



Randomized Kaczmarz Algorithms: Exact MSE Analysis and Optimal Sampling Probabilities

Yue M. Lu

Signals, Information, and Networks Group (SING)
Harvard University

Joint work with Ameya Agaskar and Chuang Wang

April 17, 2015

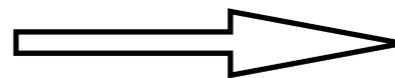
Support: The Department of Air Force under Contract #FA8721T05TCT0002
and NSF CCF-1319140

$$\begin{array}{c} \mathbf{y} \\ y_1 \\ y_2 \\ \cdot \\ \cdot \\ y_m \\ \text{observation} \end{array} = \begin{array}{c} \mathbf{A} \\ -a_1^T - \\ -a_2^T - \\ \cdot \\ \cdot \\ -a_m^T - \\ \text{system matrix} \end{array} \times \begin{array}{c} \mathbf{x} \\ x_1 \\ \cdot \\ \cdot \\ x_n \end{array} + \begin{array}{c} z_1 \\ z_2 \\ \cdot \\ \cdot \\ z_m \\ \text{noise} \end{array}$$

$$\begin{bmatrix} y_1 \\ y_2 \\ \cdot \\ \cdot \\ y_m \end{bmatrix} = \begin{bmatrix} -a_1^T - \\ -a_2^T - \\ \cdot \\ \cdot \\ -a_m^T - \end{bmatrix} \times \begin{bmatrix} x_1 \\ \cdot \\ \cdot \\ \cdot \\ x_n \end{bmatrix} + \begin{bmatrix} z_1 \\ z_2 \\ \cdot \\ \cdot \\ z_m \end{bmatrix}$$

New considerations:

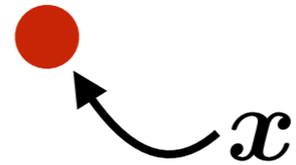
- Massive sizes
- Streaming data
- Distributed storage
- Parallel computing platform



Iterative greedy methods

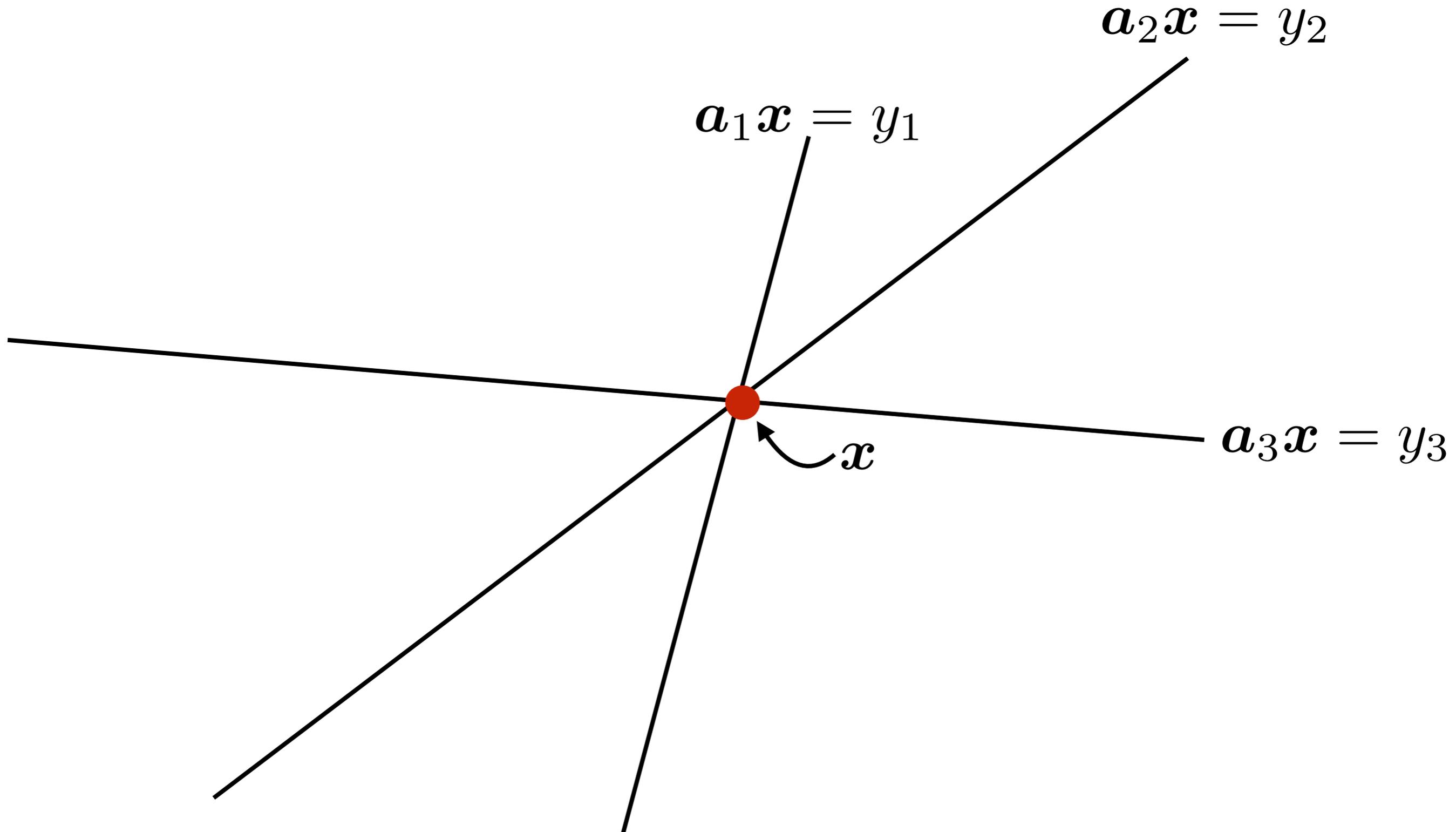
Kaczmarz Algorithm — Iterative Projections

- Noiseless case: $y = Ax$ encodes a system of m equations.



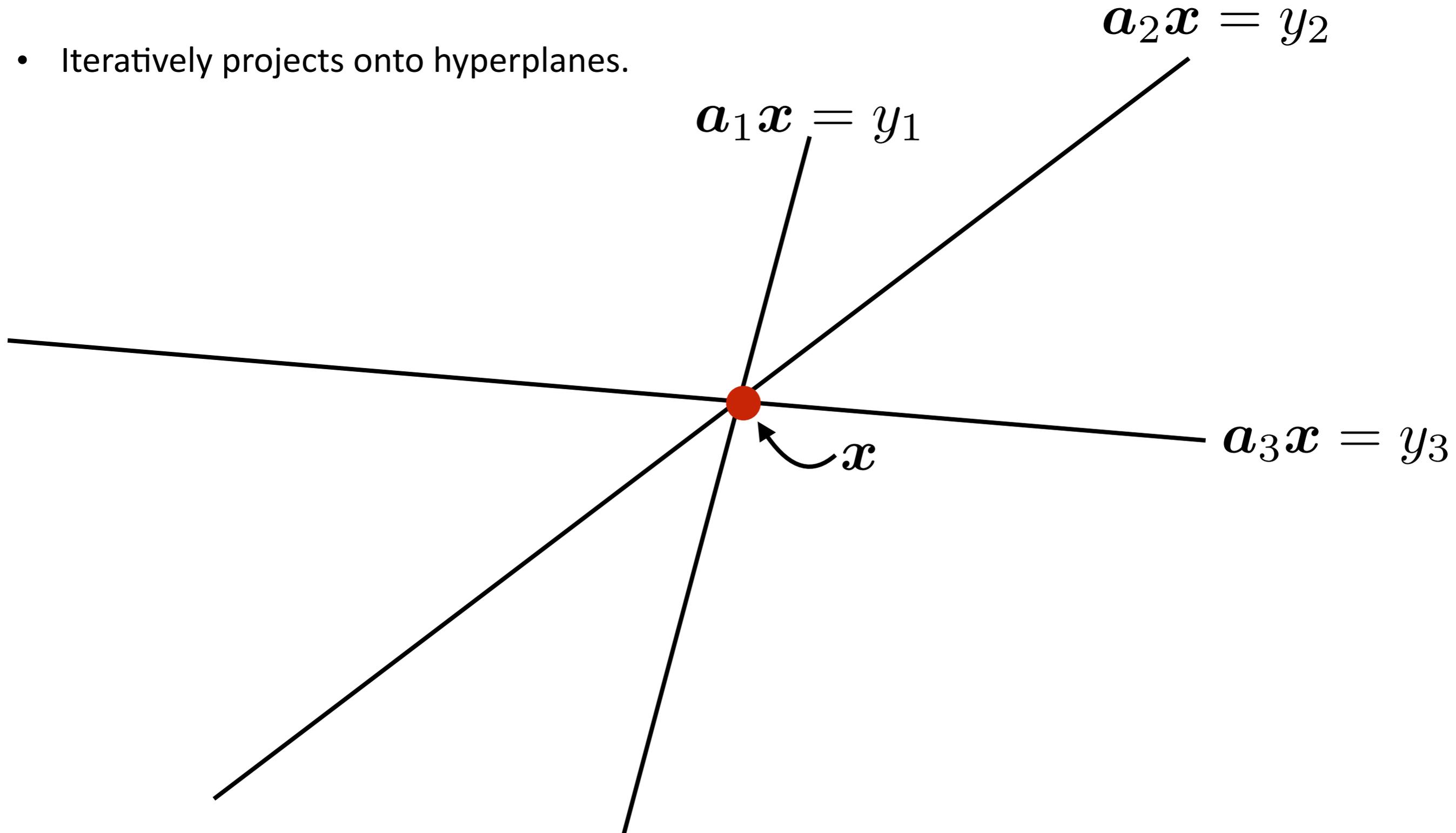
Kaczmarz Algorithm — Iterative Projections

- Noiseless case: $y = Ax$ encodes a system of m equations.



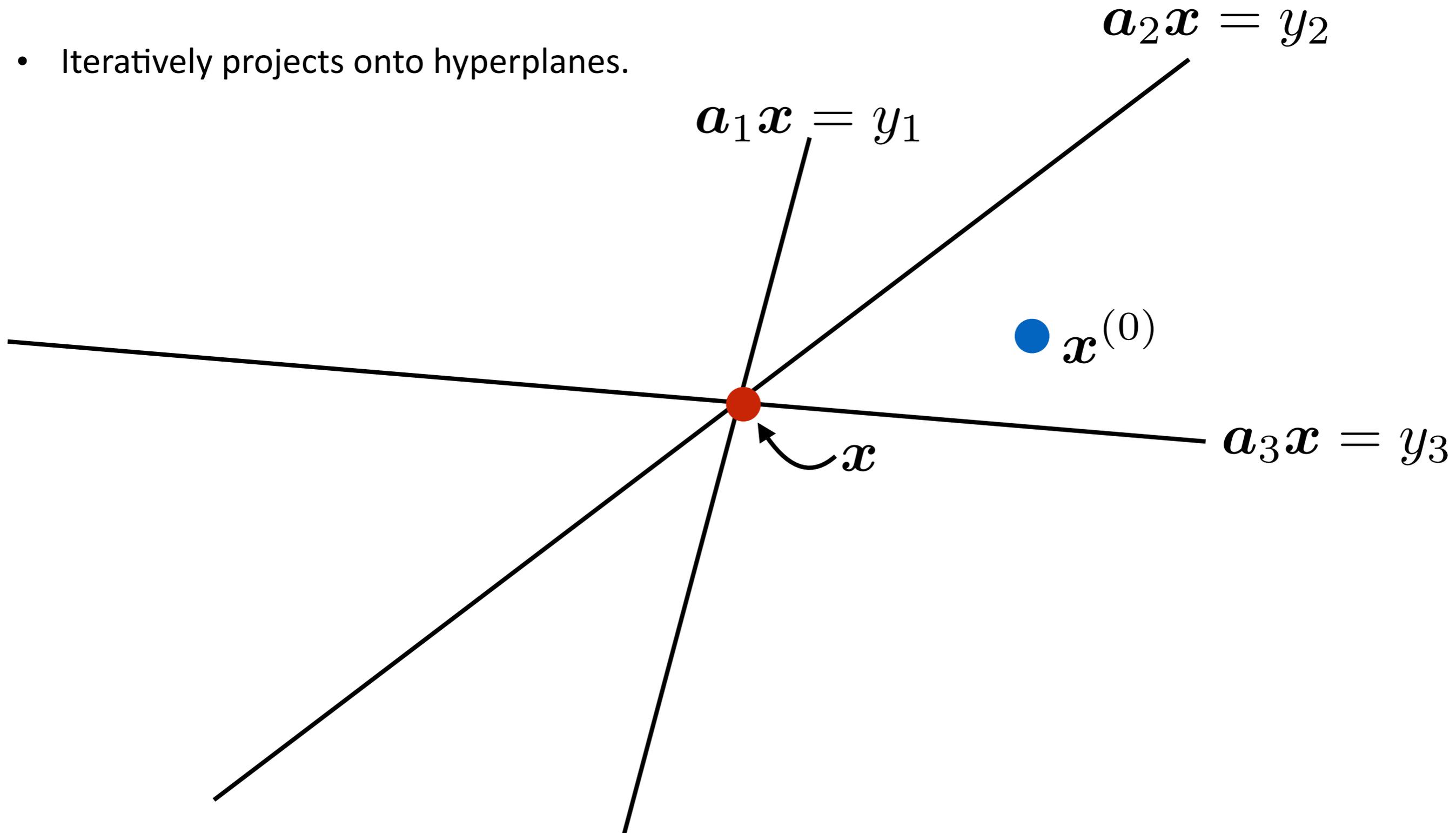
Kaczmarz Algorithm — Iterative Projections

- Noiseless case: $y = Ax$ encodes a system of m equations.
- Iteratively projects onto hyperplanes.



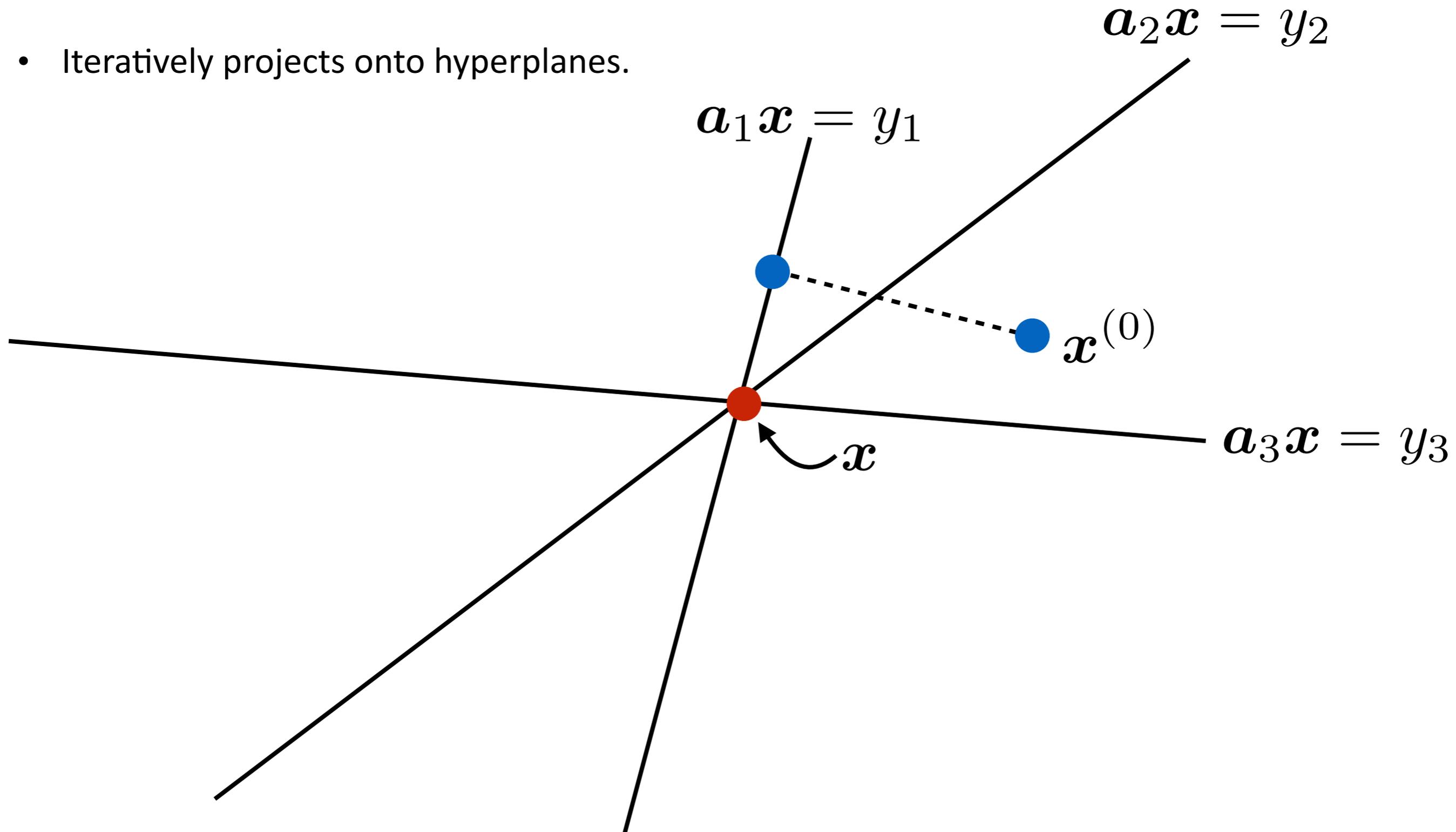
Kaczmarz Algorithm — Iterative Projections

- Noiseless case: $\mathbf{y} = \mathbf{A}\mathbf{x}$ encodes a system of m equations.
- Iteratively projects onto hyperplanes.



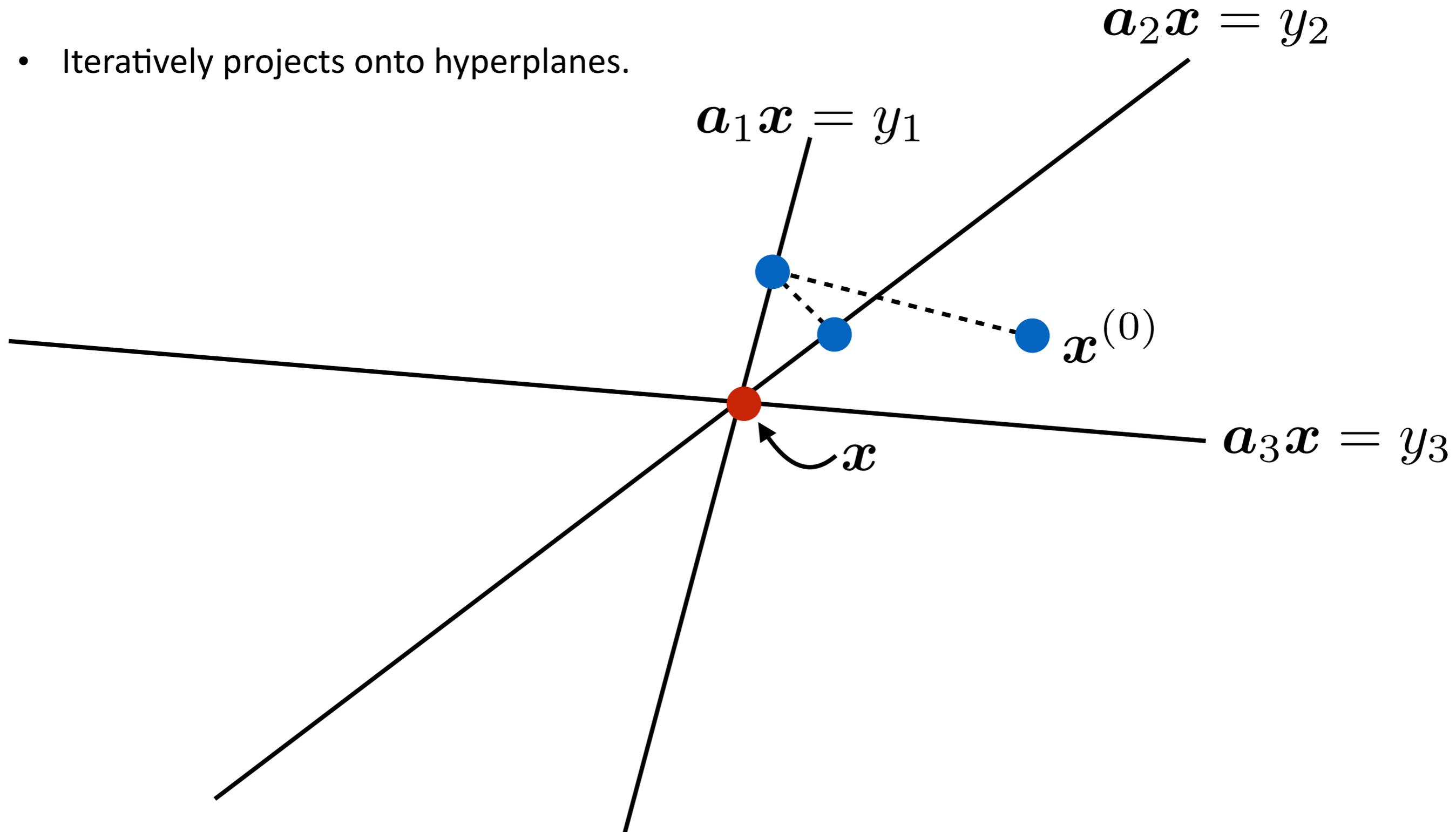
Kaczmarz Algorithm — Iterative Projections

- Noiseless case: $y = Ax$ encodes a system of m equations.
- Iteratively projects onto hyperplanes.



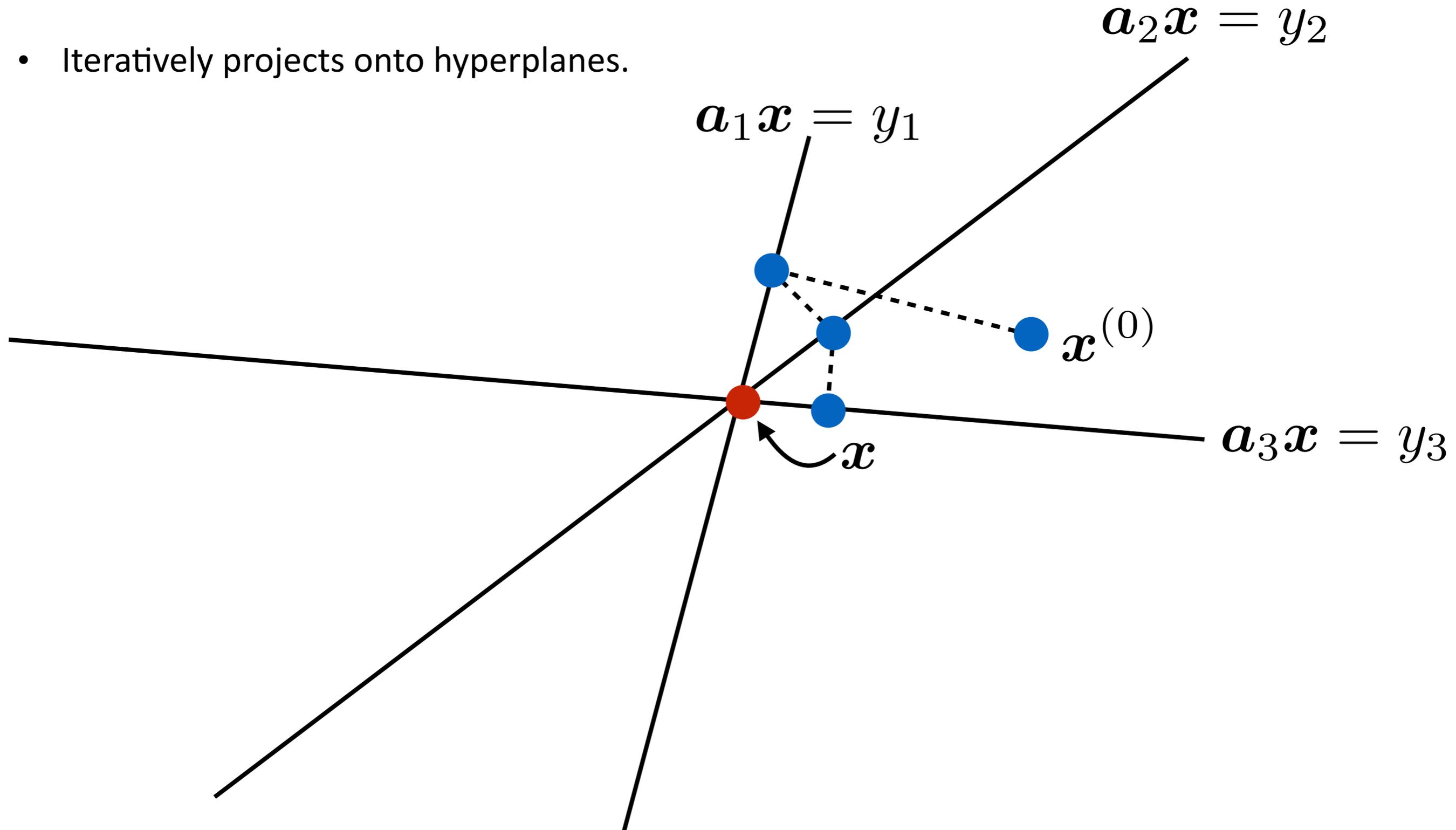
Kaczmarz Algorithm — Iterative Projections

- Noiseless case: $y = Ax$ encodes a system of m equations.
- Iteratively projects onto hyperplanes.



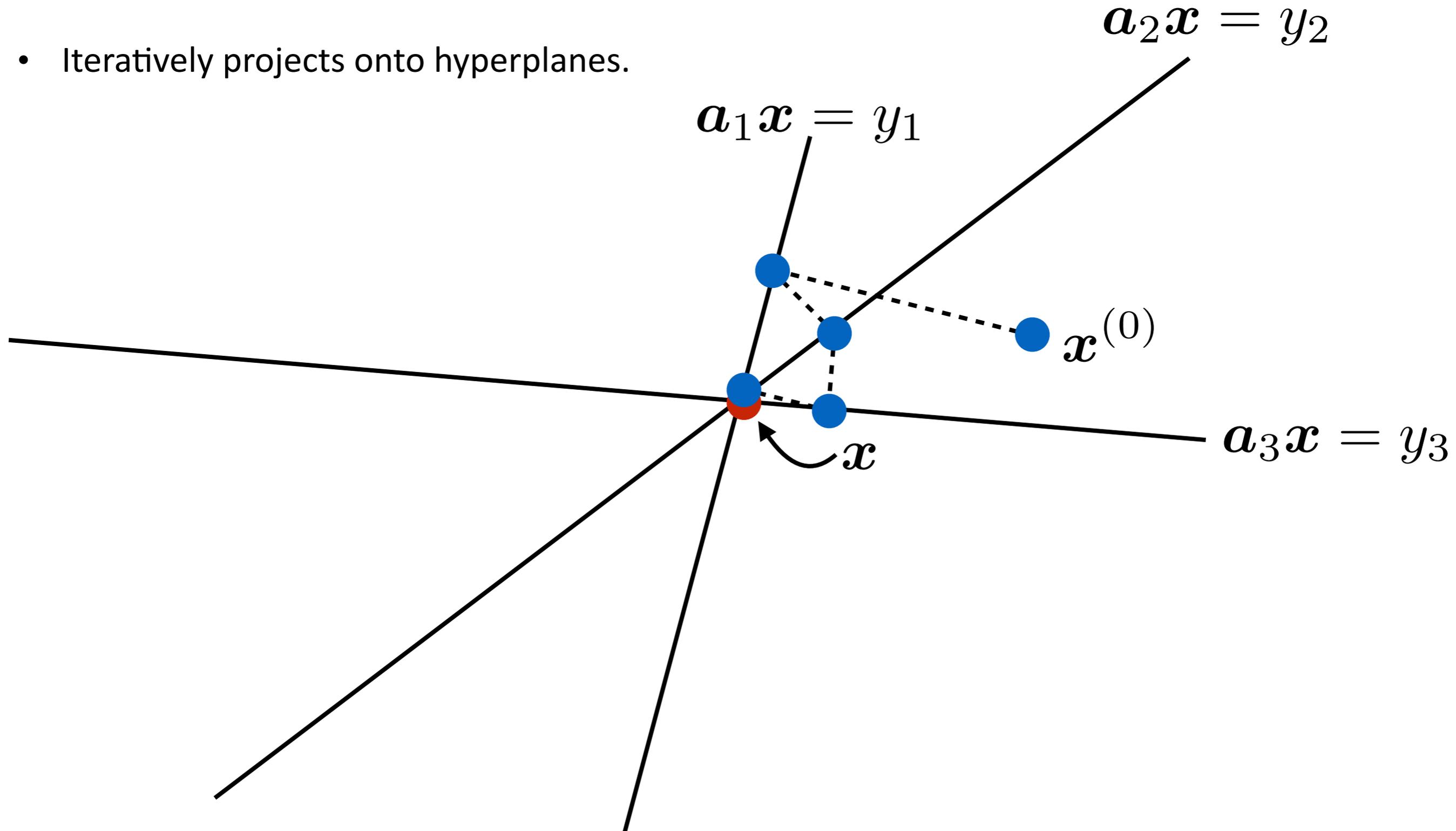
Kaczmarz Algorithm — Iterative Projections

- Noiseless case: $y = Ax$ encodes a system of m equations.
- Iteratively projects onto hyperplanes.



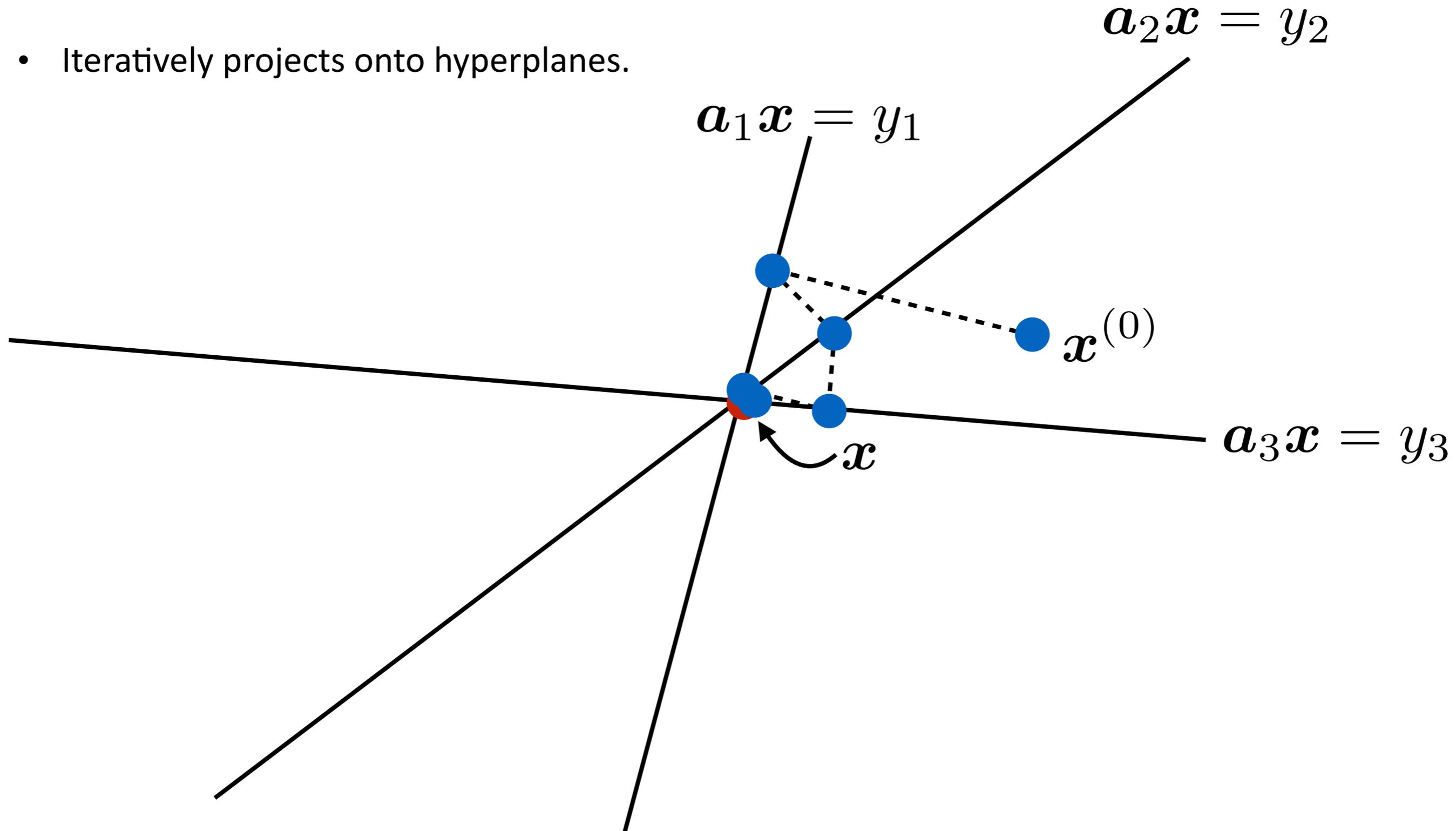
Kaczmarz Algorithm — Iterative Projections

- Noiseless case: $y = Ax$ encodes a system of m equations.
- Iteratively projects onto hyperplanes.



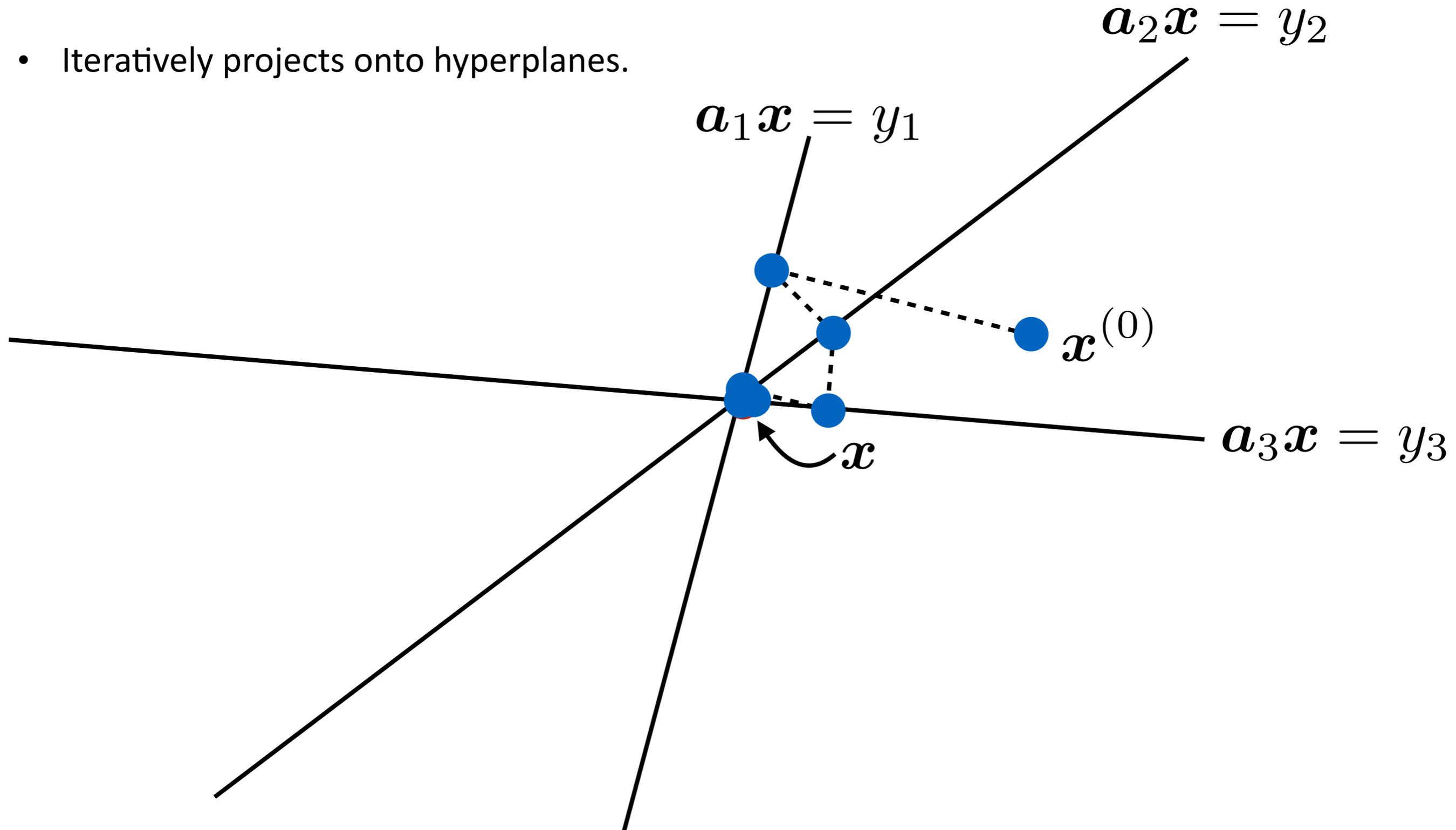
Kaczmarz Algorithm — Iterative Projections

- Noiseless case: $y = Ax$ encodes a system of m equations.
- Iteratively projects onto hyperplanes.



Kaczmarz Algorithm — Iterative Projections

- Noiseless case: $y = Ax$ encodes a system of m equations.
- Iteratively projects onto hyperplanes.



- Iterative algorithm introduced by S. Kaczmarz (1937)
- Also known as algebraic reconstruction technique (ART)
- Special case of projection onto convex sets (POCS)

Pseudocode

Initialize arbitrary $\mathbf{x}^{(0)}$

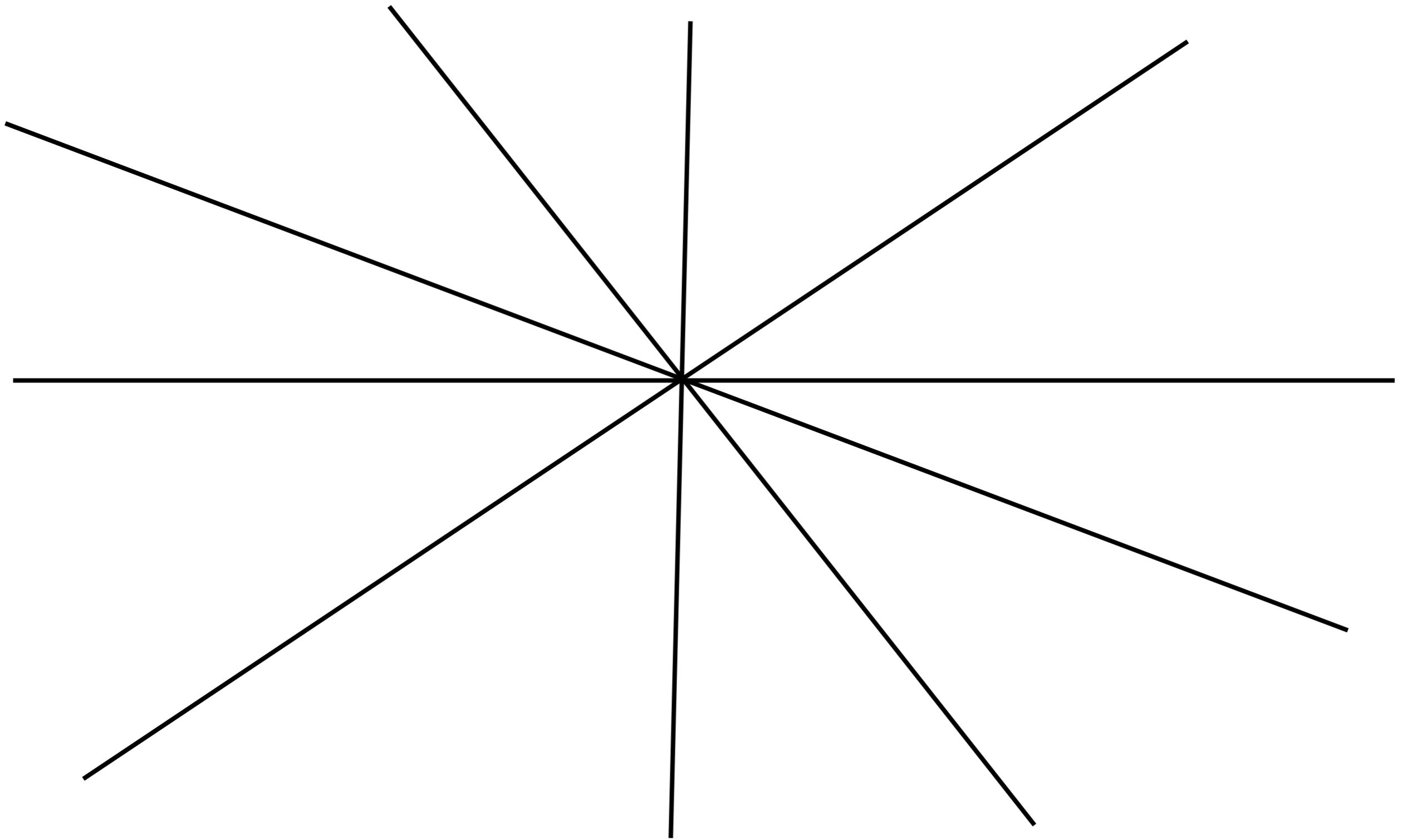
For $k = 1$ to N_{iter} :

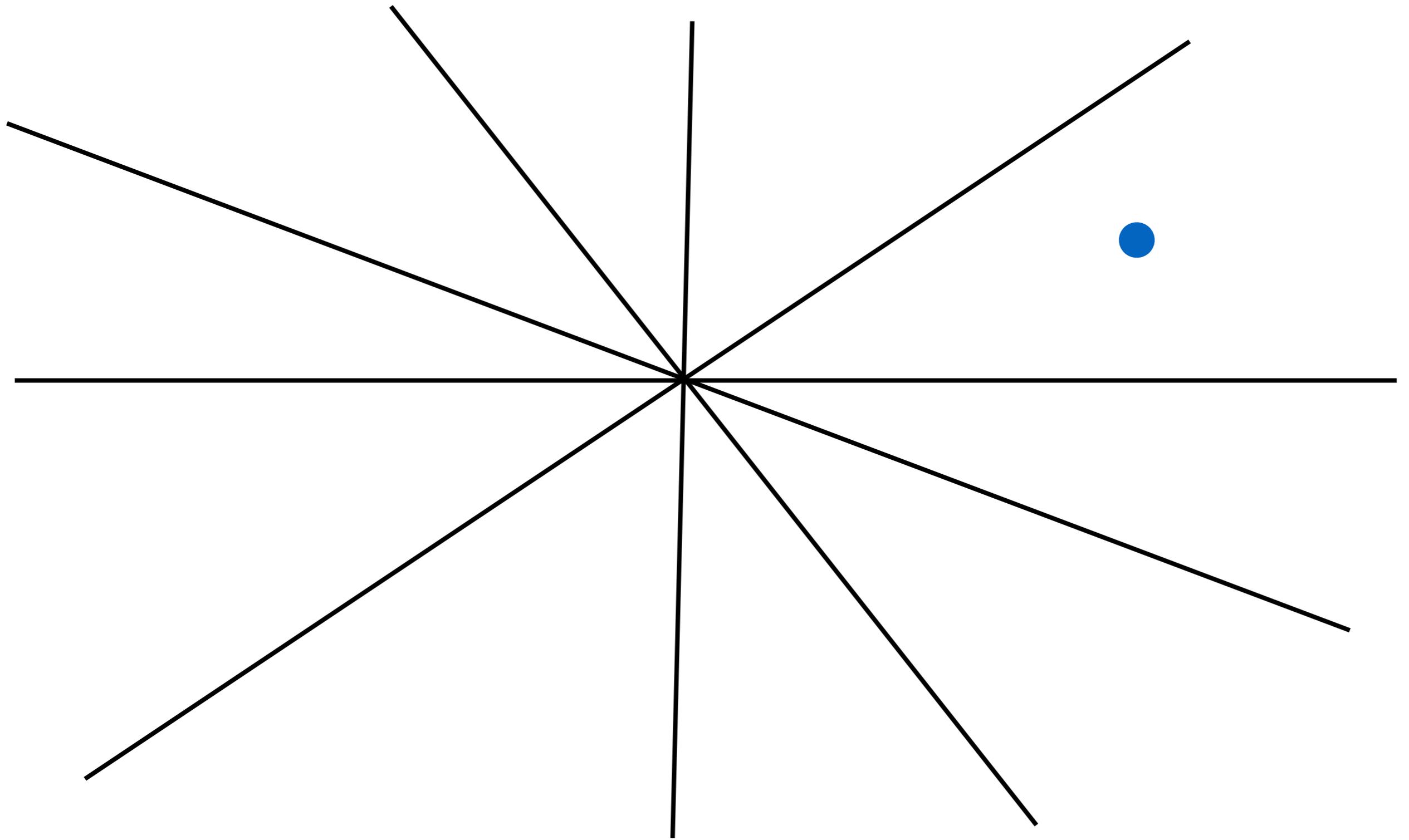
$r \leftarrow (k \bmod m) + 1$ *select the next row*

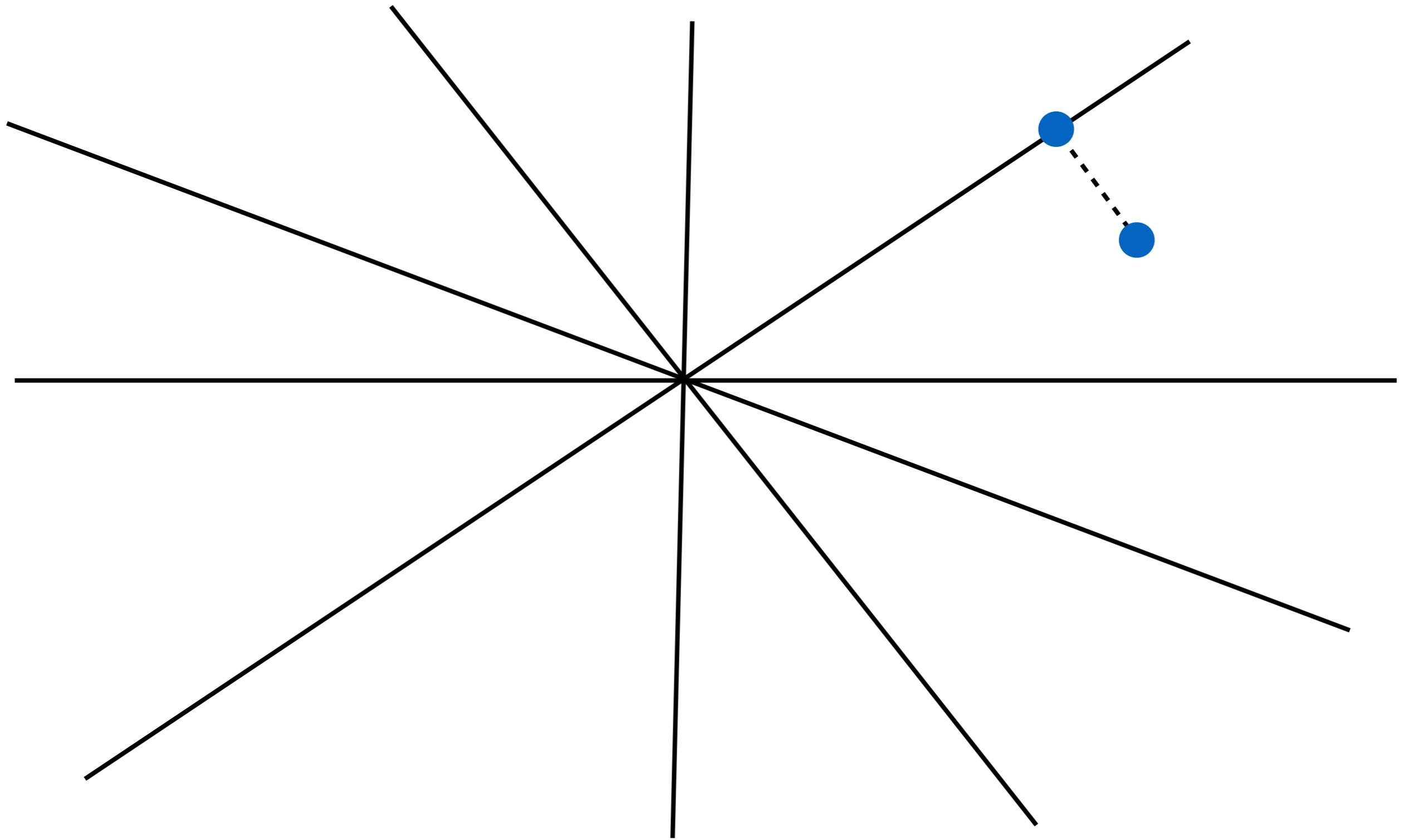
$\mathbf{x}^{(k)} \leftarrow \mathbf{x}^{(k-1)} + \frac{y_r - \mathbf{a}_r^T \mathbf{x}^{(k-1)}}{\|\mathbf{a}_r\|^2} \mathbf{a}_r$ *projection*

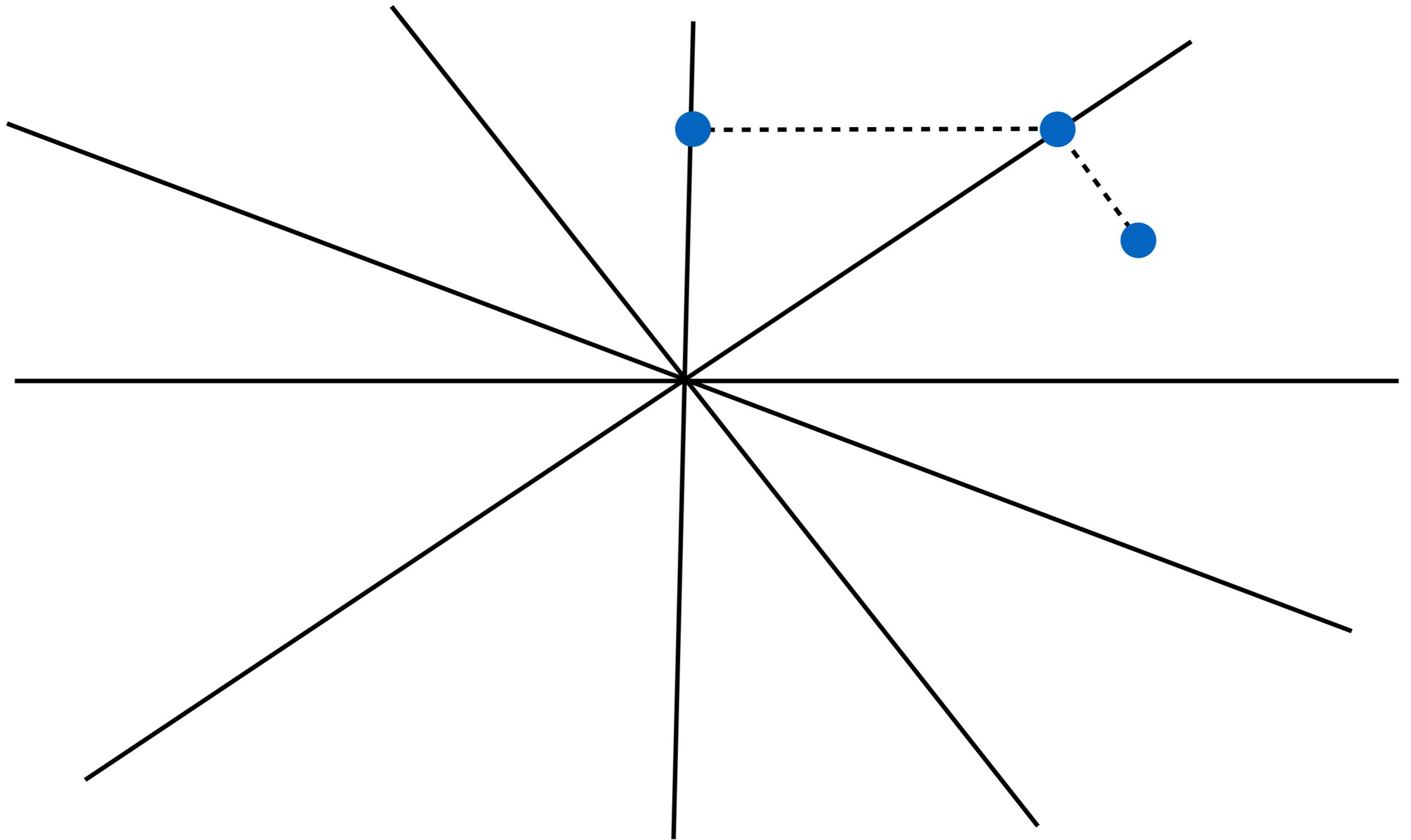
$\hat{\mathbf{x}} \leftarrow \mathbf{x}^{(N_{\text{iter}})}$

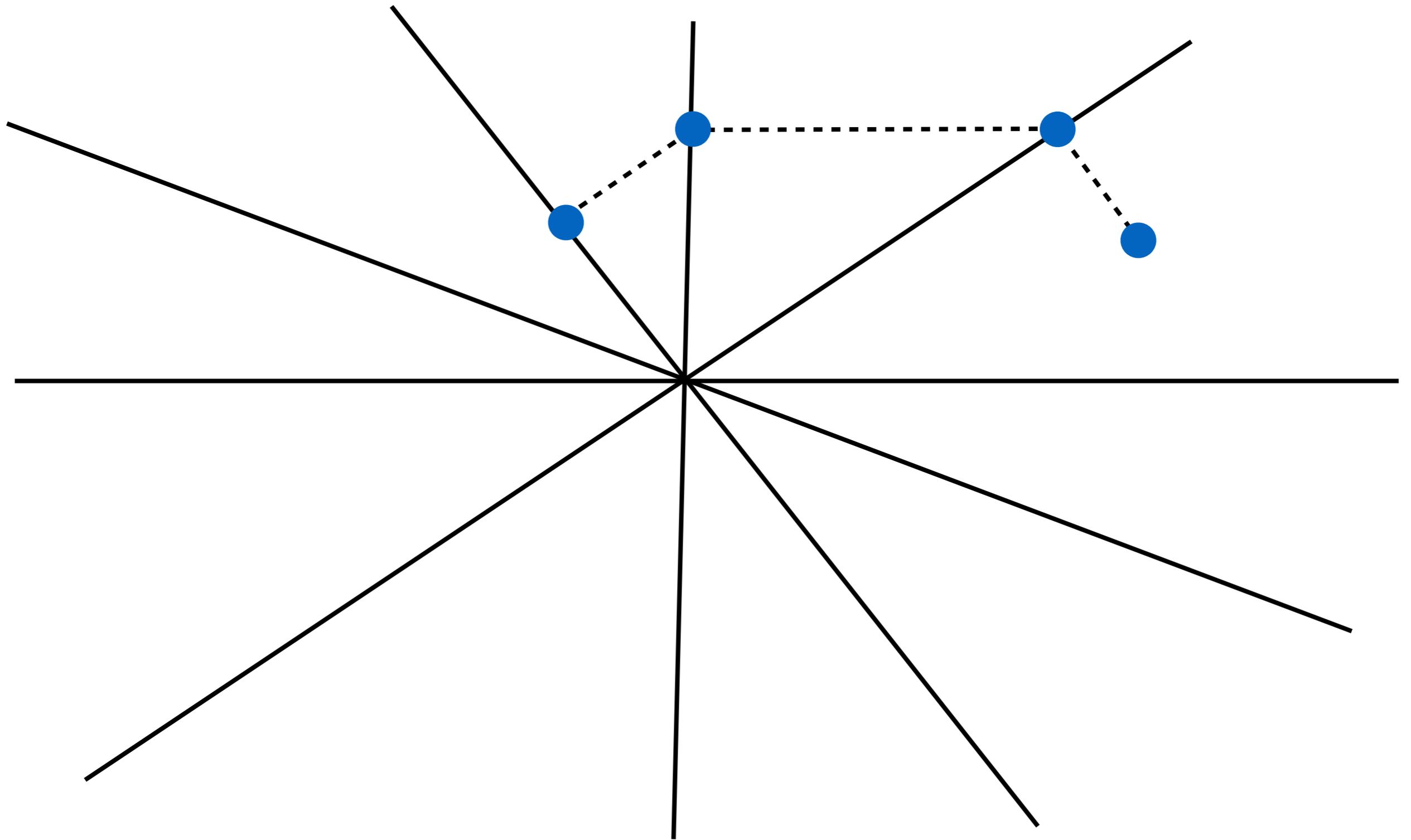
- Can be extended to find least squares estimate from noisy measurements (Zouzias & Freris 2013)

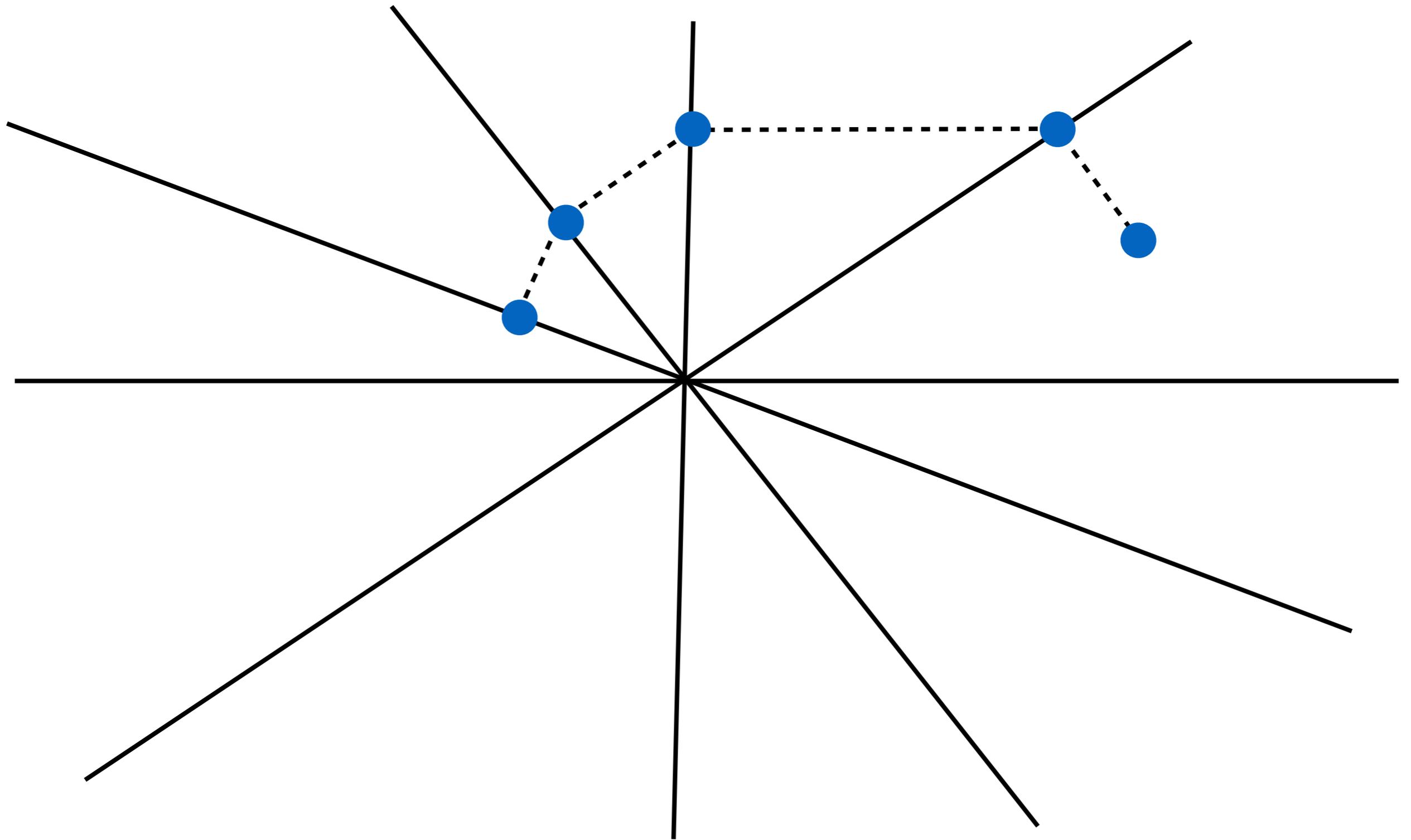


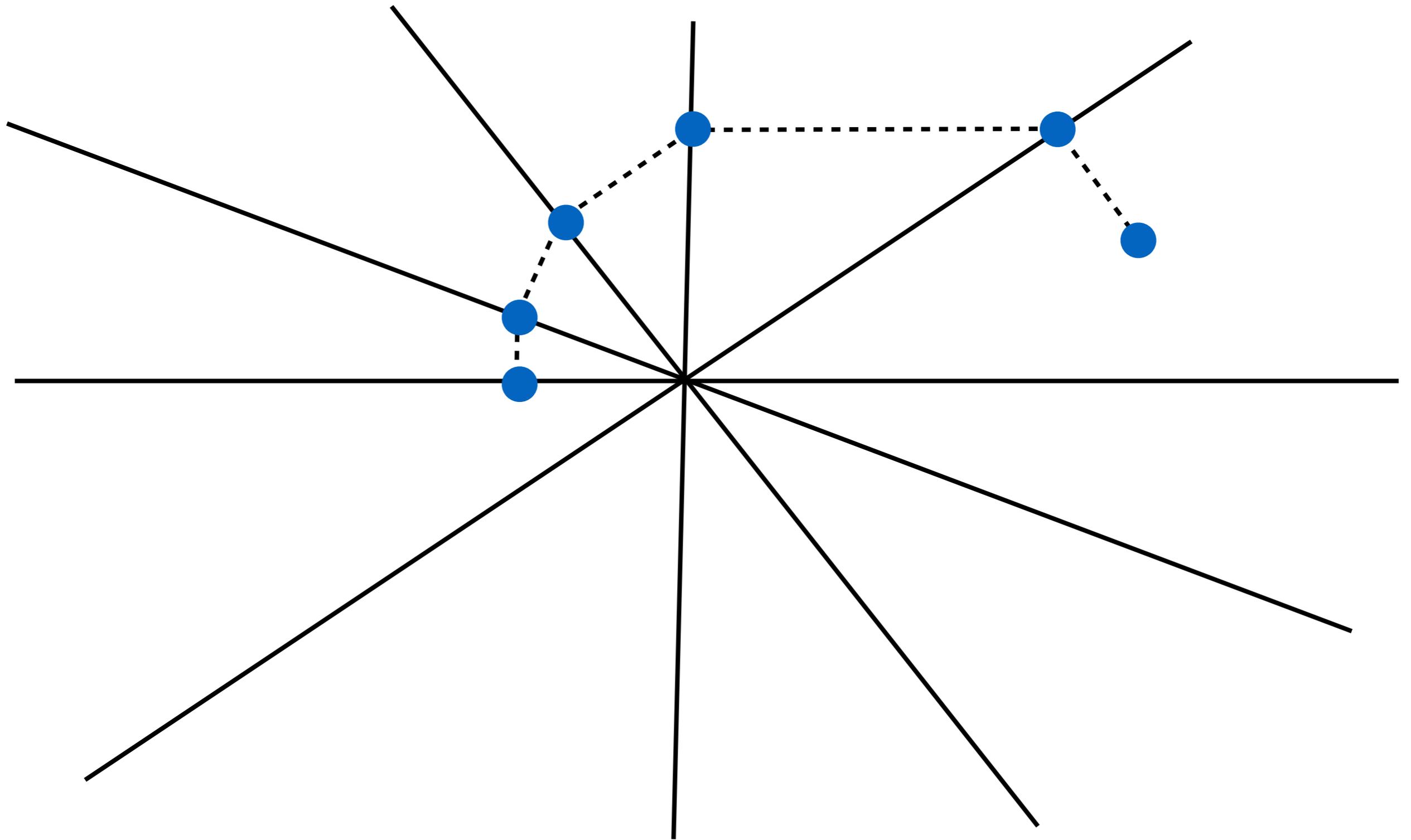


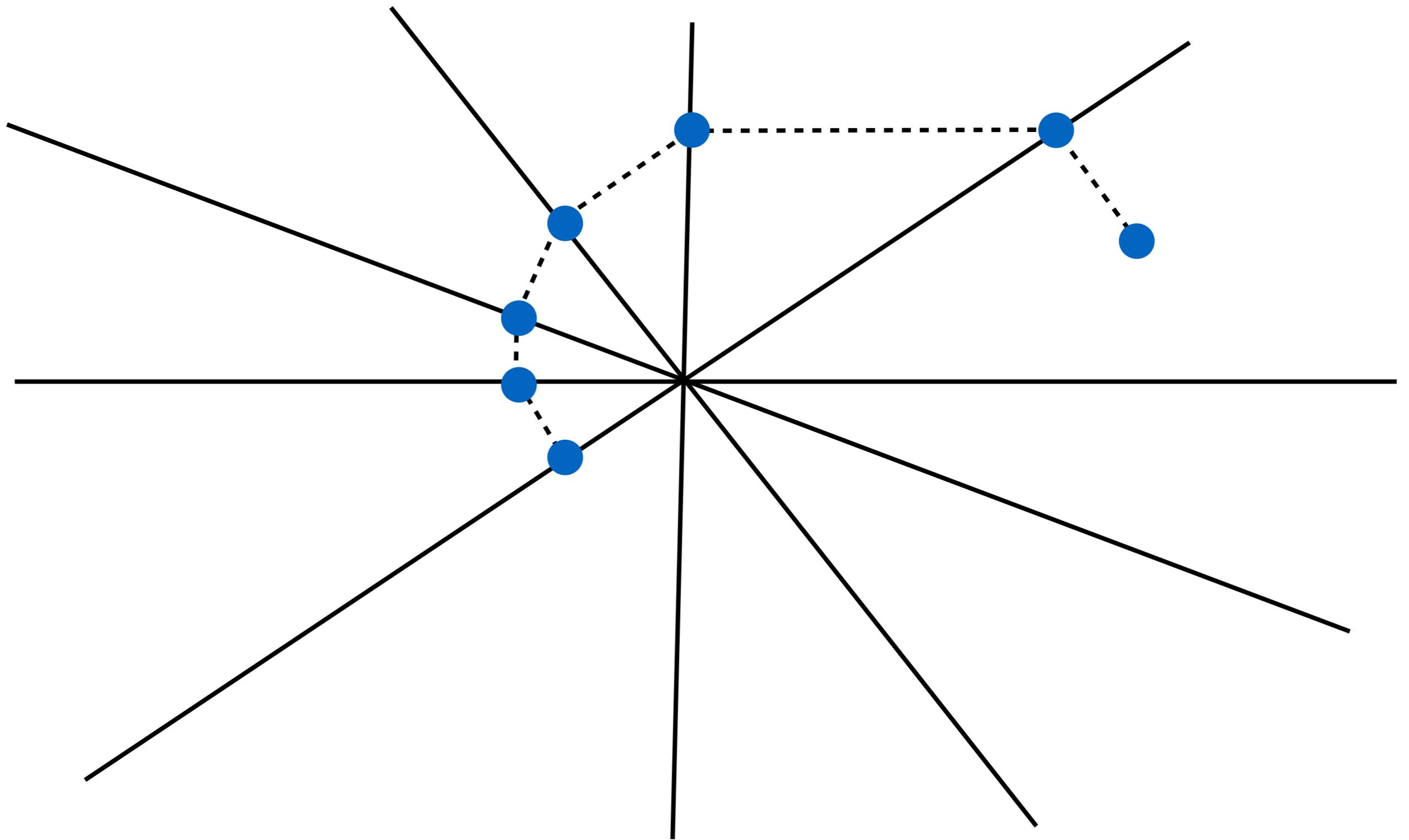


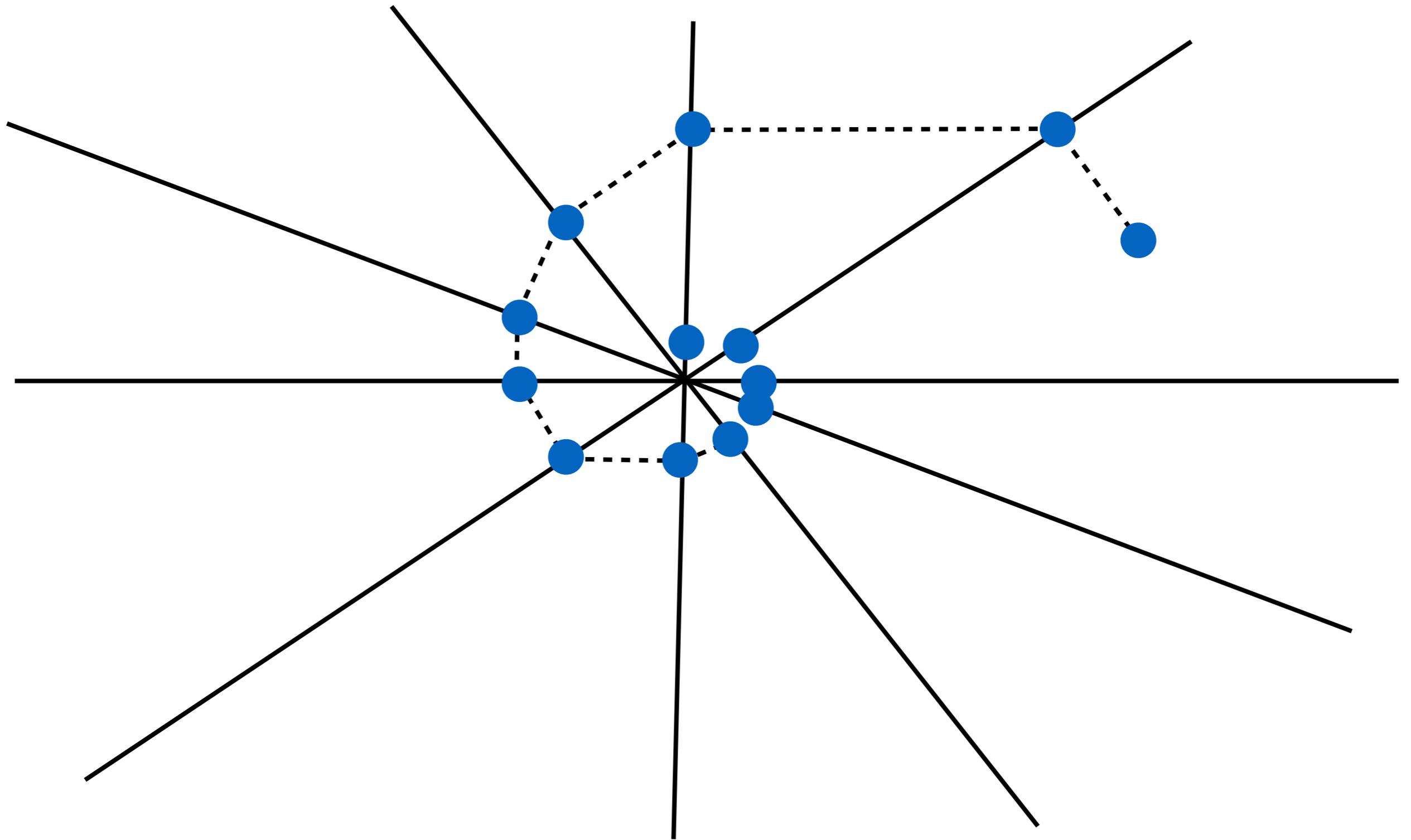


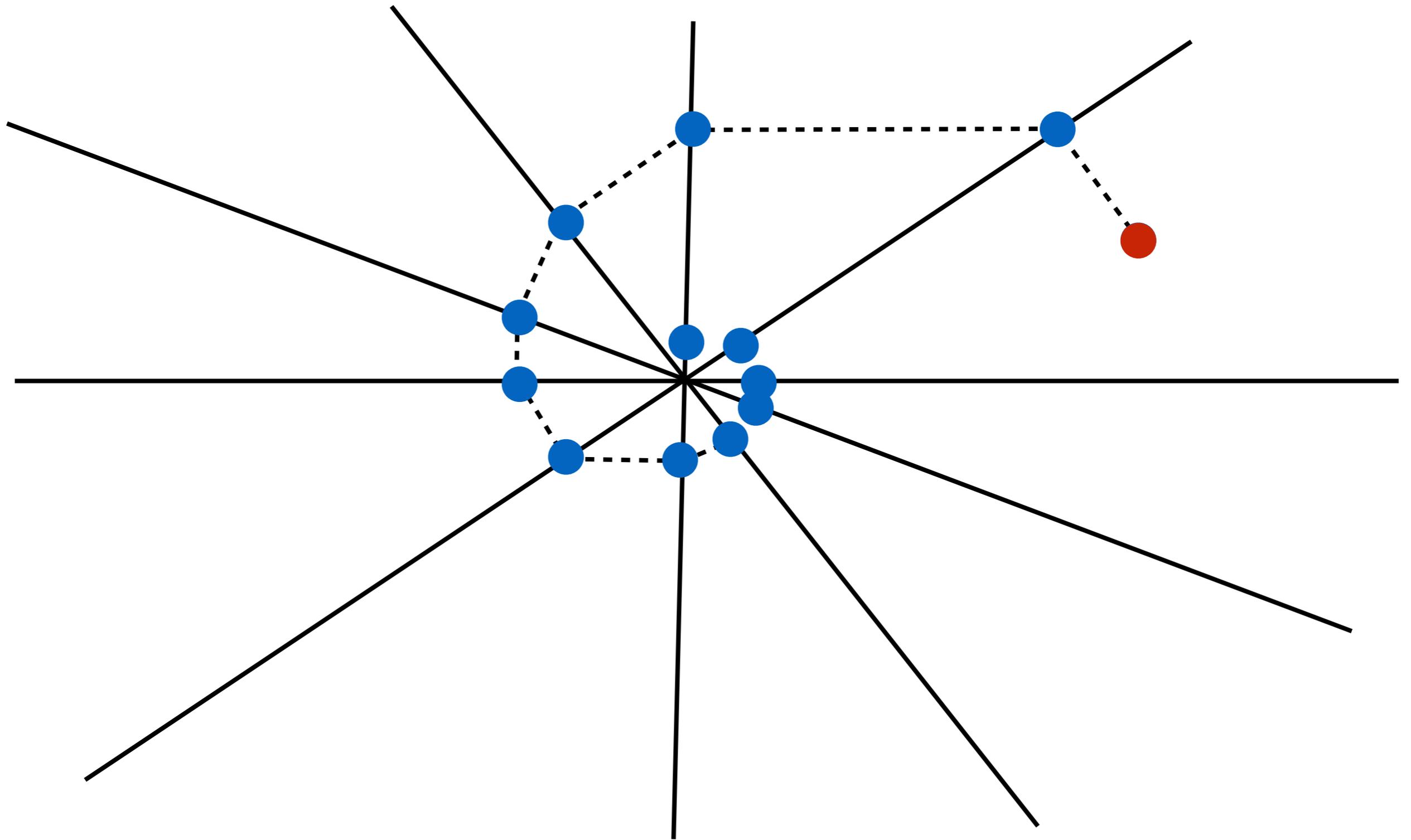


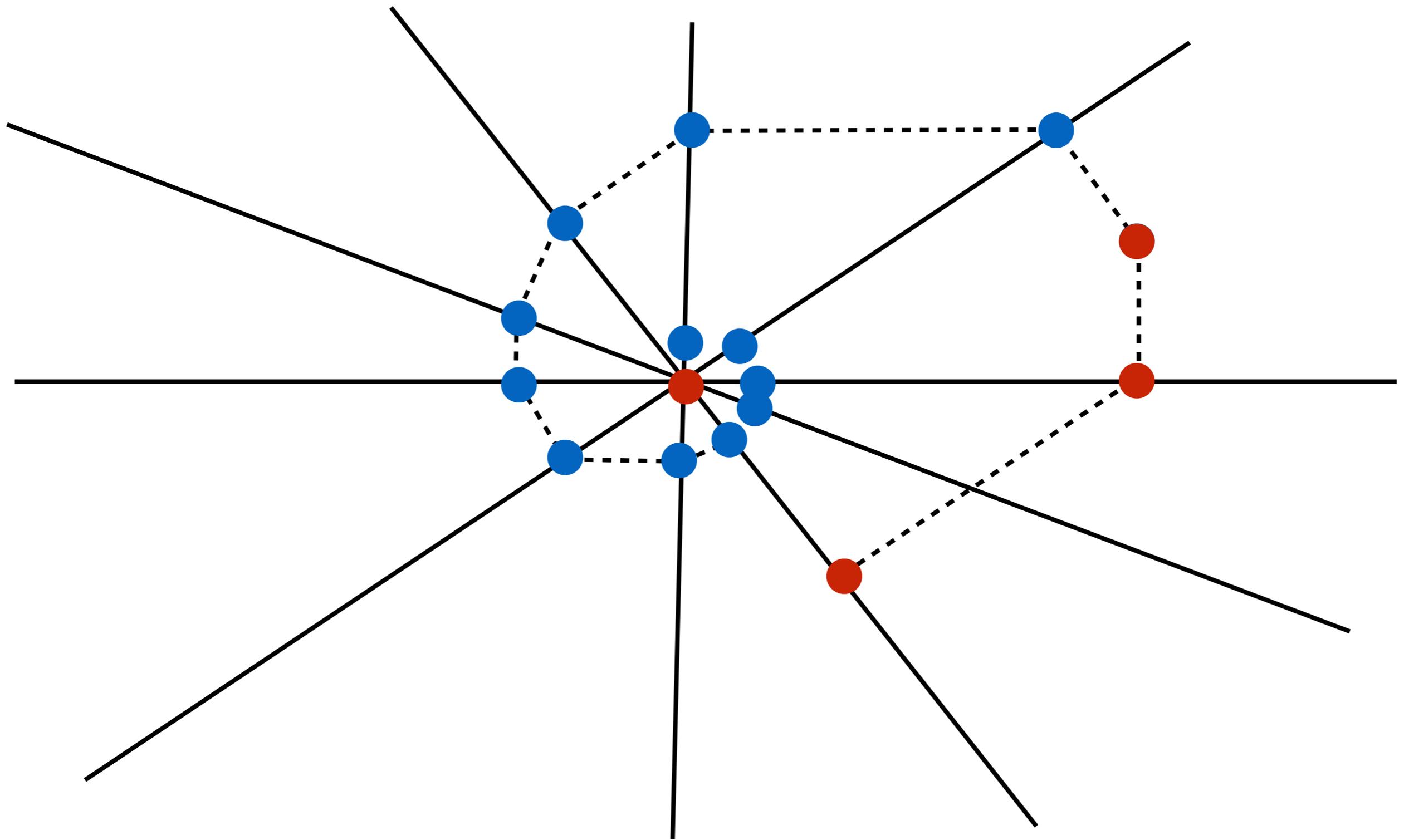












IF YOU USE THE WRONG ROW ORDER



YOU'RE GONNA HAVE A BAD TIME

- Strohmer and Vershynin (2009) proposed randomizing the order:

Choose row \mathbf{a}_i with probability proportional to $\|\mathbf{a}_i\|^2$.

- Guarantees exponential convergence:

$$\frac{\mathbb{E}\|\mathbf{x}^{(N)} - \mathbf{x}\|^2}{\|\mathbf{x}^{(0)} - \mathbf{x}\|^2} \leq \left(1 - \underbrace{\kappa(A)^{-2}}_{\|A\|_F \|A^{-1}\|_2}\right)^N$$

- Works for arbitrary probabilities by preconditioning, so we assume row i chosen with probability p_i .

- Strohmer and Vershynin (2009) proposed randomizing the order:

Choose row \mathbf{a}_i with probability proportional to $\|\mathbf{a}_i\|^2$.

- Guarantees exponential convergence:

$$\frac{\mathbb{E}\|\mathbf{x}^{(N)} - \mathbf{x}\|^2}{\|\mathbf{x}^{(0)} - \mathbf{x}\|^2} \leq \left(1 - \underbrace{\kappa(A)^{-2}}_{\|A\|_F \|A^{-1}\|_2}\right)^N$$

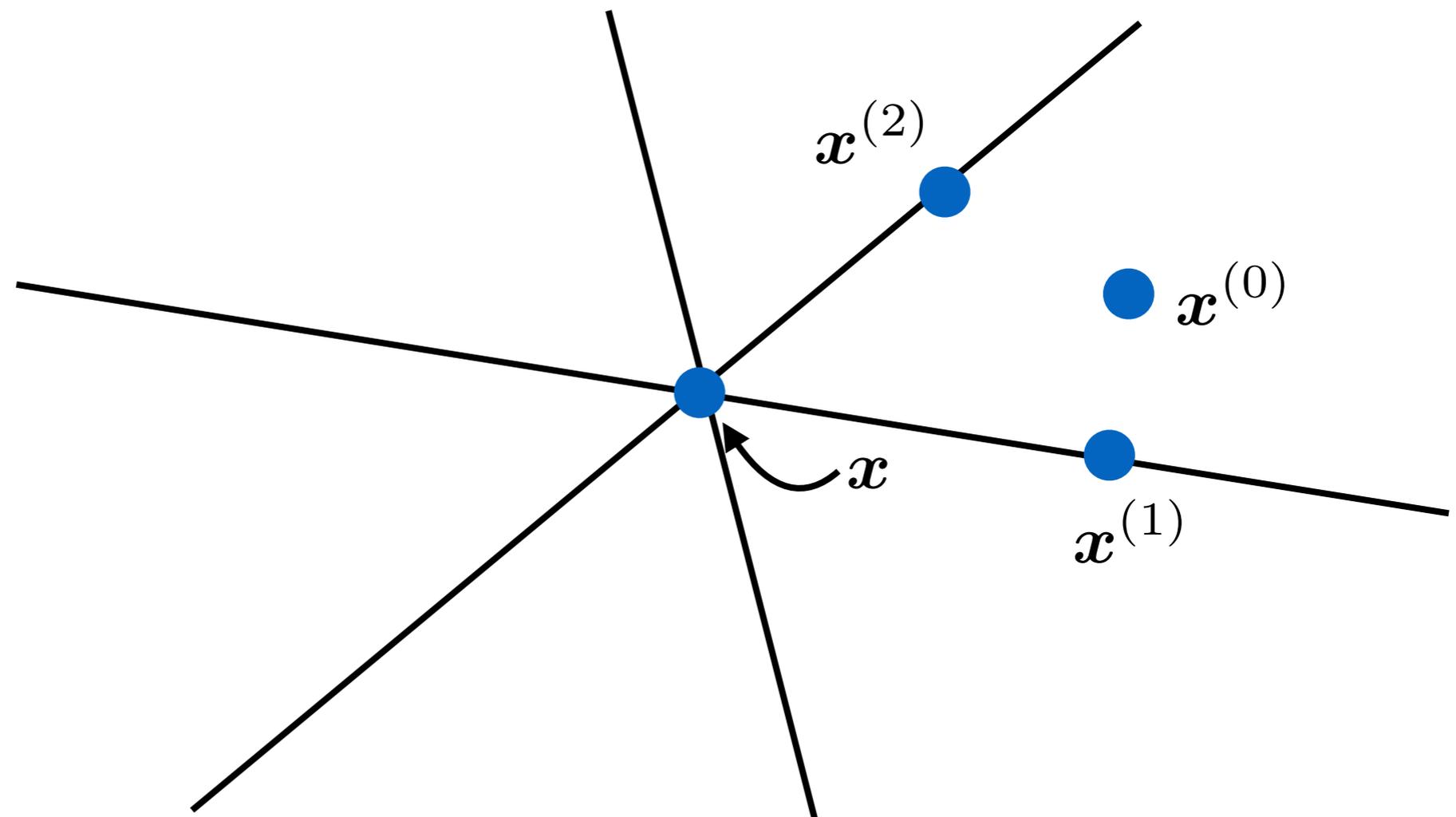
- Works for arbitrary probabilities by preconditioning, so we assume row i chosen with probability p_i .

Related work

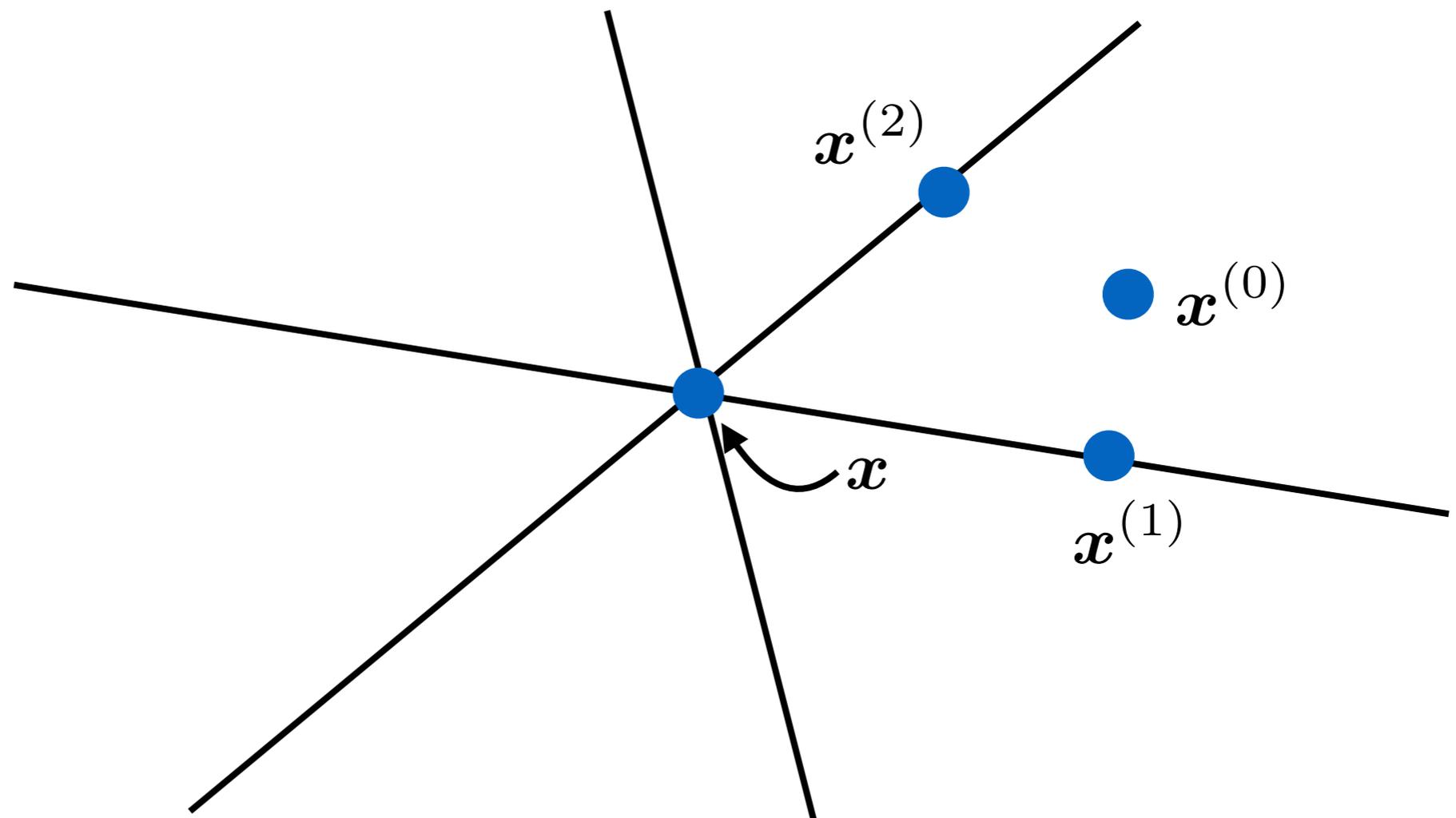
- Error bounds for ***inconsistent systems*** (Needell 2012)
- Almost-sure*** convergence (Chen & Powell 2012)
- Extension to find ***least-square solution*** in noise (Zouzias & Freris 2013)
- Block*** Kaczmarz (Needell & Tropp 2014)

1. Exact MSE formula and decay rate
2. Optimization of row selection probabilities
3. “Quenched error exponent”

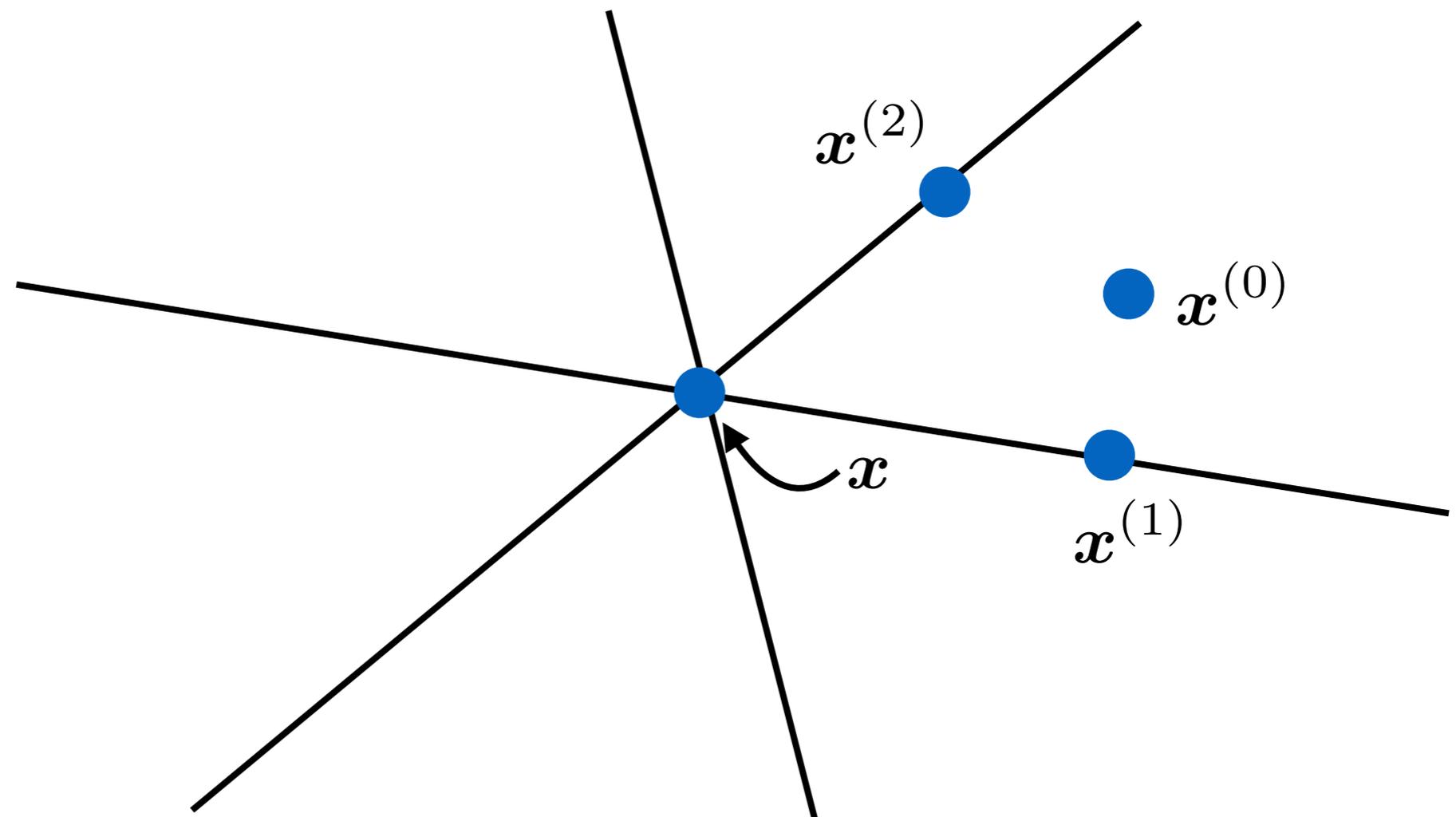
$$\mathbf{x}^{(k)} = \mathbf{x}^{(k-1)} + \frac{y_r - \mathbf{a}_r^T \mathbf{x}^{(k-1)}}{\|\mathbf{a}_r\|^2} \mathbf{a}_r$$



$$\mathbf{x}^{(k)} - \mathbf{x} = \mathbf{x}^{(k-1)} - \mathbf{x} + \frac{y_r - \mathbf{a}_r^T \mathbf{x}^{(k-1)}}{\|\mathbf{a}_r\|^2} \mathbf{a}_r$$

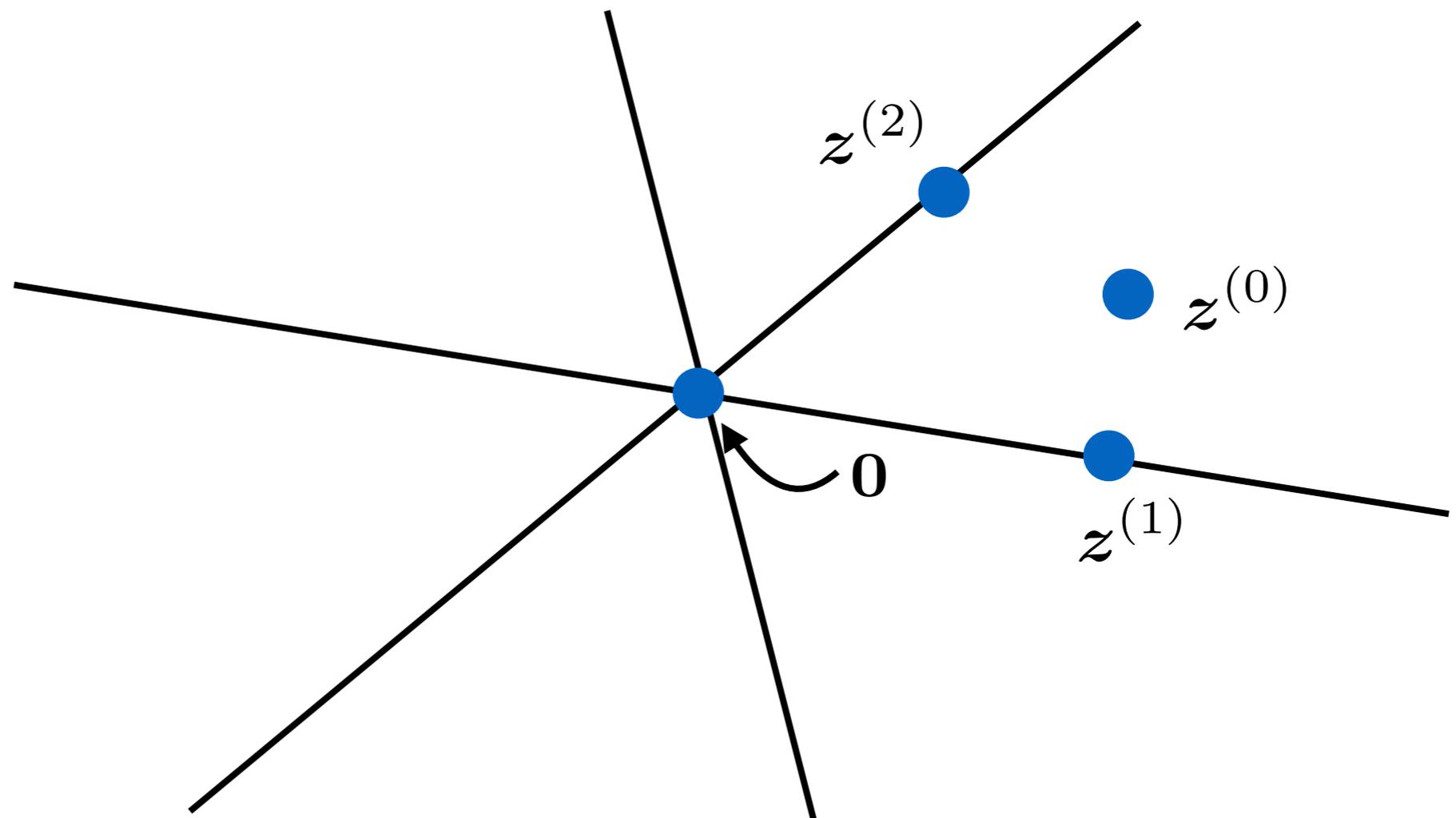


$$(\mathbf{x}^{(k)} - \mathbf{x}) = (\mathbf{x}^{(k-1)} - \mathbf{x}) - \frac{\mathbf{a}_r^T (\mathbf{x}^{(k-1)} - \mathbf{x})}{\|\mathbf{a}_r\|^2} \mathbf{a}_r$$



$$(\mathbf{x}^{(k)} - \mathbf{x}) = (\mathbf{x}^{(k-1)} - \mathbf{x}) - \frac{\mathbf{a}_r^T (\mathbf{x}^{(k-1)} - \mathbf{x})}{\|\mathbf{a}_r\|^2} \mathbf{a}_r$$

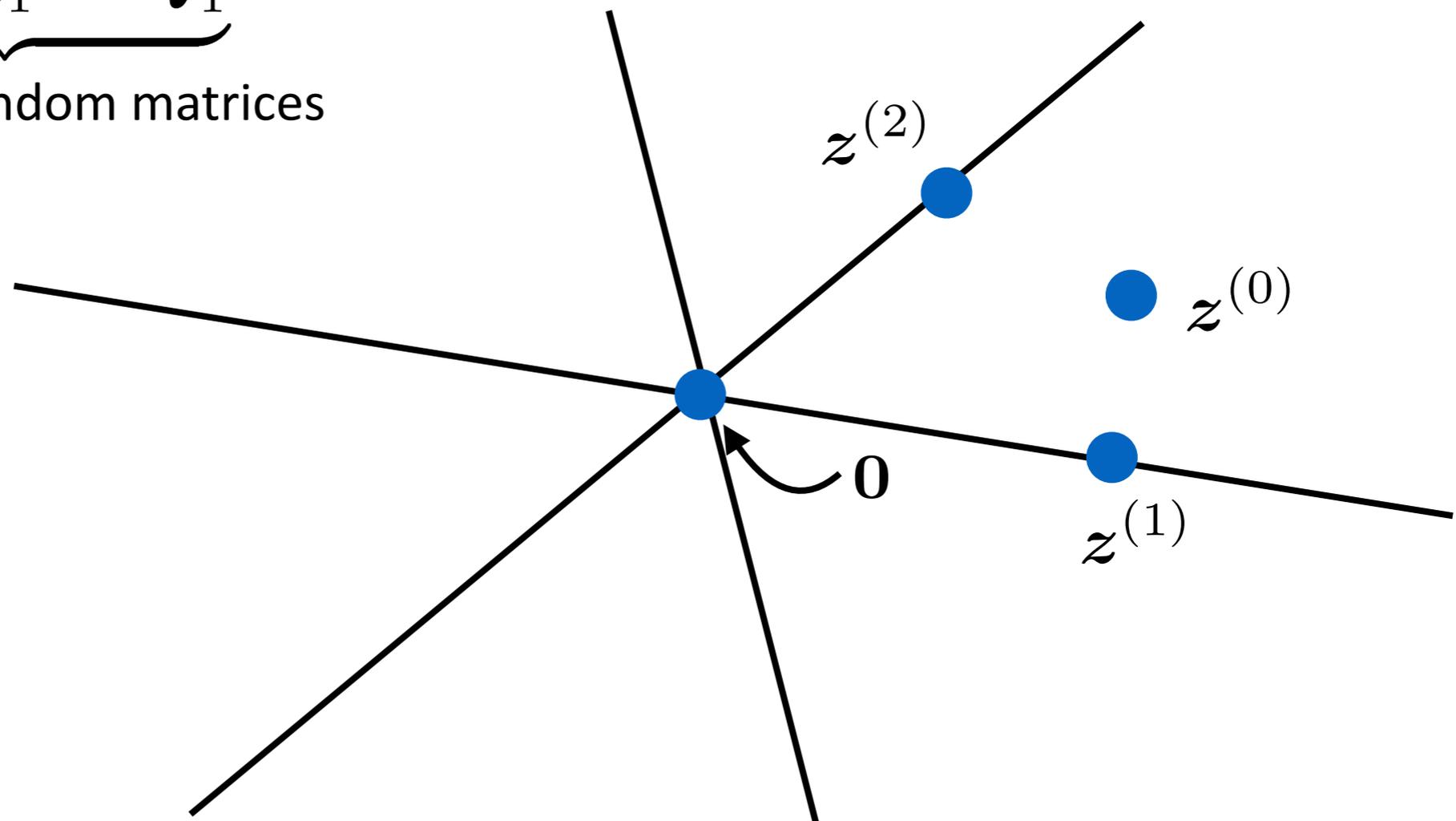
$$\mathbf{z}^{(k)} = Q_k \mathbf{z}^{(k-1)} \quad \text{where} \quad Q_k = \left(I - \frac{\mathbf{a}_r \mathbf{a}_r^T}{\|\mathbf{a}_r\|^2} \right)$$



$$(\mathbf{x}^{(k)} - \mathbf{x}) = (\mathbf{x}^{(k-1)} - \mathbf{x}) - \frac{\mathbf{a}_r^T (\mathbf{x}^{(k-1)} - \mathbf{x})}{\|\mathbf{a}_r\|^2} \mathbf{a}_r$$

$$\mathbf{z}^{(k)} = Q_k \mathbf{z}^{(k-1)} \quad \text{where} \quad Q_k = \left(I - \frac{\mathbf{a}_r \mathbf{a}_r^T}{\|\mathbf{a}_r\|^2} \right)$$

$$\mathbf{z}^{(k)} = \underbrace{Q_k Q_{k-1} \cdots Q_1}_{\text{product of random matrices}} \mathbf{z}^{(0)}$$



Proposition (A.-Wang-Lu 2014)

$$\text{MSE}_N = \text{vec } \mathbf{I}^T (\mathbb{E}\mathbf{Q} \otimes \mathbf{Q})^N \text{vec}(\mathbf{z}^{(0)} \mathbf{z}^{(0)T})$$

$$\text{where } \mathbb{E}\mathbf{Q} \otimes \mathbf{Q} = \sum_i p_i \left(\mathbf{I} - \frac{\mathbf{a}_i \mathbf{a}_i^T}{\|\mathbf{a}_i\|^2} \right)^{\otimes 2}$$

vec — vectorization operator; stack columns of matrix into vector.

\otimes — matrix Kronecker product.

$$\begin{aligned}
\underbrace{\mathbb{E} \|\mathbf{z}^{(N)}\|^2}_{\text{MSE}} &= \mathbb{E} \|\mathbf{Q}_N \mathbf{Q}_{N-1} \cdots \mathbf{Q}_1 \mathbf{z}^{(0)}\|^2 \\
&= \mathbb{E} \mathbf{z}^{(0)T} \mathbf{Q}_1 \mathbf{Q}_2 \cdots \mathbf{Q}_N \mathbf{Q}_N \cdots \mathbf{Q}_2 \mathbf{Q}_1 \mathbf{z}^{(0)} \\
&= \mathbb{E} \text{trace}(\mathbf{Q}_1 \mathbf{Q}_2 \cdots \mathbf{Q}_N \mathbf{Q}_N \cdots \mathbf{Q}_2 \mathbf{Q}_1 \mathbf{z}^{(0)} \mathbf{z}^{(0)T}) \\
&= \mathbb{E} \text{vec}(\mathbf{Q}_1 \mathbf{Q}_2 \cdots \mathbf{Q}_N \mathbf{Q}_N \cdots \mathbf{Q}_2 \mathbf{Q}_1)^T \text{vec}(\mathbf{z}^{(0)} \mathbf{z}^{(0)T})
\end{aligned}$$

Matrix identities

$$\text{vec}(ABC) = (C^T \otimes A) \text{vec}(B)$$

$$\text{trace}(AB) = \text{vec}(A)^T \text{vec}(B)$$

$$\text{trace } AB = \text{trace } BA$$

$$= \mathbb{E} \text{vec}(\mathbf{Q}_1 \mathbf{Q}_2 \cdots \mathbf{Q}_N \mathbf{Q}_N \cdots \mathbf{Q}_2 \mathbf{Q}_1)^T \text{vec}(\mathbf{z}^{(0)} \mathbf{z}^{(0)T})$$

$$\text{MSE}_N$$

$$= \mathbb{E}\{\text{vec}(\mathbf{Q}_1 \mathbf{Q}_2 \cdots \mathbf{Q}_N \mathbf{Q}_N \cdots \mathbf{Q}_2 \mathbf{Q}_1)\}^T \text{vec}(\mathbf{z}^{(0)} \mathbf{z}^{(0)T})$$

$$= \mathbb{E}\{(\mathbf{Q}_1 \otimes \mathbf{Q}_1) \text{vec}(\mathbf{Q}_2 \cdots \mathbf{Q}_N \mathbf{Q}_N \cdots \mathbf{Q}_2)\}^T \text{vec}(\mathbf{z}^{(0)} \mathbf{z}^{(0)T})$$



$$\begin{aligned}
& \text{MSE}_N \\
&= \mathbb{E} \left\{ \text{vec}(\mathbf{Q}_1 \mathbf{Q}_2 \cdots \mathbf{Q}_N \mathbf{Q}_N \cdots \mathbf{Q}_2 \mathbf{Q}_1) \right\}^T \text{vec}(\mathbf{z}^{(0)} \mathbf{z}^{(0)T}) \\
&= \left\{ (\mathbb{E} \mathbf{Q} \otimes \mathbf{Q})^N \text{vec}(\mathbf{I}) \right\}^T \text{vec}(\mathbf{z}^{(0)} \mathbf{z}^{(0)T})
\end{aligned}$$



$$\begin{aligned}
& \text{MSE}_N \\
&= \mathbb{E}\{\text{vec}(\mathbf{Q}_1 \mathbf{Q}_2 \cdots \mathbf{Q}_N \mathbf{Q}_N \cdots \mathbf{Q}_2 \mathbf{Q}_1)\}^T \text{vec}(\mathbf{z}^{(0)} \mathbf{z}^{(0)T}) \\
&= \{(\mathbb{E}\mathbf{Q} \otimes \mathbf{Q})^N \text{vec}(\mathbf{I})\}^T \text{vec}(\mathbf{z}^{(0)} \mathbf{z}^{(0)T}) \\
&= \text{vec} \mathbf{I}^T (\mathbb{E}\mathbf{Q} \otimes \mathbf{Q})^N \text{vec}(\mathbf{z}^{(0)} \mathbf{z}^{(0)T}) \blacksquare
\end{aligned}$$

MSE decays exponentially: $\mathbb{E} \|\mathbf{z}^{(N)}\|^2 = \exp(-\gamma_a N + o(N))$

Error exponent

$$\gamma_a \stackrel{\text{def}}{=} \lim_{N \rightarrow \infty} -\frac{1}{N} \log \mathbb{E} \|\mathbf{z}^{(N)}\|^2$$

We can compute the error exponent:

$$\gamma_a = -\log \lambda_{\max} \left(\sum_i p_i \left(\mathbf{I} - \frac{\mathbf{a}_i \mathbf{a}_i^T}{\|\mathbf{a}_i\|^2} \right)^{\otimes 2} \right)$$

- Can be computed in $O(mn^2)$ time
- Must be positive; exponential convergence confirmed

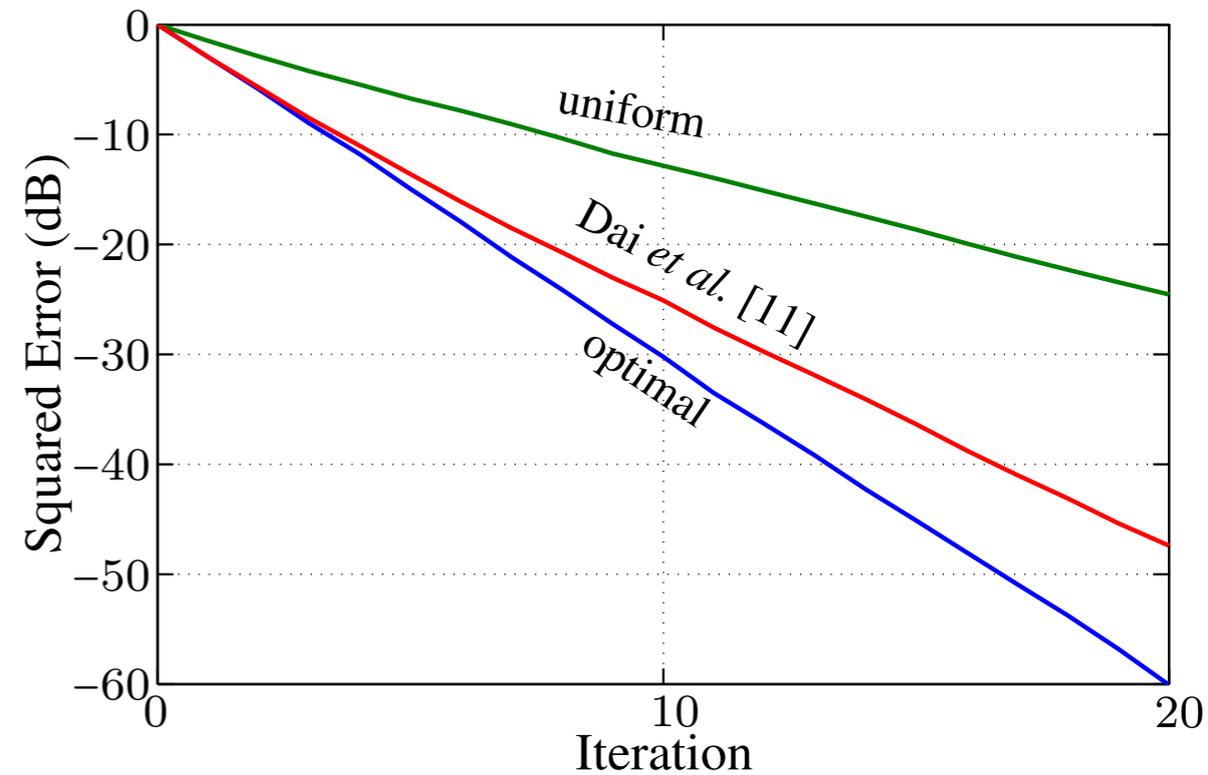
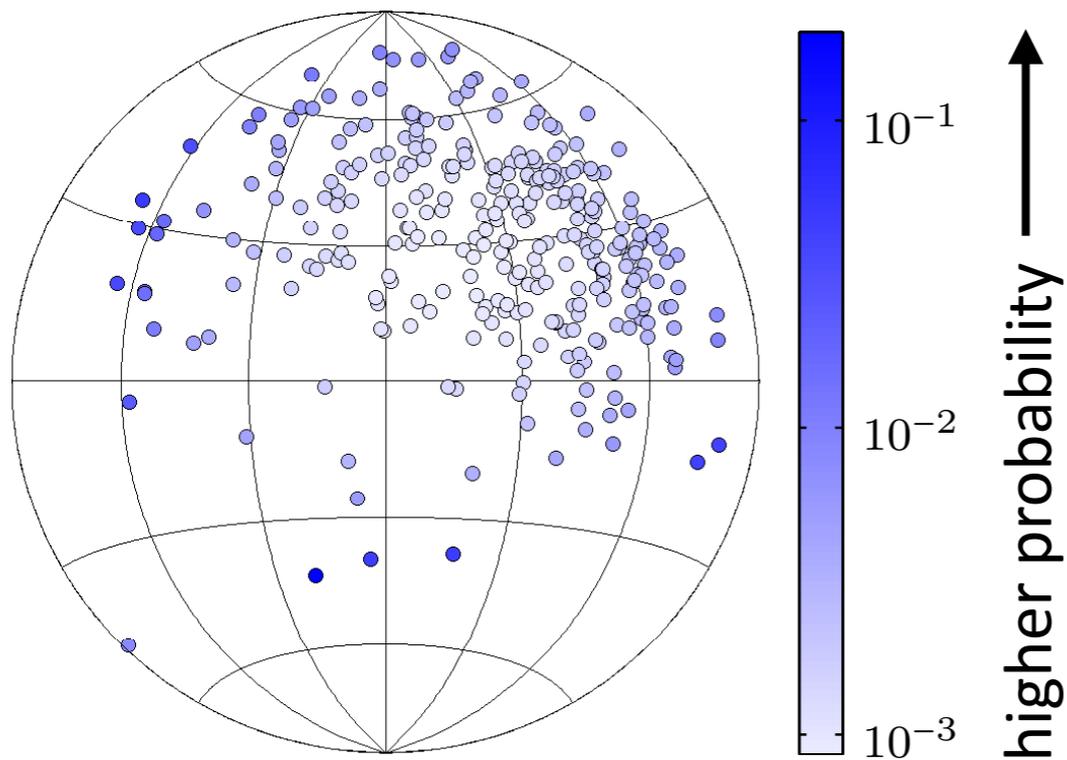
1. Exact MSE formula and decay rate
- 2. Optimization of row selection probabilities**
3. “Quenched error exponent”

Convex optimization problem: minimize error exponent.

$$(p_1, \dots, p_m) = \operatorname{argmin}_{\mathbf{p}} \lambda_{\max} \left(\sum_i p_i \left(\mathbf{I} - \frac{\mathbf{a}_i \mathbf{a}_i^T}{\|\mathbf{a}_i\|^2} \right)^{\otimes 2} \right)$$

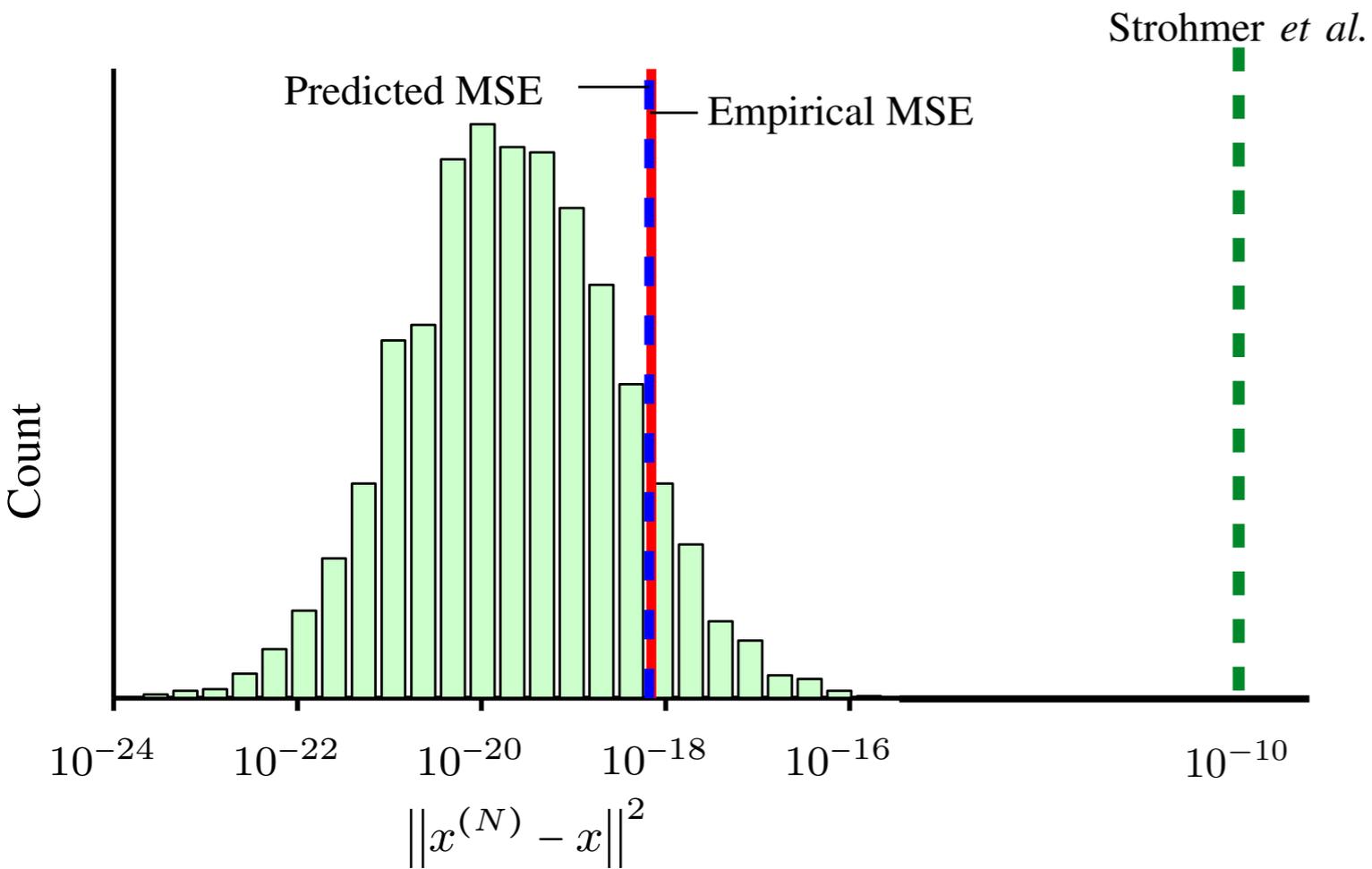
semi-definite programming

$n = 3$ lets us easily visualize the optimal probabilities

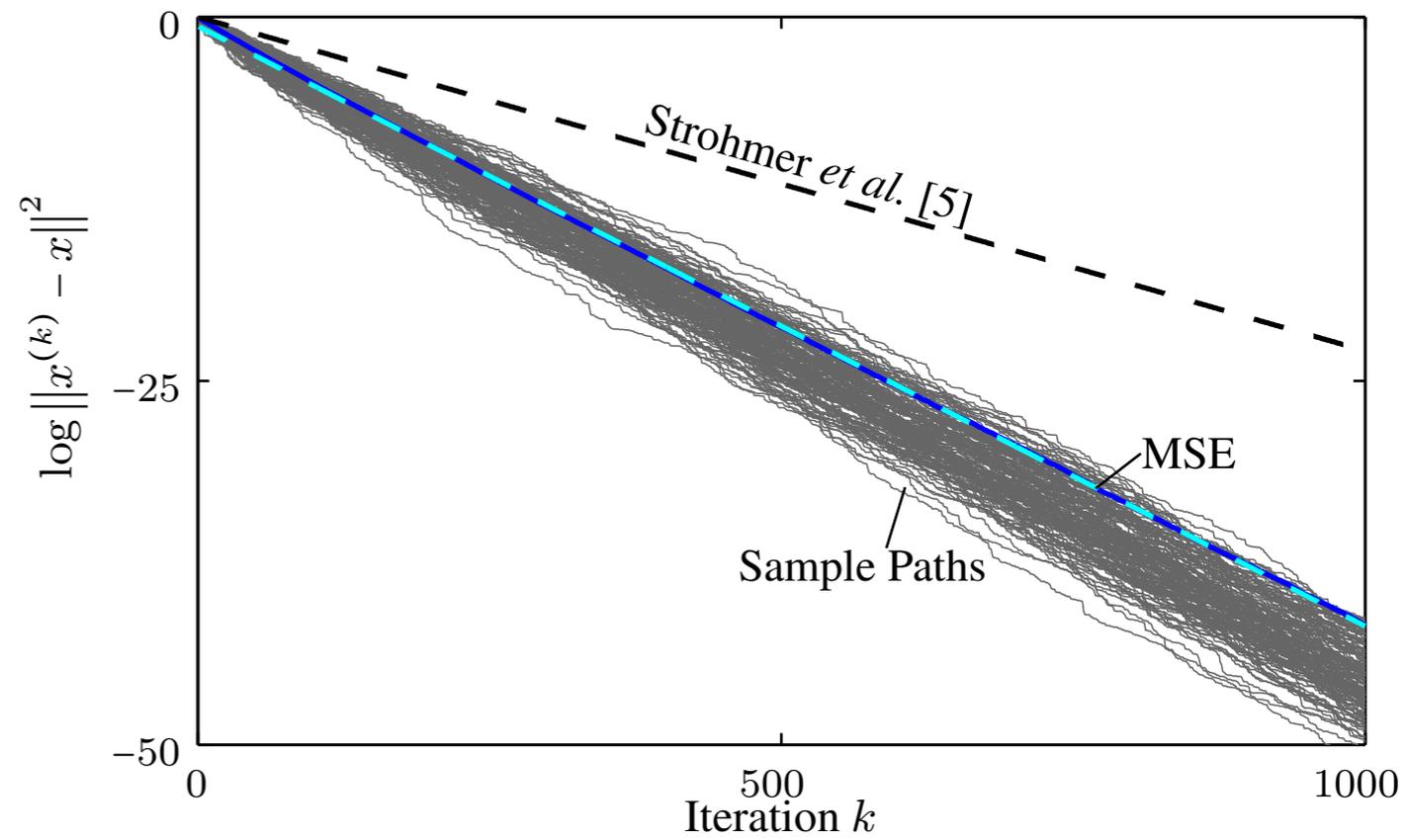
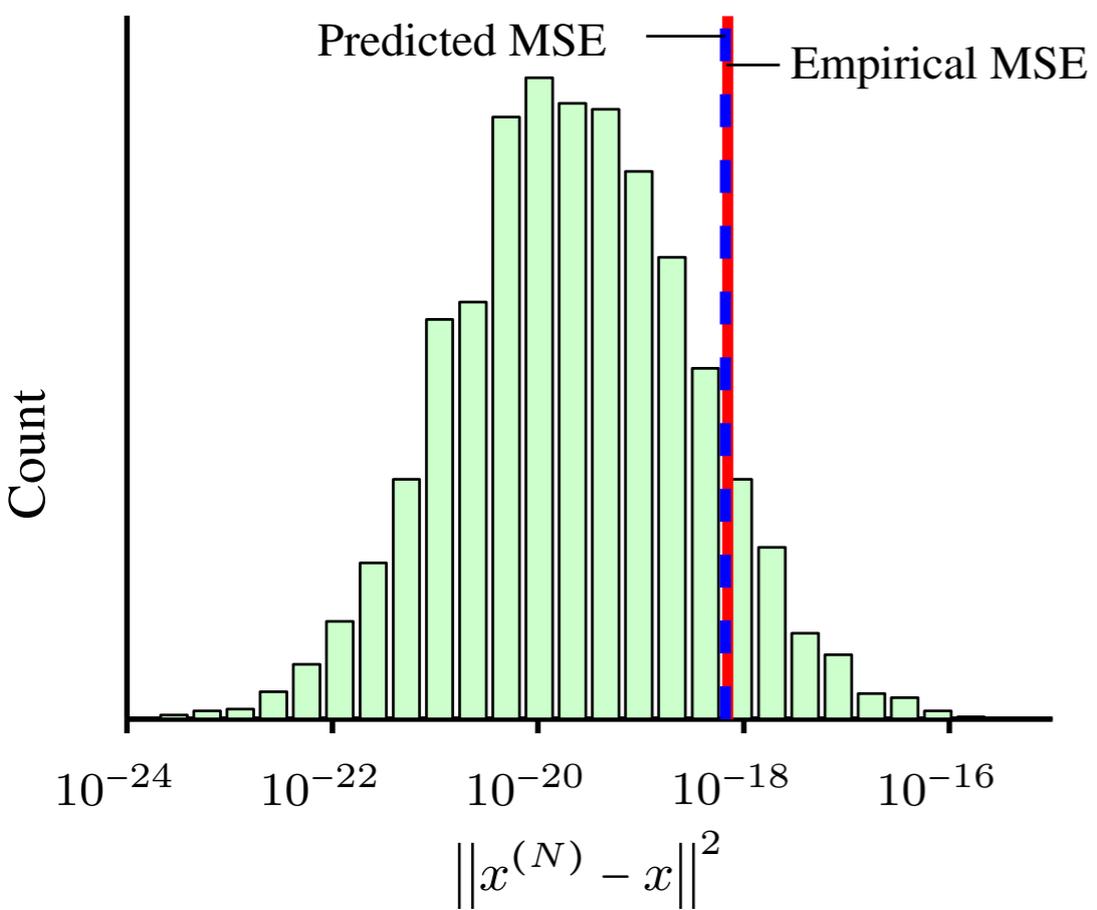


Intuition: explorers of sparsely-populated regions chosen with higher probability

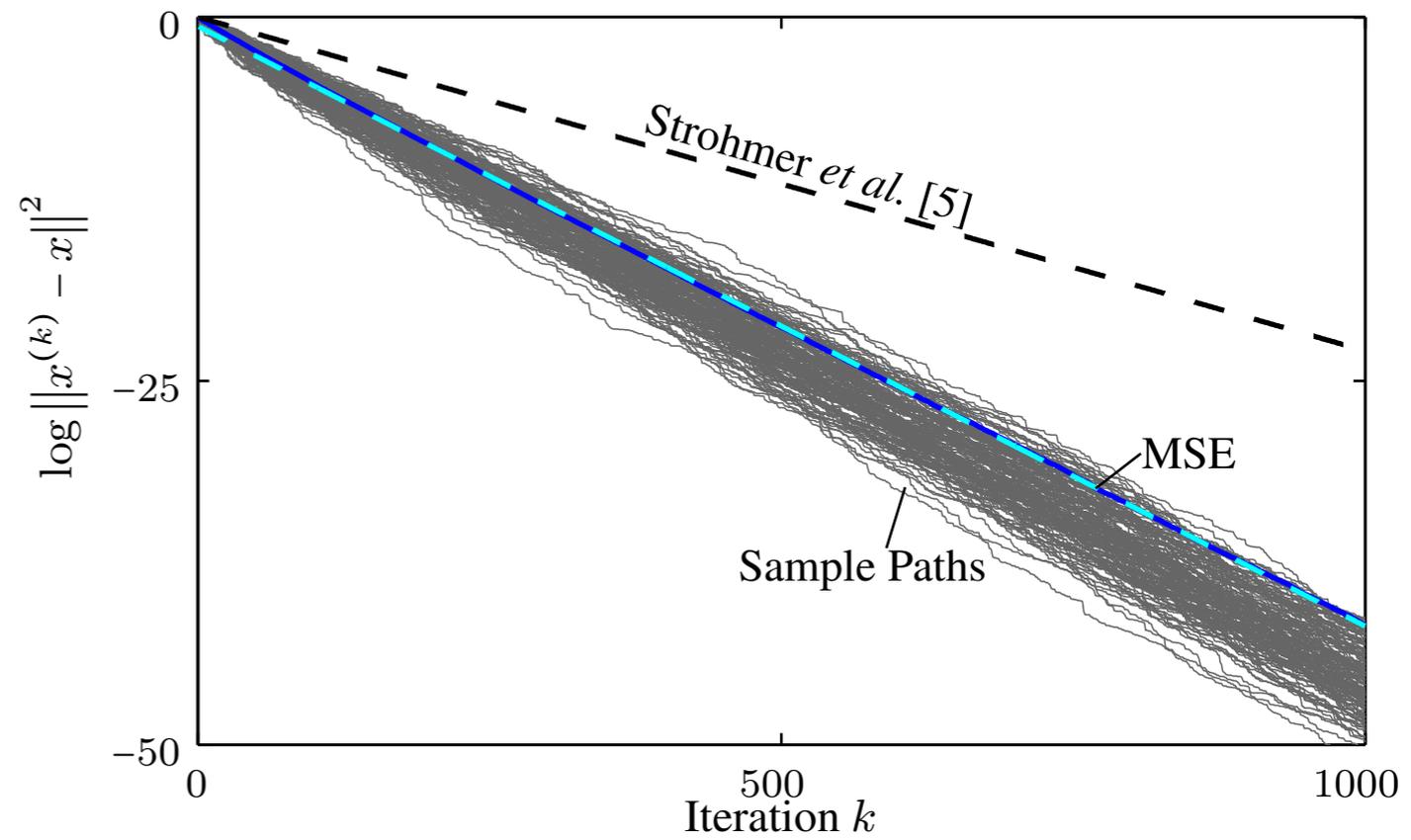
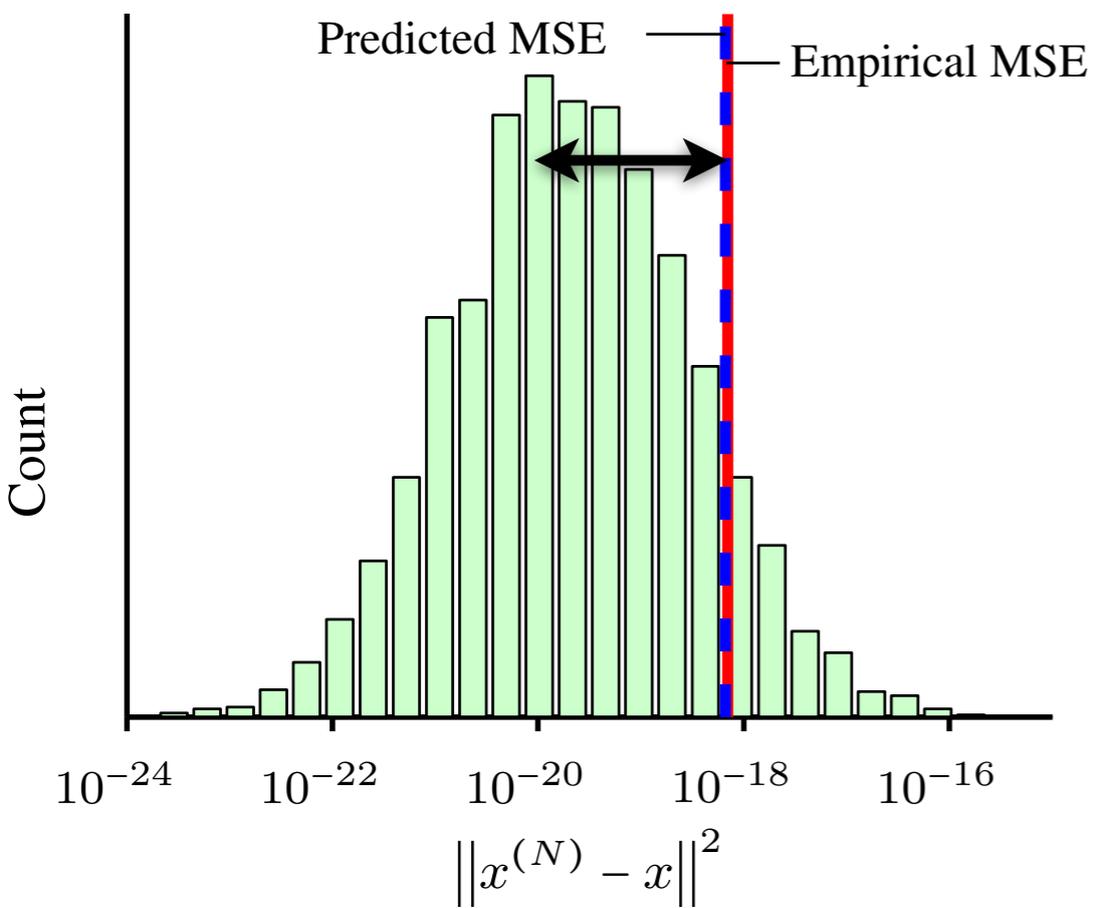
1. Exact MSE formula and decay rate
2. Optimization of row selection probabilities
- 3. “Quenched error exponent”**



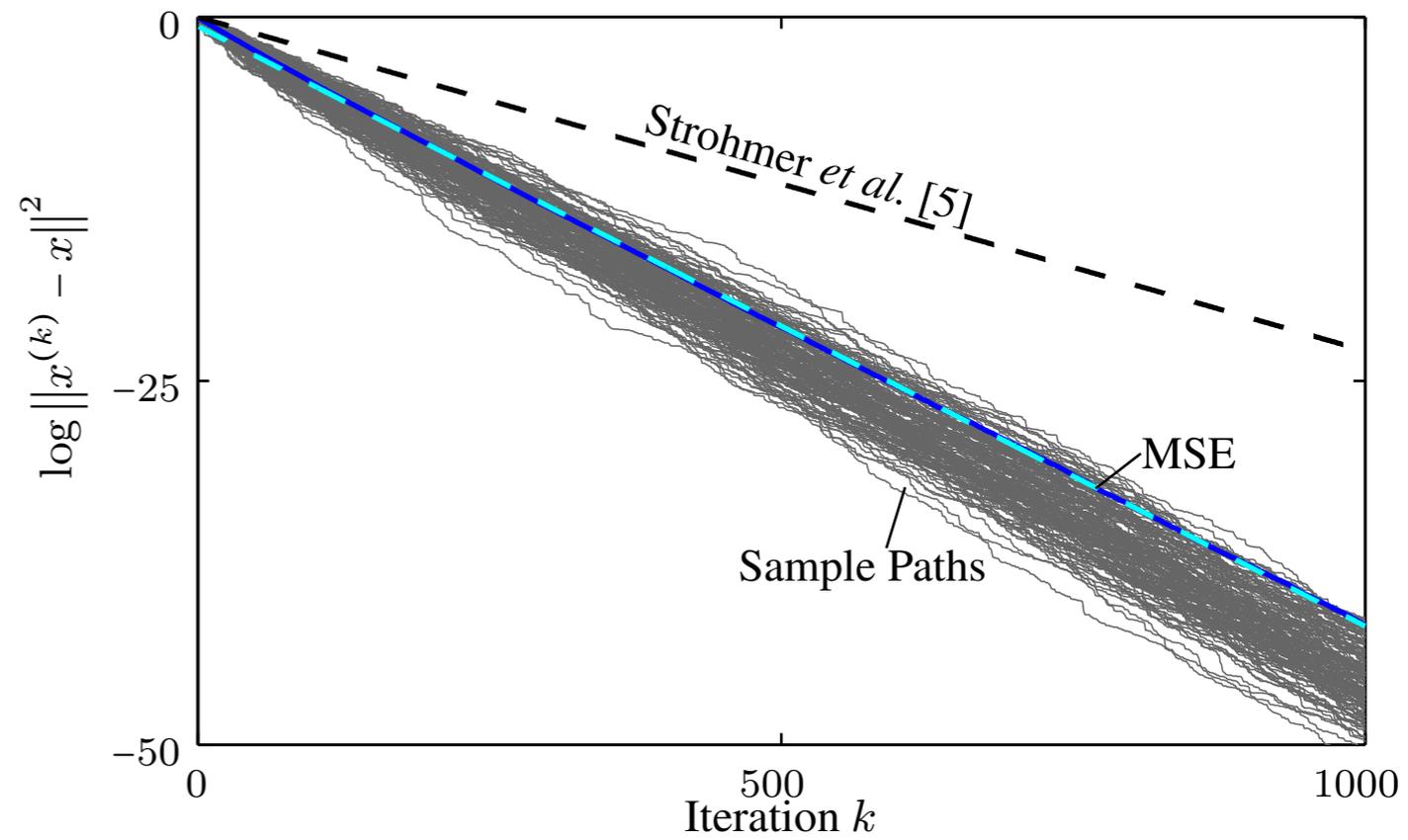
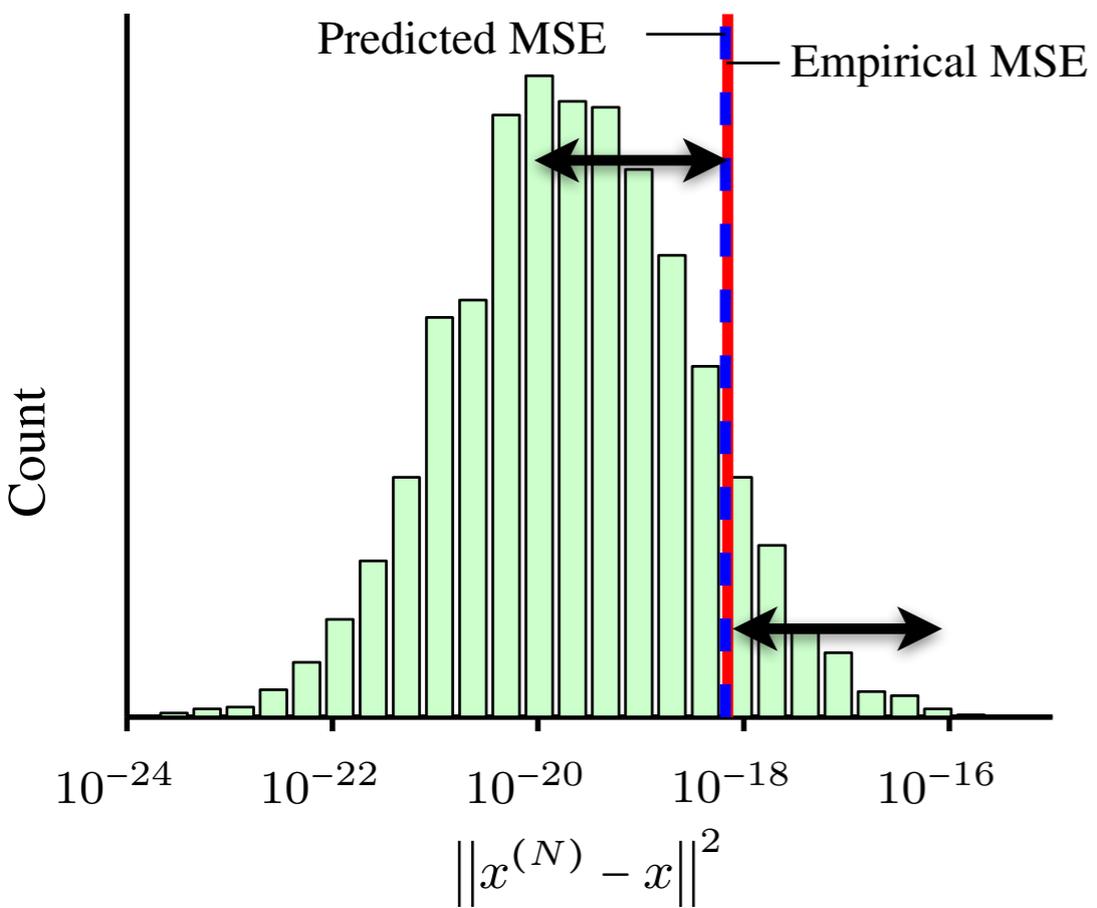
150 x 20 matrix w/ Gaussian entries.



150 x 20 matrix w/ Gaussian entries.



150 x 20 matrix w/ Gaussian entries.



150 x 20 matrix w/ Gaussian entries.

**IF YOU WANT TO MEASURE
TYPICAL PERFORMANCE**

**DON'T
USE the MSE!**

Average performance:

Annealed error exponent

$$\gamma_a \stackrel{\text{def}}{=} \lim_{N \rightarrow \infty} -\frac{1}{N} \log \mathbb{E} \|\mathbf{z}^{(N)}\|^2$$

Typical performance:

Quenched error exponent

$$\gamma_q \stackrel{\text{def}}{=} \lim_{N \rightarrow \infty} -\frac{1}{N} \mathbb{E} \log \|\mathbf{z}^{(N)}\|^2$$

- Much more difficult to analyze.
- Known to physicists as the top *Lyapunov exponent*.
- They use heuristics to solve.

Physicists have their own intuition for this trick, but we can get the same result by assuming the error is log-normal:

Assume

$$\log \|\mathbf{z}^{(N)}\|^2 \sim \mathcal{N}(N\mu, N\sigma^2).$$

Then

$$\gamma_q = -\mu$$

$$\mathbb{E} \|\mathbf{z}^{(N)}\|^2 = \exp(N[\mu + \frac{1}{2}\sigma^2])$$

$$\mathbb{E} \|\mathbf{z}^{(N)}\|^4 = \exp(N[2\mu + 2\sigma^2])$$

Naive replica method $\log Z = \lim_{n \rightarrow 0} \frac{Z^n - 1}{n}$

Solve

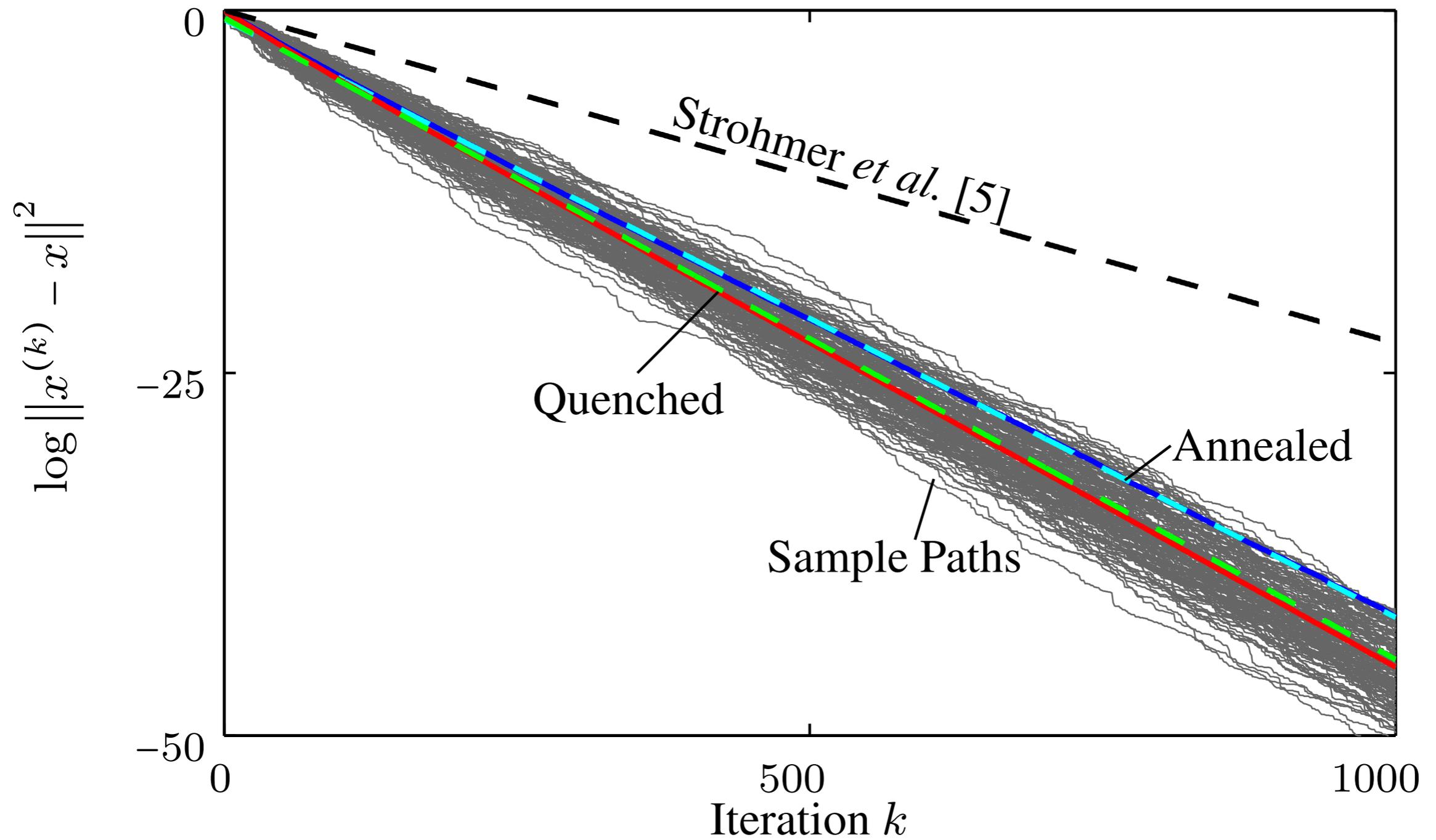
$$\mu = \frac{1}{N} \left[2 \log \mathbb{E} \|\mathbf{z}^{(N)}\|^2 - \frac{1}{2} \log \mathbb{E} \|\mathbf{z}^{(N)}\|^4 \right]$$

$$\gamma_q \approx 2\gamma_a - \frac{1}{2}\gamma_a^{(2)}$$

where $\gamma_a^{(2)} = -\log \lambda_{\max} \left(\sum_i p_i \left(\mathbf{I} - \frac{\mathbf{a}_i \mathbf{a}_i^T}{\|\mathbf{a}_i\|^2} \right)^{\otimes 4} \right)$

Quenched error exponent

$$\gamma_q \stackrel{\text{def}}{=} \lim_{N \rightarrow \infty} -\frac{1}{N} \mathbb{E} \log \|\mathbf{z}^{(N)}\|^2$$



- *Exact MSE formula for randomized Kaczmarz algorithms (and its generalizations)*

Lifting!

- *Annealed and quenched error exponents give decay rate*

Average vs. typical performance

- *Finding optimal row selection probabilities*

Convex optimization