

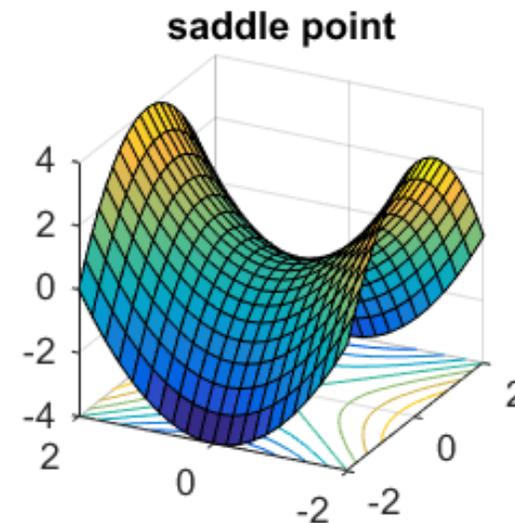
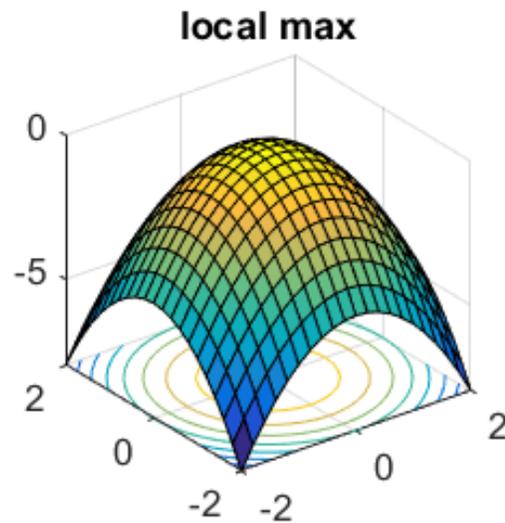
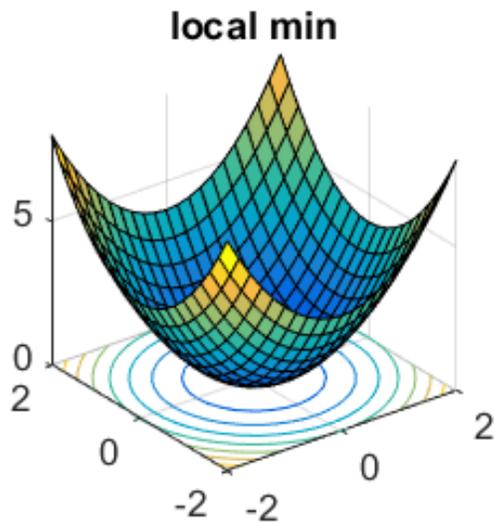


Applied Optimization: Gradients

Alexander G. Ororbia II
Introduction to Machine Learning
CSCI-335
1/26/2026 (*Mon*) & 1/30/2026 (*Fri*)

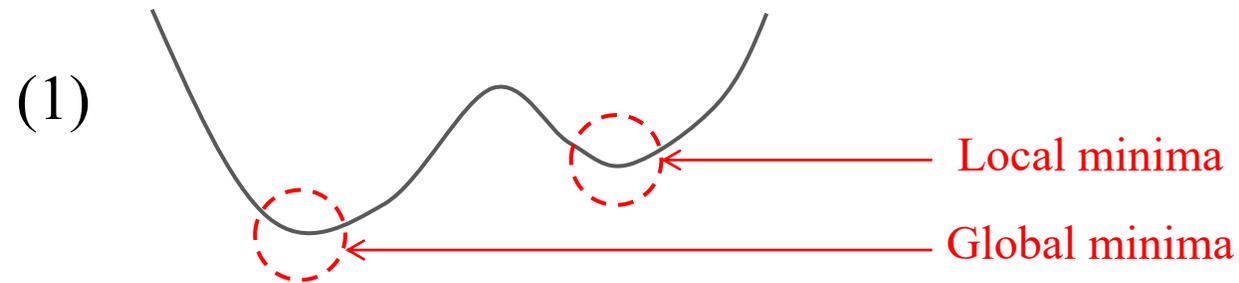
Gradients Means Derivatives; This Means Some Calculus...

- [See Crash Course on Derivatives]

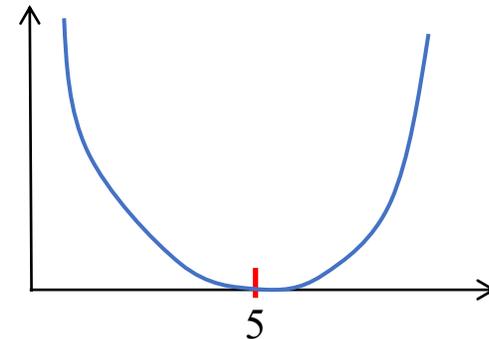


The Value of Derivatives

- How do you solve the following problems?



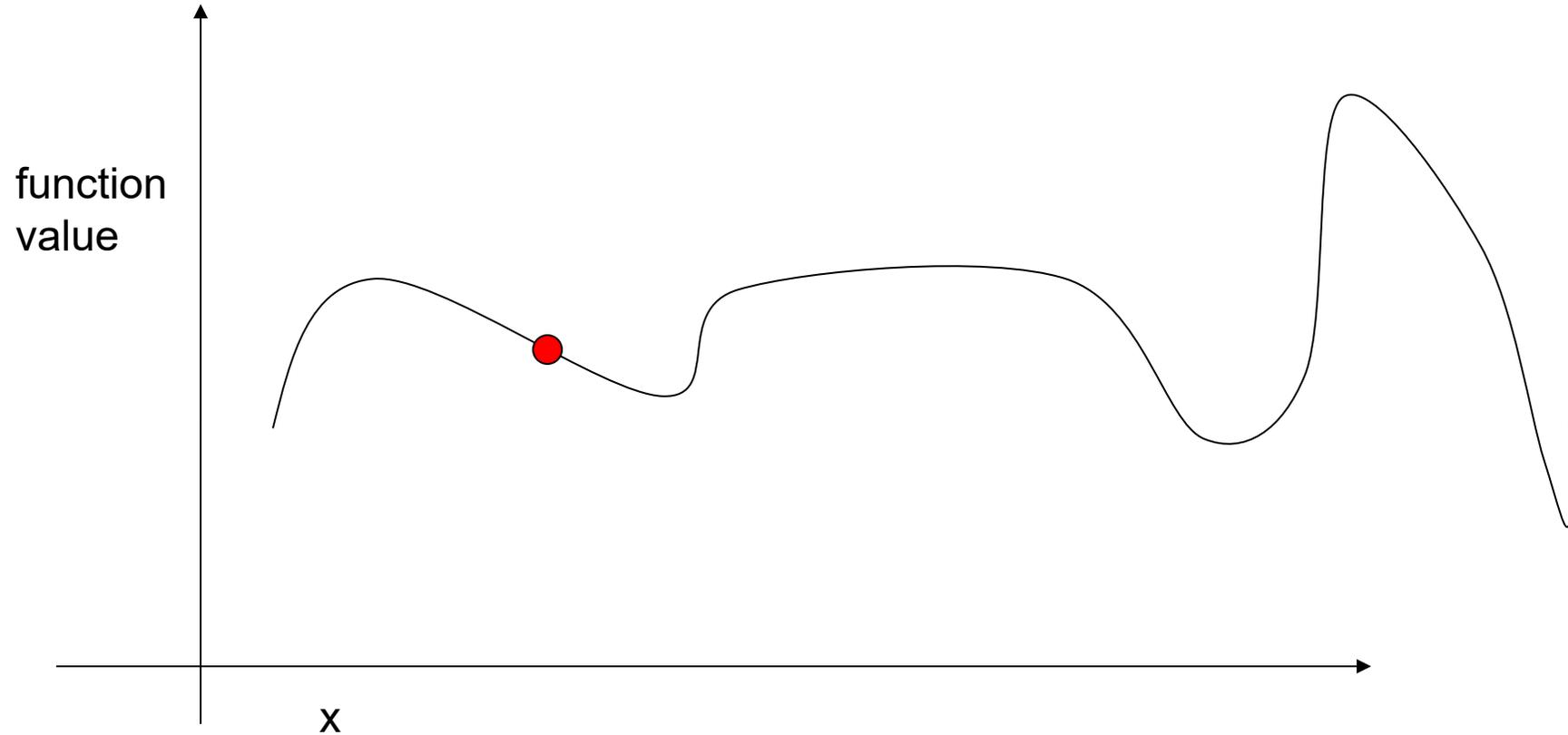
(2) $\min_w f(w) = (w-5)^2$ \implies (a) Plot



(b) Take *derivatives*, check = 0

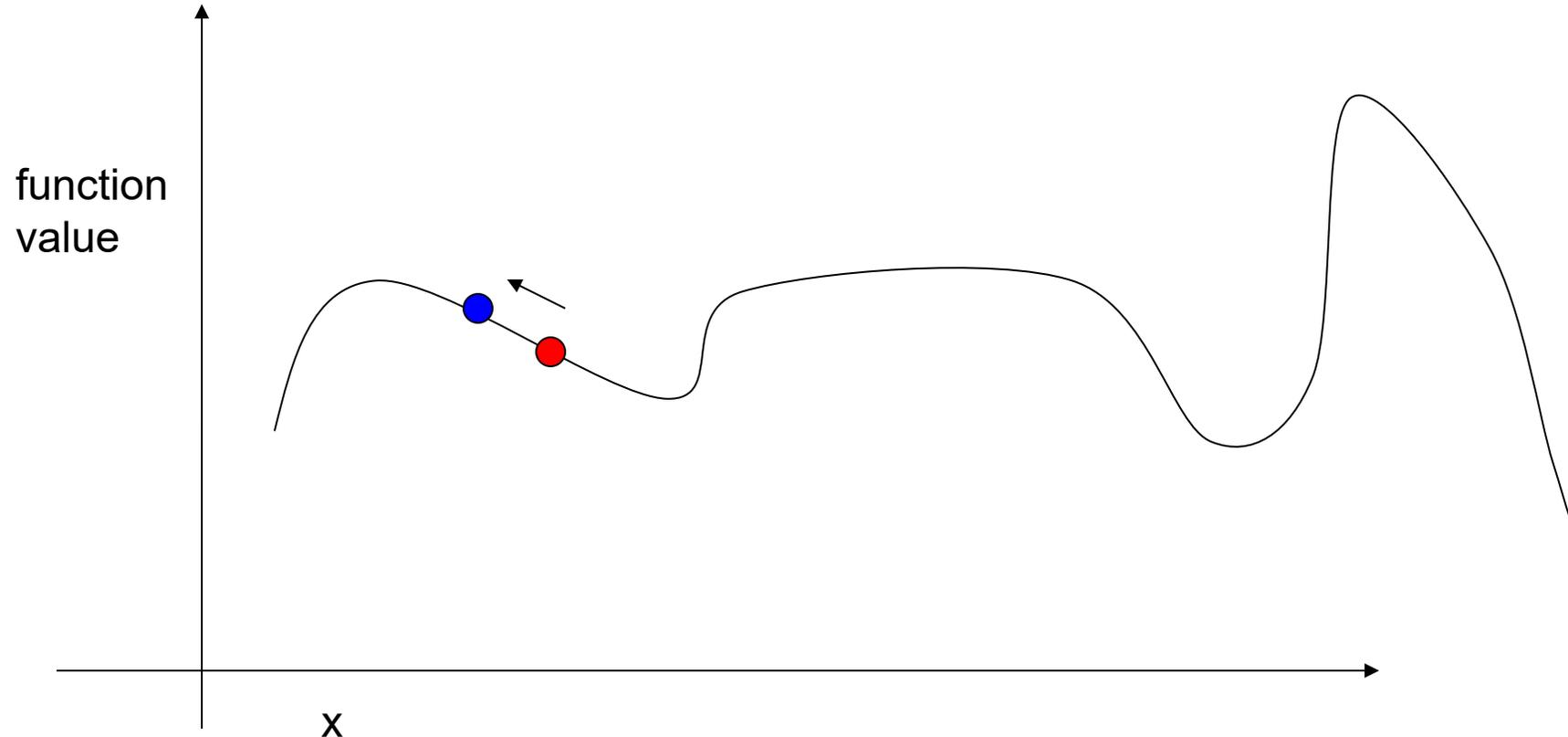
Gradient Ascent

- Random Starting Point



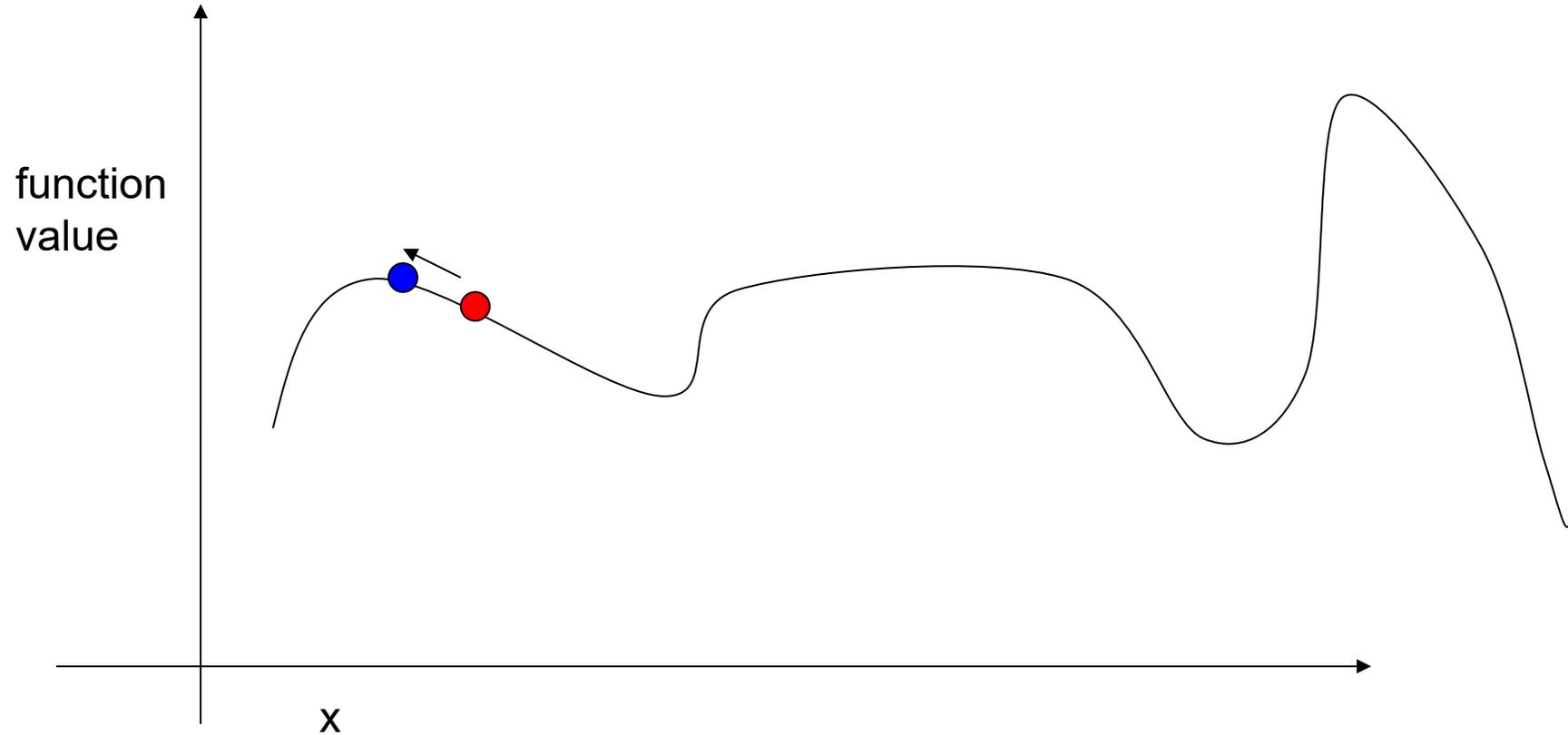
Gradient Ascent

- Take step in direction of largest increase (obvious in 1D, must be computed in higher dimensions)



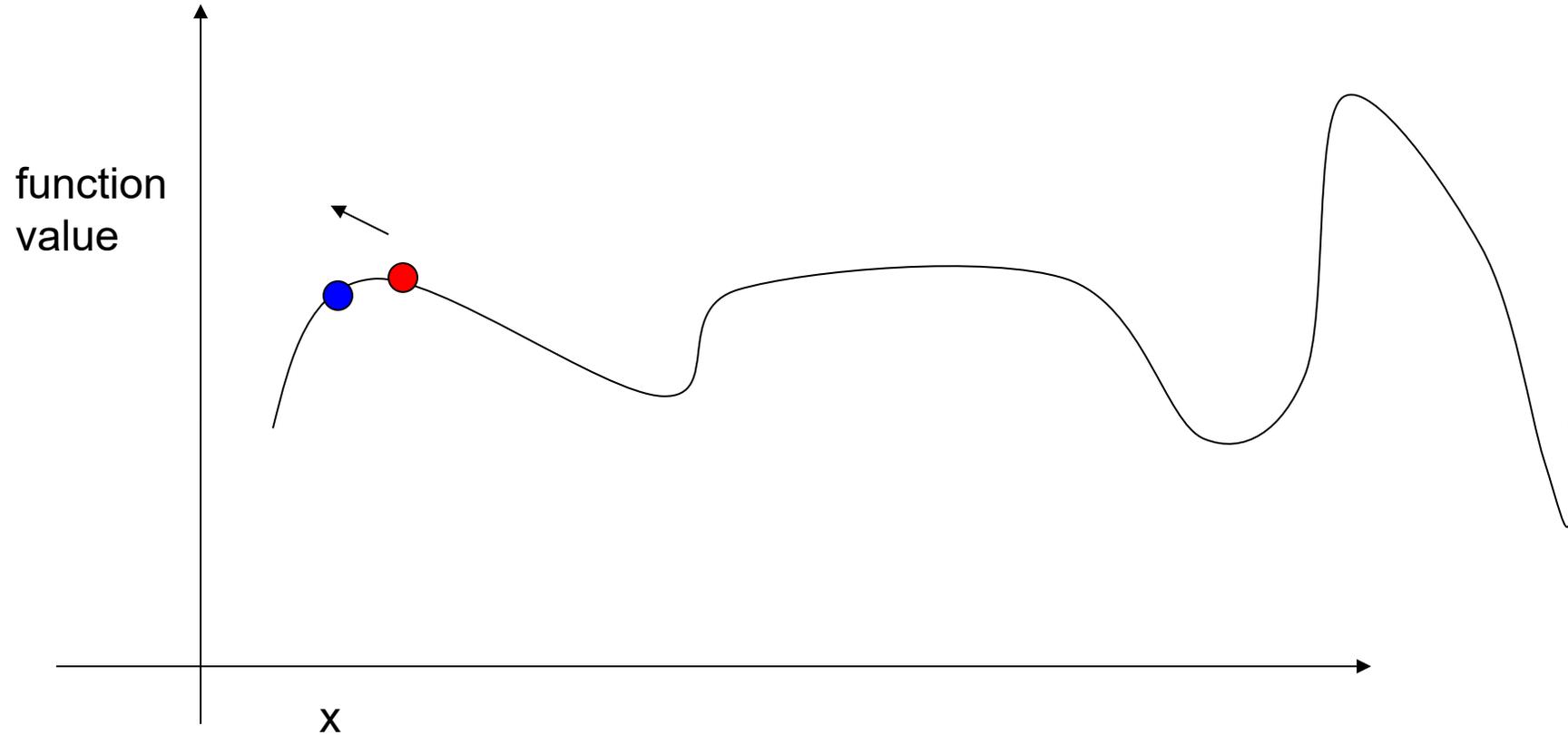
Gradient Ascent

- Repeat



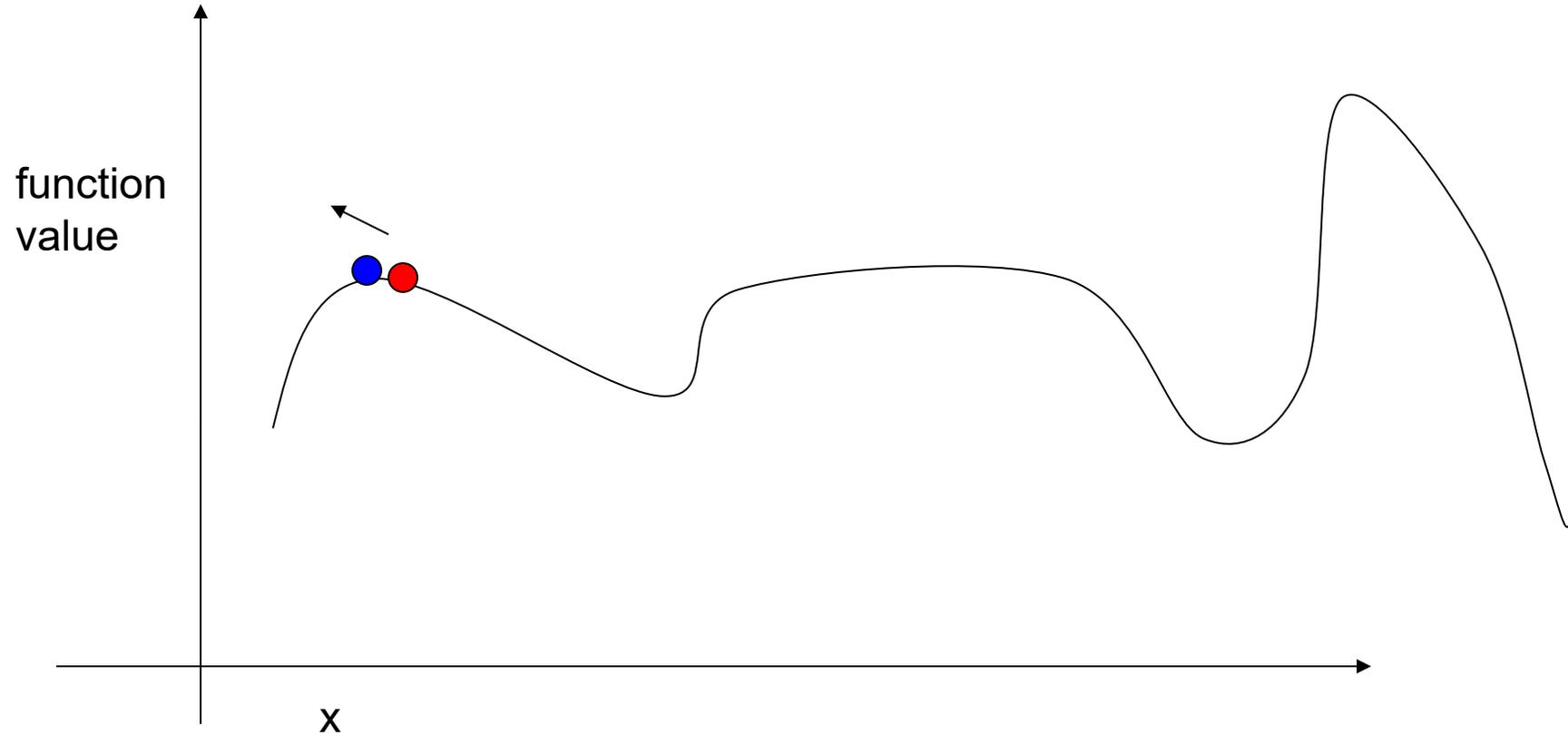
Gradient Ascent

- Next step is actually lower, so stop



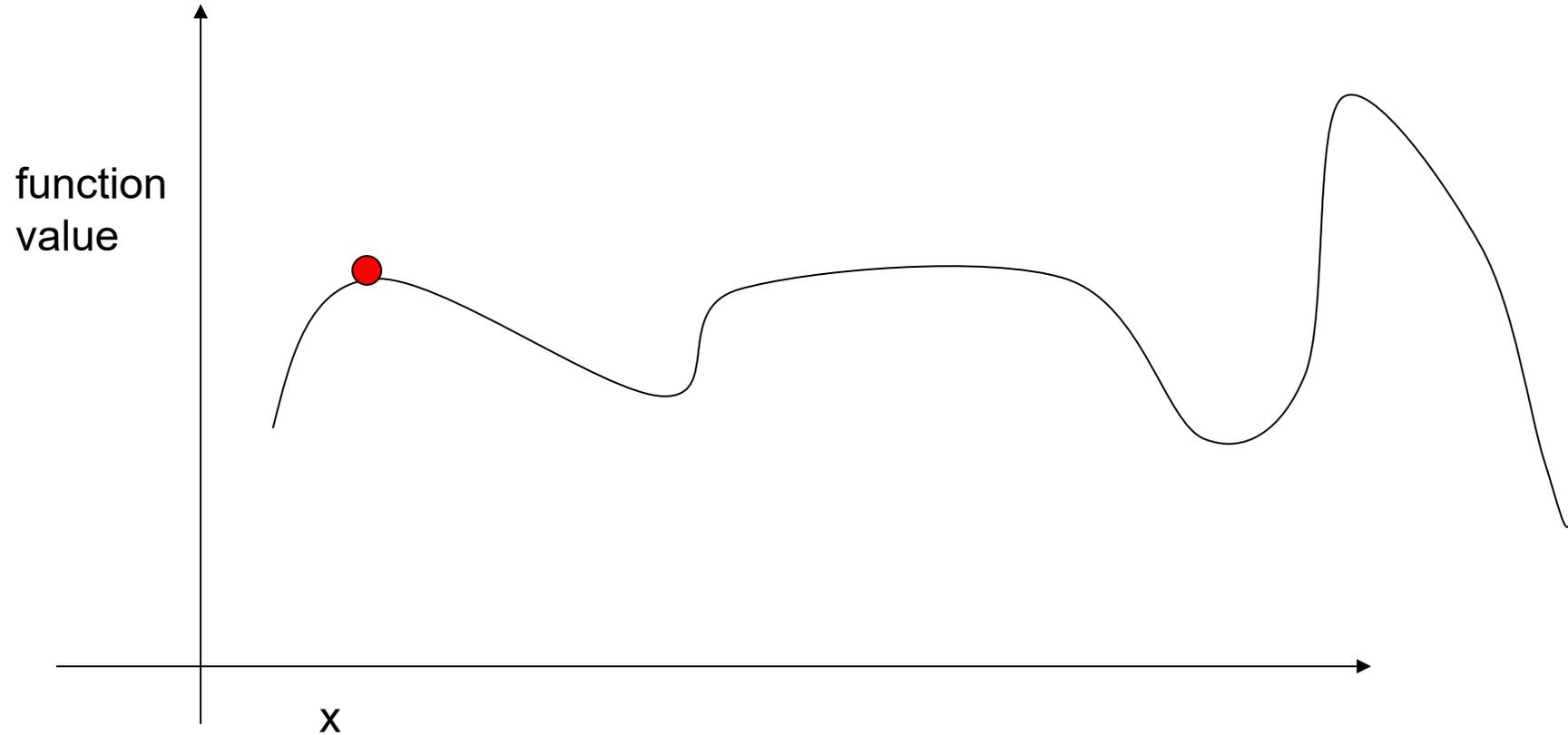
Gradient Ascent

- Could reduce step size to “hone in”



Gradient Ascent

- Converge to (local) maximum



The Finite Difference Method

The centered finite difference approximation is:

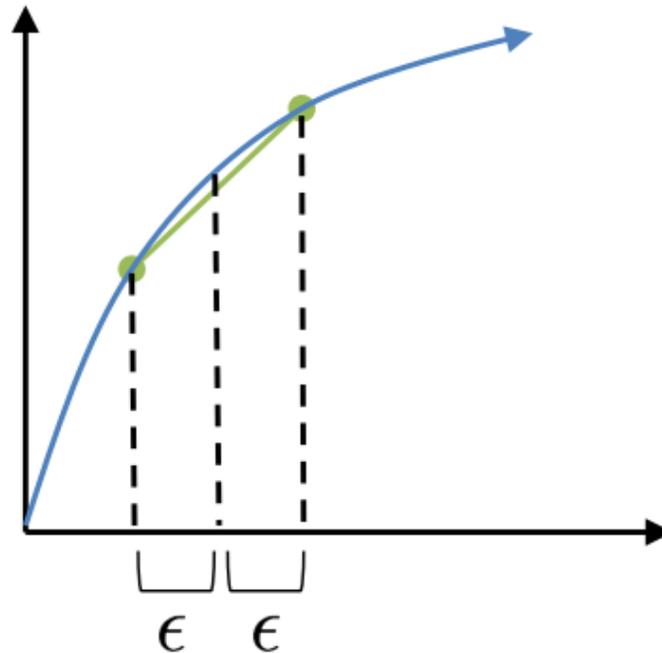
$$\frac{df(x)}{dx} = \lim_{h \rightarrow 0} \frac{f(x+h) - f(x)}{h}$$

$$\frac{\partial}{\partial \theta_i} J(\boldsymbol{\theta}) \approx \frac{(J(\boldsymbol{\theta} + \epsilon \cdot \mathbf{d}_i) - J(\boldsymbol{\theta} - \epsilon \cdot \mathbf{d}_i))}{2\epsilon}$$

where \mathbf{d}_i is a 1-hot vector consisting of all zeros except for the i th entry of \mathbf{d}_i , which has value 1.

Notes:

- Suffers from issues of floating point precision, in practice
- Typically only appropriate to use on small examples with an appropriately chosen epsilon



Ex: Gradient Descent on Least Squares

- Criterion to minimize

$$f(\mathbf{x}) = \frac{1}{2} \|A\mathbf{x} - \mathbf{b}\|^2$$

Evaluation

- Least squares regression

$$E_D(\mathbf{w}) = \frac{1}{2} \sum_{n=1}^N \{t_n - \mathbf{w}^T \phi(x_n)\}^2$$

- The gradient is

$$\nabla_{\mathbf{x}} f(\mathbf{x}) = A^T (A\mathbf{x} - \mathbf{b}) = A^T A\mathbf{x} - A^T \mathbf{b}$$

Optimization

- Gradient Descent algorithm is

1. Set step size ε , tolerance δ to small, positive nos.

2. *while* $\|A^T A\mathbf{x} - A^T \mathbf{b}\|_2 > \delta$ *do*

$$\mathbf{x} \leftarrow \mathbf{x} - \varepsilon (A^T A\mathbf{x} - A^T \mathbf{b})$$

3. *end while*

Calculus in Optimization

- Suppose we have function $y=f(x)$, x, y real nos.
 - Derivative of function denoted: $f'(x)$ or as dy/dx
 - Derivative $f'(x)$ gives the slope of $f(x)$ at point x
 - It specifies how to scale a small change in input to obtain a corresponding change in the output:

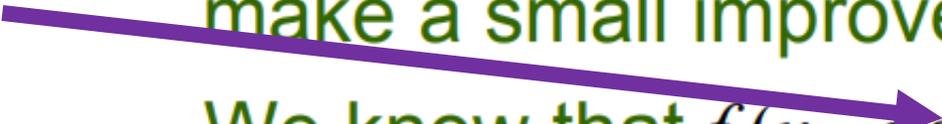
$$f(x + \varepsilon) \approx f(x) + \varepsilon f'(x)$$

– It tells how you make a small change in input to make a small improvement in y

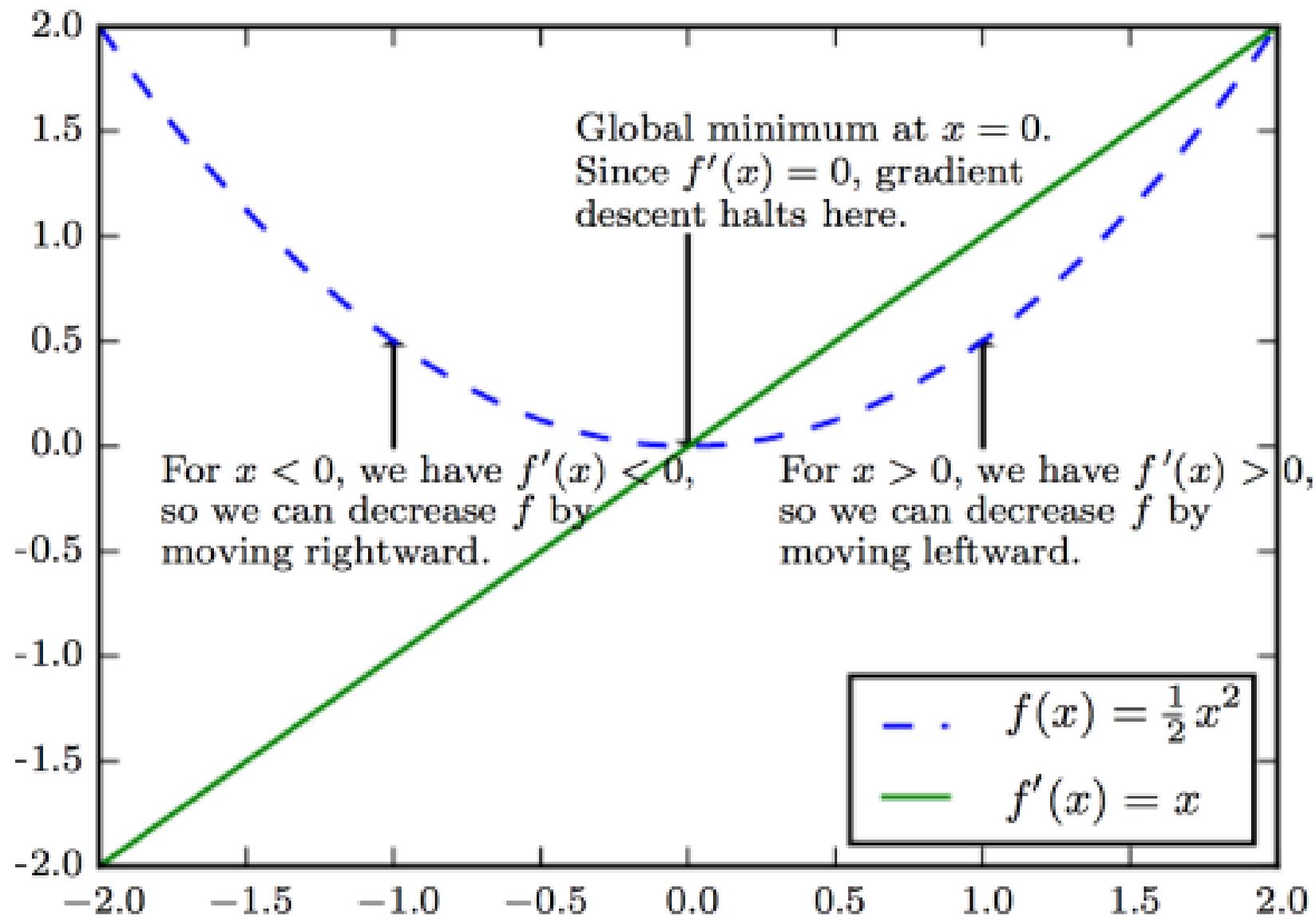
– We know that $f(x - \varepsilon \text{ sign}(f'(x)))$ is less than $f(x)$ for small ε . Thus we can reduce $f(x)$ by moving x in small steps with opposite sign of derivative

- This technique is called *gradient descent* (Cauchy 1847)

signum(v)



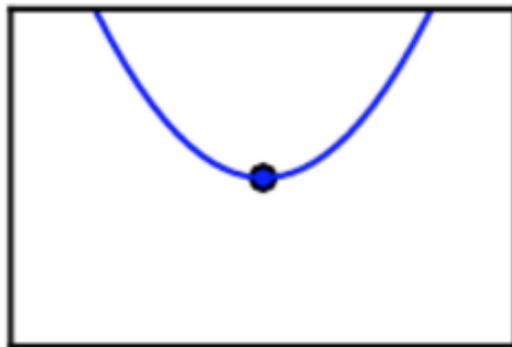
Gradient Descent Illustrated



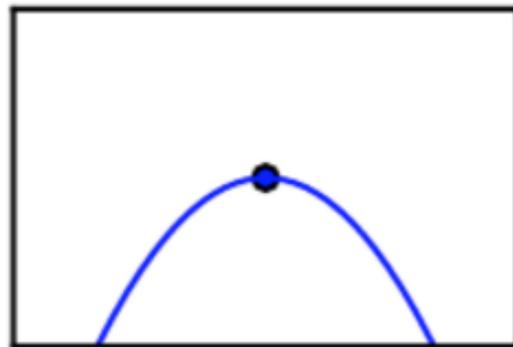
Stationary points, Local Optima

- When $f'(x)=0$ derivative provides no information about direction of move
- Points where $f'(x)=0$ are known as *stationary* or *critical points*
 - Local minimum/maximum: a point where $f(x)$ lower/higher than all its neighbors
 - Saddle Points: neither maxima nor minima

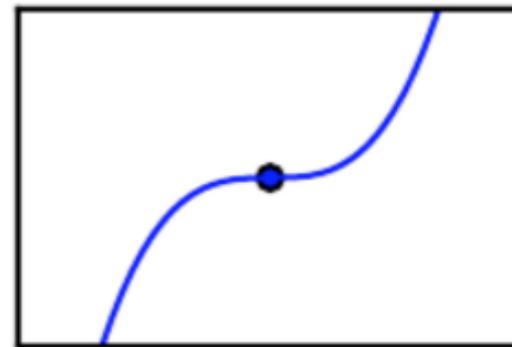
Minimum



Maximum



Saddle point



Functions with Multiple Inputs

- Need partial derivatives
- $\frac{\partial}{\partial x_i} f(\mathbf{x})$ measures how f changes as only variable x_i increases at point \mathbf{x}
- Gradient generalizes notion of derivative where derivative is wrt a vector
- Gradient is vector containing all of the partial derivatives denoted $\nabla_{\mathbf{x}} f(\mathbf{x})$
 - Element i of the gradient is the partial derivative of f wrt x_i
 - Critical points are where every element of the gradient is equal to zero

w.r.t. (with respect to)



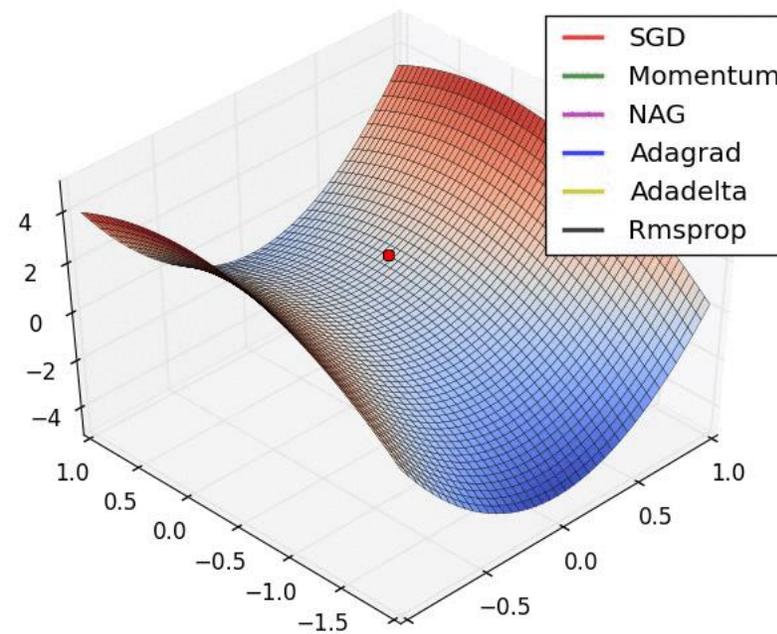
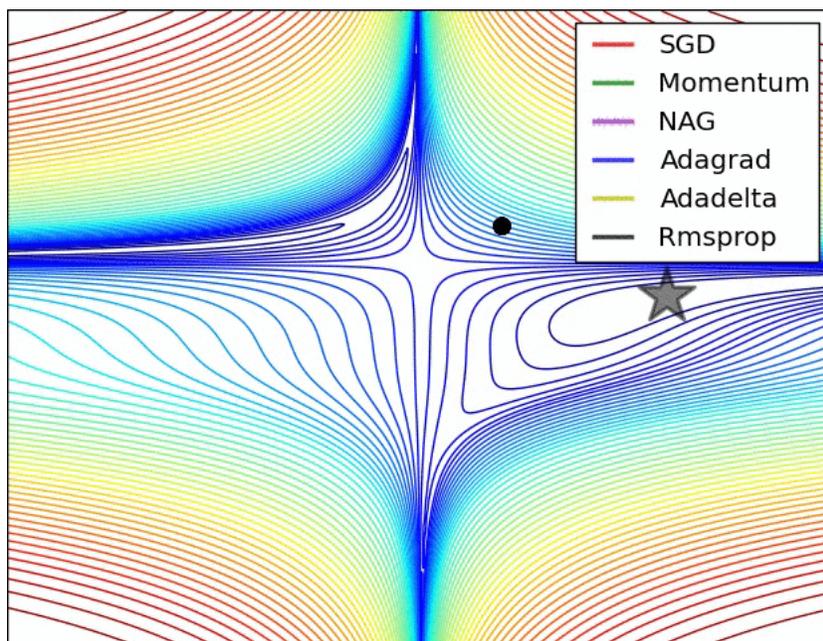
Method of Gradient Descent

- The gradient points directly uphill, and the negative gradient points directly downhill
- Thus we can decrease f by moving in the direction of the negative gradient
 - This is known as the method of steepest descent or gradient descent
- Steepest descent proposes a new point

$$\mathbf{x}' = \mathbf{x} - \varepsilon \nabla_{\mathbf{x}} f(\mathbf{x})$$

- where ε is the learning rate, a positive scalar. Set to a small constant.

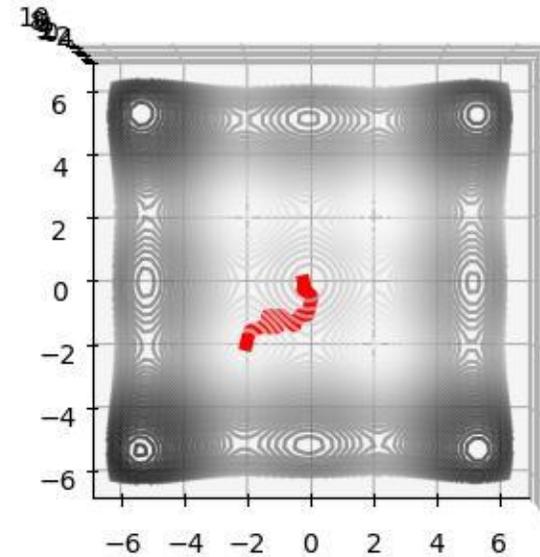
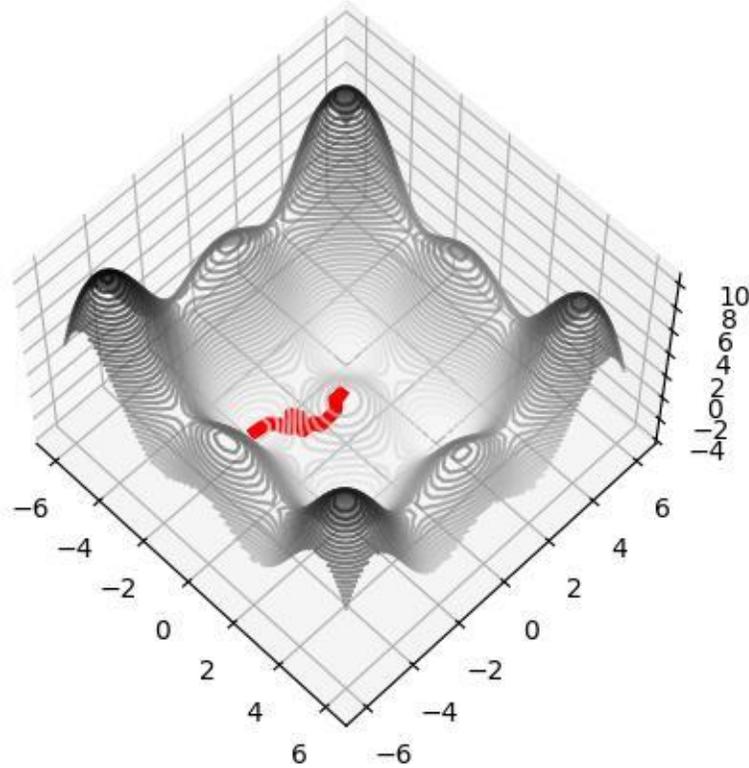
Race of the Optimizers!



<http://cs231n.github.io/neural-networks-3/#hyper>

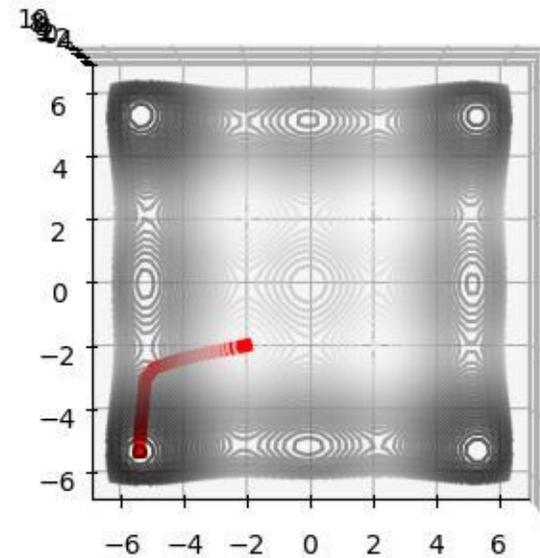
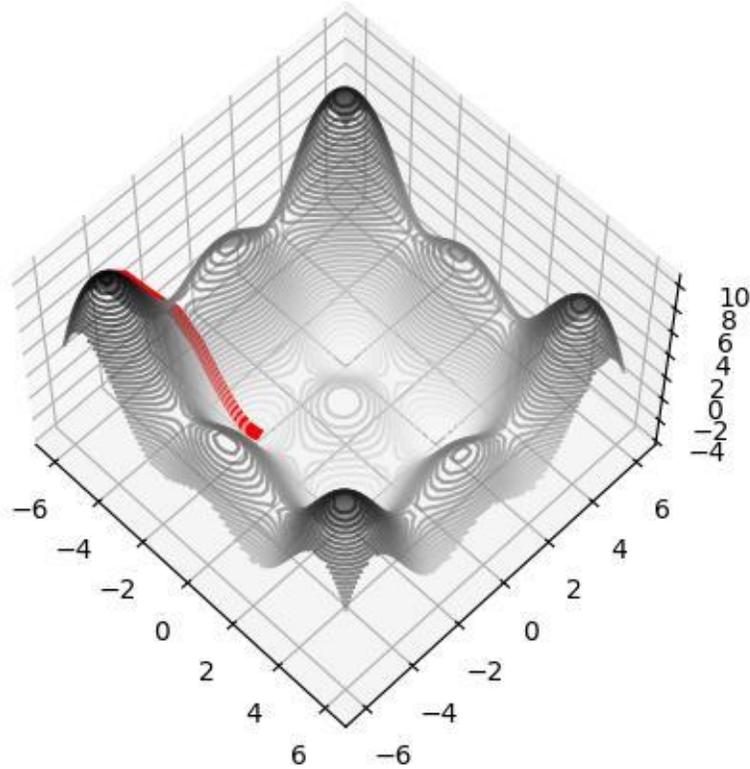
Stochastic Hill-Climbing

Negative Alpine Function $f(x) = - \sum_i (x_i \sin(x_i) + 0.1 x_i)$

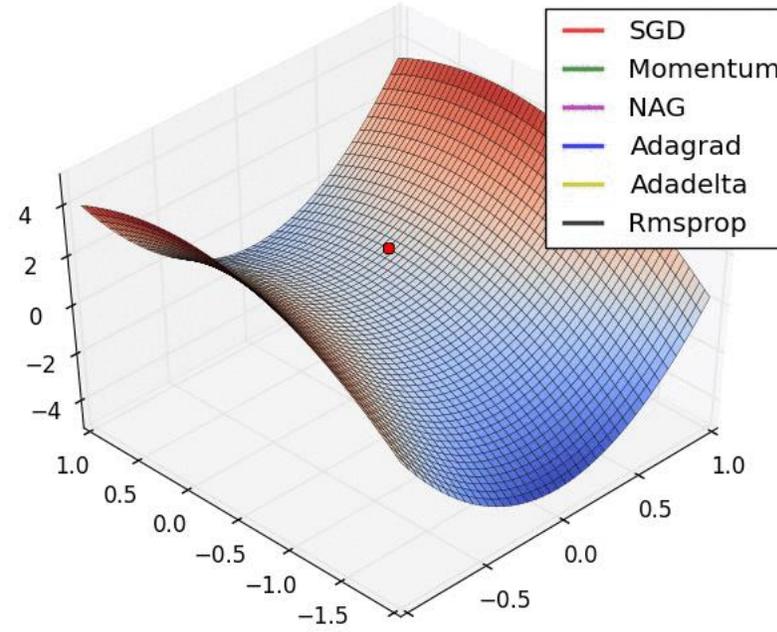
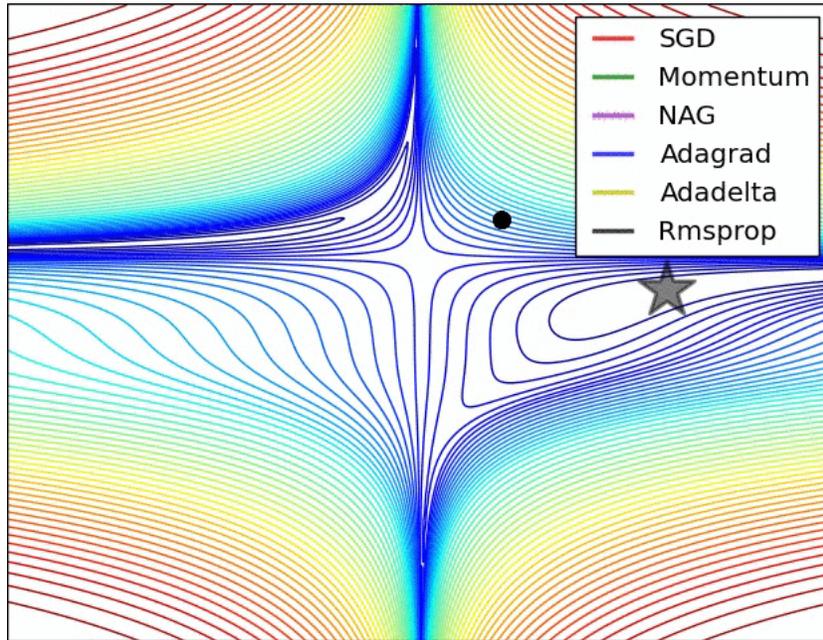


Stochastic Gradient Ascent

Negative Alpine Function $f(x) = - \sum_i (x_i \sin(x_i) + 0.1 x_i)$



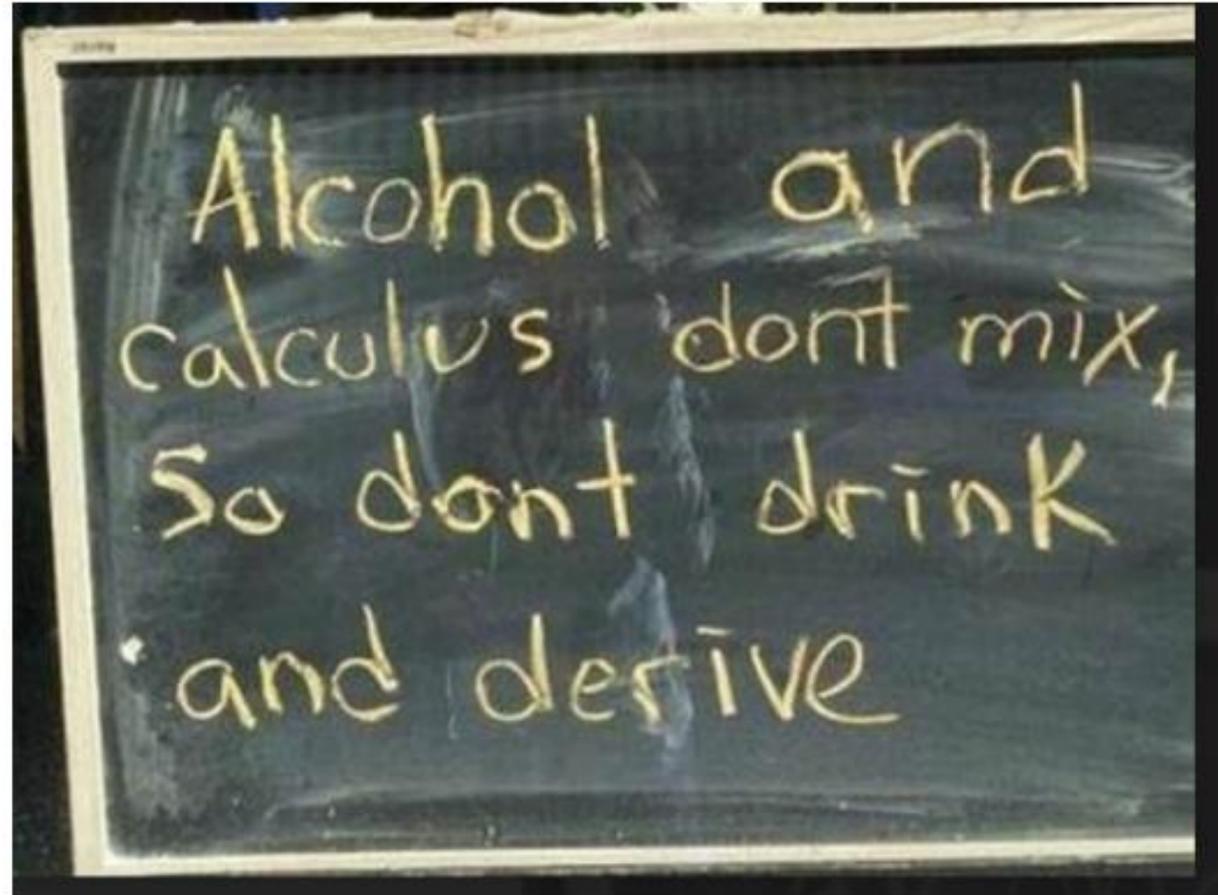
Race of the Optimizers!



<http://cs231n.github.io/neural-networks-3/#hyper>

Gradient

- Essential role of calculus



You want to build a simple univariate regression model for predicting profits y for a food truck. Furthermore, you decide to restrict yourself to a linear hypothesis space and construct a model that adheres to the following form:

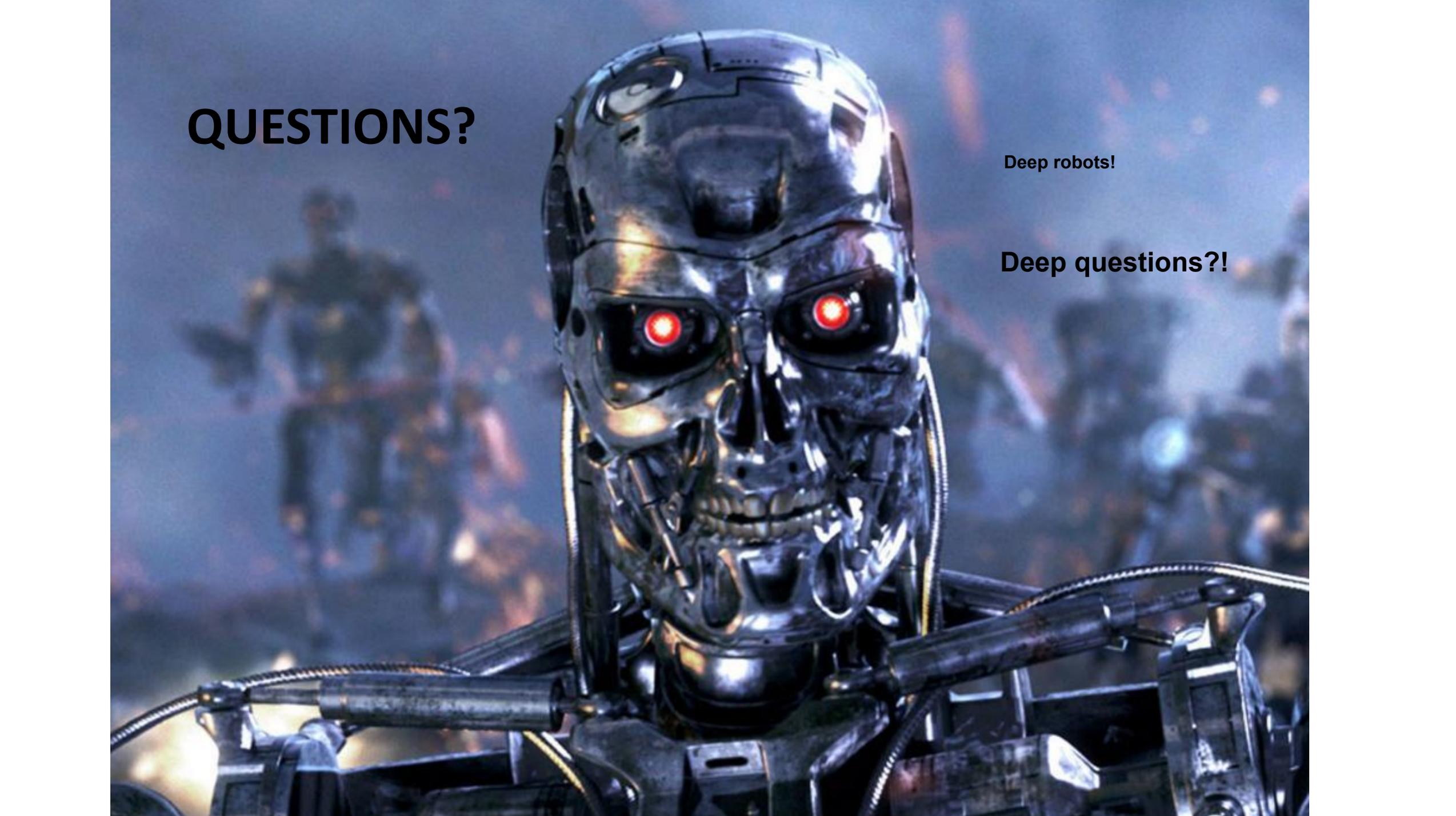
$$f_{\Theta}(x) = \theta_0 + \theta_1 x$$

Given the data you have collected, represented as a set of m complete (y, x) pairs, your goals will be to estimate the parameters of this model $\Theta = \{\theta_0, \theta_1\}$ (where $\theta_{j>0}$ is the vector of learnable coefficients that weight the observed variables, and θ_0 is a single bias coefficient) using the method of steepest gradient descent. The cost function to minimize is the well-known mean squared error (MSE) defined as follows:

$$\mathcal{J}(\Theta) = \frac{1}{2m} \sum_{i=1}^m (f_{\Theta}(x^i) - y^i)^2$$

What is a useful constrained version of this cost to get “smaller” coefficients?

(Hint: Tikhonov regularization)



QUESTIONS?

Deep robots!

Deep questions?!