# Krylov Methods

## A (very) incomplete introduction

Riccardo Rossi

# Real basics…

Let's begin with a question: imagine that i have a given vector
$$r(x) := b - Ax$$

**What does it mean to require $r(x) = 0$ ?**

# Real basics...

Let's begin with a question: imagine that i have a given vector (of size N)

$$r(x) := b - Ax$$

**What does it mean to require $r(x) = 0$ ?**

**ANSWER** – it means that all of its components are zero...

...something that we can write as

$$r(x) \cdot e_i = 0 \quad \forall i$$

**THAT IS: a vector is zero if it is ortogonal to all of the vectors of a basis of the space it lives in**

# Real basics...

But which basis? **ANY !!**

**THAT IS:** if $\{v_0 \dots v_N\}$ form a basis of $R^N$ than a vector is zero if and only if $\boldsymbol{r(x)} \cdot \boldsymbol{v_i} = \boldsymbol{0} \quad \forall \boldsymbol{i}$

**IN A NUTSHELL**, Krylov methods are all about **constructing a basis** (that grows with the number of iterations until spanning the entire $R^N$) and **making the residual to be Orthogonal to such basis**.

# Krylov Subspace

Of course we have complete freedom of "incrementally constructing" a space that eventually describes the entire $R^N$.

The space employed by Krylov techniques is known as "**Krylov Subspace**", defined as

$$K_i(\boldsymbol{A}, \boldsymbol{b}) = span\{\boldsymbol{b}, \boldsymbol{Ab}, \boldsymbol{A^2 b} \dots, \boldsymbol{A^i b}\}$$

Provided that the matrix A is invertible, and that b has a component wrt all of the eigenvectors of A, such space will eventually grow with the iterations until coinciding with R^n

Krylov techniques differe of the construction and properties of a basis of $K_i(A, b)$.

A property shared by all of the methods in the Krylov family is that **they will converge** (in exact algebra) **in at most N iterations** (although they will typically converge way before)

# Minimizing residual along a line

A tool very frequently used, is the **minimization of the residual along a given direction**. Let's imagine that we have a starting solution $x_0$ and a search direction identified by a unit vector $v$

The idea is to choose a new $x := x_0 + \alpha v$ such that $\|r(x)\|^2$ is minimal in a choosen norm. One way to accomplish this, is to make the residual to be ortogonal to the direction $v$, that is, to require that $v \cdot r(x) = 0$

Using the definition we get
$$0 = v \cdot r(x) = v \cdot (b - Ax_0 - \alpha Av) = v \cdot (r(x_0) - \alpha Av)$$

Solving for alpha
$$\alpha = \frac{v \cdot r(x_0)}{v \cdot Av}$$

# Special Case of SPD matrices

If $A$ is **SPD**, we can define a functional $\Psi(x) := x^t b - \frac{1}{2} x^t A x$ (which we will use both for CG and SD)

It is easy to see that:

1. such that $A x_{ex} = b$ is the (only) minimum of $\Psi$ (easy to prove since $A$ is SPD, hence $v^t A v > 0$ for any non zero $v$)

2. The gradient of the function is $\nabla \Psi = b - A x = r(x)$ (obviously zero in $\nabla \Psi (x_{ex}) = 0$)

# Steepest Descent
(not a member of the Krylov Family)

The "Steepest Descent" Is the first idea one may have. It Works as follows:

1 choose a starting point $x_0$ and use $\nabla\Psi$ as search direction

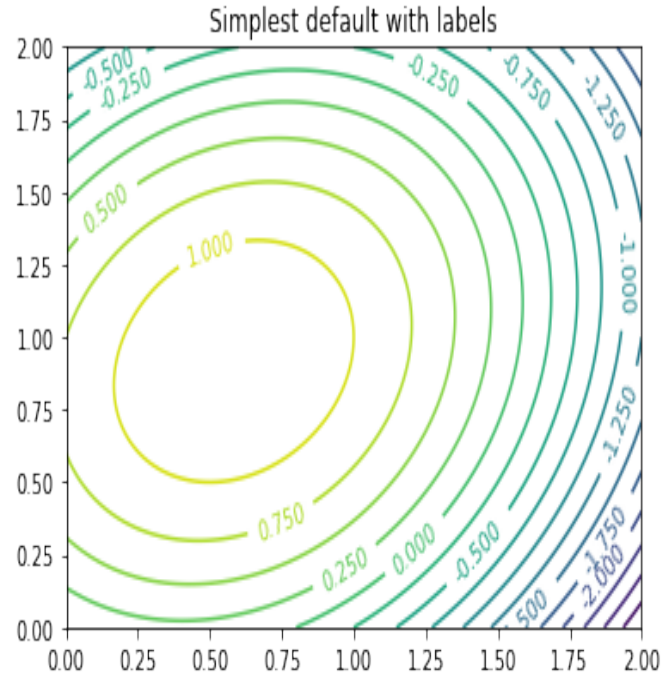$$v := \frac{\nabla\Psi\,(x_0)}{\|\nabla\Psi(x_0)\|} = \frac{r(x_0)}{\|r(x_0)\|}$$

2 evaluate the minimum in the direction of $v$ starting from $x_0$, that is, compute $x_1 = x_0 + \alpha\,v$ (with $\alpha := -\frac{v^t r(x_0)}{v^t A v}$, using the minimization formula in the previous slides)
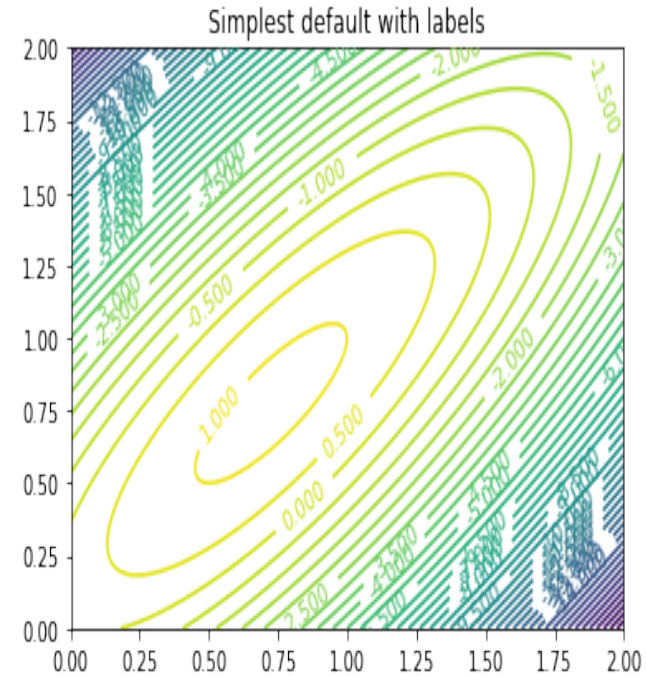
3 repeat until $\frac{r(x_i)}{\|b\|} < \epsilon$

SIMPLE but ... **may be very slow** (no guarantees it converges in k)

# Why it takes long?



Plot of $\Psi(x) = x^t b - 1/2 x^t A x$

$A = \begin{pmatrix} 3 & -1 \\ -1 & 3 \end{pmatrix}$  $b = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$  $k(A) = 2$



Plot of $\Psi(x) = x^t b - 1/2 x^t A x$

$A = \begin{pmatrix} 11 & -9 \\ -9 & 11 \end{pmatrix}$  $b = \begin{pmatrix} 1 \\ 2 \end{pmatrix}$  $k(A) = 10$
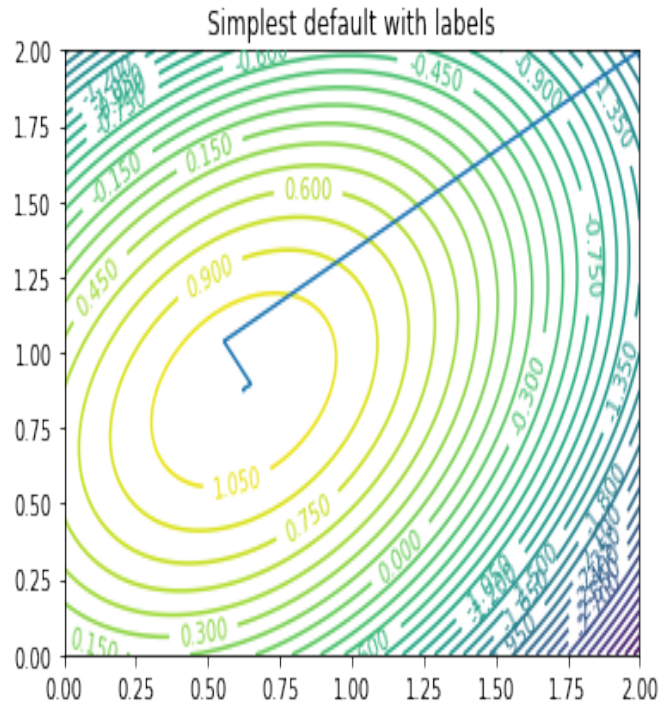
# Why it takes long?



Plot of $\Psi(x) = x^t b - 1/2 x^t A x$

$$A = \begin{pmatrix} 3 & -1 \\ -1 & 3 \end{pmatrix} \quad b = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad k(A) = 2$$
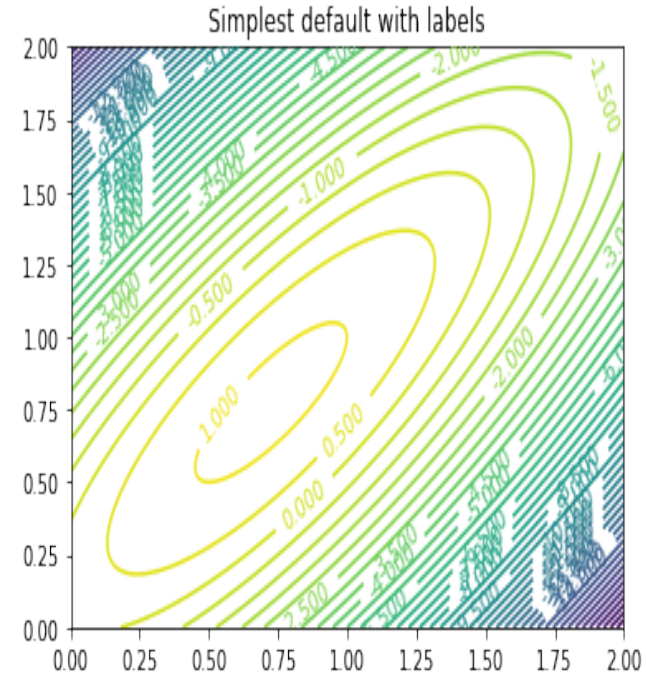
Plot of $\Psi(x) = x^t b - 1/2 x^t A x$

$$A = \begin{pmatrix} 11 & -9 \\ -9 & 11 \end{pmatrix} \quad b = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad k(A) = 10$$

# Why it takes long?



Plot of $\Psi(x) = x^t b - 1/2 x^t A x$

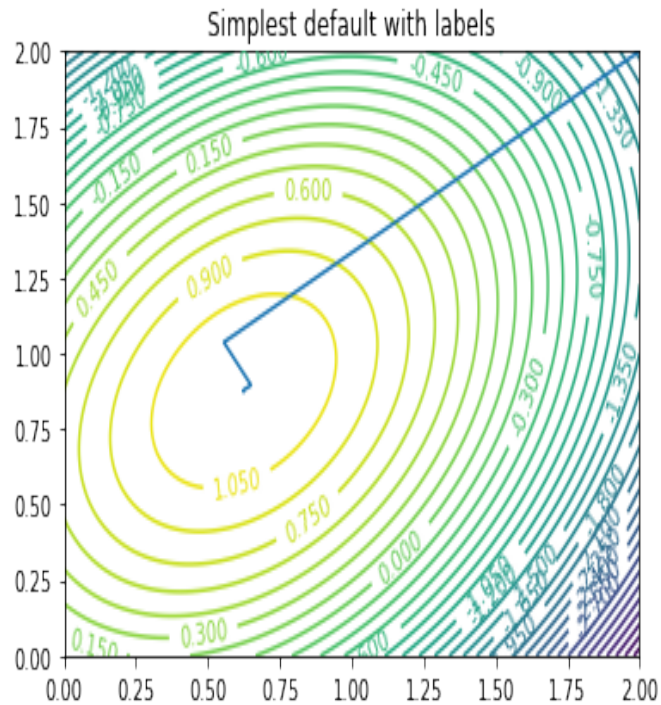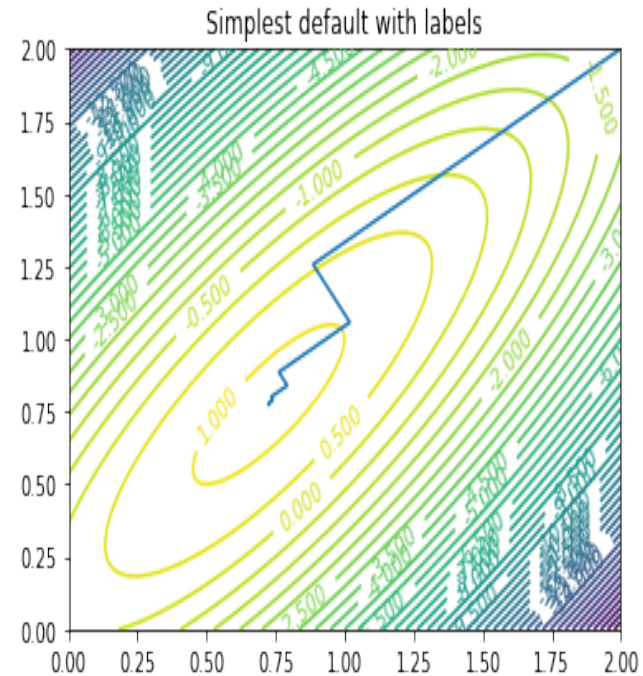$$A = \begin{pmatrix} 3 & -1 \\ -1 & 3 \end{pmatrix} \quad b = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad k(A) = 2$$



Plot of $\Psi(x) = x^t b - 1/2 x^t A x$

$$A = \begin{pmatrix} 11 & -9 \\ -9 & 11 \end{pmatrix} \quad b = \begin{pmatrix} 1 \\ 2 \end{pmatrix} \quad k(A) = 10$$

# Conjugate Gradient

The "Conjugate Gradient" method can be understood as an improvement over the "steepest descent". It Works as follows:

1 - choose a starting point $x_0$ and use $v_0 = r(x_0) = r_0$ as first search direction (**as the steepest descent!**)

2 $x_i = x_0, v_i = v_0$

3 - $x_{i+1} = x_i + \alpha \, v_i$ (with $\alpha := -\frac{v_0{}^t r(x_0)}{v_0{}^t A v_0}$)

4 - choose a new update direction as $v_{i+1} := r(x_{i+1}) + \sum_{k=0}^{i} \beta_{ik} v_k$

where the $\beta_{ik}$ are chosen so that $A v_{i+1} \cdot v_k = 0 \quad \forall k$
(**A-orthogonality instead of orthogonality!!**)

5 go back to 3, until $\frac{r(x_{i+1})}{\|b\|} < \epsilon$

**The key difference wrt the steepest descent is in Step 4**, in the choice of the

# A-orthogonality

Let's focus on the orthogonalization step, taking in mind the equation $v_{i+1} := r(x_{i+1}) + \sum_{k=0}^{i} \beta_{ik} v_k$:

1 – FIRST ITERATION $v_1 := r_1 + \beta_{00} v_0$ A-orthogonality :

$$v_0^t A v_1 = 0 \Rightarrow v_0^t A r_1 + \beta_{10} v_0^t A v_0 = 0 \Rightarrow \beta_{10} = -\frac{v_0^t A r_1}{v_0^t A v_0}$$

2 – SECOND ITERATION $v_2 := r_2 + \beta_{20} v_0 + \beta_{21} v_1$ A-orthogonality:

$$v_0^t A v_2 = 0 \Rightarrow v_0^t A r_2 + \beta_{20} v_0^t A v_0 + \beta_{21} v_0 A v_1 = 0 \Rightarrow \beta_{20} = -\frac{v_0^t A r_2}{v_0^t A v_0^t}$$

$$v_1^t A v_2 = 0 \Rightarrow v_1^t A r_2 + \beta_{20} v_1^t A v_0 + \beta_{21} v_1^t A v_1 = 0 \Rightarrow \beta_{21} = -\frac{v_1^t A r_2}{v_1^t A v_1}$$

note that we used here that $v_0 A v_1 = 0$ and that **A is symmetric**, hence $v_0 A v_1 = v_1 A v_0 = 0$.

Also, **the terms in the denominator are guaranteed to be positive since A is SPD**

3 – OTHER ITERATIONS:

$$v_{i+1} := r_{i+1} + \sum_{k=0}^{i} \beta_{ik} v_k \text{ with } \beta_{(i+1)k} = -\frac{v_k^t A r_{i+1}}{v_k^t A v_k} \forall k \le i$$

# Conjugate Gradient Magic

Let's make some observations:

- by construction $\boldsymbol{r}_{i+1} \cdot \boldsymbol{v}_i = 0$ (that's how we choose $\alpha$).

- It is easy to prove that $\boldsymbol{r}_{i+1} \cdot \boldsymbol{v}_i = 0 \Rightarrow \boldsymbol{r}_{i+1} \cdot \boldsymbol{v}_j = 0 \ \forall\, j \leq i$

- The last sentence can be paraphrased as follows:
$$\boldsymbol{r}_{i+1} \perp \text{span}\{v_0 \ldots v_i\}$$

- now, $\boldsymbol{v}_{i+1} = \boldsymbol{r}(\boldsymbol{x}_{i+1}) + \sum_{k=0}^{i} \beta_{ik} \boldsymbol{v}_k$ hence $\boldsymbol{v}_{i+1}$ **is a linear combination of the previous residuals(\*)**. It follows that
$$\text{span}\{v_0 \ldots v_i\} = \text{span}\{r_0 \ldots r_i\}$$

- But then
$$\boldsymbol{r}_{i+1} \perp \text{span}\{v_0 \ldots v_i\} \Rightarrow \boldsymbol{r}_{i+1} \perp \text{span}\{r_0 \ldots r_i\} \Rightarrow \boldsymbol{r}_{i+1} \cdot \boldsymbol{r}_j = 0 \ \forall j \leq i$$

\*we could actually show that $\text{span}\{r_0 \ldots r_i\} = K_i(A, r_0)$

# Conjugate Gradient Magic

Now the CG magic, is that taking into account that $\boldsymbol{r_{i+1} \cdot r_j} = 0$ $\forall j \le i$ we discover that **many of the $\boldsymbol{\beta_*}$ are actually 0** … hence **no need to store the search vectors**.

Proof: $\boldsymbol{r_{i+1} = b - Ax_{i+1} = b - Ax_i - \alpha Av_i = r_i + \alpha_i Av_i \Rightarrow \alpha_i Av_i = r_{i+1} - r_i}$

Hence

$$\beta_{(i+1)k} = -\frac{v_k^t A r_{i+1}}{v_k^t A v_k} = \frac{r_{i+1}^t A v_k}{v_k^t A v_k} \forall k \le i \qquad \rightarrow \qquad \beta_{(i+1)k} = \frac{r_{i+1}^t r_{k+1} - r_{i+1}^t r_k}{\alpha_i v_k^t A v_k} = \frac{-r_{i+1}^t r_{k+1}}{\alpha_i v_k^t A v_k}$$

Now $k < i \Rightarrow r_k^t r_i = 0$ it follows that the only non zero beta is for $k = i$

$$\beta_{(i+1)i} = \frac{-r_{i+1}^t r_{i+1}}{\alpha_i v_i^t A v_i} = \frac{-r_{i+1}^t r_{i+1}}{v_i^t r_{i+1} - v_i^t r_i} = \frac{r_{i+1}^t r_{i+1}}{v_i^t r_i}$$

The last step is to observe that

$$v_i = r_i + \sum_{k=0}^{i-1} \beta_{ik} v_k \Rightarrow r_i^t v_i = r_i^t r_i + \sum_{k=0}^{i-1} \beta_{ik} r_i^t v_k = r_i^t r_i$$

**Which allows to conclude that**

$$\beta_{(i+1)i} = \frac{r_{i+1}^t r_{i+1}}{r_i^t r_i} \rightarrow v_{i+1} = r_{i+1} + \frac{r_{i+1}^t r_{i+1}}{r_i^t r_i} v_i$$

# CONJUGATE GRADIENT ALGORITHM

1. $r_0 = b - Ax_o \rightarrow v_0 = r_0$

2. $\alpha_i = \dfrac{v_i^t r_i}{v_i^t \textcolor{red}{Av_i}} \rightarrow x_{i+1} = x_i + \alpha_i v_i$

3. $r_{i+1} = r_i - \alpha_i \textcolor{red}{Av_i}$ could also do $r_{i+1} = r_i - Ax_{i+1}$ but that would require one more matrix-vector product

4. $\beta_{(i+1)i} = \dfrac{r_{i+1}^t r_{i+1}}{r_i^t r_i} \rightarrow v_{i+1} = r_{i+1} + \beta_{(i+1)i} v_i$

5. Go back to 2 and loop until convergence

*COST: 1 product $\boldsymbol{Av_i}$ + a few inner products. Guaranteed to converge in N iterations, but will most likely converge much before*

# Convergence Estimates

Considering the condition number $k = k(A) = {^{\lambda_{max}}}/_{\lambda_{min}}$

Convergence estimate of **CG** is:

$$\|e\_i\| < 2\left(\frac{\textcolor{red}{\sqrt{k}} - 1}{\sqrt{k} + 1}\right)^i \|e\_0\|$$

Convergence estimate of **Steepest Descent** was:

$$\|e\_i\| < \left(\frac{\textcolor{red}{k} - 1}{k + 1}\right)^i \|e\_0\|$$

# What if the matrix is not symmetric?

CG can only be applied if the matrix is SPD, however for an arbitrary matrix $A$ (non SPD, and eventually not square) one may solve

$$A^t A x = A^t b$$

Variation of the CG algorithm exist in which $A^t A$ is never computed explicitly (good, since $A^t A$ has more nonzeros than $A$)

**PROBLEM**: the condition number of $k(A^t A) = k(A)^2$ hence the convergence is much slower.

# GMRES ALGORITHM

The most known work horse for solving non-SPD systems is the GMRES algorithm. Although we will not go in detail, the iterate of the gmres is

$$x_{i+1} = x_0 + Vy \qquad V := \begin{pmatrix} v_0 & \cdots & v_i \\ \downarrow & & \downarrow \end{pmatrix}$$

Where $y$ is chosen so to minimize

$$\|r_0 - AVy\|_2$$

The crucial difference with CG is that since A is not symmetric (nor positive definite) **we need to store all of the $v_i$**

The crucial issue, aside of the **memory occupation**, is **how to effectively perform the minimization** of $\|r_0 - AVy\|_2$

**NOTE: Implementing GMRES i way more technical**

**than CG. USE LIBRARIES!**

# Using Krylov methods as Matrix-Free

A very interesting property of laplacian methods is that the actual knowledge of the matrix entries $A_{ij}$ is not needed.

One only needs to evaluate the "action of a matrix onto a vector", that is, how to compute $Av$ for any given vector $v$.

A practical example helps in understanding this better:

$$A := \begin{pmatrix} 2 & -1 & 0 \\ -1 & 2 & -1 \\ 0 & -1 & 2 \end{pmatrix} \quad x = \begin{pmatrix} x_1 \\ x_2 \\ x_3 \end{pmatrix} \rightarrow Ax = \begin{pmatrix} 2x_1 - x_2 \\ -x_1 + 2x_2 - x_3 \\ -x_2 + 2x_3 \end{pmatrix}$$

For the CG it would be sufficient to know that the function

$$f(x) := \begin{pmatrix} 2x_1 - x_2 \\ -x_1 + 2x_2 - x_3 \\ -x_2 + 2x_3 \end{pmatrix}$$

**Can be called whenever A@x is needed**

# TODOS:

1. Implement a CG for the laplacian problems

2. In the test implementation, verify that the expected orthogonality conditions are met (use "large" laplacian matrices)

3. Implement a matrix free solution using scipy's algorithm

4. Use the matrix-free approach to impose that the solution of a Laplacian problem **without dirichlet conditions** is zero on average in the domain.

# References:

The bible of iterative methods:

https://www-users.cs.umn.edu/~saad/IterMethBook_2ndEd.pdf

An in depth dive into the CG

https://www.cs.cmu.edu/~quake-papers/painless-conjugate-gradient.pdf