# Statistics and Machine Learning

## Linear regression and beyond: bootstrapping and model complexity

Week 4 02/08 — 02/12 2001
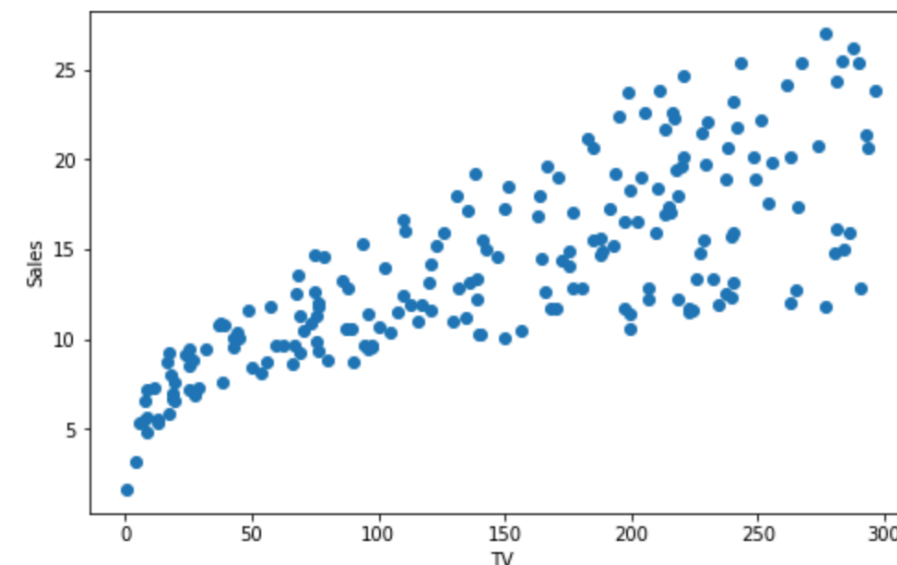
# Contents of Week 4

## Random data set and statistical meaning

- Random data set

- Interpretation of linear model

- Is 'newspaper' irrelevant to 'sales'?

- Bootstrapping data set

- Normal distribution: one sigma, two sigma, three sigma,...,
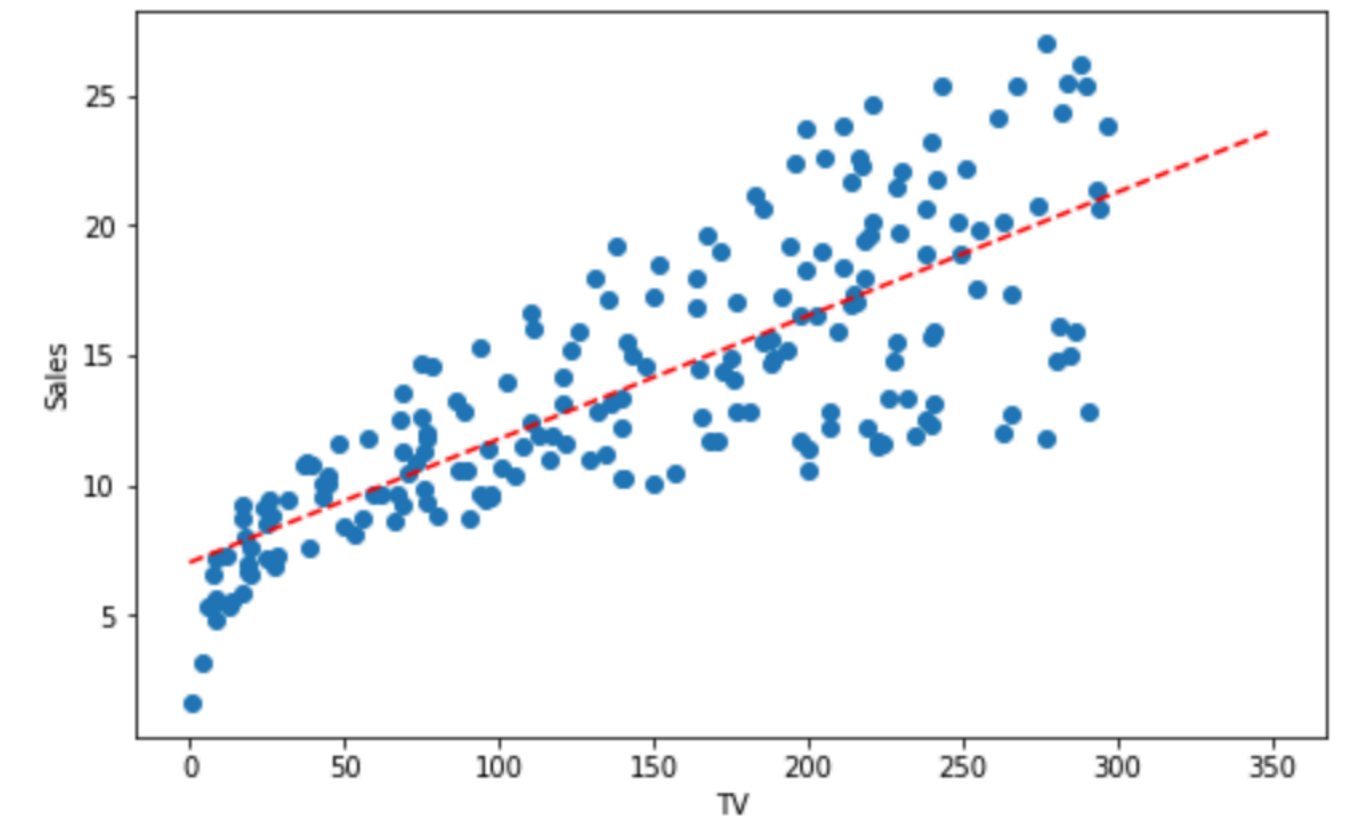
# Review of linear regression

## Finding the line which causes the minimum quadratic loss

**Data set**
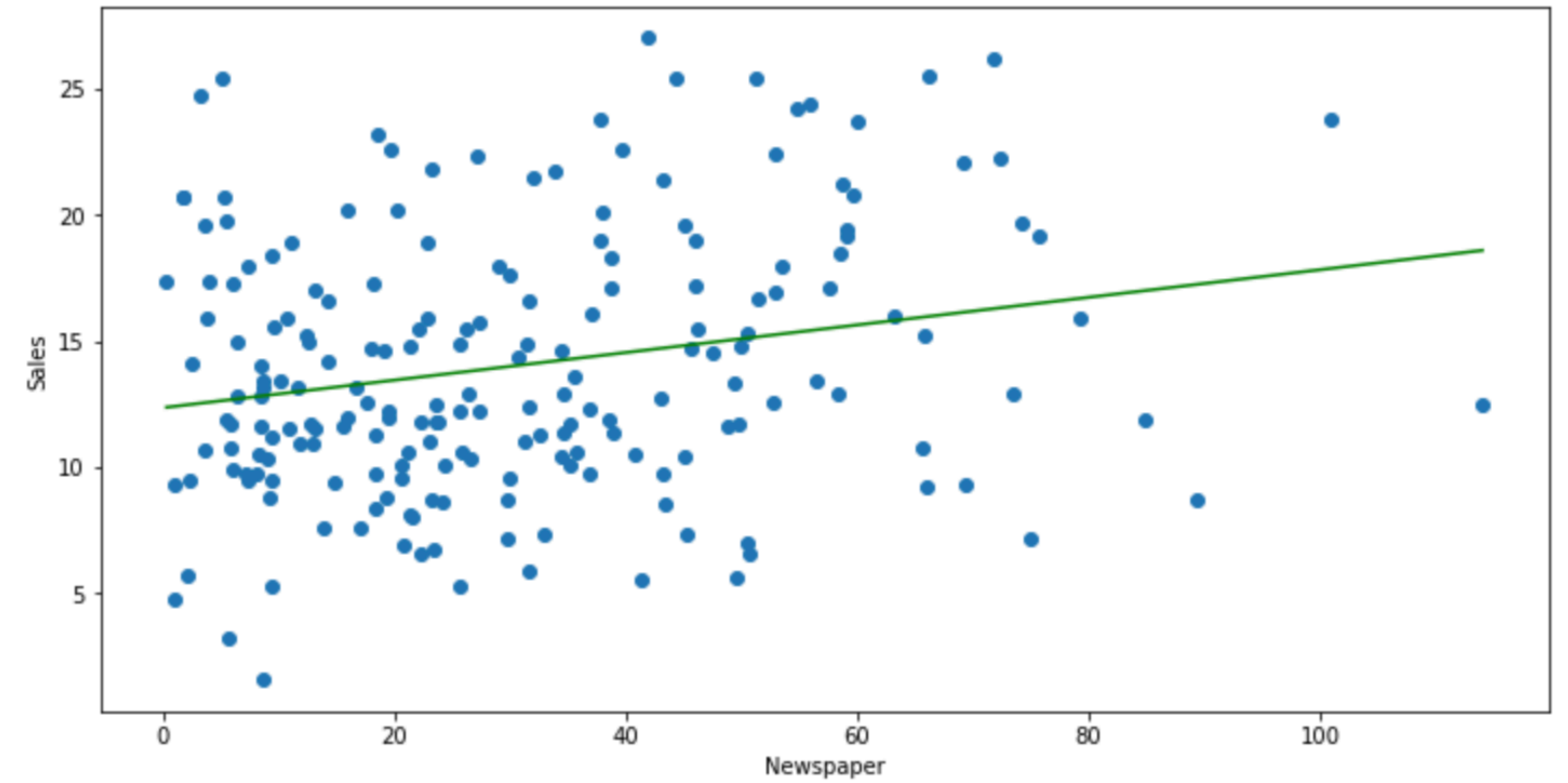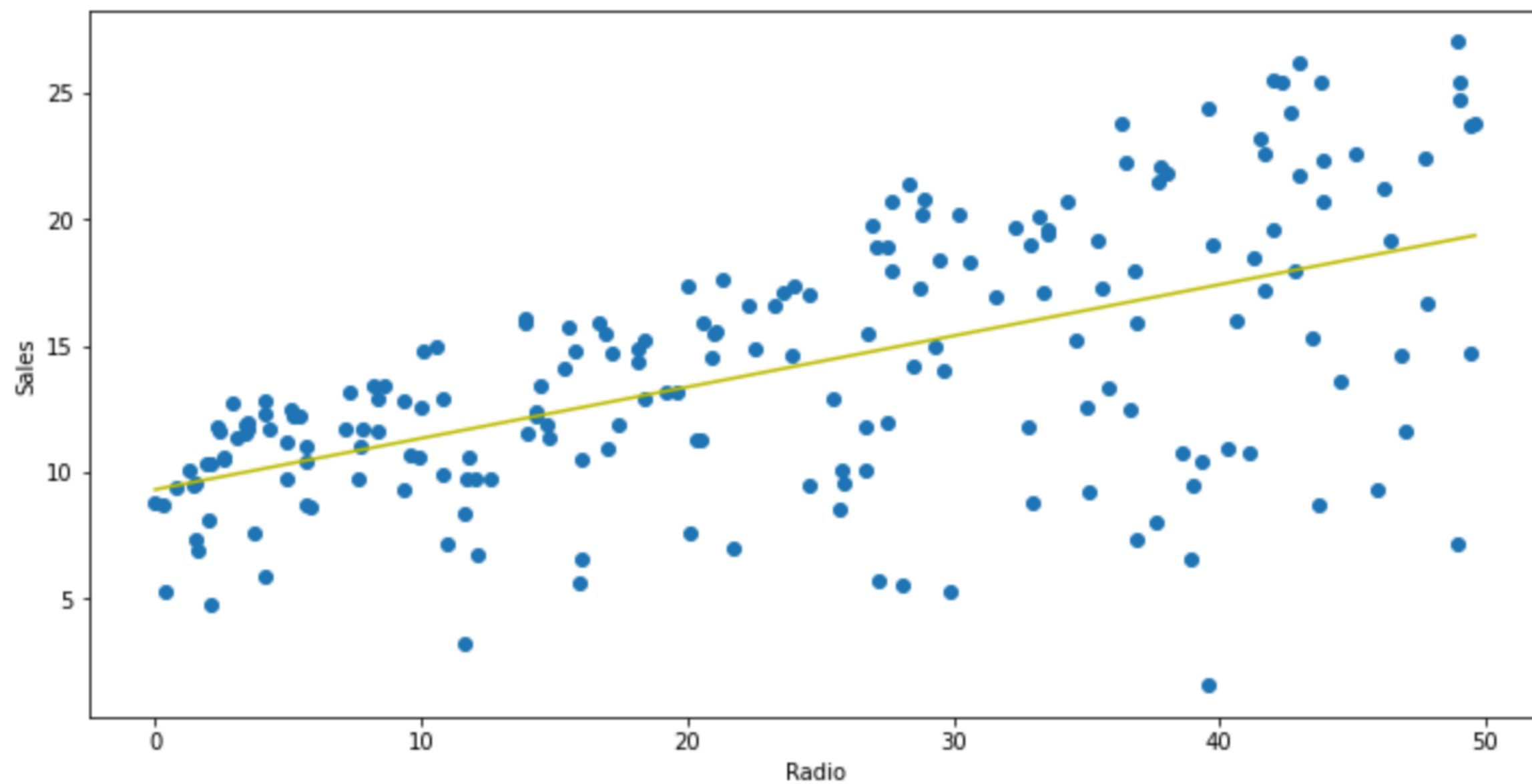
**Model**

**Linear regression**

**.fit()**
**a<−.coef_**
**b <− .intercept_**
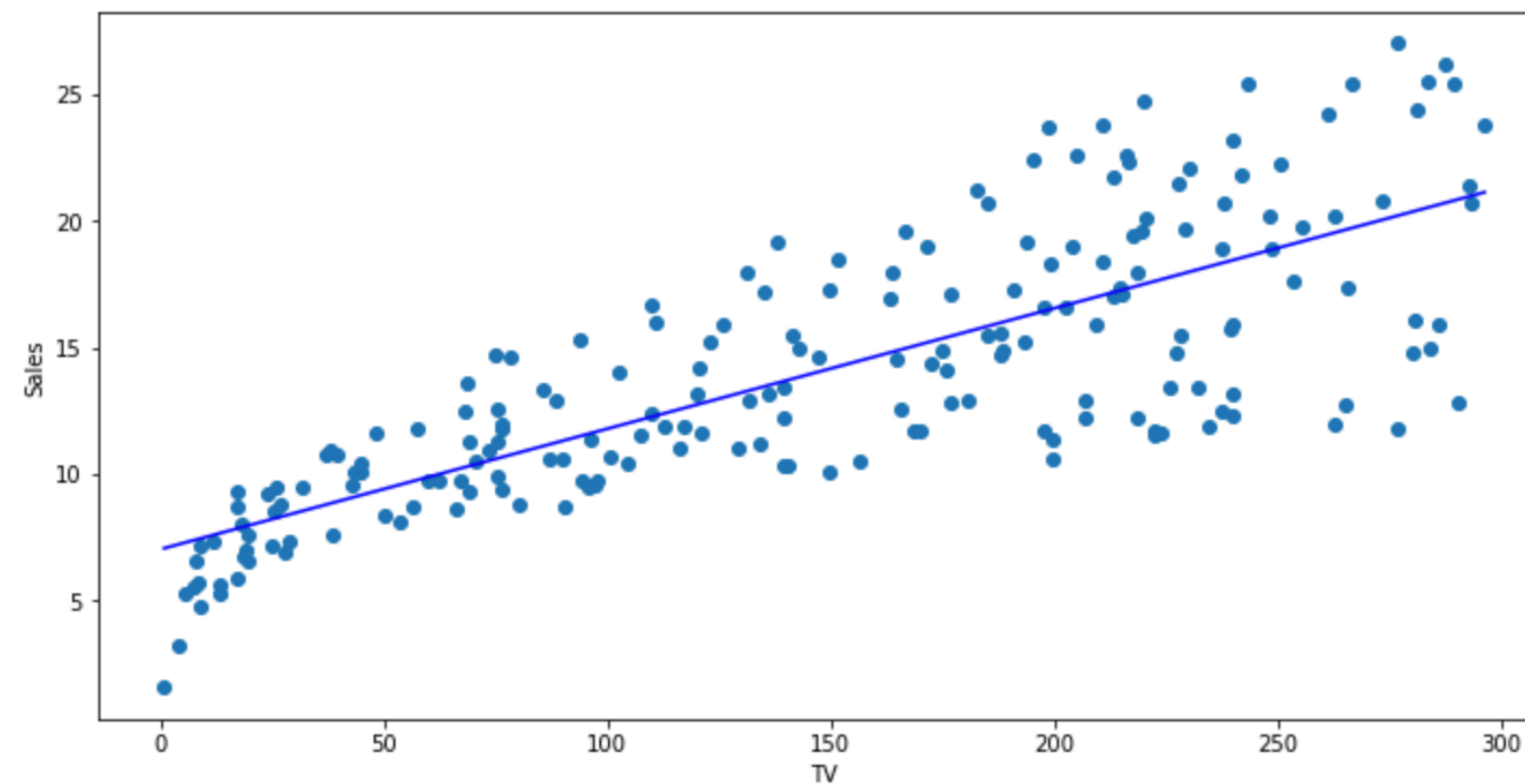
**The linear model:** $y = ax + b$

# Application to (radio, sales) and (newspaper, sales)
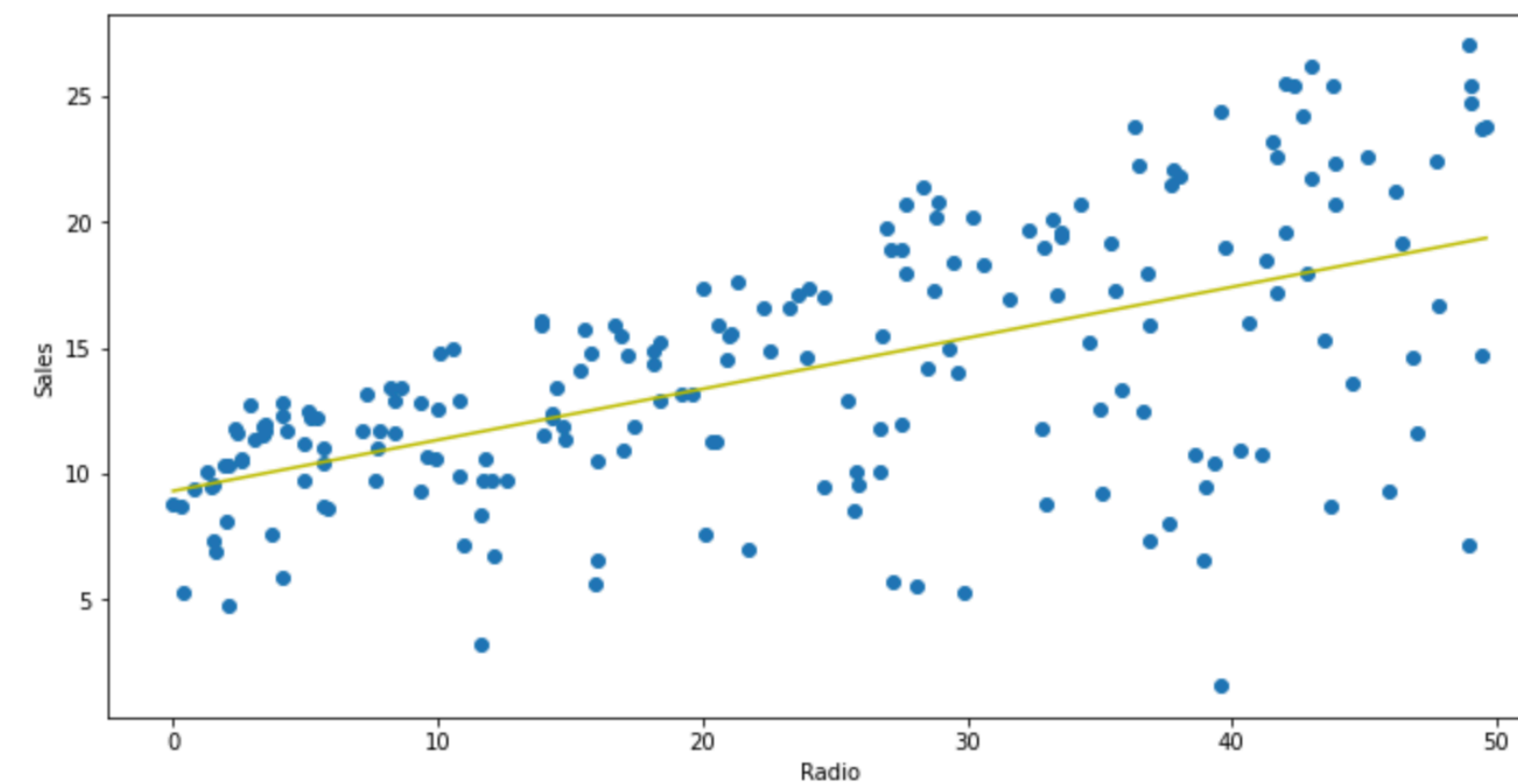
# 'TV', 'Radio', and 'Newspaper' versus 'Sales'

## Which one is the least relevant?

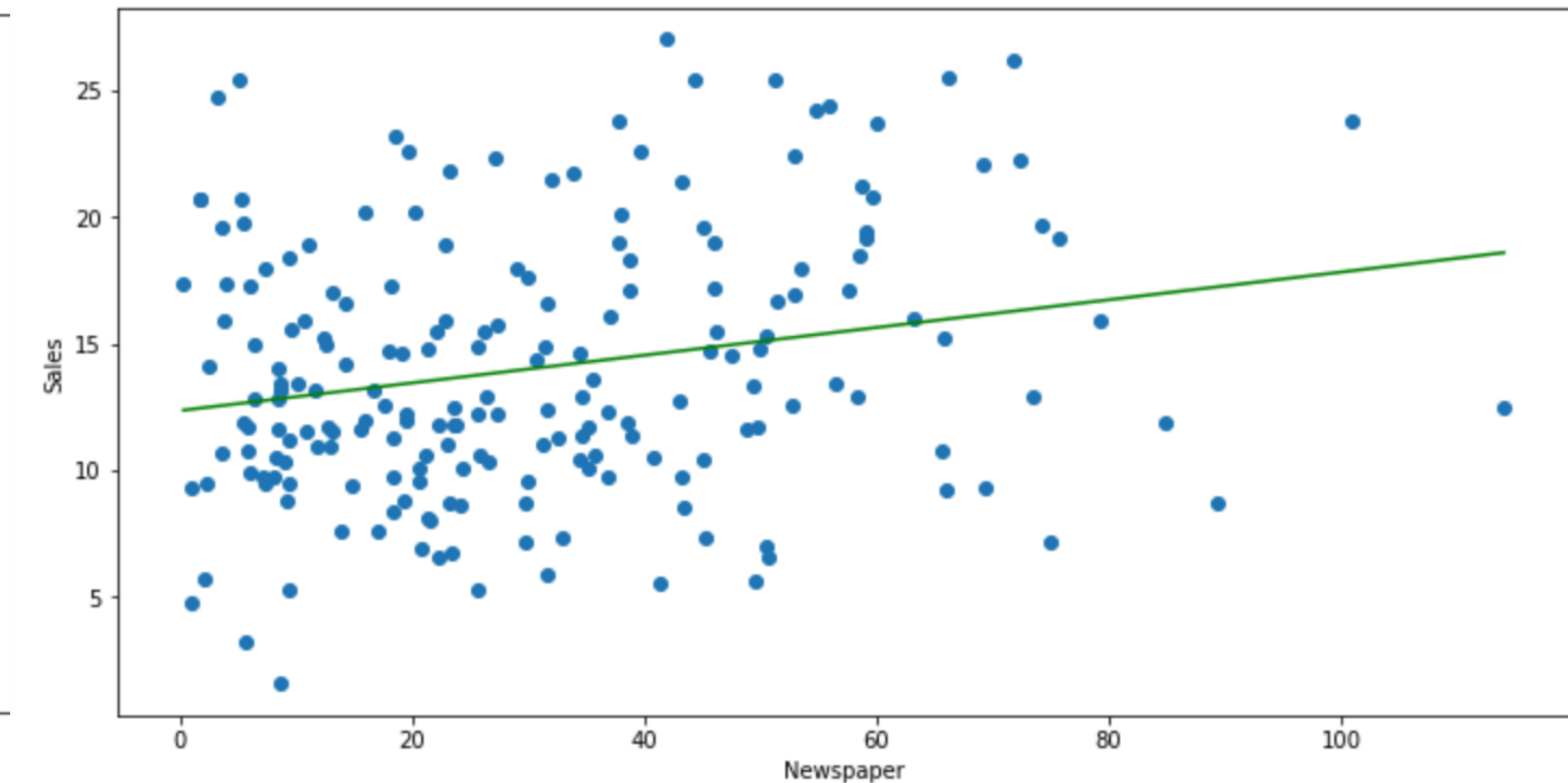Observation: if the model parameter of slope is close to zero, it means the x input is not relevant to the output!



```
regr.coef_, regr.intercept_
(array([[0.04753664]]), array([7.03259355]))
```

```
regr_radio.coef_, regr_radio.intercept_
(array([[0.20249578]]), array([9.3116381]))
```

```
regr_news.coef_, regr_news.intercept_
(array([[0.0546931]]), array([12.35140707]))
```

**Slope = 0.0475**

**Slope = 0.2**

**Slope = 0.05**

# Before we jump to conclusion,

## The caveats are

'Sales' is the result of all three attributes 'tv', 'radio', and 'newspaper'. Every time we consider one attribute and neglect the other two. It might be problematic!
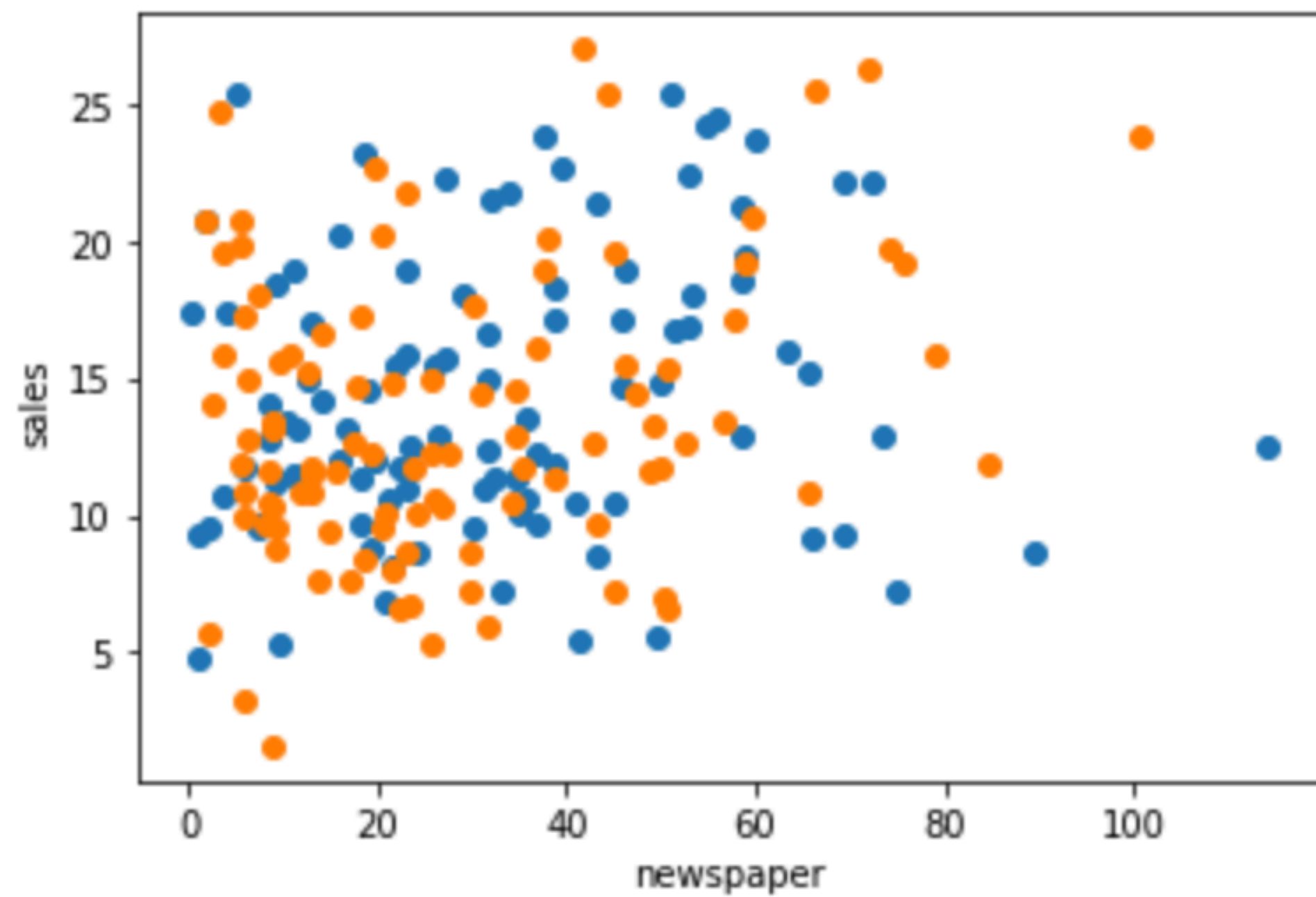
We do not know what does the unit for every attribute mean? So it might lead to comparing apple and orange.

Do we have a way out? yes.

# Different data led to different conclusions

**Example: sales versus newspaper**
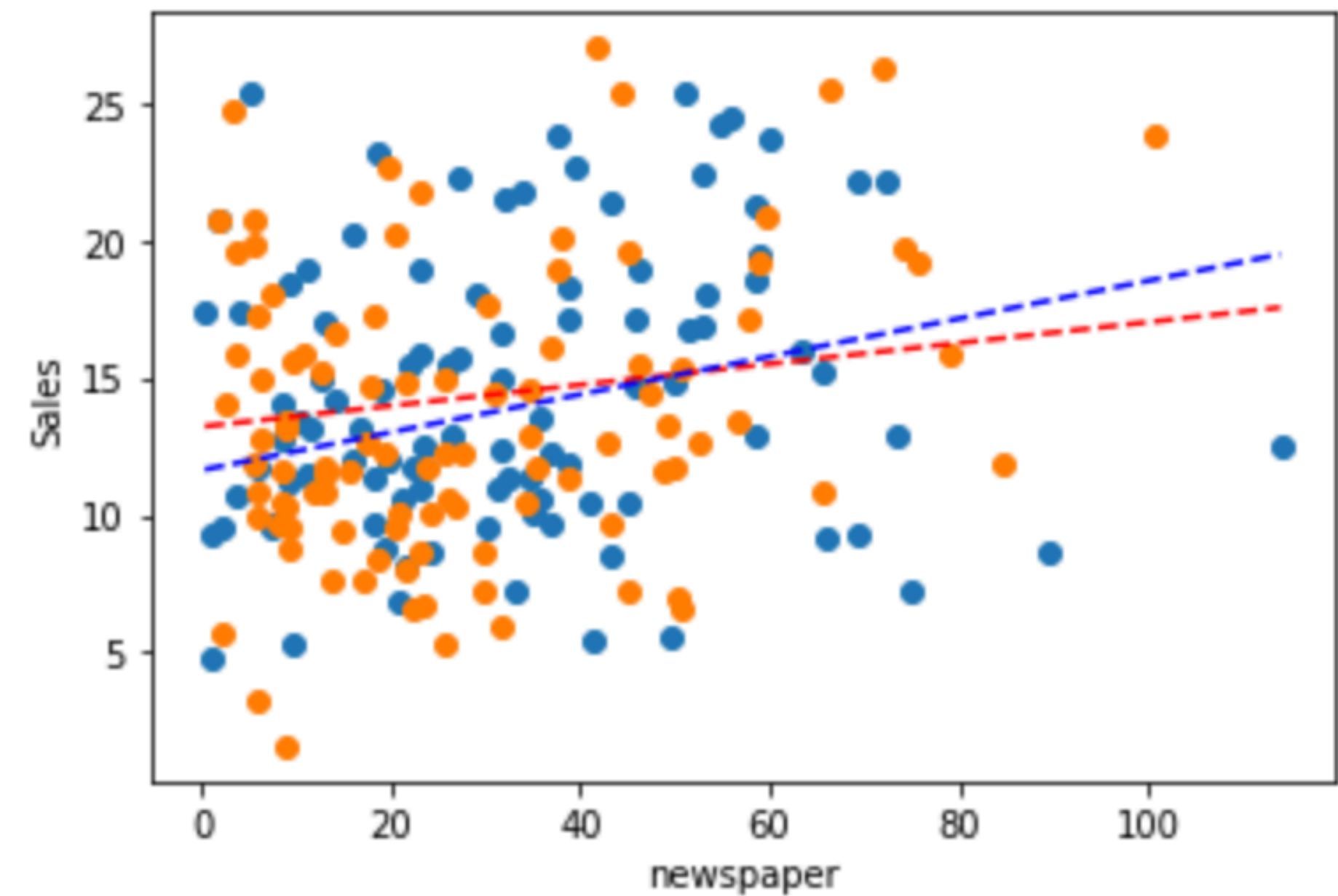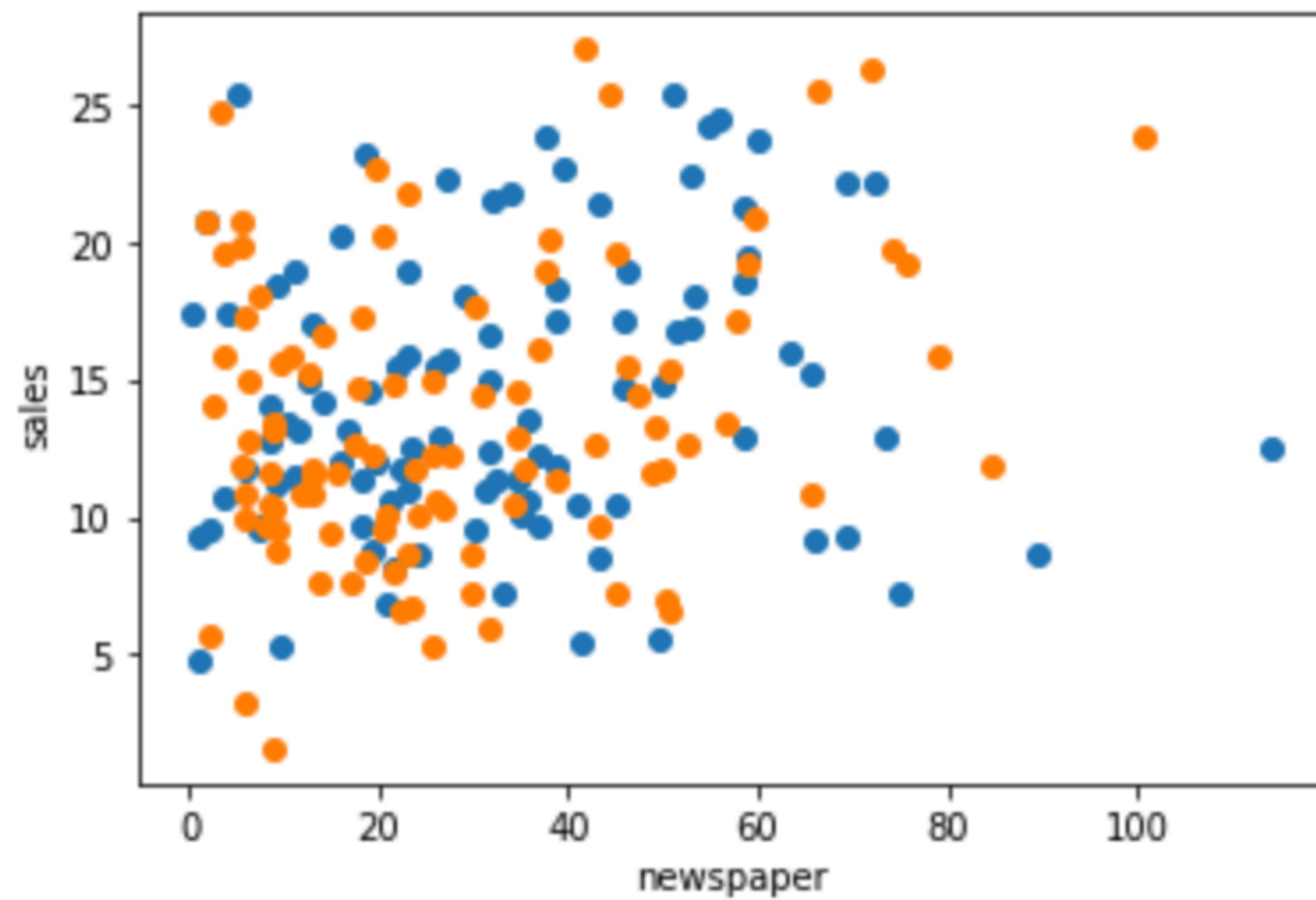
**200 data split into 100 + 100 data**

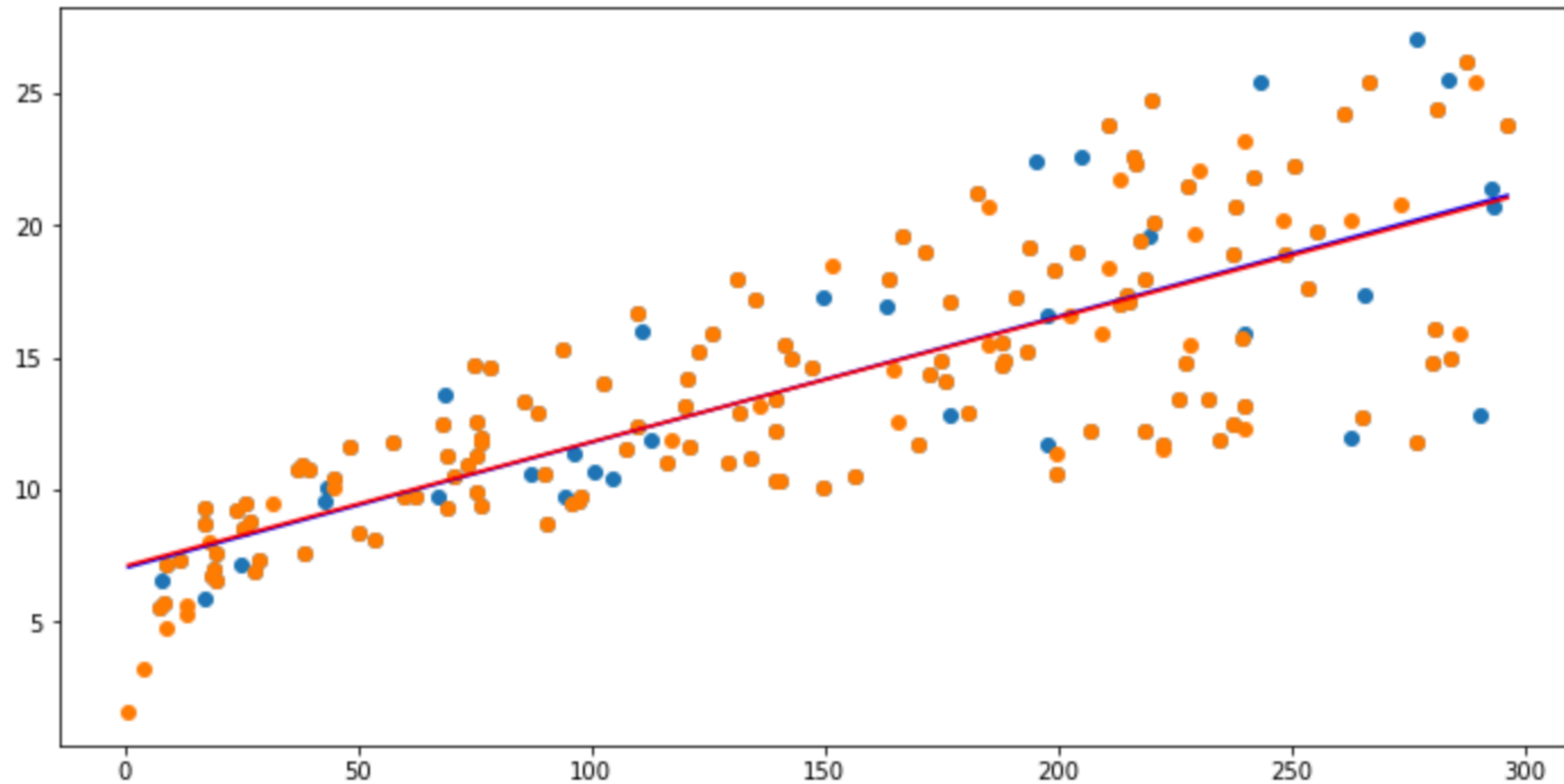# Different data led to different conclusions

**Example: sales versus newspaper**

**200 data split into 100 + 100 data**

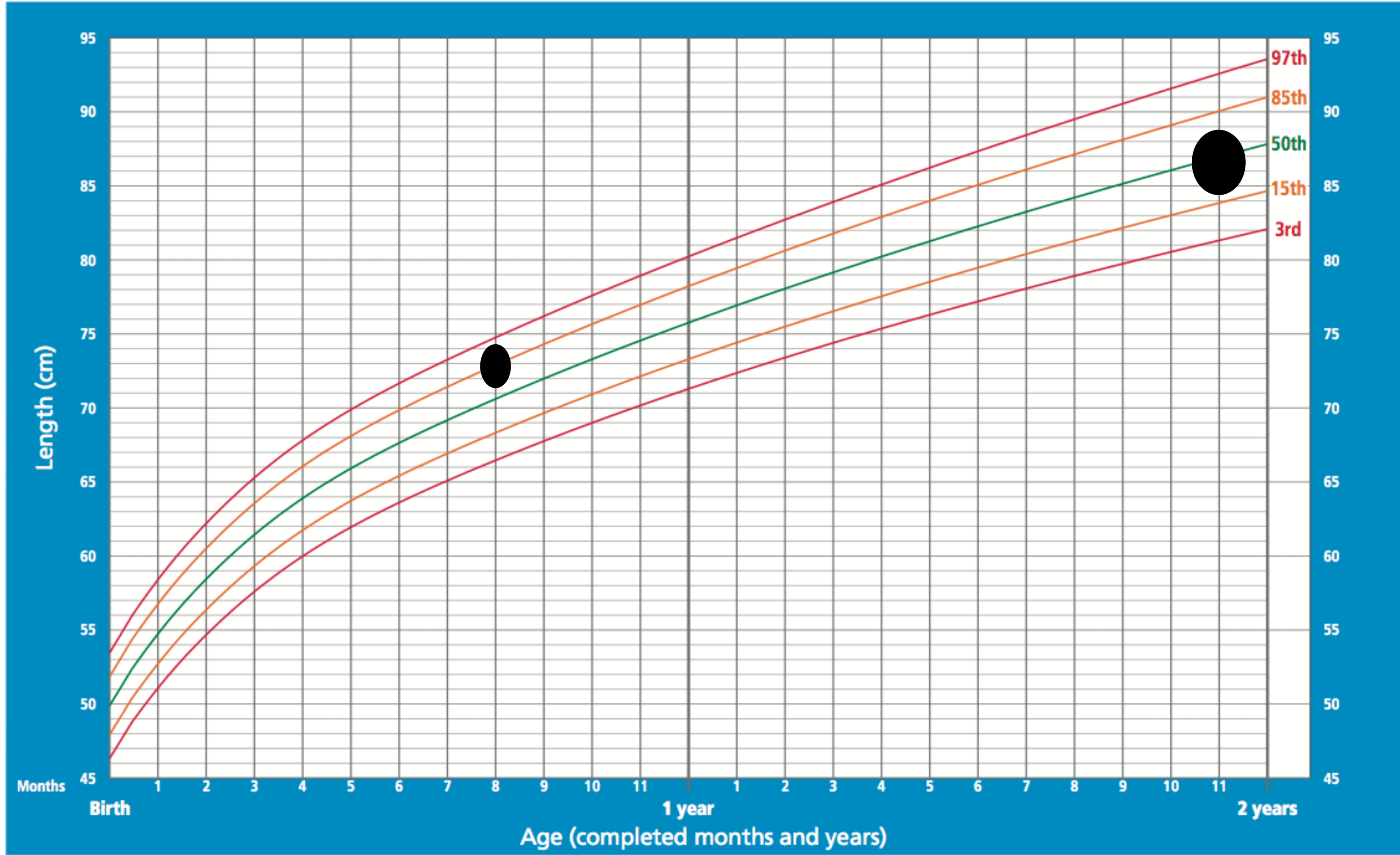# For (tv, sales), different data sets lead to almost identical linear model

# How to compare who is taller between elder and younger brothers?

## Standardized data



**Length-for-age BOYS**

Birth to 2 years (percentiles)
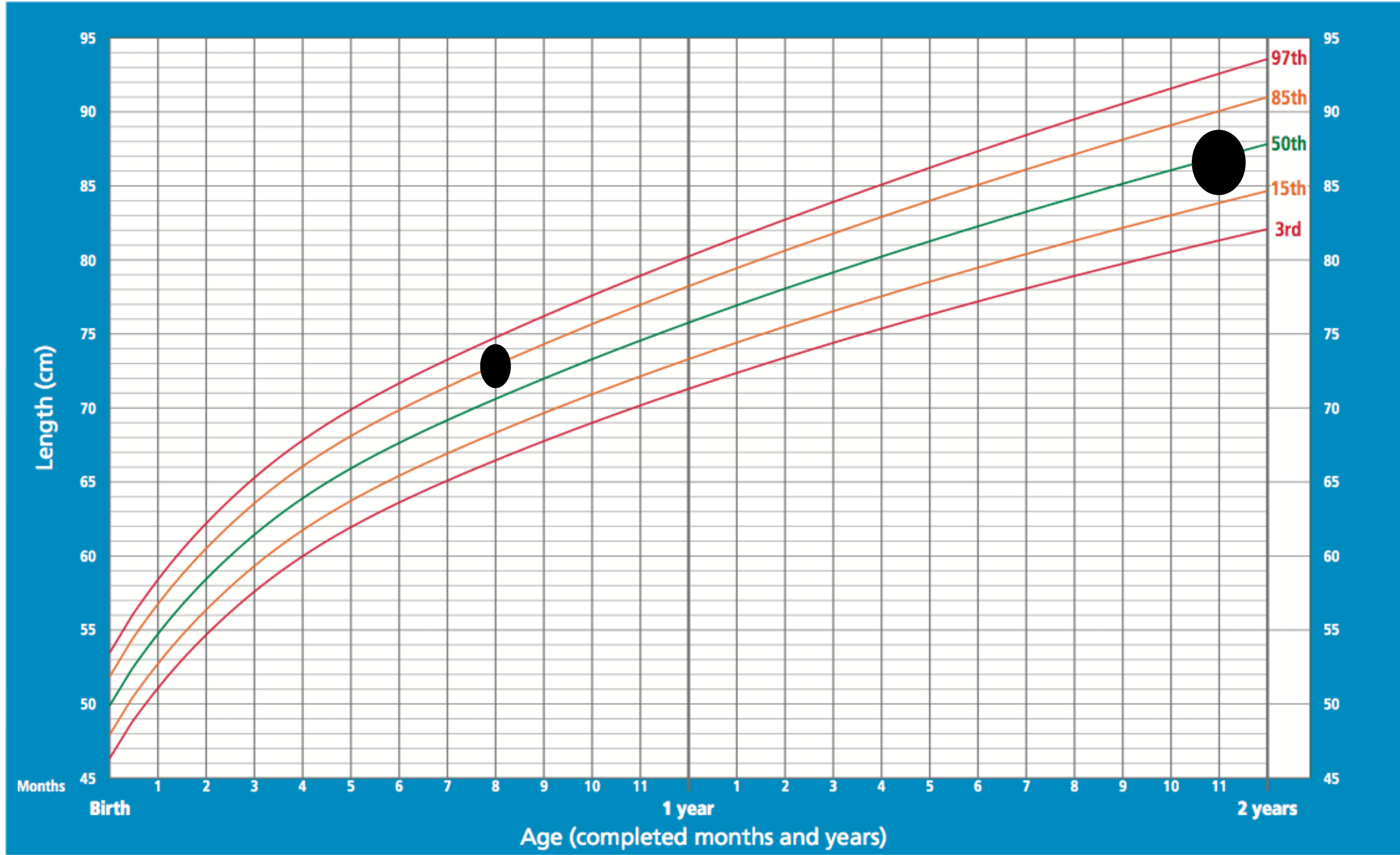
World Health Organization

WHO Child Growth Standards

# How to compare who is taller between elder and younger brothers?

## Standardized data



### Length-for-age BOYS
Birth to 2 years (percentiles)

World Health Organization
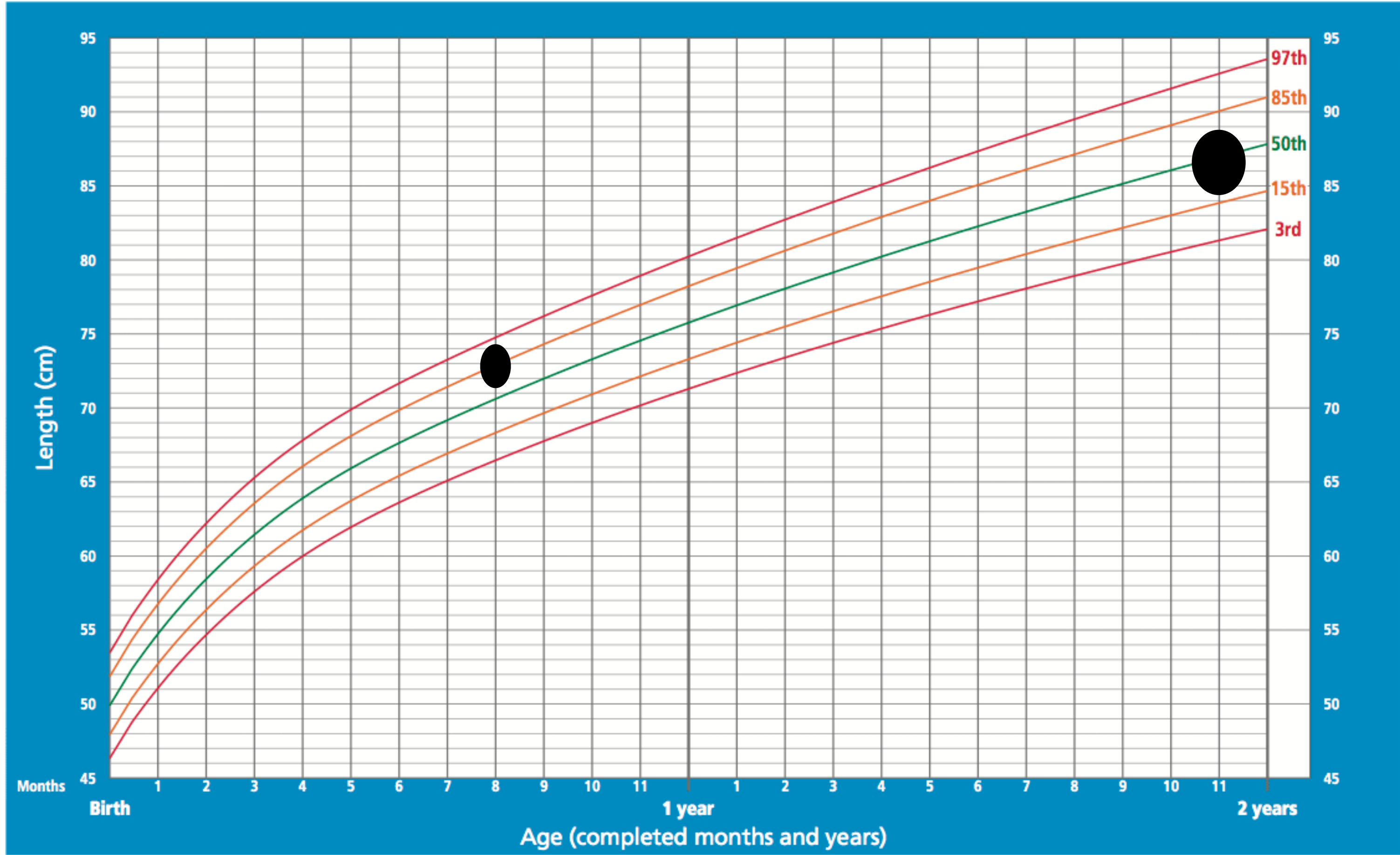
WHO Child Growth Standards

**Elder brother is 23 months old and 87 cm tall**
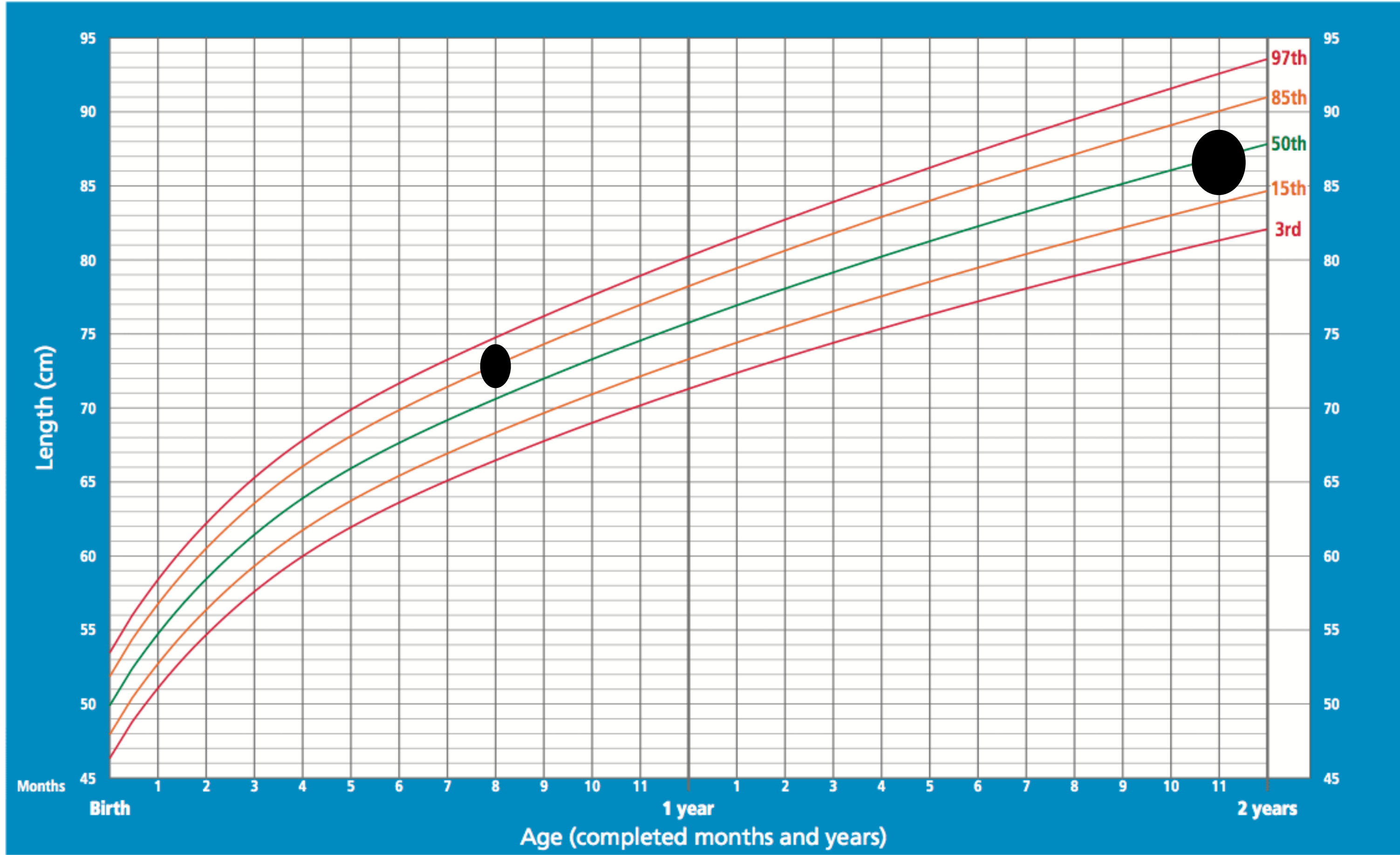**Younger brother is 8 month and 74 cm tall**

# How to compare who is taller between elder and younger brothers?

## Standardized data



Length-for-age BOYS
Birth to 2 years (percentiles)
World Health Organization

Elder brother is 23 months old and 87 cm tall
Younger brother is 8 month and 74 cm tall

Compare elder brother with his peer, he is at average.

# How to compare who is taller between elder and younger brothers?

## Standardized data



**Length-for-age BOYS**
Birth to 2 years (percentiles)
World Health Organization

WHO Child Growth Standards

Elder brother is 23 months old and 87 cm tall
Younger brother is 8 month and 74 cm tall

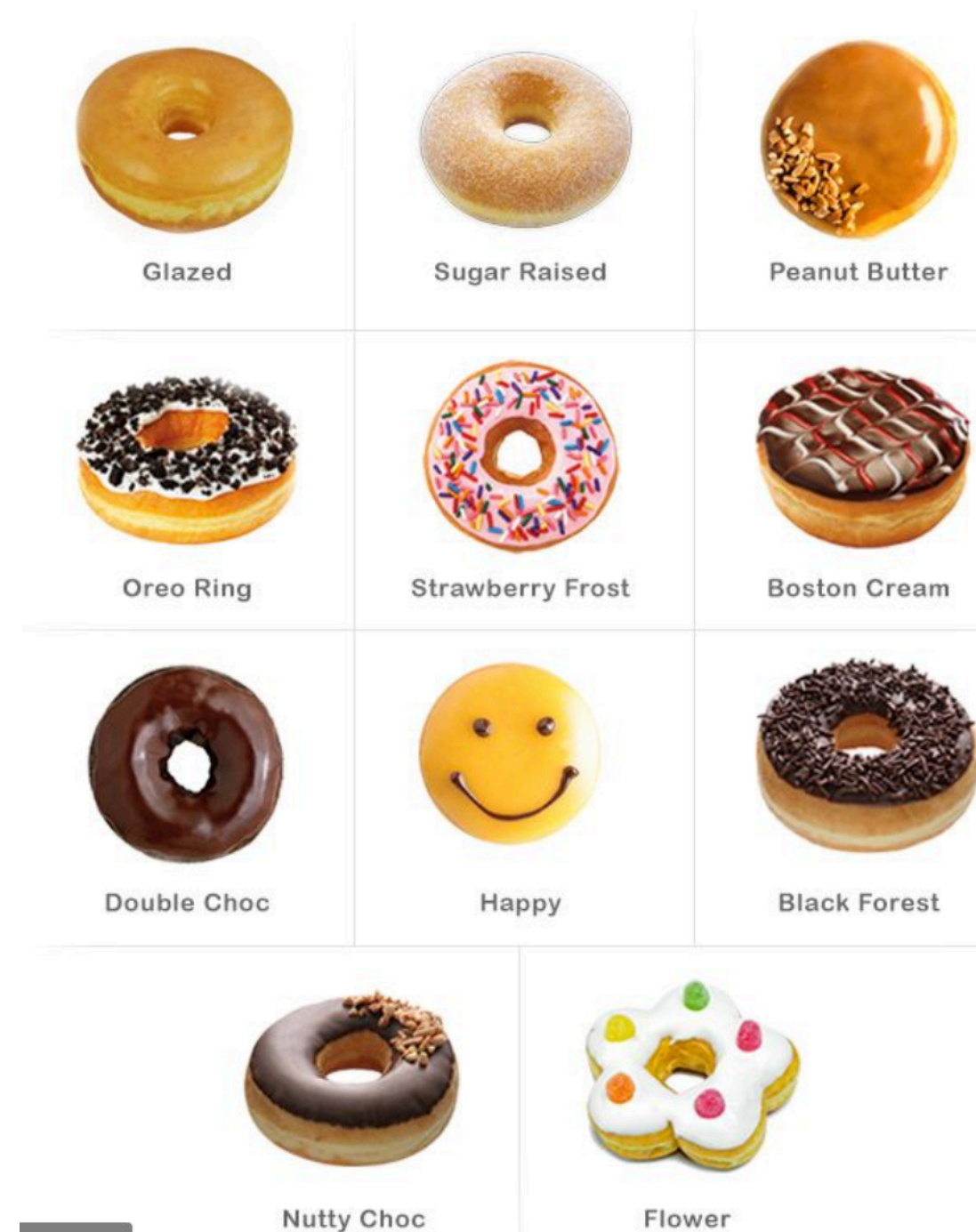Compare elder brother with his peer, he is at average.

As for younger brother, he is 1 sigma taller than his average peer.

# Bootstrapping data set

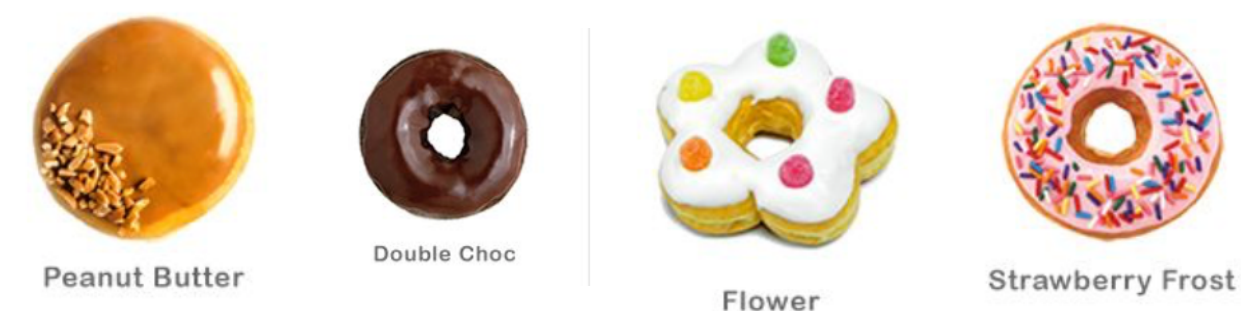## Please read the textbook from p.187 — p.190

### Menu of dunking donut



| | | |
|---|---|---|
| Glazed | Sugar Raised | Peanut Butter |
| Oreo Ring | Strawberry Frost | Boston Cream |
| Double Choc | Happy | Black Forest |
| Nutty Choc | Flower | |

**Customer 1:** Double Choc, Double Choc, Double Choc, Flower

**Customer 2:** Glazed, Double Choc, Strawberry Frost, Oreo Ring

**Customer 3:** Happy, Happy, Happy, Happy

**Customer 4:** Peanut Butter, Double Choc, Flower, Strawberry Frost

# Simulation on bootstrapping data sets

## Collecting all model parameters

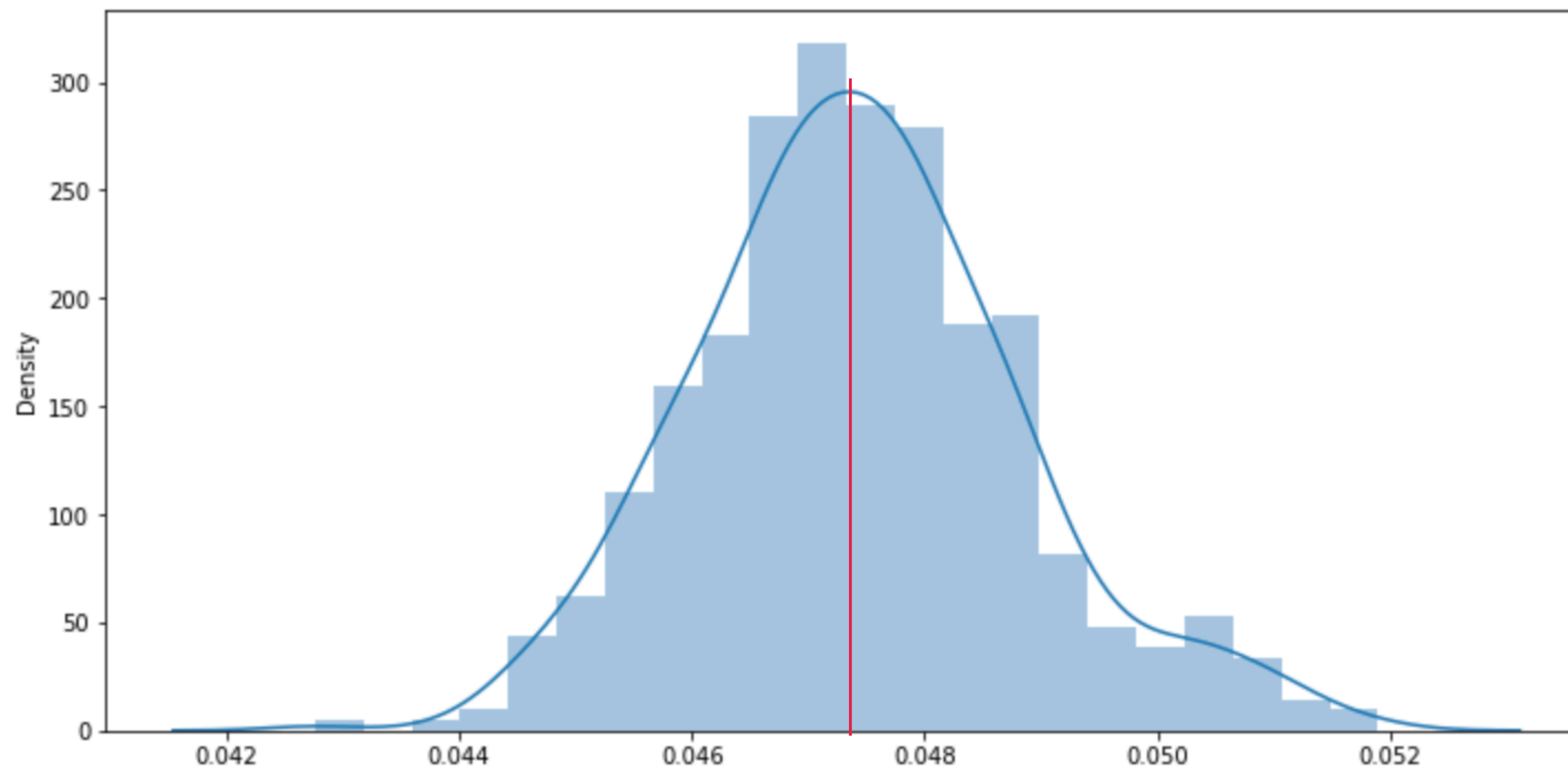# Question: where is the zero slope in tv versus sales?



0.00

```
print(np.mean(m_list), np.std(m_list))
```

0.0474348697630429 0.0014079599492549707

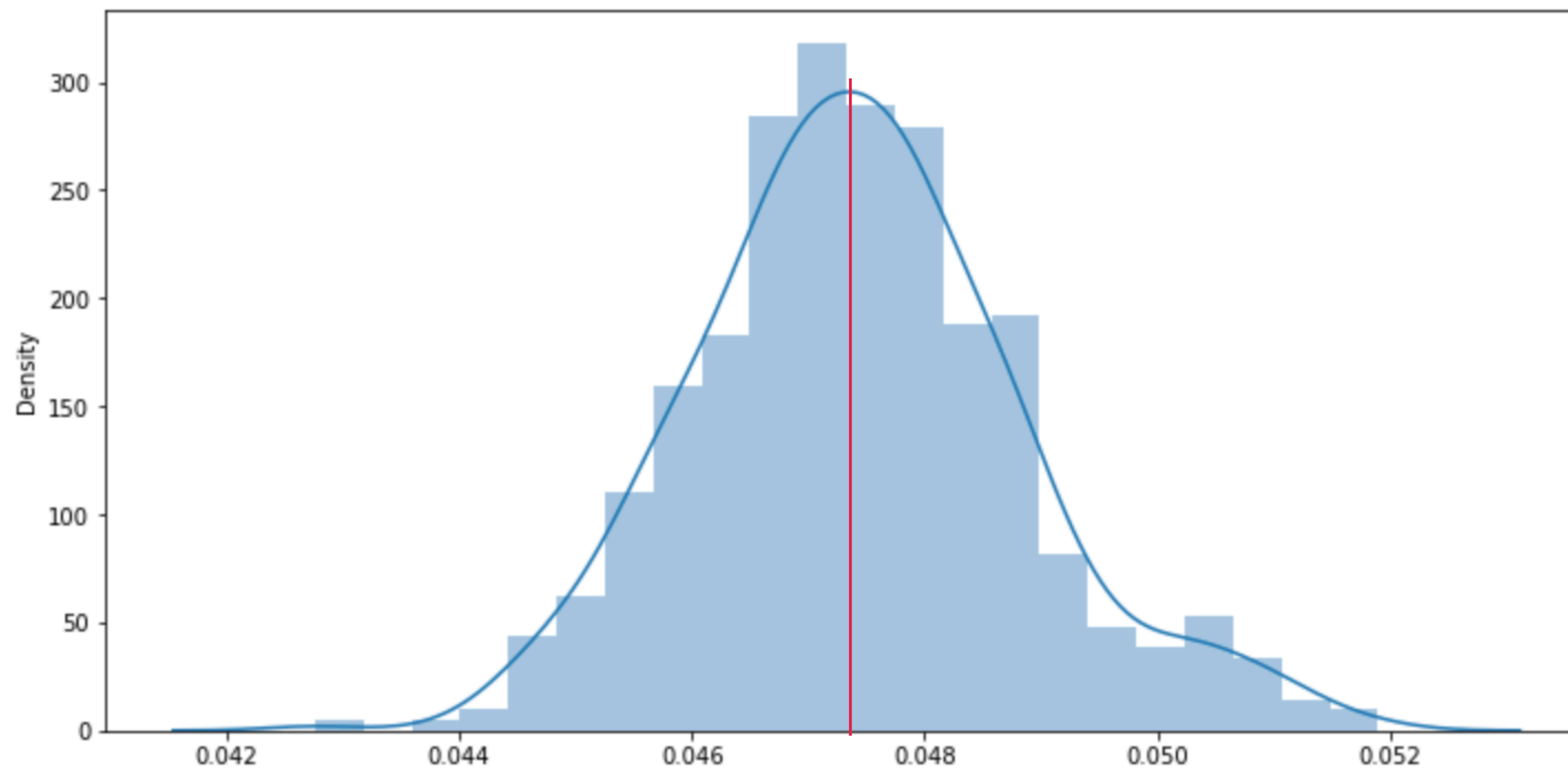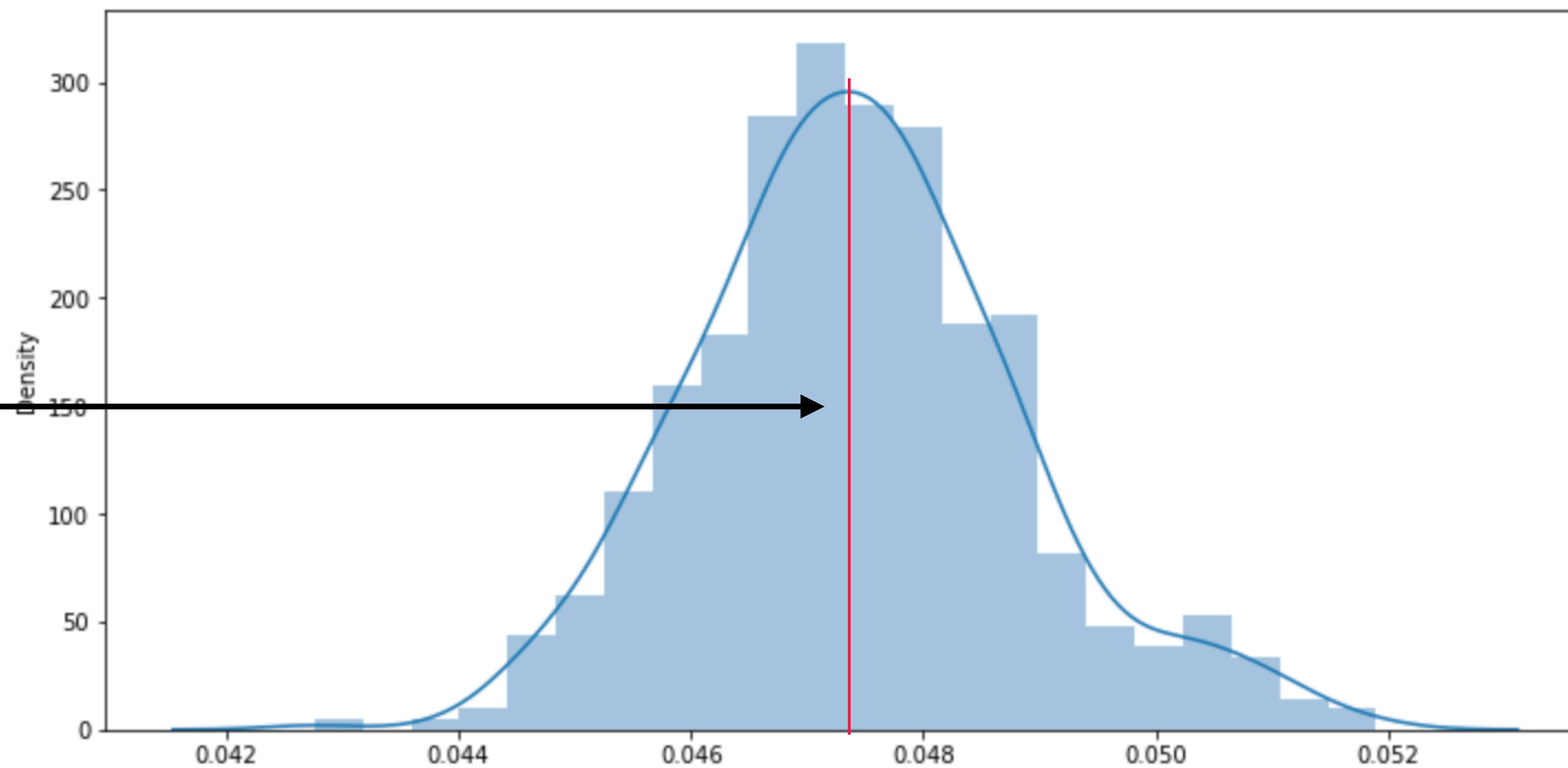# Question: where is the zero slope in tv versus sales?



0.00

```
print(np.mean(m_list), np.std(m_list))
```

0.0474348697630429 0.0014079599492549707

# Question: where is the zero slope in tv versus sales?

0.00



```
print(np.mean(m_list), np.std(m_list))
```
0.0474348697630429 0.0014079599492549707

# Question: where is the zero slope in tv versus sales?



0.00

```
print(np.mean(m_list), np.std(m_list))
```
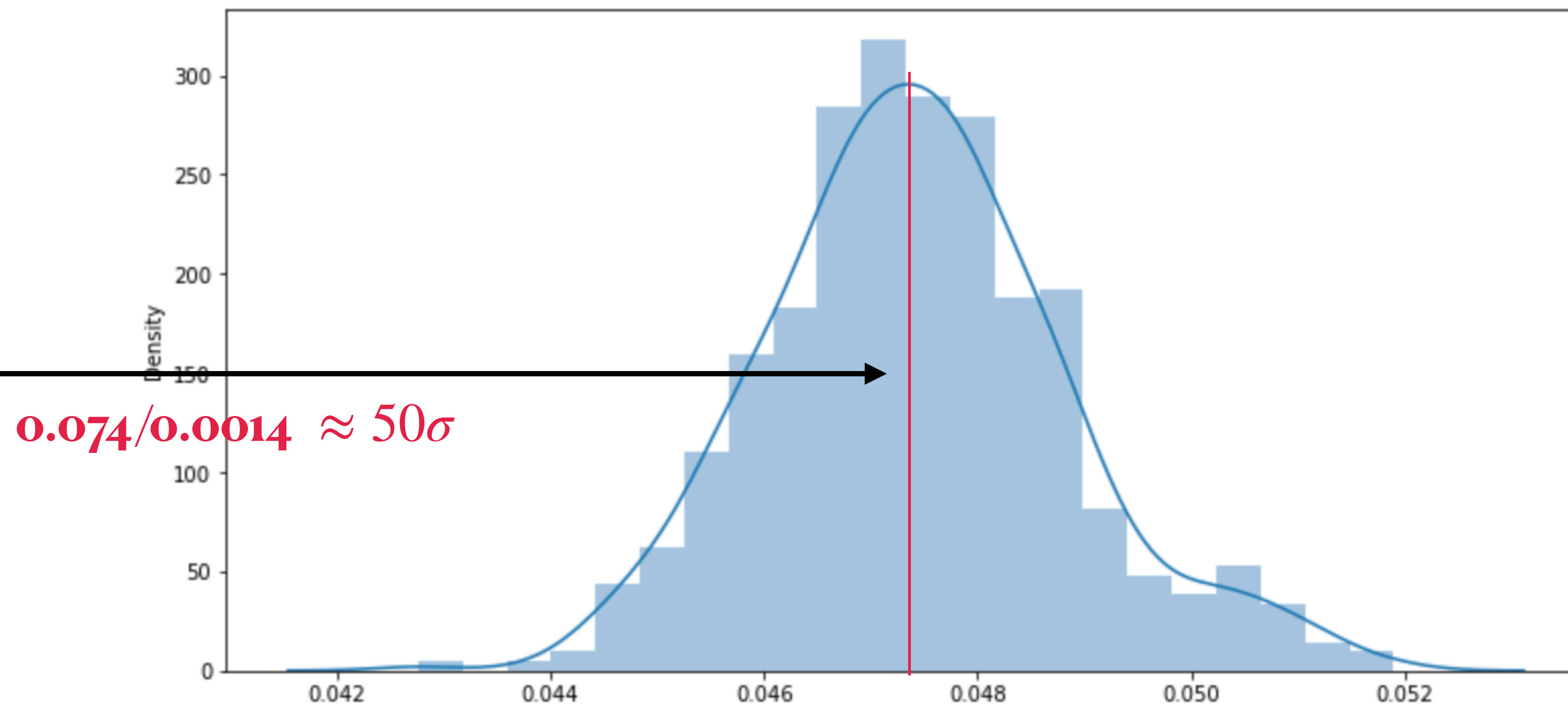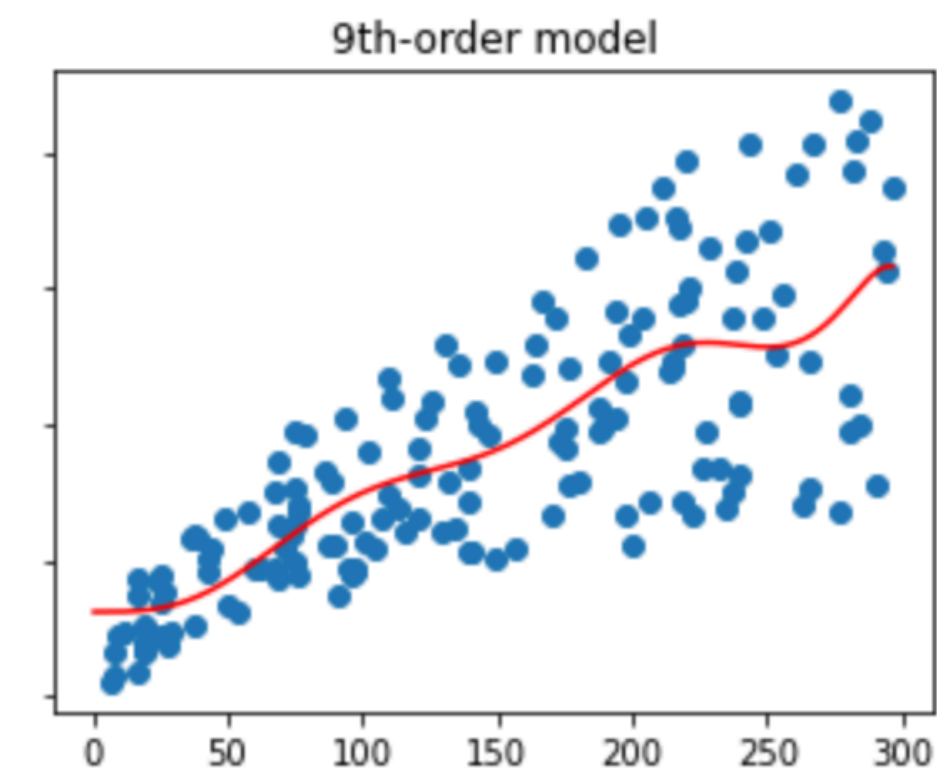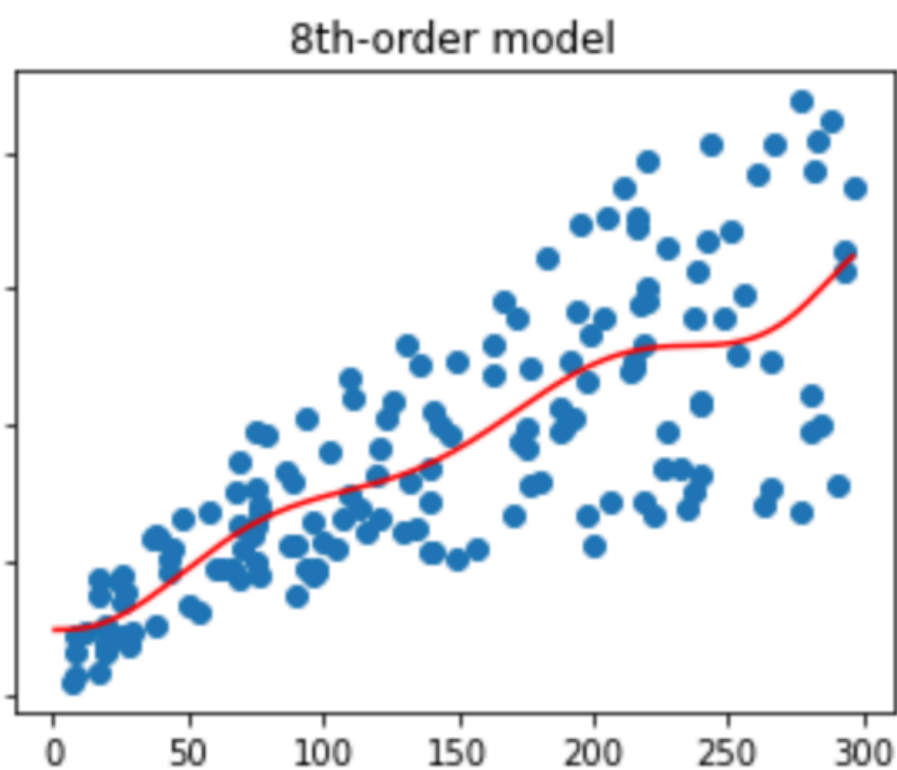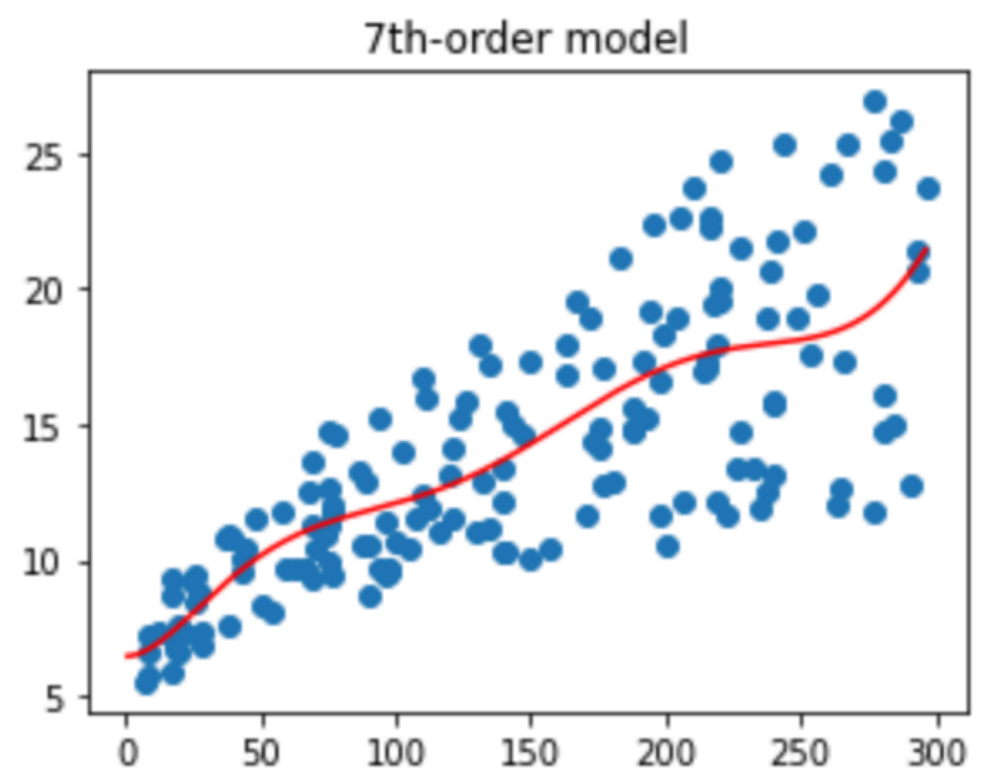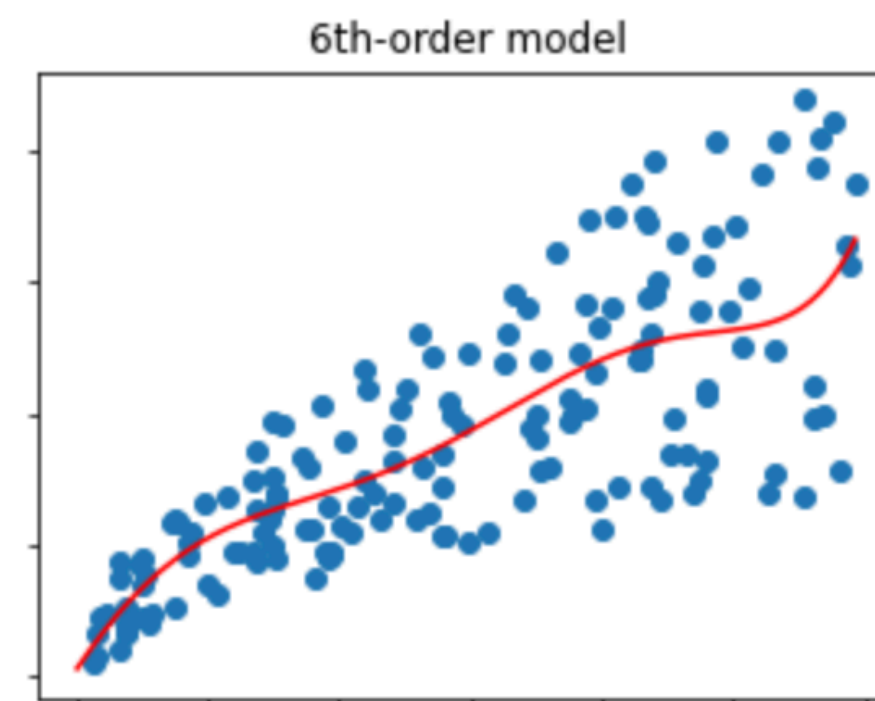
0.0474348697630429 0.0014079599492549707

# Question: where is the zero slope in tv versus sales?



$0.074/0.0014 \approx 50\sigma$
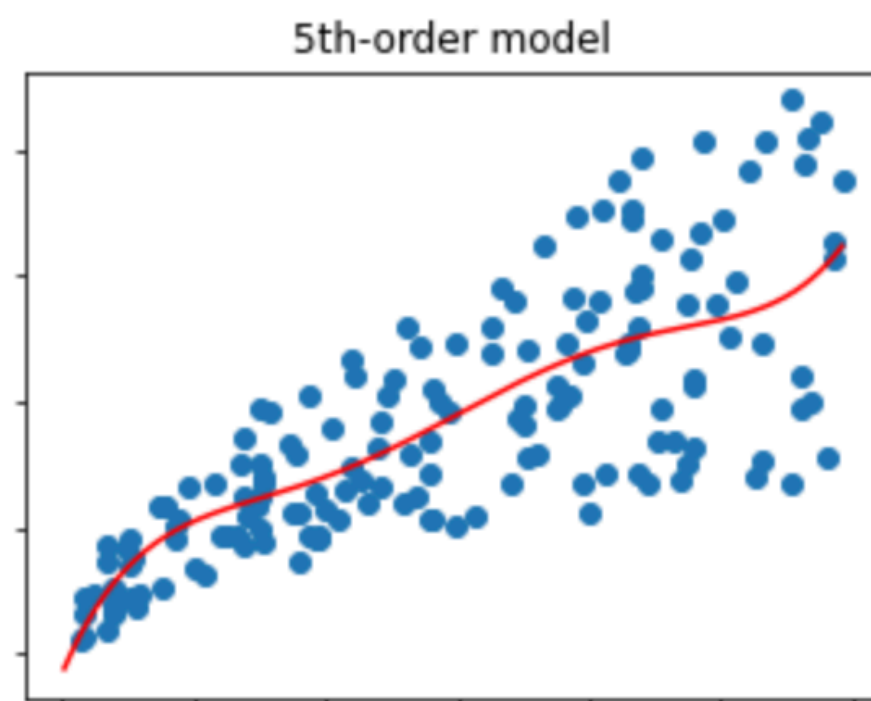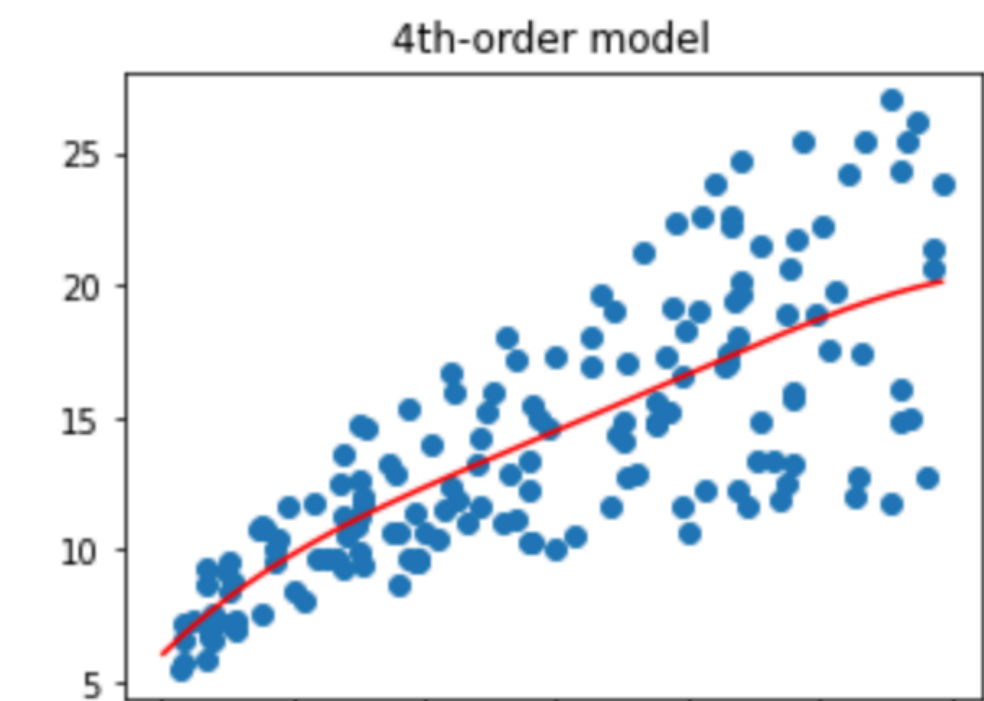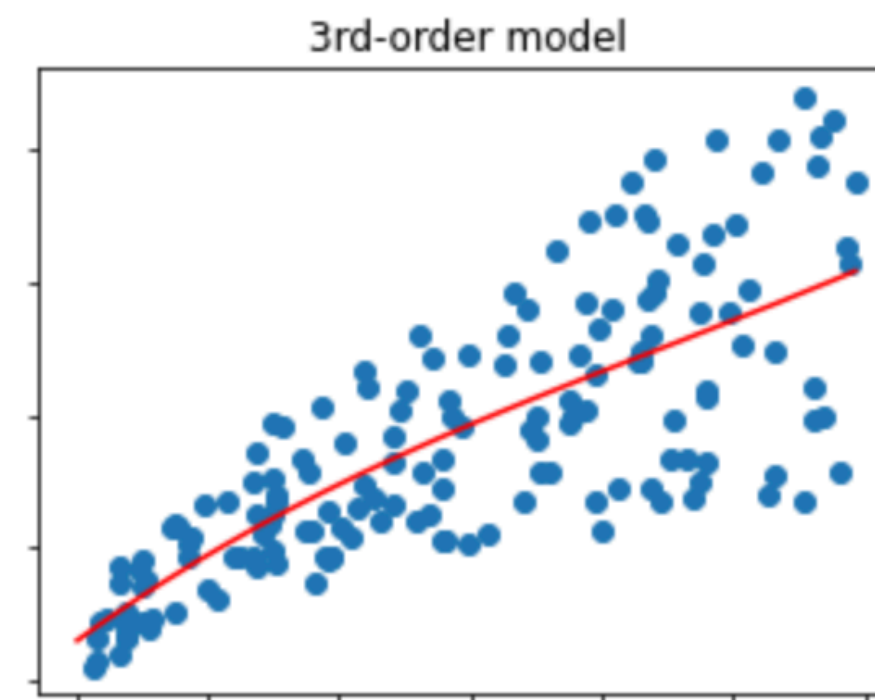
0.00

```python
print(np.mean(m_list), np.std(m_list))
```
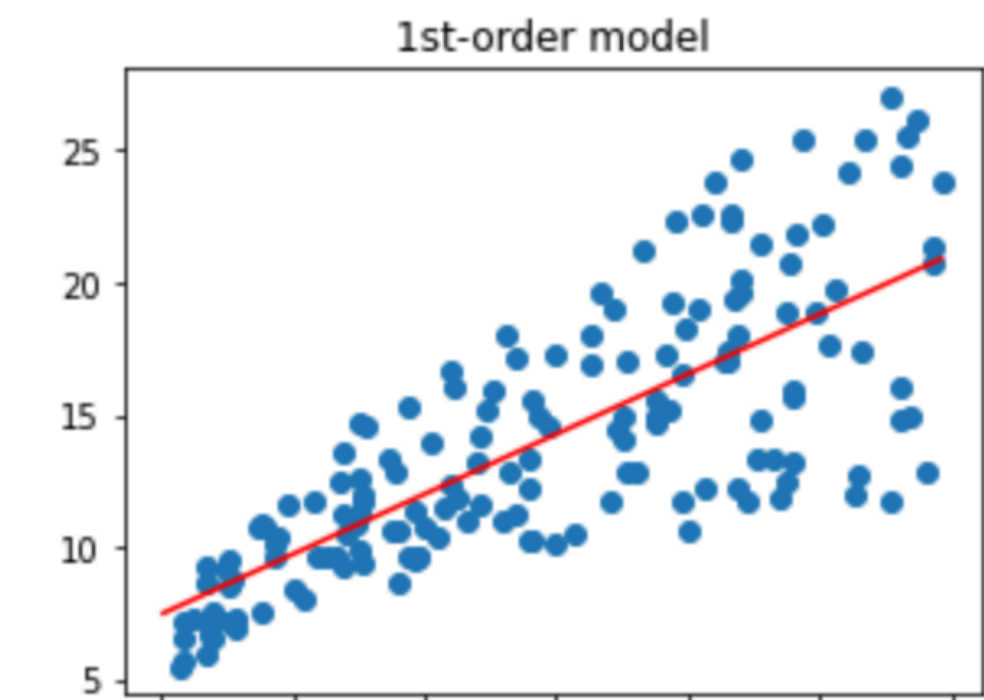0.0474348697630429 0.0014079599492549707

# You can try the same simulation to (radio,sales) and (newspaper, sales)

## It would be the next homework

# Beyond linear regression model

# How to choose which one would work best?



Error versus complexity