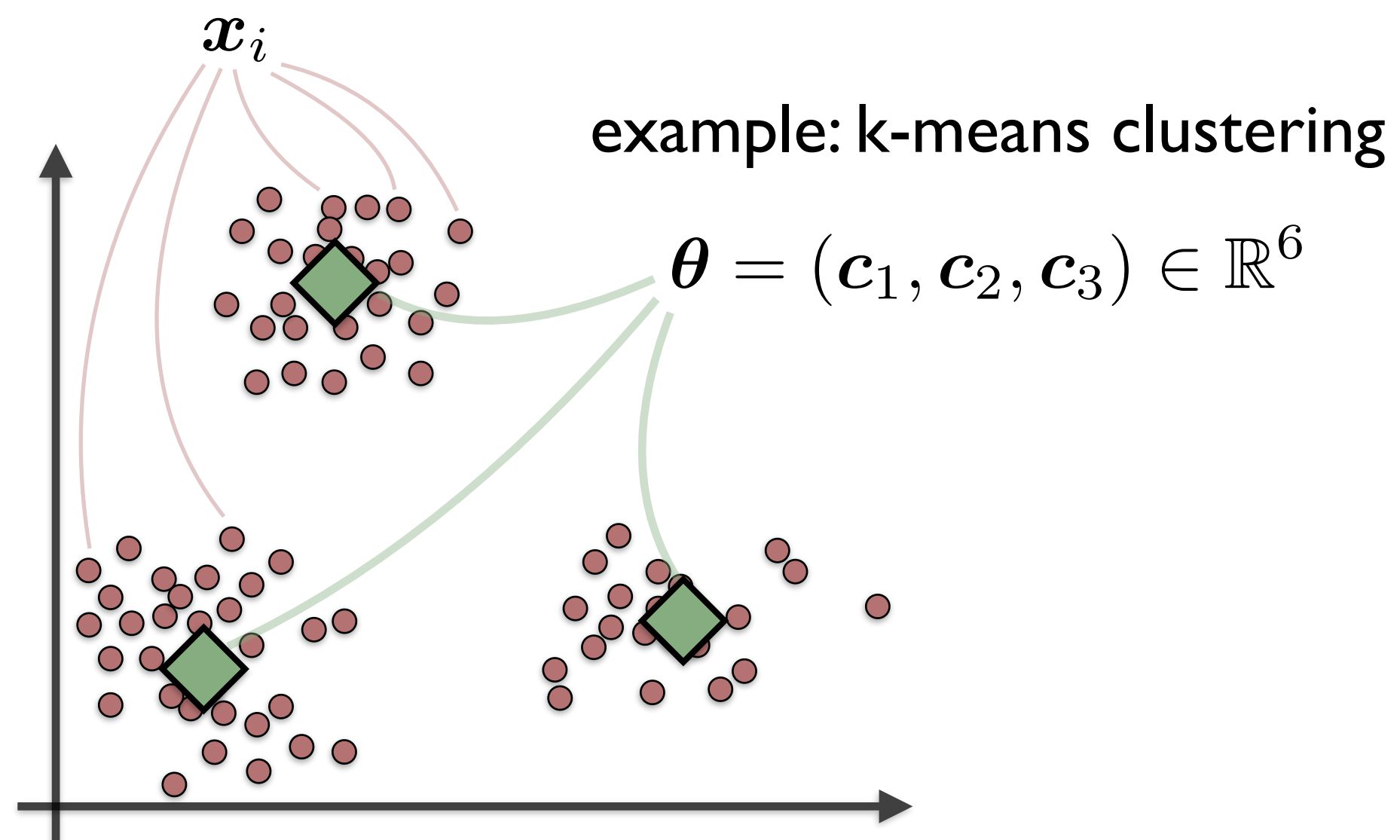
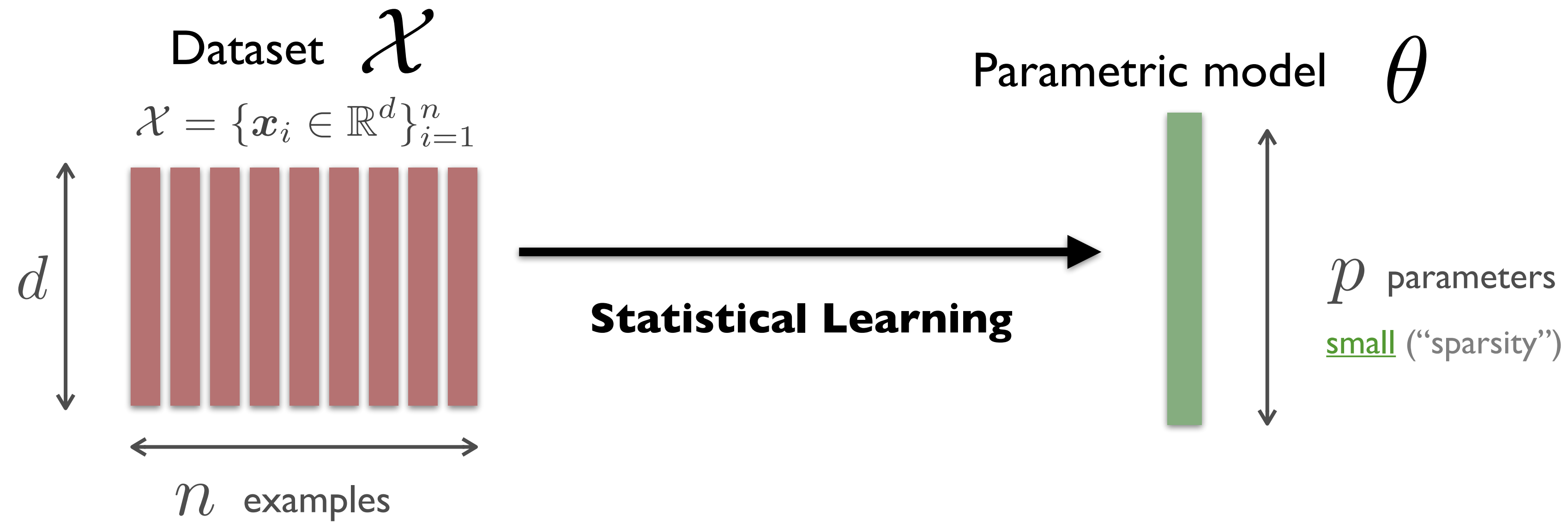


When compressive learning fails: blame the decoder or the sketch?

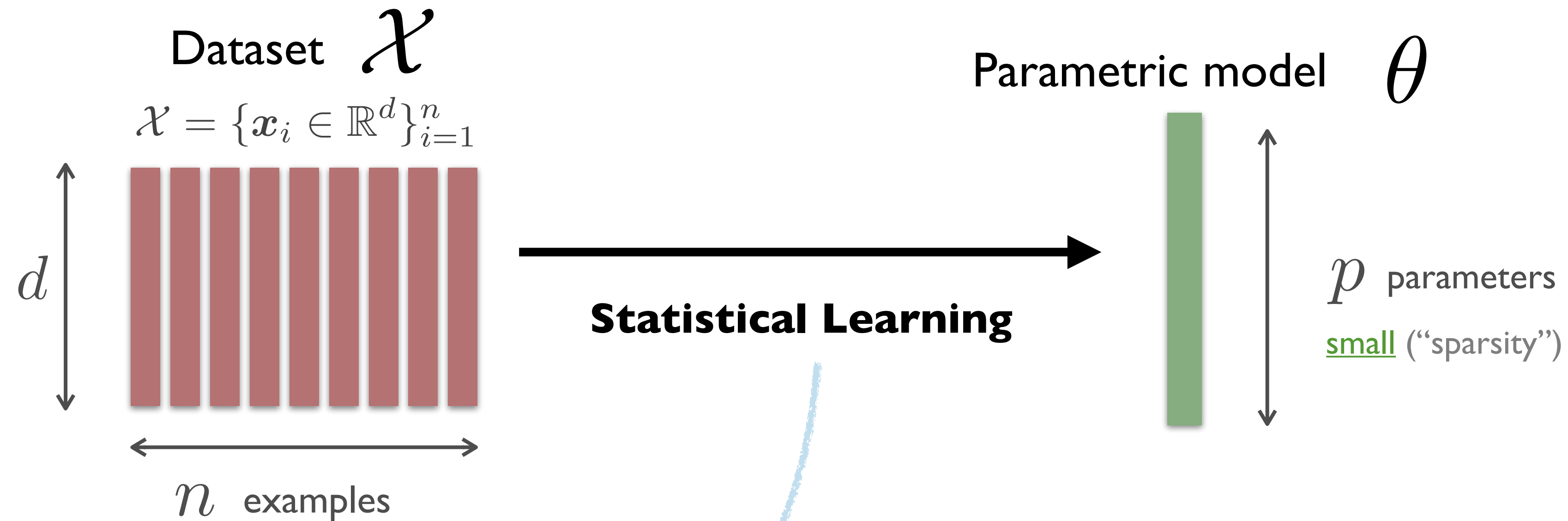
Vincent Schellekens

Laurent Jacques

Context: large-scale learning



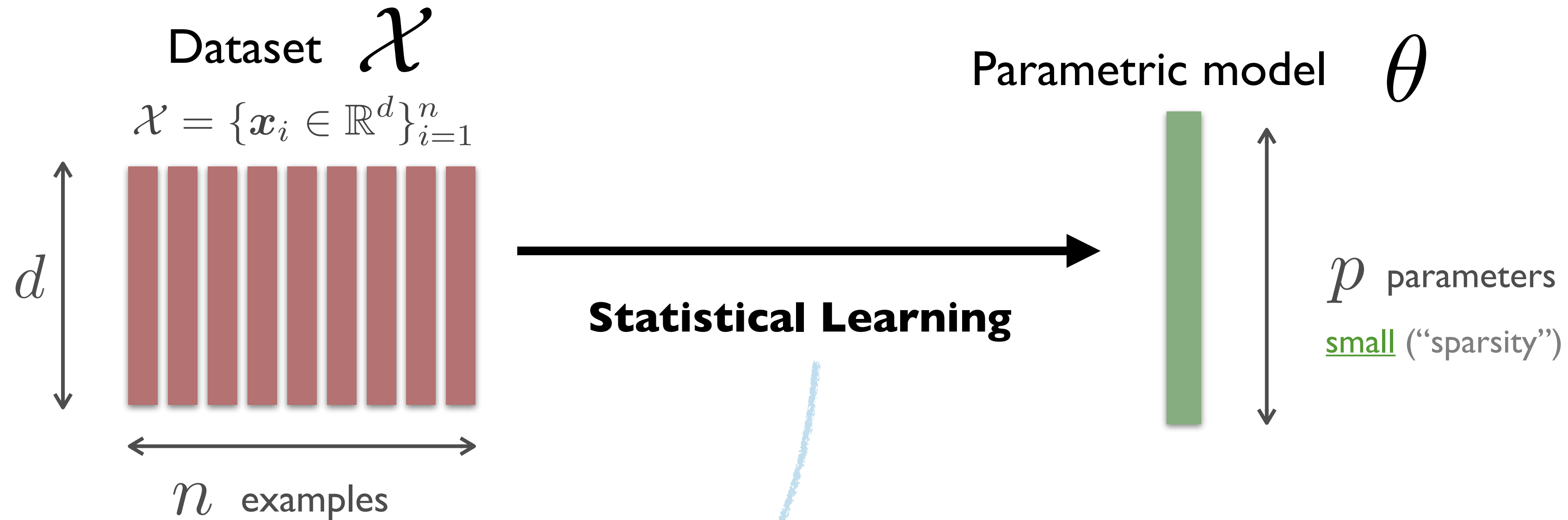
Context: large-scale learning



Choose parameters minimizing the risk

$$\hat{\theta} = \arg \min_{\theta} R(\theta | \mathcal{X}) \quad \text{with} \quad R(\theta | \mathcal{X}) = \frac{1}{n} \sum_{i=1}^n \ell(\theta, \mathbf{x}_i)$$

Context: large-scale learning



Choose parameters minimizing the risk

$$\hat{\theta} = \arg \min_{\theta} R(\theta | \mathcal{X}) \quad \text{with} \quad R(\theta | \mathcal{X}) = \frac{1}{n} \sum_{i=1}^n \ell(\theta, \mathbf{x}_i)$$

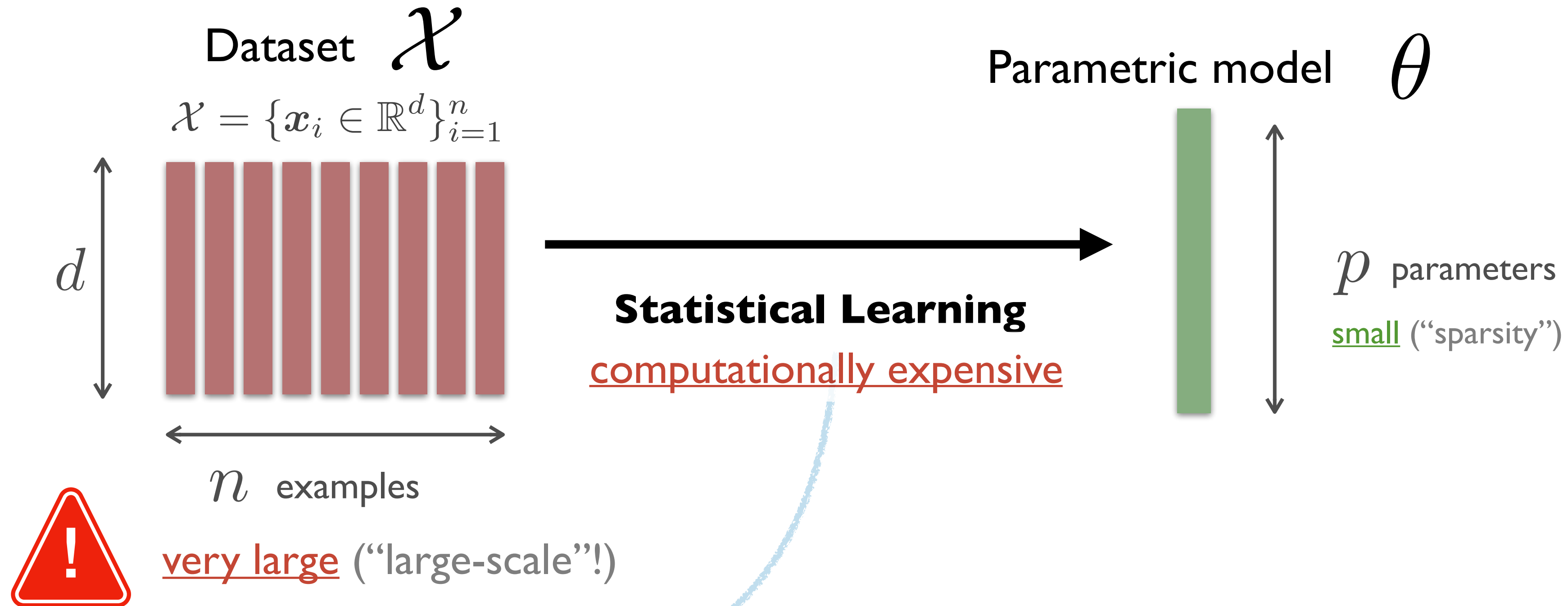
Examples:

k-means clustering:

$$R(\theta; \mathcal{X}) = \frac{1}{n} \sum_{i=1}^n \min_k \|\mathbf{c}_k - \mathbf{x}_i\|_2$$

Gaussian Mixture Model: $R(\theta; \mathcal{X}) = \frac{1}{n} \sum_{i=1}^n -\log(p_{\theta}(\mathbf{x}_i))$

Context: large-scale learning

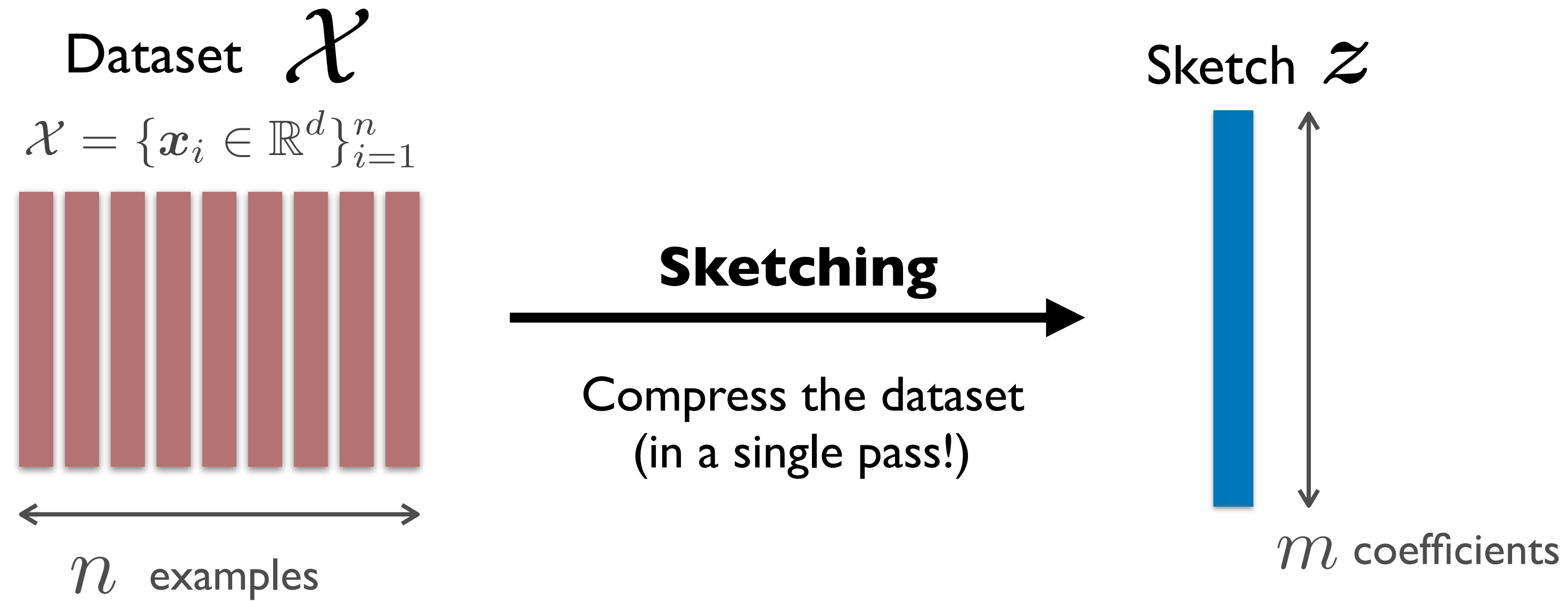


Choose parameters minimizing the risk

$$\hat{\theta} = \arg \min_{\theta} R(\theta | \mathcal{X}) \quad \text{with} \quad R(\theta | \mathcal{X}) = \frac{1}{n} \sum_{i=1}^n \ell(\theta, \mathbf{x}_i)$$

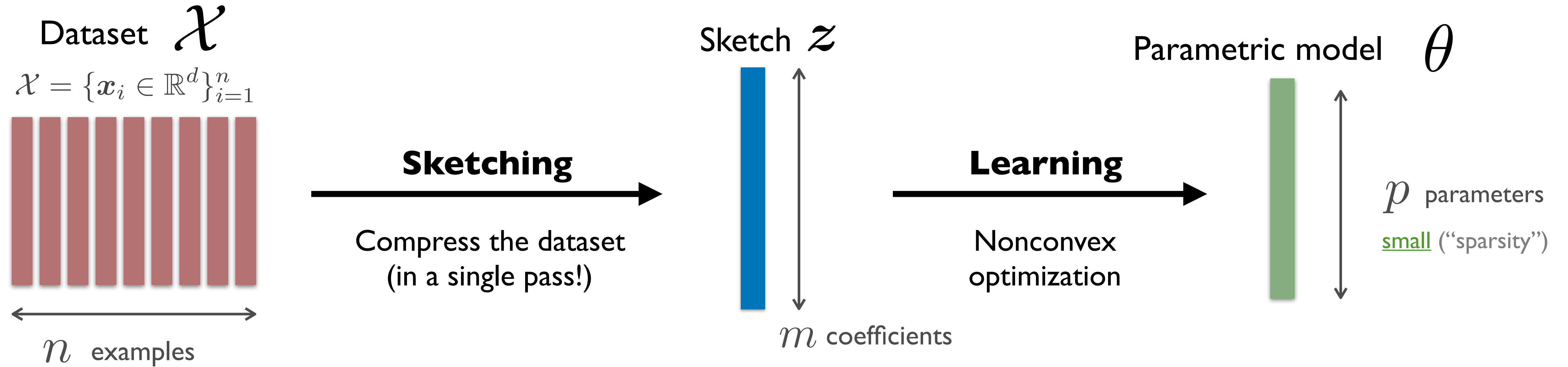
Typically, several passes over the dataset are needed

Compressive learning to the rescue



Break the training down into two “cheaper” steps

Compressive learning to the rescue



Break the training down into two “cheaper” steps

This work investigates the interplay between those steps

Sketching stage

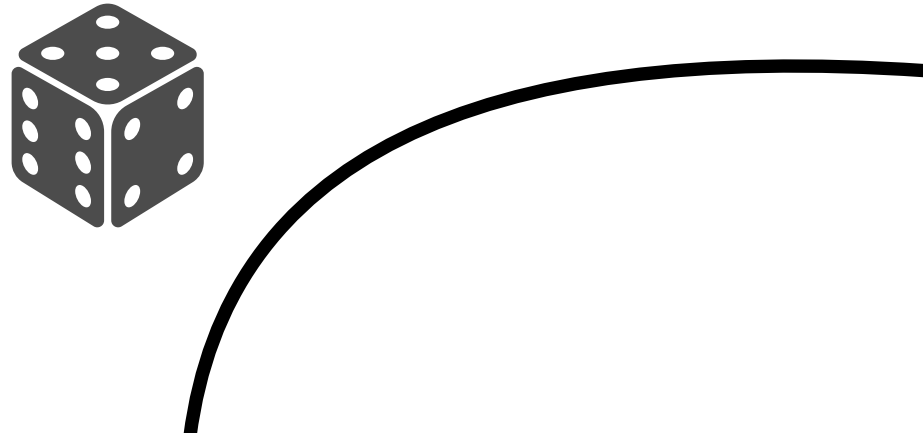
Sketch = average (one pass!) of features $\Phi : \mathbb{R}^d \rightarrow \mathbb{C}^m$

$$\mathbf{z} = \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}_i)$$

Sketching stage

Sketch = average (one pass!) of features $\Phi : \mathbb{R}^d \rightarrow \mathbb{C}^m$

Typically, random features (random linear sensing of $\hat{\mathcal{P}}_{\mathcal{X}}$)



$$\mathbf{z} = \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}_i)$$

“Low-dimensional embedding”
of the (empirical) data distribution

$$\mathbf{z} = \mathcal{A}(\hat{\mathcal{P}}_{\mathcal{X}}) = \mathcal{A}\left(\frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i}\right)$$

Sketching stage

Sketch = average (one pass!) of features $\Phi : \mathbb{R}^d \rightarrow \mathbb{C}^m$


$$\mathbf{z} = \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}_i)$$

Typically, random features (random linear sensing of $\hat{\mathcal{P}}_{\mathcal{X}}$)

In this work: random Fourier features

$$\Phi_j(\mathbf{x}) = \exp(i\boldsymbol{\omega}_j^\top \mathbf{x}), \quad j = 1, \dots, m$$

$$\boldsymbol{\omega}_j \sim i.i.d. \quad \Lambda(\boldsymbol{\omega}) = (\mathcal{F}K)(\boldsymbol{\omega})$$

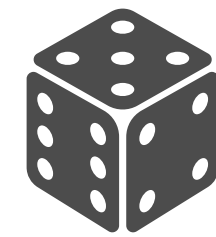
“Low-dimensional embedding”
of the (empirical) data distribution

$$\mathbf{z} = \mathcal{A}(\hat{\mathcal{P}}_{\mathcal{X}}) = \mathcal{A}\left(\frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i}\right)$$

Sketching stage

Sketch = average (one pass!) of features $\Phi : \mathbb{R}^d \rightarrow \mathbb{C}^m$

$$\mathbf{z} = \frac{1}{n} \sum_{i=1}^n \Phi(\mathbf{x}_i)$$



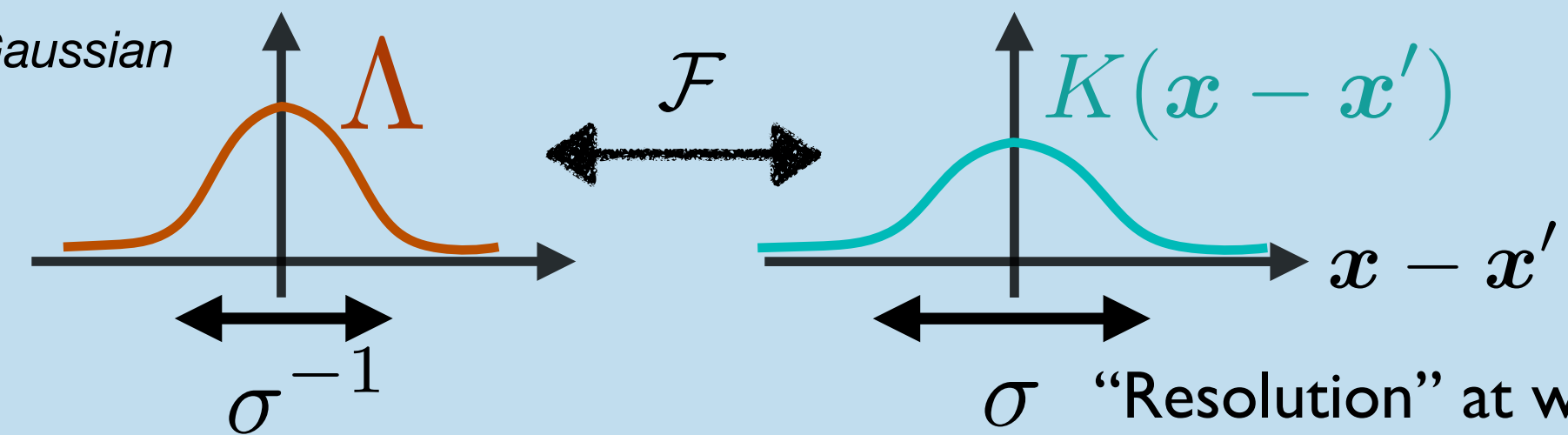
Typically, random features (random linear sensing of $\hat{\mathcal{P}}_{\mathcal{X}}$)

In this work: random Fourier features

$$\Phi_j(\mathbf{x}) = \exp(i\boldsymbol{\omega}_j^\top \mathbf{x}), \quad j = 1, \dots, m$$

$$\boldsymbol{\omega}_j \sim i.i.d. \quad \Lambda(\boldsymbol{\omega}) = (\mathcal{F}K)(\boldsymbol{\omega})$$

e.g. Gaussian



“Low-dimensional embedding”
of the (empirical) data distribution

$$\mathbf{z} = \mathcal{A}(\hat{\mathcal{P}}_{\mathcal{X}}) = \mathcal{A}\left(\frac{1}{n} \sum_{i=1}^n \delta_{\mathbf{x}_i}\right)$$

Learning stage

Optimize some “sketch matching” loss

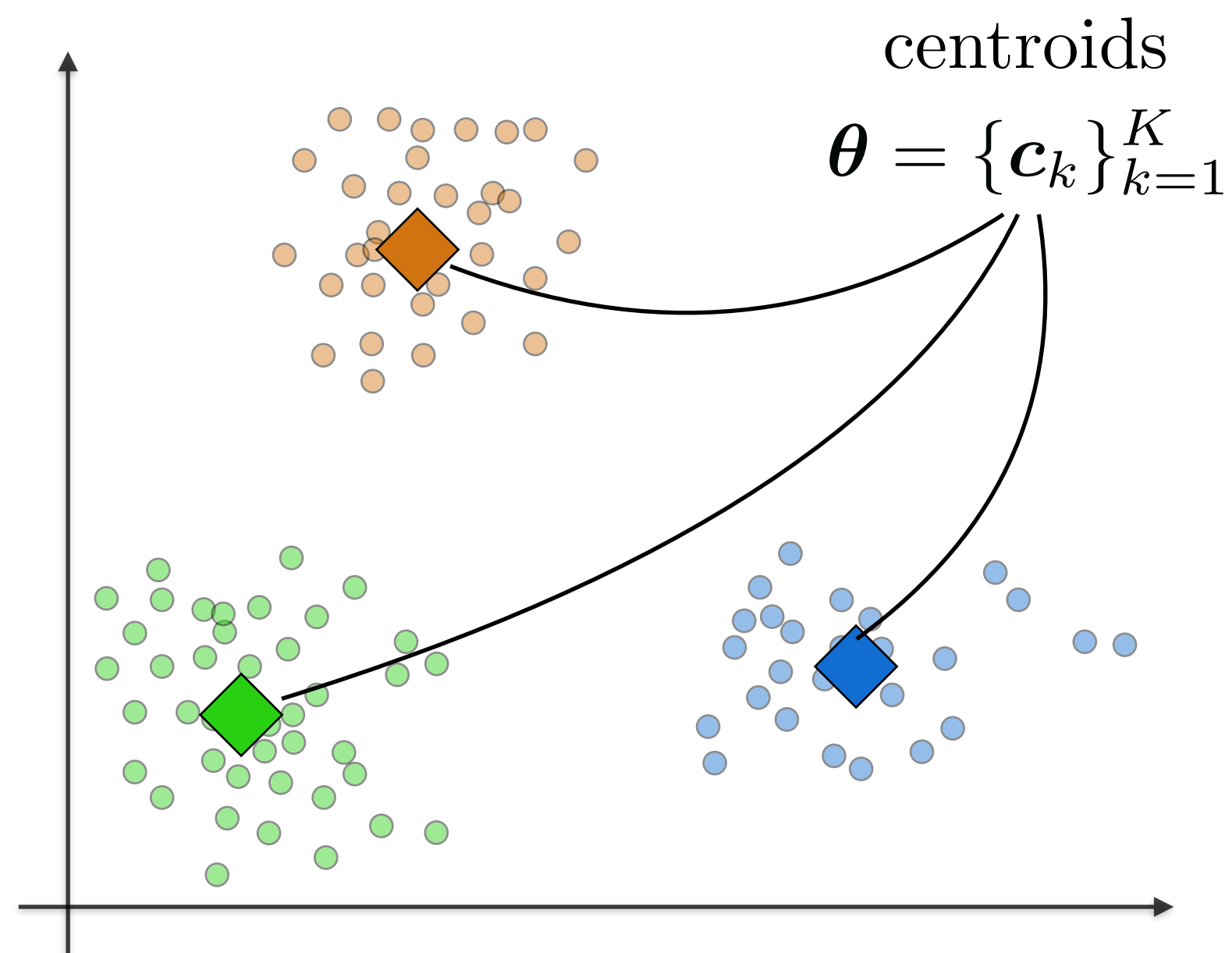
$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}; \mathbf{z})$$

Learning stage

Optimize some “sketch matching” loss

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}; \mathbf{z})$$

Example: k-means $R(\boldsymbol{\theta}; \mathcal{X}) = \sum_i \min_k \|\mathbf{c}_k - \mathbf{x}_i\|_2$

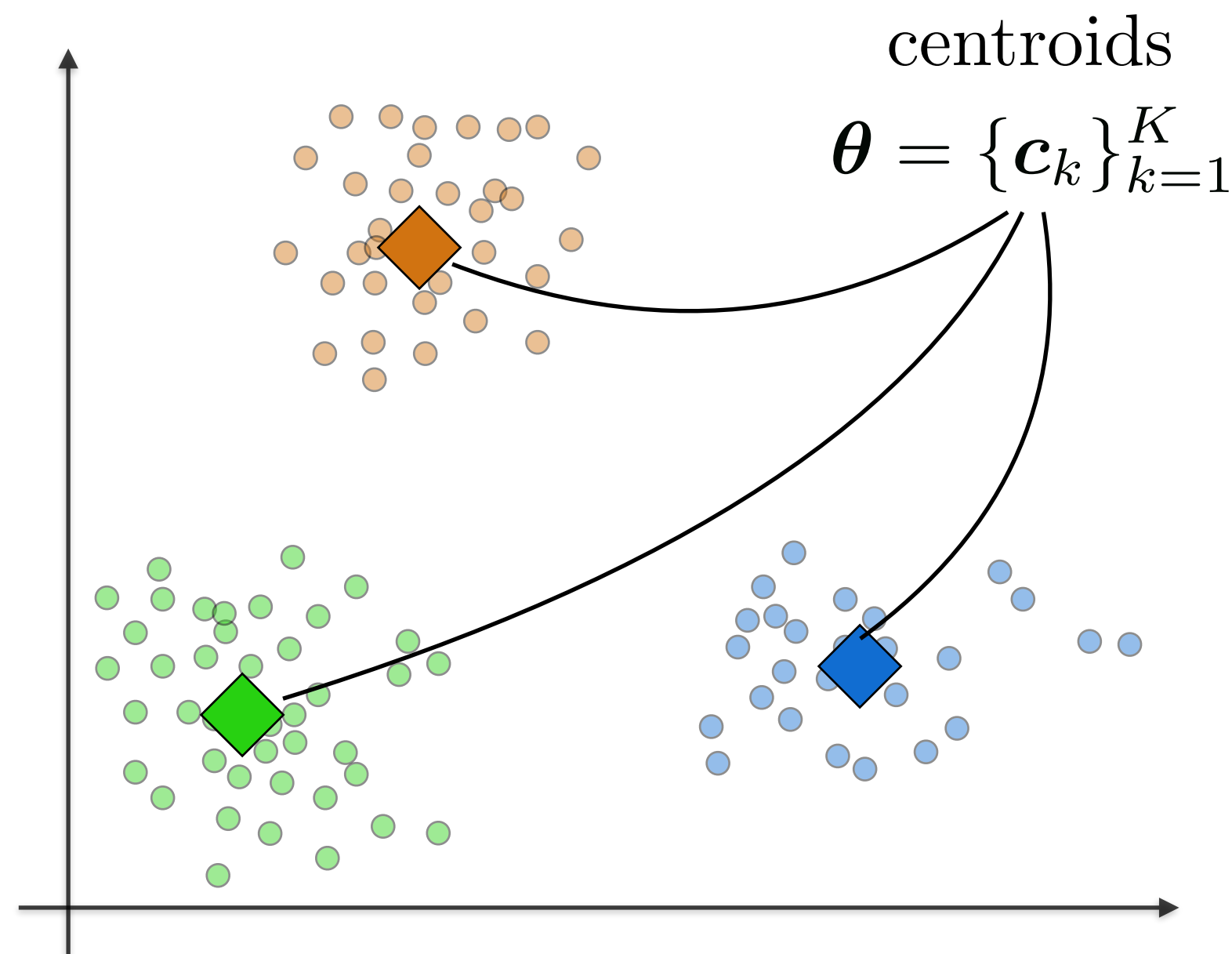


Learning stage

Optimize some “sketch matching” loss

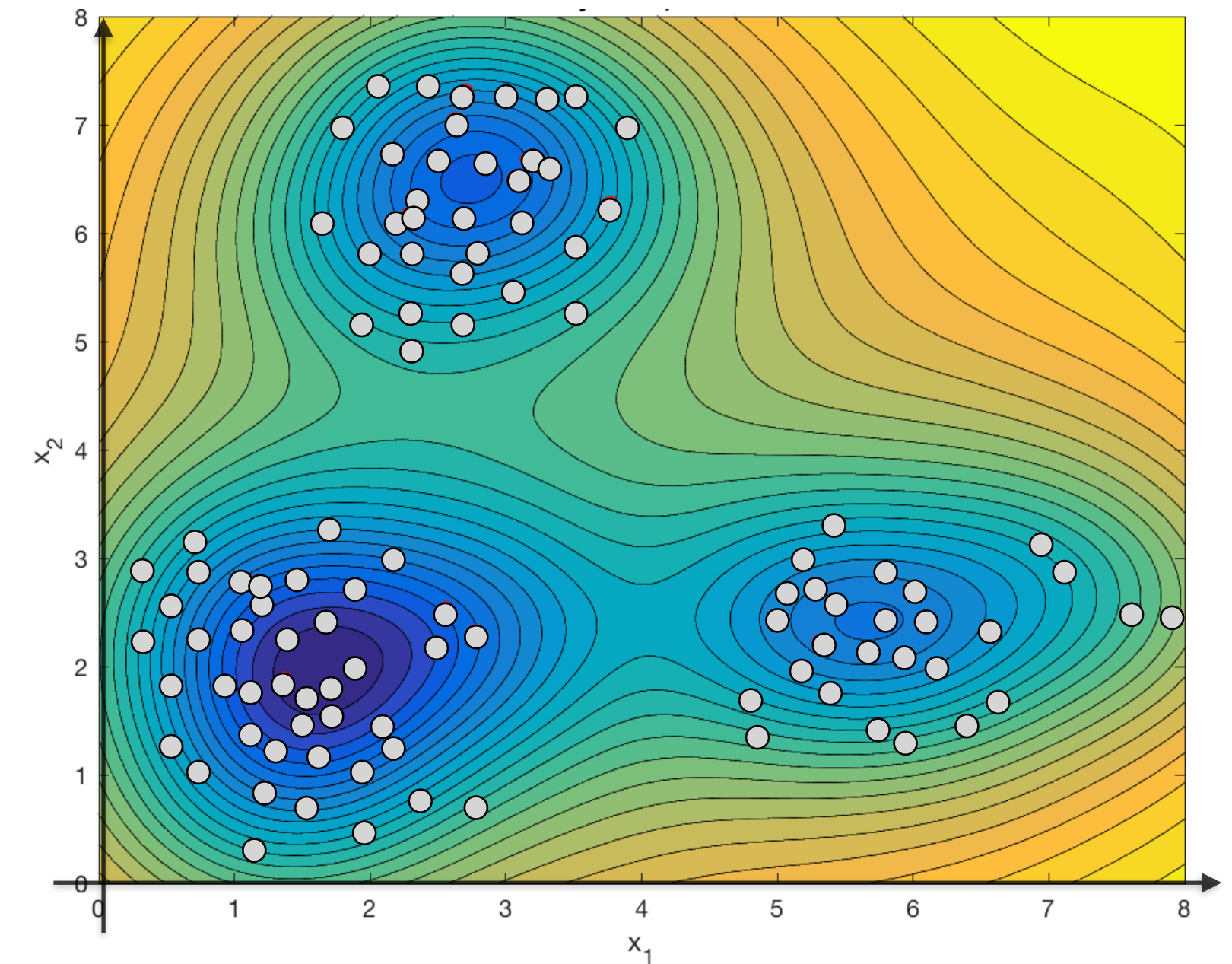
$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}(\theta; \mathbf{z})$$

Example: k-means $R(\theta; \mathcal{X}) = \sum_i \min_k \|\mathbf{c}_k - \mathbf{x}_i\|_2$



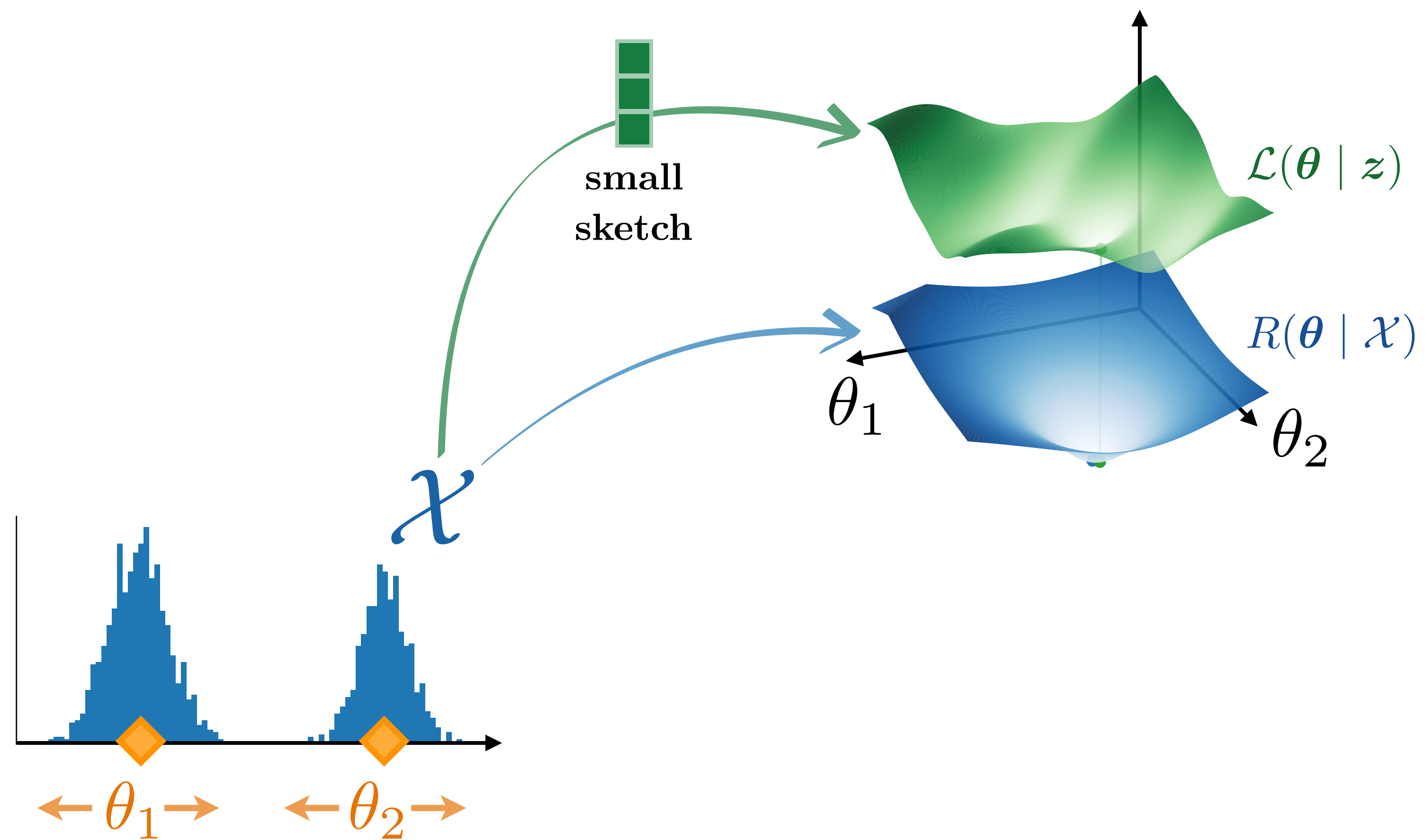
$$\mathcal{L}(\theta; \mathbf{z}) = \|\mathbf{z} - \underbrace{\frac{1}{K} \sum_{k=1}^K \Phi(\mathbf{c}_k)}_{\text{sketch of the centroids}}\|_2$$

sketch of the centroids



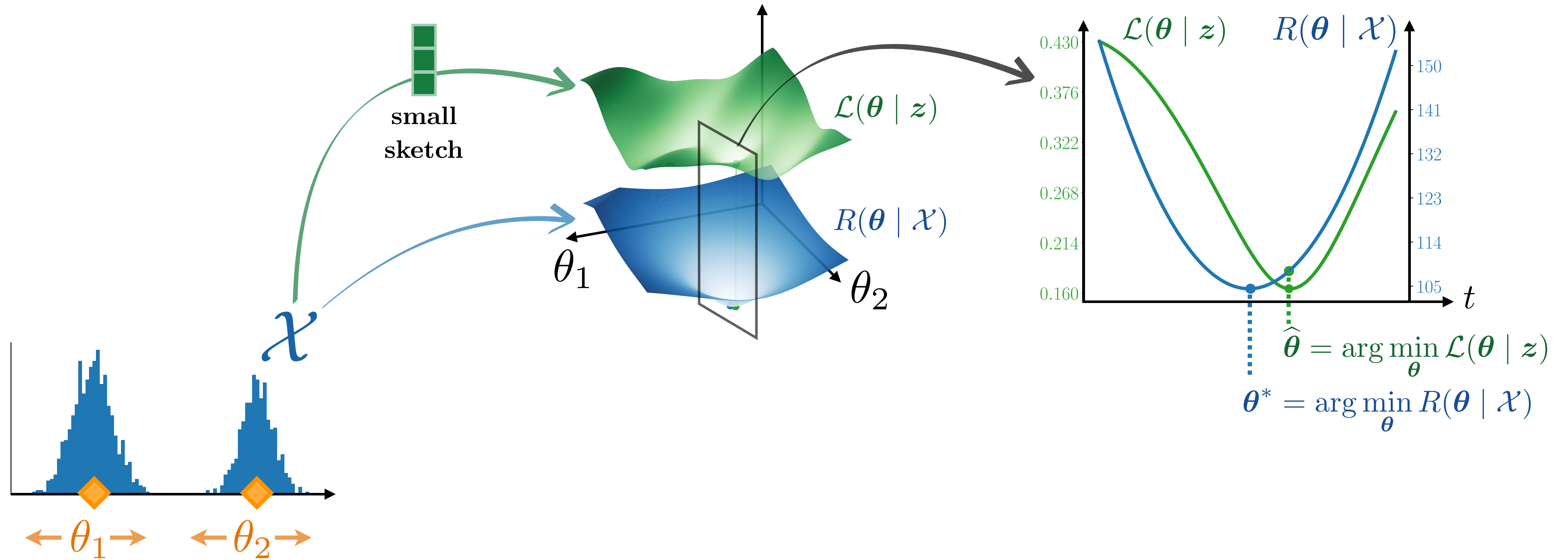
Learning stage

Cost function vs. true risk



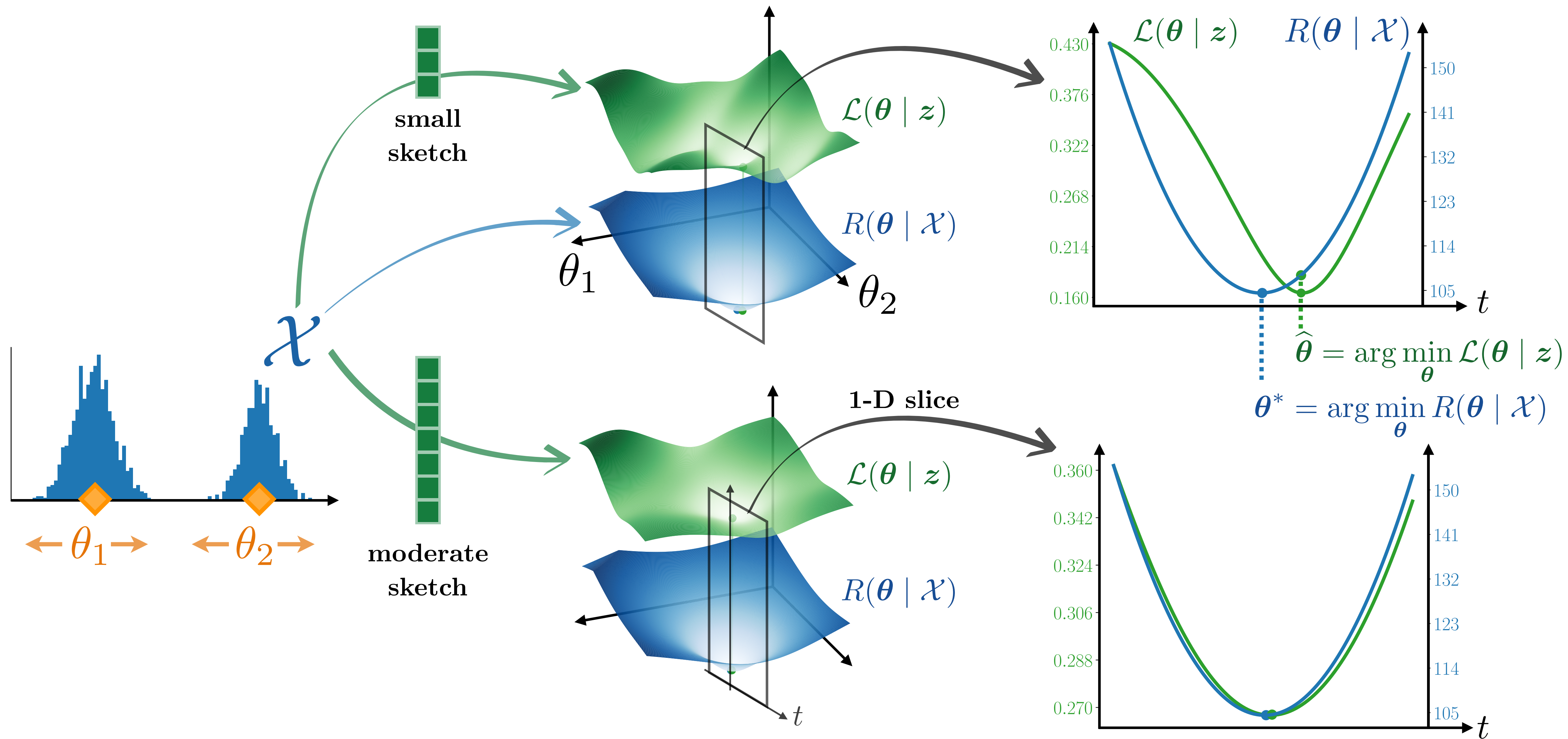
Learning stage

Cost function vs. true risk



Learning stage

Cost function vs. true risk



Learning stage

Optimize some “sketch matching” loss

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}; \mathbf{z})$$

Nonconvex! In practice, we use a “decoder”

$$\boldsymbol{\theta}_{\Delta} = \Delta[\mathbf{z}]$$

The CLOMPR decoder (based on matching pursuit) showed promising empirical success...

Learning stage

Optimize some “sketch matching” loss

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}; \mathbf{z})$$

Nonconvex! In practice, we use a “decoder”

$$\boldsymbol{\theta}_{\Delta} = \Delta[\mathbf{z}]$$

The CLOMPR decoder (based on matching pursuit) showed promising empirical success...

Our question: should we improve the decoder and/or sketch?

If so, when? What can we expect to gain?

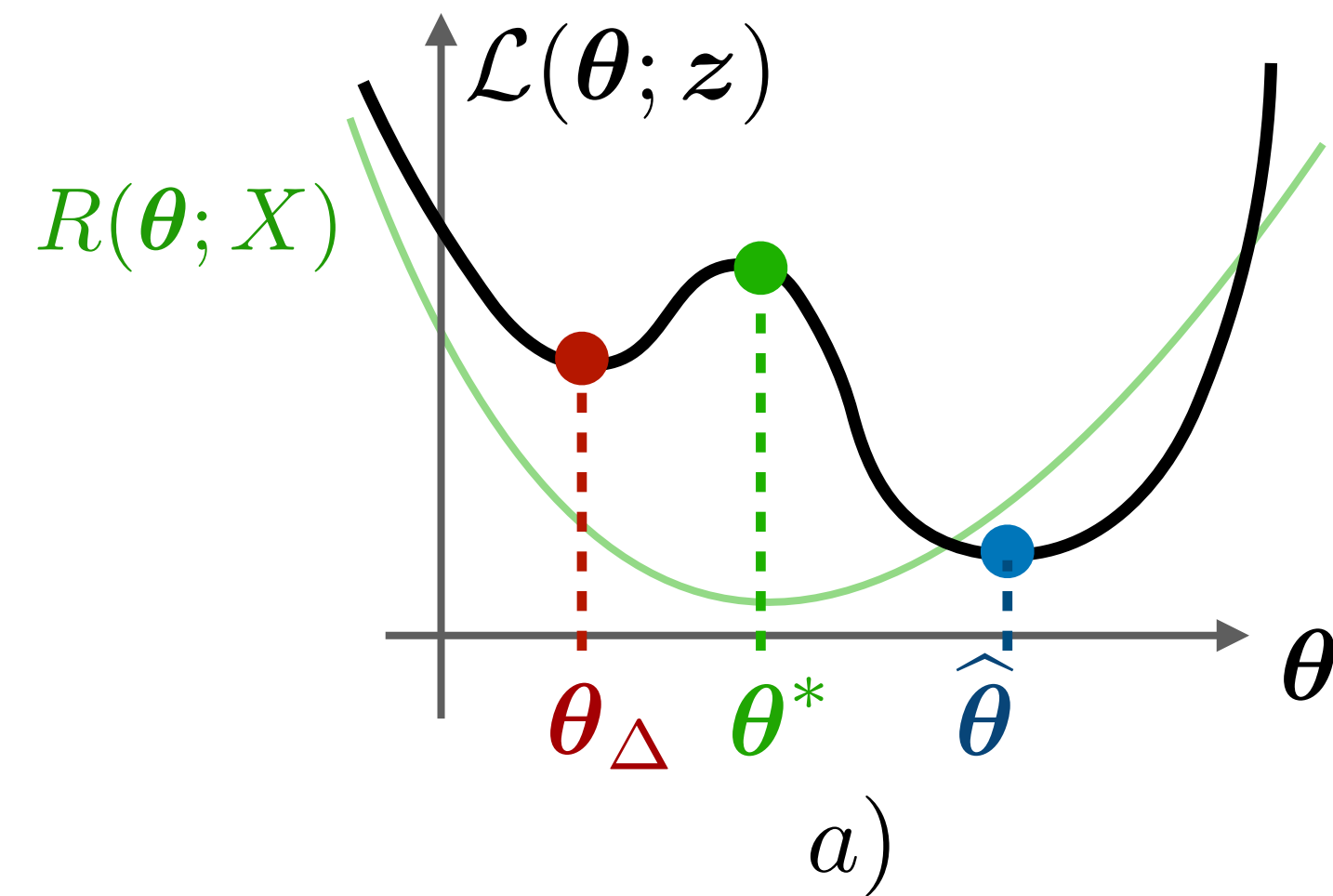
Methodology

We define 5 “failure scenarii”

$$\theta^* = \arg \min_{\theta} R(\theta; X)$$

$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}(\theta; z)$$

$$\theta_{\Delta} = \Delta[z]$$



Incorrect global minimum
Decoder fails to find it

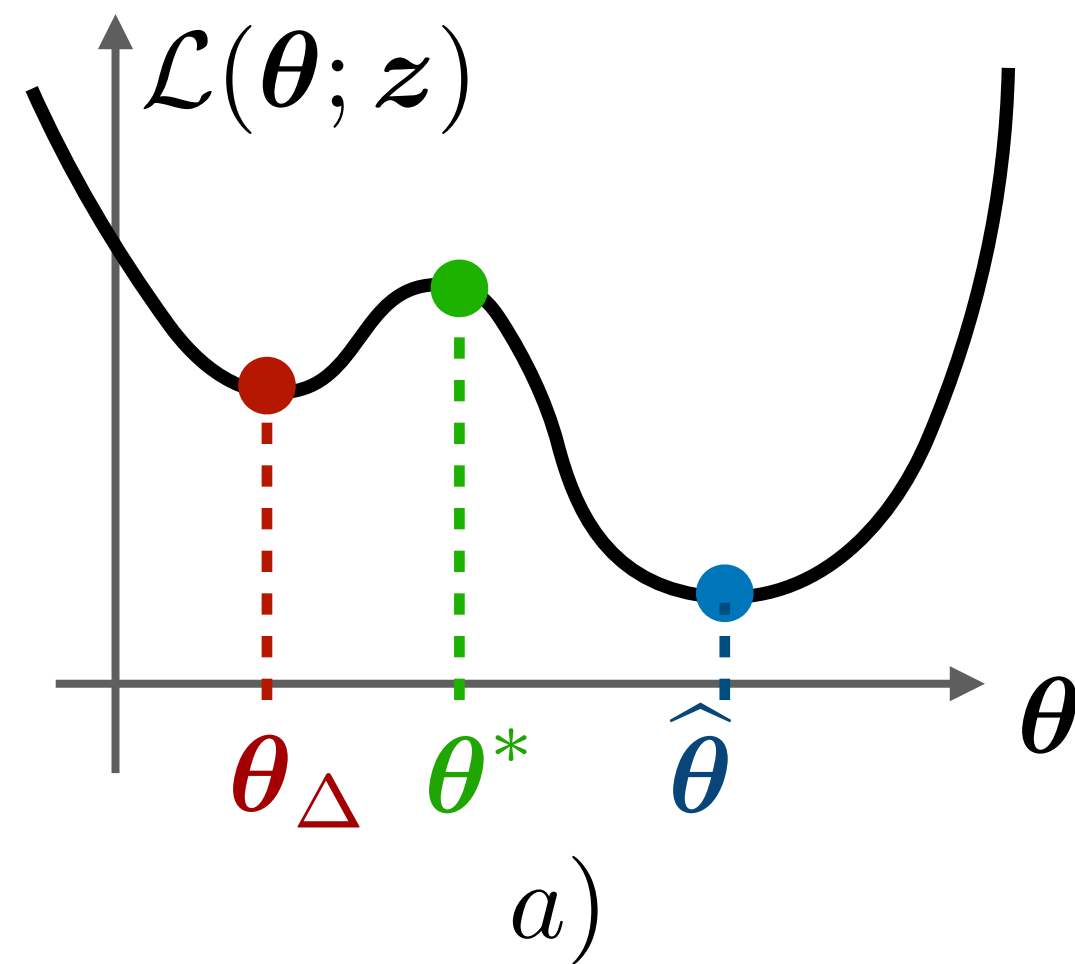
Methodology

We define 5 “failure scenarii”

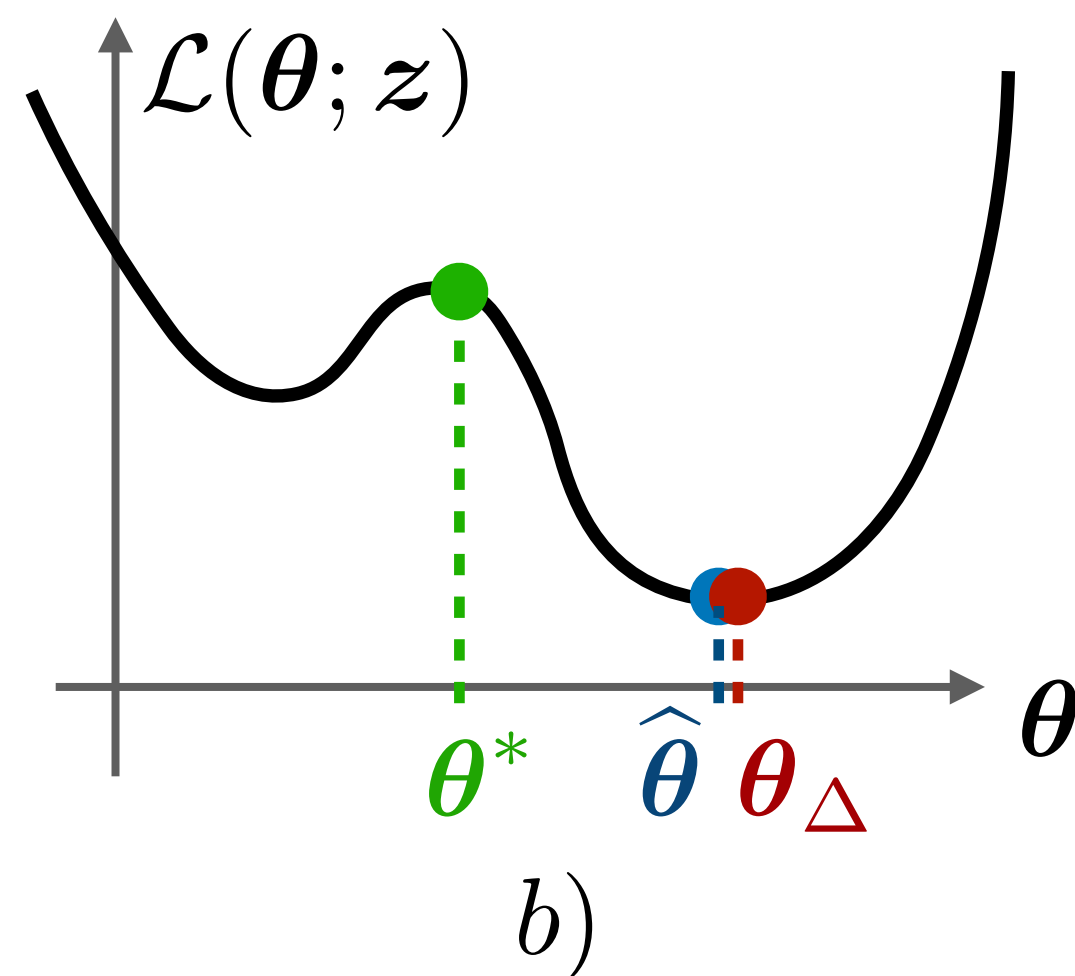
$$\theta^* = \arg \min_{\theta} R(\theta; X)$$

$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}(\theta; z)$$

$$\theta_{\Delta} = \Delta[z]$$



Incorrect global minimum
Decoder fails to find it



Incorrect global minimum
Decoder finds it

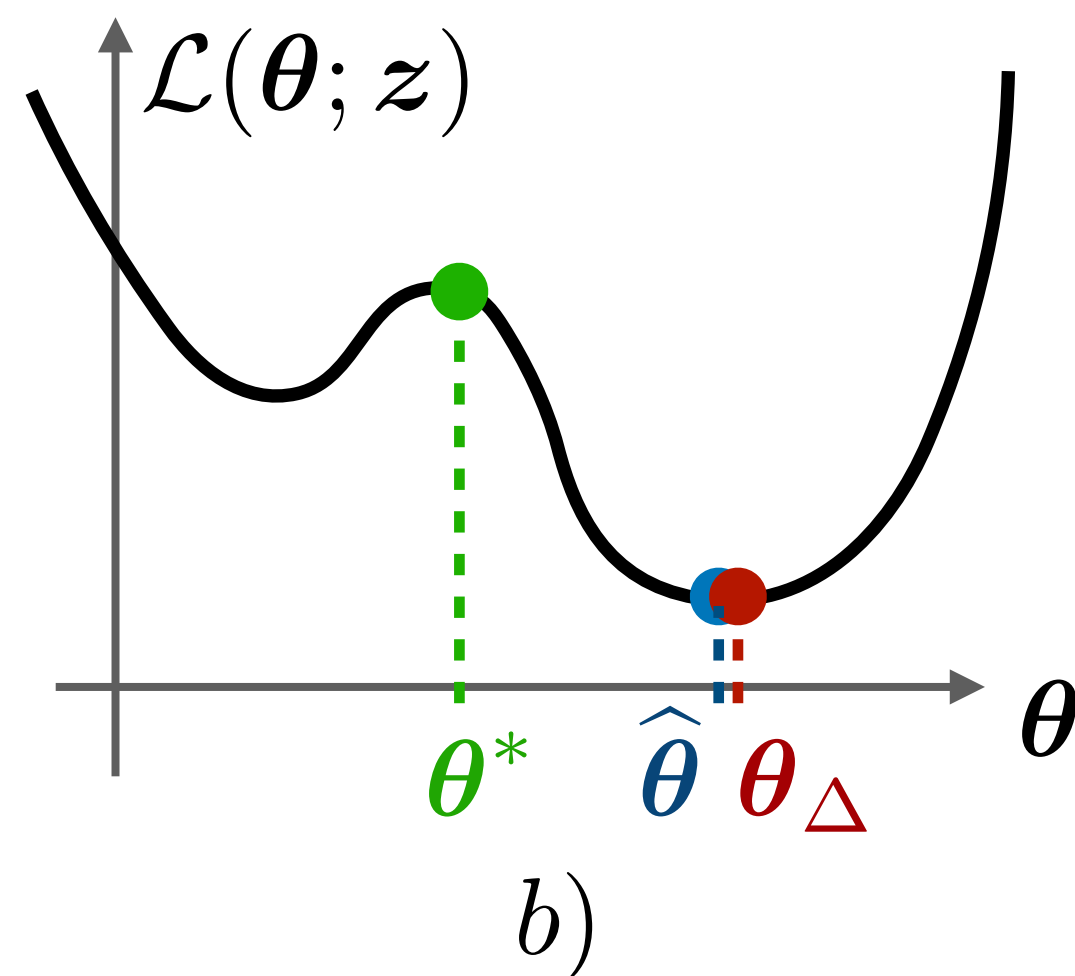
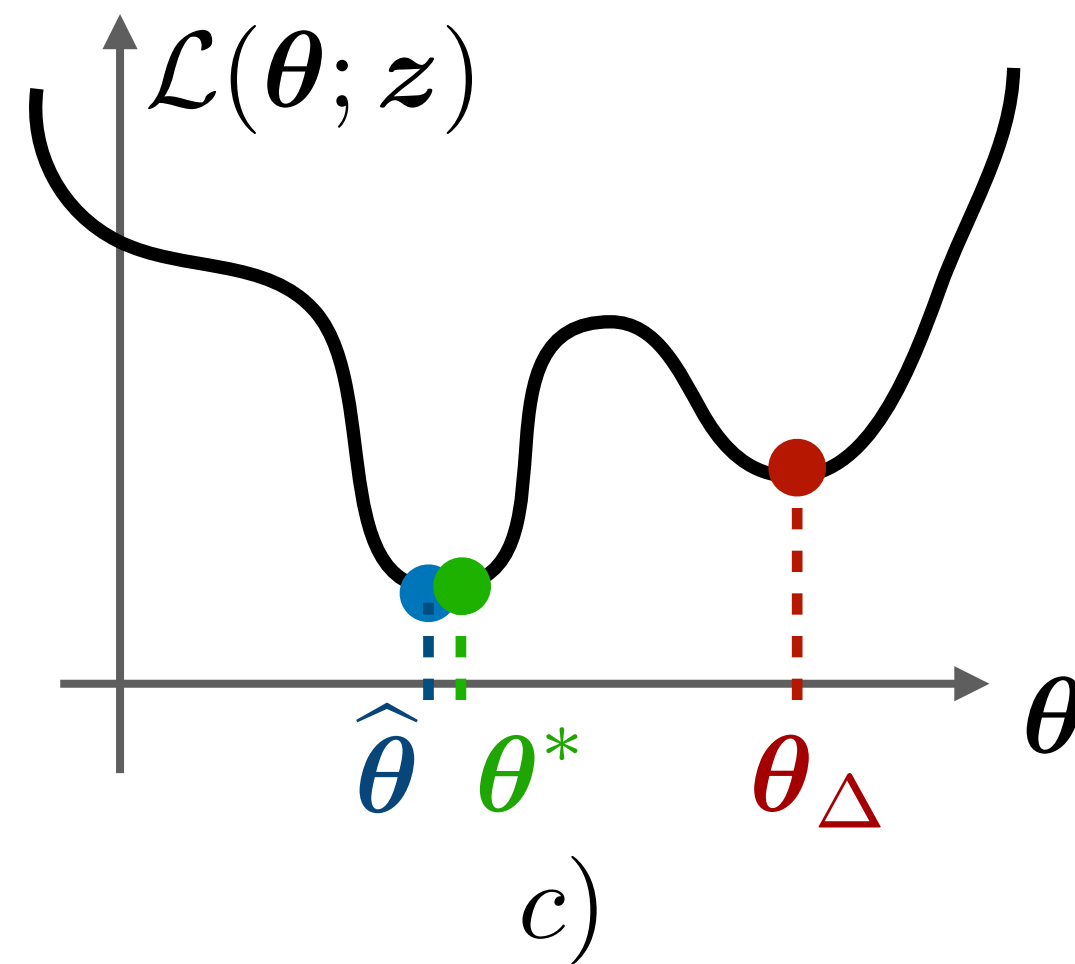
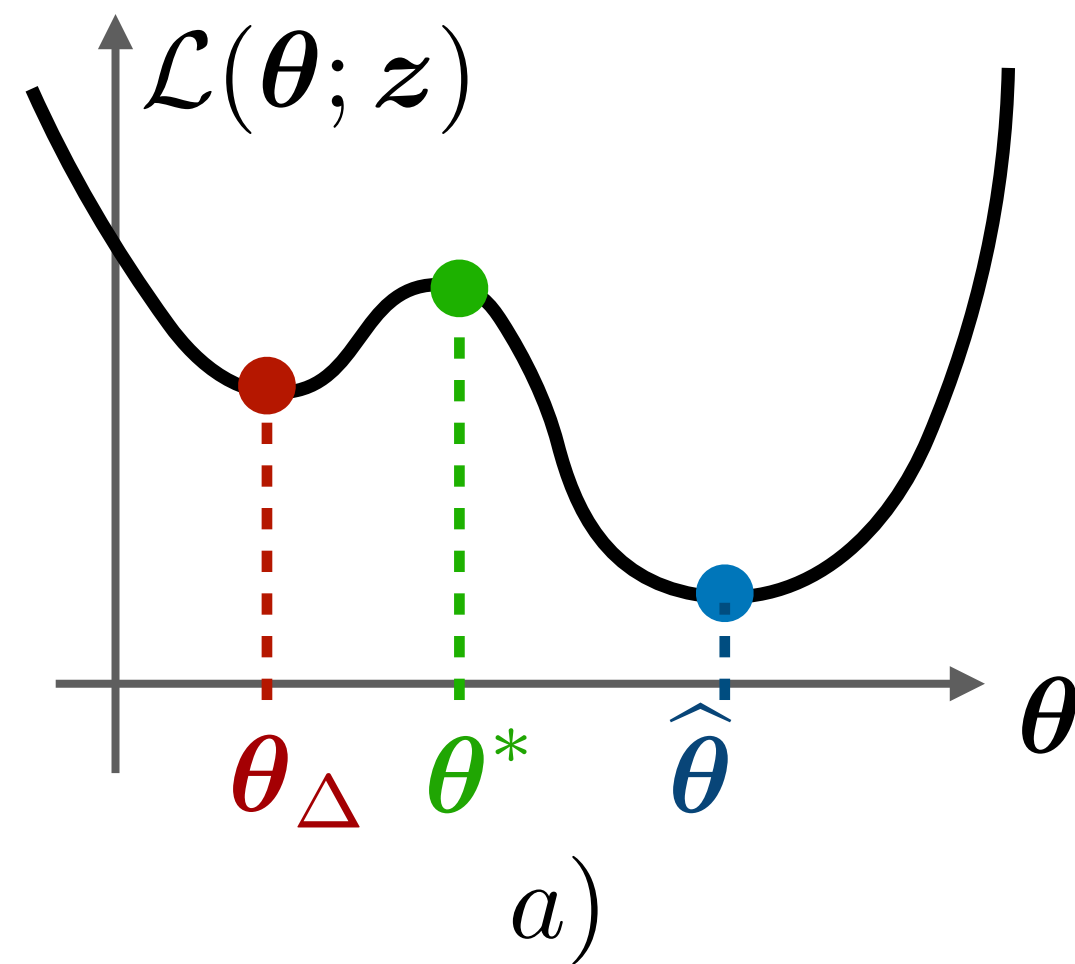
Methodology

We define 5 “failure scenarii”

$$\theta^* = \arg \min_{\theta} R(\theta; X)$$

$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}(\theta; z)$$

$$\theta_{\Delta} = \Delta[z]$$



Correct global minimum
Decoder fails to find it

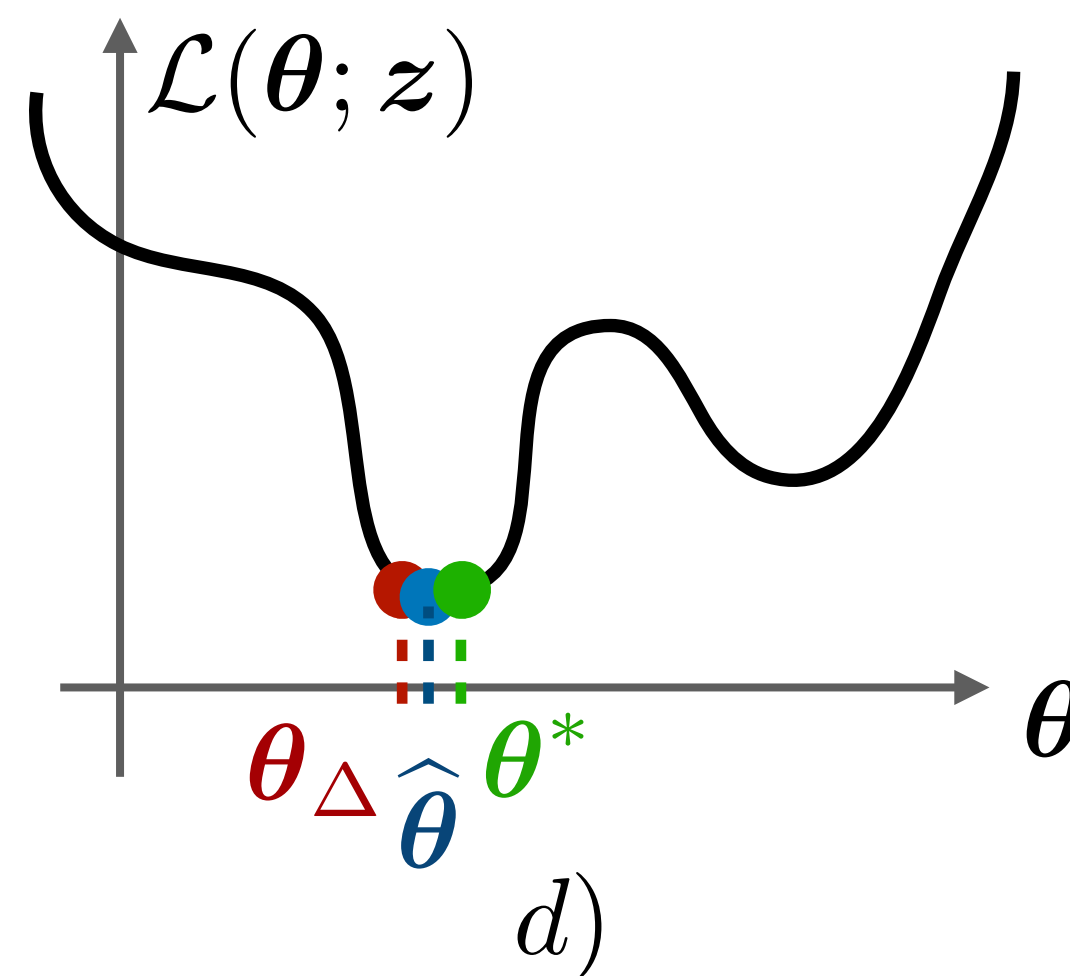
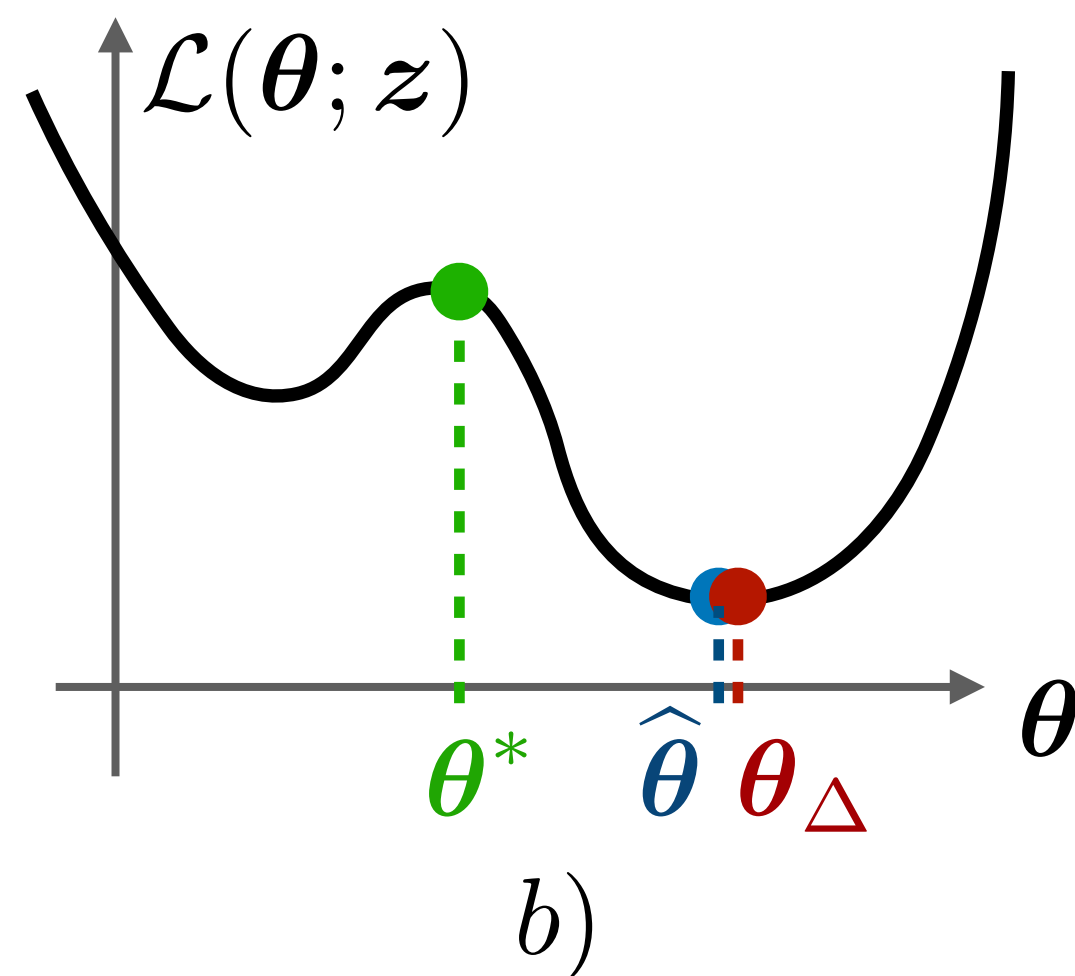
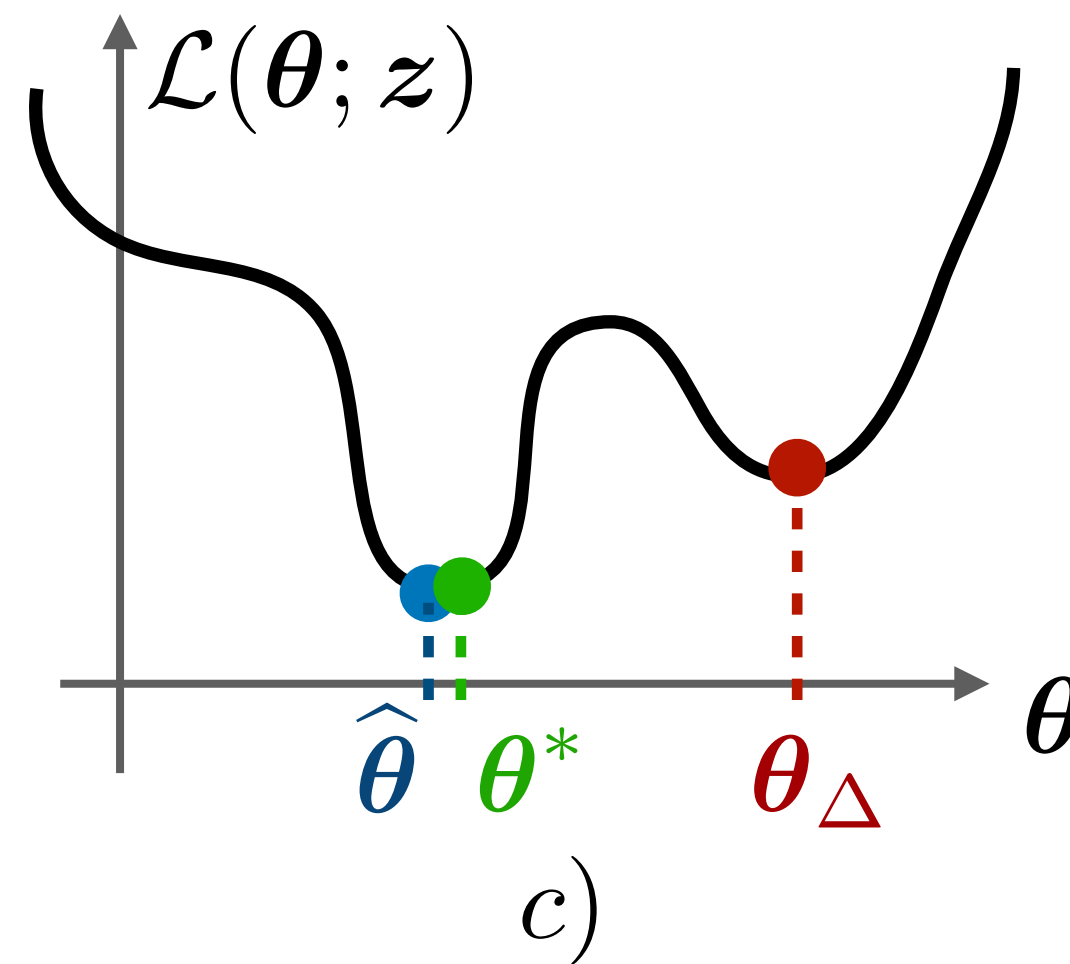
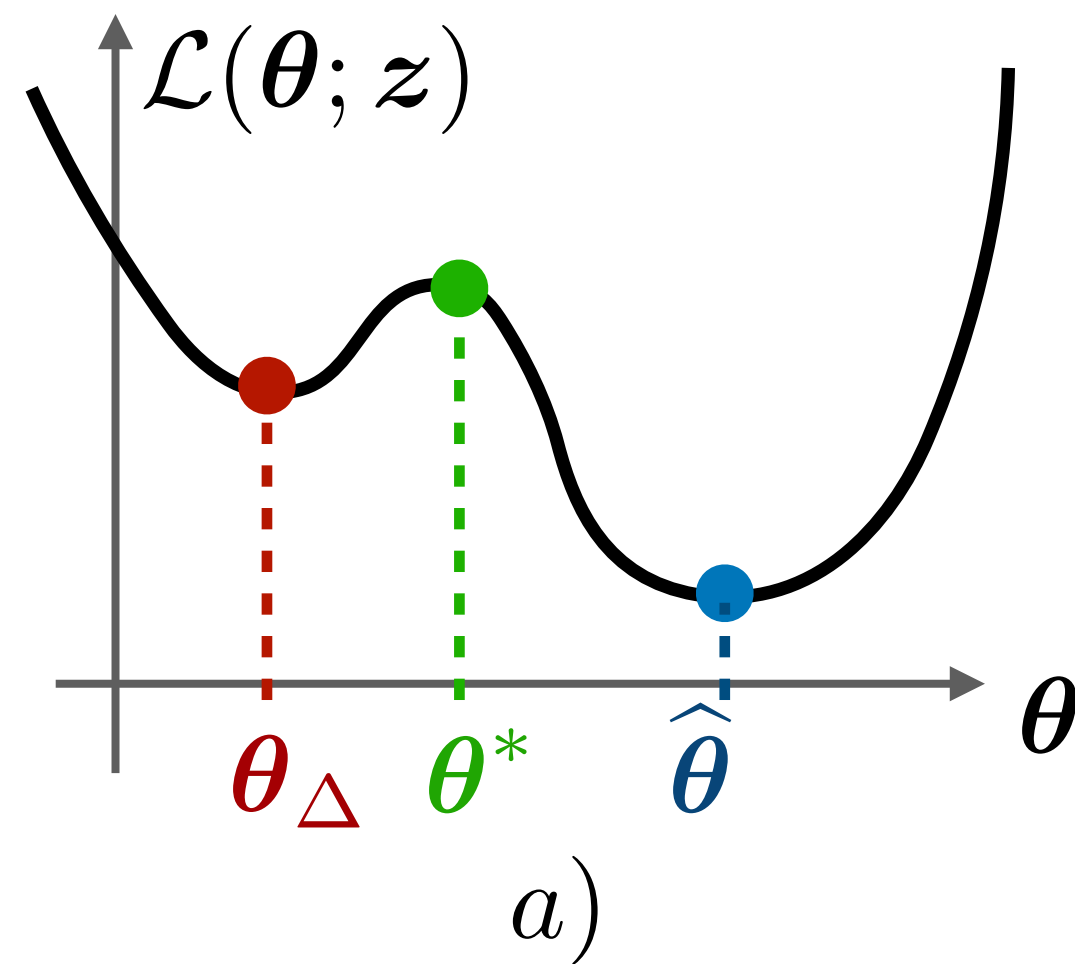
Methodology

We define 5 “failure scenarii”

$$\theta^* = \arg \min_{\theta} R(\theta; X)$$

$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}(\theta; z)$$

$$\theta_{\Delta} = \Delta[z]$$



Correct global minimum
Decoder fails to find it

Correct global minimum
Decoder finds it

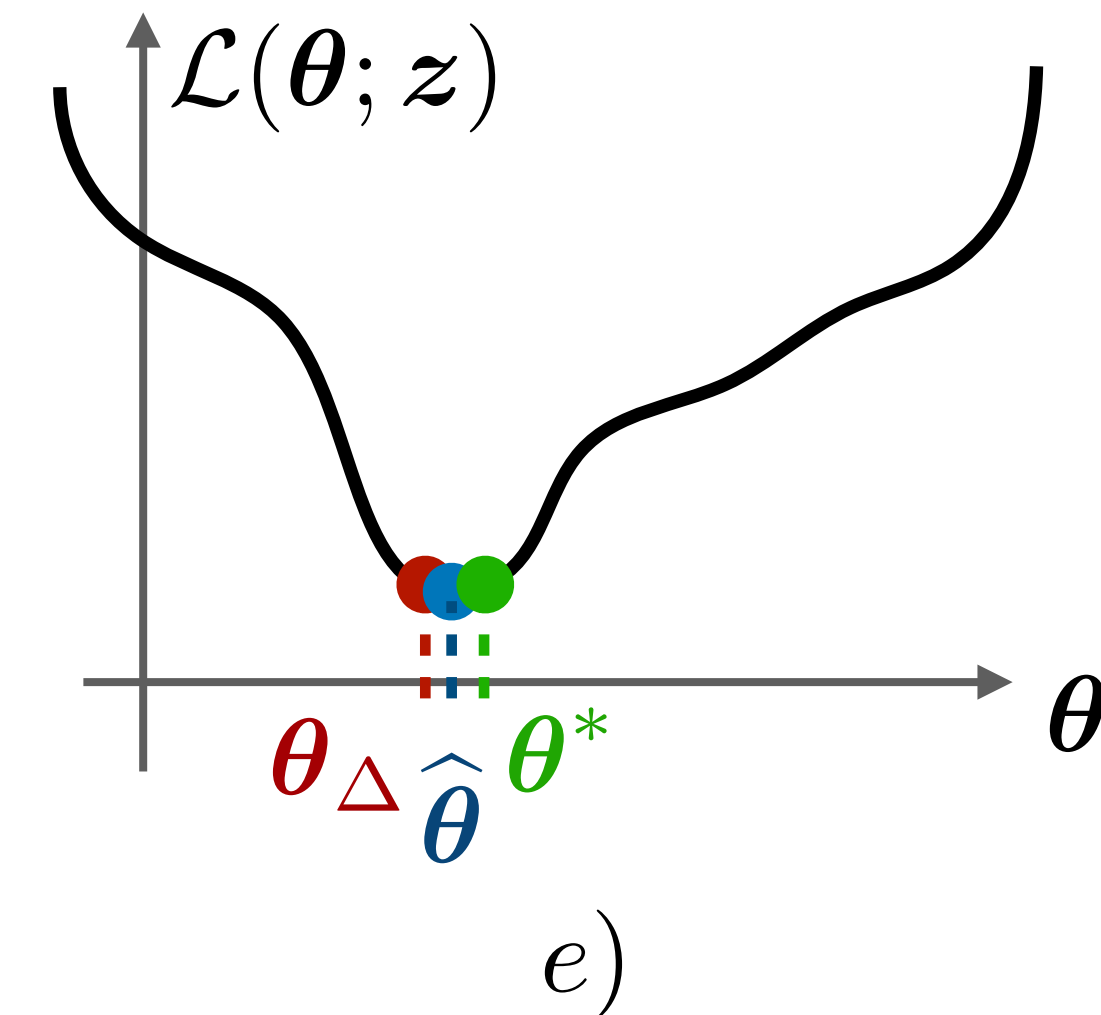
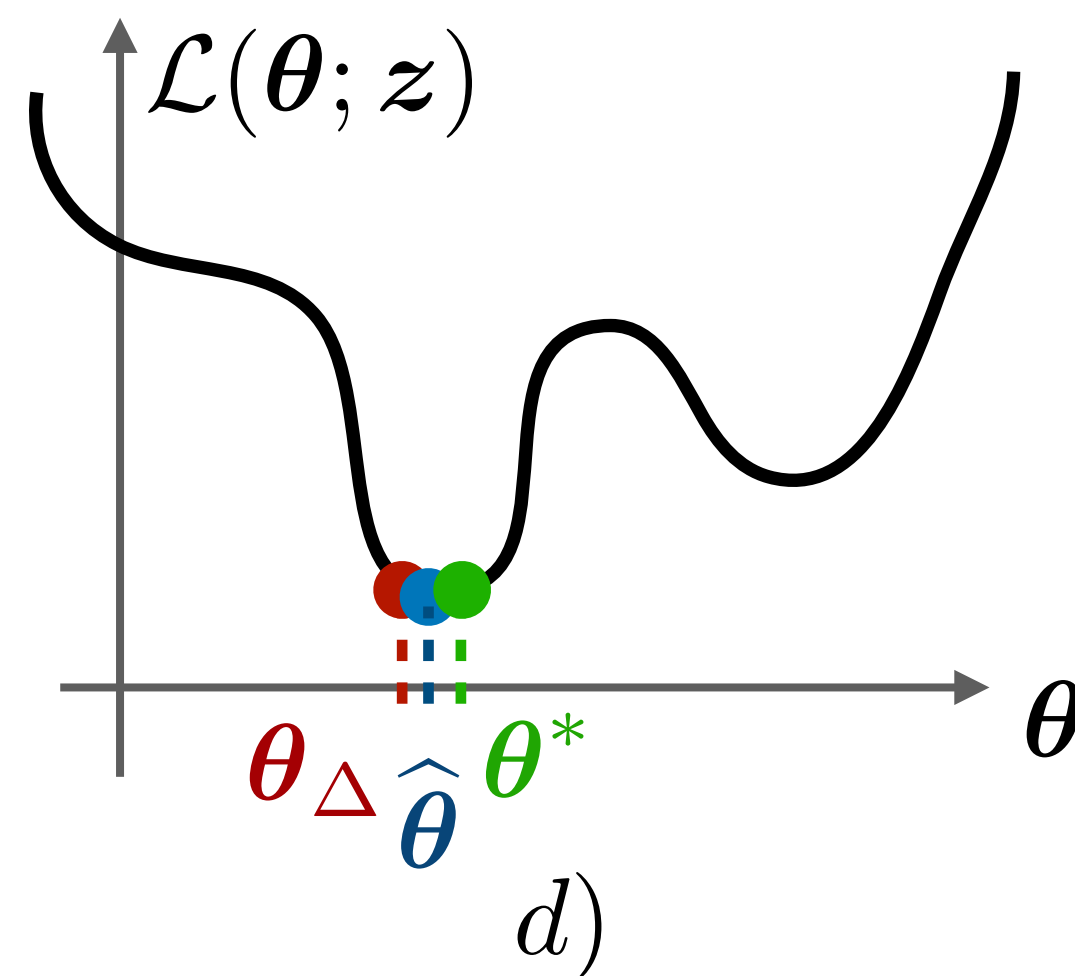
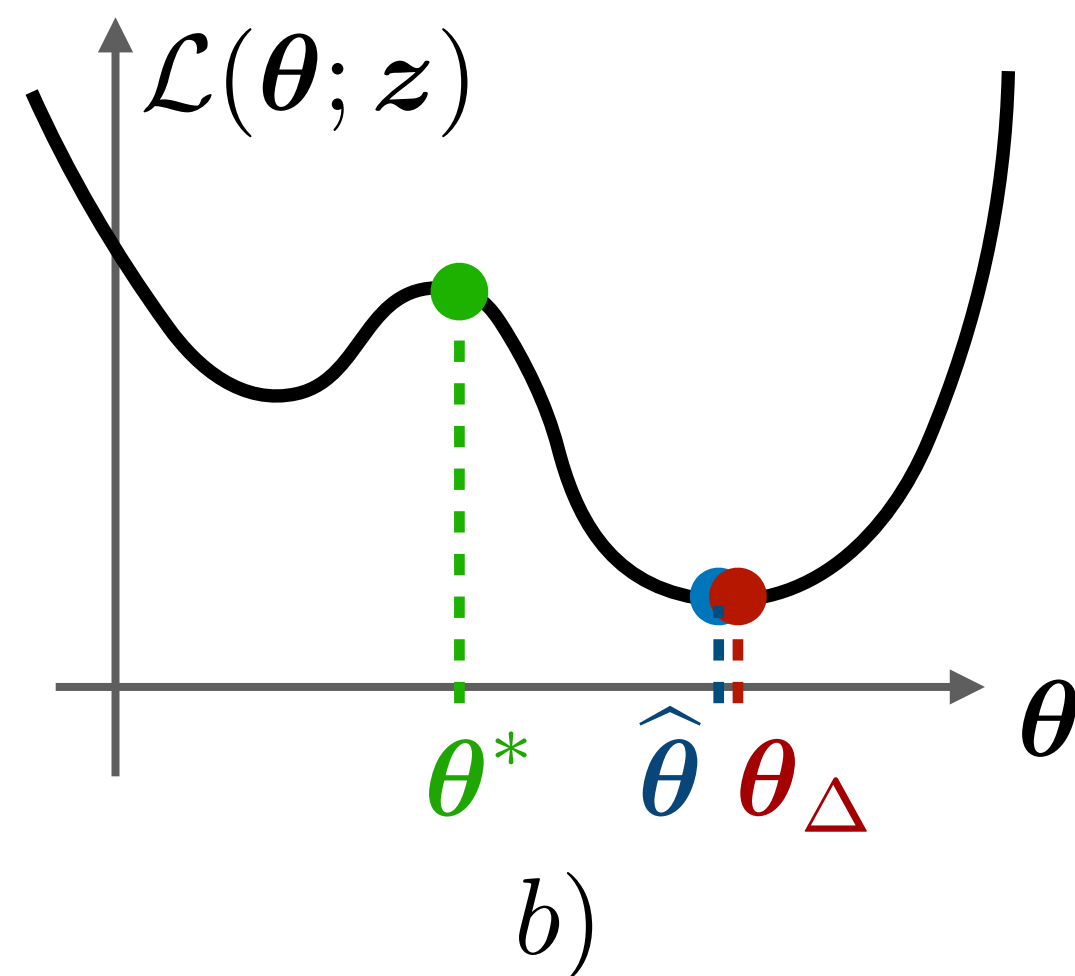
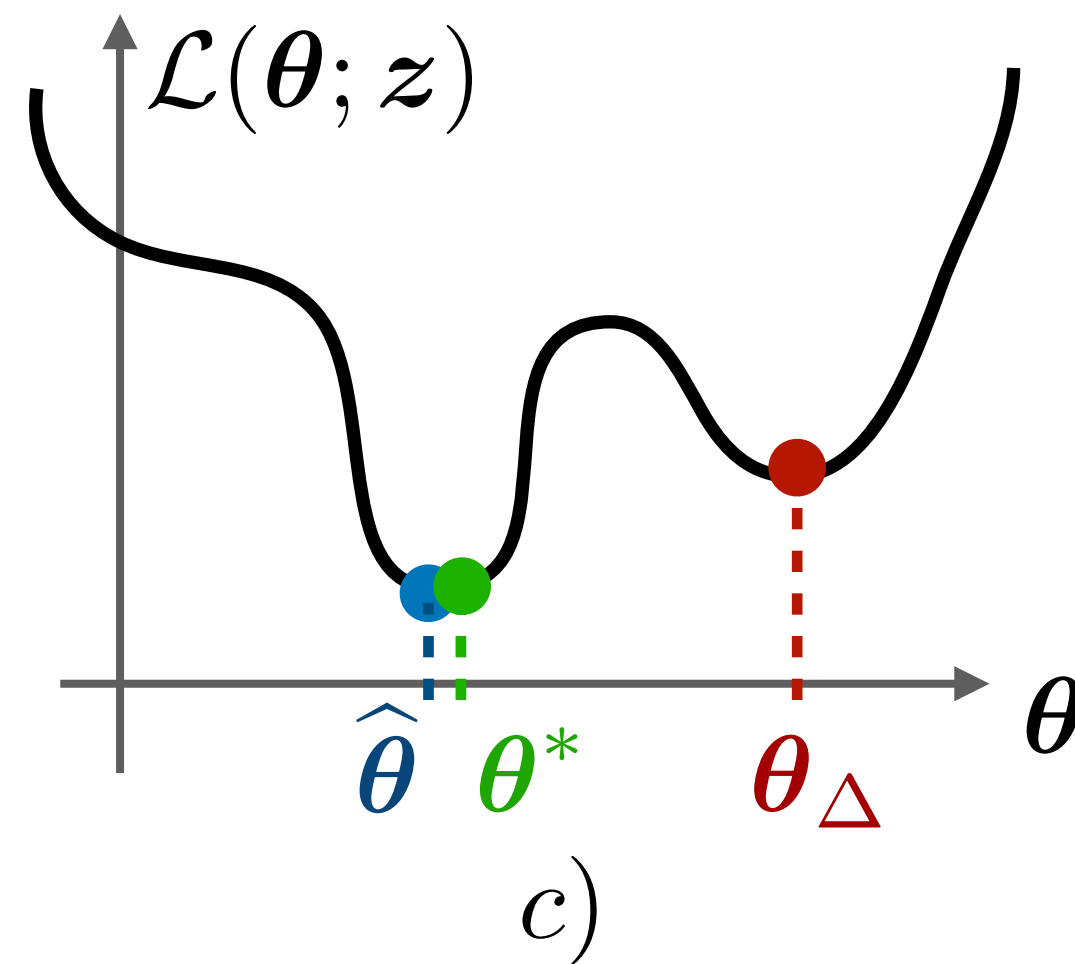
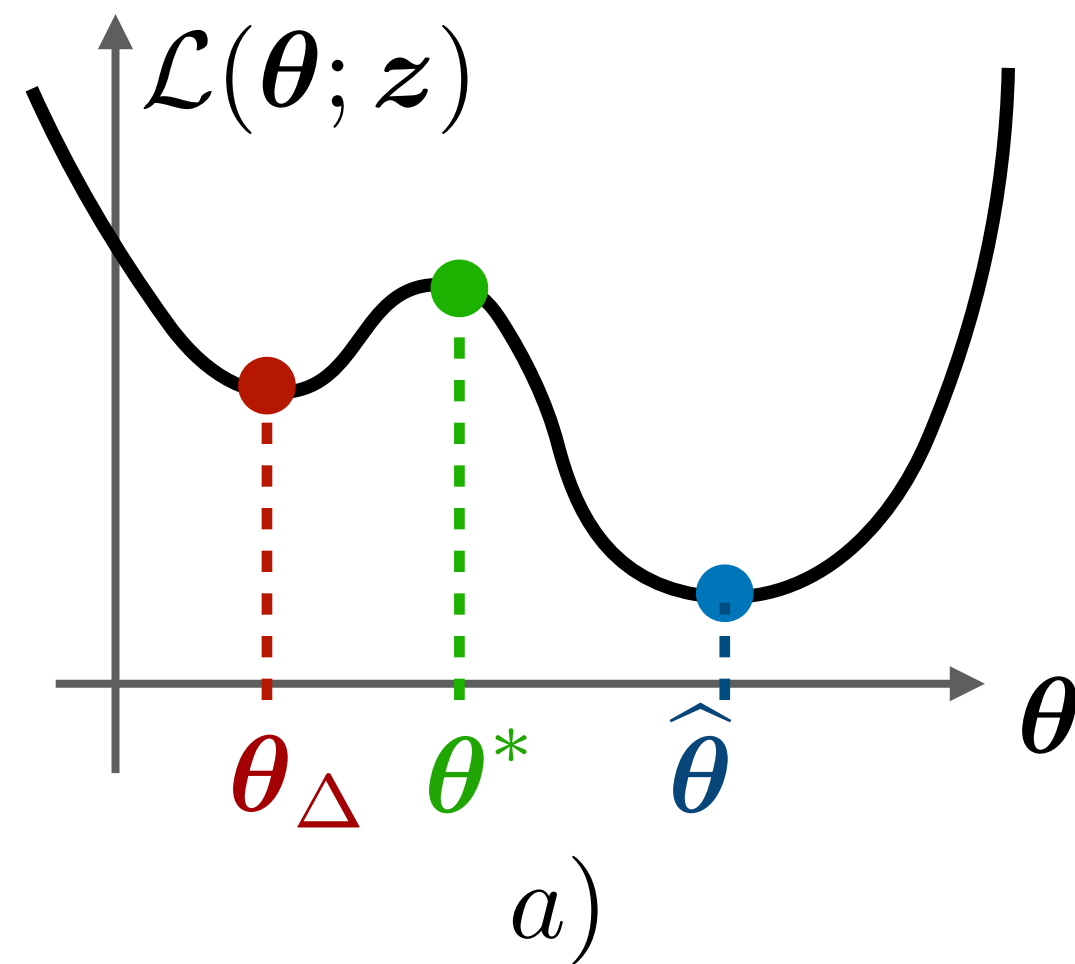
Methodology

We define 5 “failure scenarii”

$$\theta^* = \arg \min_{\theta} R(\theta; X)$$

$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}(\theta; z)$$

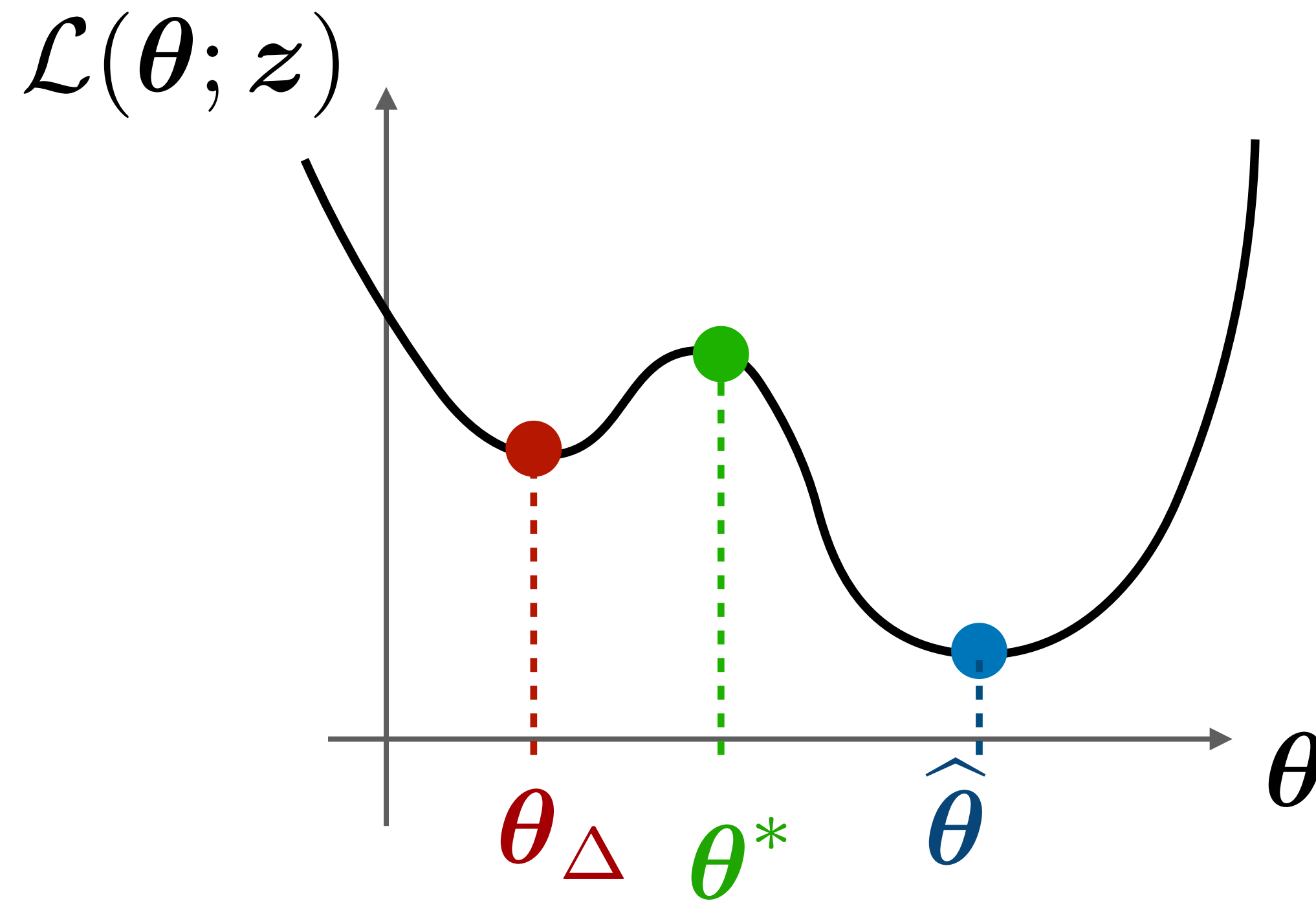
$$\theta_{\Delta} = \Delta[z]$$



Correct global minimum
Large basin of attraction,
decoder always finds it

Methodology

What do we know?



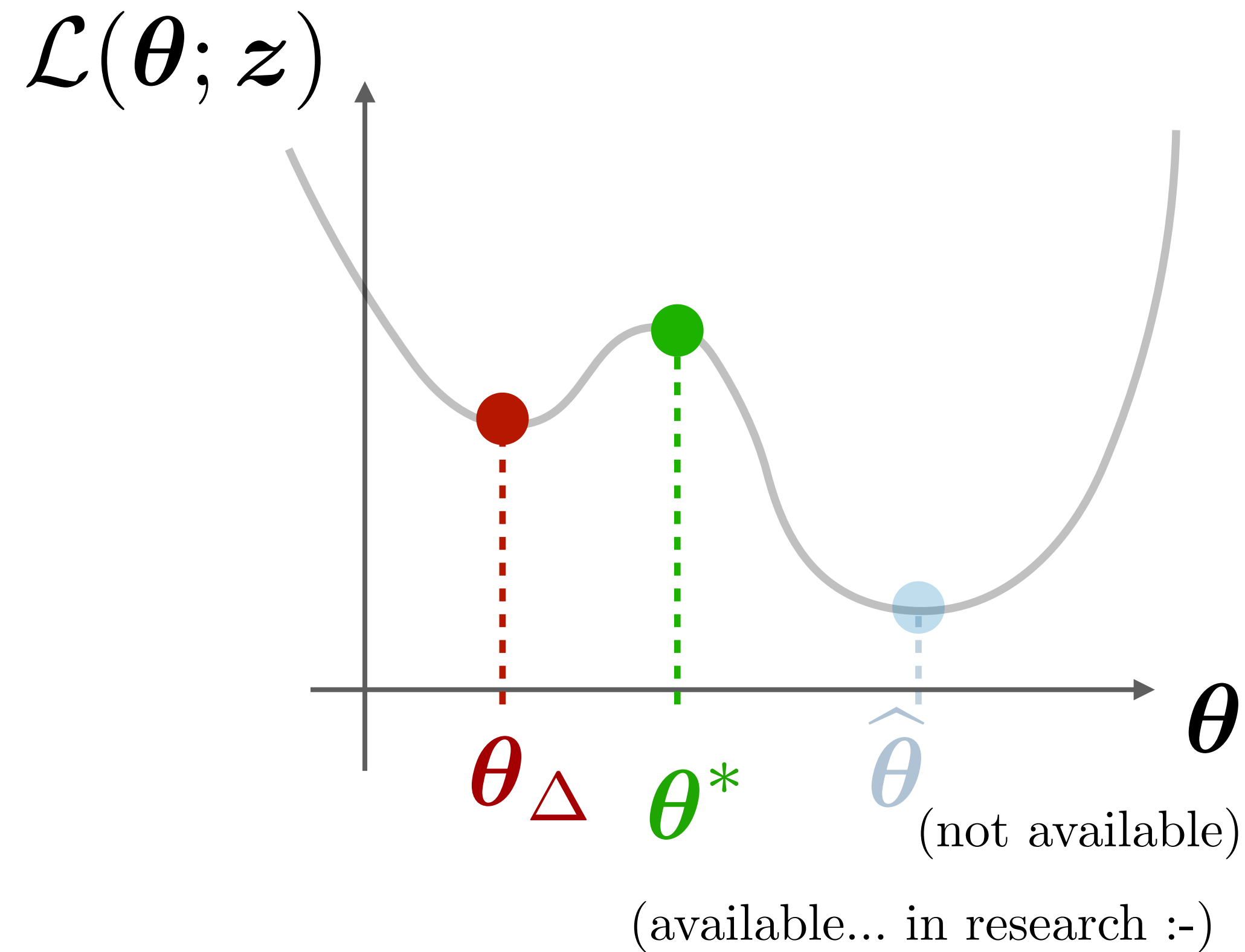
$$\theta^* = \arg \min_{\theta} R(\theta; X)$$

$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}(\theta; z)$$

$$\theta_{\Delta} = \Delta[z]$$

Methodology

What do we know?
Two points...



$$\theta^* = \arg \min_{\theta} R(\theta; X)$$

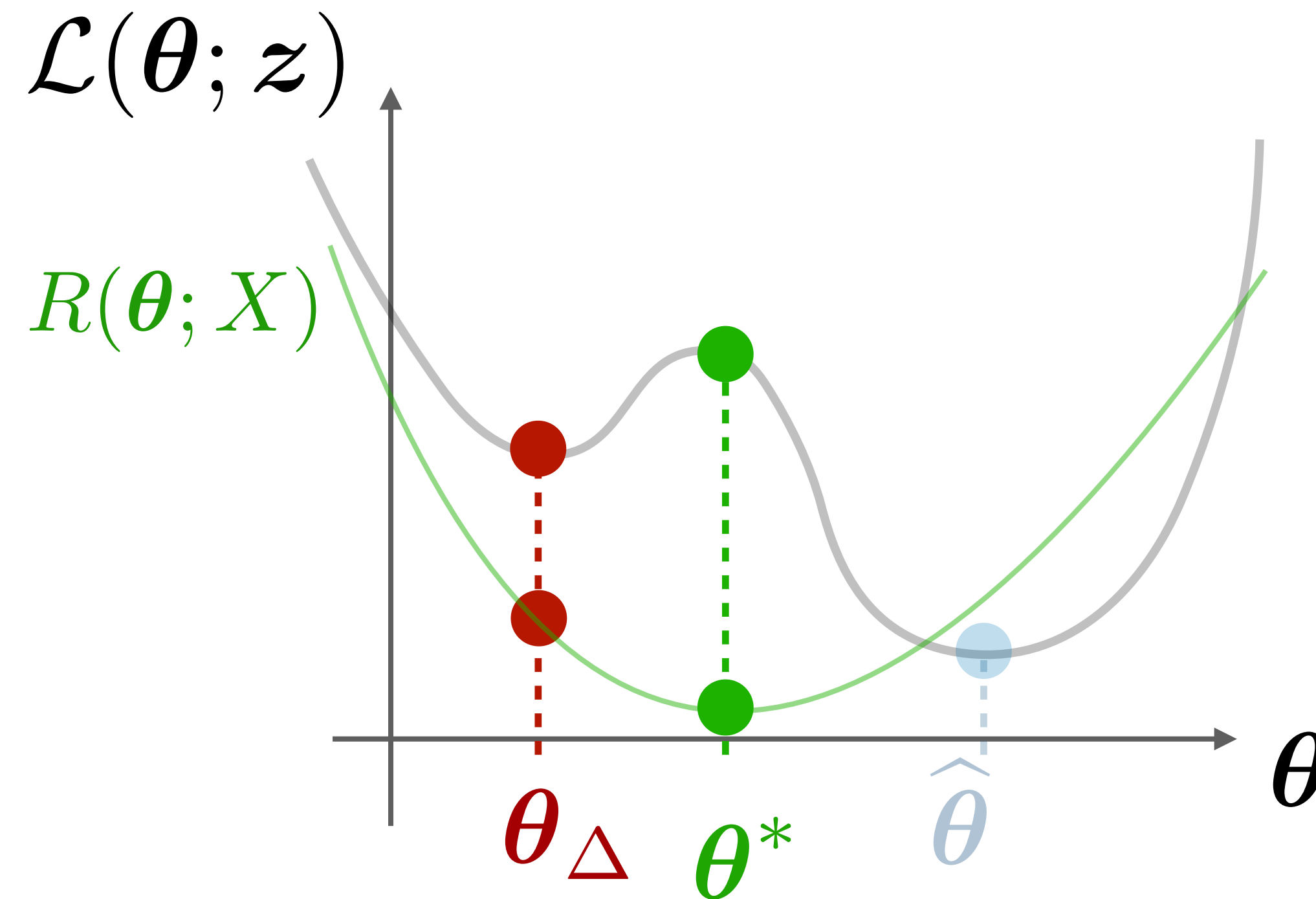
$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}(\theta; z)$$

$$\theta_{\Delta} = \Delta[z]$$

Methodology

What do we know?

Two points, two values!



$$\theta^* = \arg \min_{\theta} R(\theta; X)$$

$$\hat{\theta} = \arg \min_{\theta} \mathcal{L}(\theta; z)$$

$$\theta_\Delta = \Delta[z]$$

Methodology

What do we know?

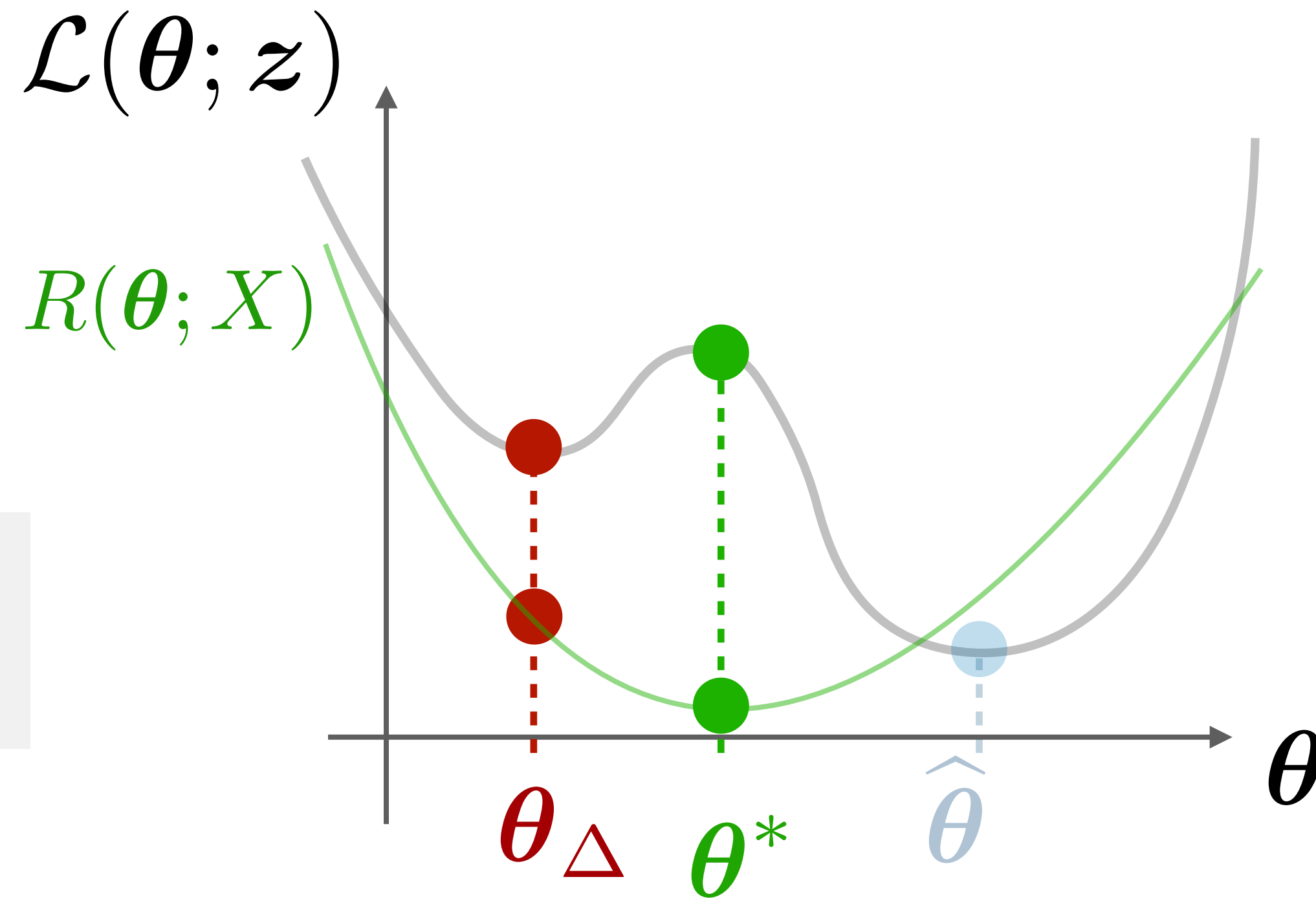
Two points, two values!

A simple test:

$$\mathcal{L}(\boldsymbol{\theta}_{\Delta}; \boldsymbol{z}) \stackrel{?}{>} \mathcal{L}(\boldsymbol{\theta}^*; \boldsymbol{z})$$

Yes

Decoder failed
(Sketch ?!)



$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} R(\boldsymbol{\theta}; X)$$

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}; \boldsymbol{z})$$

$$\boldsymbol{\theta}_{\Delta} = \Delta[\boldsymbol{z}]$$

Methodology

What do we know?

Two points, two values!

A simple test:

$$\mathcal{L}(\boldsymbol{\theta}_\Delta; \mathbf{z}) \stackrel{?}{>} \mathcal{L}(\boldsymbol{\theta}^*; \mathbf{z})$$

Yes

No

Decoder failed
(Sketch ?!)

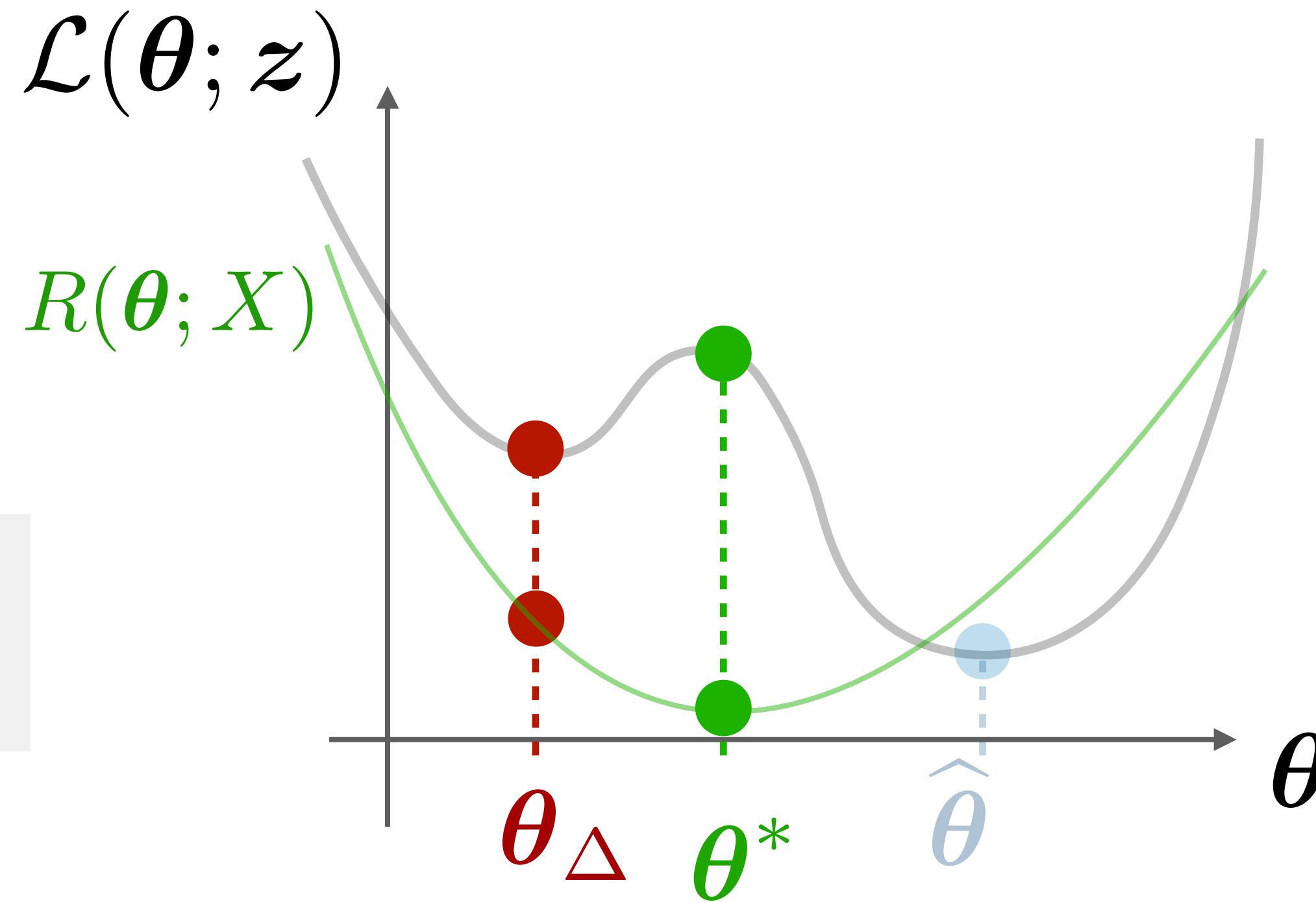
$$R(\boldsymbol{\theta}_\Delta; X) \stackrel{?}{\gg} R(\boldsymbol{\theta}^*; X)$$

Yes

No $R(\boldsymbol{\theta}_\Delta; X) \simeq R(\boldsymbol{\theta}^*; X)$

Sketch failed
(Decoder ?!)

Success :-)



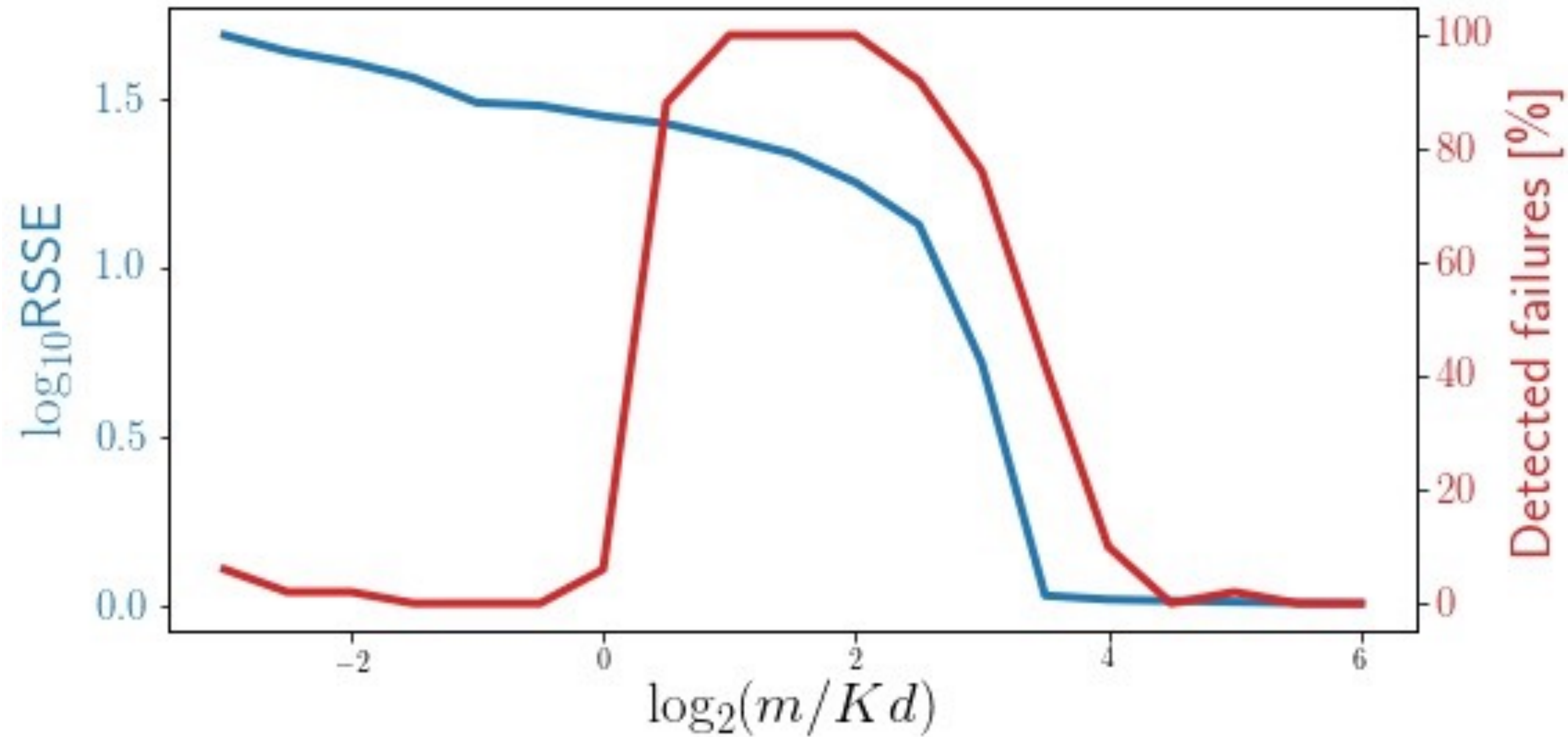
$$\boldsymbol{\theta}^* = \arg \min_{\boldsymbol{\theta}} R(\boldsymbol{\theta}; X)$$

$$\hat{\boldsymbol{\theta}} = \arg \min_{\boldsymbol{\theta}} \mathcal{L}(\boldsymbol{\theta}; \mathbf{z})$$

$$\boldsymbol{\theta}_\Delta = \Delta[\mathbf{z}]$$

Results: size of m (k-means)

Multiple sketches (50 draws)

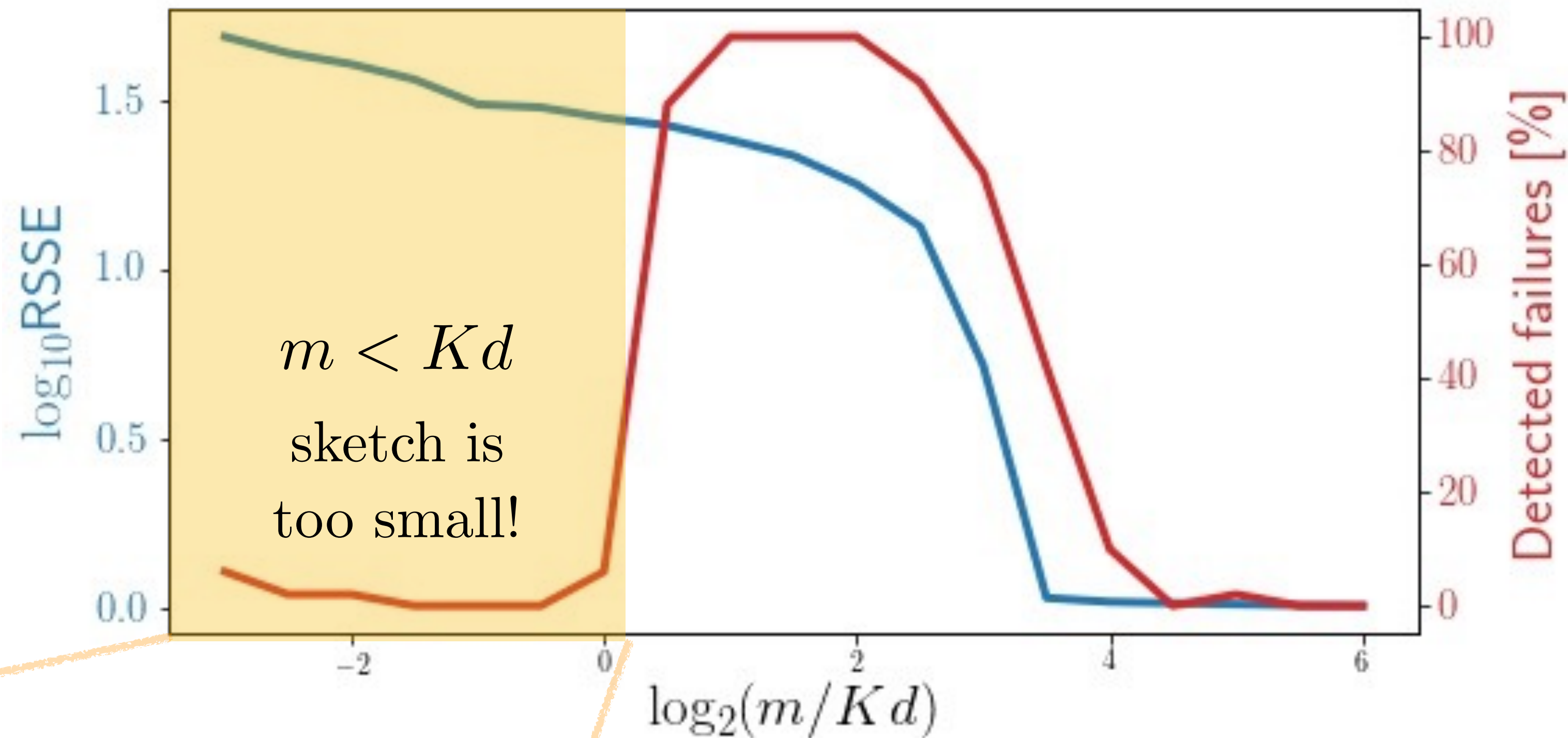


Performance
(lower is better)

Number of times
 $\mathcal{L}(\theta_{\Delta}; z) > \mathcal{L}(\theta^*; z)$

Results: size of m (k-means)

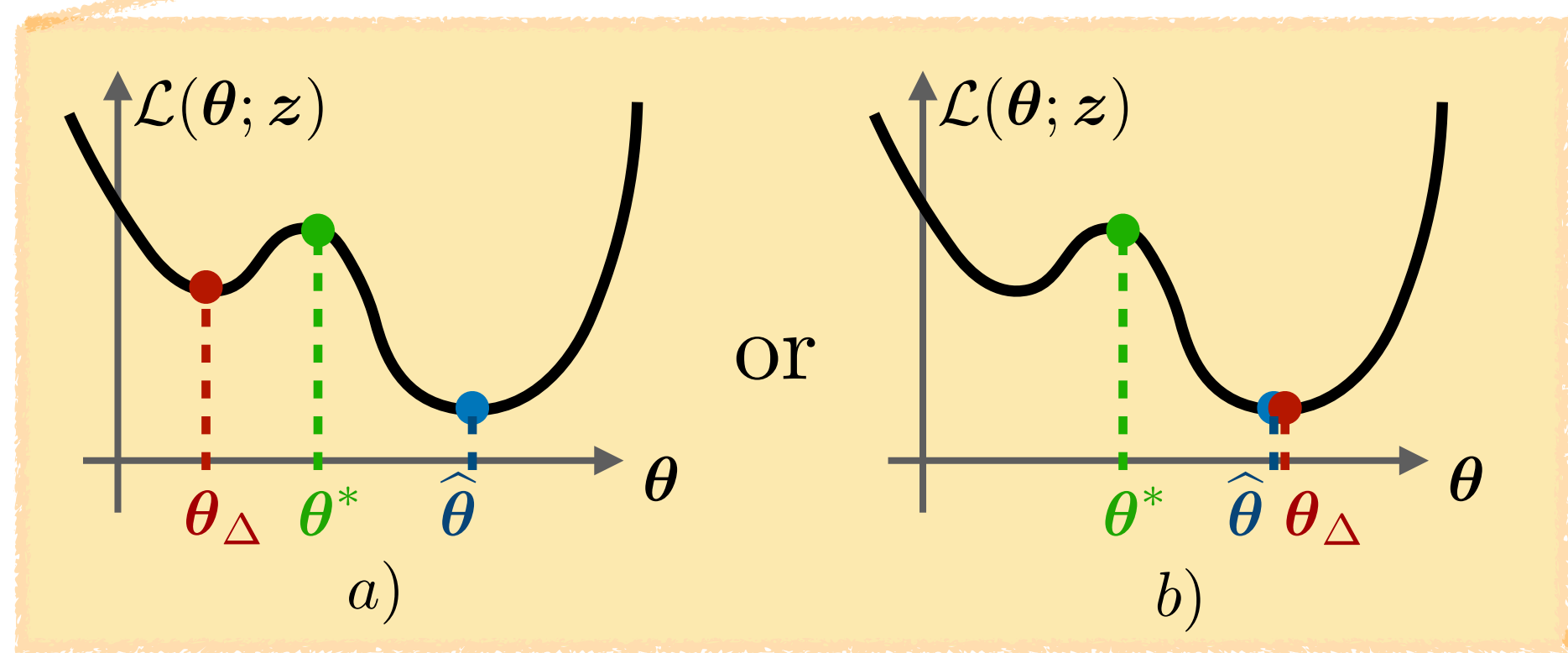
Multiple sketches (50 draws)



Performance (lower is better)

Number of times $\mathcal{L}(\theta_{\Delta}; z) > \mathcal{L}(\theta^*; z)$

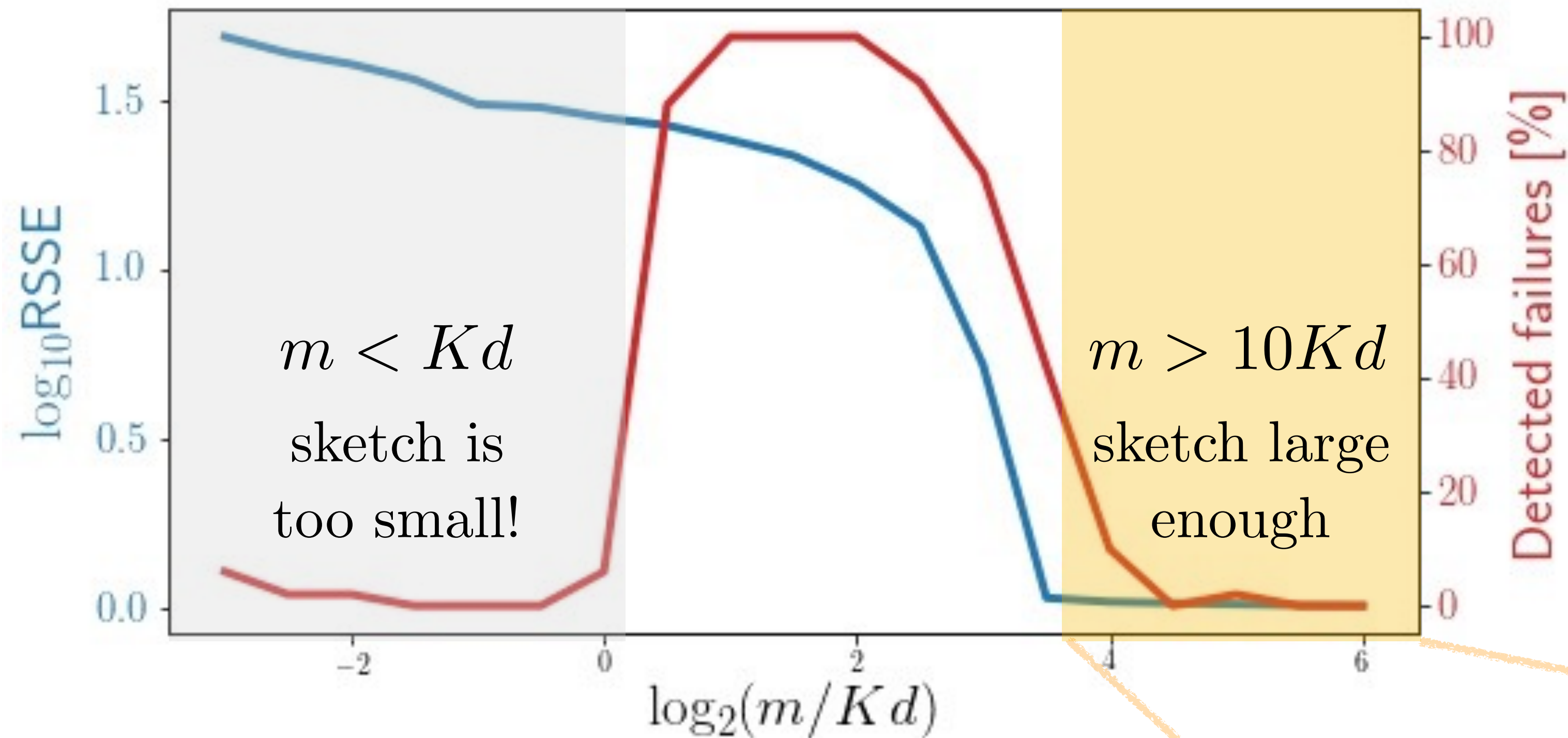
$m < Kd$
sketch is too small!



Results: size of m (k-means)

Multiple sketches (50 draws)

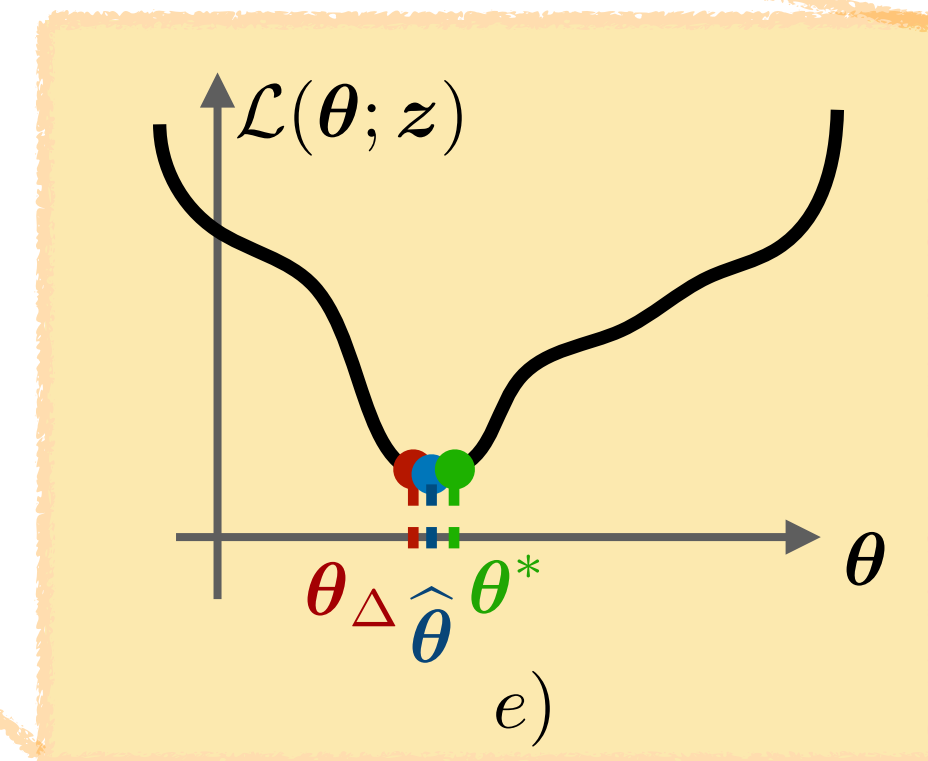
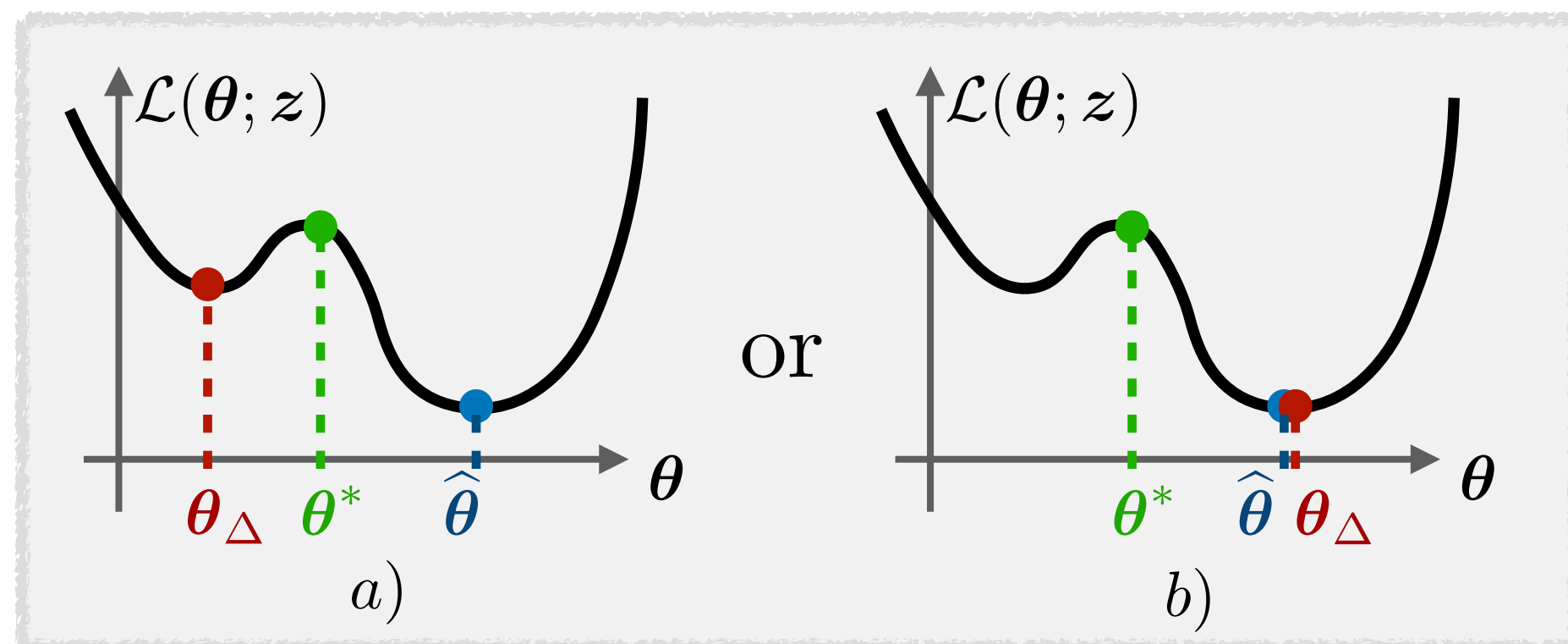
Performance
(lower is better)



$m < Kd$
sketch is
too small!

$m > 10Kd$
sketch large
enough

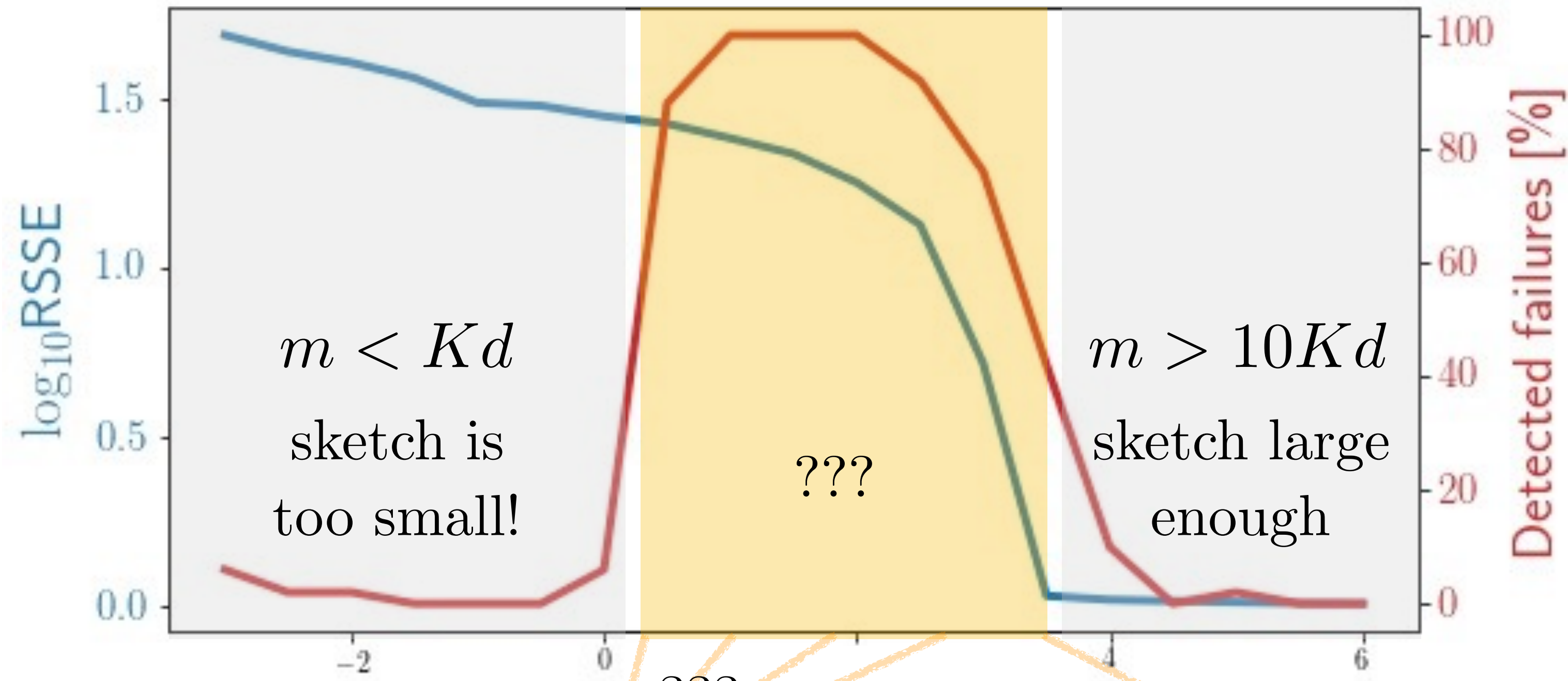
Number of times
 $\mathcal{L}(\theta_{\Delta}; z) > \mathcal{L}(\theta^*; z)$



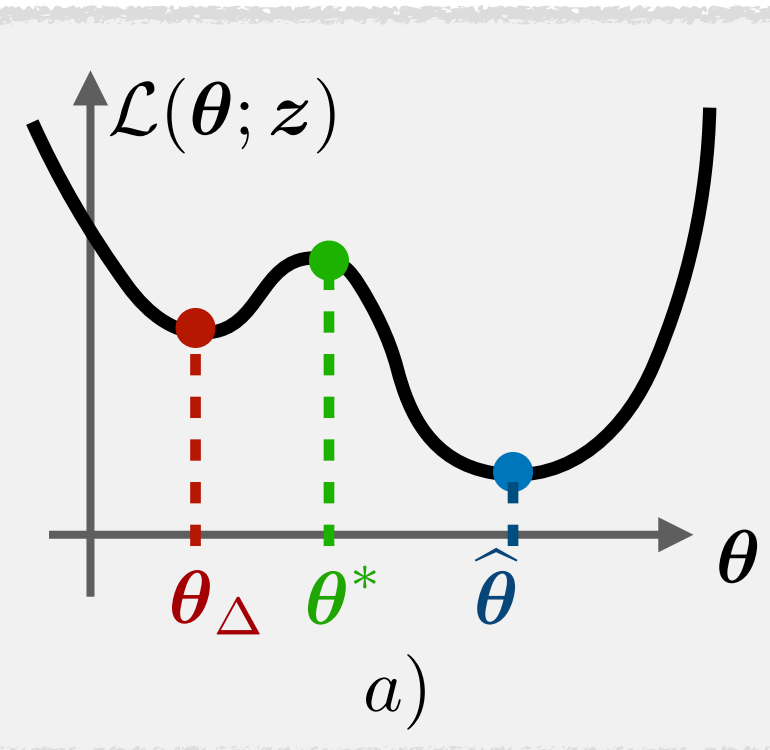
Results: size of m (k-means)

Multiple sketches (50 draws)

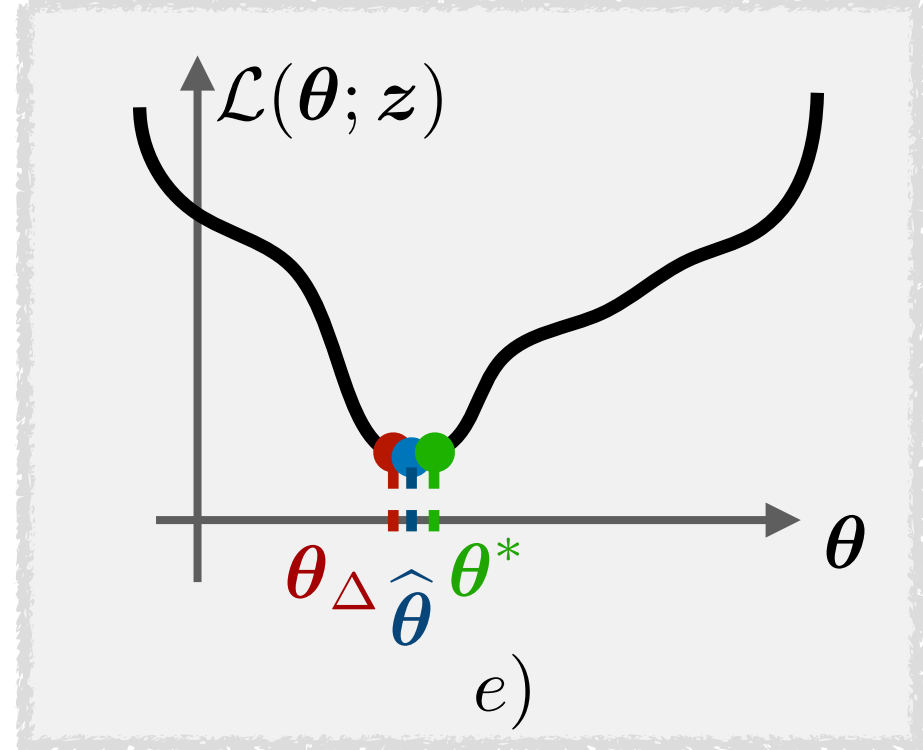
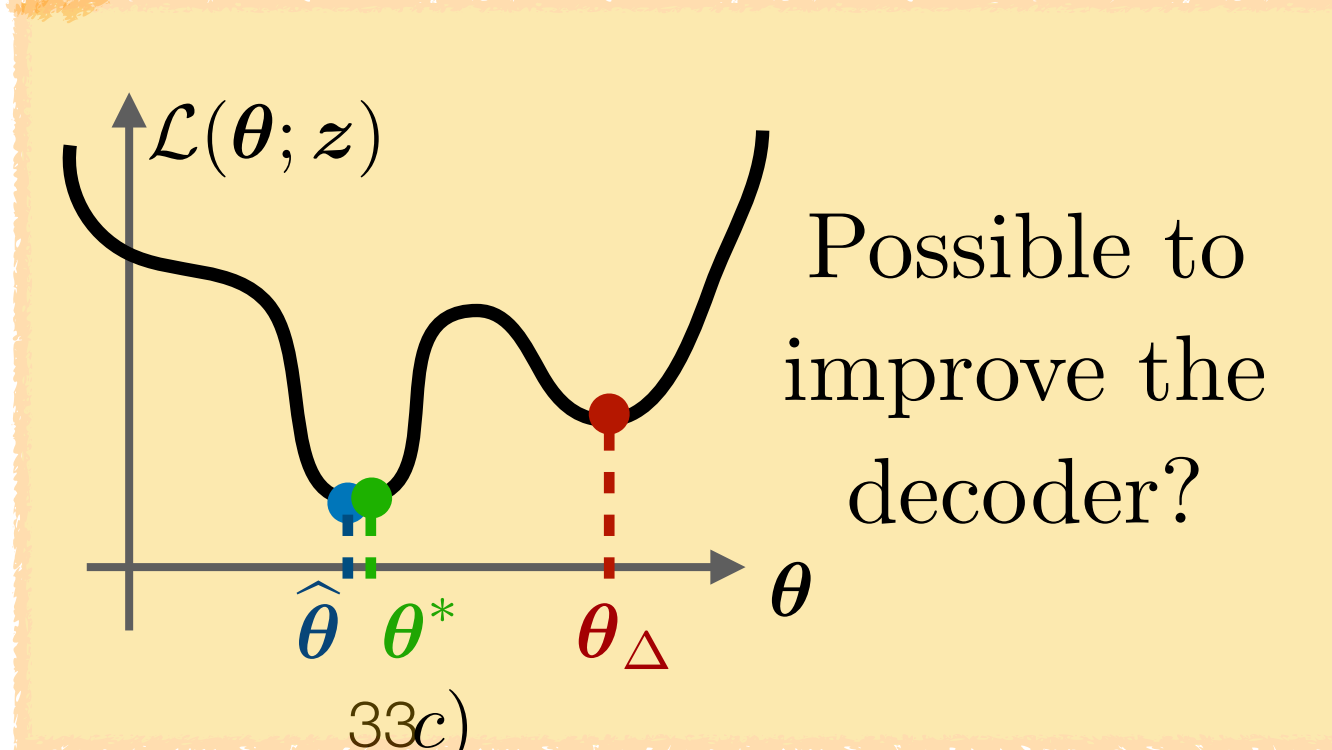
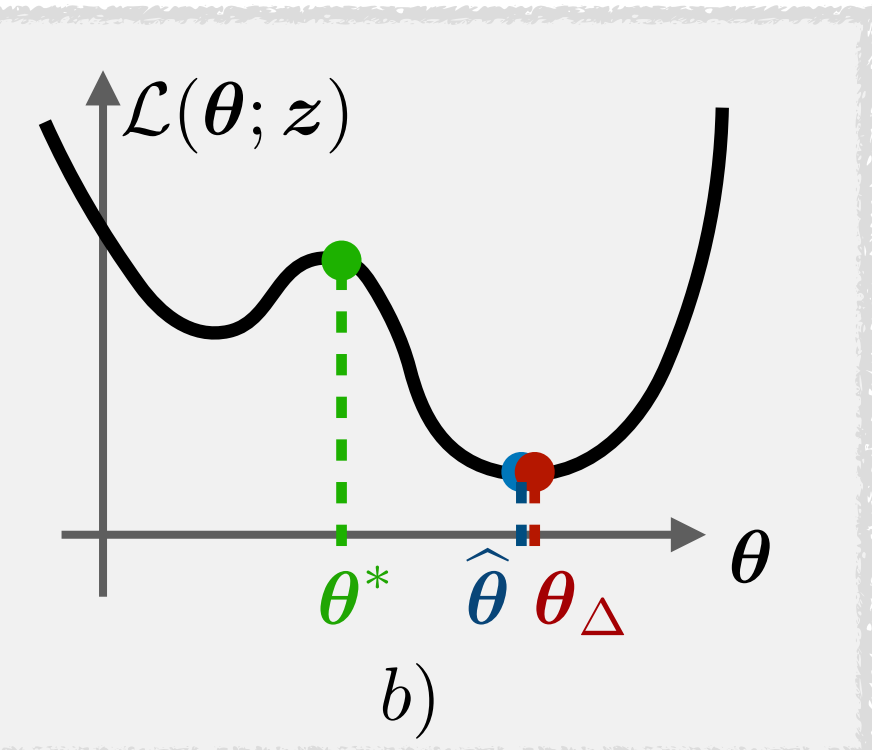
Performance
(lower is better)



Number of times
 $\mathcal{L}(\theta_{\Delta}; z) > \mathcal{L}(\theta^*; z)$



or



Results: size of m (k-means)

Can we find at least one (not necessarily efficient) decoder that succeeds where CLOMPR doesn't?

Results: size of m (k-means)

Can we find at least one (not necessarily efficient) decoder that succeeds where CLOMPR doesn't?



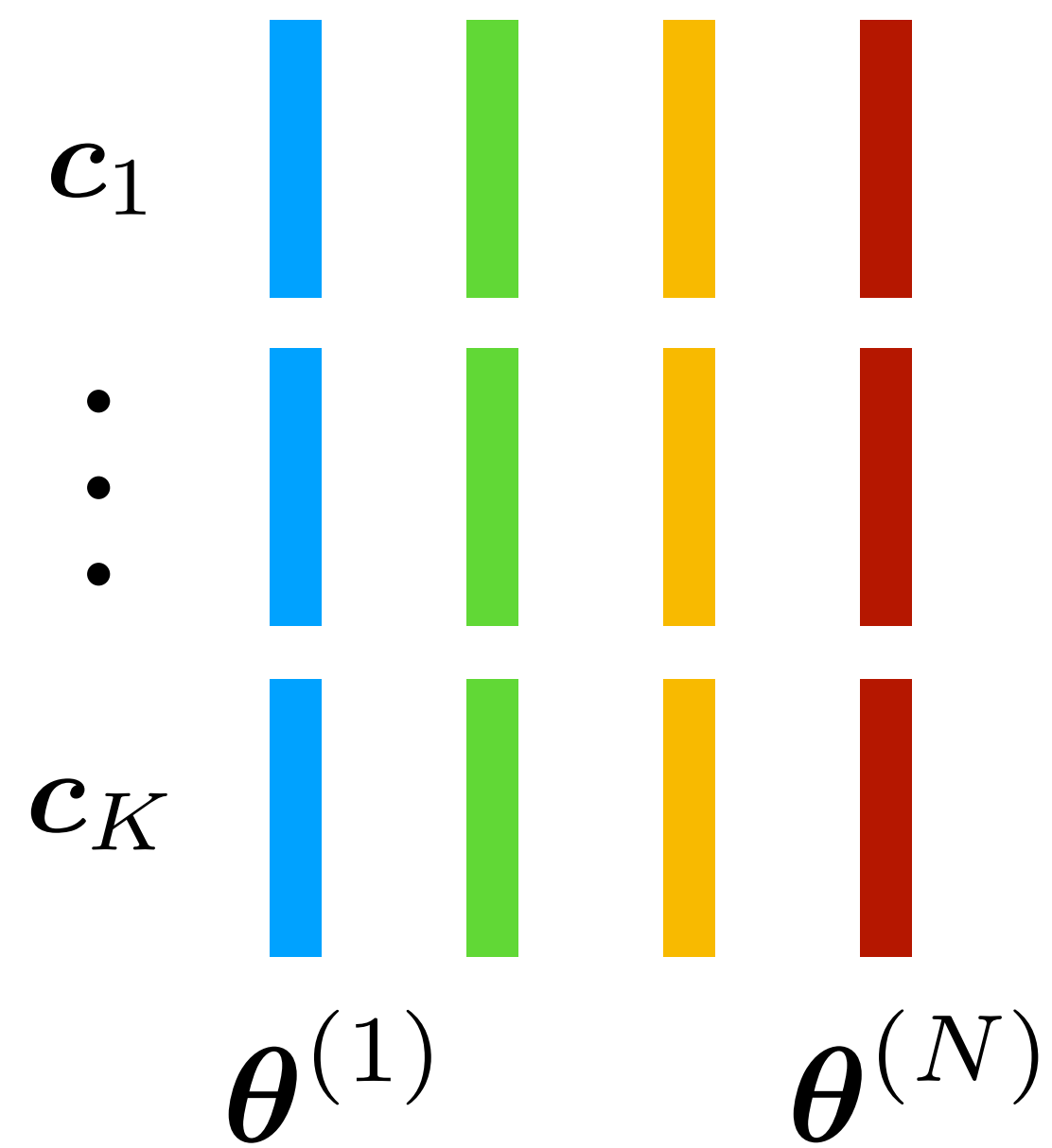
Goal: something that is *an epsilon* smarter than brute-force

For research purposes only

Results: size of m (k-means)

Can we find at least one (not necessarily efficient) decoder that succeeds where CLOMPR doesn't?

Let's explore the solution space (slightly better than) randomly: a genetic algorithm

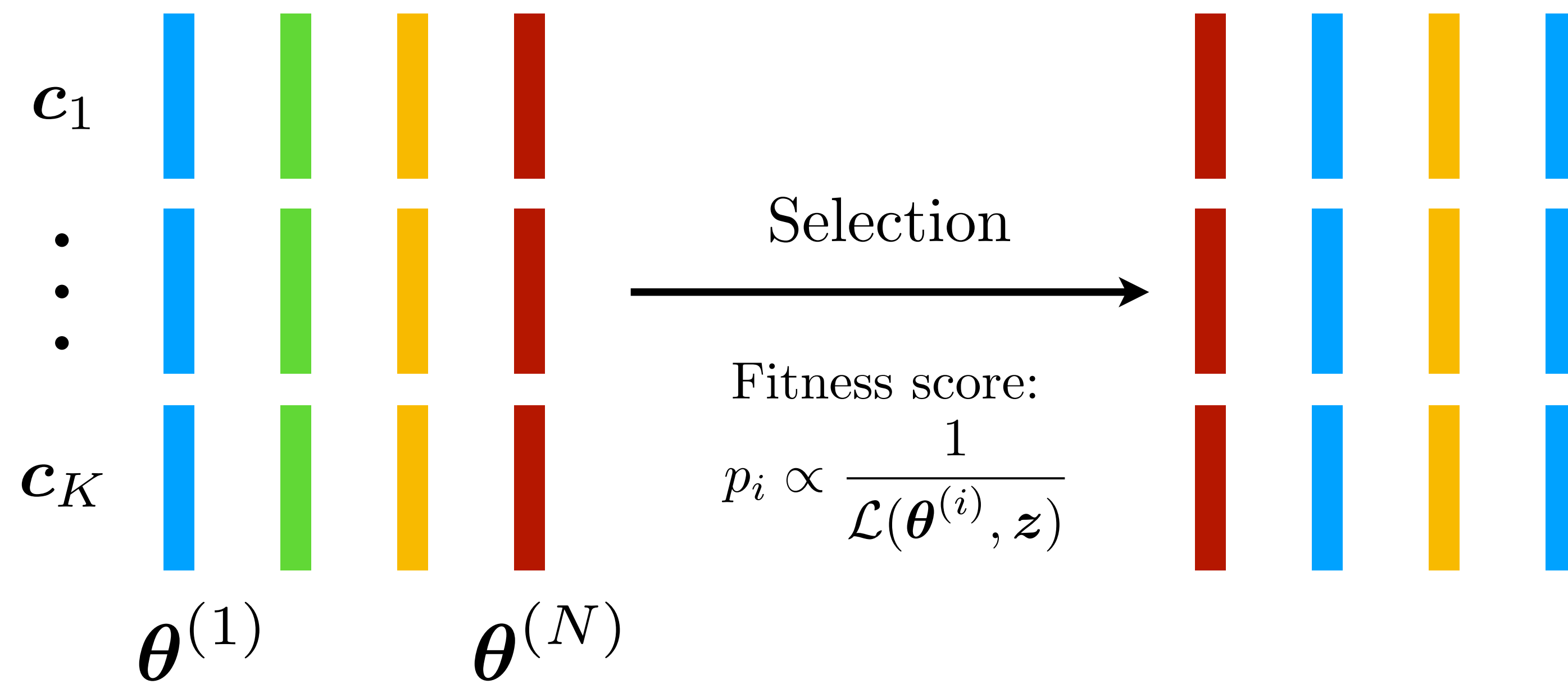


Population of (sets of)
centroids

Results: size of m (k-means)

Can we find at least one (not necessarily efficient) decoder that succeeds where CLOMPR doesn't?

Let's explore the solution space (slightly better than) randomly: a genetic algorithm

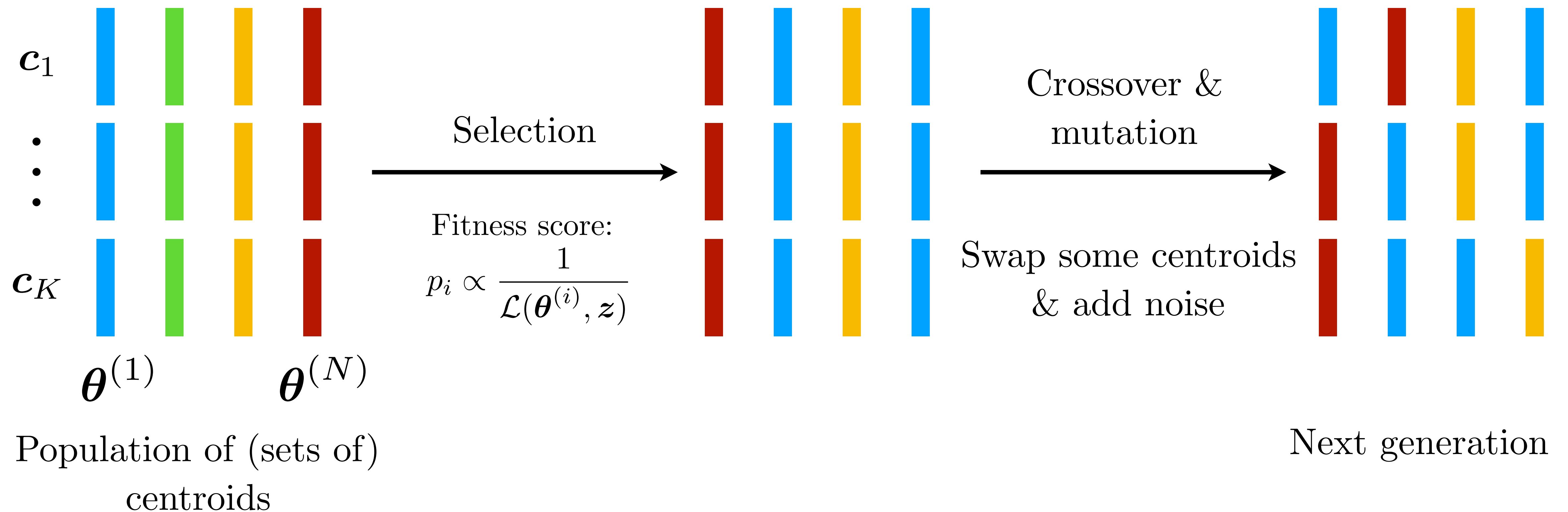


Population of (sets of)
centroids

Results: size of m (k-means)

Can we find at least one (not necessarily efficient) decoder that succeeds where CLOMPR doesn't?

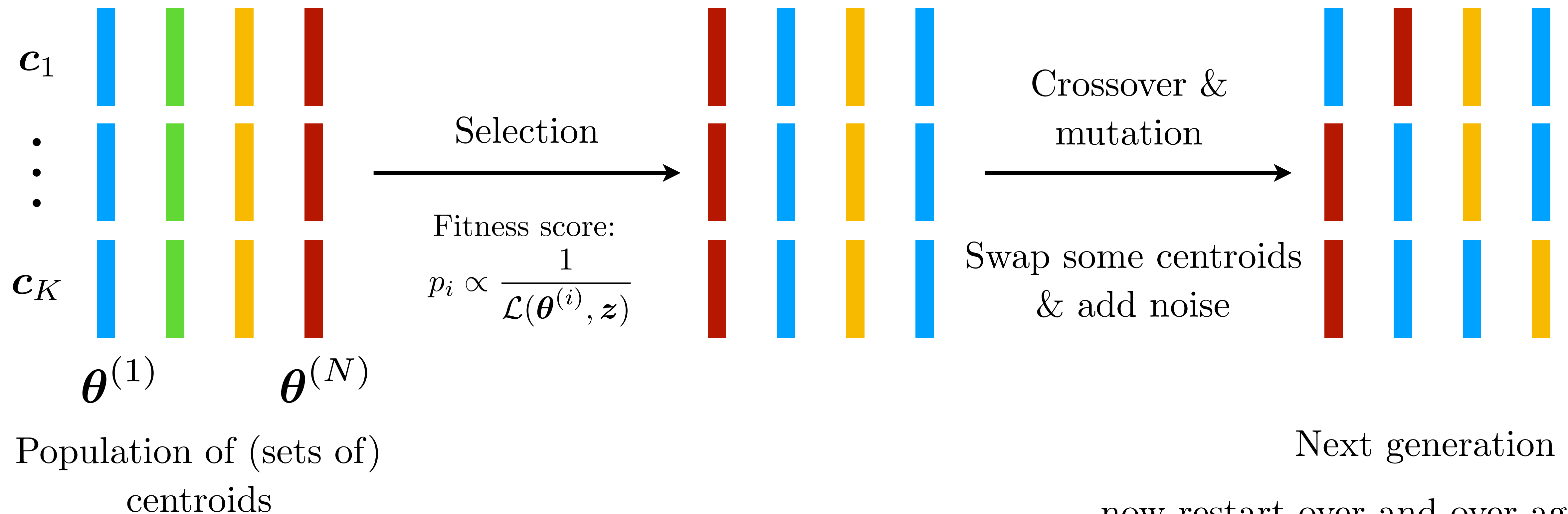
Let's explore the solution space (slightly better than) randomly: a genetic algorithm



Results: size of m (k-means)

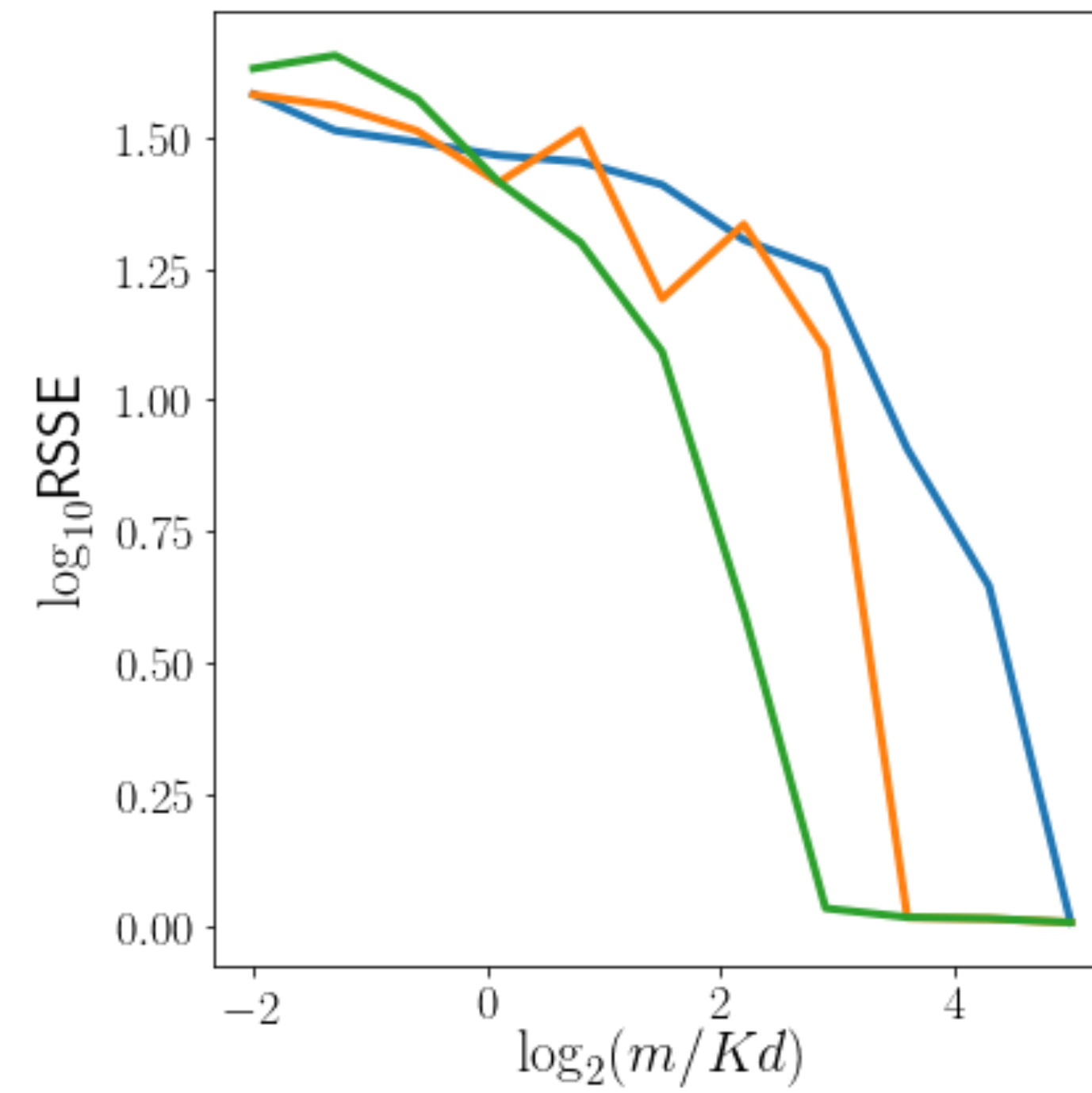
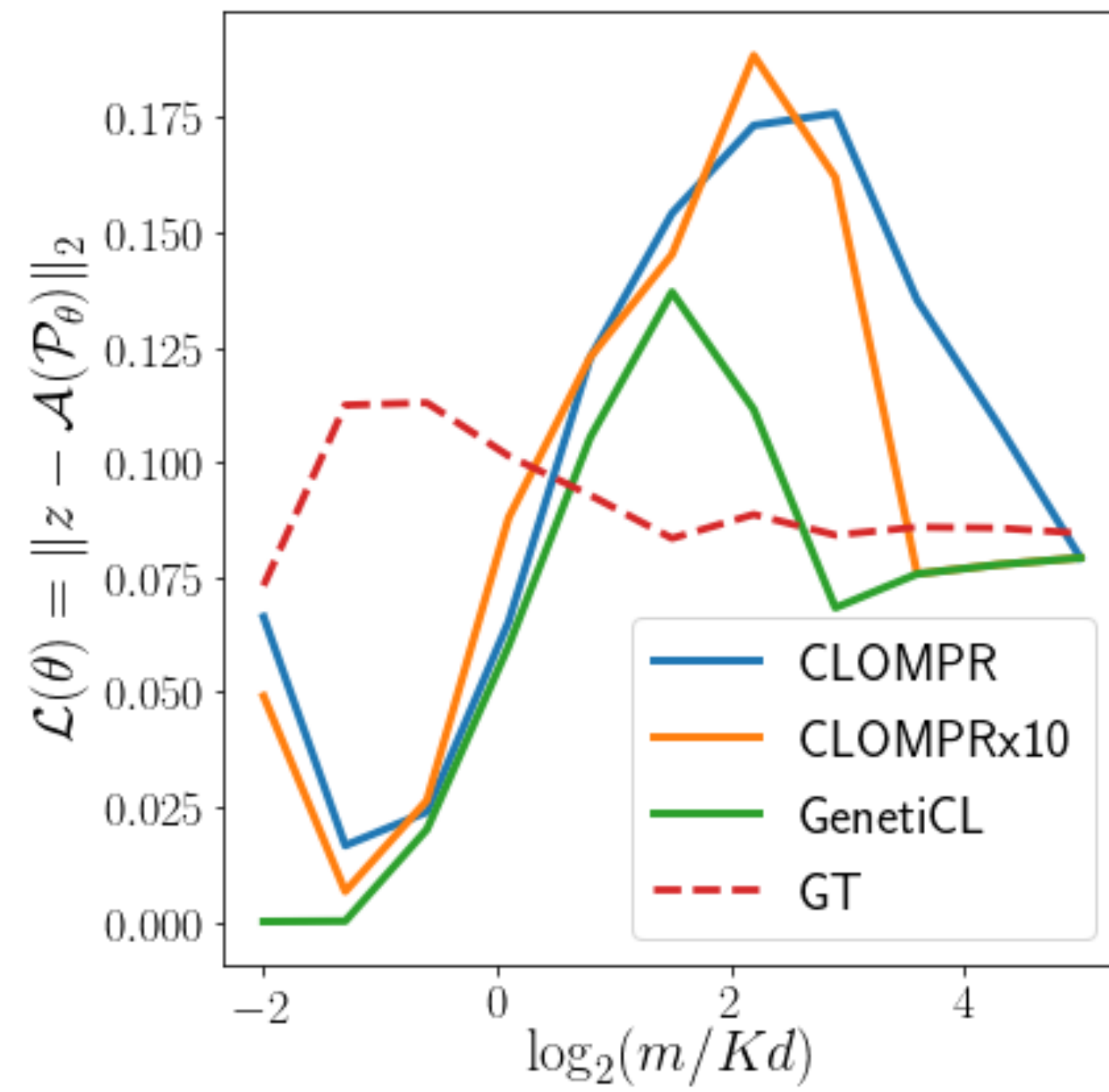
Can we find at least one (not necessarily efficient) decoder that succeeds where CLOMPR doesn't?

Let's explore the solution space (slightly better than) randomly: a genetic algorithm



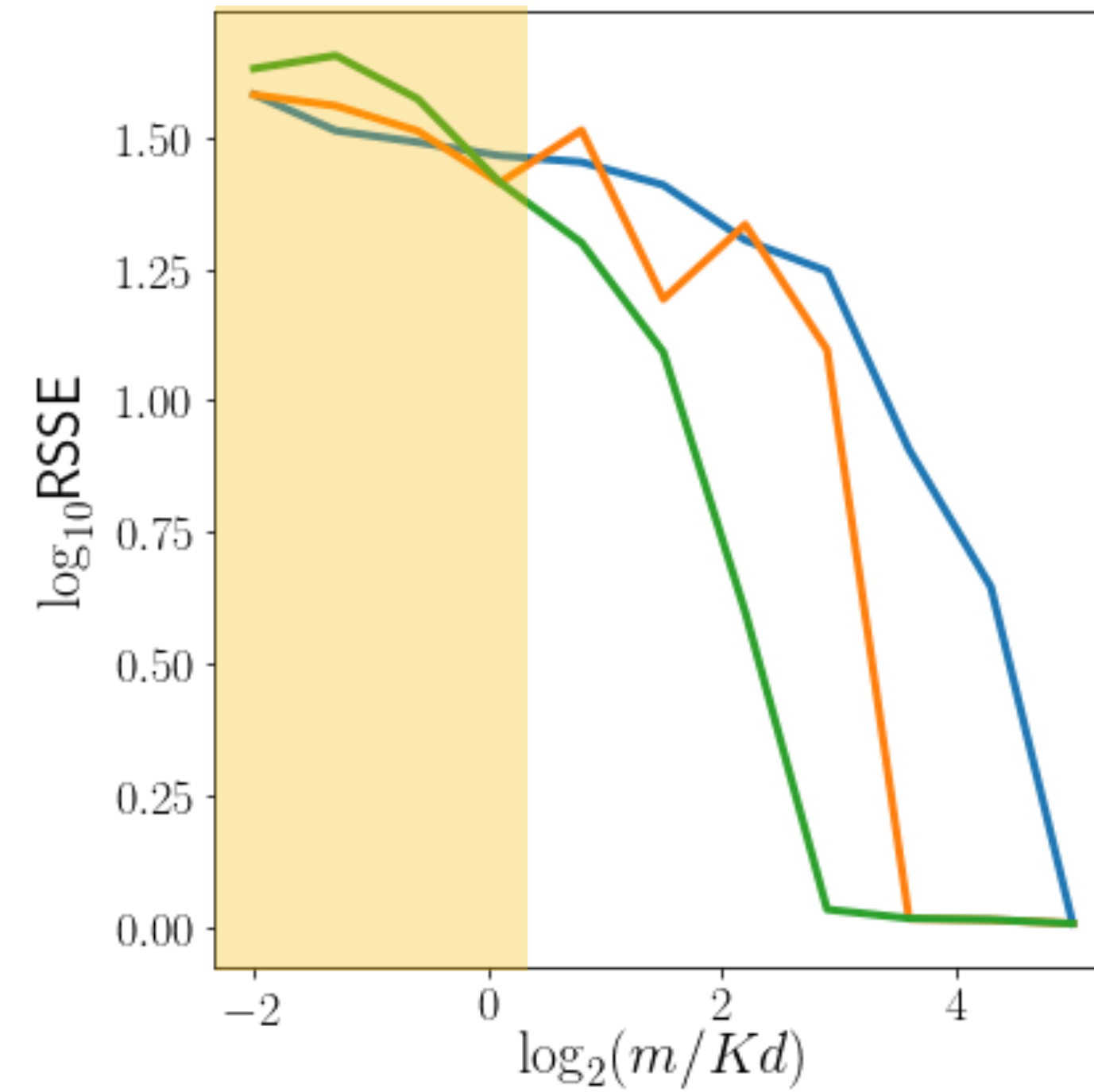
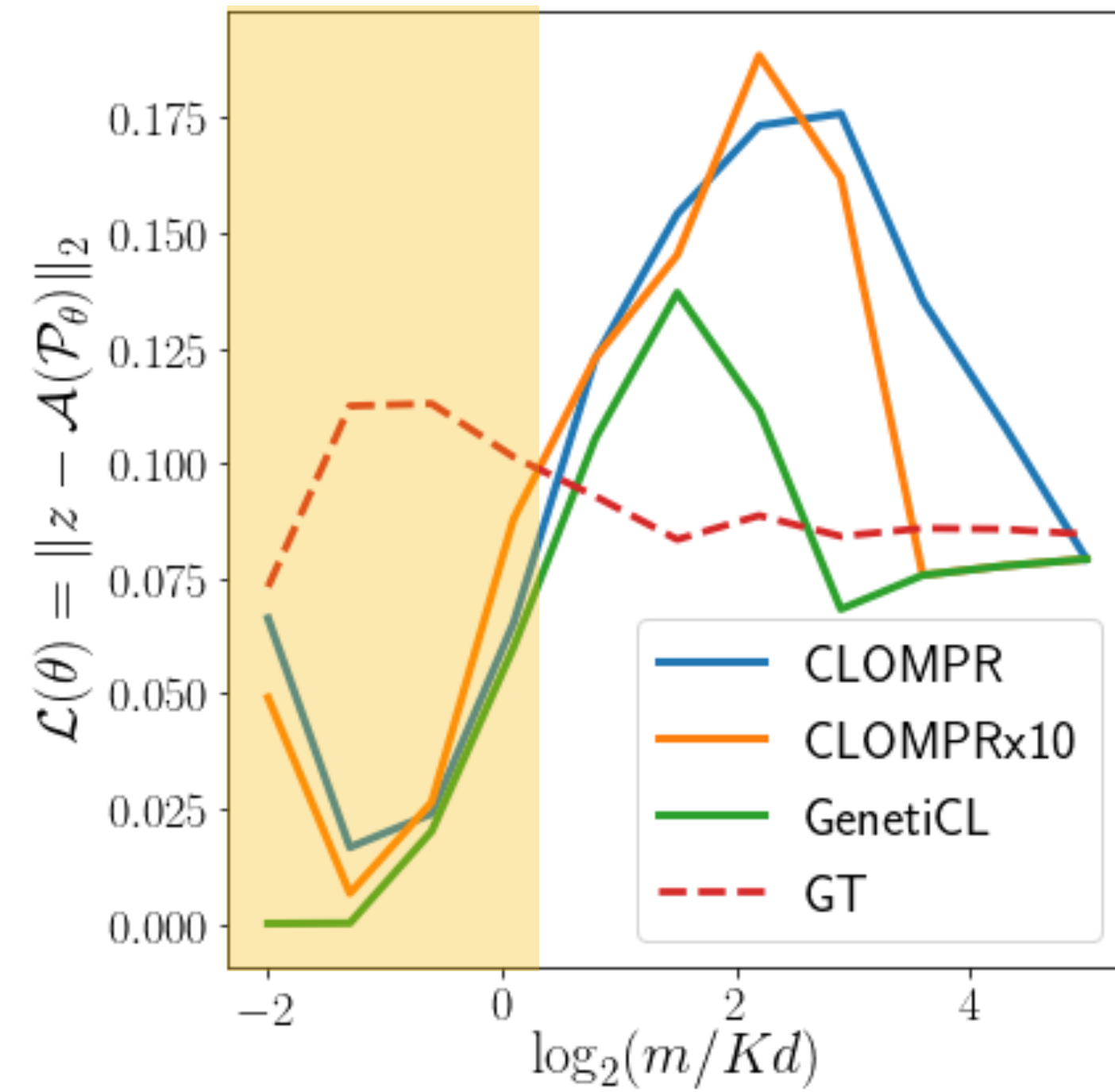
Results: size of m (k-means)

For a fixed sketch

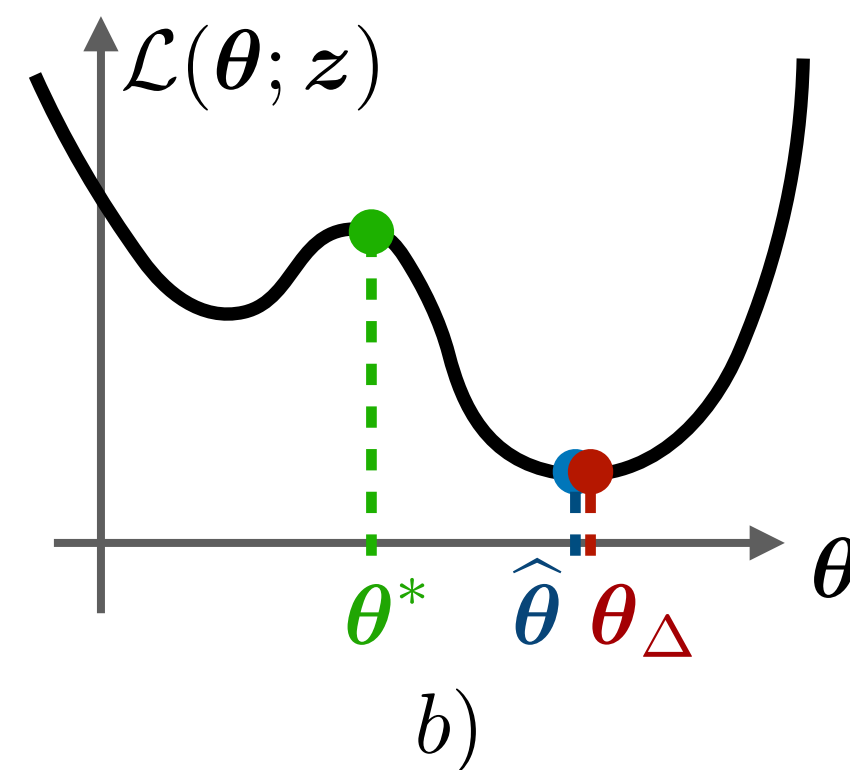


Results: size of m (k-means)

For a fixed sketch

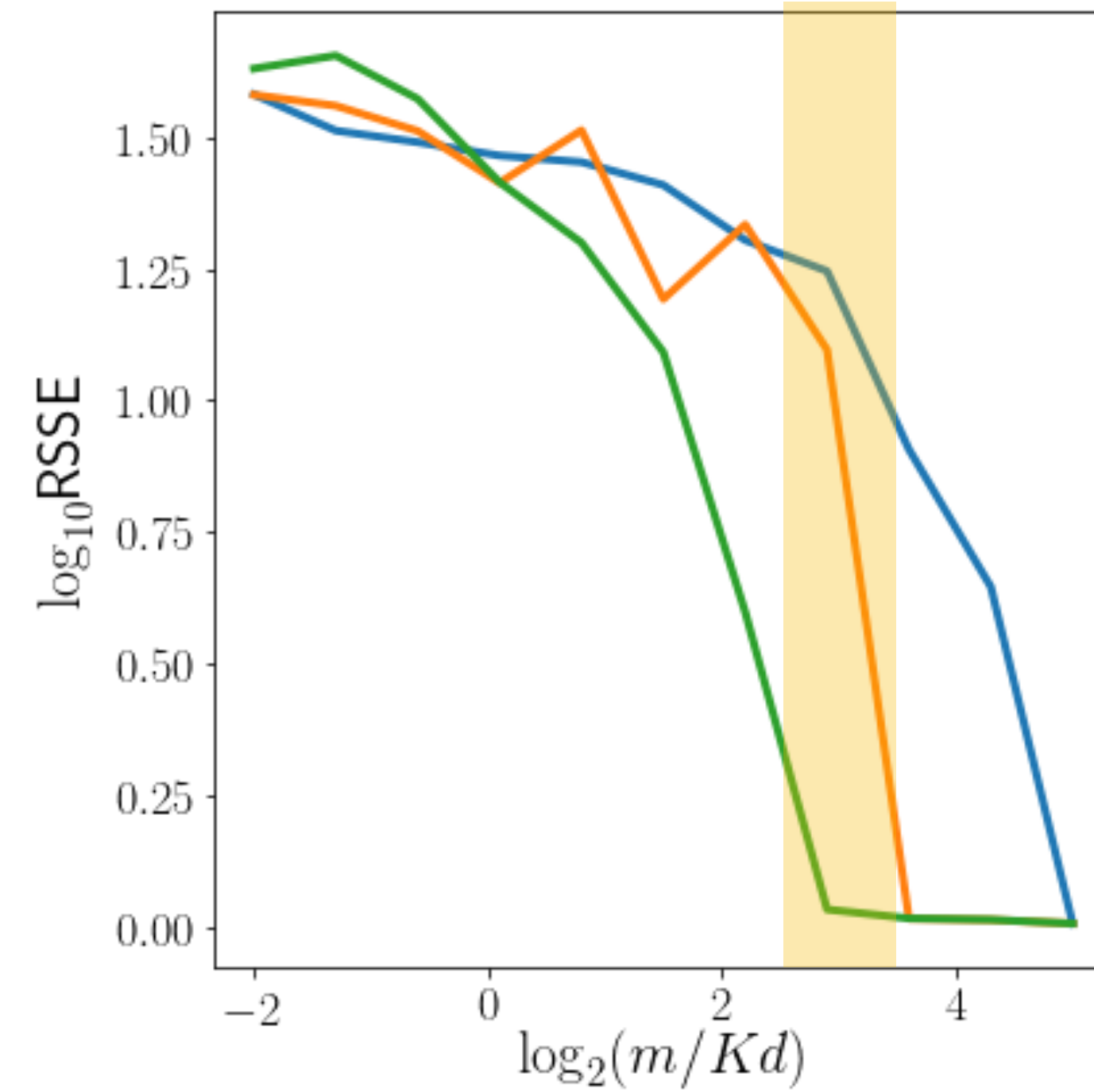
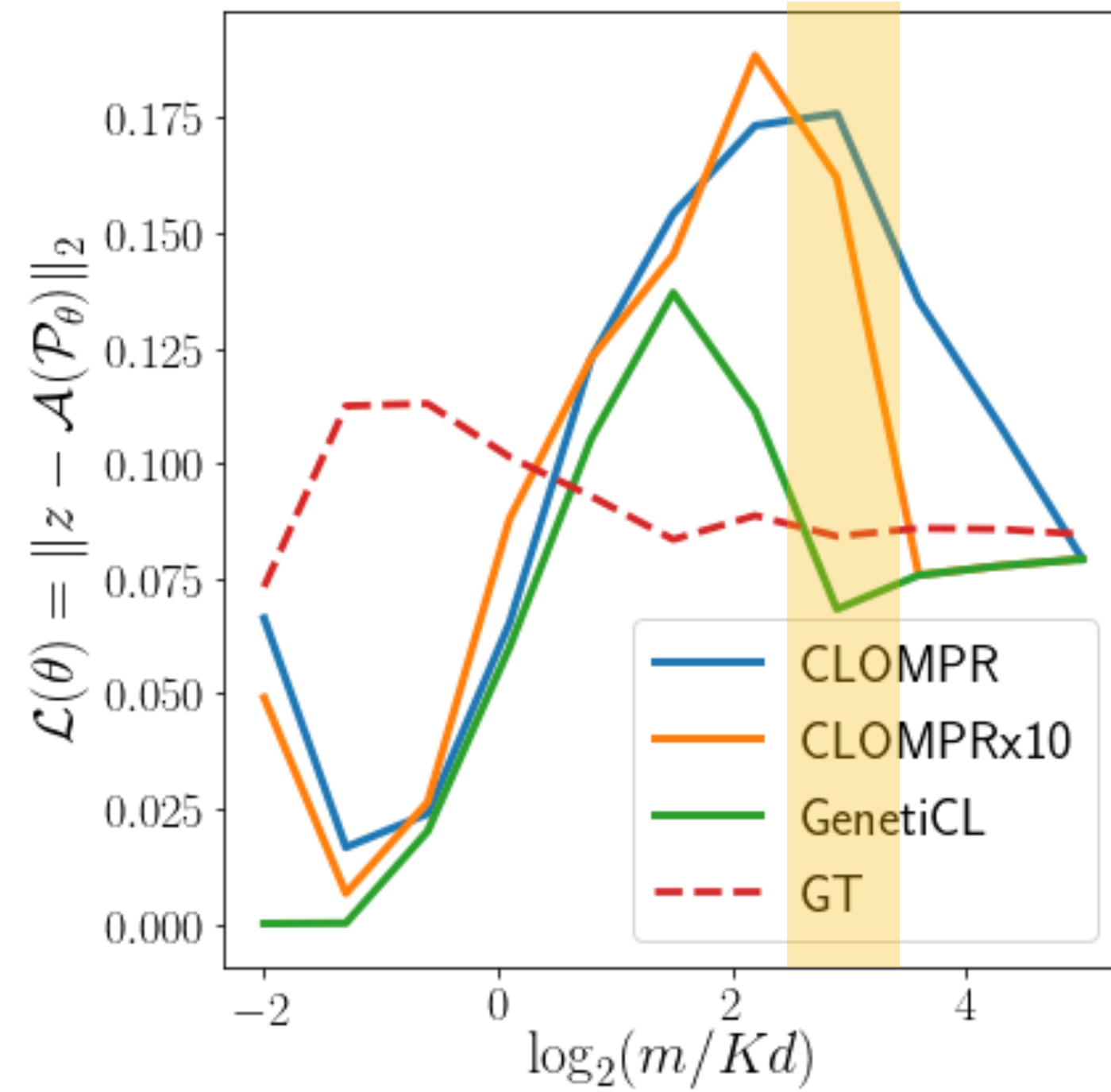


As before, loss is ill-defined when there are less measurements than parameters

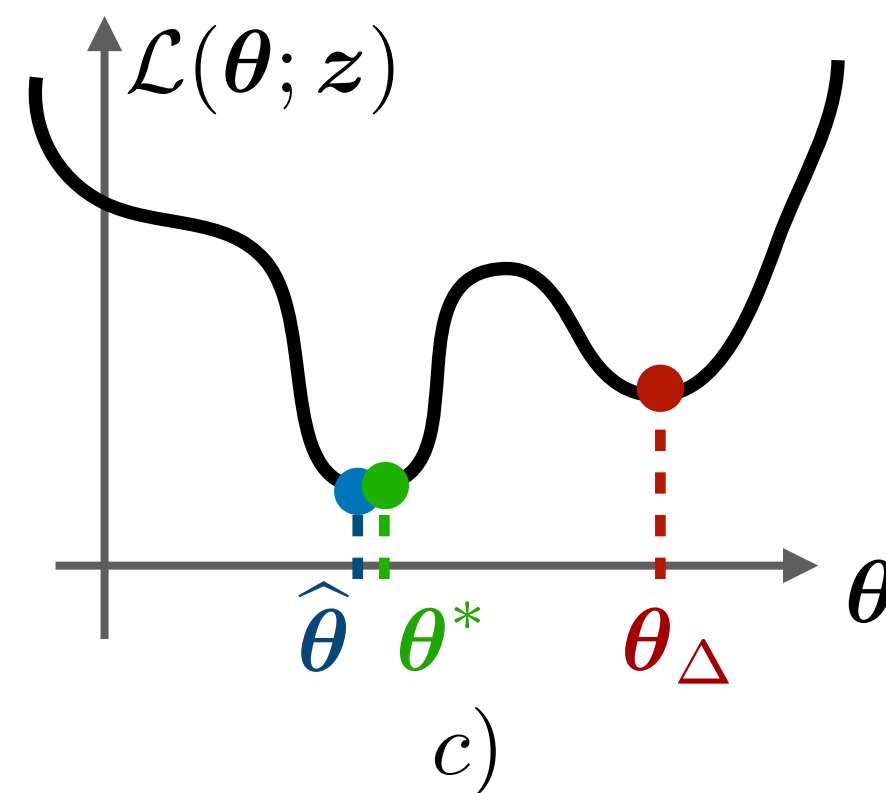


Results: size of m (k-means)

For a fixed sketch

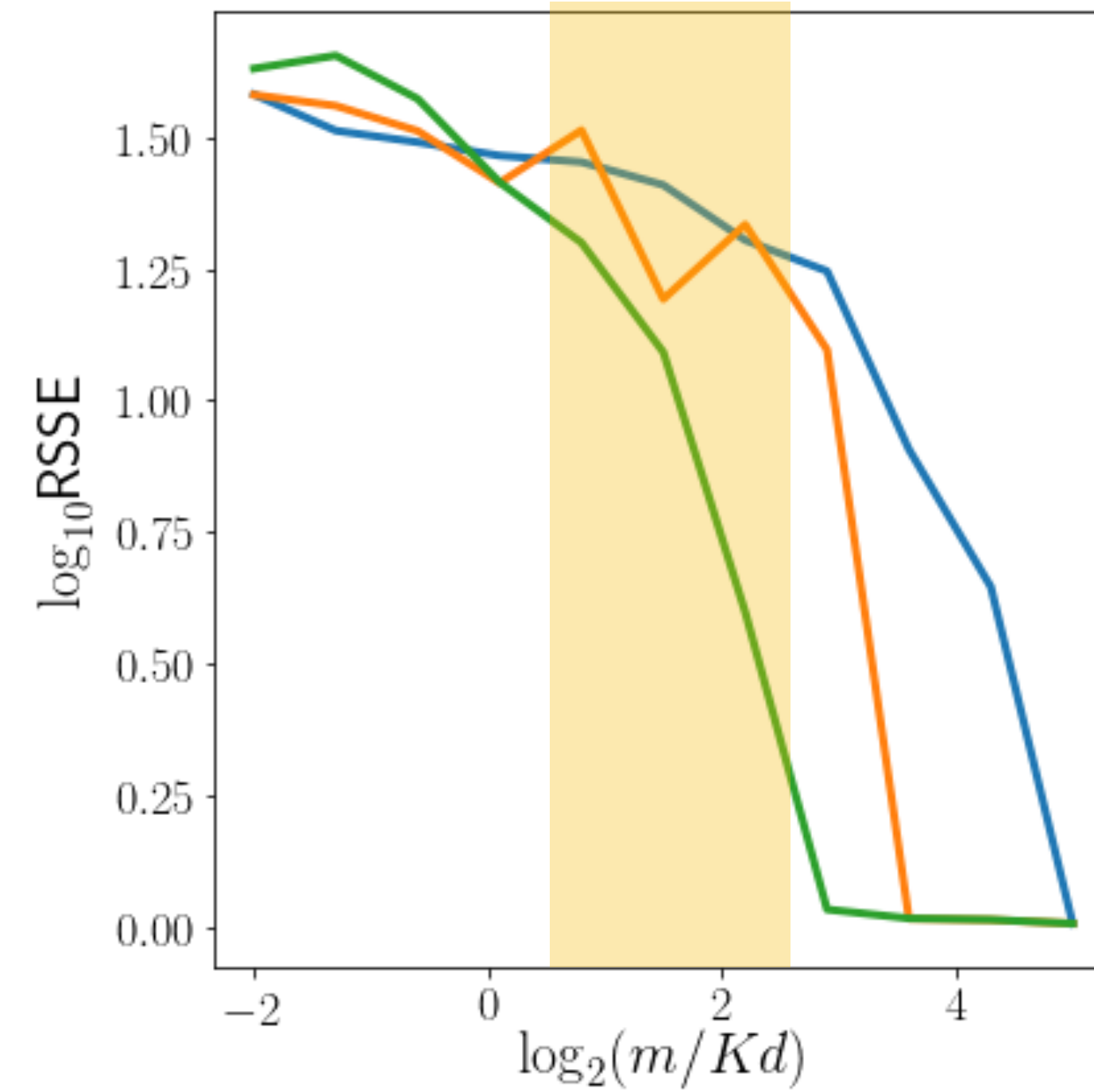
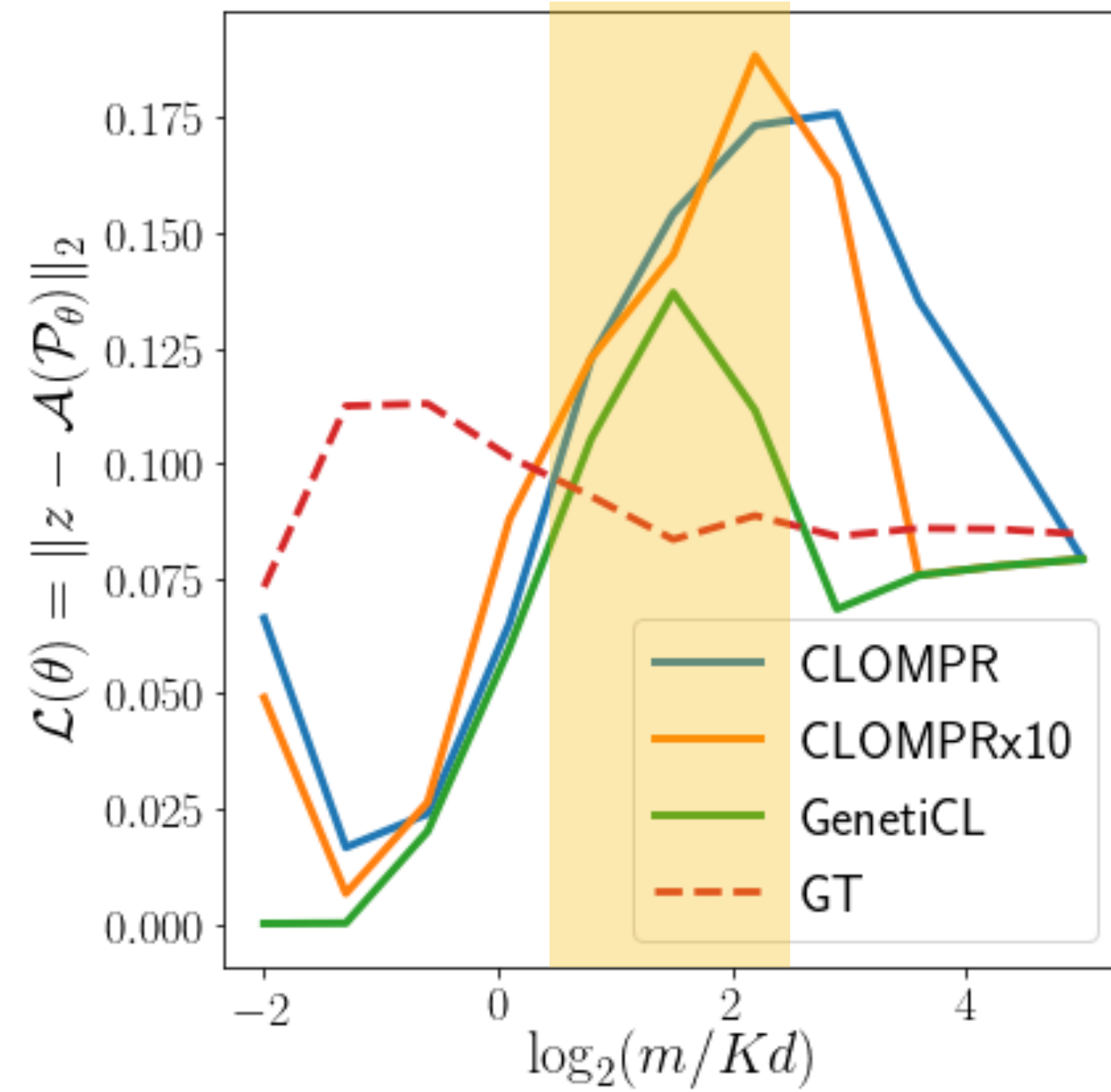


For m not too small,
GenetiCL finds a better optimum!
CLOMPR could be improved!

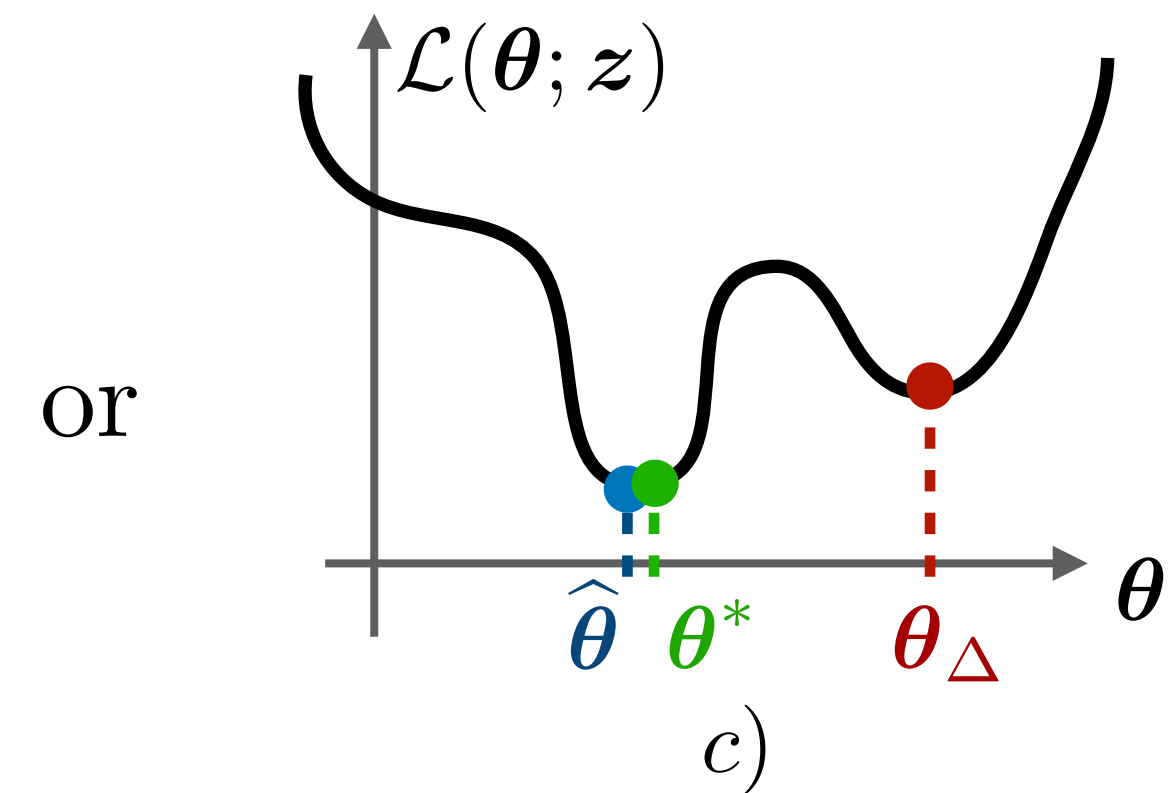
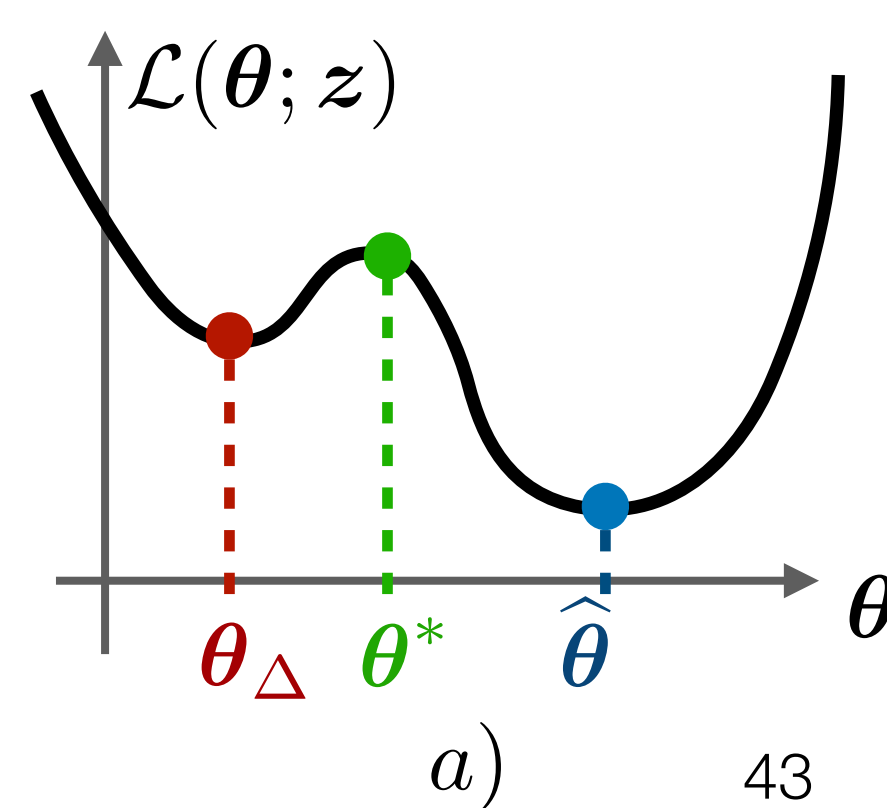


Results: size of m (k-means)

For a fixed sketch

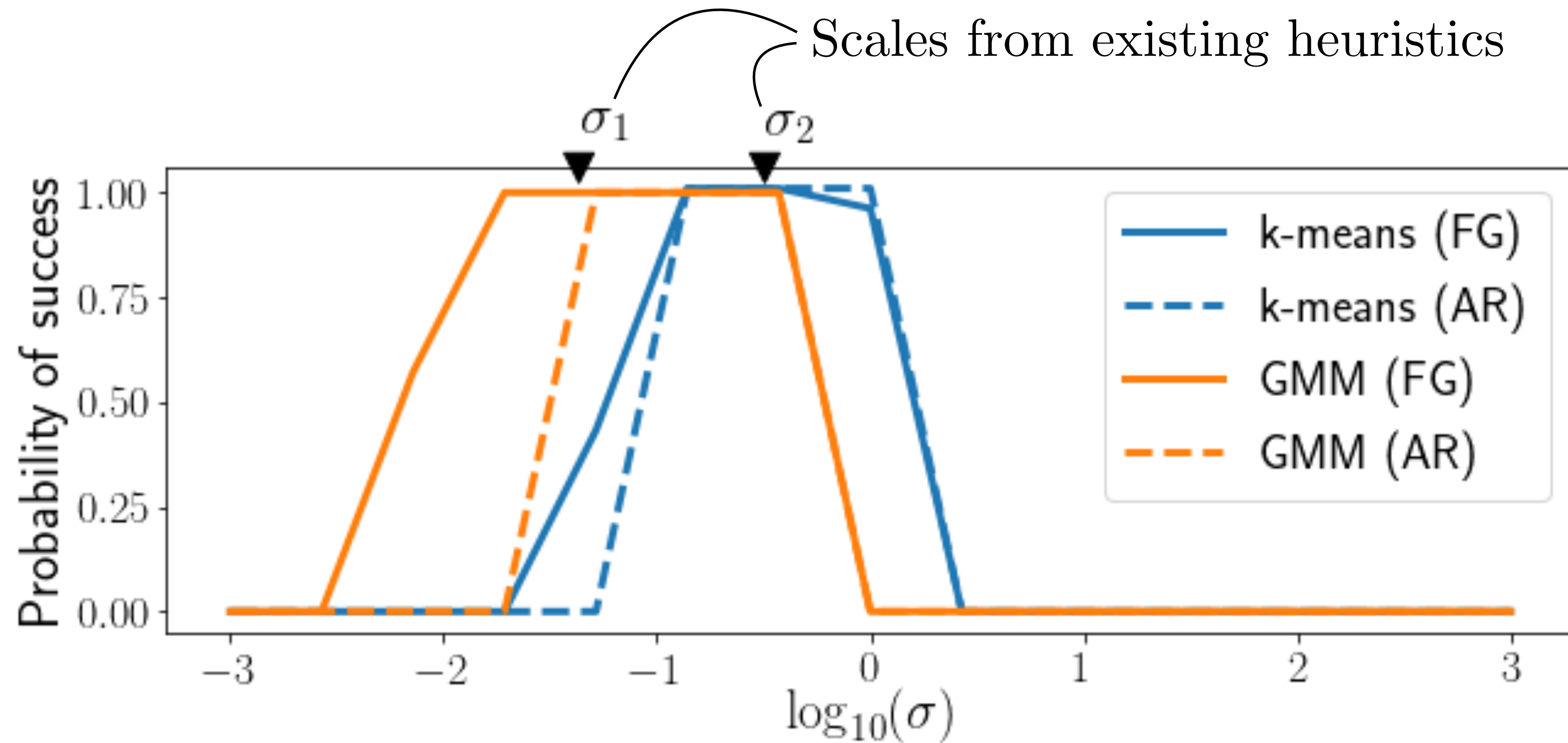


An ambiguous region remains...

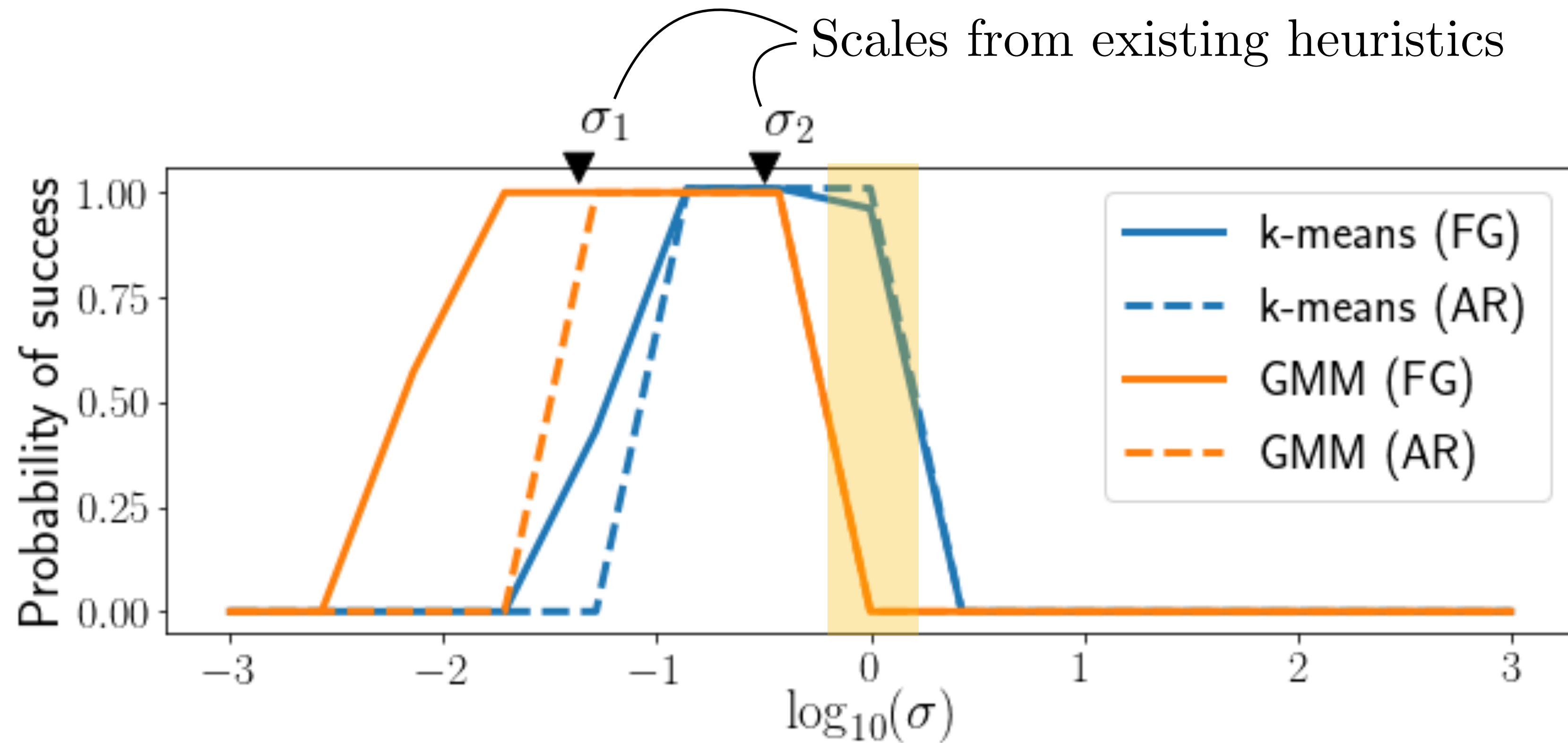


??

Results: sketch scale parameter (for different tasks)

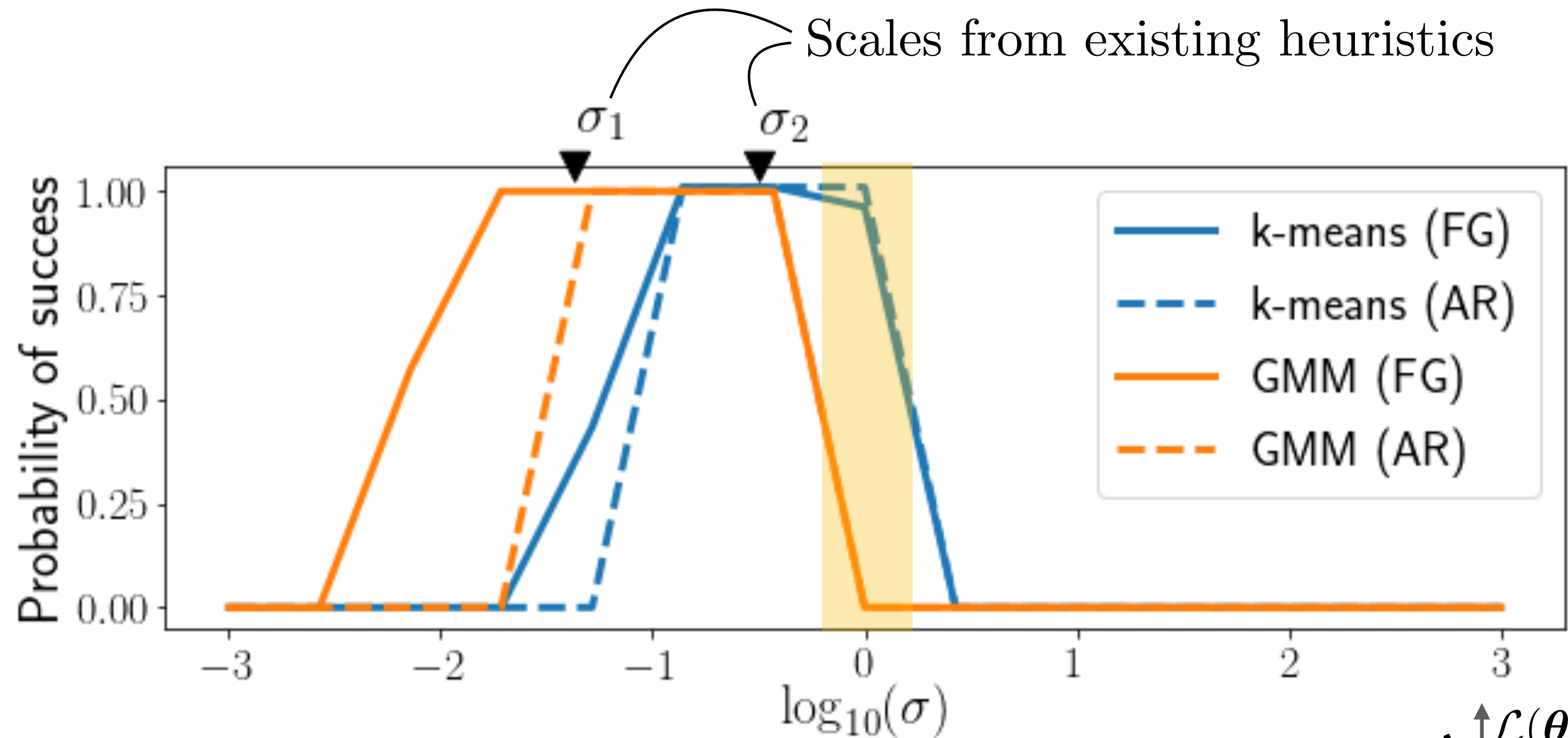


Results: sketch scale parameter (for different tasks)



At large scales, k-means succeeds but not GMM

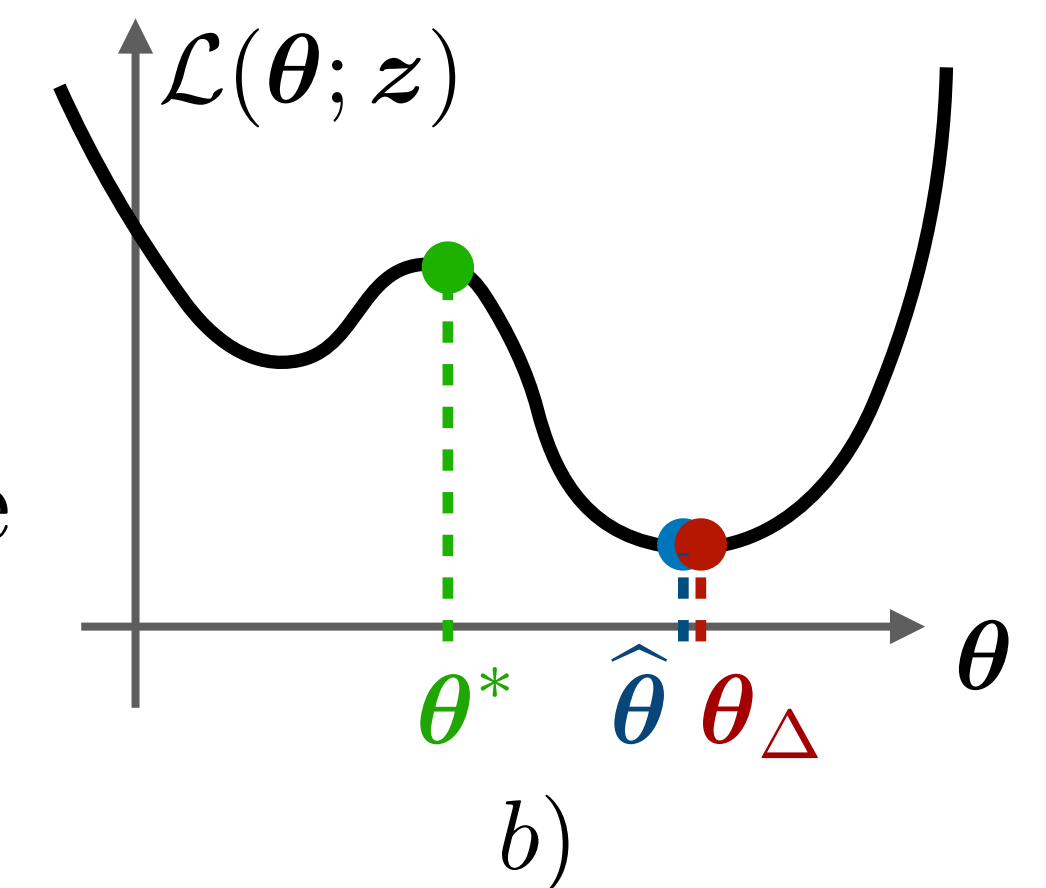
Results: sketch scale parameter (for different tasks)



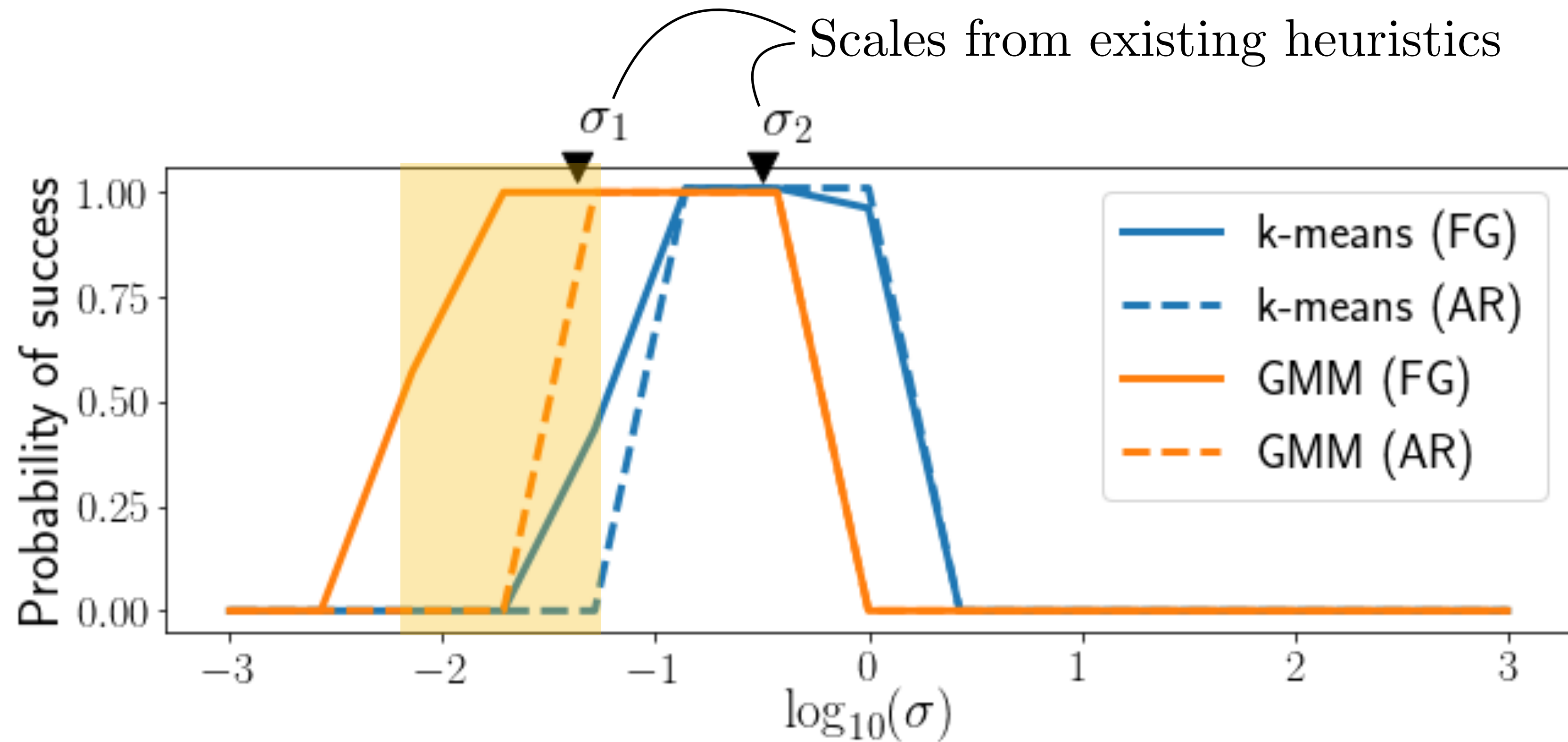
At large scales, k-means succeeds but not GMM

(sketch failures)

Sketch captures a rough (low-pass) approx. to the data distribution



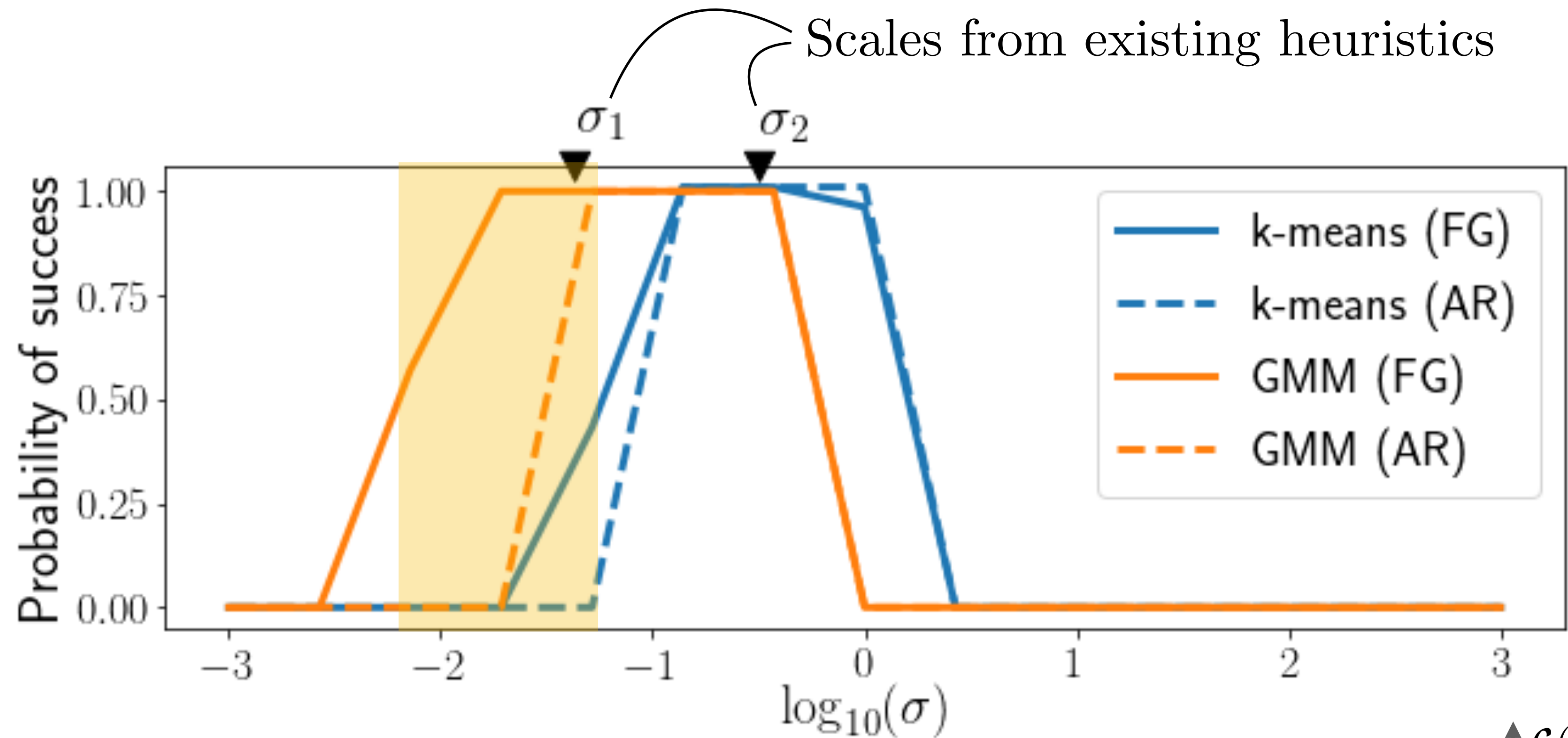
Results: sketch scale parameter (for different tasks)



At large scales, k-means succeeds but not GMM

At small scales, GMM succeeds but not k-means (decoder failures)

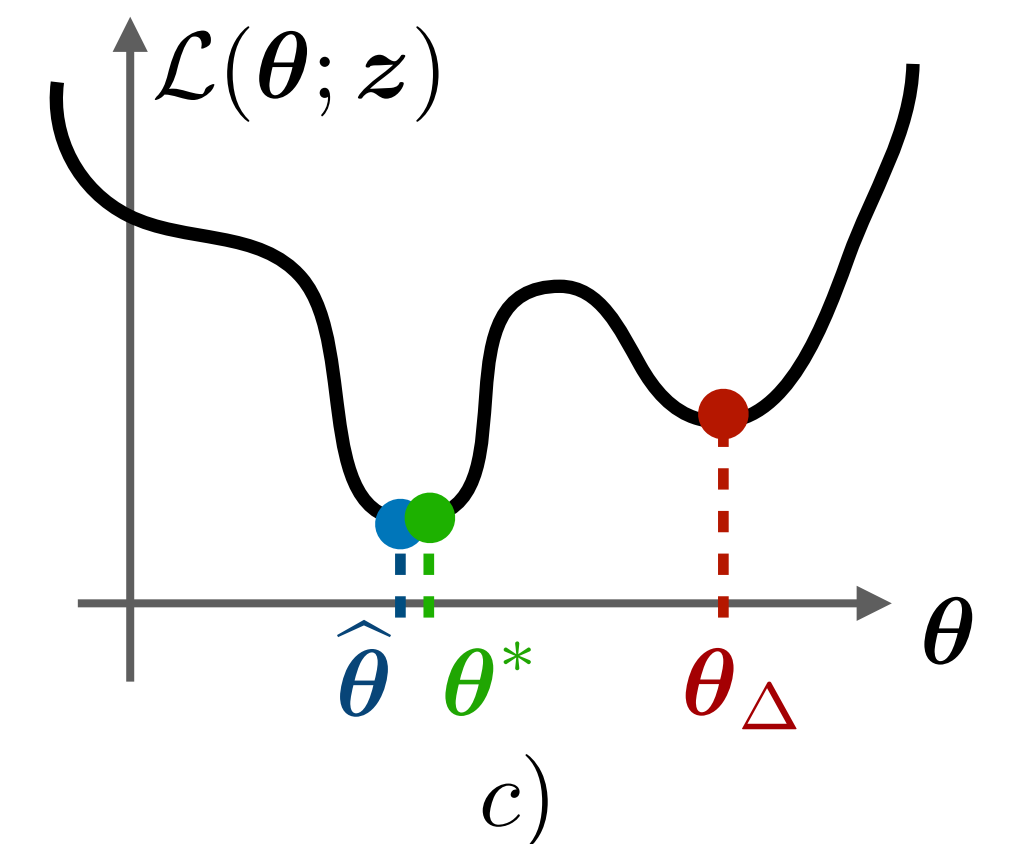
Results: sketch scale parameter (for different tasks)



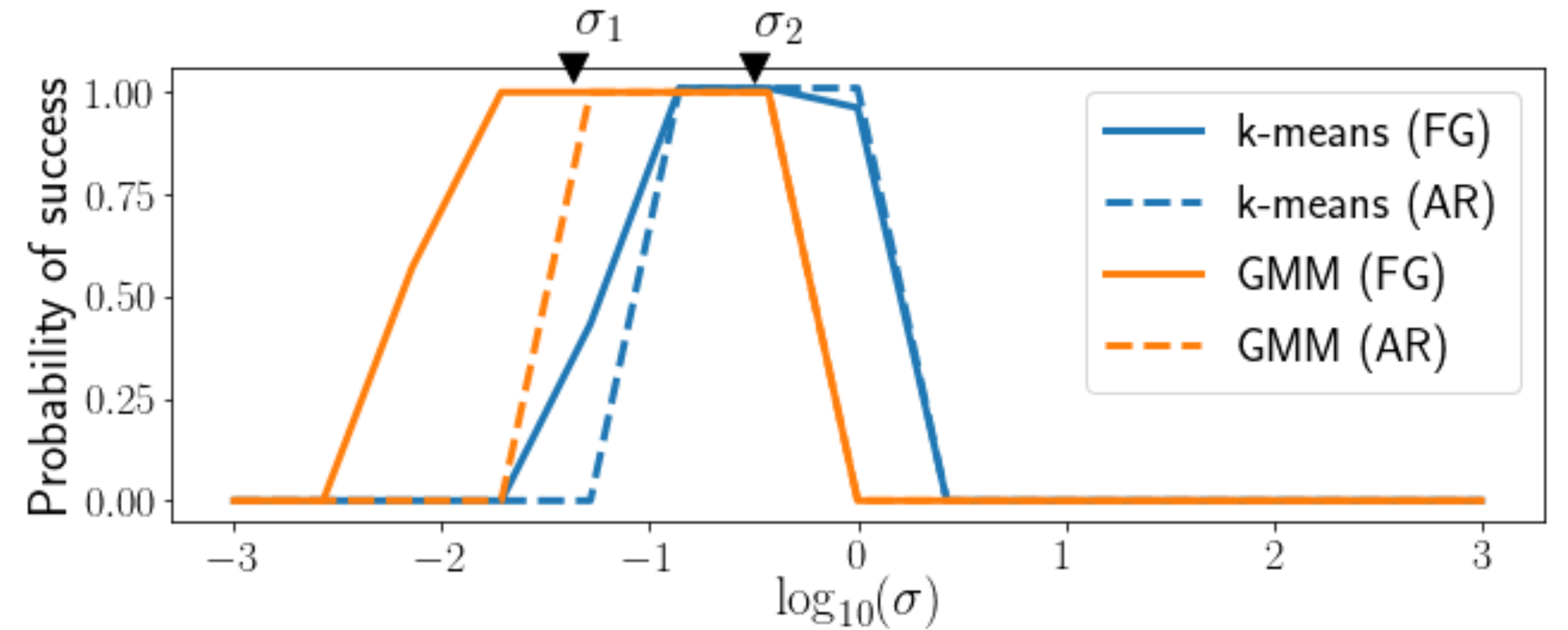
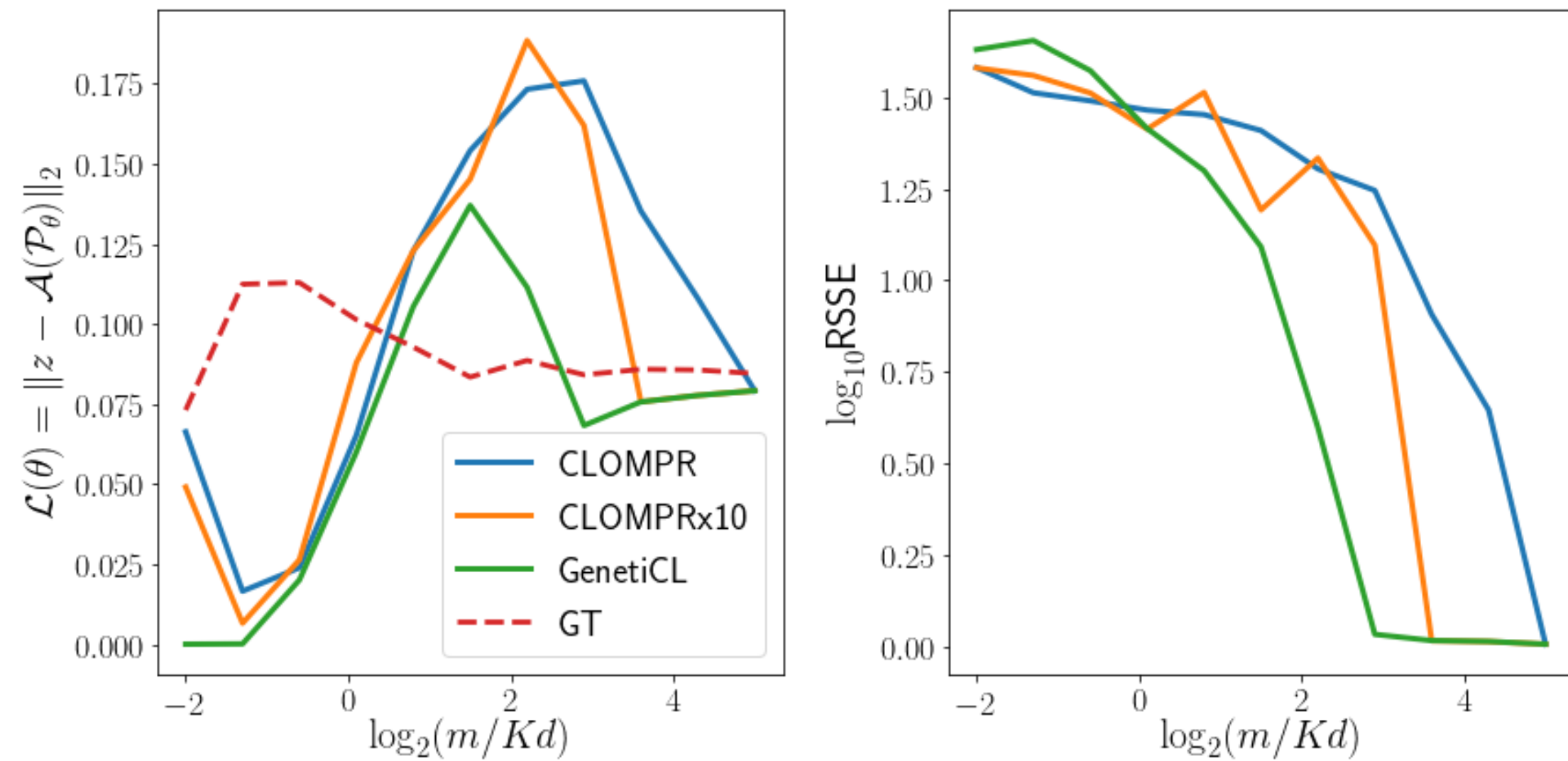
At large scales, k-means succeeds but not GMM

At small scales, GMM succeeds but not k-means

More surprising: adding variables makes the decoder behave better!



Take-home



- There are regimes where CLOMPR could be improved (small m or sigma for k-means)
- An ambiguous region remains at intermediary m (landscape probably full of local minima)
- Question for the future: can we actually do better than CLOMPR (efficiently)? Guarantees?