

Lecture 10, Feb 27, 2024

Quasi-Newton Methods (Symmetric Rank 1 (SR1))

- Recall that for quasi-Newton methods, since computing the inverse Hessian is expensive, we use approximations to speed up computation
- Idea: use an iterative update routine to approximate the inverse Hessian
- Start with a quadratic approximation of the objective function $f(\boldsymbol{\theta})$ at the current iteration, $m_k(\boldsymbol{\theta})$
- $m_k(\boldsymbol{\theta}) = f(\boldsymbol{\theta}_k) + \vec{\nabla}^T f(\boldsymbol{\theta}_k)(\boldsymbol{\theta} - \boldsymbol{\theta}_k) + \frac{1}{2}(\boldsymbol{\theta} - \boldsymbol{\theta}_k)^T \mathbf{B}_k(\boldsymbol{\theta} - \boldsymbol{\theta}_k)$
 - $\mathbf{B}_k \in \mathbb{R}^{n \times n}$ is an approximation of the inverse Hessian
 - When $\boldsymbol{\theta} = \boldsymbol{\theta}_k$ we have $m_k = f$ and $\vec{\nabla} m_k = \vec{\nabla} f$
 - These conditions are known as the *zero* and *first-order consistency conditions*
- The minimum of the quadratic model can be obtained by differentiating and setting to zero as usual
 - $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \mathbf{B}_k^{-1} \vec{\nabla} f(\boldsymbol{\theta}_k)$
- Using the new minimum, we construct a new quadratic approximation $m_{k+1}(\boldsymbol{\theta})$ using the same formula and the new approximate Hessian \mathbf{B}_{k+1}
- To update the approximate Hessian, we impose the constraint that $m_{k+1}(\boldsymbol{\theta})$ matches the gradient of $f(\boldsymbol{\theta})$ at both $\boldsymbol{\theta}_k$ and $\boldsymbol{\theta}_{k+1}$
 - We want $\vec{\nabla} m_{k+1}(\boldsymbol{\theta}_{k+1}) = \vec{\nabla} f(\boldsymbol{\theta}_{k+1}) + \mathbf{B}_{k+1}(\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}) = \vec{\nabla} f(\boldsymbol{\theta}_k)$
 - Rearrange to get $\mathbf{B}_{k+1}(\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}) = \vec{\nabla} f(\boldsymbol{\theta}_k) - \vec{\nabla} f(\boldsymbol{\theta}_{k+1})$
 - * This is known as the *secant equation*
- Since the Hessian is SPD, we want \mathbf{B}_{k+1} to also be SPD
 - Symmetry gives us $\frac{1}{2}n(n+1)$ independent entries, but the secant equation only gives a system of n equations
 - To obtain a unique solution, we impose the constraint that \mathbf{B}_{k+1} should be closest to \mathbf{B}_k
 - Therefore we use the update formula: $\mathbf{B}_{k+1} = \mathbf{B}_k + \mathbf{u}\mathbf{u}^T$
 - * $\mathbf{u}\mathbf{u}^T$ is the “symmetric rank 1” matrix
 - * This guarantees that \mathbf{B}_{k+1} is close to \mathbf{B}_k in terms of rank
- Let $\mathbf{s}_k = \boldsymbol{\theta}_{k+1} - \boldsymbol{\theta}_k$, $\mathbf{y}_k = \vec{\nabla} f(\boldsymbol{\theta}_{k+1}) - \vec{\nabla} f(\boldsymbol{\theta}_k)$
 - Rewrite the secant equation as $\mathbf{B}_{k+1}\mathbf{s}_k = \mathbf{y}_k$
 - Plugging in the SR1 update, $\mathbf{B}_k\mathbf{s}_k + \mathbf{u}\mathbf{u}^T\mathbf{s}_k = \mathbf{y}_k \implies \mathbf{u}(\mathbf{u}^T\mathbf{s}_k) = \mathbf{y}_k - \mathbf{B}_k\mathbf{s}_k$
 - * This means $\mathbf{u} = \gamma(\mathbf{y}_k - \mathbf{B}_k\mathbf{s}_k)$ where $\gamma = \mathbf{u}^T\mathbf{s}_k$ is a scalar
 - Plug this back in: $\gamma^2(\mathbf{y}_k - \mathbf{B}_k\mathbf{s}_k)(\mathbf{s}_k^T(\mathbf{y}_k - \mathbf{B}_k\mathbf{s}_k)) = \mathbf{y}_k - \mathbf{B}_k\mathbf{s}_k$
 - $\gamma^2 = \frac{1}{\mathbf{s}_k^T(\mathbf{y}_k - \mathbf{B}_k\mathbf{s}_k)}$
- The final update formula is $\mathbf{B}_{k+1} = \mathbf{B}_k + \frac{(\mathbf{y}_k - \mathbf{B}_k\mathbf{s}_k)(\mathbf{y}_k - \mathbf{B}_k\mathbf{s}_k)^T}{(\mathbf{y}_k - \mathbf{B}_k\mathbf{s}_k)^T\mathbf{s}_k}$
 - However, this only gives the approximate Hessian and not its inverse (inverting at each iteration would be too expensive)

Theorem

Sherman-Morrison-Woodbury (SMW) Formula: Given $\mathbf{A} \in \mathbb{R}^{n \times n}$, $\mathbf{u}, \mathbf{v} \in \mathbb{R}^{n \times p}$, then

$$(\mathbf{A} + \mathbf{u}\mathbf{v}^T)^{-1} = \mathbf{A}^{-1} - \mathbf{A}^{-1}\mathbf{u}(\mathbf{1} + \mathbf{v}^T\mathbf{A}^{-1}\mathbf{u})^{-1}\mathbf{v}^T\mathbf{A}^{-1}$$

- Using the SMW formula: $\mathbf{B}_{k+1}^{-1} = \mathbf{B}_k^{-1} + \frac{(\mathbf{s}_k - \mathbf{B}_k^{-1}\mathbf{y}_k)(\mathbf{s}_k - \mathbf{B}_k^{-1}\mathbf{y}_k)^T}{(\mathbf{s}_k - \mathbf{B}_k^{-1}\mathbf{y}_k)^T\mathbf{y}_k}$
 - This gives a cost of $\mathcal{O}(n^2)$ (compared to $\mathcal{O}(n^3)$ for matrix inversion)
- An alternative approach to compute \mathbf{B}_{k+1} is to formulate it as a constrained optimization problem
 - $\mathbf{B}_{k+1} = \min_{\mathbf{B}} \|\mathbf{B} - \mathbf{B}_k\|$ subject to $\mathbf{B} = \mathbf{B}^T$, $\mathbf{B}(\boldsymbol{\theta}_k - \boldsymbol{\theta}_{k+1}) = \vec{\nabla} f(\boldsymbol{\theta}_k) - \vec{\nabla} f(\boldsymbol{\theta}_{k+1})$
 - The choice of matrix norm to use leads to different variations of the method:
 - * Davidon-Fletcher-Powell (DFP): $\mathbf{B}_{k+1}^{-1} = \mathbf{B}_k^{-1} - \mathbf{A}_k + \mathbf{C}_k$

- $A_k = \frac{B_k^{-1} \mathbf{y}_k \mathbf{y}_k^T B_k^{-1}}{\mathbf{y}_k^T B_k^{-1} \mathbf{y}_k}$
- $C_k = \frac{\mathbf{s}_k \mathbf{s}_k^T}{\mathbf{s}_k^T \mathbf{y}_k}$
- This is a rank-2 method
- * Broyden–Fletcher–Goldfarb–Shanno (BFGS):
 - $B_{k+1}^{-1} = \left(\mathbf{1} - \frac{\mathbf{s}_k \mathbf{y}_k^T}{\mathbf{s}_k^T \mathbf{y}_k} \right) B_k^{-1} \left(\mathbf{1} - \frac{\mathbf{s}_k \mathbf{y}_k^T}{\mathbf{s}_k^T \mathbf{y}_k} \right) + \frac{\mathbf{s}_k \mathbf{s}_k^T}{\mathbf{s}_k^T \mathbf{y}_k}$
- Quasi-Newton methods generally have between linear and quadratic convergence; we call this *superlinear*
- In problems where n is very large such that $\mathcal{O}(n^2)$ is impractical, limiting-memory quasi-Newton methods compute the search step directly

Constrained Optimization – Penalty Methods

- Consider the problem of minimizing $f(\boldsymbol{\theta})$ subject to constraints $g_i(\boldsymbol{\theta}) \geq 0, h_j(\boldsymbol{\theta}) = 0$ where $i = 1, 2, \dots, m$ and $j = 1, 2, \dots, q$ and $\boldsymbol{\theta}_l \leq \boldsymbol{\theta} \leq \boldsymbol{\theta}_u$
- Penalty methods minimize $\pi(\boldsymbol{\theta}, \rho_k) = f(\boldsymbol{\theta}) + \rho_k \phi(\boldsymbol{\theta})$
 - $\phi(\boldsymbol{\theta})$ is the *penalty function* and ρ_k is the *penalty parameter*
 - We want $\phi(\boldsymbol{\theta})$ equal to zero when no constraints are violated and positive when constraints are violated
 - We need to ensure that the objective and the penalty function are appropriately scaled, so one doesn't dominate the other
- Quadratic penalty function: $\phi(\boldsymbol{\theta}_k) = \sum_{i=1}^m (\max(0, -g_i(\boldsymbol{\theta}_k)))^2 + \sum_{i=1}^q (h_i(\boldsymbol{\theta}_k))^2$
- Penalty methods template:
 1. Check termination conditions
 2. Minimize $\pi(\boldsymbol{\theta}, \rho_k)$ to find $\boldsymbol{\theta}_{k+1}$
 3. Increment the penalty parameter, $\rho_{k+1} > \rho_k$
 - Typically we multiply by a factor of 1.4 to 10, but this is problem dependent

Nonlinear Least Squares

- $\min_{\boldsymbol{\theta} \in \mathbb{R}^n} = \frac{1}{2} \sum_{i=1}^N r_i(\boldsymbol{\theta})^2$ where $r_i(\boldsymbol{\theta}) = \hat{f}(\mathbf{x}^{(i)}, \boldsymbol{\theta}) - y^{(i)}$
 - Assume $N > n$, i.e. we have more data points than dimensions
- Let $\mathbf{r} = \begin{bmatrix} r_1 \\ \vdots \\ r_n \end{bmatrix} \in \mathbb{R}^N$ and $f(\boldsymbol{\theta}) = \frac{1}{2} \|\mathbf{r}(\boldsymbol{\theta})\|_2^2$
- $\vec{\nabla} f(\boldsymbol{\theta}) = \sum_{j=1}^N r_j(\boldsymbol{\theta}) \vec{\nabla} r_j(\boldsymbol{\theta}) = \mathbf{J}(\boldsymbol{\theta})^T \mathbf{r}(\boldsymbol{\theta})$
 - $\mathbf{J}(\boldsymbol{\theta}) = \begin{bmatrix} \partial r_j \\ \partial r_i \end{bmatrix} \in \mathbb{R}^{N \times n}$ is the Jacobian
- $\vec{\nabla}^2 f(\boldsymbol{\theta}) = \sum_{j=1}^N \vec{\nabla} r_j(\boldsymbol{\theta}) \vec{\nabla} r_j(\boldsymbol{\theta})^T + \sum_{j=1}^N r_j(\boldsymbol{\theta}) \vec{\nabla}^2 r_j(\boldsymbol{\theta})$

$$= \mathbf{J}(\boldsymbol{\theta})^T \mathbf{J}(\boldsymbol{\theta}) + \sum_{j=1}^N r_j(\boldsymbol{\theta}) \vec{\nabla}^2 r_j(\boldsymbol{\theta})$$
 - The Jacobian is easy to compute, which gets us most of the way to the Hessian
 - Often the second term is small so we can ignore it altogether and use the Jacobian to approximate the Hessian
 - * This happens when the initial residual is small

Gauss-Newton Method

- This is similar to a modified Newton's method with line search
- Use $\vec{\nabla}^2 f(\boldsymbol{\theta}) \approx \mathbf{J}(\boldsymbol{\theta})^T \mathbf{J}(\boldsymbol{\theta})$ as an approximation of the Hessian
- Solve $\mathbf{J}(\boldsymbol{\theta})^T \mathbf{J}(\boldsymbol{\theta}) \mathbf{p}_k = -\mathbf{J}(\boldsymbol{\theta})^T \mathbf{r}(\boldsymbol{\theta}_k)$ for the search direction
 - Update $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k + \alpha_k \mathbf{p}_k$ where α_k is chosen via line search
- In the case where the initial residual is small or approximately linear in $\boldsymbol{\theta}$, the Gauss-Newton method can perform similar to the full Newton's method, despite only computing first-order derivatives
- If $\mathbf{J}(\boldsymbol{\theta}_k)$ is full-rank and $\vec{\nabla} f(\boldsymbol{\theta}_k) \neq 0$, the search direction is always a valid direction

Stochastic Gradient Descent (SGD)

- In general we have loss function $\mathcal{L}(\boldsymbol{\theta}; \mathcal{D}) = \frac{1}{N} \sum_{i=1}^N l(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) + \lambda R(\boldsymbol{\theta})$ given data $\mathcal{D} = \{(\mathbf{x}^{(i)}, y^{(i)})\}_{i=1}^N$
 - This consists of the empirical loss and a regularization term
- $\vec{\nabla} \mathcal{L} = \frac{1}{N} \sum_{i=1}^N \vec{\nabla} l(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)}) + \lambda \vec{\nabla} R(\boldsymbol{\theta})$
- Applying the steepest descent method, we get the update $\boldsymbol{\theta}_{k+1} = \boldsymbol{\theta}_k - \eta_k \vec{\nabla} \mathcal{L}(\boldsymbol{\theta}_k; \mathcal{D})$
 - η_k is the *learning rate*
 - * In classical methods we use backtracking line search to find this, but in machine learning we typically choose this heuristically, as a constant
 - * e.g. start with a sensible value like $\eta = 0.1$; take smaller steps if objective gets worse or we see oscillation; take larger steps if objective reduces too slowly
 - Since we are computing the gradient over the full dataset \mathcal{D} , this is known as *full-batch gradient descent*
- Full-batch gradient descent is typically very expensive since we need to compute the gradient over the entire dataset
- Procedure of SGD:
 1. Shuffle training indices $\{1, \dots, N\}$
 2. Initialize $\boldsymbol{\theta}_0$
 3. Repeat until we reach some convergence criteria:
 - For i from 1 to N , $\boldsymbol{\theta} \leftarrow \boldsymbol{\theta} - \eta \vec{\nabla} l(\boldsymbol{\theta}; \mathbf{x}^{(i)}, y^{(i)})$
- Each iteration of the outer loop is an *epoch*, where we loop over the full dataset
- SGD essentially uses only one datapoint at a time
 - This works because the gradient using one datapoint is an unbiased estimator of the full gradient
 - Let $\mathbf{g}_t = l(\boldsymbol{\theta}_k; \mathbf{x}^{(t)}, y^{(t)})$, we have that $\mathbb{E}[\mathbf{g}_t] = \mathcal{L}(\boldsymbol{\theta}_k; \mathcal{D})$
- Consider gradient descent over a GLM with M terms
 - The cost of full-batch gradient descent is $\mathcal{O}(NM)$, and converges in $\mathcal{O}\left(\log \frac{1}{\rho}\right)$
 - The cost of stochastic gradient descent is only $\mathcal{O}(M)$, and converges in $\mathcal{O}\left(\frac{1}{\rho}\right)$ iterations
 - * Even though SGD takes more iterations to converge (sub-linearly), it's cheaper overall when factoring in the cost per iteration
 - * Sometimes it's not practical to do full-batch gradient descent due to the size of the dataset
- In *mini-batch gradient descent* we compute the gradient over a mini-batch that is smaller than the full dataset, but more than 1 sample, in each iteration
 - This is a compromise between full-batch gradient descent and SGD
 - The larger the batch size, the closer we get to full-batch and the faster we converge (in iterations)