



The Abdus Salam
**International Centre
for Theoretical Physics**



2361-6

**School on Large Scale Problems in Machine Learning and Workshop
on Common Concepts in Machine Learning and Statistical Physics**

20 - 31 August 2012

**Large Scale Variational Bayesian Inference: Handout for Conjugate Gradients
Algorithm**

Matthias SEEGER

*Laboratory for Probabilistic Machine Learning, EPFL, CH-1015 Lausanne
Switzerland*

Large Scale Variational Bayesian Inference: Handout for Conjugate Gradients Algorithm

Matthias Seeger
Laboratory for Probabilistic Machine Learning
Ecole Polytechnique Fédérale de Lausanne

<http://lapmal.epfl.ch/>

Abstract

The idea behind the conjugate gradients algorithm, as shown in the lecture, is to maintain conjugate search directions, so that improvements realized in some iteration are not lost later on. This idea alone leads to the famous algorithm, but the derivation is too convoluted to present in a lecture. I provide it here. I'd probably not have written this note, were there a readable, yet still concise derivation out there. I took the motivation of CG from [1], but filled in the gaps.

Recall the setting. \mathbf{A} is symmetric positive definite. CG minimizes the quadratic

$$q(\mathbf{x}) = \frac{1}{2} \mathbf{x}^T \mathbf{A} \mathbf{x} - \mathbf{b}^T \mathbf{x}, \quad \mathbf{g}(\mathbf{x}) = \nabla q(\mathbf{x}) = \mathbf{A} \mathbf{x} - \mathbf{b}, \quad \mathbf{x} \in \mathbb{R}^n.$$

We start with some \mathbf{x}_0 . In iteration k , starting from \mathbf{x}_{k-1} , we pick a search direction \mathbf{d}_k , possibly using the gradient $\mathbf{g}_{k-1} = \mathbf{g}(\mathbf{x}_{k-1})$, then search along that line for the minimum point: $\mathbf{x}_k = \mathbf{x}_{k-1} + \alpha_k \mathbf{d}_k$. Be careful not to confuse search *directions* and *gradients*. They are the same only in the method of steepest descent, which CG is certainly different from (except for $\mathbf{A} = \mathbf{I}$). Annoyingly for students, CG is called *conjugate gradients*, it should be called *conjugate directions*.

Remember why steepest descent is not a good idea. Since

$$\frac{dq(\mathbf{x}_{k-1} + \alpha \mathbf{d}_k)}{d\alpha} = \mathbf{g}(\mathbf{x}_{k-1} + \alpha \mathbf{d}_k)^T \mathbf{d}_k = 0$$

at $\alpha = \alpha_k$ (line minimum), then $\mathbf{g}_k^T \mathbf{d}_k = 0$. Now, if $\mathbf{d}_k = -\mathbf{g}_{k-1}$, this means that subsequent search directions are orthogonal. And this is a really bad idea for minimizing $q(\mathbf{x})$ if \mathbf{A} has eigenvalues of widely different sizes (remember the zig-zagging problem from the lecture). A much better idea is to select \mathbf{d}_k in a way that renders *new* gradients, anywhere along the line searched, orthogonal to *old* directions:

$$0 = \mathbf{g}(\mathbf{x}_k + \alpha \mathbf{d}_{k+1})^T \mathbf{d}_k = \mathbf{g}_k^T \mathbf{d}_k + \alpha \mathbf{d}_{k+1}^T \mathbf{A} \mathbf{d}_k.$$

Since $\mathbf{g}_k^T \mathbf{d}_k = 0$ in any case (line minimization condition), we therefore need *conjugate* (rather than orthogonal) directions: $\mathbf{d}_j^T \mathbf{A} \mathbf{d}_k = 0$ for $j \neq k$.

If somebody gave us a conjugate set of directions, starting with $\mathbf{d}_1 = -\mathbf{g}_0$, we were done. However, at least to me, it is not at all obvious how to obtain such a basis in a feasible way. The conjugate gradients algorithm is a neat method for doing just that, with a *single matrix-vector multiplication* (MVM) with \mathbf{A} per iteration, requiring $O(n)$ storage only. And all that is needed in order to derive it, is stated above. It is still a bit of magic, so please read on. By the way, CG is really a cornerstone of numerical mathematics. If its three-term recurrence is interpreted differently, the Lanczos algorithm is obtained, the most important method for computing eigenvectors of very large matrices.

The idea behind the derivation of CG is to establish a number of relationships between the gradients \mathbf{g}_k , the directions \mathbf{d}_k , and the step sizes α_k (for different k values), starting from $\mathbf{g}_k^T \mathbf{d}_k = 0$ and $\mathbf{x}_k = \mathbf{x}_{k-1} + \alpha_k \mathbf{d}_k$ only, which imply that all directions are mutually conjugate, and all gradients are mutually orthogonal (yes, to get you completely confused!). There is a single *ansatz*, namely that $\mathbf{d}_{k+1} = -\mathbf{g}_k + \beta_k \mathbf{d}_k$, involving some further scalars β_k . A formally clean way is to use induction over the number of iterations k , but because it is convolved enough, I will rather proceed in the most digestible ordering. I'll also number equations excessively.

$$\mathbf{g}_k = \mathbf{A}\mathbf{x}_k - \mathbf{b} \quad (1)$$

$$\mathbf{x}_k = \mathbf{x}_{k-1} + \alpha_k \mathbf{d}_k \quad (2)$$

$$\mathbf{g}_k = \mathbf{g}_{k-1} + \alpha_k \mathbf{A}\mathbf{d}_k \quad (3)$$

$$\mathbf{g}_k^T \mathbf{d}_k = 0 \quad (4)$$

Here, (3) is just (2) times \mathbf{A} , minus \mathbf{b} , and we have discussed (4) above. We have also seen why the following holds:

$$\mathbf{g}_{k-1}^T \mathbf{d}_j = 0, \mathbf{d}_k^T \mathbf{A}\mathbf{d}_j = 0 \quad \Rightarrow \quad \mathbf{g}_k^T \mathbf{d}_j = 0 \quad (j \leq k) \quad (5)$$

In order to obtain a recurrence for the search directions, we make the ansatz

$$\mathbf{d}_{k+1} = -\mathbf{g}_k + \beta_k \mathbf{d}_k. \quad (6)$$

We would obtain steepest descent with $\beta_k = 0$, but in CG, we always have $\beta_k \neq 0$, unless we are done (we will see that below). Amazingly, this simple recurrence suffices to get all the rest. In order to determine what β_k should be, recall that we want conjugate directions.

$$0 = \mathbf{d}_{k+1}^T \mathbf{A}\mathbf{d}_k = (-\mathbf{g}_k + \beta_k \mathbf{d}_k)^T \mathbf{A}\mathbf{d}_k \quad \Rightarrow \quad \beta_k = \frac{\mathbf{g}_k^T \mathbf{A}\mathbf{d}_k}{\mathbf{d}_k^T \mathbf{A}\mathbf{d}_k}. \quad (7)$$

We will simplify this below. What about the gradients? If we assume the r.h.s. of (5) to hold, then for $j < k$:

$$\mathbf{g}_k^T \mathbf{g}_j \stackrel{(6)}{=} \mathbf{g}_k^T (\beta_j \mathbf{d}_j - \mathbf{d}_{j+1}) \stackrel{(5)}{=} 0.$$

Not relying on (5), we have

$$\mathbf{g}_k^T \mathbf{d}_j = \mathbf{g}_k^T \mathbf{d}_{j-1} = 0 \quad \Rightarrow \quad \mathbf{g}_k^T \mathbf{g}_j = 0. \quad (8)$$

OK, small break here. In steepest descent, search directions are orthogonal, and that is bad. In that method, search directions are (negative) gradients, so gradients are orthogonal. In

CG, improving on steepest descent, directions are not orthogonal, but *conjugate*. However, gradients in CG are still *orthogonal*. And to really get everybody confused, the whole method is called conjugate gradients!

We will now close the loop, by showing that orthogonality of gradients implies conjugacy of directions. Assume that for some $j < k$: $\mathbf{d}_k^T \mathbf{A} \mathbf{d}_j = 0$. Then,

$$\mathbf{d}_{k+1}^T \mathbf{A} \mathbf{d}_j \stackrel{(6)}{=} (-\mathbf{g}_k + \beta_k \mathbf{d}_k)^T \mathbf{A} \mathbf{d}_j = -\mathbf{g}_k^T \mathbf{A} \mathbf{d}_j \stackrel{(3)}{=} -\mathbf{g}_k^T (\mathbf{g}_j - \mathbf{g}_{j-1}) / \alpha_j.$$

But if gradients are orthogonal, the r.h.s. is zero. Therefore,

$$\mathbf{d}_k^T \mathbf{A} \mathbf{d}_j = 0, \mathbf{g}_k^T \mathbf{g}_j = \mathbf{g}_k^T \mathbf{g}_{j-1} = 0 \quad \Rightarrow \quad \mathbf{d}_{k+1}^T \mathbf{A} \mathbf{d}_j = 0. \quad (9)$$

The mutual orthogonality of *all* gradients and the conjugacy of *all* directions follows by running the cycle (5) \Rightarrow (8) \Rightarrow (9) \Rightarrow (5) ... Make sure you understand that this cycle is kickstarted by (4), which holds by line minimization.

We are almost done. Right now, CG would not really be elegant. While we have analytic expressions for α_k and β_k (7), they are clumsy and require more than one MVM with \mathbf{A} per iteration. By (4), (3): $0 = \mathbf{g}_k^T \mathbf{d}_k = (\mathbf{g}_{k-1} + \alpha \mathbf{A} \mathbf{d}_k)^T \mathbf{d}_k$, so that

$$\alpha_k = \frac{-\mathbf{g}_{k-1}^T \mathbf{d}_k}{\mathbf{d}_k^T \mathbf{A} \mathbf{d}_k} \stackrel{(6)}{=} \frac{-\mathbf{g}_{k-1}^T (-\mathbf{g}_{k-1} + \beta_{k-1} \mathbf{d}_{k-1})}{\mathbf{d}_k^T \mathbf{A} \mathbf{d}_k} \stackrel{(4)}{=} \frac{\|\mathbf{g}_{k-1}\|^2}{\mathbf{d}_k^T \mathbf{A} \mathbf{d}_k}. \quad (10)$$

And now, let's use (almost) everything above:

$$\beta_k \stackrel{(7)}{=} \frac{\mathbf{g}_k^T \mathbf{A} \mathbf{d}_k}{\mathbf{d}_k^T \mathbf{A} \mathbf{d}_k} \stackrel{(10)}{=} \frac{\mathbf{g}_k^T (\alpha_k \mathbf{A} \mathbf{d}_k)}{\|\mathbf{g}_{k-1}\|^2} \stackrel{(3)}{=} \frac{\mathbf{g}_k^T (\mathbf{g}_k - \mathbf{g}_{k-1})}{\|\mathbf{g}_{k-1}\|^2} \stackrel{(8)}{=} \frac{\|\mathbf{g}_k\|^2}{\|\mathbf{g}_{k-1}\|^2}. \quad (11)$$

From this equation, we see that if $\beta_k = 0$, then $\mathbf{g}_k = \mathbf{0}$, and we have reached the global minimum point of $q(\mathbf{x})$.

That's it. The CG algorithm itself is given as a lecture slide. Finally, note that (5) and the conjugacy of the \mathbf{d}_k imply that with at most n iterations, we are done. Namely, any set of conjugate directions is also linearly independent (easy exercise). But then, \mathbf{g}_n is orthogonal to all \mathbf{d}_k , $k \leq n$, which is possible only if $\mathbf{g}_n = \mathbf{0}$. For special matrices \mathbf{A} , whose characteristic polynomial has rank $< n$, this can happen earlier (a very simple example is $\mathbf{A} = \mathbf{I}$). In practice, you will probably never see this happening. Also note that (5) (new gradients orthogonal to old directions) directly implies the Krylov subspace minimization characteristic of CG that was discussed in the lecture.

References

- [1] C. Bishop. *Neural Networks for Pattern Recognition*. Oxford University Press, 1st edition, 1995.