interquartile range

How far from median?

Maternal ages



Box plot

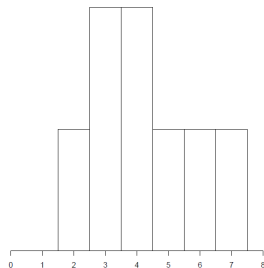


What is the typical (standard)
deviation from average?

What is the typical (standard)
deviation from average?

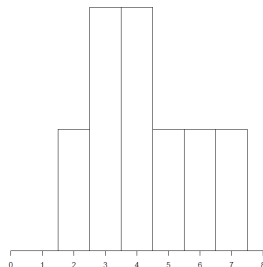
Deviation from average

What is the typical (standard) deviation from average?



Deviation from average

What is the typical (standard) deviation from average?



SD = standard deviation = směrodatná odchylka
= (square) root mean square of deviations from average
= (druhá) odmocnina průměru čtverců odchylek od průměru

variance = rozptyl
= mean square of deviations from average
= průměr čtverců odchylek od průměru

Properties of SD

Why SD is so commonly used measure of spread?

SD for given distribution measures typical distance from average.

- 1 It is non negative
- 2 It has the same units as average and the list.
- 3 It measures the average distance from the data to their mean (rms of the deviations of the data from their mean)
- 4 Chebychev inequality

Pafnuty Lvovich Chebychev (1821 - 1894)

In any list, the proportion of entries that are k or more SDs away from the average is **at most** $1/k^2$.

Pro jakoukoli číselnou řadu platí: podíl členů, které jsou od průměru vzdáleny **alespoň** k -krát SD je **nejvíce** $1/k^2$.

https://courses.edx.org/courses/BerkeleyX/Stat_2.1x/

<http://www.stat.berkeley.edu/~stark/SticiGui/Text/location.htm>