

Convex Optimization

① Convexification of the 0-1 loss

A problem for classification

$$\hat{f}_S \in \operatorname{argmin}_{f \in \mathcal{F}} \left\{ R_n(f) = \frac{1}{n} \sum_{i=1}^n \mathbb{1}\{f(x_i) \neq y_i\} \right\}$$

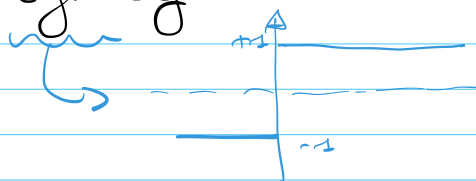
→ How to compute this?

The function $f \mapsto R_n(f)$ is discontinuous and takes $\mathbb{N}_n(x_1, \dots, x_n)$ values!

$$\binom{n}{VCC(F)} = \text{large} \quad (\text{ex: if } VCC(F) > 3)$$

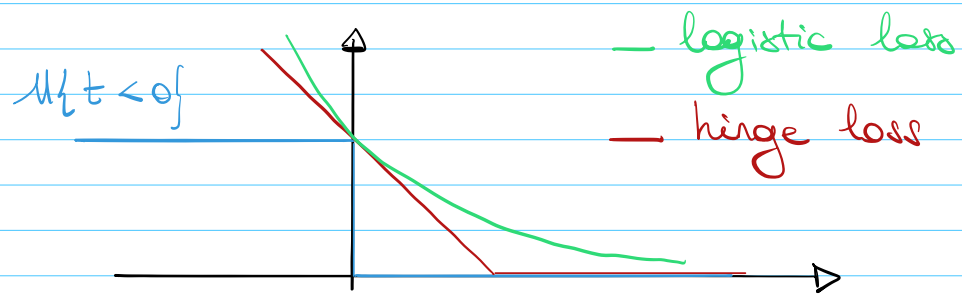
Alternative point of view:

Look for $f = \operatorname{sgn} \circ g$



$$\begin{aligned}
 \leadsto R_p(f) &= \mathbb{E}_p[\mathbb{1}\{f(x) \neq y\}] \\
 &= \mathbb{E}_p[\mathbb{1}\{\text{sgn}(g(x)) \neq y\}] \\
 &= \mathbb{E}_p[\mathbb{1}\{g(x)y < 0\}]
 \end{aligned}$$

\leadsto Still discontinuous ...



\leadsto Replace the 0-1 loss by a convex surrogate.

\hookrightarrow convex functions can be optimized efficiently.

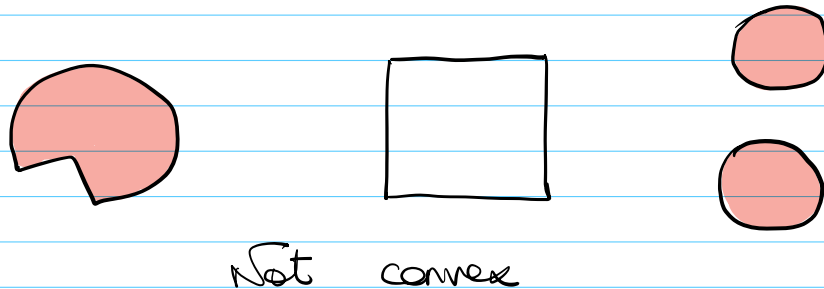
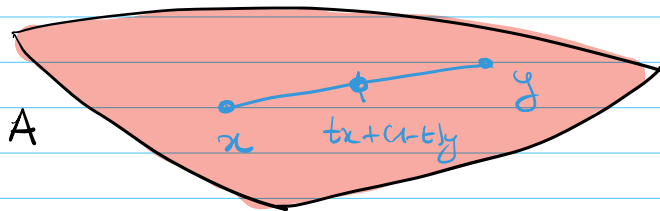
Logistic
Regression
(later today)

Support Vector
Machines
(Ch. 3)

② Convex functions:

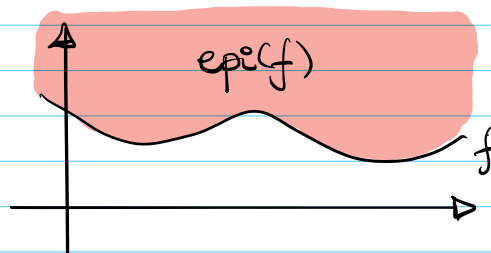
• **Convex Sets:** $A \subset \mathbb{R}^d$

$$\forall x, y \in A, t \in [0, 1], tx + (1-t)y \in A$$



• **Epigraph** of $f: A \rightarrow \mathbb{R}$

$$\text{epi}(f) = \left\{ (x, y) \in A \times \mathbb{R} : y \geq f(x) \right\} \subseteq \mathbb{R}^{d+1}$$



• **Convex function:** A convex set
 $f: A \rightarrow \mathbb{R}$

f is convex if $\text{epi}(f)$ is convex.

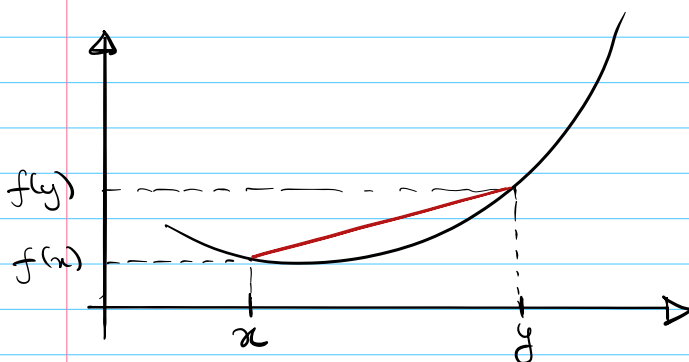
\Rightarrow We have $(x_1, f(x_1)), (x_2, f(x_2)) \in \text{epi}(f)$,

So $t(x_1, f(x_1)) + (1-t)(x_2, f(x_2)) \in \text{epi}(f)$

$\Rightarrow f(tx_1 + (1-t)x_2) \leq tf(x_1) + (1-t)f(x_2)$

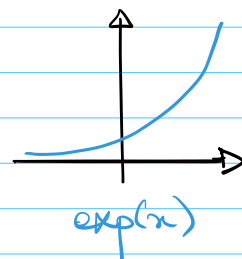
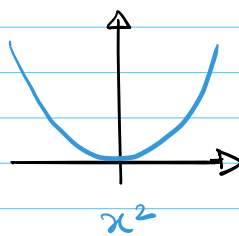
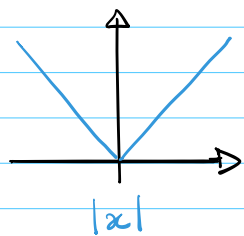
$\Leftrightarrow f$ convex if $\forall x, y \in A, t \in [0, 1]$

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



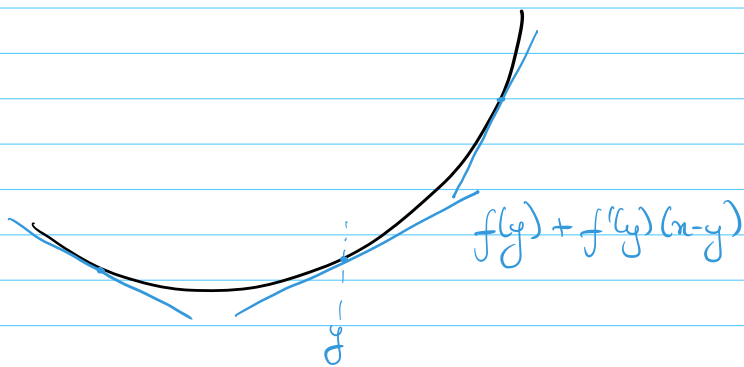
Chord is above the graph.

Examples:



Properties of convex functions:

- If f is differentiable + convex
 - $d=1$ f' is non decreasing
 - $d \geq 2$ $\langle \nabla f(x) - \nabla f(y), x - y \rangle \geq 0$and $f(x) \geq f(y) + \langle \nabla f(y), x - y \rangle$



- If f is twice differentiable + convex
 - $d=1$ $f'' \geq 0$
 - $d \geq 2$ $\forall x, \nabla^2 f(x) \succeq 0$.

symmetric positive definite matrix H

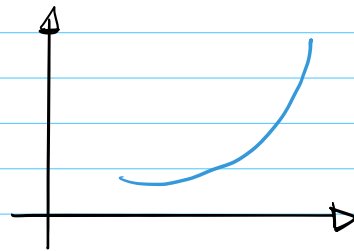
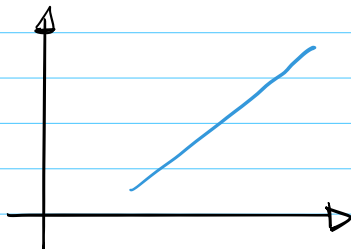
$$\forall u \in \mathbb{R}^d \quad u^T H u \geq 0$$

\leadsto Let $g_u: t \mapsto f(x+tu)$

$$g_u \text{ conv} : g_u'(t) = \langle \nabla f(x+tu), u \rangle$$

$$g_u''(t) = \langle \nabla^2 f(x+tu)u, u \rangle$$

$$\leadsto g_u''(0) = u^T \nabla^2 f(x)u \geq 0.$$



Both functions are convex ---
But one is "more" convex.

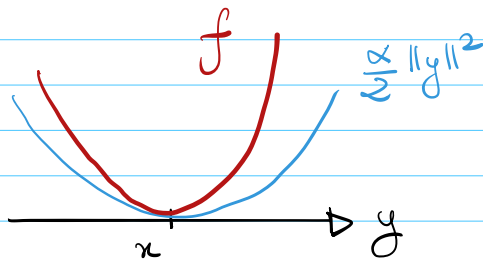
• α -strongly convex function

$$\forall x, y, t \in [0, 1]$$

$$f(tx + (1-t)y) \leq tf(x) + (1-t)f(y) - \frac{\alpha}{2} t(1-t) \|x-y\|^2$$

f differentiable \Rightarrow $\left[f(y) \geq f(x) + \langle \nabla f(x), y-x \rangle + \frac{\alpha}{2} \|x-y\|^2 \right]$

$$\begin{cases} f(x) = 0 \\ \nabla f(x) = 0 \end{cases}$$



f is sufficiently curved

f twice differentiable
 \Rightarrow

$$\nabla^2 f(x) \succeq \alpha \text{Id}$$

$$\Leftrightarrow H \succeq \alpha \text{Id} \Leftrightarrow v^T H v \geq \alpha \|v\|^2 \quad \forall v \in \mathbb{R}^d$$

$$\hookrightarrow v^T H v = \sum_i d_i \langle v, e_i \rangle^2$$

(e_1, \dots, e_d) orthonormal basis of eigenvectors
 d_1, \dots, d_d associated eigenvalues

$$\Leftrightarrow d_i \geq \alpha.$$

• β -smooth function:

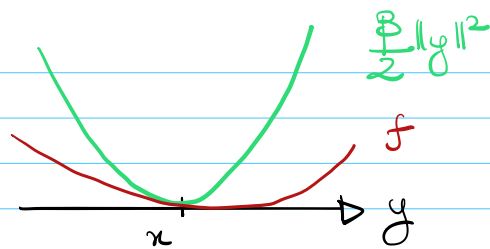
$$\forall x, y \quad \|\nabla f(x) - \nabla f(y)\| \leq \beta \|x - y\|$$

$$\Leftrightarrow \forall x, y \quad f(y) \leq f(x) + \langle \nabla f(x), y - x \rangle + \frac{\beta}{2} \|y - x\|^2$$

twice differentiable

\Leftrightarrow

$$\nabla^2 f(x) \preceq \beta \text{Id}$$



f is not too curved

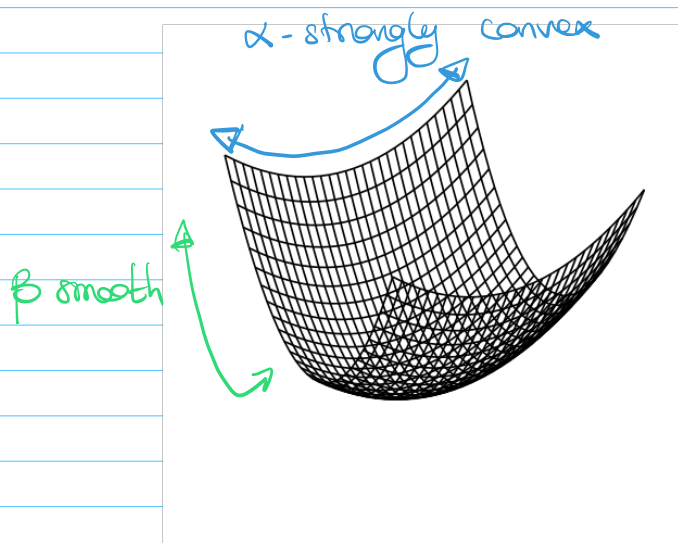
Example: A function α -strongly convex

$$\alpha \leq \beta$$

and β -smooth is

$$f_{\alpha, \beta} : x \in \mathbb{R}^2 \mapsto \frac{\alpha}{2} x_1^2 + \frac{\beta}{2} x_2^2$$

$$\nabla^2 f_{\alpha, \beta}(x) = \begin{pmatrix} \alpha & 0 \\ 0 & \beta \end{pmatrix} \begin{cases} \succeq \alpha \text{Id} \\ \preceq \beta \text{Id} \end{cases}$$



③ Gradient Descent

GOAL: Find the minimum of a convex function f .

- Start at a point x .
- Make a step of size r in some direction h

$$f(x+rh) \approx f(x) + r \langle \nabla f(x), h \rangle$$

\rightarrow largest decrease for $h = -\nabla f(x)$.

Def: $f: \mathbb{R}^d \rightarrow \mathbb{R}$
 x_0 initialization
 $s > 0$ step size

GRADIENT
DESCENT

$$x^{t+1} = x^t - s \nabla f(x^t)$$

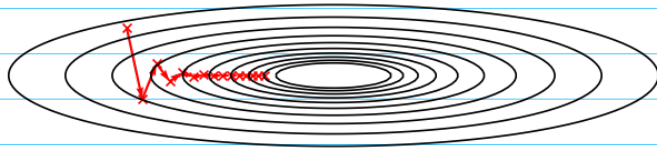
THEOREM: Let f α -strongly convex and β -smooth.

Take $s < 1/\beta$. Let $x^* = \operatorname{argmin} f$.

Gradient
 \Rightarrow
Descent

$$f(x^T) - f(x^*) \leq \underbrace{(1-\alpha s)^T}_{\leq \exp(-\alpha s T)} (f(x^0) - f(x^*))$$

Example GD on $f_{\alpha, \beta}$



proof:

$$\begin{aligned} f(x^{t+1}) &= f(x^t - s \nabla f(x^t)) \\ &\stackrel{\beta\text{-smooth}}{\leq} f(x^t) + \underbrace{\langle \nabla f(x^t), -s \nabla f(x^t) \rangle}_{-s \|\nabla f(x^t)\|^2} + s^2 \frac{\beta}{2} \|\nabla f(x^t)\|^2 \\ &= f(x^t) + s \underbrace{\left(\frac{s\beta}{2} - 1 \right)}_{< 0} \|\nabla f(x^t)\|^2 \end{aligned}$$

Poljak - Łojasiewicz inequality:

$$\textcircled{PL} \left[f(x) - f(x^*) \leq \frac{1}{2\alpha} \|\nabla f(x)\|^2 \right]$$

proof: $f(x^*) \geq f(x) + \langle \nabla f(x), x^* - x \rangle + \frac{\alpha}{2} \|x^* - x\|^2$

$$\Rightarrow f(x) - f(x^*) \leq \langle \nabla f(x), x - x^* \rangle - \frac{\alpha}{2} \|x^* - x\|^2$$

$$\Rightarrow \left[f(x^{t+1}) - f(x^*) \leq (f(x^t) - f(x^*)) \underbrace{\left(1 + 2\alpha s \left(\frac{s\beta}{2} - 1 \right) \right)}_{\leq 1 - \alpha s} \right]$$

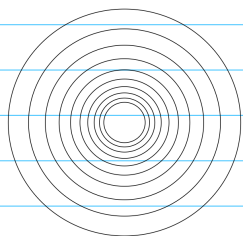
$s < \frac{1}{\beta}$ $\leq 1 - \alpha s$ \square

[GRADIENT DESCENT HAS A LINEAR RATE OF CV.]

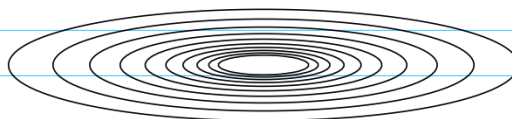
• If $s = 1/\beta$: $f(x^T) - f(x^*) \leq \exp\left(-\frac{\kappa}{\beta} T\right) (f(x^0) - f(x^*))$

$\kappa = \frac{\beta}{\alpha} \geq 1$ is the condition number.

$\kappa \gg 1$ = ill-conditioned problem
= harder to optimize



$\kappa = 1$



$\kappa \gg 1$

↳ To get a precision ϵ

$$\epsilon = \exp\left(-\frac{T}{\kappa}\right) \rightsquigarrow T = \frac{\log(\epsilon^{-1})}{\kappa} \text{ iterations}$$

Excess risk: $R_p(\hat{f}_T) - R_p(f_p^*)$ at least of order $\frac{1}{\sqrt{T}}$.

$$\rightsquigarrow \boxed{T \approx \frac{\log(\epsilon)}{\kappa} \text{ iterations.}}$$

→ What if f is only smooth? ($\alpha=0$)

THEOREM: Let f be a convex β -smooth function.

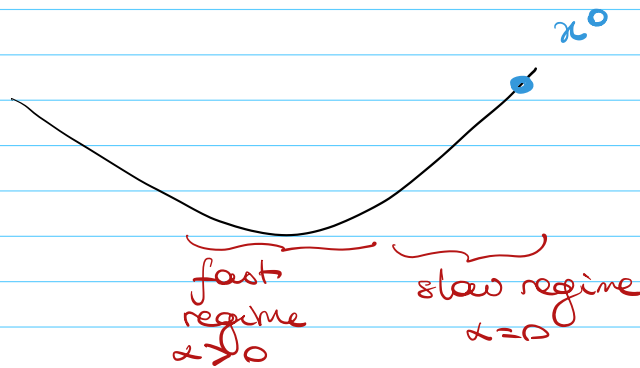
Take $s = 1/\beta$. Let $x^* = \operatorname{argmin} f$.

Gradient
Descent

$$f(x^T) - f(x^*) \leq \frac{\beta \|x^0 - x^T\|^2}{2T}$$

→ Much slower

$T = \frac{\beta}{\epsilon}$ iterations for error ϵ .



→ Gradient Descent will naturally adapt to the degree of convexity of f .

④ Newton's Method

Second order Approximation:

$$f(x+h) \approx P_{f,x}(h) = f(x) + \langle \nabla f(x), h \rangle + \frac{1}{2} h^T \nabla^2 f(x) h$$

↳ min attained at $h = (\nabla^2 f(x))^{-1} \nabla f(x)$

Def: $f: \mathbb{R}^d \rightarrow \mathbb{R}$
↳ no initialization

NEWTON'S
METHOD

$$\left[x^{t+1} = x^t - (\nabla^2 f(x^t))^{-1} \nabla f(x^t) \right]$$

THEOREM: Let f be:

- α strongly convex

- β smooth

- $\forall x, y \quad \|\nabla^2 f(x) - \nabla^2 f(y)\|_{op} \leq \gamma \|x - y\|$

- Good initialization: $\|x^0 - x^*\| \leq \alpha/\gamma$

Newton's
method

$$\left[f(x^T) - f(x^*) \leq \frac{\beta\alpha}{\gamma} 2^{-2T+1} \right]$$

[NEWTON'S METHOD HAS A QUADRATIC RATE OF CV.]

$T = O(\underbrace{\log \log \varepsilon^{-1}}_{\text{super small}})$ iterations to get an ε -approximation.

$\leadsto \log \log n$ iterations.

However: ① Compute $(\nabla^2 f(x))^{-1}$

$\rightarrow O(d^3)$ operations (naively)

Each iteration is very costly if d is large.

② Method breaks down if no α -strong convexity.

⑤ Logistic Regression

Recall: Binary classification

→ classifier of the form $f = \text{sgn} \circ g$

$$\begin{aligned} R_p(f) &= \mathbb{E}_p[\mathbb{1}\{f(x) \neq Y\}] \\ &= \mathbb{E}_p[\mathbb{1}\{\text{sgn}(g(x)) \neq Y\}] \\ &= \mathbb{E}_p[\mathbb{1}\{g(x)Y < 0\}] \end{aligned}$$

We replace $\mathbb{1}\{t < 0\}$ by a convex function.



$$\begin{aligned} \text{Logistic loss: } t &\longmapsto \log(1 + e^{-t}) \\ &= -\log(\sigma(t)) \end{aligned}$$

$$\sigma(t) = \frac{1}{1 + e^{-t}} \in [0, 1]$$

New loss: $l_{\log}(y, y') = \log(1 + e^{-yy'})$

$$\Rightarrow \tilde{R}_n(g) = \frac{1}{n} \sum_{i=1}^n \log(1 + \exp(-Y_i g(X_i)))$$

$$= \frac{1}{n} \sum_{i=1}^n (Z_i \log(1 + \exp(-g(X_i)))$$

$$Z_i = \mathbb{1}\{Y_i = 1\} + (1 - Z_i) \log(1 + \exp(g(X_i)))$$
$$= \begin{cases} 1 & \text{if } Y_i = 1 \\ 0 & \text{if } Y_i = -1 \end{cases}$$

Link with Maximum Likelihood Estimation

Statistical model: Set \mathcal{G} . Fix $g \in \mathcal{G}$.

Assume Y is obtained in the following way:

$$\begin{cases} Y = 1 & \text{with probability } \sigma(g(X)) = p_g(X, 1) \\ Y = -1 & \text{otherwise.} \end{cases} = \frac{1}{1 + e^{-g(X)}}$$

\hookrightarrow with proba $p_g(X, 0)$

$\Rightarrow P_g$ is the joint distribution of (X, Y) .

\Rightarrow We observe $(X_1, Y_1) \dots (X_n, Y_n)$ with distribution $P_{g_0} \rightarrow$ Estimate g_0 .

Likelihood of $(X_1, Y_1) \dots (X_n, Y_n)$ in g :

$$\prod_{i=1}^n p_g(X_i, Y_i) = \prod_{i=1}^n \sigma(g(X_i))^{\mathbb{1}\{Y_i=1\}} (1 - \sigma(g(X_i)))^{\mathbb{1}\{Y_i=-1\}}$$

Log-Likelihood:

$$\sum_{i=1}^n z_i \log(\sigma(g(x_i))) + (1-z_i) \log(1-\sigma(g(x_i)))$$

$$= -\tilde{R}_n(g)$$

Max Likelihood = Min Empirical Risk

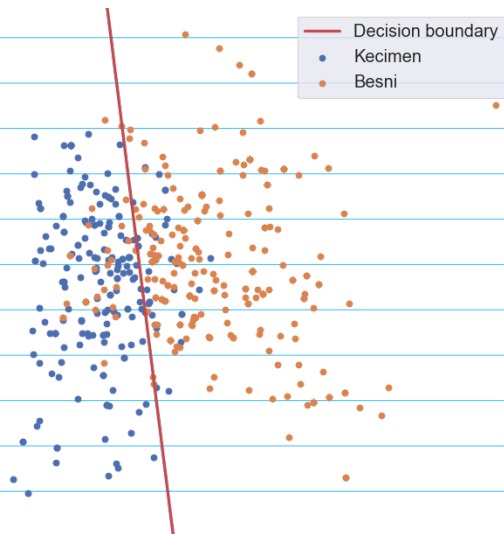
\hat{g}_{ML} satisfies strong theoretical properties.

Example: Two different variety of rainns.

Eight Geometric Features
(area, perimeter, ...)

$\mathcal{G} = \{ \text{Linear classifiers} \}$

no Optimization with gradient descent.



85% accuracy

