

## Homework 3

Assigned: 10/10/2011.

Due: 21/10/2011.

### Exercise 1. BAYESIAN UPDATE WITH CHOLESKY FACTORS

Remember the robot localization example from the lecture. The robot's state is coded in  $\mathbf{u} \in \mathbb{R}^n$ , but the knowledge of the robot about the value of  $\mathbf{u}$  is uncertain. Say, its belief in  $\mathbf{u}$  is  $N(\mathbf{u}|\boldsymbol{\mu}, \boldsymbol{\Sigma})$  (robotics people call this *belief state*). Obtaining a noisy measurement  $y$ , whose likelihood is  $P(y|\mathbf{u}) = N(\mathbf{x}^T \mathbf{u}, \sigma^2)$  for some covariate  $\mathbf{x}$  (which is known to the robot, say the specification of sensors), the robot can update its belief from prior  $P(\mathbf{u})$  to posterior  $P(\mathbf{u}|y)$ , which becomes the new belief state. This process is called *Bayesian filtering* (you'll see more of that).

From the lecture, we know that  $P(\mathbf{u}|y) = N(\boldsymbol{\mu}', \boldsymbol{\Sigma}')$ , where

$$\boldsymbol{\Sigma}' = \boldsymbol{\Sigma} - \boldsymbol{\Sigma} \mathbf{x} (\sigma^2 + \mathbf{x}^T \boldsymbol{\Sigma} \mathbf{x})^{-1} \mathbf{x}^T \boldsymbol{\Sigma}, \quad \boldsymbol{\mu}' = \boldsymbol{\mu} + \boldsymbol{\Sigma} \mathbf{x} (\sigma^2 + \mathbf{x}^T \boldsymbol{\Sigma} \mathbf{x})^{-1} (y - \mathbf{x}^T \boldsymbol{\mu}). \quad (1)$$

Revisit the material, make yourself clear why these formulae are the way they are, and that there is nothing messy about them, because you know where they come from (it's just conditioning, "product"). Also remember which type of Gaussian parameterization we are using here, that there is a different type in which all problems here would disappear, although in that case, other essential operations would be tricky (marginalization, "sum"). However, these update equations do not have good numerical properties. In other words, assume the robot got a lot of measurements, and compare the *online* setting of absorbing one measurement at a time by (1) with the *offline* setting of collecting all measurements, then doing a single update. Quite likely, you would get substantially different results. In fact, with a static robot, it would probably all work out, but once the robot moves around, increasing the uncertainty in  $P(\mathbf{u})$  due to the movement (diffusion), you'd be in trouble. In fact, whatever you do in an online scenario, you'll always have more numerical error (on average) than offline. But you can do better than (1): the topic of this exercise. This is not a numerical mathematics course, so we'll not try to find out *why* this is better. Roughly, matrices, and matrix operations, with large condition numbers have to be avoided, and the condition number of a Cholesky factor is much less than of a matrix itself (square root).

1. [1 point] Let  $\boldsymbol{\Sigma} = \mathbf{L}\mathbf{L}^T$  be the Cholesky factorization. The key is to maintain  $\mathbf{L}$  instead of  $\boldsymbol{\Sigma}$ . Also, let  $\boldsymbol{\mu} = \mathbf{L}\mathbf{a}$ . Our *belief representation* will be  $(\mathbf{L}, \mathbf{a})$ . A representation is any amount of information, fulfilling:
  - Reasonable storage requirements
  - Desired queries to  $P(\mathbf{u})$  can be computed efficiently enough
  - For Bayesian updates  $P(\mathbf{u}) \rightarrow P(\mathbf{u}|y)$ , the representation update  $(\mathbf{L}, \mathbf{a}) \rightarrow (\mathbf{L}', \mathbf{a}')$  can be computed *reliably* (with as little numerical error as possible) and efficiently enough

In short, you describe and communicate your method in terms of distributions  $P(\mathbf{u})$ , maybe graphs, but you *implement* it in terms of a representation.

For some  $\mathbf{v}$ ,  $r = \mathbf{v}^T \mathbf{u}$  ( $\mathbf{u} \sim P(\mathbf{u})$ ) is Gaussian. Show how to compute mean, variance of  $r$  most efficiently. What is the cost? Compare that to the cost for using the representation  $(\boldsymbol{\Sigma}, \boldsymbol{\mu})$  instead for this query.

2. [3 points] Now to the update, based on observation  $(\mathbf{x}, y)$ . Suppose (in this part) that you have  $\mathbf{w} = \mathbf{L}^T \mathbf{x}$ . Bring (1) into the form  $\boldsymbol{\Sigma}' = \mathbf{L}(\mathbf{I} - \mathbf{v}\mathbf{v}^T)\mathbf{L}^T$ ,  $\boldsymbol{\mu}' = \mathbf{L}\mathbf{b}$ , and show how to compute  $\mathbf{v}$ ,  $\mathbf{b}$  in  $O(n)$ .  
*Hint:* You may define another vector as a rescaled version of  $\mathbf{w}$ . If you divide by numbers, argue why they are  $\neq 0$ .
3. [1 point] Suppose I gave you some *Cholesky downdate* code, which maps  $(\mathbf{L}, \mathbf{z}, \mathbf{v}, \alpha) \rightarrow (\mathbf{L}', \mathbf{z}')$ , so that

$$\mathbf{L}'(\mathbf{L}')^T = \mathbf{L}(\mathbf{I} - \mathbf{v}\mathbf{v}^T)\mathbf{L}^T, \quad \mathbf{L}'\mathbf{z}' = \mathbf{L}(\mathbf{z} - \alpha\mathbf{v})$$

for nonzero  $\mathbf{v}$  and some  $\alpha$  of your choice. Show how to compute the representation update  $(\mathbf{L}, \mathbf{a}) \rightarrow (\mathbf{L}', \mathbf{a}')$ , if all you are allowed to do is call the downdate routine once, do a single multiplication matrix-vector multiplication with  $\mathbf{L}^T$ , and  $O(n)$  computations otherwise.

Such downdate code is contained in numerical packages. On the course website, I have linked some wrapper code of mine, as well as a report telling you how this works (if you are interested). The procedure you developed here seems a convoluted way of doing (1), but this “intuition” is wrong. It comes at exactly the same computational cost (time, memory; in fact, it is even slightly faster), and has much better numerical properties.

### Exercise 2. SPECTRAL ANALYSIS OF CONJUGATE GRADIENTS ALGORITHM

Recall the CG algorithm for approximately solving  $\mathbf{A}\mathbf{x} = \mathbf{b}$ , where  $\mathbf{A} \in \mathbb{R}^{n \times n}$  is positive definite. The algorithm minimizes the quadratic  $q(\mathbf{x}) = (1/2)\mathbf{x}^T \mathbf{A}\mathbf{x} - \mathbf{b}^T \mathbf{x}$  iteratively, constructing a sequence  $\mathbf{x}_1, \mathbf{x}_2, \dots$ , requiring a single matrix-vector multiplication with  $\mathbf{A}$  per iteration. After at most  $n$  steps, neglecting numerical errors (which, in practice, you cannot!), the exact solution  $\mathbf{x}_* = \mathbf{A}^{-1}\mathbf{b}$  is reached, in that  $\mathbf{x}_n = \mathbf{x}_*$ . Depending on  $\mathbf{A}$ , this can also happen earlier (you can use the Cayley-Hamilton theorem from linear algebra to understand this point). However, the main rationale for CG today is to *approximate*  $\mathbf{x}_*$  by  $\mathbf{x}_k$  with  $k \ll n$ . Whether  $\mathbf{x}_k$  is close to  $\mathbf{x}_*$  or not, depends on properties of  $\mathbf{A}$  and  $\mathbf{b}$  (it depends on numerical errors as well, but we ignore these in the present exercise, assuming that all computations are exact). In this exercise, we will analyze the convergence behaviour in terms of the eigenspectrum of  $\mathbf{A}$ .

Let  $\mathbf{x}_* = \mathbf{A}^{-1}\mathbf{b}$  and  $q_* = q(\mathbf{x}_*) = -(1/2)\mathbf{b}^T \mathbf{A}^{-1}\mathbf{b}$ . Then,  $q(\mathbf{x}_k)$  is nonincreasing and  $\geq q_*$ . We'll try to bound  $q(\mathbf{x}_k) - q_*$ . The eigendecomposition is  $\mathbf{A} = \mathbf{Q}\boldsymbol{\Lambda}\mathbf{Q}^T$ ,  $\mathbf{Q} \in \mathbb{R}^{n \times n}$  orthonormal ( $\mathbf{Q}^T \mathbf{Q} = \mathbf{I}$ ),  $\boldsymbol{\Lambda}$  diagonal (positive elements).

1. [2 points] Assume that we start from  $\mathbf{x}_0 = \mathbf{0}$ . The Krylov subspace  $\mathcal{K}_k$  is spanned by  $\{\mathbf{A}^j \mathbf{b} \mid j = 0, \dots, k-1\}$ . We saw that  $\mathbf{x}_k = \operatorname{argmin}_{\mathbf{x} \in \mathcal{K}_k} q(\mathbf{x})$ . For a polynomial  $P(t) = \sum_{j=0}^{k-1} \alpha_j t^j$ ,  $\alpha_j \in \mathbb{R}$ , define

$$P(\mathbf{B}) := \sum_{j=0}^{k-1} \alpha_j \mathbf{B}^j, \quad \mathbf{B} \in \mathbb{R}^{n \times n}.$$

Take care that  $P(\mathbf{B})P(\mathbf{C}) \neq P(\mathbf{C})P(\mathbf{B})$  in general if  $\mathbf{B}, \mathbf{C}$  do not commute, and recall that  $\mathbf{B}^0 := \mathbf{I}$ . Show that

$$P(\mathbf{A}) = \mathbf{Q}P(\boldsymbol{\Lambda})\mathbf{Q}^T.$$

Defining  $\mathbf{y} = \mathbf{Q}^T \mathbf{x}$ ,  $\bar{\mathbf{b}} = \mathbf{Q}^T \mathbf{b}$ , show that  $q(\mathbf{x})$  and  $q_*$  can be written in terms of  $\mathbf{y}$ ,  $\bar{\mathbf{b}}$ , and  $\{\lambda_i\}$ . *Hint:* What is an eigendecomposition of  $\mathbf{A}^{-1}$ ?

2. [4 points] Let  $\mathbf{y}_k = \mathbf{Q}^T \mathbf{x}_k$ . Show that  $\mathbf{x}_k = P_k(\mathbf{A})\mathbf{b}$  for some polynomial  $P_k(t)$  of degree  $< k$ , and that  $\mathbf{y}_k = P_k(\boldsymbol{\Lambda})\bar{\mathbf{b}}$ . In other words,  $y_{k,i} = P_k(\lambda_i)\bar{b}_i$ . By writing  $q(\mathbf{x}_k)$

in terms of  $\mathbf{y}_k$ , prove that

$$q(\mathbf{x}_k) - q_* = \min_{P_k \mid \deg(P_k) < k} (1/2) \sum_{i=1}^n (\bar{b}_i^2 / \lambda_i) (\lambda_i P_k(\lambda_i) - 1)^2.$$

Here,  $\deg(P_k)$  is the degree of  $P_k$ , the largest  $j$  such that  $t^j$  features with a nonzero coefficient.

3. **[3 points]** Argue that this means that the error  $q(\mathbf{x}_k) - q_*$  is bounded in terms of a polynomial  $P$  of degree  $\leq k$ , such that  $P(0) = -1$  (equivalent: such that  $P(0) = 1$ ). The existence of such a polynomial which is small on *all*  $\lambda_i$ , implies that the error is small. Prove that if  $\{\lambda_1, \dots, \lambda_n\} = \{\kappa_1, \dots, \kappa_k\}$ , *i.e.* if  $\mathbf{A}$  has no more than  $k$  different eigenvalues, then  $\mathbf{x}_k = \mathbf{x}_*$ .

*Remark:* Using Chebishev polynomials,  $\min_{P_k} \max_{t \in [\lambda_{\min}, \lambda_{\max}]} P_k(t)$  can be determined for  $0 < \lambda_{\min} \leq \lambda_{\max}$ , which leads to the worst-case error bound

$$q(\mathbf{x}_k) - q_* \leq \left( \frac{\rho - 1}{\rho + 1} \right)^k, \quad \rho = \sqrt{\lambda_{\max} / \lambda_{\min}}.$$

Therefore, we should aim for well-conditioned matrices  $\mathbf{A}$  ( $\lambda_{\max} / \lambda_{\min}$  small), which is what many preconditioning strategies try to do. On the other hand, if the spectrum of  $\mathbf{A}$  comes in separate clusters, this bound is overly pessimistic.