

LECTURE 7

COORDINATE DESCENT METHODS

- the basic idea is to split the vector x of optimization variables in blocks $x_i \in \mathbb{R}^{n_i}$ and solve the optimization problem for each block separately
- Advantages: cheaper computation and possibility to parallelize it
- Fixed points of the algorithms are not always optimal solutions of the original problem.
For convergence guarantees, the required assumptions depend on the algorithm and might be stronger than those for gradient methods
- They can be interpreted from the point of view of conjugate gradient methods and methods for solving linear systems
- Connections with popular methods used nowadays for large scale problems (e.g. SGD)

ITERATIVE METHODS FOR SYSTEMS OF LINEAR EQUATIONS

$$Ax = b$$

$A \in \mathbb{R}^{n \times n}$, invertible
 $a_{ii} \neq 0, \forall i$
 $b \in \mathbb{R}^n$

$$\sum_{j=1}^n a_{ij} x_j = b_i \quad (\text{i-th row})$$

Gauss-Seidel

- start with an estimate x^0
- each component is updated as

$$x_i^{k+1} = -\frac{1}{a_{ii}} \left[\sum_{j < i} a_{ij} x_j^{k+1} + \sum_{j > i} a_{ij} x_j^k - b_i \right]$$

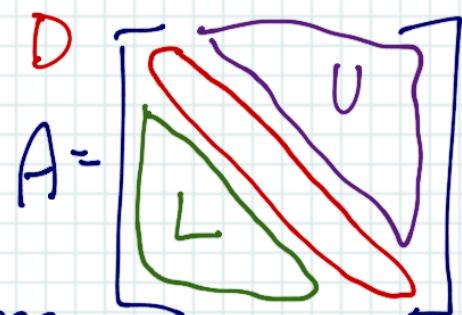
- serial computation (to compute x_i^{k+1} , we require all the previous components of x^k at iteration $k+1$, i.e. $x_1^{k+1}, \dots, x_{i-1}^{k+1}$)

Here is a little confused, why iff?

- Convergence is guaranteed if A is strictly diagonally dominant or $A \in S_{++}^n$.

- In matrix form, the iteration is:

$$x^{k+1} = -D^{-1} \left[Lx^k + Ux^k - b \right]$$



D, L, U are diagonal, strictly lower and upper parts of A , respectively.

Jacobi • start with an estimate x^0

- the component update is

$$x_i^{k+1} = -\frac{1}{a_{ii}} \left[\sum a_{ij} x_j^k - b_j \right]$$

- parallel computation (each equation is solved separately and only variables at the previous iterations are required)
- In matrix form, the iteration is:

$$x^{k+1} = -D^{-1}(L+U)x^k + D^{-1}b$$

- Convergence here is guaranteed for strictly diagonally dominant (sufficient condition), but not for every positive definite A,

Both algorithms can be written as

$$\boxed{x^{k+1} = Mx^k + Gb}$$

where $M = \begin{cases} -(I + D^{-1}L)^{-1} D^{-1}U & \text{Gauss Seidel} \\ -D^{-1}(A - D) & \text{Jacobi} \end{cases}$

- Convergence of the algorithm depends on M (note that the part depending on b has no effect on convergence properties! It can be set to 0 when studying convergence)
 - Specifically, the algorithms converge iff $\rho(M) < 1$
- spectral radius

You can interpret this from the point of view of a discrete time linear system:

$$x^t \rightarrow x^* \text{ for } t \rightarrow \infty$$

iff the state matrix of the system (which here is M) has all the eigenvalues inside the unit circle.

(this holds regardless of b , which only influences x^*)

- Moreover, the rate of convergence is governed by $\rho(M)$. The higher it is, the slower the convergence

how to actually compute the spectral radius?

The COORDINATE MINIMIZATION method generalizes the Gauss-Seidel algorithm for solving linear systems (which we know is equivalent to $\min F(x)$)

where $F(x) = \frac{1}{2} x^T A x - b^T x$, when
 $A \in S^{++}$
to a generic $F(x)$.

The PARALLEL COORDINATE MINIMIZATION is instead closely related to the Jacobi algorithm.

- Convergence guarantees and rates depend on continuity and smoothness of F .
- regularization terms can be added to F to improve convergence properties (for the parallel version, the algorithm may not converge without a regularization term)

COORDINATE GRADIENT DESCENT

gradient descent: $x^{k+1} = x^k - t^k \nabla f(x^k)$

here we take for every block-coordinate i :

$$x^{k+1} = x^k - t_i^k \begin{bmatrix} \nabla f(x^k) \end{bmatrix}_i \quad \text{e:}$$

where

- if $n_i = 1$ (1-D coordinate), then

$$e_i = \begin{bmatrix} 0 \\ \vdots \\ 0 \\ 1 \\ 0 \\ \vdots \\ 0 \end{bmatrix} \quad \text{i-th component } [\nabla f(x^k)]_i \text{ is the i-th component of the gradient}$$

- if $n_i > 1$ (block-coordinates), then the corresponding components of x will be updated.

Example:

$$f(x) = \frac{1}{2} x^T A x - b^T x, \quad x \in \mathbb{R}^6$$

We will have 3 steps:

$$\boxed{i=1} \quad x^{k+1} = x^k - t_i^k \begin{bmatrix} Q_i x^k - b_i \\ 0 \\ \vdots \\ 0 \\ 0 \\ 0 \end{bmatrix}$$

$$\boxed{i=2} \quad x^{k+1} = x^k - t_2^k \begin{bmatrix} 0 \\ Q_2 x^k - b_2 \\ 0 \\ Q_3 x^k - b_3 \\ 0 \\ Q_4 x^k - b_4 \\ 0 \\ 0 \end{bmatrix}$$

$Q_i := i\text{-th row of } A$

$$x = \begin{bmatrix} x_1 \\ x_2 \\ x_3 \\ x_4 \\ x_5 \\ x_6 \end{bmatrix} \quad \begin{array}{l} \xrightarrow{i=1} \\ \xrightarrow{i=2} \\ \xrightarrow{i=3} \end{array}$$

$$\boxed{i=3} \quad x^{k+1} = x^k - t_3^k \begin{bmatrix} 0 \\ 0 \\ 0 \\ 0 \\ 0 \\ Q_5 x^k - b_5 \\ Q_6 x^k - b_6 \end{bmatrix}$$

- t_i^k can be chosen based on L_i -smoothness
of the i -th block. In this case $t_i^k = \frac{1}{L_i}$.

All the criteria studied in gradient methods
are also valid choices in this setting.

- The coordinates can be selected:

■ in "Cyclic fashion": $i_0 = 1$,

$$i_{k+1} = \text{mod}(k, n) + 1$$

■ in "essentially cyclic fashion": for some $T \geq n$,
each coordinate is modified at least once
in every stretch of T iterations

$$\text{union-} \bigcup_{j=0}^T \{i_{k-j}\} = \{1, 2, \dots, n\} \quad \forall k \geq T$$

■ at random (not necessarily with equal
probability)

Application to a regularized problem

$$\min_x h(x)$$

$$h(x) = f(x) + \lambda \underbrace{R(x)}_{\substack{\text{regularization function} \\ (\text{assumed separable})}}$$

regularization parameter
($\lambda > 0$)

examples:

$$\cdot R(x) = \|x\|_1, R_i(x_i) = |x_i|$$

• box constraints,

$$R_i := I_{[l_i, u_i]}(x_i)$$

$$R(x) = \sum_{i=1}^n R_i(x_i)$$

The update rule has 2 steps (assume for simplicity $n_i=1$)

$$\textcircled{1} z_i^k = \underset{y}{\operatorname{argmin}} (y - x_i^k) [\nabla f(x^k)]_i + \frac{1}{2\lambda_k} \|y - x_i^k\|_2^2 + \lambda R_i(y)$$

it plays the role
of stepsize

This is a scalar subproblem aimed at computing a linear approximation of f along the i -th direction while explicitly accounting for R_i .

Notably, this problem has a closed form solution:

$$z_i^k = S_{\lambda_k} (x_i^k - \lambda_k [\nabla f(x^k)]_i)$$

$$\text{where } S_\beta(\tilde{z}) = \min_y \frac{1}{2} \|y - \tilde{z}\|_2^2 + R_i(y)$$

so no explicit search is needed.

$$\textcircled{2} x^{k+1} = x^k + (z_i^k - x_i^k) e_i$$

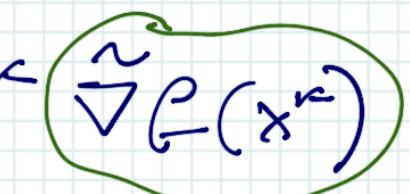
Note: if $R_i = 0$, the algorithm recovers the standard coordinate gradient descent.

STOCHASTIC GRADIENT DESCENT (SGD)

Where is the connection with coordinate methods?

idea of
SGD
algorithms

$$x^{k+1} = x^k - \epsilon^k \tilde{\nabla} f(x^k)$$



realization of a stochastic gradient approximator

Do we always have this n ? use this n to make it unbiased? actually just take one sample instead of normal

Consider for example:

This looks like coordinate gradient descent

$$\tilde{\nabla} f(x^k) = n \left[\begin{array}{c} \nabla f(x^k) \\ \vdots \\ \nabla f(x^k) \end{array} \right]_i e_i$$

direction where $p_i = p(I=i) = \frac{1}{n}$

Then:

$$E[\tilde{\nabla} f(x^k)] = \frac{1}{n} \sum_{i=1}^n \left[\nabla f(x^k) \right]_i e_i = \nabla f(x^k)$$

UNBIASED ESTIMATE

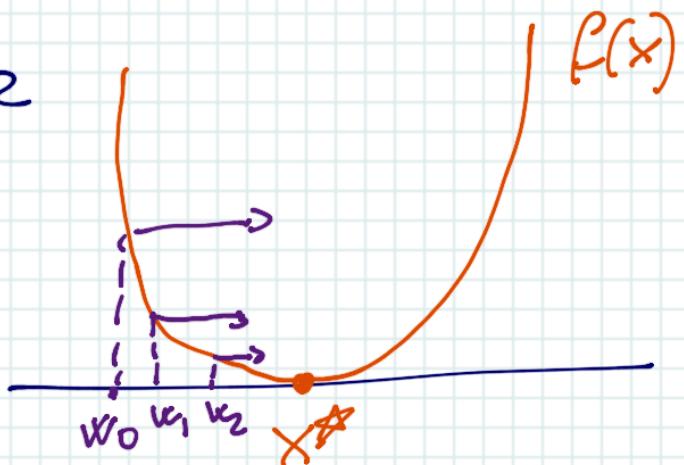
Randomized coordinate gradient descent can thus be seen as a special case of SGD methods.

However, the former has the advantage over general SGD approaches that descent in ϵ can be guaranteed at every iteration.

SGD in action

$$f(x) = \frac{1}{4} \sum_{i=1}^4 (\omega_i^T x - b_i)^2$$

Gradient descent



SGD

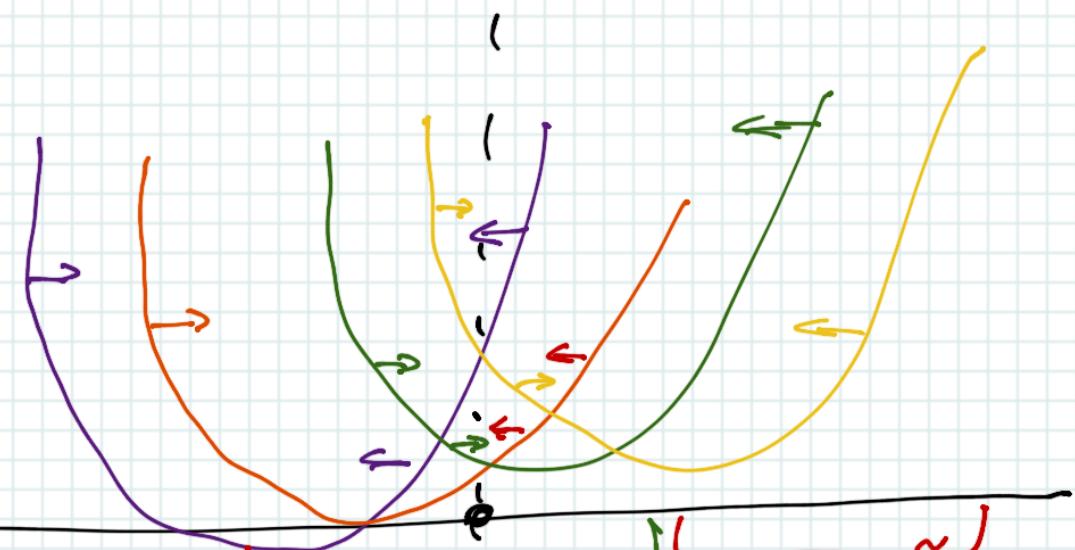
$$f_1(x) = (\omega_1^T x - b_1)^2$$

$$f_2(x) = (\omega_2^T x - b_2)^2$$

$$f_3(x) = (\omega_3^T x - b_3)^2$$

$$f_4(x) = (\omega_4^T x - b_4)^2$$

x^* SGD will minimize the average value



Every ∇f_i in this region will point towards x^* \Rightarrow

Region of confusion
Gradients here might point in both directions

every ∇f_i in this region will point towards x^*

The "confusion" is captured by the variance of the gradients:

$$\sigma^2 = \frac{1}{n} \sum_{i=1}^n \| \nabla f_i(x) - \nabla f(x) \|^2$$

If $\sigma^2 = 0 \rightarrow$ every step goes towards x^*

If $\sigma \ll 1 \rightarrow$ narrow "region of confusion", most steps go in the right direction

If $\sigma \gg 1 \rightarrow$ the algorithm operates in a region where many steps are made away from x^*

Strategies to address variance:

- adapt step-size. Specifically, we need a decreasing step size which is fast enough initially to get in the region of x^* and then slows down to be robust to variance
- minibatch \rightarrow average K local gradients when estimating $\nabla f(x^k)$. Variance is inversely proportional to the batch size.