

Optimization and the BFGS Algorithm

Jake Spiteri

December 06, 2019

Table of Contents

1 Introduction

- Unconstrained Optimization
- Motivating Example
- Intuition

2 Unconstrained Optimization

- Algorithms
 - Second-Order Methods — Newton's Method
 - Quasi-Newton Methods — The BFGS Method

The Optimization Problem

An **optimization** problem has the general form

$$\begin{array}{ll} \text{minimize} & f_0(\mathbf{x}) \\ \text{subject to} & f_i(\mathbf{x}) \leq b_i, \quad i = 1, \dots, m. \end{array}$$

We say that

- $\mathbf{x} = (x_1, \dots, x_n)$ are the *optimization variables*,
- $f_0 : \mathbb{R}^n \rightarrow \mathbb{R}$ is the *objective function*,
- $f_i : \mathbb{R}^n \rightarrow \mathbb{R}$, $i = 1, \dots, m$ are the *inequality constraint functions*.

The Optimization Problem

We say that a vector \mathbf{x}^* is optimal if it minimizes the objective function $f_0(\mathbf{x})$ and satisfies the constraint functions f_i , $i = 1, \dots, m$.

Unconstrained Problem

We minimize an objective function which is a function of real variables. We place no constraints on these values. The general problem can be written

$$\min\{f(\mathbf{x}) : \mathbf{x} \in \mathbb{R}^n\}$$

A Motivating Example — MLE

Suppose we observe data $\{x_i, y_i\}_{i=1}^n$, where our y_i is **count** data. Then we can think of Y_i as following a Poisson distribution. I.e. we can form the Poisson regression model:

$$Y_i \sim \text{Poisson}(\mu_i), \quad \mu_i = \mathbb{E}(Y_i), \quad \log(\mu_i) = \beta_0 + \beta_1 x_i$$

We want to estimate the parameters $\beta = (\beta_0, \beta_1)$ by maximum likelihood estimation.

We can write down the conditional probability and thus the log likelihood:

$$p(\mathbf{y} \mid \mathbf{x}) = \prod_{i=1}^n e^{y_i \beta \mathbf{x}_i} e^{-e^{\beta \mathbf{x}_i}} / y_i!$$
$$\ell(\beta) = \sum_{i=1}^n y_i \beta \mathbf{x}_i - e^{\beta \mathbf{x}_i} - \log(y_i!)$$

We cannot write down a closed form for β ! We require numerical optimization methods.

Intuition for Optimization Algorithms

The optimization algorithms we present are iterative processes that allow us to find the minimum of a function. Often we can only promise a local minimum.

At time step k , we have current iterate \mathbf{x}_k and we want to find the next iterate \mathbf{x}_{k+1} such that $f(\mathbf{x}_k) > f(\mathbf{x}_{k+1})$.

Our algorithms have the general form

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{d}_k,$$

where \mathbf{d}_k is a descent direction.

Intuition for Optimization Algorithms

First-order methods: Use the first derivative.

Steepest descent is a very popular first-order method with the following update step:

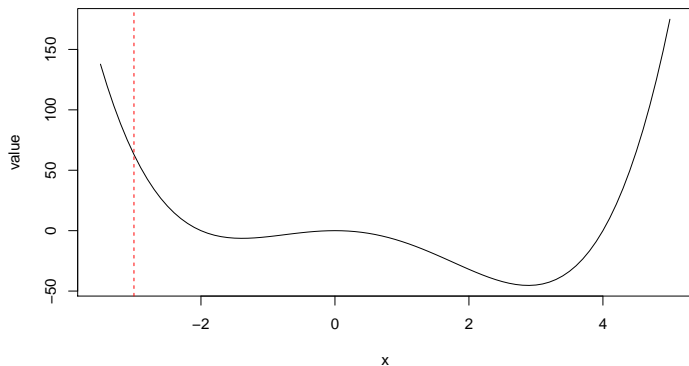
$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k).$$

Suppose we want to minimize the function $f(x) = x^4 - 2x^3 - 8x^2$.
At time step k we have $x_k = -3$. How do we choose x_{k+1} ?

Intuition for Optimization Algorithms

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)$$

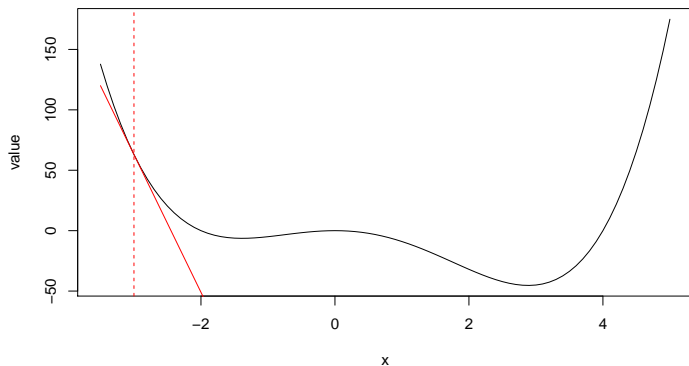
Plot of $f(x)$



Intuition for Optimization Algorithms

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)$$

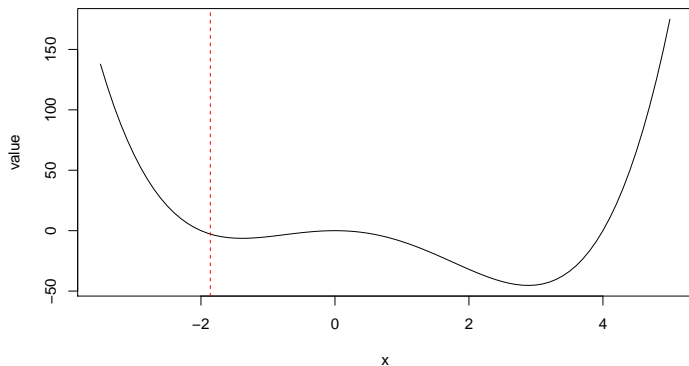
Plot of $f(x)$



Intuition for Optimization Algorithms

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \nabla f(\mathbf{x}_k)$$

Plot of $f(x)$



Intuition for Optimization Algorithms

Second-order methods: Use the first derivative and the Hessian (the matrix of second derivatives).

Newton's method is a very popular second-order method with the following update step:

$$\mathbf{x}_{k+1} = \mathbf{x}_k - (\nabla^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k).$$

Key intuition

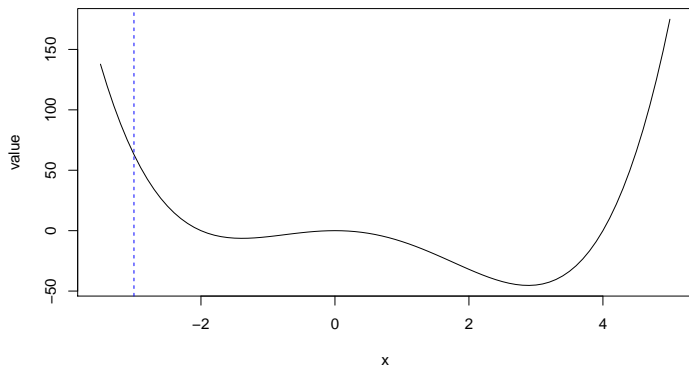
Newton's method minimizes the second-order approximation of the function.

Suppose we want to minimize the function $f(x) = x^4 - 2x^3 - 8x^2$. At time step k we have $x_k = -3$. How do we choose x_{k+1} ?

Intuition for Optimization Algorithms

$$\mathbf{x}_{k+1} = \mathbf{x}_k - (\nabla^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k)$$

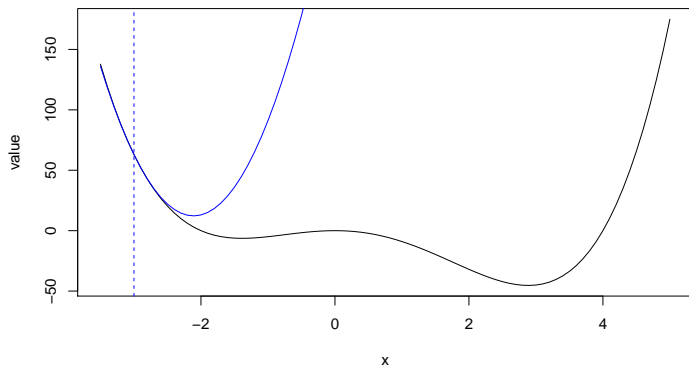
Plot of $f(x)$



Intuition for Optimization Algorithms

$$\mathbf{x}_{k+1} = \mathbf{x}_k - (\nabla^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k)$$

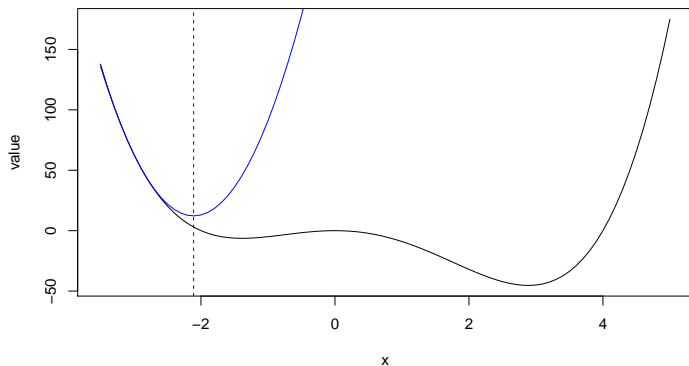
Plot of $f(x)$



Intuition for Optimization Algorithms

$$\mathbf{x}_{k+1} = \mathbf{x}_k - (\nabla^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k)$$

Plot of $f(x)$



Second-Order Methods — Approximating a function

Theorem (quadratic approximation theorem)

Let $f : U \rightarrow \mathbb{R}$ be a twice continuously differentiable function over an open set $U \subset \mathbb{R}^n$, and let $\mathbf{x} \in U, r > 0$ satisfy $B(\mathbf{x}, r) \subseteq U$. Then for any $\mathbf{y} \in B(\mathbf{x}, r)$

$$f(\mathbf{y}) = f(\mathbf{x}) + \nabla f(\mathbf{x})^T (\mathbf{y} - \mathbf{x}) + \frac{1}{2} (\mathbf{y} - \mathbf{x})^T \nabla^2 f(\mathbf{x}) (\mathbf{y} - \mathbf{x}) + o(\|\mathbf{y} - \mathbf{x}\|^2).$$

Second-Order Methods

- A second-order method uses evaluations of the Hessian.
- Newton's method is an iterative method in which each update is chosen to minimize the quadratic approximation of the objective function around \mathbf{x}_k .

Quadratic approximation of $f(\mathbf{x})$ around \mathbf{x}_k :

$$Q(\mathbf{x}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T (\mathbf{x} - \mathbf{x}_k) + \frac{1}{2} (\mathbf{x} - \mathbf{x}_k)^T \nabla^2 f(\mathbf{x}_k) (\mathbf{x} - \mathbf{x}_k) \quad \mathbf{x}_{k+1} = \operatorname{argmin}_{\mathbf{x} \in \mathbb{R}^n}$$

$$\nabla f(\mathbf{x}_k) + \nabla^2 f(\mathbf{x}_k) (\mathbf{x}_{k+1} - \mathbf{x}_k) = \mathbf{0},$$

$$\mathbf{x}_{k+1} = \mathbf{x}_k - (\nabla^2 f(\mathbf{x}_k))^{-1} \nabla f(\mathbf{x}_k)$$

Comments on Newton's Method

Newton's method converges quadratically.

Pure Newton's method may diverge if our initialization is bad. We can avoid this by modifying the Hessian at each step to make it positive (semi) definite if it is not, and ensuring that the step size is chosen such that the function decreases.

It may be difficult to evaluate the Hessian, particularly for high-dimensional objective functions.

Quasi-Newton Methods

Quasi-Newton methods use an approximation of the second derivative.

Avoids evaluating the Hessian $\nabla^2 f(\mathbf{x})$ at each iteration — an error-prone and expensive process.

Attains a **superlinear** rate of convergence, and lower computational cost ($\mathcal{O}(n^2)$ as opposed to $\mathcal{O}(n^3)$ in Newton's method).

The BFGS Method

Recall the second-order Taylor series approximation around \mathbf{x}_k ,

$$f(\mathbf{x}_k + \mathbf{d}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \nabla^2 f(\mathbf{x}_k) \mathbf{d} + o(\|\mathbf{d}\|^2)$$

We will define the approximation of the objective function at \mathbf{x}_k as

$$m_k(\mathbf{d}) = f(\mathbf{x}_k) + \nabla f(\mathbf{x}_k)^T \mathbf{d} + \frac{1}{2} \mathbf{d}^T \mathbf{B}_k \mathbf{d}$$

\mathbf{B}_k is an $n \times n$ symmetric positive definite matrix that will be updated at every iteration.

The BFGS Method

The minimizer of the convex quadratic $m_k(\mathbf{d})$, which we can write as

$$\mathbf{d} = \mathbf{B}_k^{-1} \nabla f(\mathbf{x}_k),$$

is the *search direction*. Hence, the new iterate is

$$\mathbf{x}_{k+1} = \mathbf{x}_k - \alpha_k \mathbf{B}_k^{-1} \nabla f(\mathbf{x}_k),$$

where α_k is our step size.

This iterative update is the same as Newton's method except it uses the approximate Hessian.

Instead of recomputing \mathbf{B}_k at every iteration, we iteratively update it by accounting for the curvature measured in the most recent step.

Suppose we generate a new iterate \mathbf{x}_{k+1} , then we construct the new quadratic model

$$m_{k+1}(\mathbf{d}) = f(\mathbf{x}_{k+1}) + \nabla f(\mathbf{x}_{k+1})^T \mathbf{d} + \mathbf{d}^T \mathbf{B}_{k+1} \mathbf{d}$$

In order to determine the new \mathbf{B}_{k+1} we place restrictions on $m_{k+1}(\mathbf{d})$.

We require that the gradient of m_{k+1} be equal to the gradient of f evaluated at the two most recent iterates \mathbf{x}_k and \mathbf{x}_{k+1} .

- 1 $\nabla m_{k+1}(\mathbf{0}) = \nabla f(\mathbf{x}_{k+1})$,
- 2 $\nabla m_{k+1}(-\alpha_k \mathbf{d}_k) = \nabla f(\mathbf{x}_{k+1}) - \alpha_k \mathbf{B}_{k+1} \mathbf{d}_k = \nabla f(\mathbf{x}_k)$.

Rearranging 2 gives us

$$\alpha_k \mathbf{B}_{k+1} \mathbf{d}_k = \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k)$$

We introduce the following notation:

- $\mathbf{s}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$ ($= \alpha_k \mathbf{d}_k$),
- $\mathbf{y}_k = \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k)$,

Recall from the last slide

$$\alpha_k \mathbf{B}_{k+1} \mathbf{d}_k = \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k),$$

then we can write the **secant** equation

$$\mathbf{B}_{k+1} \mathbf{s}_k = \mathbf{y}_k.$$

- Secant equation:

$$\mathbf{B}_{k+1}\mathbf{s}_k = \mathbf{y}_k$$

- BFGS update:

$$\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{B}_k^{-1} \nabla f(\mathbf{x}_k)$$

- BFGS places restrictions on the inverse of the Hessian approximation

$$\mathbf{H}_k = \mathbf{B}_k^{-1}$$

- Secant equation

$$\mathbf{H}_{k+1}\mathbf{y}_k = \mathbf{s}_k$$

Secant equation:

$$\mathbf{H}_{k+1}\mathbf{y}_k = \mathbf{s}_k$$

We require that the symmetric positive definite matrix \mathbf{H}_{k+1} maps \mathbf{y}_k to \mathbf{s}_k . This is only possible if \mathbf{y}_k and \mathbf{s}_k satisfy the curvature condition

$$\mathbf{y}_k^T \mathbf{s}_k > 0.$$

If f is strongly convex then the curvature condition is satisfied for any \mathbf{x}_k and \mathbf{x}_{k+1} .

Remark

For non-convex f we must impose this restriction explicitly. We impose the Wolfe or strong Wolfe conditions on the line search. The Wolfe conditions ensure that we take a step which produces a sufficient decrease in f , and doesn't take an unacceptably small step.

The secant equation places n restrictions on \mathbf{H}_{k+1} .

As \mathbf{H}_{k+1} is symmetric positive definite, it has $n(n+1)/2$ parameters.

\mathbf{H}_{k+1} has infinitely many solutions.

In order to determine \mathbf{H}_{k+1} uniquely, we impose the condition that among all symmetric matrices that satisfy the secant equation, \mathbf{H}_{k+1} is in some sense closest to \mathbf{H}_k .

We form the optimization problem

where \mathbf{s}_k and \mathbf{y}_k satisfy the curvature condition and \mathbf{H}_k is symmetric and positive definite.

Different matrix norms result in different quasi-Newton methods. Hence the following choices are arbitrary but they provide a nice solution and produce the BFGS algorithm update.

A norm that provides an easy solution to the optimization problem which is scale-invariant is the weighted Frobenius norm

$$\|\mathbf{A}\|_{\mathbf{W}} \equiv \|\mathbf{W}^{1/2}\mathbf{A}\mathbf{W}^{1/2}\|_F,$$

where $\|\cdot\|_F$ is defined by $\|\mathbf{A}\|_F^2 = \sum_{i=1}^n \sum_{j=1}^n a_{ij}^2$.

The matrix \mathbf{W} can be chosen to be any matrix satisfying $\mathbf{W}\mathbf{s}_k = \mathbf{y}_k$. One possible \mathbf{W} is the inverse of the average Hessian. That is,

$$\mathbf{W} = \bar{\mathbf{G}}_k^{-1}, \quad \text{where } \bar{\mathbf{G}}_k = \int_0^1 \nabla^2 f(\mathbf{x}_k + t\alpha_k \mathbf{d}_k) dt$$

Using the Taylor series approximation of $\nabla f(\mathbf{x}_{k+1})$ around \mathbf{x}_k

$$\nabla f(\mathbf{x}_{k+1}) = \nabla f(\mathbf{x}_k) + \int_0^1 \nabla^2 f(\mathbf{x}_k + t\alpha_k \mathbf{d}_k) \alpha_k \mathbf{d}_k dt,$$

we can show that $\bar{\mathbf{G}}_k$ satisfies the secant equations:

$$\mathbf{y}_k = \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k) = \bar{\mathbf{G}}_k \alpha_k \mathbf{d}_k = \bar{\mathbf{G}}_k \mathbf{s}_k,$$

The unique solution to this optimization problem is given by

$$\mathbf{H}_{k+1} = (\mathbf{I} - \rho_k \mathbf{s}_k \mathbf{y}_k^T) \mathbf{H}_k (\mathbf{I} - \rho_k \mathbf{y}_k \mathbf{s}_k^T) + \rho_k \mathbf{s}_k \mathbf{s}_k^T,$$

where $\rho_k = \frac{1}{\mathbf{y}_k^T \mathbf{s}_k}$.

[H] Given a starting point \mathbf{x}_0 , convergence tolerance $\epsilon > 0$, and inverse Hessian approximation \mathbf{H}_0 $k \leftarrow 0$ $\|\nabla f(\mathbf{x}_k)\| > \epsilon$ Compute the search direction

$$\mathbf{d}_k = -\mathbf{H}_k \nabla f(\mathbf{x}_k)$$

Set $\mathbf{x}_{k+1} = \mathbf{x}_k + \alpha_k \mathbf{d}_k$, where α_k is computed via a line search procedure to satisfy the Wolfe conditions Define $\mathbf{s}_k = \mathbf{x}_{k+1} - \mathbf{x}_k$ and

$\mathbf{y}_k = \nabla f(\mathbf{x}_{k+1}) - \nabla f(\mathbf{x}_k)$ Compute \mathbf{H}_{k+1} using

$$\mathbf{H}_{k+1} = (\mathbf{I} - \rho_k \mathbf{s}_k \mathbf{y}_k^T) \mathbf{H}_k (\mathbf{I} - \rho_k \mathbf{y}_k \mathbf{s}_k^T) + \rho_k \mathbf{s}_k \mathbf{s}_k^T \quad k \leftarrow k + 1$$

The BFGS Algorithm

Thank you!

Appendix — Optimizing the Step Size

For general functions, BFGS converges globally if the Wolfe conditions are imposed on the line search.

Definition (Wolfe Conditions)

A step length α_k satisfies the Wolfe conditions in direction \mathbf{d}_k if the following inequalities hold:

- 1 $f(\mathbf{x}_k + \alpha_k \mathbf{d}_k) \leq f(\mathbf{x}_k) + c_1 \alpha_k \nabla f(\mathbf{x}_k)^T \mathbf{d}_k,$
- 2 $\nabla f(\mathbf{x}_k + \alpha_k \mathbf{d}_k)^T \mathbf{d}_k \geq c_2 \nabla f(\mathbf{x}_k)^T \mathbf{d}_k,$

where $0 < c_1 < c_2 < 1$. We may refer to the first condition as the sufficient decrease 1, and the second condition as the curvature condition.

Appendix — Initializing the Hessian

- We often set $\mathbf{H}_0 = \beta \mathbf{I}$, for some scalar β .
- A common heuristic is to initially set $\mathbf{H}_0 = \mathbf{I}$, then set $\mathbf{H}_0 \leftarrow \frac{\mathbf{y}_k^T \mathbf{s}_k}{\mathbf{y}_k^T \mathbf{y}_k} \mathbf{I}$.

Heuristic reasoning: Makes the 'size' of \mathbf{H}_0 similar to $(\nabla^2 f(\mathbf{x}_0))^{-1}$.

Assuming the average Hessian is positive definite, we can write

$\bar{\mathbf{G}}_k = \bar{\mathbf{G}}_k^{1/2} \bar{\mathbf{G}}_k^{1/2}$. Then using $\mathbf{y}_k = \bar{\mathbf{G}}_k \mathbf{s}_k$,

$$\frac{\mathbf{y}_k^T \mathbf{s}_k}{\mathbf{y}_k^T \mathbf{y}_k} = \frac{(\bar{\mathbf{G}}_k^{1/2} \mathbf{s}_k)^T \bar{\mathbf{G}}_k^{1/2} \mathbf{s}_k}{(\bar{\mathbf{G}}_k^{1/2} \mathbf{s}_k)^T \bar{\mathbf{G}}_k \bar{\mathbf{G}}_k^{1/2} \mathbf{s}_k} = \frac{\mathbf{z}_k^T \mathbf{z}_k}{\mathbf{z}_k^T \bar{\mathbf{G}}_k \mathbf{z}_k}$$

Appendix — Solving the Minimization Problem

We have the optimization problem

$$\begin{aligned} \text{minimize } \mathbf{H} \quad & \|\mathbf{H} - \mathbf{H}_k\|, \\ \text{subject to} \quad & \mathbf{H} = \mathbf{H}^T, \\ & \mathbf{H}\mathbf{y}_k = \mathbf{s}_k, \end{aligned}$$

where \mathbf{s}_k and \mathbf{y}_k satisfy the curvature condition and \mathbf{H}_k is symmetric and positive definite.

Appendix — Solving the Minimization Problem

We introduce the following notation:

$$\begin{aligned}\widehat{\mathbf{H}} &= \mathbf{W}^{1/2} \mathbf{H} \mathbf{W}^{1/2} \\ \widehat{\mathbf{H}}_k &= \mathbf{W}^{1/2} \mathbf{H}_k \mathbf{W}^{1/2} \\ \widehat{\mathbf{s}}_k &= \mathbf{W}^{1/2} \mathbf{s}_k \\ \widehat{\mathbf{y}}_k &= \mathbf{W}^{-1/2} \mathbf{y}_k\end{aligned}$$

Then we have the minimization problem

$$\begin{aligned}\text{minimize } \widehat{\mathbf{H}} \quad & \|\widehat{\mathbf{H}} - \widehat{\mathbf{H}}_k\|_F, \\ \text{subject to} \quad & \widehat{\mathbf{H}} \widehat{\mathbf{y}}_k = \widehat{\mathbf{y}}_k (= \widehat{\mathbf{s}}_k),\end{aligned}$$

Note that $\widehat{\mathbf{y}}_k$ is an eigenvector of $\widehat{\mathbf{H}}$.

Appendix — Solving the Minimization Problem

We introduce a new orthonormal bases: $\mathbf{U} = [\mathbf{u} \mid \mathbf{u}_\perp]$. Let \mathbf{u} be the normalized eigenvector $\hat{\mathbf{y}}_k$.

Multiplications by orthogonal matrices preserve length, and the Frobenius norm is unitary invariant. Hence

$$\|\hat{\mathbf{H}} - \hat{\mathbf{H}}_k\|_F = \|\mathbf{U}^T \hat{\mathbf{H}}_k \mathbf{U} - \mathbf{U}^T \hat{\mathbf{H}} \mathbf{U}\|_F = \left\| \begin{bmatrix} * & * \\ * & \mathbf{u}_\perp^T \hat{\mathbf{H}}_k \mathbf{u}_\perp \end{bmatrix} - \begin{bmatrix} 1 & 0 \\ 0 & \mathbf{u}_\perp^T \hat{\mathbf{H}} \mathbf{u}_\perp \end{bmatrix} \right\|_F$$

We can only modify the matrix on the right, and in this the only element we can change is the lower right element. In order to minimize the norm we must set the lower left elements of the matrices to be equal.

$$\mathbf{u}_\perp^T \hat{\mathbf{H}} \mathbf{u}_\perp = \mathbf{u}_\perp^T \hat{\mathbf{H}}_k \mathbf{u}_\perp$$

Appendix — Solving the Minimization Problem

$$\mathbf{u}_{\perp}^T \widehat{\mathbf{H}} \mathbf{u}_{\perp} = \mathbf{u}_{\perp}^T \widehat{\mathbf{H}}_k \mathbf{u}_{\perp}$$

This gives the optimal solution

$$\begin{aligned}\widehat{\mathbf{H}} &= \mathbf{U} \begin{bmatrix} 1 & 0 \\ 0 & \mathbf{u}_{\perp}^T \widehat{\mathbf{H}}_k \mathbf{u}_{\perp} \end{bmatrix} \mathbf{U}^T \\ &= [\mathbf{u} \quad \mathbf{u}_{\perp}] \begin{bmatrix} 1 & 0 \\ 0 & \mathbf{u}_{\perp}^T \widehat{\mathbf{H}}_k \mathbf{u}_{\perp} \end{bmatrix} \begin{bmatrix} \mathbf{u}^T \\ \mathbf{u}_{\perp}^T \end{bmatrix} \\ &= \mathbf{u}\mathbf{u}^T + \mathbf{u}_{\perp} \mathbf{u}_{\perp}^T \widehat{\mathbf{H}}_k \mathbf{u}_{\perp} \mathbf{u}_{\perp}^T \\ &= \mathbf{u}\mathbf{u}^T + (\mathbf{I} - \mathbf{u}\mathbf{u}^T) \widehat{\mathbf{H}}_k (\mathbf{I} - \mathbf{u}\mathbf{u}^T)\end{aligned}$$