

Regression

A Practical Introduction

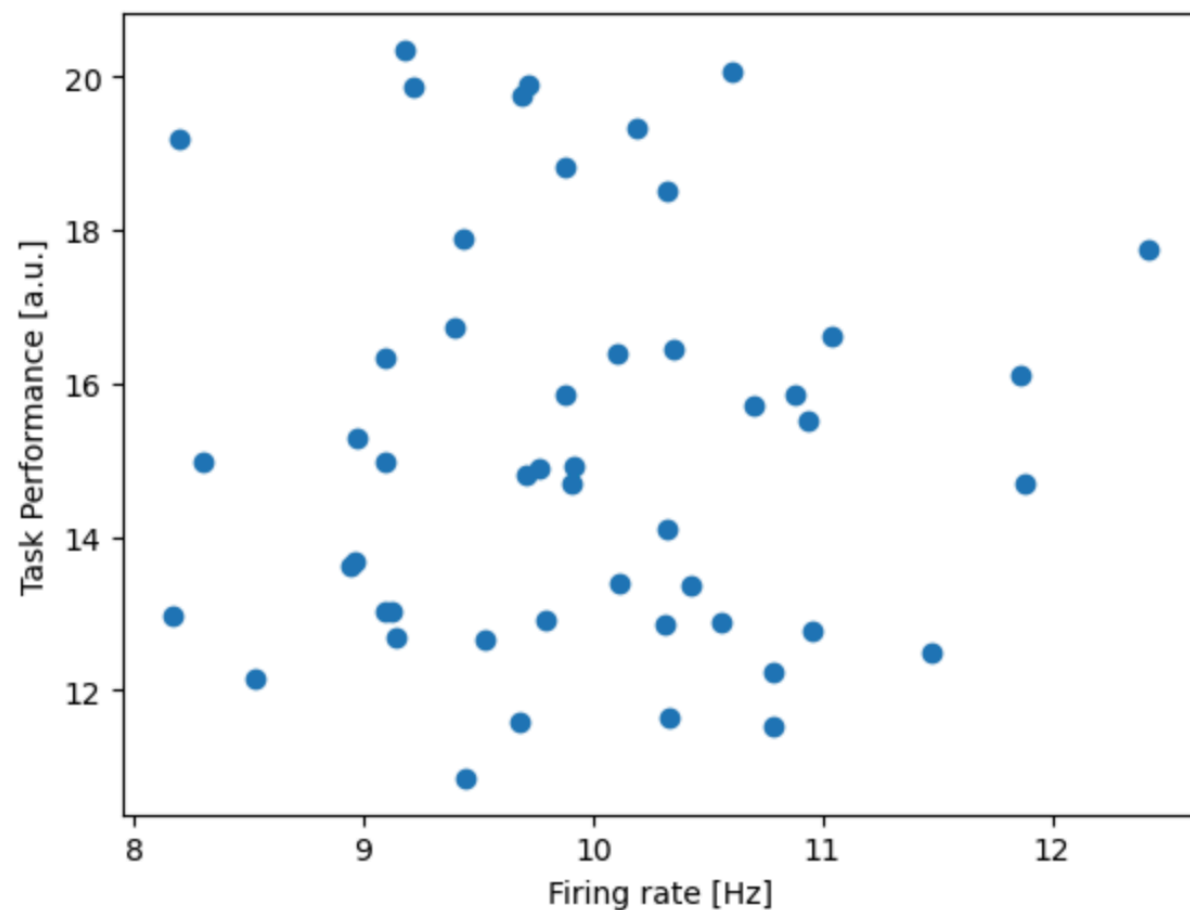
Instructor: Mark Kramer

Outline

A (very) practical introduction to linear regression

Main idea: model data as a line.

Here is my data

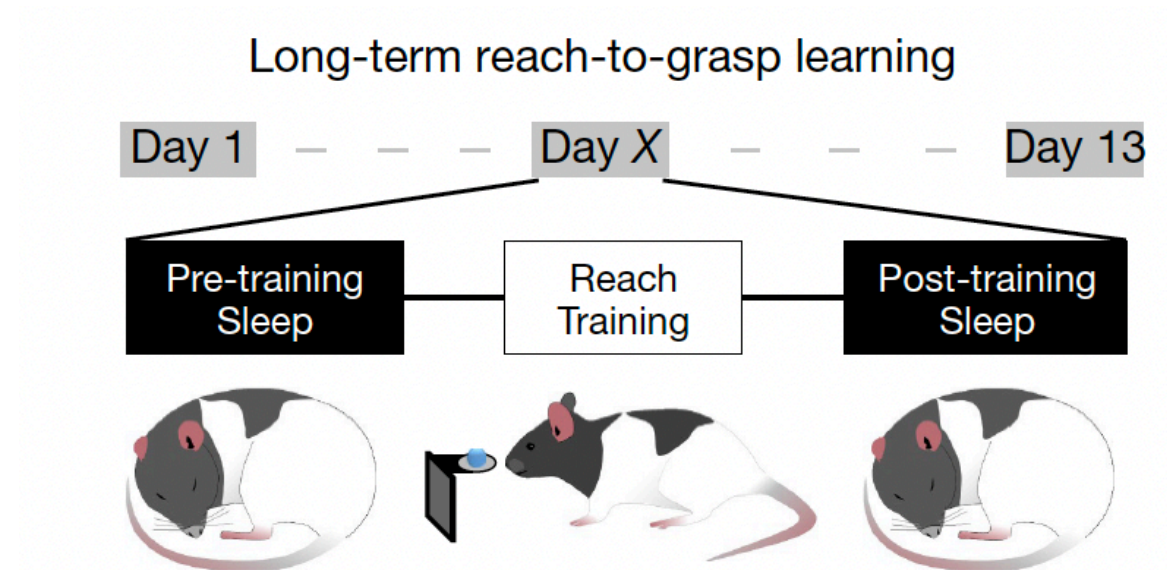
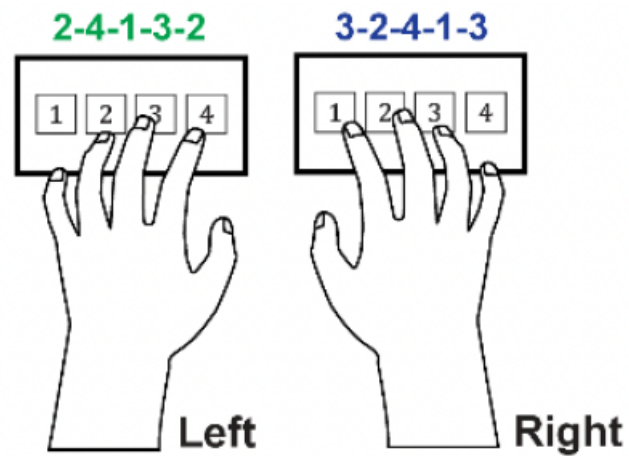


Here is my model

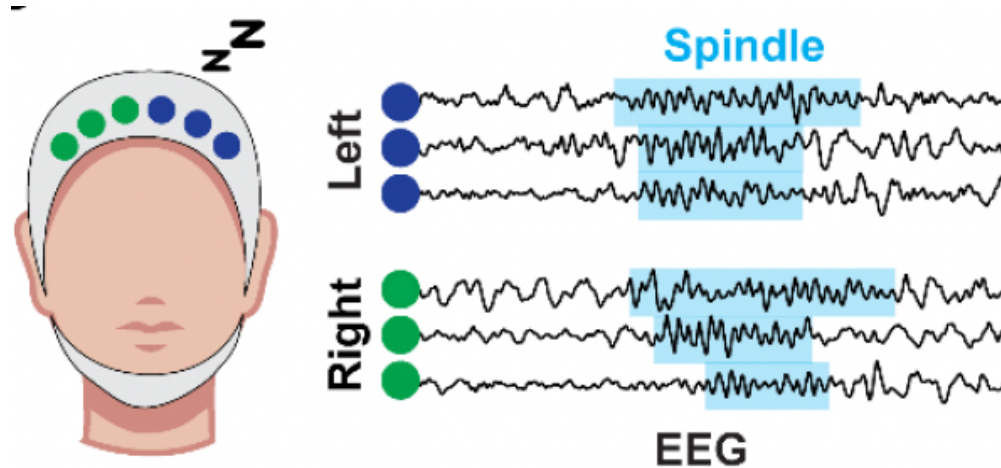
$$y = mx + b$$

Data

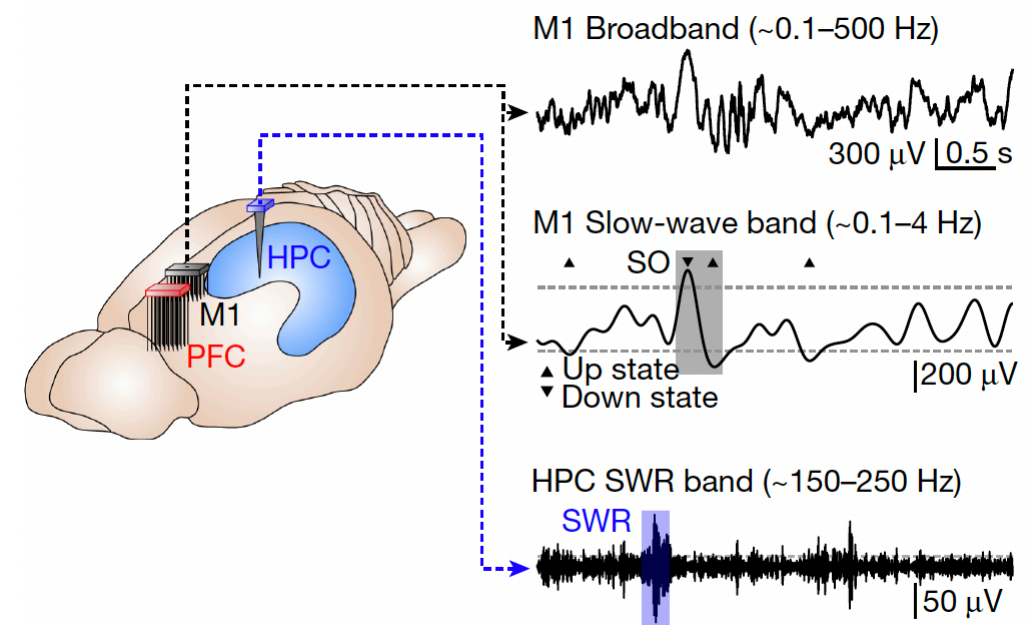
Task performance (y)



Brain activity (x)



[Kwon et al, bioRxiv, 2024]



[Kim et al, Nature, 2023]

Analyze the data (1)

Plot it ...



Python

Visual inspection:

Analyze the data (2)

Compute a statistic?

Correlation

x_n and y_n : data at index n

$$C_{xy} = \frac{1}{N} \frac{1}{\sigma_x} \frac{1}{\sigma_y} \sum_{n=1}^N (x_n - \bar{x})(y_n - \bar{y})$$

number of data points

standard deviation of x

standard deviation of y

sum from indices 1 to N



mean of x



mean of y

mean of x

$$\bar{x} = \frac{1}{N} \sum_{n=1}^N x_n$$

sum the values of x for all n indices, then divide by the total number of points summed (N)

Analyze the data (2)

Compute a statistic?

Correlation

x_n and y_n : data at index n

$$C_{xy} = \frac{1}{N} \frac{1}{\sigma_x} \frac{1}{\sigma_y} \sum_{n=1}^N (x_n - \bar{x})(y_n - \bar{y})$$

number of data points

standard deviation of x

standard deviation of y

sum from indices 1 to N

mean of x

mean of y

variance of x

$$\sigma_x^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})^2$$

characterizes the extent of fluctuations about the mean

standard deviation of x

$$\sigma_x = \sqrt{\sigma_x^2}$$

Analyze the data (2)

Compute a statistic?

Correlation

$$C_{xy} = \frac{1}{N} \frac{1}{\sigma_x} \frac{1}{\sigma_y} \sum_{n=1}^N (x_n - \bar{x})(y_n - \bar{y})$$

sum from indices 1 to N

1	2	3	...				n									N
*	*	*	*				*									*
1	2	3	...				n									N

$x - \bar{x}$

$y - \bar{y}$

then sum & scale = C_{xy}

Analyze the data (2)

Intuition

Correlation

$$C_{xy} = \frac{1}{N} \frac{1}{\sigma_x} \frac{1}{\sigma_y} \sum_{n=1}^N (x_n - \bar{x})(y_n - \bar{y})$$

Assume $\bar{x} = \bar{y} = 0$

$$C_{xy} = \frac{1}{N} \frac{1}{\sigma_x} \frac{1}{\sigma_y} \sum_{n=1}^N x_n y_n$$

Reminder:

$$\sigma_x^2 = \frac{1}{N} \sum_{n=1}^N (x_n - \bar{x})^2$$

What if x and y match?

$$C_{xy} = 1$$

What if x equals $-y$?

$$C_{xy} = -1$$

What if x and y are random?

$$C_{xy} \approx 0$$

Analyze the data (2)

Compute a statistic? Correlation

$$C_{xy} = \frac{1}{N} \frac{1}{\sigma_x} \frac{1}{\sigma_y} \sum_{n=1}^N (x_n - \bar{x})(y_n - \bar{y})$$

Python

$$C_{xy} =$$

Conclusion:

Analyze the data (3): Regression

Motivation: Characterize relationships in the data.

To do so: build a *statistical* model containing

- **systematic effects**: things we know/observe that can explain the data
- **random effects**: unknown / haphazard variations that we make no attempt to model or predict

Regression

Goal: describe succinctly the systematic variations in the data, in a way that's generalizable to other related observations (e.g., by another experimenter, at another time, in another place).

Model

$$y = \alpha + \beta x$$

random effects we don't model

+ noise

y

outcome of measured system (behavior)

x

predictor of measured system (firing rate)

α, β

parameters

Note: linear relationship

Regression

Note: we **cannot** observe y exactly ... measurement error

We observe approximately linear relationship (corrupted by noise).

Challenge: Choose values (a, b) for parameter (α, β) in our model that “best describe” the data.

We observe y_1, y_2, y_3, \dots and x_1, x_2, x_3, \dots and fit our model

$$y = \alpha + \beta x$$

to choose the values (a, b) for parameter (α, β)

Regression

If we have (a, b) , then we can compute model predictions:

$$\hat{y}_1 = a + bx_1$$

$$\hat{y}_2 = a + bx_2$$

⋮

?

Choose (a, b) to make model predictions $\hat{y}_1, \hat{y}_2, \dots$ close to the observed outcomes y_1, y_2, \dots

Note: Model predictions $\hat{y}_1, \hat{y}_2, \dots$ do **not** reproduce exactly the observed outcomes y_1, y_2, \dots

Regression

?

Choose (a, b) to make model predictions $\hat{y}_1, \hat{y}_2, \dots$ close to the observed outcomes y_1, y_2, \dots

Q: “close” ?

A: A measure of discrepancy or distance

$$S_2(y, \hat{y}) = \sum_i (y_i - \hat{y}_i)^2 \quad \text{“least squares”}$$

Choose (a, b) to minimize $S_2(y, \hat{y})$

to minimize the discrepancy between y and \hat{y}

Regression

Minimize $S_2(y, \hat{y}) = \sum_i (y_i - \hat{y}_i)^2$ assumes

1. All observation on the same physical scale (e.g., # vs % correct)
2. Observations are independent or “exchangeable”
3. Deviations $(y_i - \hat{y}_i)$ similar for different values of y
(variability independent of mean)

Regression: estimate it

Estimate the model in Python

$$y = \alpha + \beta x$$

Task performance = $\alpha + \beta$ (firing rate)

↑
intercept

↑
slope

Python

Regression: estimate it

Estimate the model in Python

$$y = \alpha + \beta x$$

Task performance = $\alpha + \beta$ (firing rate)

↑
intercept

↑
slope

OLS Regression Results				
Dep. Variable:	y	R-squared:		
Model:	OLS	Adj. R-squared:		
Method:	Least Squares	F-statistic:		
Date:	Mon, 07 Oct 2024	Prob (F-statistic):		
Time:	12:40:56	Log-Likelihood:		
No. Observations:	50	AIC:		
Df Residuals:	48	BIC:		
Df Model:	1			
Covariance Type:	nonrobust			
	coef	std err	t	P> t
Intercept	15.0190	4.037	3.720	0.001
x	0.0158	0.404	0.039	0.969
Omnibus:	4.793	Durbin-Watson:		
Prob(Omnibus):	0.091	Jarque-Bera (JB):		
Skew:	0.459	Prob(JB):		
Kurtosis:	2.153	Cond. No.		

Interpret parameters ...

Regression: plot it

Python

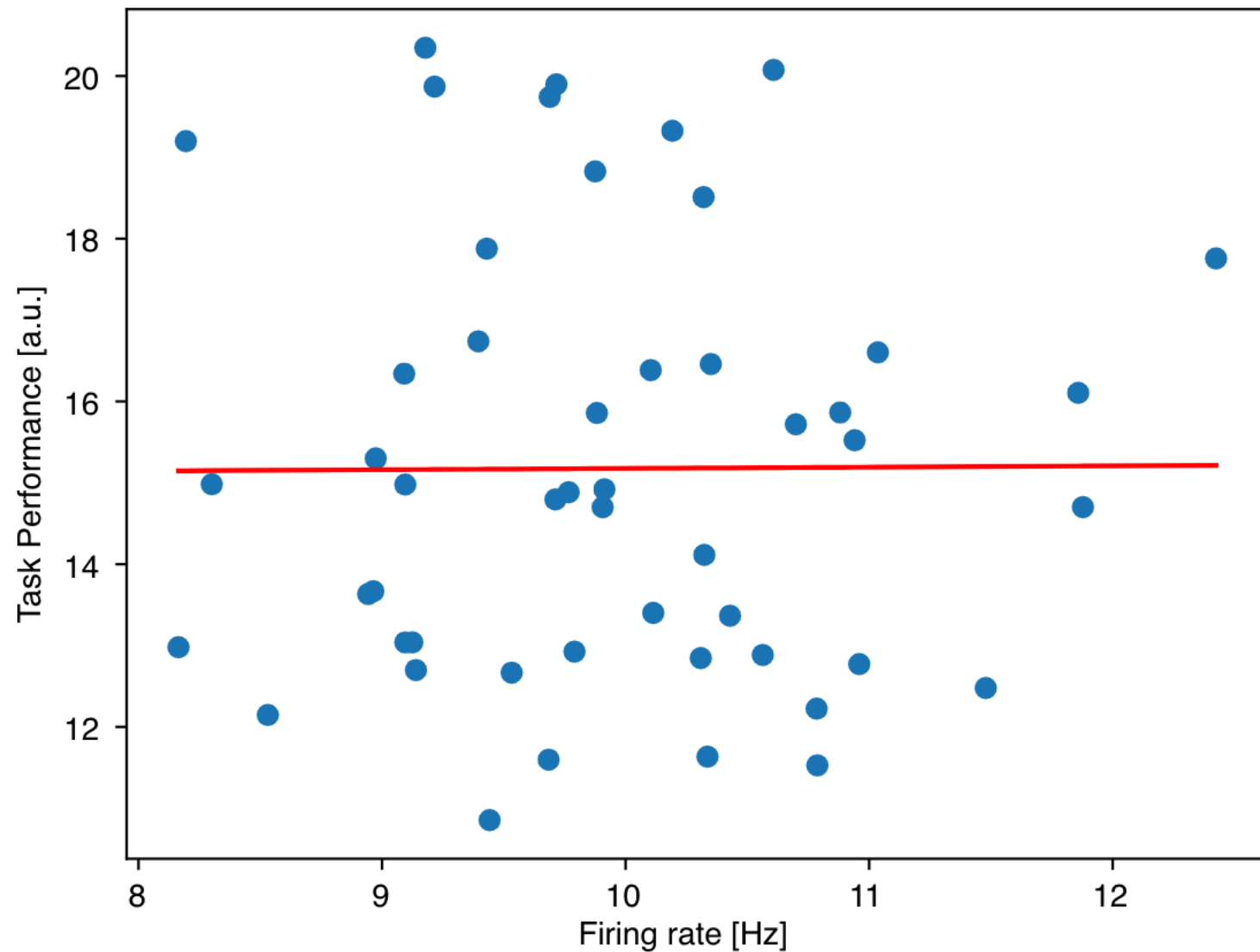
Regression: Interpret parameters

Intercept: $\alpha = 15.02$

- when firing rate (x) is 0, the task performance is ≈ 15

Slope: $\beta = 0.016$

- for each one-unit increase in firing rate, the task performance increases by 0.016.



Q: Evidence of a linear relationship between task performance and firing rate?

Regression: Interpret parameters

Q: Evidence of a linear relationship between task performance and firing rate?

A: Examine the p values

p-value: how much evidence we have to reject the null hypothesis (H_0)

Here, H_0 is that $\alpha = 0, \beta = 0$

Typically, we reject H_0 if $p < 0.05$

The probability of observing the data, or something more extreme, under the null hypothesis is less than 5%.

The observed data is unlikely to have occurred by random chance alone, assuming the null hypothesis is true.

Regression: Interpret parameters

Q: Evidence of a linear relationship between task performance and firing rate?

A: Examine the p values

Intercept: $\alpha = 15.02, p = 0.001$

- Reject H_0 that intercept = 0

Slope: $\beta = 0.016, p = 0.969$

- No evidence to reject H_0 that slope = 0.

Note: Never accept H_0 . ~~We cannot conclude slope = 0~~

Instead: “*We fail to reject the null hypothesis that slope = 0.*”

OLS Regression Results				
Dep. Variable:	y	R-squared:		
Model:	OLS	Adj. R-squared:		
Method:	Least Squares	F-statistic:		
Date:	Mon, 07 Oct 2024	Prob (F-statistic):		
Time:	12:40:56	Log-Likelihood:		
No. Observations:	50	AIC:		
Df Residuals:	48	BIC:		
Df Model:	1			
Covariance Type:	nonrobust			
	coef	std err	t	P> t
Intercept	15.0190	4.037	3.720	0.001
x	0.0158	0.404	0.039	0.969

Regression: conclusion (for now)

We considered this model:

$$\text{Task performance} = \alpha + \beta (\text{firing rate})$$

We found no evidence to reject the null hypothesis that $\beta = 0$.

We conclude that, in this model, we have no evidence of a relationship between task performance and firing rate.

Now what?

Now what?

Give up:

- we tested a hypothesis and it failed.

Confirmatory

Keep exploring:

- fishing expedition

Exploratory

Exploratory versus Confirmatory

Example

Insert tiny electrodes in the rat brain

See what happens when the rat walks, eats, drinks, grooms, sleeps, ...

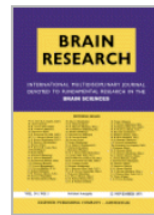
Collect data from 76 cells, then chose 8 to analyze.

Exploratory?



Brain Research

Volume 34, Issue 1, 12 November 1971, Pages 171-175



The hippocampus as a spatial map. Preliminary evidence from unit activity in the freely-moving rat

J. O'Keefe, J. Dostrovsky*

Nobel Prize in Physiology or Medicine 2014

Regression: continued

Q: Now what?

A: Look for confounds.

We learn that age impacts task performance

New variables:

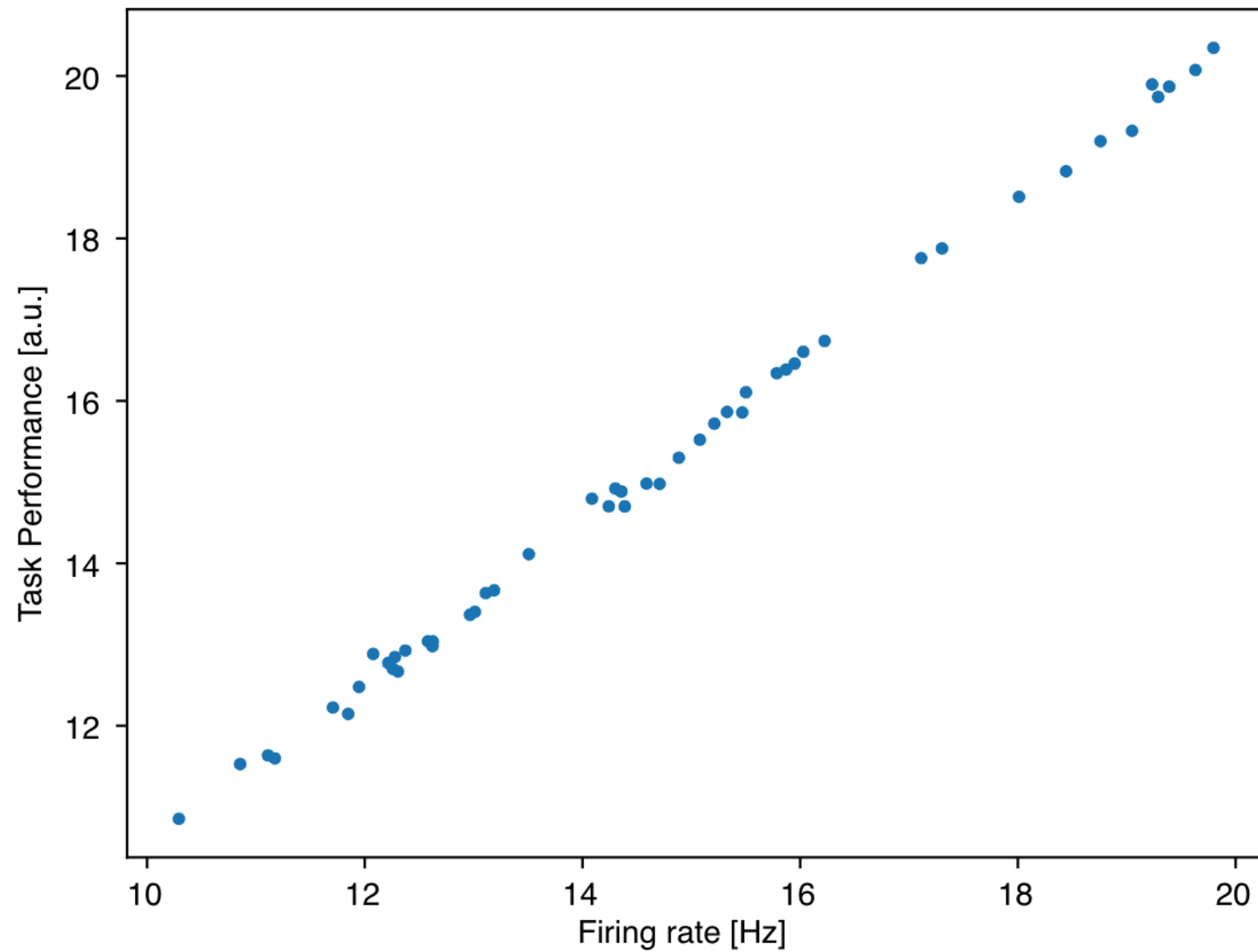
y	task performance
-----	------------------

x_1	firing rate
-------	-------------

x_2	age
-------	-----

Analyze the data (1)

Plot it task performance versus age



Visual inspection:

Analyze the data (2)

Compute the correlation between task performance and age.

Python

$$C_{xy} =$$

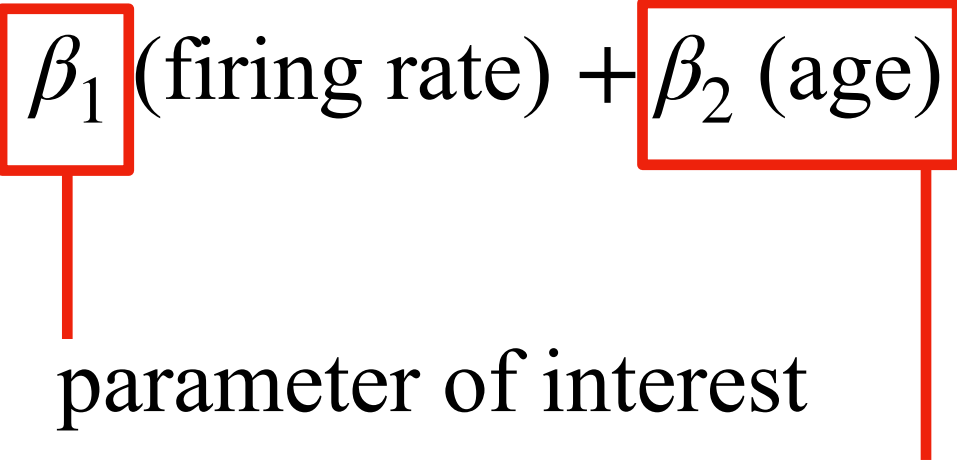
Conclusion:

Analyze the data (3): Regression

Model

$$y = \alpha + \beta_1 x_1 + \beta_2 x_2$$

Task performance = $\alpha + \beta_1$ (firing rate) + β_2 (age)



parameter of interest

confound

Q: What is the relationship between task performance (y) and firing rate (x_1) after accounting for the confound of age (x_2)?

Analyze the data (3): Regression

Python

Regression: Interpret parameters

Intercept: $\alpha =$ $p =$

Slope (firing rate): $\beta_1 =$ $p =$

Slope (age): $\beta_2 =$ $p =$

Regression: Plot the model

Python

Regression: conclusion (modified)

We considered the updated model:

$$\text{Task performance} = \alpha + \beta_1 (\text{firing rate}) + \beta_2 (\text{age})$$

We found

We conclude that

What is a “good model” ?

A: A model that makes predictions \hat{y} very close to y .

To do so, add more predictors (and parameters) to the model.

$$y = \alpha + \beta x_1 + \beta x_2 + \beta x_3 + \beta x_4 + \beta x_5 + \dots$$

No reduction in complexity.

We want a simple theoretical pattern (e.g., line) for our ragged data

parsimony of parameters (only include what we need)

What is a “good model” ?

Parsimonious model

- easier to think about
- probably makes better prediction

Modeling is an art

no formal procedure, requires imagination

‘*All models are wrong but some are useful.*’ [George Box]

eternal truth not within our grasp

use those

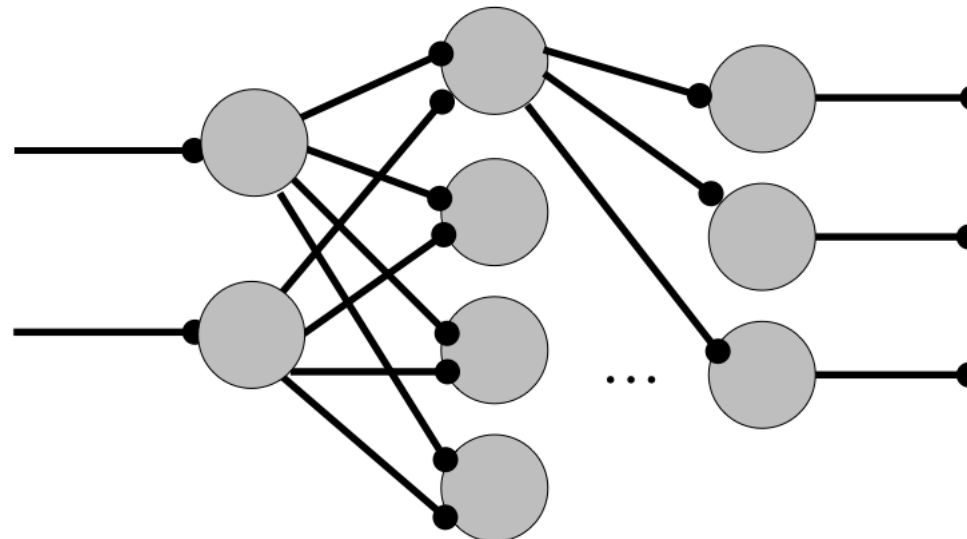
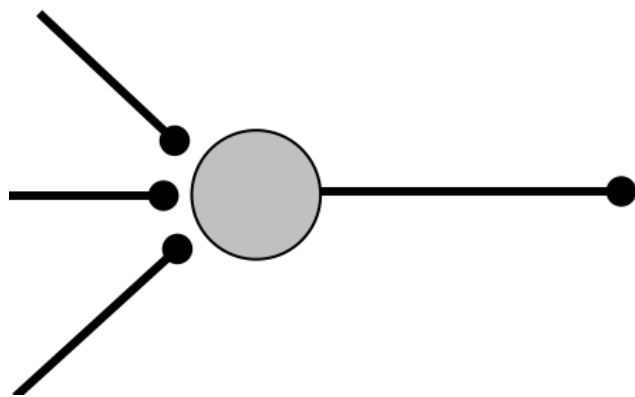
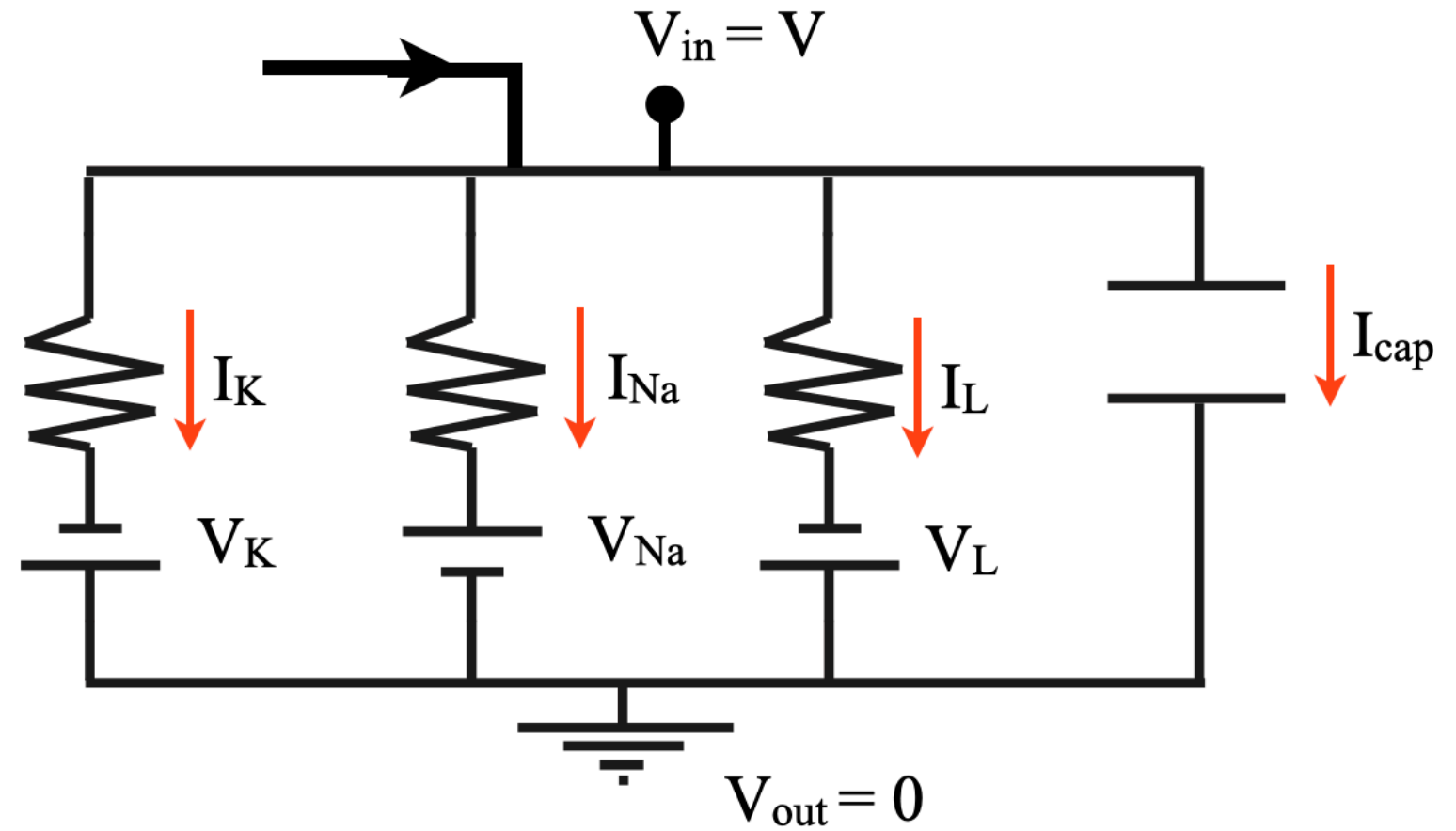
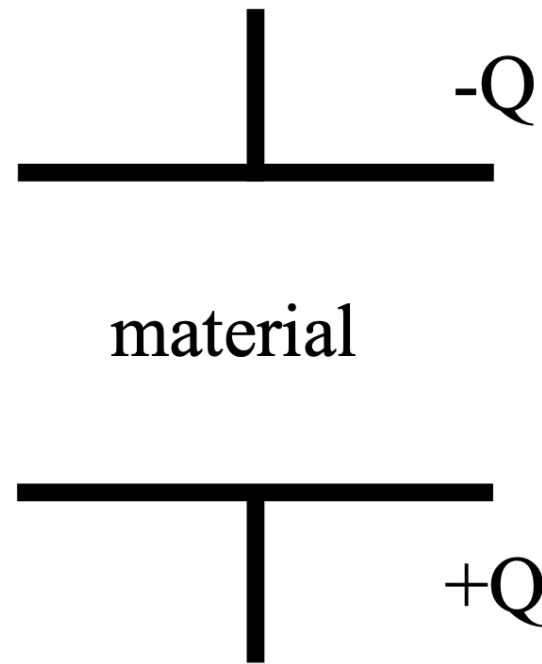
Check your model

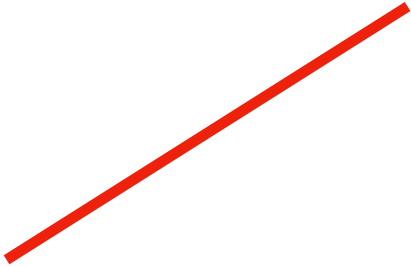
look at errors or deviations ($y_i - \hat{y}_i$)

important but not covered here

What is a model?

In MA665:




$$y = mx + b$$

What is computational neuroscience?

Mathematics:

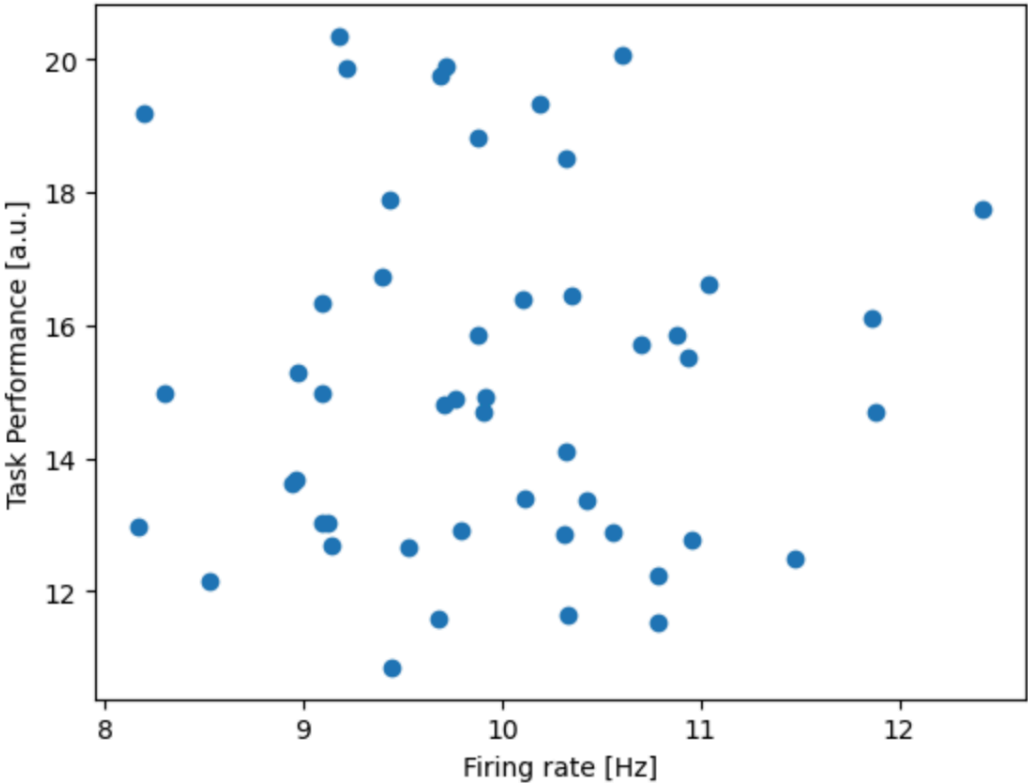
$$C \frac{dV}{dt} = I_{\text{input}}(t) - \bar{g}_K n^4 (V - V_K) - \bar{g}_{Na} m^3 h (V - V_{Na}) - \bar{g}_L (V - V_L)$$
$$\frac{dn}{dt} = - \frac{n - n_{\infty}(V)}{\tau_n(V)}$$
$$\frac{dm}{dt} = - \frac{m - m_{\infty}(V)}{\tau_m(V)}$$
$$\frac{dh}{dt} = - \frac{h - h_{\infty}(V)}{\tau_h(V)},$$

OLS Regression Results			
Dep. Variable:	y	R-squared:	0.000
Model:	OLS	Adj. R-squared:	-0.021
Method:	Least Squares	F-statistic:	0.001521
Date:	Mon, 07 Oct 2024	Prob (F-statistic):	0.969
Time:	12:40:56	Log-Likelihood:	-119.04
No. Observations:	50	AIC:	242.1
Df Residuals:	48	BIC:	245.9
Df Model:	1		
Covariance Type:	nonrobust		

Statistics:

	coef	std err	t	P> t	[0.025	0.975]
			3.720	0.001	6.901	23.137
			0.039	0.969	-0.797	0.829
.793	Durbin-Watson:					1.865
.091	Jarque-Bera (JB):					3.249
.459	Prob(JB):					0.197
.153	Cond. No.					108.

Data:



Aside: C4R

[Home](#)[About Us](#) ▾[Resources](#) ▾[Blog](#)[Contact](#)[Join](#)

Community for Rigor

Better Science Every Day



Welcome to the Community for Rigor! We are a free, open resource to help researchers of all kinds learn, practice, and promote scientific rigor.

Aside: Sample Size

<https://mark-kramer.github.io/METER-Units/>

BU METER

Sample Size - How much data is enough for your experiment?

- Interactive [notebook](#)
-

Evaluate your evaluation methods! A key to meaningful inference.

- Interactive [notebook](#)
-

Putting the p-value in context: $p < 0.05$, but what does it REALLY mean?

- Static [notebook](#)
-

Reproducible exploratory analysis: Mitigating multiplicity when mining data

- Static [notebook](#)

Aside: Sample Size

Q: Is there a relationship between x and lifespan?

A1: Do an experiment with sample size N .

A2: Fit a line...

$$lifespan = \beta_0 + \beta_1 x$$

$$\beta_1 =$$

$$p =$$

Conclusion:

Aside: Sample Size

Q: Now what?

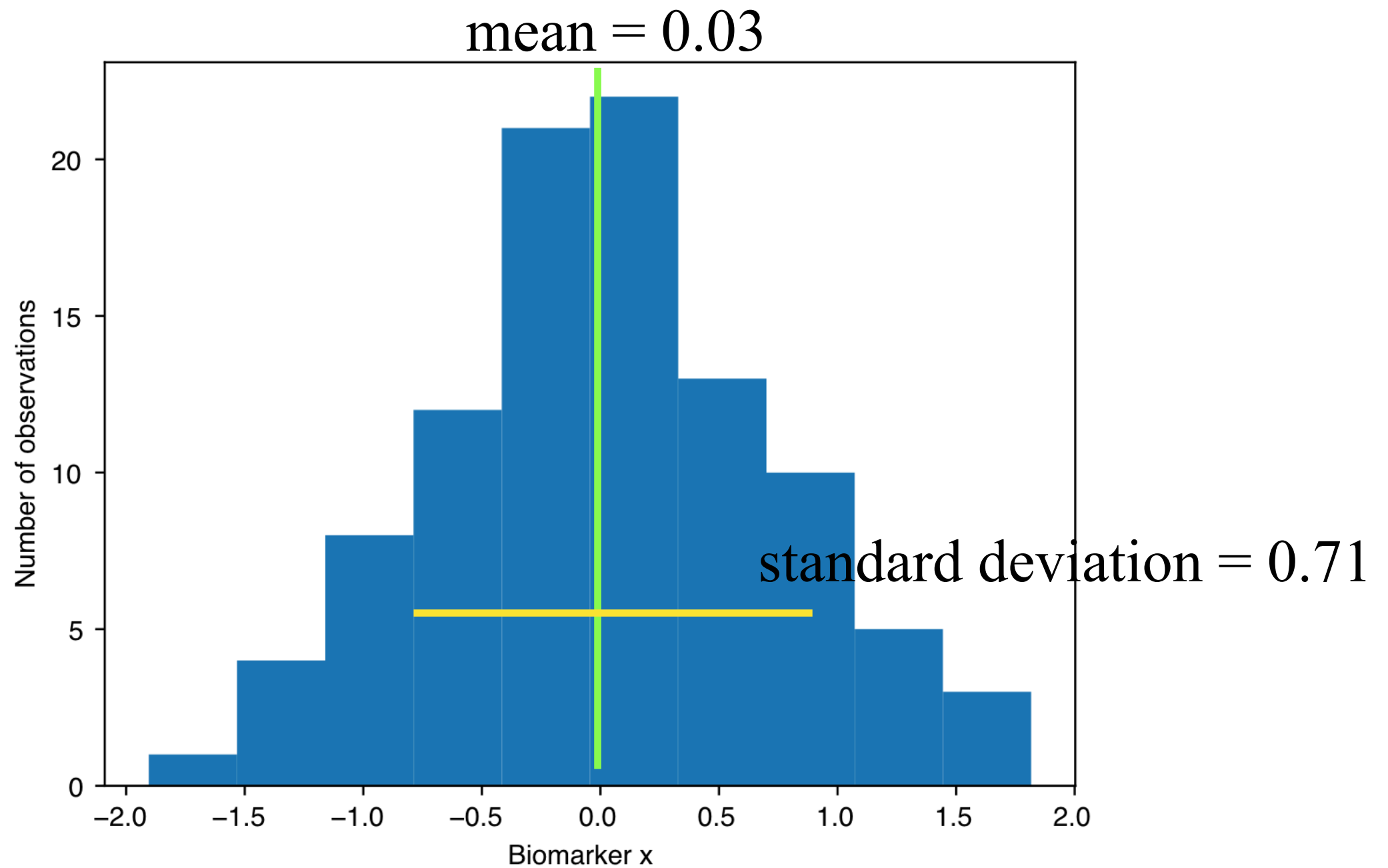
A: Maybe we failed to collect enough data to detect a relationship.

Idea:

- Reuse the data & model
- See how sample size (N) impacts conclusions.

Aside: Sample Size

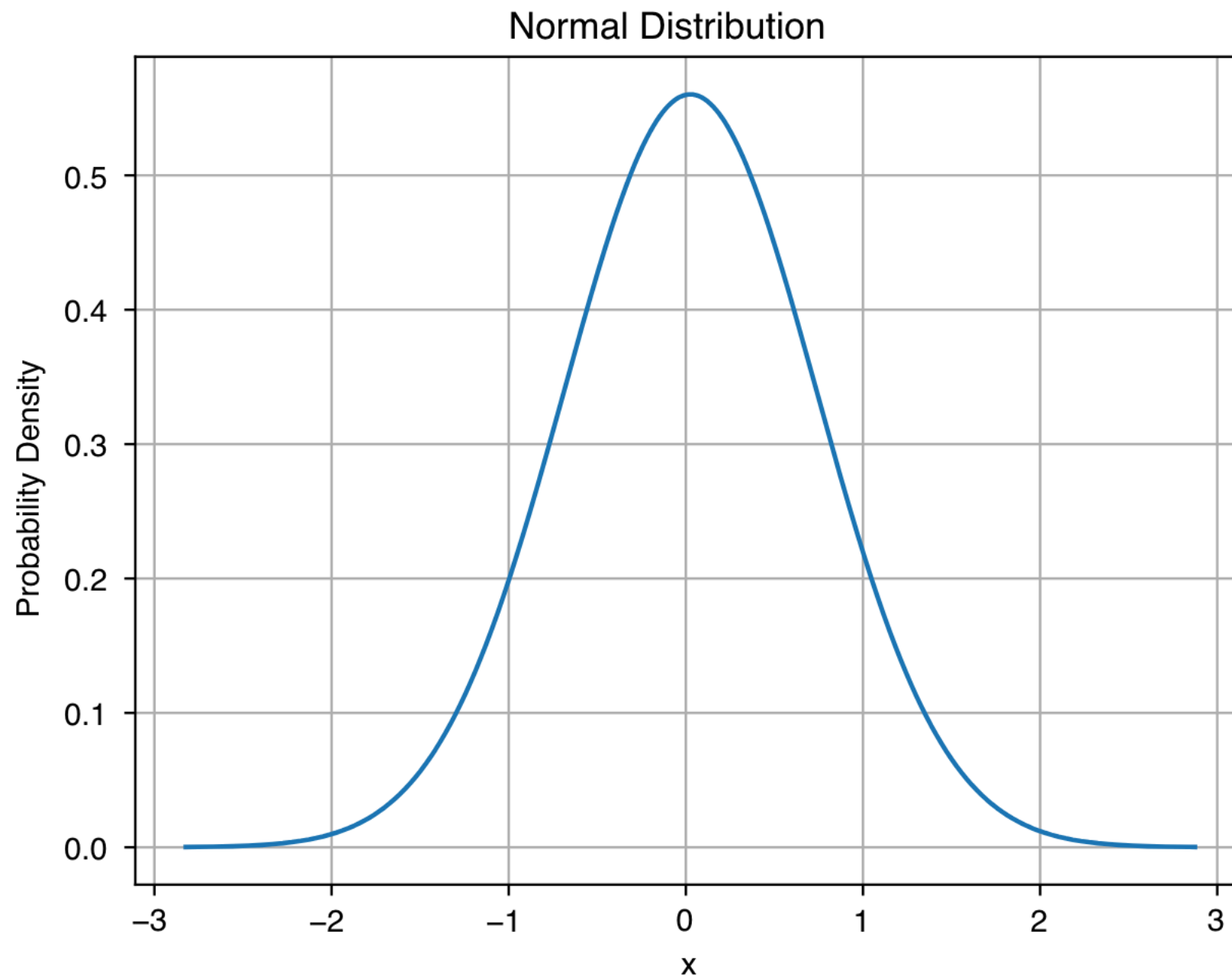
Consider biomarker x



Approximately normal

Aside: Sample Size

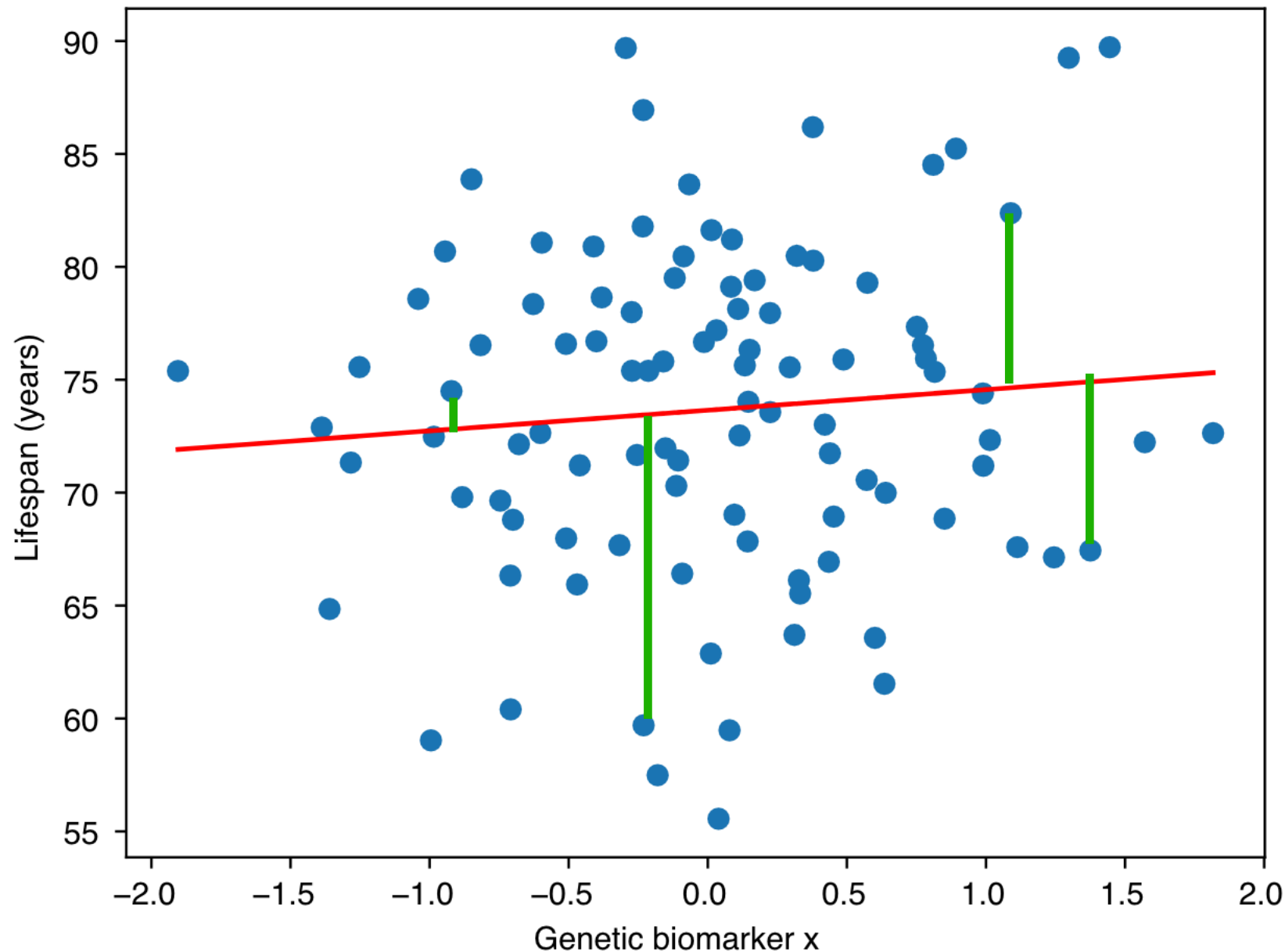
We can draw random values of x from this normal distribution



Draw 10 or 100 or 1000 or 10,000 values for x ...

Aside: Sample Size

Consider model: $lifespan = \beta_0 + \beta_1 x$



$\beta_0 = 73.65$ (intercept)

$\beta_1 = 0.91$ (slope)

There's error in our model

Normally distributed:

mean ≈ 0

stand. dev. ≈ 7

To simulate new lifespans:

- Ask the model
- Include the error

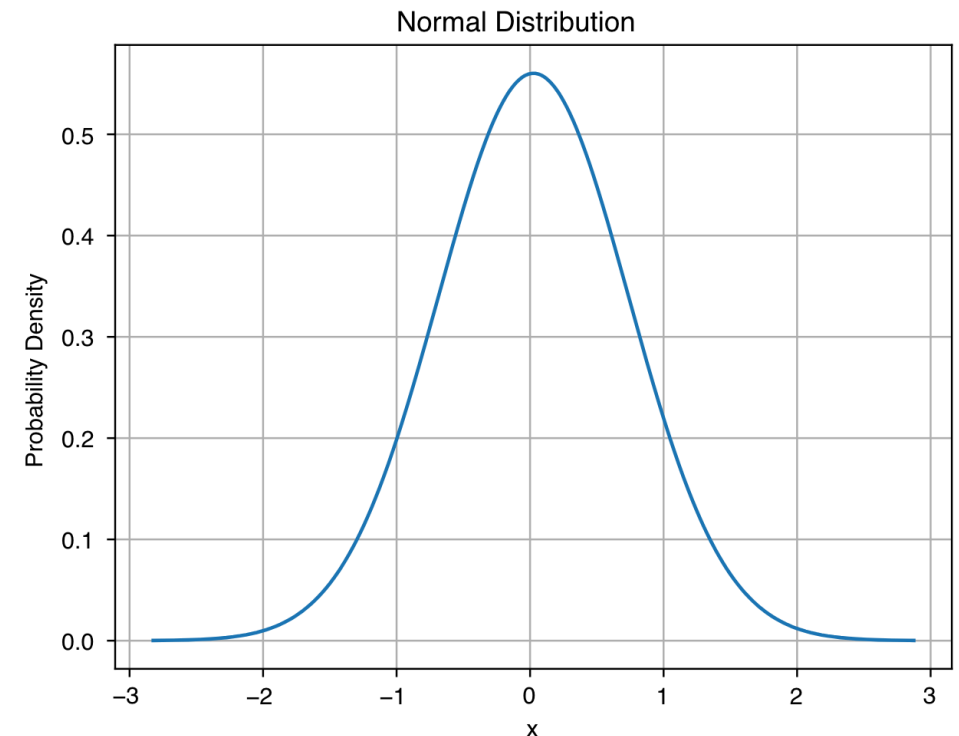
$$\text{new lifespan} = \beta_0 + \beta_1 x + \text{error}$$

Aside: Sample Size

Create new data:

- Pick new sample size N^*
- Draw new biomarkers x
- Draw new lifespans

$$\text{new lifespan} = \beta_0 + \beta_1 x + \text{error}$$



Key insight: Is there a relationship between x & lifespan in the new data?

Fit a (new) model: $\text{new lifespan} = \beta_0^* + \beta_1^* \text{new } x$

Q: At what new sample size N^* do you reliably detect a relationship?
... is $p < 0.05$ reliably.