

Acceleration of Randomized Kaczmarz Method

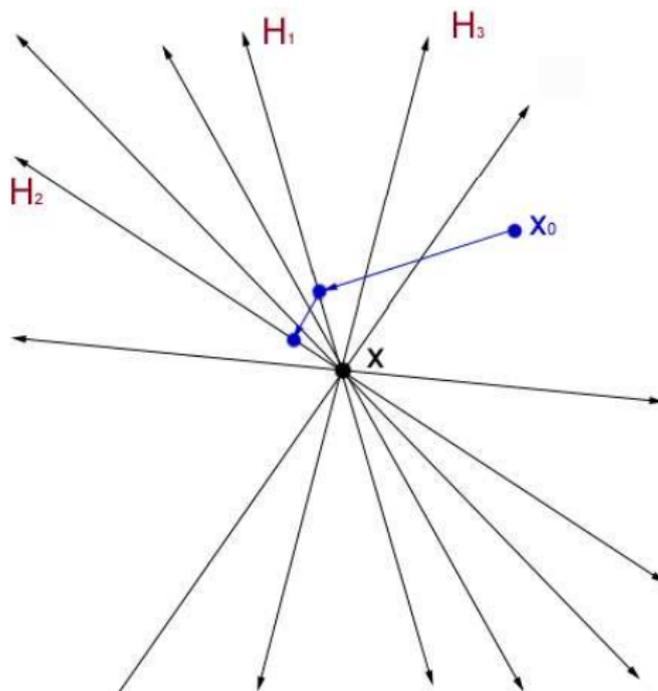
Deanna Needell [Joint work with Y. Eldar]

Stanford University

BIRS Banff, March 2011

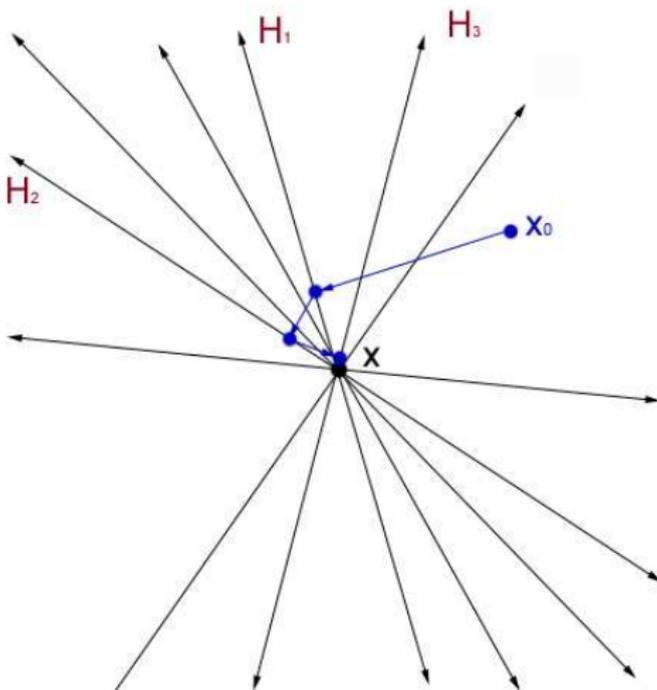
Geometrically

Denote $H_i = \{w : \langle a_i, w \rangle = b[i]\}$.



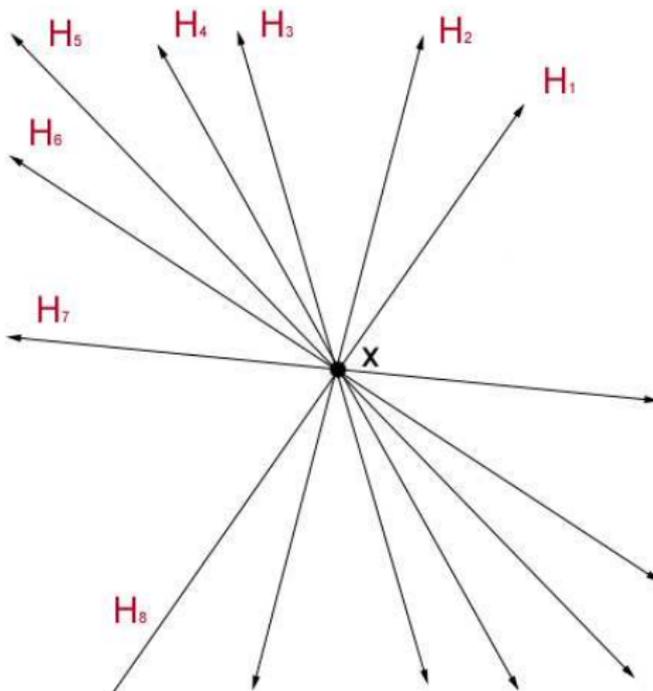
Geometrically

Denote $H_i = \{w : \langle a_i, w \rangle = b[i]\}$.



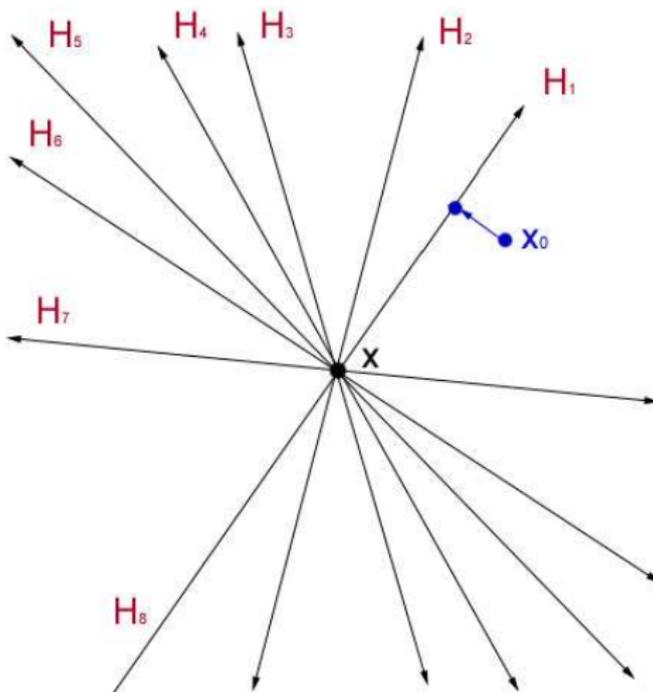
But what if...

Denote $H_i = \{w : \langle a_i, w \rangle = b[i]\}$.



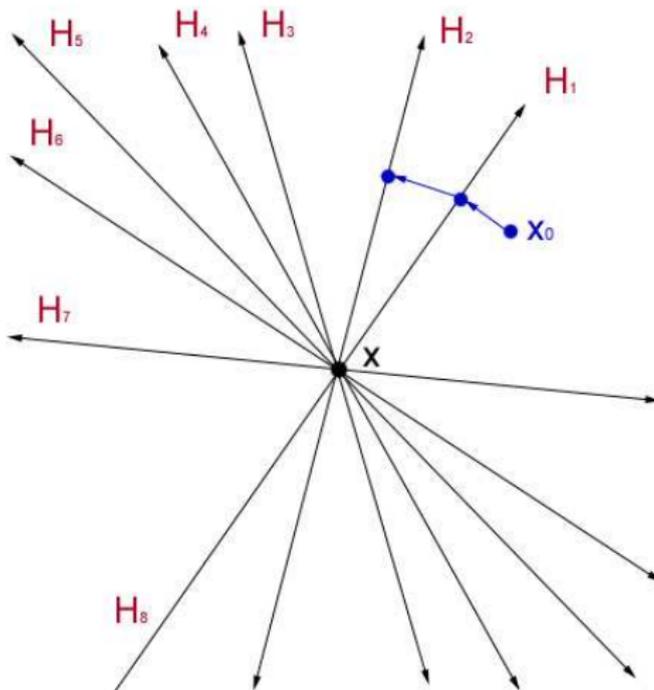
But what if...

Denote $H_i = \{w : \langle a_i, w \rangle = b[i]\}$.



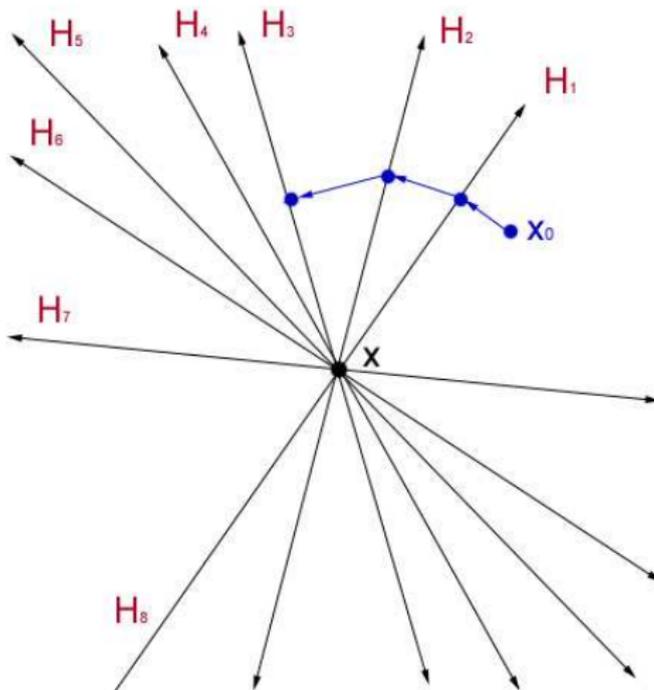
But what if...

Denote $H_i = \{w : \langle a_i, w \rangle = b[i]\}$.



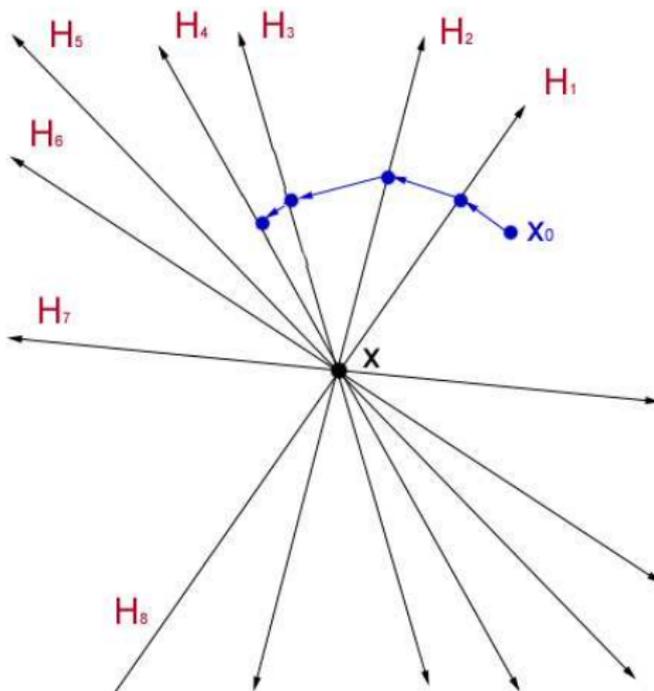
But what if...

Denote $H_i = \{w : \langle a_i, w \rangle = b[i]\}$.



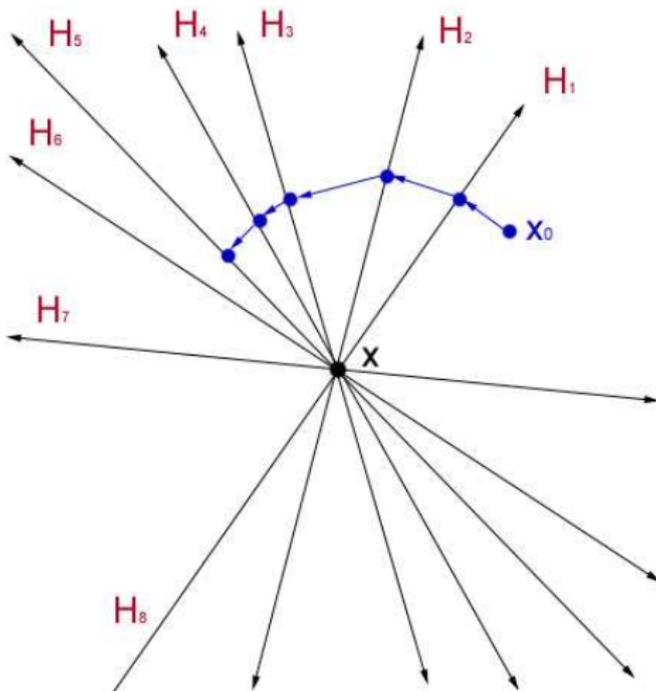
But what if...

Denote $H_i = \{w : \langle a_i, w \rangle = b[i]\}$.



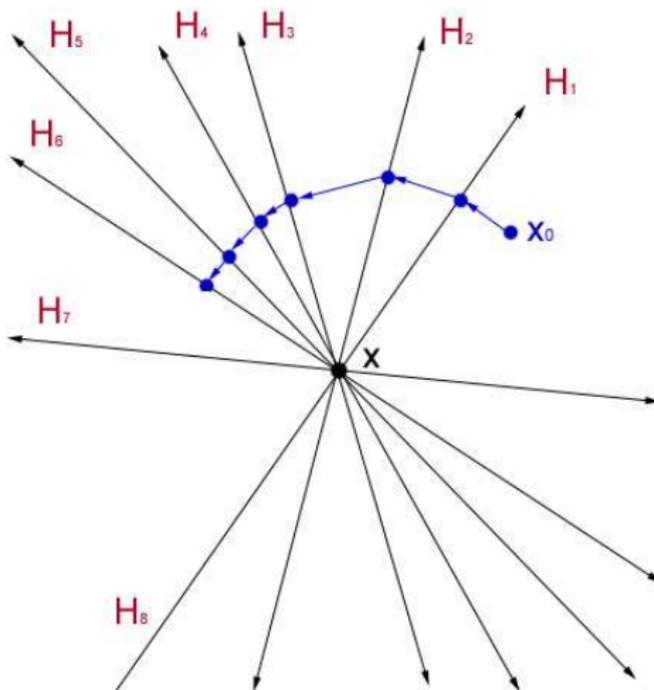
But what if...

Denote $H_i = \{w : \langle a_i, w \rangle = b[i]\}$.



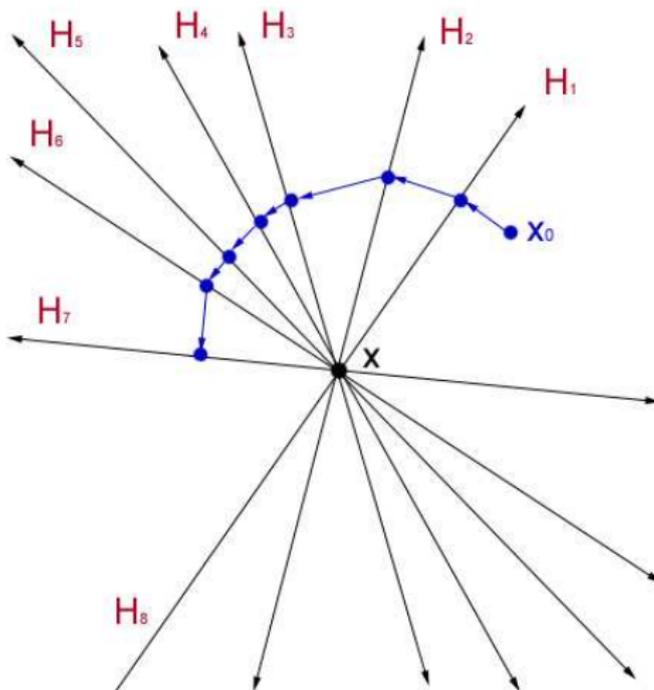
But what if...

Denote $H_i = \{w : \langle a_i, w \rangle = b[i]\}$.



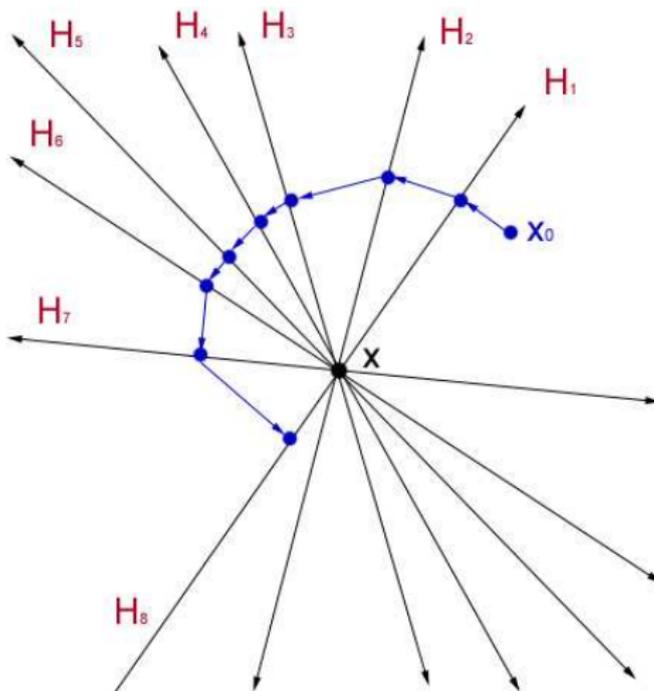
But what if...

Denote $H_i = \{w : \langle a_i, w \rangle = b[i]\}$.



But what if...

Denote $H_i = \{w : \langle a_i, w \rangle = b[i]\}$.



Randomized Kaczmarz

$$\begin{bmatrix} \text{--- } a_1 \text{ ---} \\ \text{--- } a_2 \text{ ---} \\ \vdots \\ \vdots \\ \vdots \\ \text{--- } a_m \text{ ---} \end{bmatrix} \cdot \begin{bmatrix} x \\ \vdots \\ \vdots \\ \vdots \end{bmatrix} = \begin{bmatrix} b[1] \\ b[2] \\ \vdots \\ \vdots \\ \vdots \\ b[m] \end{bmatrix}$$

Kaczmarz

- 1 Start with initial guess x_0
- 2 $x_{k+1} = x_k + \frac{b[j] - \langle a_j, x_k \rangle}{\|a_j\|_2^2} a_j$ where i is chosen *randomly*
- 3 Repeat (2)

Randomized Kaczmarz

$$\begin{bmatrix} \text{--- } a_1 \text{ ---} \\ \text{--- } a_2 \text{ ---} \\ \vdots \\ \vdots \\ \vdots \\ \text{--- } a_m \text{ ---} \end{bmatrix} \cdot \begin{bmatrix} x \\ \vdots \\ x \end{bmatrix} = \begin{bmatrix} b[1] \\ b[2] \\ \vdots \\ \vdots \\ \vdots \\ b[m] \end{bmatrix}$$

Kaczmarz

- 1 Start with initial guess x_0
- 2 $x_{k+1} = x_k + \frac{b[j] - \langle a_j, x_k \rangle}{\|a_j\|_2^2} a_j$ where i is chosen *randomly*
- 3 Repeat (2)

Randomized Kaczmarz

Strohmer-Vershynin

1 Start with initial guess x_0

2 $x_{k+1} = x_k + \frac{b_p - \langle a_p, x_k \rangle}{\|a_p\|_2^2} a_p$ where $\mathbb{P}(p = i) = \frac{\|a_i\|_2^2}{\|A\|_F^2}$

3 Repeat (2)

Randomized Kaczmarz

Strohmer-Vershynin

1 Start with initial guess x_0

2 $x_{k+1} = x_k + \frac{b_p - \langle a_p, x_k \rangle}{\|a_p\|_2^2} a_p$ where $\mathbb{P}(p = i) = \frac{\|a_i\|_2^2}{\|A\|_F^2}$

3 Repeat (2)

Randomized Kaczmarz (RK)

Strohmer-Vershynin

- Let $R = \|A^{-1}\|^2 \|A\|_F^2$ ($\|A^{-1}\| \stackrel{\text{def}}{=} \inf\{M : M\|Ax\|_2 \geq \|x\|_2 \text{ for all } x\}$)
- Then $\mathbb{E}\|x_k - x\|_2^2 \leq \left(1 - \frac{1}{R}\right)^k \|x_0 - x\|_2^2$
- Well conditioned $A \rightarrow$ Convergence in $O(n)$ iterations $\rightarrow O(n^2)$ total runtime.
- Better than $O(mn^2)$ runtime for Gaussian elimination and empirically often faster than Conjugate Gradient.

Randomized Kaczmarz (RK)

Strohmer-Vershynin

- Let $R = \|A^{-1}\|^2 \|A\|_F^2$ ($\|A^{-1}\| \stackrel{\text{def}}{=} \inf\{M : M\|Ax\|_2 \geq \|x\|_2 \text{ for all } x\}$)
- Then $\mathbb{E}\|x_k - x\|_2^2 \leq \left(1 - \frac{1}{R}\right)^k \|x_0 - x\|_2^2$
- Well conditioned $A \rightarrow$ Convergence in $O(n)$ iterations $\rightarrow O(n^2)$ total runtime.
- Better than $O(mn^2)$ runtime for Gaussian elimination and empirically often faster than Conjugate Gradient.

Randomized Kaczmarz (RK)

Strohmer-Vershynin

- Let $R = \|A^{-1}\|^2 \|A\|_F^2$ ($\|A^{-1}\| \stackrel{\text{def}}{=} \inf\{M : M\|Ax\|_2 \geq \|x\|_2 \text{ for all } x\}$)
- Then $\mathbb{E}\|x_k - x\|_2^2 \leq \left(1 - \frac{1}{R}\right)^k \|x_0 - x\|_2^2$
- Well conditioned $A \rightarrow$ Convergence in $O(n)$ iterations $\rightarrow O(n^2)$ total runtime.
- Better than $O(mn^2)$ runtime for Gaussian elimination and empirically often faster than Conjugate Gradient.

Randomized Kaczmarz (RK)

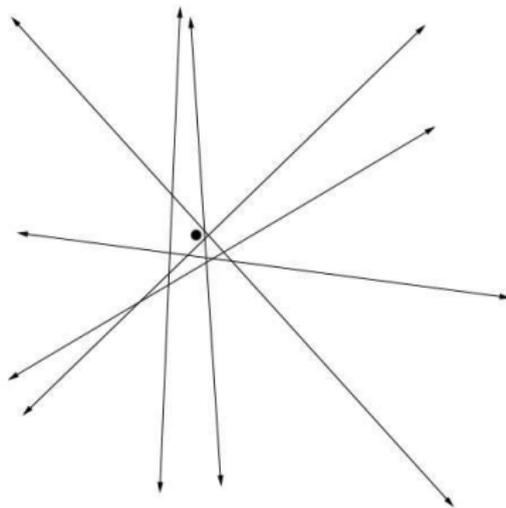
Strohmer-Vershynin

- Let $R = \|A^{-1}\|^2 \|A\|_F^2$ ($\|A^{-1}\| \stackrel{\text{def}}{=} \inf\{M : M\|Ax\|_2 \geq \|x\|_2 \text{ for all } x\}$)
- Then $\mathbb{E}\|x_k - x\|_2^2 \leq \left(1 - \frac{1}{R}\right)^k \|x_0 - x\|_2^2$
- Well conditioned $A \rightarrow$ Convergence in $O(n)$ iterations $\rightarrow O(n^2)$ total runtime.
- Better than $O(mn^2)$ runtime for Gaussian elimination and empirically often faster than Conjugate Gradient.

Randomized Kaczmarz (RK) with noise

System with noise

We now consider the consistent system $Ax = b$ corrupted by noise to form the possibly inconsistent system $Ax \approx b + z$.



Randomized Kaczmarz (RK) with noise

Theorem [N]

- Let $Ax = b$ be corrupted with noise: $Ax \approx b + z$. Then

$$\mathbb{E} \|x_k - x\|_2 \leq \left(1 - \frac{1}{R}\right)^{k/2} \|x_0 - x\|_2 + \sqrt{R}\gamma,$$

where $\gamma = \max_i \frac{|z[i]|}{\|a_i\|_2}$.

- This bound is sharp and attained in simple examples.

Randomized Kaczmarz (RK) with noise

Theorem [N]

- Let $Ax = b$ be corrupted with noise: $Ax \approx b + z$. Then

$$\mathbb{E}\|x_k - x\|_2 \leq \left(1 - \frac{1}{R}\right)^{k/2} \|x_0 - x\|_2 + \sqrt{R}\gamma,$$

where $\gamma = \max_i \frac{|z[i]|}{\|a_i\|_2}$.

- This bound is sharp and attained in simple examples.

Even better convergence? : Noiseless case revisited

- Recall $x_{k+1} = x_k + \frac{b[i] - \langle a_i, x_k \rangle}{\|a_i\|_2^2} a_i$
- Since these projections are orthogonal, the optimal projection is one that maximizes $\|x_{k+1} - x_k\|_2$.
- Therefore we choose i maximizing $\left| \frac{b[i] - \langle a_i, x_k \rangle}{\|a_i\|_2} \right|$.
- Too costly \rightarrow Project onto low dimensional subspace.
- Use the low dimensional representations to predict the optimal projection.

Even better convergence? : Noiseless case revisited

- Recall $x_{k+1} = x_k + \frac{b[i] - \langle a_i, x_k \rangle}{\|a_i\|_2^2} a_i$
- Since these projections are orthogonal, the optimal projection is one that maximizes $\|x_{k+1} - x_k\|_2$.
- Therefore we choose i maximizing $\left| \frac{b[i] - \langle a_i, x_k \rangle}{\|a_i\|_2} \right|$.
- Too costly \rightarrow Project onto low dimensional subspace.
- Use the low dimensional representations to predict the optimal projection.

JL Dimension Reduction

Moreover

- In the proof of the JL Lemma the map Φ is chosen as the projection onto a random d -dimensional subspace of \mathbb{R}^n . Now many known distributions will yield such a projection.
- Recently, transforms with fast multiplies have also been shown to satisfy the JL Lemma [Ailon-Chazelle, Hinrichs-Vybiral, Ailon-Liberty, Krahmer-Ward, ...]

Perform Reduction

Choose such a $d \times n$ projector Φ and during preprocessing set $\alpha_j = \Phi a_j$.

Runtime

- Select:
- Calculate Φx_k : In general $O(nd)$
 - Calculate γ_i for each i (of n): $O(nd)$

Test: Calculate γ_j^* and γ_l^* : $O(n)$

Project: Calculate x_{k+1} : $O(n)$

Overall Runtime

Since each iteration takes $O(nd)$, we have convergence in $O(n^2d)$.

