

Introduction to Machine Learning (IML)

LECTURE #7: OPTIMIZATION

236756 – 2022-2023 WINTER – TECHNION

LECTURER: YONATAN BELINKOV



Today

- **Part II:** *the different aspects of learning*
 1. Statistics: generalization and PAC theory
 2. Modeling:
 - Error decomposition
 - Regularization
 - Model selection
 3. Optimization: convexity, gradient descent (today)
 4. Practical aspects and potential pitfalls

Optimization

Optimizing the learning objective

- **Regularized Risk Minimization (RRM):**

$$\operatorname{argmin}_{w \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \underbrace{\ell(y_i, w^\top x_i)}_{\text{loss}} + \underbrace{\lambda R(w)}_{\text{regularization}}$$

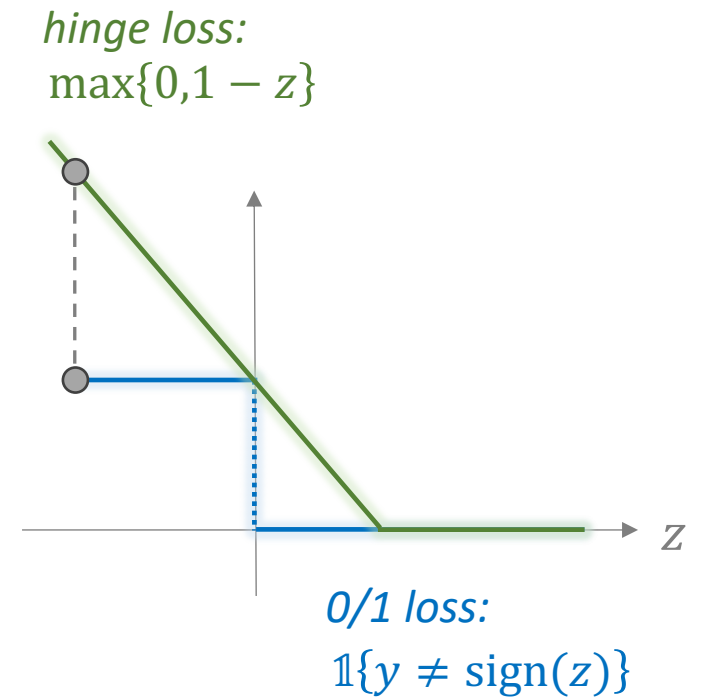
- **Recall:** ERM of 0/1 is hard (discrete)
- **Solution:** continuous proxy

Optimizing the learning objective

- **Regularized Risk Minimization (RRM):**

$$\operatorname{argmin}_{w \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \underbrace{\ell(y_i, w^\top x_i)}_{\text{loss}} + \underbrace{\lambda R(w)}_{\text{regularization}}$$

- **Recall:** ERM of 0/1 is hard (discrete)
- **Solution:** continuous proxy
- Soft SVM's *hinge loss* is one option
- **Problem:** unnecessarily large penalties for far from margin

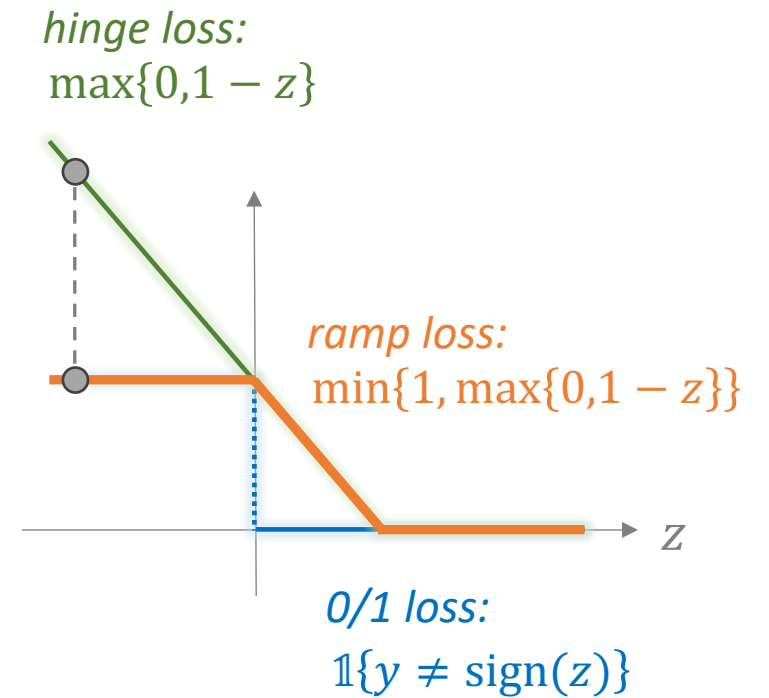


Optimizing the learning objective

- **Regularized Risk Minimization (RRM):**

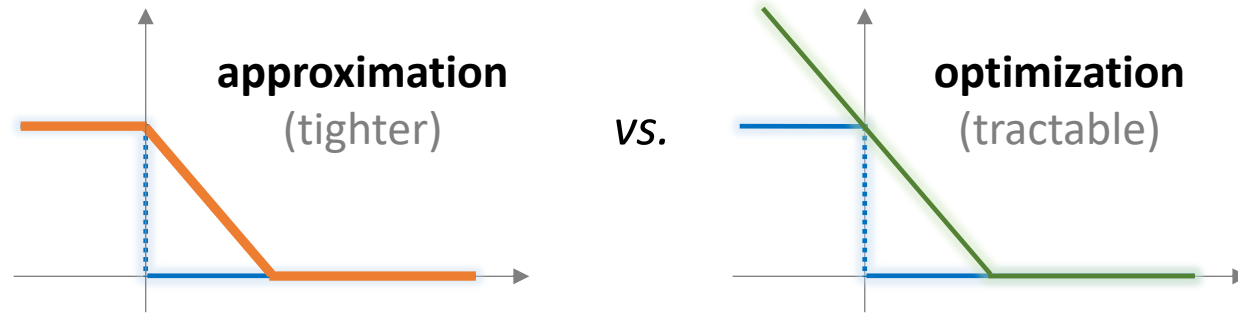
$$\operatorname{argmin}_{w \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \underbrace{\ell(y_i, w^\top x_i)}_{\text{loss}} + \underbrace{\lambda R(w)}_{\text{regularization}}$$

- **Recall:** ERM of 0/1 is hard (discrete)
- **Solution:** continuous proxy
- Soft SVM's *hinge loss* is one option
- **Problem:** unnecessarily large penalties for far from margin
- **Alternative:** *ramp loss*
- **Problem:** optimization... (today we'll see why)



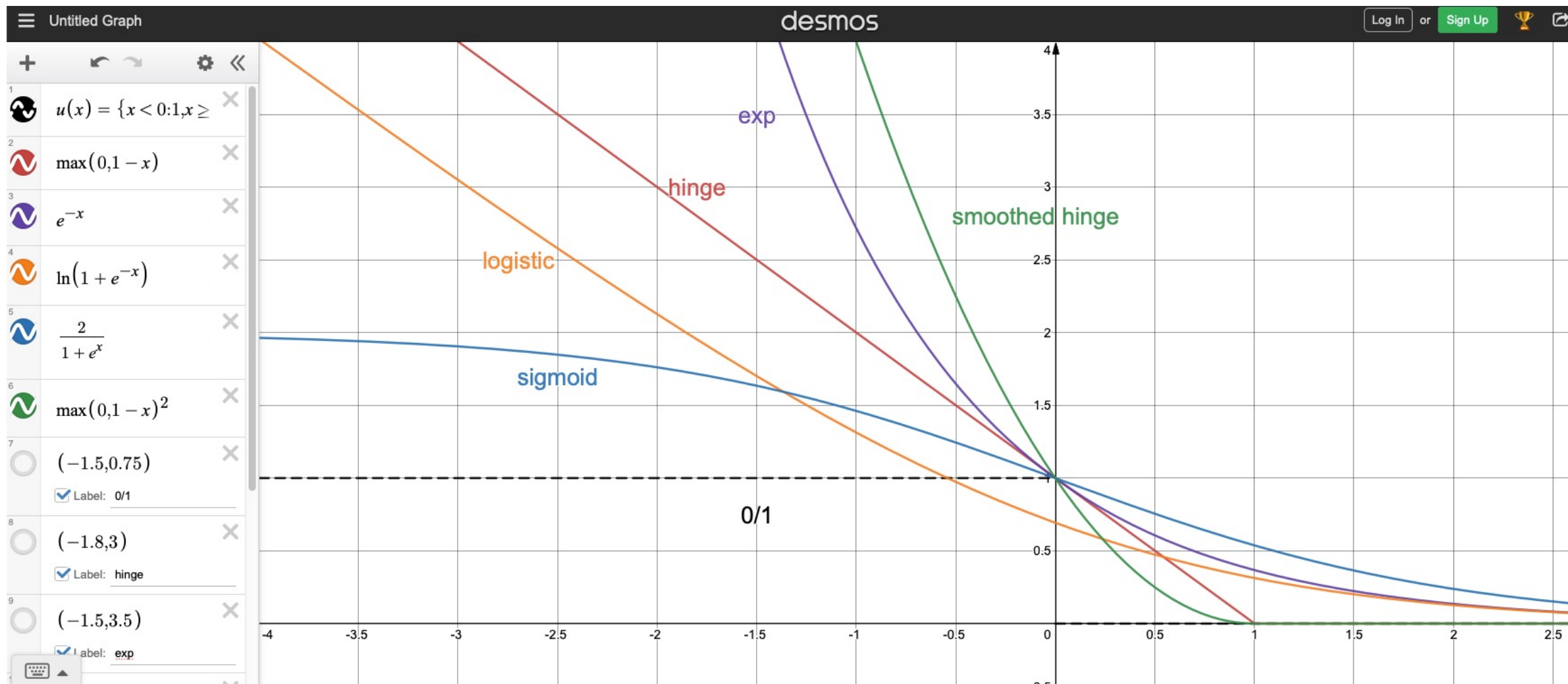
Optimizing the learning objective

- **Tradeoff:**



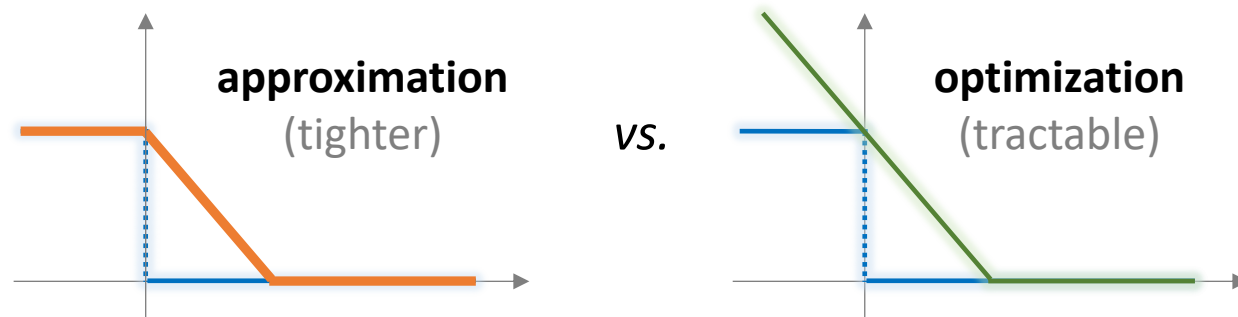
- Many proxies out there!
- [DESMOS]

Proxy losses



Optimizing the learning objective

- **Tradeoff:**



- Many proxies out there – **which is better?**
- No definite answer
- But we can say some things sometimes
- **Today:** focus on optimization (but also generalization)
- Functional properties that give optimization guarantees:

- (sub-)differentiable
- convex
- Lipschitz
- smooth
- bounded

1		$u(x) = \begin{cases} x < 0: 1, x \geq 0: 0 \end{cases}$
2		$\max(0, 1 - x)$
3		e^{-x}
4		$\ln(1 + e^{-x})$
5		$\frac{1}{1 + e^x}$
6		$\max(0, 1 - x)^2$
7		...
∞		

Gradient descent



One algorithm to rule them all

- **Goal:** solve $\min_{w \in \mathbb{R}^d} f(w)$

- (in learning: f is the learning objective, w are model parameters)

- **Approach:** Gradient Descent (GD) [\sim Cauchy 1847]

$$w_{t+1} = w_t - \eta \nabla f(w_t)$$

- “First order iterative method”
- Forms the basis for many other optimization algorithms
(we’ll see some today)

Gradient descent

- **Definition:**

The *gradient* of f w.r.t. w is a vector function of partial derivatives,

$$\nabla f = \begin{pmatrix} \frac{\partial f}{\partial w_1} \\ \vdots \\ \frac{\partial f}{\partial w_d} \end{pmatrix}$$

- *Gradient evaluation* at $\bar{w} \in \mathbb{R}^d$:

$$\nabla f(\bar{w}) \in \mathbb{R}^d$$

Gradient Descent (GD)

- Initialize w_0 (e.g., $w_0 = \vec{0}$)
- Repeat:
 - $w_{t+1} = w_t - \eta \nabla f(w_t)$
- Until convergence
(e.g., $\|w_{t+1} - w_t\| \leq \epsilon$)

Gradient descent

- **View I:** fog in mountains
- **View II:** local greedy descent
- **View III:** optimize simple approximations

$$@t: f(w) \approx f(w_t) + \nabla f(w_t)^\top (w - w_t)$$

- **View IV:** Taylor expansion + tradeoff
 - $f(w_t + \delta) \approx f(w_t) + \nabla f(w_t)^\top \delta, \quad \delta \in \mathbb{R}^d$
 - $w_{t+1} = w_t + \delta$ close to w_t , i.e., $\|\delta\|_2$ is small
 - $w_{t+1} = \operatorname{argmin}_w \frac{1}{2} \|w - w_t\|_2^2 + \eta (f(w_t) + \nabla f(w_t)^\top (w - w_t))$
 - $\partial = 0 \Rightarrow w_{t+1} = w_t - \eta \nabla f(w_t)$

Gradient Descent (GD)

- Initialize w_0 (e.g., $w_0 = \vec{0}$)
- Repeat:
 - $w_{t+1} = w_t - \eta \nabla f(w_t)$
- Until convergence
(e.g., $\|w_{t+1} - w_t\| \leq \epsilon$)

Gradient descent

- **View I:** fog in mountains
- **View II:** local greedy descent
- **View III:** optimize simple approximations

$$@t: f(w) \approx \nabla f(w_t)^\top w$$

- **View IV:** Taylor expansion + tradeoff
 - $f(w_t + \delta) \approx f(w_t) + \nabla f(w_t)^\top \delta, \delta \in \mathbb{R}^d$
 - $w_{t+1} = w_t + \delta$ close to w_t , i.e., $\|\delta\|_2$ is small
 - $w_{t+1} = \operatorname{argmin}_w \frac{1}{2} \|w_t - w\|_2^2 + \eta (f(w_t) + \nabla f(w_t)^\top (w - w_t))$
 - $\partial = 0 \Rightarrow w_{t+1} = w_t - \eta \nabla f(w_t)$

Gradient Descent (GD)

- Initialize w_0 (e.g., $w_0 = \vec{0}$)
- Repeat:
 - $w_{t+1} = w_t - \eta \nabla f(w_t)$
- Until convergence (e.g., $\|w_{t+1} - w_t\| \leq \epsilon$)



Need to determine:

1. Step size or learning rate η
2. Stopping criterion

Gradient descent

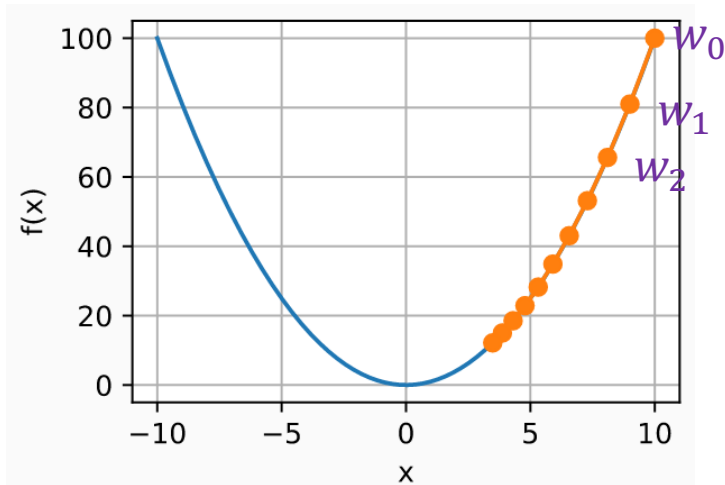
- **Definition:** Let $B_\epsilon(w) = \{v : \|w - v\| \leq \epsilon\}$ be a ball of radius ϵ around w . Then w is a **local minimum** of f if

$$\exists \epsilon > 0 \forall v \in B_\epsilon(w), f(w) \leq f(v)$$

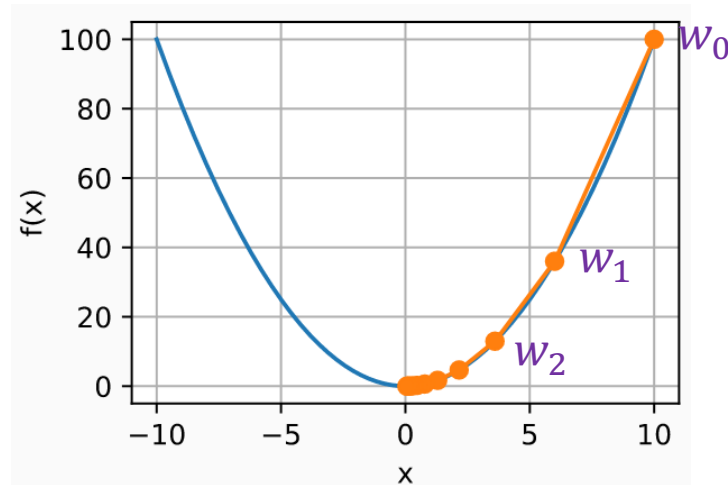
- **Claim:** if η is small enough, then $f(w_{t+1}) < f(w_t)$ for any non-minimum w_t (this is the “descent” part of “gradient descent”)
- **Proof:** consider non-minimum w_t and negate
- **Corollary:** GD converges* to a local minimum (or saddle point!)
* for an appropriate choice of η
- **Trick:** use time-varying learning rate η_t , for example, $\eta_t = 1/t$
- **Note:** if w_t is a local minimum then $\nabla f(w_t) = 0$ (but not iff!)

Learning rate

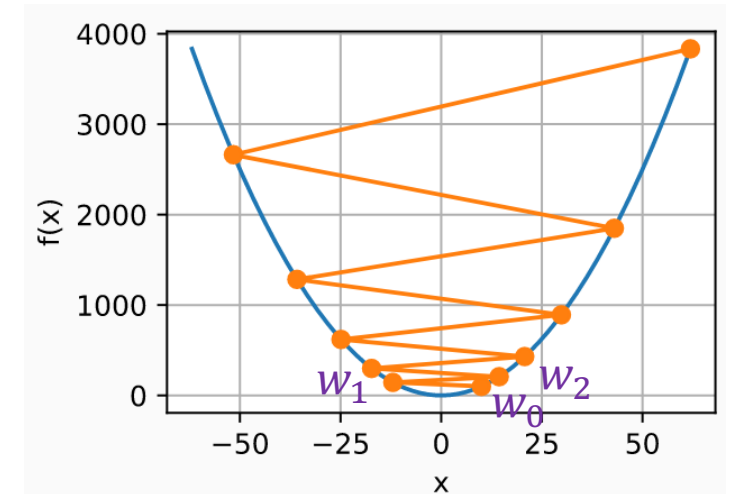
- **Q1:** So GD converges for learning rate $\eta_t = 1/t$. Are we done?
- **A1:** No!



$\eta_t = 1/t$
slow convergence
(too slow!)



$\eta_t = 1/\sqrt{t}$
"just right"
(how can we know?)

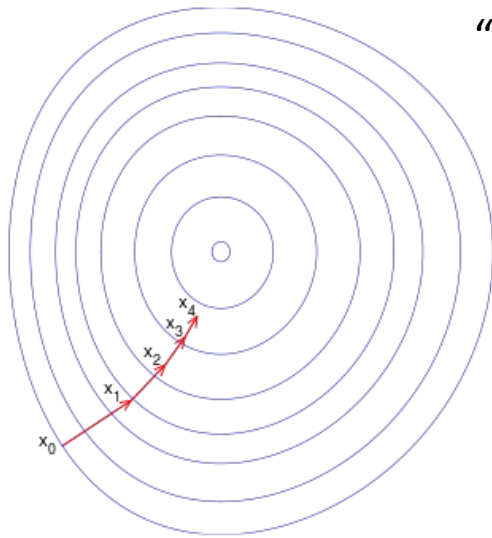


$\eta_t = 10$
divergence
(too fast!)

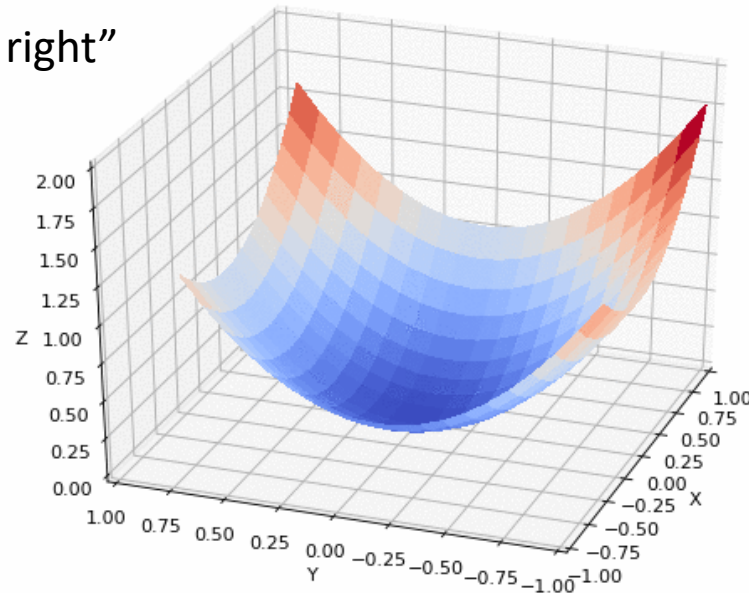
(illustrative)

Learning rate

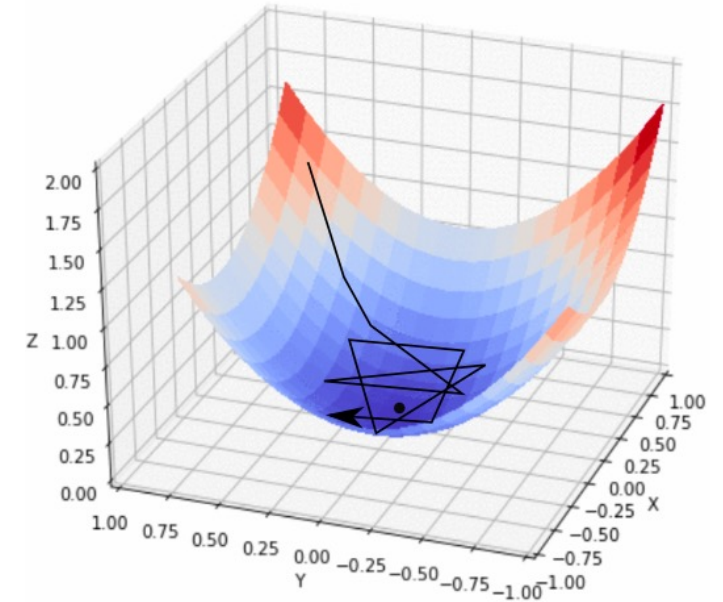
- **Q1:** So GD converges for learning rate $\eta_t = 1/t$. Are we done?
- **A1:** No!



“just right”



divergence?
slow convergence?



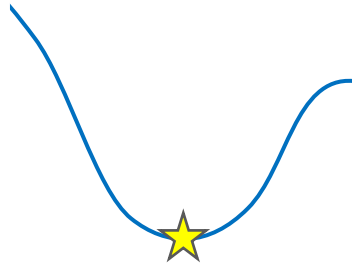
- Choosing the learning rate can be tricky – we will return to this.

When to stop

- **Q2:** So GD converges in theory. But in practice, how do we know when (and if) it did?
- **A2:** Can't "know", but can estimate.
- Set criteria for stopping, for example:
 - stop when $\|w_{t+1} - w_t\| \leq \epsilon$
 - or when $\|f(w_{t+1}) - f(w_t)\| \leq \epsilon$
 - or when $\|\nabla f(w_t)\| \leq \epsilon$
 - or ...
- Actually, knowing when to stop is (also) tricky – and has statistical implications!
We will return to this as well.

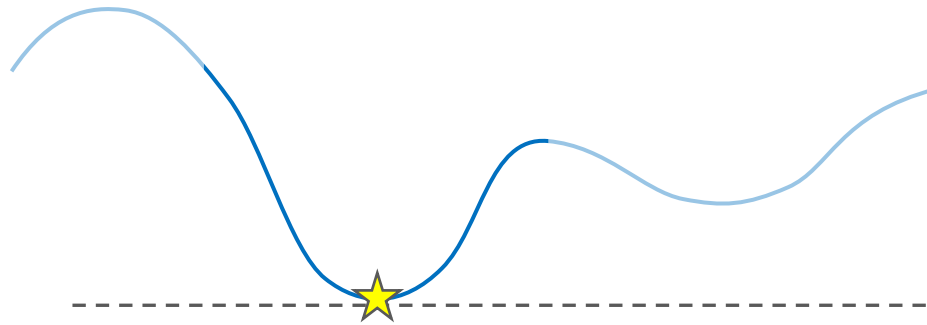
Local minima

- **Q3:** So GD can converge to a local minimum. Should we be happy?
- **A3:** It depends...



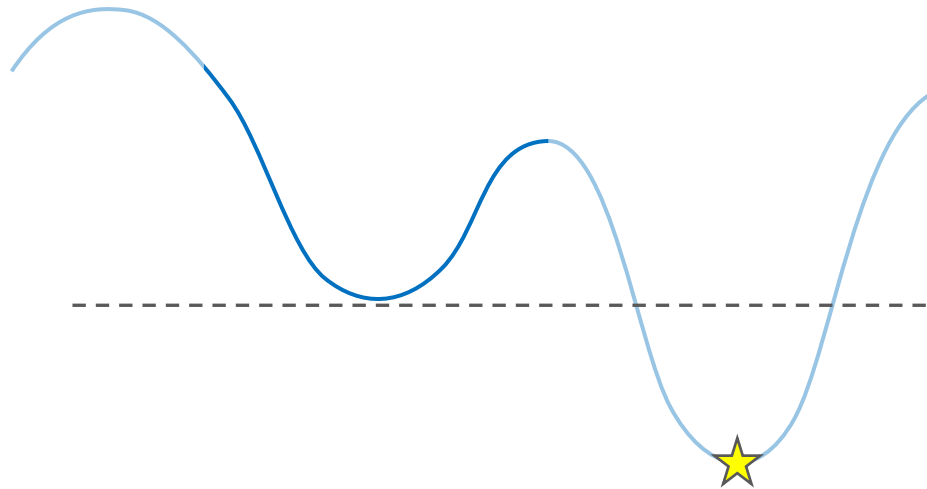
Local minima

- **Q3:** So GD can converge to a local minimum. Should we be happy?
- **A3:** It depends...



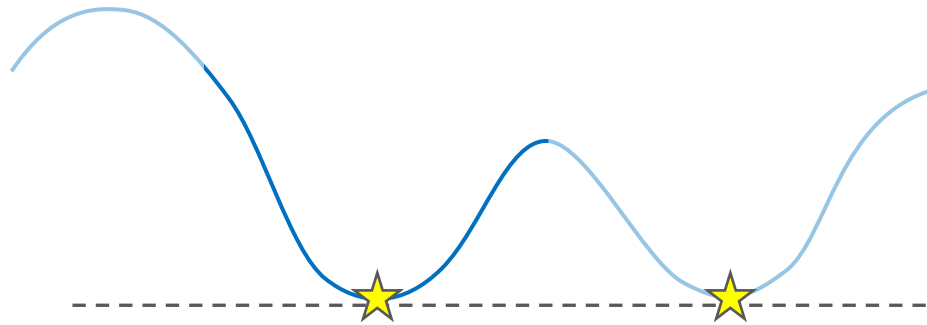
Local minima

- **Q3:** So GD can converge to a local minimum. Should we be happy?
- **A3:** It depends...



Local minima

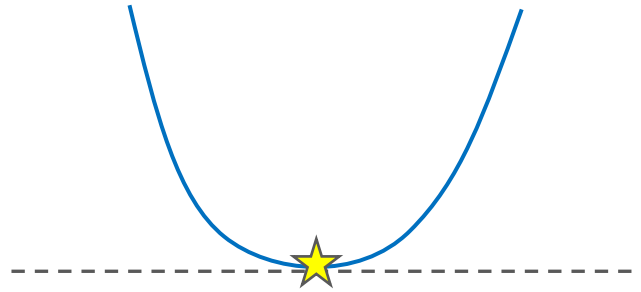
- **Q3:** So GD can converge to a local minimum. Should we be happy?
- **A3:** It depends...



- **Goal:** find conditions on f for which GD is guaranteed to work well

Local minima

- **Q3:** So GD can converge to a local minimum. Should we be happy?
- **A3:** It depends...



- **Goal:** find conditions on f for which GD is guaranteed to work well

Convexity

Descending with guarantees

- Recall foggy mountain story
- When will the “going down” approach help you reach your cabin?
- When mountain range is:
 1. continuous
 2. smooth
 3. **has a single valley**
- Enter **convexity**.
- [on board]



Convex sets

- **Definition:** *convex combination*:

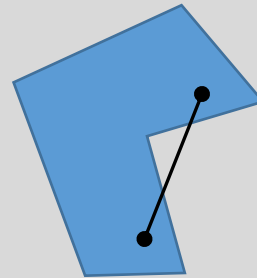
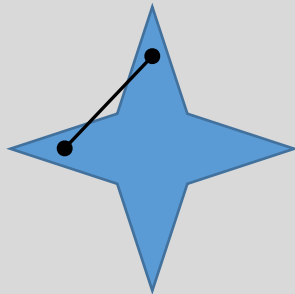
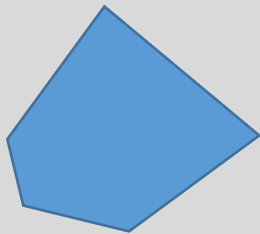
$$\alpha u + (1 - \alpha)v \quad \text{for } u, v \in \mathbb{R}^d, \alpha \in [0,1]$$

- As interpolation – point on line between u, v

- **Definition:** C is a *convex set* if

$$\forall u, v \in C, \text{ any convex combination of } u, v \text{ is also in } C \quad \mathbb{R}^d$$

- **Examples:**



$$H_{w,b} = \{x : w^\top x \geq b\}$$

- **Definition:** For C convex set, $f: C \rightarrow \mathbb{R}$ is *convex function* if $\forall u, v \in C, \alpha \in [0,1]$:
$$f(\alpha u + (1 - \alpha)v) \leq \alpha f(u) + (1 - \alpha)f(v)$$

GD convergence rates

- β -Smooth: $f(x_t) - f(x^*) \leq \frac{2\beta\|x_0 - x^*\|}{t+4} = O\left(\frac{\beta}{t}\right)$
 - Lower bound: $\Omega\left(\frac{\beta}{t^2}\right)$
- L -Lipschitz: $f(x_t) - f(x^*) = O\left(\frac{L}{\sqrt{t}}\right)$
- σ -Strongly convex: $f(x_t) - f(x^*) = O(e^{-2\sigma t})$

- **Depends on:**

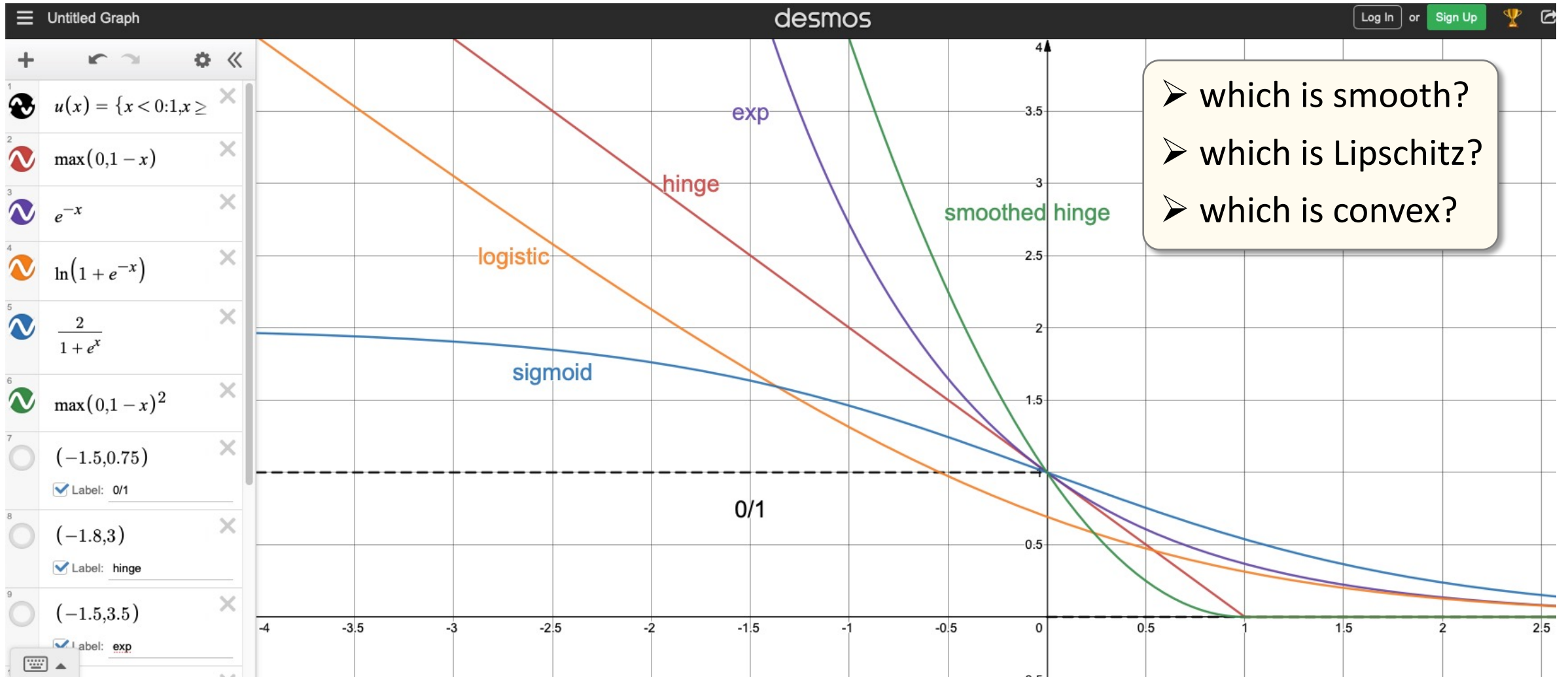
- Type
- Parameter
- Initial guess!

f is β -Smooth if
$$\|\nabla f(u) - \nabla f(w)\| \leq \beta\|u - w\|$$

f is L -Lipschitz if
$$\|f(u) - f(w)\| \leq L\|u - w\|$$

f is σ -Strongly convex if
$$f(\alpha w + (1 - \alpha)u) \leq \alpha f(w) + (1 - \alpha)f(u) - \frac{\sigma}{2}\alpha(1 - \alpha)\|u - w\|^2$$

GD convergence rates



- which is smooth?
- which is Lipschitz?
- which is convex?

Optimization for learning

Convex learning problems

- Recall ERM/RLM:

$$\operatorname{argmin}_{\theta} \frac{1}{m} \sum_{i=1}^m \ell(y_i, f_{\theta}(x_i)) + \lambda R(\theta)$$

- How can we tell if the learning objective is convex?

Composition rules

- **Convexity-preserving operations:**

1. **scaling:** f convex, $\alpha \geq 0 \Rightarrow \alpha f$ convex [easy]
2. **sum:** f, g convex $\Rightarrow f + g$ convex [tirgul]
3. **noneg. weighted sum:** f_i convex, $\alpha_i \geq 0 \Rightarrow \sum_i \alpha_i f_i$ convex
4. **composition:** f convex, g linear $\Rightarrow f \circ g$ convex [on board]
5. ...

- (these are sufficient conditions, but not necessary)

Composition of convex and linear is convex

- **Claim:** if g is convex and h is linear, then $f = g \circ h$ is convex
- **Proof:**
$$\begin{aligned} f(\alpha u + (1 - \alpha)v) &= g(h(\alpha u + (1 - \alpha)v)) \stackrel{h \text{ linear}}{=} \\ g((\alpha u + (1 - \alpha)v)^T x + b) &= g(\alpha u^T x + (1 - \alpha)v^T x + b) = \\ g(\alpha u^T x + (1 - \alpha)v^T x + (\alpha + 1 - \alpha)b) &= \\ g(\alpha(u^T x + b) + (1 - \alpha)(v^T x + b)) &= g(\alpha h(u) + (1 - \alpha)h(v)) \\ \leq_{g \text{ convex}} \alpha g(h(u)) + (1 - \alpha)g(h(v)) &= \alpha f(u) + (1 - \alpha)f(v) \end{aligned}$$

Convex learning problems

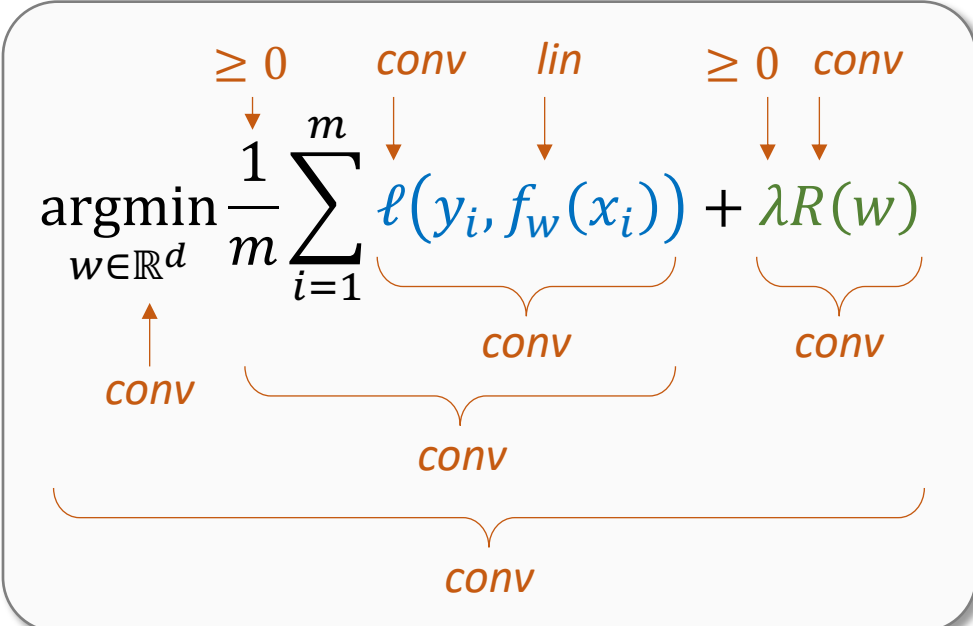
- **Claim:** The following conditions are sufficient for the learning objective to be convex (in w):
 1. $\ell(\cdot, \cdot)$ is convex (in its second argument)
 2. $R(\cdot)$ is convex
 3. f_w is linear (in w , i.e., $f_w(x) = w^\top x$)

Convex learning problems

- **Claim:** The following conditions are sufficient for the learning objective to be convex (in w):

1. $\ell(\cdot, \cdot)$ is convex (in its second argument)
2. $R(\cdot)$ is convex
3. f_w is linear (in w , i.e., $f_w(x) = w^\top x$)

- **Let's try** 


$$\operatorname{argmin}_{w \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \ell(y_i, f_w(x_i)) + \lambda R(w)$$

Annotations in the diagram:

- ≥ 0 (above $\frac{1}{m}$)
- conv (above ℓ)
- lin (above f_w)
- ≥ 0 (above λ)
- conv (above R)
- conv (below ℓ)
- conv (below R)
- conv (below the sum)
- conv (below the entire expression)

Convex learning problems

- **Claim:** The following conditions are sufficient for the learning objective to be convex (in w):

1. $\ell(\cdot, \cdot)$ is convex (in it's second argument)
2. $R(\cdot)$ is convex
3. f_w is linear (in w , i.e., $f_w(x) = w^\top x$)

- **Let's try** →

- **Claim:** The Soft SVM objective is convex

- **Need to prove:**

1. $\max\{0, 1 - z\}$ is convex [ex]
2. ℓ_2 -norm squared is convex [tirgul]

- **Corollary:** can solve with GD!

$$\operatorname{argmin}_{w \in \mathbb{R}^d} \underbrace{\frac{1}{m} \sum_{i=1}^m \underbrace{\ell(y_i, f_w(x_i))}_{\text{conv}}}_{\text{conv}} + \underbrace{\lambda R(w)}_{\text{conv}}$$

Annotations: ≥ 0 (above $\frac{1}{m}$), conv (above ℓ), lin (above f_w), ≥ 0 (above λ), conv (above R). Brackets indicate the convexity of the summands and the overall expression.

$$\operatorname{argmin}_{w \in \mathbb{R}^d} \frac{1}{m} \sum_{i=1}^m \max\{0, 1 - y_i \cdot w^\top x_i\} + \lambda \|w\|_2^2$$

Variations, extensions,
and beyond

The computational cost of GD

- For running GD we need to compute gradients of the learning objective.

- **Note:**
$$\nabla \left(\frac{1}{m} \sum_{i=1}^m \ell(y_i, f_w(x_i)) \right) = \frac{1}{m} \sum_{i=1}^m \underbrace{\nabla \ell(y_i, f_w(x_i))}_{\nabla_i}$$

- **For Soft SVM:**
 - Each $\nabla \max\{0, 1 - y_i \cdot w^\top x_i\}$ costs: $O(d)$
 - Overall: $O(dm)$
- GD can be costly!
- **Solution:** use *approximate* gradients

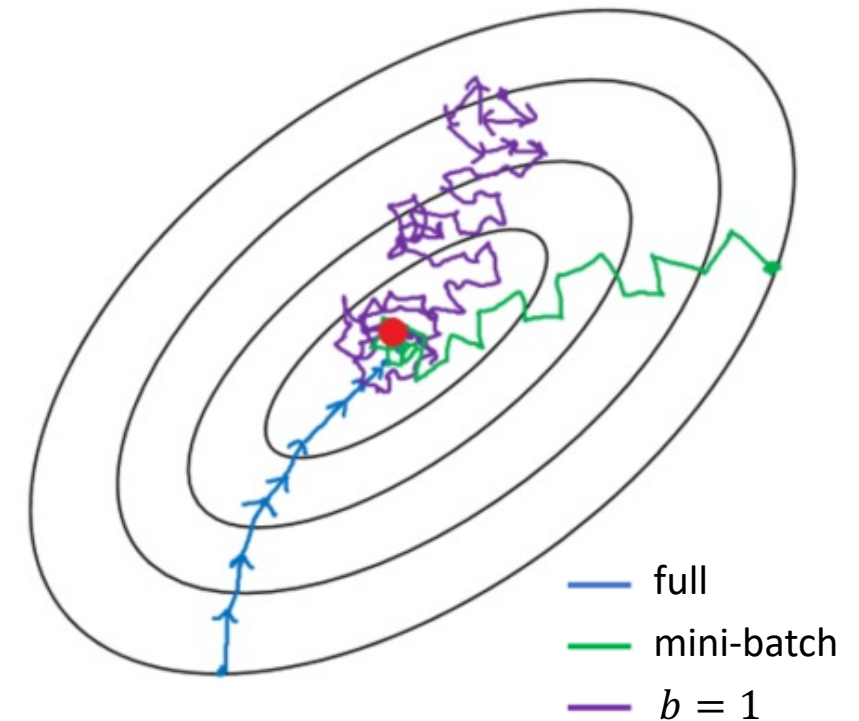
Taking approximate gradient steps

- **Observation:** sampling average is unbiased (but noisy) estimator
- **Idea:** replace full gradient with average of small random set of examples

- **Stochastic Gradient Descent (SGD):***

1. Sample small random “mini-batch” $B \subset S$ of size b
2. Compute average gradient $\bar{\nabla} = \frac{1}{b} \sum_{i \in B} \nabla_i$
3. Apply *approximate* gradient step $x_{t+1} = x_t - \eta \bar{\nabla}$

- **Pros:** reduces compute time (significantly!)
- **Cons:** adds noise (but often worth it)
- Hyper parameter b trades off compute time with noise



* The name “SGD” typically refers to the $b = 1$ case

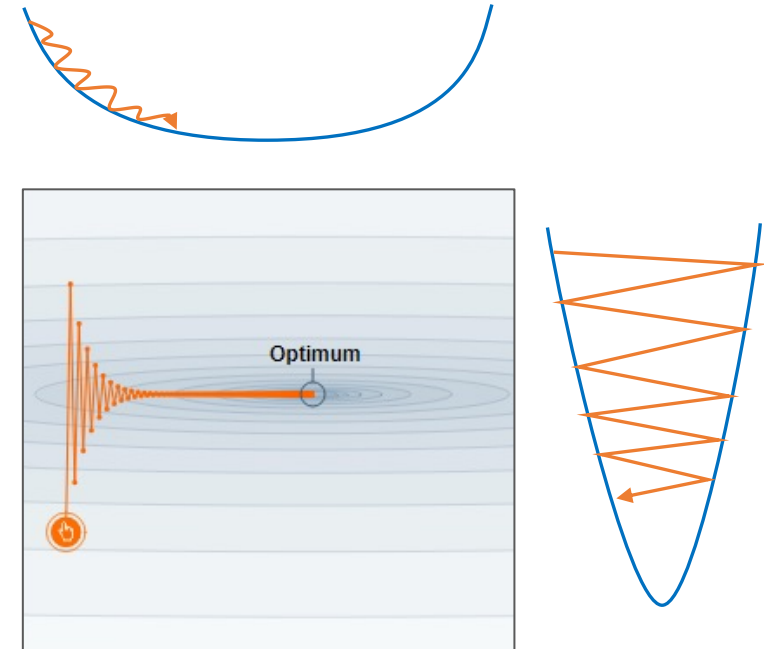
Looking back

- Consider 2D quadratic with highly varying curvatures
 - On low-curvature dimension, GD crawls hesitantly
 - On high-curvature dimension, GD oscillates frantically
- Adding “momentum” sorts this out:

$$\begin{aligned}v_{t+1} &= \beta v_t - \eta \nabla f(w_t) && \rightarrow \text{velocity/“memory”} \\w_{t+1} &= w_t + v_{t+1} && \rightarrow \text{position} \\&= w_t - \eta \nabla f(w_t) + \beta v_t\end{aligned}$$

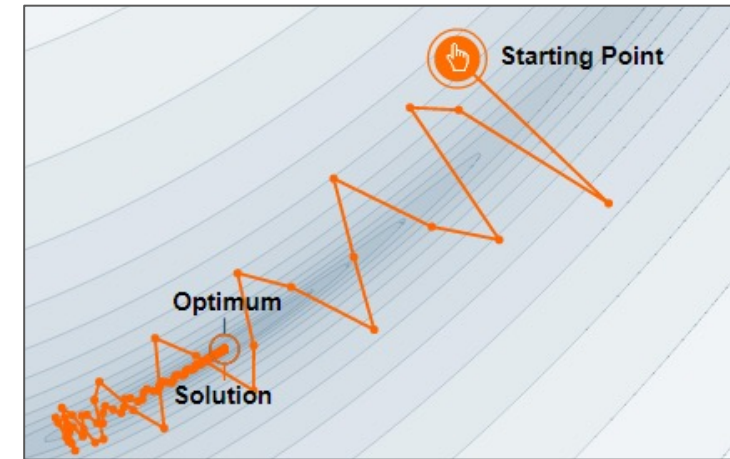
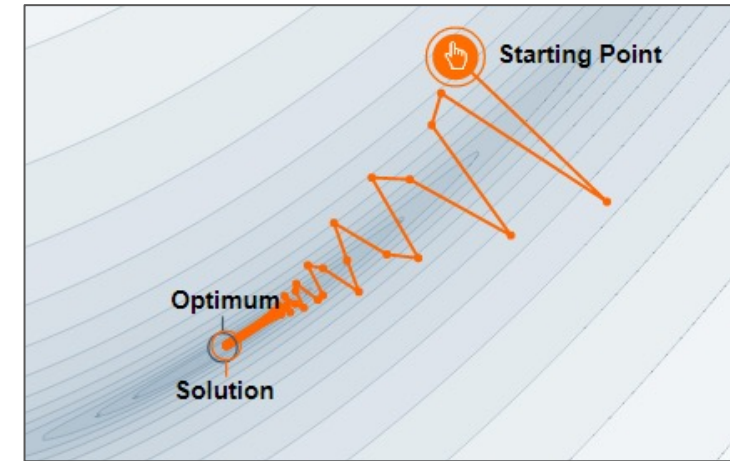
(imagine a rolling ball)

- Effects of momentum:
 - *increases* in dimensions where gradient preserves direction
 - *decreases* in dimensions where gradient direction varies
- Typical β values: 0.9, 0.95, 0.99



Looking forward

- Momentum is great, but can overshoot
- **Solution:** look into *future*
- **Idea:** *imagine* momentum has been applied, then compute gradient



Looking forward

- Momentum is great, but can overshoot
- **Solution:** look into *future*
- **Idea:** *imagine* momentum has been applied, then compute gradient

- **Nesterov's accelerated gradient:**

$$\tilde{w}_{t+1} = w_t + \beta v_t$$

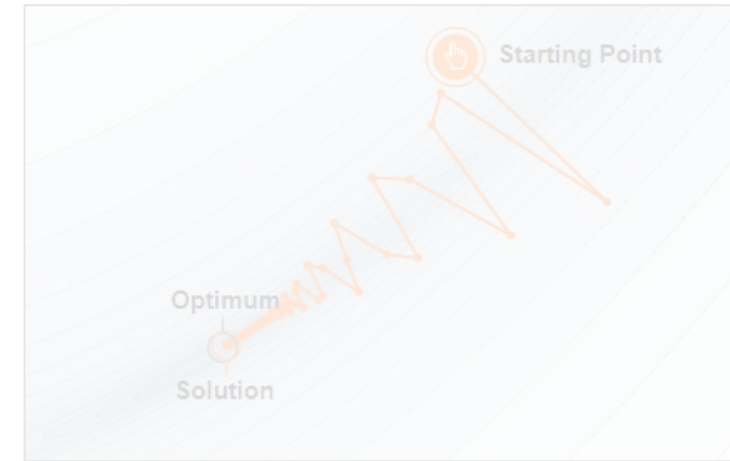
$$v_{t+1} = \beta v_t - \eta \nabla f(\tilde{w}_{t+1})$$

$$w_{t+1} = w_t + v_{t+1}$$

- **Equivalent:**

$$v_{t+1} = \beta v_t - \eta \nabla f(w_t)$$

$$w_{t+1} = w_t + (1 + \beta)v_{t+1} - \beta v_t$$

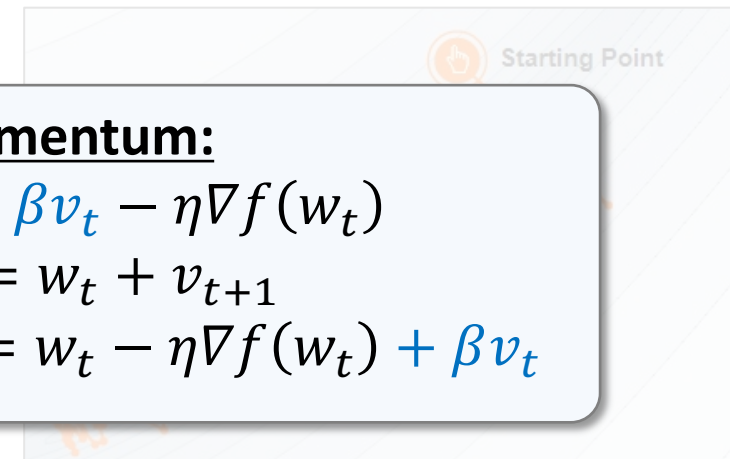


Vs. momentum:

$$v_{t+1} = \beta v_t - \eta \nabla f(w_t)$$

$$w_{t+1} = w_t + v_{t+1}$$

$$= w_t - \eta \nabla f(w_t) + \beta v_t$$



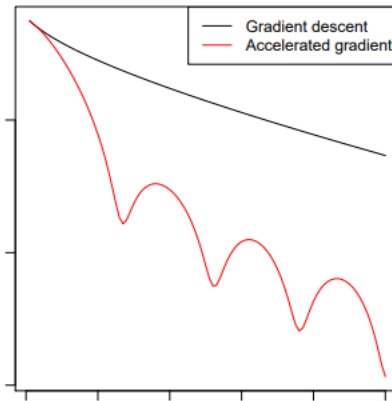
Looking forward

- **Nesterov's accelerated gradient:**

$$v_{t+1} = \beta v_t - \eta \nabla f(w_t)$$

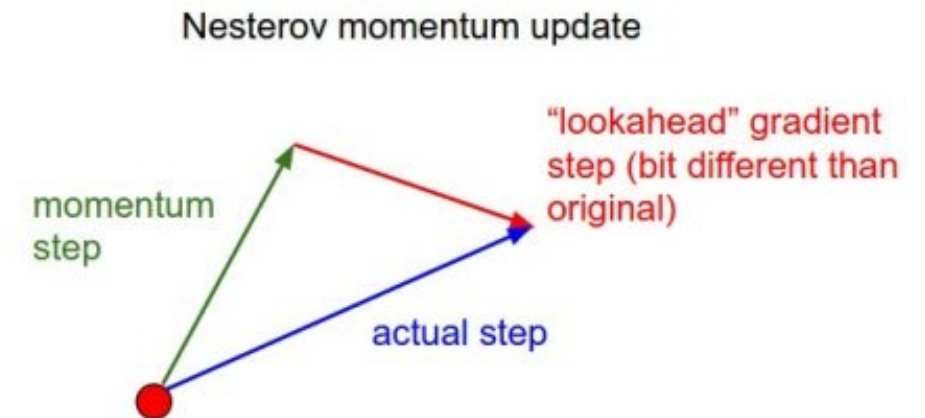
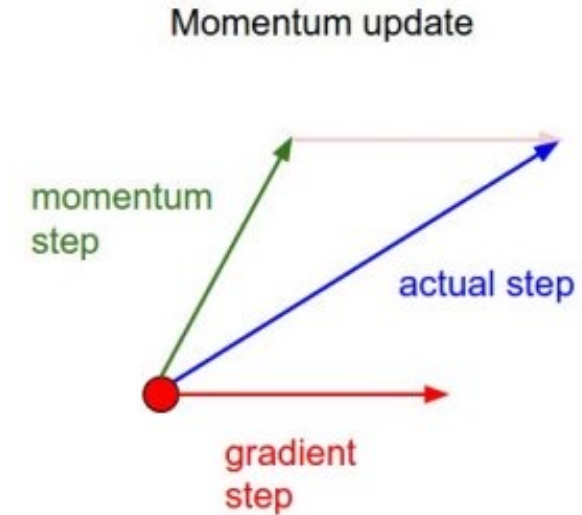
$$w_{t+1} = w_t + (1 + \beta)v_{t+1} - \beta v_t$$

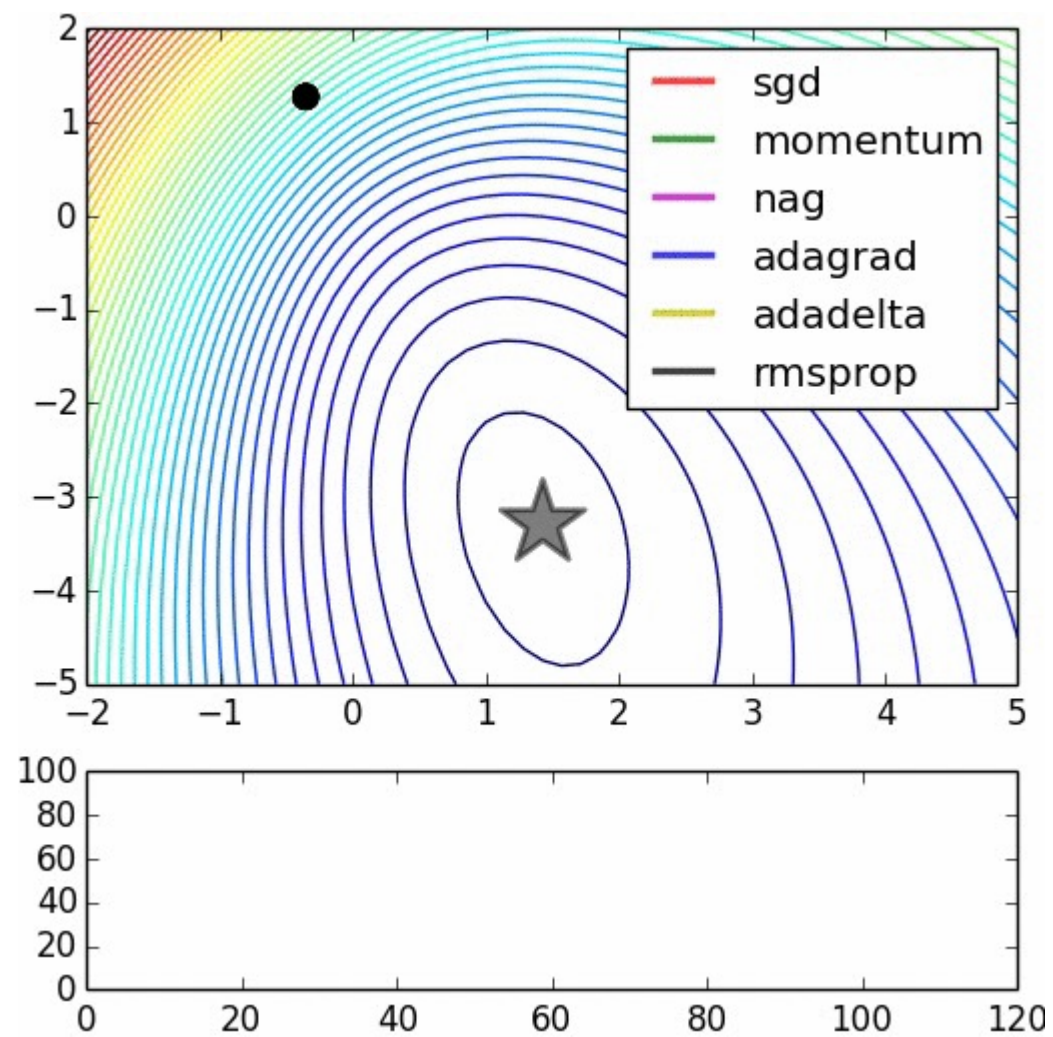
- This is not a descent method!



- Nonetheless, converges faster:

$$f(x_t) - f(x^*) \leq \frac{2\beta \|x_0 - x^*\|}{t^2} = o\left(\frac{\beta}{t^2}\right) \quad \leftarrow \text{matches lower bound!}$$

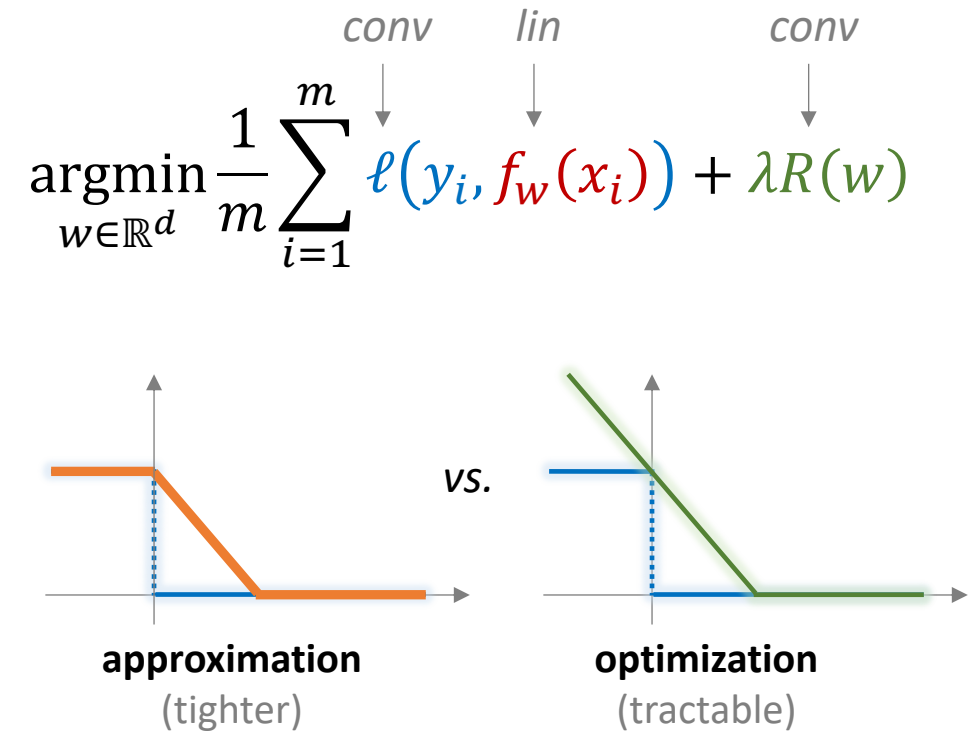




Beyond convexity

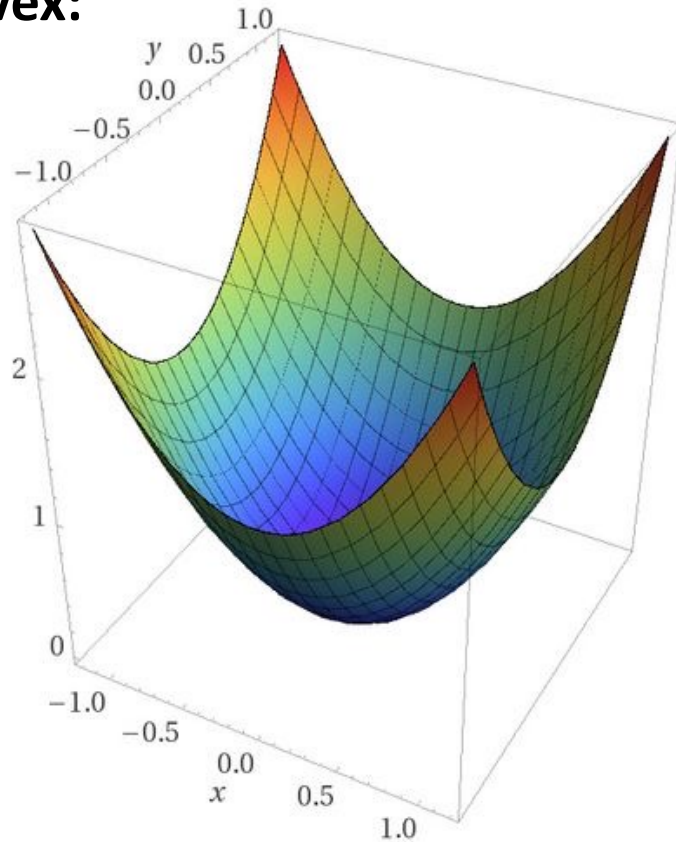
Back to modelling

- **Recall:**
 - Really want to optimize 0/1 loss
 - Instead optimize a continuous proxy
 - Proxies trade off in approximation vs. optimization
 - GD provides strong guarantees for *convex* objectives
 - **But can still be applied to non-convex objectives!**
 - non-convex losses
 - non-convex regularizers
 - non-convex predictive models
- (just need differentiability)
- Finds only local minimum, but many “tricks” for ending up at good local minimum

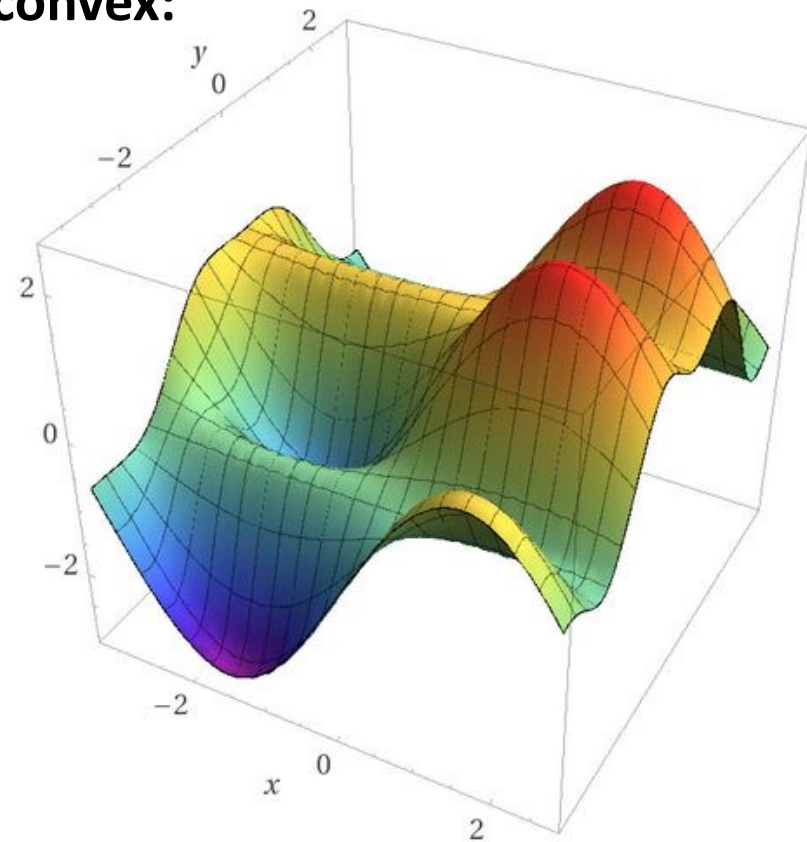


Optimization landscape

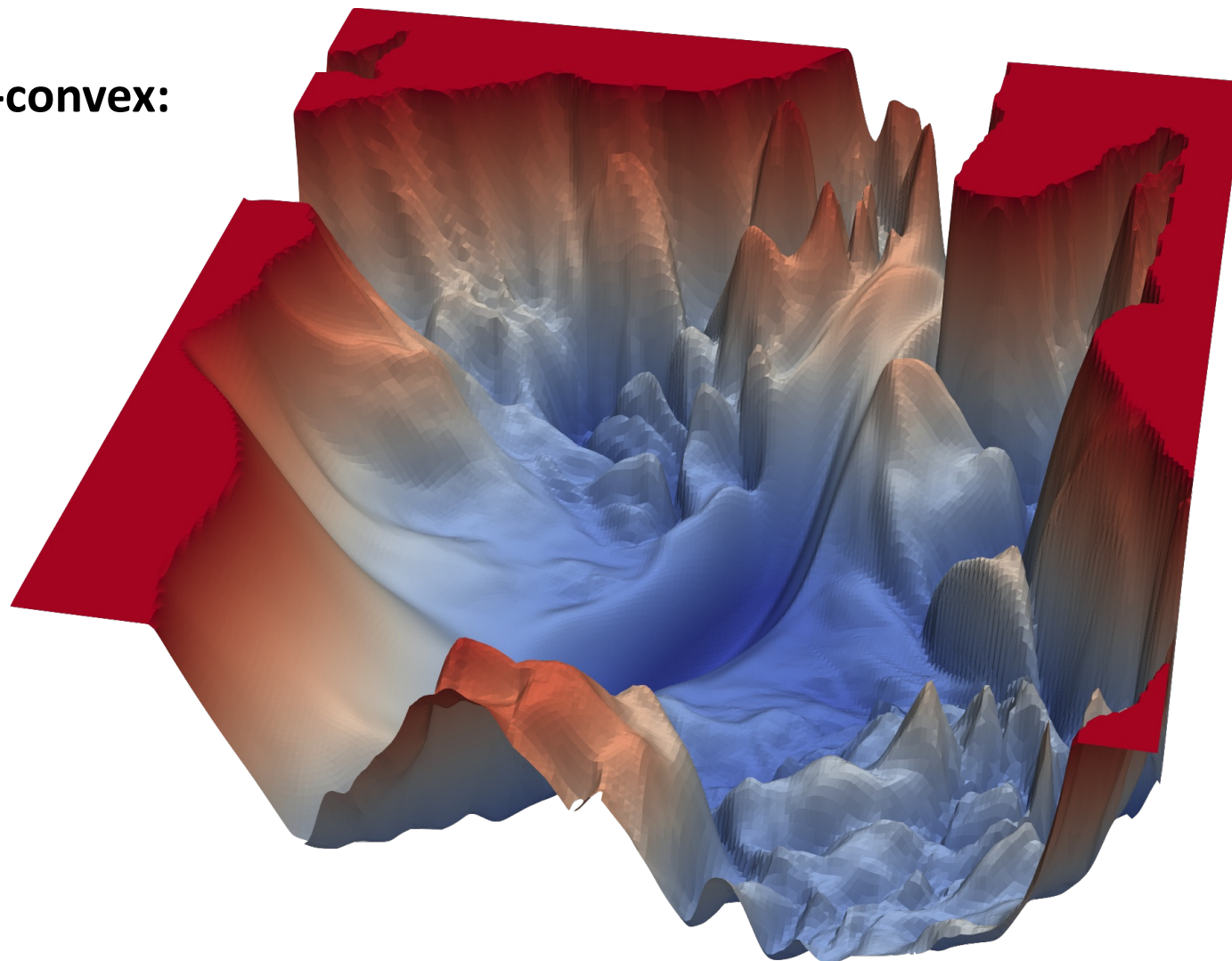
convex:

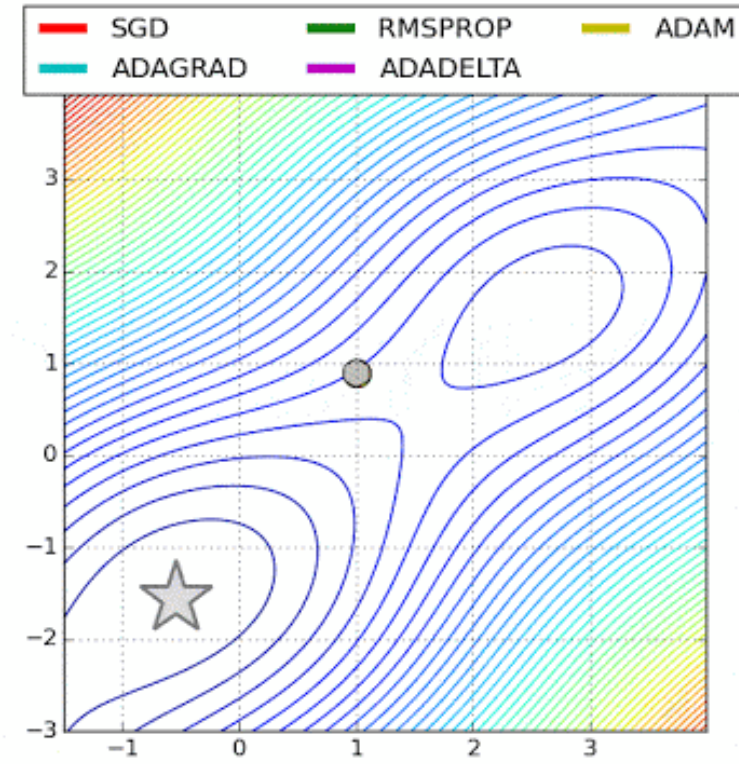
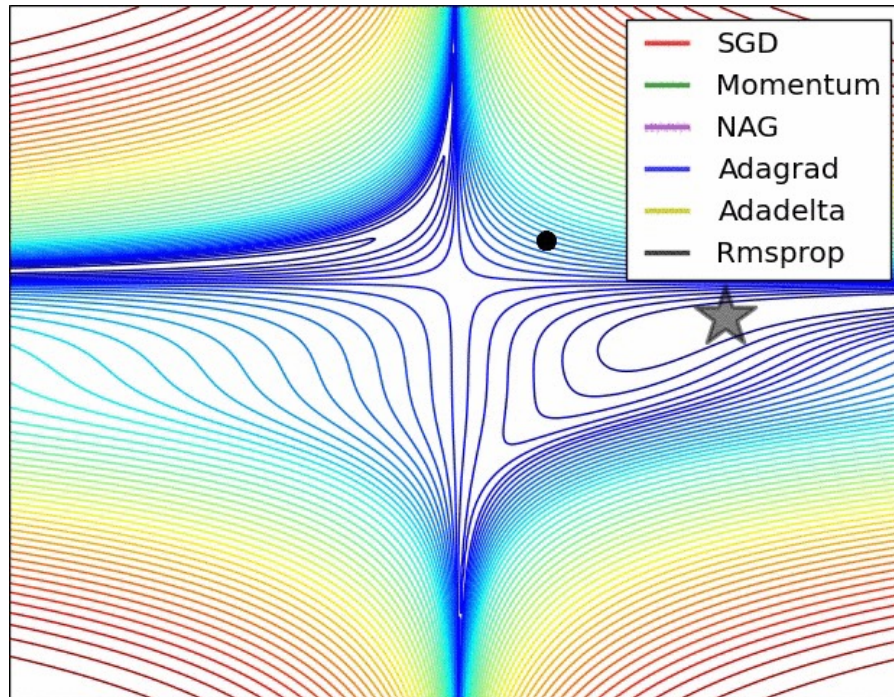


non-convex:

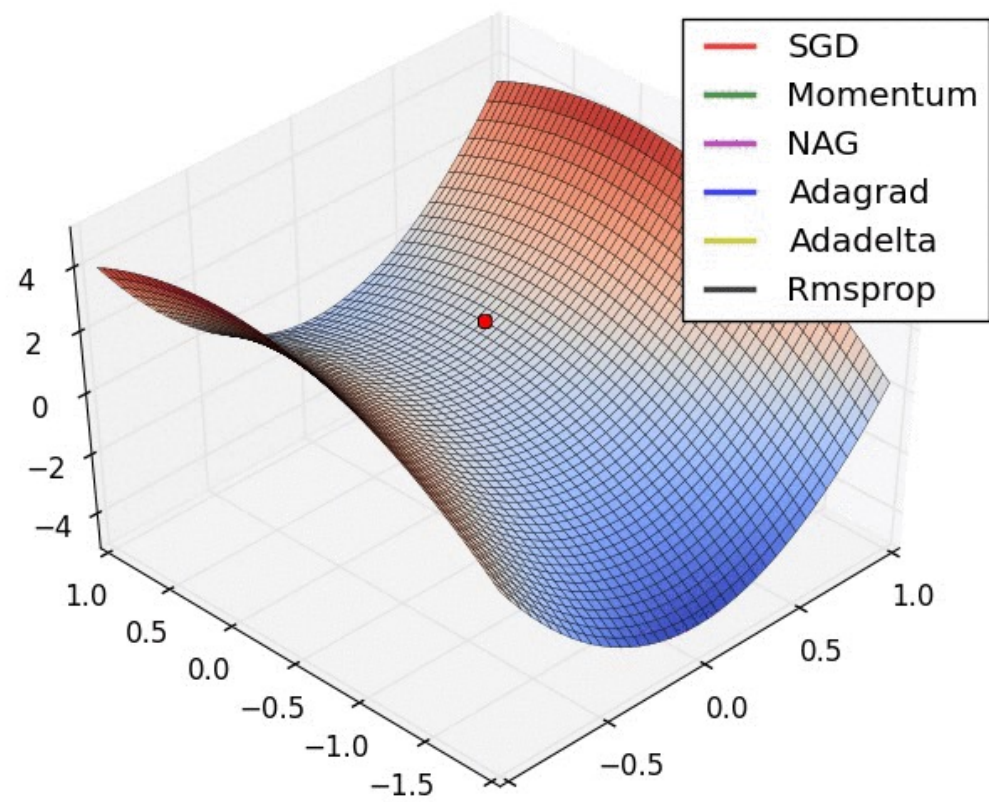
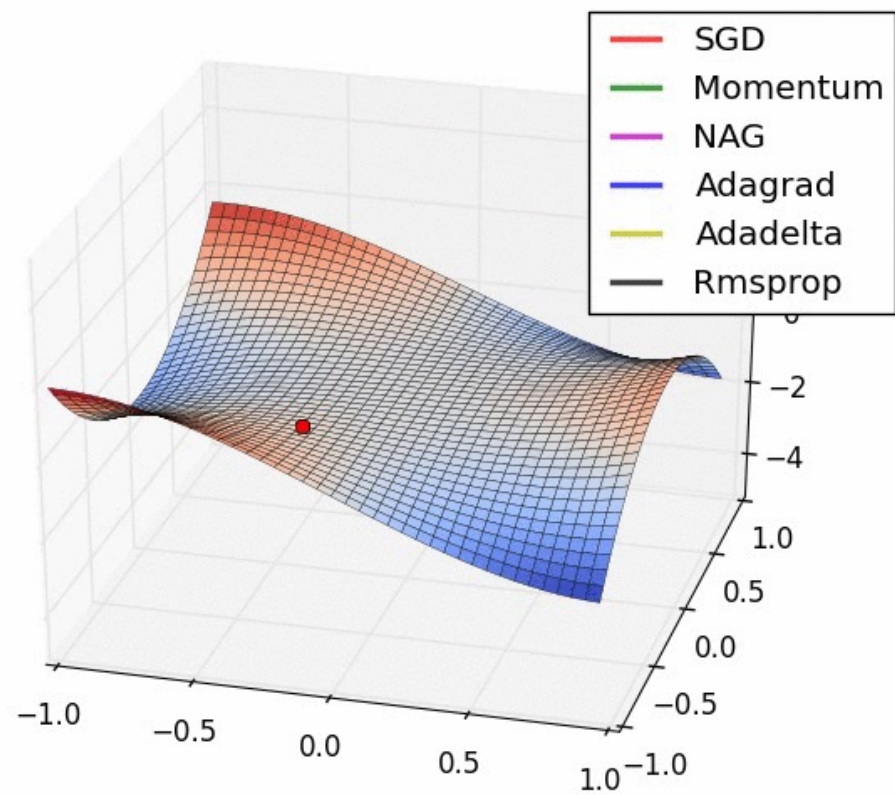


highly non-convex:





Can still apply gradient methods!



Automatic differentiation

- GD requires access to gradients
- **Stone age**: had to compute gradients by hand

$$P_k = \frac{e^{f_k}}{\sum_j e^{f_j}} \quad L_i = -\sum_k p_{i,k} \log P_k \quad f_m = (x_i W)_m$$

$$\text{when } k = m, \quad \frac{\partial P_k}{\partial f_m} = \frac{e^{f_k} \sum_j e^{f_j} - e^{f_k} \cdot e^{f_k}}{(\sum_j e^{f_j})^2} = P_k(1 - P_k)$$

$$\text{when } k \neq m, \quad \frac{\partial P_k}{\partial f_m} = -\frac{e^{f_k} e^{f_m}}{(\sum_j e^{f_j})^2} = -P_k P_m$$

then:

$$\begin{aligned} \frac{\partial L_i}{\partial f_m} &= -\sum_k p_{i,k} \frac{\partial \log P_k}{\partial f_m} \\ &= -\sum_k p_{i,k} \frac{1}{P_k} \frac{\partial P_k}{\partial f_m} \\ &= -\sum_{k=m} p_{i,k} \frac{1}{P_k} P_k(1 - P_k) + \sum_{k \neq m} p_{i,k} \frac{1}{P_k} P_k P_m \\ &= \sum_{k \neq m} p_{i,k} P_m - \sum_{k=m} p_{i,k} (1 - P_k) \\ &= \begin{cases} P_m & , \quad m \neq y_i \\ P_m - 1 & , \quad m = y_i \end{cases} \\ &= P_m - p_{i,m} \end{aligned}$$

Last:

$$\frac{\partial L_i}{\partial W_k} = \frac{\partial L_i}{\partial f_m} \frac{\partial f_m}{\partial W_k} = x_i^T (P_m - p_{i,m})$$

$$\nabla_{W_k} L = -\frac{1}{N} \sum_i x_i^T (p_{i,m} - P_m) + 2\lambda W_k$$

Automatic differentiation

- GD requires access to gradients
- **Stone age**: had to compute gradients by hand
- **Modern age**: automatic differentiation (AutoDiff)

forward
 $(y, \text{dy/dx}) = \text{foo}(x)$
backward

- Gradient computation completely abstracted away
- Building blocks + composition = **differentiable programs**
- We'll return to this when we discuss deep learning

$$P_k = \frac{e^{f_k}}{\sum_j e^{f_j}} \quad L_i = -\sum_k p_{i,k} \log P_k \quad f_m = (x_i W)_m$$

$$\text{when } k = m, \quad \frac{\partial P_k}{\partial f_m} = \frac{e^{f_k} \sum_j e^{f_j} - e^{f_k} \cdot e^{f_k}}{(\sum_j e^{f_j})^2} = P_k(1 - P_k)$$

$$\text{when } k \neq m, \quad \frac{\partial P_k}{\partial f_m} = -\frac{e^{f_k} e^{f_m}}{(\sum_j e^{f_j})^2} = -P_k P_m$$

then:

$$\begin{aligned} \frac{\partial L_i}{\partial f_m} &= -\sum_k p_{i,k} \frac{\partial \log P_k}{\partial f_m} \\ &= -\sum_k p_{i,k} \frac{1}{P_k} \frac{\partial P_k}{\partial f_m} \\ &= -\sum_{k=m} p_{i,k} \frac{1}{P_k} P_k(1 - P_k) + \sum_{k \neq m} p_{i,k} \frac{1}{P_k} P_k P_m \\ &= \sum_{k \neq m} p_{i,k} P_m - \sum_{k=m} p_{i,k} (1 - P_k) \\ &= \begin{cases} P_m & , \quad m \neq y_i \\ P_m - 1 & , \quad m = y_i \end{cases} \\ &= P_m - p_{i,m} \end{aligned}$$

Last:

$$\frac{\partial L_i}{\partial W_k} = \frac{\partial L_i}{\partial f_m} \frac{\partial f_m}{\partial W_k} = x_i^T (P_m - p_{i,m})$$

$$\nabla_{W_k} L = -\frac{1}{N} \sum_i x_i^T (p_{i,m} - P_m) + 2\lambda W_k$$

<http://blog.jasouris.com/2012/02/02/>

Up next

- **Part II:** *the different aspects of learning*
 1. Statistics: generalization and PAC theory
 2. Modeling:
 - Error decomposition
 - Regularization
 - Model selection
 3. Optimization: convexity, gradient descent
 4. Practical aspects and potential pitfalls

Perceptron

input: sample set $S = \{(x_i, y_i)\}_{i=1}^m$

algorithm:

- initialize $w_0 = \vec{0}$
- for $t = 1, 2, \dots$
 - if $\exists i \in [m]$ s.t. $y_i w_t^\top x_i \leq 0$ #wrong classification
 - $w_{t+1} = w_t + y_i x_i$
 - else
 - return w_t

