

# Prediction With Linear Regression

POLS 602

Dr. Mike Burnham  
Texas A&M Political Science

# Review

1. You have a bunch of data. The only thing you know about the data is that the mean of your sample is 50. A, observation from the sample is chosen at random and you have to guess its value. What is your best guess?
2. You have a bunch of data. The only thing you know about the data is that the mean of your sample is 50 and the median is 60. A, observation from the sample is chosen at random and you have to guess its value. What is your best guess?
3. In plain english, what is the standard deviation?
4. In plain english, what is Z in this formula?  $Z = \frac{x - \mu}{\sigma}$
5. Why would we want to use units of standard deviations to communicate how far something is from the mean?

# Review

$$Z = \frac{x - \mu}{\sigma}$$

$$Z = \frac{x - \mu}{\sigma}$$

$$Z_i^X = \frac{(X_i - \bar{X})}{sd(X)}$$

$$Z_{\{i\}^{\wedge}\{X\}} = \frac{(X_i - \bar{X})}{sd(X)}$$

# Predictive Modeling

# Why Predictive Modeling

- Uncover new causal mechanisms and generate new hypotheses
- Discover new measures and compare different operationalizations
- Improvement to existing explanatory models
- Assessing the distance between theory and practice
- Compare competing theories
- Quantifying the level of predictability of measurable phenomena by creating benchmarks

Source:

To Explain or to Predict? Galit Shmueli, 2010

What Can We Learn from Predictive Modeling? Cranmer and Desmarais, 2017

# Notation and Vocabulary

- **X**: Predictors, features
- **Y**: Outcomes, observed outcomes
- **Model**: a mathematical expression that represents a relationship between data
- **Data Generating Process**: The mechanisms or processes in the real world that produce the data we observe

# M System Metro

wmata.com  
Customer Information Service: 202 637-7000  
TTY Phone: 202 638-3788  
Metro Transit Police: 202 960-2121

**Legend**

- Red Line • Glenmont / Shady Grove
- Orange Line • New Carrollton / Vienna
- Blue Line • Franconia-Springfield / Largo Town Center
- Green Line • Branch Ave / Greenbelt
- Yellow Line • Huntington / Fort Totten
- Silver Line • Wiehle-Reston East / Largo Town Center

**Station Features**

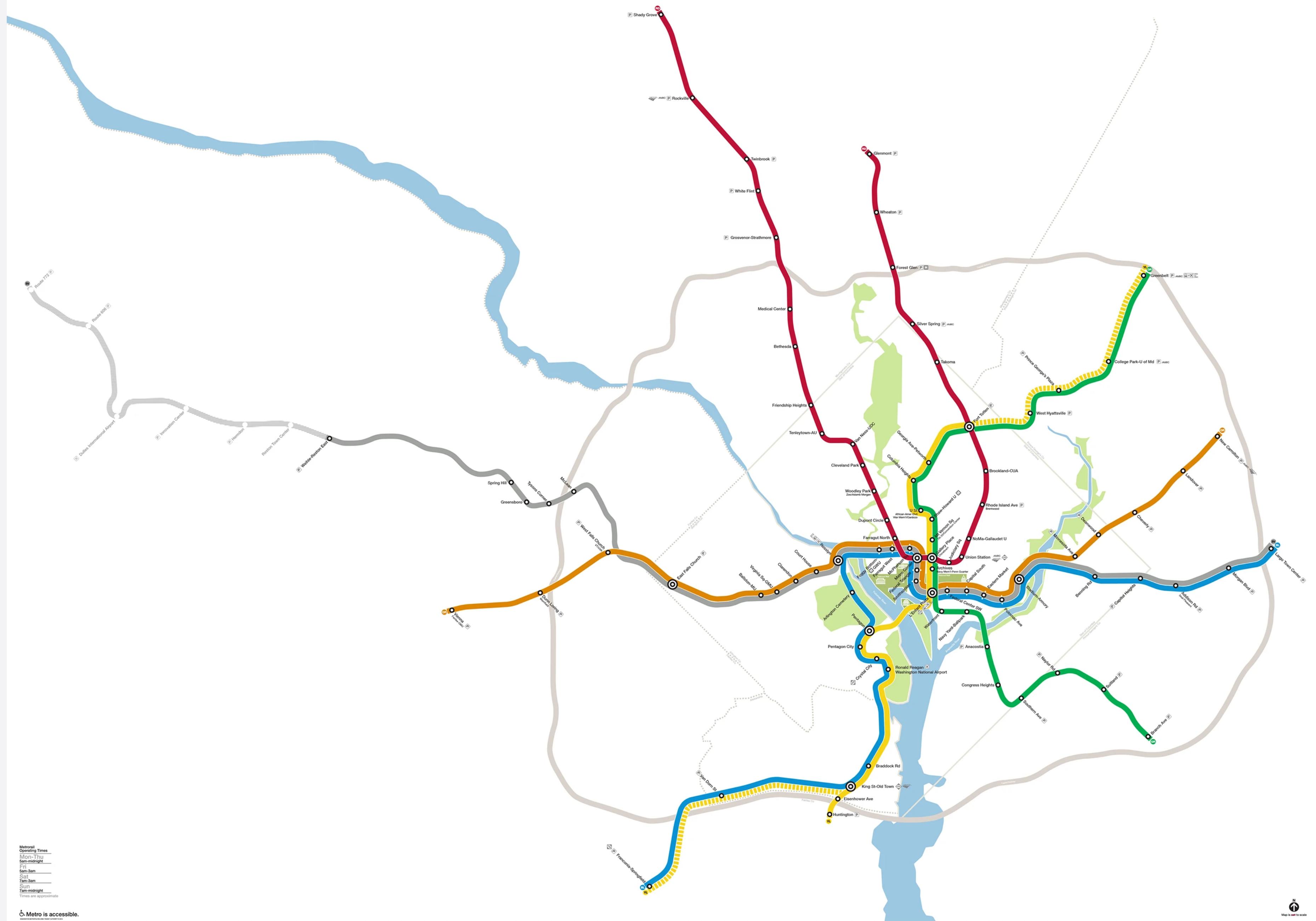
- Bus Stop
- Parking
- Hospital
- Airport

**Connecting Rail Systems**

- Amtrak
- Metrorail

**Tram/Streetcar**

- Under Construction
- Full-Time Service
- Rush-Only Service: Monday-Friday  
6:30am - 9:00am 3:30pm - 6:00pm
- Station in Service



N

Map is not to scale

# M System Map

wmata.com  
Information: 202-GO-METRO | TTY: 202-962-2033  
Metro Transit Police: 202-962-2121 | Text: MYMTPD (696873)

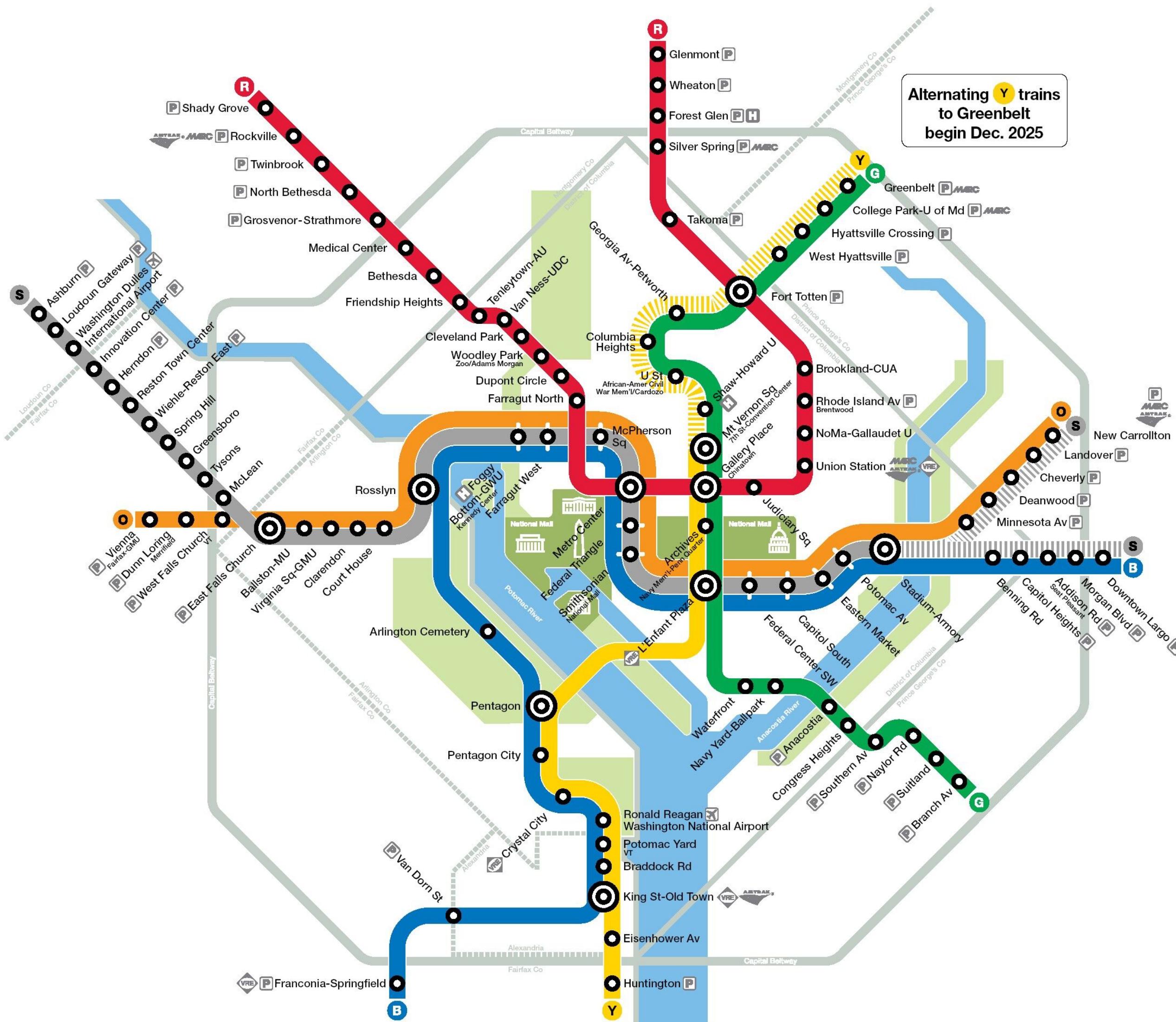
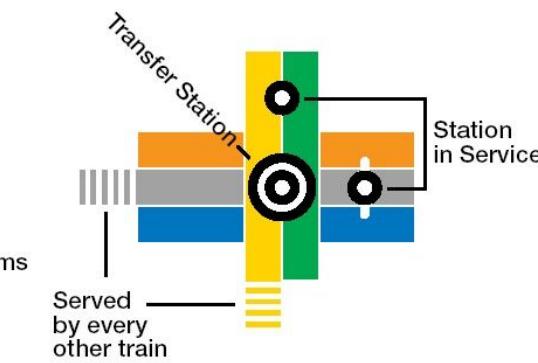
## Terminal stations

- R Red Line • Glenmont / Shady Grove
- O Orange Line • New Carrollton / Vienna
- B Blue Line • Franconia-Springfield / Downtown Largo
- G Green Line • Branch Av / Greenbelt
- Y Yellow Line • Huntington / Greenbelt
- S Silver Line • Ashburn / Downtown Largo & New Carrollton

## Station Features

- Parking
- Hospital
- Airport

## Connecting Rail Systems

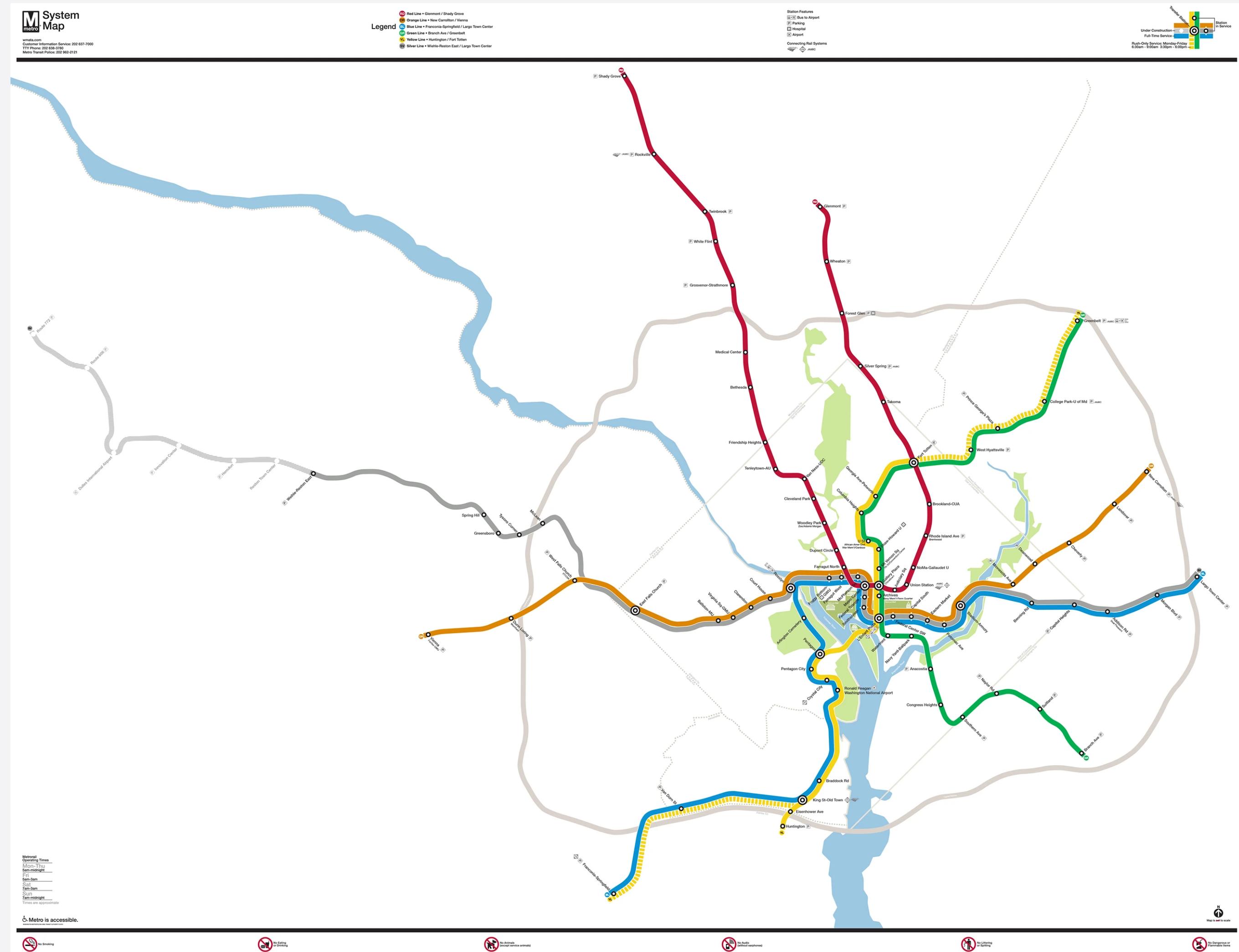
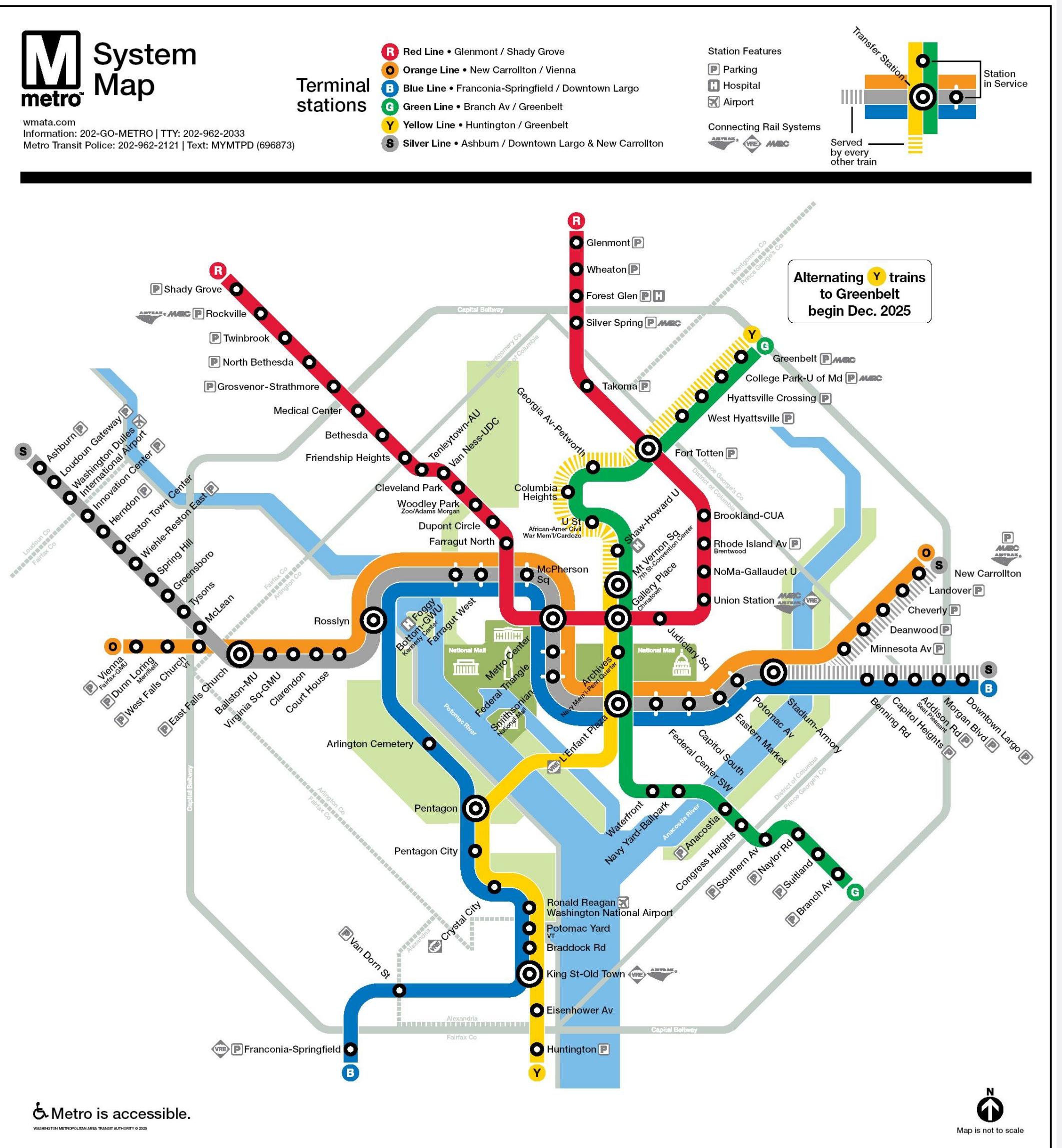


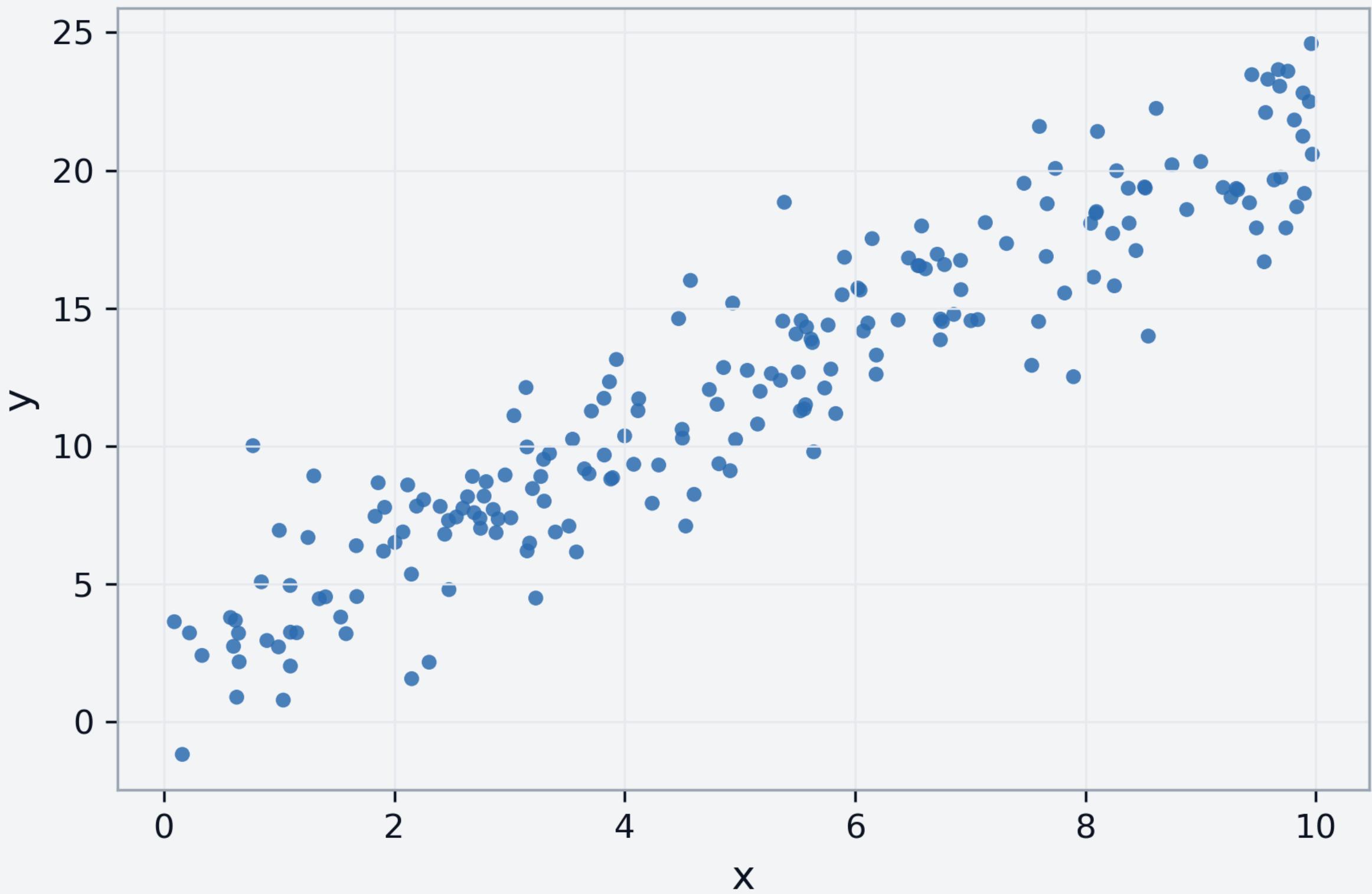
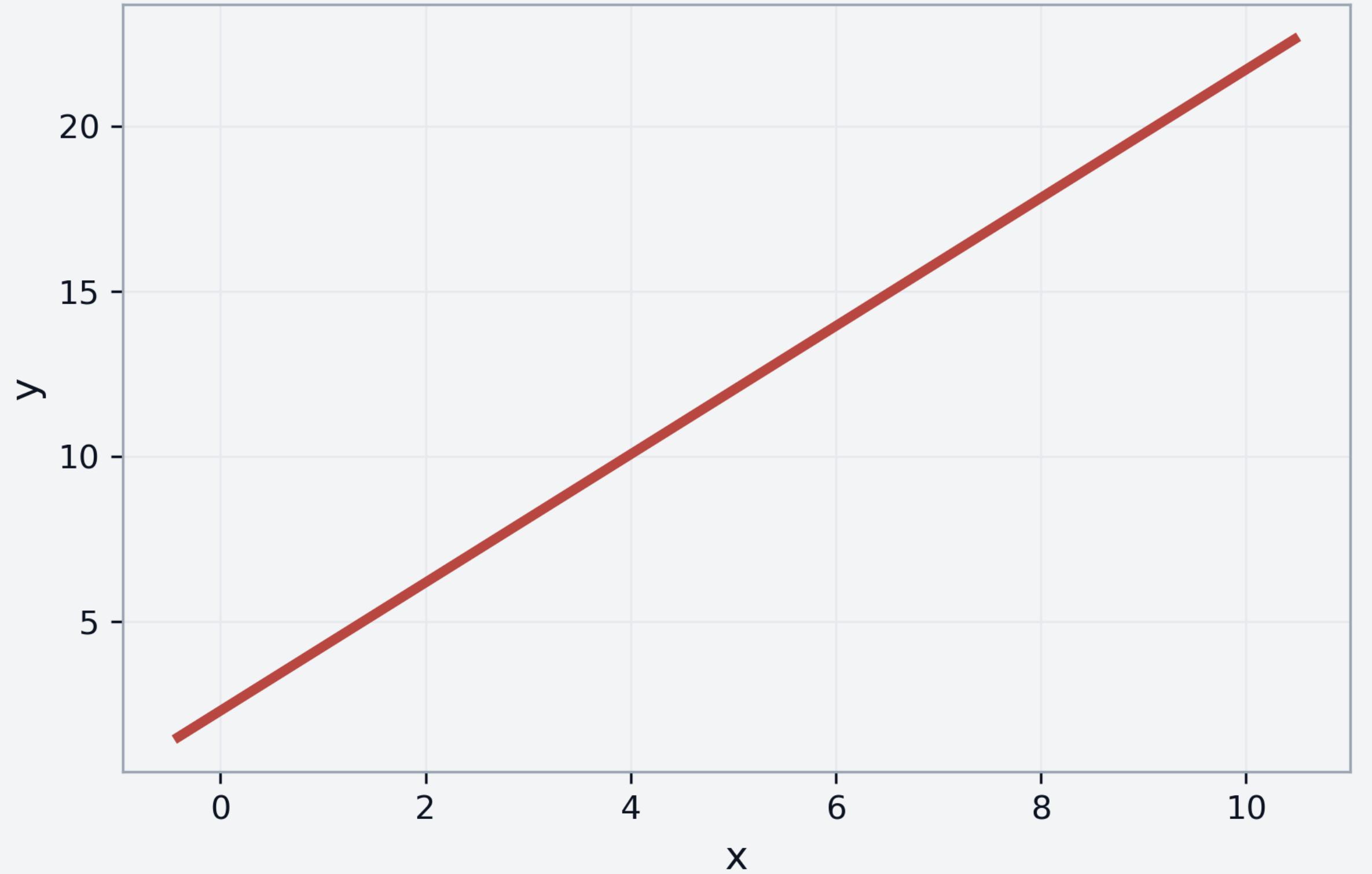
Metro is accessible.

WASHINGTON METROPOLITAN AREA TRANSIT AUTHORITY © 2025



Map is not to scale





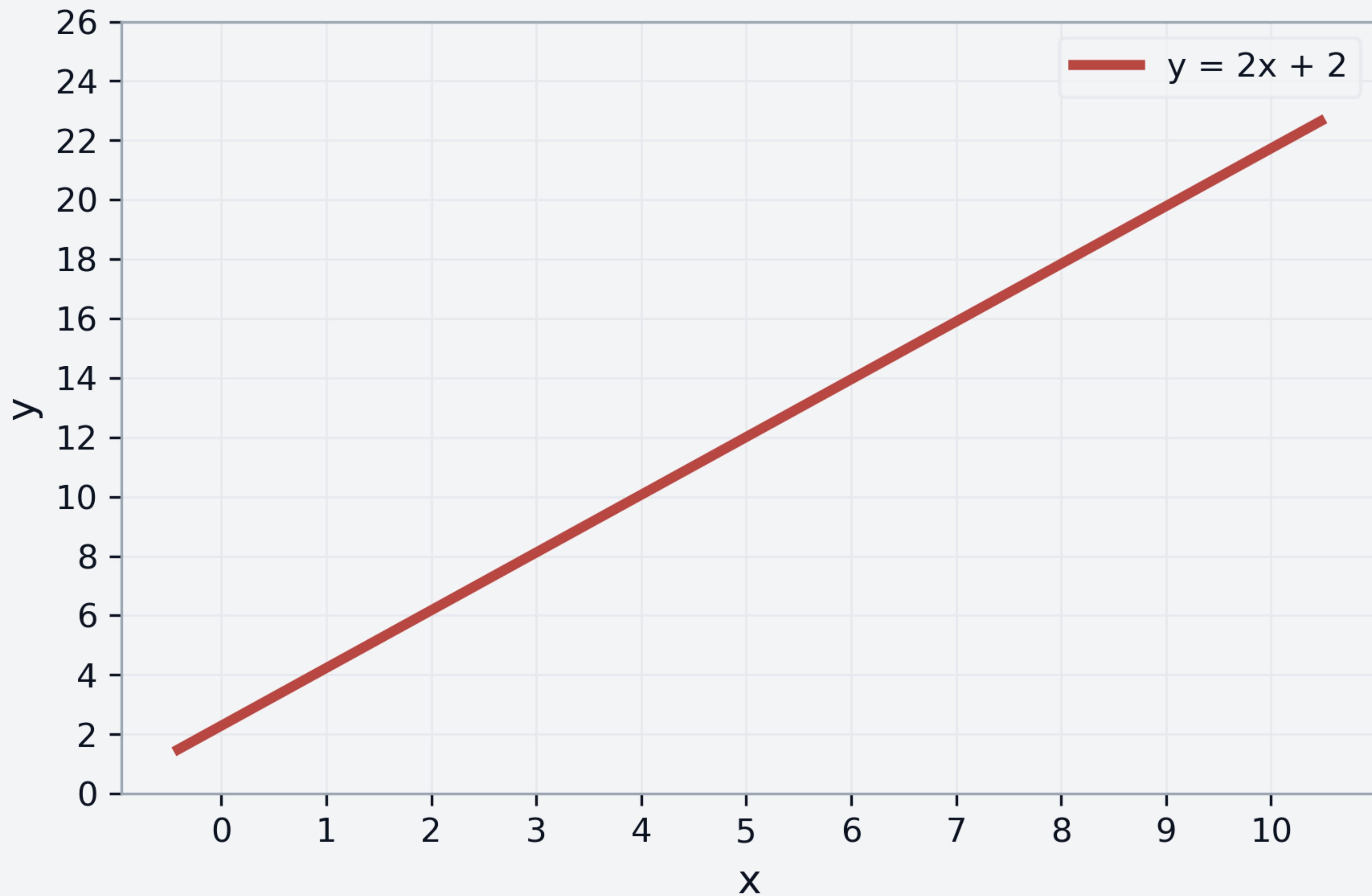
# Linear Regression

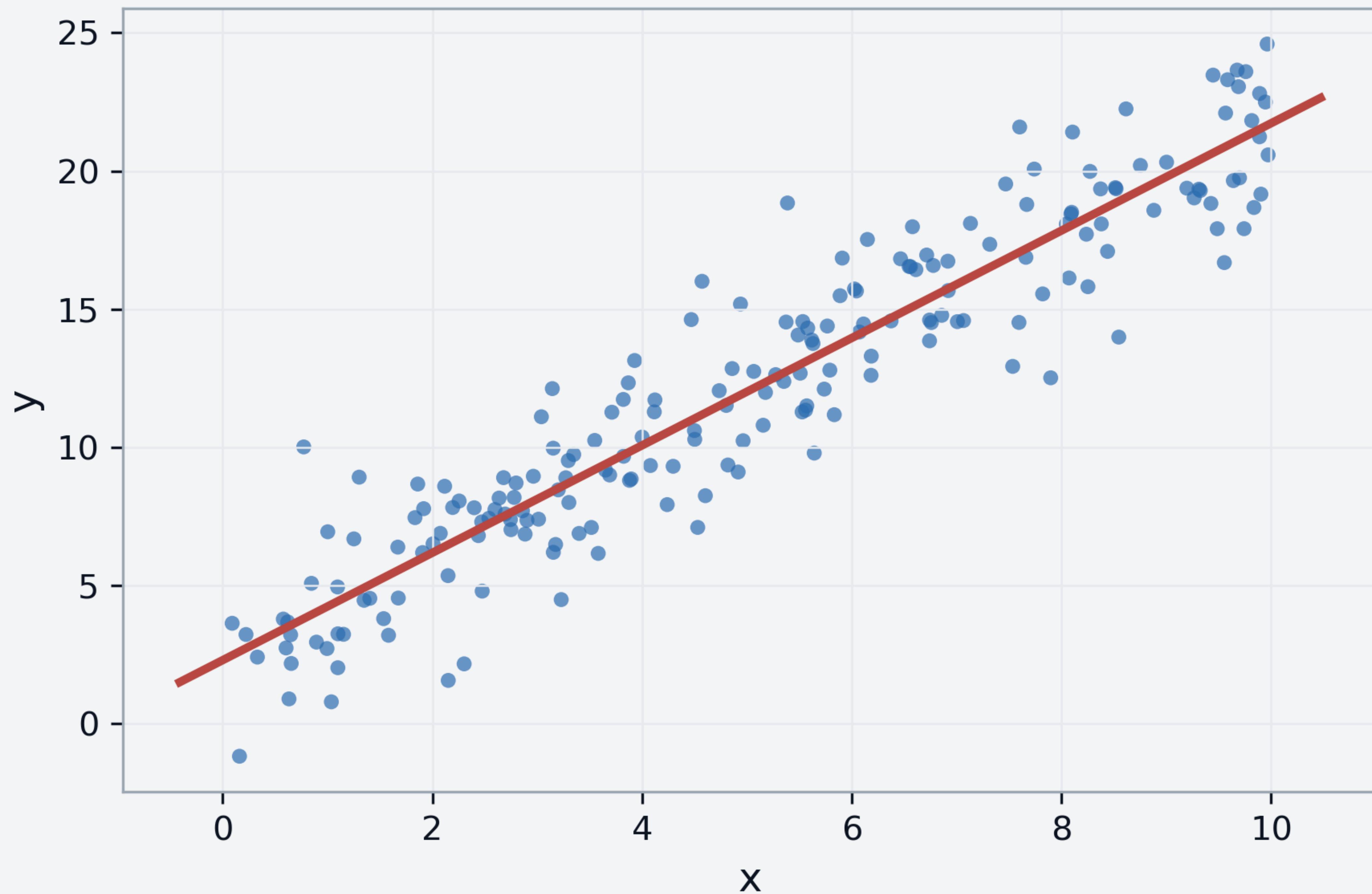
$$y = mx + b$$

$$y = mx + b$$

Slope

Intercept





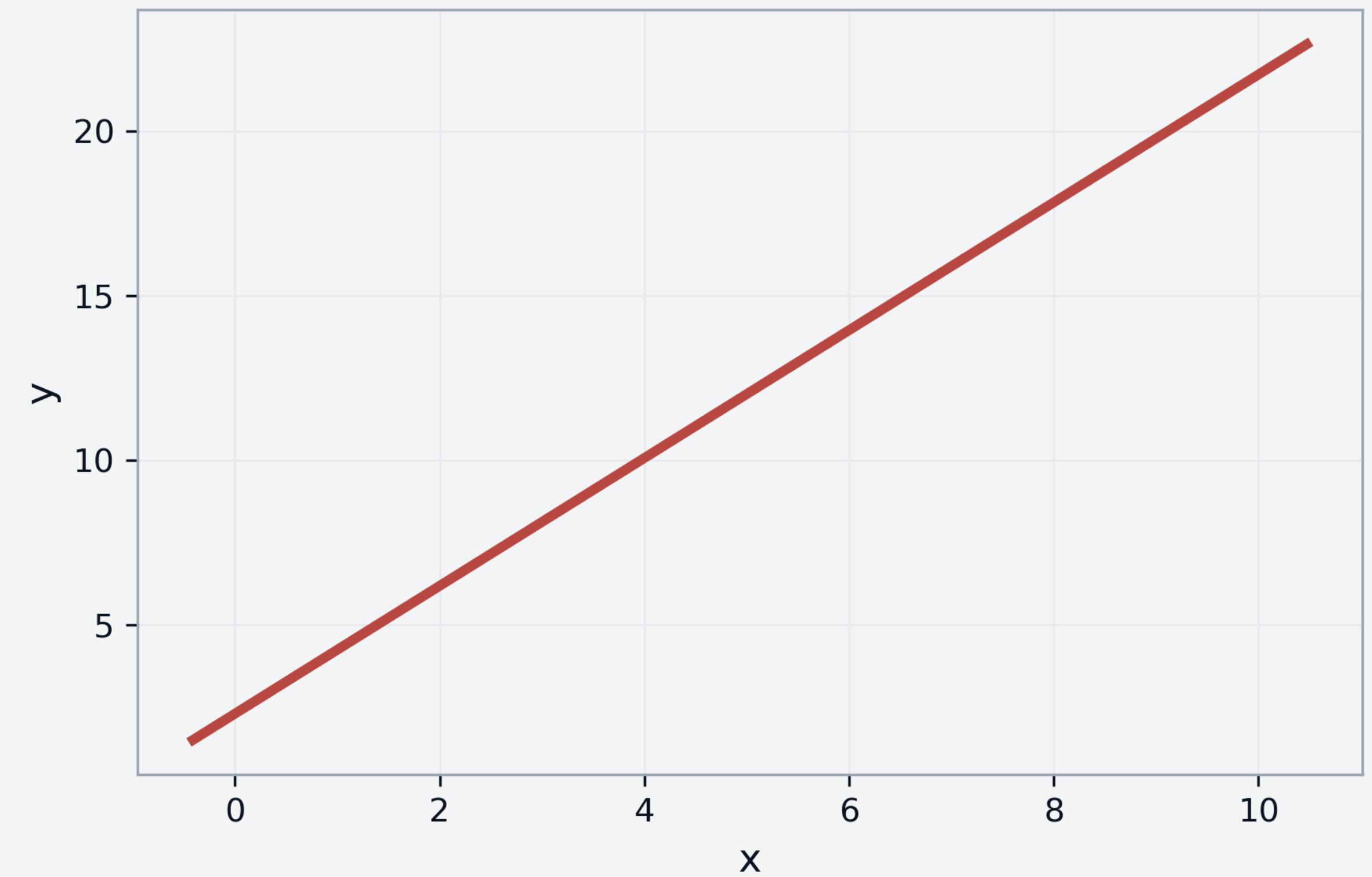
# Estimated Values

$$\hat{Y}_i = \hat{\alpha} + \hat{\beta}X_i$$

$$\hat{Y}_i = \hat{\beta}_0 + \hat{\beta}_1 X_i$$

- $\hat{Y}$ : Predicted values, predicted outcomes, fitted values. The values of Y predicted by a model. \hat{Y}
- $\hat{\alpha}$ : The intercept
- $\hat{\beta}$ : The slope coefficient

# Estimated Values



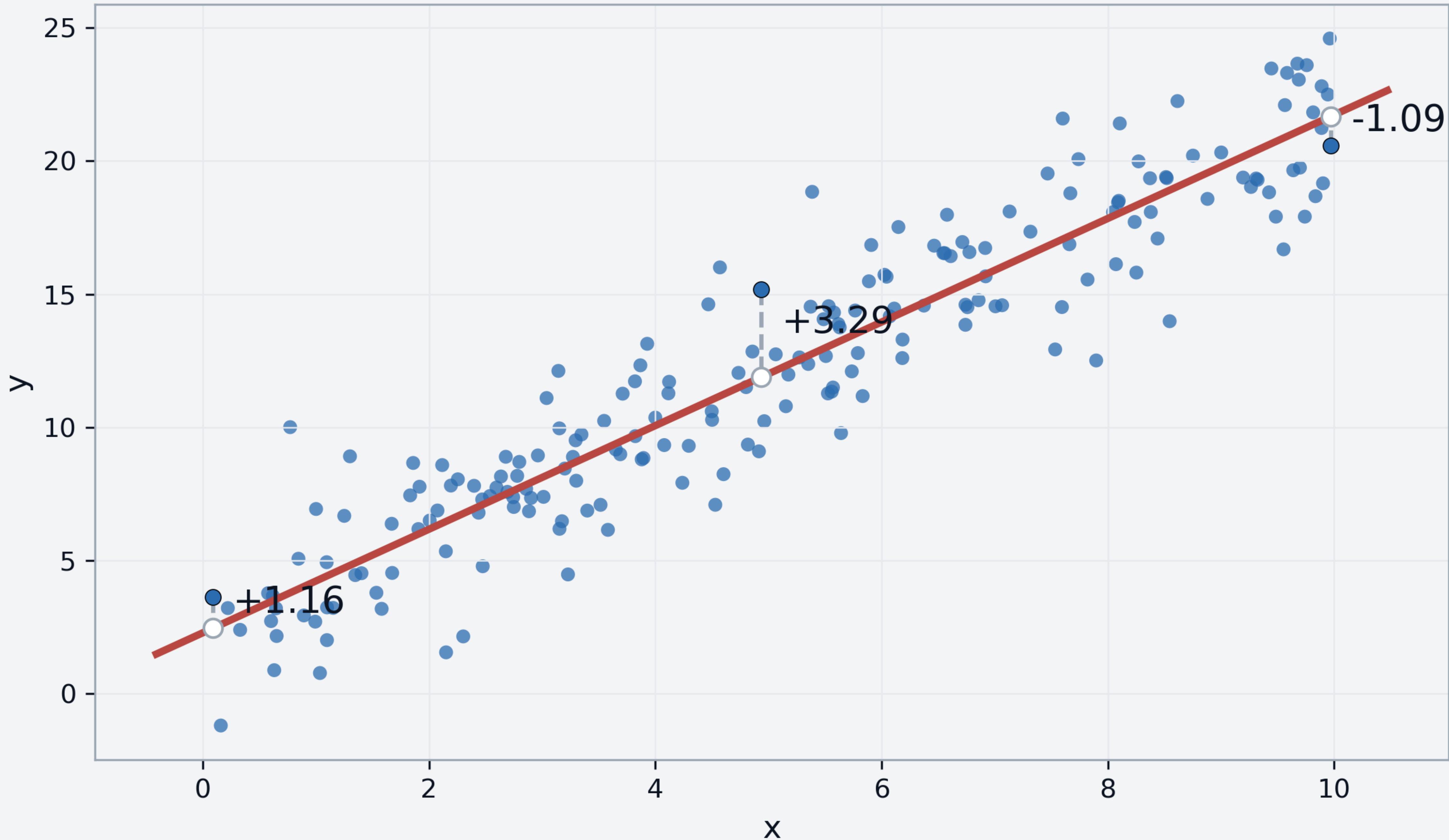
$$\hat{Y}_i = 2 + 2X_i$$

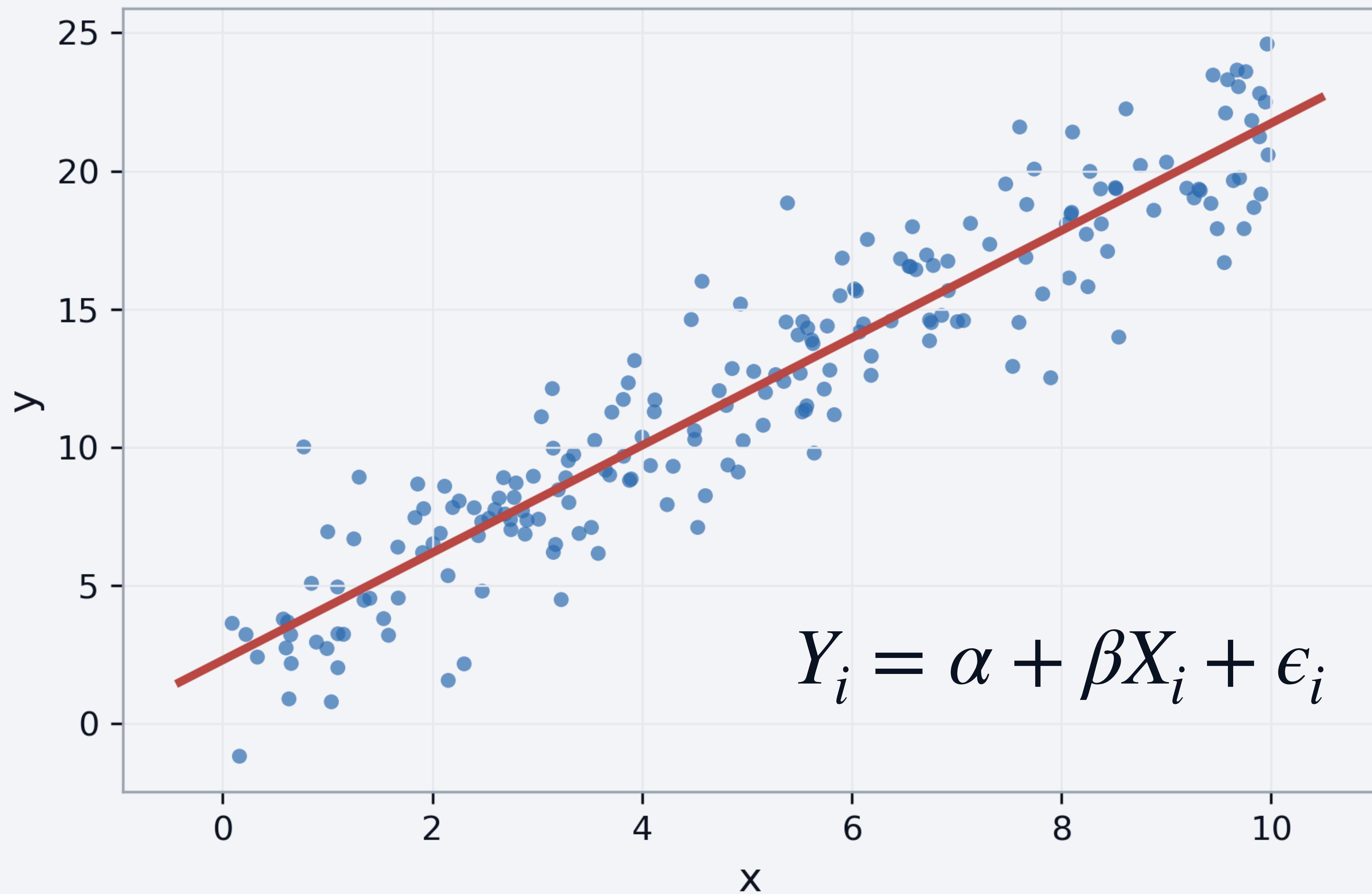
# The Linear Regression Model

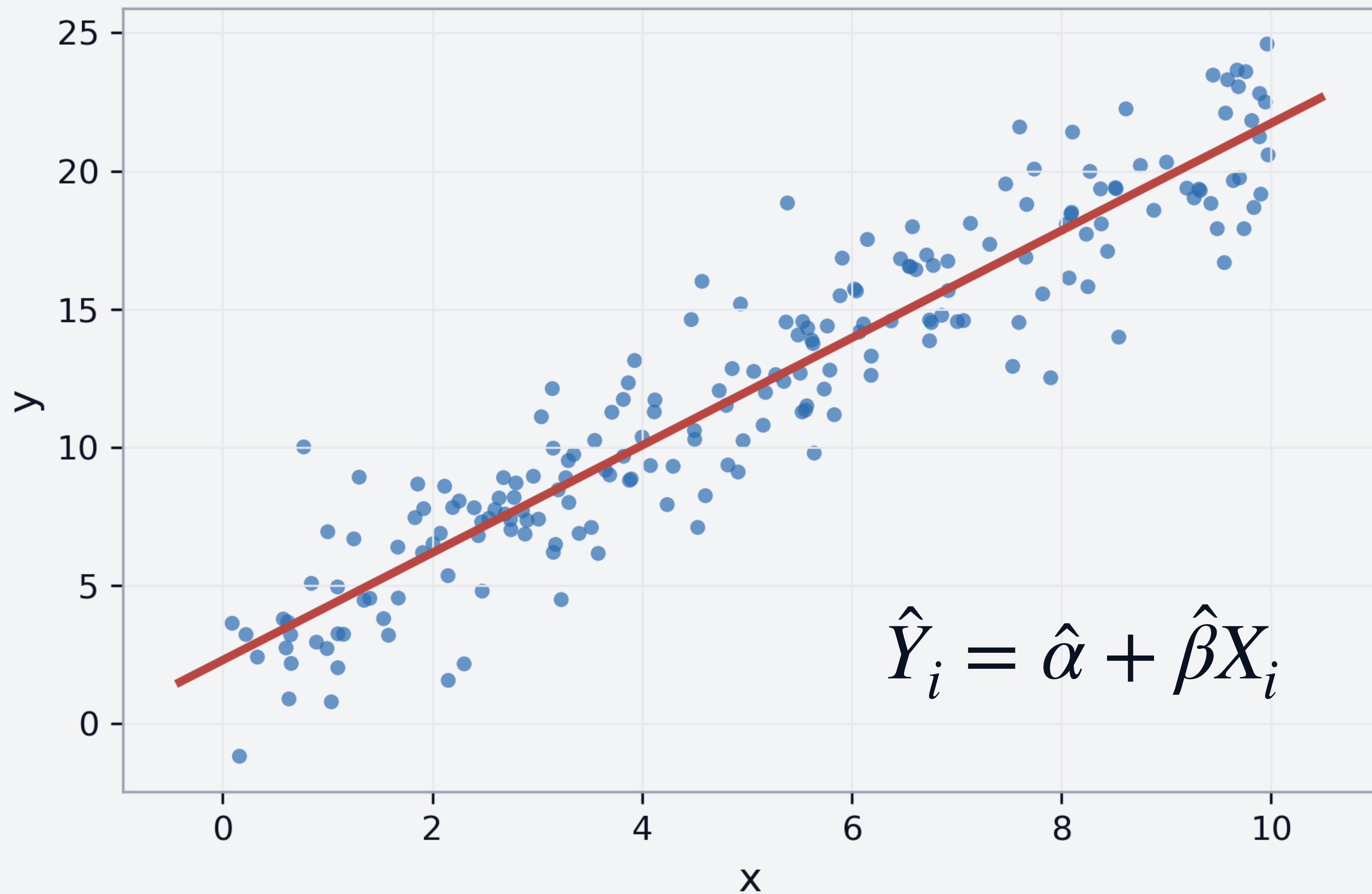
$$Y_i = \alpha + \beta X_i + \epsilon_i$$

- Notice no hats!
- The new term  $\epsilon_i$  is the prediction error, the error term, or the residuals
- $\epsilon_i = Y_i - \hat{Y}_i$

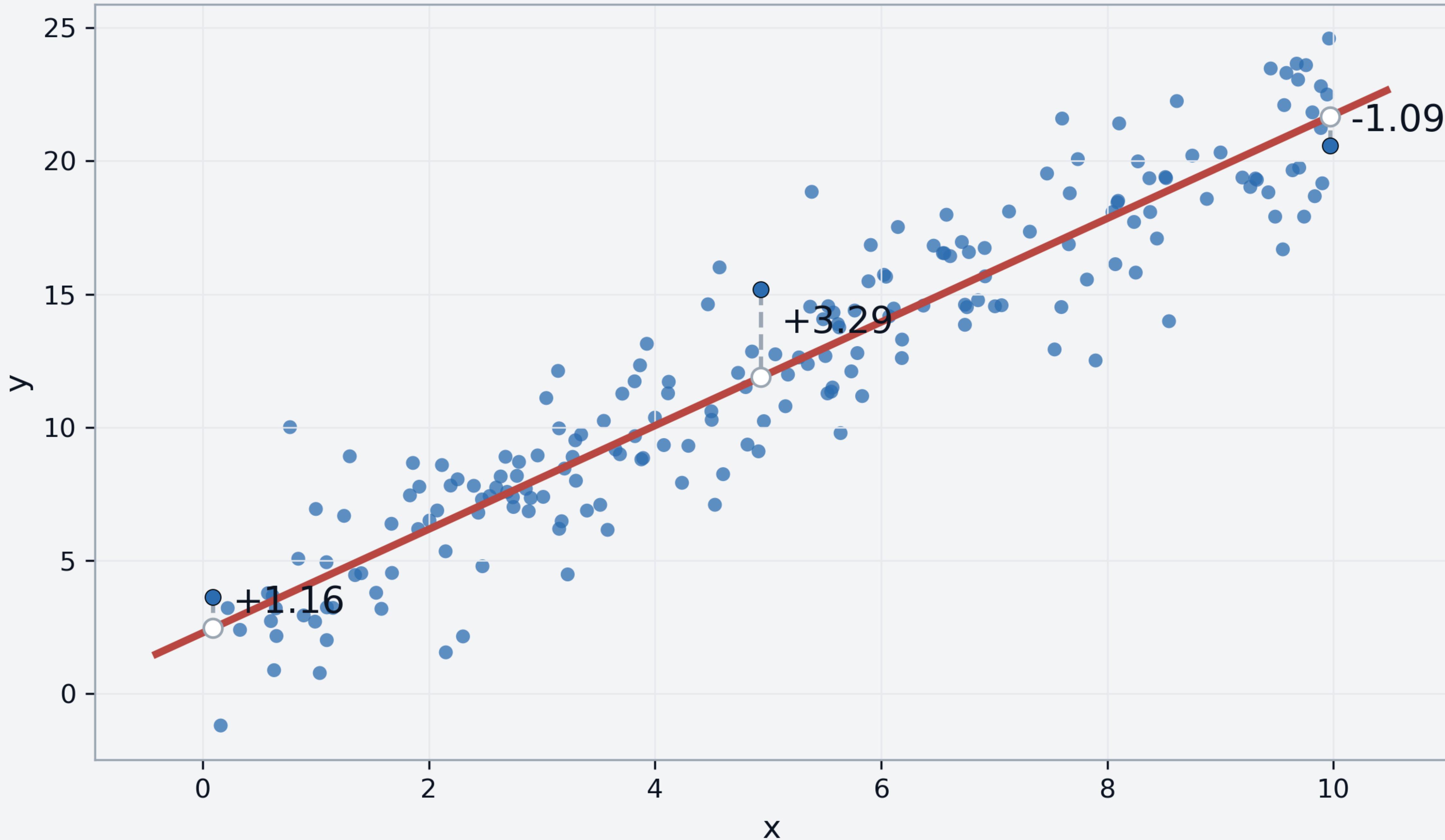
# Regression Errors: dashed lines show residuals ( $y - \hat{y}$ )



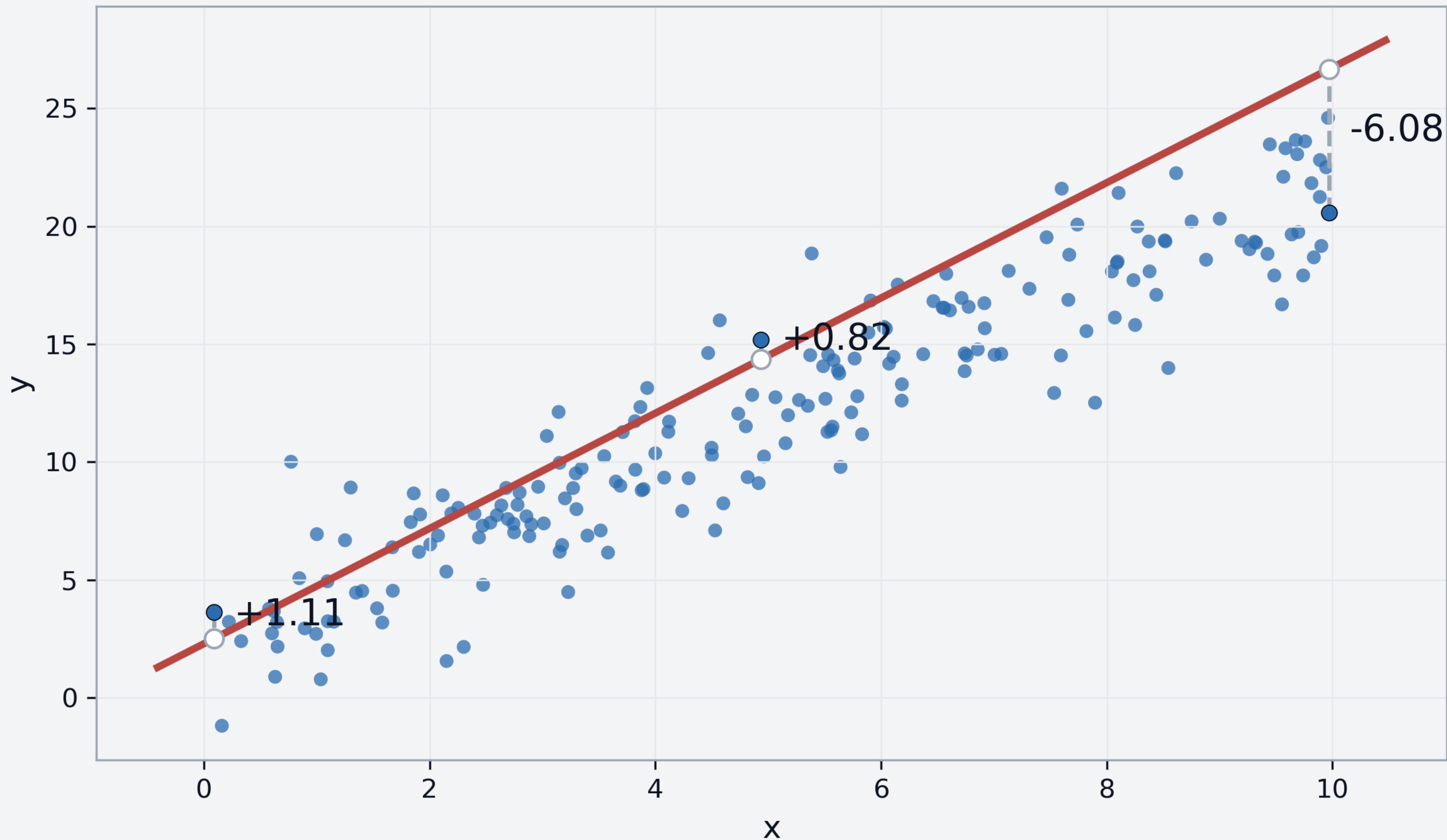




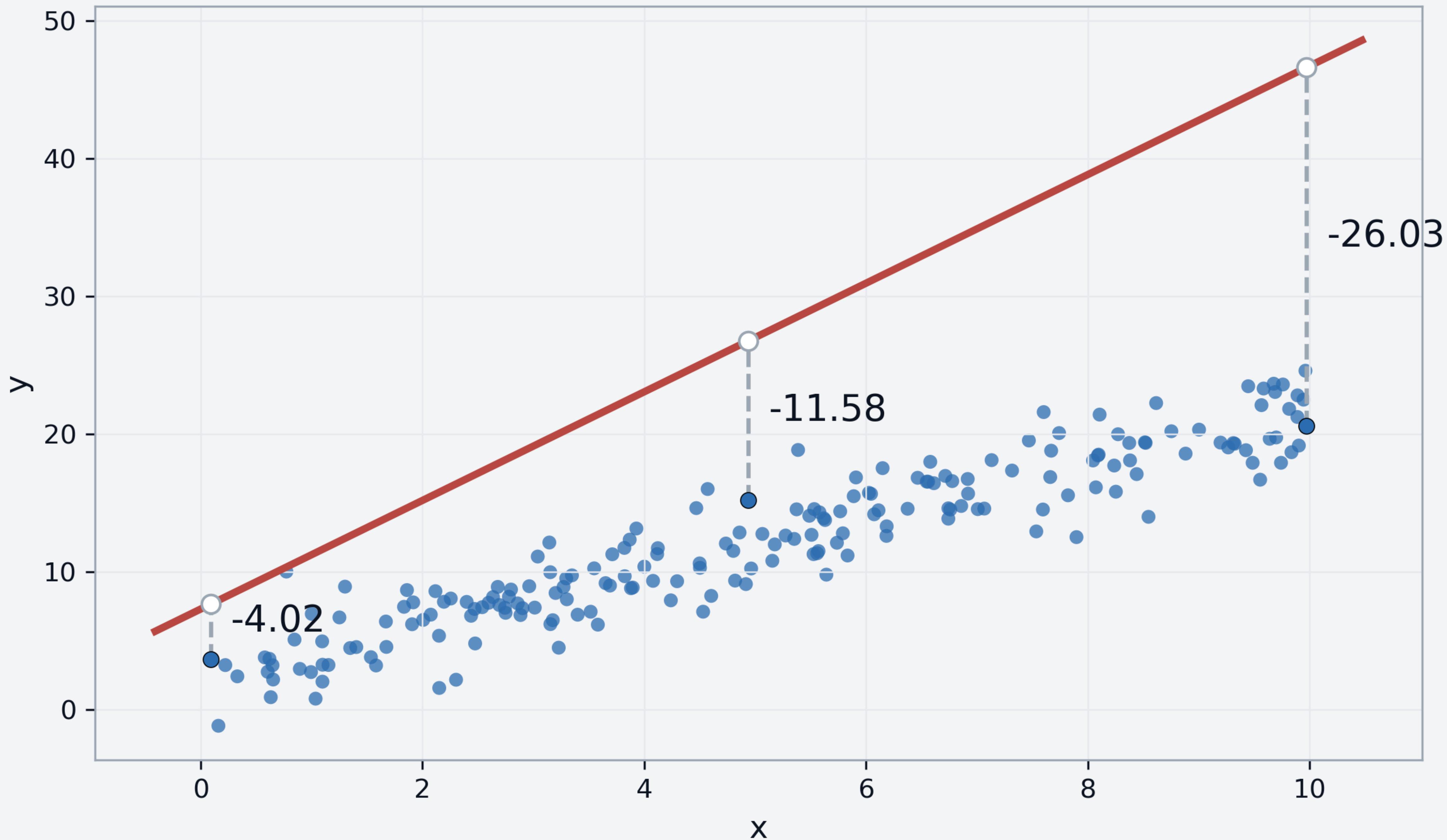
# Regression Errors: dashed lines show residuals ( $y - \hat{y}$ )



# Regression Errors: dashed lines show residuals ( $y - \hat{y}$ )



# Regression Errors: dashed lines show residuals ( $y - \hat{y}$ )



# Regression Errors: dashed lines show residuals ( $y - \hat{y}$ )

