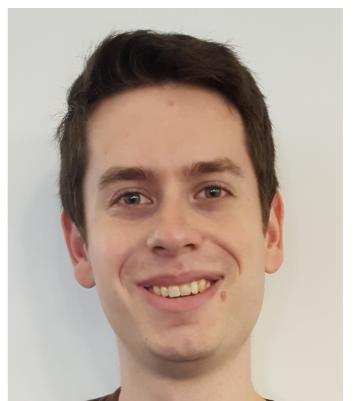


Approximate Cross Validation for Large Data and High Dimensions

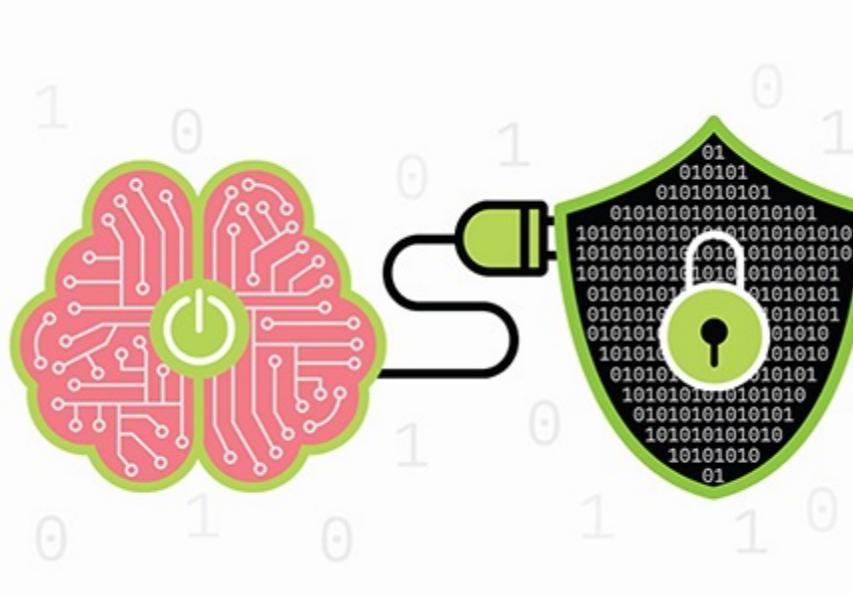
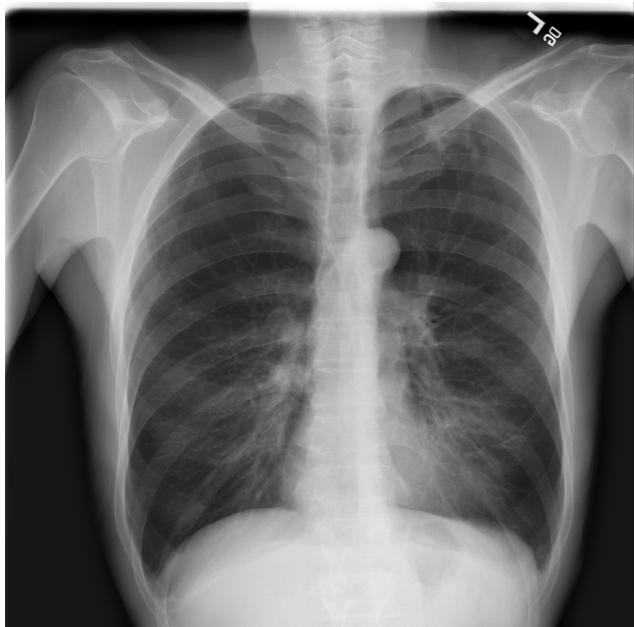
Tamara Broderick

Associate Professor,
MIT

Ryan Giordano, William T. Stephenson

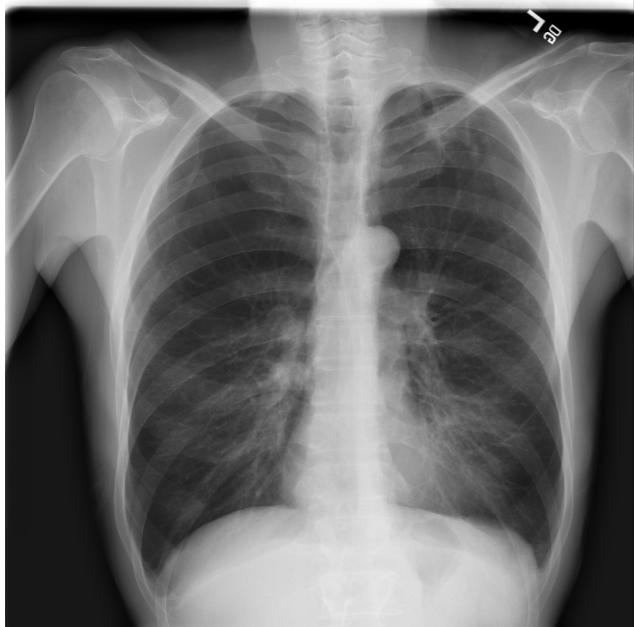


- Better computation → increasingly complex and expensive data analysis in life-changing application areas



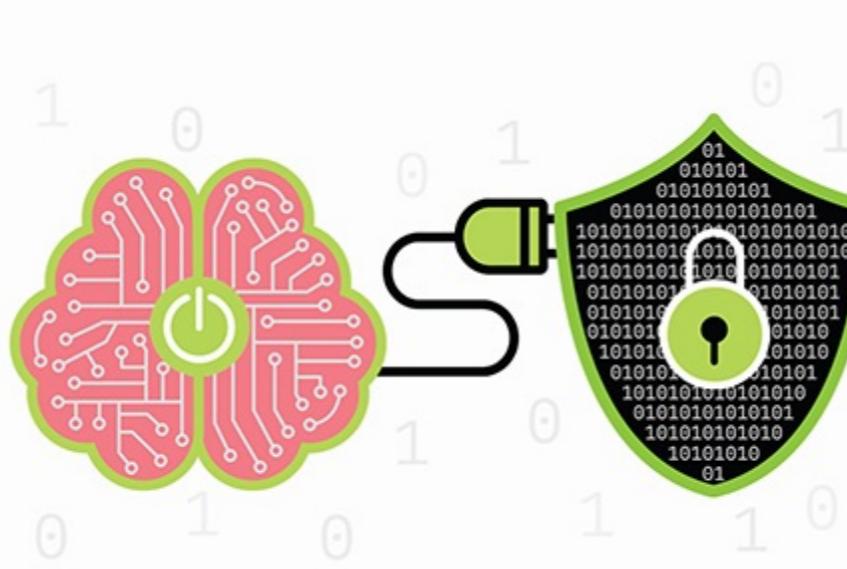
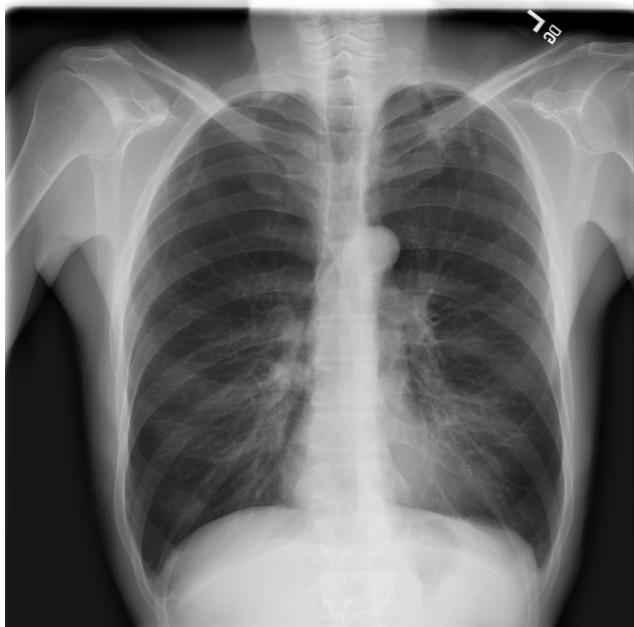
Evaluation

- Better computation → increasingly complex and expensive data analysis in life-changing application areas



Evaluation

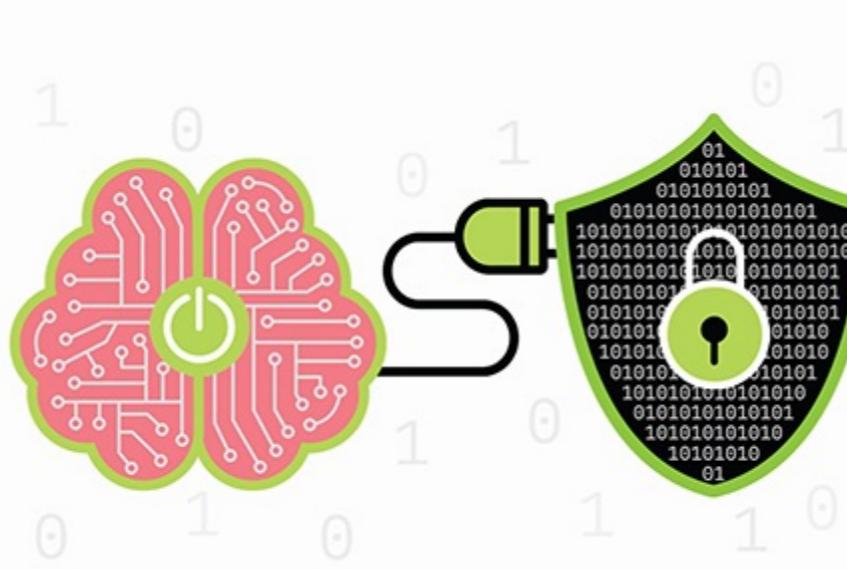
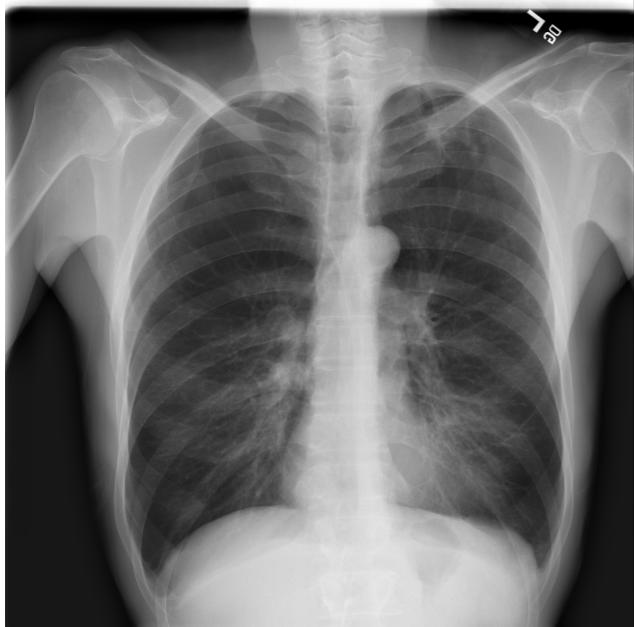
- Better computation → increasingly complex and expensive data analysis in life-changing application areas



- How well does the data analysis work?

Evaluation

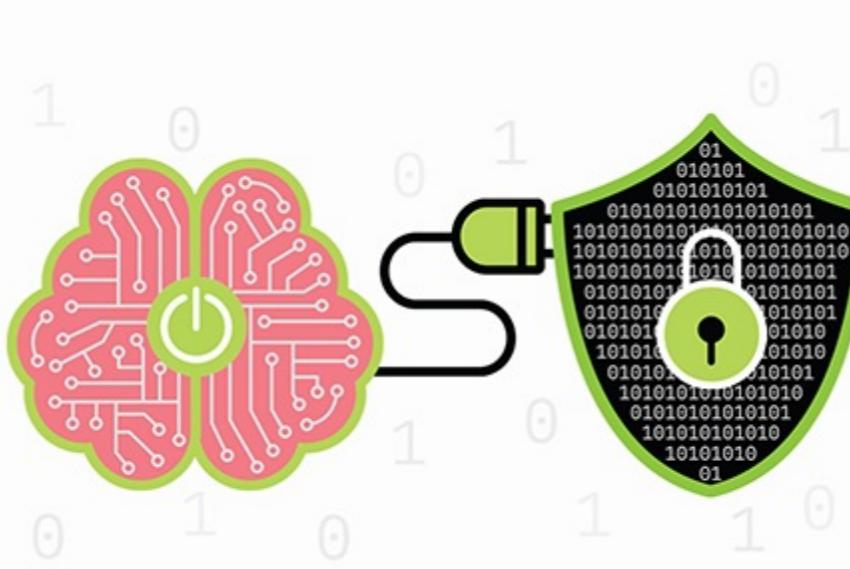
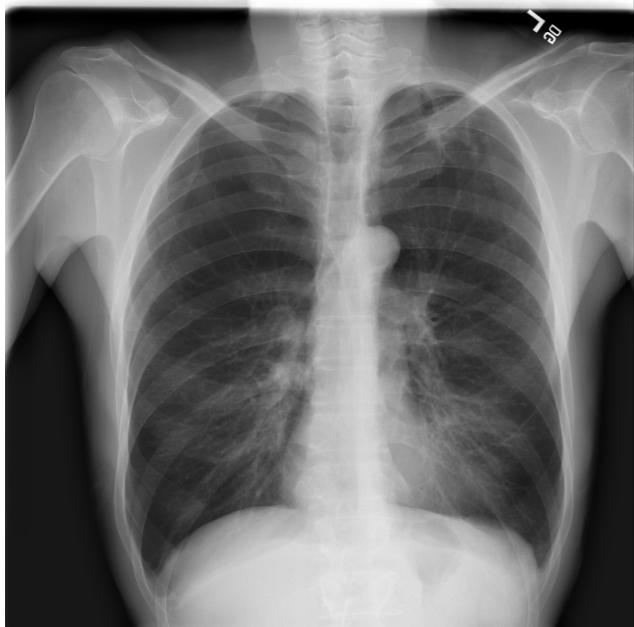
- Better computation → increasingly complex and expensive data analysis in life-changing application areas



- How well does the data analysis work? Evaluation concerns: accuracy, compute time, & user time

Evaluation

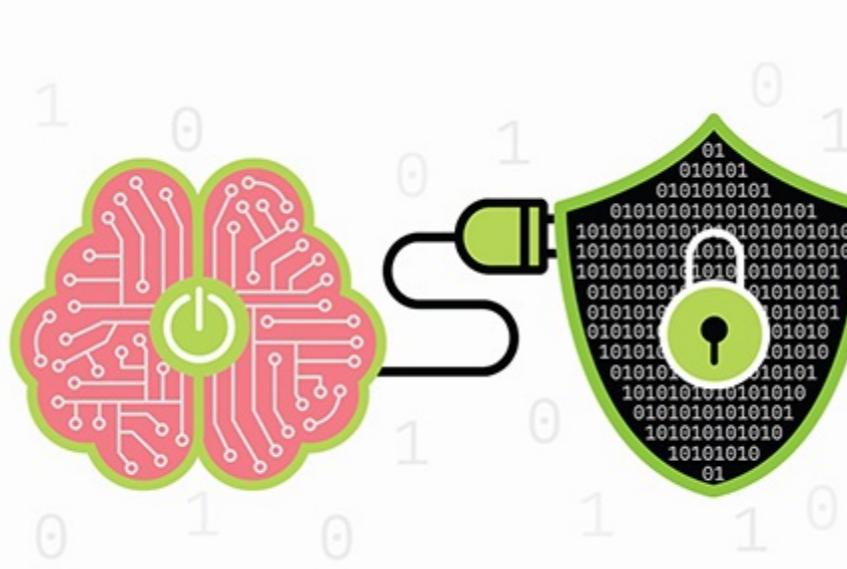
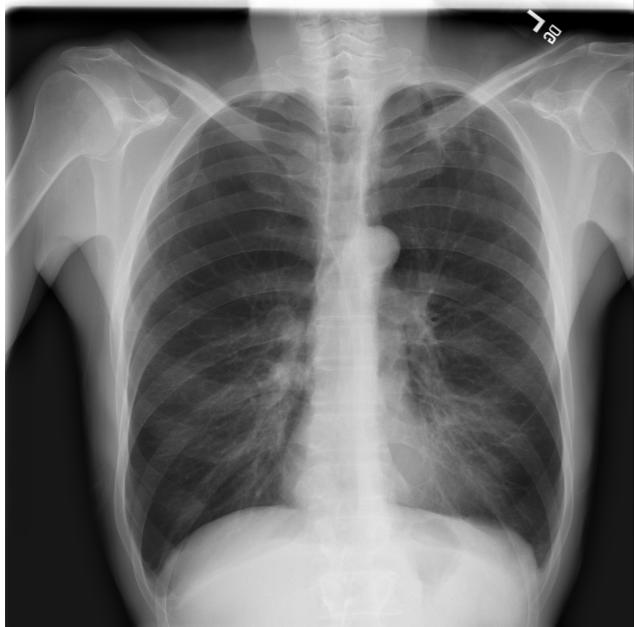
- Better computation → increasingly complex and expensive data analysis in life-changing application areas



- How well does the data analysis work? Evaluation concerns: accuracy, compute time, & user time
- Cross validation seems to work well, is model agnostic, but can be expensive (multiple analysis runs)

Evaluation

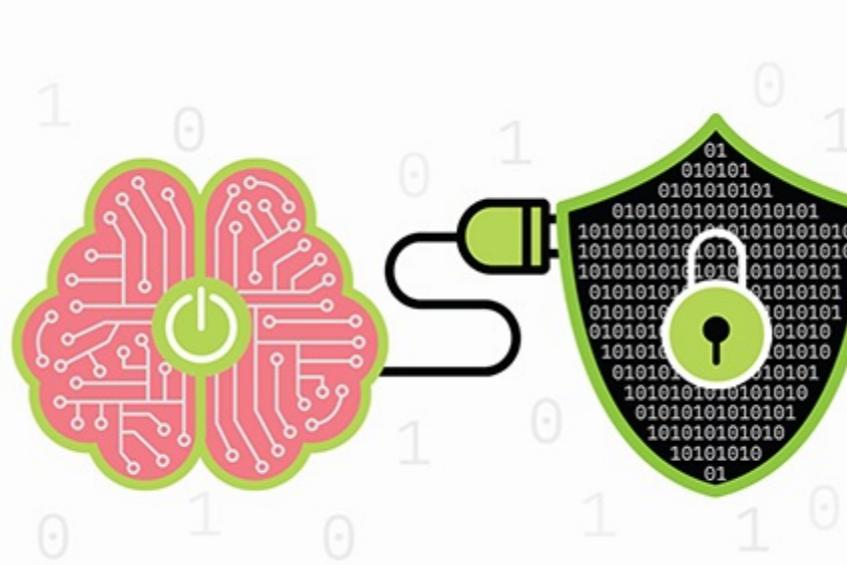
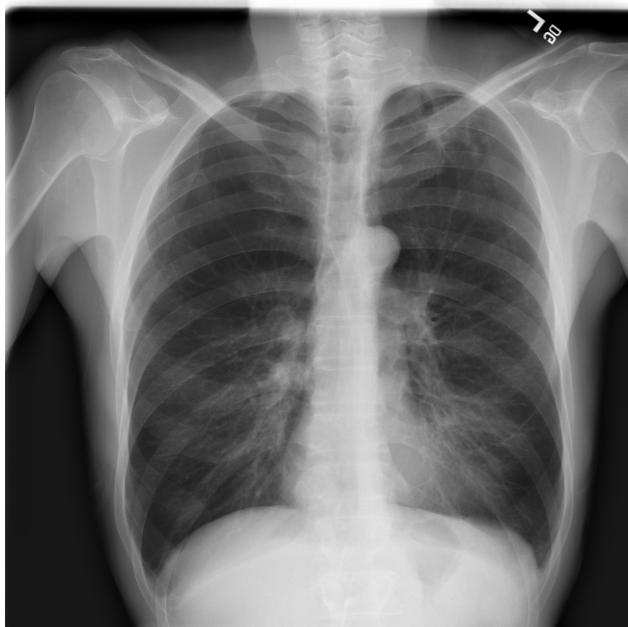
- Better computation → increasingly complex and expensive data analysis in life-changing application areas



- How well does the data analysis work? Evaluation concerns: accuracy, compute time, & user time
- Cross validation seems to work well, is model agnostic, but can be expensive (multiple analysis runs)
- Can *approximate* CV: Taylor expand + reduce dimension
 - We show: is **fast**, **automatic**, and provably **high-quality**

Evaluation

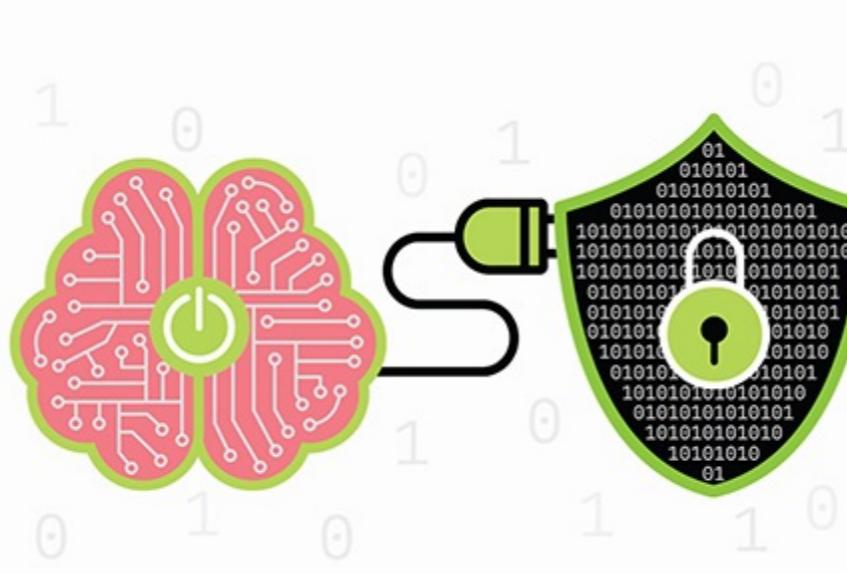
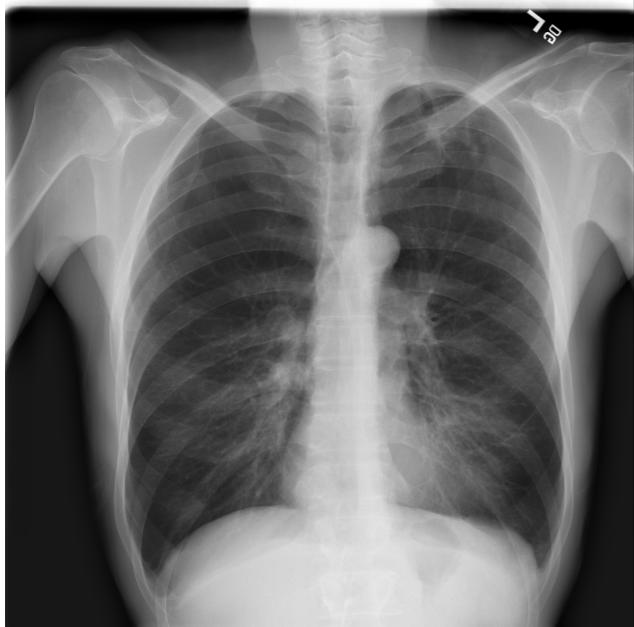
- Better computation → increasingly complex and expensive data analysis in life-changing application areas



- How well does the data analysis work? Evaluation concerns: accuracy, compute time, & user time
- Cross validation seems to work well, is model agnostic, but can be expensive (multiple analysis runs)
- Can *approximate* CV
 - We show: is **fast**, **automatic**, and provably **high-quality**

Evaluation

- Better computation → increasingly complex and expensive data analysis in life-changing application areas



- How well does the data analysis work? Evaluation concerns: accuracy, compute time, & user time
- Cross validation seems to work well, is model agnostic, but can be expensive (multiple analysis runs)
- Can *approximate* general resampling/reweighting
 - We show: is **fast**, **automatic**, and provably **high-quality**

Roadmap

Roadmap

- Cross validation (CV) setup

Roadmap

- Cross validation (CV) setup
- CV Challenge #1: an expensive data analysis

Roadmap

- Cross validation (CV) setup
- CV Challenge #1: an expensive data analysis
- Solution #1: a Taylor series approximation

Roadmap

- Cross validation (CV) setup
- CV Challenge #1: an expensive data analysis
- Solution #1: a Taylor series approximation
("Swiss Army infinitesimal jackknife")

Roadmap

- Cross validation (CV) setup
- CV Challenge #1: an expensive data analysis
- Solution #1: a Taylor series approximation (“Swiss Army infinitesimal jackknife”)
- CV Challenge #2: high dimensions

Roadmap

- Cross validation (CV) setup
- CV Challenge #1: an expensive data analysis
- Solution #1: a Taylor series approximation (“Swiss Army infinitesimal jackknife”)
- CV Challenge #2: high dimensions
- Solution #2: using sparsity (carefully)

Roadmap

- Cross validation (CV) setup
 - CV Challenge #1: an expensive data analysis
 - Solution #1: a Taylor series approximation (“Swiss Army infinitesimal jackknife”)
- CV Challenge #2: high dimensions
- Solution #2: using sparsity (carefully)

Cross validation setup

Cross validation setup

- A data analysis:

Cross validation setup

- A data analysis:

parameters
 θ

Cross validation setup

- A data analysis:

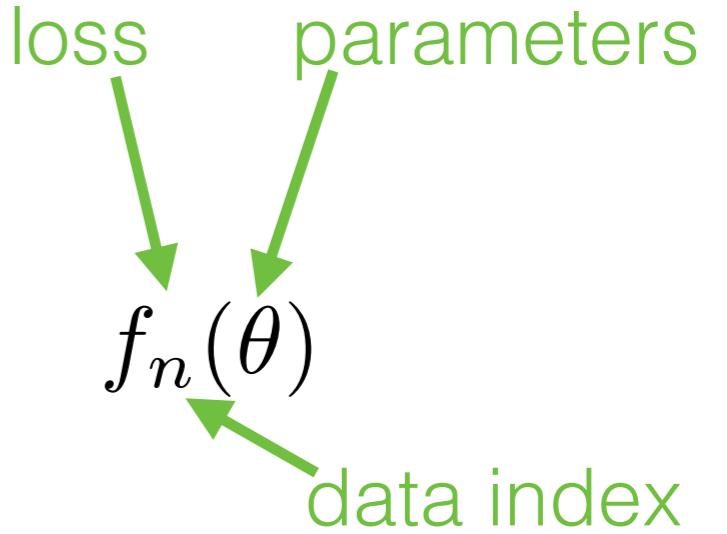
$$f_n(\theta)$$

parameters
data index

```
graph TD; A[parameters] --> B["f_n(<math>\theta</math>)"]; C[data index] --> B
```

Cross validation setup

- A data analysis:



Cross validation setup

- A data analysis:

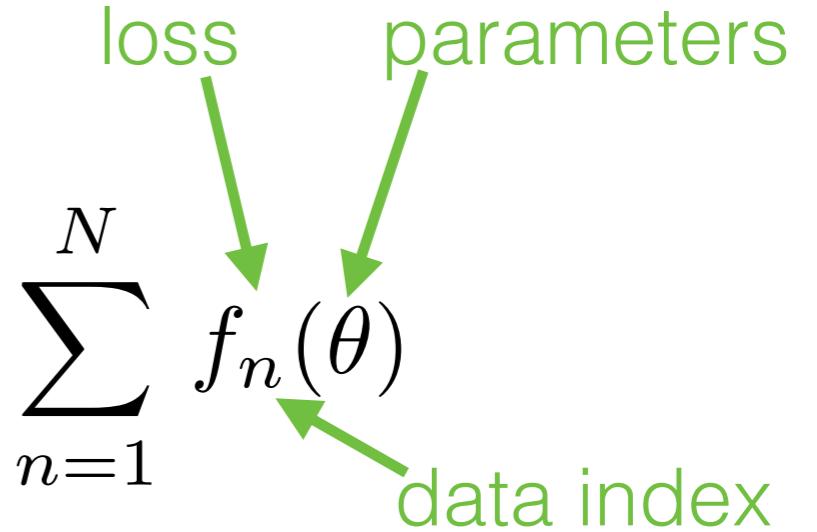
$$\sum_{n=1}^N f_n(\theta)$$

loss parameters
data index

Cross validation setup

- A data analysis:

$$\operatorname{argmin}_{\theta \in \Theta}$$



Cross validation setup

- A data analysis:

$$\operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N f_n(\theta)$$

loss
parameters
data index

Cross validation setup

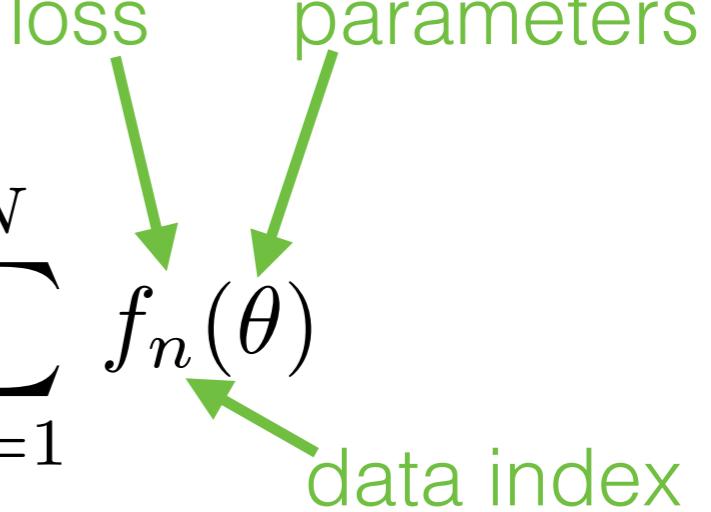
- A data analysis:

$$\hat{\theta} := \underset{\text{estimator}}{\arg\min_{\theta \in \Theta}} \frac{1}{N} \sum_{n=1}^N f_n(\theta)$$

The diagram illustrates the components of the cross-validation setup. It shows three green arrows pointing to the function $f_n(\theta)$: one from the word "loss", one from the word "parameters", and one from the word "data index". The "loss" arrow points to the term f_n , the "parameters" arrow points to the parameter θ , and the "data index" arrow points to the index n in the summation.

Cross validation setup

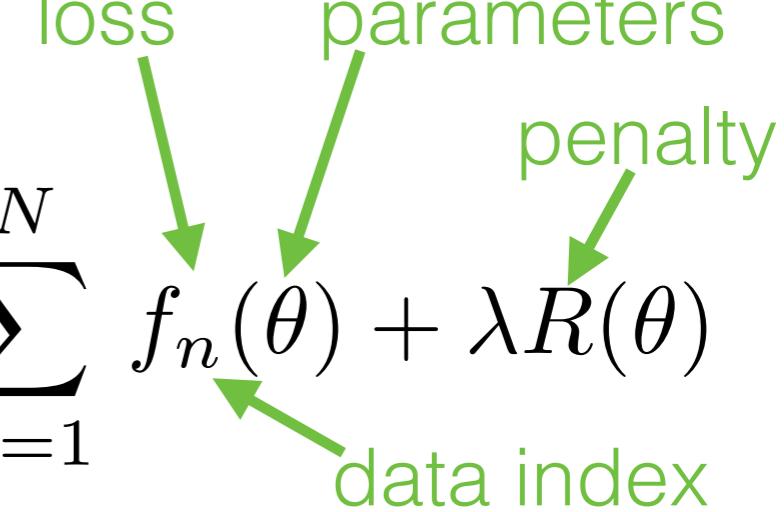
- A data analysis:

$$\hat{\theta} := \underset{\text{estimator}}{\operatorname{argmin}_{\theta \in \Theta}} \frac{1}{N} \sum_{n=1}^N f_n(\theta)$$


- E.g. max likelihood, min loss, M-estimation

Cross validation setup

- A data analysis:

$$\hat{\theta} := \underset{\text{estimator}}{\operatorname{argmin}_{\theta \in \Theta}} \frac{1}{N} \sum_{n=1}^N f_n(\theta) + \lambda R(\theta)$$


- E.g. max likelihood, min loss, M-estimation

Cross validation setup

- A data analysis:

$$\hat{\theta} := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$$

estimator

loss
parameters
penalty
data index

- E.g. max likelihood, min loss, M-estimation

Cross validation setup

- A data analysis:
estimator $\hat{\theta} := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- E.g. max likelihood, min loss, M-estimation

loss
parameters
penalty
data index

The diagram shows four green arrows pointing from labels to specific terms in the mathematical expression. The first arrow points from 'loss' to the term $w_n f_n(\theta)$. The second arrow points from 'parameters' to the term $R(\theta)$. The third arrow points from 'penalty' to the same term $R(\theta)$. The fourth arrow points from 'data index' to the term n in the summation index $n=1$.

Cross validation setup

- A data analysis: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
 - E.g. max likelihood, min loss, M-estimation

loss
parameters
penalty
data index

estimator

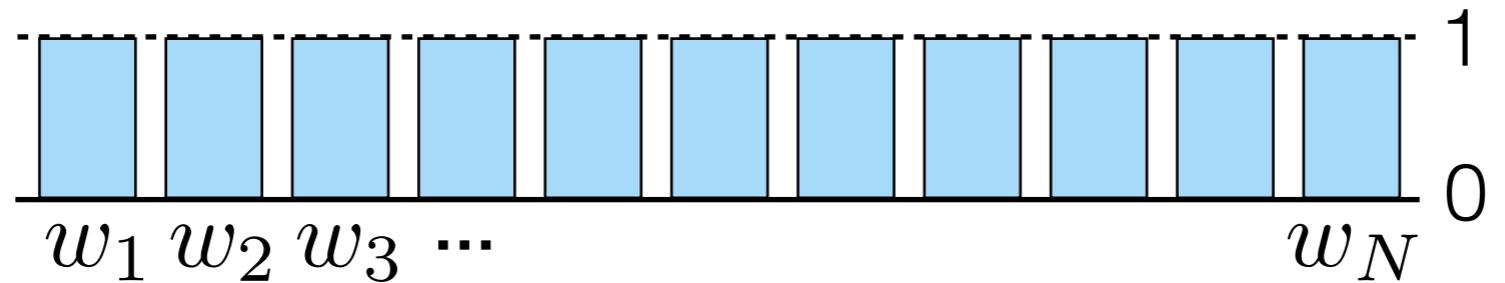
Cross validation setup

- A data analysis: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
 - E.g. max likelihood, min loss, M-estimation

loss
parameters
penalty
data index

Cross validation setup

- A data analysis: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
 - E.g. max likelihood, min loss, M-estimation
 - Original problem: $w = 1_N = (1, \dots, 1, 1, 1, \dots, 1)$



loss
parameters
penalty
data index

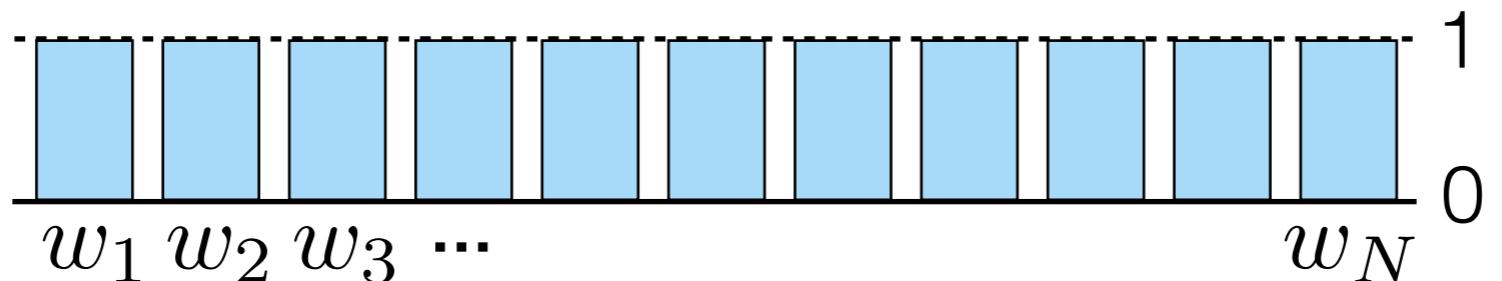
Cross validation setup

- A data analysis: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$

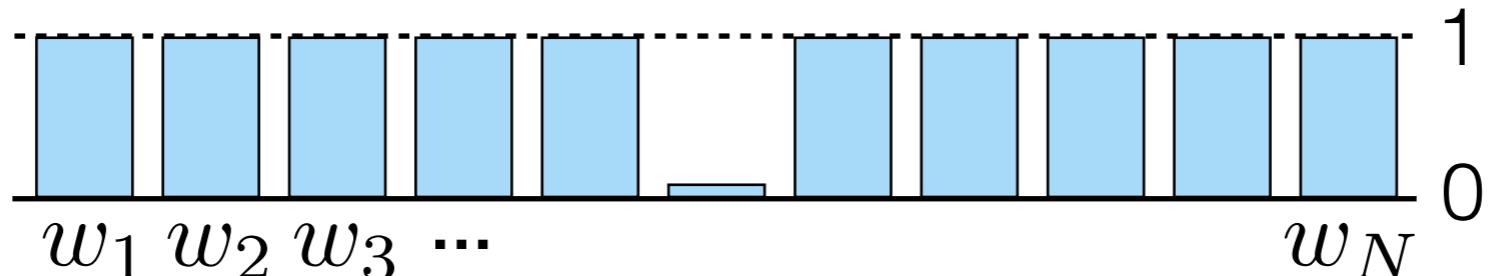
loss
parameters
penalty
data index

- E.g. max likelihood, min loss, M-estimation

- Original problem: $w = 1_N = (1, \dots, 1, 1, 1, \dots, 1)$



- Cross validation (CV): $w = (1, \dots, 1, 0, 1, \dots, 1)$



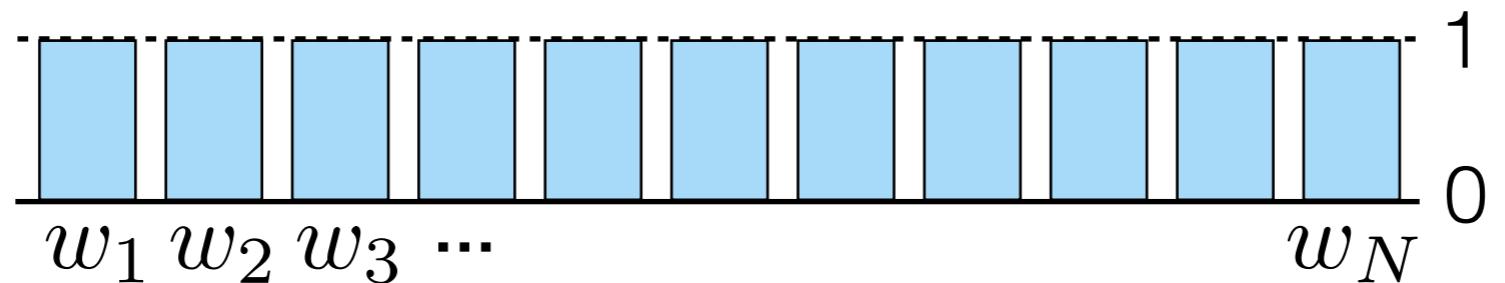
Cross validation setup

- A data analysis: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$

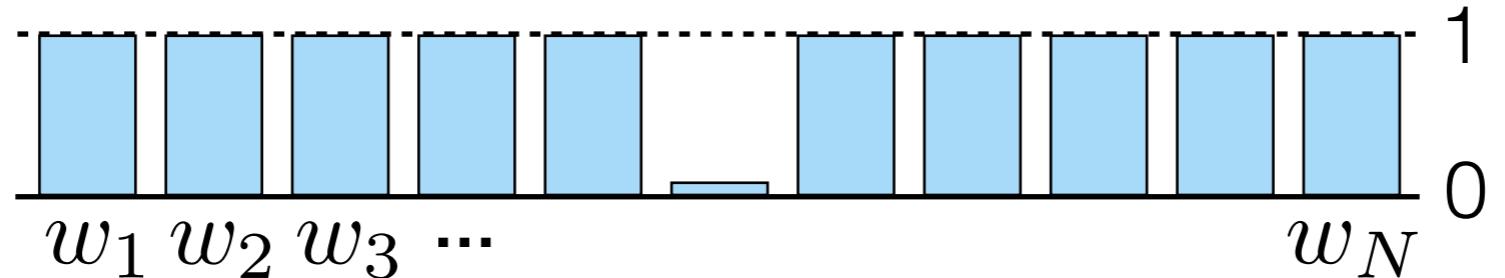
loss
parameters
penalty
data index

- E.g. max likelihood, min loss, M-estimation

- Original problem: $w = 1_N = (1, \dots, 1, 1, 1, \dots, 1)$



- Cross validation (CV): $w = (1, \dots, 1, 0, 1, \dots, 1)$



- CV user time cost: easy to use

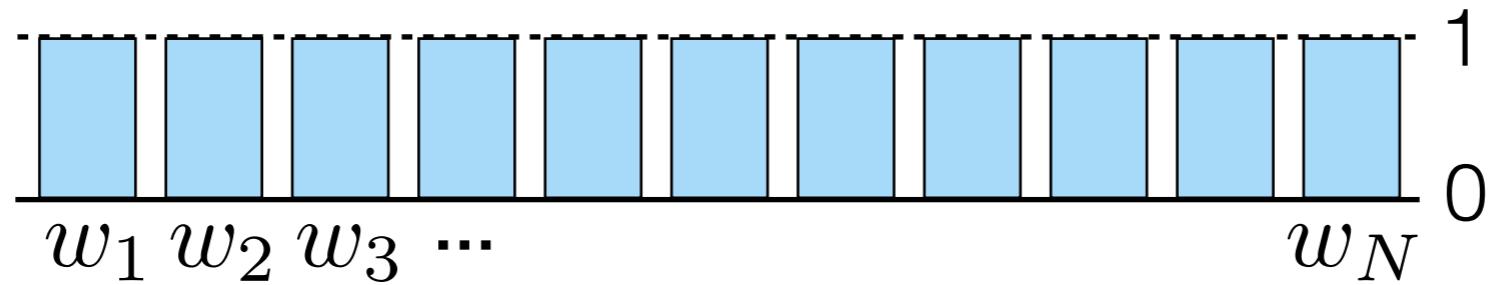
Cross validation setup

• A data analysis: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$

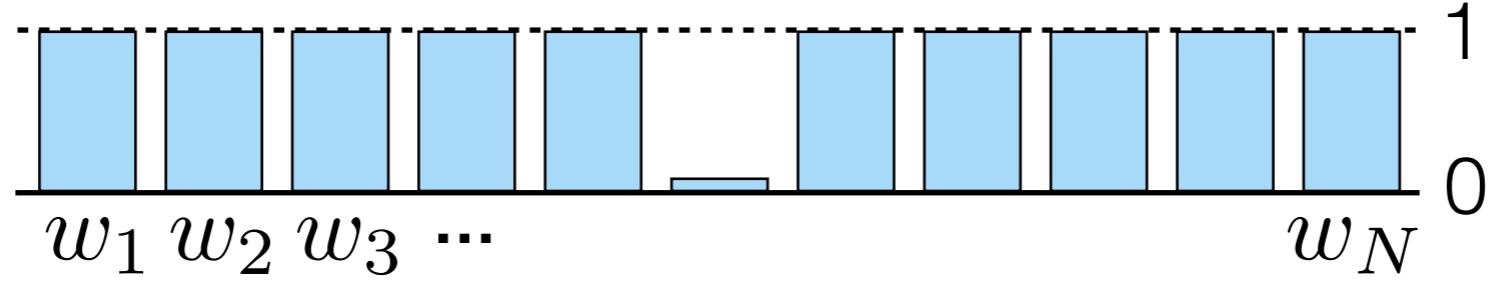
estimator

loss
parameters
penalty
data index

- E.g. max likelihood, min loss, M-estimation
- Original problem: $w = 1_N = (1, \dots, 1, 1, 1, \dots, 1)$



- Cross validation (CV): $w = (1, \dots, 1, 0, 1, \dots, 1)$



- CV user time cost: easy to use
- CV compute time cost: many \times (cost of 1 data analysis)

Approximation for large data

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$

Approximation for large data

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Idea: (first-order) Taylor expansion around $w = 1_N$

Approximation for large data

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Idea: (first-order) Taylor expansion around $w = 1_N$

$$\hat{\theta}_{IJ}(w) := \hat{\theta}(1_N) + \left. \frac{d\hat{\theta}(w)}{dw^T} \right|_{w=1_N} (w - 1_N)$$

Approximation for large data

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Idea: (first-order) Taylor expansion around $w = 1_N$

$$\begin{aligned}\hat{\theta}_{IJ}(w) &:= \hat{\theta}(1_N) + \left. \frac{d\hat{\theta}(w)}{dw^T} \right|_{w=1_N} (w - 1_N) \\ &= \hat{\theta}(1_N) - H(1_N)^{-1} G(1_N) (w - 1_N)\end{aligned}$$

Approximation for large data

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Idea: (first-order) Taylor expansion around $w = 1_N$

$$\begin{aligned}\hat{\theta}_{IJ}(w) &:= \hat{\theta}(1_N) + \left. \frac{d\hat{\theta}(w)}{dw^T} \right|_{w=1_N} (w - 1_N) \\ &= \hat{\theta}(1_N) - H(1_N)^{-1} G(1_N)(w - 1_N)\end{aligned}$$

*H: Hessian of
the objective*

$$H_{d_1 d_2} = \frac{\partial^2}{\partial \theta_{d_1} \partial \theta_{d_2}} \text{objective}$$

Approximation for large data

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Idea: (first-order) Taylor expansion around $w = 1_N$

$$\begin{aligned}\hat{\theta}_{IJ}(w) &:= \hat{\theta}(1_N) + \left. \frac{d\hat{\theta}(w)}{dw^T} \right|_{w=1_N} (w - 1_N) \\ &= \hat{\theta}(1_N) - H(1_N)^{-1} G(1_N)(w - 1_N)\end{aligned}$$

H : Hessian of
the objective G : matrix of
 gradients

$$H_{d_1 d_2} = \frac{\partial^2}{\partial \theta_{d_1} \partial \theta_{d_2}} \text{objective} \quad G_{dn} = \frac{\partial^2}{\partial \theta_d \partial w_n} \text{objective}$$

Approximation for large data

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Idea: (first-order) Taylor expansion around $w = 1_N$

$$\begin{aligned}\hat{\theta}_{IJ}(w) &:= \hat{\theta}(1_N) + \left. \frac{d\hat{\theta}(w)}{dw^T} \right|_{w=1_N} (w - 1_N) \\ &= \hat{\theta}(1_N) - H(1_N)^{-1} G(1_N)(w - 1_N)\end{aligned}$$

H : Hessian of
the objective G : matrix of
gradients

$$H_{d_1 d_2} = \frac{\partial^2}{\partial \theta_{d_1} \partial \theta_{d_2}} \text{objective} \quad G_{dn} = \frac{\partial^2}{\partial \theta_d \partial w_n} \text{objective}$$

$$G = \frac{1}{N} \begin{pmatrix} & & & & \\ & | & & & | \\ \nabla_{\theta} f_1 & \cdots & & & \nabla_{\theta} f_N \\ & | & & & | \end{pmatrix}$$

Approximation for large data

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Idea: (first-order) Taylor expansion around $w = 1_N$

$$\begin{aligned}\hat{\theta}_{IJ}(w) &:= \hat{\theta}(1_N) + \left. \frac{d\hat{\theta}(w)}{dw^T} \right|_{w=1_N} (w - 1_N) \\ &= \hat{\theta}(1_N) - H(1_N)^{-1} G(1_N) (w - 1_N)\end{aligned}$$

H: Hessian of *G*: matrix of
the objective gradients

- Just one data analysis

Approximation for large data

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Idea: (first-order) Taylor expansion around $w = 1_N$

$$\begin{aligned}\hat{\theta}_{IJ}(w) &:= \hat{\theta}(1_N) + \left. \frac{d\hat{\theta}(w)}{dw^T} \right|_{w=1_N} (w - 1_N) \\ &= \hat{\theta}(1_N) - H(1_N)^{-1} G(1_N) (w - 1_N)\end{aligned}$$

*H: Hessian of G: matrix of
the objective gradients*

- Just one data analysis; requires 2x differentiability

Approximation for large data

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Idea: (first-order) Taylor expansion around $w = 1_N$

$$\begin{aligned}\hat{\theta}_{IJ}(w) &:= \hat{\theta}(1_N) + \left. \frac{d\hat{\theta}(w)}{dw^T} \right|_{w=1_N} (w - 1_N) \\ &= \hat{\theta}(1_N) - H(1_N)^{-1} G(1_N)(w - 1_N)\end{aligned}$$

*H: Hessian of G: matrix of
the objective gradients*

- Just one data analysis; requires 2x differentiability
- Local sensitivity analysis / influence scores

Approximation for large data

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Idea: (first-order) Taylor expansion around $w = 1_N$

$$\begin{aligned}\hat{\theta}_{IJ}(w) &:= \hat{\theta}(1_N) + \left. \frac{d\hat{\theta}(w)}{dw^T} \right|_{w=1_N} (w - 1_N) \\ &= \hat{\theta}(1_N) - H(1_N)^{-1} G(1_N)(w - 1_N)\end{aligned}$$

H: Hessian of *G*: matrix of
the objective gradients

- Just one data analysis; requires 2x differentiability
- Local sensitivity analysis / influence scores
- Other linear approximations

Approximation for large data

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Idea: (first-order) Taylor expansion around $w = 1_N$

$$\begin{aligned}\hat{\theta}_{IJ}(w) &:= \hat{\theta}(1_N) + \left. \frac{d\hat{\theta}(w)}{dw^T} \right|_{w=1_N} (w - 1_N) \\ &= \hat{\theta}(1_N) - H(1_N)^{-1} G(1_N)(w - 1_N)\end{aligned}$$

*H: Hessian of G: matrix of
the objective gradients*

- Just one data analysis; requires 2x differentiability
- Local sensitivity analysis / influence scores
- Other linear approximations
- What about higher orders? See our latest on arXiv

[Giordano, Stephenson, Liu, Jordan, Broderick 2019; Giordano, Broderick, Jordan 2018;
Giordano, Jordan, Broderick 2019;

Approximation for large data

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Idea: (first-order) Taylor expansion around $w = 1_N$

$$\hat{\theta}_{IJ}(w) := \hat{\theta}(1_N) + \left. \frac{d\hat{\theta}(w)}{dw^T} \right|_{w=1_N} (w - 1_N)$$

IJ: infinitesimal
jackknife

$$= \hat{\theta}(1_N) - H(1_N)^{-1} G(1_N)(w - 1_N)$$

H : Hessian of G : matrix of
the objective gradients

- Just one data analysis; requires 2x differentiability
- Local sensitivity analysis / influence scores
- Other linear approximations
- What about higher orders? See our latest on arXiv

[Giordano, Stephenson, Liu, Jordan, Broderick 2019; Giordano, Broderick, Jordan 2018;
Giordano, Jordan, Broderick 2019; Jaeckel 1972; Clarke 1983;

Approximation for large data

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Idea: (first-order) Taylor expansion around $w = 1_N$

$$\hat{\theta}_{IJ}(w) := \hat{\theta}(1_N) + \left. \frac{d\hat{\theta}(w)}{dw^T} \right|_{w=1_N} (w - 1_N)$$

IJ: infinitesimal
jackknife

$$= \hat{\theta}(1_N) - H(1_N)^{-1} G(1_N)(w - 1_N)$$

H : Hessian of G : matrix of
the objective gradients

Our contributions:

Approximation for large data

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Idea: (first-order) Taylor expansion around $w = 1_N$

$$\hat{\theta}_{IJ}(w) := \hat{\theta}(1_N) + \left. \frac{d\hat{\theta}(w)}{dw^T} \right|_{w=1_N} (w - 1_N)$$

IJ: infinitesimal
jackknife

$$= \hat{\theta}(1_N) - H(1_N)^{-1} G(1_N)(w - 1_N)$$

H : Hessian of G : matrix of
the objective gradients

Our contributions:

- Practical automatic differentiation wrappers

Approximation for large data

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Idea: (first-order) Taylor expansion around $w = 1_N$

$$\hat{\theta}_{IJ}(w) := \hat{\theta}(1_N) + \left. \frac{d\hat{\theta}(w)}{dw^T} \right|_{w=1_N} (w - 1_N)$$

IJ: infinitesimal
jackknife

$$= \hat{\theta}(1_N) - H(1_N)^{-1} G(1_N)(w - 1_N)$$

H : Hessian of G : matrix of
the objective gradients

Our contributions:

- Practical automatic differentiation wrappers
- Practical theory assumptions (vs. e.g. term-wise bounded gradients)

Approximation for large data

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Idea: (first-order) Taylor expansion around $w = 1_N$

$$\hat{\theta}_{IJ}(w) := \hat{\theta}(1_N) + \left. \frac{d\hat{\theta}(w)}{dw^T} \right|_{w=1_N} (w - 1_N)$$

IJ: infinitesimal
jackknife

$$= \hat{\theta}(1_N) - H(1_N)^{-1} G(1_N)(w - 1_N)$$

H : Hessian of G : matrix of
the objective gradients

Our contributions:

- Practical automatic differentiation wrappers
- Practical theory assumptions (vs. e.g. term-wise bounded gradients)
- Finite-sample theoretical bounds on quality

A “Swiss Army infinitesimal jackknife”

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Idea: (first-order) Taylor expansion around $w = 1_N$

$$\hat{\theta}_{IJ}(w) := \hat{\theta}(1_N) + \left. \frac{d\hat{\theta}(w)}{dw^T} \right|_{w=1_N} (w - 1_N)$$

IJ: infinitesimal
jackknife

$$= \hat{\theta}(1_N) - H(1_N)^{-1} G(1_N)(w - 1_N)$$

H : Hessian of G : matrix of
the objective gradients

Our contributions:

- Practical automatic differentiation wrappers
- Practical theory assumptions (vs. e.g. term-wise bounded gradients)
- Finite-sample theoretical bounds on quality

Swiss Army IJ: Easy to use

- In Python:
 - **autograd** for autodiff <https://github.com/HIPS/autograd>
 - **vittles** for sensitivity <https://github.com/rgiordan/vittles>

Swiss Army IJ: Easy to use

- In Python:
 - `autograd` for autodiff <https://github.com/HIPS/autograd>
 - `vittles` for sensitivity <https://github.com/rgiordan/vittles>

```
import autograd.numpy as np
import vittles

def objective_fun(theta, w):
    # Your objective here

# optimize...
theta_opt = optimize(lambda theta: objective_fun(theta,
w_ones))

predictor =
vittles.HyperparameterSensitivityLinearApproximation(
objective_fun, theta_opt, w_ones)

theta_IJ = predictor.predict_opt_par_from_hyper_par(w_new)
```

Theory assumptions (no *term-wise* bounded gradients)

Theory assumptions (no *term-wise* bounded gradients)

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$

Theory assumptions (no *term-wise* bounded gradients)

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Consider:

Theory assumptions (no *term-wise* bounded gradients)

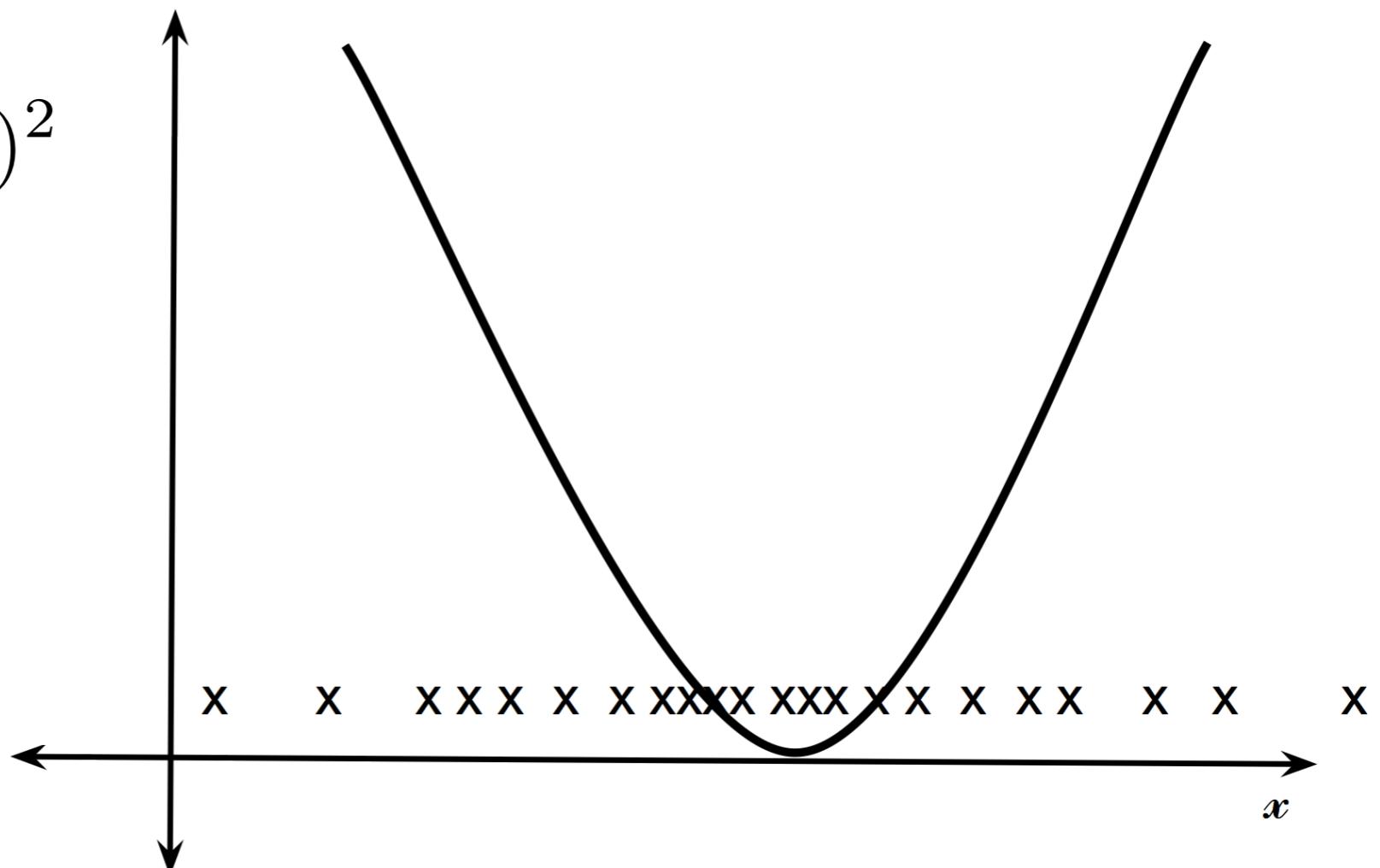
- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Consider:

$$f_n(\theta) = \frac{1}{2} (x_n - \theta)^2$$

Theory assumptions (no *term-wise* bounded gradients)

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Consider:

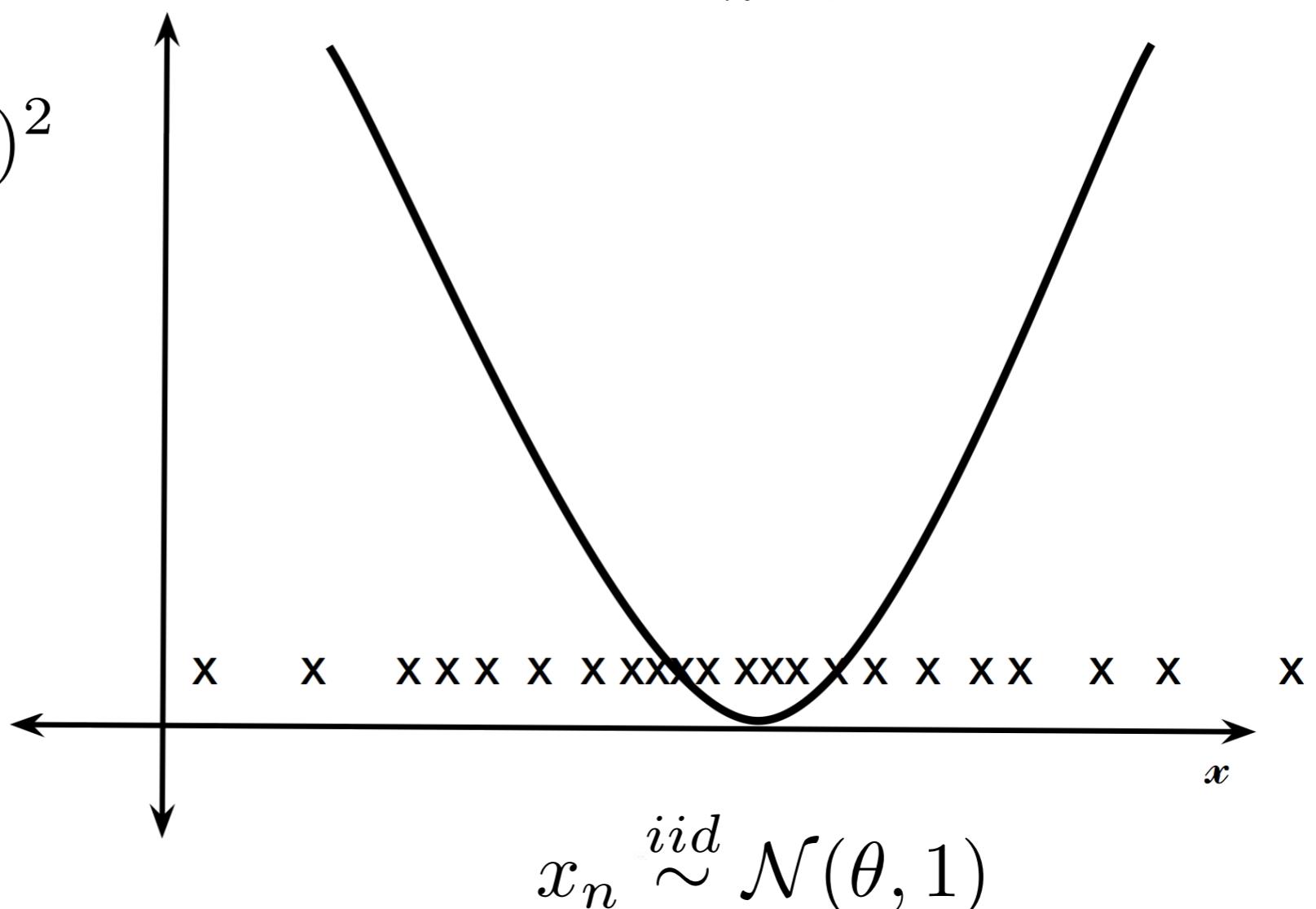
$$f_n(\theta) = \frac{1}{2} (x_n - \theta)^2$$



Theory assumptions (no *term-wise* bounded gradients)

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Consider:

$$f_n(\theta) = \frac{1}{2} (x_n - \theta)^2$$

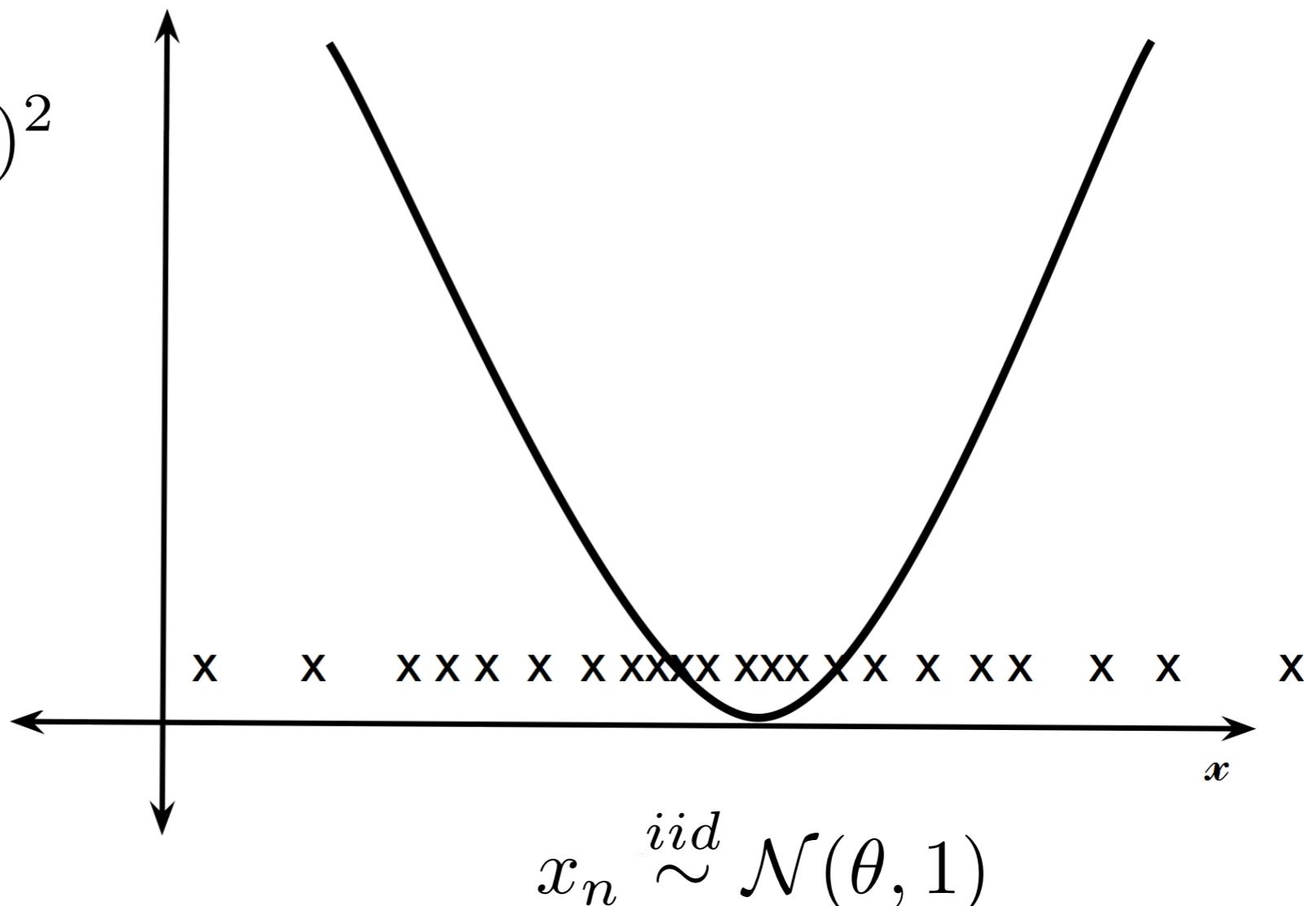


Theory assumptions (no *term-wise* bounded gradients)

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Consider:

$$f_n(\theta) = \frac{1}{2}(x_n - \theta)^2$$

$$\nabla_\theta f_n(\theta) = x_n - \theta$$



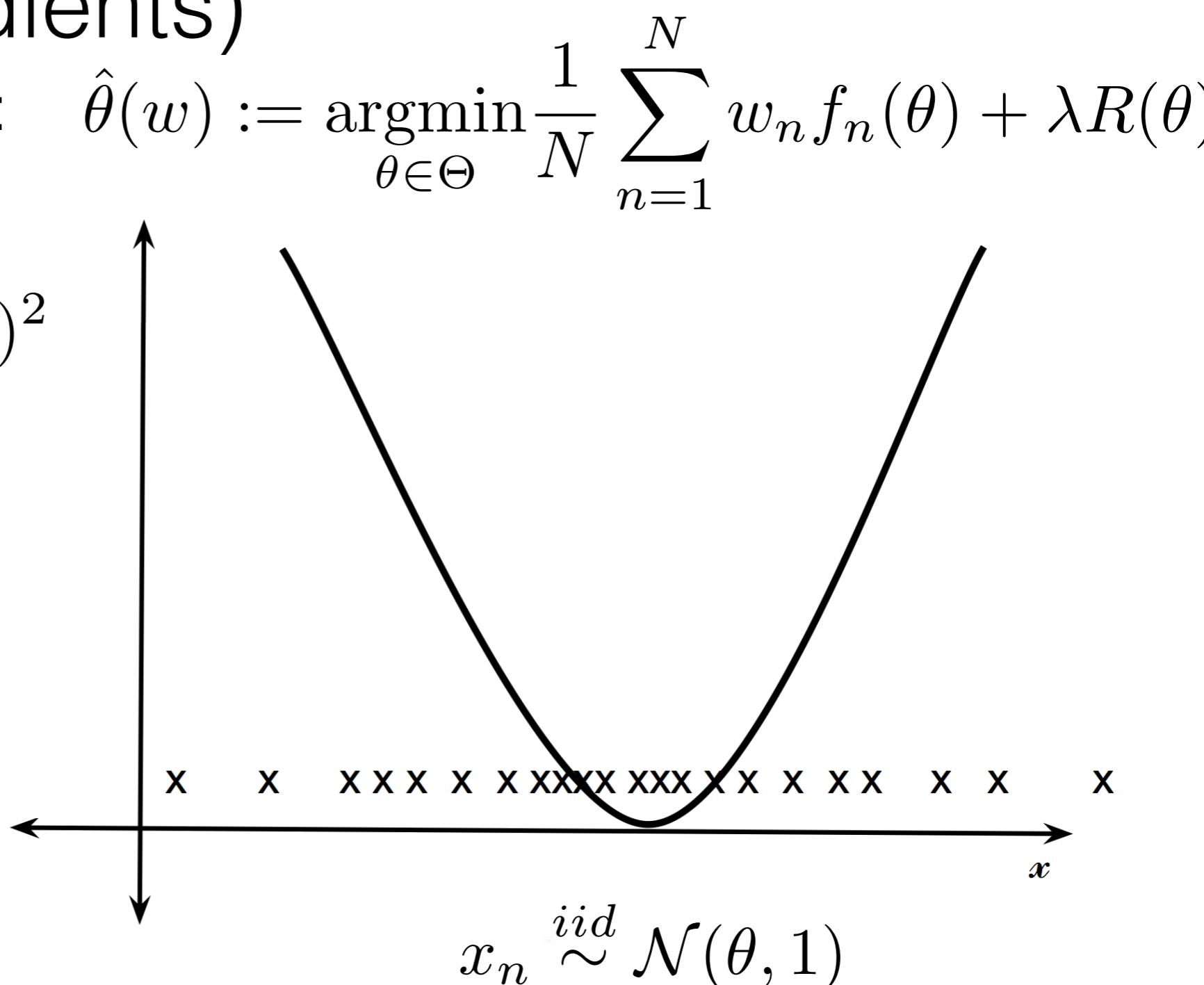
Theory assumptions (no *term-wise* bounded gradients)

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Consider:

$$f_n(\theta) = \frac{1}{2} (x_n - \theta)^2$$

$$\nabla_\theta f_n(\theta) = x_n - \theta$$

$$\max_{n \in 1:N} |x_n - \theta|$$



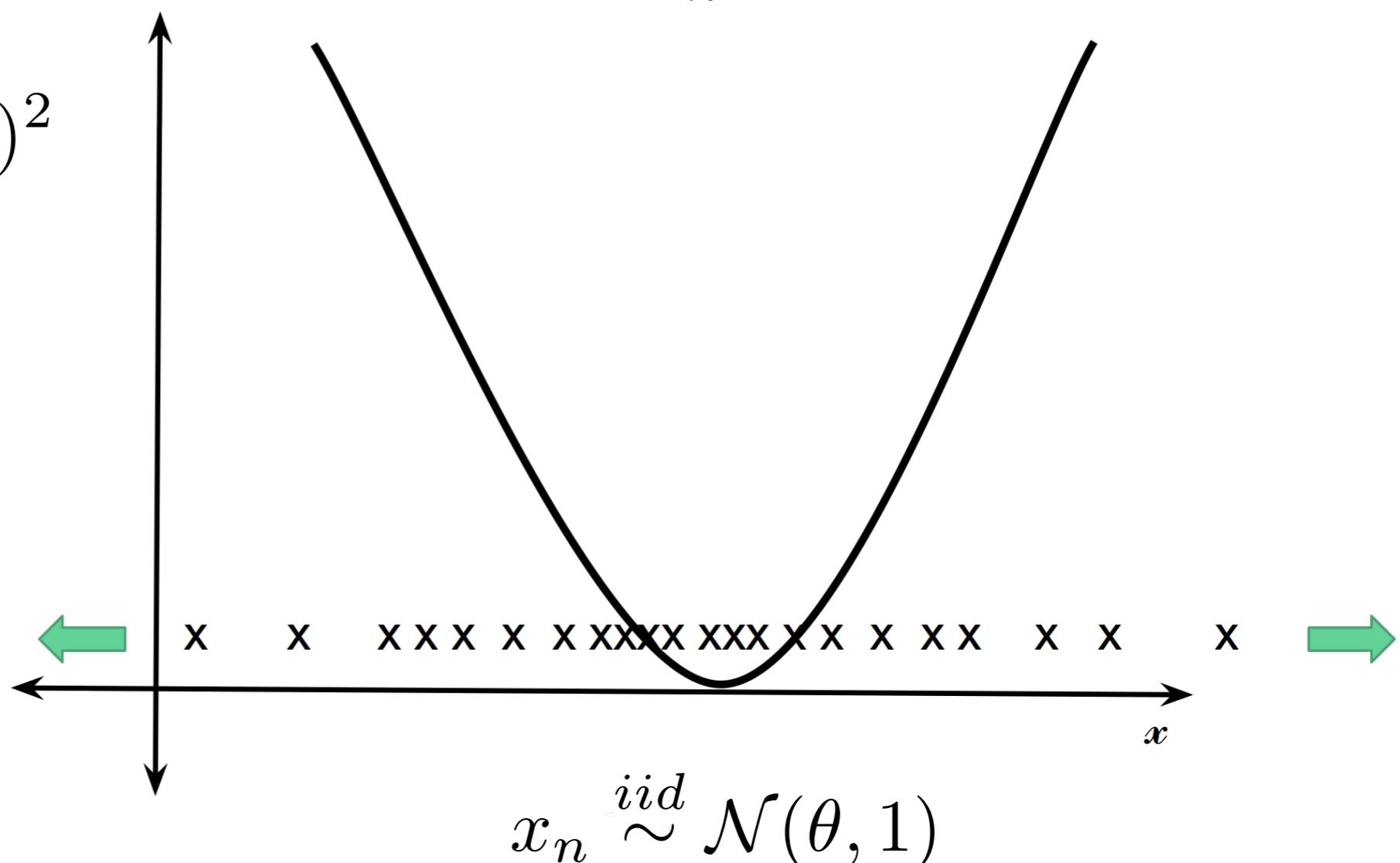
Theory assumptions (no *term-wise* bounded gradients)

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Consider:

$$f_n(\theta) = \frac{1}{2} (x_n - \theta)^2$$

$$\nabla_\theta f_n(\theta) = x_n - \theta$$

$$\max_{n \in 1:N} |x_n - \theta|$$



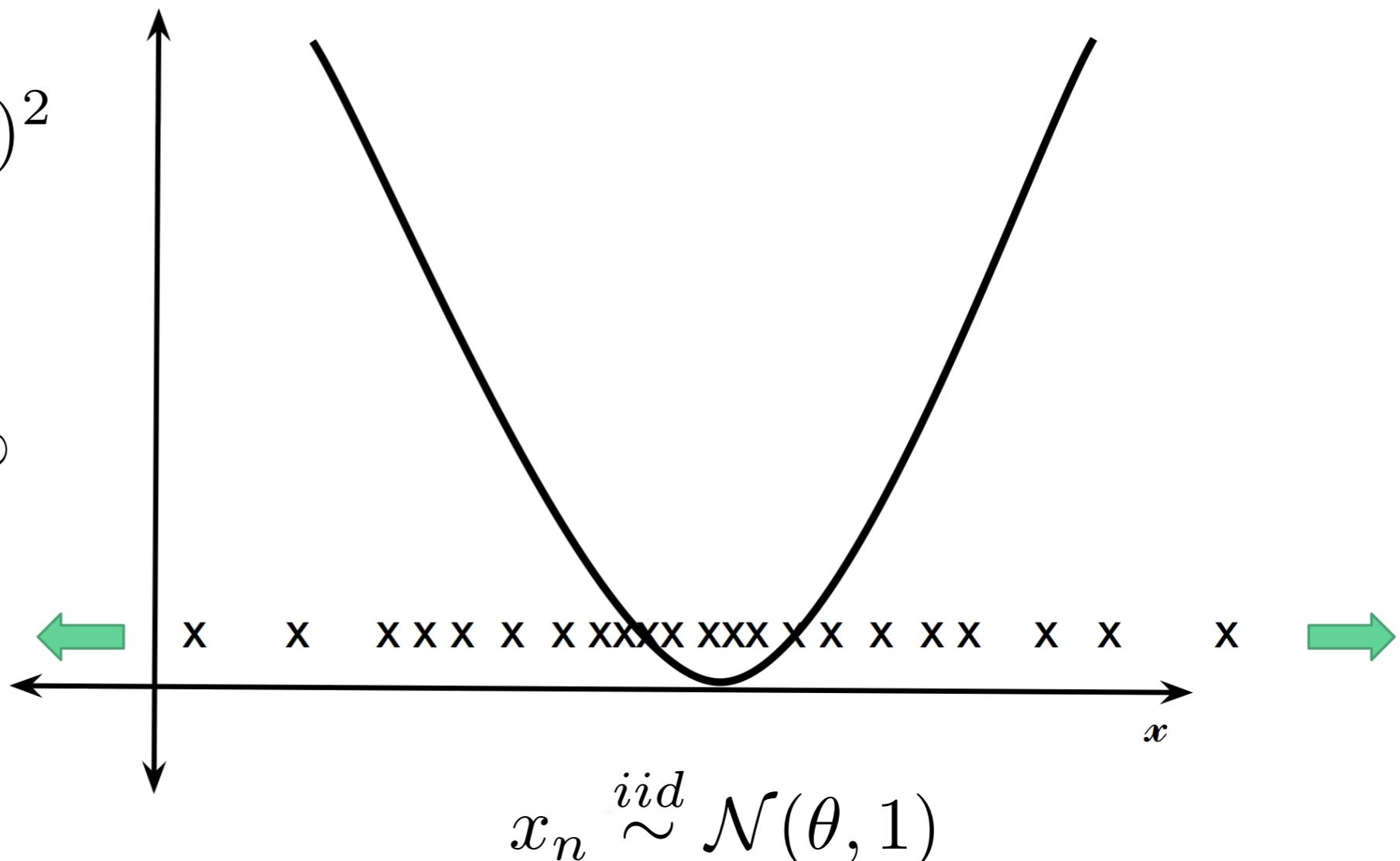
Theory assumptions (no *term-wise* bounded gradients)

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Consider:

$$f_n(\theta) = \frac{1}{2}(x_n - \theta)^2$$

$$\nabla_\theta f_n(\theta) = x_n - \theta$$

$$\max_{n \in 1:N} |x_n - \theta| \rightarrow \infty$$



Theory assumptions (no *term-wise* bounded gradients)

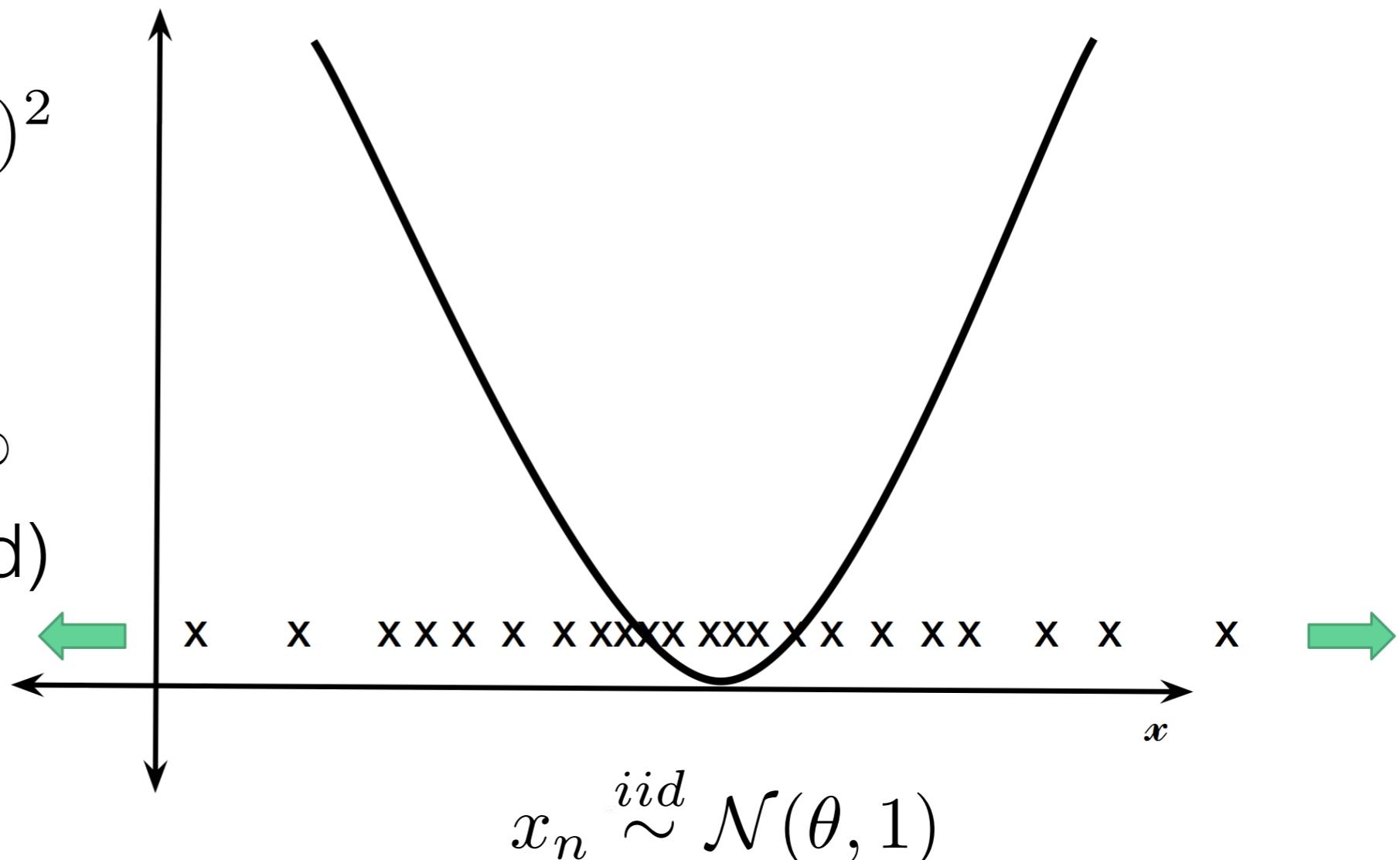
- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Consider:

$$f_n(\theta) = \frac{1}{2} (x_n - \theta)^2$$

$$\nabla_\theta f_n(\theta) = x_n - \theta$$

$$\max_{n \in 1:N} |x_n - \theta| \rightarrow \infty$$

(not bounded)



Theory assumptions (no *term-wise* bounded gradients)

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Consider:

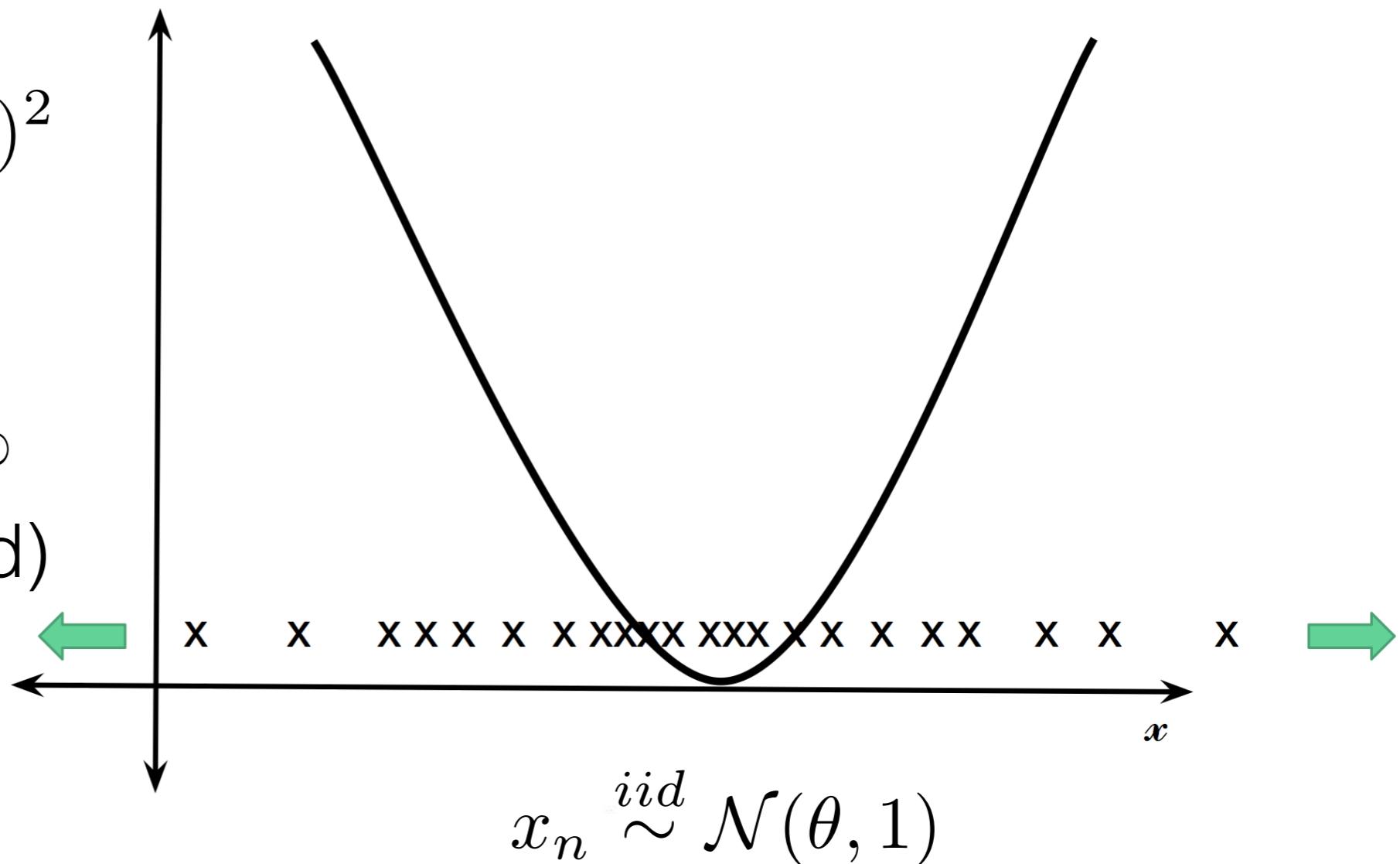
$$f_n(\theta) = \frac{1}{2}(x_n - \theta)^2$$

$$\nabla_{\theta} f_n(\theta) = x_n - \theta$$

$$\max_{n \in 1:N} |x_n - \theta| \rightarrow \infty$$

(not bounded)

$$\frac{1}{N} \sum_{n=1}^N (x_n - \theta)^2$$



Theory assumptions (no term-wise bounded gradients)

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Consider:

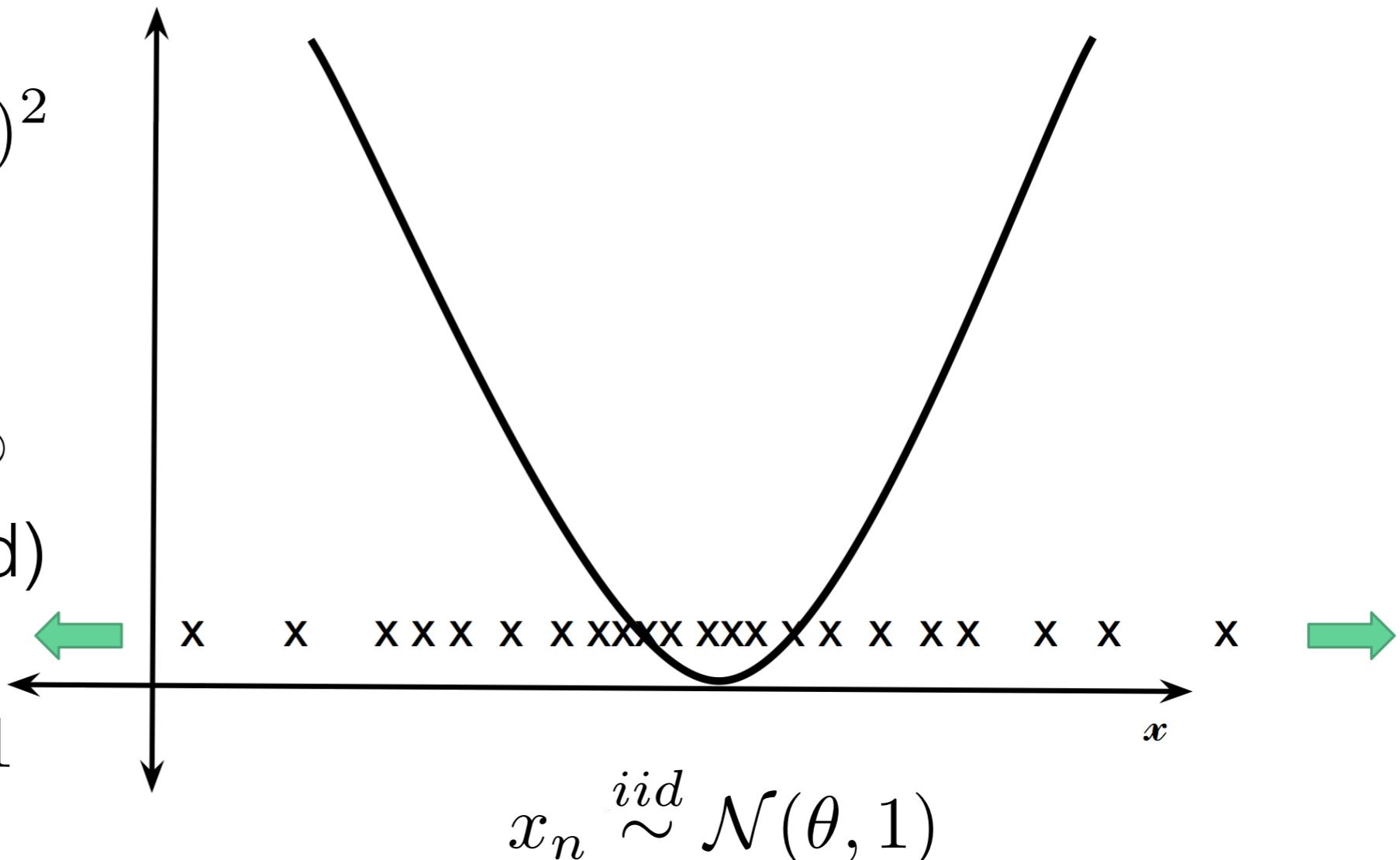
$$f_n(\theta) = \frac{1}{2}(x_n - \theta)^2$$

$$\nabla_\theta f_n(\theta) = x_n - \theta$$

$$\max_{n \in 1:N} |x_n - \theta| \rightarrow \infty$$

(not bounded)

$$\frac{1}{N} \sum_{n=1}^N (x_n - \theta)^2 \rightarrow 1$$



Theory assumptions (no term-wise bounded gradients)

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Consider:

$$f_n(\theta) = \frac{1}{2}(x_n - \theta)^2$$

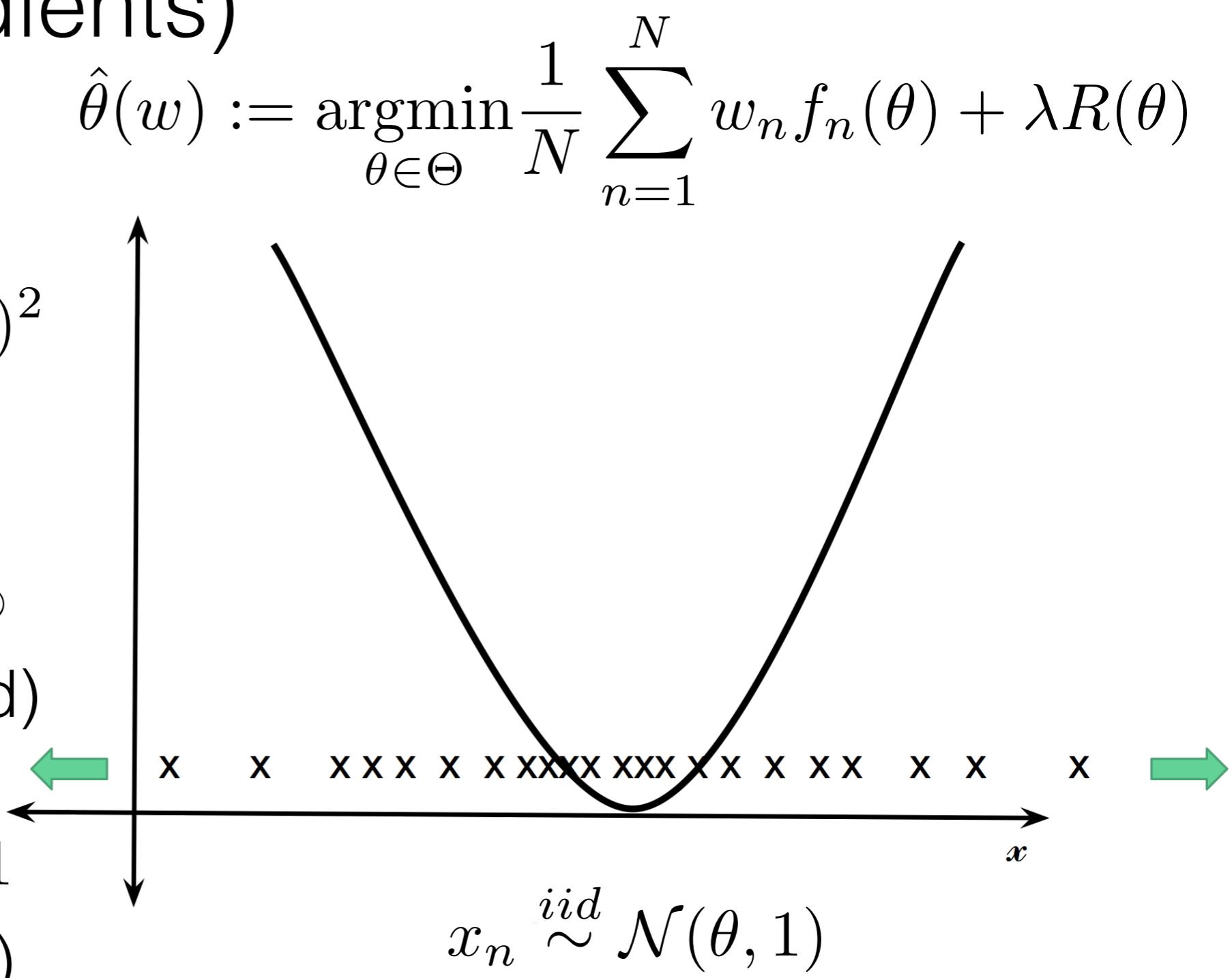
$$\nabla_{\theta} f_n(\theta) = x_n - \theta$$

$$\max_{n \in 1:N} |x_n - \theta| \rightarrow \infty$$

(not bounded)

$$\frac{1}{N} \sum_{n=1}^N (x_n - \theta)^2 \rightarrow 1$$

(is bounded)



Theory assumptions (no *term-wise* bounded gradients)

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$
- Consider:

$$f_n(\theta) = \frac{1}{2}(x_n - \theta)^2$$

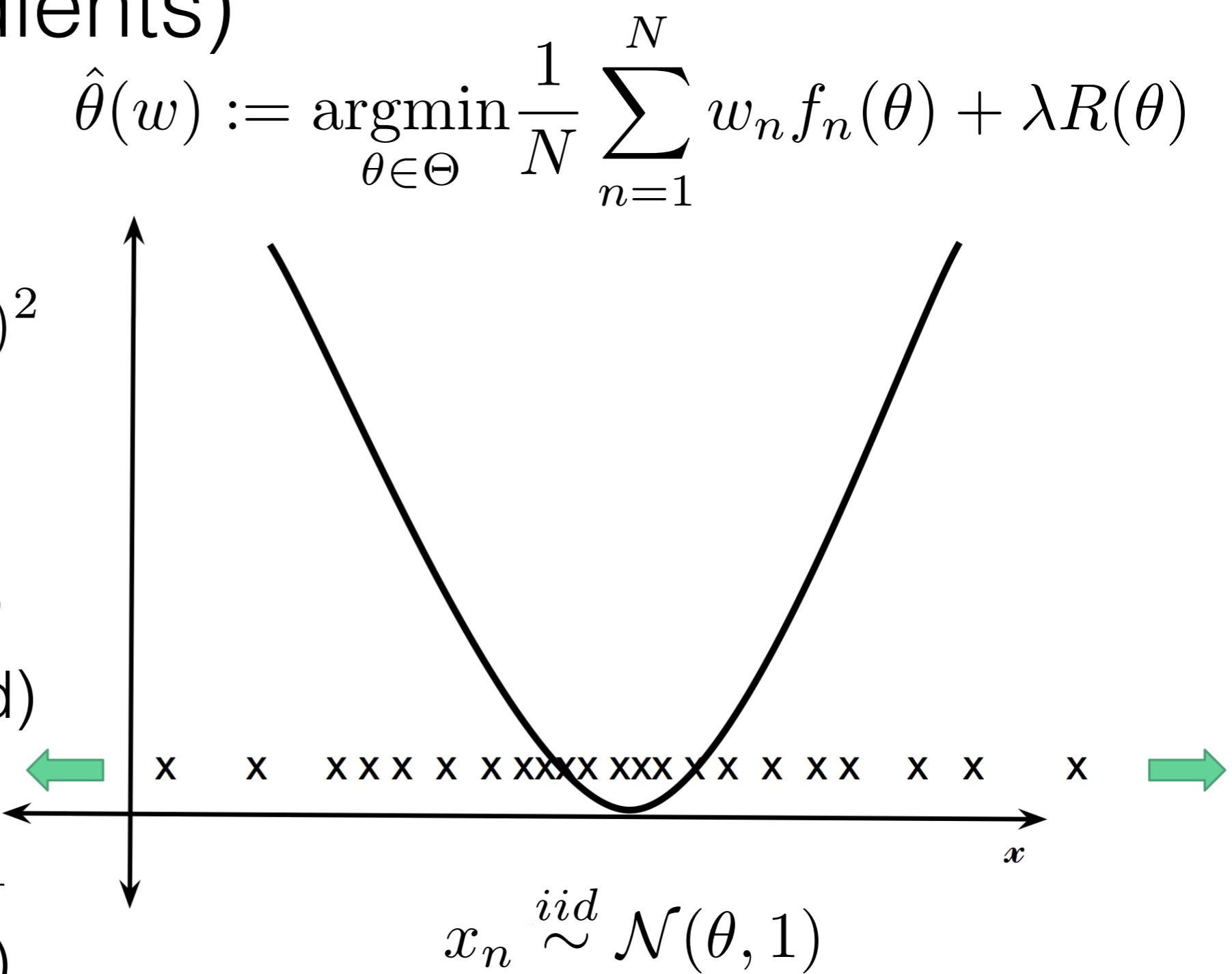
$$\nabla_\theta f_n(\theta) = x_n - \theta$$

$$\max_{n \in 1:N} |x_n - \theta| \rightarrow \infty$$

(not bounded)

$$\frac{1}{N} \sum_{n=1}^N (x_n - \theta)^2 \rightarrow 1$$

(is bounded)



- Our assumption: bounded *sample variances* of gradients and Hessians

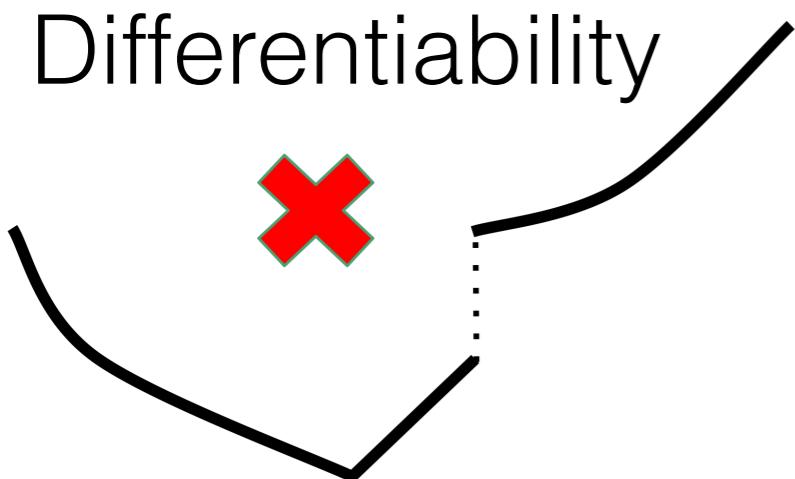
More theory assumptions

- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$

More theory assumptions

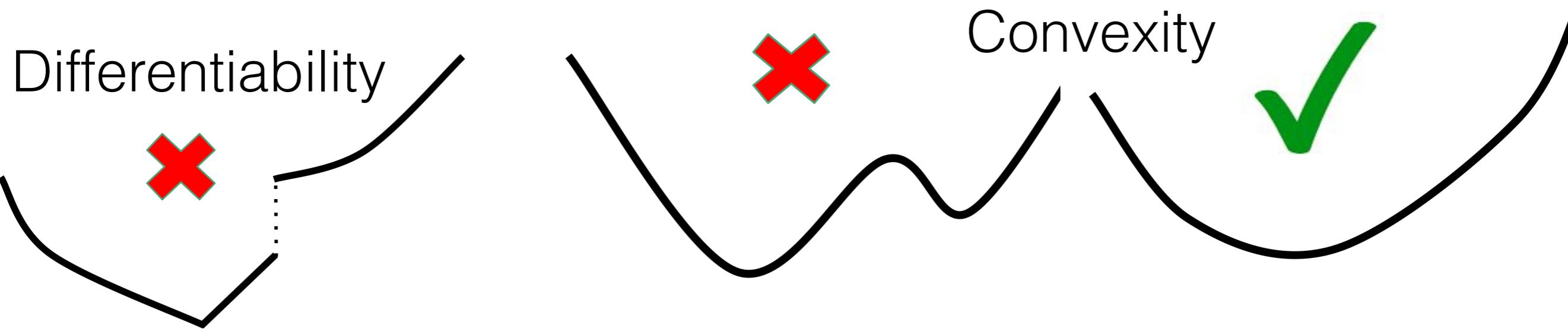
- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$

Differentiability



More theory assumptions

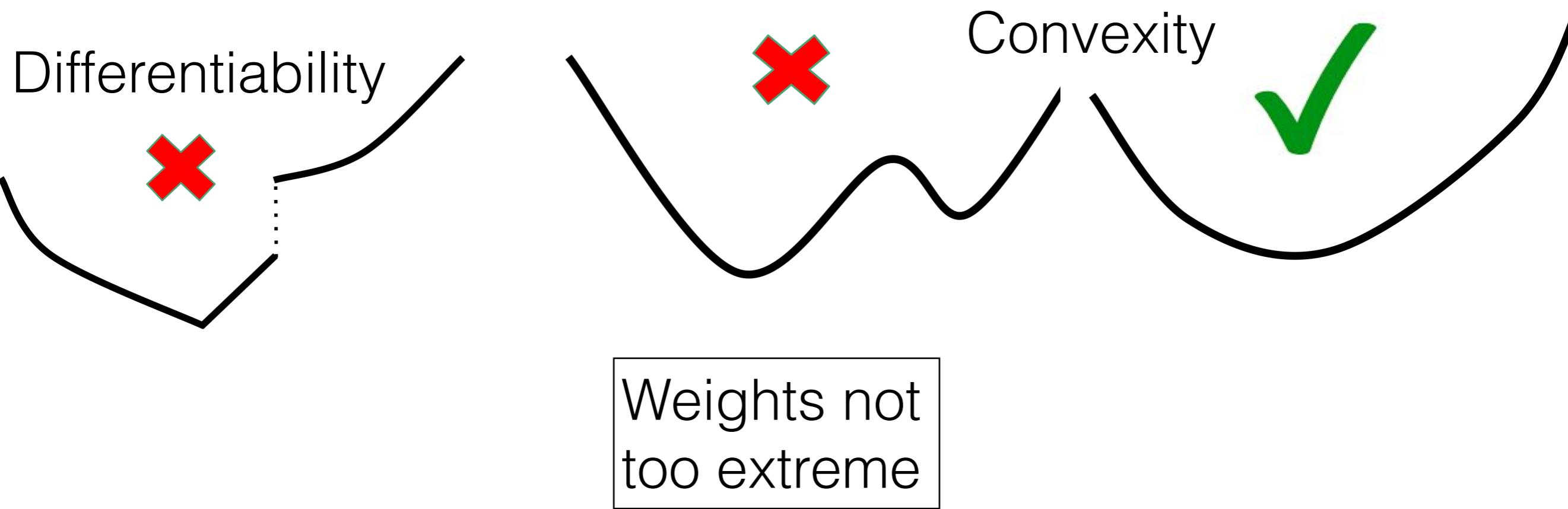
- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$



More theory assumptions

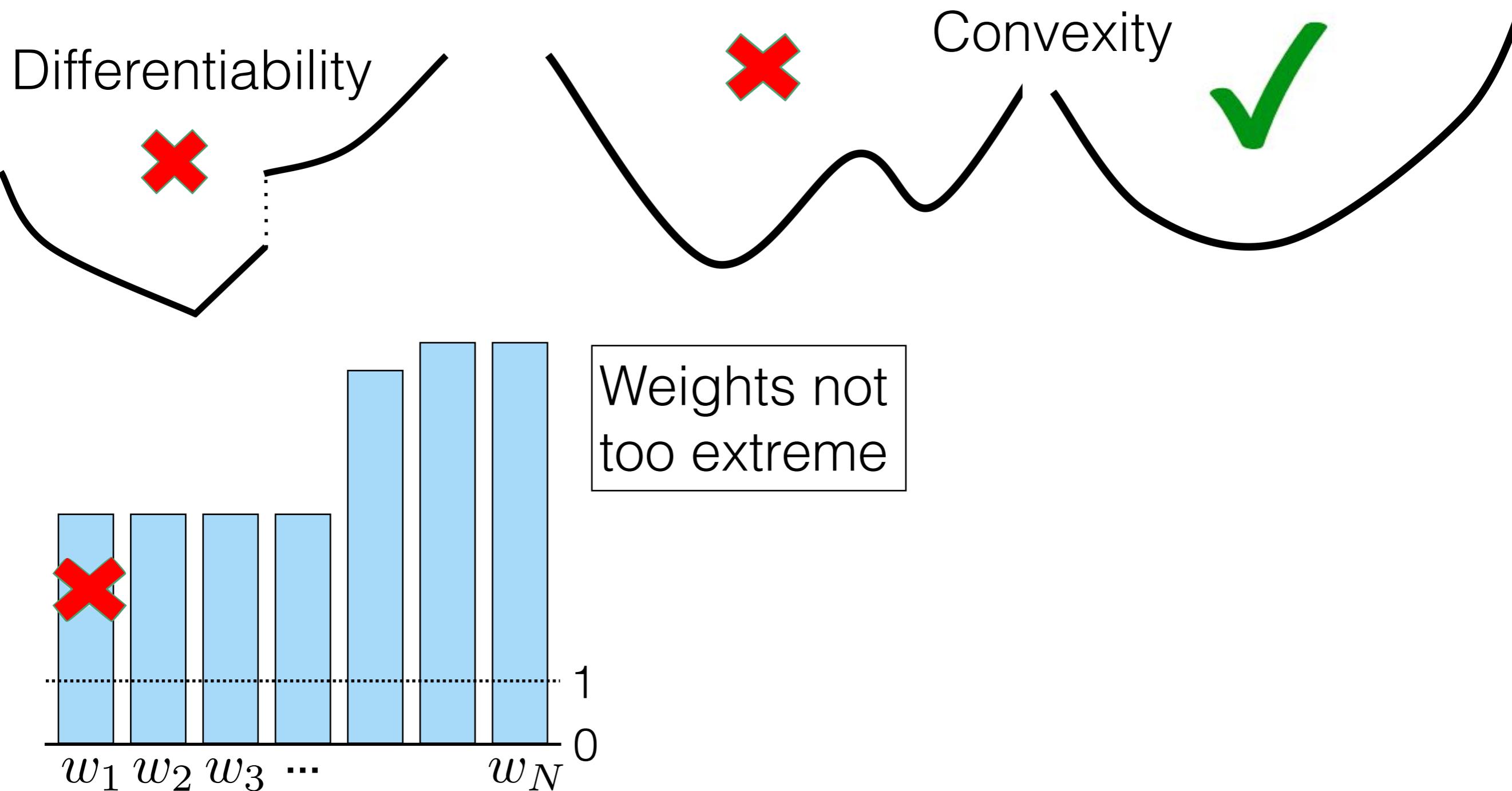
- Exact CV round:

$$\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$$



More theory assumptions

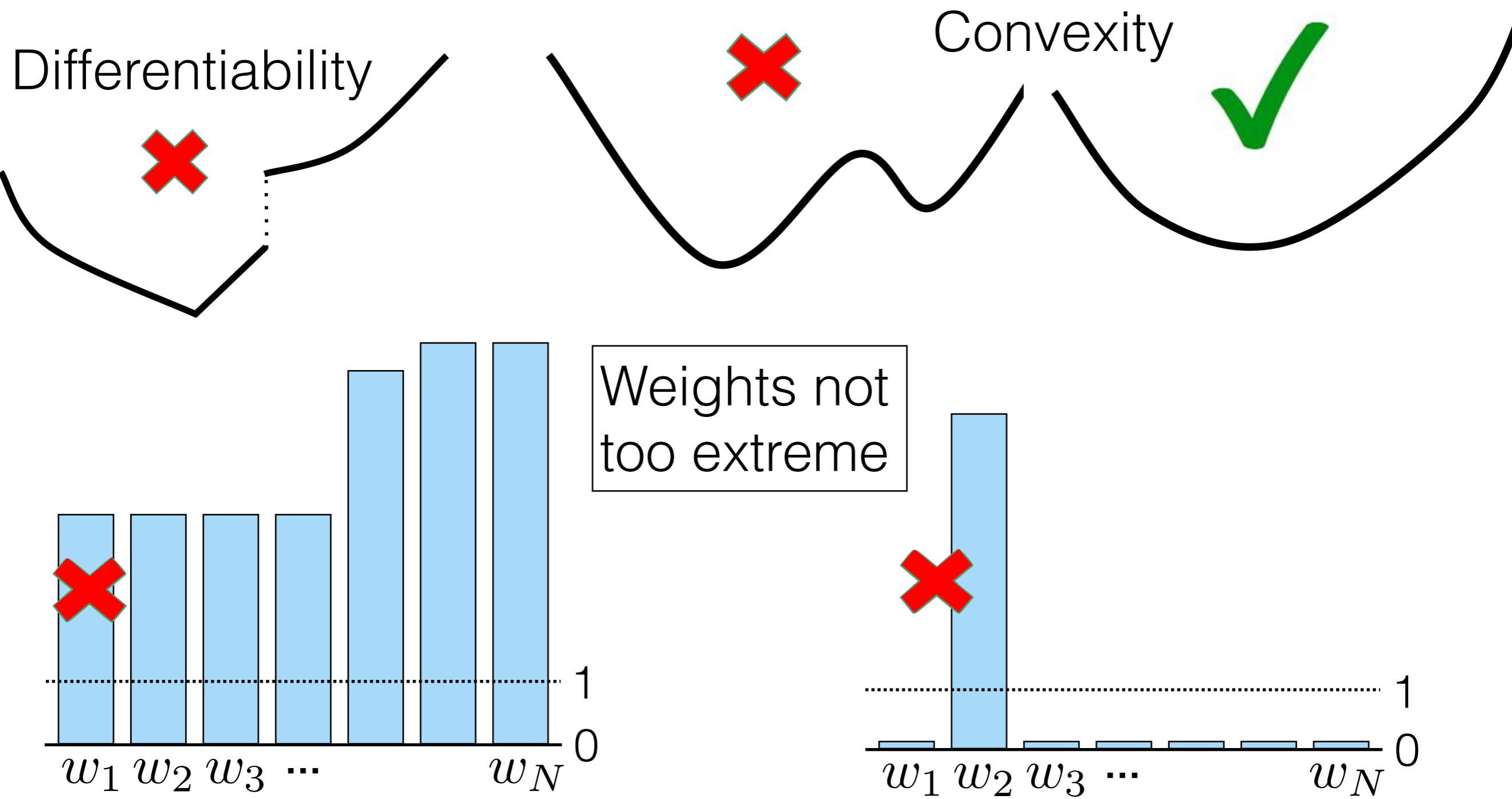
- Exact CV round: $\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$



More theory assumptions

- Exact CV round:

$$\hat{\theta}(w) := \operatorname{argmin}_{\theta \in \Theta} \frac{1}{N} \sum_{n=1}^N w_n f_n(\theta) + \lambda R(\theta)$$



Theory condition

Theory condition

- A set complexity condition:

$$\max_{w \in W_\delta} \sup_{\theta} \left\| \frac{1}{N} \sum_{n=1}^N (w_n - 1) \nabla_{\theta} f_n(\theta) \right\|_1 \leq \delta$$

Theory condition

- IJ approximation:

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- A set complexity condition:

$$\max_{w \in W_\delta} \sup_{\theta} \left\| \frac{1}{N} \sum_{n=1}^N (w_n - 1) \nabla_{\theta} f_n(\theta) \right\|_1 \leq \delta$$

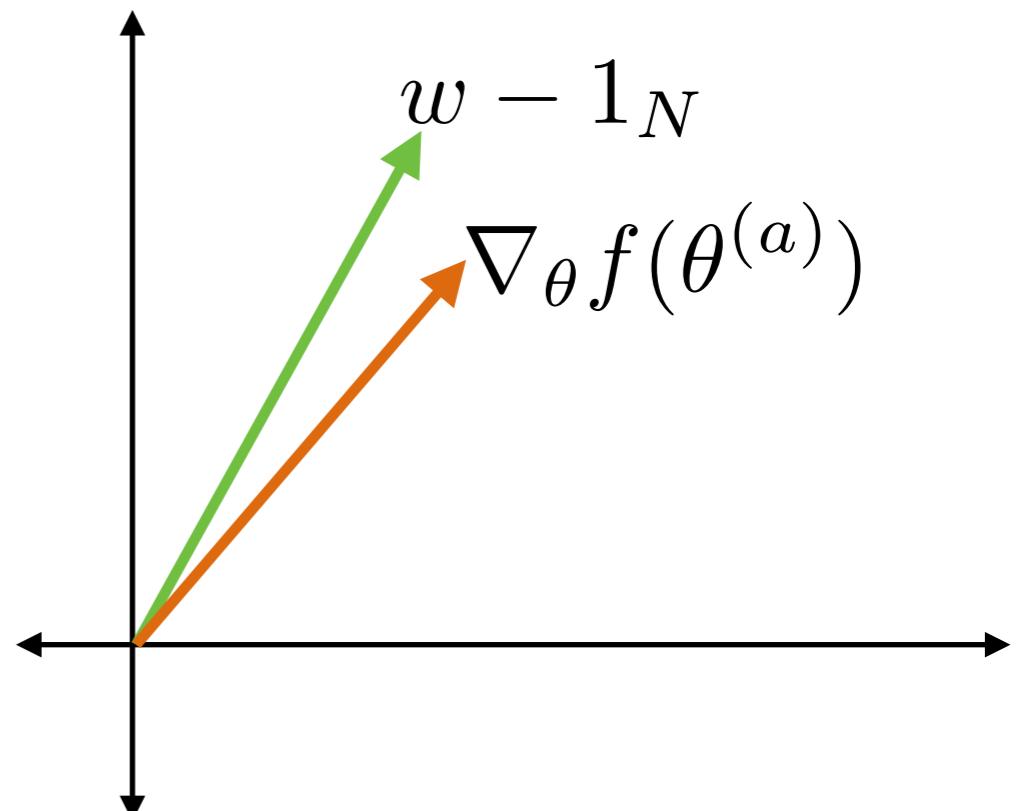
Theory condition

- IJ approximation:

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- A set complexity condition:

$$\max_{w \in W_\delta} \sup_{\theta} \left\| \frac{1}{N} \sum_{n=1}^N (w_n - 1) \nabla_{\theta} f_n(\theta) \right\|_1 \leq \delta$$



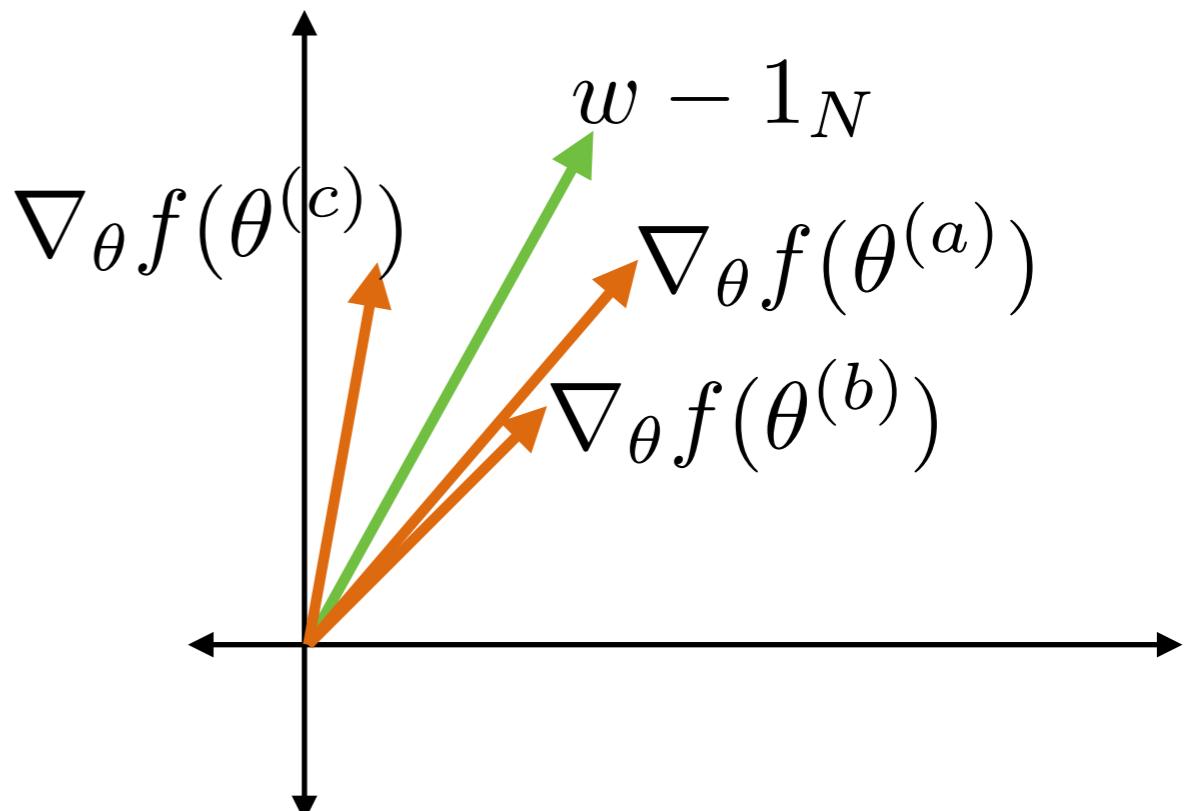
Theory condition

- IJ approximation:

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- A set complexity condition:

$$\max_{w \in W_\delta} \sup_{\theta} \left\| \frac{1}{N} \sum_{n=1}^N (w_n - 1) \nabla_{\theta} f_n(\theta) \right\|_1 \leq \delta$$



Theory condition

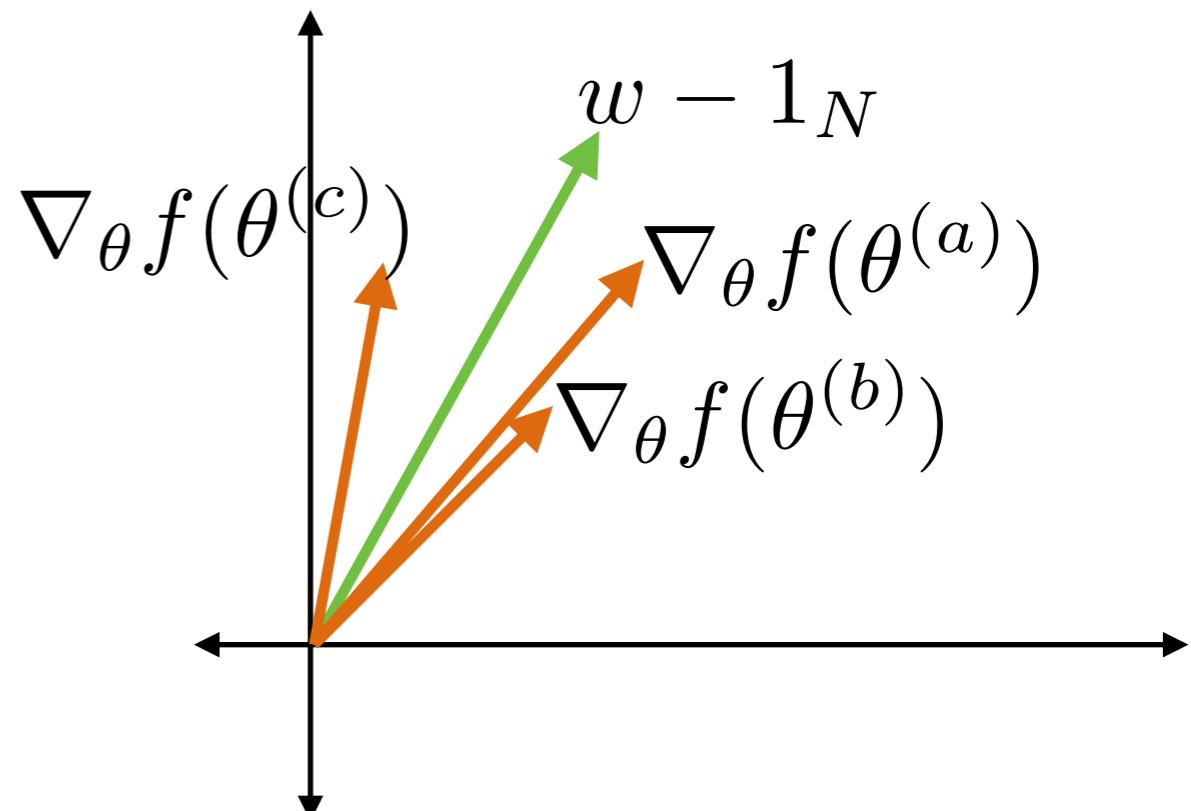
- IJ approximation:

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- A set complexity condition:

$$\max_{w \in W_\delta} \sup_{\theta} \left\| \frac{1}{N} \sum_{n=1}^N (w_n - 1) \nabla_{\theta} f_n(\theta) \right\|_1 \leq \delta$$

- + Analogous assumption on the Hessian terms



Theory condition

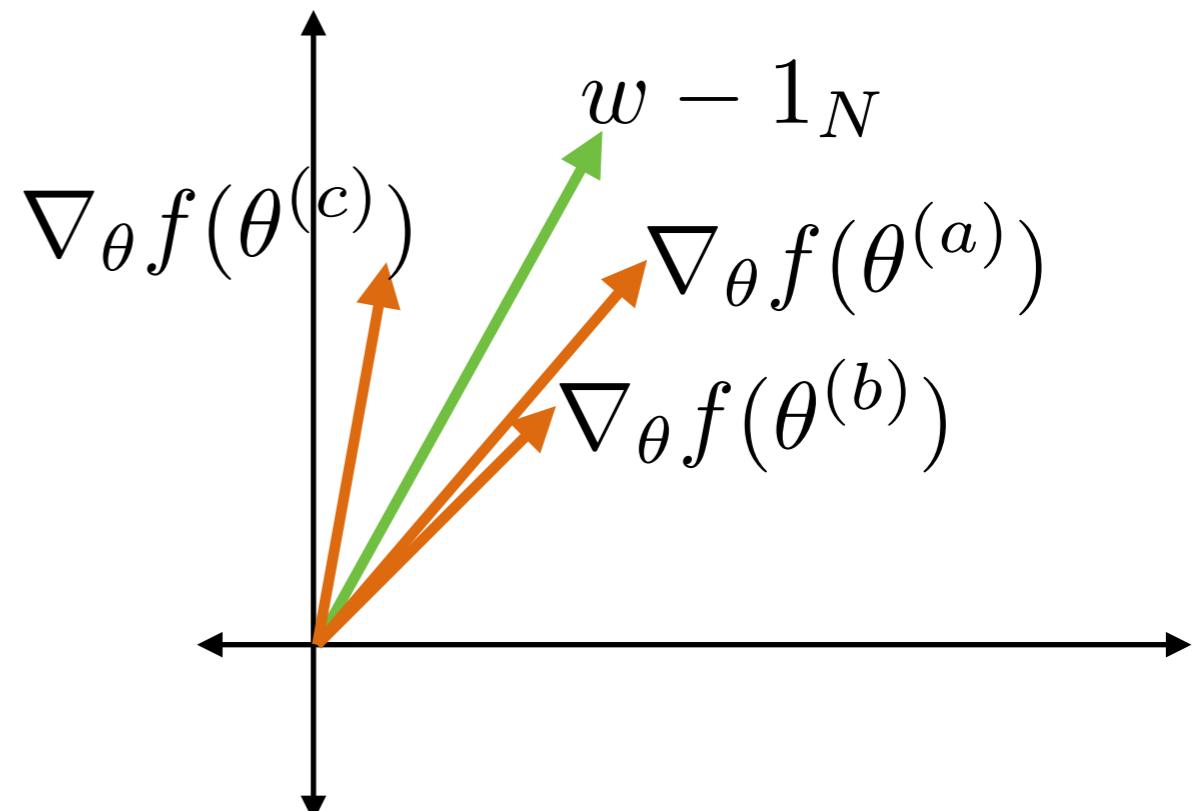
- IJ approximation:

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- A set complexity condition:

$$\max_{w \in W_\delta} \sup_{\theta} \left\| \frac{1}{N} \sum_{n=1}^N (w_n - 1) \nabla_{\theta} f_n(\theta) \right\|_1 \leq \delta$$

- + Analogous assumption on the Hessian terms
- Always a non-empty weight set for any delta



Theory condition

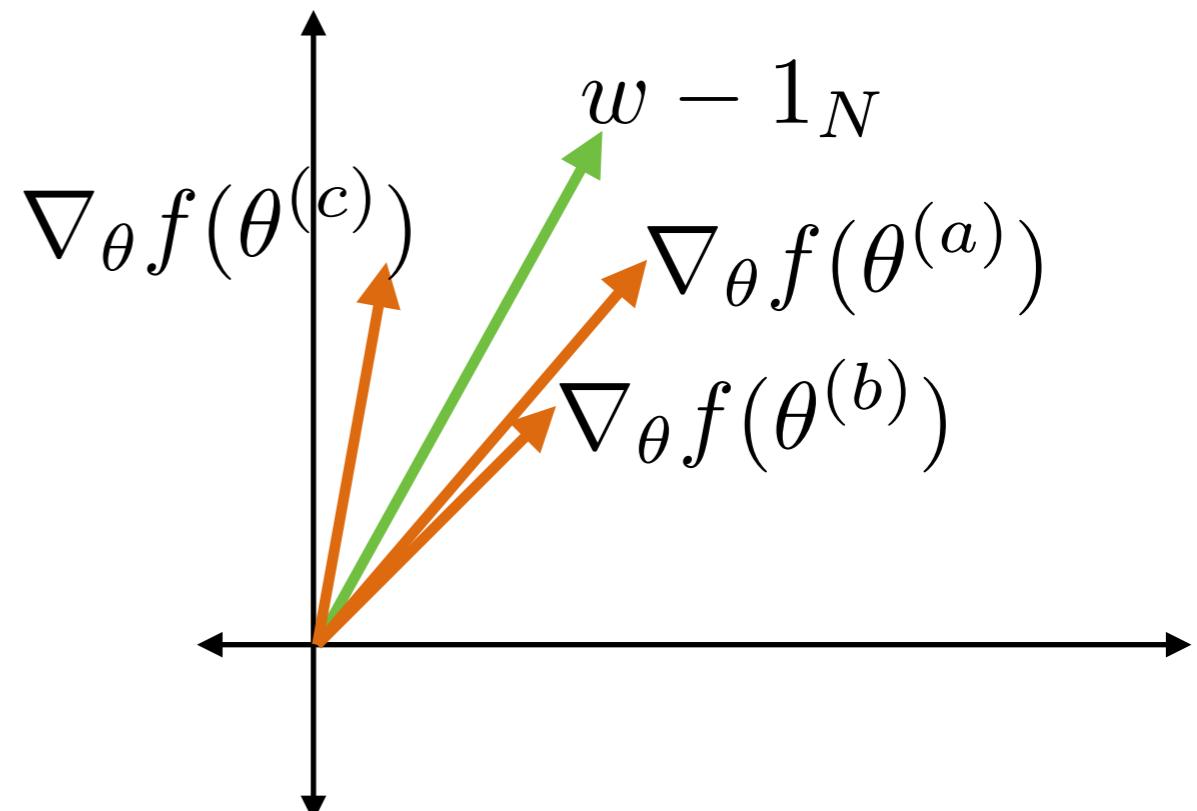
- IJ approximation:

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- A set complexity condition:

$$\max_{w \in W_\delta} \sup_{\theta} \left\| \frac{1}{N} \sum_{n=1}^N (w_n - 1) \nabla_{\theta} f_n(\theta) \right\|_1 \leq \delta$$

- + Analogous assumption on the Hessian terms
- Always a non-empty weight set for any delta
- CV works via Holder's inequality



Theory condition

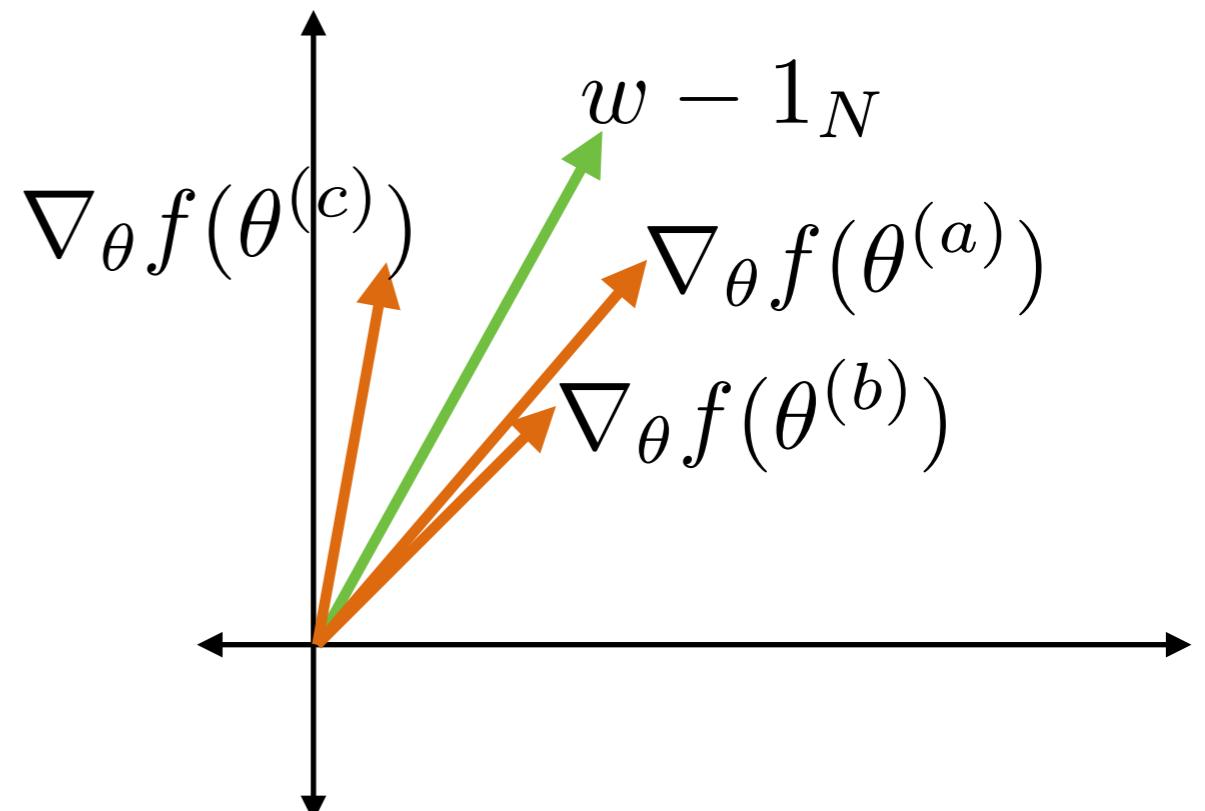
- IJ approximation:

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- A set complexity condition:

$$\max_{w \in W_\delta} \sup_{\theta} \left\| \frac{1}{N} \sum_{n=1}^N (w_n - 1) \nabla_{\theta} f_n(\theta) \right\|_1 \leq \delta$$

- + Analogous assumption on the Hessian terms
- Always a non-empty weight set for any delta
- CV works via Holder's inequality



Theory condition

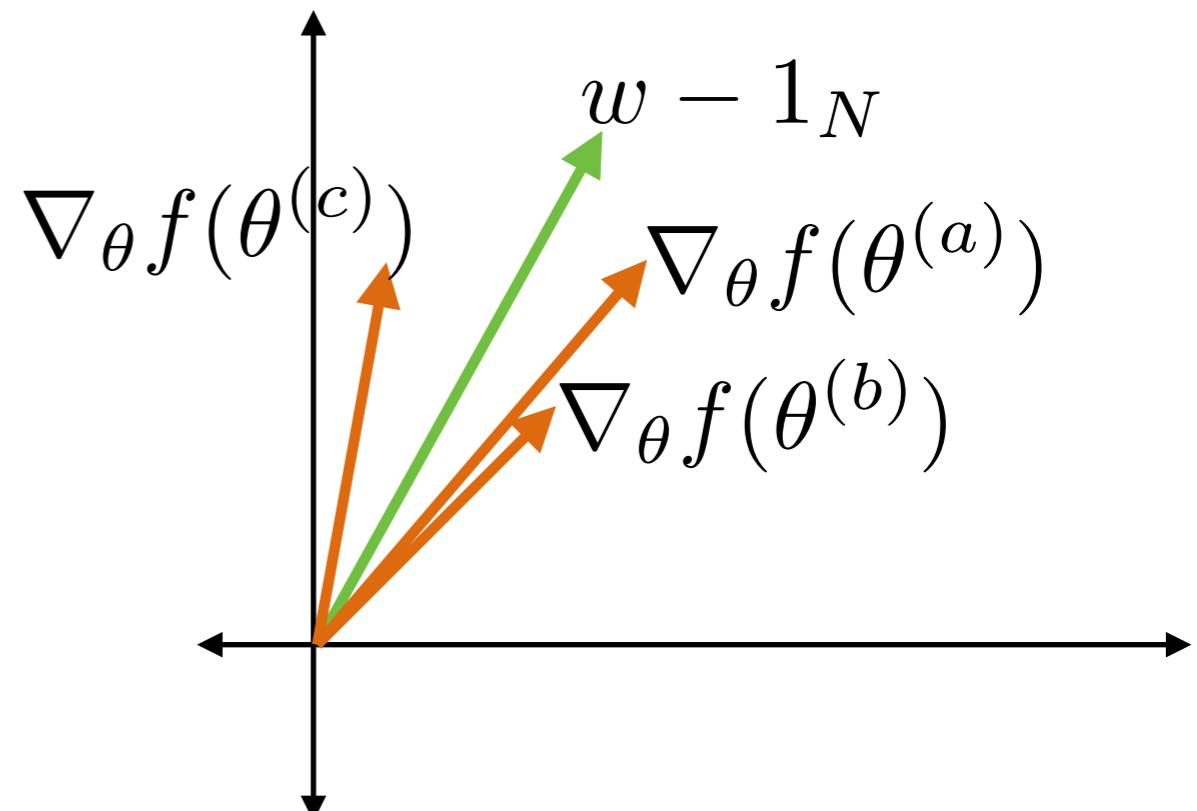
- IJ approximation:

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- A set complexity condition:

$$\max_{w \in W_\delta} \sup_{\theta} \left\| \frac{1}{N} \sum_{n=1}^N (w_n - 1) \nabla_{\theta} f_n(\theta) \right\|_1 \leq \delta$$

- + Analogous assumption on the Hessian terms
- Always a non-empty weight set for any delta
- CV works via Holder's inequality



Theory

Theorem: Under our assumptions,

$$\delta \leq \Delta \Rightarrow \max_{w \in W_\delta} \left\| \hat{\theta}_{IJ}(w) - \hat{\theta}(w) \right\|_2 \leq C\delta^2$$

Explicit functions of constants from our assumptions

Theory

Theorem: Under our assumptions,

$$\delta \leq \Delta \Rightarrow \max_{w \in W_\delta} \left\| \hat{\theta}_{IJ}(w) - \hat{\theta}(w) \right\|_2 \leq C\delta^2$$

Explicit functions of constants from our assumptions

Corollary: For leave- k -out CV,

$$\max_{w \in W^{(k)}} \left\| \hat{\theta}_{IJ}(w) - \hat{\theta}(w) \right\|_2 \leq \frac{C_1 \|\nabla_{\theta} f\|_{\infty}^2}{N^2} \leq C_2 N^{-1}$$

Swiss Army IJ: Experiments

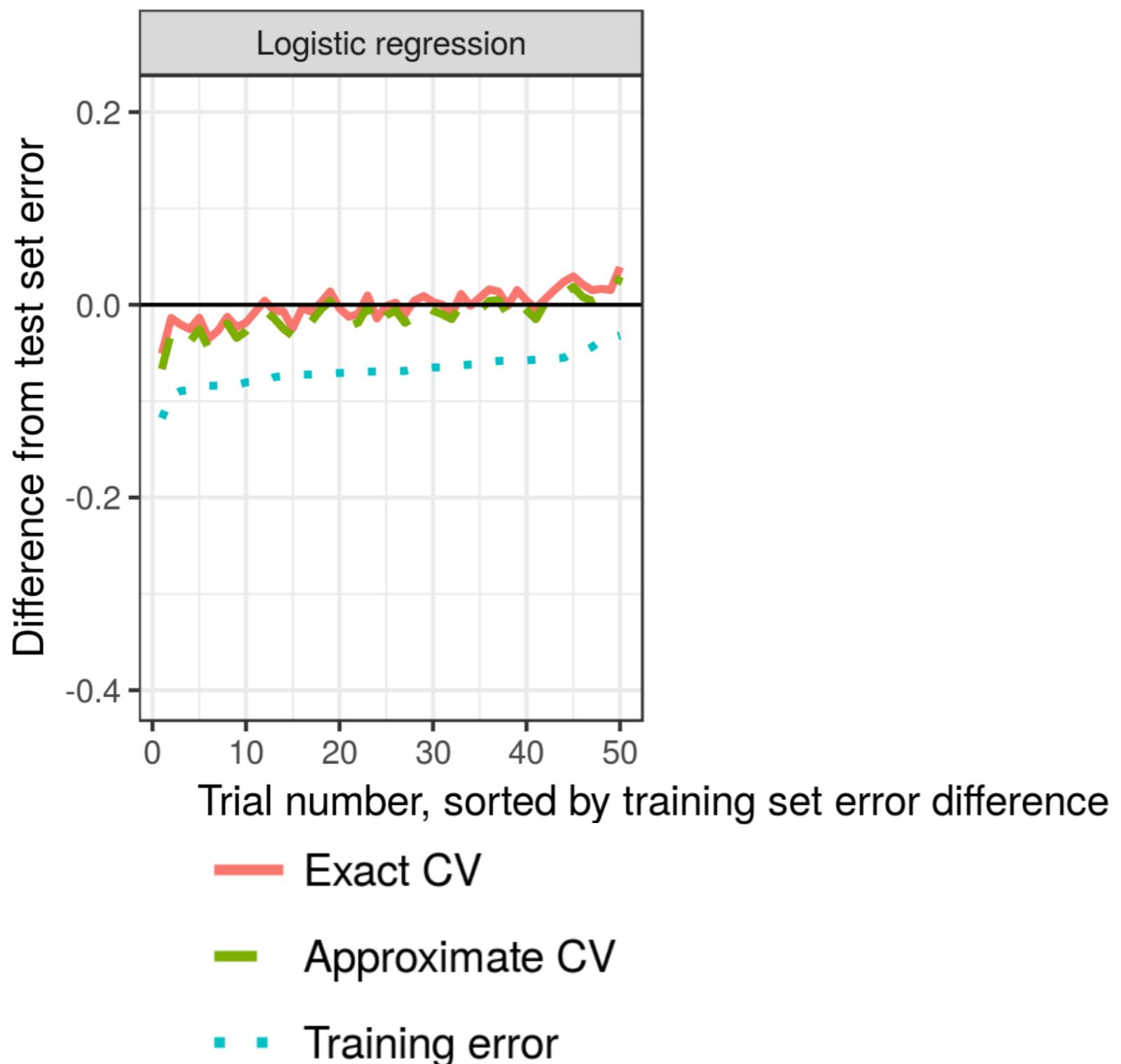
Swiss Army IJ: Experiments

- Simulated data

Swiss Army IJ: Experiments

- Simulated data

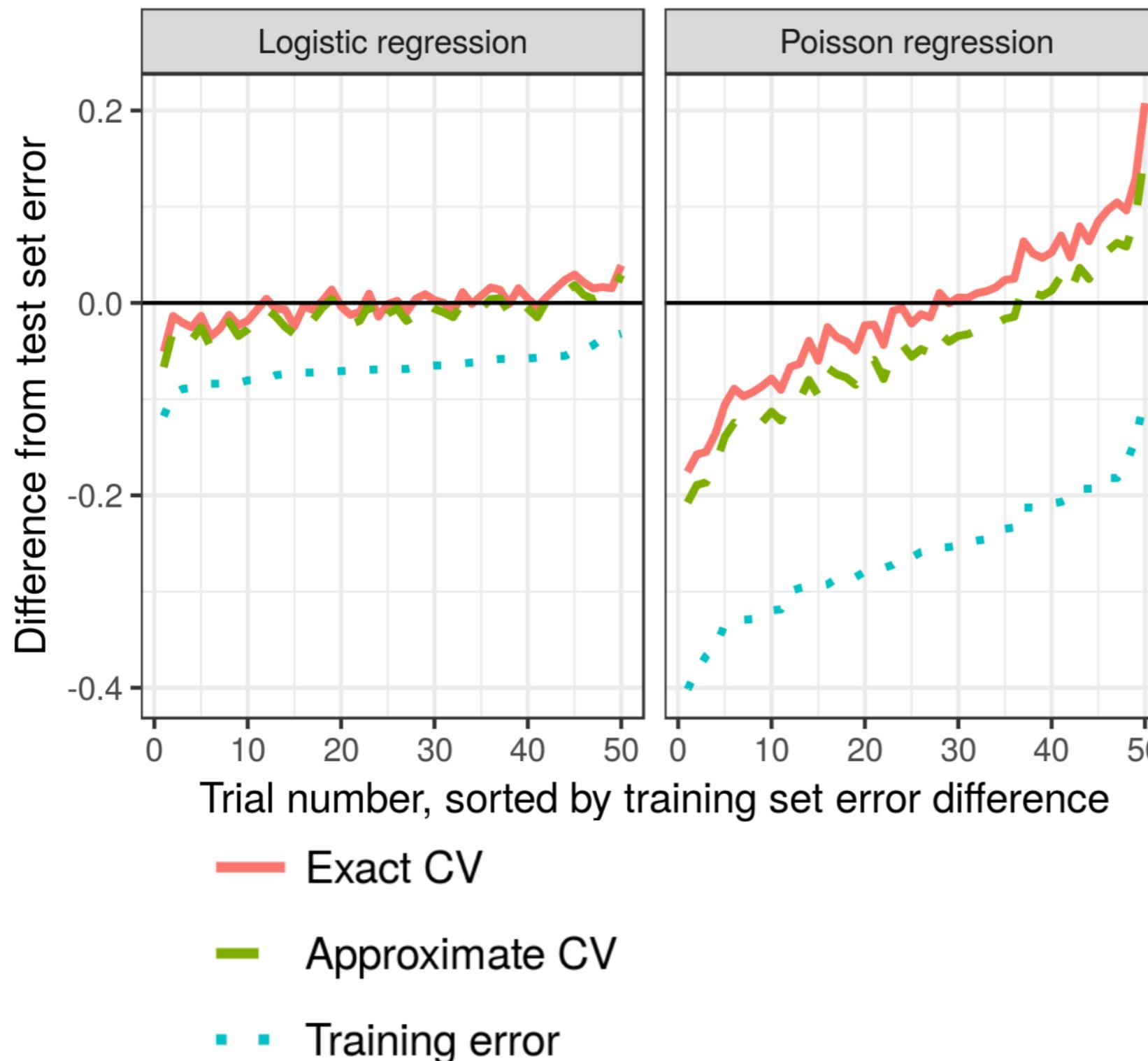
Accuracy comparison



Swiss Army IJ: Experiments

- Simulated data

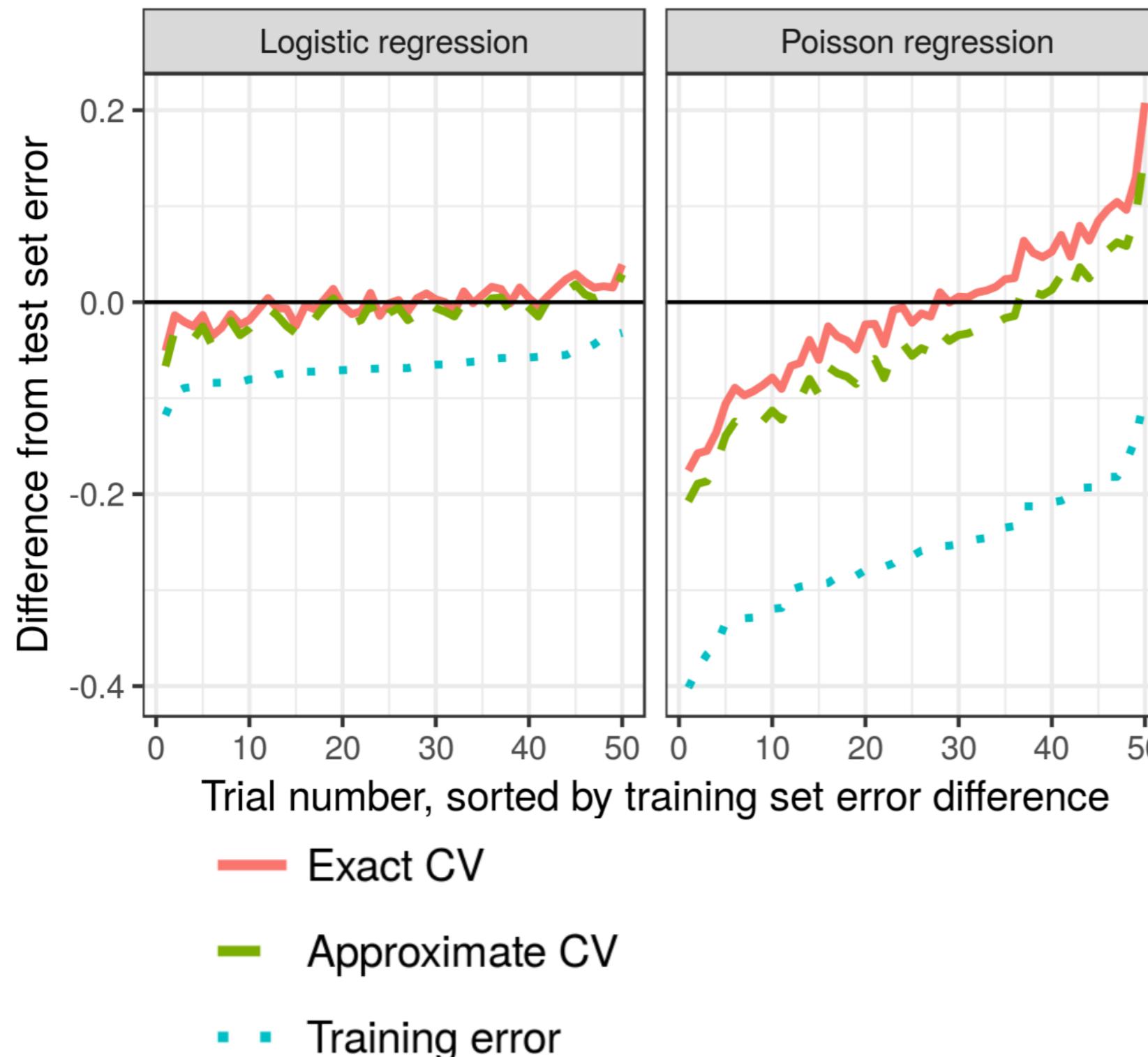
Accuracy comparison



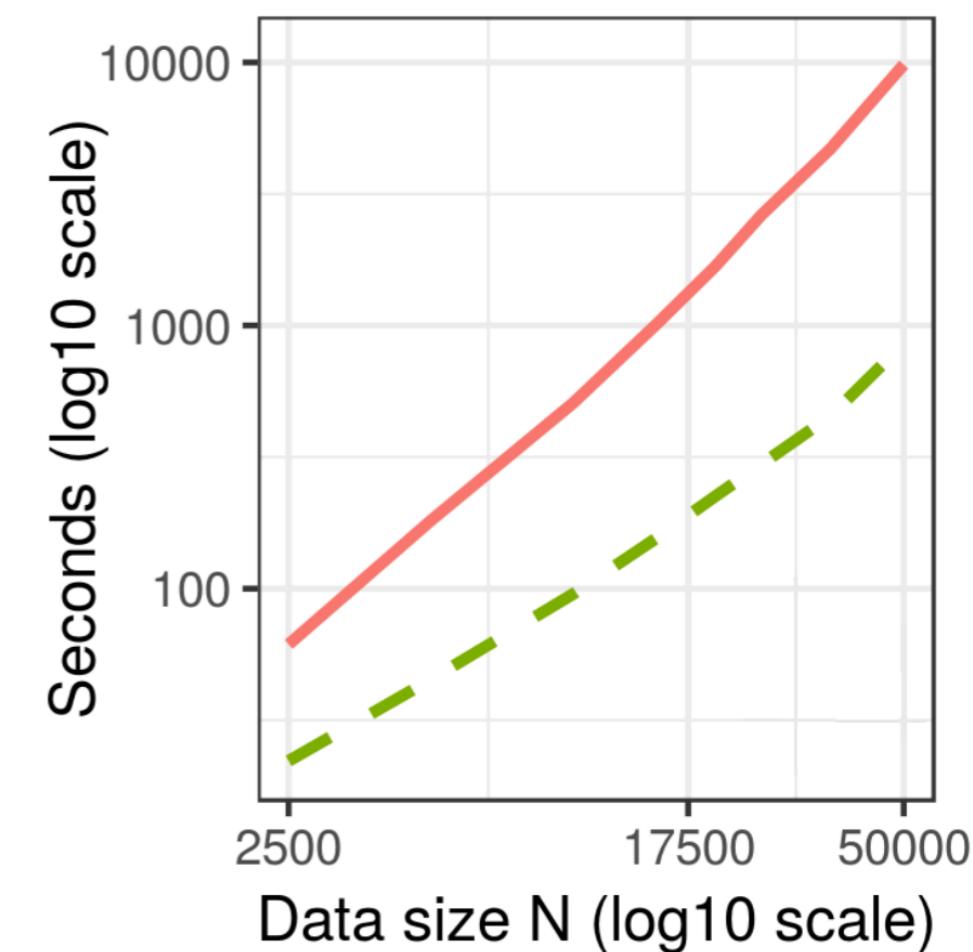
Swiss Army IJ: Experiments

- Simulated data

Accuracy comparison



Timing comparison



Swiss Army IJ: Experiments

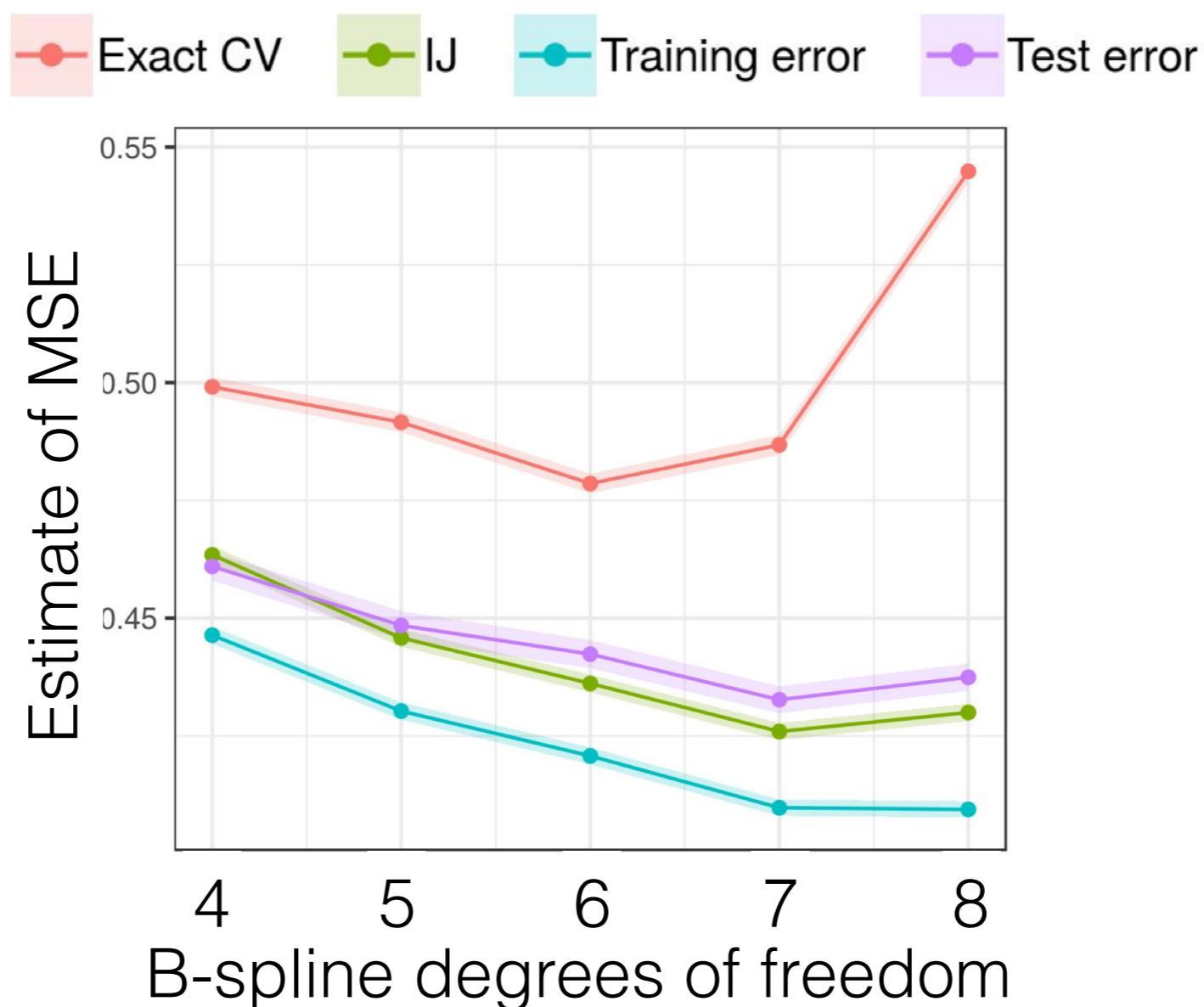
- Mice gene expression ($N=1000$) time series, after influenza infection

Swiss Army IJ: Experiments

- Mice gene expression ($N=1000$) time series, after influenza infection
- Two-stage analysis: Regress gene expression on B-spline basis, then cluster (transformed) fits

Swiss Army IJ: Experiments

- Mice gene expression ($N=1000$) time series, after influenza infection
- Two-stage analysis: Regress gene expression on B-spline basis, then cluster (transformed) fits



Roadmap

- Cross validation (CV) setup
- CV Challenge #1: an expensive data analysis
- Solution #1: a Taylor series approximation (“Swiss Army infinitesimal jackknife”)
- CV Challenge #2: high dimensions
- Solution #2: using sparsity (carefully)

Roadmap

- Cross validation (CV) setup
- CV Challenge #1: an expensive data analysis
- Solution #1: a Taylor series approximation (“Swiss Army infinitesimal jackknife”)
- CV Challenge #2: high dimensions
- Solution #2: using sparsity (carefully)

Approximation for high dimensions

- IJ approximation:

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

Approximation for high dimensions

- IJ approximation:

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- Cost: (1 data analysis)

Approximation for high dimensions

- IJ approximation: H : Hessian of the objective G : matrix of gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- Cost: (1 data analysis)

Approximation for high dimensions

- IJ approximation: H : Hessian of the objective G : matrix of gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- Cost: (1 data analysis)
- Recall:

Approximation for high dimensions

- IJ approximation: H : Hessian of the objective G : matrix of gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- Cost: (1 data analysis)
- Recall:

$$H_{d_1 d_2} = \frac{\partial^2}{\partial \theta_{d_1} \partial \theta_{d_2}} \text{objective}$$

Approximation for high dimensions

- IJ approximation: H : Hessian of the objective G : matrix of gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- Cost: (1 data analysis)
- Recall:

$$H_{d_1 d_2} = \frac{\partial^2}{\partial \theta_{d_1} \partial \theta_{d_2}} \text{objective}$$

$D \times D$

Approximation for high dimensions

- IJ approximation: H : Hessian of the objective G : matrix of gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- Cost: (1 data analysis)
- Recall:

$$H_{d_1 d_2} = \frac{\partial^2}{\partial \theta_{d_1} \partial \theta_{d_2}} \text{objective}$$

$D \times D$

$$G_{dn} = \frac{\partial^2}{\partial \theta_d \partial w_n} \text{objective}$$

Approximation for high dimensions

- IJ approximation:

H : Hessian of
the objective

G : matrix of
gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- Cost: (1 data analysis)

- Recall:

$$H_{d_1 d_2} = \frac{\partial^2}{\partial \theta_{d_1} \partial \theta_{d_2}} \text{objective}$$

$D \times D$

$$G_{dn} = \frac{\partial^2}{\partial \theta_d \partial w_n} \text{objective}$$

$D \times N$

Approximation for high dimensions

- IJ approximation:

H : Hessian of
the objective

G : matrix of
gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- Cost: (1 data analysis) + $O(D^3)$

- Recall:

$$H_{d_1 d_2} = \frac{\partial^2}{\partial \theta_{d_1} \partial \theta_{d_2}} \text{objective}$$

$D \times D$

$$G_{dn} = \frac{\partial^2}{\partial \theta_d \partial w_n} \text{objective}$$

$D \times N$

Approximation for high dimensions

- IJ approximation: H : Hessian of the objective G : matrix of gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- Cost: (1 data analysis) + $O(D^3)$

Approximation for high dimensions

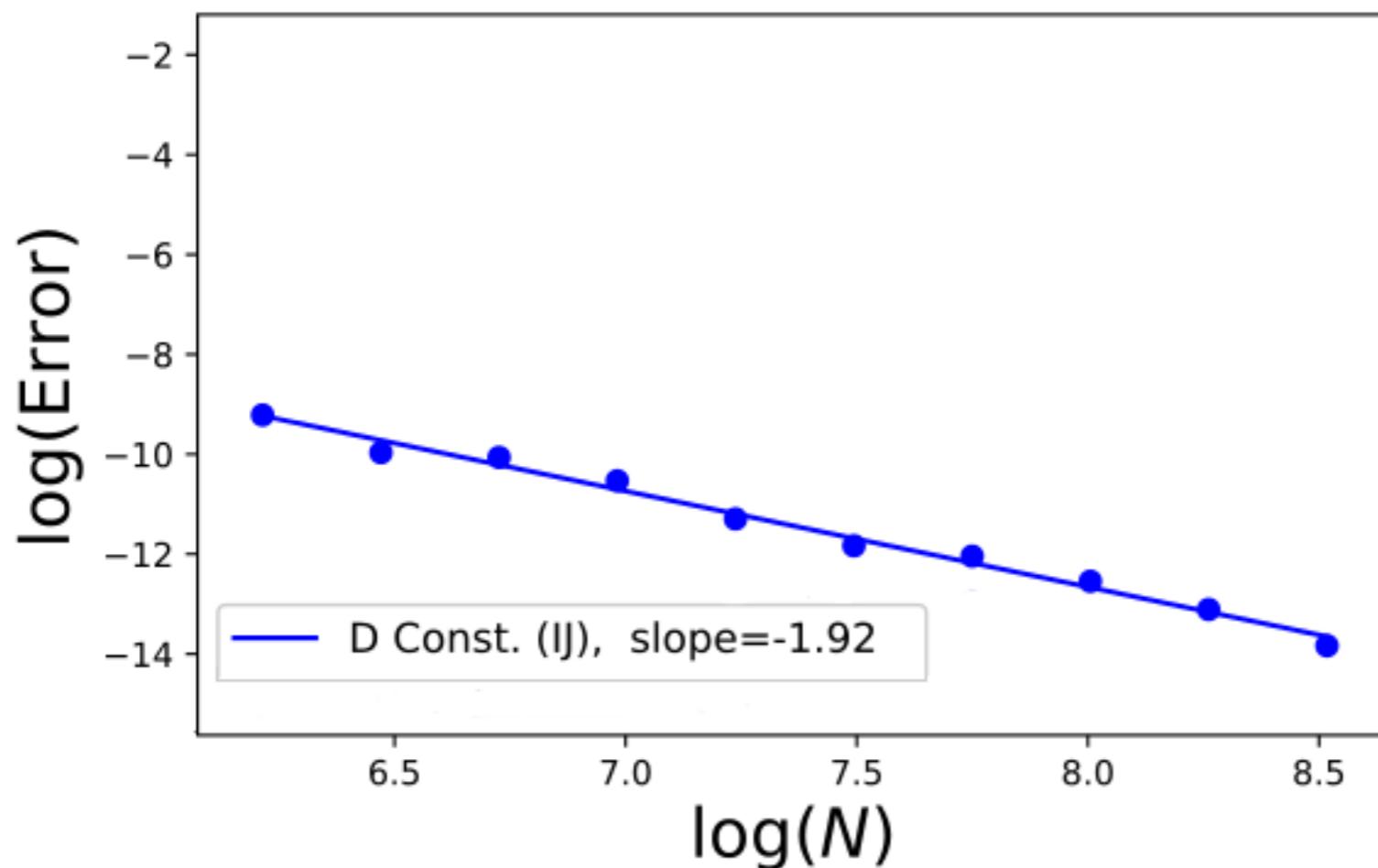
- IJ approximation: H : Hessian of the objective G : matrix of gradients
$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$
 - Cost: (1 data analysis) + $O(D^3)$
-
- Accuracy also bad in high dims

Approximation for high dimensions

- IJ approximation: H : Hessian of the objective G : matrix of gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- Cost: (1 data analysis) + $O(D^3)$



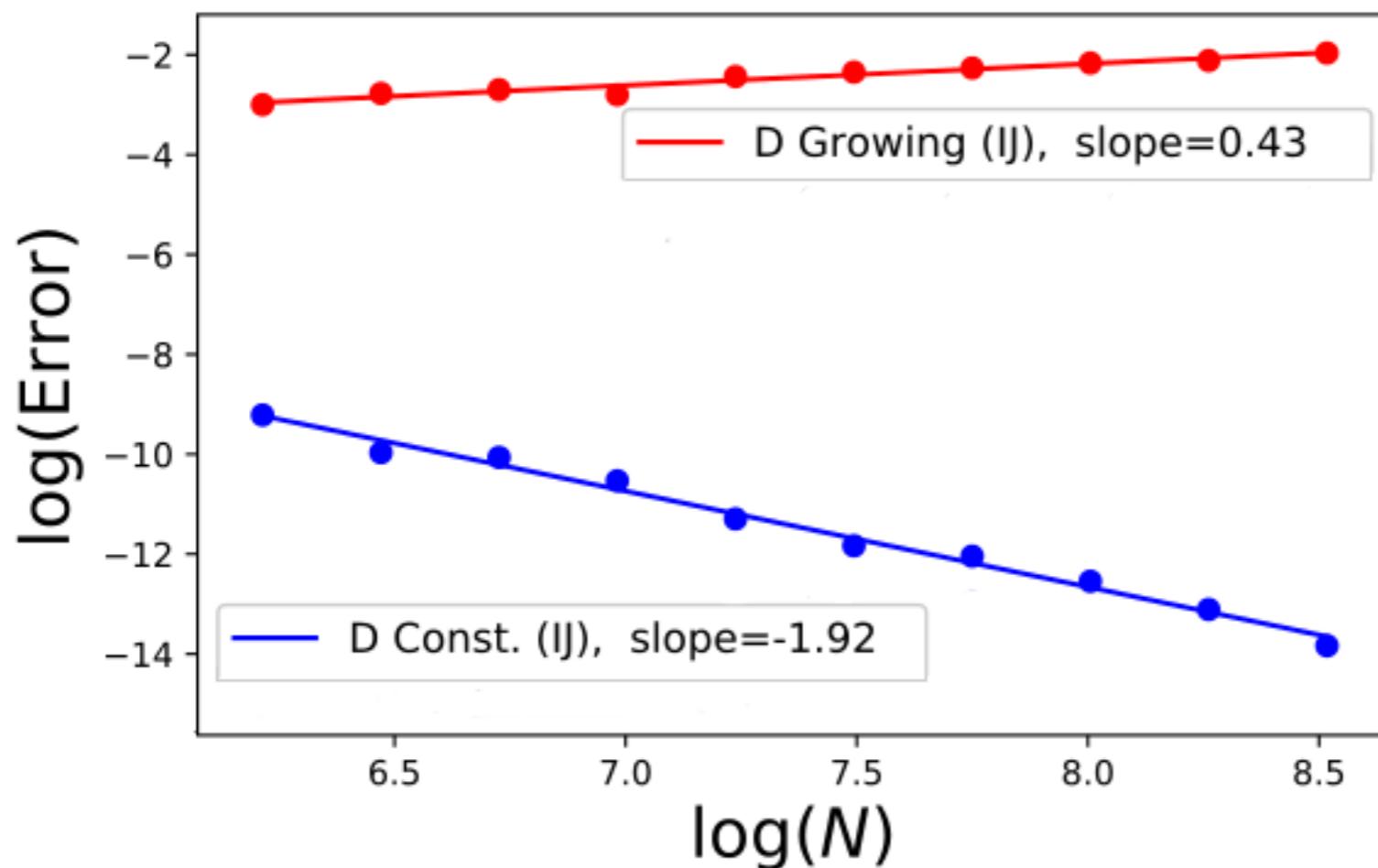
- Accuracy also bad in high dims

Approximation for high dimensions

- IJ approximation: H : Hessian of the objective G : matrix of gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- Cost: (1 data analysis) + $O(D^3)$



- Accuracy also bad in high dims

Approximation for high dimensions

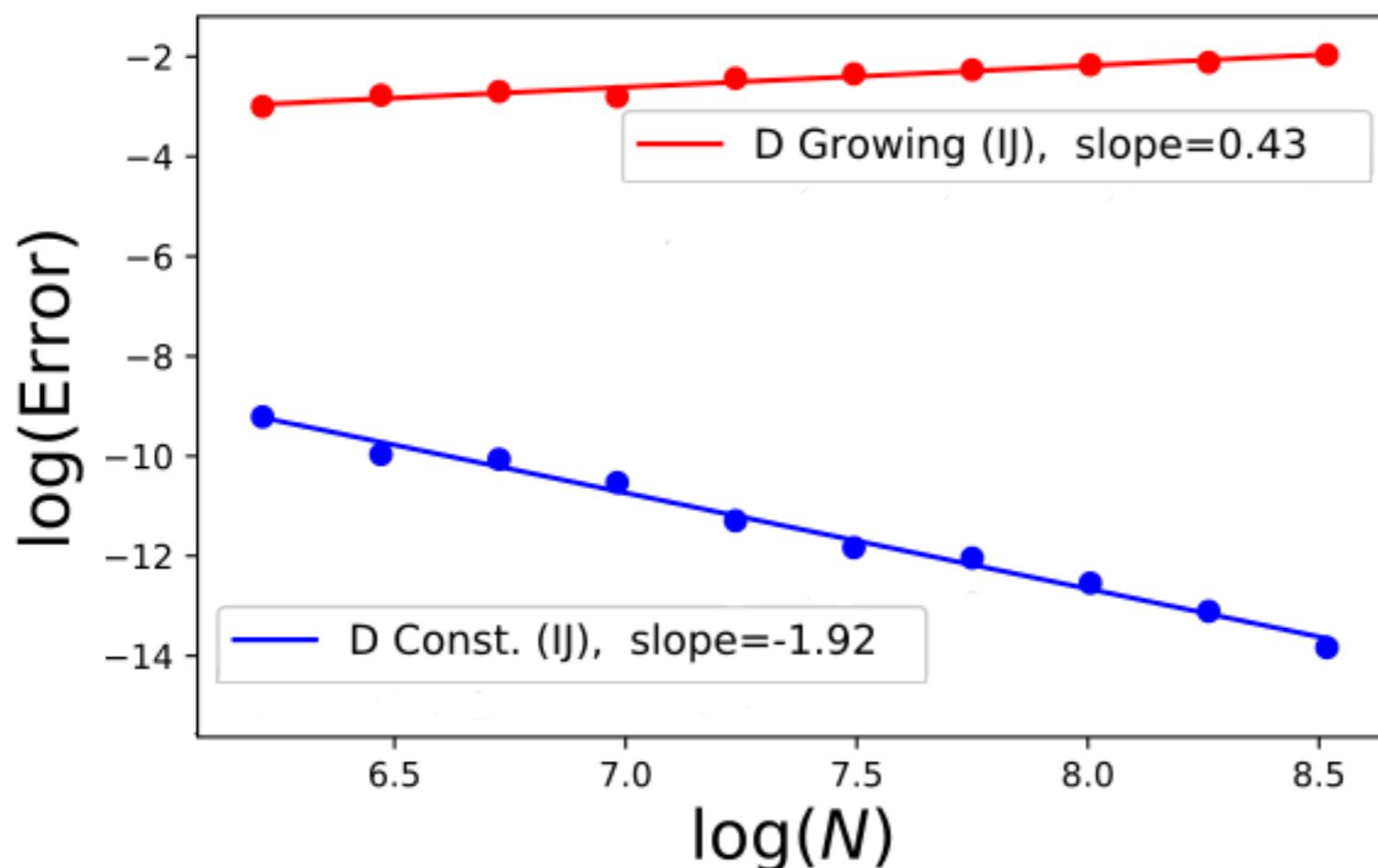
- IJ approximation:

H : Hessian of
the objective

G : matrix of
gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- Cost: (1 data analysis) + $O(D^3)$



- Accuracy also bad in high dims
- Problem doesn't go away on its own when there's a smaller "effective dimension" D_{eff}

Approximation for high dimensions

- IJ approximation: H : Hessian of the objective G : matrix of gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- What to do?

Approximation for high dimensions

- IJ approximation: H : Hessian of the objective G : matrix of gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- What to do?

Ideas:

Does it work?

Approximation for high dimensions

- IJ approximation: H : Hessian of the objective G : matrix of gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- What to do?

Ideas:

- Further approximation via randomized solvers

Does it work?

Approximation for high dimensions

- IJ approximation:

H : Hessian of
the objective

G : matrix of
gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- What to do?

Ideas:

- Further approximation via randomized solvers

Does it work?

- Can be fast but additional error

Approximation for high dimensions

- IJ approximation:

H : Hessian of
the objective

G : matrix of
gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- What to do?

Ideas:

- Further approximation via randomized solvers
- Use Swiss Army IJ with L1 regularization

Does it work?

- Can be fast but additional error

Approximation for high dimensions

- IJ approximation:

H : Hessian of
the objective

G : matrix of
gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- What to do?

Ideas:

- Further approximation via randomized solvers
- Use Swiss Army IJ with L1 regularization

Does it work?

- Can be fast but additional error
- L1 penalty not differentiable

Approximation for high dimensions

- IJ approximation:

H : Hessian of
the objective

G : matrix of
gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- What to do?

Ideas:

- Further approximation via randomized solvers
- Use Swiss Army IJ with L1 regularization
- Use Swiss Army IJ with smooth approximation to L1

Does it work?

- Can be fast but additional error
- L1 penalty not differentiable

Approximation for high dimensions

- IJ approximation:

H : Hessian of
the objective

G : matrix of
gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- What to do?

Ideas:

- Further approximation via randomized solvers
- Use Swiss Army IJ with L1 regularization
- Use Swiss Army IJ with smooth approximation to L1

Does it work?

- Can be fast but additional error
- L1 penalty not differentiable
- We show is inaccurate, slow

Approximation for high dimensions

- IJ approximation:

H : Hessian of
the objective

G : matrix of
gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- What to do?

Ideas:

- Further approximation via randomized solvers
- Use Swiss Army IJ with L1 regularization
- Use Swiss Army IJ with smooth approximation to L1
- Run L1 regularization *then* IJ on the (small) support

Does it work?

- Can be fast but additional error
- L1 penalty not differentiable
- We show is inaccurate, slow

Approximation for high dimensions

- IJ approximation:

H : Hessian of
the objective

G : matrix of
gradients

$$\hat{\theta}_{IJ}(w) = \hat{\theta}(1_N) - H(1_N)^{-1}G(1_N)(w - 1_N)$$

- What to do?

Ideas:

- Further approximation via randomized solvers
- Use Swiss Army IJ with L1 regularization
- Use Swiss Army IJ with smooth approximation to L1
- Run L1 regularization *then* IJ on the (small) support

Does it work?

- Can be fast but additional error
- L1 penalty not differentiable
- We show is inaccurate, slow
- Yes! Accurate, fast, supported by our theory

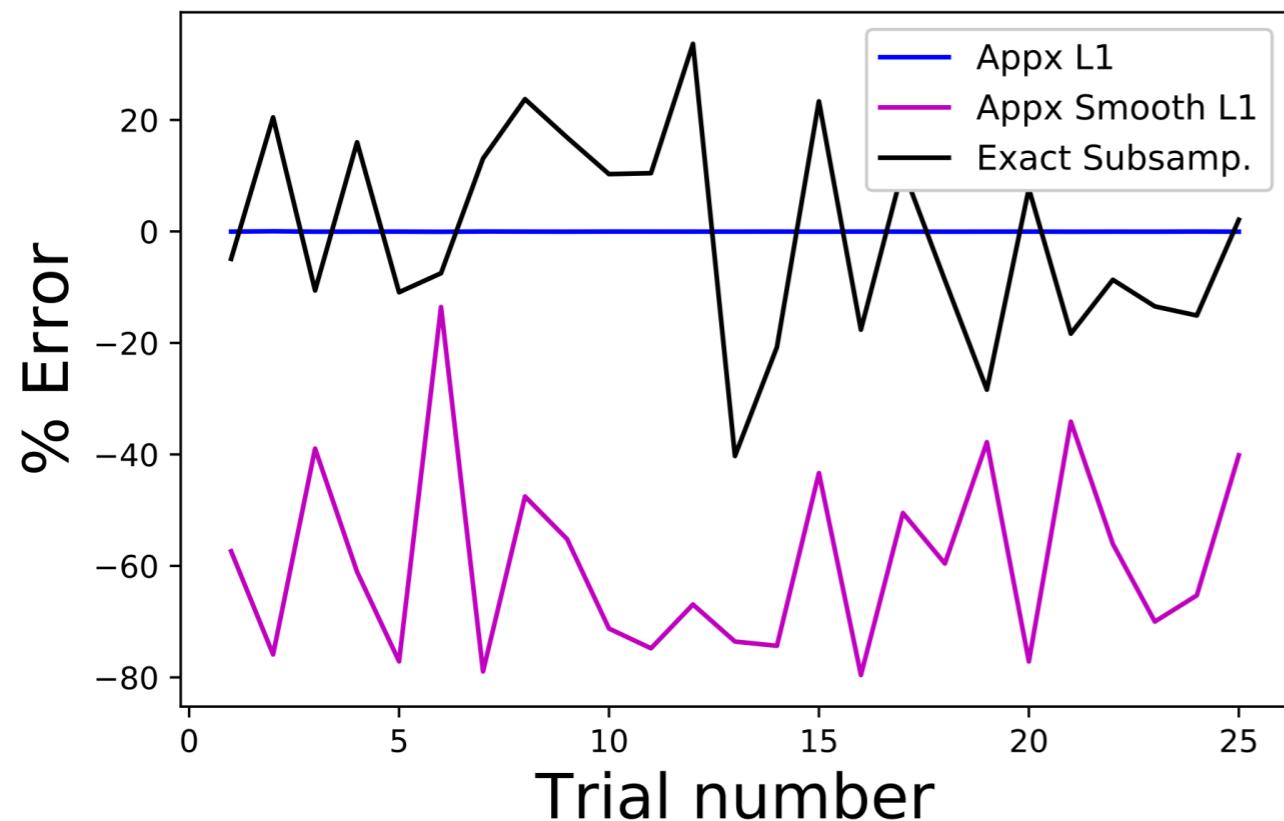
High dimensions: Experiments

High dimensions: Experiments

- Simulated data ($N=500$, $D=40,000$)

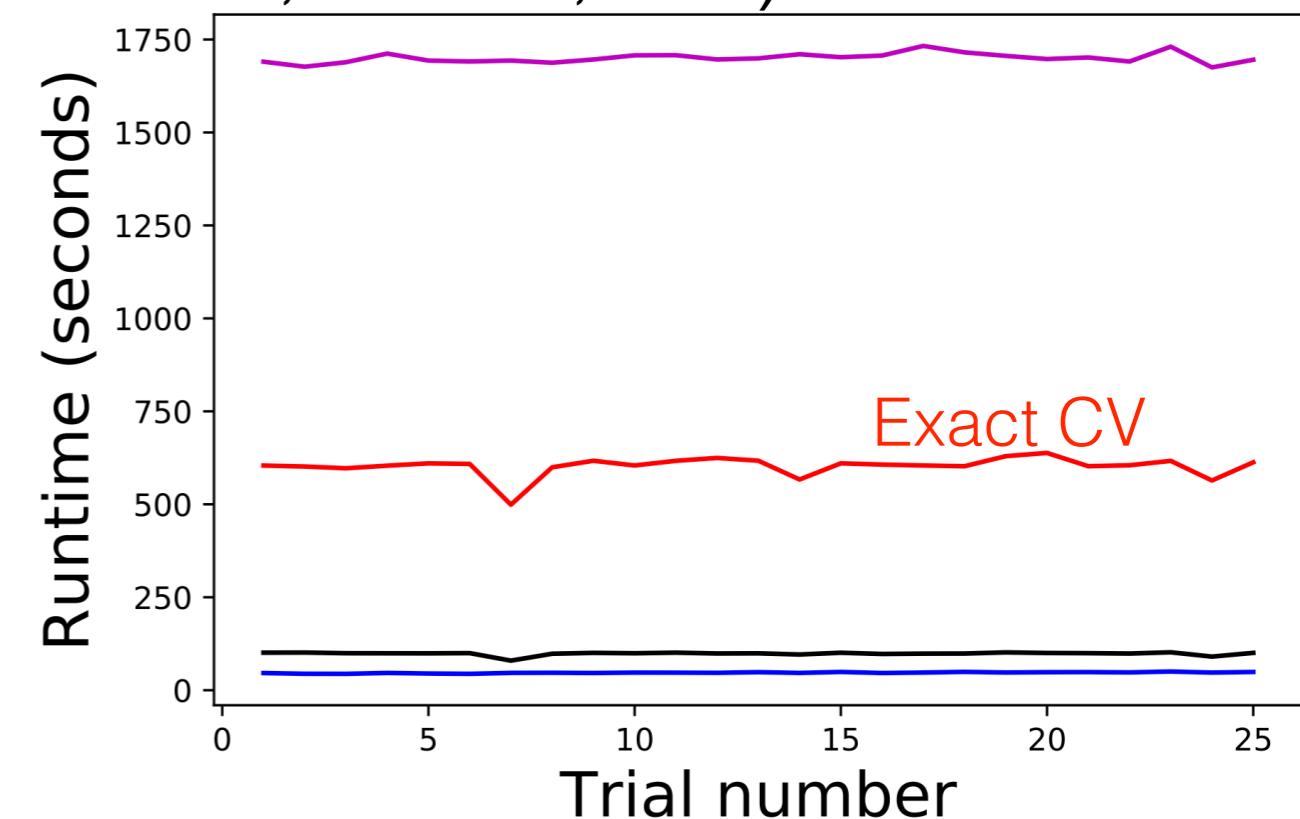
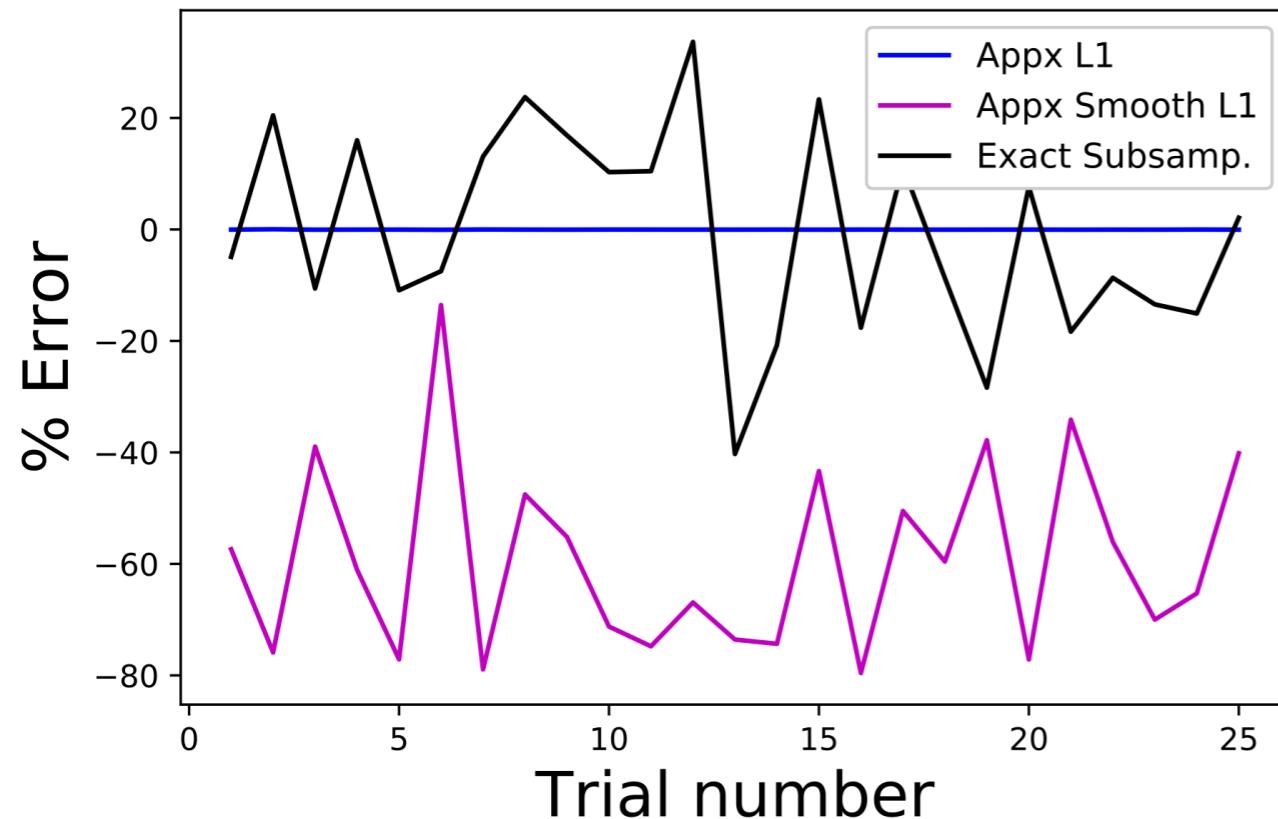
High dimensions: Experiments

- Simulated data ($N=500$, $D=40,000$)



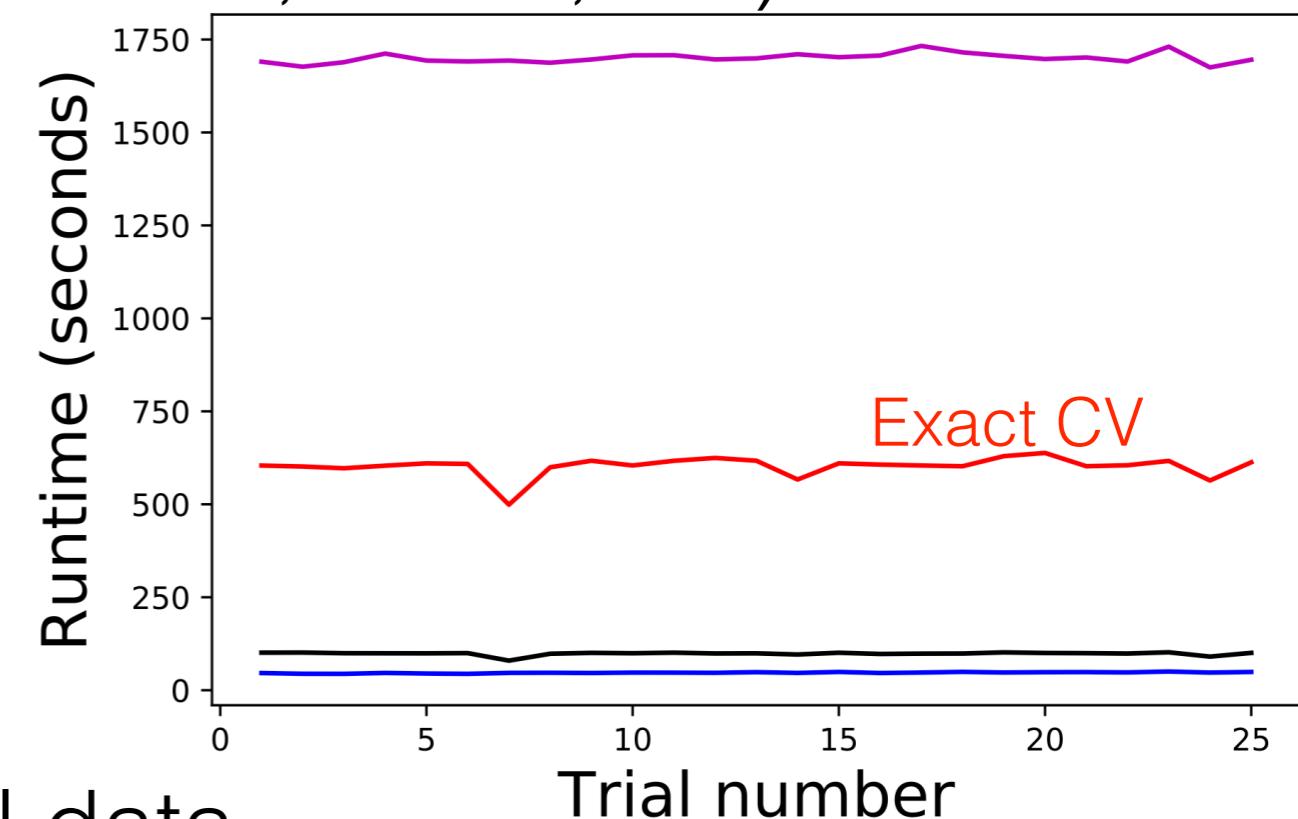
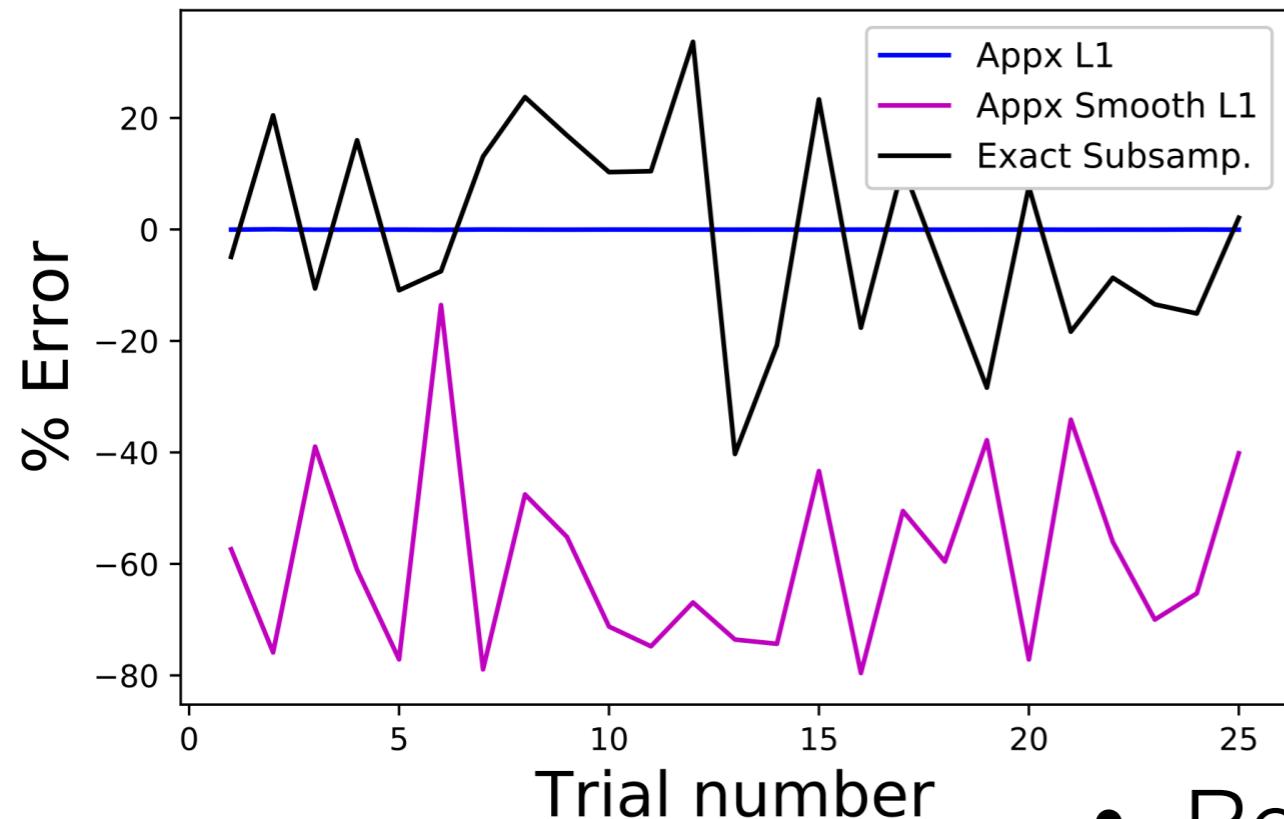
High dimensions: Experiments

- Simulated data ($N=500$, $D=40,000$)



High dimensions: Experiments

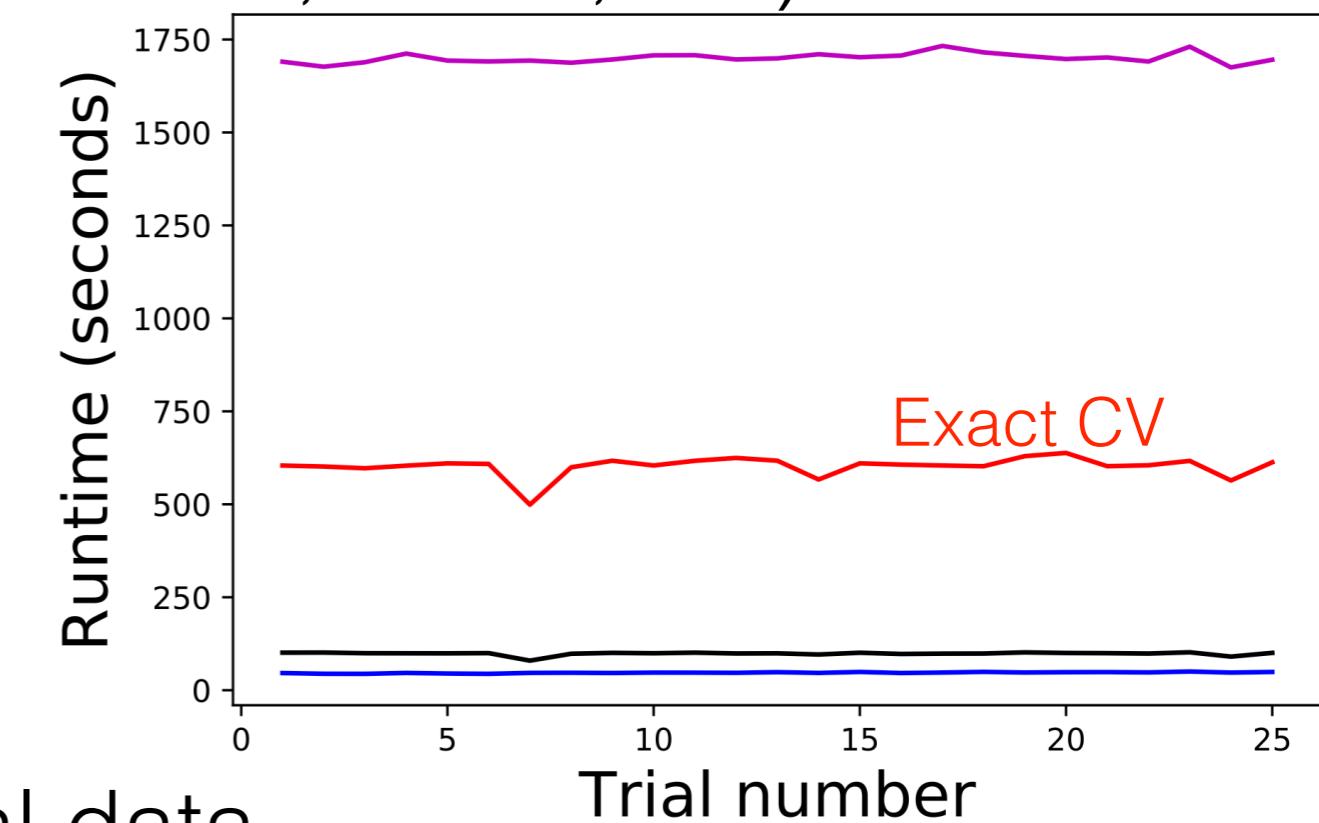
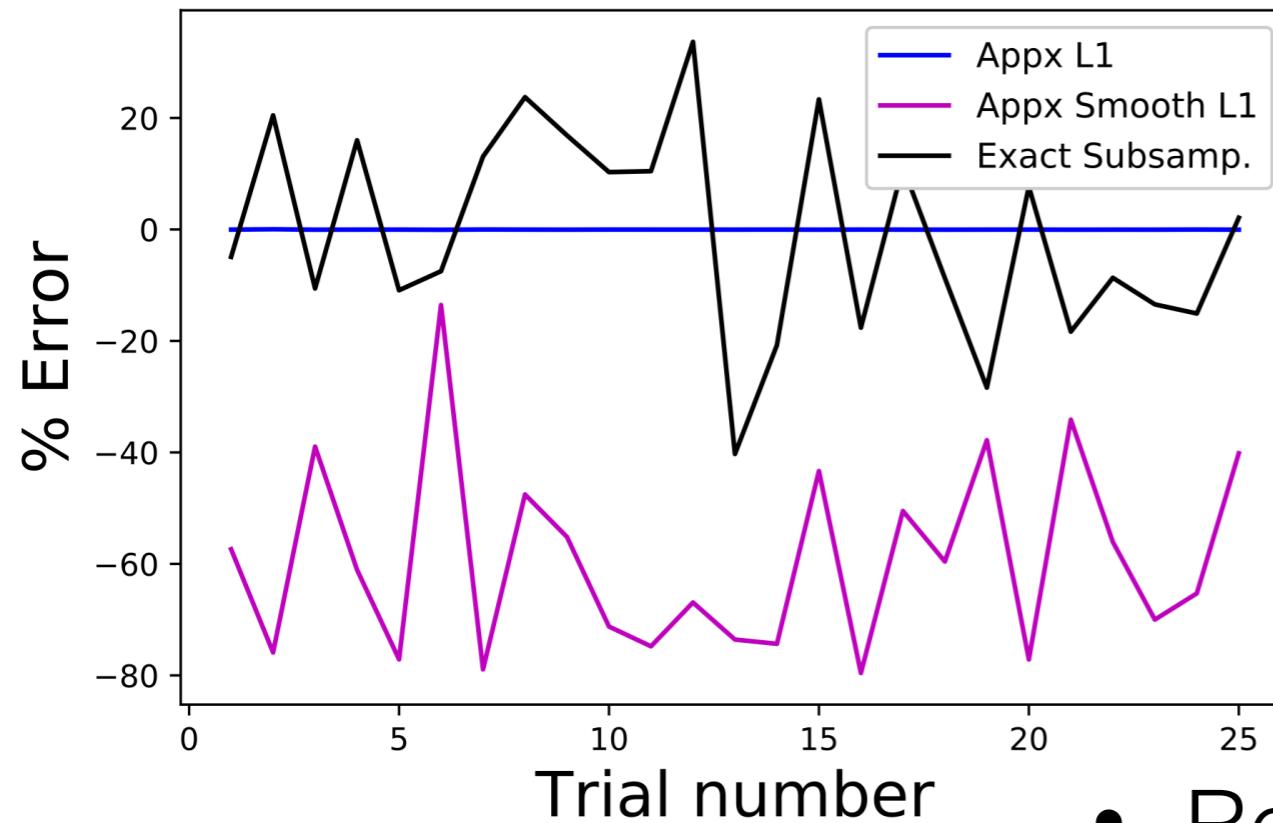
- Simulated data ($N=500$, $D=40,000$)



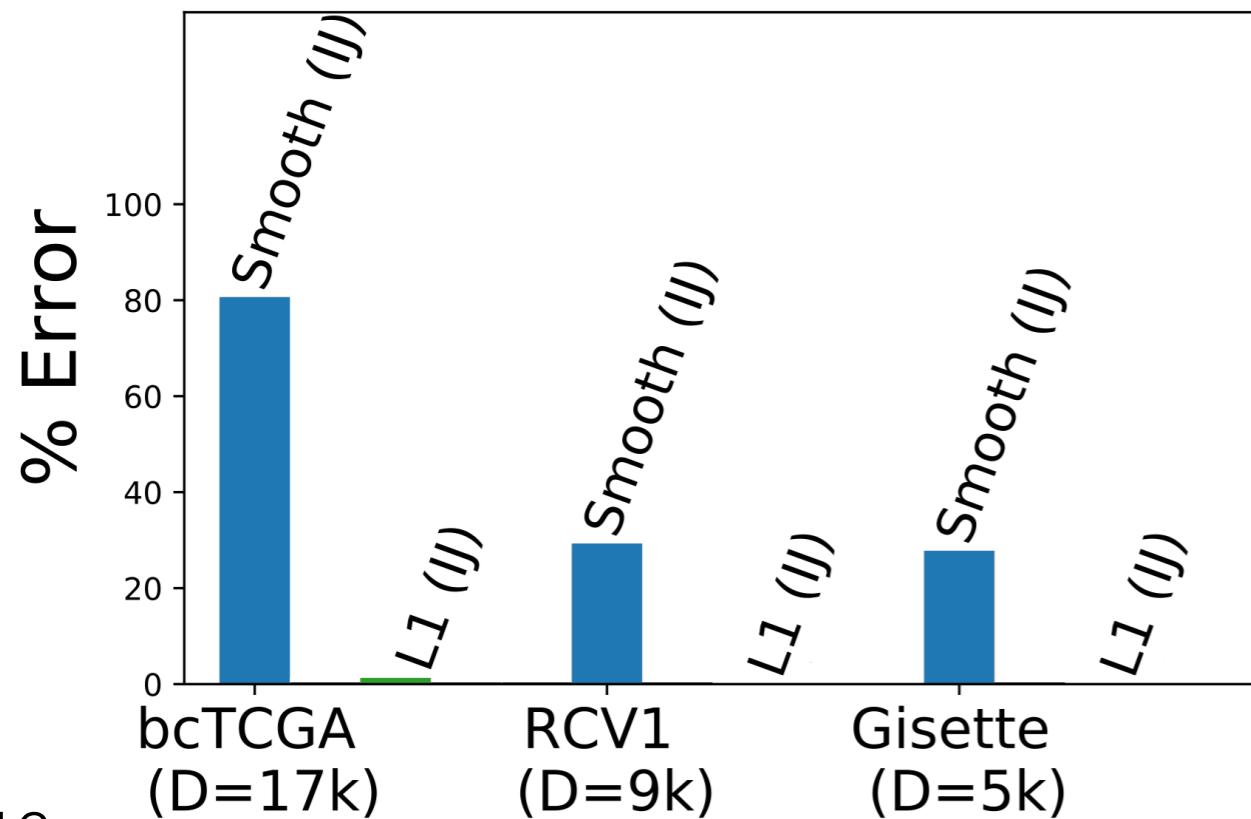
- Real data

High dimensions: Experiments

- Simulated data ($N=500$, $D=40,000$)

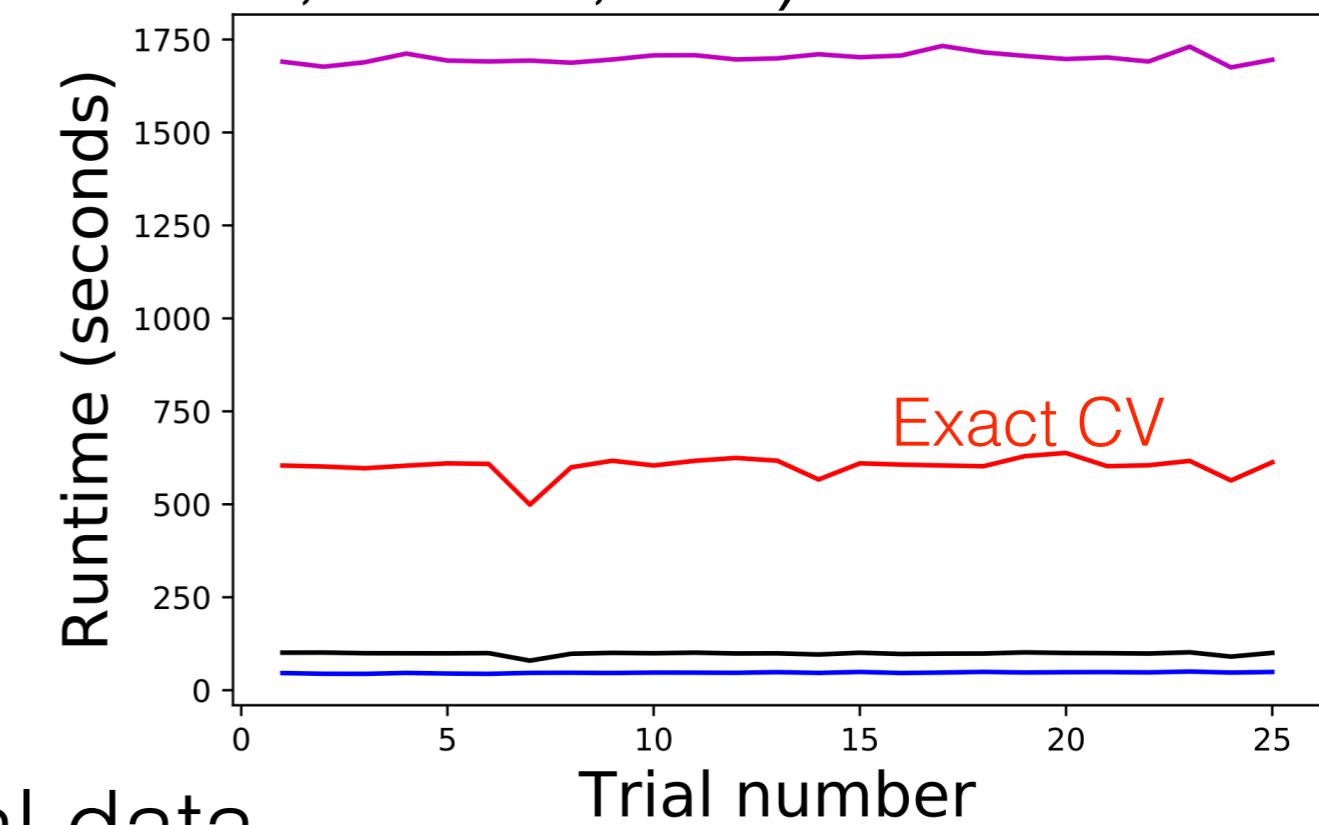
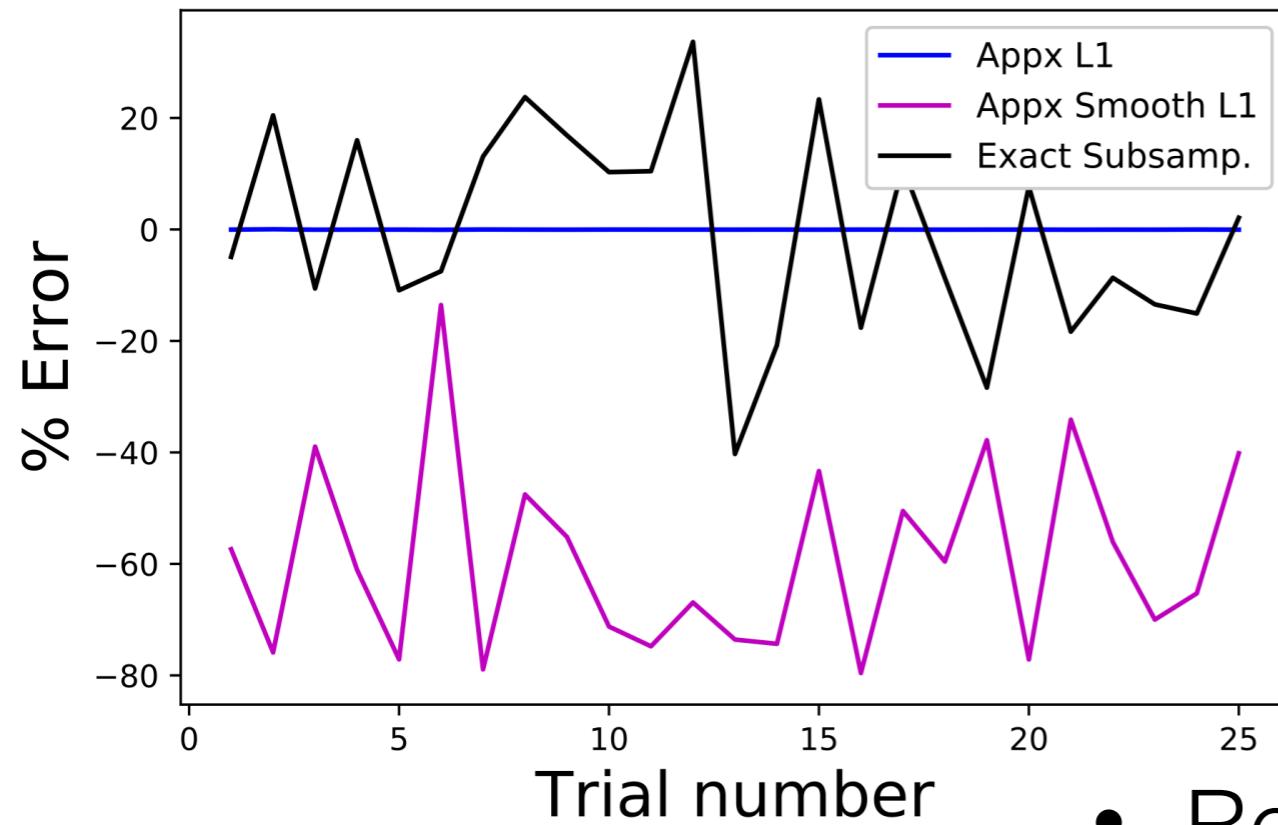


- Real data

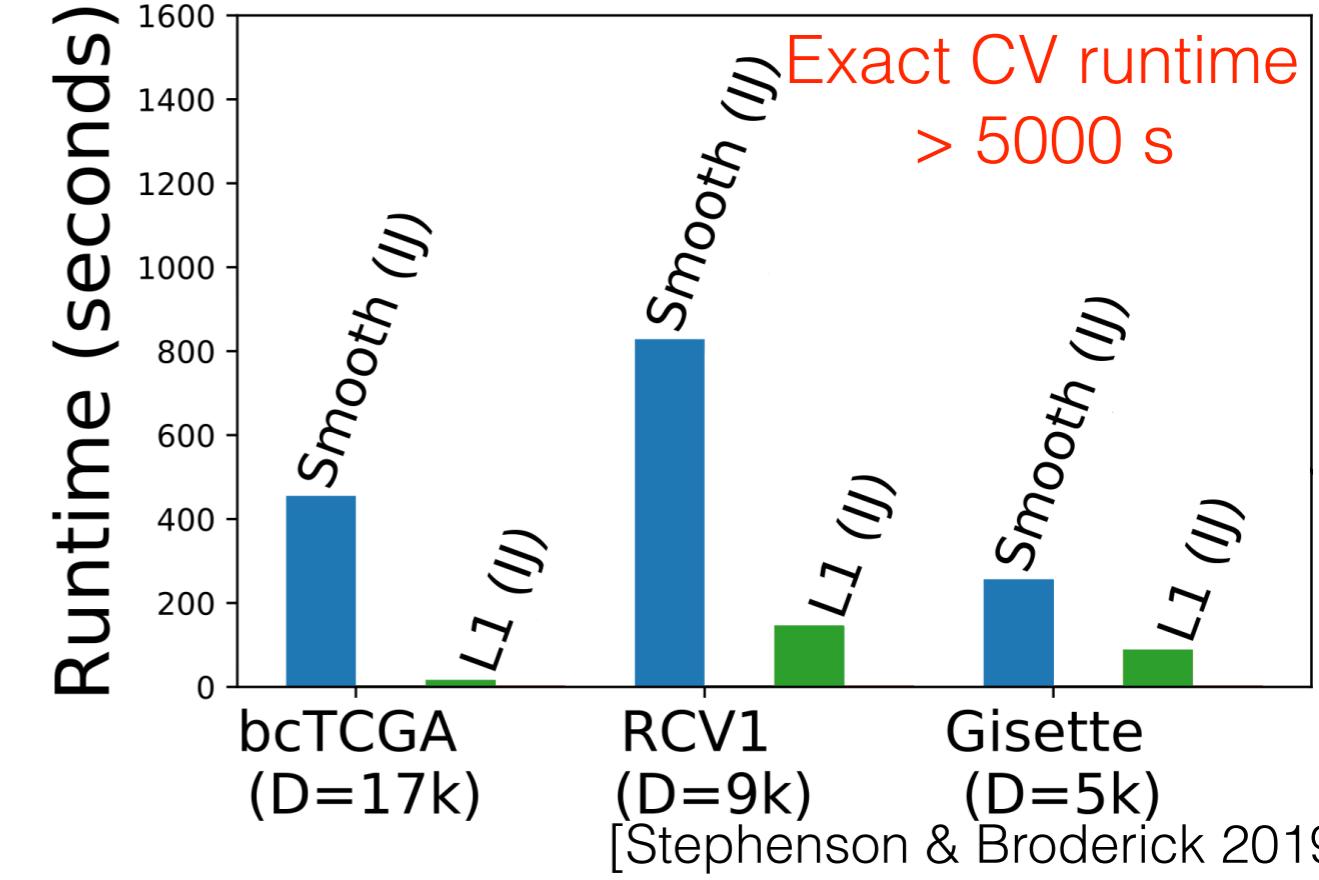
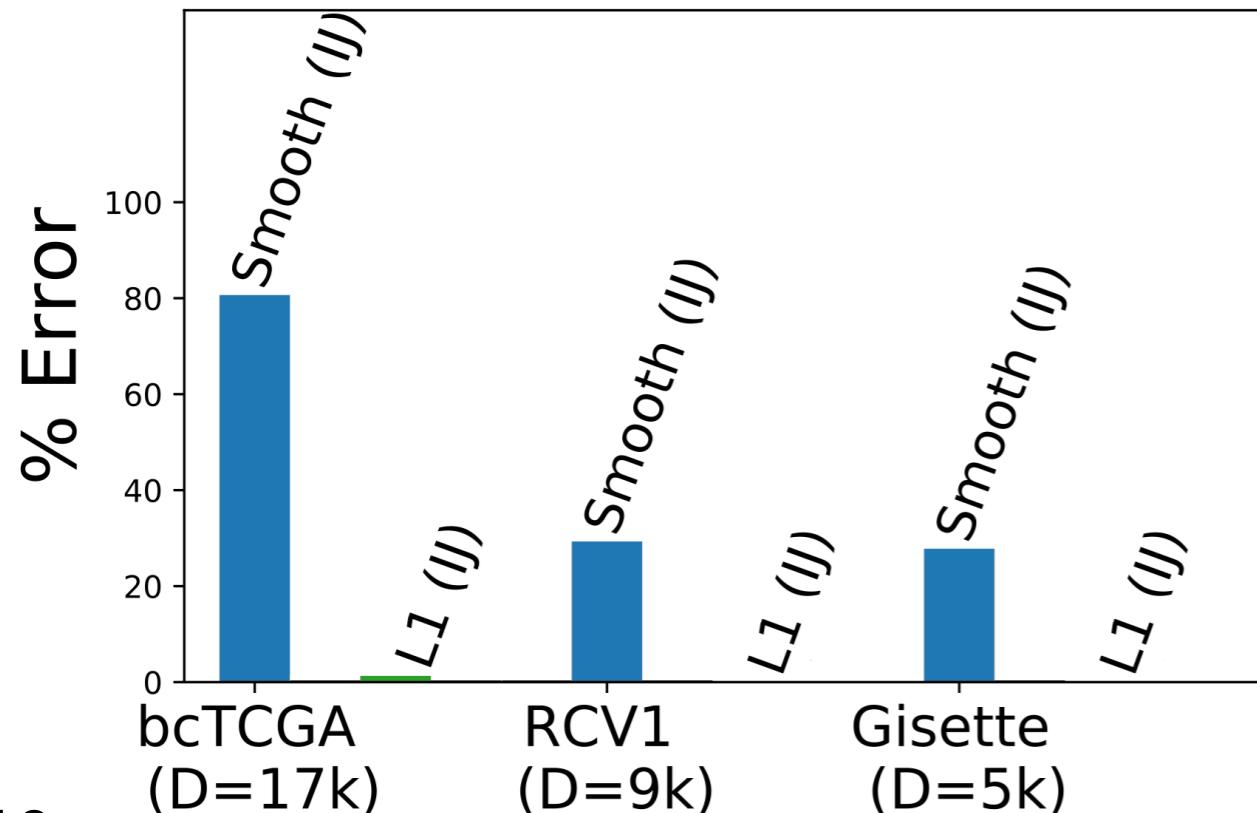


High dimensions: Experiments

- Simulated data ($N=500$, $D=40,000$)



- Real data



Conclusions

Conclusions

- Note: On a larger RCV1 data set, we estimate full CV takes > 2 weeks. Our method takes 3 minutes.

Conclusions

- Note: On a larger RCV1 data set, we estimate full CV takes > 2 weeks. Our method takes 3 minutes.
- Our CV approximations require *just 1 algorithm run*

Conclusions

- Note: On a larger RCV1 data set, we estimate full CV takes > 2 weeks. Our method takes 3 minutes.
- Our CV approximations require *just 1 algorithm run*
 - Taylor expansion + L1-regularization

Conclusions

- Note: On a larger RCV1 data set, we estimate full CV takes > 2 weeks. Our method takes 3 minutes.
- Our CV approximations require *just 1 algorithm run*
 - Taylor expansion + L1-regularization
 - We show: Fast, automatic, and high-quality

Conclusions

- Note: On a larger RCV1 data set, we estimate full CV takes > 2 weeks. Our method takes 3 minutes.
- Our CV approximations require *just 1 algorithm run*
 - Taylor expansion + L1-regularization
 - We show: Fast, automatic, and high-quality

**R Giordano, W Stephenson, R Liu, MI Jordan, and T Broderick.
A Swiss army infinitesimal jackknife. AISTATS 2019. (Notable
Paper Award winner)**

- <https://github.com/rgiordan/AISTATS2019SwissArmyIJ>
- <https://github.com/rgiordan/vittles>

**W Stephenson and T Broderick. Approximate cross validation
in high dimensions with guarantees. AISTATS 2020, to appear.
ArXiv:1905.13657.**

References

R Giordano, W Stephenson, R Liu, MI Jordan, and T Broderick.
A Swiss army infinitesimal jackknife. *AISTATS* 2019. (Notable Paper Award winner)

- <https://github.com/rgiordan/AISTATS2019SwissArmyIJ>
- <https://github.com/rgiordan/vittles>

W Stephenson and T Broderick. Approximate cross validation in high dimensions with guarantees. *AISTATS* 2020, to appear.
ArXiv:1905.13657.

References

- R Agrawal, T Campbell, JH Huggins, and T Broderick. Data-dependent compression of random features for large-scale kernel approximation. *AISTATS* 2019.
- T Campbell and T Broderick. Automated scalable Bayesian inference via Hilbert coresets. *Journal of Machine Learning Research*, 2019.
- T Campbell and T Broderick. Bayesian Coreset Construction via Greedy Iterative Geodesic Ascent. *ICML* 2018.
- JH Huggins, T Campbell, and T Broderick. Coresets for scalable Bayesian logistic regression. *NeurIPS* 2016.

**R Giordano, W Stephenson, R Liu, MI Jordan, and T Broderick.
A Swiss army infinitesimal jackknife. AISTATS 2019. (Notable
Paper Award winner)**

- <https://github.com/rgiordan/AISTATS2019SwissArmyIJ>
- <https://github.com/rgiordan/vittles>

**W Stephenson and T Broderick. Approximate cross validation
in high dimensions with guarantees. AISTATS 2020, to appear.
ArXiv:1905.13657.**