

COS 324: Lecture 9

Stochastic Gradient Descent

Elad Hazan Yoram Singer



Admin

- HW4

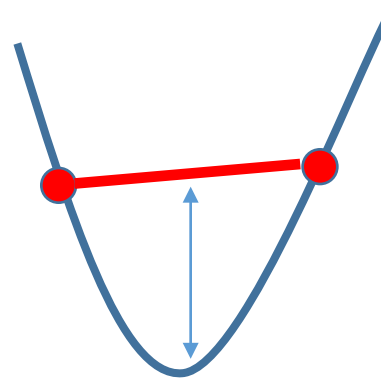
Recap

- Started convex analysis
- Gradient descent + analysis
- Today: faster optimization through randomization: stochastic gradient descent
 - Part 1: how to deal with constraints?
 - Part 2: SGD

Convex Functions and Sets

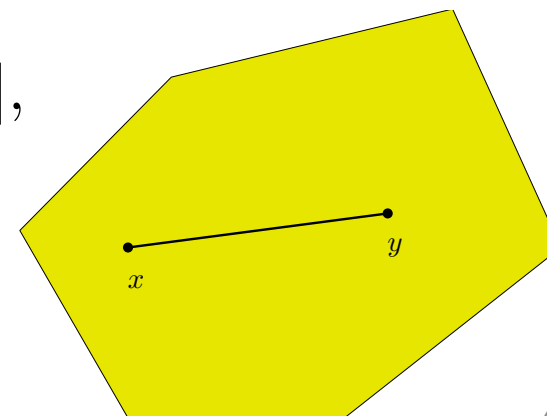
A function $f : \mathbb{R}^n \rightarrow \mathbb{R}$ is convex if for $x, y \in \text{dom } f$ and any $a \in [0, 1]$,

$$f(ax + (1 - a)y) \leq af(x) + (1 - a)f(y)$$



A set $C \subseteq \mathbb{R}^n$ is convex if for $x, y \in C$ and any $a \in [0, 1]$,

$$ax + (1 - a)y \in C$$

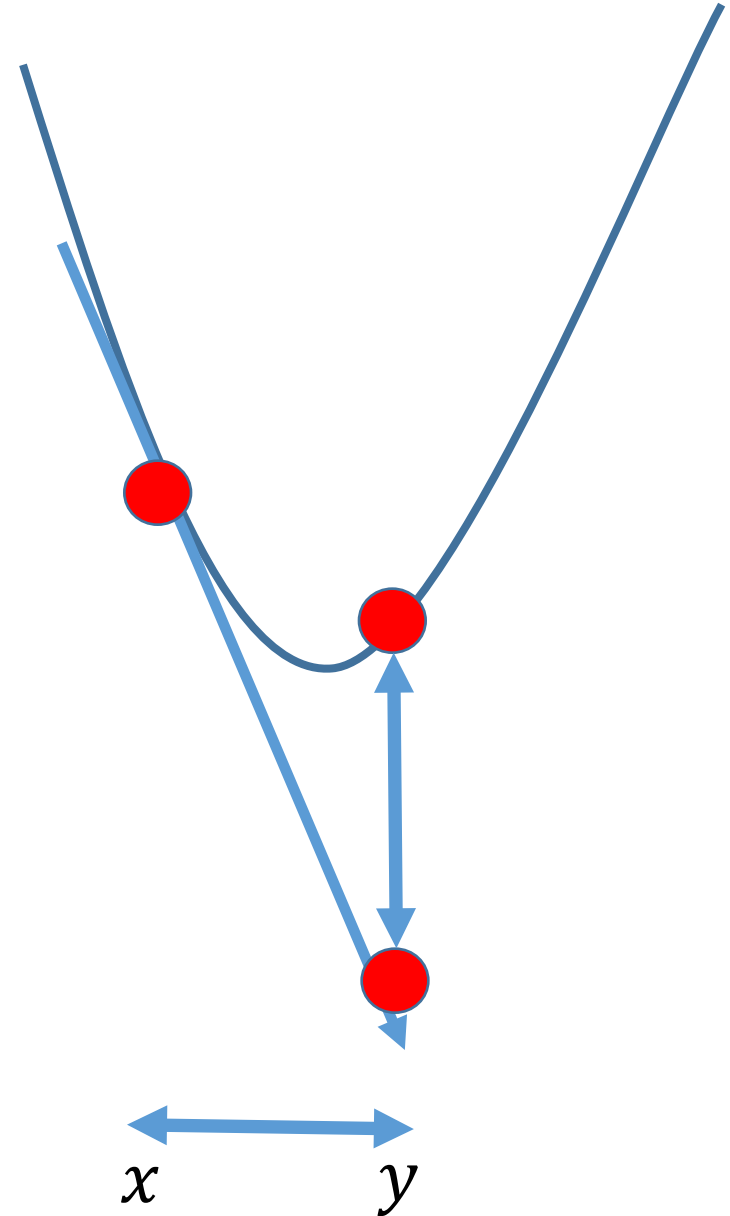


Convexity: local \rightarrow global

Convexity

- Alternative definition:

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x)$$



Lipschitzness

f is G -Lipschitz if for every $x, y \in K$, we have

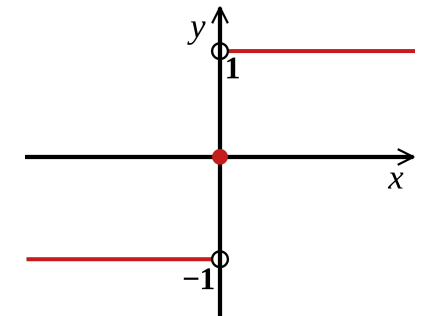
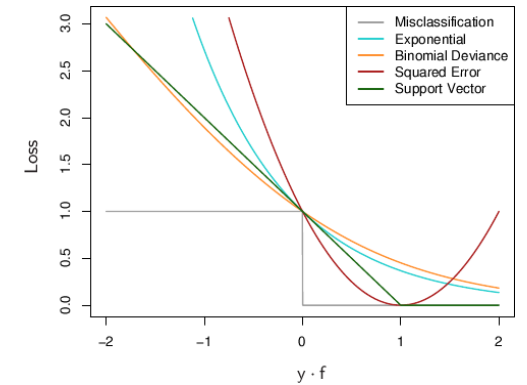
$$|f(x) - f(y)| \leq G |x - y|$$

Note: for convex functions, suffices that the gradient is bounded (why?)

$$\forall x \in K \quad |\nabla f(x)| \leq G$$

Recall for convex functions:

$$f(x) - f(y) \leq \nabla f(x)^\top (x - y) \leq |\nabla f(x)| |x - y| \leq G |x - y|$$



Optimality conditions

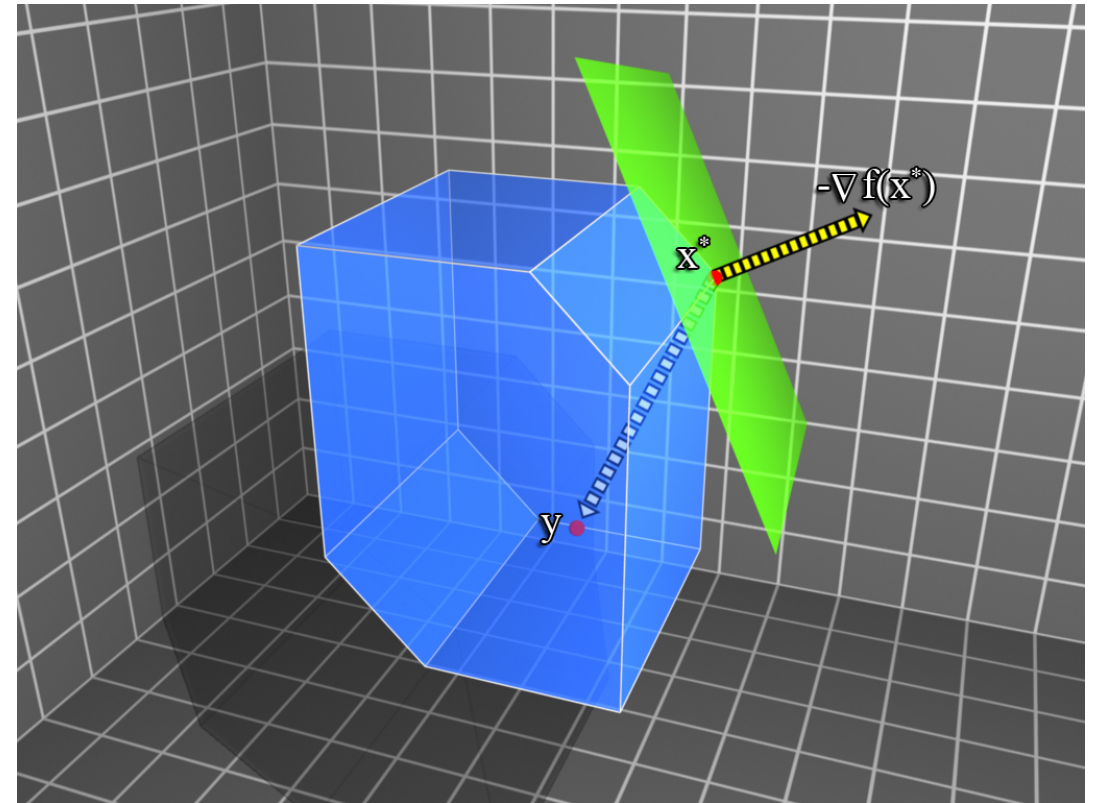
x^* is the minimum of convex function f iff
 $|\nabla f(x^*)| = 0$

If we have a constrained set K , then it is optimum iff

$$\prod_K [x^* - \nabla f(x^*)] = x^*$$

Here \prod_K denotes the projection operation, defined as:

$$\prod_K [y] = \arg \min_{x \in K} |x - y|$$



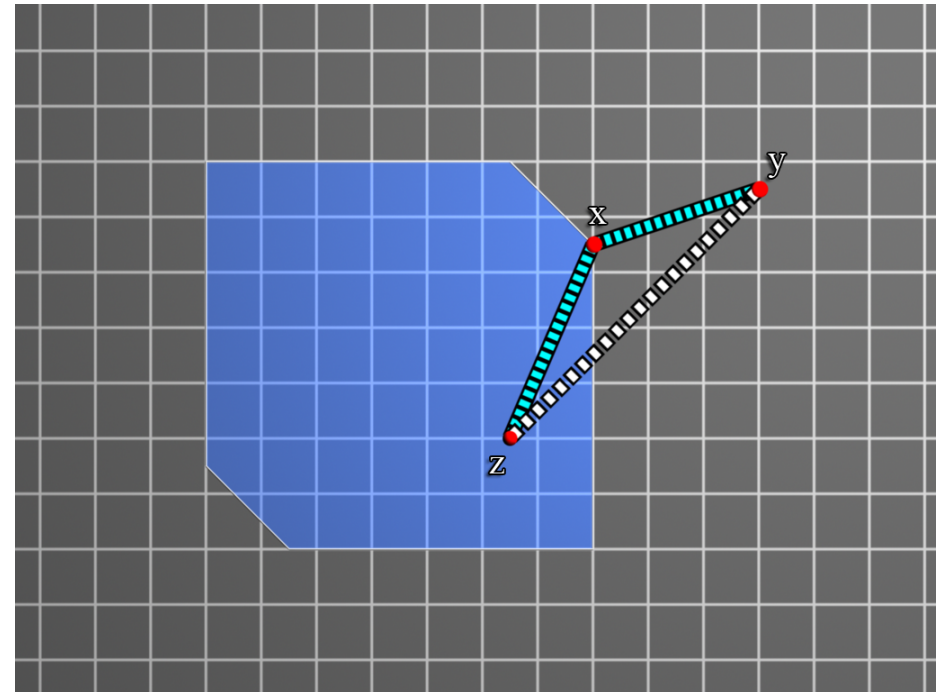
Projections

For projections over convex sets, defined as

$$\prod_K [y] = \arg \min_{x \in K} |x - y|$$

We have the Pythagorean theorem:

$$|y - x|^2 \leq |y - z|^2$$



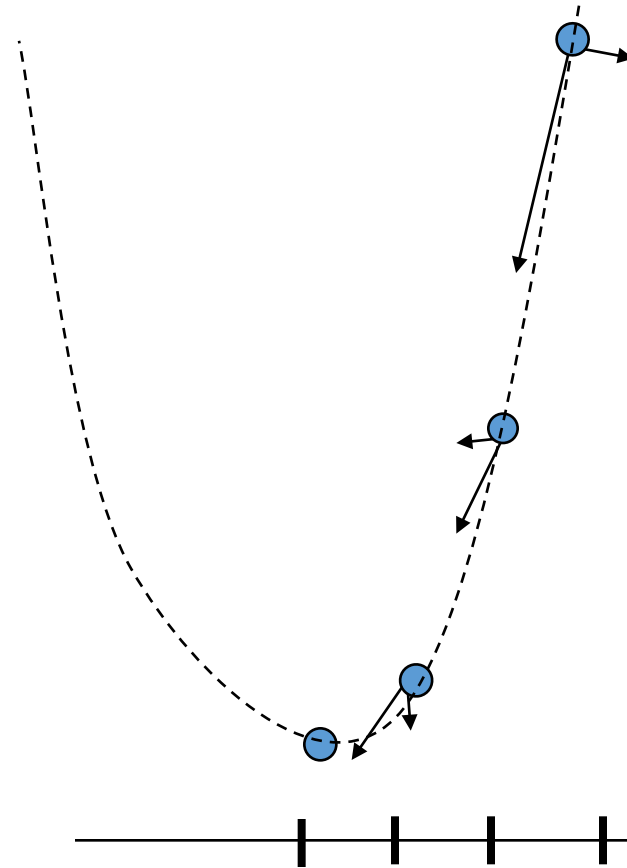
Greedy optimization: gradient descent

- Move in the direction of steepest descent, which is:

$$-[\nabla f(x)]_i = -\frac{\partial}{\partial x_i} f(x)$$

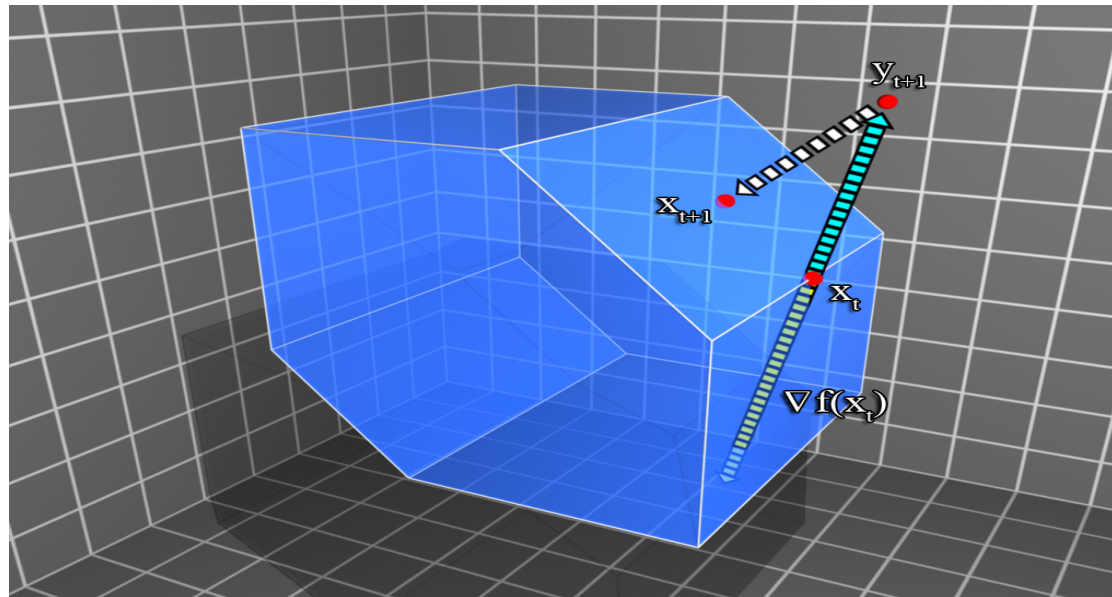
$$x_{t+1} \leftarrow x_t - \eta \nabla f(x_t)$$

“step size” or “Learning rate”



gradient descent – constrained set

$$y_{t+1} \leftarrow x_t - \eta \nabla f(x_t)$$
$$x_{t+1} = \arg \min_{x \in K} |y_{t+1} - x|$$



gradient descent – constrained set

Let:

- G = upper bound on norm of gradients

$$|\nabla f(x_t)| \leq G$$

$$y_{t+1} \leftarrow x_t - \eta \nabla f(x_t)$$
$$x_{t+1} = \arg \min_{x \in K} |y_{t+1} - x|$$

- D = diameter of constraint set

$$\forall x, y \in K \quad |x - y| \leq D$$

Theorem: for step size $\eta = \frac{D}{G\sqrt{T}}$

$$f\left(\frac{1}{T} \sum_t x_t\right) \leq \min_{x^* \in K} f(x^*) + \frac{DG}{\sqrt{T}}$$

online gradient descent

Sequence of functions f_1, f_2, \dots, f_T . Let:

- G = upper bound on norm of gradients

$$|\nabla f_t(x_t)| \leq G$$

$$y_{t+1} \leftarrow x_t - \eta \nabla f_t(x_t)$$
$$x_{t+1} = \arg \min_{x \in K} |y_{t+1} - x|$$

- D = diameter of constraint set

$$\forall x, y \in K \quad |x - y| \leq D$$

Theorem: for step size $\eta = \frac{D}{G\sqrt{T}}$

$$\sum_t f_t(x_t) \leq \min_{x^* \in K} \sum_t f_t(x^*) + 2DG\sqrt{T}$$

Proof:

1. Observation 1:

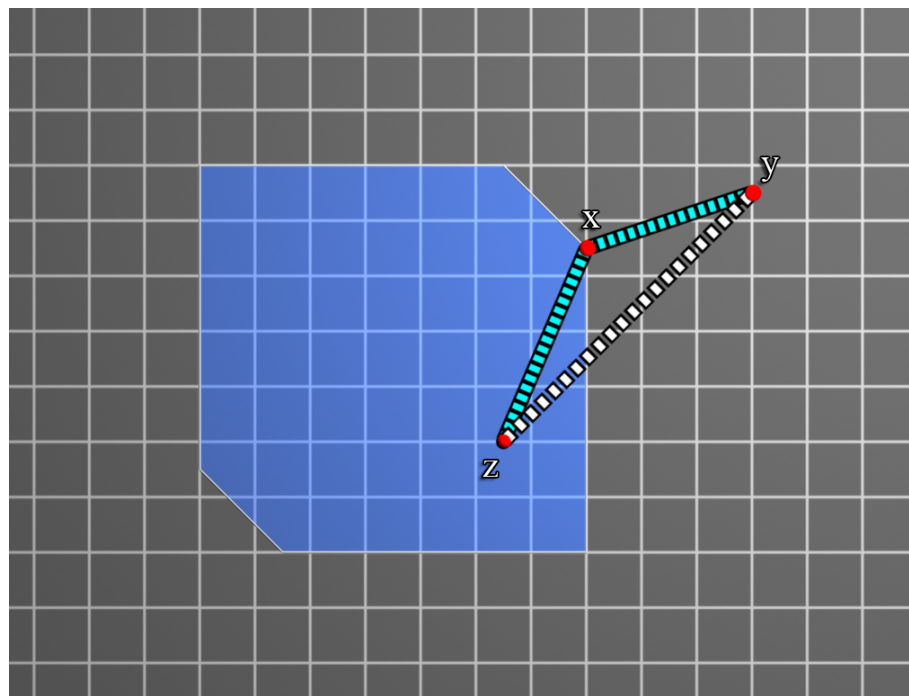
$$|x^* - y_{t+1}|^2 = |x^* - x_t|^2 - 2\eta \nabla f_t(x_t)(x_t - x^*) + |\nabla f_t(x_t)|^2$$

2. Observation 2:

$$|x^* - x_{t+1}|^2 \leq |x^* - y_{t+1}|^2$$

$$y_{t+1} \leftarrow x_t - \eta \nabla f_t(x_t)$$
$$x_{t+1} = \arg \min_{x \in K} |y_{t+1} - x|$$

This is the Pythagorean theorem:



Proof:

1. Observation 1:

$$|\mathbf{x}^* - \mathbf{y}_{t+1}|^2 = |\mathbf{x}^* - \mathbf{x}_t|^2 - 2\eta \nabla f_t(\mathbf{x}_t)(\mathbf{x}_t - \mathbf{x}^*) + \eta^2 |\nabla f_t(\mathbf{x}_t)|^2$$

2. Observation 2:

$$|\mathbf{x}^* - \mathbf{x}_{t+1}|^2 \leq |\mathbf{x}^* - \mathbf{y}_{t+1}|^2$$

Thus:

$$|\mathbf{x}^* - \mathbf{x}_{t+1}|^2 \leq |\mathbf{x}^* - \mathbf{x}_t|^2 - 2\eta \nabla f_t(\mathbf{x}_t)(\mathbf{x}_t - \mathbf{x}^*) + \eta^2 G^2$$

And hence:

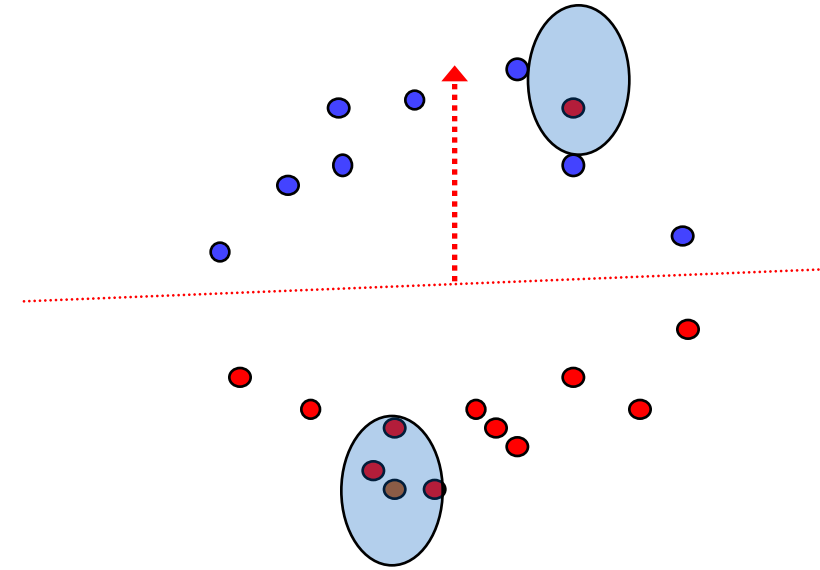
$$\begin{aligned} \sum_t f_t(\mathbf{x}_t) - \sum_t f_t(\mathbf{x}^*) &\leq \sum_t \nabla f_t(\mathbf{x}_t)(\mathbf{x}_t - \mathbf{x}^*) \\ &\leq \sum_t \frac{1}{2\eta} (|\mathbf{x}^* - \mathbf{x}_t|^2 - |\mathbf{x}^* - \mathbf{x}_{t+1}|^2) + \frac{\eta}{2} TG^2 \\ &\leq \frac{1}{2\eta} D^2 + \frac{\eta}{2} TG^2 \leq DG\sqrt{T} \end{aligned}$$

$$\begin{aligned} \mathbf{y}_{t+1} &\leftarrow \mathbf{x}_t - \eta \nabla f_t(\mathbf{x}_t) \\ \mathbf{x}_{t+1} &= \arg \min_{\mathbf{x} \in K} |\mathbf{y}_{t+1} - \mathbf{x}| \end{aligned}$$

GD for linear classification

$$\min_{|w| \leq 1} \frac{1}{m} \sum_i \ell(w^\top x_i, y_i)$$

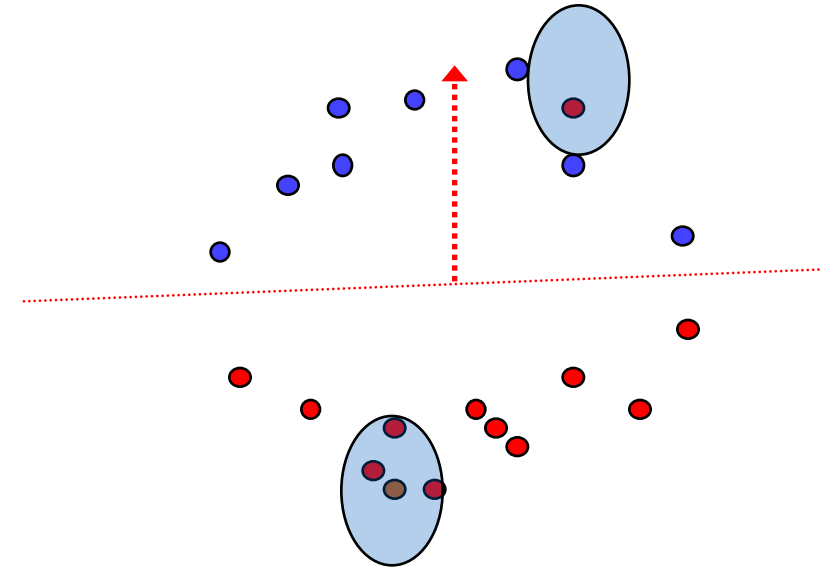
$$w_{t+1} = w_t - \eta \frac{1}{m} \sum_i \ell'(w_t^\top x_i, y_i) x_i$$



- Complexity? $\frac{1}{\epsilon^2}$ iterations, each taking \sim linear time in data set
- Overall $O\left(\frac{md}{\epsilon^2}\right)$ running time, $m = \#$ of examples in \mathbb{R}^d
- Can we speed it up??

GD for linear classification

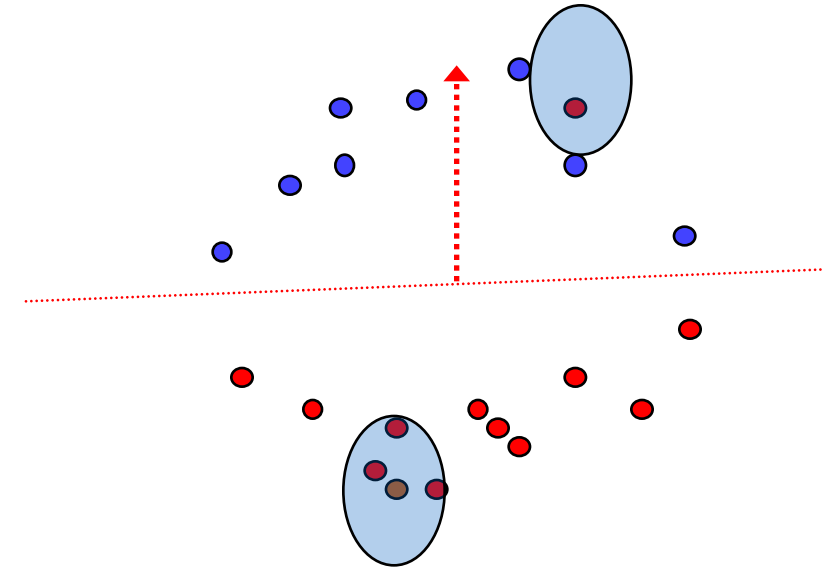
- What if we take a single example, and compute gradient only w.r.t it's loss??
- Which example?
- --> uniformly at random...
- Why would this work?



SGD for linear classification

$$\min_{|w| \leq 1} \frac{1}{m} \sum_i \ell(w^\top x_i, y_i)$$

$$w_{t+1} = w_t - \eta \ell'(w_t^\top x_{i_t}, y_{i_t}) x_{i_t}$$



- Uniformly at random?! $i_t \sim U[1, \dots, m]$

Has expectation = full gradient

- Each iteration is much faster $O(md) \rightarrow O(d)$, convergence??

Crucial for SGD: linearity of expectation and derivatives

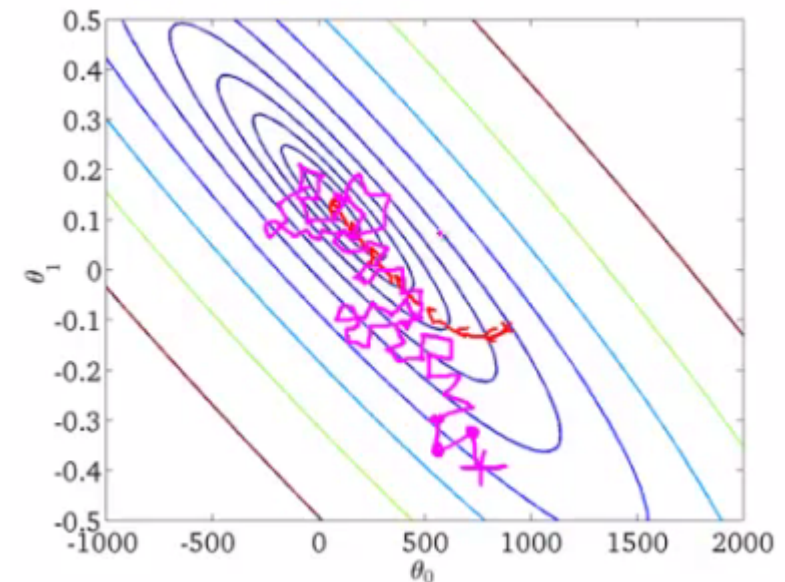
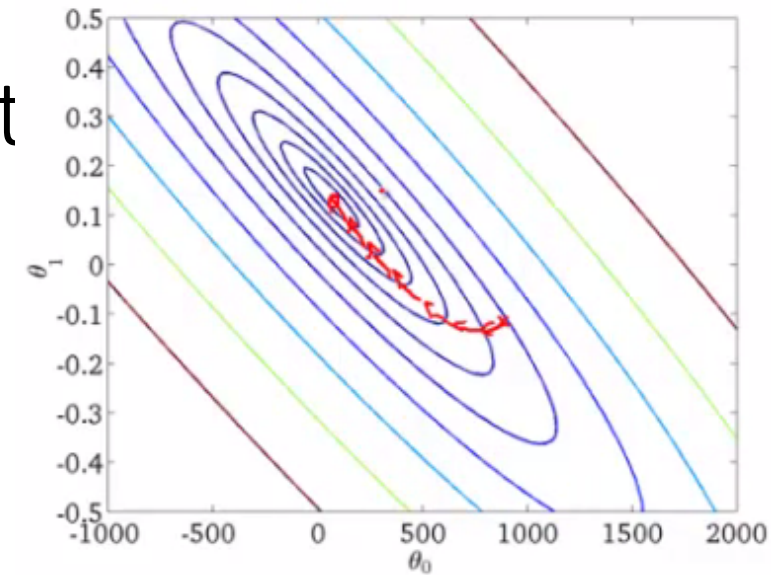
Let $f(w) = \frac{1}{m} \sum_i \ell_i(w)$, then for $i_t \sim U[1, \dots, m]$ chosen uniformly at random, we have

$$\mathbb{E}[\nabla \ell_{i_t}(w)] = \sum_{i=1}^m \frac{1}{m} \nabla \ell_i(w) = \nabla \frac{1}{m} \sum_i \ell_i(w) = \nabla \ell(w)$$

Greedy optimization: gradient descent

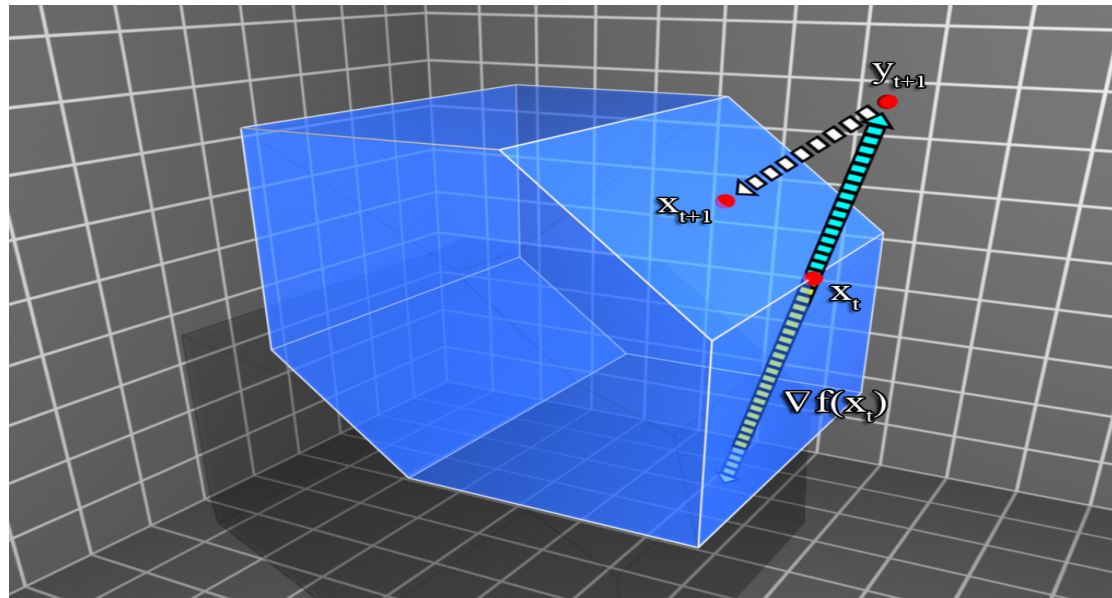
- Move in a random direction, whose expectation is the steepest descent:
- Denote by $\widetilde{\nabla}f(w)$ a vector random variable whose expectation is the gradient,
$$E[\widetilde{\nabla}f(w)] = \nabla f(w)$$

$$x_{t+1} \leftarrow x_t - \eta \widetilde{\nabla}f(x_t)$$



Stochastic gradient descent – constrained case

$$y_{t+1} \leftarrow x_t - \eta \widetilde{\nabla f(x_t)} \quad , \quad \mathbb{E}[\widetilde{\nabla f(x_t)}] = \nabla f(x_t)$$
$$x_{t+1} = \arg \min_{x \in K} |y_{t+1} - x|$$



Stochastic gradient descent – constrained set

Let:

- G = upper bound on norm of gradient **estimators**

$$|\nabla \widetilde{f}(x_t)| \leq G$$

- D = diameter of constraint set

$$\forall x, y \in K \quad |x - y| \leq D$$

Theorem: for step size $\eta = \frac{D}{G\sqrt{T}}$

$$\mathbb{E}[f\left(\frac{1}{T}\sum_t x_t\right)] \leq \min_{x^* \in K} f(x^*) + \frac{DG}{\sqrt{T}}$$

$$\begin{aligned} y_{t+1} &\leftarrow x_t - \eta \nabla \widetilde{f}(x_t) \\ \mathbb{E}[\nabla \widetilde{f}(x_t)] &= \nabla f(x_t) \\ x_{t+1} &= \arg \min_{x \in K} |y_{t+1} - x| \end{aligned}$$

$$y_{t+1} \leftarrow x_t - \eta \nabla \widetilde{f}(x_t)$$

$$\mathbb{E}[\nabla \widetilde{f}(x_t)] = \nabla f(x_t)$$

$$x_{t+1} = \arg \min_{x \in K} |y_{t+1} - x|$$

Proof:

- We have proved: (for any sequence of f_t)

$$\left(\frac{1}{T} \sum_t f_t(x_t) \right) \leq \min_{x^* \in K} \frac{1}{T} \sum_t f_t(x^*) + \frac{DG}{\sqrt{T}}$$

- Let $f_t(x) = \nabla \widetilde{f}(x_t)^\top x = \tilde{v}_t^\top x$, then

$$\left(\frac{1}{T} \sum_t \tilde{v}_t^\top x_t \right) \leq \min_{x^* \in K} \frac{1}{T} \sum_t \tilde{v}_t^\top x^* + \frac{DG}{\sqrt{T}}$$

- By property of expectation:

$$\mathbb{E}\left[f\left(\frac{1}{T} \sum_t x_t\right) - \min_{x^* \in K} f(x^*) \right] \leq \mathbb{E}\left[\frac{1}{T} \sum_t \nabla \widetilde{f}(x_t)^\top (x_t - x^*) \right] \leq \frac{DG}{\sqrt{T}}$$

Recap

$$y_{t+1} \leftarrow x_t - \eta \nabla \widetilde{f}(x_t)$$
$$\mathbb{E}[\nabla \widetilde{f}(x_t)] = \nabla f(x_t)$$
$$x_{t+1} = \arg \min_{x \in K} |y_{t+1} - x|$$

- SGD convergence:

$$\mathbb{E}[f\left(\frac{1}{T} \sum_t x_t\right) - \min_{x^* \in K} f(x^*)] \leq \frac{DG}{\sqrt{T}}$$

- Faster per-iteration step $O(d)$ instead of $O(md)$
- Faster convergence? (in terms of T)

Convexity

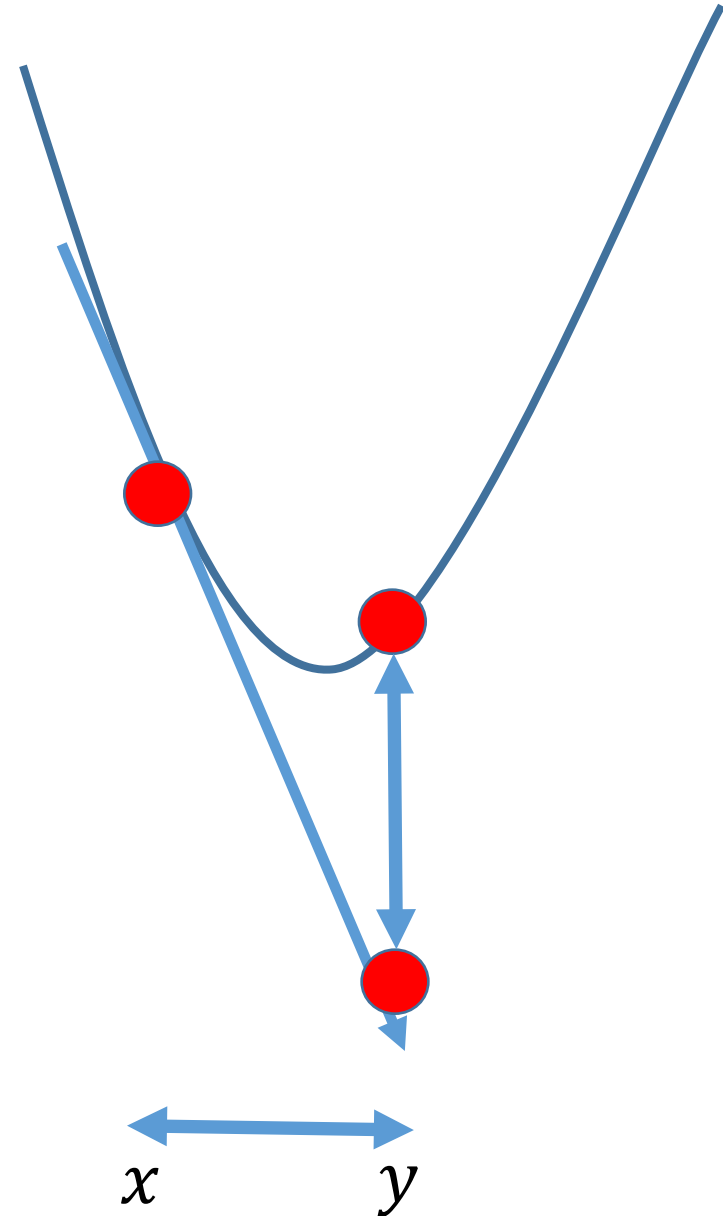
- definition:

$$f(y) \geq f(x) + \nabla f(x)^\top (y - x)$$

- if twice differentiable:
second derivative is non-negative (in 1D)

- Higher dimensions:

$$\nabla^2 f(x) \succcurlyeq 0$$



Strong convexity and smoothness

Convexity:

$$\nabla^2 f(x) \succeq 0$$

Strong convexity:

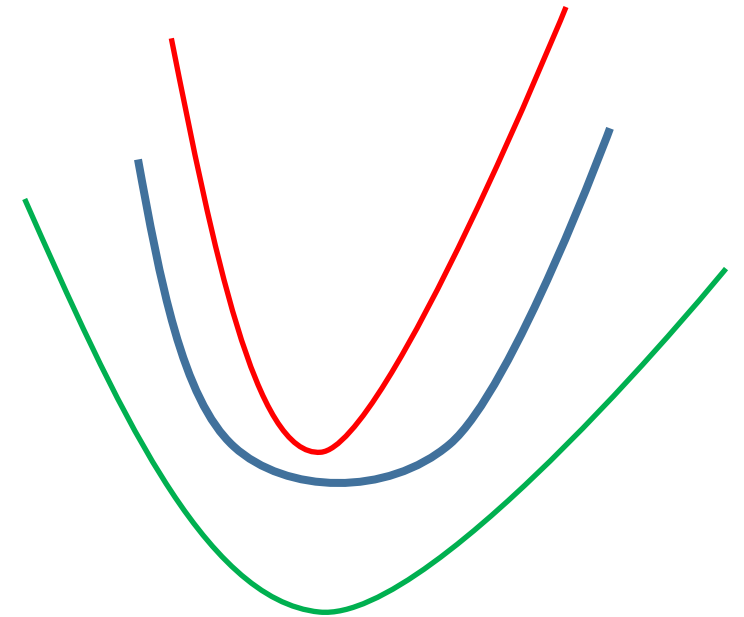
$$\nabla^2 f(x) \succeq \alpha I$$

Smoothness:

$$\nabla^2 f(x) \preceq \beta I$$

GD convergence:

$$f(x_T) - \min_{x^* \in K} f(x^*) \leq O\left(e^{-\frac{\alpha}{\beta} t}\right) \text{ (as opposed to } \frac{DG}{\sqrt{T}})$$



Summary

- (online) gradient descent algorithm with constraints
- Stochastic GD
- Conditions that imply faster convergence