

Stats + Python



Week 1 - Day 04

Descriptive Statistics

Vs.

Inferential Statistics

Basic statistics

Mean/Average

Median

Mode

Min/Max

Quartiles/Percentiles

Values are sorted from min to max

Median - split in 2

Quartile - split in 4

1st quartile

2nd quartile = median

3rd quartile

Median - split in 2

Quartile - split in 4

Percentile - split in 100

1st quartile = 25th percentile

2nd quartile = 50th percentile = median

3rd quartile = 75th percentile

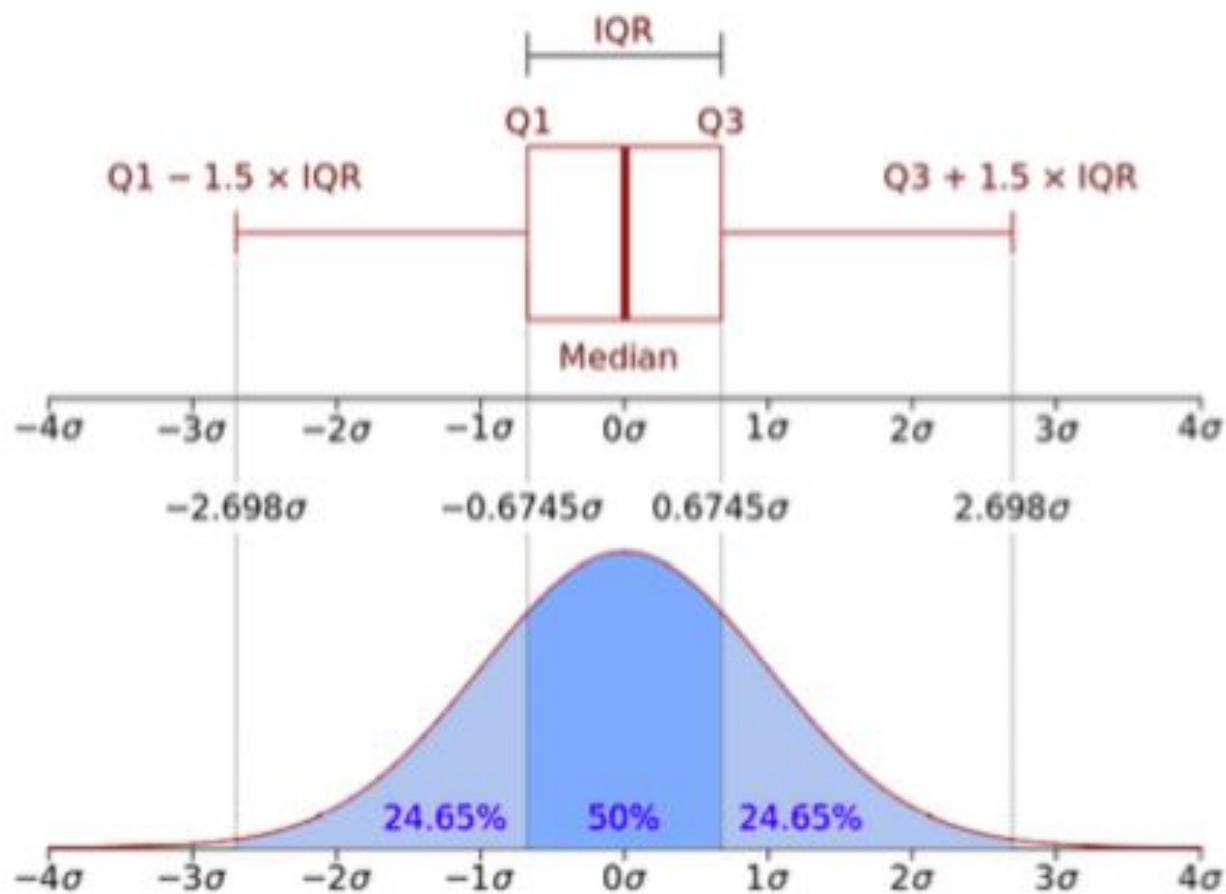
What's the value of the 5th percentile?

What's the value of the 95th percentile?

InterQuartile Range (IQR)

=

3th quartile - 1st quartile



Measures of Dispersion

How spread is our data?

[98, 99, 101, 102]

Vs.

[1, 2, 198, 199]

Range = max-min

Variance = var

Standard deviation = std = sqrt(var)

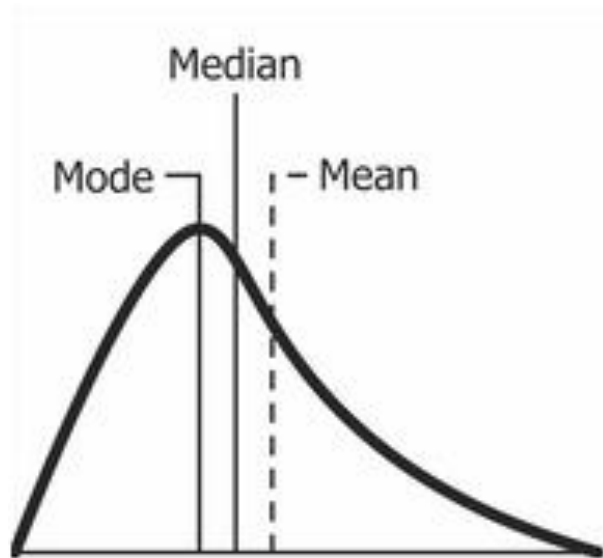
[98, 99, 101, 102], mean = 100

[-2, 1, 1, 2]

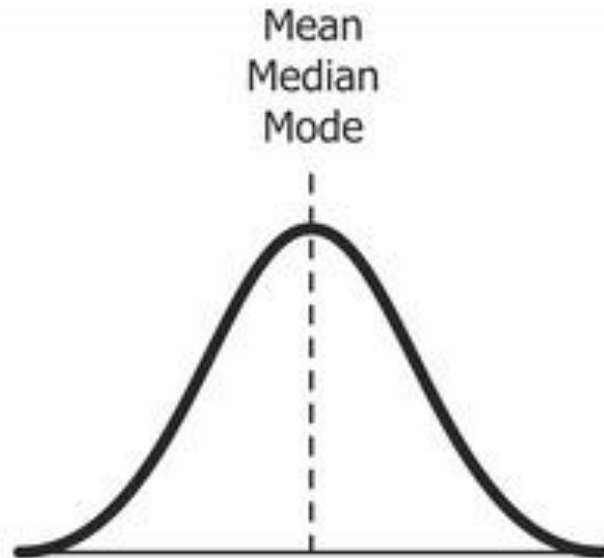
[4, 1, 1, 4]

$\text{sum}([4, 1, 1, 4]) / \text{len}([4, 1, 1, 4]) = 2.4$

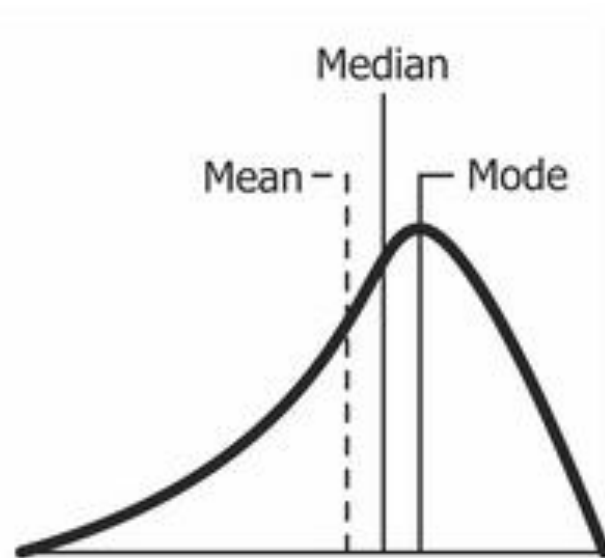
$\text{sqrt}(2.4) = 1.58$



Positive
Skew

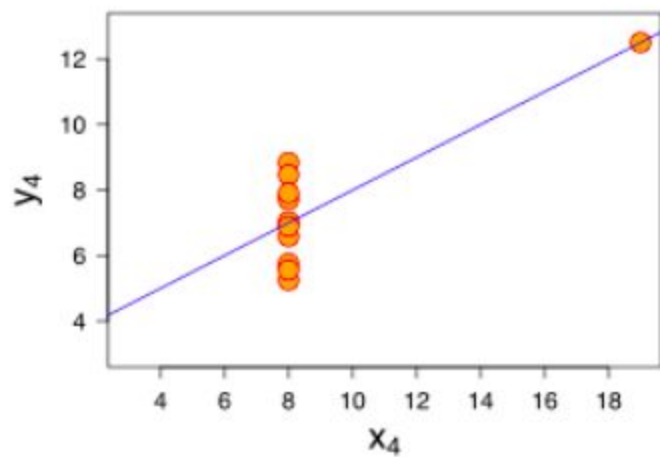
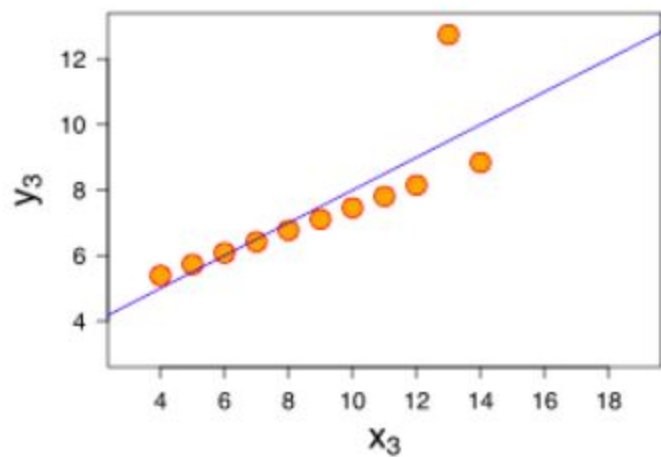
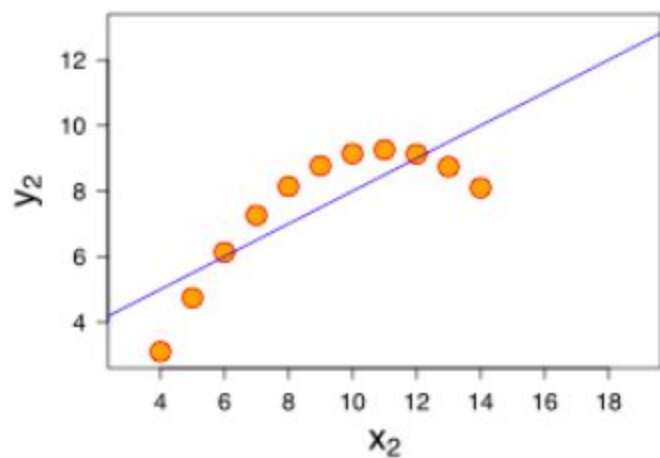
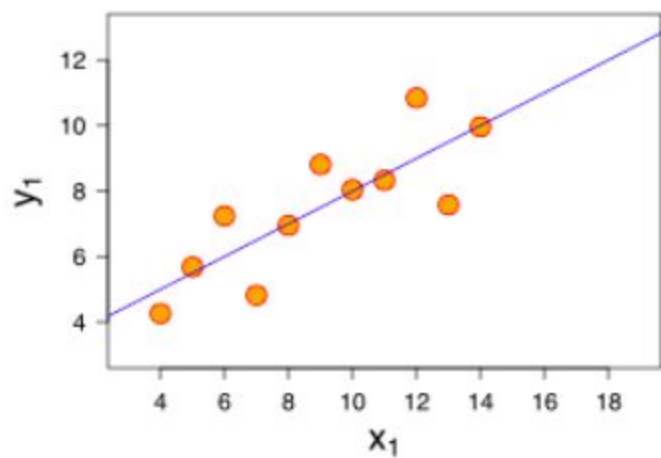


Symmetrical
Distribution



Negative
Skew

**What's wrong with
summary statistics?**



Property	Value
Mean of x	9
Sample variance of x	11
Mean of y	7.50
Sample variance of y	4.125
Correlation between x and y	0.816
Linear regression line	$y = 3.00 + 0.500x$
Coefficient of determination of the linear regression	0.67

Mean

What would my
starting salary be?



I'll put it this way:
our average starting
salary is \$80,000!



you → \$ 30,000

all your coworkers {
\$ 30,000
\$ 30,000
\$ 30,000
\$ 30,000
\$ 30,000
\$ 30,000

CEO's son → \$ 430,000

Average: \$80,000.



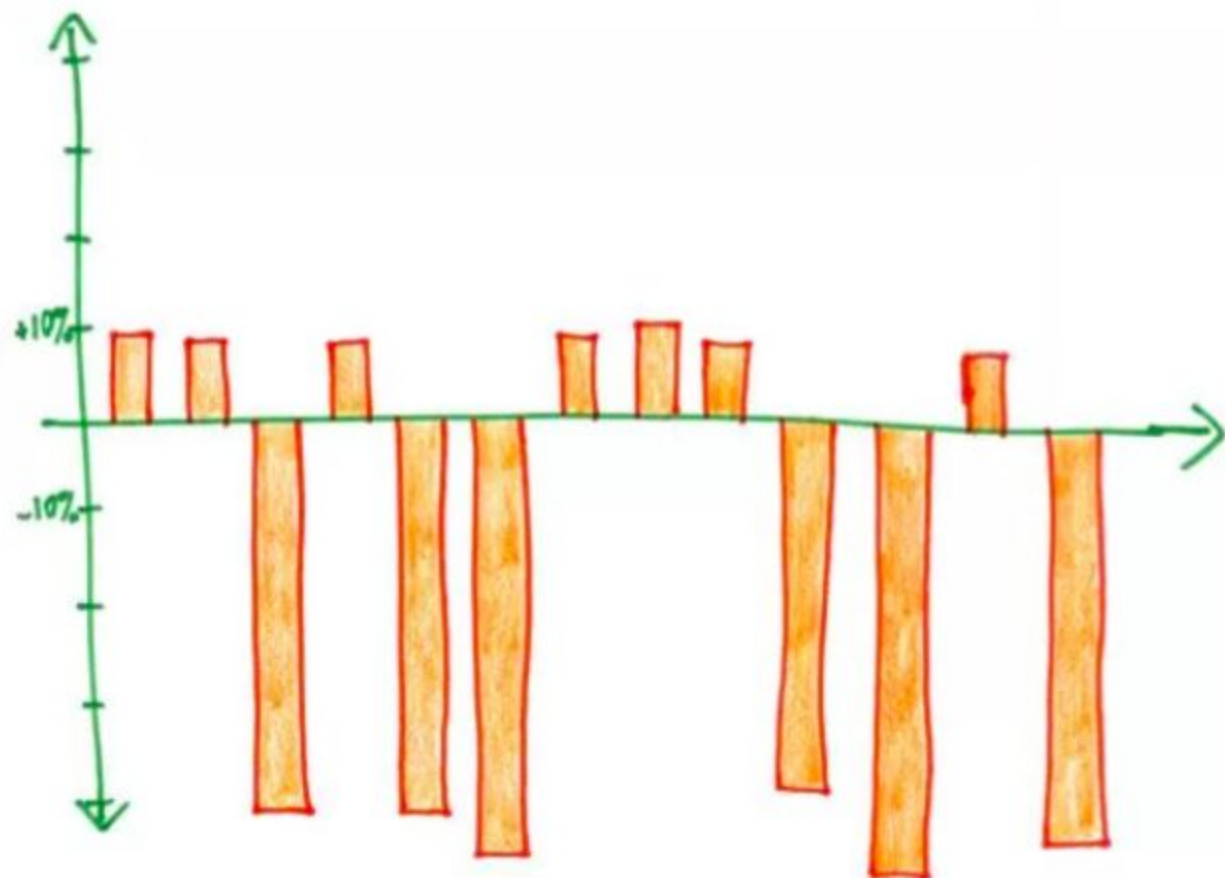
Median

So, why should I
invest with you?



Well, not to brag, but
my fund has a median
gain of 8% per year!





Mode

How are you doing
on your tests?



My modal category
is 70-80%!



please don't ask
about the mean...

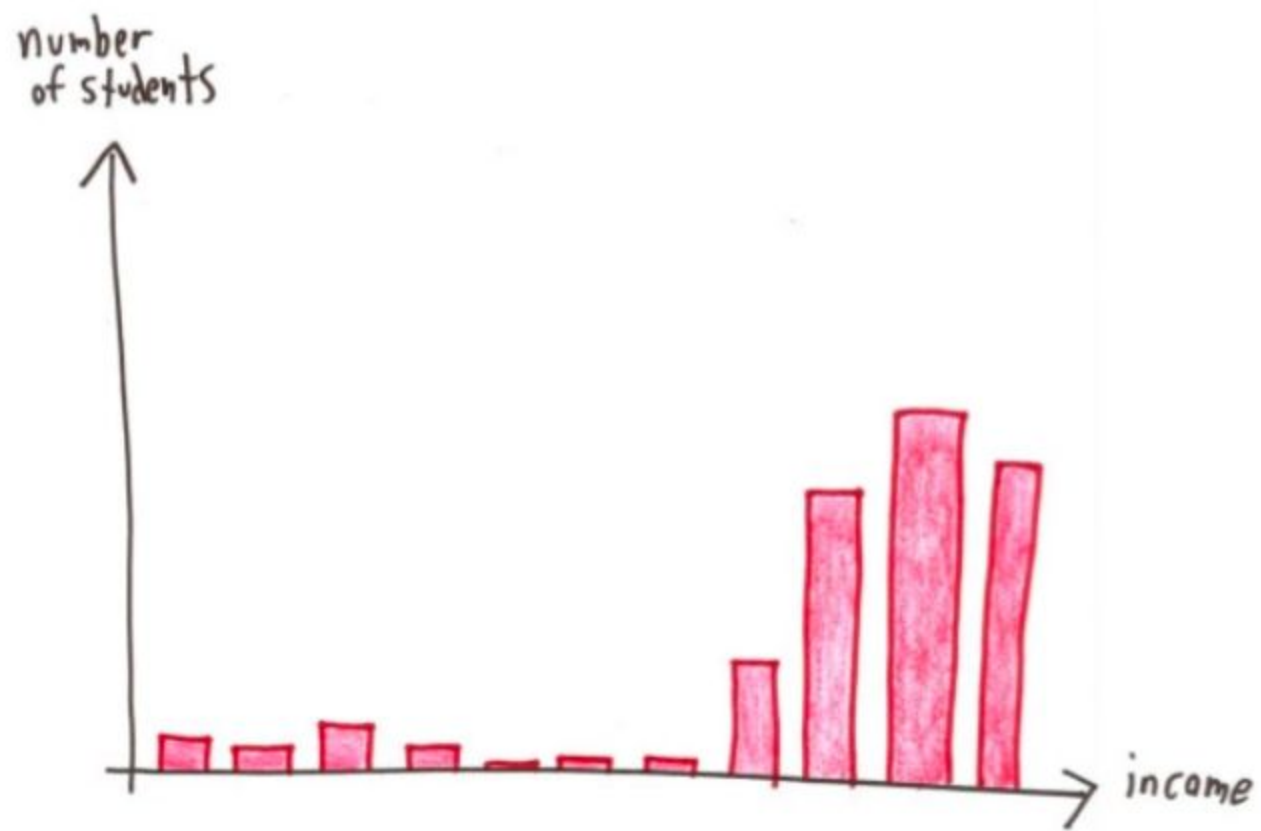


Score Category	Number of Tests
90s	0
80s	0
70s	2
60s	1
50s	1
40s	1
30s	1
20s	1

Range

Our students come from a
wide range of
socioeconomic
backgrounds...





Variance

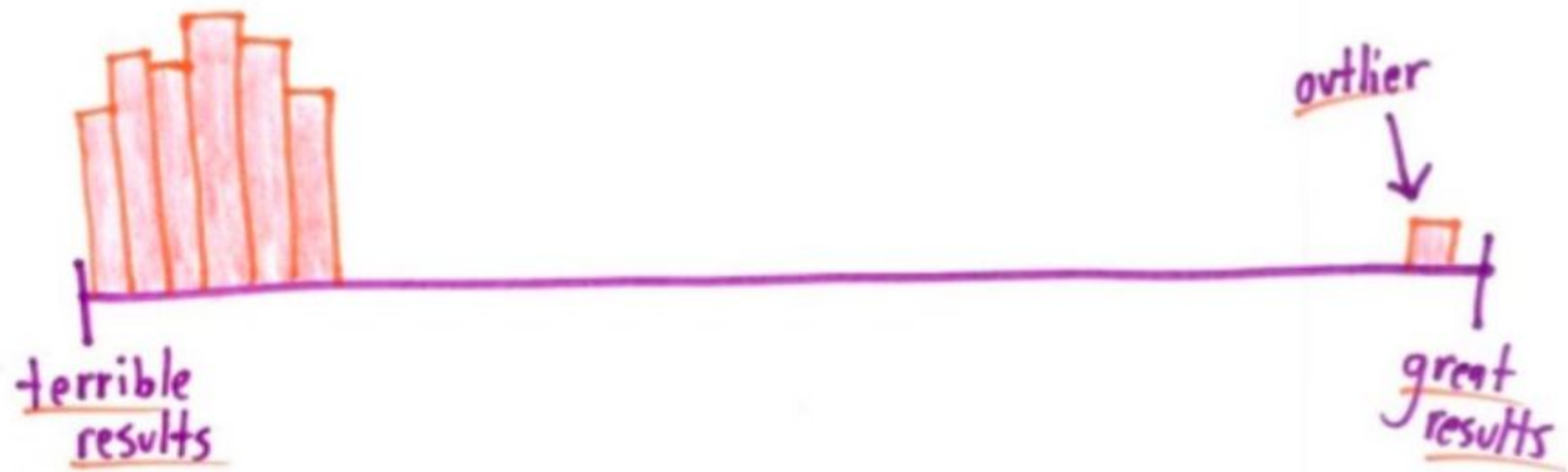
These results are
a disaster!



Sure, they look bad,
but there's a lot of variance!

Don't rush
to judgment.





Solution 1 = Data Visualization

Solution 2 = Use 3+ metrics

**Stats + Python =
Numpy**

```
import numpy as np
```

```
np.mean([1,2,3])
```

```
np.median([1,2,3])
```

```
np.std([1,2,3])
```

```
np.percentile([1,2,3], 50)
```