# Chapter 8

# Optimization basics

## Contents (class version)

---

## 8.0 Introduction

Many of the previous topics have involved **optimization** formulations: linear LS, Procrustes, low-rank approximation, multidimensional scaling. In all these cases we derived analytical solutions, like the pseudo-inverse for minimum-norm LS problems and the truncated SVD for low-rank approximation.

But often we need iterative optimization algorithms, *e.g.*,
• if no closed-form minimizer exists, or
• if the analytical solution requires too much computation and/or memory, *e.g.*, SVD for large problems.

To solve a problem like $\hat{\boldsymbol{x}} = \arg\min_{\boldsymbol{x}} f(\boldsymbol{x})$ via an iterative method, we start with some initial guess $\boldsymbol{x}_0$, and then the algorithm produces a sequence $\{\boldsymbol{x}_k\}$ where hopefully the sequence **converges** to $\hat{\boldsymbol{x}}$, meaning $\|\boldsymbol{x}_k - \hat{\boldsymbol{x}}\| \to 0$ for some norm $\|\cdot\|$ as $k \to \infty$, as discussed in Ch. 5.

The homework has introduced a few such optimization algorithms. This chapter elaborates on those. Along the way we introduce several new matrix concepts: **matrix square root**, **matrix powers**, more about **positive (semi)definite matrices**, **commuting matrices**, and **majorization**.

---

### 8.1 Preconditioned gradient descent (PGD) for LS

For the (possibly regularized) LS problem with $A \in \mathbb{F}^{M \times N}$ the **cost function** is **quadratic**:

$$\hat{x} = \arg\min_{x} f(x), \quad \underbrace{f(x) = \frac{1}{2} \|Ax - y\|_2^2}_{f : \mathbb{F}^N \mapsto \mathbb{R}}, \quad \underbrace{\nabla f(x) = A'(Ax - y)}_{\nabla f : \mathbb{F}^N \mapsto \mathbb{F}^N}.$$

The homework problems focused on the classical **gradient descent** (**GD**) iterative method:

$$x_{k+1} = x_k - \alpha \nabla f(x_k) = x_k - \alpha A'(Ax_k - y),$$

where convergence of the sequence $\{x_k\}$ is ensured (according to HW) by choosing the **step size** $\alpha$ to satisfy

$$0 < \alpha < \frac{2}{\sigma_1(A'A)} = \frac{2}{\sigma_1^2(A)}. \tag{8.1}$$

Now consider a generalization called **preconditioned gradient descent** (**PGD**):

$$x_{k+1} = x_k - P\nabla f(x_k) = x_k - PA'(Ax_k - y), \tag{8.2}$$

where $P$ is a $N \times N$ **preconditioning matrix**. Classical gradient descent simply uses $P = \alpha I_N$. But a single step size $\alpha$ may be suboptimal (and hard to choose), especially if different elements of $x$ have different units.

What conditions on $P$ ensure convergence? Will some other $P \neq \alpha I$ provide faster convergence? Why does (8.1) suffice? Answering these questions will use many matrix methods!

## Tool: Matrix square root

We need another linear algebra tool first. (And more will come before we finish this topic.)

Define. We call $S$ a (matrix) square root of a square matrix $A$ iff $SS = A$.

It is not unique: if $S$ is a square root of $P$, then $-S$ is also square root of $P$.

Example. For the rotation matrix $R = \begin{bmatrix} \cos\phi & \sin\phi \\ -\sin\phi & \cos\phi \end{bmatrix}$ a square root is $S = \begin{bmatrix} \cos(\phi/2) & \sin(\phi/2) \\ -\sin(\phi/2) & \cos(\phi/2) \end{bmatrix}$ because $SS = R$. This example has an intuitive geometric interpretation.

If $P$ is positive semidefinite, i.e., $P \succeq 0$, then $P$ has a positive semidefinite square root.
Proof: $P \succeq 0 \Longrightarrow P = V\Lambda V'$, where the eigenvalues are all real and nonnegative,
so let $S = V\Lambda^{1/2}V'$ where $\Lambda^{1/2} \triangleq \mathrm{diag}\{\sqrt{\lambda_i}\}$. Clearly $SS = P$ and $S \succeq 0$.

If $P$ is positive definite, $P \succ 0$, then $P$ has a positive definite (and hence invertible) square root.

Proof: Same as above only now the eigenvalues are all positive.

Define. We write $P^{1/2}$ to denote the positive (semi)definite square root of $P$; we use this notation only if $P$ is positive (semi)definite. When being careful, we call $P^{1/2}$ the principal square root, but often we just call it "the square root" of $P$, just like people say "the square root of 9 is 3."

Every **diagonalizable** matrix has a **matrix square root**.
A: True                                    B: False                    ??

Example. If $A = \begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$ then $S = \begin{bmatrix} 1 & 1/2 \\ 0 & 1 \end{bmatrix}$ is a matrix square root of $A$.

If a square matrix has a **matrix square root**, then it is **diagonalizable**.
A: True                                    B: False                    ??

More generally, let $A$ be any $2 \times 2$ matrix. Even if $A$ is not diagonalizable, we can still factor it using the **Jordan normal form**: as follows:

$$A = V \begin{bmatrix} \lambda & 1 \\ 0 & \lambda \end{bmatrix} V^{-1}$$

where $V$ is an invertible matrix consisting of **generalized eigenvectors**. One can verify that if $\lambda \neq 0$, then a square root of this matrix is

$$S = V \begin{bmatrix} \sqrt{\lambda} & \frac{1}{2\sqrt{\lambda}} \\ 0 & \sqrt{\lambda} \end{bmatrix} V^{-1}.$$

This **Jordan form** approach to find a matrix square root generalizes to any square matrix having positive eigenvalues. But not every matrix has a square root.

Example. The matrix $\begin{bmatrix} 0 & 1 \\ 0 & 0 \end{bmatrix}$ has no square root.

Example.

For $x \neq 0$, consider $A = xx' \succeq 0$. Which of the following is a matrix square root of $A$?

A: $xx'$      B: $xx' / \|x\|_2$      C: $xx' / \|x\|_2^2$      D: $xx' \|x\|_2$      E: $xx' \|x\|_2^2$

??

Example.

If $Q$ is a unitary matrix, then $Q$ has a matrix square root.

A: True                 B: False          ??

??

**Practical implementation of matrix square root** _____

For a square matrix, the command for finding a **matrix square root** is:      `sqrt(A)`

This operation differs *completely* from element-wise root:      `sqrt.(A)`

In MATLAB (and old JULIA versions) the operation is      `sqrtm(A)`

Caution. If $A$ is square but not diagonalizable then the output of `sqrt` can be meaningless.

Try `sqrt([0 1; 0 0])` in JULIA or `sqrtm([0 1; 0 0])` in MATLAB.

Hereafter we assume the **preconditioner** $P$ in PGD (8.2) is **positive definite**, *i.e.*, $P \succ 0$, so that $P^{1/2}$ exists and is invertible. We denote its inverse by $P^{-1/2}$.

### Convergence rate analysis of PGD: first steps

Let $\hat{x}$ denote any minimizer of $f(x)$. That $\hat{x}$ must satisfy the normal equations: $A'A\hat{x} = A'y$. For analysis purposes, replace $A'y$ in (8.2) with $A'A\hat{x}$:

$$\begin{aligned}
x_{k+1} = x_k - PA'\left(Ax_k - y\right) &= x_k - P\left(A'Ax_k - A'y\right) \\
&= x_k - P\left(A'Ax_k - A'A\hat{x}\right) = x_k - PA'A\left(x_k - \hat{x}\right) \\
\implies x_{k+1} - \hat{x} &= x_k - \hat{x} - PA'A\left(x_k - \hat{x}\right) \\
&= \rule{4cm}{0.4cm} \\
\implies \underbrace{P^{-1/2}\left(x_{k+1} - \hat{x}\right)}_{\delta_{k+1}} &= P^{-1/2}\left(I - PA'A\right)\left(x_k - \hat{x}\right) \\
&= \left(I - P^{1/2}A'AP^{1/2}\right)\underbrace{P^{-1/2}\left(x_k - \hat{x}\right)}_{\delta_k} \\
\implies \delta_{k+1} = \left(I - P^{1/2}A'AP^{1/2}\right)\delta_k &= G\delta_k, \qquad \begin{aligned} &\delta_k \triangleq P^{-1/2}\left(x_k - \hat{x}\right) \text{ (mod. error vector)} \\ &G \triangleq I - P^{1/2}A'AP^{1/2} \end{aligned} \\
\implies \delta_k = \rule{2cm}{0.4cm}\,.
\end{aligned}$$

The matrix $G$ "governs" the convergence of PGD.

---

**Tool: Matrix powers**

Now we need **matrix powers** to proceed, another easy byproduct of **eigendecomposition**.

Here $G = G'$, so $G = V\Lambda V'$ for some unitary $V$ and $\Lambda$ is real (but *not* nonnegative in general).

So $G^2 = GG = V\Lambda V'V\Lambda V' = V\Lambda^2 V'$ and more generally $G^k = V\Lambda^k V'$, $\forall k \in \mathbb{N}$.

If $A$ is **diagonalizable** and **invertible**, is $A^k$ diagonalizable and invertible for all $k \in \mathbb{N}$?

A: Yes                          B: No                          C: Depends                          ??

**PGD convergence condition using eigenvalues** _____

So we can express the modified error vector $\boldsymbol{\delta}_k$ in terms of the initial error:

$$\boldsymbol{\delta}_k = \phantom{xxxxxxx}$$

Define $\tilde{\boldsymbol{\delta}}_k \triangleq \boldsymbol{V}'\boldsymbol{\delta}_k$ to be the error coordinates in the $\boldsymbol{V}$ basis, then

$$\tilde{\boldsymbol{\delta}}_k = \phantom{xxxxxxx}$$

We have reduced the question of convergence of PGD down to seeing whether $\{\lambda_i^k\}$ is a decreasing geometric series. We need $-1 < \lambda_i < 1$ for all $i = 1, \ldots, N$ to ensure that all error components diminish to zero.

Put concisely, a necessary and sufficient condition to ensure convergence of PGD from any initializer $\boldsymbol{x}_0$ is:

$$\rho(\boldsymbol{G}) = \rho\left(\boldsymbol{I} - \boldsymbol{P}^{1/2}\boldsymbol{A}'\boldsymbol{A}\boldsymbol{P}^{1/2}\right) < 1, \quad i.e., \phantom{xxxxxxxxxxxx} \tag{8.3}$$

**Classical GD: step size bounds**

Consider the classical case where $P = \alpha I$. Then the above condition (8.3) simplifies to

$$-1 < \text{eig}\{I - \alpha A'A\} < 1, \quad \Longleftrightarrow \quad \boxed{\phantom{xxxxxxxxxx}} \quad \Longleftrightarrow \quad \boxed{\phantom{xxxx}}$$

where $\{\sigma_i\}$ denotes the singular values of $A$.

The lower bound condition $0 < \alpha \sigma_i^2$ holds for all $i$ if $0 < \alpha$ and if $A$ has full column rank (*e.g.*, when Tikhonov regularization is used). So even though the analysis did not assume full rank at the start, if we want to make sure that all error modes diminish to zero, we must have full (column) rank matrix $A$.

The upper bound condition $\alpha \sigma_i^2 < 2$ holds for all $i$ if $\alpha \sigma_1^2 < 2$ because $\sigma_1$ is the largest, so we need $\alpha < 2/\sigma_1^2$.

In summary, we derived the step-size condition (8.1) for convergence of classical GD for LS problem:

$$0 < \alpha < \frac{2}{\sigma_1^2(A)} = \frac{2}{\sigma_1(A'A)}.$$

## Optimal step size for GD

For classical GD, where $\boldsymbol{P} = \alpha \boldsymbol{I}$, the *optimal* step size (for fastest convergence) is:                    (HW)

$$\alpha_* = \frac{2}{\sigma_1^2(\boldsymbol{A}) + \sigma_N^2(\boldsymbol{A})}.$$                    (8.4)

This step size makes $\rho(\boldsymbol{I} - \alpha \boldsymbol{A}'\boldsymbol{A})$ as small as possible, *i.e.*, so that its largest absolute eigenvalue is as close to zero as possible. That largest eigenvalue will be the mode that converges to zero the slowest, so it will ultimately govern the convergence rate.

However, the step-size conditions (8.1) and (8.4) are unsatisfying practically because if a LS problem is so large that we cannot use the SVD to compute the pseudo-inverse, then probably we do not really want to use an SVD (or `svds`) to find $\sigma_1(\boldsymbol{A})$ and perhaps also $\sigma_N(\boldsymbol{A})$.

If $\boldsymbol{A}$ is unitary, what is the best step size $\alpha$ when $\boldsymbol{P} = \alpha \boldsymbol{I}$?
A: 0               B: 1               C: 1.99               D: 2               E: None of these               ??

Suppose $\boldsymbol{P} = \alpha(\boldsymbol{A}'\boldsymbol{A})^{-1}$ for some $\alpha \geq 0$. What is the best value of $\alpha$, so that the eigenvalues of $\boldsymbol{G}$ are as small (in magnitude) as possible, so PGD converges as fast as possible?
A: 0               B: 1/2               C: 1               D: $\sqrt{2}$               E: $2 - \epsilon$               ??

## Practical step size for GD

One option for avoiding the SVD is to use the bound from HW5:

$$\sigma_1(\boldsymbol{A}) = \|\boldsymbol{A}\|_2 \le \|\boldsymbol{A}\|_F.$$

So choosing $0 < \alpha < 2/\|\boldsymbol{A}\|_F^2$ always provides a valid step size. This is useful because it is easy to compute $\|\cdot\|_F$ because no SVD is needed. However, the upper bound above is often quite loose. If $\boldsymbol{A} = \boldsymbol{I}_N$, then $\|\boldsymbol{A}\|_2 = 1$ but $\|\boldsymbol{A}\|_F = \sqrt{N}$. So using this bound, the step size $\alpha$ would be much smaller than necessary.

Another option is to recall that because $\boldsymbol{A}'\boldsymbol{A}$ is symmetric:

$$\|\boldsymbol{A}'\boldsymbol{A}\|_2 = \sigma_1(\boldsymbol{A}'\boldsymbol{A}) = \rho(\boldsymbol{A}'\boldsymbol{A}) \le \|\boldsymbol{A}'\boldsymbol{A}\|,$$

for any matrix norm $\|\cdot\|$. In particular, as discussed in Ch. 5, $\|\boldsymbol{A}'\boldsymbol{A}\|_1 \le \|\boldsymbol{A}'\|_1 \|\boldsymbol{A}\|_1 = \|\boldsymbol{A}\|_\infty \|\boldsymbol{A}\|_1$.
So convergence is always ensured by choosing:

$$0 < \alpha < \frac{2}{\|\boldsymbol{A}'\boldsymbol{A}\|_1} \quad \text{or} \quad 0 < \alpha < \frac{2}{\|\boldsymbol{A}\|_\infty \|\boldsymbol{A}\|_1}.$$

Why did I choose $\|\boldsymbol{A}'\boldsymbol{A}\|_1$ instead of $\|\boldsymbol{A}'\boldsymbol{A}\|_\infty$ ?

Which norm is bigger:
A: $\|\boldsymbol{A}'\boldsymbol{A}\|_1$ usually   B: $\|\boldsymbol{A}'\boldsymbol{A}\|_1$ always   C: $\|\boldsymbol{A}'\boldsymbol{A}\|_\infty$ usually   D: $\|\boldsymbol{A}'\boldsymbol{A}\|_\infty$ always   E: They are the same.
??

### Ideal preconditioner for PGD

Before looking at more general choices for the preconditioner $P$, we first explore the "ideal" $P$.

The ideal preconditioner is $P_{\text{ideal}} = (A'A)^{-1}$,
because for that choice:

$$G = I - P_{\text{ideal}}^{1/2} A'A P_{\text{ideal}}^{1/2} = $$

so $\delta_{(1)} = G^1 \delta_0 = 0 \delta_0 = 0$, *i.e.*, PGD converges in one iteration:

$$x_1 = x_0 - P_{\text{ideal}} A'(Ax_0 - y) = x_0 - (A'A)^{-1} A'(Ax_0 - y) = (A'A)^{-1} A'y = A^+y.$$

Why not always use this ideal preconditioner?
Inverting $A'A$ is expensive: it is as hard as the original LS problem. So we look for other preconditioners. Specifically we would like to find a preconditioner $P \approx (A'A)^{-1}$ or $P^{-1} \approx A'A$ but where the inverse is much easier to compute.
For that we need another tool.

## Tool: Positive (semi)definiteness properties

Recall that we define positive (semi)definiteness only for Hermitian matrices and:
- **positive semidefinite**: $A \succeq 0 \iff x'Ax \geq 0, \forall x$
- **positive definite**: $A \succ 0 \iff x'Ax > 0, \forall x \neq 0$

Properties ($A$ and $B$ are all Hermitian here):

- $X'X \succeq 0$ for *any* matrix $X$, even non-square (shown previously)
- $A \succeq 0 \implies X'AX \succeq 0$, for any matrix $X$ with `size(A,2) == size(X,1)`
- $A \succeq 0$ and $\alpha \geq 0 \implies \alpha A \succeq 0$
- $A \succeq 0 \implies$ all eigenvalues of $A$ are real and nonnegative
  Proof 1: if $x$ is an eigenvector of $A$ then $0 \leq x'(Ax) = x'\lambda x = \lambda \|x\|_2^2 \implies \lambda \geq 0$
  Proof 2: $A$ Hermitian implies it has an orthogonal eigendecomposition: $A = V\Lambda V'$ so $V\Lambda V' \succeq 0$.
  Using 2nd property above: $\implies V'(V\Lambda V')V \succeq 0 \implies \Lambda \succeq 0 \implies e_j'\Lambda e_j \geq 0 \implies \lambda_j \geq 0$.
- $A \succ 0 \implies A$ invertible and $A^{-1} \succ 0$
- $B \succeq A \iff B - A \succeq 0$ (trivial from definition of $\succeq$)
- $B \succ A \iff B - A \succ 0$
- $B \succeq A \succ 0 \implies A^{-1} \succeq B^{-1} \succ 0$
- $B \succeq A \succ 0 \implies \gamma B \succ A, \forall \gamma > 1$
- $A \succ 0, B \succ 0 \implies BAB \succ 0$ (HW)

Notice the pattern: start with a definition, then develop properties (especially addition and multiplication).

The above properties relate to multiplication; now consider **matrix addition**.

If $A \succeq 0$ and $B \succeq 0$ have the same size, then $A + B \succeq 0$. (?)

A: True                                                B: False        ??

If $A \succ 0$ and $B \succeq 0$ have the same size, then $A + B \succ 0$. (?)

A: True                                                B: False        ??

The above properties are all useful and important for GD for LS because:
• we choose $P \succ 0$
• and we assume $A'A \succ 0$.

**General preconditioners for PGD**

Repeating (8.3), for convergence we want

$$-1 < \text{eig}\{I - P^{1/2}A'AP^{1/2}\} < 1. \tag{8.5}$$

Now the RHS of (8.5) is easy; assuming $A'A$ and $P$ are positive definite:

$$\text{eig}\{I - P^{1/2}A'AP^{1/2}\} = 1 - \underbrace{\text{eig}\{P^{1/2}A'AP^{1/2}\}}_{> 0} < 1.$$

For the LHS of (8.5), we want $-1 < \text{eig}\{I - P^{1/2}A'AP^{1/2}\} = 1 - \text{eig}\{P^{1/2}A'AP^{1/2}\}$, *i.e.*,  (HW)

$$\text{eig}\{P^{1/2}A'AP^{1/2}\} < 2 \iff 2I \succ P^{1/2}A'AP^{1/2} \iff 2P^{-1} \succ A'A.$$

**Matrix majorizers** _____

Define. We say that a (Hermitian) symmetric matrix $A$ is a **majorizer** of (or **majorizes**) a (Hermitian) symmetric matrix $B$ iff $A \succeq B$, *i.e.*, iff $A - B$ is positive semidefinite.

Summarizing: if $\beta P^{-1}$ **majorizes** $A'A$ for some $0 < \beta < 2$, then convergence of PGD is ensured, because $\beta P^{-1} \succeq A'A \implies 2P^{-1} \succ A'A$.

---

**Diagonal majorizer**

Fact (diagonal majorizer). If $\boldsymbol{B}$ is (Hermitian) symmetric and $N \times N$, then:     (HW)

$$\operatorname{diag}\{|\boldsymbol{B}|\,\mathbf{1}_N\} \succeq \boldsymbol{B}, \tag{8.6}$$

where here $|\boldsymbol{B}|$ denotes the element-wise absolute value of $\boldsymbol{B}$, *i.e.*, `abs.(B)` in JULIA.

Thus, a valid preconditioner for PGD is $\boldsymbol{P}^{-1} = \alpha^{-1} \operatorname{diag}\{|\boldsymbol{A}'\boldsymbol{A}|\,\mathbf{1}_N\}$ for any $0 < \alpha < 2$, leading to the **diagonally preconditioned GD** iteration:

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - \alpha \boldsymbol{D}\boldsymbol{A}'\left(\boldsymbol{A}\boldsymbol{x}_k - \boldsymbol{y}\right), \quad \boldsymbol{D} \triangleq \operatorname{diag}\{|\boldsymbol{A}'\boldsymbol{A}|\,\mathbf{1}_N\}^{-1}, \quad 0 < \alpha < 2. \tag{8.7}$$

Typically we use $1 \le \alpha \le 1.99$ and it is easy to adjust $\alpha$ in this range to find good value.

If $\boldsymbol{A} = \begin{bmatrix} 1 & 0 \\ 0 & 2 \\ 1 & -1 \end{bmatrix}$ then what is the diagonal preconditioner $\boldsymbol{D}$ in (8.7) ?

A: $\begin{bmatrix} 1/3 & 0 \\ 0 & 1/6 \end{bmatrix}$     B: $\begin{bmatrix} 1/2 & 0 \\ 0 & 1/3 \end{bmatrix}$     C: $\begin{bmatrix} 1/2 & 0 \\ 0 & 1/5 \end{bmatrix}$     D: $\begin{bmatrix} 2 & 0 \\ 0 & 5 \end{bmatrix}$     E: None of these.

??

The "normalization" by the inverse matrix $D$ in (8.7) makes finding $\alpha$ much easier than in classical GD where changing the units of the problem (units of $x$ or $y$ and hence $A$) necessitate finding a new $\alpha$ value.

The following code compares eig$\{I - P^{-1/2}A'AP^{-1/2}\}$ for
- classical GD with $P = \alpha_* I$ for the optimal step size $\alpha_*$,
- versus diagonally preconditioned GD with $\alpha = 1.3$ chosen empirically.

The former has eigenvalues $\pm 0.5151$, and the latter has eigenvalues $(-0.3, 0.35)$ so at least in this example the diagonal preconditioner accelerates convergence, without requiring an eigendecomposition to find $P$.

```julia
using LinearAlgebra
A = [1 0; 0 2; 1 -1]
@show A'A
D = Diagonal(1 ./ (abs.(A'A) * ones(2))) # diagonal preconditioner
alpha1 = 1.3
P1 = alpha1 * D
G1 = I - sqrt(P1) * (A' * A) * sqrt(P1)
@show eigvals(G1)

alphabest = 2 / (sum(svdvals(A'A)[[1, end]])) # optimal GD step size
G2 = I - alphabest * (A' * A)
@show eigvals(G2)
```
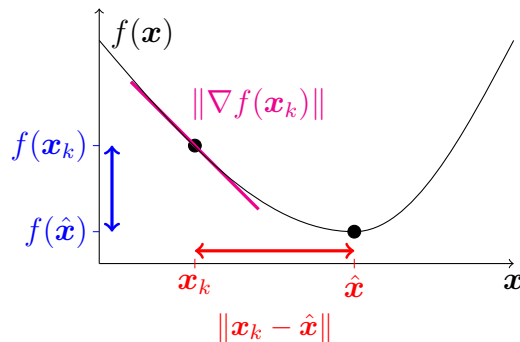
We have focused on the quadratic problem here, but the principles of diagonal majorization apply to using GD (and accelerated GD) in general for convex optimization of cost functions having Lipschitz gradients [1].

### Convergence rates

There are many ways to assess the convergence rate of an iterative algorithm like PGD. Researchers study:

- $f(\boldsymbol{x}_k) \to f(\hat{\boldsymbol{x}})$
- $\|\nabla f(\boldsymbol{x}_k)\| \to 0$
- $\|\boldsymbol{x}_k - \hat{\boldsymbol{x}}\| \to 0$
- $\|\boldsymbol{\delta}_k\| \to 0, \quad \boldsymbol{\delta}_k \triangleq \boldsymbol{P}^{-1/2}(\boldsymbol{x}_k - \hat{\boldsymbol{x}})$
- $\|\tilde{\boldsymbol{\delta}}_k\| \to 0, \quad \tilde{\boldsymbol{\delta}}_k \triangleq \boldsymbol{V}'\boldsymbol{\delta}_k$



Quantifying bounds on the rates of decrease of these quantities is an active research area. Even classical GD has relatively recent results [2] that tighten up the traditional bounds. The tightest possible worst-case bound for GD for the decrease of the cost function (with a fixed step size $\alpha = 1/L$) is $O(1/k)$:

$$f(\boldsymbol{x}_k) - f(\hat{\boldsymbol{x}}) \leq \frac{\|\boldsymbol{x}_0 - \hat{\boldsymbol{x}}\|_2^2}{L(4k+2)},$$
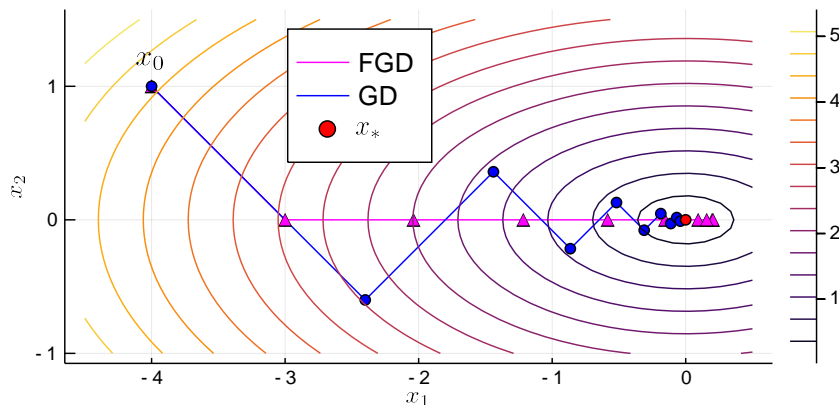
where $L$ is the Lipschitz constant of the gradient $\nabla f(\boldsymbol{x})$.

In contrast, Nesterov's fast gradient method has a worst-case cost function decrease at rate at least $O(1/k^2)$, which can be improved (and has) by only a constant factor [3].

Example. The following figure illustrates how slow GD can converge for a simple LS problem
with $A = \begin{bmatrix} 1 & 0 \\ 0 & 2 \end{bmatrix}$ and $y = 0$. This case used the optimal step size $\alpha_*$ for illustration.
This slow convergence has been the impetus of thousands of papers on faster algorithms!



The ellipses show the contours of the LS cost function $\|Ax - y\|$.
Also shown is Nesterov's fast gradient descent (FGD) method for comparison.

What is the optimal step size $\alpha_*$ here?

A: 1/5          B: 2/5          C: 3/5          D: 1/3          E: 2/3          ??

**Tool: commuting (square) matrices**

Let $X$ and $Y$ denote square matrices of the same size.
Which of the following conditions guarantee that $X$ and $Y$ commute, *i.e.*, $XY = YX$.
Choose the most general correct combination of condition(s).
- 1 $X$ diagonalizable
- 2 $Y$ diagonalizable
- 3 $X$ and $Y$ have same eigenvectors (*i.e.*, are simultaneously diagonalizable)
- 4 $X$ (Hermitian) symmetric
- 5 $Y$ (Hermitian) symmetric

A: 1 & 2
B: 1 & 2 & 3
C: 1 & 2 & 3 & (4 or 5)
D: 1 & 2 & 3 & 4 & 5
E: None of these suffice.

??

All $N \times N$ **companion** matrices commute. (?)
A: True                                    B: False                    ??
??

All $N \times N$ **circulant** matrices commute. (?)
A: True                                    B: False                    ??

If $X$ and $Y$ are positive semidefinite and are simultaneously diagonalizable then

$$XY = YX = X^{1/2}YX^{1/2} = Y^{1/2}XY^{1/2}.$$

(?)
A: True                                    B: False                    ??

Being simultaneously diagonalizable is a sufficient condition. A necessary (but not sufficient) condition for **commuting matrices** is that they are **simultaneously triangularizable**.

Example. The matrices $\begin{bmatrix} 1 & 1 \\ 0 & 1 \end{bmatrix}$ and $\begin{bmatrix} 1 & 0 \\ 0 & 1 \end{bmatrix}$ commute, but the first is not diagonalizable.

**Monotonicity**

Some algorithms get closer to the solution every iteration, and/or decrease the cost function every iteration, whereas other algorithms have no such guarantee.

For PGD, there are many quantities we can evaluate to see if they do (or do not) **decrease monotonically**, where to be precise we really mean "non-increasing" but usually we say "decreasing" for brevity:

- $f(\boldsymbol{x}_{k+1}) \overset{?}{\leq} f(\boldsymbol{x}_k)$ cost function
- $\|\nabla f(\boldsymbol{x}_{k+1})\| \overset{?}{\leq} \|\nabla f(\boldsymbol{x}_k)\|$ gradient norm
- $\|\boldsymbol{x}_{k+1} - \hat{\boldsymbol{x}}\| \overset{?}{\leq} \|\boldsymbol{x}_k - \hat{\boldsymbol{x}}\|$ distance to solution in natural coordinates
- $\|\boldsymbol{\delta}_{k+1}\| \overset{?}{\leq} \|\boldsymbol{\delta}_k\|, \quad \boldsymbol{\delta}_k \triangleq \boldsymbol{P}^{-1/2}(\boldsymbol{x}_k - \hat{\boldsymbol{x}})$ distance to solution with a $\boldsymbol{P}$-related coordinate transform
- $\left\|\tilde{\boldsymbol{\delta}}_{k+1}\right\| \overset{?}{\leq} \left\|\tilde{\boldsymbol{\delta}}_k\right\|, \quad \tilde{\boldsymbol{\delta}}_k \triangleq \boldsymbol{V}'\boldsymbol{\delta}_k$ distance to solution under the $\boldsymbol{V}$ transform

When one of the above inequalities hold (typically with strict inequality for $\boldsymbol{x}_k \neq \hat{\boldsymbol{x}}$ we say the algorithm is **monotonic** or **monotone** but to be complete we should also qualify that by saying in which sense it is monotone (*e.g.*, cost function, gradient norm, distance to solution).

The easiest of the above list to analyze for PGD is the last one.

Recall $\tilde{\boldsymbol{\delta}}_{k+1} = \boldsymbol{\Lambda}\tilde{\boldsymbol{\delta}}_k$. If $\tilde{\boldsymbol{\delta}}_k \neq \boldsymbol{0}$, then because we choose $\boldsymbol{P}$ so that $-1 < \lambda_n < 1$ per (8.5), the distance to the solution diminishes monotonically in these coordinates:

$$\left\|\tilde{\boldsymbol{\delta}}_{k+1}\right\|_2 = \left\|\boldsymbol{\Lambda}\tilde{\boldsymbol{\delta}}_k\right\|_2 \leq \|\boldsymbol{\Lambda}\|_2 \left\|\tilde{\boldsymbol{\delta}}_k\right\|_2 < \left\|\tilde{\boldsymbol{\delta}}_k\right\|_2.$$

If $\boldsymbol{P} \succ \boldsymbol{0}$ and $2\boldsymbol{P}^{-1} \succ \boldsymbol{A}'\boldsymbol{A} \succ \boldsymbol{0}$, does $\|\boldsymbol{\delta}_k\|_2$ decrease monotonically?

A: Yes.          B: Not necessarily.       ??

Under the same conditions, when $x_k \neq \hat{x}$ does the distance $\|x_k - \hat{x}\|_2$ decrease monotonically? Choose most general correct answer.

  A: Yes.
  B: Yes, if the right singular vectors of $\boldsymbol{A}$ are eigenvectors of $\boldsymbol{P}$.
  C: Yes, if the left singular vectors of $\boldsymbol{A}$ are eigenvectors of $\boldsymbol{P}$.
  D: Yes, if $\boldsymbol{P} = \alpha\boldsymbol{I}$ for suitably chosen $\alpha$.
  E: No.

??

Hint: $\|x_{k+1} - \hat{x}\|_2 = \|(\boldsymbol{I} - \boldsymbol{P}\boldsymbol{A}'\boldsymbol{A})(x_k - \hat{x})\|_2 \leq \|\boldsymbol{I} - \boldsymbol{P}\boldsymbol{A}'\boldsymbol{A}\|_2 \|x_k - \hat{x}\|_2$

## 8.2 Preconditioned steepest descent

(Read)

Instead of using GD with a fixed step size $\alpha$, an alternative is to do a **line search** to find the best step size *at each iteration*. This variation is called **steepest descent** (or GD with a line search). Here is how **preconditioned steepest descent** for a linear LS problem works:

$$
\begin{aligned}
\boldsymbol{d}_k &= -\boldsymbol{P}\nabla f(\boldsymbol{x}_k) = -\boldsymbol{P}\boldsymbol{A}'(\boldsymbol{A}\boldsymbol{x}_k - \boldsymbol{y}) \quad &\text{search direction (negative preconditioned gradient)} \\
\alpha_k &= \arg\min_{\alpha} f(\boldsymbol{x}_k + \alpha\boldsymbol{d}_k) \quad &\text{step size} \\
\boldsymbol{x}_{k+1} &= \boldsymbol{x}_k + \alpha_k\boldsymbol{d}_k \quad &\text{update.}
\end{aligned}
$$

- Finding $\alpha_k$ is a HW problem
- By construction, this iteration is guaranteed to decrease the cost function monotonically: $f(\boldsymbol{x}_{k+1}) \leq f(\boldsymbol{x}_k)$ with strict decrease unless $\boldsymbol{x}_k$ is already a minimizer, provided the preconditioner $\boldsymbol{P}$ is **positive definite**.
- Computing $\alpha_k$ takes some extra work, especially for non-quadratic problems. Often Nesterov's fast gradient method or the **optimized gradient method** (**OGM**) [3] are preferable because they do not require a line search (if the Lipschitz constant is available).

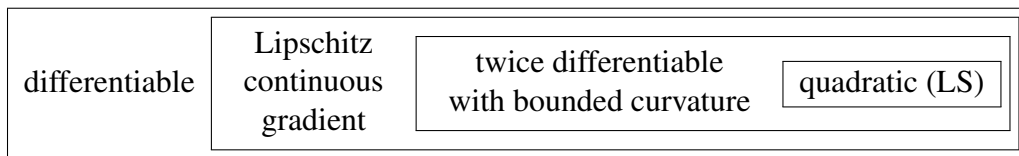## 8.3 Gradient descent for smooth convex functions

So far we used (preconditioned) **gradient descent** (**PGD**) only for LS problems.
We often need to solve more general (non LS) unconstrained minimization problems:

$$\hat{\boldsymbol{x}} = \arg\min_{\boldsymbol{x}\in\mathbb{R}^N} f(\boldsymbol{x}).$$

What algorithm we use depends in part on the properties of the cost function $f : \mathbb{R}^N \mapsto \mathbb{R}$.

| Venn diagram for convex functions: | differentiable | Lipschitz continuous gradient | twice differentiable with bounded curvature | quadratic (LS) |
|---|---|---|---|---|

GD is applicable to the broader family of **convex** functions having **gradients** that are **Lipschitz continuous**.
We call these **smooth convex** functions.

Define. A differentiable function $f(\boldsymbol{x})$ has a **Lipschitz continuous gradient** if there exists $L < \infty$ such that
$$\|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{z})\|_2 \leq L \|\boldsymbol{x} - \boldsymbol{z}\|_2, \quad \forall \boldsymbol{x}, \boldsymbol{z} \in \mathbb{R}^N.$$

Example. For $f(\boldsymbol{x}) = \frac{1}{2} \|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{y}\|_2^2$ we have $\|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{z})\|_2 = \|\boldsymbol{A}'(\boldsymbol{A}\boldsymbol{x} - \boldsymbol{y}) - \boldsymbol{A}'(\boldsymbol{A}\boldsymbol{z} - \boldsymbol{y})\|_2$
$= \|\boldsymbol{A}'\boldsymbol{A}(\boldsymbol{x} - \boldsymbol{z})\|_2 \leq \|\boldsymbol{A}'\boldsymbol{A}\|_2 \|\boldsymbol{x} - \boldsymbol{z}\|_2$.
So the Lipschitz constant of $\nabla f$ is $L = \|\boldsymbol{A}'\boldsymbol{A}\|_2 = \|\boldsymbol{A}\|_2^2 = \sigma^2(\boldsymbol{A}) = \rho(\boldsymbol{A}'\boldsymbol{A})$.

Fact. If $f(x)$ is **twice differentiable** and if there exists $L < \infty$ such that its **Hessian matrix** has a bounded **spectral norm**:
$$\left\|\nabla^2 f(x)\right\|_2 \leq L, \quad \forall x \in \mathbb{R}^N,$$
then $f(x)$ has a **Lipschitz continuous gradient** with Lipschitz constant $L$.

So twice differentiability with **bounded curvature** is sufficient, but not necessary, for a function to have Lipschitz continuous gradient.

Proof. Using **Taylor's theorem** and the **triangle inequality** and the definition of **spectral norm**:
$$
\begin{aligned}
\|\nabla f(x) - \nabla f(z)\|_2 &= \left\| \int_0^1 \nabla^2 f(x + \tau(z - x)) \, \mathrm{d}\tau \, (x - z) \right\|_2 \\
&\leq \left( \int_0^1 \left\| \nabla^2 f(x + \tau(z - x)) \right\|_2 \mathrm{d}\tau \right) \|z - z\|_2 \\
&\leq \left( \int_0^1 L \, \mathrm{d}\tau \right) \|x - z\|_2 = L \|x - z\|_2.
\end{aligned}
$$

Example. $f(x) = \frac{1}{2} \|Ax - y\|_2^2 \implies \nabla^2 f = A'A$ so $\|\nabla^2 f\|_2 = \|A'A\|_2$.

The Lipschitz constant for the gradient of $f(x) \triangleq x' \begin{bmatrix} 1 & 2 \\ 2 & 4 \end{bmatrix} x$ is:

A: 1     B: 2     C: 4     D: 5     E: 10     ??

**Convexity and Hessian**

If $f(x)$ is twice differentiable, then $f(x)$ is **convex** iff its Hessian $\nabla^2 f(x)$ is **positive semidefinite** for all $x$.

Example. $f(x) = \frac{1}{2} \|Ax - y\|_2^2 \implies \nabla^2 f(x) = A'A \succeq 0 \implies f(x)$ is convex for any $A$.

One can also show convexity of this $f(x)$ directly from the definition of convex functions.

If $f(x)$ is twice differentiable, then $f(x)$ is **strictly convex** if its Hessian $\nabla^2 f(x)$ is **positive definite** for all $x$.

Why not only if? $f(x) = x^4$ is strictly convex but $\ddot{f}(0) = 0$.

Example. If $A$ has full column rank, then $A'A$ is positive definite so $f(x) = \frac{1}{2} \|Ax - y\|_2^2$ is strictly convex.

Suppose $A$ is a wide matrix and consider the regularized LS cost function $f(x) \triangleq \frac{1}{2} \|Ax - y\|_2^2 + \beta \|x\|_2^2$ where $\beta > 0$. Then $f$ is a **strictly convex** function. (?)

A: True                                      B: False                            ??

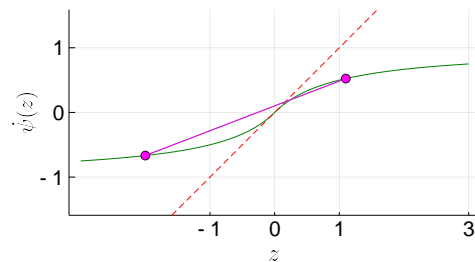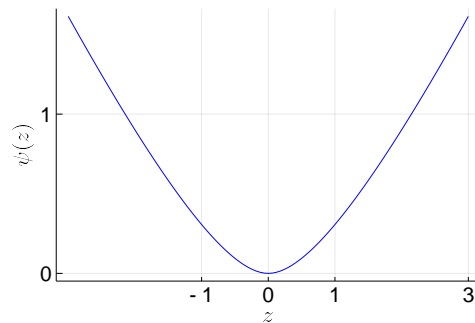Example. The **Fair potential** used in many imaging applications [4] [5] is

$$\psi(x) = |x| - \log(1 + |x|),$$

and has the property of being roughly quadratic for $x \approx 0$ and roughly like $|x|$ for $|x| \gg 0$.
For this function:

$$\dot{\psi}(x) = \frac{x}{1 + |x|} \text{ and } \ddot{\psi}(x) = \frac{1}{(1 + |x|)^2} \leq 1,$$

so the Lipschitz constant of the derivative of $\psi(\cdot)$ is 1.
Furthermore, its second derivative is nonnegative so it is a convex function.

Is the Fair potential $\psi$ itself **Lipschitz continuous**?
A: No          B: Yes with $L = 1/2$          C: Yes with $L = 1$          D: Yes with $L = 2$          ??

**GD convergence theorem**

- If **convex** function $f(\boldsymbol{x})$ has a (not necessarily unique) minimizer $\hat{\boldsymbol{x}}$ for which

$$-\infty < f(\hat{\boldsymbol{x}}) \leq f(\boldsymbol{x}), \quad \forall \boldsymbol{x} \in \mathbb{R}^N,$$

- the gradient of $f(\boldsymbol{x})$ is **Lipschitz continuous**, *i.e.*,

$$\|\nabla f(\boldsymbol{x}) - \nabla f(\boldsymbol{z})\|_2 \leq L \|\boldsymbol{x} - \boldsymbol{z}\|_2, \forall \boldsymbol{x}, \boldsymbol{z} \in \mathbb{R}^N,$$

- the **step size** $\alpha$ is chosen such that

$$0 < \alpha < 2/L,$$

then the GD iteration

$$\boldsymbol{x}_{k+1} = \boldsymbol{x}_k - \alpha \nabla f(\boldsymbol{x}_k)$$

has the following convergence properties [6, p. 207]:
- the cost function is non-increasing: $f(\boldsymbol{x}_{k+1}) \leq f(\boldsymbol{x}_k)$
- for any minimizer $\hat{\boldsymbol{x}}$, the distance to that minimizer is non-increasing: $\|\boldsymbol{x}_{k+1} - \hat{\boldsymbol{x}}\| \leq \|\boldsymbol{x}_k - \hat{\boldsymbol{x}}\|$
- the sequence $\{\boldsymbol{x}_k\}$ **converges** to a minimizer of $f(\cdot)$.
- For $0 < \alpha \leq 1/L$, the cost function decrease is bounded by [2]:

$$f(\boldsymbol{x}_k) - f(\hat{\boldsymbol{x}}) \leq \frac{L \|\boldsymbol{x}_0 - \hat{\boldsymbol{x}}\|_2^2}{2} \max\left(\frac{1}{2k\alpha + 1}, (1 - \alpha)^{2k}\right).$$

This upper bound is conjectured to also hold for $1/L < \alpha < 2/L$ [7].

**Convex sets** _____

Define.  A nonempty set $\mathcal{C}$ in a vector space $\mathcal{V}$ is a **convex set** iff

$$\boldsymbol{x}, \boldsymbol{z} \in \mathcal{C} \implies \alpha \boldsymbol{x} + (1 - \alpha)\boldsymbol{z} \in \mathcal{C}, \quad \forall 0 \leq \alpha \leq 1.$$

The linear combination $\alpha \boldsymbol{x} + (1 - \alpha)\boldsymbol{z}$ for any $0 \leq \alpha \leq 1$ is a **convex combination**.

Example. Any subspace $\mathcal{S}$ in a vector space $\mathcal{V}$ is convex.

Example. The **nonnegative orthant** $\mathbb{R}^N_+ \triangleq \left\{ \boldsymbol{x} \in \mathbb{R}^N : \boldsymbol{x} \geq \boldsymbol{0} \right\}$ is convex.

Example. For any norm, the ball of radius $r > 0$ $\left\{ \boldsymbol{x} : \|\boldsymbol{x}\| \leq r \right\}$ is convex. (HW)

Finding the nearest point in a convex set is an important operation. $\mathcal{P}_{\mathcal{C}}(\cdot)$ denotes the "**projection**" of its argument onto the closest point in $\mathcal{C}$:

$$\mathcal{P}_{\mathcal{C}}(\boldsymbol{z}) \triangleq \min_{\boldsymbol{x} \in \mathcal{C}} \|\boldsymbol{x} - \boldsymbol{z}\|$$

For $\mathcal{C} = \left\{ \boldsymbol{x} \in \mathbb{R}^N : \|\boldsymbol{x}\|_\infty \leq 5 \right\}$, which of these is the projection of a point $\boldsymbol{z} \in \mathbb{R}^N$ onto $\mathcal{C}$?
A: `min.(x,5)`   B: `min.(abs.(x),5)`   C: `min.(abs.(x),5).*sign(x)`   D: None of these
??

**Gradient projection method** _____

In many applications, we seek the minimizer of a convex function $f(\boldsymbol{x})$ over a convex set $\mathcal{C}$:

$$\hat{\boldsymbol{x}} = \arg\min_{\boldsymbol{x}\in\mathcal{C}} f(\boldsymbol{x}).$$

This is called **constrained optimization**.

> Fact. The conclusions about convergence for GD given above also hold for the more **gradient projection** (**GP**) method
> $$\boldsymbol{x}_{k+1} = \mathcal{P}_{\mathcal{C}}(\boldsymbol{x}_k - \alpha \nabla f(\boldsymbol{x}_k)).$$

<u>Example.</u> (HW) NNLS, where $f(\boldsymbol{x}) = \frac{1}{2}\|\boldsymbol{A}\boldsymbol{x} - \boldsymbol{y}\|_2^2$ and $\mathcal{C} = \mathbb{R}_+^N$.

<u>Example.</u> If $f(\boldsymbol{x}) = \frac{1}{2}\left\| \begin{bmatrix} 4 & 0 \\ 0 & 2 \end{bmatrix}\boldsymbol{x} - \begin{bmatrix} 4 \\ -2 \end{bmatrix}\right\|_2^2$ and we apply GP with $\alpha = 1/L$ and $\boldsymbol{x}_0 = \boldsymbol{0}$ then $\boldsymbol{x}_1 = ?$

A: $(0,0)$　　　　　B: $(1,0)$　　　　　C: $(1,-1)$　　　　　D: $(0,-1)$　　　　　E: $(1,1)$　　　　　??
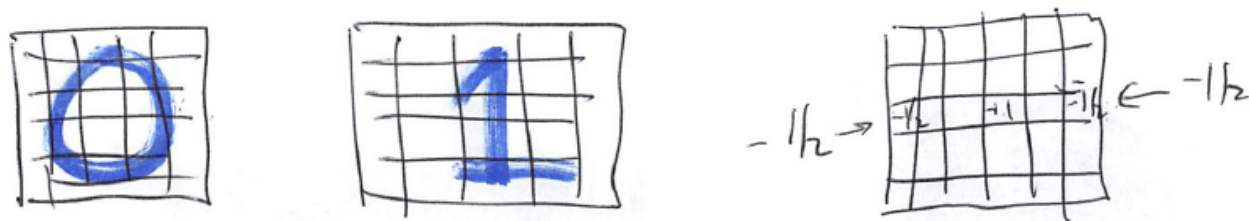
## 8.4 Machine learning via logistic regression for binary classification

In a **binary classification** problem we are given training feature vectors $\{\boldsymbol{v}_i\} \in \mathbb{R}^N$ and corresponding binary labels $\{y_i = \pm 1 : i = 1, \ldots, M\}$. Given these training pairs $(\boldsymbol{v}_i, y_i)$, the training goal is to learn weights $\boldsymbol{x} \in \mathbb{R}^N$ such that

- $\langle \boldsymbol{x}, \boldsymbol{v}_i \rangle > 0$ if $y_i = +1$ and
- $\langle \boldsymbol{x}, \boldsymbol{v}_i \rangle < 0$ if $y_i = -1$,

*i.e.*, such that $\langle \boldsymbol{x}, y_i \boldsymbol{v}_i \rangle > 0$, so that a reasonable classifier is: $\operatorname{sign}(\langle \boldsymbol{x}, \boldsymbol{v} \rangle)$.

Example. Classifying gray-scale images of hand-written digits 0 & 1.



In this example, we "hand crafted" the weights $\boldsymbol{x}$, hoping that $\operatorname{sign}(\langle \boldsymbol{x}, \boldsymbol{v} \rangle)$ will be correct for a test image $\boldsymbol{v}$.
For good statistical performance, we want to learn $\boldsymbol{x}$ from training data.
Learning is especially important for harder problems like classifying 5 and 8.

To learn the weights $x$, we can minimize a cost function with a regularization parameter $\beta \geq 0$:

$$\hat{x} = \arg\min_{x} f(x), \qquad f(x) = \sum_{i=1}^{M} h(y_i \langle x,\ v_i \rangle) + \beta \frac{1}{2} \|x\|_2^2 = \mathbf{1}'_M h.(Ax) + \beta \frac{1}{2} \|x\|_2^2,$$

where the $m$th row of the $M \times N$ matrix $A$ is $y_m v_m^T$, i.e., $A \triangleq \begin{bmatrix} y_1 v_1^T \\ \vdots \\ y_M v_M^T \end{bmatrix}$.

A regularization term like $\|x\|_2^2$ is especially important in the typical case where the feature vector dimension $N$ is large relative to the sample size $M$. The 1-norm is also often used; see EECS 598.

We want the function $h$ to discourage incorrect classification of the training data.

- The 0-1 loss function $h(z) = \mathbb{I}_{\{z \leq 0\}}$ is natural because it counts how many training samples are misclassified, but it is nonconvex and nondifferentiable, so very difficult to use for optimization. Instead one usually uses **surrogate loss functions**.
- The **hinge loss** function $h(z) = \max(1 - z, 0)$ is related to the soft-margin **support-vector machine** (**SVM**) and is convex, but non-differentiable. (See EECS 598.)
- The **logistic** loss function is convex and has a **Lipschitz continuous** derivative:
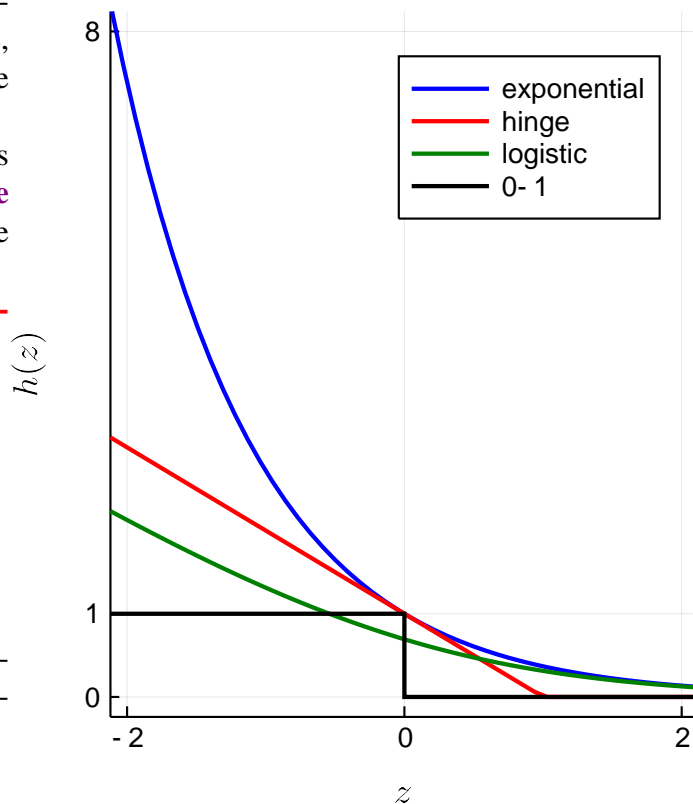
$$h(z) = \log\left(1 + e^{-z}\right)$$
$$\dot{h}(z) = -1/\left(e^z + 1\right)$$
$$\ddot{h}(z) = \frac{e^z}{\left(e^z + 1\right)^2} \in \left(0, \frac{1}{4}\right].$$

It is suitable for gradient-based methods.
- The exponential loss function is convex and differentiable, but its derivative is not Lipschitz continuous. We will not consider it further.

## Loss functions (surrogates)

For the logistic loss, the cost function is not quadratic, but it does have a Lipschitz continuous gradient. For gradient-based optimization, we need the cost function gradient:

$$\underbrace{\nabla f(\boldsymbol{x})}_{N \times 1} = \nabla \left( \boldsymbol{1}_M' h.(\boldsymbol{A}\boldsymbol{x}) + \beta \frac{1}{2} \|\boldsymbol{x}\|_2^2 \right) = \left( \sum_{m=1}^{M} \nabla h(\boldsymbol{A}_{m,:}\boldsymbol{x}) \right) + \beta \boldsymbol{x}$$

$$= \left( \sum_{m=1}^{M} \boldsymbol{A}_{m,:}' \dot{h}(\boldsymbol{A}_{m,:}\boldsymbol{x}) \right) + \beta \boldsymbol{x} = \boldsymbol{A}' \dot{h}.(\boldsymbol{A}\boldsymbol{x}) + \beta \boldsymbol{x}.$$

The cost function **Hessian matrix** is:

$$\nabla^2 f(\boldsymbol{x}) = \nabla^T \nabla f(\boldsymbol{x}) = \nabla^T \left( \sum_{m=1}^{M} \boldsymbol{A}_{m,:}' \dot{h}(\boldsymbol{A}_{m,:}\boldsymbol{x}) + \beta \boldsymbol{x} \right) = \sum_{m=1}^{M} \boldsymbol{A}_{m,:}' \ddot{h}(\boldsymbol{A}_{m,:}\boldsymbol{x}) \boldsymbol{A}_{m,:}' + \beta \boldsymbol{I}$$

$$= \boldsymbol{A}' \operatorname{diag}\left\{ \ddot{h}.(\boldsymbol{A}\boldsymbol{x}) \right\} \boldsymbol{A} + \beta \boldsymbol{I} = \underbrace{\boldsymbol{A}' \boldsymbol{D}(\boldsymbol{x}) \boldsymbol{A}}_{} + \beta \underbrace{\boldsymbol{I}}_{}, \qquad \boldsymbol{D}(\boldsymbol{x}) \triangleq \operatorname{diag}\left\{ \ddot{h}.(\boldsymbol{A}\boldsymbol{x}) \right\} \quad \rule{2em}{0pt}.$$

$f$ is a **strictly convex** function when $\psi$ is the **logistic** loss function and $\beta > 0$. (?)
A: True                                   B: False                                   ??

A Hessian majorizer leads to Lipschitz constant bound for $\nabla f(\boldsymbol{x})$ that is useful for many iterative algorithms: (HW)

$$\nabla^2 f(\boldsymbol{x}) = \boldsymbol{A}'\boldsymbol{D}\boldsymbol{A} + \beta\boldsymbol{I} \preceq \boldsymbol{A}'\left(\frac{1}{4}\boldsymbol{I}\right)\boldsymbol{A} + \beta\boldsymbol{I} = \frac{1}{4}\boldsymbol{A}'\boldsymbol{A} + \beta\boldsymbol{I} \preceq \frac{1}{4}\|\boldsymbol{A}'\boldsymbol{A}\|_2\boldsymbol{I} + \beta\boldsymbol{I} = \left(\frac{1}{4}\|\boldsymbol{A}\|_2^2 + \beta\right)\boldsymbol{I}.$$

Thus, a Lipschitz constant bound that depends on the training data $\boldsymbol{A}$ and the regularization parameter $\beta$ is:

$$L \triangleq \frac{1}{4}\|\boldsymbol{A}\|_2^2 + \beta.$$

Notice the process here: typically we do not try to find $L$ numerically. Instead we analyze either $\nabla f$ or $\nabla^2 f$ and use properties and inequalities (analytically, on paper) to find a suitable $L$ expressed in terms of the problem variables (in this case $\boldsymbol{A}$ and $\beta$).
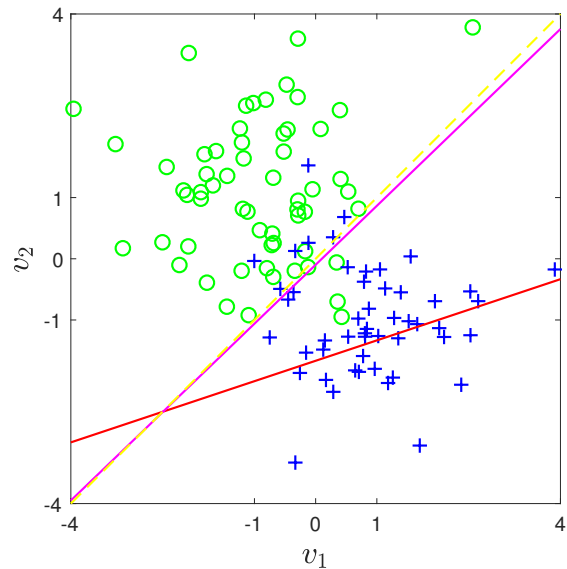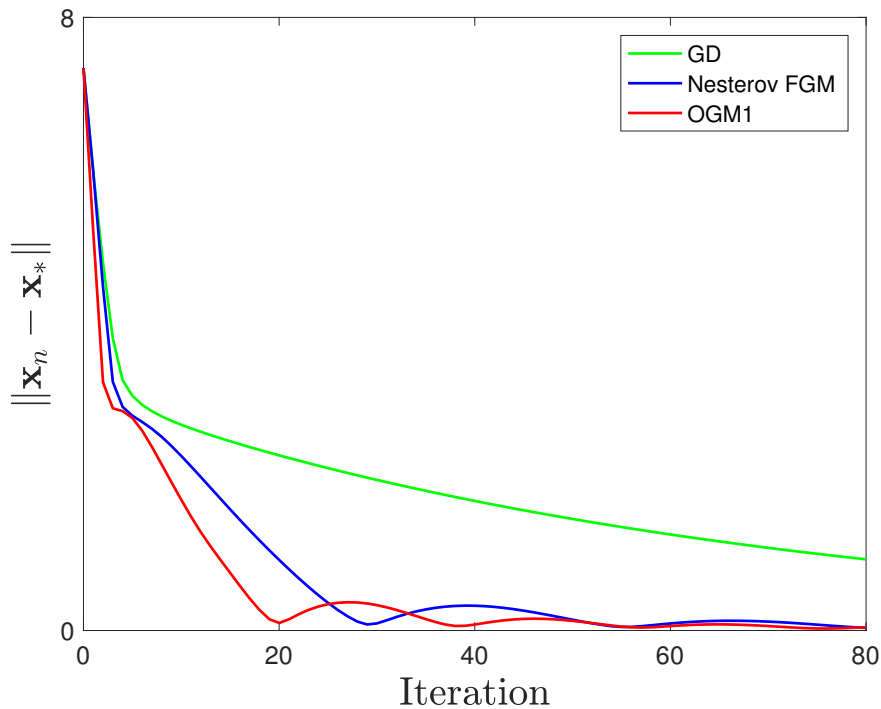
Practical implementation:
- Normalizing each row of $\boldsymbol{A}$ to unit norm can help keep $e^z$ from overflowing.
- Tuning $\beta$ should use **cross validation** or other such tools from machine learning.
- The cost function is convex with Lipschitz gradient, so it is well-suited for Nesterov's fast gradient method (or OGM).
- When feature dimension $N$ is very large, seeking a sparse weight vector $\boldsymbol{x}$ may be preferable. For that, replace the Tikhonov regularizer $\|\boldsymbol{x}\|_2^2$ with $\|\boldsymbol{x}\|_1$ and then use FISTA (or POGM [8]) for optimization.

This **logistic regression** approach to classification will be explored in a task sheet.
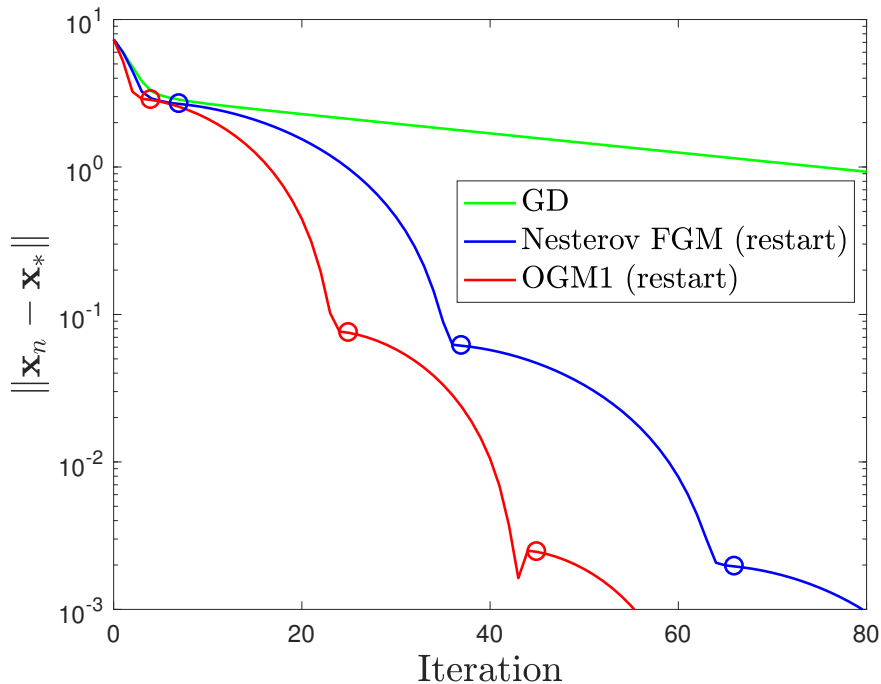
## Numerical Results: logistic regression

Labeled training data (green and blue points);
initial decision boundary (red);
final decision boundary (magenta);
ideal boundary (yellow).

**Numerical Results: convergence rates**



OGM faster than FGM in early iterations...

**Adaptive restart of accelerated GD** _____



FGM restart, O'Donoghue & Candès, 2015. [9]
OGM restart [10]

---

## 8.5 Summary

This chapter barely scratches the surface of the field of optimization algorithms, but it illustrates how crucial matrix methods are to that field.

### Bibliography

[1]     W. Zuo and Z. Lin. "A generalized accelerated proximal gradient approach for total-variation-based image restoration". In: *IEEE Trans. Im. Proc.* 20.10 (Oct. 2011), 2748–59 (cit. on p. 8.18).

[2]     Y. Drori and M. Teboulle. "Performance of first-order methods for smooth convex minimization: A novel approach". In: *Mathematical Programming* 145.1-2 (June 2014), 451–82 (cit. on pp. 8.19, 8.30).

[3]     D. Kim and J. A. Fessler. "Optimized first-order methods for smooth convex minimization". In: *Mathematical Programming* 159.1 (Sept. 2016), 81–107 (cit. on pp. 8.19, 8.25).

[4]     R. C. Fair. "On the robust estimation of econometric models". In: *Ann. Econ. Social Measurement* 2 (Oct. 1974), 667–77 (cit. on p. 8.29).

[5]     K. Lange. "Convergence of EM image reconstruction algorithms with Gibbs smoothing". In: *IEEE Trans. Med. Imag.* 9.4 (Dec. 1990). Corrections, T-MI, 10:2(288), June 1991., 439–46 (cit. on p. 8.29).

[6]     B. T. Polyak. *Introduction to optimization*. New York: Optimization Software Inc, 1987 (cit. on p. 8.30).

[7]     A. B. Taylor, J. M. Hendrickx, and Francois Glineur. "Smooth strongly convex interpolation and exact worst-case performance of first- order methods". In: *Mathematical Programming* 161.1 (Jan. 2017), 307–45 (cit. on p. 8.30).

[8]     A. B. Taylor, J. M. Hendrickx, and Francois Glineur. "Exact worst-case performance of first-order methods for composite convex optimization". In: *SIAM J. Optim.* 27.3 (Jan. 2017), 1283–313 (cit. on p. 8.37).

[9]     B. O'Donoghue and E. Candes. "Adaptive restart for accelerated gradient schemes". In: *Found. Comp. Math.* 15.3 (June 2015), 715–32 (cit. on p. 8.40).

[10]   D. Kim and J. A. Fessler. "Adaptive restart of the optimized gradient method for convex optimization". In: *J. Optim. Theory Appl.* 178.1 (July 2018), 240–63 (cit. on p. 8.40).