

Modern Optimization Methods for Big Data Problems

MATH11146
The University of Edinburgh
Spring 2017

Peter Richtárik



31 / 78

Modern Optimization Methods for Big Data Problems

Lecture 3

Randomized Methods for Solving Linear Systems:
Equivalence and Exactness

January 23, 2017



32 / 78

Projection and Pseudoinverse



33 / 78

Projection Operators and Pseudoinverse Matrices - I

Definition 7

The **B-pseudoinverse** of a matrix \mathbf{M} , is defined as

$$\mathbf{M}^{\dagger_B} \stackrel{\text{def}}{=} \mathbf{B}^{-1} \mathbf{M}^{\top} (\mathbf{M} \mathbf{B}^{-1} \mathbf{M}^{\top})^{\dagger}, \quad (16)$$

where \dagger denotes the standard pseudoinverse.

Exercise 4

Show that

- (i) $\mathbf{A}^{\dagger} \mathbf{A}$ is a symmetric matrix
- (ii) $\mathbf{A}^{\top} (\mathbf{A} \mathbf{A}^{\top})^{\dagger} = \mathbf{A}^{\dagger}$
- (iii) The **I**-pseudoinverse is the standard Moore-Penrose pseudoinverse.



34 / 78

Projection Operators and Pseudoinverse Matrices - II

Lemma 8

The projection onto $\mathcal{L} = \{x : \mathbf{A}x = b\}$ is given by

$$\Pi_{\mathcal{L}}^{\mathbf{B}}(x) = x - \mathbf{B}^{-1}\mathbf{A}^{\top}(\mathbf{A}\mathbf{B}^{-1}\mathbf{A}^{\top})^{\dagger}(\mathbf{A}x - b) \stackrel{(16)}{=} x - \mathbf{A}^{\dagger\mathbf{B}}(\mathbf{A}x - b). \quad (17)$$

Proof.

Do it yourself. □

Exercise 5

Show that \mathbf{B} -pseudoinverse satisfies

$$\mathbf{A}^{\dagger\mathbf{B}}b = \Pi_{\mathcal{L}}^{\mathbf{B}}(0) = \arg \min_x \{\|x\|_{\mathbf{B}} : \mathbf{A}x = b\}.$$



35 / 78

Equivalence of Algorithms



36 / 78

Gradient and Hessian of $f_S(x)$ - I

In order to keep the expressions as brief as possible throughout, it will be useful to define

$$\mathbf{Z} \stackrel{\text{def}}{=} \mathbf{A}^\top \mathbf{H} \mathbf{A} \stackrel{(13)}{=} \mathbf{A}^\top \mathbf{S} (\mathbf{S}^\top \mathbf{A} \mathbf{B}^{-1} \mathbf{A}^\top \mathbf{S})^\dagger \mathbf{S}^\top \mathbf{A}. \quad (18)$$

Lemma 9

$\mathbf{B}^{-1} \mathbf{Z}$ is the projection, in the \mathbf{B} -norm, onto $\text{Range}(\mathbf{B}^{-1} \mathbf{A}^\top \mathbf{S})$. In particular,

$$(\mathbf{B}^{-1} \mathbf{Z})^2 = \mathbf{B}^{-1} \mathbf{Z} \quad \text{and} \quad \mathbf{Z} \mathbf{B}^{-1} \mathbf{Z} = \mathbf{Z}. \quad (19)$$

Recall from (3) that $f_S(x) \stackrel{\text{def}}{=} \frac{1}{2} \|\mathbf{A}x - b\|_{\mathbf{H}}^2 = \frac{1}{2} (\mathbf{A}x - b)^\top \mathbf{H} (\mathbf{A}x - b)$. By combining this with (18), this can be also written in the compact form

$$f_S(x) = \frac{1}{2} (x - x_*)^\top \mathbf{Z} (x - x_*), \quad (20)$$

where x_* is any point in \mathcal{L} .



37 / 78

Gradient and Hessian of $f_S(x)$ - II

Lemma 10

For each $x, h \in \mathbb{R}^n$ we have the expansion

$$f_S(x + h) = f_S(x) + \langle \nabla f_S(x), h \rangle_{\mathbf{B}} + \frac{1}{2} \langle (\nabla^2 f_S) h, h \rangle_{\mathbf{B}},$$

where

$$\nabla f_S(x) \stackrel{\text{def}}{=} \mathbf{B}^{-1} \mathbf{A}^\top \mathbf{H} (\mathbf{A}x - b) \quad \text{and} \quad \nabla^2 f_S \stackrel{\text{def}}{=} \mathbf{B}^{-1} \mathbf{Z} \quad (21)$$

are the gradient and Hessian of f_S with respect to the \mathbf{B} -inner product, respectively.²

In view of (18) and (21), the gradient can also be written as

$$\nabla f_S(x) = \mathbf{B}^{-1} \mathbf{Z} (x - x_*), \quad x \in \mathbb{R}^n, \quad x_* \in \mathcal{L}. \quad (22)$$

²If $\mathbf{B} = \mathbf{I}$, then $\langle \cdot, \cdot \rangle_{\mathbf{B}}$ is the standard Euclidean inner product, and we recover formulas for the standard gradient and Hessian. Note that $\mathbf{B}^{-1} \mathbf{Z}$ is both self-adjoint and positive semidefinite with respect to the \mathbf{B} -inner product. Indeed, for all $x, y \in \mathbb{R}^n$ we have $\langle \mathbf{B}^{-1} \mathbf{Z} x, y \rangle_{\mathbf{B}} = \langle \mathbf{Z} x, y \rangle_{\mathbf{I}} = \langle x, \mathbf{Z} y \rangle_{\mathbf{I}} = \langle x, \mathbf{B}^{-1} \mathbf{Z} y \rangle_{\mathbf{B}}$, and $\langle \mathbf{B}^{-1} \mathbf{Z} x, x \rangle_{\mathbf{B}} = \langle \mathbf{Z} x, x \rangle_{\mathbf{I}} \geq 0$.



38 / 78

Useful Identities Involving $f_{\mathbf{S}}(x)$

Lemma 11

For all $x \in \mathbb{R}^n$, we have

$$\begin{aligned}\nabla f_{\mathbf{S}}(x) &= (\nabla^2 f_{\mathbf{S}}) \nabla f_{\mathbf{S}}(x) = (\nabla^2 f_{\mathbf{S}})^{\dagger \mathbf{B}} \nabla f_{\mathbf{S}}(x) \\ &= x - \Pi_{\mathcal{L}_{\mathbf{S}}}^{\mathbf{B}}(x) = \mathbf{B}^{-1} \mathbf{A}^{\top} \mathbf{H}(\mathbf{A}x - b).\end{aligned}\tag{23}$$

Moreover,

$$f_{\mathbf{S}}(x) = \frac{1}{2} \|\nabla f_{\mathbf{S}}(x)\|_{\mathbf{B}}^2.\tag{24}$$

Finally, if $\mathcal{L}_{\mathbf{S}}$ is the set of minimizers of $f_{\mathbf{S}}$, then $\mathcal{L} \subseteq \mathcal{L}_{\mathbf{S}}$, and

- (i) $\mathcal{L}_{\mathbf{S}} = \{x : f_{\mathbf{S}}(x) = 0\} = \{x : \nabla f_{\mathbf{S}}(x) = 0\}$
- (ii) $\mathcal{L}_{\mathbf{S}} = x_* + \text{Null}(\mathbf{B}^{-1} \mathbf{Z})$ for all $x_* \in \mathcal{L}$
- (iii) $\mathcal{L}_{\mathbf{S}} = \{x : \mathbf{B}^{-1} \mathbf{A}^{\top} \mathbf{H} \mathbf{A} x = \mathbf{B}^{-1} \mathbf{A}^{\top} \mathbf{H} b\}$ (see (4))
- (iv) $\mathcal{L}_{\mathbf{S}} = \{x : \mathbf{S}^{\top} \mathbf{A} x = \mathbf{S}^{\top} b\}$ (see (6))



39 / 78

Some Consequences of Lemma 11

- The identity $(\nabla^2 f_{\mathbf{S}}) \nabla f_{\mathbf{S}}(x) = \nabla f_{\mathbf{S}}(x)$ means that the stochastic gradients of $f_{\mathbf{S}}$ are eigenvectors of the stochastic Hessian $\nabla^2 f_{\mathbf{S}}$, corresponding to eigenvalue one.
- The identity $(\nabla^2 f_{\mathbf{S}})^{\dagger \mathbf{B}} \nabla f_{\mathbf{S}}(x) = \nabla f_{\mathbf{S}}(x)$ means that the stochastic gradients of $f_{\mathbf{S}}$ are eigenvectors of the \mathbf{B} -pseudoinverse of the stochastic Hessian $\nabla^2 f_{\mathbf{S}}$, corresponding to eigenvalue one.
- Function f can be represented in multiple ways:

$$f(x) = \frac{1}{2} \mathbb{E} [\|x - \Pi_{\mathcal{L}_{\mathbf{S}}}^{\mathbf{B}}(x)\|_{\mathbf{B}}^2] = \frac{1}{2} \mathbb{E} [\|\nabla f_{\mathbf{S}}(x)\|_{\mathbf{B}}^2].\tag{25}$$

- The gradient and Hessian of f (with respect to the \mathbf{B} -inner product) are given by

$$\nabla f(x) = \mathbf{B}^{-1} \mathbb{E}[\mathbf{Z}](x - x_*), \quad \text{and} \quad \nabla^2 f = \mathbf{B}^{-1} \mathbb{E}[\mathbf{Z}],\tag{26}$$

respectively, where x_* is any point in \mathcal{L} .



40 / 78

Equivalence of Algorithms

Theorem 12

Algorithm 1 (Basic Method) can be equivalently written as stochastic gradient descent (8), stochastic Newton method (9), stochastic fixed point method (10), and stochastic projection method (11).

Proof.

This follows from identities (23) in Lemma 11. □



41 / 78

Proof of Lemma 11 - I

Pick any $x_* \in \mathcal{L}$. First, we have

$$\Pi_{\mathcal{L}_S}^{\mathbf{B}}(x) \stackrel{(17)}{=} x - \mathbf{B}^{-1} \mathbf{A}^\top \mathbf{H}(\mathbf{A}x - b) \stackrel{(21)}{=} x - \nabla f_S(x).$$

To establish (23), it now only remains to consider the two expressions involving the Hessian. We have

$$\nabla^2 f_S \nabla f_S(x) \stackrel{(21)+(22)}{=} \mathbf{B}^{-1} \mathbf{Z} \mathbf{B}^{-1} \mathbf{Z}(x - x_*) \stackrel{(19)}{=} \mathbf{B}^{-1} \mathbf{Z}(x - x_*) \stackrel{(22)}{=} \nabla f_S(x),$$

and

$$\begin{aligned} (\nabla^2 f_S)^\dagger \nabla f_S(x) &\stackrel{(16)}{=} \mathbf{B}^{-1} (\nabla^2 f_S)^\top ((\nabla^2 f_S) \mathbf{B}^{-1} (\nabla^2 f_S)^\top)^\dagger \nabla f_S(x) \\ &\stackrel{(21)}{=} \mathbf{B}^{-1} (\mathbf{B}^{-1} \mathbf{Z})^\top ((\mathbf{B}^{-1} \mathbf{Z}) \mathbf{B}^{-1} (\mathbf{B}^{-1} \mathbf{Z})^\top)^\dagger \mathbf{B}^{-1} \mathbf{Z}(x - x_*) \\ &= \mathbf{B}^{-1} \mathbf{Z} \mathbf{B}^{-1} (\mathbf{B}^{-1} \mathbf{Z} \mathbf{B}^{-1} \mathbf{Z} \mathbf{B}^{-1})^\dagger \mathbf{B}^{-1} \mathbf{Z}(x - x_*) \\ &\stackrel{(19)}{=} (\mathbf{B}^{-1} \mathbf{Z} \mathbf{B}^{-1}) (\mathbf{B}^{-1} \mathbf{Z} \mathbf{B}^{-1})^\dagger (\mathbf{B}^{-1} \mathbf{Z} \mathbf{B}^{-1}) \mathbf{B}(x - x_*) \\ &= \mathbf{B}^{-1} \mathbf{Z}(x - x_*) \\ &\stackrel{(22)}{=} \nabla f_S(x). \end{aligned}$$



42 / 78

Proof of Lemma 11 - II

Identity (24) follows from

$$\frac{1}{2} \|\nabla f_{\mathbf{S}}(x)\|_{\mathbf{B}}^2 \stackrel{(22)}{=} \frac{1}{2} (x - x_*)^\top \mathbf{Z} \mathbf{B}^{-1} \mathbf{Z} (x - x_*) \stackrel{(19)}{=} \frac{1}{2} (x - x_*)^\top \mathbf{Z} (x - x_*) \stackrel{(20)}{=} f_{\mathbf{S}}(x).$$

If $x \in \mathcal{L}$, then by picking $x_* = x$ in (22), we see that $x \in \mathcal{L}_{\mathbf{S}}$.

It remains to show that the sets defined in (i)–(iv) are identical.

- ▶ Equivalence between (i) and (ii) follows from (22).
- ▶ Now consider (ii) and (iii). Any $x_* \in \mathcal{L}$ belongs to the set defined in (iii), which follows immediately by substituting $b = \mathbf{A}x_*$. The rest follows after observing the nullspaces are identical.
- ▶ In order to show that (iii) and (iv) are equivalent, it suffices to compute $\Pi_{\mathcal{L}_{\mathbf{S}}}^{\mathbf{B}}(x)$ and observe that $\Pi_{\mathcal{L}_{\mathbf{S}}}^{\mathbf{B}}(x) = x$ if and only if x belongs to the set defined in (iii).



43 / 78

Equivalence of 4 Stochastic Reformulations



44 / 78

Equivalence of the Stochastic Formulations

The below theorem says that the solution sets of the four stochastic problems (2), (4), (5), and (6) are identical. **In this sense, the four stochastic problems are equivalent.**

Theorem 13 (Equivalence of stochastic formulations)

Let $x_* \in \mathcal{L}$. The following sets are identical:

- (i) $\mathcal{X} = \arg \min f(x) = \{x : f(x) = 0\} = \{x : \nabla f(x) = 0\} \rightarrow (2)$
- (ii) $\mathcal{X} = \{x : \mathbf{B}^{-1} \mathbf{A}^\top \mathbb{E}[\mathbf{H}] \mathbf{A}x = \mathbf{B}^{-1} \mathbf{A}^\top \mathbb{E}[\mathbf{H}] b\} = x_* + \text{Null}(\mathbb{E}[\mathbf{Z}]) \rightarrow (4)$
- (iii) $\mathcal{X} = \{x : \mathbb{E}[\Pi_{\mathcal{L}_S}^{\mathbf{B}}(x)] = x\} \rightarrow (5)$
- (iv) $\mathcal{X} = \{x : \text{Prob}(x \in \mathcal{L}_S) = 1\} \rightarrow (6)$

Moreover, \mathcal{X} does not depend on \mathbf{B} .



45 / 78

Proof of Theorem 13 - Part I

As f is convex, nonnegative and achieving the value of zero (since $\mathcal{L} \neq \emptyset$), the sets in (i) are all identical. We shall now show that the sets defined in (ii)–(iv) are equal to that defined in (i).

(i) \leftrightarrow (ii): Using the formula for the gradient from (26), we see that

$$\begin{aligned} \{x : \nabla f(x) = 0\} &= \{x : \mathbf{B}^{-1} \mathbb{E}[\mathbf{Z}](x - x_*) = 0\} \\ &= \{x : \mathbb{E}[\mathbf{Z}](x - x_*) = 0\} \\ &= x_* + \{h : \mathbb{E}[\mathbf{Z}]h = 0\} \\ &= x_* + \text{Null}(\mathbb{E}[\mathbf{Z}]), \end{aligned}$$

which shows that (i) and (ii) are the same.

(i) \leftrightarrow (iii): Equivalence of (i) and (iii) follows by taking expectations in (23) to obtain

$$\nabla f(x) = \mathbb{E}[\nabla f_S(x)] \stackrel{(23)}{=} \mathbb{E}[x - \Pi_{\mathcal{L}_S}^{\mathbf{B}}(x)].$$



46 / 78

Proof of Theorem 13 - Part II

(i) \leftrightarrow (iv): It remains to establish equivalence between (i) and (iv). Let

$$\mathcal{X} = \{x : f(x) = 0\} \stackrel{(25)}{=} \left\{x : E \left[\|x - \Pi_{\mathcal{L}_S}^{\mathbf{B}}(x)\|_{\mathbf{B}}^2 \right] = 0 \right\} \quad (27)$$

and let \mathcal{X}' be the set from (iv).

We need to show that $\mathcal{X}' = \mathcal{X}$. For easier reference, let

$$\xi_S(x) \stackrel{\text{def}}{=} \|x - \Pi_{\mathcal{L}_S}^{\mathbf{B}}(x)\|_{\mathbf{B}}^2.$$

Note that the following three probabilistic events are identical:

$$[x \in \mathcal{L}_S] = [x = \Pi_{\mathcal{L}_S}^{\mathbf{B}}(x)] = [\xi_S(x) = 0]. \quad (28)$$

We first show that $\mathcal{X}' \subseteq \mathcal{X}$.

In view of (28), if $x \in \mathcal{X}'$, then the random variable $\xi_S(x)$ is equal to zero with probability 1, which implies $E[\xi_S(x)] = 0$, whence $x \in \mathcal{X}$.



47 / 78

Proof of Theorem 13 - Part III

Let us now show that $\mathcal{X} \subseteq \mathcal{X}'$.

Let $1_{[\xi_S(x) \geq t]}$ be the indicator function of the event $[\xi_S(x) \geq t]$. Note that since $\xi_S(x)$ is a nonnegative random variable, for all $t \in \mathbb{R}$ we have the inequality

$$\xi_S(x) \geq t 1_{[\xi_S(x) \geq t]}. \quad (29)$$

Now take $x \in \mathcal{X}$ and consider $t > 0$. By taking expectations in (29), we obtain

$$0 = E[\xi_S(x)] \geq E[t 1_{[\xi_S(x) \geq t]}] = t E[1_{[\xi_S(x) \geq t]}] = t \text{Prob}(\xi_S(x) \geq t),$$

which implies that $\text{Prob}(\xi_S(x) \geq t) = 0$. Now choose $t_i = 1/i$ for $i = 1, 2, \dots$ and note that the event $[\xi_S(x) > 0]$ can be written as

$$[\xi_S(x) > 0] = \bigcup_{i=1}^{\infty} [\xi_S(x) \geq t_i].$$



48 / 78

Proof of Theorem 13 - Part IV

Therefore, by the union bound,

$$\text{Prob}(\xi_{\mathbf{s}}(x) > 0) \leq \sum_{i=1}^{\infty} \text{Prob}(\xi_{\mathbf{s}}(x) \geq t_i) = 0,$$

which immediately implies that $\text{Prob}(\xi_{\mathbf{s}}(x) = 0) = 1$. From (28) we conclude that $x \in \mathcal{X}'$.

Independence on \mathbf{B} . Since characterization (iv) of \mathcal{X} does not depend on \mathbf{B} , we conclude that \mathcal{X} does not depend on \mathbf{B} .



49 / 78

Exactness of the Reformulations



50 / 78

Rangespace and Nullspace of a Matrix - I

Let $\mathbf{M} \in \mathbb{R}^{m \times n}$.

Definition 14 (Rangespace of a matrix)

By $\text{Range}(\mathbf{M})$ we mean the **rangespace of matrix \mathbf{M}** . This is the linear subspace of \mathbb{R}^m generated by the columns of \mathbf{M} :

$$\text{Range}(\mathbf{M}) \stackrel{\text{def}}{=} \{\mathbf{M}\mathbf{x} : \mathbf{x} \in \mathbb{R}^n\} = \left\{ \sum_j \mathbf{M}_{:j} x_j, \quad \mathbf{x} \in \mathbb{R}^n \right\}.$$

Definition 15 (Nullspace of a matrix)

By $\text{Null}(\mathbf{M})$ we mean the **nullspace of matrix \mathbf{M}** . This is the linear subspace of \mathbb{R}^n formed by the vectors orthogonal (under standard Euclidean inner product) to all rows of \mathbf{M} :

$$\text{Null}(\mathbf{M}) \stackrel{\text{def}}{=} \{\mathbf{x} \in \mathbb{R}^n : \mathbf{M}\mathbf{x} = \mathbf{0}\} = \{\mathbf{x} \in \mathbb{R}^n : \langle \mathbf{M}_{i:}^\top, \mathbf{x} \rangle = 0 \quad \forall i\}.$$



51 / 78

Rangespace and Nullspace of a Matrix - II

Definition 16 (Orthogonal complement)

Let X be a subspace of a vector space. The **orthogonal complement** of X is the linear subspace $X^\perp \stackrel{\text{def}}{=} \{\mathbf{y} : \langle \mathbf{y}, \mathbf{x} \rangle = 0 \quad \forall \mathbf{x} \in X\}$.

Here we collect some useful identities involving rangespaces and nullspaces of a matrix:

Fact 17

For any $\mathbf{M} \in \mathbb{R}^{m \times n}$, we have

- (i) $\text{Range}(\mathbf{M}^\top) = \text{Null}(\mathbf{M})^\perp$
- (ii) $\text{Range}(\mathbf{M}^\top)^\perp = \text{Null}(\mathbf{M})$
- (iii) If $\mathbf{G} \succ 0$, then $\text{Null}(\mathbf{M}^\top \mathbf{G} \mathbf{M}) = \text{Null}(\mathbf{M})$



52 / 78

Exactness

Key Question: When are the stochastic formulations (2), (4), (5), (6) equivalent to the linear system (1)? That is, when is their set of solutions \mathcal{X} identical to solution set of the linear system \mathcal{L} ?

This leads to the concept of **exactness**:

Assumption 3 (Exactness)

*Stochastic reformulations (2), (4), (5), (6) of problem (1) are exact.
That is, $\mathcal{X} = \mathcal{L}$.*

In what follows, we will

- ▶ Give **sufficient**, and **necessary & sufficient** conditions for exactness.
- ▶ Use this assumption to **prove convergence** of the algorithms to a specific point in \mathcal{L} .



53 / 78

Necessary and Sufficient Conditions for Exactness

Theorem 18 (\Leftrightarrow Conditions for exactness)

The following statements are equivalent:

- (i) *Assumption 3 (exactness) holds*
- (ii) $\text{Null}(\mathbb{E}[\mathbf{Z}]) = \text{Null}(\mathbf{A})$
- (iii) $\text{Null}(\mathbf{B}^{-1/2}\mathbb{E}[\mathbf{Z}]\mathbf{B}^{-1/2}) = \text{Null}(\mathbf{A}\mathbf{B}^{-1/2})$
- (iv) $\text{Range}(\mathbf{A}) \cap \text{Null}(\mathbb{E}[\mathbf{H}]) = \{0\}$



54 / 78

Proof of Theorem 18 - I

(i) \leftrightarrow (ii): Choose any $x_* \in \mathcal{L}$. We know that $\mathcal{L} = x_* + \text{Null}(\mathbf{A})$. On the other hand, Theorem 13 says that $\mathcal{X} = x_* + \text{Null}(\mathbb{E}[\mathbf{Z}])$.

(ii) \leftrightarrow (iii): If (ii) holds, then

$$\text{Null}(\mathbf{A}) = \text{Null}(\mathbb{E}[\mathbf{Z}]) = \text{Null}(\mathbf{B}^{-1/2}\mathbb{E}[\mathbf{Z}]),$$

and (iii) follows. If (iii) holds, then

$$\text{Null}(\mathbf{A}) = \text{Null}(\mathbf{B}^{-1/2}\mathbb{E}[\mathbf{Z}]) = \text{Null}(\mathbb{E}[\mathbf{Z}]),$$

proving (ii).



55 / 78

Proof of Theorem 18 - II

(ii) \leftrightarrow (iv): First, note that

$$\mathbb{E}[\mathbf{Z}] = \mathbf{A}^\top (\mathbb{E}[\mathbf{H}])^{1/2} (\mathbb{E}[\mathbf{H}])^{1/2} \mathbf{A}.$$

In view of Fact 17, for any matrix \mathbf{M} we have $\text{Null}(\mathbf{M}^\top \mathbf{M}) = \text{Null}(\mathbf{M})$. Therefore,

$$\text{Null}(\mathbb{E}[\mathbf{Z}]) = \text{Null}\left((\mathbb{E}[\mathbf{H}])^{1/2} \mathbf{A}\right).$$

Moreover, we know that

(a) $\text{Null}((\mathbb{E}[\mathbf{H}])^{1/2} \mathbf{A}) = \text{Null}(\mathbf{A})$ if and only if $\text{Range}(\mathbf{A}) \cap \text{Null}((\mathbb{E}[\mathbf{H}])^{1/2}) = \{0\}$, and

(b) $\text{Null}((\mathbb{E}[\mathbf{H}])^{1/2}) = \text{Null}(\mathbb{E}[\mathbf{H}])$ (see Fact 17).

It remains to combine these observations.



56 / 78

Sufficient Conditions for Exactness

We now list some sufficient conditions for exactness.

Lemma 19 (Sufficient conditions for exactness)

Any of these conditions implies that Assumption 3 is satisfied:

- (i) $E[\mathbf{H}] \succ 0$
- (ii) $\text{Null}(E[\mathbf{H}]) \subseteq \text{Null}(\mathbf{A}^\top)$

Proof.

If (i) holds, then $\text{Null}(E[\mathbf{Z}]) = \text{Null}(\mathbf{A}^\top E[\mathbf{H}] \mathbf{A}) = \text{Null}(\mathbf{A})$, where the last equality follows from Fact 17. Exactness now follows by applying Theorem 18.

On the other hand, in view of Fact 17, (ii) implies statement (iv) in Theorem 18, and hence exactness follows. □



57 / 78

Condition Number



58 / 78

Spectral Decomposition of the Hessian of f

Recall that the **Hessian of f** is given by

$$\nabla^2 f = \mathbb{E}_{\mathbf{s} \sim \mathcal{D}} [\nabla^2 f_{\mathbf{s}}] = \mathbf{B}^{-1} \mathbb{E} [\mathbf{Z}]. \quad (30)$$

Lemma 20

Matrices $\mathbf{B}^{-1} \mathbb{E} [\mathbf{Z}]$ and $\mathbf{B}^{-1/2} \mathbb{E} [\mathbf{Z}] \mathbf{B}^{-1/2}$ have the same eigenvalues.

Proof.

It is known that for any $\mathbf{X}, \mathbf{Y} \in \mathbb{R}^{n \times n}$, the matrices \mathbf{XY} and \mathbf{YX} have the same eigenvalues. It only remains to apply this to $\mathbf{X} = \mathbf{B}^{-1/2}$ and $\mathbf{Y} = \mathbf{B}^{-1/2} \mathbb{E} [\mathbf{Z}]$. □

The above result allows us to study spectral properties of the Hessian $\nabla^2 f$ through the **eigenvalue decomposition** of the symmetric positive definite matrix $\mathbf{B}^{-1/2} \mathbb{E} [\mathbf{Z}] \mathbf{B}^{-1/2}$.



59 / 78

Eigenvalues of the Hessian of f

Let

$$\mathbf{W} \stackrel{\text{def}}{=} \mathbf{B}^{-1/2} \mathbb{E} [\mathbf{Z}] \mathbf{B}^{-1/2} = \mathbf{U} \mathbf{\Lambda} \mathbf{U}^\top = \sum_{i=1}^n \lambda_i \mathbf{u}_i \mathbf{u}_i^\top \quad (31)$$

be the **eigenvalue decomposition of \mathbf{W}** , where

$$\mathbf{U} = [\mathbf{u}_1, \dots, \mathbf{u}_n] \in \mathbb{R}^{n \times n}$$

is an orthonormal matrix composed of **eigenvectors** (i.e., we have $\mathbf{U} \mathbf{U}^\top = \mathbf{U}^\top \mathbf{U} = \mathbf{I}$), and

$$\mathbf{\Lambda} = \mathbf{Diag}(\lambda_1, \lambda_2, \dots, \lambda_n)$$

is a diagonal matrix of **eigenvalues**. Assume without loss of generality that the eigenvalues are ordered from largest to smallest:

$$\lambda_{\max} \stackrel{\text{def}}{=} \lambda_1 \geq \lambda_2 \geq \dots \geq \lambda_n \stackrel{\text{def}}{=} \lambda_{\min}.$$



60 / 78

All Eigenvalues of \mathbf{W} are Between 0 and 1

Lemma 21

$0 \leq \lambda_i \leq 1$ for all i .

Proof.

Since $\mathbf{B}^{-1/2}\mathbf{Z}\mathbf{B}^{-1/2}$ is symmetric positive semidefinite, so is its expectation \mathbf{W} , implying that $\lambda_i \geq 0$ for all i .

Further, note that $\mathbf{B}^{-1/2}\mathbf{Z}\mathbf{B}^{-1/2}$ is a projection matrix. Indeed, it is the projection (in the standard \mathbf{I} -norm) onto $\text{Range}(\mathbf{B}^{-1/2}\mathbf{A}^\top \mathbf{S})$. Therefore, its eigenvalues are all zeros or ones. Since the map $\mathbf{X} \mapsto \lambda_{\max}(\mathbf{X})$ is convex, by **Jensen's inequality** we get

$$\lambda_{\max}(\mathbf{W}) = \lambda_{\max}\left(\mathbb{E}\left[\mathbf{B}^{-1/2}\mathbf{Z}\mathbf{B}^{-1/2}\right]\right) \leq \mathbb{E}\left[\lambda_{\max}(\mathbf{B}^{-1/2}\mathbf{Z}\mathbf{B}^{-1/2})\right] \leq 1.$$

□



61 / 78

Smallest Nonzero Eigenvalue

Lemma 22

If Assumption 3 (exactness) holds, then $\lambda_{\max} > 0$.

Proof.

Assume, by contradiction, that $\lambda_i = 0$ for all i . Then from Theorem 18 and the fact that $\text{Null}(\mathbf{W}) = \text{Range}(u_i : \lambda_i = 0)$ we conclude that $\text{Null}(\mathbf{A}\mathbf{B}^{-1/2}) = \mathbb{R}^n$, which in turn implies that $\text{Null}(\mathbf{A}) = \mathbb{R}^n$. This can only happen if $\mathbf{A} = 0$, which contradicts with our assumption on \mathbf{A} . □

Now, let j be the largest index for which $\lambda_j > 0$. This identifies the **smallest nonzero eigenvalue of \mathbf{W}** , which we shall denote as

$$\lambda_{\min}^+ = \lambda_j.$$

If all eigenvalues $\{\lambda_i\}$ are positive, then $j = n$.



62 / 78

Condition Number

Definition 23

The **condition number** associated with the four stochastic reformulations is the quantity³

$$\zeta(\mathbf{A}, \mathbf{B}, \mathcal{D}) = \zeta \stackrel{\text{def}}{=} \|\mathbf{W}\| \|\mathbf{W}^\dagger\| = \frac{\lambda_{\max}}{\lambda_{\min}^+}. \quad (32)$$

Remark:

- ▶ As we shall see, convergence rate of the Basic Method is described by ζ .
- ▶ As one varies the parameters defining the reformulation (i.e., \mathcal{D} and \mathbf{B}), ζ changes. As a general rule of thumb, simple distributions will lead to reformulations with a small condition number. For instance, choosing $\mathbf{S} = \mathbf{I}$ with probability one gives $\zeta = 1$. However, in such a case each step of the Basic Method is very expensive. One needs to strike the right balance.

³ $\|\mathbf{X}\|$ denotes the **spectral norm** of \mathbf{X} . In general, $\|\mathbf{X}\| = (\lambda_{\max}(\mathbf{X}^\top \mathbf{X}))^{1/2}$. If \mathbf{X} is symmetric positive semidefinite, then $\|\mathbf{X}\|^2 = \lambda_{\max}(\mathbf{X}^\top \mathbf{X}) = \lambda_{\max}(\mathbf{X}^2) = (\lambda_{\max}(\mathbf{X}))^2$. Therefore, $\|\mathbf{X}\| = \lambda_{\max}(\mathbf{X})$.

