

Lecture 10

Feature Engineering, Gradient Descent

DSC 40A, Spring 2024

Agenda

- Feature engineering and transformations.
- Minimizing functions using gradient descent.

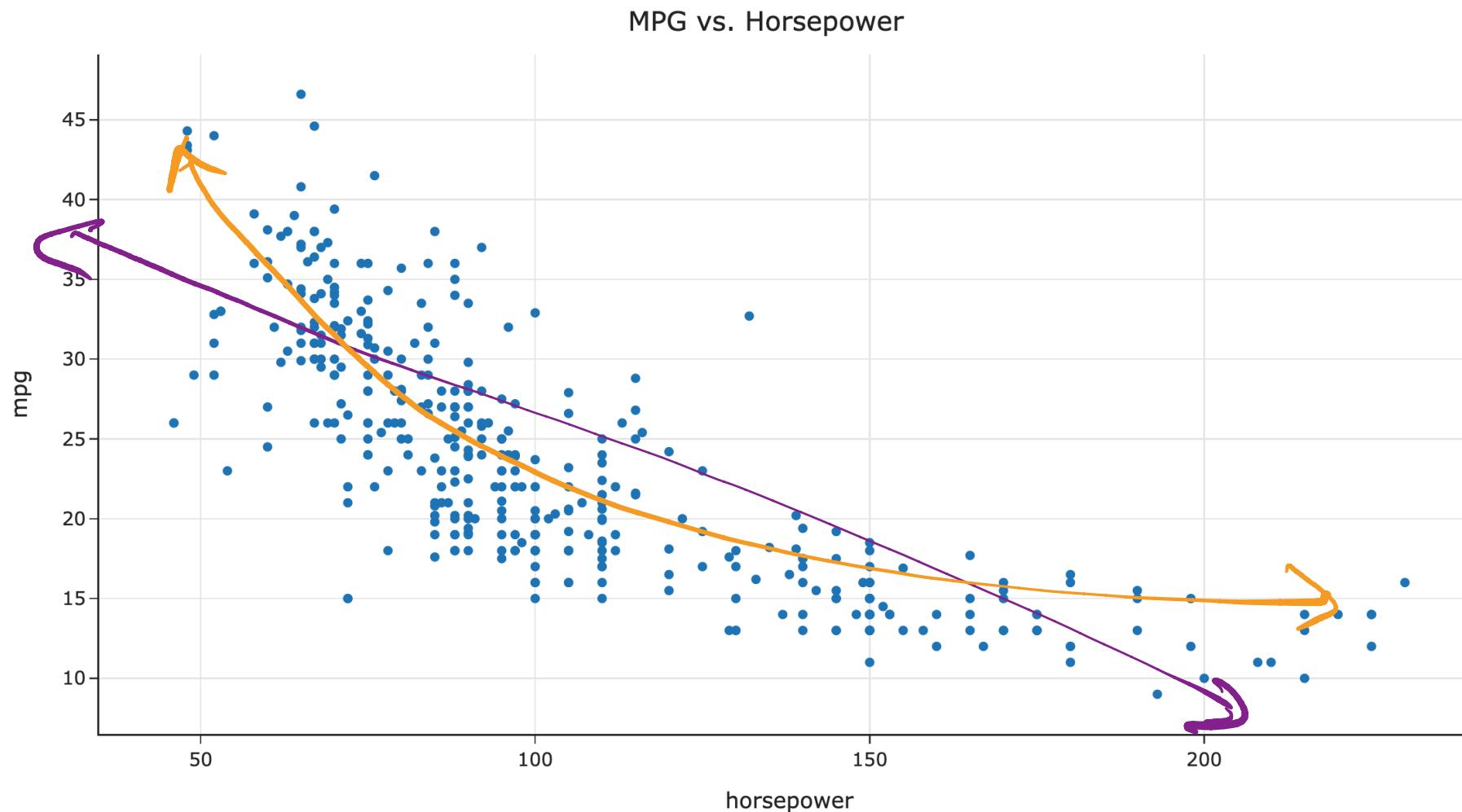
Question 🤔

Take a moment to pause and reflect...

If you have any questions please post online to our forms/Q&A site.

Course staff will answer them ASAP!

Feature engineering and transformations



Question: Would a linear hypothesis function work well on this dataset?

? $\xrightarrow{\text{This!!}}$ $w_0 \underline{?} + w_1 \underline{?} + w_2 \underline{?} + \dots + w_n \underline{?}$

Linear in the parameters

- We can fit rules like:

$$\vec{y} = \vec{X} \vec{w}$$

good ✓

$$w_0 + w_1 x + w_2 x^2$$

$$w_1 e^{-x^{(1)^2}} + w_2 \cos(x^{(2)} + \pi) + w_3 \frac{\log 2x^{(3)}}{x^{(2)}}$$

- This includes arbitrary polynomials.
- These are all linear combinations of (just) features.

Features

- We can't fit rules like:

nope!

$$w_0 + e^{w_1 x}$$

$$w_0 + \sin(w_1 x^{(1)} + w_2 x^{(2)})$$

- These are **not** linear combinations of just features! nope

- We can have any number of parameters, as long as our hypothesis function is **linear in the parameters**, or linear when we think of it as a function of the parameters.

Example: Amdahl's Law

- Amdahl's Law relates the runtime of a program on p processors to the time to do the sequential and nonsequential parts on one processor.

$$H(p) = t_S + \frac{t_{NS}}{p}$$

- Collect data by timing a program with varying numbers of processors:

Processors	Time (Hours)
1	8
2	4
4	3

We care about X , \vec{w} , \vec{y}

Example: Fitting $H(x) = w_0 + w_1 \cdot \frac{1}{x}$

Processors	Time (Hours)
1	8
2	4
4	3

X y

$$X = \begin{bmatrix} 1 & \frac{1}{1} \\ 1 & \frac{1}{2} \\ 1 & \frac{1}{4} \end{bmatrix}_{3 \times 2}$$

$$\vec{w} = \begin{bmatrix} w_0 \\ w_1 \end{bmatrix}_{2 \times 1}$$

$$\vec{y} = \begin{bmatrix} 8 \\ 4 \\ 3 \end{bmatrix}_{3 \times 1}$$

$$\vec{h} = X\vec{w}$$

Goal is to find w_0^* and w_1^* by:

$$X^T X \vec{w}^* = X^T \vec{y}$$

"System of two equations"
two unknowns!

$$\vec{w}^* = (X^T X)^{-1} X^T \vec{y}$$

$\Rightarrow X^T X$ is invertible (full rank & X 's cols are linearly independent)

How do we fit hypothesis functions that aren't linear in the parameters?

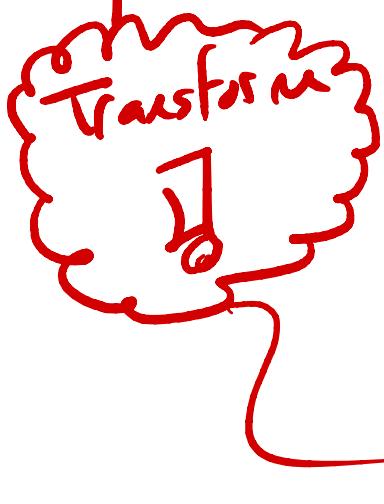
- Suppose we want to fit the hypothesis function:

$$H(x) = w_0 e^{w_1 x}$$

not
linear!

needs to change!

- This is **not** linear in terms of w_0 and w_1 , so our results for linear regression don't apply.
- Possible solution: Try to apply a transformation.



Goal is to get it in the form :

$$w_0 + w_1 \cdot \underline{?} + w_n \cdot \underline{?}$$

Transformations

$$y = H(x)$$

very close!
 $w_0 + w_1 \dots$?

- Question: Can we re-write $H(x) = w_0 e^{w_1 x}$ as a hypothesis function that is linear in the parameters?

☞ Use the log

$$H(x) = w_0 e^{w_1 x}$$

$$\log(y) = \log(w_0 e^{w_1 x})$$

$$\log(y) = \log(w_0) + \log(e^{w_1 x})$$

$$\log(y) = \log(w_0) + w_1 x$$

$$z = b_0 + b_1 x$$

$$z = \log(y)$$

$$b_0 = \log(w_0)$$

$$\Rightarrow w_0^* = e^{b_0}$$

Prop. of logs

$$\log(ab) = \log(a) + \log(b)$$

Conversions
we need!

Transformations

Provided us a way to solve

$$\mathbf{X}^T \mathbf{X} \vec{b}^* = \mathbf{X}^T \vec{z}$$

\vec{b} = parameter vector
 \vec{z} = new observation vector

- **Solution:** Create a new hypothesis function, $T(x)$, with parameters b_0 and b_1 , where $T(x) = b_0 + b_1 x$.
- This hypothesis function is related to $H(x)$ by the relationship $T(x) = \log H(x)$.
- \vec{b} is related to \vec{w} by $b_0 = \log w_0$ and $b_1 = w_1$.

- Our new observation vector, \vec{z} , is
$$\begin{bmatrix} \log y_1 \\ \log y_2 \\ \vdots \\ \log y_n \end{bmatrix}.$$

- $T(x) = b_0 + b_1 x$ is linear in its parameters, b_0 and b_1 .
- Use the solution to the normal equations to find \vec{b}^* , and the relationship between \vec{b} and \vec{w} to find \vec{w}^* .

Non-linear hypothesis functions in general

- Sometimes, it's just not possible to transform a hypothesis function to be linear in terms of some parameters.
- In those cases, you'd have to resort to other methods of finding the optimal parameters.
 - For example, $H(x) = w_0 \sin(w_1 x)$ can't be transformed to be linear.
 - But, there are other methods of minimizing mean squared error:
$$R_{\text{sq}}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - w_0 \sin(w_1 x))^2$$
- One method: **gradient descent**, the topic we're going to look at next!
- Hypothesis functions that are linear in the parameters are much easier to work with.

Question 🤔

Σ w. ?

Pause the video and try to answer the question...

Which hypothesis function is **not** linear in the parameters?

- ✓ • A. $H(\vec{x}) = \boxed{w_1}(x^{(1)}x^{(2)}) + \boxed{w_2} \sin(x^{(2)})$ → $\boxed{w_2} \cdot \frac{\sin(x^{(2)})}{x^{(1)}}$
- ✗ • B. $H(\vec{x}) = 2^{w_1}x^{(1)}$
- ✓ • C. $H(\vec{x}) = \vec{w} \cdot \text{Aug}(\vec{x})$ → $w_1(x^{(1)}) + w_2(x^{(2)})$
- ✓ • D. $H(\vec{x}) = \boxed{w_1} \cos(x^{(1)}) + \boxed{w_2} 2^{x^{(2)}} \log x^{(3)}$
- ✓ • E. More than one of the above.

The modeling recipe

given some Data & targets "y's" How do we make good Pred.

1. Choose a model.

(I.) $H(x) = h$ constant

(II.) $H(x) = \omega_0 + \omega_1 x$
Simple Linear Reg.

multi linear Regression

(III.) $H(\vec{x}) = \omega_0 + \omega_1 x^{(1)} + \omega_2 x^{(2)} + \dots + \omega_d x^{(d)}$

- Single Pred: $\vec{\omega} \cdot \text{Aug}(\vec{x})$
- All Preds: $\vec{h} = X \vec{\omega}$

2. Choose a loss function.

(A.) Squared loss: $(y_i - H(x_i))^2$ (C.) 0-1 loss

(B.) Absolute loss: $|y_i - H(x_i)|$

3. Minimize average loss to find optimal model parameters.

Constant Model w/ Squared loss:

$$R_{\text{sq}}(h) = \frac{1}{n} \sum_{i=1}^n (y_i - h)^2 \Rightarrow h^* = \text{Mean}(y_1, y_2, \dots, y_n)$$

Multiple lin. Reg w/ Squared loss:

$$R_{\text{sq}}(\vec{\omega}) = \frac{1}{n} \| \vec{y} - X \vec{\omega} \|^2$$

Minimizing functions using gradient descent

Minimizing empirical risk

- Repeatedly, we've been tasked with **minimizing** the value of empirical risk functions.
 - Why? To help us find the **best** model parameters, h^* or w^* , which help us make the **best** predictions!
- We've minimized empirical risk functions in various ways.

$$\circ R_{\text{sq}}(h) = \frac{1}{n} \sum_{i=1}^n (y_i - h)^2 \rightarrow \text{calculus to solve these equations}$$
$$\circ R_{\text{abs}}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n |y_i - (w_0 + w_1 x)| \rightarrow \text{By hand w/ Brute force approach!}$$
$$\circ R_{\text{sq}}(\vec{w}) = \frac{1}{n} \|\vec{y} - X\vec{w}\|^2 \rightarrow \text{Principles of Linear Algebra!}$$

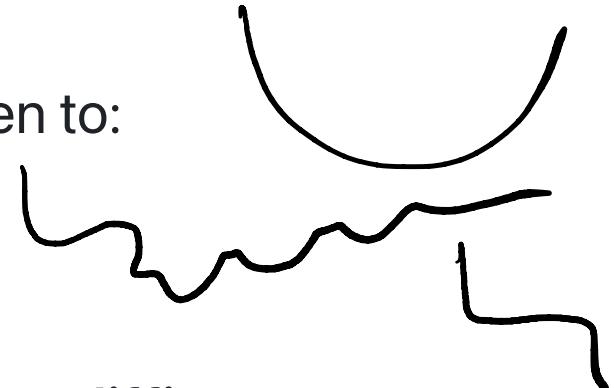
Minimizing arbitrary functions

a derivative does exist

- Assume $f(t)$ is some **differentiable** single-variable function.
- When tasked with minimizing $f(t)$, our general strategy has been to:
 - i. Find $\frac{df}{dt}(t)$, the derivative of f .
 - ii. Find the input t^* such that $\frac{df}{dt}(t^*) = 0$.
- However, there are cases where we can find $\frac{df}{dt}(t)$, but it is either difficult or impossible to solve $\frac{df}{dt}(t^*) = 0$.

$$f(t) = 5t^4 - t^3 - 5t^2 + 2t - 9$$

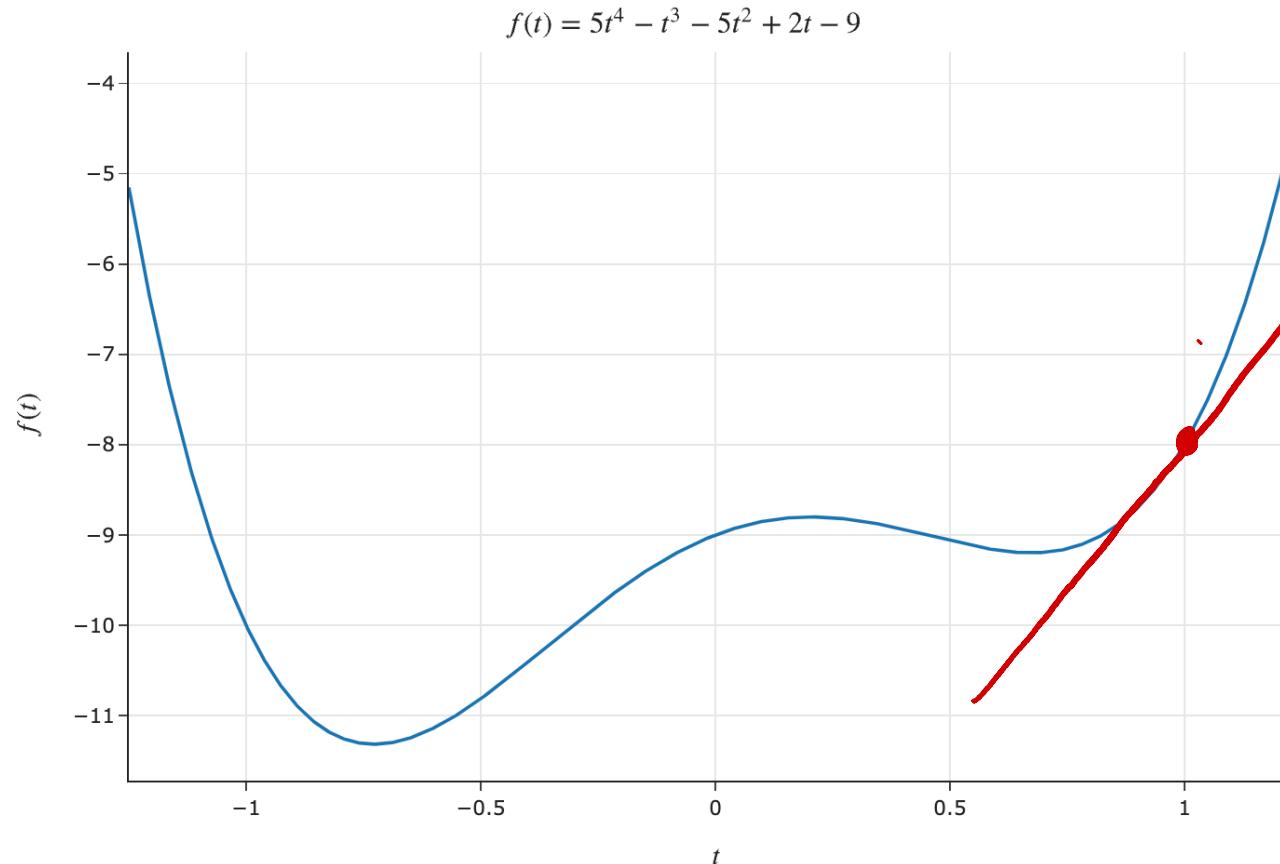
$$\frac{df}{dt}(t) = 20t^3 - 3t^2 - 10t + 2$$



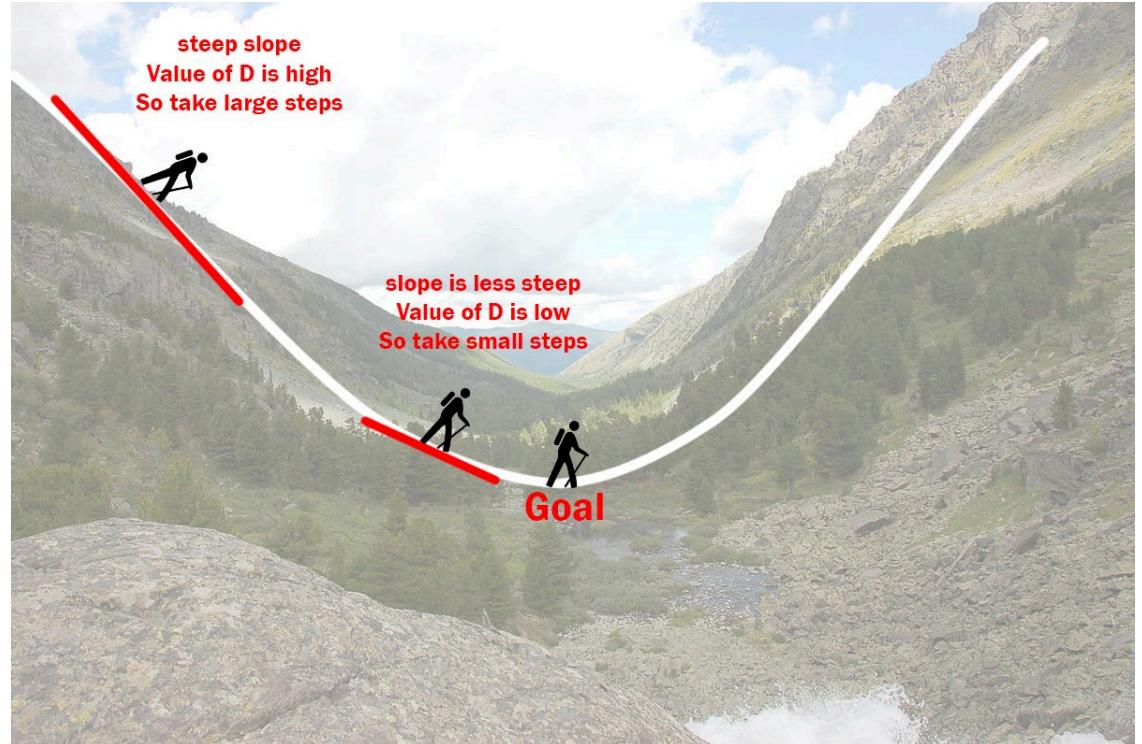
- Then what?

What does the derivative of a function tell us?

- **Goal:** Given a **differentiable** function $f(t)$, find the input t^* that minimizes $f(t)$.
- What does $\frac{d}{dt} f(t)$ mean?

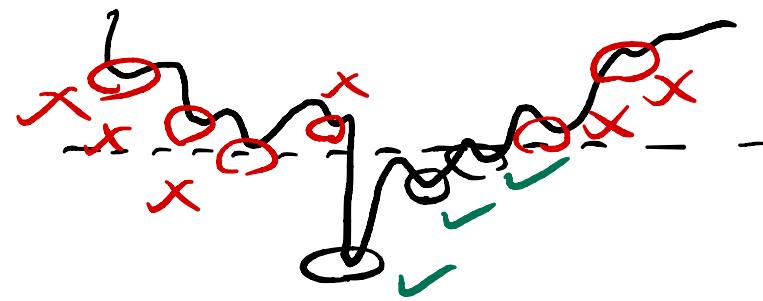
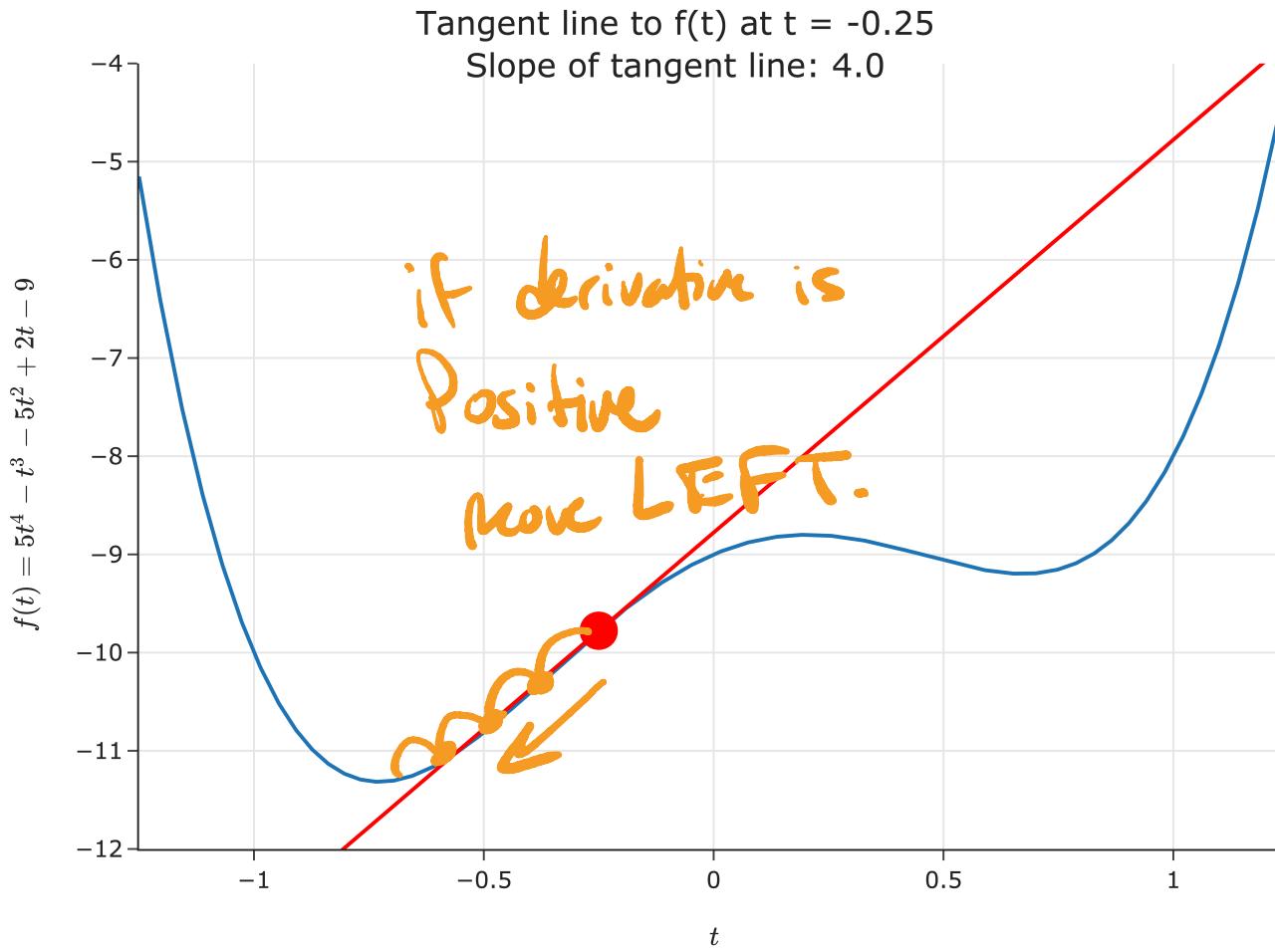


Let's go hiking!

- Suppose you're at the top of a mountain  and need to get **to the bottom**.
- Further, suppose it's really cloudy 

The image shows a mountainous landscape with a winding path through a valley. Two hikers are visible on the path. A red line highlights a steep section of the path, and a white line highlights a less steep section. Text overlays provide instructions for each section:
steep slope
Value of D is high
So take large steps
slope is less steep
Value of D is low
So take small steps
Goal

Searching for the minimum

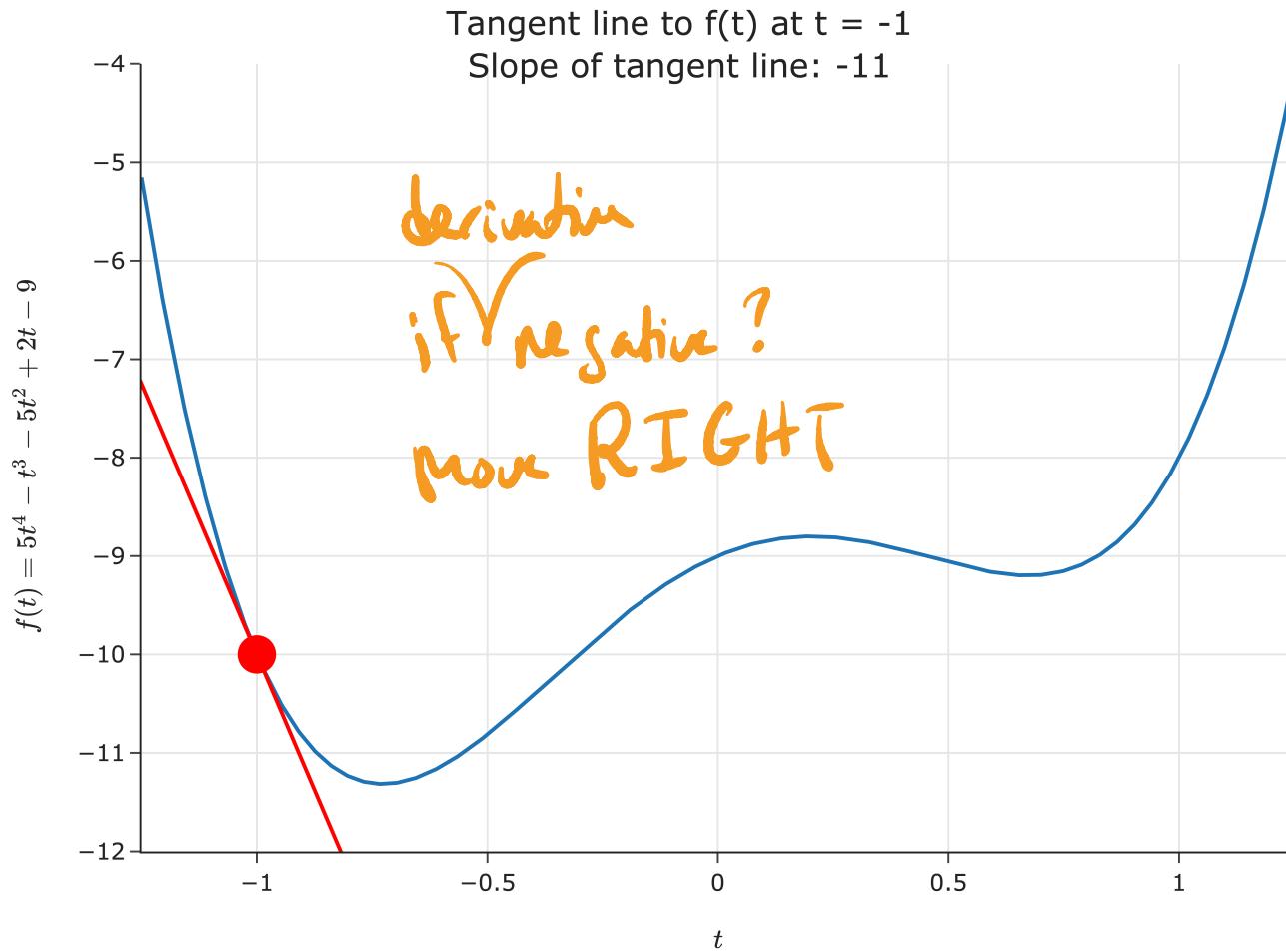


Suppose we're given an initial guess for a value of t that minimizes $f(t)$.

If the **slope of the tangent line at $f(t)$ is positive ↗**:

- Increasing t increases f .
- This means the minimum must be to the **left** of the point $(t, f(t))$.
- Solution: **Decrease t** ⬇.

Searching for the minimum



Suppose we're given an initial guess for a value of t that minimizes $f(t)$.

If the **slope of the tangent line at $f(t)$ is negative** 📉:

- Increasing t **decreases** f .
- This means the minimum must be to the **right** of the point $(t, f(t))$.
- Solution: **Increase t** ⬆.

Intuition

- To minimize $f(t)$, start with an initial guess t_0 .
- Where do we go next?
 - If $\frac{df}{dt}(t_0) > 0$, decrease t_0 .
 - If $\frac{df}{dt}(t_0) < 0$, increase t_0 .
- One way to accomplish this:

There exists some t^* that minimizes $f(t)$ and $t_0, t_1, t_2, \dots, t_n$ are guesses for t^* we improve t_n via this

Grad. Descent
update Rule

$$t_1 = t_0 - \frac{df}{dt}(t_0)$$

opposite direction
of the Derivative

Gradient descent

To minimize a **differentiable** function f :

- Pick a positive number, α . This number is called the **learning rate**, or **step size**.
- Pick an **initial guess**, t_0 .
- Then, repeatedly update your guess using the **update rule**:

$$t_{i+1} = t_i - \alpha \frac{df}{dt}(t_i)$$

Step Sizes

- Repeat this process until **convergence** – that is, when t doesn't change much.
- This procedure is called **gradient descent**.

What is gradient descent?

- Gradient descent is a numerical method for finding the input to a function f that minimizes the function.
- Why is it called **gradient** descent?
 - The gradient is the extension of the derivative to functions of multiple variables.
 - We will see how to use gradient descent with multivariate functions next class.
- What is a **numerical** method?
 - A numerical method is a technique for approximating the solution to a mathematical problem, often by using the computer.
- Gradient descent is **widely used** in machine learning, to train models from linear regression to neural networks and transformers (including ChatGPT)!

Lingering questions

Next class, we'll explore the following ideas:

- When is gradient descent *guaranteed* to converge to a global minimum?
 - What kinds of functions work well with gradient descent?
- How do I choose a step size?
- How do I use gradient descent to minimize functions of multiple variables, e.g.:

$$R_{\text{sq}}(w_0, w_1) = \frac{1}{n} \sum_{i=1}^n (y_i - (w_0 + w_1 x_i))^2$$