

Gradient descent : theory and practice

Pierre Ablin

Slides courtesy of Robert Gower



Machine learning task

Finite Sum Training Problem

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n f_i(w) =: f(w)$$

Today: assume that f is differentiable and L -smooth

! ∇f exists

Iterative minimization

Finite Sum Training Problem

$$\min_{w \in \mathbb{R}^d} \frac{1}{n} \sum_{i=1}^n f_i(w) =: f(w)$$

Usually cannot solve this in closed form : $w^* = \dots$

Idea: start from initial guess w^0 and try to find a new, better point. Iterative process $w^0 \rightarrow w^1 \rightarrow \dots$

Gradient descent : basic idea

Given w^0 , look for w^1 as $w^1 = w^0 + d$ where d is a small displacement.

Ideally: $d \in \arg \min_{d \in \mathbb{R}^p} f(w^0 + d)$

Just as hard as the original problem :(

Solution: $d \in \arg \min_{\|d\| \leq \varepsilon} f(w^0 + d)$

Q: as ε goes to 0, what is the limit of d ?

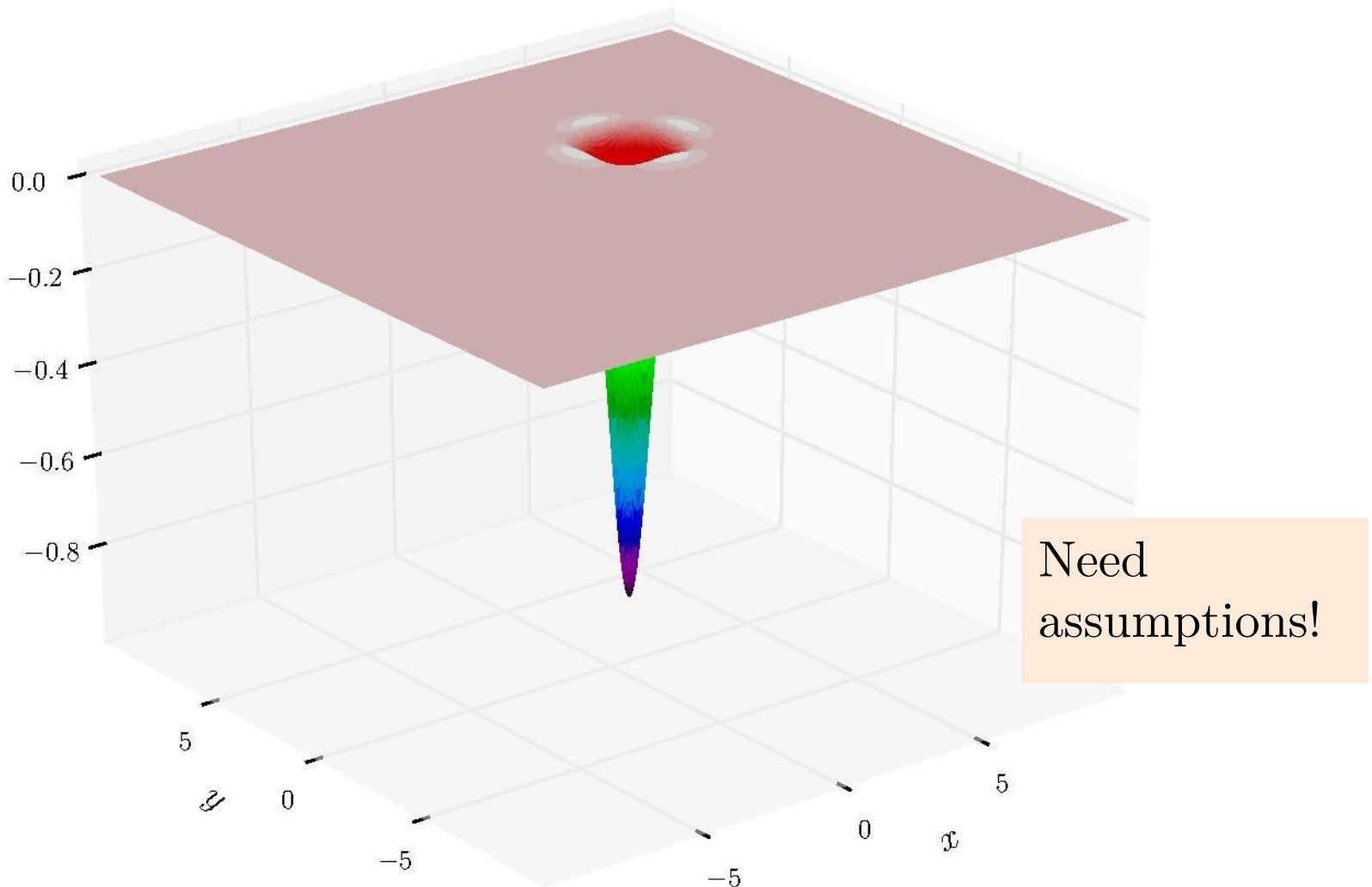
Gradient descent algorithm

Init : Select initial guess w^0
For $t = 0, 1, \dots, T$:
 - Select a step size $\rho^t \in \mathbb{R}^+$
 - Update $w^{t+1} = w^t - \rho^t \nabla f(w^t)$
Return : w^{T+1}

Questions :

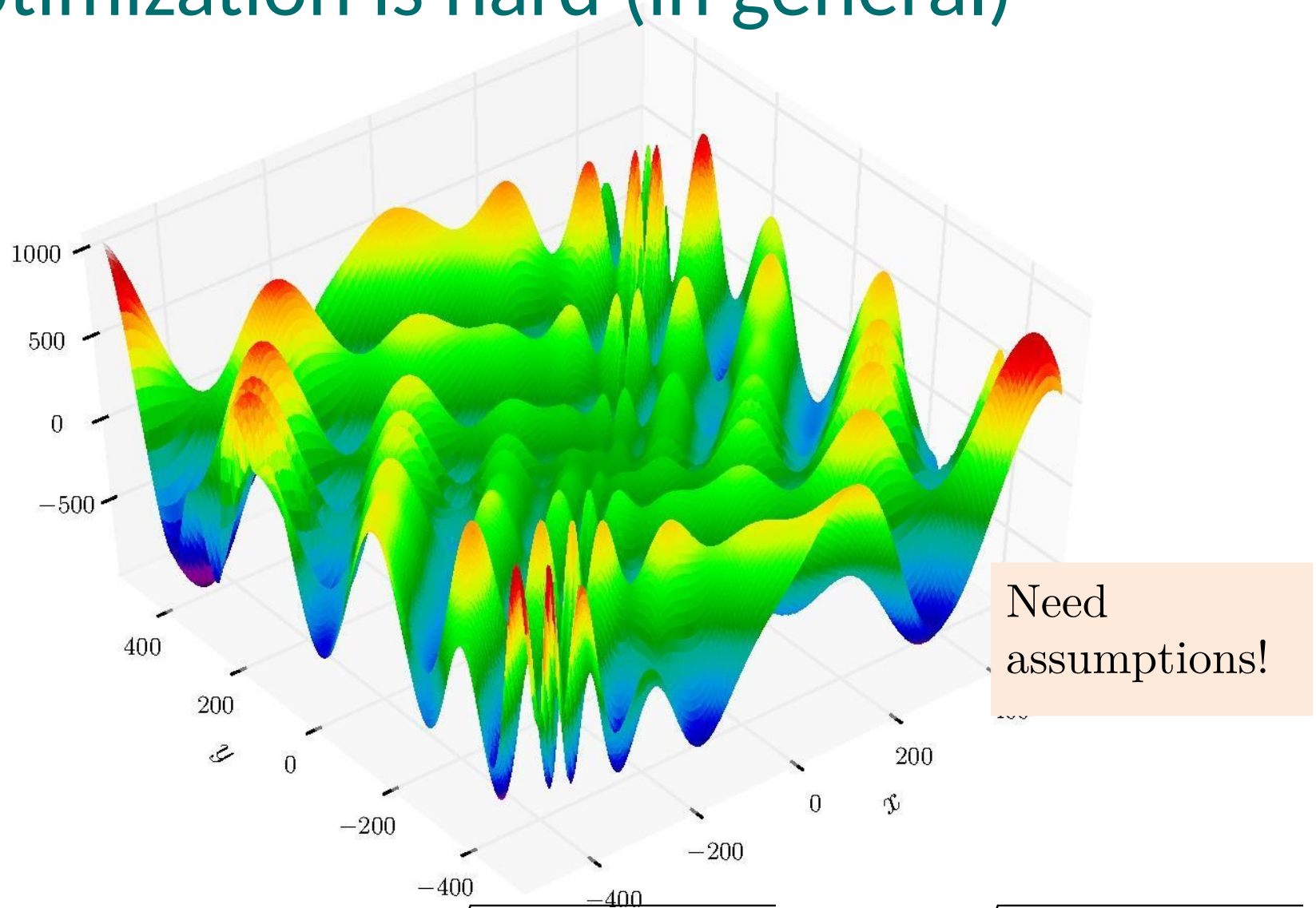
- Does it converge? In which sense?
- At which speed?
- Choice of ρ^t ?

Optimization is hard (in general)



$$f(x, y) = -\cos(x) \cos(y) \exp \left(-(x - \pi)^2 - (y - \pi)^2 \right)$$

Optimization is hard (in general)



$$f(x, y) = -(y + 47) \sin \sqrt{\left| \frac{x}{2} + (y + 47) \right|} - x \sin \sqrt{\left| \frac{x}{2} - (y + 47) \right|}$$

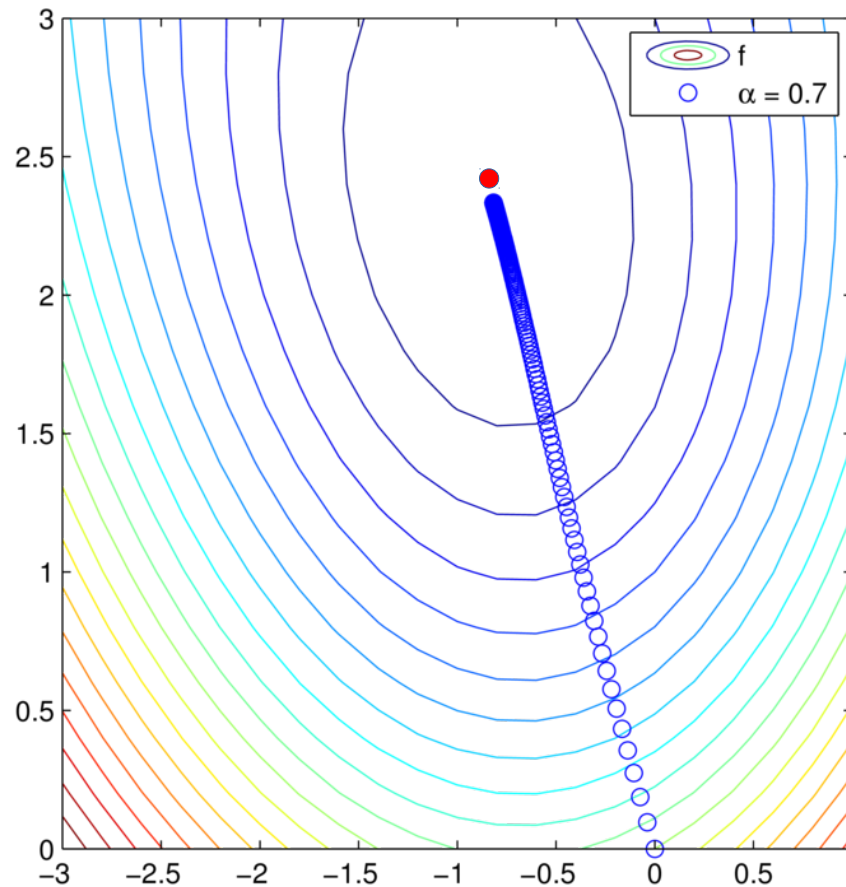
Gradient Descent Example

A Logistic Regression problem using the fourclass labelled data from LIBSVM

$(n, d) = (862, 2)$

Logistic Regression

$$\min_{w \in \mathbf{R}^d} \frac{1}{n} \sum_{i=1}^n \ln(1 + e^{-y^i \langle w, x^i \rangle}) + \lambda \|w\|_2^2$$



Can we prove that this always works?

Gradient Descent Example

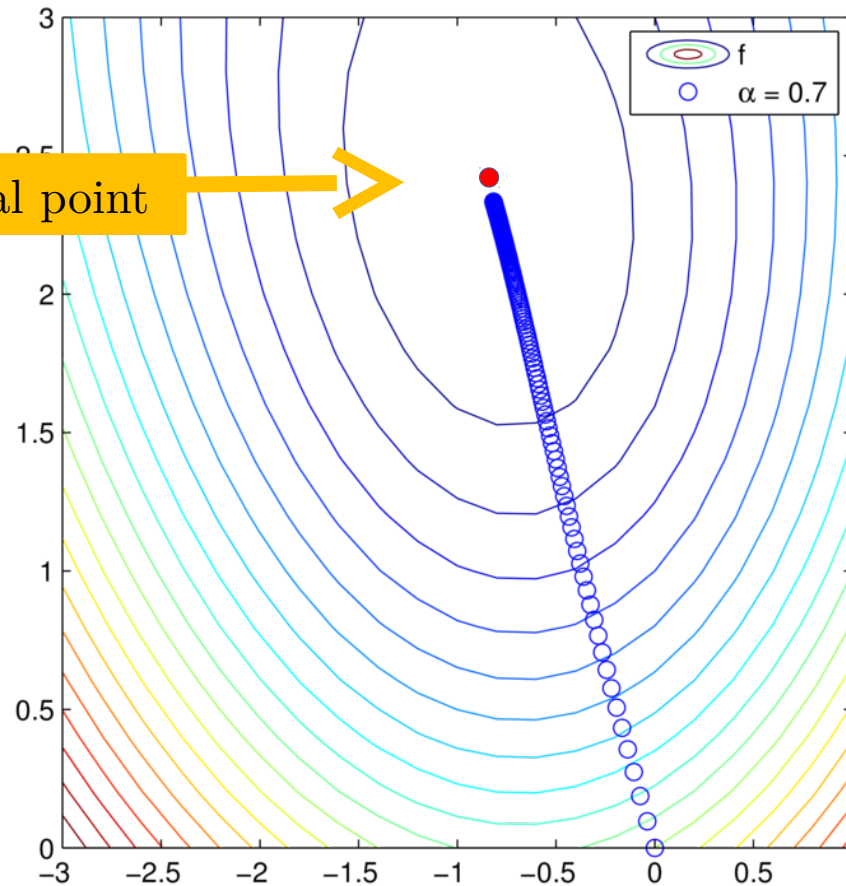
Optimal point

A Logistic Regression problem using the fourclass labelled data from LIBSVM

$(n, d) = (862, 2)$

Logistic Regression

$$\min_{w \in \mathbf{R}^d} \frac{1}{n} \sum_{i=1}^n \ln(1 + e^{-y^i \langle w, x^i \rangle}) + \lambda \|w\|_2^2$$



Can we prove that this always works?

Gradient Descent Example

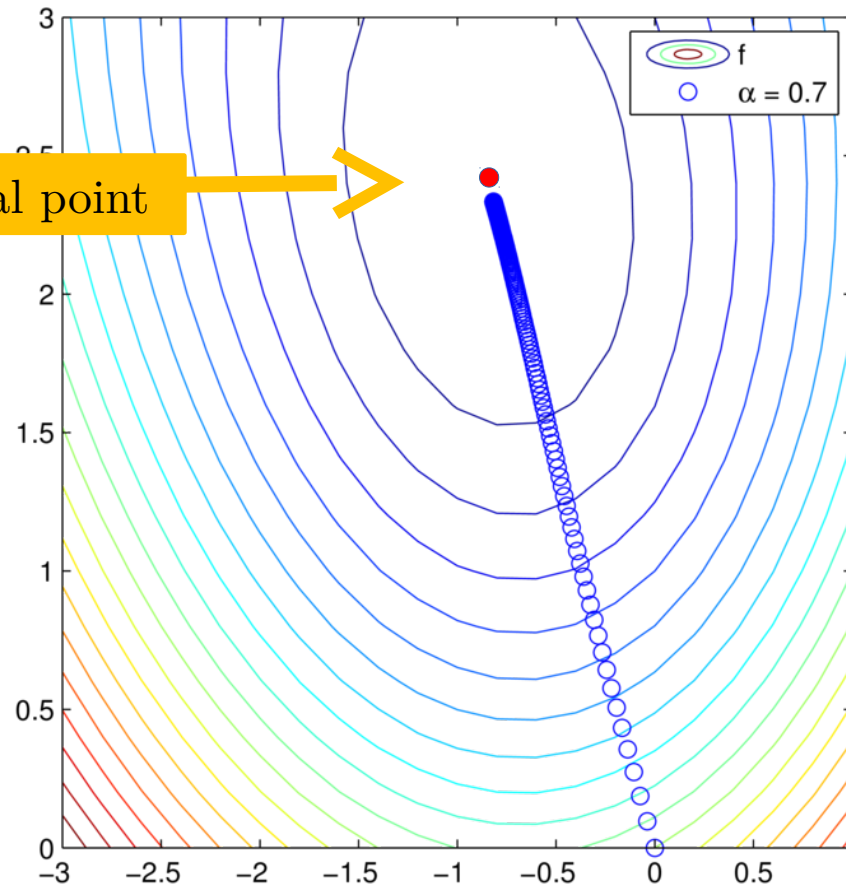
Optimal point

A Logistic Regression problem using the fourclass labelled data from LIBSVM

$(n, d) = (862, 2)$

Logistic Regression

$$\min_{w \in \mathbf{R}^d} \frac{1}{n} \sum_{i=1}^n \ln(1 + e^{-y^i \langle w, x^i \rangle}) + \lambda \|w\|_2^2$$



Can we prove that this always works?

No! There is no universal optimization method. The “no free lunch” of Optimization

Gradient Descent Example

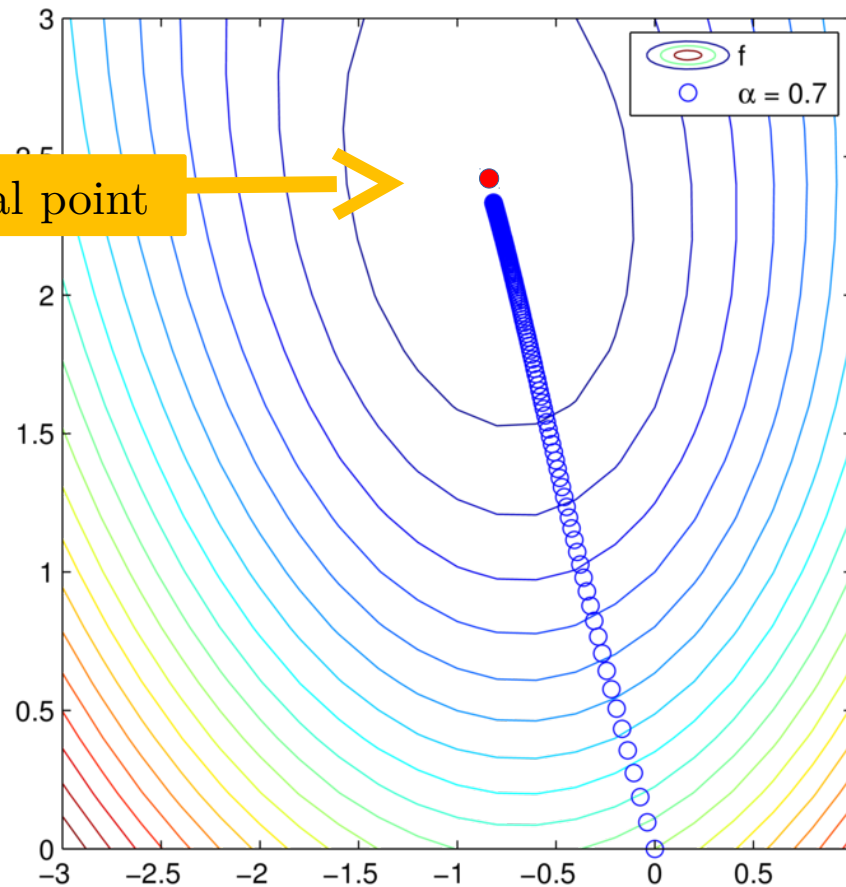
Optimal point

A Logistic Regression problem using the fourclass labelled data from LIBSVM

$(n, d) = (862, 2)$

Logistic Regression

$$\min_{w \in \mathbf{R}^d} \frac{1}{n} \sum_{i=1}^n \ln(1 + e^{-y^i \langle w, x^i \rangle}) + \lambda \|w\|_2^2$$



Can we prove that this always works?

No! There is no universal optimization method. The “no free lunch” of Optimization

Specialize

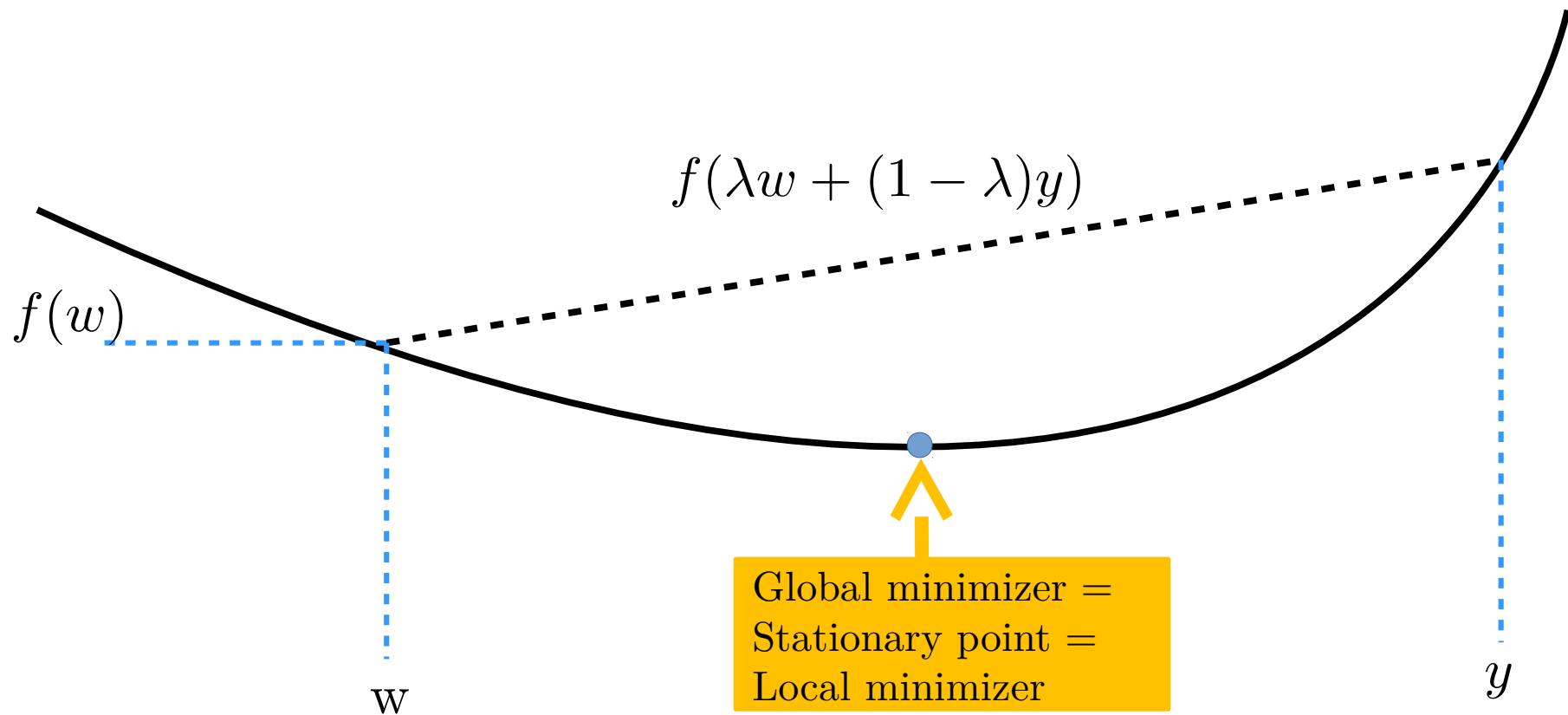


Convex and smooth training problems

Convexity

We say $f : \text{dom}(f) \subset \mathbb{R}^p \rightarrow \mathbb{R}$ is convex if $\text{dom}(f)$ is convex and

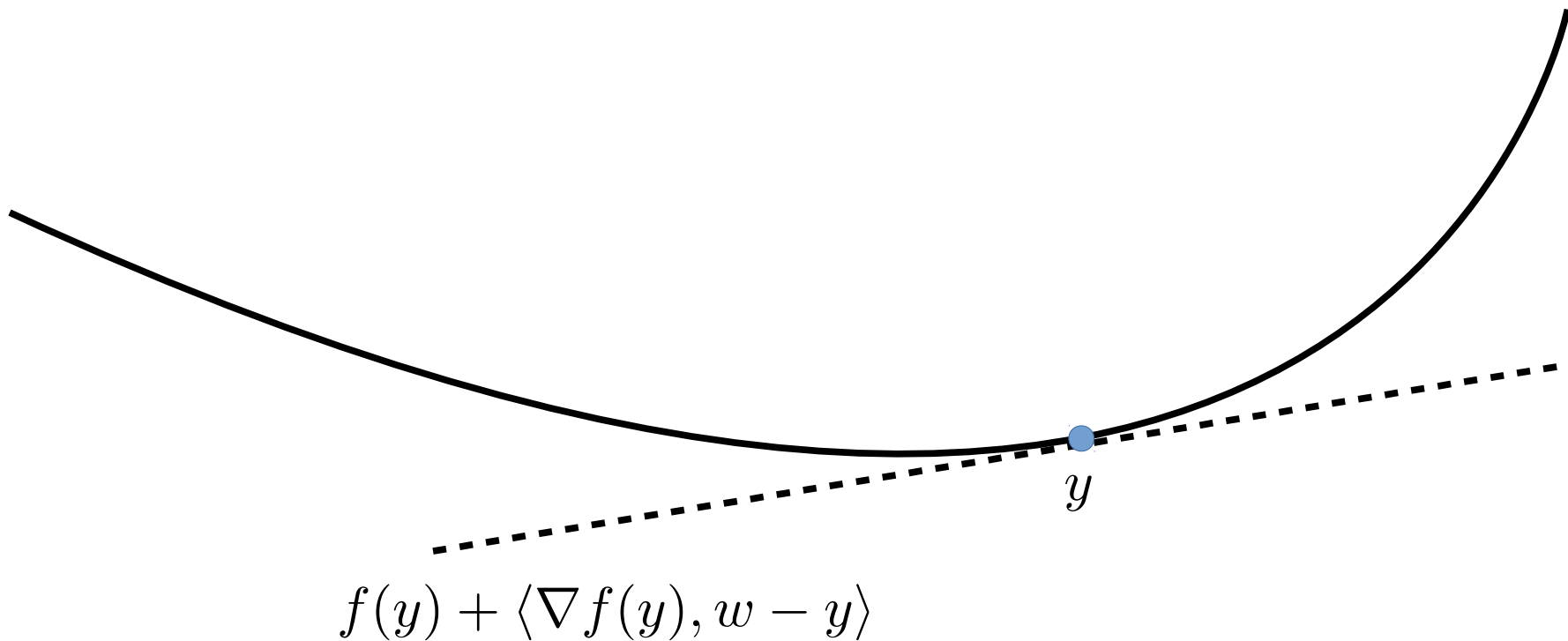
$$f(\lambda w + (1 - \lambda)y) \leq \lambda f(w) + (1 - \lambda)f(y), \quad \forall w, y \in C, \lambda \in [0, 1]$$



Convexity

A differentiable function $f : \text{dom}(f) \subset \mathbb{R}^p \rightarrow \mathbb{R}$ is convex iff

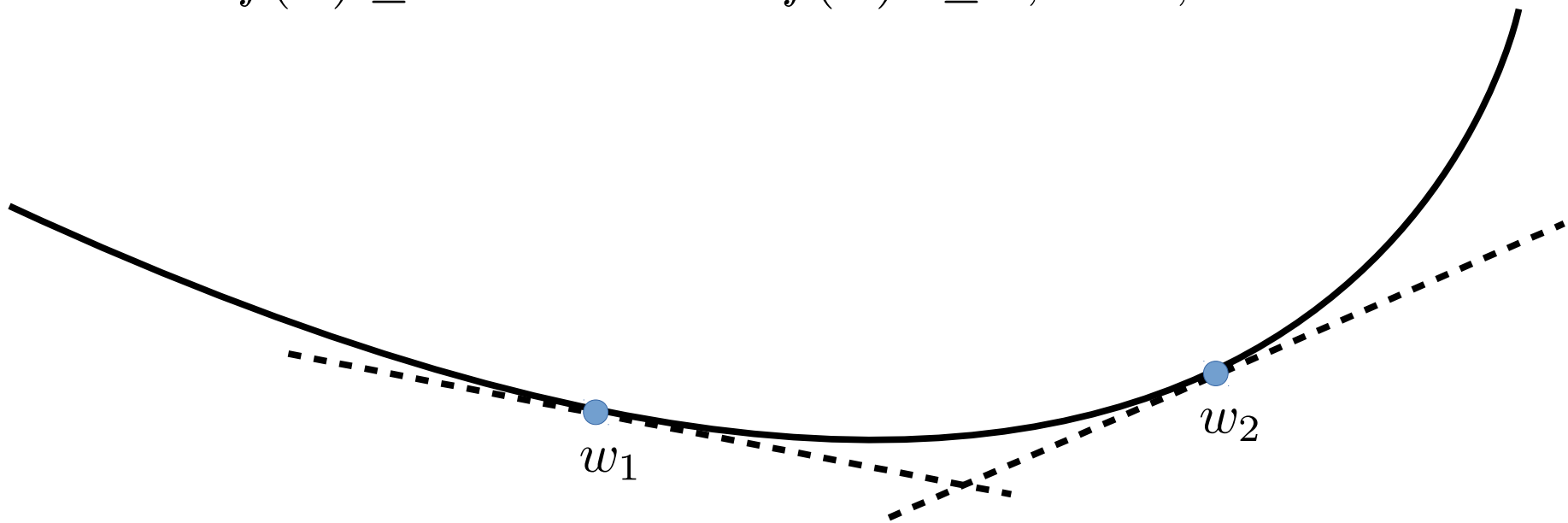
$$f(w) \geq f(y) + \langle \nabla f(y), w - y \rangle$$



Convexity

A twice differentiable function $f : \text{dom}(f) \subset \mathbb{R}^p \rightarrow \mathbb{R}$ is convex iff

$$\nabla^2 f(w) \succeq 0 \quad \Leftrightarrow \quad v^\top \nabla^2 f(w) v \geq 0, \quad \forall w, v \in \mathbb{R}^p$$



$$w_1 \leq w_2 \quad \Rightarrow \quad f'(w_1) \leq f'(w_2)$$

Main Advantage of Convexity

Nice Property

If $\nabla f(w^*) = 0$ then $f(w^*) \leq f(w), \quad \forall w \in \mathbb{R}^d$

Main Advantage of Convexity

Nice Property

If $\nabla f(w^*) = 0$ then $f(w^*) \leq f(w), \quad \forall w \in \mathbb{R}^d$

All stationary points are
global minima

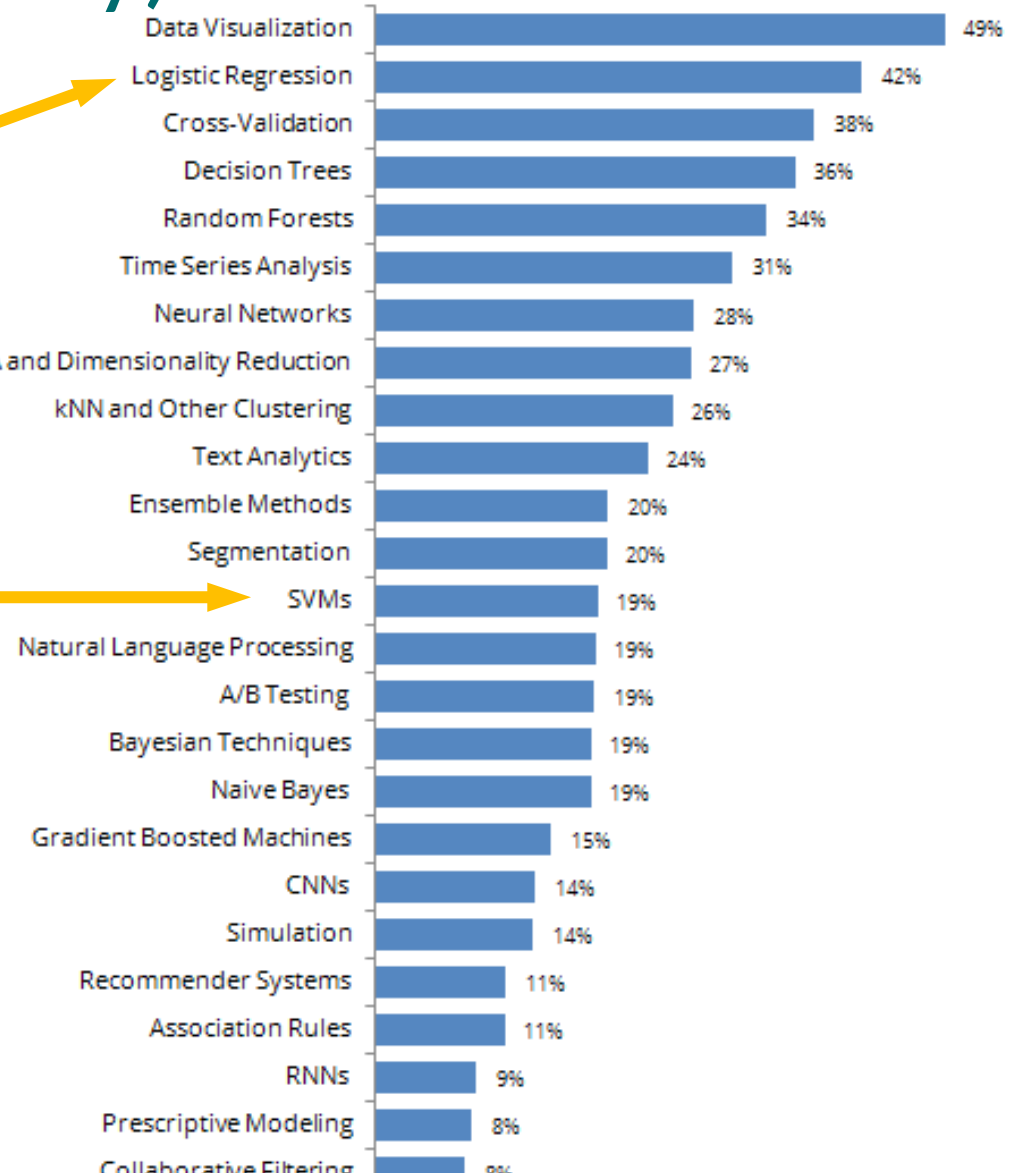
Lemma: Convexity \Rightarrow Nice property

If $f(w) \geq f(y) + \langle \nabla f(y), w - y \rangle, \quad \forall w, y \in \mathbb{R}^d$
then Nice Property holds

PROOF: Choose $y = w^*$

Data science methods most used (Kaggle 2017 survey)

Convex
Optimization
problems



Convexity: Examples

Extended-value extension:

$$f : \mathbb{R}^p \rightarrow \mathbb{R} \cup \{\infty\}$$

$$f(x) = \infty, \quad \forall x \notin \text{dom}(f)$$

Norms and squared norms:

$$x \mapsto \|x\|$$

$$x \mapsto \|x\|^2$$

Negative log and logistic:

$$x \mapsto -\log(x)$$

$$x \mapsto \log \left(1 + e^{-y \langle a, x \rangle} \right)$$

Hinge loss

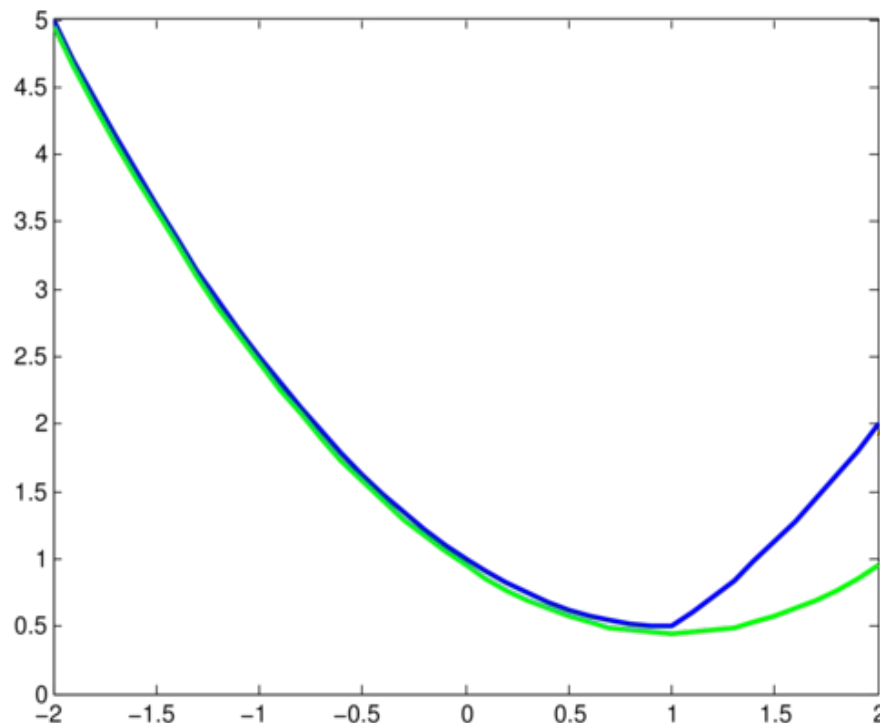
$$x \mapsto \max\{0, 1 - yx\}$$

Negatives log determinant, exponentiation ... etc

Strong convexity

We say $f : \mathbb{R}^p \rightarrow \mathbb{R} \cup \{\infty\}$ is μ -strongly convex if

$$f(w) \geq f(y) + \langle \nabla f(y), w - y \rangle + \frac{\mu}{2} \|w - y\|^2, \quad \forall w, y \in \mathbb{R}^p$$



Hinge loss + L2
 $\max\{0, 1 - w\} + \frac{1}{2} \|w\|_2^2$

Quadratic lower bound

Smoothness

We say $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ is smooth if

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|, \quad \forall x, y \in \mathbb{R}^p$$

Smoothness

We say $f : \mathbb{R}^p \rightarrow \mathbb{R} \cup \{\infty\}$ is smooth if

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|, \quad \forall x, y \in \mathbb{R}^p$$

If a twice differentiable $f : \mathbb{R}^p \rightarrow \mathbb{R} \cup \{\infty\}$ is L -smooth then

$$1) \quad d^\top \nabla^2 f(x) d \leq L \cdot \|d\|_2^2, \quad \forall x, d \in \mathbb{R}^p$$

$$2) \quad f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{L}{2} \|x - y\|^2, \quad \forall x, y \in \mathbb{R}^p$$

Smoothness

We say $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ is smooth if

$$\|\nabla f(x) - \nabla f(y)\| \leq L\|x - y\|, \quad \forall x, y \in \mathbb{R}^n$$

If a twice differentiable $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ is L -smooth then

$$1) \quad d^\top \nabla^2 f(x) d \leq L \cdot \|d\|_2^2, \quad \forall x, d \in \mathbb{R}^n$$

$$2) \quad f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{L}{2} \|x - y\|^2, \quad \forall x, y \in \mathbb{R}^n$$

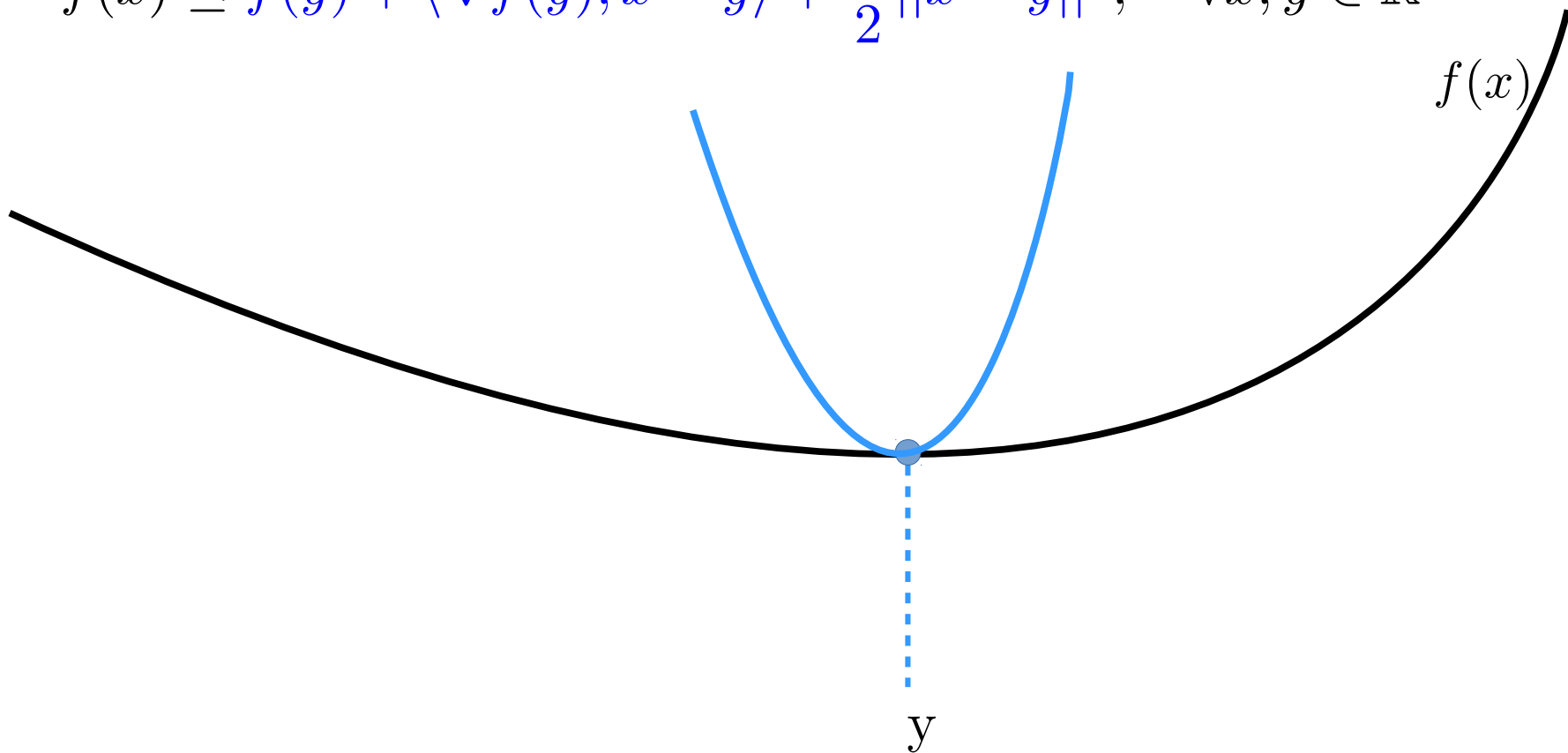
EXE: determine the strong convexity / smoothness constants of

$$f(w) := \frac{1}{2} \|X^\top w - b\|_2^2 \text{ for } X \in \mathbb{R}^{n \times p}, \quad b \in \mathbb{R}^n$$

Important consequences of Smoothness

If $f : \mathbb{R}^n \rightarrow \mathbb{R} \cup \{\infty\}$ is L -smooth then

$$f(x) \leq f(y) + \langle \nabla f(y), x - y \rangle + \frac{L}{2} \|x - y\|^2, \quad \forall x, y \in \mathbb{R}^n$$



Smoothness: Examples

Convex quadratics:

$$x \mapsto x^\top Ax + b^\top x + c$$

Logistic:

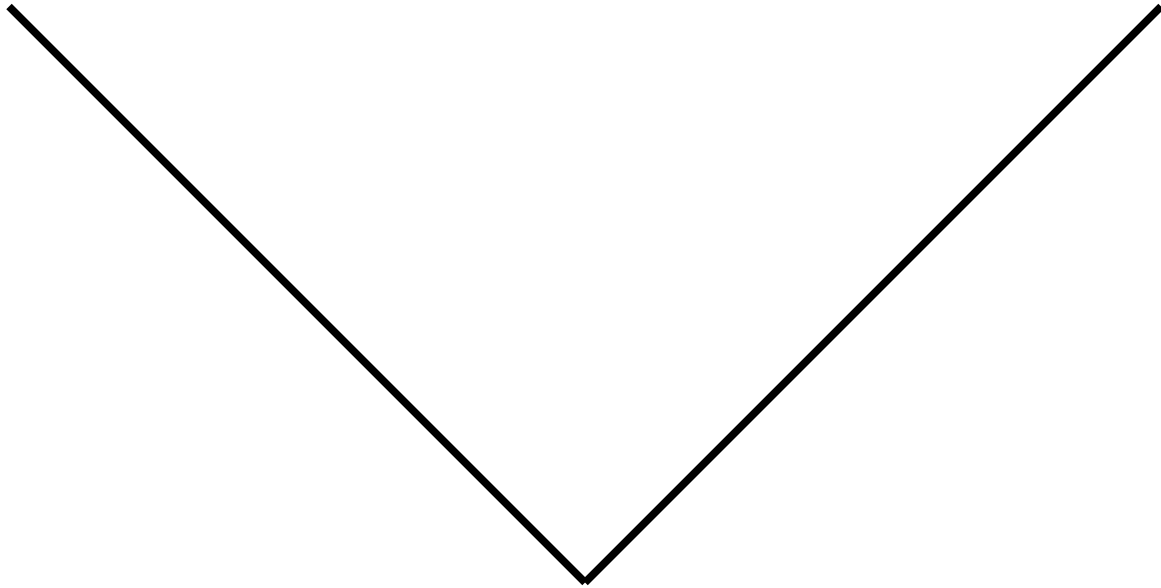
$$x \mapsto \log \left(1 + e^{-y \langle a, x \rangle} \right)$$

Trigonometric:

$$x \mapsto \cos(x), \sin(x)$$

Smoothness: Convex counter-example

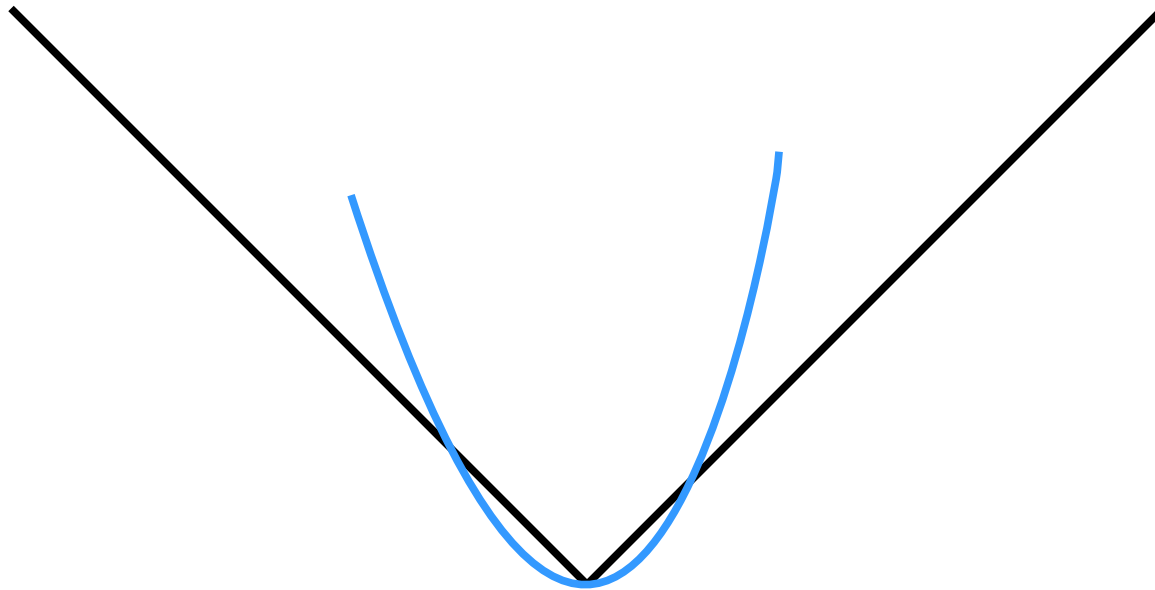
$$f(w) = ||w||_1 = \sum_{i=1}^n |w_i|$$



We'll see how to handle
this problem next class

Smoothness: Convex counter-example

$$f(w) = ||w||_1 = \sum_{i=1}^n |w_i|$$



We'll see how to handle this problem next class

Does not fit
Not smooth

Insight into Gradient Descent using Smoothness

$$f(w) \leq f(w^0) + \langle \nabla f(w^0), w - w^0 \rangle + \frac{L}{2} \|w - w^0\|^2$$

Q: what is the minimizer of the upper bound in w ?

Insight into Gradient Descent using Smoothness

$$f(w) \leq f(w^0) + \langle \nabla f(w^0), w - w^0 \rangle + \frac{L}{2} \|w - w^0\|^2$$

Minimizing the upper bound in w we get:

$$\nabla_w \left(f(w^0) + \langle \nabla f(w^0), w - w^0 \rangle + \frac{L}{2} \|w - w^0\|^2 \right) = \nabla f(w^0) + L(w - w^0)$$



A gradient
descent step !

$$w = w^0 - \frac{1}{L} \nabla f(w^0)$$

Insight into Gradient Descent using Smoothness

$$f(w) \leq f(w^0) + \langle \nabla f(w^0), w - w^0 \rangle + \frac{L}{2} \|w - w^0\|^2$$

Minimizing the upper bound in w we get:

$$\nabla_w \left(f(w^0) + \langle \nabla f(w^0), w - w^0 \rangle + \frac{L}{2} \|w - w^0\|^2 \right) = \nabla f(w^0) + L(w - w^0)$$

Smoothness Lemma (EXE):

If f is L -smooth, show that

$$f\left(y - \frac{1}{L} \nabla f(y)\right) - f(y) \leq -\frac{1}{2L} \|\nabla f(y)\|_2^2, \quad \forall y$$

$$f(w^*) - f(w) \leq -\frac{1}{2L} \|\nabla f(w)\|_2^2, \quad \forall w \in \mathbb{R}^n$$

$$\text{where } f(w^*) \leq f(w), \quad \forall w \in \mathbb{R}^n$$



A gradient
descent step !

$$w = w^0 - \frac{1}{L} \nabla f(w^0)$$

Convergence rates: smooth case

We have
$$f\left(w - \frac{1}{L} \nabla f(w)\right) - f(w) \leq -\frac{1}{2L} \|\nabla f(w)\|^2$$

Gradient descent with step $\rho^t = \frac{1}{L}$:

$$f(w^{t+1}) - f(w^t) \leq -\frac{1}{2L} \|\nabla f(w^t)\|^2$$

$$\sum_{t=0}^T \|\nabla f(w^t)\|^2 \leq 2L(f(w^0) - f(w^*)), \quad \forall T > 0$$

Q: what does it mean?

Convergence rates: smooth case

Theorem : if f is L -smooth, the iterates of gradient descent verify

$$\nabla f(w^t) \rightarrow 0$$

$$\inf_{t \leq T} \|\nabla f(w^t)\|^2 \leq \frac{2L}{T} (f(w^0) - f(w^*))$$



Convergence speed

Slow convergence

Say $2L(f(w^0) - f(w^*)) = 1$

In order to have $\inf_{t \leq T} \|\nabla f(w^t)\|^2 \leq 10^{-4}$

Need 10^4 iterations...

Convergence GD strongly convex

Theorem

Let f be μ -strongly convex and L -smooth.

$$\|w^t - w^*\|_2^2 \leq \left(1 - \frac{\mu}{L}\right)^t \|w^1 - w^*\|_2^2$$

Where

$$w^{t+1} = w^t - \frac{1}{L} \nabla f(w^t), \quad \text{for } t = 1, \dots, T$$

$$\Rightarrow \text{for } \frac{\|w^T - w^*\|_2^2}{\|w^1 - w^*\|_2^2} \leq \epsilon \text{ we need } T \geq \frac{L}{\mu} \log \left(\frac{1}{\epsilon} \right) = O \left(\log \left(\frac{1}{\epsilon} \right) \right)$$

Convergence GD strongly convex

Theorem

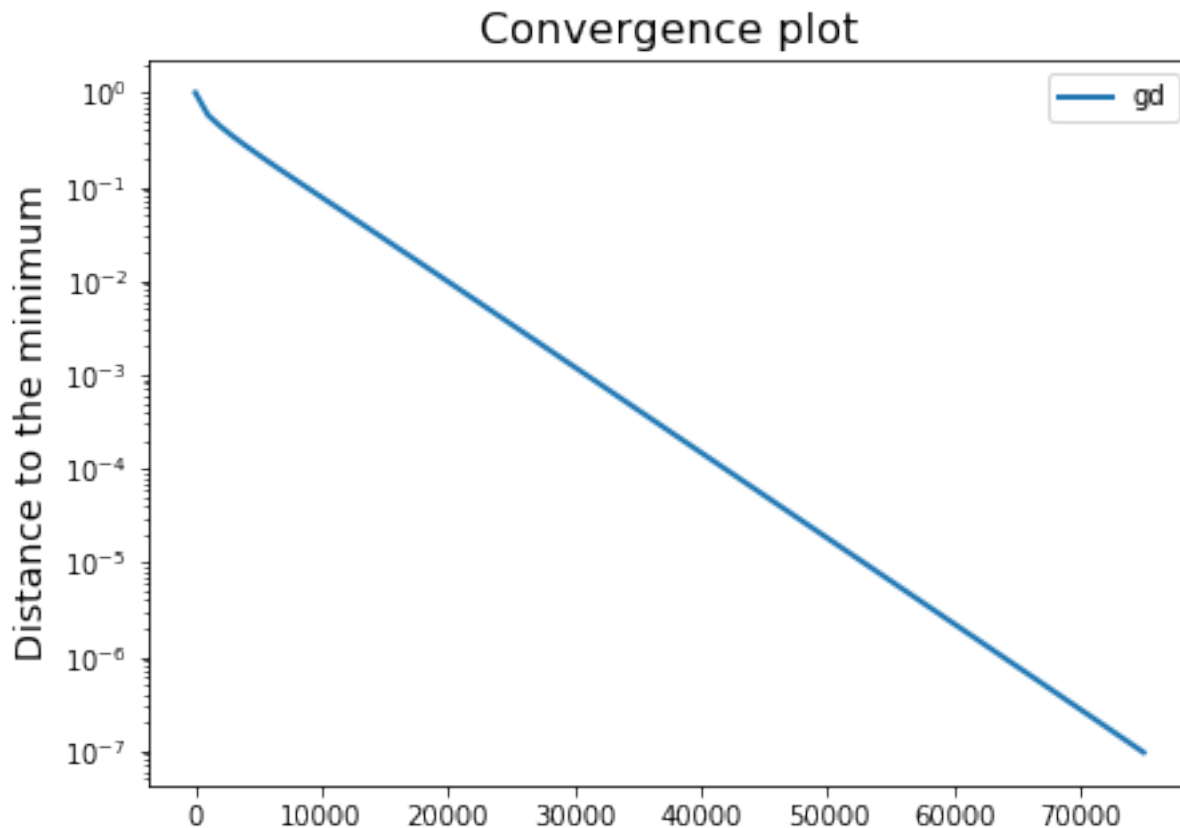
Let f be μ -strongly convex and L -smooth.

$$f(w^t) - f(w^*) \leq \left(1 - \frac{\mu}{L}\right)^t (f(w^0) - f(w^*))$$

Where

$$w^{t+1} = w^t - \frac{1}{L} \nabla f(w^t), \quad \text{for } t = 1, \dots, T$$

Gradient Descent Example: logistic regression



$$y\text{-axis} = \frac{\|w^t - w^*\|_2^2}{\|w^1 - w^*\|_2^2} \quad \longrightarrow \quad \log \left(\frac{\|w^t - w^*\|_2^2}{\|w^1 - w^*\|_2^2} \right) \leq t \log \left(1 - \frac{\mu}{L} \right)$$

Proof Convergence GD strongly convex + smooth

Smoothness



Proof:

$$\begin{aligned} f(w^{t+1}) &\leq f(w^t) + \langle \nabla f(w^t), w^{t+1} - w^t \rangle + \frac{L}{2} \|w^{t+1} - w^t\|^2 \\ &= f(w^t) - \frac{1}{2L} \|\nabla f(w^t)\|^2 \end{aligned}$$

Proof Convergence GD strongly convex + smooth

Smoothness



Proof:

$$\begin{aligned} f(w^{t+1}) &\leq f(w^t) + \langle \nabla f(w^t), w^{t+1} - w^t \rangle + \frac{L}{2} \|w^{t+1} - w^t\|^2 \\ &= f(w^t) - \frac{1}{2L} \|\nabla f(w^t)\|^2 \end{aligned}$$

Polyak-Lojasiewicz (PL) inequality :

Q: show that strong convexity \Rightarrow PL

$$\|\nabla f(w)\|^2 \geq 2\mu(f(w) - f(w^*)), \quad \forall w$$

Proof Convergence GD strongly convex + smooth

Smoothness




Proof:

$$\begin{aligned} f(w^{t+1}) &\leq f(w^t) + \langle \nabla f(w^t), w^{t+1} - w^t \rangle + \frac{L}{2} \|w^{t+1} - w^t\|^2 \\ &= f(w^t) - \frac{1}{2L} \|\nabla f(w^t)\|^2 \end{aligned}$$

Polyak-Lojasiewicz (PL) inequality :

$$\|\nabla f(w)\|^2 \geq 2\mu(f(w) - f(w^*)), \quad \forall w$$

$$f(w^{t+1}) \leq f(w^t) - \frac{\mu}{L} (f(w^t) - f(w^*))$$


$$f(w^t) - f(w^*) \leq \left(1 - \frac{\mu}{L}\right)^t (f(w^0) - f(w^*))$$

Examples of smooth machine learning problems

Least squares

Data: $x_1, \dots, x_n \in \mathbb{R}^p$, and $y_1, \dots, y_n \in \mathbb{R}$

Assumption: There exists w^* such that

$$y_i \simeq \langle x_i, w^* \rangle$$

Optimization problem: $\min_w f(w) = \frac{1}{n} \sum_{i=1}^n (\langle x_i, w \rangle - y_i)^2$

Q: show that we can rewrite $f(w) = \frac{1}{n} \|Xw - y\|^2$

Is the problem convex, smooth? Compute the associated constants

Ridge regression

Problem :

$$\min_w f(w) = \frac{1}{n} \sum_{i=1}^n (\langle x_i, w \rangle - y_i)^2$$

Has infinitely many solutions when $n < p$. Bad conditioning, and very sensitive to X .

Solution : regularize !

$$\min_w f(w) = \frac{1}{n} \sum_{i=1}^n (\langle x_i, w \rangle - y_i)^2 + \frac{\lambda}{2} \|w\|^2$$

Q: Is the problem convex, smooth? Compute the associated constants

Logistic regression

Data: $x_1, \dots, x_n \in \mathbb{R}^p$, and $y_1, \dots, y_n \in \{-1, +1\}$

Assumption: There exists w^* such that

$$y_i \simeq \text{sign}(\langle x_i, w^* \rangle)$$

Optimization problem:

$$\min_w f(w) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i \langle x_i, w \rangle))$$

Q: Is the problem convex, smooth? Compute the associated constants

Regularized logistic regression

Data: $x_1, \dots, x_n \in \mathbb{R}^p$, and $y_1, \dots, y_n \in \{-1, +1\}$

Assumption: There exists w^* such that

$$y_i \simeq \text{sign}(\langle x_i, w^* \rangle)$$

Optimization problem:

$$\min_w f(w) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-y_i \langle x_i, w \rangle)) + \frac{\lambda}{2} \|w\|^2$$

Q: Is the problem convex, smooth? Compute the associated constants