

L665 Optimization

Damir Cavar
Indiana University
Feb. 2018

Agenda

Learning

Optimization

Example: K-means

Gradient Descent

Machine Learning as Optimization

Machine Learning Algorithms in Groups:

- Supervised learning
- Unsupervised learning
- Reinforcement learning

Supervised Learning

Learning using annotations in data, linguistic corpora

- Brown corpus
- Penn treebank
- Etc.

Learning algorithms learn:

- associations between data and target responses

Supervised Learning

Responses could be

- the annotations themselves, or
- secondary meta-information associated with the annotated data sets

Good examples via training set.

Generalization and application to test set or unknown data.

Supervised Learning

Problems addressed:

- Regression problem: target is a numerical value
- Classification problem: target is a qualitative variable (PoS-tag, sentiment label, etc.)

Unsupervised Learning

Data to learn from does not contain annotations or class meta-information.

Mapping of data to:

- Some features (also vectors) that then could represent classes or some numerical value.

Reinforcement Learning

As in Unsupervised Learning:

- Only unlabeled data is available.

Provide feedback about the classification that the unsupervised algorithm proposes:

- Negative or positive feedback
- Trial and error learning

Reinforcement Learning

Example:

- Google DeepMind lets an algorithm learn how to play Atari video games:
 - <https://www.youtube.com/watch?v=V1eYniJ0Rnk>

Logic of Machine Learning Algorithms

Assumption:

- There is an unknown mathematical function or formulation and data.
- The mathematical function or formulation can be guessed from data.

Supervised Example:

- Data and features (annotations) with a class assignment.

Logic of Machine Learning Algorithms

Training:

- See many examples of feature and data combinations for a given class and optimize to identify correctly unknown data to belong to any of the classes with some probability.

Classifier in ML:

- Building a *target function* that can express or identify the characteristics of some class given the observed data and its features.
- *Mapping* is the discovery of such a function by observing the outputs.

Optimization

Parameters of the learning algorithm:

- Model (distributions, associations of input output tuples, vectors, metrics)

Internal parameters have dimensions, type and so on:

- The dimensions determine the target function.

Optimization:

- Assuming some initial state of parameters.
- Changing parameters or identifying the best combination that maps input with output.

Optimization

Hypothesis space:

- Set of all potential target functions (parameter set values or properties) that the learning algorithm can identify for a given data set.

Non learnable parameters:

- Called *hyper-parameters*, are provided by the ML engineer.

Cost-functions (loss function, objective function, scoring function, error function):

- Measure for how well a learning algorithm maps the target function.

Optimization

Motivation:

- Parameters (or coefficients) of a cost function cannot be computed analytically.
- Must be searched for via optimization.

Optimization

Most algorithms:

- Optimize a given cost function.

Some algorithms:

- Provide a set of possible cost functions.

There might be a mismatch between learning algorithm optimization objectives and desired objectives wrt. performance:

- Often expressed by using *error* or *loss function*

Optimization

For algorithms with fixed cost functions:

- One can tweak the hyper-parameters and select specific features to drive the result towards a desired objective.

Cost functions:

- Guiding learning algorithm in choice of features that lead to better predictions of the model or parameter choice.
- Optimization results from iterative application of learning algorithms and cost function guided improvements of the parameters.

Optimization Strategies

Examples:

- Clustering using K-means
- Gradient Descent

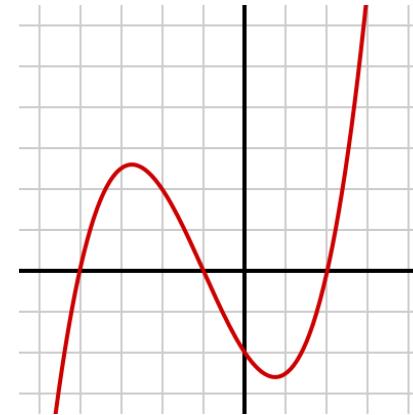
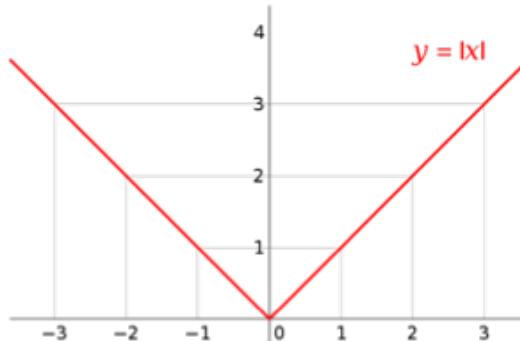
Gradient Descent

- See tutorial and slides by Kris Hauser
http://homes.soic.indiana.edu/classes/spring2012/csci/b553-hauserk/gradient_descent.pdf
- <http://homes.soic.indiana.edu/classes/spring2012/csci/b553-hauserk/05gradientdescent.pptx>

Gradient Descent

Given a differentiable scalar field $f(x)$ (Graphs from Wikipedia)

- A field associates every scalar value to every point in the space, and addition, subtraction, multiplication, division are defined.
- A differentiable function is one whose derivative exists at each point in its domain (it has no breaks, bends, cusps).



Gradient Descent

Initialization:

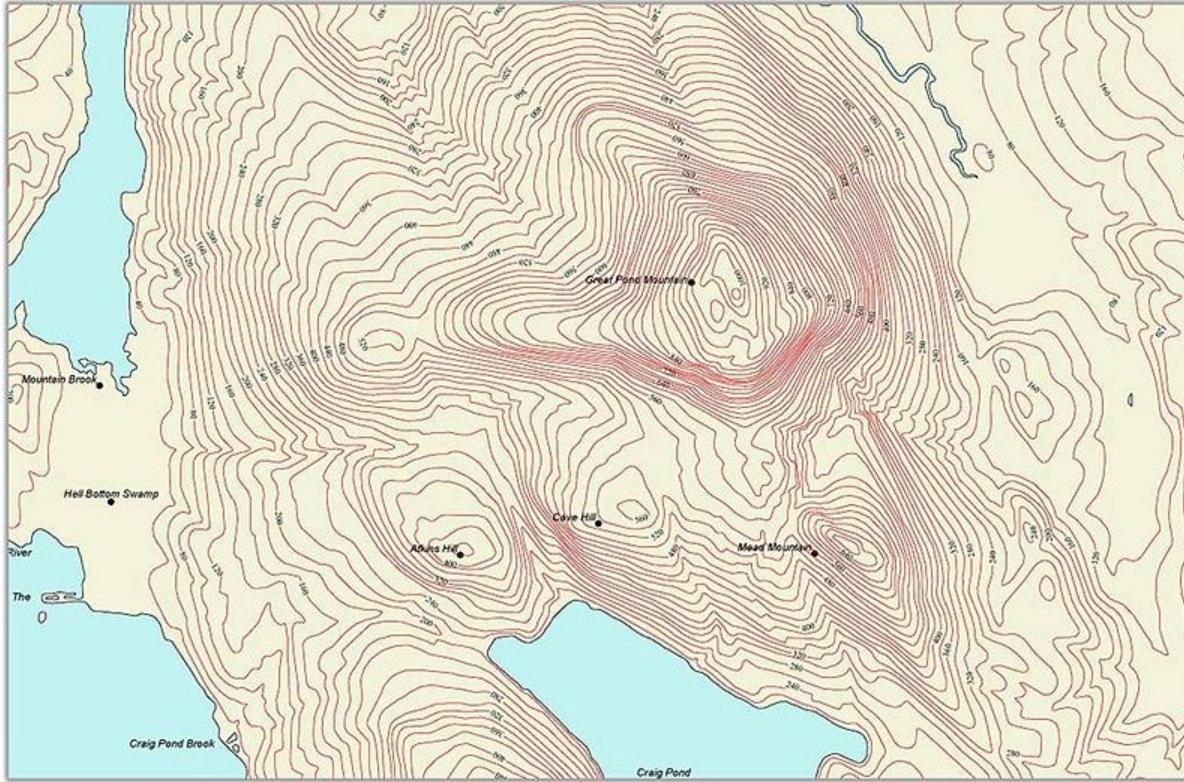
- Initial guess x_1

GD iteratively moves the guess towards lower values of f

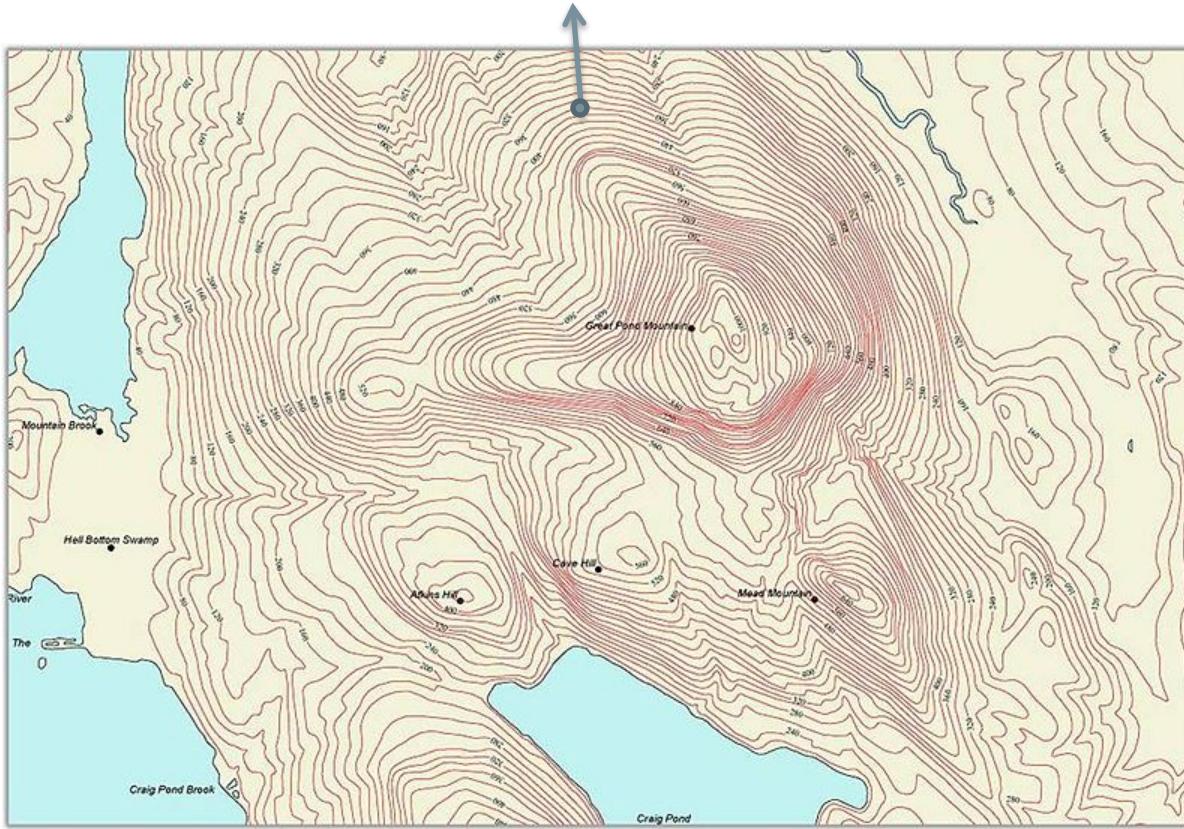
Taking steps in the direction of the negative gradient $-\nabla f(x)$

Locally the negative gradient is the steepest descent direction

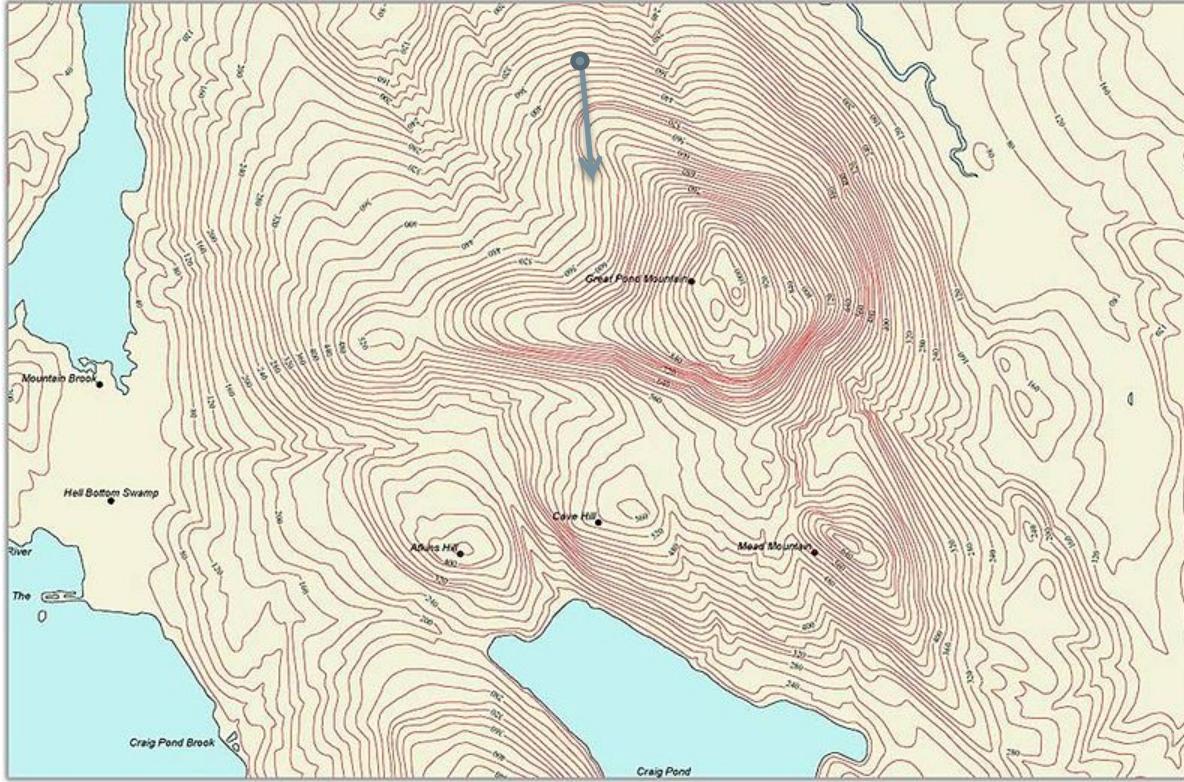
Following images from Kris Hauser's slides, see previous slides:



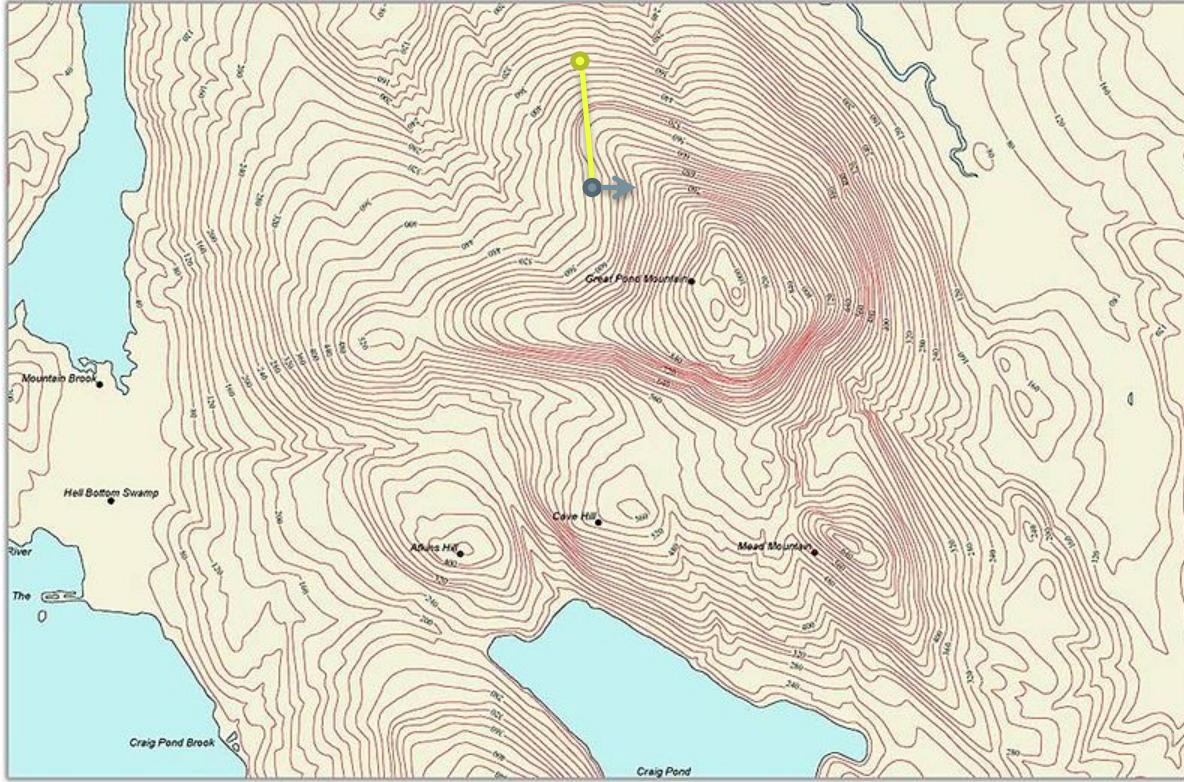
Gradient direction ∇f is orthogonal to the level sets (contours) of f , points in direction of steepest increase



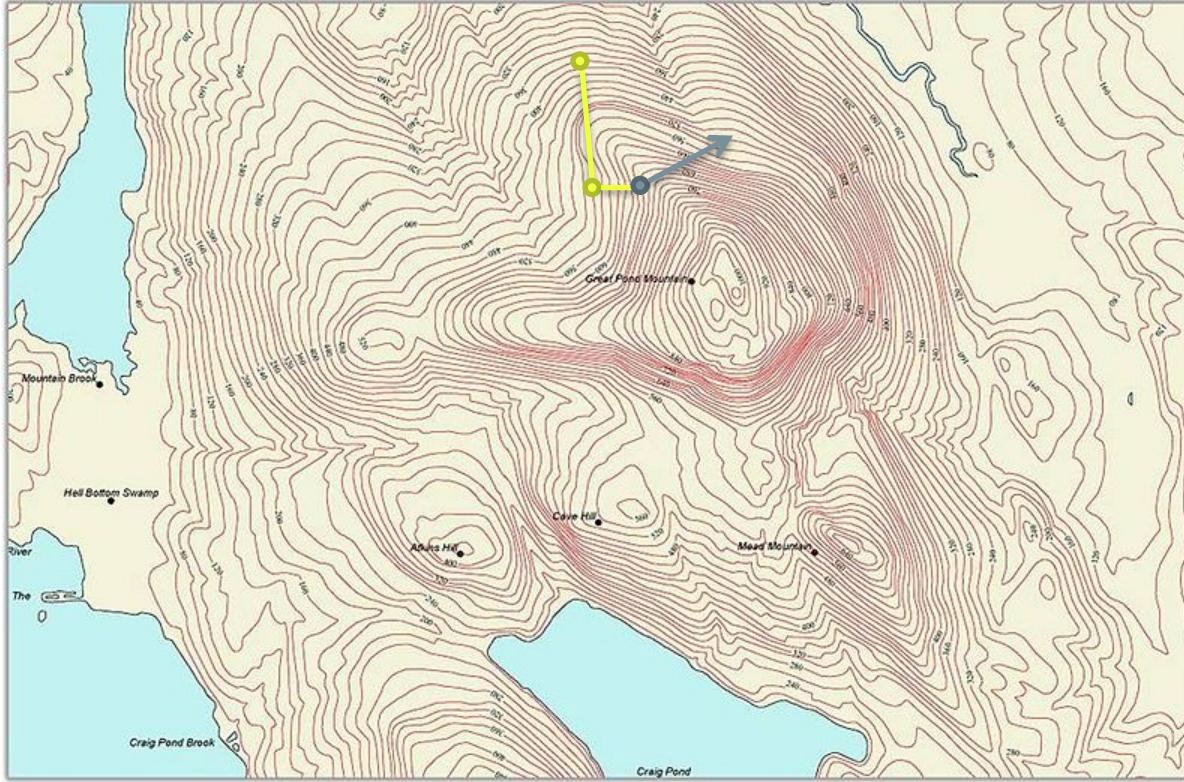
Gradient direction ∇f is orthogonal to the level sets (contours) of f ,
points in direction of steepest increase



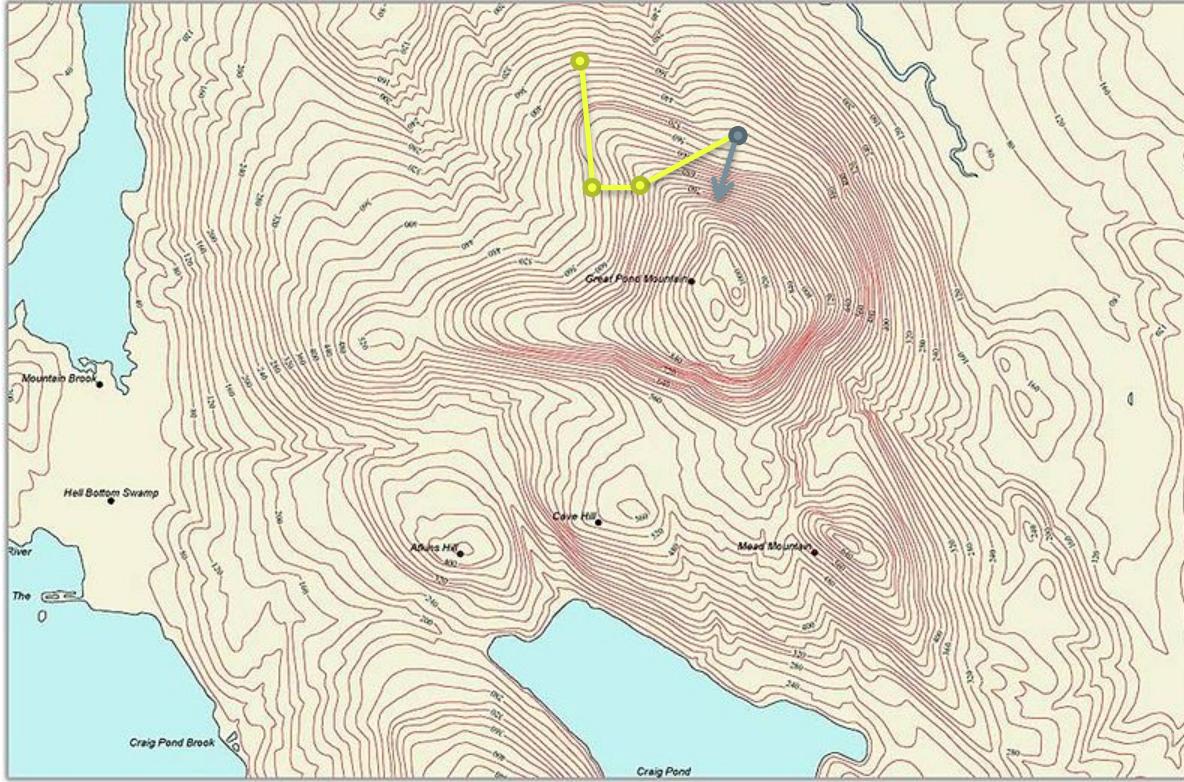
Gradient descent: iteratively move in direction $-\nabla f$



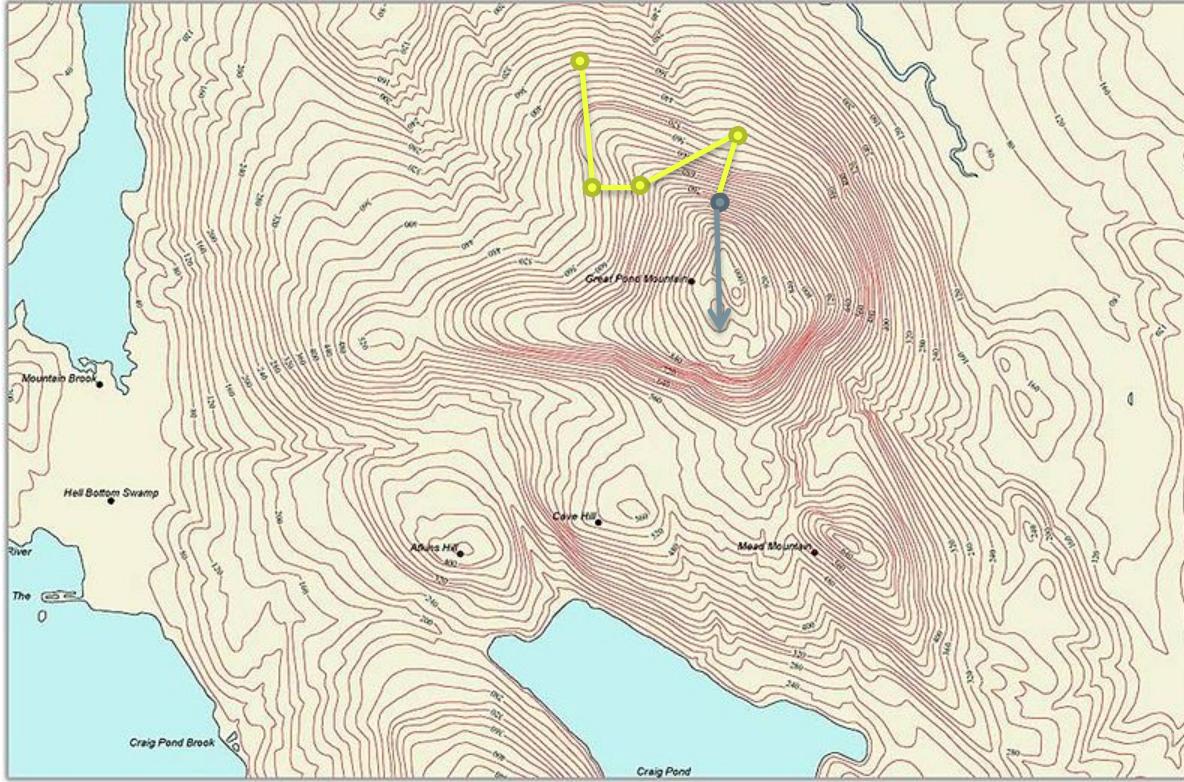
Gradient descent: iteratively move in direction $-\nabla f$



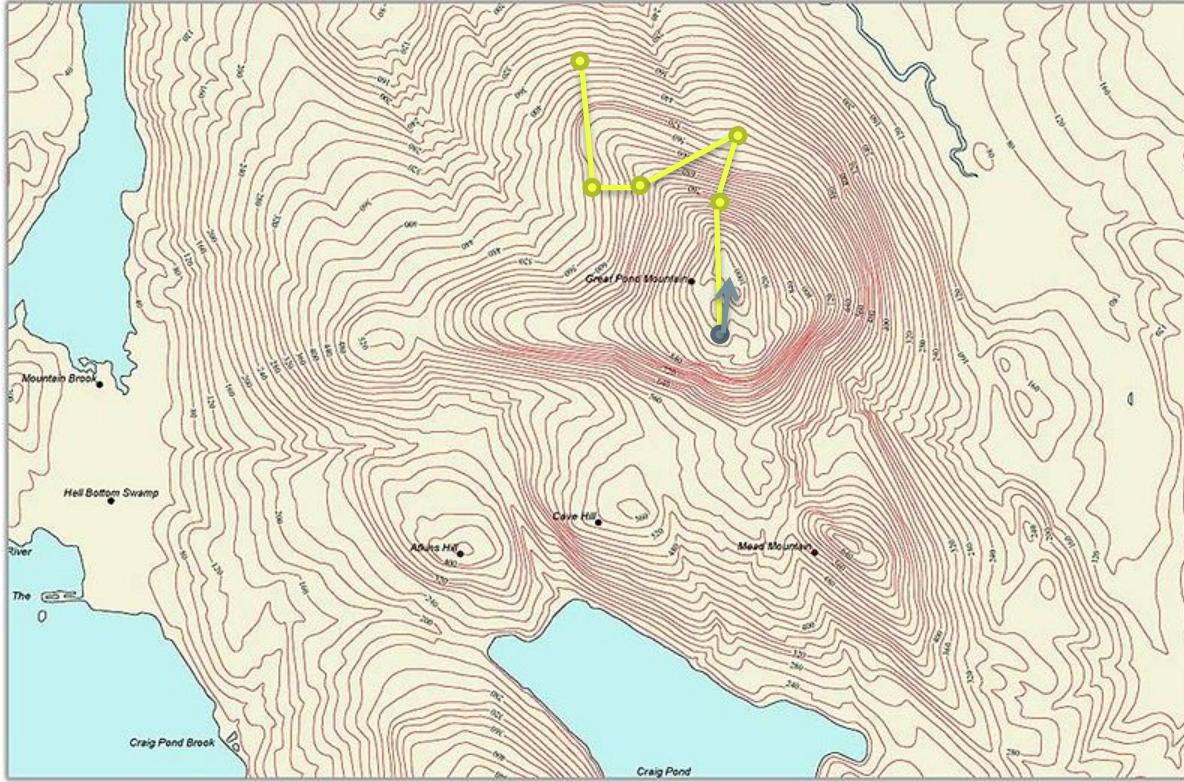
Gradient descent: iteratively move in direction $-\nabla f$



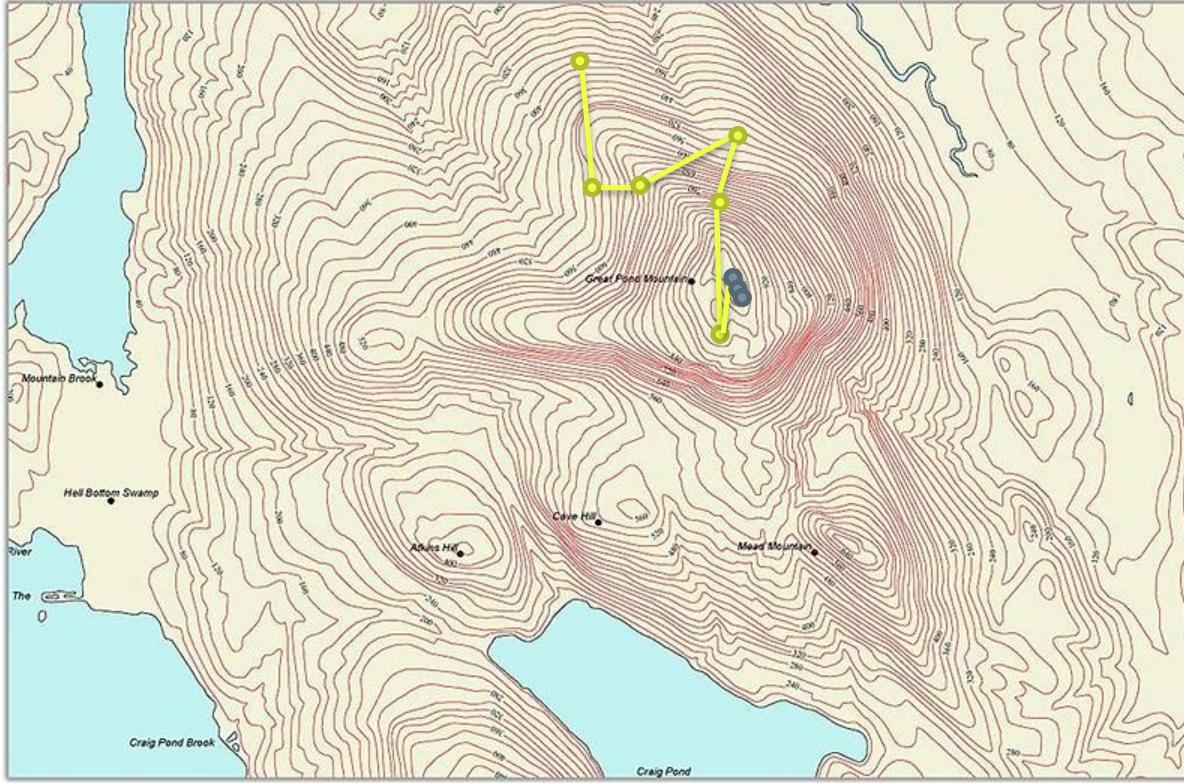
Gradient descent: iteratively move in direction $-\nabla f$



Gradient descent: iteratively move in direction $-\nabla f$



Gradient descent: iteratively move in direction $-\nabla f$



Gradient descent: iteratively move in direction $-\nabla f$

Gradient Descent

Norm tolerance: ϵ_g

If the algorithm has arrived at critical point.

Step tolerance: ϵ_x

If significant progress has been made.

Step size: α_t

Should be chosen to maintain a balance between convergence speed and avoiding divergence. Note that α_t may depend on the step t .

Gradient Descent

Input: f , starting value x_1 , termination tolerances

- For $t=1,2,\dots,N_{\max}$:
 - Compute the search direction $d_t = -\nabla f(x_t)$
 - If $\|d_t\| < \varepsilon_g$ then:
 - return “Converged to critical point”, output x_t
 - Find α_t so that $f(x_t + \alpha_t d_t) < f(x_t)$ using line search
 - If $\|\alpha_t d_t\| < \varepsilon_x$ then:
 - return “Converged in x ”, output x_t
 - Let $x_{t+1} = x_t + \alpha_t d_t$
- Return “Max number of iterations reached”, output $x_{N_{\max}}$

Gradient Descent

Line Search

- Input
 - function f , an initial point \mathbf{x} and direction \mathbf{d}
 - initial candidate step size γ_1
- Output
 - Step size $\gamma > 0$ such that $f(\mathbf{x} + \gamma\mathbf{d}) < f(\mathbf{x})$, or failure if the step size falls below some threshold.
- First step: check whether $\mathbf{x} + \gamma_1\mathbf{d}$ decreases the value of f .
 - If true, then:
 - terminate the line search, or
 - search for a larger valid γ by repeated doubling: $\gamma_n = \gamma_1 2^n$

Gradient Descent

- If a step of size γ_1 increases f , search through repeated halving: $\gamma_n = \gamma_1 / 2^n$.
- Continues until a reduction in f is found, or minimum step size tolerance is reached.
- Why does line search not find the optimal step size rather than any one that reduces the function value?
 - In practice: not worthwhile the extra work to obtain an optimal line search.

