

Machine Learning

Lecture 5: Optimization

Prof. Dr. Stephan Günnemann

Data Mining and Analytics
Technical University of Munich

19.11.2018

Reading material

Reading material

- Boyd - Convex Optimization: chapters 2.1 – 2.3, 3.1, 3.2, 4.1 – 4.4, 9
 - free pdf version online
- Sebastian Ruder - An overview of gradient descent optimization algorithms
 - <https://arxiv.org/abs/1609.04747>

Motivation

- Many data mining/machine learning tasks are optimization problems
- Examples we've already seen:
 - Linear Regression $\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} \frac{1}{2} (\mathbf{X}\mathbf{w} - \mathbf{y})^T (\mathbf{X}\mathbf{w} - \mathbf{y})$
 - Logistic Regression $\mathbf{w}^* = \underset{\mathbf{w}}{\operatorname{argmin}} -\ln p(\mathbf{y}|\mathbf{w}, \mathbf{X})$
- Other examples:
 - Support Vector Machines: find hyperplane that separates the classes with a maximum margin
 - k-means: find clusters and centroids such that the squared distances is minimized
 - Matrix Factorization: find matrices that minimize the reconstruction error
 - Neural networks: find weights such that the loss is minimized
 - And many more...

General task

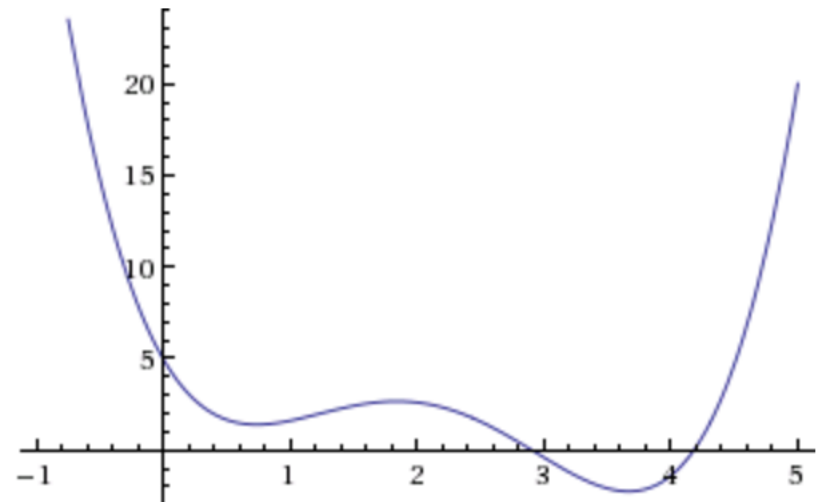
- Let θ denote the variables/parameters of our problem we want to learn
 - e.g. $\theta = \mathbf{w}$ in Logistic Regression
- Let \mathcal{X} denote the domain of θ ; the set of valid instantiations
 - constraints on the parameters!
 - e.g. \mathcal{X} = set of (positive) real numbers
- Let $f(\theta)$ denote the **objective function**
 - e.g. f is the negative log likelihood
- Goal: Find solution θ^* minimizing function f : $\theta^* = \operatorname{argmin}_{\theta \in \mathcal{X}} f(\theta)$
 - find a global minimum of the function f !
 - similarly, for some problems we are interested in finding the maximum

Introductory example

- Goal: Find minimum of function

$$f(\theta) = 0.6 * \theta^4 - 5 * \theta^3 + 13 * \theta^2 - 12 * \theta + 5$$

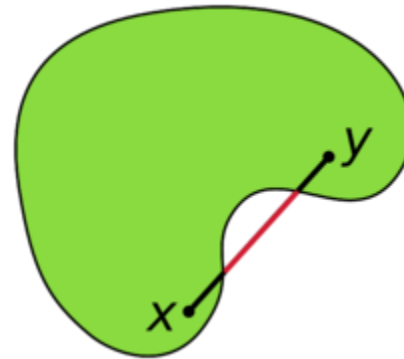
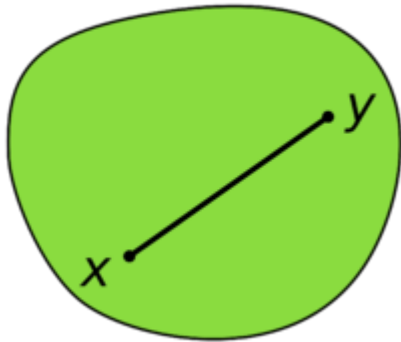
- Unconstrained optimization + differentiable function
- Necessary condition for minima
 - Gradient = 0
 - Sufficient?



- General challenge: multiple local minima possible

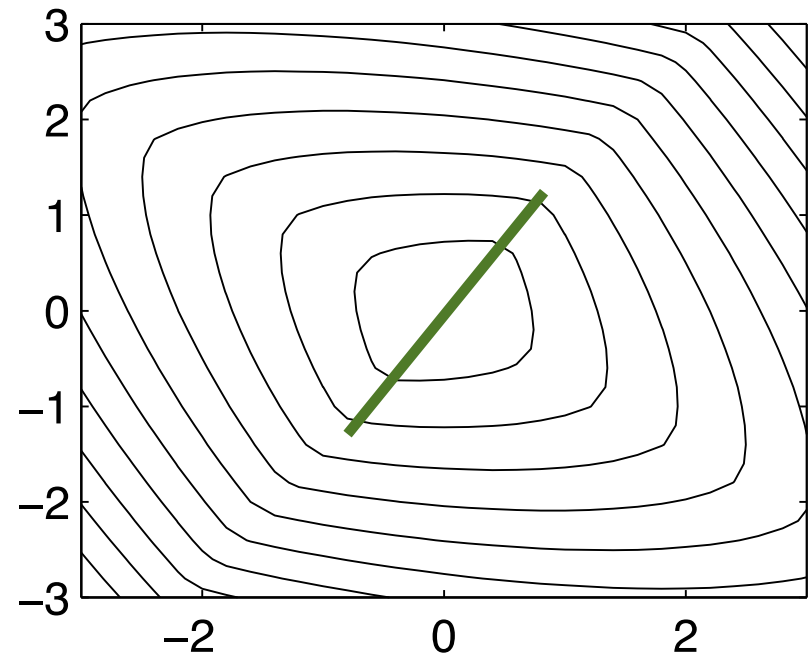
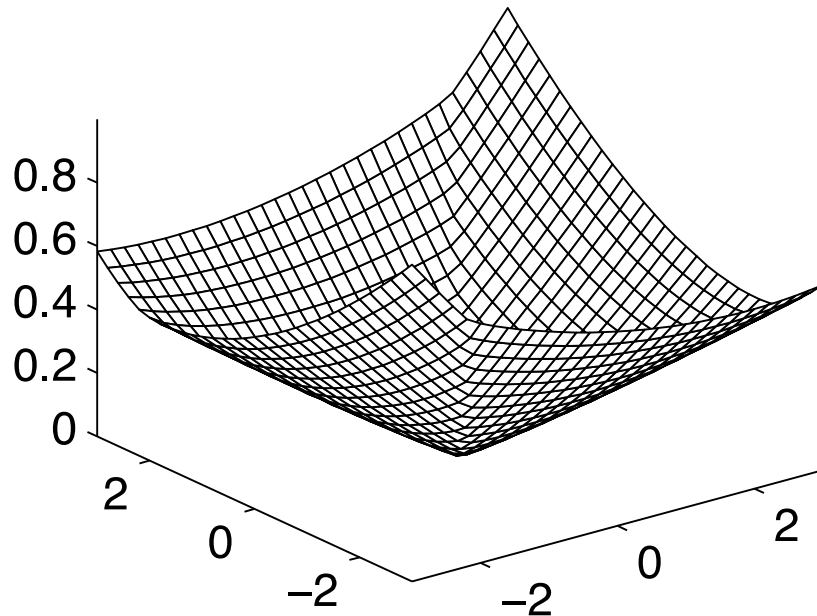
Convexity: Sets

- X is a convex set
iff
for all $x, y \in X$ it follows that $\lambda x + (1 - \lambda)y \in X$ for $\lambda \in [0,1]$



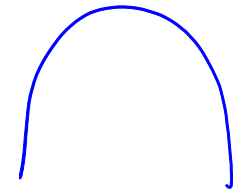
Convexity: Functions

- $f(\mathbf{x})$ is a convex function on a convex set X
iff
for all $\mathbf{x}, \mathbf{y} \in X$: $\lambda f(\mathbf{x}) + (1 - \lambda)f(\mathbf{y}) \geq f(\lambda\mathbf{x} + (1 - \lambda)\mathbf{y})$ for $\lambda \in [0,1]$

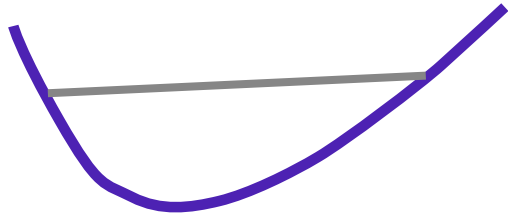


Convexity and *minimization* problems

CONCAVE



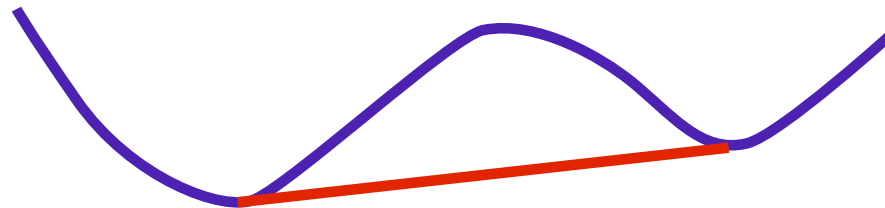
- Region **above** a convex function is convex



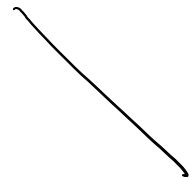
$$f(\lambda \mathbf{x} + (1 - \lambda) \mathbf{y}) \leq \lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y})$$

hence $\lambda f(\mathbf{x}) + (1 - \lambda) f(\mathbf{y}) \in X$ for $\mathbf{x}, \mathbf{y} \in X$

- Convex functions don't have local minima
 - Proof by contradiction - linear interpolation breaks local minimum condition



- Each local minimum is a global minimum
 - zero gradient implies (local) minimum for convex functions
 - if f_0 is a convex function and $\nabla f_0(\boldsymbol{\theta}^*) = 0$ then $\boldsymbol{\theta}^*$ is a global minimum
 - minimization becomes "relatively easy"



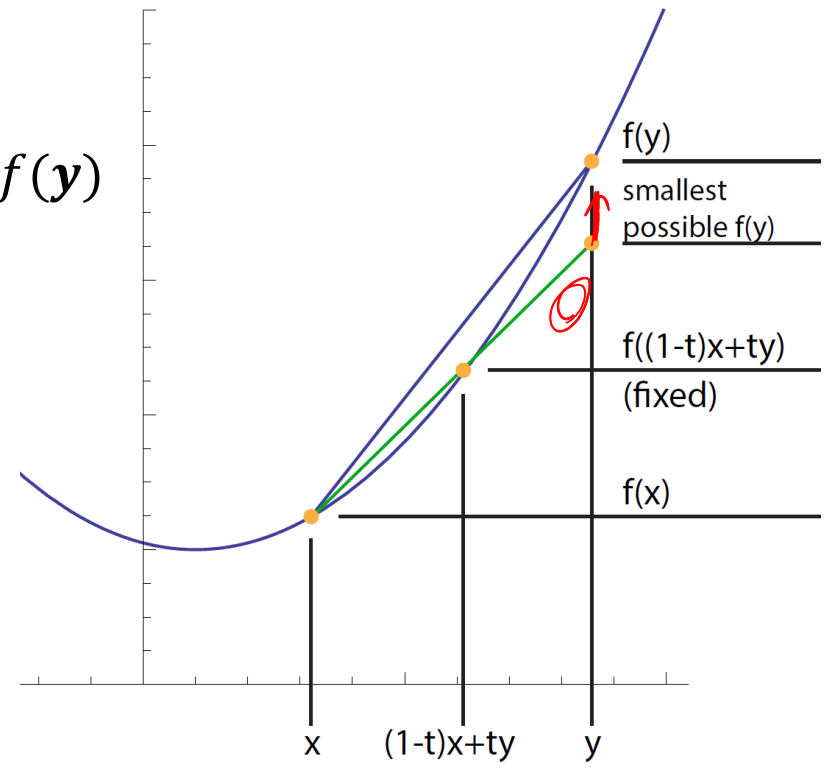
First order convexity conditions (I)

- Convexity imposes a rate of rise on the function

- $$f((1-t)x + ty) \leq (1-t)f(x) + tf(y)$$

- $$f(y) - f(x) \geq \frac{f((1-t)x + ty) - f(x)}{t}$$

- Difference between $f(y)$ and $f(x)$ is bounded by function values between x and y



First order convexity conditions (II)

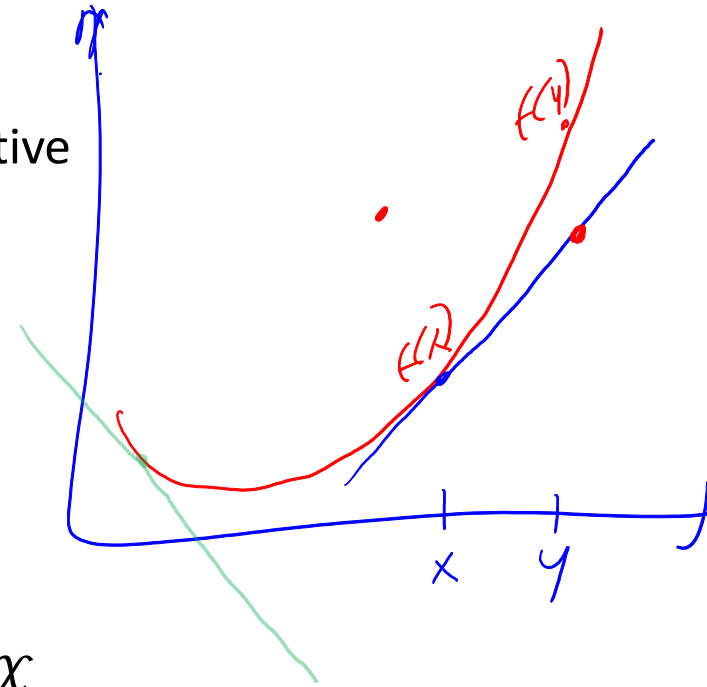
- $f(\mathbf{y}) - f(\mathbf{x}) \geq \frac{f((1-t)\mathbf{x}+t\mathbf{y}) - f(\mathbf{x})}{t}$
- Let $t \rightarrow 0$ and apply the definition of the derivative
- $f(\mathbf{y}) - f(\mathbf{x}) \geq (\mathbf{y} - \mathbf{x})^T \nabla f(\mathbf{x})$

- Theorem:

Suppose $f: \mathcal{X} \rightarrow \mathbb{R}$ is a differentiable function and \mathcal{X} is convex. Then f is convex iff for $\mathbf{x}, \mathbf{y} \in \mathcal{X}$

$$f(\mathbf{y}) \geq f(\mathbf{x}) + (\mathbf{y} - \mathbf{x})^T \nabla f(\mathbf{x})$$

- Proof. See Boyd p.70



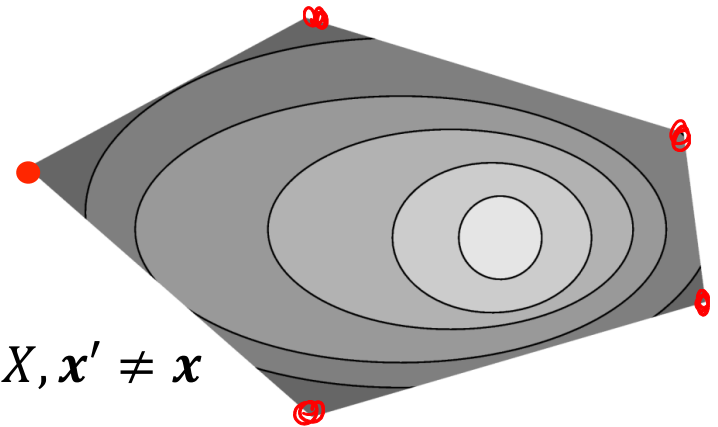
Convexity: Vertices & Convex Hull



- Given a (closed) convex set X .
 $\mathbf{x} \in X$ is a **vertex** of the convex set if it cannot be extrapolated within the convex set:

$$(\lambda \mathbf{x} + (1 - \lambda) \mathbf{x}') \notin X \text{ for } \lambda > 1 \text{ for all } \mathbf{x}' \in X, \mathbf{x}' \neq \mathbf{x}$$

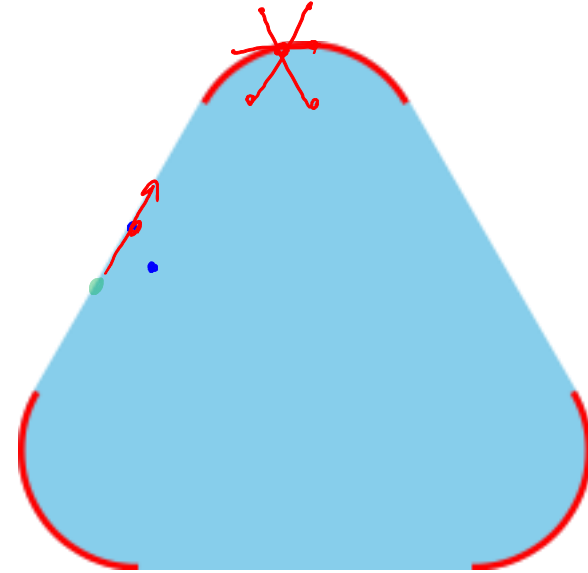
- We denote with $\text{Ve}(X)$ the set of all vertices of X



- Convex hull: Given a set of points $X \subseteq \mathbb{R}^d$, the convex hull is defined as
$$\text{Conv}(X) := \left\{ \sum_{i=1}^n \alpha_i \cdot \mathbf{x}_i \mid \mathbf{x}_i \in X, n \in \mathbb{N}, \sum \alpha_i = 1, \alpha_i \geq 0 \right\}$$
- Convex hull of a set is a convex set

Example: Vertices & Convex Hull

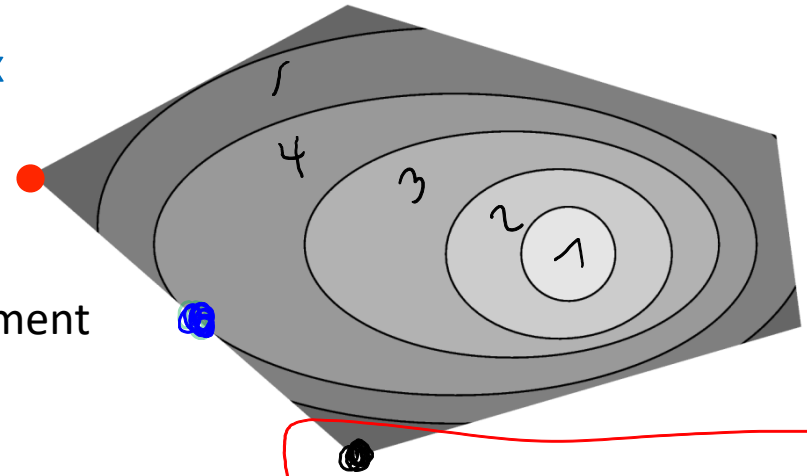
- Red = set X
- Red + Blue = convex hull of X
- Here: set X equals to the vertices of the convex hull
- In general:
 - $\text{Ve}(\text{Conv}(X)) \subseteq X$



Convexity and *maximization* problems (I)

$$\max_{\Theta \in \mathcal{F}} f(\Theta)$$

- Maximum over a convex function on a convex set is obtained on a vertex
- Proof:
 - Assume that maximum inside line segment
 - Then function cannot be convex
 - Hence it must be on vertex
- We only need to test the vertices to find the maximum
- In some cases this set is finite (see figure)



$$f(\bullet) \geq f(\bullet)$$
$$f(\bullet) \geq f(\bullet)$$

Convexity and *maximization* problems (II)

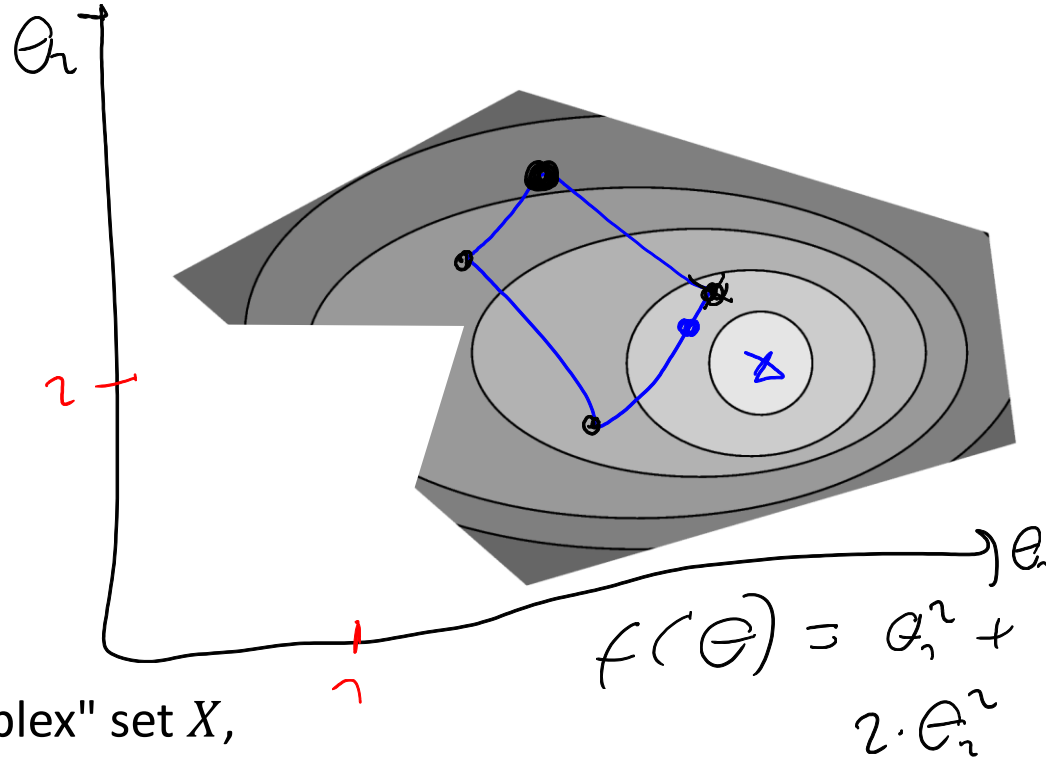
f is CONVEX !!!

- Supremum on convex hull

$$\sup_{x \in X} f(x) = \sup_{x \in \text{Conv}(X)} f(x)$$

- Proof: by contradiction

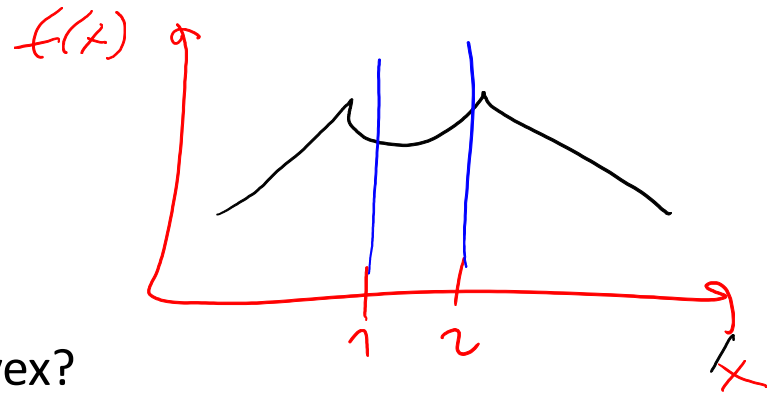
- Instead of working with a "complex" set X , operate with the easier, i.e. convex, set $\text{Conv}(X)$
- One might also simply focus on $\text{Ve}(\text{Conv}(X))$



$$\theta_1 + \theta_2 \leq 10$$

$$\theta_1 \geq 1 \text{ if } \theta_2 \leq 2$$

Verifying convexity (I)



- Convexity makes optimization "easier"
- How to verify whether a function is convex?
- For example: $e^{x_1+2x_2} + x_1 - \log(x_2)$ convex on $[1, \infty) \times [1, \infty)$?

1. Prove whether the definition of convexity holds (See slide 7)

2. Exploit special results

- First order convexity (See slide 10)
- Example: A twice differentiable function of one variable is convex on an interval if and only if its **second-derivative is non-negative** on this interval
- More general: a twice differentiable function of several variables is convex (on a convex set) if and only if its **Hessian matrix is positive semidefinite** (on the set)

$$x^T A x \geq 0$$

Verifying convexity (II)

3. Show that the function can be obtained from simple convex functions by operations that preserve convexity

a) Start with simple convex functions, e.g.

- $f(x) = \text{const}$ and $f(\mathbf{x}) = \mathbf{x}^T \cdot \mathbf{b}$ (these are also concave functions)
- $f(x) = e^x$

b) Apply "construction rules" (next slide)

Convexity preserving operations

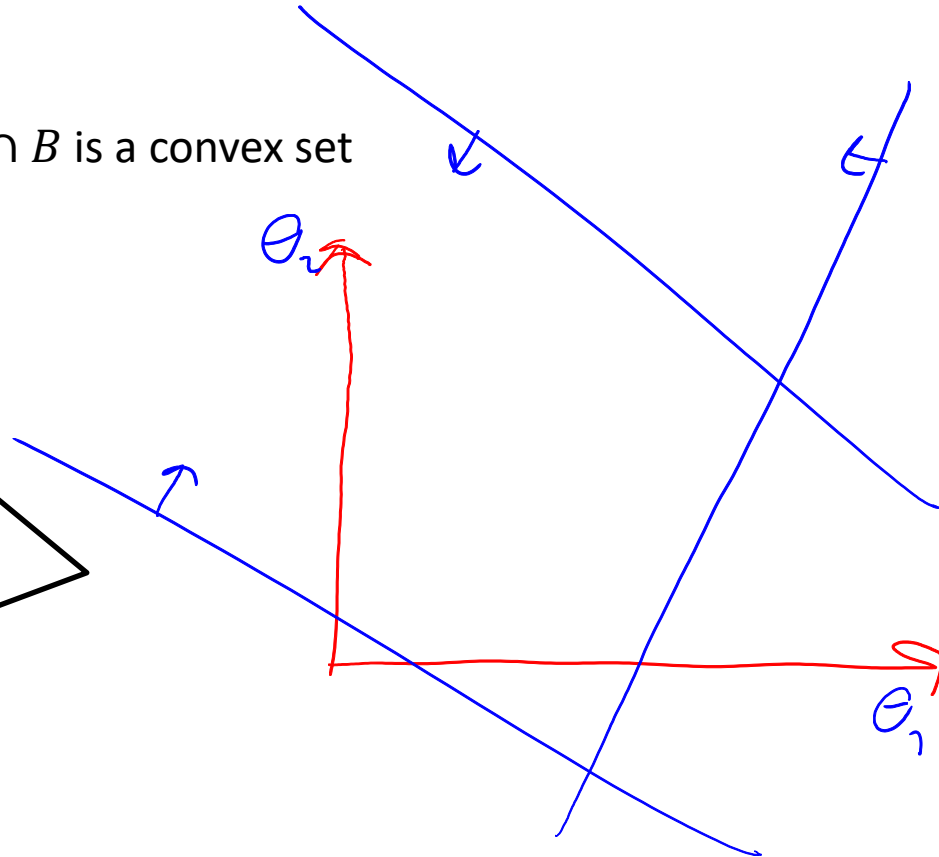
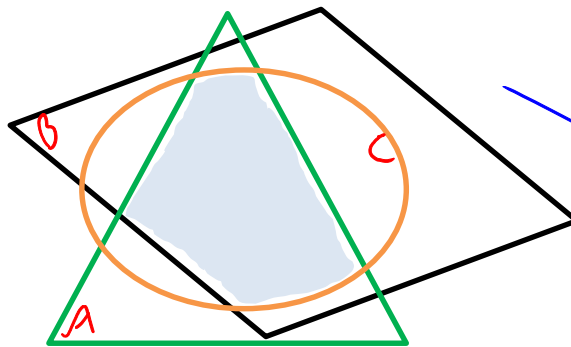
\<1
>1

- Let $f_1: \mathbb{R}^d \rightarrow \mathbb{R}$ and $f_2: \mathbb{R}^d \rightarrow \mathbb{R}$ be **convex** functions, and $g: \mathbb{R}^d \rightarrow \mathbb{R}$ be a **concave** function, then
 - $h(\mathbf{x}) = f_1(\mathbf{x}) + f_2(\mathbf{x})$ is convex
 - $h(\mathbf{x}) = \max\{f_1(\mathbf{x}), f_2(\mathbf{x})\}$ is convex
 - $h(\mathbf{x}) = c \cdot f_1(\mathbf{x})$ is convex if $c \geq 0$
 - $h(\mathbf{x}) = c \cdot g(\mathbf{x})$ is convex if $c \leq 0$
 - $h(\mathbf{x}) = f_1(A\mathbf{x} + \mathbf{b})$ is convex (A matrix, b vector)
 - $h(\mathbf{x}) = m(f_1(\mathbf{x}))$ is convex if $m: \mathbb{R} \rightarrow \mathbb{R}$ is convex and nondecreasing
- Example: $e^{x_1+2x_2} + x_1 - \log(x_2)$ is convex on, e.g., $[1, \infty) \times [1, \infty)$



Verifying convexity for sets

1. Prove definition
 - often easier for sets than for functions
2. Apply intersection rule
 - Let A and B be convex sets, then $A \cap B$ is a convex set

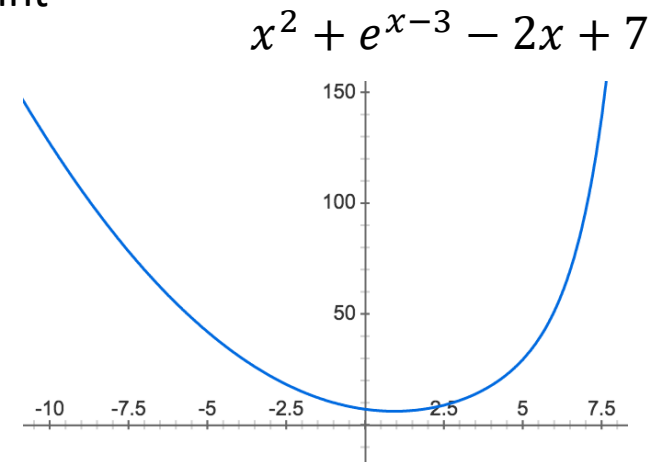


An easy problem

$$\theta = \dots$$

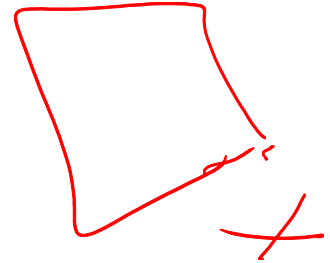
Convex objective function f

- Objective function differentiable on its whole domain
 - i.e. we are able to compute gradient f' at every point
- We can solve $f'(\theta) = 0$ for θ analytically
 - i.e. solution for θ where gradient = 0 is known
- Unconstrained minimization
 - i.e. above computed solution for θ is valid
- We are done!
- Example: Ordinary Least Squares Regression

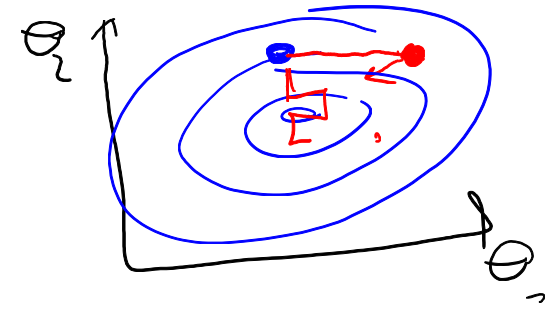


Outlook

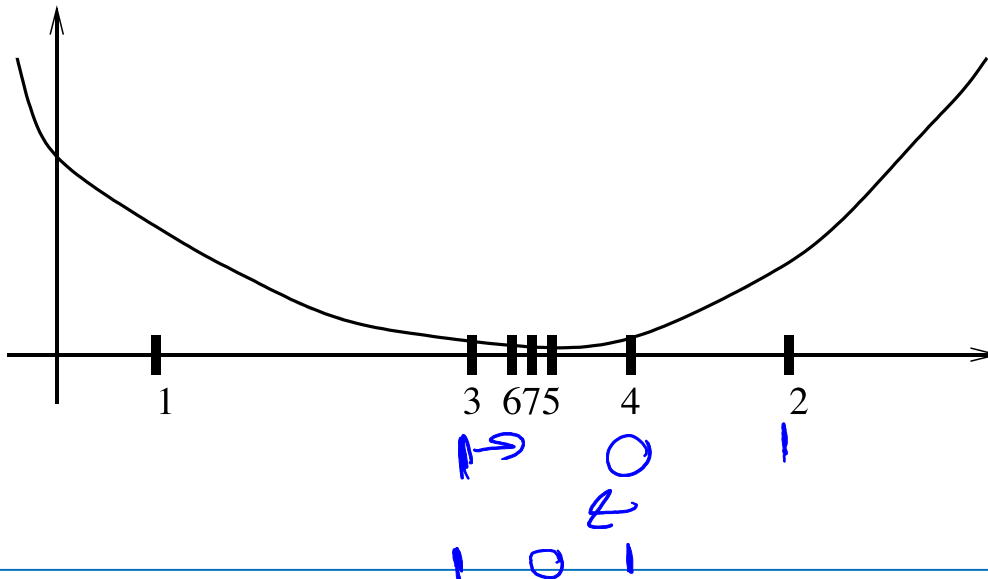
- Unfortunately, many problems are harder...
- No analytical solution for $f'(\boldsymbol{\theta}) = 0$
 - e.g. Logistic Regression
 - Solution: try numerical approaches, e.g. gradient descent
- Constraints on $\boldsymbol{\theta}$
 - e.g. $f'(\boldsymbol{\theta}) = 0$ only holds for points outside the domain
 - Solution: constrained optimization
- f not differentiable on whole domain
 - Potential solution: subgradients; or is it a discrete optimization problem?
- f not convex
 - Potential solution: convex relaxations; convex in some variables?



One-dimensional problems



- Key idea
 - For differentiable f search for θ with $\nabla f(\theta) = 0$
 - Interval bisection (derivative is monotonic)
- Can be extended to nondifferentiable problems
 - exploit convexity in upper bound and keep 5 points



Require: a, b , Precision ϵ

Set $A = a, B = b$

repeat

if $f'(\frac{A+B}{2}) > 0$ **then**

$B = \frac{A+B}{2}$

else

$A = \frac{A+B}{2}$

end if

until $(B - A) \min(|f'(A)|, |f'(B)|) \leq \epsilon$

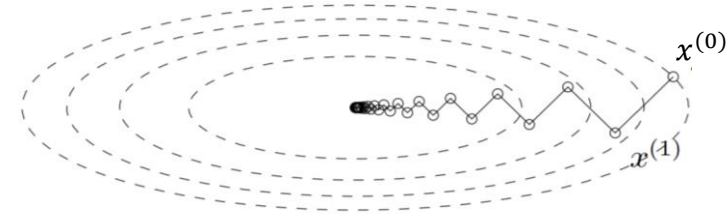
Output: $x = \frac{A+B}{2}$

*solution on
the left*

Gradient Descent

- Key idea

- Gradient points into steepest ascent direction
- Locally, the gradient is a good approximation of the objective function



- GD with Line Search

- Get descent direction, then unconstrained line search



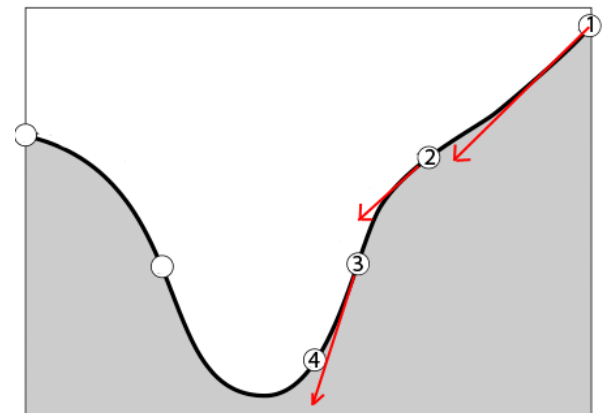
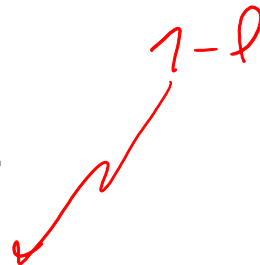
- Turn a multidimensional problem into a one-dimensional problem that we already know how to solve

given a starting point $\theta \in \text{dom}(f)$.

repeat

1. $\Delta\theta := -\nabla f(\theta)$
2. Line search. $t = \arg \min_{t>0} f(\theta + t \cdot \Delta\theta)$
3. Update. $\theta := \theta + t\Delta\theta$

until stopping criterion is satisfied.



Gradient Descent convergence

- Let p^* be the optimal value, $\boldsymbol{\theta}^*$ be the minimizer – the point where the minimum is obtained, and $\boldsymbol{\theta}^{(0)}$ be the starting point
- The residual error ρ , for the k-th iteration is (for strongly convex f):
$$\rho = f(\boldsymbol{\theta}^{(k)}) - p^* \leq c^k (f(\boldsymbol{\theta}^{(0)}) - p^*), \quad c < 1$$

 $f(\boldsymbol{\theta}^{(k)})$ converges to p^* as $k \rightarrow \infty$
- We must have $f(\boldsymbol{\theta}^{(k)}) - p^* \leq \epsilon$ after at most $\frac{\log((f(\boldsymbol{\theta}^{(0)}) - p^*)/\epsilon)}{\log(1/c)}$ iterations
- Linear convergence for strongly convex objective
 - $k \sim \log(\rho^{-1})$ // k = number of iterations, ρ
- Attention: linear convergence = exponentially fast
 - i.e. linear when plotting on a log scale – old statistics terminology

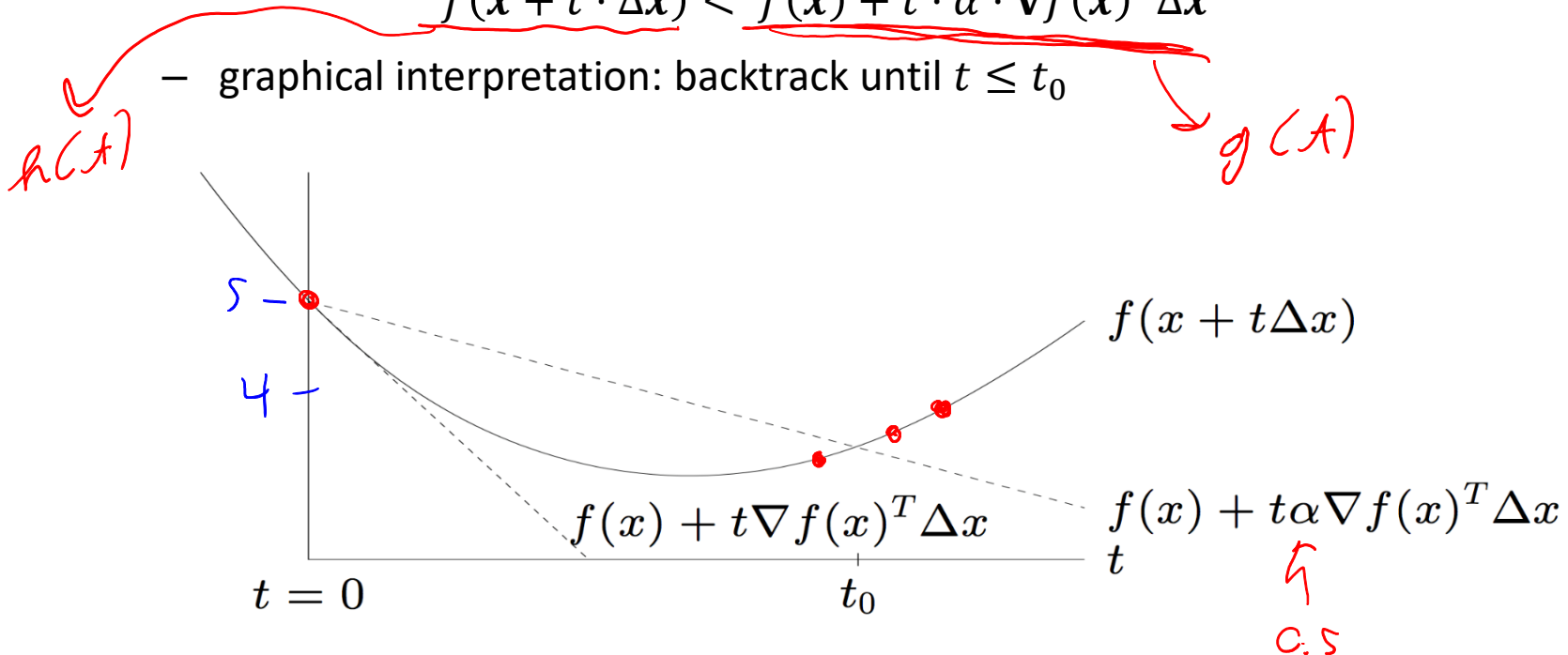
Line search types

$$f(\theta_{b+n}) < f(\theta_{\#})$$

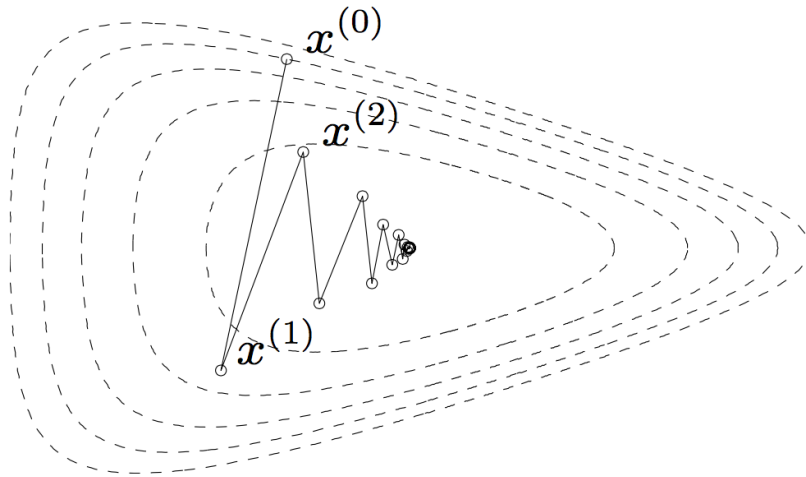
- Exact line search: $t = \arg \min_{t>0} f(\mathbf{x} + t \cdot \Delta \mathbf{x})$
- Backtracking line search: (with parameters $\alpha \in (0, 1/2), \beta \in (0, 1)$)
 - starting at $t = 1$, repeat $t := \beta t$ until

$$f(\mathbf{x} + t \cdot \Delta \mathbf{x}) < f(\mathbf{x}) + t \cdot \alpha \cdot \nabla f(\mathbf{x})^T \Delta \mathbf{x}$$

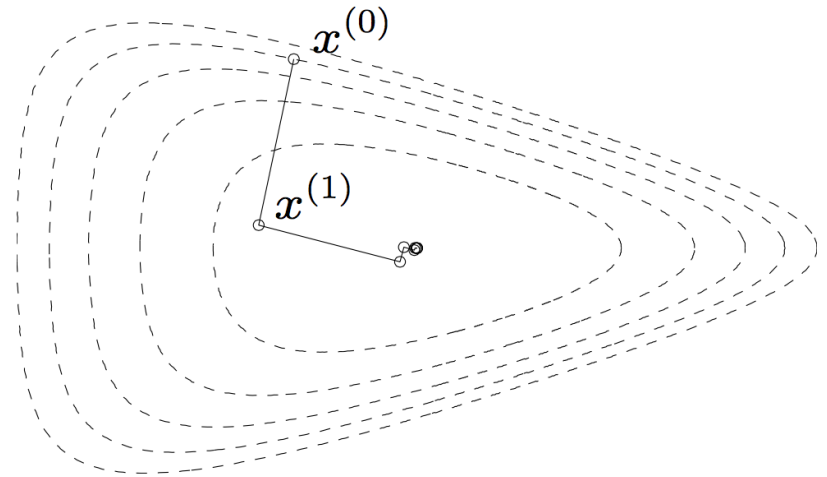
- graphical interpretation: backtrack until $t \leq t_0$



Backtracking vs. exact line search



backtracking line search



exact line search

from Boyd & Vandenberghe

Distributed/Parallel implementation

- Often problems are of the form
 - $f(\boldsymbol{\theta}) = \sum_i L_i(\boldsymbol{\theta}) + g(\boldsymbol{\theta})$
 - where i iterates over, e.g., each data instance
- Example OLS regression: // with regularization
 - $L_i(\mathbf{w}) = (\mathbf{x}_i^T \mathbf{w} - y_i)^2$ $g(\mathbf{w}) = \lambda \cdot \|\mathbf{w}\|_2^2$
- Gradient can simply be decomposed based on the sum rule
- Easy to parallelize/distribute

Basic steps (I)

given a starting point $\theta \in \text{Dom}(f)$.

repeat

1. $\Delta\theta := -\nabla f(\theta)$
2. Line search. $t = \operatorname{argmin}_{t>0} f(\theta + t \cdot \Delta\theta)$
3. Update. $\theta := \theta + t\Delta\theta$

until stopping criterion is satisfied.

easy parallel
computation

- Distribute data over several machines
- Compute partial gradients (on each machine in parallel)
- Aggregate the partial gradients to the final one
- Communicate the final gradient back to all machines

Basic Steps (II)

given a starting point $\theta \in \text{Dom}(f)$.

repeat

1. $\Delta\theta := -\nabla f(\theta)$
2. Choose t via exact or backtracking line search.
3. Update. $\theta := \theta + t\Delta\theta$

until stopping criterion is satisfied.

update value in search direction and feed back
(might be done multiple times: expensive!)

communicate final step
size to each machine

- Line search is expensive
 - for each tested step size: scan through all datapoints

Scalability analysis

- + Linear time in number of instances
- + Linear memory consumption in problem size (not data)
- + Logarithmic time in accuracy
- + 'Perfect' scalability
- Multiple passes through dataset for each iteration

A faster algorithm

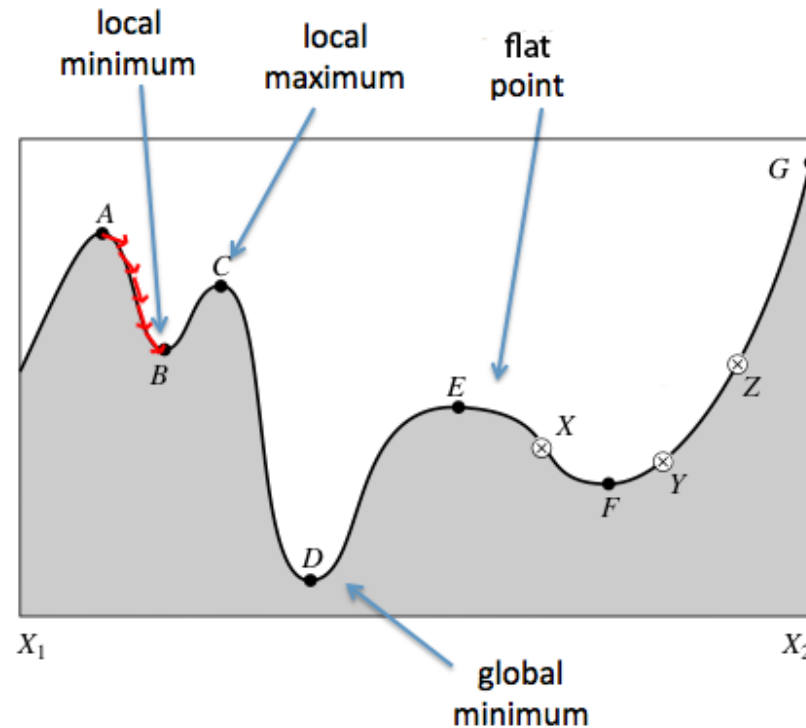
- Avoid the line search; simply pick update

$$\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t - \tau \cdot \nabla f(\boldsymbol{\theta}_t)$$

- τ is often called the **learning rate**
- Only single pass through data per iteration
- Logarithmic iteration bound (as before)
 - if learning rate is chosen "correctly"
- How to pick the learning rate?
 - too small: slow convergence
 - too high: algorithm might oscillate, no convergence
- Interactive tutorial on optimization
 - <http://www.benfrederickson.com/numerical-optimization/>

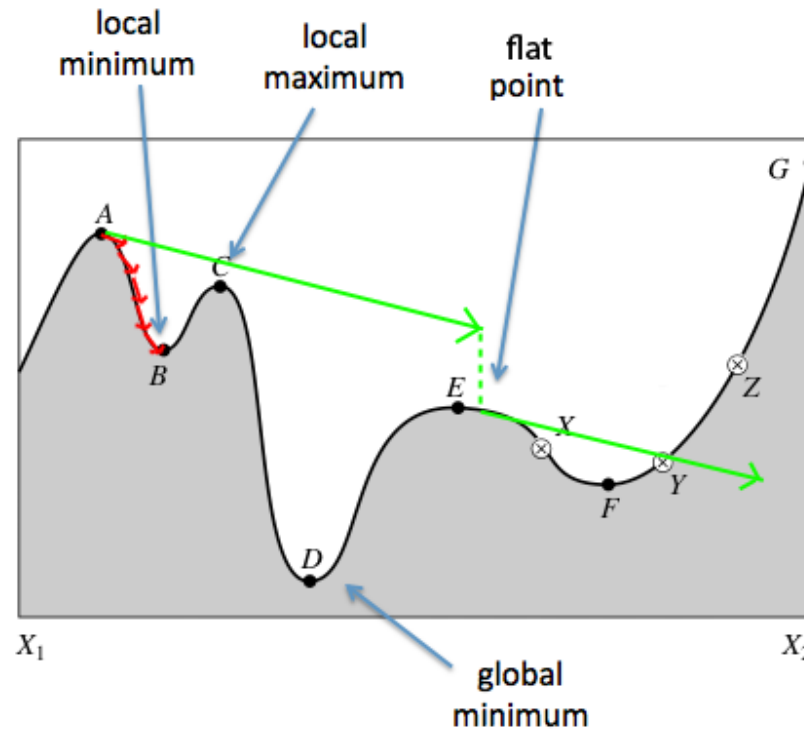
The value of τ

- A too **small** value for τ has two drawbacks
 - We find the minimum more slowly
 - We end up in local minima or saddle/flat points



The value of τ

- A too **large** value for τ has one drawback
 - You may never find a minimum; oscillations usually occur
- We only need 1 steps to overshoot



Learning rate adaptation

- Simple solution: let the learning rate be a decreasing function τ_t of the iteration number t
 - so called **learning rate schedule**
 - first iterations cause large changes in the parameters; later do fine-tuning
 - convergence easily guaranteed if $\lim_{t \rightarrow \infty} \tau_t = 0$
 - example: $\tau_{t+1} \leftarrow \alpha \cdot \tau_t$ for $0 < \alpha < 1$

Learning rate adaptation

- Other solutions: Incorporate "history" of previous gradients
- **Momentum:**
 - $\mathbf{m}_t \leftarrow \tau \cdot \nabla f(\boldsymbol{\theta}_t) + \gamma \cdot \mathbf{m}_{t-1}$ // often $\gamma = 0.5$
 - $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t - \mathbf{m}_t$
 - As long as gradients point to the same direction, the search accelerates
- **AdaGrad:**
 - different learning rate per parameter
 - learning rate depends inversely on accumulated "strength" of all previously computed gradients
 - large parameter updates ("large" gradients) lead to small learning rates

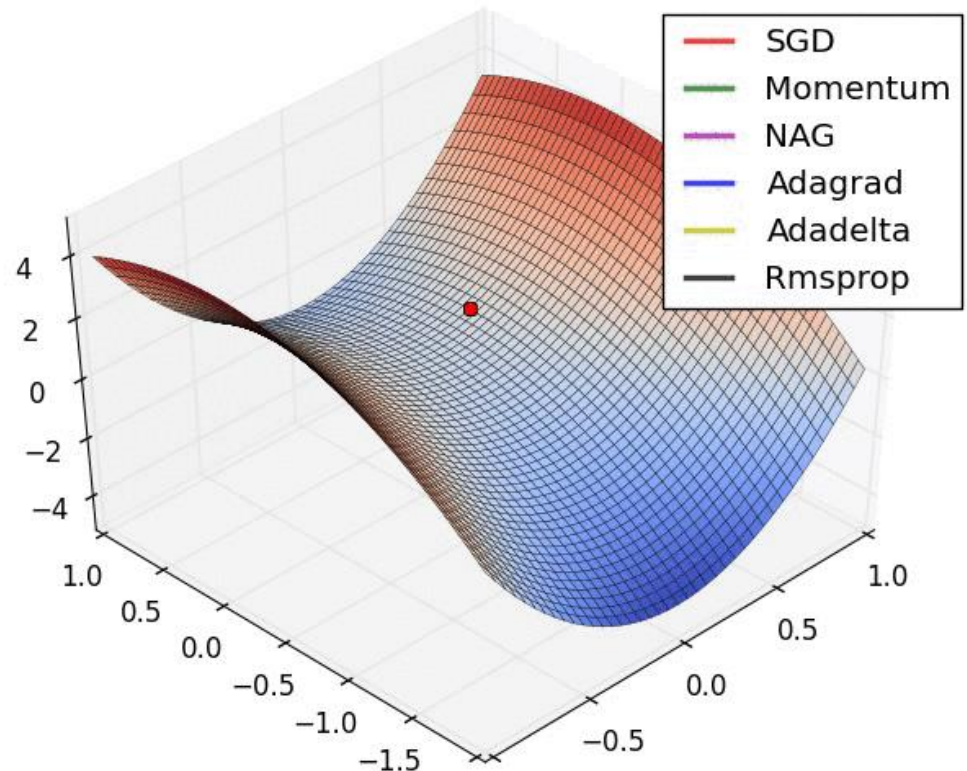
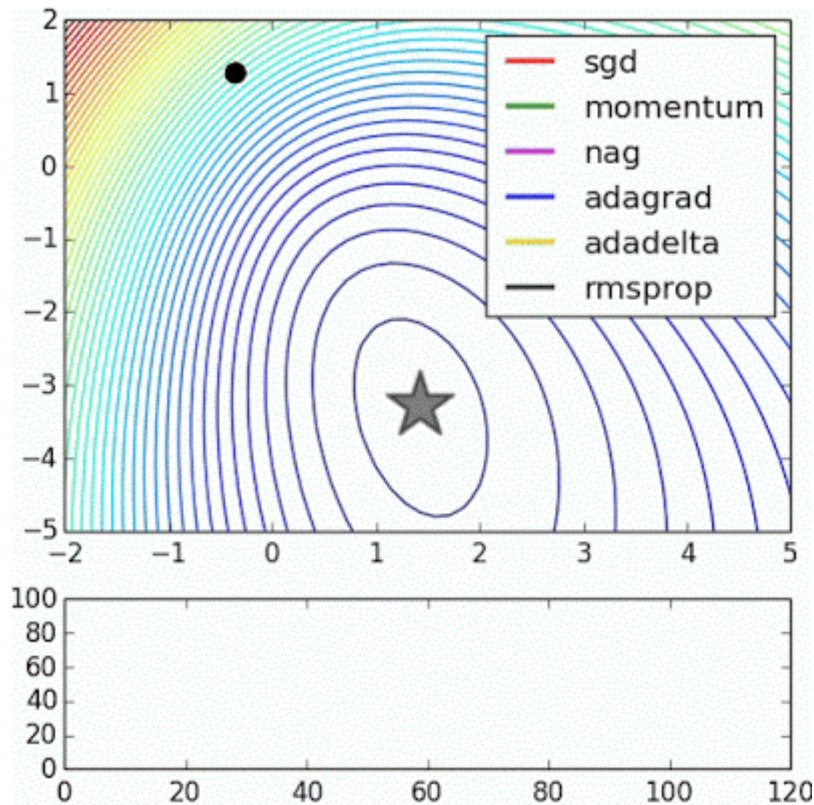
Adaptive moment estimation (Adam)

- $\mathbf{m}_t = \beta_1 \mathbf{m}_{t-1} + (1 - \beta_1) \nabla f(\boldsymbol{\theta}_t)$
 - estimate of the first moment (mean) of the gradient
 - Exponentially decaying average of past gradients m_t (similar to momentum)
- $\mathbf{v}_t = \beta_2 \mathbf{v}_{t-1} + (1 - \beta_2) (\nabla f(\boldsymbol{\theta}_t))^2$
 - estimate of the second moment (uncentered variance) of the gradient
 - Exponentially decaying average of past squared gradients v_t
- To avoid bias towards zero (due to 0's initialization) use bias-corrected version instead:
 - $\widehat{\mathbf{m}}_t = \frac{\mathbf{m}_t}{1 - \beta_1^t} \quad \widehat{\mathbf{v}}_t = \frac{\mathbf{v}_t}{1 - \beta_2^t}$
- Finally, the Adam update rule for parameters θ :
 - $\boldsymbol{\theta}_{t+1} = \boldsymbol{\theta}_t - \frac{\tau}{\sqrt{\widehat{\mathbf{v}}_t} + \epsilon} \widehat{\mathbf{m}}_t$
- Default values: $\beta_1 = 0.9, \beta_2 = 0.999, \epsilon = 10^{-8}$

Visualizing gradient descent variants

- AdaGrad and variants
 - often have faster convergence
 - might help to escape saddlepoints

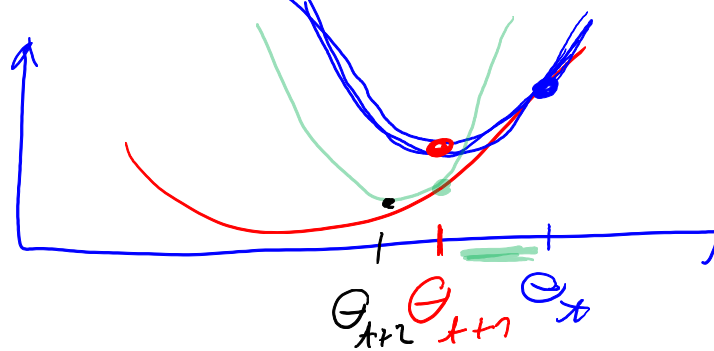
<http://sebastianruder.com/optimizing-gradient-descent/>



Discussion

- Gradient descent and similar techniques are called first-order optimization techniques
 - they only exploit information of the gradients (i.e. first order derivative)
- Higher-order techniques use higher-order derivatives
 - e.g. second-order = Hessian matrix
 - Example: Newton Method

Newton method



- Convex objective function f
- Nonnegative second derivative: $\nabla^2 f(\theta) \succcurlyeq 0$ // Hessian matrix

- Taylor expansion of f at point θ_t

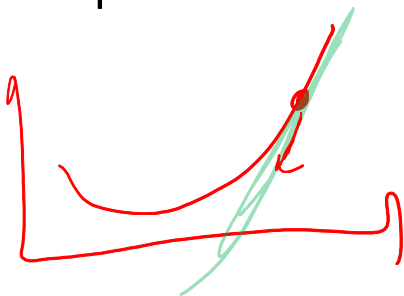
$$f(\theta_t + \delta) = f(\theta_t) + \delta^T \cdot \nabla f(\theta_t) + \frac{1}{2} \delta^T \nabla^2 f(\theta_t) \delta + O(\delta^3)$$

$h(\delta) \approx h'(\delta) \leftarrow \text{Approx}$

- Minimize approximation: leads to

$$\theta_{t+1} \leftarrow \theta_t - [\nabla^2 f(\theta_t)]^{-1} \nabla f(\theta_t) = \theta_t + \delta$$

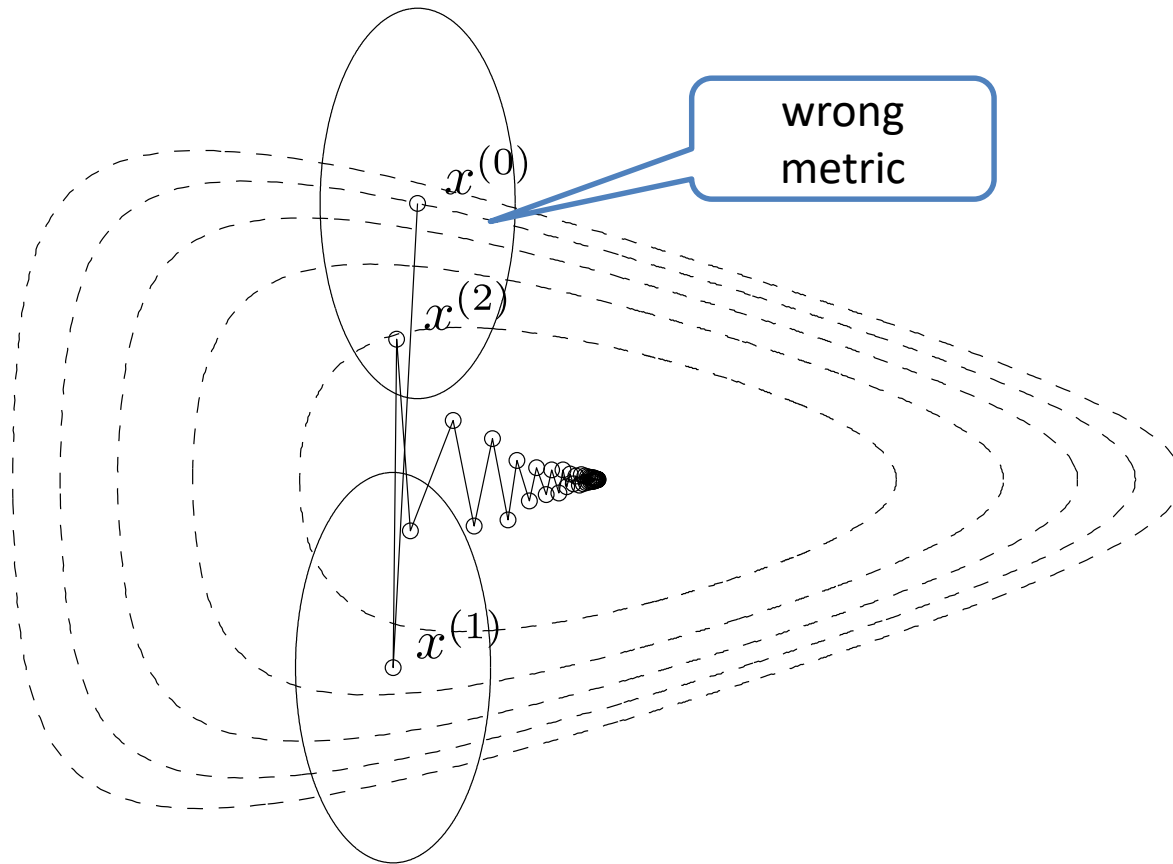
- Repeat until convergence



$$\nabla h'(\delta) \stackrel{!}{=} 0 \Leftrightarrow \delta =$$

Rescaling of space

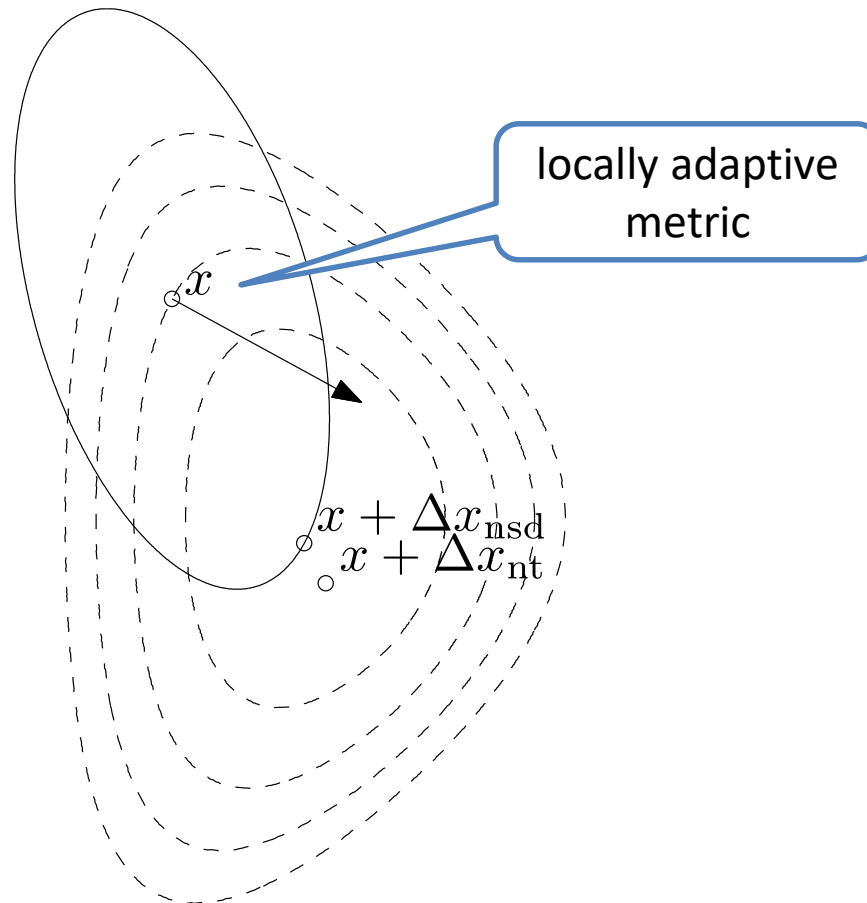
- Newton method rescales the space



from Boyd & Vandenberghe

Rescaling of space

- Newton method rescales the space



Parallel Newton method

- + Good rate of convergence
- + Few passes through data needed
- + Parallel aggregation of gradient and Hessian
- + Gradient requires $O(d)$ data
- Hessian requires $O(d^2)$ data
- Update step is $O(d^3)$ & nontrivial to parallelize
- Use it only for low dimensional problems!

Large scale optimization

$$f(\theta) = \sum_{i=1}^n \mathcal{L}_i(\theta)$$

- Higher-order techniques have nice properties (e.g. convergence) but they are prohibitively expensive for high dimensional problems
- For large scale data / high dimensional problems use first-order techniques
 - i.e. variants of gradient descent
- But for real-world large scale data even first-order methods are too costly
- Solution: **Stochastic optimization!**

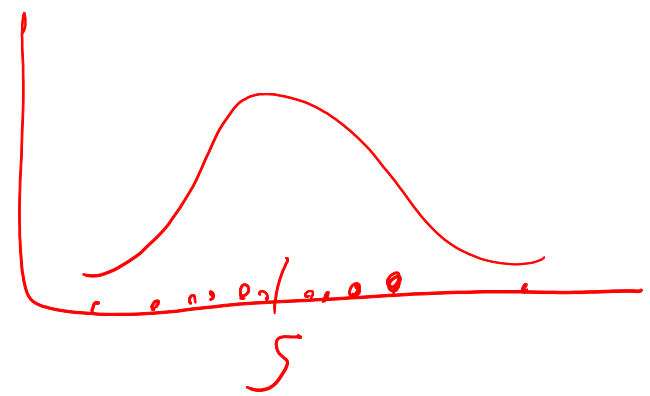
Motivation: Stochastic Gradient Descent

- Goal: minimize $f(\boldsymbol{\theta}) = \sum_{i=1}^n L_i(\boldsymbol{\theta})$ + potential constraints
- For very large data: even a single pass through the data is very costly
- Lots of time required to even compute the very first gradient
- Is it possible to update the parameters more frequently/faster?

Stochastic Gradient Descent

- Consider the task as empirical risk minimization

$$\frac{1}{n} (\sum_{i=1}^n L_i(\boldsymbol{\theta})) = \mathbb{E}_{i \sim \{1, \dots, n\}} [L_i(\boldsymbol{\theta})]$$



- (Exact) expectation can be approximated by smaller sample:
- $\mathbb{E}_{i \sim \{1, \dots, n\}} [L_i(\boldsymbol{\theta})] \approx \frac{1}{|S|} \sum_{j \in S} (L_j(\boldsymbol{\theta}))$ // with $S \subseteq \{1, \dots, n\}$

or equivalently: $\sum_{i=1}^n L_i(\boldsymbol{\theta}) \approx \frac{n}{|S|} \sum_{j \in S} L_j(\boldsymbol{\theta})$

Stochastic Gradient Descent

- Intuition: Instead of using "exact" gradient, compute only a **noisy (but still unbiased) estimate** based on smaller sample
- Stochastic gradient decent:
 1. randomly pick a (small) subset S of the points \rightarrow so called mini-batch
 2. compute gradient based on mini-batch
 3. update: $\boldsymbol{\theta}_{t+1} \leftarrow \boldsymbol{\theta}_t - \tau \cdot \frac{n}{|S|} \cdot \sum_{j \in S} \nabla L_j(\boldsymbol{\theta}_t)$
 4. pick a new subset and repeat with 2
- "Original" SGD uses mini-batches of size 1
 - larger mini-batches lead to more stable gradients (i.e. smaller variance in the estimated gradient)

Example: Perceptron

$$w_{x+1} \leftarrow w_x - 0$$

- Simple linear binary classifier:

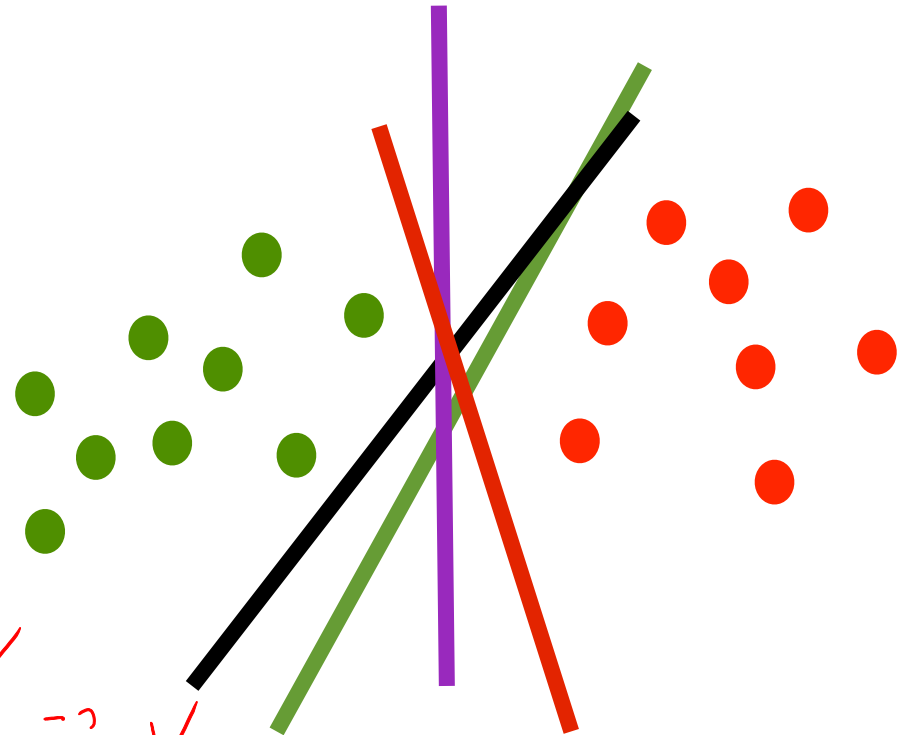
$$\delta(\mathbf{x}) = \begin{cases} 1 & \text{if } \mathbf{w}^T \mathbf{x} + b > 0 \\ -1 & \text{else} \end{cases}$$

- Learning task:

Given $(\mathbf{x}_1, y_1), \dots, (\mathbf{x}_n, y_n)$ find

$$\min_{\mathbf{w}, b} \sum_i L(y_i, \mathbf{w}^T \mathbf{x}_i + b)$$

+1 ✓
+1 ✓
-1 ✓
-1 ✓



- L is the loss function

$$- \text{e.g. } L(u, v) = \max(0, -u \cdot v) = \begin{cases} -uv & \text{if } uv < 0 \\ 0 & \text{else} \end{cases} \quad \begin{array}{l} \leftarrow \text{incorrect classification} \\ \leftarrow \text{correct classification} \end{array}$$

$$-(y_i \cdot (\mathbf{w}^T \mathbf{x}_i + b)) =$$

Example: Perceptron

- Let's solve this problem via SGD
- Result:

```
initialize  $w = 0$  and  $b = 0$   
repeat  
  if  $y_i \cdot (w^T x_i + b) \leq 0$  then  
     $w \leftarrow w + \tau \cdot n \cdot y_i \cdot x_i$  and  $b \leftarrow b + \tau \cdot n \cdot y_i$   
  end if  
until all classified correctly
```

- Note: Nothing happens if classified correctly
 - gradient is zero
- Does this remind you of the original learning rules for perceptron?

Optimizing Logistic Regression

- Recall we wanted to solve $\mathbf{w}^* = \arg \min_{\mathbf{w}} E(\mathbf{w})$
- $E(\mathbf{w}) = -\ln p(\mathbf{y} \mid \mathbf{w}, \mathbf{X})$
$$= -\sum_{i=1}^N y_i \ln \sigma(\mathbf{w}^T \mathbf{x}_i) + (1 - y_i) \ln(1 - \sigma(\mathbf{w}^T \mathbf{x}_i))$$
- Closed form solution does not exist
- Solution:
 - Computed the gradient $\nabla E(\mathbf{w})$
 - Find \mathbf{w}^* using gradient descent
- Is $E(\mathbf{w})$ convex?
- Can you use SGD?
- How can you choose the learning rate?
- What changes if we add regularization, i.e. $E_{reg}(\mathbf{w}) = E(\mathbf{w}) + \lambda \|\mathbf{w}\|_2^2$?

Convergence in expectation

- Subject to relatively mild assumptions, stochastic gradient descent converges **almost surely** to a global minimum when the objective function is convex
 - almost surely to a local minimum for non-convex functions

- The expectation of the residual error decreases with speed

$$\mathbb{E}[\rho] \sim t^{-1} \quad // \text{ i.e. } t \sim \mathbb{E}[\rho]^{-1}$$

- Note: Standard GD has speed $t \sim \log \rho^{-1}$
 - faster convergence speed; but each iteration takes longer

Summary

- General task: Find solution θ^* minimizing function f
- Convex sets & functions
 - Global vs. local minimum
 - Convex hull $\text{Conv}(X)$
 - Verifying convexity: Definition, special results (first-order convexity, 2nd derivative), convexity-preserving operations
- Gradient descent: $\theta := \theta - t \nabla f(\theta)$
 - Line search: How to choose t ? E.g. backtracking, exact
 - Learning rate: Fix $t = \tau$, change via learning rate adaptation (momentum, AdaGrad)
 - Faster methods: Adam, Newton method
 - Stochastic gradient descent (SGD): Only use part of data (mini-batches) at each step